

UC Berkeley

UC Berkeley Electronic Theses and Dissertations

Title

Essays in Labor Economics

Permalink

<https://escholarship.org/uc/item/7cv7f50r>

Author

Scuderi, Benjamin Michael

Publication Date

2022

Peer reviewed|Thesis/dissertation

Essays in Labor Economics

by

Benjamin Scuderi

A dissertation submitted in partial satisfaction of the

requirements for the degree of

Doctor of Philosophy

in

Economics

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor Patrick Kline, Chair

Professor Christopher Walters

Professor Matthew Backus

Spring 2022

Essays in Labor Economics

Copyright 2022
by
Benjamin Scuderi

Abstract

Essays in Labor Economics

by

Benjamin Scuderi

Doctor of Philosophy in Economics

University of California, Berkeley

Professor Patrick Kline, Chair

This dissertation develops and applies econometric methods to understand key issues in labor economics. In particular, this dissertation investigates methods to transparently measure quantities of central importance to understanding equity and efficiency in labor markets, like the productivity of workers and firms, the nature of worker’s preferences and firm’s strategic behavior, and the size of wage markdowns.

The first chapter studies the nature and implications of firm wage-setting conduct on a large online job board for full-time U.S. tech workers. Utilizing granular data on the choice sets and decisions of firms and job seekers, I first develop and implement a novel estimator of worker preferences that accounts for both the vertical and horizontal differentiation of firms. The average worker is willing to pay 14% of their salary for a standard deviation increase in firm amenities. However, at the average firm, the standard deviation of valuations of that firm’s amenities across coworkers is also equivalent to 14% of their salaries, indicating that preferences are not well mdescribed by a single ranking of firms. Following the modern Industrial Organization literature, I use the labor supply estimates to compute the wage markdowns implied by a series of models of firm conduct that vary in the degree to which worker preference heterogeneity gives rise to market power. I then formulate a testing procedure that can discriminate between these models. Oligopsonistic models of wage setting are rejected in favor of monopsonistic models exhibiting near uniform markdowns of roughly 18%. Relative to a competitive benchmark, imperfect competition substantially exacerbates gender gaps in both wages and welfare. However, blinding employers to the gender of job candidates would have negligible effects on wage inequality.

The second chapter proposes a novel framework for conducting causal inference when researchers wish to compare a large number of treatments, as in studies of value-added that aim to quantify heterogeneity in skill, productivity, or preferences across workers, decisionmakers or service providers. Rather than apply parametric assumptions about the data-generating process, the framework I propose leverages only the common assumption that assignment of observations to treatments is unconfounded, and as such leads to “design-based” inferences of causal effects (in contrast to “model-based” approaches). I first illustrate identification of the causal effects of interest when the mechanism governing assignment – the propensity score –

is known. I then propose a method for estimating the features of the assignment mechanism when the true propensity scores are unknown and must be estimated. In settings with a large number of treatments (e.g. teachers, judges, or firms), the standard overlap assumption that all observations face a strictly positive probability of assignment to every treatment is likely to fail. I therefore propose a propensity score estimator that allows for structural failures of overlap, provide a computational guarantee for the estimation algorithm, and develop a finite-sample bound on the error of the estimator that holds with high probability. Finally, I provide an algorithm for using the estimated propensity scores to optimally trim the sample, such that a traditional notion of overlap is likely to hold on the resultant subsample and treatments can be reliably compared.

The third chapter applies the econometric framework of the second chapter to understand the distribution of productivity in a particular setting: legal defense for indigent individuals. This chapter quantifies the extent to which variation in case outcomes across indigent criminal defendants can be attributed to variation in the quality of their assigned counsel. Applying my estimation framework to data on case outcomes from three Texas counties that assign cases through conditionally randomized “wheel” systems, I find that attorney quality is highly variable. For defendants in felony cases, a one-standard-deviation decrease in attorney quality is associated with a 5.6 percentage-point increase in the probability of incarceration. These findings suggest that outcomes in criminal cases are driven in a nontrivial way by the luck of the draw, undermining the extent to which the criminal legal system can achieve traditional notions of fairness and efficiency. Using estimates of attorney quality, I evaluate the effects of a program that allowed defendants to choose attorneys. Perhaps because attorney quality is difficult to predict using observable characteristics, the program had essentially no effects on aggregate case outcomes, although it did significantly shift the burden of caseloads across attorneys.

In loving memory of my friend, Austin Hudson-Lapore.

Acknowledgments

I am deeply indebted to my advisor, Patrick Kline. As an advisor, Pat is astonishingly generous. He is not only generous temporally – he devotes a huge amount of time to his students – but also intellectually, by which I mean that he readily shares his energy for thinking through difficult problems with others. I am grateful to Pat for the tremendous amount he has taught me, but I am even more grateful for the fact that in being so generous intellectually, Pat instills a belief in others that they, too, are up to the challenge of tackling even the most difficult intellectual tasks. I am also indebted to my committee members, Christopher Walters and Matthew Backus. I have learned a great deal from both Chris and Matt, who have offered invaluable support, guidance, and encouragement along the way. I am thankful for the many Berkeley faculty who have given support and feedback at various points during my time in graduate school: David Card, Bryan Graham, Hilary Hoynes, Emmanuel Saez, Jesse Rothstein, Steven Raphael, and Jonah Gelbach. I am especially lucky to have been able to teach for, and get to know, David and Christy Romer, who are two of the kindest people I know.

I am also thankful for the wonderful friends and colleagues I have met at Berkeley. I feel extremely lucky to be able call Nina Roussille both a friend and co-author. Our weekly zoom check-ins kept me sane during the pandemic, and I am so grateful to be the beneficiary of her overwhelming brilliance and kindness. I cannot imagine what Berkeley would have been like without Arlen Guarin, Kaveh Danesh, Pedro Pires, and Priscila de Oliveira, each of whom made the the otherwise crushing load of graduate school bearable. Finally, I am thankful for Yotam Shem-Tov, who has been both a fantastic mentor and friend.

Beyond graduate school, I am grateful for the support of my family and friends. I thank Jimmy Biblarz, Nick Fandos, Amy Weiss-Meyer, and Michael Clegg for their humor, kindness, and unwavering encouragement. I thank my brother Louis Scuderi, and my parents Joan Drake and Louis Scuderi, for both believing in me and having the patience to put up with me through the ups and downs of the past six years.

Last, I am thankful for my wife and best friend, ImeIme Umana. You inspire me every day with your brilliance, courage, and deeply moral sense of justice. You make me laugh like no one else can. You push me to do better. Your love makes everything possible.

Contents

| | |
|--|------------|
| Contents | iii |
| List of Figures | v |
| List of Tables | vi |
| 1 Bidding for Talent: Equilibrium Wage Dispersion on a High-Wage Online Job Board | 1 |
| 1.1 Introduction | 1 |
| 1.2 Setting and Data | 6 |
| Market description | 6 |
| Sample restrictions: connected set | 8 |
| Descriptive Statistics | 9 |
| 1.3 Model | 10 |
| Setup | 10 |
| Labor Supply | 11 |
| Labor Demand | 13 |
| 1.4 Econometric Framework | 17 |
| Candidate Preferences | 17 |
| Labor Demand | 23 |
| Discriminating Between Non-Nested Models of Conduct | 28 |
| 1.5 Model Estimates | 33 |
| Labor Supply | 33 |
| Labor Demand | 39 |
| 1.6 Counterfactual Simulations of Bidding Behavior | 43 |
| Scenarios of interest | 43 |
| Computing new equilibria | 44 |
| Simulation Results | 44 |
| 1.7 Conclusion | 46 |
| 2 A Framework for Design-Based Inference of Many Treatment Effects | 66 |
| 2.1 Introduction | 66 |
| 2.2 Setup | 69 |

| | |
|---|------------|
| Variables and Notation | 69 |
| Potential Outcomes Model | 70 |
| 2.3 Estimation with a Known Assignment Mechanism | 72 |
| Consistency of the Fixed Effects Estimator | 72 |
| Estimating the Distribution of Treatment Effects | 73 |
| 2.4 Modeling the Assignment Process | 76 |
| Maximum Likelihood Estimation via the EM Algorithm | 77 |
| Estimation under Exact Rank Constraints | 78 |
| 2.5 A Nuclear Norm Regularized Estimator of Propensity Scores | 80 |
| Estimation Algorithm and Computational Guarantees | 81 |
| Statistical Guarantees | 85 |
| 2.6 An Algorithm for Sample Selection | 87 |
| 2.7 Conclusion | 89 |
| 3 Spinning the Wheel: Heterogeneity and Choice in the Provision of In- | |
| digent Defense | 90 |
| 3.1 Introduction | 90 |
| 3.2 Institutional Background | 92 |
| Indigent Defense in Bexar County | 92 |
| Indigent Defense in Comal County & the Client Choice Program | 93 |
| 3.3 Data and Summary Statistics | 94 |
| Data Sources | 94 |
| Summary Statistics | 95 |
| 3.4 Results | 96 |
| Variation in Treatment Effects Across Attorneys | 96 |
| The Full Distribution of Attorney Treatment Effects and Policy Simulations | 97 |
| Evaluating the Client Choice Program | 99 |
| 3.5 Conclusion | 101 |
| Bibliography | 118 |
| A Appendix to Chapter 1 | 130 |
| A.1 Additional Figures | 130 |
| A.2 Additional Tables | 135 |
| A.3 Illustration of conceptual framework | 137 |
| A.4 Details of EM algorithm | 138 |
| A.5 Properties of bidding strategies | 139 |
| A.6 Proof of the consistency of \hat{c}_j^m | 140 |

List of Figures

| | |
|--|-----|
| 1.1 Timeline of the Recruitment Process on Hired.com | 47 |
| 1.2 Distribution of Fraction of Interview Requests Accepted Across Firms | 48 |
| 1.3 Empirical Patterns in Bid and Ask Strategies | 49 |
| 1.4 Bids are Sticky in Expectation | 50 |
| 1.5 Interview Rejection Reasons as a Function of Firm Rankings | 51 |
| 1.6 Differentiation between and within firms | 52 |
| 1.7 First Stage | 53 |
| 1.8 Predicted Markdowns | 54 |
| 1.9 Visualizing the Vuong Test | 55 |
| 1.10 Relationship between Productivity and Amenity Values | 56 |
| 3.1 Flowchart of Case Assignment Process in Bexar County | 102 |
| 3.2 Maximum Likelihood Estimates of Distribution of Attorney Effects for Pleas | 103 |
| 3.3 Maximum Likelihood Estimates of Distribution of Attorney Effects for Incarceration | 104 |
| 3.4 Simulated Reduction in Incarceration Rate from a Layoff Policy | 105 |
| 3.5 Simulated Reduction in Incarceration Rate from a Retention Policy | 106 |
| 3.6 Time series of HHI of Attorney Representation | 107 |
| 3.7 Event Study: Effect on Distribution of Attorney Effects on Probability to Plea | 108 |
| 3.8 Event Study: Effect on Distribution of Attorney Effects on Probability of Incarceration | 109 |
| A.1 Mandatory features of a candidate profile, at the time of the study | 130 |
| A.2 Typical interview request message sent by a company to a candidate, at the time of the study | 131 |
| A.3 Model Fit: Labor Supply | 132 |
| A.4 Relationship between bids and systematic component of valuations, $\gamma_j(x_i)$ | 133 |
| A.5 Summary Statistics of Benefits listed by Firms | 134 |

List of Tables

| | |
|--|-----|
| 1.1 Summary Statistics for Candidate Characteristics | 57 |
| 1.2 Summary Statistics: Job Search and Job Finding | 58 |
| 1.3 Candidate Preference Model Goodness-of-Fit | 59 |
| 1.4 Which Firm Characteristics are Correlated with Amenity Values? | 60 |
| 1.5 Oaxaca-Blinder Decompositions of Components of Utility | 61 |
| 1.6 Non-Nested Model Comparison Tests | 62 |
| 1.7 (Subset of) Labor Demand Parameters Γ : $\log(\varepsilon_{ij}) = z'_j \Gamma x_i + \nu_{ij}$ | 63 |
| 1.8 Variance Decomposition of Bids | 64 |
| 1.9 Counterfactual Simulations | 65 |
| 3.1 Defendant Summary Statistics | 110 |
| 3.2 Attorney Summary Statistics | 111 |
| 3.3 Case Summary Statistics | 112 |
| 3.4 Validation Exercise | 113 |
| 3.5 Estimated Attorney Effect Variances | 114 |
| 3.6 Correlation Between Estimated Effects \hat{p}_j and Attorney Characteristics | 115 |
| 3.7 Comparison of Moments of Estimated $G(\cdot)$ s | 116 |
| 3.8 Estimates of Defendant Preferences over Attorney Characteristics | 117 |
| A.1 Comparison of data sources | 135 |
| A.2 Match productivity estimates: $\gamma_j(x_i) = z'_j \Gamma x_i$ | 136 |

Chapter 1

Bidding for Talent: Equilibrium Wage Dispersion on a High-Wage Online Job Board

This chapter is coauthored with Nina Roussille.

1.1 Introduction

How should economists interpret the empirical regularity that observably similar workers often receive markedly different wages across firms (Card et al. 2018)? A large literature has explored a variety of factors that can explain this heterogeneity: productivity (Abowd, Kramarz, and Margolis 1999; Gibbons et al. 2005; Faggio, Salvanes, and Reenen 2010; Dunne et al. 2004; Barth et al. 2016), compensating differentials (Rosen 1986; Hamermesh 1999; Pierce 2001; Mas and Pallais 2017a; Wiswall and Zafar 2018; Taber and Vejlin 2020; Sorkin 2018), and, more recently, imperfect competition (Manning 2011; Lamadon, Mogstad, and Setzler 2022; Berger, Herkenhoff, and Mongey 2017; Jarosch, Nimczik, and Sorkin 2021). Because most studies of the relative contributions of each of these factors use data on equilibrium matches, they generally rely on strong assumptions about the nature of the process by which workers and firms meet and by which wages are formed. For instance, a form of random matching is often assumed: given a set of equilibrium wages, workers have no control over the vacancies they are matched up with. An assumption of this kind is necessary when the menu of jobs workers choose from (their “choice set”) is not measured, but instead must be inferred. However, erroneous inference of these choice sets can introduce substantial bias (Barseghyan et al. 2021).

A particularly important assumption for any analysis of equilibrium wage dispersion regards the nature of firm wage-setting *conduct*: how firms determine which workers to hire, and how much to pay them. Despite the recent surge in interest in imperfect competition, little attention has been paid to testing which of the many possible models of conduct best describes firms’ observed behavior. Typically, existing analyses either propose a reduced-

form test of a particular imperfect-competition alternative relative to a perfect-competition null, or simply assume a single form of firm conduct. In practice, this means that prior studies make untested assumptions about key aspects of firm behavior, like whether firms behave strategically or the extent to which firms know workers' preferences. These assumptions then become key ingredients in the estimation of the size of markdowns and the distribution of welfare. Yet, different modes of conduct imply markedly different conclusions about the sources of wage dispersion and the extent of firms' market power. For example, models with strategic interactions predict more substantial markdowns at larger firms, implying that observed firm size-wage gradients are indicative of even steeper gradients in unobserved productivity. In contrast, models without strategic interactions need not imply differential markdowns by firm size, *ceteris paribus* (Boal and Ransom 1997). More broadly, erroneous assumptions about the form of conduct can lead to severely biased inferences about welfare and efficiency (Berger, Herkenhoff, and Mongey 2017).

This paper provides direct evidence about the nature of firms' wage-setting behavior by developing a testing procedure to adjudicate between non-nested models of conduct in the labor market. In particular, we focus on two sets of alternatives relevant to ongoing debates in the labor literature: first, whether firms compete strategically (Berger, Herkenhoff, and Mongey 2017; Jarosch, Nimczik, and Sorkin 2021), and second, whether firms tailor wage offers to workers' outside options (Caldwell and Harmon 2019; Flinn and Mullins 2021). We overcome the data limitations of previous studies by using detailed information from a large, high-stakes online job board on the choice sets and decisions of candidates and firms. On the platform, workers do not directly apply to jobs—rather, firms looking to fill vacancies submit “bids” on workers. Each bid must include an initial indication of the salary the firm is willing to pay (hereafter “the bid salary”), as well as a description of the job they are trying to fill, both of which may be individually tailored to each candidate. Because candidates can only enter the recruitment process at firms that bid on them, we are able to measure the full set of options they choose from. And, since the platform records whether candidates accept or reject firms' initial bids, we can cleanly infer candidates' revealed preferences over firms. Further, our data on bids reveal detailed variation in firms' willingness to pay for candidates that extends beyond just those the firm ultimately hires. These features of the data allow us to disentangle workers' selection into firms (labor supply) from firms' preference over workers (labor demand).

Armed with these data, our paper develops and implements a new framework for analyzing worker preferences over firms and the wage-setting conduct of those firms. In a first step, we propose a novel method for estimating the amenity values candidates associate with firms. Because we fully observe candidates' choice sets, we can cleanly infer a partial ordering of options for every candidate—our estimator ranks firms by aggregating those revealed preferences. The logic of our estimator is recursive, like that of Sorkin (2018), in that the estimated amenity value of any firm depends on the estimated amenity values of the firms it was revealed-preferred to: conditional on the bid salary, firms that offer good amenities will be revealed-preferred to other firms that offer good amenities. Importantly, our estimator flexibly models both the vertical differentiation (between-firm differences in amenity values common to all candidates) and horizontal differentiation (within-firm differences in amenity

values across candidates) of firms. In contrast to existing estimates of amenity values, we neither assume that all candidates share the same (mean) ranking of amenities, nor that candidates' (mean) rankings are a deterministic function of their demographics. Instead, we describe candidates' preferences as a mixture over types, each with a unique mean ranking of firms, where the distribution of types can depend upon candidate characteristics. Our estimator incorporates another unique feature of our data: candidates must publicly list the salary they wish to make at their next job (what we call the *ask* salary). To match reduced form evidence from both our setting and similar settings (e.g. Hall and Mueller 2018), we model preferences as reference-dependent: the labor supply function is kinked at the ask salary, which is analogous to an older tradition in IO where firms conjecture kinked product demand curves (Sweezy 1939; Bhaskar, Machin, and Reid 1991; Camerer et al. 1997; Farber 2015).

Next, we propose a general blueprint for analyzing labor demand that allows us to adjudicate between many non-nested models of firm wage-setting conduct. The fundamental intuition of our test is that if labor supply can be identified in a first step, applying an assumption about firm conduct immediately reveals implied equilibrium markdowns and therefore firms' valuations of candidates' labor (or, interchangeably, candidates' productivity) (see e.g. Berry and Haile 2014). Model-implied estimates of the valuations can then be used to test between modes of conduct via exclusion restrictions: instrumental variables that are excluded from the determinants of labor productivity should not be correlated with model-implied valuations. The logic of our procedure builds on the modern Industrial Organization literature studying product markets, beginning with Bresnahan (1987) and recently reviewed by Gandhi and Nevo (2021). Importantly, this empirical strategy avoids the endogeneity issues associated with relating variation in prices to variation in measures of market structure (like the Herfindahl–Hirschman Index) across markets, as in the “Structure-Conduct-Performance” (SCP) paradigm (Robinson 1933; Chamberlain and Robinson 1933; Bain 1951).

We translate this logic to the labor market setting: given our estimates of candidate preferences, we compute the wage markdowns implied by a set of non-nested models of firm wage-setting conduct. In order to adapt models of conduct to our data, we analogize the behavior of firms on the platform to that of bidders in a large online auction marketplace: just as in an auction market, firms compete against each other by bidding for workers' talent. We draw upon insights from the empirical auction literature (e.g. Guerre, Perrigne, and Vuong 2000; Backus and Lewis 2020) to define an equilibrium concept, establish the identification of markdowns, and propose a method for estimating those markdowns. To test between the various models of conduct, we implement the Vuong non-nested model comparison test (Vuong 1989; Rivers and Vuong 2002). The logic of the Vuong test is simple: when comparing two alternative models, the one that is closer to the truth should fit better. Following Berry and Haile (2014), Backus, Conlon, and Sinkinson (2021) and Duarte et al. (2021), we ensure that our test has power to discriminate between alternatives by using instruments that shift predicted markdowns but are excluded from productivity.

Our initial set of findings focuses on the labor supply. We document substantial vertical differentiation of firms on the platform: the average worker is willing to pay 14% of her

desired salary to enjoy a standard deviation increase in firm amenities. However, horizontal variation is just as important—the average standard deviation in valuations of amenities across coworkers at the same firm is also 14%. Our preferred estimates of labor supply describe preferences as a mixture over three types of workers. While preferences vary on a number of axes, the three groups can roughly be distinguished by preferences over firm size: some workers strongly prefer larger, more established firms, while others prefer smaller firms. Because the platform focuses on tech jobs, we loosely interpret these differences as differences in candidates’ risk tolerance. Finally, there is a residual gender gap in welfare, even conditional on the gender gap in bid salaries. This finding contrasts with other settings in which gender gaps in compensation have been shown to be driven in part by differences in preferences over working conditions (e.g. Bolotnyy and Emanuel 2022).

We then use those estimates to implement our procedure for comparing models of firm behavior. As a baseline, we are able to resoundingly reject the perfect competition model against all possible imperfect competition alternatives. However, in every version of our test, models that assume firms ignore strategic interactions in wage setting significantly outperform models that incorporate strategic interactions. This finding has significant implications for our conclusions about the size of wage markdowns—under the preferred model, we find markdowns of 18.2% on average, while models with strategic firms would have implied average markdowns of 25.8%. We also find evidence that firms do not actively tailor wage offers to candidates on the basis of predictable horizontal variation in preferences. In other words, our tests suggest that firms do not take advantage of predictable variation in firm-specific labor supply when making hiring decisions, which may lead to substantial misallocation in equilibrium. This finding is especially striking in the context of online labor markets which ostensibly seek to reduce information frictions in the search and matching process.

To quantify the impacts of imperfect competition on welfare, we use labor demand estimates from the preferred model to compute counterfactual equilibria under a range of conduct assumptions. Relative to a price-taking baseline, we find that firms make significantly more offers under the preferred model, but that the wages firms attach to those offers are lower. On net, this change leads to meaningful welfare losses. Relative to the preferred model, however, the average value of bids and the total number of bids are significantly lower in simulations of strategic firms, substantially decreasing overall welfare. We also find that the form of conduct has important implications for gender gaps: relative to men, women receive significantly fewer bids when firms predict horizontal preference variation than when they do not. Imperfect competition exacerbates gender gaps relative to the price-taking baseline. Finally, we find that blinding employers to the gender of candidates may lead to modest reductions in gender gaps.

This paper contributes to several strands of literature. First, our paper is most directly related to a growing literature that employs tools from industrial organization to study the role of firms in labor market inequality. Studies in this literature typically assume a single model of firm conduct, which they estimate using matched employer-employee data. Card et al. (2018) and Lamadon, Mogstad, and Setzler (2022) consider models in which firms are assumed to be monopsonistically competitive: that is, firms internalize upward-sloping labor supply, but do not act strategically. Berger, Herkenhoff, and Mongey (2017)

and Jarosch, Nimczik, and Sorkin (2021), on the other hand, write down models of non-atomistic firms that compete in local oligopolies. Our study departs from this prior work by explicitly formulating a testing procedure for discriminating *between* different modes of firm conduct, rather than assuming a single mode of conduct, more closely mirroring the industrial organization literature on estimating supply and demand and testing between models of conduct in product markets (Bresnahan 1989; Nevo 2001; Berry and Haile 2014, 2020; Backus, Conlon, and Sinkinson 2021; Gandhi and Nevo 2021). Second, because our data records not only equilibrium matches, but also the full set of offers made by firms to candidates (both accepted and rejected), we are able to separate the estimation of supply and demand. Finally, we focus on a single labor market in which it is likely that conduct of all firms is well-approximated by a single model, rather than applying our model to a national labor market defined by regional sub-markets. In this way, our study is related to a long tradition of single-industry studies in labor economics (Freeman 1976; Lipsky and Farber 1976; Staiger, Spetz, and Phibbs 2010; Goldin and Katz 2016).

Our paper more broadly contributes to a large literature exploring imperfect competition in labor markets (Boal and Ransom 1997; Bhaskar and To 1999; Bhaskar, Manning, and To 2002; Bhaskar and To 2003; Manning 2005, 2011). We adapt models of imperfect labor market competition to our setting, which combines the characteristics of online auction markets and terrestrial labor markets. In a similar context, Azar, Berry, and Marinescu (2019) gauge the potential market power of employers by estimating labor supply to individual firms on a large, online labor market using modern discrete choice methods. Our paper extends their analysis by characterizing both the nature of horizontal differentiation and the nature of firm conduct. A number of recent studies have examined the relationship between measures of market structure—typically, concentration measures like the Herfindahl–Hirschman Index (HHI)—and wages across markets in order to gauge the importance of imperfect competition (Azar et al. 2020; Schubert, Stansbury, and Taska 2021; Arnold 2021; Macaluso, Hershbein, and Yeh 2021). Since wages and market concentration are joint outcomes in models of labor markets, and finding excludable instruments for market structure is challenging (Berry 2021; Schmalensee 1989). In testing whether firms’ wage offers depend upon workers’ preference types, our study also relates to a line of research that connects heterogeneity in wages to outside options and the mode of wage determination (Hall and Krueger 2012; Caldwell and Harmon 2019; Lachowska et al. 2021).

Next, our paper relates to the literature on the estimation of non-wage amenities and their role in wage dispersion (Rosen 1986). Recent contributions in this area include Sorkin (2018) and Taber and Vejlin (2020) who use matched employer-employee data to identify search models that incorporate dispersion in non-wage amenities of firms. Because these studies use data on equilibrium matches, they infer amenity values from flows of workers across firms. By contrast, we observe the full set of options available to each worker on the platform, and therefore estimate amenity values by aggregating candidates’ revealed preferences over these options. In providing estimates of amenity values and exploring the relationship between those values and candidate characteristics, our paper also relates to a large literature on estimating heterogeneity in amenity values, e.g. Mas and Pallais (2017b) and Wiswall and Zafar (2018). In contrast to these studies, which are primarily carried out

in lab or experimental settings, we study the career decisions of workers in a high-stakes environment.

Finally, our paper contributes to strands of the literature in labor and industrial organization on the nature of competition on online markets. Using experiments, Dube et al. (2020) and Dube, Manning, and Naidu (2020) demonstrate the importance of monopsony in online labor markets for task work, and conclude that the presence of monopsony power in markets that are specifically designed to reduce search frictions suggests that imperfect competition may be pervasive in other “putatively thick” markets. Our paper more broadly relates to others describing the behavior of firms and workers in online labor markets. For instance, a recent study by Horton, Johari, and Kircher (2021) on the informative content of cheap talk about wages in online labor markets. We similarly find that cheap talk on Hired.com—in the form of firms’ initial offers and workers’ desired salaries—is an important signalling mechanism.

1.2 Setting and Data

Market description

As illustrated in Appendix Table A.1, a key limitation of the literature estimating revealed preferences from worker flows is that workers’ choice sets are rarely observed, and almost never available in a high-stakes, real-world environment. Because of this, existing estimates of worker preferences are either computed in surveys and lab environments (e.g., Wiswall and Zafar (2018), Mas and Pallais (2017b)), or reliant on strong assumptions applied to observational data. In survey or experimental settings, sample sizes and external validity to more traditional labor markets can be limited. In observational settings, estimates may be confounded by differences in choice sets or erroneous inference of workers’ options.

Two features of the recruitment process on Hired.com allow us to overcome this limitation. First, wage bargaining on Hired.com is high-stakes: the modal candidate on the platform is a software engineer in San Francisco looking for a full-time job with a salary of about \$120,000. Second, the recruitment process on Hired.com allows us to cleanly identify the choice set of candidates deciding which firms to interview with as well as the full set of observable profile characteristics firms have access to when deciding to send interview requests to a candidates. We explore these distinctive features below.

On the candidate side, Hired.com mostly serves candidates looking for full-time, high-wage engineering jobs based in the U.S. Table 1.1 shows that, on Hired.com, candidates are highly educated: 87.2 % of them have at least a bachelor’s degree and 40.3% have at least a master’s degree. Accordingly, the average salary offered by firms on the platform is high (\$114,505). Candidates on Hired.com are broadly comparable to those listed on other recruitment platforms for similar careers. For instance, the most common profile on Hired.com is a software engineer in San Francisco. As of April 2020, the average salary of candidates with this profile was \$119,488 on Glassdoor and \$132,000 on Paysa.¹ Hired’s

¹ Paysa is a personalized career service offering salary compensation and job matching for corporate employees. It is a useful reference for comparing employee salaries in the tech industry.

average salary for such profiles is \$129,783, which is between Glassdoor’s (lower bound) and Payscale’s (upper bound) salaries. The Hired.com sample also features profiles with different levels of seniority. For instance, among SF software engineers, 6% have 0-2 years of experience in software engineering, 22% have 2-4 years of experience, 22% have 4-6 years of experience, 33% have 6-10 years of experience, 8% have 10-15 years of experience, and 7% have more than 15 years of experience. This distribution is similar to the one reported by Payscale for this combination of job and location.² On the firm side, companies hiring on the platform are representative of the tech ecosystem: a mix of early stage firms, more mature start-ups (e.g. Front, Agolia), and larger, more established firms (e.g. Zillow, Toyota). With more than 13,000 candidates and jobs in our analysis sample, the market we study should be thought of as a large, high stakes job board for well-qualified candidates.

Our ability to cleanly identify the choice sets of candidates deciding which firms to interview with emerges from the unique chronology of hiring on the platform. On a traditional job board, firms post a job description and then candidates apply to each posted job separately. By contrast, on Hired.com, companies apply to candidates based on their profiles, and candidates decide whether or not to interview with companies based on the job descriptions and bid salaries they receive. Importantly, candidates have no way to directly view and apply to job postings without receiving an interview request. As a result, for each candidate on Hired.com, we know their consideration set (the set of all the firms that apply to them), and their choices (whether or not they decided to interview with any given firm in the consideration set).

Formally, the recruitment process can be divided into the following three sequential steps, also described in Figure 1.1:

Supply side: Candidates create a profile that contains standardized resume entries (education, past experience, etc.) and, crucially, the salary that the candidate would prefer to make. We call this the *ask salary*. Appendix Figure A.1 is a screenshot of a typical candidate’s profile. In short, every profile includes the current and desired location(s) of the candidate, their desired job title (software engineering, web design, product management, etc.), their experience (in years) in this job, their top skills (mostly coding languages such as R or Python), their education (degree and institution), their work history (i.e., firms they worked at), their contract preferences (remote or on-site, contract or full-time, and visa requirements), as well as their search status, which describes whether the candidate is ready to interview and actively searching or simply exploring new opportunities. Importantly, the ask salary is prominently featured on all profiles since it is a required field.

Demand side: Firms get access to candidate profiles that match standard requirements for the job they want to fill (i.e., job title, experience, and location). To apply for an interview with a candidate, the company sends them a message—the *interview request*—that typically contains a basic description of the job as well as, crucially, the salary at which they would be willing to hire the candidate. We call this the *bid salary*. Appendix Figure A.2 is a screen-

² Payscale’s page for SF software engineer profiles can be found [here](#).

shot of a typical message sent to a candidate by a company. The bid salary is prominently featured in the subject line of the message and is required to be able to send the message. The equity field also exists but is optional.

Demand meets supply: Hired.com records whether the candidate accepts or rejects the interview request. While interviews are conducted outside of the platform, Hired.com gathers information on whether the company makes a final offer of employment to the candidate and at what salary. We refer to this as the *final salary*. It is important to note that the bid salary is non-binding, so the final salary can differ from the bid. Finally, we observe whether the candidate accepts the final salary offer, in which case the candidate is hired. Given these three steps of the recruitment process and the nature of candidates and jobs on the platform, our setting combines a high stakes environment with clean identification of the consideration set of each candidate and their decisions at the interview stage. One a priori caveat is that, while the consideration set is comprehensive—that is, we observe all the firms that the candidate considers on the platform—it is not exogenous, as firms select into sending an interview request to candidates. However, the fact that we observe all information about candidates available to firms at the time they decide to send an interview request allows us to circumvent this issue.³

Sample restrictions: connected set

As we explain below, we can only estimate amenity values for firms that are members of a connected set. To be a member of this set, a firm must have been both revealed-preferred to at least one member of the set, and have been revealed-dispreferred to at least one member of the set. While several job titles and locations are represented on Hired.com, the candidate market is highly skewed towards software engineers in San Francisco: 60.1% of the candidates are software engineers and 31.1% live in the Bay Area. In addition, the jobs on the platform are even more concentrated in these profiles: 76% of interview requests go to software engineers in the Bay Area. Therefore, while the average number of interview requests on the platform is 4.5, the average number of interview requests received by a software engineer in the Bay Area is 11.2. For these reasons, we zoom in on the highly connected market of San Francisco software engineers. Table 1.2 provides simple descriptive statistics on the sample sizes, for the full dataset, for the subset of jobs in the San Francisco Bay Area and finally for the connected set of firms within that market. The full sample includes 7,877 companies that sent 856,665 requests for 64,539 different jobs to 224,499 candidates. While the average number of bids sent per job is 13.3, the median is 5.0, suggesting large differences in the extent to which companies reach out to candidates. More than a fourth (n=16,907) of all jobs on Hired.com in the full sample are based in the SF Bay area. For these jobs, 2,121 companies sent out 267,940 interview requests to 44,321 candidates, averaging 15.8 bids per job (median 5 bids) and 4.1 bids per candidate. The average probability of accepting a bid remains almost constant between 60% and 62.5 % in both sets. 1,649 companies meet the

³ Assumption 1.1 in Section 1.4 formalises this argument.

requirements to qualify for the connected set. Companies in this sample are more targeted when approaching candidates, sending on average only 9.5 bids (median 4 bids) for 13,072 different jobs to 14,344 candidates. However, the average number of bids per person is with 4.8 around 37% higher than in the full sample and candidates accept only 56.4% of received interview requests.

Descriptive Statistics

As noted above, we can only estimate the amenity values of firms that have both been accepted and rejected by at least one candidate. This implies that candidates must necessarily incur an interview cost, such that they would not accept all the interview requests they receive. Figure 1.2 empirically tests this assumption by displaying the distribution of the share of bids accepted for a given firm. It first shows that firms are frequently rejected by candidates: on average, candidates only accept 60.5% of the interview requests they receive. In addition, there is significant heterogeneity across companies in the likelihood that an interview request is accepted: while the mean share of bids accepted is 60.5%, 10.2% of the firms see less than 40% of their interview requests accepted, while 16.2% of the firms see more than 75% of their interview requests accepted.

Figure 1.3 further illustrates several empirical patterns that are the foundation of our modelling strategy. Figure 1.3a plots the probability of acceptance of an interview request against the ratio of the bid to ask salary. The first fact is that higher bids are associated with a higher acceptance probability: when the bid salary matches the ask salary, the acceptance probability is 62%. When the ratio is 1.2 or more, the acceptance probability goes to 73%, whereas when it is 0.8 or less it averages 36%. The second notable pattern is that there is a clear discontinuity of the probability of acceptance in the neighborhood of $\frac{bid}{ask} = 1$. In particular, while the probability of acceptance is 52% when $\frac{bid}{ask} = 0.95$, it jumps to 62% when the ratio is 1.⁴ Figure 1.3b shows the relationship between the probability that the bid is, respectively, less than, equal to, or greater than the ask, and the level of the ask salary. First, across all levels of ask salary, the probability that the bid is exactly equal to the ask is very high, averaging 76.5%. A second, intuitive, observation is that the probability that the ask is lower than the bid increases with the level of the ask from virtually 0% at the lowest levels of ask salary to just shy of 40% for the highest levels of ask salary. Symmetrically, the probability that the bid is greater than the ask decreases from around 20% to 0%. This empirical pattern provides strong suggestive evidence that the asked wage serves as a behavioral reference point in the formation of the bid salary. Figure 1.3c shows the relationship between the bid premium - the difference between bid and ask salaries - and the within-job deviation of the log salary. This figure illustrates the fact that there is large heterogeneity of bid salaries for the same job. Indeed, if the data were on the -45 red line, firms' bids for the same job would remain constant, independent of the candidates' ask salaries. Empirically, we observe that the slope of the relationship is dramatically flatter

⁴ Leveraging a survey of 6,000 job seekers in New Jersey, Figure 3 in Hall and Mueller (2018) shows the job offer acceptance frequency as a function of the difference between the log hourly offered wage and the log hourly reservation wage. A clear kink is observed at offered wage = reservation wage.

than this “full compression” line: changes in the ask are almost entirely offset by changes in the bid - indicating that, even for a given job, firms increase their bids almost one for one with the asks. In fact, only 1.4% of jobs offer the same bid salary to all candidates, and the within-job variation in salaries is substantial: the average standard deviation of offers for a given job is \$23,041.

The bid salary is what firms declare they are willing to pay the candidate solely based on their profile, before any interaction with them. The final salary is offered to a candidate at the hiring stage. Given that companies are by no means contractually bound by their bids, final salaries may differ from bids. Given our focus in this paper on the interview stage of the process, it is important to point out that firms effectively commit to making final offers that are close to the bids. Figure 1.4 shows the relationship between the bid and final offer for the subset of candidates that receive one. Strikingly, this relationship is very linear, with a slope close to one. Additionally, 31% of all final offers are identical to the bid and 72% of all final offers are within 10% of the bid.

1.3 Model

Setup

This section describes our model of the recruitment process on the platform. We index candidates by $i = 1, \dots, N$ and firms by $j = 1, \dots, J$. Firms encounter a candidate pool, I_j , the size and composition of which varies depending on the time period of the firm’s search. Likewise, candidates encounter a time-specific firm pool J_i .⁵ We denote the observable characteristics of firms by z_j (which includes a constant), and let $j = 0$ denote an outside option. Candidates post resume information x_i , which includes their asked salary a_i (and a constant), before interacting with firms on the platform. Firms browse active candidate profiles and decide whether to send each candidate an interview request, and if so, how much to bid. As stated above, firms’ bids are made before the firm has had any interaction with the candidate, on the basis of the observable candidate characteristics x_i alone. We denote the bid of firm j on candidate i by b_{ij} , and let the indicator variable B_{ij} equal one if firm j sends a bid to candidate i . After a candidate receives an interview request, she decides whether to accept and thereby move forward with the recruitment process, or to reject the offer. We let the indicator variable D_{ij} equal one if candidate i accepts firm j ’s interview request. After the interview process is complete, the firm can make a final offer of employment to the candidate. We let B_{ij}^f equal one if j makes a final offer to i , and we denote the salary attached with that final offer by b_{ij}^f . Finally, we let D_{ij}^f equal one if i accepts j ’s final offer of employment.

Our analysis focuses on the initial stages of the recruitment process. In order to specify a tractable model of firm and candidate behavior at the initial stages, we make several simplifying assumptions about the later stages of the process. In particular, we assume firms are risk neutral, and that firms do not treat bids as cheap talk – rather, we assume that firms

⁵ We assume that agents’ beliefs are stationary, such that they behave as if they are in a steady state, as in Backus and Lewis (2020). We defer consideration of dynamics for future research.

credibly expect to pay their bids, should they decide to make a final offer. In practice, this assumption is an accurate description of firm behavior: the correlation between initial bids b_{ij} and final offers b_{ij}^f is 0.86 (see Figure 1.4). Second, we assume that candidates' choices at the interview request and final offer stages are governed by the same basic preference structure. While our framework is consistent with certain forms of preference updating on the part of candidates after interviews take place, we remain agnostic about those mechanisms here. These assumptions allow us to model the bid determination process straightforwardly: when a firm encounters a candidate, the firm decides to bid on that candidate by maximizing the ex-ante option value associated with an interview request. The option value is determined by the firm's forecast of the candidate's marginal revenue product, net of the bid, and the probability that the candidate would accept a final offer of employment, given the bid.

Labor Supply

The first component of our model is a labor supply system. In our model, candidates' asked wages a_i play two important roles. First, motivated by the visual evidence in Figure 1.3, we assume that the asked wage acts as a behavioral reference point: the elasticity of labor supply may be relatively larger when firms offer less than the asked wage than when they offer more than the asked wage. This feature is a potential mechanism driving the bunching of *offered* wages at exactly the asked wage, even conditional on detailed candidate-specific controls. Second, we assume that the asked wage serves as a sufficient statistic for the monetary component of utility associated with candidates' outside options, up to an additive constant. For the large fraction of workers on the platform engaging in on-the-job search, this assumption can easily be justified if candidates formulate asked wages as a function of their current wage. Workers searching from unemployment post lower asked wages even conditional on a rich set of covariates (conditional on other profile characteristics, employed candidates ask for \$8,366 more than unemployed candidates), suggesting that asked wages of unemployed candidates indeed reflect the relatively worse outside options available to those workers. We therefore normalize the "bid" associated with the outside option as $b_{i0} = a_i$.

We model the utility candidate i associates with option j at bid b_{ij} as additively separable:

$$V_{ij} = u(b_{ij}, a_i) + \Xi_{ij},$$

where the function $u(b_{ij}, a_i)$ is the monetary component of utility and Ξ_{ij} is the non-monetary component of utility that candidate i associates with option j . Because only relative utilities matter for choices, we normalize $u(a, a) = 0$ without loss of generality. The utility of the outside option is therefore given by:

$$V_{i0} = \Xi_{i0}.$$

We assume that $u(b, a)$ is continuous, strictly increasing, and twice continuously differentiable in its first argument, except at the point $b = a$, where $\lim_{b \rightarrow a^-} \partial u(b, a) / \partial b > \lim_{b \rightarrow a^+} \partial u(b, a) / \partial b$. This assumption encodes reference-dependence around the asked wage: utility decreases relatively more quickly for every dollar below the asked wage than it increases for every dollar above the asked wage.

The non-monetary component of utility can be further decomposed into a systematic *amenity value* and an idiosyncratic *taste shock*:

$$\Xi_{ij} = A_{ij} + \xi_{ij}.$$

We assume that the idiosyncratic preference shocks ξ_{ij} are independent and identically-distributed draws from a probability distribution, $\xi_{ij} \stackrel{iid}{\sim} F_\xi(\cdot)$, where F_ξ admits a continuous, log-concave density $f_\xi(\cdot)$ with support on the full real line.⁶ Preference shocks ξ_{ij} are private information: they are observed by workers, but not by firms. Further, the distribution of preference shocks is independent of x_i : $F_{\xi|x} = F_\xi$.

The amenity value candidate i associates with option j is determined by i 's *latent preference type*, which we denote by Q_i :

$$A_{ij} = A_j(Q_i).$$

Candidates i and ℓ with $Q_i = Q_\ell$ share a common mean valuation of amenities at all firms. We assume that candidates' preference types are not directly observable by recruiters, but that the distribution of preference types F_Q may depend non-trivially on candidates' observable resume characteristics x_i : $F_{Q|x} \neq F_Q$. In this sense, A_{ij} is not purely the private information of the candidate, but instead may be forecast by firms on the basis of the observables available on candidate profiles.

We assume that a candidate accepts an interview request if and only if the utility associated with that requests exceeds that of her outside option:

$$D_{ij} = B_{ij} \times \mathbf{1}[V_{ij} \geq V_{i0}].$$

Likewise, let V_{ij}^f denote the utility level i associated with a final offer of b_{ij}^f from j . Candidates pick the top choice among all final offers, such that:

$$D_{ij}^f = \mathbf{1} \left[V_{ij}^f \geq V_{ik}^f \quad \forall k \quad \text{s.t.} \quad B_{ij}^f = 1 \right].$$

For simplicity's sake, we model the utility candidates associate with final offers as $V_{ij}^f = u(b_{ij}^f, a_i) + \Xi_{ij}$, such that the same utility shocks that enter into candidates' interview offer decisions also govern candidates' final job choice. Because we focus mainly on the ex-ante perspective of firms formulating bids, we view this assumption as a simplifying abstraction that may be relaxed in future work.

⁶ A function f_ξ is log-concave if:

$$f_\xi(\lambda y + (1 - \lambda)x) \geq f_\xi(y)^\lambda f_\xi(x)^{1-\lambda} \quad \forall x, y \in \mathbb{R}, \lambda \in [0, 1].$$

A large number of common probability distributions admit log-concave densities, including but not limited to the normal, logistic, extreme value, and Laplace distributions. Log-concave probability distributions are commonly used in models of search (Bagnoli and Bergstrom 2005), and possess a number of desirable qualities. Among other things, log-concavity of f_ξ implies that F_ξ and $1 - F_\xi = \bar{F}_\xi$ are also log-concave, that f_ξ/F_ξ is monotone decreasing, and that f_ξ/\bar{F}_ξ is monotone increasing.

Labor Demand

A General Bidding Framework

We next write down a general framework for rationalizing firms' bidding behavior. Firms are risk neutral and equally well informed. Firms do not observe candidates' latent types Q_i , but rather can form predictions over those types using the available candidate characteristics x_i . For each candidate i it encounters, firm j formulates an optimal bid b_{ij}^* to maximize the expected option value of making an interview request given that candidate's observables. This is given by maximizing an *expected option value* function $\pi_{ij}(b)$:

$$b_{ij}^* = \arg \max_b \pi_{ij}(b).$$

Firms decide to bid on candidates if the maximized value of the expected option value function surpasses an interview cost threshold c_j :

$$B_{ij} = \mathbf{1} \left[\pi_{ij}(b_{ij}^*) \geq c_j \right].$$

We may therefore write realized bids as:

$$b_{ij} = B_{ij} \times b_{ij}^*.$$

We use the shorthand $b_{ij} = 0$ to indicate the event $B_{ij} = 0$.

The option value of an interview request to a particular candidate depends upon both her labor supply decision and her productivity. Define the potential outcome:

$$D_{ij}^\circ(b) \triangleq \mathbf{1} \left[i \text{ would accept } j\text{'s offer of employment} \mid b_{ij} = b \right],$$

which encodes candidate i 's final labor supply decision, given the firm's choice of bid b . We refer to $\pi_{ij}(b)$ as an expected *option value* function because even if the event $D_{ij}^\circ(b) = 1$ is realized, the firm may choose not to hire i (for instance, if a candidate the firm prefers over i would also accept its offer). Denote the ex-post productivity of a match between candidate i and firm j as ε_{ij}° . Given these definitions, the expected option value/profit function can then be written:

$$\pi_{ij}(b) = \mathbb{E}_{ij} \left[D_{ij}^\circ(b_{ij}) \times (\varepsilon_{ij}^\circ - b_{ij}) \mid b_{ij} = b \right]$$

where \mathbb{E}_{ij} denotes expectation taken over the information set of firm j when it evaluates candidate i , and so implicitly conditions on firm, candidate, and market-level variables. The connection between this representation of the firm's problem and the objective function of a bidder in a standard first-price auction is immediate: indeed, the problems are nearly identical. In a first-price auction, a bidder's objective is simply to maximize her expected utility, where her bid affects both the net payoff should she win ($\varepsilon_{ij}^\circ - b$) and the probability that she wins the auction (the distribution of $D_{ij}^\circ(b)$). In a standard auction, the win probability depends only upon the monetary values of the competing bids – the bidder who submits the highest bid wins. In our setting, horizontal differentiation weakens this relation:

the firm that submits the highest monetary bid is not guaranteed to be the candidate's top-ranked choice.

Conditional on the firm's information set, we assume that potential outcomes $D_{ij}^o(b)$ and ex-post marginal revenue products ε_{ij}^o are independent. Further, conditional on the information known to the firm at the time it bids, ε_{ij}^* is independent of the firm's choice of bid b_{ij} . The first of these assumptions rules out, among other things, scenarios in which the event of winning the "auction" for candidate i reveals information about other firms' productivity forecasts that is relevant to j 's forecast (sometimes called the "winner's curse"). Since all firms must bid on candidates before productivity is revealed, this assumption essentially establishes the sufficiency of the observables available to the firm for forecasting productivity. The second assumption rules out behavioral effects of increasing compensation (e.g. efficiency wages). Together, they imply:

$$\pi_{ij}(b) = \Pr_{ij}(D_{ij}^*(b) = 1) \times (\mathbb{E}_{ij}[\varepsilon_{ij}^*] - b).$$

The first term in the above expression is j 's forecast of i 's labor supply decision, which we denote by:

$$\Pr_{ij}(D_{ij}^*(b) = 1) \triangleq G_{ij}(b).$$

Firms' forecasts of ex-post productivity, which we denote by ε_{ij} , are functions of a systematic component (determined by candidate covariates) and an idiosyncratic component:

$$\mathbb{E}_{ij}[\varepsilon_{ij}^*] \triangleq \varepsilon_{ij} = \gamma_j(x_i, \nu_{ij}).$$

We further assume $\nu_{ij} \stackrel{iid}{\sim} F_\nu(\cdot)$, and that ν_{ij} is independent of x_i , z_j , and market-level variables. The function $\gamma_j(x, \cdot)$ encodes the systematic component of productivity shared by all candidates with observables $x_i = x$ at firm j . We impose the normalization $\mathbb{E}[\nu_{ij}] = 0$ without loss of generality. Substituting these definitions into the expected option value function gives:

$$\pi_{ij}(b) = G_{ij}(b) \times (\varepsilon_{ij} - b) = G_{ij}(b) \times (\gamma_j(x_i, \nu_{ij}) - b).$$

Given the parallels between our setting and the auction setting, we refer to ε_{ij} as either j 's *valuation* for i or i 's (ex-ante) productivity at j , and $G_{ij}(b)$ as either j 's win probability for i or i 's labor supply to j . Firms' strategies are described by an optimal bidding function that maps valuations into actions:

$$b_{ij}(\varepsilon) = \begin{cases} \arg \max_b G_{ij}(b) \times (\varepsilon - b) & \text{if } \max_b G_{ij}(b) \times (\varepsilon - b) \geq c_j \\ 0 & \text{otherwise.} \end{cases}$$

To close the model, we define a notion of equilibrium. In a standard Bayes-Nash equilibrium, players' actions are best responses given their beliefs, which are themselves consistent with equilibrium play. In the subsequent analysis, we test models of firm behavior in which

firms' forecasts of candidates' labor supply decisions may not fully incorporate the relevant available information. In order to accommodate these models, we modify the standard definition of equilibrium as follows. Denote the maximum utility level offered to i by V_i^1 , and let Λ_i be a random variable that governs the distribution of V_i^1 . We assume that beliefs are consistent *conditional on the information firms use to construct those beliefs*. In particular, let $\Omega_{ij} = \{\omega_{ij}^\Lambda, \omega_{ij}^Q\}$ encode the information j uses to forecast Λ_i and Q_i , respectively, and let $F_{\Lambda,Q}(\lambda, q \mid \Omega)$ denote the population joint CDF of Λ_i and Q_i , conditional on Ω_{ij} . We may now define equilibrium as follows:

Definition 1.1 (Equilibrium). *Conditional on an information structure $\{\Omega_{ij}\}_{i=1, j=1}^{N,J}$, a pure strategy equilibrium is a set of tuples $\{b_{ij}(\cdot), G_{ij}(\cdot)\}_{i=1, j=1}^{N,J}$ such that:*

(Optimality) $b_{ij}(\varepsilon)$ is j 's best response for valuation ε given beliefs $G_{ij}(b)$.

(Consistency) Conditional on the information Ω_{ij} , firm j 's beliefs obey:

$$G_{ij}(b) = \iint \Pr(V_{ij} = V_i^1 \mid \Lambda_i = \lambda, Q_i = q, b_{ij} = b) \times dF_{\Lambda,Q}(\lambda, q \mid \Omega_{ij}).$$

In the classic first-price auction setting, the function $G_{ij}(b)$ is nonparametrically identified by the observed distribution of bids: the seller accepts the highest bid, and so (under the assumption that bidders have rational expectations) an estimate of $G_{ij}(b)$ can be constructed by calculating the empirical CDF of winning bids. This argument is the basic intuition of the approach of 2000 (GPV). In our setting, the win probability $G_{ij}(b)$ depends not only upon the monetary value of the bid a firm submits, but also the non-monetary components $A_j(Q_i) + \xi_{ij}$. Despite this difference, we adopt the basic logic of GPV in our estimation strategy, which we detail below: given estimates of the labor supply parameters and the assumption of rational expectations, the empirical distribution of inclusive values for each candidate can be used in combination with an assumption on firm conduct (where various models of conduct are indexed by m) to construct estimates of $G_{ij}^m(b)$ – the conditional win probability under model m .

Defining Firm Conduct

Given the framework of the previous section, we next consider various modes of firm conduct. We operationalize our notion of conduct in this setting as sets of assumptions on the information firms use to forecast candidates' labor supply decisions. In practice, that means specifying which variables are included in the components of Ω_{ij} . This notion of conduct is not the only interesting feature of firm behavior in wage setting, and indeed there are many potentially interesting questions about the ways firms behave in labor markets that we do not test. However, our setting – one in which firms have the ability to offer fully individualized wages to each candidate – is particularly well-suited for thinking about how firms incorporate information about the distribution of preferences into their recruitment decisions. In Appendix A.3, we illustrate the implications of our conduct assumptions, and how the conceptual framework of our study differs from those that relate measures of

market structure to wages, via a simplified model similar to that of Bhaskar, Manning, and To (2002).

We first consider a model of “perfect competition” in which firms are assumed to bid their valuations: $b_{ij}(\varepsilon) = \varepsilon$. In this model, interview costs c_j are normalized to 0 without loss of generality. This model does not fit cleanly into the framework of the previous section – to rationalize bidding at exactly its valuation, a firm must believe that there always exists a competitor with a valuation arbitrarily close to its own valuation. Even so, the perfect competition model we estimate serves as a useful baseline against which we can compare more complicated models of conduct that incorporate additional sources of wage dispersion beyond differences in the marginal revenue product of labor (MRPL).

In order to specify additional conduct assumptions of interest, we decompose the joint CDF of Λ_i and Q_i given Ω_{ij} as:

$$F_{\Lambda,Q}(\lambda, q \mid \Omega_{ij}) = F_{\Lambda|Q}(\lambda \mid Q_i = q, \omega_{ij}^\Lambda) \times F_Q(q \mid \omega_{ij}^Q).$$

The first conduct assumption we test concerns the information firms use to forecast types. We specify two alternatives – firms are assumed to be either:

- **Type Predictive:** $\omega_{ij}^Q = x_i$, such that $F_Q(q \mid \omega_{ij}^Q) = F_{Q|x}(q \mid x_i)$, or
- **Not Predictive:** ω_{ij}^Q is empty, such that $F_Q(q \mid \omega_{ij}^Q) = F_Q(q)$.

This assumption governs how firms internalize horizontal differentiation: do firms engage in what is sometimes called *direct segmentation*? Our model allows for the possibility that workers who have the same level of productivity at a particular firm may belong to different preference types. Variation in preference types can itself be partially predicted by candidate characteristics, raising the possibility that type-predictive firms might offer different wages to candidates with identical productivity levels. Non-predictive conduct implies that firms make fewer offers than under an efficient allocation, although workers may capture a larger share of the surplus. Type-predictive conduct implies less misallocation, but potentially at the cost of workers’ share of the surplus. How firms do or do not use information has been a matter of debate in the labor literature. For instance, Burdett and Mortensen (1998) assume that firms are not type-predictive, leading to efficiency losses that they show can be reduced by the introduction of a minimum wage. On the other hand, Postel-Vinay and Robin (2002) assume that firms are not just type-predictive, but fully informed about the types of workers they meet, allowing them to engage in classic first-degree price discrimination. More recently, Postel-Vinay and Robin (2004) and Flinn and Mullins (2021) analyze models in which firms differ in whether they commit to posted wages (akin to non-predictive conduct) or negotiate wages in response to outside offers (akin to type-predictive conduct). Similarly, whether firms use information on within-firm variation in price elasticities has been the subject of interest in the industrial organisation literature on uniform pricing (DellaVigna and Gentzkow 2019).

The second conduct assumption we test concerns the nature of interactions between vertically-differentiated firms. Again, we specify two alternatives – firms are assumed to

be either:

- **Monopsonistically Competitive:** ω_{ij}^Λ omits j 's bid as a (direct) determinant of Λ_i ,
or
- **Oligopsonists:** ω_{ij}^Λ includes j 's bid as a (direct) determinant of Λ_i .

In a monopsonistically competitive model, firms are differentiated, but view themselves as atomistic relative to the market: they ignore the effect of their behavior on the distribution of options available to each candidate. This assumption is maintained in a number of studies, including Card et al. (2018) and Lamadon, Mogstad, and Setzler (2022), among others. When firms are oligopsonists, on the other hand, they actively incorporate the effects of their behavior on the distribution of options available to each candidate into their wage-setting decisions. In this way, models of oligopsony incorporate strategic interactions between firms. Berger, Herkenhoff, and Mongey (2017) and Jarosch, Nimczik, and Sorokin (2021) estimate models that include strategic interactions of this form. Berger, Herkenhoff, and Mongey (2017) note that, under oligopsony, structural labor supply elasticities to the firm are not equal to reduced-form elasticities, as they are under monopsonistic competition. Under oligopsony, these elasticities depend upon the value of the firms' own amenities, in addition to competitor's amenities (and bids). Importantly, our definition of oligopsonistic behavior encompasses multiple mechanisms that have been explored separately in prior work (for instance, our framework subsumes both size- and differentiation-based mechanisms by which oligopsonists generate wage markdowns).

1.4 Econometric Framework

Candidate Preferences

Identification

We first consider identification of the preference structure from choice data. Our principal identification assumption is that firms do not directly observe Q_i , but rather predict type membership on the basis of observable characteristics. This implies that, given a vector of characteristics x_i , the probability that candidate i receives offer set $\mathcal{B}_i = \{b_{ij}, B_{ij}\}_{j=0}^J$ is independent of i 's true type membership Q_i :

Assumption 1.1. (Conditional Independence) *Firms do not observe Q_i , and so only make decisions about whether and how much to bid on the basis of x_i . This implies that, conditional on posted resume characteristics x_i , firms' bids are independent of candidates' latent preference types Q_i :*

$$\Pr(\mathcal{B}_i \mid Q_i = q, x_i) = \Pr(\mathcal{B}_i \mid x_i).$$

An immediate consequence of Assumption 1.1 is that the distribution of candidate types conditional on received bids \mathcal{B}_i and characteristics x_i is equal to the distribution of types conditional on x_i alone:

$$\Pr(Q_i = q \mid \mathcal{B}_i, x_i) = \frac{\Pr(\mathcal{B}_i \mid Q_i = q, x_i) \Pr(Q_i = q \mid x_i)}{\Pr(\mathcal{B}_i \mid x_i)} = \Pr(Q_i = q \mid x_i).$$

In administrative data, like linked employer-employee records, assumptions similar to Assumption 1.1 are highly implausible due to the various selection mechanisms at play in the formation of equilibrium matches. By contrast, our data contains not only the final matches between firms and candidates, but also the full *distribution* of bids candidates receive. Further, the rules of contact on the platform require firms to make initial bids on the basis of candidate profiles alone, before they have the chance to interact with candidates (and thereby update their forecasts of candidate preferences). Since we observe the same profile information that firms do (x_i), we are able to closely approximate the information set available to firms when forming bids. This feature is one of the advantages of using data from online hiring platforms and has been recognized in other studies. For instance, Hangartner, Kopp, and Siegenthaler (2021) study discrimination in hiring on a large online job board. Because they observe all variables visible to employers on the site, they argue that they are able to control for all relevant confounds.

We next formalize additional assumptions about the structure of preferences implicit in the model of labor supply specified in the previous section. Denote the set of bids that i accepts by \mathcal{B}_i^1 , and likewise denote the set of bids i rejects by $\mathcal{B}_i^0 = \mathcal{B}_i \setminus \mathcal{B}_i^1$. Given a set of bids \mathcal{B}_i , we let $\mathcal{B}_i^1 \succ \mathcal{B}_i^0$ denote the event $\min_{j \in \mathcal{B}_i^1} V_{ij} \geq \max_{k \in \mathcal{B}_i^0} V_{ik}$: every option in i 's accepted set is revealed-preferred to every option in i 's rejected set. We refer to $\mathcal{B}_i^1 \succ \mathcal{B}_i^0$ as a partial ordering over options.

Assumption 1.2. (Mixture Model) *The probability of observing any partial ordering is described by a finite mixture model over latent preference types:*

- a) (Finite Support)** *The support of the distribution of latent types is finite – without loss of generality, we restrict the support of Q_i to the integers $1, \dots, Q$. The conditional probability of type membership is denoted by:*

$$\Pr(Q_i = q \mid x_i) \triangleq \alpha_q(x_i).$$

- b) (Exclusion Restriction)** *Conditional on a candidate's latent type and offer set, the probability of observing any partial ordering is independent of x_i :*

$$\Pr(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid Q_i = q, x_i) = \Pr(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid Q_i = q) \triangleq \mathcal{P}_q(\mathcal{B}_i^1 \succ \mathcal{B}_i^0).$$

Assumption 1.2a is a modelling choice about the form of unobserved heterogeneity in preferences over firms. Assumption 1.2b is an exclusion restriction that governs how preferences are related to individual characteristics: the variables in x_i shift the distribution of

types, but provide no additional information about preferences conditional on those types. Importantly, Assumption 1.2b is an implication of the labor supply model we specified in the previous section.

Combining Assumptions 1.1 and 1.2, we may express the likelihood of the partial ordering $\mathcal{B}_i^1 \succ \mathcal{B}_i^0$, given an option set \mathcal{B}_i and profile characteristics x_i , as:

$$\begin{aligned} \Pr(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid \mathcal{B}_i, x_i) &= \sum_{q=1}^Q \Pr(Q_i = q \mid x_i) \times \Pr(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid Q_i = q) \\ &= \sum_{q=1}^Q \alpha_q(x_i) \times \mathcal{P}_q(\mathcal{B}_i^1 \succ \mathcal{B}_i^0). \end{aligned}$$

Mixtures of random utility models (RUMs) of this form have been studied in both econometrics and computer science/machine learning. In particular, Soufiani et al. (2013) establish identifiability of a finite-mixture-of-types RUM for which the idiosyncratic error components follow a log-concave distribution, as assumed in our model. As in Sorkin (2018), we can only rank firms that are members of a connected set: to be a member of the set, a firm must have been both revealed-preferred to at least one member of the set, and have been revealed-dispreferred to a at least one member of the set. This identification condition is identical to that of conditional logit models that require variation in binary outcomes for every unit.

Estimation

We produce estimates of the labor supply parameters using a two-step procedure. In the first step, we estimate β and a transformation of the amenity values A_{qj} . To do so, we maximize the likelihood of each candidate’s revealed preference ranking over firms *for which they received identical wage offers*.⁷ Once we have obtained first step estimates, we use them in a second step to estimate the remaining labor supply parameters. In particular, we estimate those parameters in a generalized method of moments procedure in which we specify conditional moment restrictions on the interview acceptance probability.

Parameterization. In order to estimate preferences, we first specify a tractable parameterization of the labor supply model. The monetary component utility function is assumed to be continuous, with a kink at the point at which the bid salary equals the ask salary. We

⁷ Typically, exact matching of observations on a continuous covariate is extremely challenging. In our case, however, the overwhelming bunching of wage offers at ask (in addition to additional bunching of wage offers at round numbers) means that we may still use the majority of observations for estimation of amenity values and the distribution of unobserved heterogeneity.

write this function as:

$$\begin{aligned} u(b, a) &= \theta_0 \cdot [\log(b) - \log(a)] + \theta_1 \cdot [\log(b) - \log(a)]_- \\ &= (\theta_0 + \theta_1 \cdot \mathbf{1}[b < a]) \cdot \log(b/a) \\ &= \begin{cases} \theta_0 \cdot \log(b/a) & \text{if } b \geq a \\ (\theta_0 + \theta_1) \cdot \log(b/a) & \text{if } b < a, \end{cases} \end{aligned}$$

where $[x]_- = x \cdot \mathbf{1}[x < 0]$ denotes the negative part of x . Note that we have defined $u(b, a)$ relative to the outside option: when $b = a$, $\log(b/a) = \log(1) = 0$, and so $u(b, a)$ is continuous at $b = a$.⁸ Under monopsonitic competition, the structural labor supply elasticity parameters θ_0 and θ_1 coincide with the elasticities of labor supply to individual firms, and markdowns only vary based upon whether bids are above or below ask. Under oligopsony, the elasticity of labor supply to each firm depends additionally on the amenity value of the firm, and therefore varies both across firms and within firms between workers of different preference types. When oligopsonistic firms are not type-predictive, they only exploit across-firm differences in average labor supply elasticities, while type-predictive oligopsonists exploit both between- and within-firm differences in labor supply elasticities.

We let \mathbf{Q}_i denote a $Q \times 1$ vector of mutually exclusive and exhaustive indicators Q_{iq} for membership in type q ($Q_{iq} = 1$ if $Q_i = q$). We specify the distribution of types as a multinomial logit in profile characteristics x_i :

$$\Pr(Q_{iq} = 1 \mid x_i) = \alpha_q(x_i \mid \beta) = \frac{\exp(x'_i \beta_q)}{\sum_{q'=1}^Q \exp(x'_i \beta_{q'})}.$$

We additionally let $A_j(Q_i) = \mathbf{Q}'_i \mathbf{A}_j$, where \mathbf{A}_j is a $Q \times 1$ vector of type-specific mean amenity values at firm j with q -th component A_{qj} . Finally, we assume that the distribution of taste shocks is extreme value type 1:

$$\xi_{ij} \stackrel{iid}{\sim} EV_1,$$

and so the particular labor supply system we estimate is a discrete mixed-logit random utility model.

First Step. The first step of our procedure is to estimate the distribution of preference types and (a transformation of) the type-specific mean amenity valuations, or rankings, for each firm. Our estimation strategy is based on a simple observation: if candidate i accepts an offer from j and rejects an offer from k when $b_{ij} = b_{ik}$, then by revealed preference:

$$\mathbf{Q}'_i (\mathbf{A}_j - \mathbf{A}_k) \geq \xi_{ik} - \xi_{ij}.$$

Candidates often have several offers at the same bid salary – most often at exactly their ask, but also often at round numbers. Because exact matching of offers at the same salary is

⁸ To make comparisons of utility between candidates, we add back the monetary component associated with the outside option: $u(b, a) + \theta_0 \cdot \log(a)$.

possible in our setting, we subset to sets of offers made to candidates at the same bid salary for the purpose of estimating amenity values.

In order to model the joint probability of the full set of choices candidates make, we must derive the probability of observing an arbitrary partial ordering of firms, $\mathcal{P}_q(\mathcal{B}_i^1 \succ \mathcal{B}_i^0)$. Define the re-parameterization:

$$\rho_{qj} = \frac{\exp(A_{qj})}{\sum_{k=1}^J \exp(A_{qk})},$$

and let $\sigma(\cdot) : \{1, \dots, J\} \rightarrow \{1, \dots, J\}$ denote a linear order or ranking of all J alternatives. A multinomial logit model over rankings of alternatives is sometimes called a Plackett-Luce (Plackett 1975; Luce 1959) model, or an exploded logit. Given this notation, the likelihood of observing any full ranking of alternatives is given by:

$$\Pr(\sigma(\cdot) \mid \boldsymbol{\rho}_q) = \prod_{r=1}^J \frac{\rho_{q\sigma^{-1}(r)}}{\sum_{s=r}^J \rho_{q\sigma^{-1}(s)}}.$$

Unlike the standard Plackett-Luce/exploded logit setting, we only observe candidates' partial orderings of firms. Following Allison and Christakis (1994), we could compute the probability of observing any particular partial ordering of preferences by summing over all linear orders that are consistent with that partial ordering. Even with a small number of alternatives, however, this strategy is computationally intractable: the number of concordant linear orders grows exponentially in the number of alternatives. Simulation methods that sample linear orders (e.g. Liu et al. 2019) are likely to be slow, and introduce additional sources of noise.

We circumvent this issue by implementing a novel numerical approximation to the partial order likelihood that greatly reduces the computational burden of estimation. Our strategy relies on the well known fact that the maximum of independent EV_1 random variables is also distributed EV_1 :

$$\Pr\left(\max_{k \in \mathcal{B}_i^0} \log(\rho_{qk}) + \xi_{ik} < v\right) = F_\xi\left(v - \log\left(\sum_{k \in \mathcal{B}_i^0} \rho_{qk}\right)\right),$$

where $F_\xi(x) = \exp(-\exp(-x))$ is the EV_1 CDF. Using this observation, in combination with a simple change of variables argument, we can re-write the probability of the partial ordering

$\mathcal{B}_i^1 \succ \mathcal{B}_i^0$, conditional on preference parameters $\boldsymbol{\rho}_q$, as:

$$\begin{aligned}
\mathcal{P}(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid \boldsymbol{\rho}_q) &= \Pr \left(\min_{j \in \mathcal{B}_i^1} \log(\rho_{qj}) + \xi_{ij} > \max_{k \in \mathcal{B}_i^0} \log(\rho_{qk}) + \xi_{ik} \mid \boldsymbol{\rho}_q \right) \\
&= \int_{-\infty}^{\infty} \prod_{j \in \mathcal{B}_i^1} (1 - F_{\xi}(v - \log(\rho_{qj}))) \times dF_{\xi} \left(v - \log \left(\sum_{k \in \mathcal{B}_i^0} \rho_{qk} \right) \right) \\
&= \int_{-\infty}^{\infty} \prod_{j \in \mathcal{B}_i^1} \left(1 - F_{\xi} \left(v - \log \left(\sum_{k \in \mathcal{B}_i^0} \rho_{qk} \right) \right)^{\rho_{qj} / \sum_{k \in \mathcal{B}_i^0} \rho_{qk}} \right) \\
&\hspace{20em} \times dF_{\xi} \left(v - \log \left(\sum_{k \in \mathcal{B}_i^0} \rho_{qk} \right) \right) \\
&= \int_0^1 \prod_{j \in \mathcal{B}_i^1} \left(1 - u^{\rho_{qj} / \sum_{k \in \mathcal{B}_i^0} \rho_{qk}} \right) du.
\end{aligned}$$

The second line uses the independence of ξ_{ij} and the distribution of $\max_{k \in \mathcal{B}_i^0} \log(\rho_{qk}) + \xi_{ik}$, the third line uses the fact that $F_{\xi}(x - \log(a)) = F_{\xi}(x - \log(b))^{a/b}$, and the fourth line substitutes $u = F_{\xi}(v - \log(\sum_{k \in \mathcal{B}_i^0} \rho_{qk}))$. This expression, and its derivatives, can be quickly and accurately approximated by numerical quadrature. The log-integrated likelihood of i 's revealed partial order is therefore given by:

$$\mathcal{L}(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid x_i, \boldsymbol{\beta}, \boldsymbol{\rho}) = \log \left(\sum_{q=1}^Q \alpha_q(x_i \mid \boldsymbol{\beta}) \times \mathcal{P}(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid \boldsymbol{\rho}_q) \right).$$

We estimate $\boldsymbol{\beta}$ and $\boldsymbol{\rho}$ via a first-order generalized EM-algorithm. Details of the estimation procedure are given in Appendix A.4.

While our estimation procedure differs in several ways from those of existing studies, the logic of the ranking methodology is similar to that of Sorkin (2018) and Avery et al. (2013). As in those studies, the estimated rank of firm j depends not on j 's raw acceptance probability, but the composition of firms to which j was revealed preferred. Sorkin (2018) summarizes this property as a recursion: highly-ranked firms are those that are revealed-preferred to other highly-ranked firms. Avery et al. (2013) note that producing rankings in this way is robust to potential strategic manipulations of the units being ranked – a key property in our setting. While we do not present a formal proof of consistency here, parameter consistency of the MLE for similar models has been established under sequences in which the number of items to be ranked (here, the number of firms J) grows asymptotically, avoiding the usual incidental parameters problem (Neyman and Scott 1948). Simons and Yao (1999) established the consistency and asymptotic normality of the maximum likelihood estimator of the parameters of Bradley-Terry models of paired comparisons (a special case of Plackett-Luce) under asymptotics that hold fixed the number of comparisons available between each pair of choices, but let the number of choices tend to infinity. Yan, Yang, and Xu (2012) and Han, Xu, and Chen (2020) generalized this result to sparse comparison matrices in which not all choices are compared and the numbers of available comparisons for each pair of choices are random variables. Graham (2020) develops similar results for logistic regression under sparse network asymptotics.

Second Step. The second step of our procedure requires estimating the labor supply elasticity parameters (θ_0, θ_1) , outside option values (\mathbf{A}_0) , and scaling factors $(\boldsymbol{\sigma})$, which we carry out by GMM. We form moment conditions around the model-implied probability of accepting an interview request, given our first-step estimates $\hat{\boldsymbol{\beta}}$ and $\hat{\boldsymbol{\rho}}$ and the remaining parameters $\Theta = \{\theta_0, \theta_1, \mathbf{A}_0, \boldsymbol{\sigma}\}$. This probability is given by:

$$\Pr(D_{ij} = 1 \mid b_{ij}, x_i) = \sum_{q=1}^Q \alpha_q(x_i \mid \hat{\boldsymbol{\beta}}) \times \Lambda\left((\theta_0 + \theta_1 \cdot \mathbf{1}[b_{ij} < a_i]) \cdot \log(b_{ij}/a_i) + \sigma_q \times \log(\hat{\rho}_{qj}) - A_{q0}\right),$$

where the function $\Lambda(x) = (1 + \exp(-x))^{-1}$ is the logistic CDF. Let $m(b_{ij}, x_i \mid \Theta)$ denote this model-based estimate of $\Pr(D_{ij} = 1 \mid b_{ij}, x_i)$ evaluated at the parameters Θ . We specify conditional moment conditions of the form:

$$\mathbb{E}[x_i \cdot (D_{ij} - m(b_{ij}, x_i \mid \Theta))] = 0 \quad \text{and} \quad \mathbb{E}[z_j \cdot (D_{ij} - m(b_{ij}, x_i \mid \Theta))] = 0.$$

We compute the sample analogues of these moment conditions and stack them in the vector $\widehat{m}(\Theta)$. We estimate the components of Θ by minimizing:

$$\widehat{\Theta} = \arg \min_{\Theta} \widehat{m}(\Theta)' \mathbf{W} \widehat{m}(\Theta)$$

for a symmetric, positive-semidefinite weighting matrix \mathbf{W} . In practice, we use an efficient two-step GMM procedure, in which we produce an initial estimate $\widehat{\Theta}^0$ with \mathbf{W}^0 set equal to an identity matrix. We construct an updated weighting matrix \mathbf{W} by computing the inverse of the covariance matrix of the moment conditions evaluated at the initial estimate $\widehat{\Theta}^0$, which we then use to construct an efficient estimate $\widehat{\Theta}$.

Labor Demand

Preliminaries: Construction of $G_{ij}^m(b)$

Before we can implement the estimation and testing procedure outlined below, we must first produce approximations to firms' beliefs for each combination of conduct assumptions. Definition 1.1 specified a general form for beliefs in equilibrium. Beliefs depend upon the probability that candidates will rank a firm's bid highest among all available options, and that probability conditions on a random variable Λ_i which summarizes the distribution of the maximum of the utilities available to i . In our multinomial logit setting, we take Λ_i to be the *inclusive value* of the full set of bids offered to i :

$$\Lambda_i = \log\left(\sum_{k: B_{ik}=1} \exp\left(u(b_{ik}, a_i) + Q'_i A_k\right)\right).$$

Given Λ_i , the probability that i ranks j 's bid highest can be written:

$$\Pr(V_{ij} = V_i^1 \mid \Lambda_i, b_{ij} = b) = \exp\left(u(b, a_i) + Q'_i A_j\right) / \exp\left(\Lambda_i\right).$$

Using this expression, we re-write firms' beliefs as:

$$G_{ij}(b) = \sum_{q=1}^Q \left(\int \left[\exp(u(b, a_i) + A_{qj}) / \exp(\lambda) \right] dF_{\Lambda|Q}(\lambda \mid Q_{iq} = 1, \omega_{ij}^{\Lambda}) \right) \times \alpha_q(\omega_{ij}^Q).$$

We construct approximations to $G_{ij}(b)$ under two alternative conduct assumptions about firms' beliefs about the distribution of Λ_i :

Monopsonistic Competition: Under the monopsonistic competition alternative, firms do not take into account the contribution of their own bid on the inclusive value Λ_i – in other words, $b_{ij} \notin \omega_{ij}^{\Lambda}$. Let Λ_{iq} denote the inclusive value of i 's offer set, conditional on $Q_{iq} = 1$. Under this assumption, the expression for firms' beliefs simplifies to:

$$G_{ij}(b) = \sum_{q=1}^Q \left(\exp(u(b, a_i) + A_{qj}) \times \mathbb{E} \left[\exp(-\Lambda_{iq}) \mid \omega_{ij}^{\Lambda} \right] \right) \times \alpha_q(\omega_{ij}^Q).$$

Since firms are assumed to have rational expectations conditional on the information ω_{ij}^{Λ} , the quantity $\mathbb{E} \left[\exp(-\Lambda_{iq}) \mid \omega_{ij}^{\Lambda} \right]$ can be approximated by regressing $\exp(-\Lambda_{iq})$ on a flexible function of the variables contained in ω_{ij}^{Λ} (which include x_i, z_j). This argument mirrors the intuition of Guerre, Perrigne, and Vuong (2000): in a rational expectations equilibrium, bidders' beliefs are consistent with the true distribution of winning bids in an auction, and so beliefs (and therefore markdowns) can be approximated by the empirical distribution of winning bids. Given estimates of the labor supply parameters and $\mathbb{E} \left[\exp(-\Lambda_{iq}) \mid \omega_{ij}^{\Lambda} \right]$, the beliefs of monopsonistically-competitive firms can be written as: $G_{ij}^m(b) = (b/a_i)^{\theta_0 + \theta_1} \mathbf{1}_{[b < a_i]} \times C_{ij}^m$, where C_{ij}^m is a model-specific constant. This implies that markdowns are a constant fraction of the wage on either side of $b_{ij} = a_i$: $\frac{\theta_0}{1 + \theta_0}$ when $b_{ij} > a_i$, and $\frac{\theta_0 + \theta_1}{1 + \theta_0 + \theta_1}$ when $b_{ij} < a_i$. When $b_{ij} = a_i$, we have that $\mu_{ij}^m = a_i / \varepsilon_{ij} \in \left[\frac{\theta_0}{1 + \theta_0}, \frac{\theta_0 + \theta_1}{1 + \theta_0 + \theta_1} \right]$.

Oligopsony: Under the oligopsony alternative, firms do take into account the contribution of their own bid on the inclusive value Λ_{iq} – in other words, $b_{ij} \in \omega_{ij}^{\Lambda}$. In this case, we have that:

$$\Lambda_{iq} \mid b_{ij} \sim \exp(u(b_{ij}, a_i) + A_{qj}) + \exp(\Lambda_{iq}^{-j}),$$

where $\Lambda_{iq}^{-j} = \log \left(\sum_{k \neq j: B_{ik} = 1} \exp(u(b_{ik}, a_i) + Q'_i A_k) \right)$ is the leave- j -out inclusive value. Denote the probability distribution of Λ_{iq}^{-j} by $F_{\Lambda_q^{-j}}$. Under this assumption, firms' beliefs can be written:

$$G_{ij}(b) = \sum_{q=1}^Q \int \left(\frac{\exp(u(b, a_i) + A_{qj})}{\exp(u(b_{ij}, a_i) + A_{qj}) + \exp(\lambda)} \times dF_{\Lambda_q^{-j}}(\lambda \mid \omega_{ij}^{\Lambda}) \right) \times \alpha_q(\omega_{ij}^Q).$$

Again, since firms' beliefs are assumed to be consistent, $F_{\Lambda_q^{-j}}(\lambda \mid \omega_{ij}^{\Lambda})$ can be approximated by computing the distribution of leave-one-out inclusive values in the sample – for instance,

by computing a series of quantile regressions of Λ_{iq}^{-j} on a flexible function of the variables contained in ω_{ij}^Λ . We can then use these estimates to construct a numerical approximation to the integral over the distribution of leave- j -out inclusive values. Unlike monopsonistic competition, there is no simple closed-form expression for markdowns in the oligopsony case.

In order to approximate $G_{ij}(b)$, we must also specify how firms forecast candidate preferences. We consider two alternatives for assumptions about firms' beliefs about the distribution of Q_i :

Type Predictive: Under the type-predictive alternative, firms predict candidate types given observed profile characteristics x_i ($\omega_{ij}^Q = x_i$). In this case, we approximate these predictions using the estimated prior over types, $\alpha_q(\omega_{ij}^Q) = \alpha_q(x_i | \hat{\beta})$.

Not Predictive: Under the not-predictive alternative, firms do not predict candidate types given observed profile characteristics x_i ($\omega_{ij}^Q = \emptyset$). In this case, we assume that firms weight type-specific win probabilities by the average probability of type membership, $\alpha_q(\omega_{ij}^Q) = \bar{\alpha}_q = \frac{1}{N} \sum_{i=1}^N \alpha_q(x_i | \hat{\beta})$.

We produce approximations to $G_{ij}(b)$ under all four combinations of these conduct assumptions. In addition, we consider a baseline **Perfect Competition** case, in which firms are assumed to bid their valuations.

Identification and Estimation in the General Model

Next, we consider identification and estimation in our general framework for labor demand. Let m denote a choice of model, as specified by a combination of conduct assumptions. Each model m is associated with a particular belief about the population win probability $G_{ij}(b)$, which we denote by $G_{ij}^m(b)$. To illustrate the intuition of our estimation procedure, assume for the moment that $G_{ij}^m(b)$ is differentiable, and denote the derivative of $G_{ij}^m(b)$ with respect to b as $g_{ij}^m(b)$. Under this assumption, bids must satisfy the first-order condition:

$$\underbrace{b + \frac{G_{ij}^m(b)}{g_{ij}^m(b)}}_{=\varepsilon_{ij}^m(b)} = \gamma_j^m(x_i, \nu_{ij}^m),$$

where $\varepsilon_{ij}^m(b)$ is the inverse bidding function under model m ($b = b_{ij}^m(\varepsilon_{ij}^m(b))$).⁹ Crucially, the inverse bidding function is *known* once we have specified a set of conduct assumptions

⁹ Labor economists may be more familiar with the equivalent formulation of the firms' first-order condition in terms of a multiplicative markdown $\mu_{ij}^m(b)$ expressed as a function of the elasticity of labor supply to the firm, $\eta_{ij}^m(b)$, evaluated at the optimal bid:

$$\mu_{ij}^m(b) = \frac{b \times g_{ij}^m(b) / G_{ij}^m(b)}{1 + b \times g_{ij}^m(b) / G_{ij}^m(b)} = \frac{\eta_{ij}^m(b)}{1 + \eta_{ij}^m(b)}.$$

m and plugged in labor supply parameters estimated in a previous step: in a Bayes-Nash Equilibrium, productivity is “revealed” by the bid. If the function $\varepsilon_{ij}^m(\cdot)$ is an injection, then a unique implied valuation $\varepsilon_{ij}^m = \varepsilon_{ij}^m(b_{ij})$ can be inferred for every bid b_{ij} . Given conditional moment restrictions of the form $\mathbb{E}[\nu_{ij}^m \mid \Omega_{ij}] = 0$ (arising, for instance, from exclusion restrictions), we could estimate the productivity function $\gamma_j^m(x_i, \nu_{ij})$ by regressing ε_{ij}^m on (flexible functions of) the determinants of productivity under a functional form assumption. By standard arguments, the parameters that govern $\gamma_j^m(\cdot, \cdot)$ are identified given sufficient variation in model-implied markdowns and the covariates. This approach is taken by Backus, Conlon, and Sinkinson (2021) in their analysis of the common-ownership hypothesis in product markets. Our setting differs from this example in two important ways, both of which motivate the maximum likelihood framework we adopt.

First, we explicitly model labor supply as a kinked function of the bid. This implies that $G_{ij}^m(b)$ is not differentiable at $b = a$, and so the first-order condition for pricing does not hold in general. In Appendix A.5, we establish that bidding strategies $b_{ij}^m(\cdot)$ and option values $\pi_{ij}^{m*}(\cdot)$ are continuous, monotonic functions of firms’ valuations ε_{ij} as a consequence of the log-concavity of F_ξ and the shape restrictions we place on $u(b, a)$. In particular, we show that $b_{ij}^m(\cdot)$ is a strictly-increasing function of ε_{ij} outside an interval $[\varepsilon_{ij}^{m-}, \varepsilon_{ij}^{m+}]$, and is equal to a_i when ε_{ij} is inside that interval. We also show that $\pi_{ij}^{m*}(\cdot)$ is strictly increasing over all valuations. This implies that bids partially identify valuations (and therefore option values) in each model: bids not equal to ask map to a unique valuation, while bids equal to ask map to an interval of possible valuations $[\varepsilon_{ij}^{m-}, \varepsilon_{ij}^{m+}]$. This motivates our use of a Tobit-style maximum likelihood procedure that incorporates a mass point of bids made exactly at ask.

Second, selection into bidding is a key feature of our setting: firms only bid on candidates for whom the maximized option value exceeds a threshold c_j . This implies that the conditional moment restriction $\mathbb{E}[\nu_{ij}^m \mid \Omega_{ij}] = 0$ does not hold in general, but rather that $\mathbb{E}[\nu_{ij}^m \mid \Omega_{ij}] > 0$ in the sample for which bids are observed. While selection poses an estimation challenge, it also provides an opportunity for an additional source of differentiation between models: different conduct assumptions lead to different predictions about the option value of each bid, and thereby imply different patterns of selection which may or may not be reflected in the data. We deal with selection by leveraging a feature of our models of bidding: under each conduct assumption, firms’ bids reveal not only their valuations, but also the maximized value of their objective functions. For every bid made not at ask, we can construct the option value implied by the model, and for every bid made at ask, we can construct an upper bound on the option value implied by the model. We denote these values by $\hat{\pi}_{ij}^{m*}$, and use them to construct a consistent estimate of each firm’s interview cost threshold (under the assumptions of model m) by taking the minimum among all bids made by that firm:

$$\hat{c}_j^m = \min_{i: B_{ij}=1} \hat{\pi}_{ij}^{m*} \xrightarrow{\text{a.s.}} c_j^m.$$

The consistency of our estimate of c_j necessarily depends upon the number of observations per firm growing without bound. See Appendix A.6 for a proof of this result.¹⁰

¹⁰ Our proof of the consistency of \hat{c}_j^m for each firm j (and model m) closely follows the proof of Lemma 1 (ii) of Donald and Paarsch (2002).

Using this estimate, we can compute a lower bound on the valuation associated with each bid, which we use to implement a selection correction. Because $\pi_{ij}^{m*}(\cdot)$ is a strictly increasing function, there is a unique lower-bound valuation $\underline{\varepsilon}_{ij}^m$ at which firm j is indifferent between bidding and not bidding on candidate i . This lower bound controls the selection into bidding: employer j must draw a valuation of at least $\underline{\varepsilon}_{ij}^m$ to make a bid on candidate i , and so the distribution of valuations is censored from below by $\underline{\varepsilon}_{ij}^m$. Given our estimate of c_j , we construct candidate-specific lower bounds by numerically inverting the option value function; $\widehat{\underline{\varepsilon}}_{ij}^m$ is the number that sets:

$$\pi_{ij}^{m*}(\widehat{\underline{\varepsilon}}_{ij}^m) = \widehat{c}_j^m.$$

We use these lower bound estimates to construct the likelihood contribution of each bid, which is given by:

$$\begin{aligned} \mathcal{L}_{ij}^m(\Psi^m) &= \Pr\left(\varepsilon_{ij} = \varepsilon_{ij}^m(b_{ij}) \mid \varepsilon_{ij} \geq \widehat{\underline{\varepsilon}}_{ij}^m, \Psi^m\right)^{\mathbf{1}[b_{ij} \neq a_i]} \\ &\quad \times \Pr\left(\varepsilon_{ij} \in [\varepsilon_{ij}^{m-}, \varepsilon_{ij}^{m+}] \mid \varepsilon_{ij} \geq \widehat{\underline{\varepsilon}}_{ij}^m, \Psi^m\right)^{\mathbf{1}[b_{ij} = a_i]} \\ &= \left(\frac{f_\varepsilon(\varepsilon_{ij}^m(b_{ij}); \Psi^m)}{1 - F_\varepsilon(\widehat{\underline{\varepsilon}}_{ij}^m; \Psi^m)}\right)^{\mathbf{1}[b_{ij} \neq a_i]} \times \left(\frac{F_\varepsilon(\varepsilon_{ij}^{m+}; \Psi^m) - F_\varepsilon(\max(\varepsilon_{ij}^{m-}, \widehat{\underline{\varepsilon}}_{ij}^m); \Psi^m)}{1 - F_\varepsilon(\widehat{\underline{\varepsilon}}_{ij}^m; \Psi^m)}\right)^{\mathbf{1}[b_{ij} = a_i]}, \end{aligned}$$

where Ψ^m denotes the parameters for model m , $f_\varepsilon(\cdot; \Psi^m)$ is the density of ε_{ij} given parameters Ψ^m , $F_\varepsilon(\cdot; \Psi^m)$ is the CDF of ε_{ij} given parameters Ψ^m , $\varepsilon_{ij}^m(\cdot)$ is the inverse bidding function for model m , and ε_{ij}^{m+} and ε_{ij}^{m-} are the model-implied upper and lower bounds on ε_{ij} when $b_{ij} = a_i$.¹¹

Parameterization: In order to estimate the distribution of valuations under each set of conduct assumptions, we make assumptions about the functional forms of $\gamma_j(x_i, \nu_{ij})$ and the distribution of ν_{ij} , F_ν . We parameterize $\gamma_j(x_i, \nu_{ij})$ as log-linear in the sum of ν_{ij} and a bi-linear form in candidate and firm characteristics, as in Lindenlaub and Postel-Vinay (2021):

$$\begin{aligned} \gamma_j(x_i, \nu_{ij}) &= \exp\left(z_j' \Gamma x_i + \nu_{ij}\right) \\ z_j' \Gamma x_i &= \sum_k \sum_\ell \gamma_{k\ell} z_{jk} x_{i\ell}, \end{aligned}$$

where both x_i and z_j include a constant. We further assume:

$$\nu_{ij} \stackrel{iid}{\sim} N(0, \sigma_\nu).$$

¹¹ The approach we take here – concentrating the c_j parameters out of the likelihood by computing the minimum order statistic – is similar to that of Donald and Paarsch (1993, 1996, 2002), who consider models in the classic procurement auction setting. However, because the thresholds c_j are not functions of any of the other parameters of the model, our estimation procedure yields a proper likelihood (unlike some of the cases they consider).

For each model m , we construct estimates $\hat{\Gamma}^m$ and $\hat{\sigma}_\nu^m$ by maximizing the log-likelihood of the complete set of bids for all companies in the connected set (this includes bids on all candidates, not just those in the connected set).

Discriminating Between Non-Nested Models of Conduct

We next turn to our testing procedure. Given sets of parameter estimates for each model, our objective is to determine which of those models is closest to the true data-generating process. The models we consider are *non-nested*: “Broadly speaking, two models (or hypotheses) are said to be ‘non-nested’ if neither can be obtained from the other by the imposition of appropriate parametric restrictions or as a limit of a suitable approximation; otherwise they are said to be ‘nested’” Pesaran (1990). In our setting, models are non-nested as long as they 1) generate distinct combinations of markdowns and selection corrections, and 2) those markdowns and selection corrections are not co-linear with the determinants of productivity (the elements of z_j and x_i and their interactions).

To provide intuition for our testing procedure, consider again the simpler case in which $G_{ij}(b)$ is assumed to be differentiable. Under our functional form assumptions and the true conduct assumption, we may write:

$$\log(\varepsilon_{ij}(b_{ij})) = z'_j \Gamma x_i + \nu_{ij}.$$

This equation includes only one source of error: the idiosyncratic component of firms’ valuations, ν_{ij} , which are assumed to be independent of both x_i and z_j , in addition to market-level variables. Of course, the true model of conduct is unknown, so in practice we must substitute the true inverse bidding function $\varepsilon_{ij}(\cdot)$ with our approximation under conduct assumption m , $\varepsilon_{ij}^m(\cdot)$.¹² If model m is mis-specified, then using $\varepsilon_{ij}^m(\cdot)$ in place of $\varepsilon_{ij}(\cdot)$ introduces a mis-specification error:

$$\log(\varepsilon_{ij}^m(b_{ij})) = z'_j \Gamma x_i + \nu_{ij} + \zeta_{ij}^m.$$

The presence of mis-specification error suggests two rather intuitive conclusions. First, models that are further from the truth should perform worse on standard goodness-of-fit metrics, since the residual variance combines the contributions of both ν_{ij} and ζ_{ij}^m . Second, if labor supply responses (and therefore markdowns) are determined in part by variables that are excluded from the productivity function, then the estimated residuals of models that are far from the truth should be strongly correlated with those excluded variables.

This is the basic logic of Berry and Haile (2014). They establish the necessity of instruments that shift demand (analogous to labor supply in our setting), but that are excluded from the marginal cost function (analogous to valuations or productivity in our setting), for identification in the product market setting with data only on market shares. Such variation, they note, is particularly important for testing between models of conduct. Following this logic, Backus, Conlon, and Sinkinson (2021) implement a test of conduct that formalizes the second conclusion above: under true conduct assumptions, instruments that affect markups

¹² Keeping in mind, under assumption m , we may treat $\varepsilon_{ij}^m(b_{ij})$ as data.

(markdowns) but do not affect marginal costs (valuations) should not be correlated with recovered idiosyncratic cost shocks (ν_{ij}).

Our setting, and the nature of the data we use, differs in several key ways from that of Berry and Haile (2014). The most basic difference is that we have access to micro data on individual choices, rather than market-level data. Berry and Haile (2020) consider identification of differentiated products demand using micro data on individual choices, and demonstrate that access to micro data significantly reduces reliance on instruments. Our use of micro data in the form of multiple choices for each candidate, combined with our ability to condition on all information available to firms when they bid, allowed us to identify candidate preferences without requiring additional instruments for prices (bids). A second major difference between our setting and that of Berry and Haile (2014) is that we analyze *individualized* bids rather than uniform market prices. Bids are made before any negotiation has taken place and without direct knowledge of the competition, and so they do not have to satisfy a market clearing condition. Rather, we assume that firms' behavior must satisfy a conditional form of rational expectations about competition. Given this assumption, our identification arguments follow those of the empirical auction literature, like Guerre, Perrigne, and Vuong (2000) or Backus and Lewis (2020).

Despite the relatively less stringent requirements for instruments to identify labor supply in our setting, the power of our testing procedure to discriminate between models of conduct still depends upon using additional sources of variation in markdowns that are independent of the determinants of firms' valuations. Without such variation, our ability to discriminate between models of conduct may be severely limited. In other words: without an instrument, our ability to discriminate between models will be driven by differences in functional form.

Instrumenting Labor Supply with Market Tightness

To obviate these concerns, and thereby increase the power of our testing procedure, we use relative market tightness as an instrument for firms' expectations about competing bids. Our use of market tightness as an instrument mirrors the arguments of papers studying auctions with entry that use variation in the potential number of entrants to identify models of auctions with selective entry (e.g. Gentry and Li (2014)). We define tightness as the number of active candidates in a particular experience, occupation, and two-week period cell divided by the number of firms searching for candidates in that experience, occupation, and two-week period cell.¹³ For every candidate, we define the variables n_{iw}^J , n_{iw}^I as the number of firms searching for i 's experience level and occupation during two-week period w and the number of candidates with active profiles in i 's experience level and occupation during two-week period w , respectively. Market tightness is the ratio of the two counts:

$$t_{iw} = n_{iw}^I / n_{iw}^J,$$

where the prevailing level of tightness at the time j bids on i is denoted t_{ij} (similarly define n_{ij}^I and n_{ij}^J). We define tightness within occupation and experience bins because those

¹³ Technically, our instrument is the inverse of the usual definition of market tightness, which is the ratio of vacancies to the level of unemployment. The particular form of instrument does not matter for our analysis.

categories are the primary search fields recruiters use when browsing candidates. Further, we define tightness within two-week periods because that is the default length of time a candidate’s profile will remain active, and therefore variation in tightness between periods is driven primarily by the rate of flow of new candidates onto the platform.

We assume that labor market tightness does not affect firm valuations, but does affect firms’ expectations about competition for i as encoded by Λ_i . The intuition is simple: the more active firms there are per active candidate, the more bids those candidates can expect to receive. We formalize this assumption as:

Assumption 1.3. (Instrument Exogeneity) *Labor market tightness is independent of idiosyncratic determinants of labor demand:*

$$t_{ij} \perp\!\!\!\perp \nu_{ij} \mid x_i, z_j.$$

We incorporate variation in tightness by including t_{ij} (and n_{ij}^I , n_{ij}^J , and occupation, experience, and two-week period dummies) in the set of variables firms use to predict inclusive values, ω_{ij}^Λ (which also includes x_i and z_j). Variation in tightness thereby drives variation in predicted markdowns that is independent of firms’ valuations. We propose two non-nested model comparison tests that leverage this exclusion restriction in different, but complementary, ways.

Option 1: The Vuong (1989) Likelihood Ratio Test

Because we estimate models by maximum likelihood, a natural first option for our test of conduct is a straightforward application of the Vuong (1989) likelihood ratio test. The Vuong (1989) test is a pairwise, rather than ensemble, testing procedure: rather than explicitly identifying the “best” model among a set of alternatives, the test considers each pair of models in turn and asks whether one of those models is closer to the truth than the other. In the likelihood setting, the “better” of two models is the one with greatest goodness-of-fit, as measured by the maximized log-likelihoods.¹⁴

Let $s = |ij : B_{ij} = 1|$ denote the sample size. For a pair of models m_1 and m_2 , denote the maximized sample log-likelihoods by $\mathcal{L}_s^{m_1}$ and $\mathcal{L}_s^{m_2}$, respectively, where:

$$\mathcal{L}_s^m = \max_{\Psi} \sum_{ij: B_{ij}=1} \log \left(\mathcal{L}_{ij}^m(\Psi) \right),$$

and Ψ^m denotes the arg max. The null hypothesis of our test is that m_1 and m_2 are equally close to the truth, or *equivalent*. In this case, the population expectation of the difference in log likelihoods is zero. There are two one-sided alternative hypotheses: that m_1 is closer to the truth than m_2 , and vice versa. When m_1 is closer to the true data-generating process, the population expectation of the likelihood ratio $\mathbb{E}^0[\log(\mathcal{L}_{ij}^{m_1}(\Psi^{m_1})/\mathcal{L}_{ij}^{m_2}(\Psi^{m_2}))]$ is greater than zero. Vuong (1989) shows that when m_1 and m_2 are non-nested, an appropriately-scaled

¹⁴ The population expectation of the log-likelihood measures the distance, in terms of the Kullback-Liebler Information Criterion (KLIC), between the model and the true data generating process.

version of the sample likelihood ratio is asymptotically normal under the null that the two models are equivalent:

$$Z_s^{m_1, m_2} = \frac{\mathcal{L}_s^{m_1} - \mathcal{L}_s^{m_2}}{\sqrt{s} \cdot \hat{\omega}_s^{m_1, m_2}} \xrightarrow{D} \mathcal{N}(0, 1),$$

where $\hat{\omega}_s^{m_1, m_2}$ is the square root of a consistent estimate of the asymptotic variance of the likelihood ratio, $\omega_*^{2m_1, m_2}$. We set:

$$\hat{\omega}_s^{m_1, m_2} = \left(\frac{1}{s} \sum_{ij: B_{ij}=1} \log \left(\frac{\mathcal{L}_{ij}^{m_1}(\Psi^{m_1})}{\mathcal{L}_{ij}^{m_2}(\Psi^{m_2})} \right)^2 \right)^{1/2}.$$

We construct test statistics $Z_s^{m_1, m_2}$ for every pair of models we estimate. Given a significance level α with critical value c_α , we reject the null hypothesis that m_1 and m_2 are equivalent in favor of the alternative that m_1 is better than m_2 when $Z_s^{m_1, m_2} > c_\alpha$, and vice versa if $Z_s^{m_1, m_2} < -c_\alpha$. If $|Z_s^{m_1, m_2}| \leq c_\alpha$, the test cannot discriminate between the two models.

How does variation in the instrument increase the power of the test? The answer depends on the relevance of the instrument for predicting markdowns. Returning to the simplified example above, we may write the mis-specification error as:

$$\zeta_{ij}^m = \log(\varepsilon_{ij}^m(b_{ij})) - \log(\varepsilon_{ij}(b_{ij})).$$

To the extent that variation in tightness drives variation in markdowns under the true model, variation in tightness will also generate variation in ζ_{ij}^m if the assumed model m is mis-specified. This implies that relatively more mis-specified models will imply valuations that are more difficult to explain using observables than those that are closer to the truth.

Option 2: The Rivers and Vuong (2002) Test

Rivers and Vuong (2002) proposed a generalization of the Vuong (1989) testing procedure that extended the logic of that test to a much wider class of objective functions. In their analysis of firm conduct, Backus, Conlon, and Sinkinson (2021) implement a version of the Rivers and Vuong (2002) test by specifying a single moment condition involving the residuals of fitted models and excluded instruments. We propose a variant of that test using the generalized residuals associated with the likelihood we estimate. Gourieroux et al. (1987) define generalized residuals and explicate their use in testing. In the context of maximum likelihood estimation, the generalized residuals are defined by the scores of the likelihood. Let $s_{ijk\ell}^m(\Psi) = \partial \mathcal{L}_{ij}^{m_1}(\Psi) / \partial \psi_{k\ell}^m$ denote the k, ℓ -th component of the score vector for observation ij . The scores may be written as $s_{ijk\ell}^m(\Psi) = h_{ij}^m(\Psi) \cdot z_{jk} \cdot x_{i\ell}$, where $h_{ij}^m(\Psi)$ is the generalized residual for observation ij under model m and parameters Ψ . The maximum likelihood estimate $\hat{\Psi}^m$ is the vector that sets the mean of the scores to zero:

$$\sum_{ij: B_{ij}=1} s_{ijk\ell}^m(\hat{\Psi}^m) = \sum_{ij: B_{ij}=1} h_{ij}^m(\hat{\Psi}^m) \cdot z_{jk} \cdot x_{i\ell} = 0 \quad \forall k, \ell,$$

and so generalized residuals are constrained to be orthogonal to covariates. The generalized residuals for each model can be easily computed by taking the derivative of the individual likelihood contributions.

We form the generalized residuals for each model, and use them to compute the scalar moment/lack-of-fit measure:

$$Q_s^m = \left(\frac{1}{s} \sum_{ij: B_{ij}=1} h_{ij}^m(\widehat{\Psi}^m) \cdot t_{ij} \right)^2.$$

Q_s^m measures the covariance between the generalized residuals of each model and the excluded instrument t_{ij} . Under proper specification, the influence of the instrument on markdowns is completely summarized by the inverse bidding function, and so there should be zero correlation between the instrument and the generalized residual. A separate way to motivate the lack-of-fit measure Q_s^m is as an unscaled version of the score test statistic for testing against the null hypothesis that the coefficient on t_{ij} in the labor demand equation is zero.

Following Backus, Conlon, and Sinkinson (2021),¹⁵ we formulate a pairwise test statistic for testing between models m_1 and m_2 as an appropriately-scaled difference between $Q_s^{m_1}$ and $Q_s^{m_2}$, which Rivers and Vuong (2002) show to be asymptotically normal:

$$T_s^{m_1, m_2} = \frac{Q_s^{m_1} - Q_s^{m_2}}{\widehat{\sigma}_s^{m_1, m_2} / \sqrt{s}} \xrightarrow{D} \mathcal{N}(0, 1),$$

where $\widehat{\sigma}_s^{m_1, m_2}$ is an estimate of the population variance of $Q_s^{m_1} - Q_s^{m_2}$. We compute an estimate of $\widehat{\sigma}_s^{m_1, m_2} / \sqrt{s}$ as the variance of $Q_s^{m_1} - Q_s^{m_2}$ across bootstrap replications. Given a significance level α with critical value c_α , we reject the null hypothesis that m_1 and m_2 are equivalent in favor of the alternative that m_1 is better than m_2 when $T_s^{m_1, m_2} < c_\alpha$, and vice versa if $T_s^{m_1, m_2} > c_\alpha$. If $|T_s^{m_1, m_2}| \leq c_\alpha$, the test cannot discriminate between the two models.

The intuition for this test is relatively more straightforward than for the first test: the lack-of-fit measures each pairwise test compares can themselves be interpreted as test statistics associated with a score test of the exclusion restriction. In some ways, this feature makes the test relatively more appealing than the first option. However, the power of the test depends entirely on the ability of the instrument to predict differential markdowns and selection corrections, which is not the case for our first test (see Duarte et al. (2021) for a discussion of weak instruments problems in conduct testing). For these reasons, we present the results of both tests and view the two procedures as complementary.

¹⁵ Backus, Conlon, and Sinkinson (2021) formulate their moment-based test statistic by interacting residuals with an appropriate function of both the excluded instrument and all other exogenous variables, and connect their choice of that function to the literature on optimal instruments (Chamberlain 1987). In our setting, the formulation of an appropriate function that combines the instrument and other exogenous variables is complicated by the issues of selection and partial identification we previously highlighted. While we do not pursue it here, the formulation of such a function is a focus of future work.

1.5 Model Estimates

Labor Supply

Model selection and validation

Before describing the estimated preference orderings and group structures, we must settle on a baseline version of the model. In particular, we need to specify the number of latent preference classes Q , and we need to specify how class membership is related to candidate observables. To that effect, for each pair of models – a given number of ladders and a given set of observables used to define group membership –, we calculate a standard likelihood ratio statistic and compute the appropriate χ^2 p -value. In addition to formal likelihood ratio (LR) statistics, we also compute a more directly-interpretable “goodness-of-fit” (GoF) statistic for each model. The statistic is simply the fraction of pairwise revealed-preference comparisons that are concordant with the estimated rankings for each model. Specifically, we define:

$$\text{GoF} = \frac{1}{N_{pw}} \sum_{i=1}^N \sum_{q=1}^Q \alpha_q(x_i | \hat{\beta}) \times \left(\sum_{j \in \mathcal{B}_i^1} \sum_{k \in \mathcal{B}_i^0} \mathbf{1} \left[\hat{A}_{qj} \geq \hat{A}_{qk} \right] \right),$$

where N_{pw} is the total number of pairwise comparisons implied by revealed preference.

Table 1.3 reports these goodness-of-fit statistics for several versions of our labor supply model. Each row corresponds to a given number of ladders (from one to four) and each column corresponds to the observables leveraged to construct class membership. In the first row, we estimate the model with a single preference group ($Q = 1$), such that there is no additional preference heterogeneity for a given firm aside from variation in idiosyncratic preference shocks ξ_{ij} . In the second row, we estimate a model with two preference groups. The first column allows men and women to have different rankings of firms, and the second columns splits candidates between above- and below-median experience. The last column leverages all the observables we access for the candidates to define latent preference groupings. In particular, we estimate the prior probability of group membership $\alpha_q(x_i)$ concurrently with the preference orderings themselves. We then refer to each preference class as a separate job ladder.

A model that assigns random numbers for each A_{qj} would in expectation yield a GoF statistic of 0.5. As reported in the first row of Table 1.3, the single-ladder model, in which there is common mean ranking of firms for all candidates, increases goodness-of-fit over that baseline to 0.67.¹⁶ Table 1.3 second finding concerns the comparison of goodness-of-fit between the single-ladder model and the two models that split candidates into preference groups based on observable characteristics. In Column 1, allowing women and men to have distinct rankings of firms on the second row has no additional explanatory power for the

¹⁶ The goodness of fit measure varies slightly across the three columns because the estimation samples are different. For instance, to be ranked in the model that splits the ladder by gender, a firm needs to have been accepted once and rejected once by candidates of both genders. The resulting sample will differ from the model that splits by experience, where to be in the connected set, a firm needs to have been accepted once and rejected once by candidates of all experience levels.

revealed preferences in the data, in comparison to the single ladder model from the first row: the GoF statistic increases imperceptibly (from 0.672 to 0.680), and the formal LR test fails to reject the null that the two-ladder model is equivalent to the single-ladder model ($p = 0.27$). The finding that men and women have very similar mean preference orderings over firms mirrors that of Sorkin (2017), who also finds that the implied preference orderings of men and women over firms are extremely similar. Splitting by experience does only marginally better: while the LR test can reject the null that the two-ladder model is equivalent to the single-ladder model ($p < 0.001$), the GoF statistic only increases by 1.6 percentage point. Our third finding is that the model using the full set of observables to define the clusters performs markedly better than the gender- and experience-split models. For the same number of ladders (two), the GoF statistics for the model-based clustering is 0.744, that is 10.7 percentage points higher than the gender or experience splits. Our final finding concerns the number of ladders: sequential LR tests between the one- and two-ladder models and two- and three-ladder models both reject the null that the more-complex models are equivalent to the simpler models ($p < 0.001$). In addition to the two- and three-ladder models, we estimated a model with four preference groupings, but were unable to reject the null that this model was equivalent to the three-ladder alternative. We therefore adopt the three-ladder model as our baseline model of candidate preferences. Plugging in those estimated rankings into our second-step GMM procedure yields the following labor supply elasticity parameter estimates:

$$u(b_{ij}, a_i) = \left[\begin{array}{c} 4.05 \\ (0.33) \end{array} + \begin{array}{c} 1.58 \\ (0.28) \end{array} \cdot \mathbf{1}[b < a_i] \right] \cdot \log(b/a_i).$$

These estimates are similar to others in the literature – for instance, Berger, Herkenhoff, and Mongey (2017) report an estimate of 3.74 for this parameter (what they call the within-market substitutability parameter), while Azar et al. (2020) report an estimate of 5.8.¹⁷

In order to validate the estimated rankings, we take advantage of the fact that candidates may sometimes provide reasons for rejecting an interview request. While the platform does not require candidates to list a reason, 58% of them do. When providing a rejection reason, candidates select from a list of options that includes reasons like “company culture”, “firm size”, and “poor timing”, among others. We divide the list into two categories: personal reasons that should correspond to a low draw of ξ_{ij} and job-related reasons that should correspond to a low value of A_{qj} . If the model provides a good fit to the data, then we should find that candidates are more likely to reject highly-ranked firms for personal reasons than job-related reasons relative to lower-ranked firms. Figure 1.5 plots the probability that a firm was rejected for a job-related reason as a function of firms’ ordinal rankings (where lower ranks are better) – we indeed find that workers are significantly less likely to reject the most-preferred companies for job-related reasons than they are for lower-ranked companies. Appendix Figure A.3 provides additional evidence of the quality of the fit of the preferred 3-type model. For every bid, we compute the model-implied probability that the bid

¹⁷ Note that, in contrast with other studies, our model allows for kinked labor supply and therefore our estimates of the parameter is 5.63 below the kink, i.e. when $b < w_i$, and 4.05 above the kink.

will be accepted. Appendix Figure A.3 plots the relationship between those model-implied probabilities and the empirical acceptance probability – the model-implied probabilities are extremely close to the actual probability of acceptance throughout the range of the data.

Characterizing the distribution of amenity values

Figure 1.6 illustrates the scale of vertical and horizontal differentiation of firms implied by our preferred model estimates. To understand the relative importance of the amenity values workers attach to firms, we compute a willingness-to-accept statistic (WTA) for every firm. The statistic is equal to the fraction of a candidate’s ask salary that the model implies a firm must offer to make that candidate indifferent between accepting or rejecting an interview request, on average. Specifically, we compute WTA_{qj} as the number that solves:

$$\left(4.05 + 1.58 \times \mathbf{1}[WTA_{qj} < 1]\right) \times \log(WTA_{qj}) + \hat{A}_{qj} - \hat{A}_{q0} = 0.$$

where A_{q0} is the q -th component of the vector of type-specific mean amenity values at the outside option.

Panel (a) of Figure 1.6 plots the distribution of the mean WTA at each firm, averaging over the population probabilities of each type: $WTA_j = \sum_{q=1}^3 \bar{\alpha}_q \times WTA_{qj}$. The average mean WTA is 0.99, indicating that candidates are willing to accept roughly 1% less than their ask at the average firm. The standard deviation of mean WTA across firms is 0.14, which suggests a large range of variability in the amenity values candidates attach to firms. Indeed, there are a nontrivial number of firms for which the average candidate would be willing to accept less than 80% of their ask, and an even larger number of firms for which candidates demand over 120% of their ask. Panel (b) illustrates the systematic component of horizontal differentiation. Here, we plot the within-firm standard-deviation of WTA_{qj} across preference types. The mean within-firm SD of WTA is 0.14, suggesting that the horizontal differentiation is about as important as vertical differentiation. The implication of these estimates is that there is large scope for firms to exercise market power in the ways we have specified: the significant horizontal differentiation suggests that firms may stand to gain significantly from accurately predicting which candidates are in which preference groups, while the significant vertical differentiation suggests that firms with high rankings can afford to mark down wages significantly (assuming they act strategically). Given the significant scope for wage markdown based on preference heterogeneity, assessing whether firms are able to predict the types is crucial to the understanding of their ability to offer type-specific marked down wages. Section 1.5 explores whether firms are type predictive.

What firm characteristics are associated with higher amenity values? To partially answer this question, we report regressions of (standardized) estimates of A_{qj} on firm covariates z_j in the sample for which those covariates are available in Table 1.4. Here, larger values of A_{qj} correspond to better rankings. These covariates represent only a small fraction of the potential relevant characteristics candidates may consider when they choose among job offers – importantly, the (“all-in”) amenity values we estimate do not depend upon exhaustive knowledge of what candidates value. Even with the relatively coarse covariates available, some clear patterns are evident. In particular, the basic evidence in Table 1.4 suggests

a loose classification of groups as “baseline” (group 2), “risk-averse” (group 3), and “risk-loving” (group 1). Relative to baseline, members of group 3 are more interested in working at larger, established firms for which there may be less employment risk, while members of group 1 are more interested in working at the smallest firms that may be more risky bets.

How do worker characteristics shift the probability of preference group membership? In our preliminary goodness-of-fit exercise, we found that explicitly splitting candidates by gender or experience only marginally improved our ability to explain choices – does that result carry over to the more flexible group membership model we estimated? In order to more concretely gauge the associate between covariates and preference types, we compute the model-implied posterior probabilities of type membership for every candidate and correlate those probabilities with candidate characteristics (our discussion of the EM algorithm in Appendix A.4 covers the construction of these probabilities). We find that women are 7 percentage points more likely to be in the risk-averse group and 7 percentage points less likely to be in the risk-loving group, while candidates with above-median experience are 10 percentage points less likely to be in the risk-averse group and 9 percentage points more likely to be in the risk-loving group. While there is significant residual variation in preferences conditional on covariates, our preferred model estimates suggest that covariates are indeed predictive of preference type.

Decomposing group differences in welfare

Given our estimates of amenity values and labor supply parameters, we may fully characterize the utility value candidates associate with the portfolios of bids they receive. Importantly, this allows us to ask whether observable differences in average bids between groups are reflective of underlying differences in welfare. We decompose mean differences in welfare using the Oaxaca-Blinder (OB) decomposition (Oaxaca 1973; Blinder 1973). The OB decomposition posits that variable Y_{ig} corresponding to individual i in group $g = 0, 1$ can be written:

$$Y_{ig} = X'_{ig}\beta_g + \epsilon_{ig},$$

where X_{ig} are covariates measured for all individuals and $\mathbb{E}(\epsilon_{ig}) = 0$. The average value of Y_{ig} in group g is therefore given by $\bar{Y}_g = \bar{X}'_g\beta_g$. We can decompose the difference in the average value of Y_{ig} between groups $g = 1$ and $g = 0$ as:

$$\bar{Y}_1 - \bar{Y}_0 = \underbrace{\bar{X}'_1\beta_1 - \bar{X}'_0\beta_0}_{\text{endowments}} = \underbrace{(\bar{X}_1 - \bar{X}_0)'\beta_0}_{\text{endowments}} + \underbrace{\bar{X}'_0(\beta_1 - \beta_0)}_{\text{coefficients/returns}} + \underbrace{(\bar{X}_1 - \bar{X}_0)'(\beta_1 - \beta_0)}_{\text{interactions}}.$$

The classic OB decomposition apportions the difference in the mean of a variable between two groups into components due to: 1) differences between those groups in *endowments*, or the distribution of relevant covariates; 2) differences between those groups in *coefficients* or *returns* associated with those covariates; and 3) the *interactions* between coefficient and endowment differences.¹⁸ Roughly speaking, the greater the share of the mean difference the

¹⁸ Note that the OB decomposition is not unique – an equivalent “reverse” decomposition may be obtained by replacing β_0 with β_1 in the first term, \bar{X}_0 with \bar{X}_1 in the second term, and flipping the sign of the third term.

OB decomposition apportionments to endowments relative to returns, the more we can conclude that a difference in means is driven by differences in characteristics between those groups, and not how those groups are treated conditional on those characteristics values (differential returns to characteristics). The OB decompositions we present should be interpreted as purely descriptive (Guryan and Charles 2013). Importantly, we exclude the asked salary as an explanatory variable in our OB decompositions of welfare, because candidates formulate their asks as endogenous functions of all of their other characteristics (including gender). The endogeneity of the ask greatly complicates the interpretation of decompositions that include the asked salary: if asks themselves are functions of gender, then gender differences in asks may not be appropriately interpreted as reflecting differing endowments.¹⁹

We report decompositions of welfare-relevant quantities in Table 1.5. The utility associated with each portfolio of bids depends both upon the number of bids received and the composition of those bids. In order to gauge the relative importance of quantity and quality, we compute the total number of bids received by each candidate, as well as the mean values of the components of utility associated with the bids each candidate received. We calculate the monetary component of utility for each bid as:

$$\bar{u}(b_{ij}, a_i) = \left(4.05 + 1.58 \cdot \mathbf{1}[b_{ij} < a_i]\right) \cdot \log(b_{ij}/a_i) + 4.05 \cdot \left(\log(a_i) - \overline{\log(a_i)}\right),$$

where we subtract the (grand) mean of the log of the ask salary ($\overline{\log(a_i)}$) without loss of generality, since the absolute level of utility is not identified. We also compute the mean amenity values associated with each bid, which we decompose into two parts: a common component of amenity valuations shared by all workers, and the worker-specific deviation from that common component: $A_{ij} = \bar{A}_j + \Delta A_{ij}$. The common component is the average candidates' amenity valuation: $\bar{A}_j = \sum_{q=1}^Q \bar{\alpha}_q \cdot \hat{A}_{qj}$ (where $\bar{\alpha}_q$ is the population share of type q). The candidate-specific deviation is the difference between candidate i 's amenity valuation and the average amenity valuation: $\Delta A_{ij} = \sum_{q=1}^Q (\alpha_q(x_i | \hat{\beta}) - \bar{\alpha}_q) \cdot \hat{A}_{qj}$.

To understand how these differences map into welfare, we compute the (expected) inclusive value of every offer set:

$$\Lambda_i^* = \sum_{q=1}^Q \alpha_q(x_i | \hat{\beta}) \cdot \log \left(\sum_{j \in \mathcal{B}_i} \exp(\bar{u}(b_{ij}, a_i) + \hat{A}_{qj}) \right).$$

We decompose (expected) inclusive values into a monetary component and an amenity component. We compute the monetary component of the inclusive value by setting $\hat{A}_{qj} = 0$ for all q and j :

$$\Lambda_i^b = \log \left(\sum_{j \in \mathcal{B}_i} \exp(\bar{u}(b_{ij}, a_i)) \right).$$

¹⁹ Because we omit the ask salary from these decompositions, the effect of the ask salary will be apportioned between the endowments and coefficients components. Any differential patterns in the relationship between characteristics and asks will be reflected in the coefficients component, while mean differences in asks are reflected in the endowments component.

We compute the amenity component of the inclusive value by setting $\bar{u}(b_{ij}, a_i) = 0$ for all i and j . We further decompose the amenity portion into a common component:

$$\bar{\Lambda}_i^A = \sum_{q=1}^Q \bar{\alpha}_q \cdot \log \left(\sum_{j \in \mathcal{B}_i} \exp(\hat{A}_{qj}) \right),$$

and a candidate-specific deviation:

$$\Delta \Lambda_i^A = \sum_{q=1}^Q \left(\alpha_q(x_i | \hat{\beta}) - \bar{\alpha}_q \right) \cdot \log \left(\sum_{j \in \mathcal{B}_i} \exp(\hat{A}_{qj}) \right).$$

Because the inclusive value is a nonlinear function, the relative contributions of each component will not sum to one.

Panel A of Table 1.5 reports decompositions of mean gaps in these quantities by gender (here, the reference group corresponds to women, so positive differences correspond to larger values for men). Column 1 decomposes the gap in the number of bids received by men and women: on average, women receive fewer bids than men. However, slightly more than 100% of this raw gap is driven by differences in endowments: conditional on covariates, women and men receive nearly the same number of bids. Column 2 reports the decomposition of the mean gap in the monetary component of utility: the average monetary value of bids is significantly lower for women than for men. This result is driven by the fact that women ask for less (see Table 1.1), and therefore receive less, conditional on other characteristics—but as discussed above, the ask is an endogenous function of gender. Our decomposition, which excludes the ask as an explanatory variable, suggests that differences in characteristics between men and women can only explain about 1/3 of the raw gap in monetary values, with the rest explained by differential returns. Column 3 decomposes the mean difference in the common component of amenity values. Unconditionally, the bids men receive are from firms with better amenities than the bids women receive. Differences in the returns to characteristics, representing differential selection of firms into bidding by gender, explain 1/3 of this gap. In other words, even conditional on covariates, women receive bids from firms the average worker values relatively less than those that bid on men.

Column 4 decomposes differences in candidate-specific components of the amenity valuation. Here, we find a (small) reverse gap: women value the amenities associated with the bids they receive relatively more than the average worker would, and do so to a greater degree than men. What might be driving this pattern? Without knowing how firms behave, we cannot discriminate between possible explanations. One possibility is that the pattern is driven by differences in the degree of assortative matching of firms to male and female candidates—that is, firms' valuations over candidates might be more correlated with the preference of female candidates than male candidates. Another possibility is that firms are type-predictive and better at targeting offers to female candidates relative to male candidates, all else equal.²⁰ These qualitative patterns are reflected in the decompositions of

²⁰ Evidence from Section 1.5 that firms are in fact not type-predictive suggests the former explanation is more likely than the latter.

components of inclusive values, reported in columns 5-8. Taken together, these results suggest that the large observed gender gap in bids is reflective of a large gender gap in welfare. Unconditionally, the gap in welfare between men and women is exacerbated by differences in the amenity values of the bids they receive. However, differences in covariates between men and women account for most of the unconditional gap.

Panel B of Table 1.5 reports decompositions of mean gaps in welfare by education level, where the reference group is candidates without a graduate degree. Here, we find that candidates with graduate degrees receive slightly fewer bids than those without graduate degrees, but that the average quality of those bids is higher along all components. Again, differences in the monetary component of utility are driven by the fact that candidates with graduate degrees ask for more than those without on average (candidates without graduate degrees ask for \$10,800 less than those with graduate degrees). This differential is reflected in the share of the gap explained by returns, which explain about 40% of the raw gap. Unlike with gender, we find that differential returns do not explain differences in the common component of amenity valuations between education levels, although we do find that differences in returns explain nearly all the difference in candidate-specific components of valuations. Again, the evidence we find in these decompositions is consistent with either assortative matching between workers and firms (candidates with high productivity at firm j also value the amenities of firm j), or the effective targeting of firms' bids to the candidates most likely to accept those bids.

Labor Demand

Testing between models

We next describe the results of implementing our estimation and testing framework for labor demand. As a preliminary matter, Figure 1.7 plots the “first stage” relationship between the model-implied inclusive values Λ_i and Λ_i^{-j} and the instrumental variable t_{ij} , conditional on firm and candidate covariates and two-week period dummies. Intuitively, the fewer candidates there are relative to firms (low t_{ij}), the more offers those candidates should receive, and the larger the inclusive values associated with their offer sets should be. This intuition is borne out in Figure 1.7: both full- and leave-one-out inclusive values are strongly negatively related to labor market tightness. As described in the previous sections 1.4, 1.5, and 1.5, we estimate the distribution of full- and leave-one-out inclusive values conditional on all firm covariates, candidate covariates, and instruments, and use those estimated distributions to construct approximations to firms' beliefs under each combination of conduct assumptions.

Figure 1.8 plots the distributions of predicted markdowns in dollars under both the monopsonistic competition and oligopsony alternatives. We compute markdowns as the difference between the model-implied firm valuation and the observed bid: $\varepsilon_{ij}^m - b_{ij}$. In cases where the implied valuation is not point identified (the bid is equal to ask), we take the midpoint of the model-implied range of valuations: $(\varepsilon_{ij}^{m+} + \varepsilon_{ij}^{m-})/2 - b_{ij}$. The two alternatives predict markedly different distributions of markdowns. Under the monopsonistic competition alternative, the average predicted markdown is \$30,503, with a standard deviation of \$6,658.

Further, the distribution of markdowns is relatively symmetric—the mean and median of the distribution are separated by less than \$300, and the skewness of the distribution of markdowns is just 0.35. By contrast, the oligopsony model predicts uniformly larger markdowns than the monopsonistic competition alternative: the mean model-implied markdown under oligopsony is \$43,385. Further, the distribution of markdowns under oligopsony is significantly more variable, with a standard deviation of \$16,357. Finally, the distribution of markdowns under oligopsony is highly skewed: the mean markdown is \$4,000 larger than the median markdown, and the skewness of the distribution is just over 2. The two sets of markdowns are positively correlated, with a correlation coefficient of 0.42. The large differences highlighted by Figure 1.8 illustrate the importance of understanding which form of conduct best describes firm behavior—different assumptions about the presence or absence of strategic interactions lead to strikingly different conclusions about the size of wage markdowns.

Table 1.6 reports the results of implementing our pairwise testing procedure on the five models we estimated, using both the likelihood-based and moment-based versions of the Vuong test. The test statistics we report suggest that we can resoundingly reject the null hypothesis of model equivalence in most cases, and both versions of the test yield remarkably similar conclusions. The “Perfect Competition” model unambiguously performs the worst of all the models we tested. Among the remaining alternatives, the two monopsonistic competition models outperform the two oligopsony models, with the not-predictive monopsonistic competition alternative performing best. We visualize these results in Figure 1.9, which plots generalized residuals for two alternative models against the excluded instrument. Under proper specification, the generalized residuals should not be correlated with the instrument – the further a model’s generalized residuals are from the x-axis, the greater the degree of mis-specification. In the figure, the generalized residuals for the monopsonistic competition alternative are closely aligned with the x-axis, while the generalized residuals for the oligopsony alternative are strongly negatively related to tightness.

Our tests therefore suggest that models of firm behavior in which firms ignore strategic interactions in wage setting are closer approximations to firms’ true bidding behavior on the platform than are models in which firms act strategically. Additionally, while we cannot reject the null hypothesis that the two monopsonistic competition models are equivalent in the likelihood-based test, the moment-based version of the test strongly rejects the type-predictive alternative relative to the not-predictive alternative. The weight of the evidence therefore suggests that firms are not actively type-predictive: in the context of the monopsonistic competition model selected by our procedure, firms do not appear to target their offers to the candidates who are most willing to accept those offers, conditional on productivity. In the following analysis, we adopt the not-predictive monopsonistic competition model as our baseline.

Markdowns and valuations in the preferred model

Given the results of our testing procedure, we next characterize the distribution of valuations implied by the preferred model. Table 1.7 reports a subset of the estimated matrix

of coefficients $\hat{\Gamma}$ that govern labor demand, $\gamma_j(x_i, \nu_{ij}) = \exp(z_j' \Gamma x_i + \nu_{ij})$. The full set of coefficient estimates are reported in Appendix Table A.2. Each cell of Table 1.7 reports the coefficient on the interaction of the variables specified in the corresponding row and column. Column variables are candidate characteristics (x_i), and row variables are firm characteristics (z_j). We normalize the log ask salary by subtracting the log of the unconditional mean asked salary (equivalently, by taking the log of the ratio of ask to mean ask), such that the constant term reflects productivity at the mean ask. The second, third, and fourth rows correspond to dummies for firm size categories, such that the omitted category (subsumed into the constant, the first row of the table) corresponds to the smallest firms (between one and fifteen employees). The remaining three rows correspond to non-exclusive sector dummies. The implied R^2 of the observed determinants of productivity is 0.89, suggesting that the bilinear form we adopted provides a close approximation to the data.

Column 1 of Table 1.7 reports the main effects of each firm characteristic. Interestingly, there at first appears to be essentially no firm size-productivity gradient: small and large firms tend to pay roughly equivalent salaries, all else equal. The apparent lack of a strong relationship between firm size and productivity disappears, however, when we consider the interaction of candidate ask salaries and firm characteristics in Column 2. As first suggested by Roussille (2021), the ask salary is a powerful predictor of productivity: the elasticity of valuations with respect to the asked salary is 0.795. This elasticity is strongly increasing in firm size: workers that are more productive everywhere (on the basis of their ask) are even more productive at larger firms. The next three Columns (3-5) report the main and interaction effects of dummy variables recording gender (= 1 if female), current employment (= 1 if currently employed), and education (= 1 if candidate has at least one graduate degree). In Column 3, we find evidence of a small residual gender gap in firms' valuations: the main effect of the female dummy is a 0.8% reduction in valuations, with some heterogeneity by firm size and industry. Importantly, this residual gender gap is conditional on the level of the ask salary: Roussille (2021) previously documented a statistically- and economically-meaningful gender gap in ask salaries. In Column 4, we find no evidence of any difference in labor demand between employed and unemployed candidates, all else equal. This result is somewhat surprising in light of Kroft, Lange, and Notowidigdo (2013) and Jarosch and Pilossoph (2018), who find that employers screen out unemployed candidates. It may be the case that in our setting, the rich profile information available to employers and the information encoded in the ask salary provide more informative signals of quality than current employment status. Finally, in Column 5, we report estimates of the main and interaction effects of holding a graduate degree. While the main effect is positive, we find a reverse firm size gradient: larger firms value graduate degrees relatively less, all else equal. To assess model fit, in Appendix Figure A.4, we plot the relationship between observed bids and the systematic component of valuations $\gamma_j(x_i)$. The two are very strongly and positively correlated.

How much does variation in observable determinants of demand contribute to overall variation in bids? Given our labor demand parameter estimates and the estimated markdowns for the preferred model, we can decompose variation in bids across firms and candidates to gauge the relative contributions of markdowns, systematic components of valuations, and

idiosyncratic components of valuations. We define markdowns here as the ratio of the observed bid and the model-implied productivity level $\hat{\varepsilon}_{ij}$ ²¹: $\log(\mu_{ij}) = \log(b_{ij}) - \log(\hat{\varepsilon}_{ij})$. The idiosyncratic component of the valuation is therefore given by $\hat{\nu}_{ij} = \hat{\varepsilon}_{ij} - z_j \hat{\Gamma} x_i$. We can then write:

$$\log(b_{ij}) = \underbrace{\log(\mu_{ij})}_{\text{markdown}} + \underbrace{z_j \hat{\Gamma} x_i}_{\text{systematic comp.}} + \underbrace{\hat{\nu}_{ij}}_{\text{idiosyncratic comp.}}.$$

We compute a simple decomposition of the variance of bids by taking the covariance of each side of the above equation with the bid, yielding:

$$\text{Var}(b_{ij}) = \text{Cov}(\log(b_{ij}), \log(\mu_{ij})) + \text{Cov}(\log(b_{ij}), z_j \hat{\Gamma} x_j) + \text{Cov}(\log(b_{ij}), \hat{\nu}_{ij}).$$

Dividing each side of the decomposition by $\text{Var}(\log(b_{ij}))$ yields a simple representation of the relative importance of each factor.²² Individual components of variance are reported in Table 1.8, for both the (preferred) monopsonistic competition/not predictive model as well as the (dispreferred) oligopsony/not predictive model. Under monopsonistic competition (Panel A) markdowns are nearly constant across candidates, such that variation in components of firms' valuations account for 100% of the variation in log bids. The intuition for this is simple: when firms are monopsonistically competitive, they view the structural labor supply elasticity (governed by θ_0 and θ_{a1}) as the elasticity of labor supply to the firm, and so there is no (perceived) variation in labor supply elasticities across firms. (Variation in elasticities around the kink accounts for the small extent of variation in markdowns.) 91% of that variation can be attributed to systematic components of valuations, while the remainder is accounted for by idiosyncratic components. As an illustration of the implications of incorrect assumptions about the form of firm conduct, Panel B reports the variance decomposition under the oligopsony model. Under oligopsony, markdowns account for 10% of the variation in log bids, while systematic components of valuations account for 78% and idiosyncratic components account for 12%. Relative to monopsonistic competition, interpreting variation in bids under the assumption that firms act strategically implies that firms mark down wages much more steeply, and that valuations themselves are more variable (conditional on candidate x 's).

How do our estimates relate to models of additive worker and firm effects (Abowd, Kramarz, and Margolis 1999)? Our model of productivity includes both firm-specific contributions (here captured by z_j), worker-specific contributions (captured by x_i), and the interactions of firm- and worker-specific covariates. Tables 1.7 and A.2 provide evidence that

²¹ Again taking the midpoint of the implied interval of productivity levels when bid equals ask $\hat{\varepsilon}_{ij} = (\hat{\varepsilon}_{ij}^+ + \hat{\varepsilon}_{ij}^-)/2$

²² A second decomposition may be computed by taking the variance of both sides:

$$\begin{aligned} \text{Var}(\log(b_{ij})) &= \text{Var}(\log(\mu_{ij})) + \text{Var}(z_j \hat{\Gamma} x_j) + \text{Var}(\hat{\nu}_{ij}) - 2 \cdot \text{Cov}(\log(b_{ij}), z_j \hat{\Gamma} x_j) \\ &\quad - 2 \cdot \text{Cov}(\log(b_{ij}), \hat{\nu}_{ij}) + 2 \cdot \text{Cov}(z_j \hat{\Gamma} x_j, \hat{\nu}_{ij}). \end{aligned}$$

interactions of worker and firm factors are statistically meaningful determinants of productivity. However, the interaction effects we estimate are generally small, which suggests that additive models might well approximate productivity. To explore this, we regress bids, predicted ε_{ij} , and the predicted systematic component of productivity $\exp(z_j'\hat{\Gamma}x_i)$ on all candidate and firm characteristics, without including interactions. Consistent with Card, Heining, and Kline (2013)'s informal assessment of the log-additivity of wages using mean residuals from Abowd, Kramarz, and Margolis (1999) regressions, we find that the main effects of worker and firm characteristics separately explain the vast majority of variation in bids and productivity, as reflected in uniformly high (adjusted) R^2 values: 0.924 for bids, 0.905 for ε_{ij} , and 0.999 for $\exp(z_j'\hat{\Gamma}x_i)$. In the context of the near-constant markdowns our preferred model implies, this further suggests that additive models of worker and firm effects provide good approximations to log wages.

Finally, how do our estimates of productivity relate to amenities? To explore this question, we compute regression-adjusted averages of amenities and productivity within firm types defined by combinations of size and industry. We regress the model-implied amenity and productivity values on the (log) ask salary, and an exhaustive set of fixed effects for combinations of all other worker characteristics x_i , and dummies for each firm type. Figure 1.10 plots the relationship between (average) firm amenity values and (average) components of productivity, as measured by the estimated firm-type fixed effects. Like Lagos (2021), we find that the highest-amenity firms also tend to be the highest-productivity firms. The story is different for low-productivity firms, where there is a negative relationship between amenities and productivity. These patterns are broadly consistent with a model of endogenous amenities in which firms do not invest in amenities before they reach a certain productivity level. Because wage markdowns are a near-constant fraction of productivity in the preferred model, Figure 1.10 suggests that there may be compensating differentials between low-amenity firms at the competitive fringe of the labor market for tech workers, but not between high-amenity firms.

1.6 Counterfactual Simulations of Bidding Behavior

Scenarios of interest

To better understand the implications of imperfect competition for welfare, we use our supply and demand estimates to simulate bidding outcomes under all four conduct scenarios: $\{\text{monopsonistic competition, oligopsony}\} \times \{\text{not predictive, type-predictive}\}$. To gauge the losses due to imperfect competition, we define a new form of conduct, which we term **price taking**. Under the price taking conduct alternative, firms have no discretion over the wages they offer. Instead, firms are constrained to offer a prevailing market wage, as if set by a Walrasian auctioneer. In our price-taking alternative, we set the equilibrium wage equal to the systematic component of firms' valuations, $b_{ij} = \exp(z_j'\Gamma x_i)$. Given this set of wages, the only decision firms have to make is whether to bid on each candidate. Because firms are price takers in this scenario, we assume that they view themselves as atomistic, as in monopsonistic

competition.²³ In addition to these simulations, we also simulate the effects of a simple policy meant to reduce gender disparities in wages: blinding employers to candidates' gender. This counterfactual entails replacing gender-specific estimates of labor demand with cross-gender averages, and doing the same for estimates of labor supply.

Computing new equilibria

In order to compute counterfactuals, we randomly select 500 firms and 500 candidates from the universe of firms and candidates in the analysis sample. For each firm-candidate pair, we compute the model-implied systematic component of firm valuations using our preferred estimates of labor demand parameters, $\exp(z'_j \widehat{\Gamma} x_i)$. Under a particular conduct assumption, equilibrium is determined by a set of beliefs over the distribution of the utility afforded by the best option in each candidates' offer set. The inclusive value is itself a sufficient statistic for the distribution of the maximum utility option for each candidate. At an equilibrium, firms' beliefs about inclusive values must be consistent with the true distribution of inclusive values generated by the bidding behavior of competing firms. We make the assumption that those beliefs depend only upon the expected value of the inclusive value to simplify our calculations here.

To compute new equilibria, we first conjecture an initial set of (expected) inclusive values Λ_i^1 . We then iterate the following steps:

1. At iteration t , take *iid* draws from a normal distribution with mean zero and standard deviation $\widehat{\sigma}_\nu$ to produce a new set of idiosyncratic components of firms' valuations, ν_{ij}^t . Use these draws, plus the systematic components of valuations $z'_j \widehat{\Gamma} x_i$, to compute ε_{ij}^t .
2. Given ε_{ij}^t and Λ_i^t , compute b_{ij}^t as firm j 's best response (under the assumed form of conduct). If there is no number b such that $G_{ij}^m(b)(\varepsilon_{ij} - b) \geq \widehat{c}_j$, then set $B_{ij}^t = 0$.
3. Given firms' best responses b_{ij}^t and B_{ij}^t , calculate the realized inclusive value for each candidate, $\Lambda_i^{t*} = \mathbb{E}[\log(\sum_{j: B_{ij}^t=1} \exp(u(b_{ij}^t, a_i) + A_{ij}))]$. Compute the vector of expected inclusive values at the next iteration by taking a step $\alpha^t \in [0, 1]$ towards Λ_i^{t*} :

$$\Lambda_i^{t+1} = \alpha^t \Lambda_i^{t*} + (1 - \alpha^t) \Lambda_i^t.$$

We iterate this procedure until the distribution of inclusive values converges. We then use the equilibrium distribution of inclusive values to compute mean counterfactual outcomes by constructing the average across 50 simulations of firm bidding decisions.

Simulation Results

Table 1.9 reports the results of our simulations. For each scenario, we compute the average bid, ratio of bid to ask, markdown, and number of bids received per candidate. We

²³ Because bids vary even conditional on our detailed controls, we automatically ruled out this form of price taking as a potential mode of conduct to describe firms' actual bidding behavior on the platform.

also compute the averages of (scaled) components of utility associated with each candidates' portfolio of bids. The absolute magnitudes of these components of utility do not have a direct interpretation, but relative differences across scenarios are meaningful.

The unconditional means of each of these variables across simulation repetitions are reported in Panel A of Table 1.9. We first consider scenarios in which firms are assumed to be not predictive. Unsurprisingly, average bids are higher (\$169k vs \$145k), and markdowns are lower (10% vs 18%), in the price taking model (column 1) relative to the preferred monopsonistic competition model (column 2). Additionally, candidates receive markedly fewer bids (20 vs 43) under price taking than under monopsonistic competition, reflecting the increased labor costs under price taking. Even though they receive fewer bids under price taking, the increased monetary value of bids more than makes up for the substantial drop in the number offers: the average candidates' expected utility is higher under price taking than it is under monopsonistic competition. On the other hand, candidates fare far worse when firms act strategically (column 3): under oligopsony, candidates receive even fewer bids than when firms are price-takers (13.5), and the monetary value of those bids is even lower than under monopsonistic competition (\$139k). As a result, candidates' expected utilities are lowest under oligopsony. Interestingly, switching to modes of conduct in which firms are assumed to be type-predictive does little to change the unconditional means of each of the variables we summarize here (columns 4-6).

The lack of a difference between the type-predictive and not-predictive alternatives in unconditional mean outcomes obscures substantial differences in outcomes between men and women when firms are type-predictive relative to when they are not predictive. We report differences in mean outcomes across simulations between women and men in panel B of Table 1.9. Across all conduct assumptions, women receive fewer bids than men (note, however, that this difference is not conditional on other characteristics). In absolute terms, the largest gender gaps in bids and welfare are predicted by the monopsonistic competition model, although these differences are partly driven by the fact that firms unconditionally make more bids under monopsonistic competition than they do under the other alternatives. Relative to the unconditional average, women receive 8-10% fewer offers when firms are not type predictive. The gap widens to 12-18% when firms are assumed to be type-predictive, and the oligopsony model predicts the largest relative gaps. Female candidates' expected utility also drops, although to only a relatively small degree. The upshot of these simulations is that firms have significant ability to exercise market power in ways that expand gender gaps, as first posited by Robinson (1933).

Can a simple policy that blinds employers to the gender of the candidates they consider narrow these gaps? Panel C reports differences between mean outcomes for men and women across simulation draws in which firms are constrained to no longer observe the candidate gender. The results from our simulations suggest that the efficacy of such a policy is relatively limited. Across all conduct possibilities, the policy is predicted to marginally increase the expected utility of female candidates relative to their male counterparts—across conduct scenarios, blinding employers to gender lowers the gender gap in expected utilities by 6-9.5%. Interestingly, while blinding not-predictive firms to gender modestly increases the number of offers women receive relative to men, the opposite is true when firms are type-predictive.

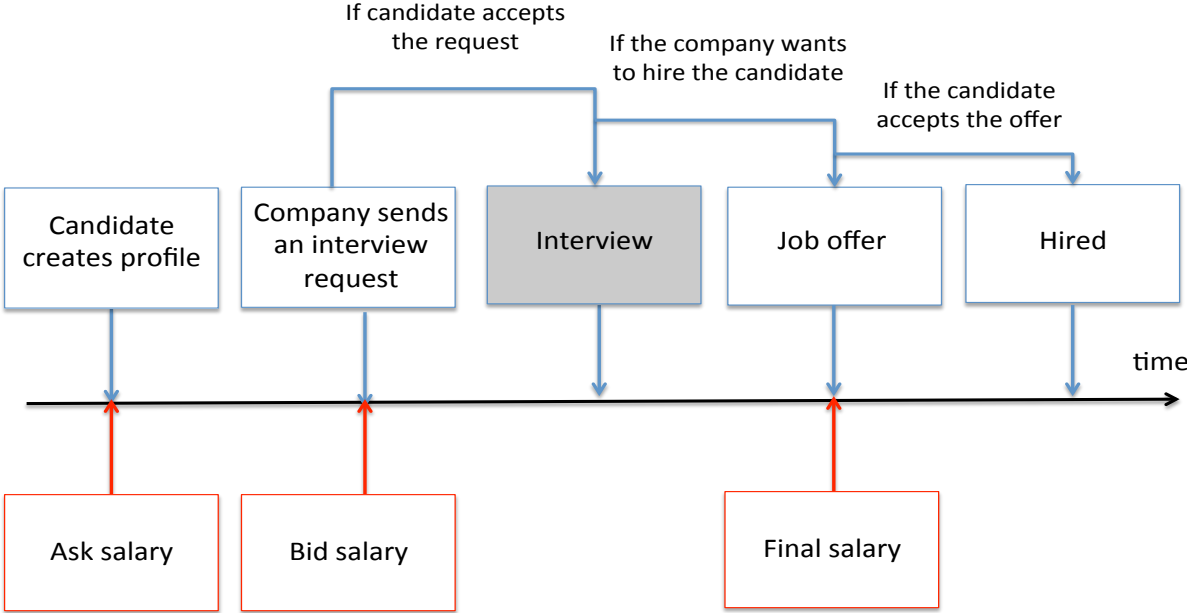
1.7 Conclusion

This paper provides direct evidence about the nature of firms' wage-setting behavior by developing a testing procedure to adjudicate between many non-nested models of conduct in the labor market. In particular, we focus on two sets of alternatives relevant to ongoing debates in the labor literature: first, whether firms compete strategically (Berger, Herkenhoff, and Mongey 2017; Jarosch, Nimczik, and Sorkin 2021), and second, whether firms tailor wage offers to workers' outside options (Caldwell and Harmon 2019; Flinn and Mullins 2021). Applying our testing procedure, we find evidence against strategic interactions in wage setting as well as against the tailoring of offers to workers of different types. Although we study a specific labor market, these findings suggest that the relatively simple model of wage determination posited by Card et al. (2018) provides a reasonable approximation to firm wage-setting conduct in labor markets where many employers are competing for workers. Importantly, we find that incorrect conduct assumptions can lead to substantial biases: in our preferred model, wages are marked down by 18.2% on average, while an oligopsonistic model predicts average markdowns of 25.8%.

Finally, we explore simulations of alternative conduct assumptions to quantify the impact of imperfect competition on welfare. Relative to a price-taking baseline, we find that firms make significantly more offers under the preferred model, but that the wages firms attach to those offers are lower. Relative to the preferred model, however, the average value of bids, the total number of bids, and welfare are significantly lower in simulated equilibria with strategic interactions. We also find that the form of conduct has important implications for gender gaps: relative to men, women receive significantly fewer bids when firms predict horizontal preference variation than when they do not. Imperfect competition exacerbates gender gaps relative to the price-taking baseline. Finally, we find that blinding employers to the gender of candidates generates only modest reductions in gender gaps.

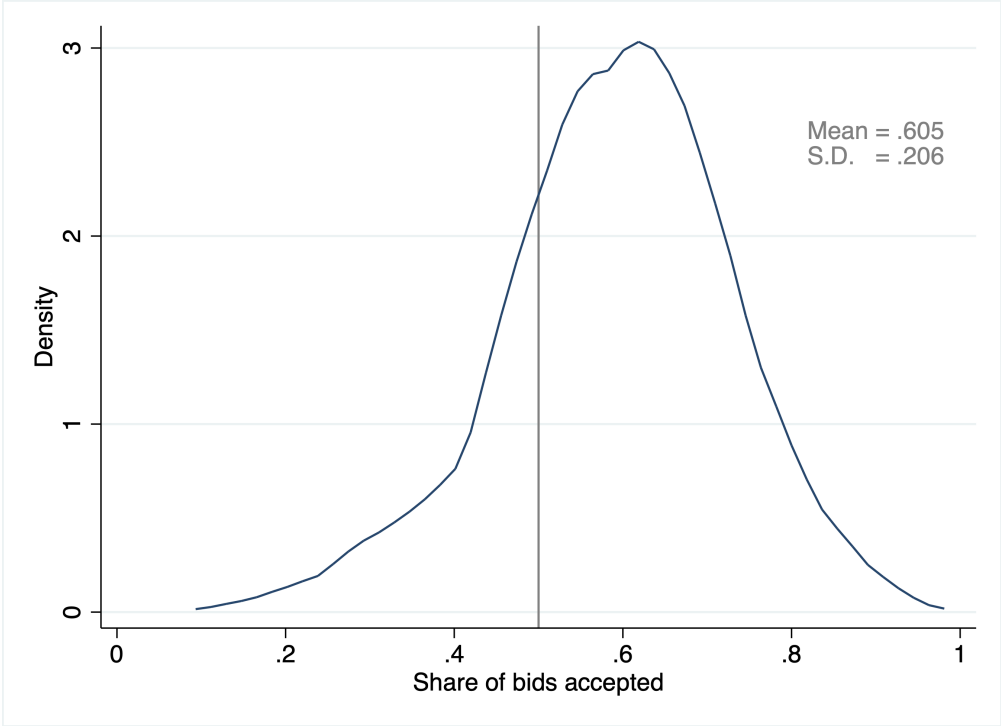
Figures

Figure 1.1: Timeline of the Recruitment Process on Hired.com



Note: This Figure shows the timeline of a recruitment on Hired.com. In red boxes are the different salaries that are captured on the platform. The blue boxes describe all the steps of a recruitment on the platform, from profile creation to hiring. The grey shading for the interview stage indicates that we do not have meta data from companies about their interview process. In green are the classification of the recruitment process between labor demand side (companies) and labor supply side (candidates).

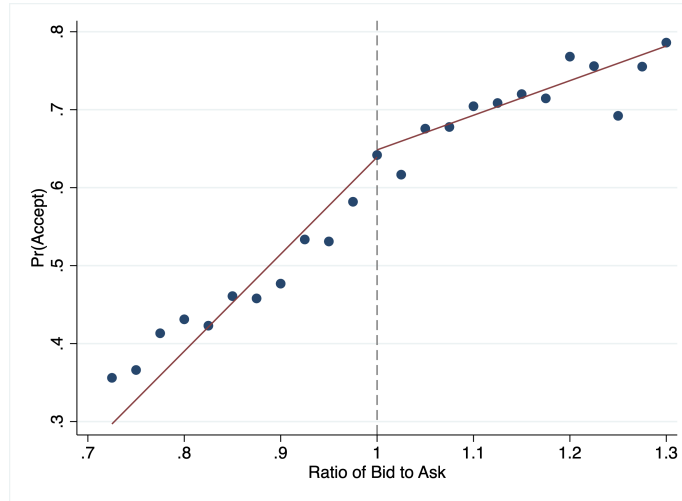
Figure 1.2: Distribution of Fraction of Interview Requests Accepted Across Firms



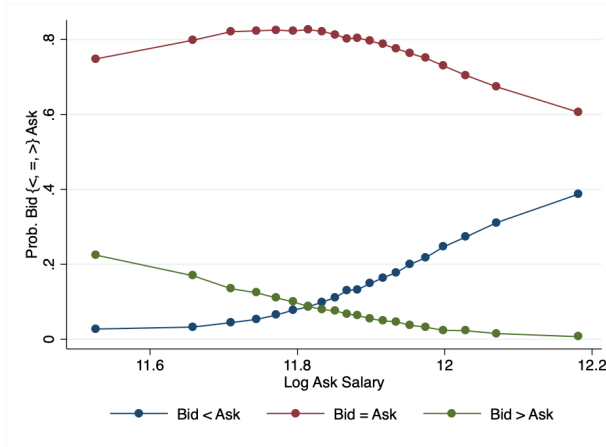
Note: This Figure shows the distribution of the share of accepted interview requests for a given firm. Firms interview requests are frequently rejected by candidates. On average, an interview request by a firm is only accepted 60.5% (SD .206) of the time. For 10.2% of the firms the likelihood that their interview is accepted is less than 40% , while 16.2% of the firms see more than 75% of their interview requests accepted.

Figure 1.3: Empirical Patterns in Bid and Ask Strategies

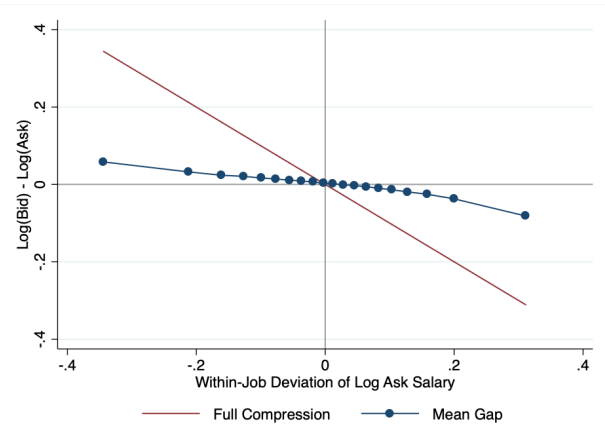
(a) Kink at Bid = Ask



(b) Bids often match Ask

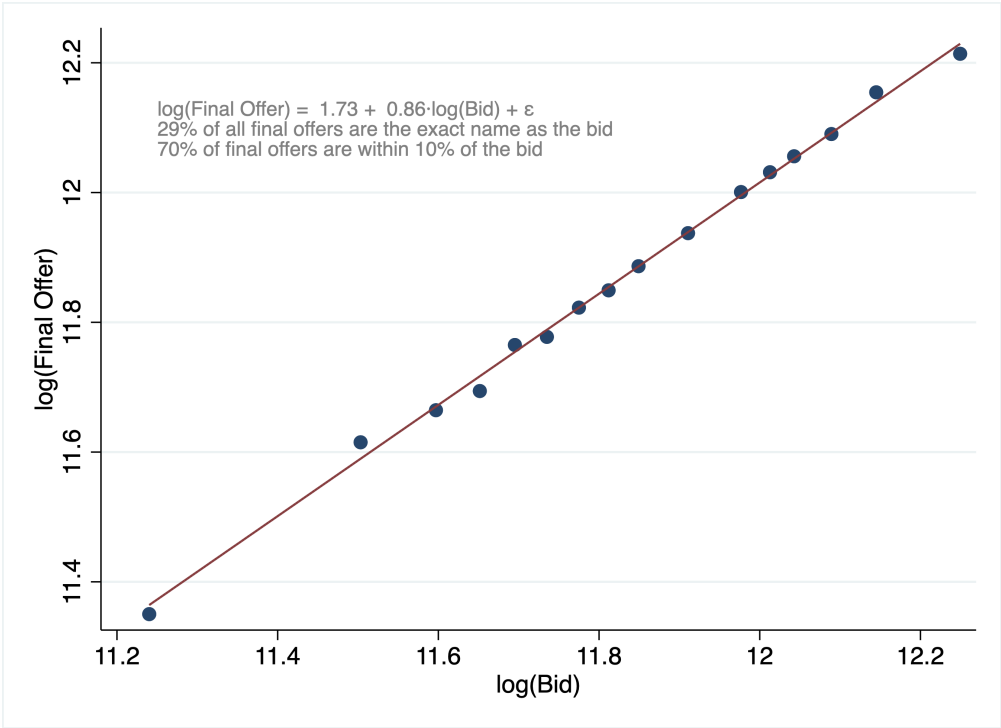


(c) Large range of bid salaries for same job



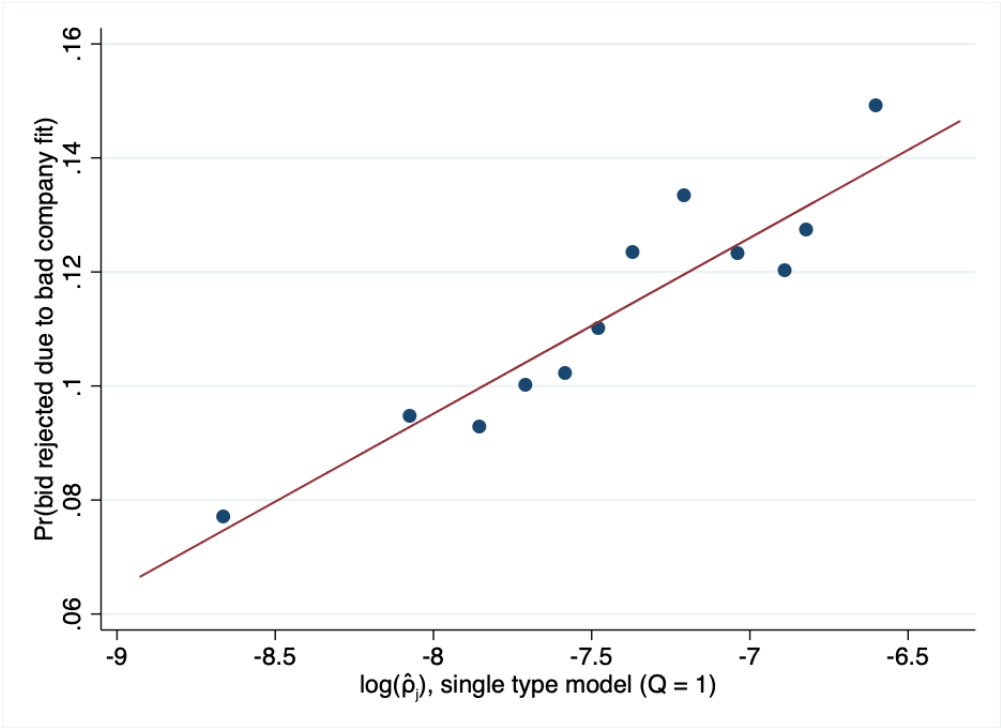
Note: This Figure illustrates several empirical patterns in the relationship between bid and ask salaries. Panel (a) plots the average probability that a candidate accepts an interview request by the company against the ratio of the bid to ask salary in the analysis sample. The slope of the regression line for a bid ask ration of less than one is 1.304 (SE .022), while the slope of the regression line for values greater or equal to 1 is 0.546 (SE .030). Panel (b) shows the relationship between the probability that the bid is, respectively, less, the same or more than the ask, and the level of the (log) ask salary. Panel (c) plots the relationship between the premium – the difference between (log) bid and ask salary – and the within-job deviation of the (log) ask salary.

Figure 1.4: Bids are Sticky in Expectation

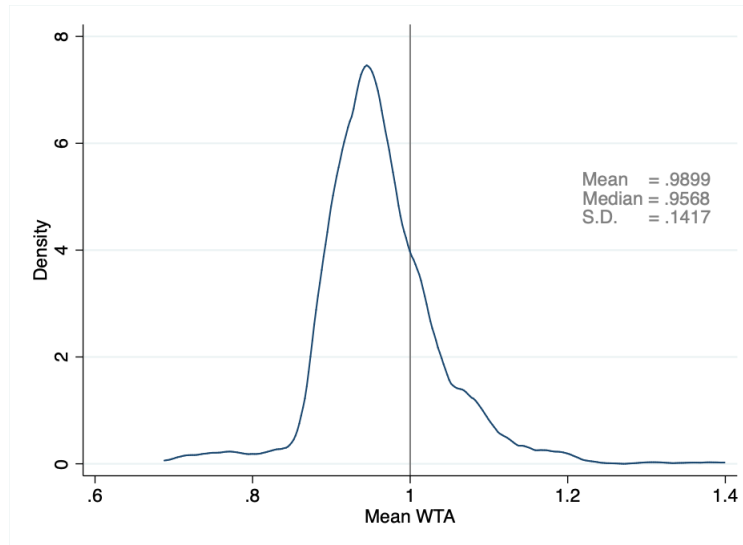
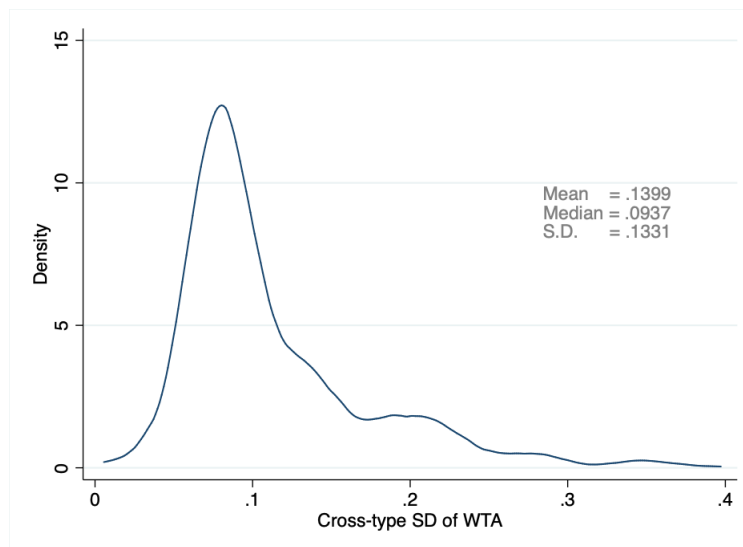


Note: This Figure illustrates the relationship between the initial bid salary sent by a company and the final offer of candidates that are hired for the subset of the analysis sample. The correlation between log bid and log final salary is 0.86 (SE .458). 29% of all final offers in this subset are identical to the bid and 70% of all final offers are within 10% of the initial bid salary.

Figure 1.5: Interview Rejection Reasons as a Function of Firm Rankings

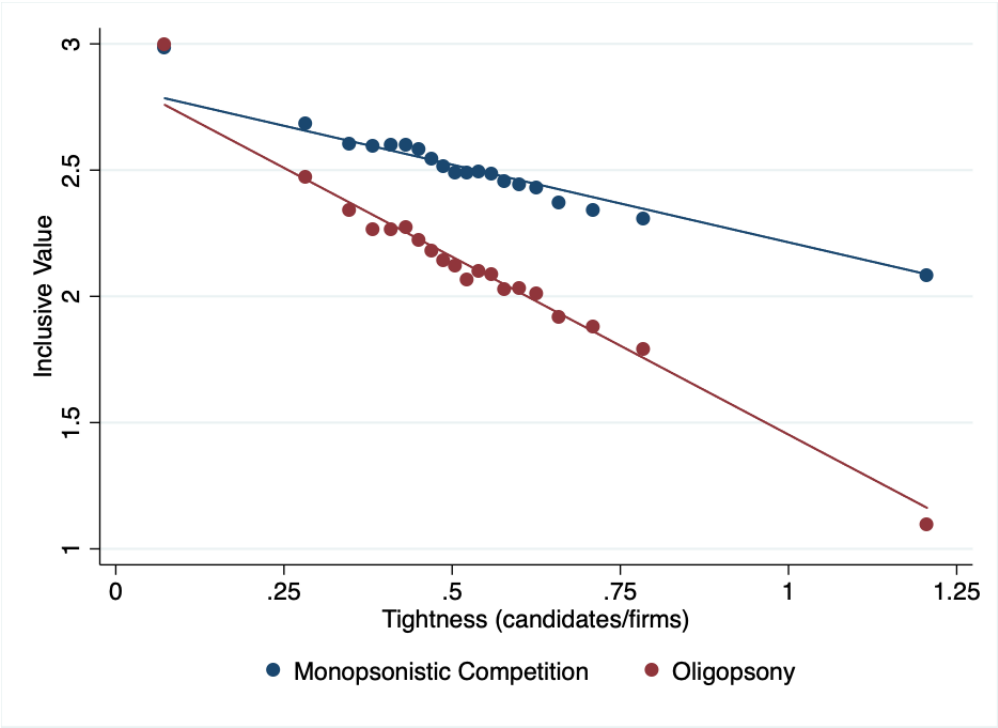


Note: This Figure plots the probability that a firm was rejected for a job-related reason as a function of firms’ ordinal rankings (where lower ranks are better) for the analysis sample. When a candidate receives a bid, she can decide to reject it, that is she can refuse to interview with the company. For a sub-sample (57%) of these rejections, candidates opted to provide a justification. They can choose from justifications such as “company size”, “insufficient compensation” or “company culture”. The latter is the justification we label as “bad company fit”. We plot the probability of rejection due to bad company fit against estimated rankings from the single-type model.

Figure 1.6: Differentiation between and within firms**(a)** Vertical Differentiation**(b)** Horizontal Differentiation

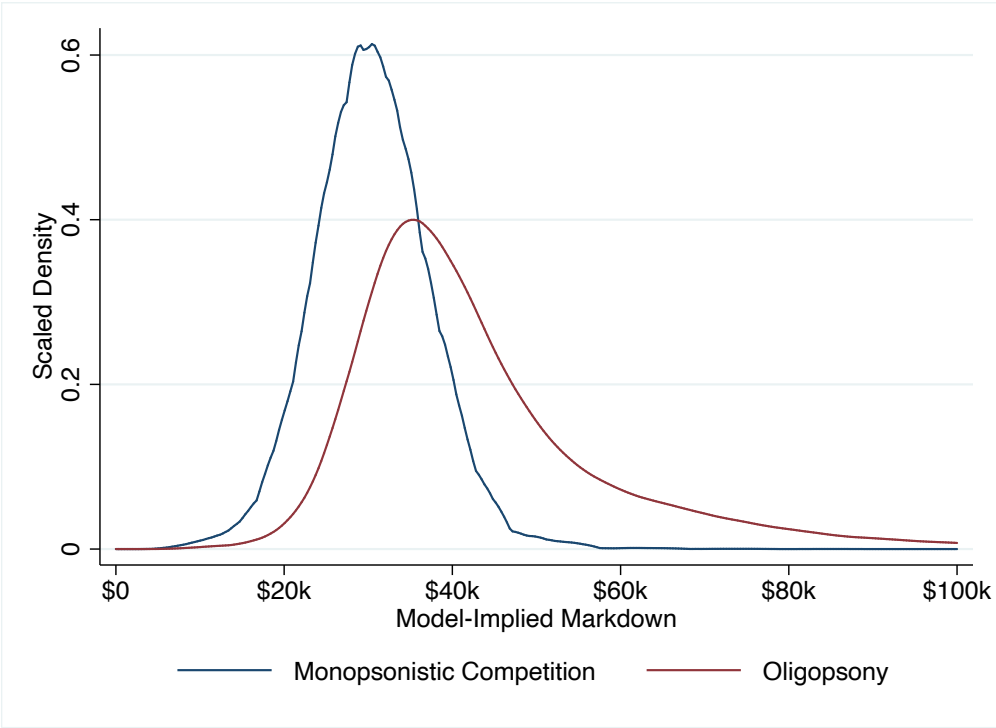
Note: This Figure illustrates the scale of vertical and horizontal differentiation of firms implied by our preferred model estimates. Willingness to Accept (WTA) is equal to the fraction of a candidate's ask salary that the model implies a firm must offer to make that candidate indifferent between accepting or rejecting an interview request, on average. Panel (a) plots the distribution of the mean Willingness to Accept (WTA) at each firm, averaging over the population probabilities of each type. Panel (b) illustrates the systematic component of horizontal differentiation, plotting the distribution of the within-firm standard-deviation of (WTA) across preference types.

Figure 1.7: First Stage



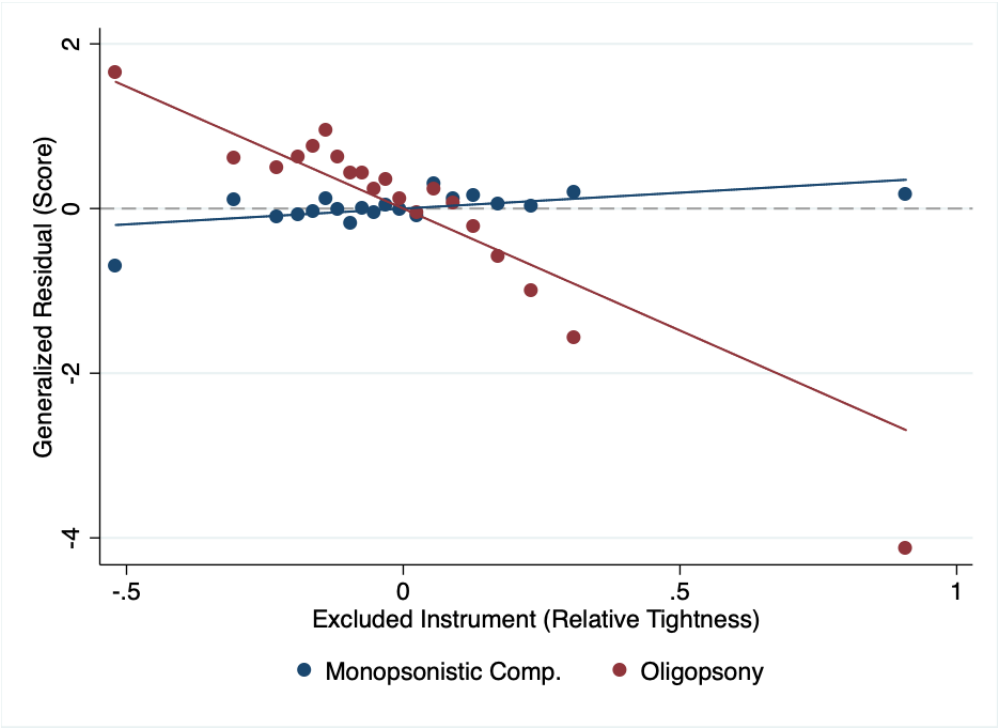
Note: This figure plots the “first stage” relationship between the model-implied inclusive values Λ_i and Λ_i^{-j} and the instrumental variable t_{ij} , conditional on firm covariates z_j and candidate covariates x_i and two-week period dummies.

Figure 1.8: Predicted Markdowns



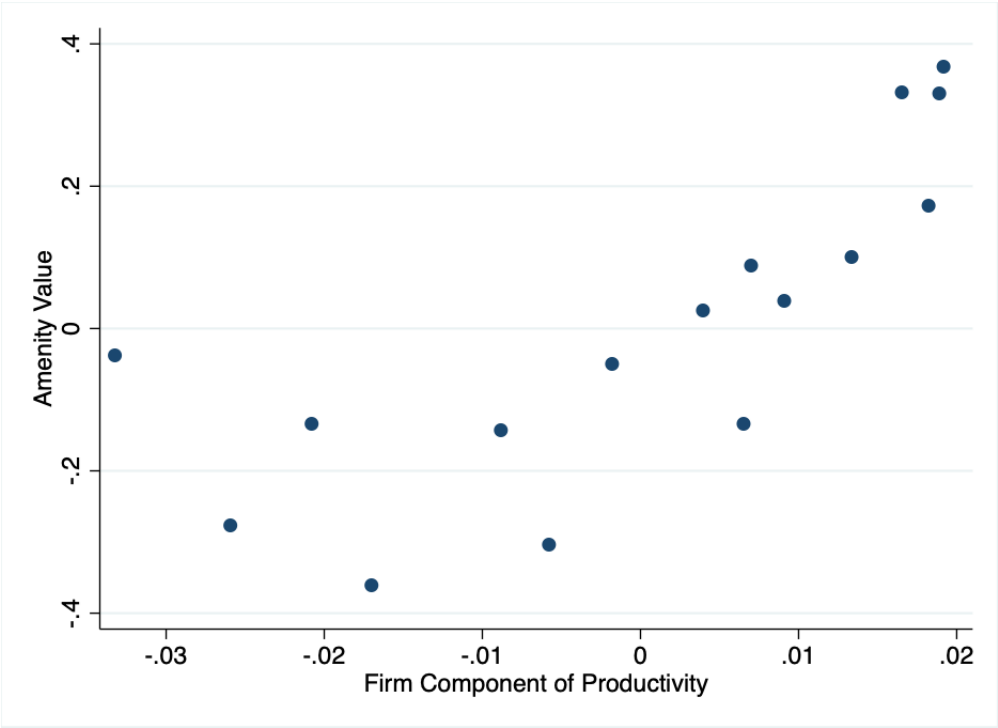
Note: This Figure plots the distribution of predicted markdowns under monopsonistic competition and oligopsony alternatives (in both cases, assuming firms are not type-predictive). For observations with bid equal to ask, we take the midpoint of the possible range of markdowns: $(\varepsilon_{ij}^+ + \varepsilon_{ij}^-)/2 - a_i$.

Figure 1.9: Visualizing the Vuong Test



Note: This Figure plots the relationship between generalized residuals and the excluded instrument (labor market tightness) for the non-predictive monopsonistic competition and oligopsony models. Under proper specification, the correlation of the generalized residuals and the excluded instrument should be zero (the dashed line). The larger the deviation from zero, the greater the degree of mis-specification of the model.

Figure 1.10: Relationship between Productivity and Amenity Values



Note: This Figure plots regression-adjusted measures of the average firm component of amenity values against the average firm component productivity for 16 categories of firms defined by combinations of firm size and industry. We compute regression-adjusted firm-type averages as the coefficients on a set of fixed effects in bid-level regressions of model-implied amenity and productivity values on $\log(\text{ask})$, an exhaustive set of fixed effects for combinations of other worker characteristics x_i , and dummies for firm type.

Tables

Table 1.1: Summary Statistics for Candidate Characteristics

| Variable (mean) | (1) All ($n = 43630$) | (2) Female (19%) | (3) Male (81%) |
|---------------------|----------------------------|---------------------|-------------------|
| Salary | | | |
| Ask/Expectation | \$137k | \$126k | \$140k |
| Education | | | |
| Has a BA+ | 0.872 | 0.913 | 0.862 |
| Has an MA+ | 0.403 | 0.437 | 0.395 |
| Has a CS degree | 0.629 | 0.558 | 0.645 |
| Attended an IvyPlus | 0.154 | 0.185 | 0.147 |
| Work History | | | |
| Years of experience | 11.3 | 10.1 | 11.6 |
| Software engineer | 0.684 | 0.512 | 0.724 |
| Worked at a FAANG | 0.108 | 0.097 | 0.111 |
| Employed | 0.748 | 0.719 | 0.755 |

Note: This Table reports summary statistics for the subset of candidates in the connected set, in particular, candidates' posted ask salary, education and previous work history. We report statistics both pooled and by gender. Previous work history is reported in years, ask/expectation salary in dollars, and all other statistics in percentages.

Table 1.2: Summary Statistics: Job Search and Job Finding

| Variable | (1) Full Sample | (2) Analysis Sample | (3) Connected Set |
|---------------------------------------|--------------------|------------------------|----------------------|
| Company Side | | | |
| Number of companies | 7,877 | 2,121 | 1,649 |
| Number of jobs | 64,539 | 16,907 | 13,072 |
| Number of interview requests sent | 856,665 | 267,940 | 124,075 |
| Average number of bids sent | 13.3 | 15.8 | 9.5 |
| Median number of bids sent | 5.0 | 6.0 | 4.0 |
| Candidate side | | | |
| Number of candidates | 224,499 | 44,321 | 14,344 |
| Average number of bids received | 3.5 | 4.1 | 4.8 |
| Probability of accepting a bid (in %) | 60.2 | 62.5 | 56.4 |

Note: This Table reports summary statistics for three increasingly restrictive samples of the data. The full sample includes the universe of entries on the platform. The analysis sample contains all candidates who had been contacted by a job that listed SF as the job location. The connected set includes all companies that can be ranked. The average and median number of bids sent statistics are calculated within job.

Table 1.3: Candidate Preference Model Goodness-of-Fit

| | | (1) | (2) | (3) |
|-----------------------|--------|-----------------|---------------------|----------------------|
| | | Split on Gender | Split on Experience | Model-Based Clusters |
| One | Log. L | -43,463 | -45,184 | -47,155 |
| Ladder | GOF | 0.672 | 0.673 | 0.677 |
| Two | Log. L | -42,962 | -44,535 | -45,558 |
| Ladders | GOF | 0.680 | 0.684 | 0.744 |
| | p(2,1) | 0.271 | <0.001 | <0.001 |
| Three | Log. L | - | - | -44,594 |
| Ladders | GOF | - | - | 0.779 |
| | p(3,2) | - | - | <0.001 |
| Four | Log. L | - | - | -43,857 |
| Ladders | GOF | - | - | 0.808 |
| | p(4,3) | - | - | >0.999 |
| Number of Firms | | 975 | 1,128 | 1,649 |
| Number of Candidates | | 13,658 | 13,830 | 14,344 |
| Number of Comparisons | | 209,934 | 222,935 | 235,827 |

Note: This Table reports goodness-of-fit (GOF) measures and p -values to adjudicate between labor supply models with different numbers of ladders (rows). Each column represents a different way to split candidates into preference types. The GOF statistic is calculated as the fraction of pairwise comparisons correctly predicted by the model, $\mathbb{E}[(\hat{A}_{qj} > \hat{A}_{qk}) \times (j \succ_i k)]$, and p -values are calculated via the likelihood ratio. Each column corresponds to a different sample determined by (overlapping, if relevant) connected sets.

Table 1.4: Which Firm Characteristics are Correlated with Amenity Values?

| | (1) \hat{A}_{1j} | (2) \hat{A}_{2j} | (3) \hat{A}_{3j} |
|------------------|-----------------------|-----------------------|-----------------------|
| Year Founded | 0.00521 (0.00374) | 0.00641 (0.00385) | -0.00502 (0.00358) |
| 15-50 Employees | -0.0836 (0.0881) | 0.114 (0.0907) | 0.105 (0.0843) |
| 50-500 Employees | -0.0531 (0.0829) | 0.222** (0.0853) | 0.337*** (0.0793) |
| 500+ Employees | -0.00169 (0.0993) | 0.287** (0.102) | 0.640*** (0.0950) |
| Finance | 0.0153 (0.0694) | 0.0474 (0.0715) | -0.105 (0.0664) |
| Tech | -0.0179 (0.0567) | -0.0312 (0.0584) | -0.0594 (0.0543) |
| Health | 0.0174 (0.0911) | 0.117 (0.0938) | -0.0778 (0.0872) |
| adj. R^2 | -0.004 | 0.009 | 0.085 |
| N | 913 | 913 | 913 |

Note: This Table reports regressions of standardized estimates of firm amenity values, \hat{A}_{qj} , on basic firm characteristics z_j . The omitted category for the number of employees is 0-15. Standard errors in parentheses, constant not reported. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Table 1.5: Oaxaca-Blinder Decompositions of Components of Utility

| | (1) Quantity | (2) | (3) Composition | (4) | (5) | (6) Inclusive Value | (7) | (8) |
|---------------------------|---------------------|------------------------|----------------------|----------------------|----------------------|------------------------|----------------------|----------------------|
| | # Bids | $\bar{u}(b_{ij}, a_i)$ | \bar{A}_j | ΔA_{ij} | Λ_i^b | $\bar{\Lambda}_i^A$ | $\Delta \Lambda_i^A$ | Λ_i^* |
| <i>Panel A: Gender</i> | | | | | | | | |
| Men | 4.854*** (0.037) | -0.265*** (0.006) | 0.336*** (0.002) | 0.006*** (0.001) | 0.806*** (0.009) | 1.434*** (0.006) | 0.003*** (0.001) | 1.199*** (0.010) |
| Women | 4.348*** (0.069) | -0.667*** (0.012) | 0.315*** (0.003) | 0.014*** (0.002) | 0.332*** (0.018) | 1.348*** (0.012) | 0.017*** (0.002) | 0.721*** (0.020) |
| Difference | 0.507*** (0.078) | 0.402*** (0.014) | 0.021*** (0.003) | -0.008*** (0.002) | 0.474*** (0.021) | 0.085*** (0.013) | -0.013*** (0.002) | 0.478*** (0.022) |
| Endowments | 0.577*** (0.045) | 0.151*** (0.009) | 0.018*** (0.002) | 0.025*** (0.002) | 0.243*** (0.013) | 0.111*** (0.008) | 0.024*** (0.002) | 0.287*** (0.015) |
| Coefficients | -0.083 (0.074) | 0.242*** (0.013) | 0.007* (0.004) | -0.033*** (0.002) | 0.215*** (0.020) | -0.026* (0.013) | -0.037*** (0.002) | 0.181*** (0.021) |
| Interaction | 0.012 (0.044) | 0.010 (0.008) | -0.005* (0.002) | -0.001 (0.001) | 0.017 (0.012) | 0.001 (0.008) | -0.000 (0.001) | 0.010 (0.012) |
| <i>Panel B: Education</i> | | | | | | | | |
| No Grad School | 4.943*** (0.045) | -0.478*** (0.007) | 0.320*** (0.002) | 0.000 (0.001) | 0.596*** (0.011) | 1.424*** (0.007) | -0.004*** (0.001) | 0.969*** (0.012) |
| Grad School | 4.489*** (0.046) | -0.140*** (0.007) | 0.349*** (0.002) | 0.017*** (0.001) | 0.892*** (0.012) | 1.408*** (0.008) | 0.020*** (0.001) | 1.312*** (0.013) |
| Difference | 0.454*** (0.065) | -0.338*** (0.010) | -0.029*** (0.003) | -0.017*** (0.002) | -0.296*** (0.016) | 0.016 (0.011) | -0.023*** (0.002) | -0.343*** (0.017) |
| Endowments | -0.039 (0.041) | -0.101*** (0.007) | -0.017*** (0.002) | 0.001 (0.001) | -0.132*** (0.011) | -0.047*** (0.007) | -0.001 (0.001) | -0.149*** (0.012) |
| Coefficients | 0.554*** (0.071) | -0.137*** (0.010) | -0.001 (0.003) | -0.013*** (0.002) | -0.057*** (0.017) | 0.080*** (0.012) | -0.017*** (0.002) | -0.073*** (0.018) |
| Interaction | -0.062 (0.053) | -0.100*** (0.008) | -0.011*** (0.002) | -0.005*** (0.001) | -0.107*** (0.013) | -0.016 (0.009) | -0.006*** (0.001) | -0.121*** (0.014) |
| <i>N</i> | 38,231 | 38,231 | 38,231 | 38,231 | 38,231 | 38,231 | 38,231 | 38,231 |

Note: This Table reports Oaxaca-Blinder decompositions of components of utility. Panel A reports decompositions by gender. Panel B reports decompositions by education. Column 1 decomposes the gap in the number of bids. Column 2 decomposes the mean gap in the monetary component of utility. Column 3 decomposes the mean difference in the common component of amenity values. Column 4 decomposes differences in candidate-specific components of the amenity valuation. Columns 5-8 decompose components of the inclusive value. The Endowments, Coefficients, and Interaction rows sum to the Difference row in every column. Robust standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 1.6: Non-Nested Model Comparison Tests

| Model | (1) | (2) | (3) | (4) |
|---|----------------|-----------------|----------------|-----------------|
| | Monopsonistic | | Oligopsony | |
| | Not Predictive | Type Predictive | Not Predictive | Type Predictive |
| <i>Panel A: Likelihood-Based Test (Vuong (1989))</i> | | | | |
| Perfect Competition | -237.57 | -237.67 | -156.16 | -154.34 |
| Monopsonistic, Not Predictive | – | 1.28 | 90.17 | 90.39 |
| Monopsonistic, Type Predictive | | – | 88.45 | 89.81 |
| Oligopsony, Not Predictive | | | – | 6.88 |
| Oligopsony, Type Predictive | | | | – |
| <i>Panel B: Moment-Based Test (Rivers and Vuong (2002))</i> | | | | |
| Perfect Competition | -54.84 | -54.40 | -39.92 | -39.91 |
| Monopsonistic, Not Predictive | – | 7.83 | 3.98 | 2.69 |
| Monopsonistic, Type Predictive | | – | 2.77 | 1.54 |
| Oligopsony, Not Predictive | | | – | -3.67 |
| Oligopsony, Type Predictive | | | | – |

Note: This Table reports test statistics from the Vuong (1989) non-nested model comparison procedure. We implement the testing procedure for each pair of the five models we estimated, using both the likelihood-based test (Panel A) and the moment-based test (Panel B). Positive values imply the row model is preferred to the column model. Under the null of model equivalence, the test statistics are asymptotically normal with mean zero and unit variance.

Table 1.7: (Subset of) Labor Demand Parameters Γ : $\log(\varepsilon_{ij}) = z_j' \Gamma x_i + \nu_{ij}$

| | (1) | (2) | (3) | (4) | (5) |
|--|------------------------|-----------------------|---------------------------------------|---------------------|------------------------|
| | Constant | log(Ask) | Female | Employed | Grad School |
| Constant | 11.9897*** (0.0523) | 0.7954*** (0.0046) | -0.0079*** (0.0025) | -0.0014 (0.0040) | 0.0094*** (0.0021) |
| 16-50 Employees | 0.0305 (0.0448) | 0.0814*** (0.0039) | 0.0046 (0.0027) | 0.0006 (0.0044) | -0.0022 (0.0023) |
| 51-500 Employees | 0.0503 (0.0510) | 0.0832*** (0.0045) | -0.0010 (0.0025) | 0.0037 (0.0041) | -0.0069*** (0.0022) |
| 501+ Employees | 0.0612 (0.0516) | 0.1073*** (0.0045) | -0.0009 (0.0026) | 0.0011 (0.0043) | -0.0090*** (0.0022) |
| Finance | -0.0008 (0.0526) | 0.0156*** (0.0046) | 0.0055*** (0.0016) | 0.0024 (0.0028) | 0.0022 (0.0013) |
| Tech | 0.0052 (0.0314) | 0.0166*** (0.0027) | 0.0043*** (0.0013) | -0.0028 (0.0023) | -0.0001 (0.0011) |
| Health | -0.0028 (0.0462) | 0.0011 (0.0040) | 0.0009 (0.0022) | -0.0006 (0.0037) | -0.0004 (0.0017) |
| Std. Dev. of ν_{ij} ($\hat{\sigma}_\nu$) | 0.0743 | (0.0001) | $N = 181,927$, Implied $R^2 = 0.888$ | | |

Note: This table reports a subset of maximum likelihood parameter estimates from our preferred model. The parameters relate combinations of candidate and firm characteristics to the distribution of firms' valuations over each candidate (or, the ex-ante productivity of that candidate at that firm). The log of productivity/valuations is modelled as normally distributed, with mean $z_j' \Gamma x_i$ and variance σ_ν . Each cell reports the coefficient on the interaction of the variables specified in the corresponding row and column. Column variables are candidate characteristics (x_i), and row variables are firm characteristics (z_j). The second, third, and fourth rows correspond to dummies for firm size categories, such that the omitted category (subsumed into the constant, the first row of the table) corresponds to the smallest firms (between one and fifteen employees). The remaining three rows correspond to non-exclusive sector dummies. Column 1 reports the main effects of each firm characteristic. Column 2 reports the main effects and interactions for the log ask salary, where the log ask salary has been de-measured. Columns 3-5 report coefficients on dummies recording whether the candidate is female, was employed, or has received at least a master's degree. Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 1.8: Variance Decomposition of Bids

| | (1) | (2) | (3) | (4) |
|--|----------------|------------------|-------------------|------------|
| | Bids | Markdowns | Valuations | |
| | $\log(b_{ij})$ | $\log(\mu_{ij})$ | $z'_j \Gamma x_i$ | ν_{ij} |
| <i>Panel A: Monopsonistic Competition</i> | | | | |
| $\log(b_{ij})$ | 1.000 | -0.001 | 0.910 | 0.091 |
| $\log(\mu_{ij})$ | | 0.03 | 0.007 | -0.011 |
| $z'_j \Gamma x_i$ | | | 0.897 | 0.006 |
| ν_{ij} | | | | 0.097 |
| <i>Panel B: Oligopsony</i> | | | | |
| $\log(b_{ij})$ | 1.000 | 0.101 | 0.777 | 0.122 |
| $\log(\mu_{ij})$ | | 0.133 | 0.080 | -0.113 |
| $z'_j \Gamma x_i$ | | | 0.680 | 0.016 |
| ν_{ij} | | | | 0.219 |
| Standard Deviation of $\log(b_{ij}) = 0.221$. | | | | |

Note: This Table describes the variance decomposition of log bids. Each cell reports the covariance of the row and column variables, standardized (divided) by the overall variance of log bids. Panel A is computed using estimates from the preferred model, monopsonistic competition/not predictive conduct. Panel B is computed using the dis-preferred oligopsony/type-predictive conduct model

Table 1.9: Counterfactual Simulations

| | (1) | (2) | (3) | (4) | (5) | (6) |
|--|----------------|--------|--------|-----------------|--------|--------|
| <i>Panel A: Unconditional Means</i> | | | | | | |
| | Not Predictive | | | Type-Predictive | | |
| | PT | MC | OG | PT | MC | OG |
| Bid, b_{ij} | \$169k | \$145k | \$139k | \$169k | \$145k | \$139k |
| Ratio of Bid/Ask, b_{ij}/a_i | 1.196 | 1.024 | 0.979 | 1.196 | 1.025 | 0.978 |
| Markdown, $1 - b_{ij}/\varepsilon_{ij}$ | 0.099 | 0.182 | 0.182 | 0.099 | 0.183 | 0.183 |
| # Bids Received/Candidate | 20.1 | 43.2 | 13.5 | 19.6 | 42.0 | 13.2 |
| Inclusive Value, Λ_i^* | 0.930 | 0.886 | 0.822 | 0.932 | 0.888 | 0.822 |
| Monetary Component, Λ_i^b | 0.033 | 0.015 | 0.000 | 0.033 | 0.016 | 0.000 |
| Common Amenity Comp., $\bar{\Lambda}_i^A$ | 0.282 | 0.357 | 0.315 | 0.281 | 0.355 | 0.314 |
| Type-Specific Amenity Comp., $\Delta\Lambda_i^A$ | 0.002 | 0.004 | 0.004 | 0.005 | 0.008 | 0.007 |
| <i>Panel B: Differences, Women - Men</i> | | | | | | |
| | Not Predictive | | | Type-Predictive | | |
| | PT | MC | OG | PT | MC | OG |
| # Bids Received/Candidate | -1.830 | -3.793 | -1.434 | -2.411 | -5.681 | -2.529 |
| Inclusive Value, Λ_i^* | -0.053 | -0.069 | -0.019 | -0.056 | -0.070 | -0.019 |
| Monetary Component, Λ_i^b | -0.026 | -0.052 | -0.016 | -0.027 | -0.051 | -0.016 |
| Common Amenity Comp., $\bar{\Lambda}_i^A$ | -0.003 | -0.005 | -0.003 | -0.004 | -0.007 | -0.004 |
| Type-Specific Amenity Comp., $\Delta\Lambda_i^A$ | 0.005 | 0.010 | 0.013 | 0.003 | 0.010 | 0.011 |
| <i>Panel C: Differences, Women - Men, Gender Blind Firms</i> | | | | | | |
| | Not Predictive | | | Type-Predictive | | |
| | PT | MC | OG | PT | MC | OG |
| # Bids Received/Candidate | -1.652 | -3.749 | -1.529 | -2.776 | -6.162 | -2.549 |
| Inclusive Value, Λ_i^* | -0.050 | -0.066 | -0.018 | -0.053 | -0.068 | -0.019 |
| Monetary Component, Λ_i^b | -0.025 | -0.051 | -0.016 | -0.027 | -0.050 | -0.016 |
| Common Amenity Comp., $\bar{\Lambda}_i^A$ | -0.003 | -0.005 | -0.003 | -0.002 | -0.006 | -0.002 |
| Type-Specific Amenity Comp., $\Delta\Lambda_i^A$ | 0.004 | 0.011 | 0.013 | 0.005 | 0.009 | 0.011 |

Note: This Table reports results of counterfactual simulations under various conduct assumptions. Columns labelled PT refer to the price-taking model of conduct, columns labelled MC refer to the monopsonistic competition model of conduct, and columns labelled OG refer to the oligopsony model of conduct. Each cell reports the average of the statistic over 50 simulation draws. In each simulation draw, we sample from the distribution of valuations for a set of 500 firms considering 500 workers (a single sample of workers and firms is used for all simulations). Panel A reports the unconditional means of various statistics. Panel B reports differences in means between women and men. Panel C reports differences in means between women and men for simulations in which firms are constrained to be gender blind.

Chapter 2

A Framework for Design-Based Inference of Many Treatment Effects

2.1 Introduction

Social scientists are often interested in estimating the distribution of an unobserved, or partially-observed, attribute among a population. In economics, latent attributes of interest are commonly measures of productivity or preferences that are allowed to vary across units (e.g. people, firms, institutions). Differences in productivity or preferences across units implies variation in the counterfactual potential outcomes of individuals assigned to those units, such that each unit is associated with a unique treatment effect. Analysis of unit-specific treatment effects reflects the basic fact that treatments of interest to social scientists are rarely delivered uniformly: for instance, economists cannot (yet) manufacture pills to deliver uniform doses of human capital to students in a randomized controlled trial. Instead, human capital is “delivered” to students by individual teachers, each of whom may vary in skill or teaching practices. When policymakers say they want to improve the quality of education, for instance, they are necessarily referring to the productivity of the teachers who actually do the work of educating students. In order to craft policies that improve the quality of education, then, it is important to understand the extent to which teachers differ, if at all, in their ability to they have on their students (see, e.g. Chetty, Friedman, and Rockoff 2014). The same logic applies when considering policies to improve health outcomes when patients are treated by individual doctors, or when considering policies to reduce bias in legal outcomes when cases are handled by individual judges, among other examples.

There are (broadly) two alternative modes of analysis social scientists have adopted when they wish to quantify variation in latent treatment effects across units. The first alternative is what might be called a direct, structural, or supervised approach, in which researchers specify a list of measurable characteristics of units and assume that latent attributes that determine treatment effects are shared within groups of units that share those characteristics. The effect of this assumption is to partially reveal those latent attributes, allowing for direct estimation of treatment effects that compare individuals assigned to units in each group. For instance, we might assume that teacher productivity is a function of teacher experience and education,

such that all teachers with the same combinations of experience and education have the same mean productivity. This approach is clearly attractive if appropriate assumptions can be maintained, since it allows researchers to pool data across units to estimate treatment effects precisely.

However, the direct approach has clear drawbacks if observable proxies for latent attributes are incomplete, difficult to measure, or only weakly correlated with the latent attribute of interest. In that case, the direct approach may at best deliver an uninformative lower bound on the variance of treatment effects. At worst, the direct approach may deliver parameter estimates polluted by omitted variables bias, since it is generally not possible to know whether measurable characteristics of units are correlated meaningfully with other unobserved characteristics, and how those characteristics combine to determine the outcome of interest. Importantly, this is true even under pure random assignment of individuals to units, since identification of the mechanisms that drive productivity requires random assignment of *attributes* across units (e.g. random assignment of training to teachers, such that training is independent of other teacher attributes). Further, it is often extremely burdensome, if not impossible, to measure all of the possible characteristics that may be relevant to the latent attribute.

The second alternative is an outcome-oriented, reduced-form, or unsupervised approach, in which researchers specify a statistical model that relates the distribution of outcomes to the latent attribute in which no structural relationship between the unobserved attribute and measurable inputs is assumed – in other words, each unit is allowed to have a completely unique value of the latent attribute. The outcome-oriented approach circumvents the problem of direct measurement of inputs by inferring variation in those inputs through modeling the distribution of readily-measurable outputs. One could think of the this approach as measuring the variation in a sufficient statistic for the outcome that collapses all possible latent factors of units. In some cases, this is a strength of the approach, since the outcome that is measured is one that may have direct welfare implications (for instance, in the analysis of variation in physician “productivity” as measured by patient mortality). In other applications, the lack of a structural interpretation is a weakness. Because there are typically many more units than there are measurable attributes of those units, moving from the direct approach to the outcome-oriented approach is fundamentally about accepting increased uncertainty around estimates of variation in attributes in exchange for a reduction in the bias of those estimates.

This paper is concerned with the second approach to estimation of distributions of attributes, which has grown increasingly popular among economists. The popularity of the outcome-oriented approach is in part due to the fact that it requires placing fewer restrictions on the data generating process than the direct approach: both require the absence of systematic sorting of individuals to units based on unobserved factors, but the direct approach additionally requires that measurable inputs are randomly assigned to units, or that all relevant inputs can be measured. The assumption that there is no sorting of individuals to units on the basis of unobserved factors is key for identification of unit-specific attributes, and is typically formalized as an assumption of unconfoundedness: conditional on a set of observable variables, the actual assignment of individuals to units is independent of indi-

viduals’ vector of potential outcomes under counterfactual assignment to each unit. This assumption is usually motivated by institutional knowledge of the assignment process.

In practice, however, researchers who adopt the outcome-oriented approach almost always apply additional functional form and distributional assumptions which may not be motivated by any particular knowledge of the data-generating process. For instance, it is common to model student test scores as the sum of a linear combination of measured confounds, a teacher effect, and an idiosyncratic individual effect, where the teacher and individual effects are assumed to be drawn from normal distributions (this is the basic modeling assumption of Chetty, Friedman, and Rockoff (2014), for instance). By specifying a complete model of the data generating process, the inferences these studies make about the distribution of treatment effects are in a sense *model-based* (Card 2012; Sterba 2009). Under correct specification, model-based estimators of treatment effects can perform optimally (in the sense of bias and variance) relative to alternative estimators. However, a mis-specified outcome model may lead to severely biased estimates of treatment effects.

This paper draws upon the key insight that the modeling assumptions used in studies that adopt the outcome-oriented approach are unnecessary when researchers have already assumed that assignment of individuals to units is unconfounded, and the object of interest is a distribution of average treatment effects across units. Indeed, if assignment to units is unconfounded conditional on a set of measured variables X , all one needs to estimate the average treatment effect of each unit is the probability of assignment to treatment given X – the propensity score (Rosenbaum and Rubin 1983). Given a model of the propensity score, there is no need to specify a parametric model for the outcome. Rather, taking averages of each unit’s outcomes weighted by the inverse of the propensity score will deliver estimates of the average treatment effects associated with each unit (e.g. Hirano, Imbens, and Ridder 2003). In contrast to the model-based approach to causal inference outlined above, this is a fundamentally *design-based* approach to causal inference: the propensity score fully summarizes the sampling design that gave rise to the observed data (Card 2012; Sterba 2009).

There are advantages to the design-based approach relative to the model-based approach beyond eschewing potentially mis-specified functional form and distributional assumptions. By specifying the sampling mechanism, the design-based approach formalizes both the notion of which causal quantity is the target parameter of interest (Rubin, Stuart, and Zanutto 2004) and the population of individuals who are actually at risk for assignment to each treatment. In so doing, the design-based approach makes explicit researcher’s assumptions about the extent to which units are reliably comparable without out-of-sample extrapolation, while those assumptions remain implicit (and often unrecognized) in studies that adopt the model-based approach to causal inference. Despite the advantages of design-based approaches to causal inference about distributions of treatment effects over model-based approaches, design-based approaches are rarely used in these settings (a notable recent exceptions is Angrist et al. 2020).

This paper proposes a framework for design-based inference for many unit-specific treatment effects and the distribution from which those treatment effects are drawn. Importantly, for the purposes of this paper, “design-based inference” refers primarily to the method by

which causal quantities are estimated, and not necessarily the mode of statistical inference used to make probability statements about those casual quantities. In this sense, our use of the terms “design-based” and “model-based” mirror their use in Card (2012) and Card (2022): identification of the causal effect(s) of interest arises from manipulation of assignments of individuals to treatments. While the design-based approach offers many potential advantages over the model-based approach, implementing the design-based approach in a setting with a large number of treatments presents unique econometric challenges. First, the standard overlap assumptions that are maintained in settings where treatment is binary are very likely to fail in settings where the number of treatments is large. Second, the dimensionality of the problem is greatly increased relative to the binary treatment case. In order to grapple with the dual issues of limited overlap and high dimensionality, we develop a regularized propensity score estimator that allows for structural failure in overlap. We then use tools from the literature on matrix completion to analyze the properties of this estimator. Given estimates of the assignment model, we then propose a novel sample trimming routine that selects the largest subset of the sample for which a traditional notion of overlap is likely to hold. Finally, we illustrate a method for estimating the distribution from which treatment effects are drawn via an inverse propensity score weighted nonparametric maximum likelihood routine.

2.2 Setup

Variables and Notation

Consider collecting data on a sample $i = 1, \dots, N$ of the form:

$$Y_i \in \{0, 1\}, D_i = \left(D_{ij} \in \{0, 1\} \right)_{j=1}^J, S_i = \left(S_{ir} \in \{0, 1\} \right)_{r=1}^R, X_i' \in \mathbb{R}^K.$$

Here, we have assumed that Y_i is a binary outcome of interest, although the analysis below generalizes to other types of outcomes. D_i is a J -dimensional vector that encodes the assignment of individual i to unit (treatment arm) j . S_i is an R -dimensional vector encoding the assignment of individual i to “randomization stratum” r . The variable S_i is an exhaustive discretization of the full set of measured confounds that determine the probability of assignment to treatment. Assuming S_i takes on a finite number of values is typically without loss of generality, since the variables that govern assignment probabilities in usual applications are themselves discrete (e.g. age, location, etc.). Finally, X_i is a K -dimensional vector of pre-determined characteristics of unit i . The ultimate quantity of interest is the set of treatment effects associated with assignment of D_i to levels $j = 1, \dots, J$ on outcome Y_i . The observed variables X_i are a subset of all factors aside from D_i and S_i that affect the outcome Y_i . We denote all unobserved factors by X_i^u and denote the union of these sets of random variables by $U_i = \{X_i, X_i^u\}$. Both S_i and U_i are “pre-treatment” variables, in that they are determined before treatment assignments are made, and cannot be affected by treatment assignments.

Potential Outcomes Model

We next specify a framework for defining causal treatment effects along the lines of (Rosenbaum and Rubin 1983). Under the Stable Unit Treatment Value Assumption (SUTVA), we may write the potential outcome of individual i assigned to unit j as:

$$Y_{ij}^* = Y_j^*(S_i, U_i),$$

such that the observed outcome can be written:

$$Y_i = \sum_{j=1}^J D_{ij} Y_{ij}^*.$$

We next state the first of two sets of assumptions necessary for identification of treatment effects:

Assumption 2.1. (Unconfoundedness)

a) *Conditional on S_i , assignment to treatment is independent of U_i :*

$$U_i \perp D_i \mid S_i.$$

b) *Potential outcomes are iid Bernoulli conditional on U_i and S_i :*

$$Y_j^*(S_i, U_i) \stackrel{iid}{\sim} \text{Bernoulli}(p_j(U_i, S_i)).$$

Assumption 2.1a is sometimes called selection-on-observables: conditional on randomization stratum, assignments of individuals to treatments are orthogonal to the remaining determinants of outcomes (U_i). Assumption 2.1b states that individuals in the same stratum S_i and with the same characteristics U_i face identical potential outcome distributions. Together, these two assumptions imply that assignment to treatment is *weakly unconfounded* (Imbens 2000):

Definition 2.1. (Weak Unconfoundedness) *Assignment to treatment is weakly unconfounded, given randomization stratum S_i , if:*

$$Y_{ij}^* \perp D_i \mid S_i.$$

What does assumption 2.1 imply? Denote the distribution of U_i conditional on S_i by:

$$U_i \mid S_i = e_r \sim F_{U|S}(\cdot \mid e_r),$$

where e_r is the r -th standard basis vector. Next, define:

$$p_j(e_r) = \int p_j(u, e_r) dF_{U|S}(u \mid e_r).$$

Given this notation, we have:

$$\begin{aligned} Y_i \mid D_{ij} = 1, S_{ir} = 1 &\stackrel{d}{=} Y_{ij}^* \mid D_{ij} = 1, S_{ir} = 1 \\ &\stackrel{d}{=} Y_{ij}^* \mid S_{ir} = 1 \\ &\sim \text{Bernoulli}(p_j(e_r)). \end{aligned}$$

The first equality follows from the definition of the potential outcomes model, while the second equality follows from assumption 2.1. Therefore, the quantities $p_j(e_r)$ are identified from observable data. If each unit was observed many times in each stratum, the unit j - stratum r cell mean would converge to $p_j(e_r)$. Finally, define:

$$p_j = \sum_{r=1}^R p_j(e_r) dF_S(e_r),$$

where $dF_S(e_r)$ is the unconditional distribution of individuals across randomization strata. The unconditional distribution of potential outcomes is therefore given by:

$$Y_{ij}^* \sim \text{Bernoulli}(p_j).$$

The distribution of p_j across the J treatment arms is the primary object of interest. Variation in p_j across units reflects differences in the average treatment effects of those units.

Next, we introduce the second set of assumptions necessary for identification of treatment effects. While Assumption 2.1 is relatively standard, these assumptions are non-standard:

Assumption 2.2. (Overlap)

a) Conditional on a latent variable $Z_{jr} \in \{0, 1\} \perp S_i$:

$$\begin{aligned} 0 < \Pr(D_{ij} = 1 \mid S_{ir} = 1, Z_{jr} = 1) < 1, \text{ and} \\ \Pr(D_{ij} = 1 \mid S_{ir} = 1, Z_{jr} = 0) = 0. \end{aligned}$$

b) Conditional on a latent variable $V_i = \sum_{r=1}^R S_{ir} V_r$, where $V_r \in \mathbb{R}^Q$, with $Q < \min(J, R)$:

$$\Pr(D_{ij} = 1 \mid S_{ir} = 1, V_i = v, Z_{jr} = z) = \Pr(D_{ij} = 1 \mid V_i = v, Z_{jr} = z)$$

Assumption 2.2a is a weakening of the typical overlap assumption, which requires $\Pr(D_{ij} = 1 \mid S_{ir} = 1) > 0$ for all j and r . The latent variables Z_{jr} encode whether individuals in stratum r could have been assigned to treatment j . When $Z_{jr} = 0$, then there is zero probability of assignment to j in stratum r , and overlap fails. When $Z_{jr} = 1$ then there is a positive probability of assignment to j in stratum r . Assumption 2.2b implies that the matrix of propensity scores (or an appropriate transformation of those propensity scores) is of low rank. In particular, all R of the J -dimensional vectors of strata-specific assignment probabilities are generated by taking combinations of a small number (Q) of baseline selection regimes. Further, assumption 2.2b states that randomization strata which combine the underlying selection regimes in the same way ($V_r = V_{r'}$) have identical vectors of propensity scores.

2.3 Estimation with a Known Assignment Mechanism

Consistency of the Fixed Effects Estimator

To estimate the full set of treatment effects, consider the following weighted likelihood:

$$L(\mathbf{p}) = \prod_{i=1}^N \prod_{j=1}^J \left(p_j^{Y_i} (1 - p_j)^{1 - Y_i} \right)^{w_{ij} D_{ij}},$$

where w_{ij} is a set of weights (to be specified). Importantly, the assumption that potential outcomes are distributed according to a Bernoulli distribution is a trivial implication of the assumption that the outcome Y_i is binary, and so specifying a full likelihood for the data does not impose any substantive restrictions. Maximizing the log-weighted-likelihood with respect to the full parameter vector $\mathbf{p} = (p_1, \dots, p_J)'$ yields the following fixed-effects estimator for each p_j :

$$\hat{p}_j^{\text{fe}} = \frac{\frac{1}{N} \sum_{i=1}^N w_{ij} D_{ij} Y_i}{\frac{1}{N} \sum_{i=1}^N w_{ij} D_{ij}}.$$

In this section, we assume that the assignment mechanism is known, by which we mean that Z_{jr} , V_r , and $\Pr(D_{ij} = 1 \mid V_i = v, Z_{ij} = 1)$ are known, where $Z_{ij} = \sum_{r=1}^R S_{ir} Z_{jr}$. When the assignment mechanism is known, and given a particular choice of weighting function, the \hat{p}_j^{fe} is a consistent estimator of the average treatment effect of unit j for the population of individuals in strata for which $Z_{jr} = 1$. Let \mathcal{O} denote a set of unit and strata indices (j and r) such that $Z_{jr} = 1$ for all $j, r \in \mathcal{O}$. Define the average treatment effect of unit j in subsample \mathcal{O} as:

$$p_j^{\mathcal{O}} = \frac{\sum_{r \in \mathcal{O}} p_j(e_r) dF_S(e_r)}{\sum_{r \in \mathcal{O}} dF_S(e_r)}.$$

Given these definitions, we now state the basic result:

Theorem 2.1. (Consistency of \hat{p}_j^{fe} , Known Assignment Mechanism) *Assume the conditions of assumptions 2.1 and 2.2 hold, and additionally assume that the assignment mechanism is known (such that all Z_{jr} , V_r , and $\Pr(D_{ij} = 1 \mid V_i = v, Z_{ij} = 1)$ are observed). Let \mathcal{O} denote a set of unit (j) and randomization strata (r) indices such that $Z_{jr} = 1$ for all $j, r \in \mathcal{O}$. Define the weighting function:*

$$w_{ij} = w_j(v) = \frac{\mathbf{1}[r(i) \in \mathcal{O}]}{\Pr(D_{ij} = 1 \mid V_i = v, Z_i = 1)},$$

where $r(i)$ returns the index of the randomization stratum of individual i . Then:

$$\hat{p}_j^{\text{fe}} \xrightarrow{P} p_j^{\mathcal{O}}.$$

Proof. Without loss of generality, we assume $Z_{jr} = 1$ for all j and r and suppress dependence on the overlap set \mathcal{O} (otherwise, we could select an overlap set and re-define all indices). We also suppress i subscripts for brevity. As N tends to infinity, the numerator of \hat{p}_j^{fe} converges

in probability to $\Pr(D = e_j) \times \mathbb{E}[w_j(V)Y \mid D = e_j]$, and the denominator converges to $\Pr(D = e_j) \times \mathbb{E}[w_j(V) \mid D = e_j]$. We have:

$$\begin{aligned} \mathbb{E}[w_j(V) Y \mid D = e_j] &= \\ &= \int_v \sum_{r=1}^R \mathbb{E}[w_j(v) Y \mid V = v, S = e_r, D = e_j] dF_{S|V,D}(e_r \mid v, e_j) dF_{V|D}(v \mid e_j). \end{aligned}$$

By Bayes' rule and Assumption 2, we may write:

$$dF_{S|V,D} = \frac{dF_{D|V,S} dF_{S|V}}{dF_{D|V}} = \frac{dF_{D|V} dF_{S|V}}{dF_{D|V}} = dF_{S|V}.$$

Additionally, we may write:

$$\mathbb{E}[w_j(v) Y \mid V = v, S = e_r, D = e_j] = w_j(v) \mathbb{E}[Y_j^* \mid V = v, S = e_r].$$

This implies:

$$\begin{aligned} \mathbb{E}[w_j(V) Y \mid D = e_j] &= \int_v w_j(v) \left[\sum_{r=1}^R \mathbb{E}[Y_j^* \mid V = v, S = e_r] dF_{S|V}(e_r \mid v) \right] dF_{V|D}(v \mid e_j) \\ &= \int_v w_j(v) \mathbb{E}[Y_j^* \mid V = v] dF_{V|D}(v \mid e_j) \\ &= \int_v \mathbb{E}[Y_j^* \mid V = v] dF_{V|D}(v \mid e_j) \frac{1}{dF_{D|V}(e_j \mid v)} \frac{dF_{D|V}(e_j \mid v) dF_V(v)}{dF_D(e_j)} \\ &= \Pr(D = e_j)^{-1} \times \int_v \mathbb{E}[Y_j^* \mid V = v] dF_V(v) \\ &= \Pr(D = e_j)^{-1} \times \mathbb{E}[Y_j^*]. \end{aligned}$$

The numerator clearly converges in probability to $\mathbb{E}[Y_j^*] = p_j$. Similarly, the denominator converges in probability to 1. And so, by Slutsky's Theorem:

$$\hat{p}_j^{\text{fe}} \xrightarrow{p} p_j.$$

□

Estimating the Distribution of Treatment Effects

Although the Maximum Likelihood fixed effects estimator of \mathbf{p} is consistent, there are several practical downsides -- namely, fixed effect estimates tend to have poor out-of-sample predictive power. In order to address this concern, we next consider estimation of the distribution of the p_j , which is denoted by $G(\cdot)$. For the purposes of this section, we again assume the weights w_{ij} and the Z_i are known, and implicitly condition on $Z_i = 1$ without loss of generality.

We now formally assume:

Assumption 2.3. (Random Sampling)

$$p_j \stackrel{iid}{\sim} G(\cdot).$$

Under assumption 2.3, moments of up to order J are identified (as in Kline and Walters 2021). In particular, we may write:

$$\mathbb{E} \left[\frac{1}{J} \sum_{j=1}^J \left(\widehat{p}_j^{\text{fe}} \right)^m \right] = \sum_{i=0}^m \omega_{im} \mu_i,$$

where the ω_{im} are known constants that are functions of the weights w_{ij} . In particular, the identified moments of $G(\cdot)$ can be used to construct an unbiased estimate of the variance of treatment effects across units:

$$\widehat{\text{Var}}(p_j) = \frac{1}{J-1} \sum_{j=1}^J \left(\widehat{p}_j^{\text{fe}} - \frac{1}{J} \sum_{j=1}^J \widehat{p}_j^{\text{fe}} \right)^2 - \frac{1}{J} \sum_{j=1}^J \text{SE} \left(\widehat{p}_j^{\text{fe}} \right),$$

where $\text{SE} \left(\widehat{p}_j^{\text{fe}} \right)$ is an unbiased estimate of the standard error of $\widehat{p}_j^{\text{fe}}$. It can be shown that $\widehat{\text{Var}}(p_j)$ is a U-statistic, and therefore:

$$\sqrt{J} \left(\widehat{\text{Var}}(p_j) - \text{Var}(p_j) \right) \rightsquigarrow N(0, V_\sigma/J).$$

We can construct an unbiased estimate of V_σ following the method of Wang and Lindsay (2014), and use that estimate to test against the null hypothesis of no heterogeneity.

Beyond estimating individual moments of the distribution of $G(\cdot)$, we may be interested in producing estimates of the entire distribution itself. To do so, re-arrange the weighted likelihood:

$$L(\mathbf{p}) = \prod_{j=1}^J \prod_{i:D_{ij}=1} \left(p_j^{Y_i} (1-p_j)^{1-Y_i} \right)^{w_{ij}}.$$

Now, let

$$f_j(p) = f(p \mid \{Y_i, S_i\}_{i:D_{ij}=1}) = \prod_{i:D_{ij}=1} \left(p^{Y_i} (1-p)^{1-Y_i} \right)^{w_{ij}}$$

for $j = 1, \dots, J$ denote the weighted likelihood of the observed vector of outcomes for each unit j . Define the integrated likelihood as:

$$L^*(G) = \int \cdots \int \prod_{j=1}^J \prod_{i:D_{ij}=1} \left(p_j^{Y_i} (1-p_j)^{1-Y_i} \right)^{w_{ij}} dG(p_1, \dots, p_J)$$

By Tonelli's theorem, the integrated likelihood and its logarithm can be written:

$$L^*(G) = \prod_{j=1}^J \int \prod_{i:D_{ij}=1} \left(p_j^{Y_i} (1-p_j)^{1-Y_i} \right)^{w_{ij}} dG(p_j)$$

$$\ell^*(G) = \sum_{j=1}^J \log \left(\int f_j(p | Y_j, Z_j) dG(p) \right).$$

We may produce estimates of the distribution of treatment effects across units by maximizing $\ell^*(G)$ with respect to $G(\cdot)$. Because $G(\cdot)$ is infinite dimensional, and the log-integrated likelihood is nonconvex, this optimization problem poses a challenge. One way to simplify the problem is to assume that $G(\cdot)$ is a member of a parametric family of distributions. A natural choice here is a Beta(α, β) distribution. Estimation of α and β can be achieved via the EM algorithm.

A second option, proposed by Koenker and Mizera (2014), is a nonparametric alternative to the parametric method described above. We assume that $G(\cdot)$ takes on the form:

$$G(p) = \sum_{k=0}^K g_k \times \mathbf{1}[p \leq k/K], \text{ with } G(1) = 1, g_k \geq 0 \forall k,$$

for some K relatively large. This assumption restricts the support of $G(\cdot)$ to a fine grid of points. Given this assumption, the likelihood becomes:

$$\begin{aligned} \ell^*(\theta) &= \sum_{j=1}^J \log \left(\sum_{k=0}^K f_j(p_k | Y_j, Z_j) g_k \right) \\ &= \sum_{j=1}^J \log \left(\sum_{k=0}^K f_{jk} g_k \right) \\ &= \sum_{j=1}^J \log (F_j' G), \end{aligned}$$

where $f_{jk} = f_j(p_k | Y_j, Z_j)$, $F_j = (f_{jk})_{k=1}^K$, and $G = (g_k)_{k=1}^K$. Let $F = (F_1, \dots, F_J)$. Estimation of G proceeds via nonlinear convex programming:

$$\min_G -\mathbf{1}'_J \log(F'G) \text{ s.t. } \mathbf{1}'_K G = 1, G \geq 1.$$

The distributions produced by this routine are “spiky,” with estimated $g_k > \epsilon \approx 0$ for approximately $\log(J)$ points only. Efron (2016) proposed an empirical bayes deconvolution estimator that is essentially a smoothed version of Koneker and Mizera's estimator that imposes additional (exponential family) structure on the g_k . Specifically, Efron sets:

$$g_k = g_k(\alpha) = \exp(Q'_k \alpha - \phi(\alpha)), \text{ with } \phi(\alpha) = \log \left(\sum_{k=0}^K \exp(Q'_k \alpha) \right).$$

Here, α is a p -dimensional parameter vector and $Q = (Q_0 \dots Q_K)$ is a known $p \times K + 1$ design matrix. The full procedure specifies $Q'\alpha$ as a spline in p_k and imposes a penalty function on α .

2.4 Modeling the Assignment Process

In most practical applications, $w_j(v) = \mathbf{1}[r(i) \in \mathcal{O}] / \Pr(D_{ij} = 1 \mid V_i = v, Z_i = 1)$ is unknown and must be estimated. Estimation of the $J \times R$ matrix of propensity scores is complicated by the fact that in many real-world applications, a large share of unit-by-strata cells will contain exactly zero observations. This problem is exacerbated when the assumption of unconfoundedness only holds conditional on a very large set of factors, such that the dimension R grows large relative to the overall sample size. Traditionally, the problem of limited overlap is overcome by trimming on the estimated propensity score: excluding observations with estimated propensity scores above or below cutoff values (e.g Crump et al. 2009). Because the data is sparse and treatment is high-dimensional rather than binary, trimming on the estimated propensity score poses a fundamental problem: we need to infer whether a unit/stratum cell contains zero observations because there was actually no chance that an observation could have been assigned to that unit in that stratum, or rather that the cell contained zero observations by chance.

For each unit $j = \{1, \dots, J\}$, and randomization strata $r = \{1, \dots, R\}$ cell, we observe a count n_{jr} which measures the number of times unit j was assigned to an individual in strata r . Each randomization strata r is associated with a marginal count $n_r = \sum_{j=1}^J n_{jr}$. We form the J -row-by- R column matrix

$$\mathbf{N} = \left(n_{jr} \right)_{j=1, r=1}^{J, R},$$

which encodes the cell counts. We model the cell counts as draws from a mixture distribution:

Assumption 2.4. (Zero-Inflated Poisson (ZIP) Model) n_{jr} , the number of times unit j appears in stratum r , follows a mixture distribution:

$$n_{jr} \sim Z_{jr} \times M_{jr},$$

where the terms of the mixture are independent and distributed:

$$\begin{aligned} Z_{jr} &\sim \text{Bernoulli}(\pi_{jr}), \text{ and} \\ M_{jr} &\sim \text{Poisson}(\lambda_{jr}), \end{aligned}$$

where $\pi_{jr} \in (0, 1)$ and $\lambda_{jr} > 0 \forall jr$.

Let $\mathbf{\Lambda} = \left(\lambda_{jr} \right)_{j=1, r=1}^{J, R}$, $\mathbf{\Pi} = \left(\pi_{jr} \right)_{j=1, r=1}^{J, R}$, and $\mathbf{Z} = \left(Z_{jr} \right)_{j=1, r=1}^{J, R}$ denote the $J \times R$ matrices collecting the Poisson rate parameters, Bernoulli success probabilities, and latent overlap indicators, respectively. The ZIP model specified by assumption 2.4 is standard for estimating sparse count models. The model allows for over-dispersion (excess zeros) in the observed counts n_{jr} when the sample is derived from a (potentially) sparse data generating process. A key feature of this model is that the random variables Z_{jr} , which determine overlap, are actually partially observed. When $n_{jr} > 0$, we can infer that $Z_{jr} = 1$. However, when $n_{jr} = 0$, we can only imperfectly predict Z_{jr} .

Maximum Likelihood Estimation via the EM Algorithm

Denote the parameters of the full model by $\boldsymbol{\theta} = \{\boldsymbol{\Pi}, \boldsymbol{\Lambda}\}$. The complete-data log-likelihood for the above model is:

$$\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{Z}, \mathbf{N}) = \sum_{j=1}^J \sum_{r=1}^R \log \left((1 - Z_{jr})(1 - \pi_{jr})\mathbf{1}[n_{jr} = 0] + Z_{jr}\pi_{jr}\lambda_{jr}^{n_{jr}} \exp(-\lambda_{jr})/(n_{jr}!) \right).$$

We form the log integrated likelihood by replacing the expression inside the logarithm with its expectation, conditional on the observed counts:

$$\mathcal{L}^*(\boldsymbol{\theta} \mid \mathbf{N}) = \sum_{j=1}^J \sum_{r=1}^R \log \left((1 - \pi_{jr})\mathbf{1}[n_{jr} = 0] + \pi_{jr}\lambda_{jr}^{n_{jr}} \exp(-\lambda_{jr})/(n_{jr}!) \right).$$

The log integrated likelihood is highly nonconvex, rendering estimation by direct maximization infeasible. We therefore maximize this function via the EM algorithm. The E-step entails forming the function:

$$\mathcal{E}_{\mathbf{N}}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) = \sum_{j=1}^J \sum_{r=1}^R \left[(1 - \tilde{\pi}_{jr}^t) \log \left((1 - \pi_{jr})\mathbf{1}[n_{jr} = 0] \right) + \tilde{\pi}_{jr}^t \log \left(\pi_{jr}\lambda_{jr}^{n_{jr}} \exp(-\lambda_{jr})/(n_{jr}!) \right) \right],$$

where $\boldsymbol{\theta}^t$ denotes the parameter values at the t -th iteration of the algorithm, and $\tilde{\pi}_{jr}^t$ denotes the posterior probability of the event $Z_{jr} = 1$ conditional on the observed values of n_{jr} and those parameters. Suppressing t superscripts, these probabilities are given by:

$$\begin{aligned} \tilde{\pi}_{jr} &= \Pr(Z_{jr} = 1 \mid n_{jr}; \boldsymbol{\theta}) \\ &= \begin{cases} \frac{\pi_{jr} \exp(-\lambda_{jr})}{\pi_{jr} \exp(-\lambda_{jr}) + (1 - \pi_{jr})} & \text{if } n_{jr} = 0, \text{ and} \\ 1 & \text{otherwise.} \end{cases} \end{aligned}$$

Let $\tilde{\boldsymbol{\Pi}}$ denote a matrix collecting all $J \times R$ posterior probabilities $\tilde{\pi}_{jr}$. Note that the event $n_{jr} > 0$ implies $\tilde{\pi}_{jr} = 1$. We adopt the information-theoretic convention that $0 \log(0) = \lim_{p \downarrow 0} p \log(p) = 0$, which implies that $(1 - \tilde{\pi}_{jr}^t) \log \left((1 - \pi_{jr})\mathbf{1}[n_{jr} = 0] \right) = 0$ when $n_{jr} > 0$, and more generally that the expression $(1 - \tilde{\pi}_{jr}^t) \log \left((1 - \pi_{jr})\mathbf{1}[n_{jr} = 0] \right)$ is always equal to $(1 - \tilde{\pi}_{jr}^t) \log(1 - \pi_{jr})$. Importantly, the $\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t, \mathbf{N})$ function is separable in its arguments:

$$\begin{aligned} \mathcal{E}_{\mathbf{N}}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) &= \mathcal{E}_{\mathbf{N}}^0(\boldsymbol{\Pi} \mid \boldsymbol{\theta}^t) + \mathcal{E}_{\mathbf{N}}^1(\boldsymbol{\Lambda} \mid \boldsymbol{\theta}^t) + C^t, \text{ with} \\ \mathcal{E}_{\mathbf{N}}^0(\boldsymbol{\Pi} \mid \boldsymbol{\theta}^t) &= \sum_{j=1}^J \sum_{r=1}^R \left[(1 - \tilde{\pi}_{jr}^t) \log(1 - \pi_{jr}) + \tilde{\pi}_{jr}^t \log(\pi_{jr}) \right], \text{ and} \\ \mathcal{E}_{\mathbf{N}}^1(\boldsymbol{\Lambda} \mid \boldsymbol{\theta}^t) &= \sum_{j=1}^J \sum_{r=1}^R \left[\tilde{\pi}_{jr}^t (n_{jr} \log(\lambda_{jr}) - \lambda_{jr}) \right], \end{aligned}$$

and where C^t is a constant. The M-step entails maximizing $\mathcal{E}_{\mathbf{N}}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t)$ with respect to the elements of $\boldsymbol{\theta}$, which can be achieved by maximizing each part separately. In our implementations, both maximization sub-tasks involve approximating nonnegative $J \times R$ matrices with nonnegative matrices of (much) lower rank. We consider two possible alternatives for estimating the parameters of this model.

Estimation under Exact Rank Constraints

The model we have specified is heavily overparameterized: every observed count n_{jr} is assumed to be a draw from a two-parameter distribution, where each jr cell may have unique values of those parameters. In order to identify the model, we impose constraints on the structure of the parameter matrices $\mathbf{\Lambda}$ and $\mathbf{\Pi}$. These constraints more explicitly formalize 2.2b. We first consider estimating the ZIP model of counts under exact constraints on the rank of both $\mathbf{\Lambda}$ and $\mathbf{\Pi}$:

Assumption 2.5. (Exact Rank Constraints)

$$\text{rank}(\mathbf{\Lambda}) = Q < \min(J, R) \text{ and } \text{rank}(\mathbf{\Pi}) = 1.$$

Under assumption 2.5, we may write $\mathbf{\Lambda} = \mathbf{U}\mathbf{V}'$ with \mathbf{U} a $J \times Q$ matrix and \mathbf{V} an $R \times Q$ matrix, and $\mathbf{\Pi} = \boldsymbol{\mu}\boldsymbol{\nu}'$ with $\boldsymbol{\mu}$ a $J \times 1$ vector and $\boldsymbol{\nu}$ an $R \times 1$ vector. Note that the choice to model $\mathbf{\Pi}$ as a rank-1 matrix is arbitrary, although there are some intuitive features of this representation. Algorithms to optimize objectives like $\mathcal{E}_{\mathbf{N}}^0(\mathbf{\Pi} \mid \boldsymbol{\theta}^t)$ or $\mathcal{E}_{\mathbf{N}}^1(\mathbf{\Lambda} \mid \boldsymbol{\theta}^t)$ under exact rank constraints have been called “Non-negative Matrix Factorization” (NMF) routines (Lee and Seung 2000). NMF methods have existed under various names since at least the 1990s, but gained popularity after they were successfully implemented in the Netflix Prize competition, which asked entrants to create an algorithm to predict which movies Netflix users would rate highly given data on their ratings of prior movies. In the machine learning literature, these models are often referred to as “recommender systems” or “collaborative filtering” algorithms. NMF methods have been shown to be equivalent to several well-known statistical modeling procedures when appropriate constraints are applied, including spectral clustering, K-means clustering, and Probabilistic Latent Semantic Indexing (a popular method for text analysis).

We parameterize the Poisson rate parameters as follows:

$$\lambda_{jr} = \sum_{q=1}^Q u_{jq}v_{rq}, \text{ such that } \sum_{j=1}^J u_{jq} = 1 \forall q, \quad u_{jq} \geq \tau > 0 \forall jq, \text{ and } v_{rq} \geq 0 \forall rq.$$

We collect the u_{jq} into Q separate $J \times 1$ column vectors $\mathbf{u}_q = (u_{1q}, \dots, u_{Jq})'$ (*factors*), and the v_{rq} into R separate $Q \times 1$ column vectors $\mathbf{v}_r = (v_{r1}, \dots, v_{rQ})'$ (*factor loadings*). We may then write:

$$\begin{aligned} \mathbf{U} &= [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_Q] \in [\tau, 1 - \tau(J - 1)]^{J \times Q}, \text{ and} \\ \mathbf{V} &= [\mathbf{v}^1 \ \mathbf{v}^2 \ \dots \ \mathbf{v}^R]' \in \mathbb{R}_+^{R \times Q}. \end{aligned}$$

The requirement that the elements of each of the q factors sum to one is not a substantive restriction on its own (dividing a column of \mathbf{U} by a scalar is equivalent to multiplying the corresponding column of \mathbf{V} by that scalar). However, imposing this restriction ensures identifiability and interpretability of the parameters. The restriction that the factors \mathbf{u}_q and factor loadings \mathbf{v}^r are nonnegative *is* a substantive restriction. These restrictions ensure

that the Poisson rate parameters are well-defined (these parameters cannot be negative). Finally, the restriction that the factor values all be greater than τ serves two purposes. First, this restriction helps ensure the model is identified. If the u_{jq} were allowed to equal zero identically, then some of the λ_{jr} could be set to zero identically. In such a case, there is no meaningful distinction between the events $Z_{jr} = 1$ and $Z_{jr} = 0$. Second, this restriction in combination with the summing-to-one restrictions motivates an interpretation of the model as determining the subsample of the data for which overlap in the distributions of assignments is most likely to be weak. The tuning parameter τ corresponds to the researcher's prior about level of overlap in the sample.

We now present an algorithm for estimating the model parameters. First, consider the parameters $\boldsymbol{\mu}$ and $\boldsymbol{\nu}$. Denote the Karush–Kuhn–Tucker (KKT) multipliers that encode the constraints $1 - \mu_j > 0$ and $1 - \nu_r > 0$ by γ_j and γ_r , respectively (the positivity constraints are trivially non-binding). The FOCs for these parameters can be written:

$$\begin{aligned}\mu_j &: \sum_{r=1}^R \nu_r \left(\frac{\tilde{\pi}_{jr}}{\mu_j \nu_r} - \frac{1 - \tilde{\pi}_{jr}}{1 - \mu_j \nu_r} \right) - \gamma_j = 0, \text{ and} \\ \nu_r &: \sum_{j=1}^J \mu_j \left(\frac{\tilde{\pi}_{jr}}{\mu_j \nu_r} - \frac{1 - \tilde{\pi}_{jr}}{1 - \mu_j \nu_r} \right) - \gamma_r = 0.\end{aligned}$$

Denote $\frac{\tilde{\pi}_{jr}}{\mu_j \nu_r}$ by ψ_{jr}^1 and $\frac{1 - \tilde{\pi}_{jr}}{1 - \mu_j \nu_r}$ by ψ_{jr}^0 , which are collected in the corresponding $J \times R$ matrices $\boldsymbol{\Psi}^1$ and $\boldsymbol{\Psi}^0$. Eliminating the KKT multipliers and re-arranging yields the following multiplicative update rules, which guarantee at each step that the parameters satisfy the constraints and that the objective function weakly increases:

$$\begin{aligned}\boldsymbol{\mu} &\leftarrow \mathbf{1}_J - (\mathbf{1}_J - \boldsymbol{\mu}) \circ \frac{\boldsymbol{\Psi}^0 \boldsymbol{\nu}}{\boldsymbol{\Psi}^1 \boldsymbol{\nu}}, \text{ and} \\ \boldsymbol{\nu} &\leftarrow \mathbf{1}_R - (\mathbf{1}_R - \boldsymbol{\nu}) \circ \frac{\boldsymbol{\Psi}^{0'} \boldsymbol{\mu}}{\boldsymbol{\Psi}^{1'} \boldsymbol{\mu}},\end{aligned}$$

where $\mathbf{1}_J$ and $\mathbf{1}_R$ are conformable column vectors of ones, \circ denotes the element-wise (Hadamard) product of matrices, and division is also element-wise. Iterating over these updates (which involves re-calculating the $\boldsymbol{\Psi}$ matrices after each step) is a gradient descent procedure with a fixed step size. To avoid numerical issues, a small constant may be added to both the numerator and denominator of each expression.

Next, consider the parameters \mathbf{U} and \mathbf{V} . Denote the KKT multipliers that encode the constraints $u_{jq} - \tau \geq 0$ and $v_{rq} \geq 0$ by γ_{jq} and γ_{rq} , respectively. In addition, denote the Lagrange multipliers that encode the constraints $\sum_{q=1}^Q u_{jq} = 1$ by γ_j . The FOCs for these parameters are then:

$$\begin{aligned}u_{jq} &: \sum_{r=1}^R \tilde{\pi}_{jr} v_{rq} \left(\frac{n_{jr}}{\lambda_{jr}} - 1 \right) + \gamma_{jq} - \gamma_j = 0, \text{ and} \\ v_{rq} &: \sum_{j=1}^J \tilde{\pi}_{jr} u_{jq} \left(\frac{n_{jr}}{\lambda_{jr}} - 1 \right) + \gamma_{rq} = 0.\end{aligned}$$

Denote $\frac{\tilde{\pi}_{jr}n_{jr}}{\lambda_{jr}}$ by ω_{jr} , which are collected in the corresponding $J \times R$ matrix $\mathbf{\Omega}$. Eliminating the KKT multipliers and summing the FOCs for u_{jq} over $j = 1, \dots, J$ gives:

$$\gamma_q = \frac{(\mathbf{u}_q - \tau)' (\mathbf{\Omega} - \tilde{\mathbf{\Pi}}) \mathbf{v}_q}{1 - J\tau},$$

where τ is a conformable matrix with every entry equal to the scalar τ , \mathbf{u}_q is the q -th factor (column) of the factor matrix \mathbf{U} and \mathbf{v}_q is the q -th column of the factor loading matrix \mathbf{V} . Let $\gamma_{\tilde{\pi}} = \text{diag}((\mathbf{U} - \tau)' \tilde{\mathbf{\Pi}} \mathbf{V}) / (1 - J\tau)$ and $\gamma_{\omega} = \text{diag}((\mathbf{U} - \tau)' \mathbf{\Omega} \mathbf{V}) / (1 - J\tau)$, where $\text{diag}(\cdot)$ is a function that returns the central diagonal of its argument as a column vector. Eliminating the KKT multipliers and simplifying gives the following update rules:

$$\begin{aligned} \mathbf{U} &\leftarrow \tau + (\mathbf{U} - \tau) \circ \frac{\mathbf{\Omega} \mathbf{V} + \gamma'_{\tilde{\pi}}}{\tilde{\mathbf{\Pi}} \mathbf{V} + \gamma'_{\omega}}, \text{ and} \\ \mathbf{V} &\leftarrow \mathbf{V} \circ \frac{\mathbf{\Omega}' \mathbf{U}}{\tilde{\mathbf{\Pi}}' \mathbf{V}}. \end{aligned}$$

Iterating over the full procedure (the E and M steps, with iterative maximization of each term of \mathcal{E} in the M step) yields estimates of the model parameters $\theta = \{\boldsymbol{\mu}, \boldsymbol{\nu}, \mathbf{U}, \mathbf{V}\}$.

2.5 A Nuclear Norm Regularized Estimator of Propensity Scores

While the the rank-constrained EM algorithm presented in the prior section is relatively simple to implement, it is highly nonconvex. So, while the algorithm might converge to a local stationary point, there is no guarantee that that stationary point represents a global, or even local, maximum of the (constrained) integrated likelihood. Recognizing this drawback, it is common to replace the exact constraint on the rank of \mathbf{A} with the *convex relaxation* of that constraint. The rank of a matrix \mathbf{X} can be written as:

$$\text{rank}(\mathbf{X}) = \sum_{i=1}^{\min(J,R)} \mathbf{1}[\sigma_i(\mathbf{X}) > 0],$$

where $\sigma_i(\mathbf{X})$ is the i -th singular value of \mathbf{X} . The Nuclear Norm of \mathbf{X} is the convex relaxation of the rank constraint:

$$\|\mathbf{X}\|_* = \sum_{i=1}^{\min(J,R)} \sigma_i(\mathbf{X}).$$

In this section, we propose an alternative estimator for the model parameters under a relaxed version of the rank constraint. In particular, we assume:

Assumption 2.6. (Nuclear Norm Constraints)

$$\|\mathbf{A}\|_* \leq M_{\lambda} \text{ and } \|\mathbf{\Pi}\|_* \leq M_{\pi}.$$

In practice, we implement these constraints by regularizing the likelihood, since selecting an appropriate choice of regularization parameters is equivalent to enforcing the constraints. We form the regularized negative log-integrated likelihood as:

$$\mathcal{L}_\tau^*(\boldsymbol{\theta} \mid \mathbf{N}) = -\mathcal{L}^*(\boldsymbol{\theta} \mid \mathbf{N}) + \underbrace{\tau_1 \|\boldsymbol{\Pi}\|_* + \tau_2 \|\boldsymbol{\Lambda}\|_*}_{=\mathcal{R}_\tau(\boldsymbol{\theta})},$$

where τ_1 and τ_2 are regularization parameters. Estimation proceeds by minimizing this function. For the remainder of this section, we will implicitly condition on the sample, such that $\mathcal{L}_\tau^*(\boldsymbol{\theta})$, $\mathcal{L}^*(\boldsymbol{\theta})$, and $\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t)$ denote $\mathcal{L}_\tau^*(\boldsymbol{\theta} \mid \mathbf{N})$, $\mathcal{L}^*(\boldsymbol{\theta} \mid \mathbf{N})$ and $\mathcal{E}_{\mathbf{N}}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t)$, respectively.

Estimation Algorithm and Computational Guarantees

We now outline an algorithm for estimating the model parameters. To maximize the regularized log integrated likelihood, we adapt the analysis of Lu, Freund, and Nesterov (2018) to our setting. As a preliminary matter, we introduce the notion of *Bregman Divergences*, which are a generalization of norms:

Definition 2.2. (Bregman Divergence) *Let $h(\cdot)$ be a strictly convex and continuously differentiable function defined on a closed convex set Ω . Then the h -Bregman divergence between any two points $x, y \in \Omega$ is:*

$$D_h(x, y) = h(x) - h(y) - \langle \nabla h(y), x - y \rangle.$$

When considering matrices or vectors, we write $D_h(\mathbf{X}, \mathbf{Y})$ to denote the sum of the Bregman divergences between all individual elements of \mathbf{X} and \mathbf{Y} : if \mathbf{X} and \mathbf{Y} are $J \times R$ matrices, then $D_h(\mathbf{X}, \mathbf{Y}) = \sum_{j=1}^J \sum_{r=1}^R D_h(x_{jr}, y_{jr})$ (this is sometimes called the collective Bregman divergence). Like norms, Bregman divergences $D_h(x, y)$ are strictly positive when $x \neq y$, are equal to zero when $x = y$. Unlike norms, Bregman divergences are not symmetric in general: $D_h(x, y)$ need not be equal to $D_h(y, x)$.

Typically, algorithms to minimize complicated functions rely on an assumption that objectives are *smooth* and *strongly convex*. Together, these assumptions require that, for any given point x in the domain of the objective $f(\cdot)$, the difference between the value of $f(\cdot)$ at any other point y ($f(y)$) and the Taylor series approximation to $f(y)$ around x ($f(x) + \langle \nabla f(x), y - x \rangle$) can be upper- and lower-bounded (respectively) by scalar multiples of the squared Euclidean distance between x and y . The objective $\mathcal{L}_\tau^*(\cdot)$ does not obey either of these conditions, and so traditional results about the convergence of optimization algorithms do not necessarily apply.

Lu, Freund, and Nesterov (2018) consider problems in which the objective does not obey standard notions of smoothness and strong convexity, and provide computation guarantees for an alternative optimization algorithm that accommodates such settings. To do so, they define alternative notions of smoothness and strong convexity relative to an appropriately-chosen Bregman divergence:

Definition 2.3. (Relative Smoothness and Strong Convexity)

a) A function $f(\cdot)$ is L -smooth relative to $h(\cdot)$ on a set Ω if for any $x, y \in \Omega$, there is a scalar L for which

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + LD_h(y, x).$$

b) A function $f(\cdot)$ is μ -strongly convex relative to $h(\cdot)$ on a set Ω if for any $x, y \in \Omega$, there is a scalar $\mu \geq 0$ for which

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \mu D_h(y, x).$$

Given this notation, we now state Theorem 3.1 of Lu, Freund, and Nesterov (2018), which provides a computation guarantee for a primal gradient scheme to minimize a function obeying relative smoothness and strong convexity conditions:

Theorem 2.2. (Computation Guarantee for Primal Gradient Scheme) *Let $f(\cdot)$, $h(\cdot)$, L , μ , and Ω satisfying Definition 2.3 be given. Let $x^1 \in \Omega$ denote an initialization point of our algorithm. At each iteration, perform an update that sets:*

$$x^{t+1} = \arg \min_{x \in \Omega} \{f(x^t) + \langle \nabla f(x^t), x - x^t \rangle + LD_h(x, x^t)\}.$$

Then, for all $t \geq 1$ and $x \in \Omega$, the sequence $\{f(x^t)\}$ is monotonically decreasing, and the following inequality holds:

$$f(x^{t+1}) - f(x) \leq \frac{\mu D_h(x, x^1)}{(1 + \frac{\mu}{L-\mu})^t - 1} \leq \frac{L - \mu}{t} D_h(x, x^1),$$

where, in the case when $\mu = 0$, the middle expression is defined in the limit as $\mu \rightarrow 0+$, which is equal to L/t .

We now prove the following proposition about the regularized log integrated likelihood:

Proposition 2.1. (Smoothness of \mathcal{L}_τ^*) *Define $m = \max_{j,r} n_{jr}$. The regularized log integrated likelihood function is 1-smooth relative to a function $h(\cdot)$:*

$$\mathcal{L}_\tau^*(\boldsymbol{\theta}) \leq \mathcal{L}_\tau^*(\boldsymbol{\theta}^t) + \langle \nabla \mathcal{L}_\tau^*(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_h(\boldsymbol{\theta}, \boldsymbol{\theta}^t),$$

where the function $h(\cdot)$ is given by:

$$h(\boldsymbol{\theta}) = - \sum_{j=1}^J \sum_{r=1}^R \left[\log(\pi_{jr}) + \log(1 - \pi_{jr}) + m \log(\lambda_{jr}) \right] + \mathcal{R}_\tau(\boldsymbol{\theta}).$$

Proof. First, note that the log integrated likelihood can be written as the sum of the EM proxy function and an entropy term (Dempster, Laird, and Rubin 1977):

$$\mathcal{L}^*(\boldsymbol{\theta}) = \mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) + \mathcal{H}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t),$$

where

$$\mathcal{H}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) = - \sum_{j=1}^J \sum_{r=1}^R \sum_{z=0}^1 \log \left(\Pr(Z_{jr} = z \mid n_{jr}, \boldsymbol{\theta}) \right) \Pr(Z_{jr} = z \mid n_{jr}, \boldsymbol{\theta}^t),$$

and $\mathcal{H}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t) \geq \mathcal{H}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t)$ for all $\boldsymbol{\theta}$. This implies that we may write:

$$\mathcal{L}_\tau^*(\boldsymbol{\theta}) - \mathcal{L}_\tau^*(\boldsymbol{\theta}^t) = - \left[\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) - \mathcal{E}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t) \right] - \left[\mathcal{H}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) - \mathcal{H}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t) \right] + \left[\mathcal{R}_\tau(\boldsymbol{\theta}) - \mathcal{R}_\tau(\boldsymbol{\theta}^t) \right].$$

The condition $\mathcal{H}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) \geq \mathcal{H}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t)$ implies :

$$\mathcal{L}_\tau^*(\boldsymbol{\theta}) - \mathcal{L}_\tau^*(\boldsymbol{\theta}^t) \leq - \left[\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) - \mathcal{E}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t) \right] + \left[\mathcal{R}_\tau(\boldsymbol{\theta}) - \mathcal{R}_\tau(\boldsymbol{\theta}^t) \right].$$

Next, note that $\nabla \mathcal{L}_\tau^*(\boldsymbol{\theta}^t) = -\nabla \mathcal{L}^*(\boldsymbol{\theta}^t) + \nabla \mathcal{R}_\tau(\boldsymbol{\theta}^t)$. Simple calculations show that:

$$\frac{\partial \mathcal{L}^*(\boldsymbol{\theta})}{\partial \pi_{jr}} = \frac{\tilde{\pi}_{jr}}{\pi_{jr}} - \frac{1 - \tilde{\pi}_{jr}}{1 - \pi_{jr}} \quad \text{and} \quad \frac{\partial \mathcal{L}^*(\boldsymbol{\theta})}{\partial \lambda_{jr}} = \tilde{\pi}_{jr} \left(\frac{n_{jr}}{\lambda_{jr}} - 1 \right),$$

where $\tilde{\pi}_{jr}$ is the posterior probability that $Z_{jr} = 1$ given n_{jr} and parameters $\boldsymbol{\theta}$ (defined above). Next, rearranging terms gives:

$$\begin{aligned} & - \left[\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) - \mathcal{E}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t) \right] = \\ & - \sum_{j=1}^J \sum_{r=1}^R \left[\tilde{\pi}_{jr}^t \log \left(\frac{\pi_{jr}}{\tilde{\pi}_{jr}^t} \right) + (1 - \tilde{\pi}_{jr}^t) \log \left(\frac{1 - \pi_{jr}}{1 - \tilde{\pi}_{jr}^t} \right) + \tilde{\pi}_{jr}^t \left(n_{jr} \log \left(\frac{\lambda_{jr}}{\lambda_{jr}^t} \right) - (\lambda_{jr} - \lambda_{jr}^t) \right) \right]. \end{aligned}$$

To simplify this expression, we use the identity: $\log(x/y) = (x - y)/y - D_{-\log}(x, y)$, where $D_{-\log}(x, y)$ is the Bregman divergence associated with the convex function $-\log(\cdot)$. Making this substitution and collecting terms, we may write:

$$\begin{aligned} & - \left[\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) - \mathcal{E}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t) \right] = \langle -\nabla \mathcal{L}^*(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle \\ & + \sum_{j=1}^J \sum_{r=1}^R \left[\tilde{\pi}_{jr}^t D_{-\log}(\pi_{jr}, \tilde{\pi}_{jr}^t) + (1 - \tilde{\pi}_{jr}^t) D_{-\log}(1 - \pi_{jr}, 1 - \tilde{\pi}_{jr}^t) + \tilde{\pi}_{jr}^t n_{jr} D_{-\log}(\lambda_{jr}, \lambda_{jr}^t) \right]. \end{aligned}$$

Since $D_{-\log}(x, y)$ is always positive, $0 \leq \tilde{\pi}_{jr}^t \leq 1$, and $\tilde{\pi}_{jr}^t n_{jr} \leq m$, we have:

$$\begin{aligned} & - \left[\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t) - \mathcal{E}(\boldsymbol{\theta}^t \mid \boldsymbol{\theta}^t) \right] \leq \langle -\nabla \mathcal{L}^*(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle \\ & + \underbrace{\sum_{j=1}^J \sum_{r=1}^R \left[D_{-\log}(\pi_{jr}, \tilde{\pi}_{jr}^t) + D_{-\log}(1 - \pi_{jr}, 1 - \tilde{\pi}_{jr}^t) + m D_{-\log}(\lambda_{jr}, \lambda_{jr}^t) \right]}_{=D_{\mathcal{B}_m}(\boldsymbol{\theta}, \boldsymbol{\theta}^t)}. \end{aligned}$$

The function $\mathcal{B}_m(\boldsymbol{\theta}) = -\sum_{j=1}^J \sum_{r=1}^R [\log(\pi_{jr}) + \log(1 - \pi_{jr}) + m \log(\lambda_{jr})]$ is the first term of the reference function: $h(\cdot) = \mathcal{B}_m(\cdot) + \mathcal{R}_\tau(\cdot)$. $\mathcal{B}_m(\cdot)$ and can be thought of as a “barrier” function that enforces the constraints that $\pi_{jr} \in (0, 1)$ and $\lambda_{jr} > 0$. Since Bregman divergences are linear, we have that:

$$D_h(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = D_{\mathcal{B}_m}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) + D_{\mathcal{R}_\tau}(\boldsymbol{\theta}, \boldsymbol{\theta}^t).$$

To complete the proof, we substitute the above inequality into the prior inequality, and add and subtract $\langle \nabla \mathcal{R}_\tau(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle$:

$$\begin{aligned} \mathcal{L}_\tau^*(\boldsymbol{\theta}) - \mathcal{L}_\tau^*(\boldsymbol{\theta}^t) &\leq \langle -\nabla \mathcal{L}^*(\boldsymbol{\theta}^t) + \nabla \mathcal{R}_\tau(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_{\mathcal{B}_m}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) + D_{\mathcal{R}_\tau}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \\ &= \langle \nabla \mathcal{L}_\tau^*(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_h(\boldsymbol{\theta}, \boldsymbol{\theta}^t). \end{aligned}$$

□

For the remainder of this section we maintain the following assumption, variants of which are standard the literature analyzing properties of EM algorithms:

Assumption 2.7. (0-Relative Strong Convexity) *Define $h(\cdot)$ as in Proposition 2.1, let $\boldsymbol{\theta}^*$ denote the minimum of $\mathcal{L}_\tau^*(\cdot)$, and let $\boldsymbol{\theta}^1$ denote an initialization point of our algorithm. The function $\mathcal{L}_\tau^*(\cdot)$ is 0-strongly convex relative to $h(\cdot)$ in a region containing $\boldsymbol{\theta}^*$, and all initialization points $\boldsymbol{\theta}^1$ are also contained in this region.*

We are now prepared to describe our estimation algorithm, which is a modification of the EM algorithm:

Modified EM Algorithm: Initialize the algorithm at point $\boldsymbol{\theta}^1$. At iteration $t+1$, compute $\nabla \mathcal{L}^*(\boldsymbol{\theta}^t)$, the gradient of the unregularized log-likelihood. Then update according to:

$$\begin{aligned} \boldsymbol{\theta}^{t+1} &= \arg \min_{\boldsymbol{\theta}} \left\{ \mathcal{L}_\tau^*(\boldsymbol{\theta}^t) + \langle \nabla \mathcal{L}_\tau^*(\boldsymbol{\theta}^t), \boldsymbol{\theta} - \boldsymbol{\theta}^t \rangle + D_h(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \right\} \\ &= \arg \min_{\boldsymbol{\theta}} \left\{ -\langle \nabla \mathcal{L}^*(\boldsymbol{\theta}^t), \boldsymbol{\theta} \rangle + \mathcal{R}_\tau(\boldsymbol{\theta}) + D_{\mathcal{B}_m}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \right\}, \end{aligned}$$

where the second line drops constants and rearranges terms. Optimization at each step can be achieved efficiently using modern convex programming solvers. Given Proposition 2.1 and Assumption 2.7, the following Corollary is an immediate consequence of Theorem 2.2:

Corollary 2.1. (Convergence of the Modified EM Algorithm) *Define $h(\cdot)$ as in Proposition 2.1, and define $\boldsymbol{\theta}^*$ and $\boldsymbol{\theta}^1$ as in Assumption 2.7. If updates are computed as:*

$$\boldsymbol{\theta}^{t+1} = \arg \min_{\boldsymbol{\theta}} \left\{ -\langle \nabla \mathcal{L}^*(\boldsymbol{\theta}^t), \boldsymbol{\theta} \rangle + \mathcal{R}_\tau(\boldsymbol{\theta}) + D_{\mathcal{B}_m}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \right\},$$

then the sequence $\{\mathcal{L}_\tau^(\boldsymbol{\theta}^t)\}$ is monotonically decreasing, and*

$$\mathcal{L}_\tau^*(\boldsymbol{\theta}^T) - \mathcal{L}_\tau^*(\boldsymbol{\theta}^*) \leq \frac{D_h(\boldsymbol{\theta}^*, \boldsymbol{\theta}^1)}{T+1}.$$

While Corollary 2.1 establishes convergence of the modified EM algorithm, the rate of convergence is sublinear. Exploring improved computation guarantees is an important topic for further work.

Statistical Guarantees

We now consider the calculation of finite-sample bounds on the discrepancy between our estimates and the true parameters, $\boldsymbol{\theta}^0 = \{\boldsymbol{\Lambda}^0, \boldsymbol{\Pi}^0\}$. For the purposes of this paper, we will only consider the error associated with the Poisson rate parameters $\boldsymbol{\Lambda}$. We reserve analysis of finite-sample bounds on the discrepancy between $\boldsymbol{\Pi}^*$ and $\boldsymbol{\Pi}^0$ for future work, although similar arguments will apply in both cases.

For this analysis, we assume convergence of the modified EM algorithm to a stationary point $\boldsymbol{\theta}^*$. Remember that the EM proxy function $\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^t)$ can be written as the sum of two components, such that each component depends only upon $\boldsymbol{\Lambda}$ or $\boldsymbol{\Pi}$ separately. Additionally, note that at the stationary point, $\mathcal{L}_\tau^*(\boldsymbol{\theta}^*) = -\mathcal{E}(\boldsymbol{\theta}^* \mid \boldsymbol{\theta}^*) + \mathcal{R}_\tau(\boldsymbol{\theta}^*)$, and $\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} -\mathcal{E}(\boldsymbol{\theta} \mid \boldsymbol{\theta}^*) + \mathcal{R}_\tau(\boldsymbol{\theta})$. These facts allow us to analyze the finite sample error of $\boldsymbol{\Lambda}$ separately from that of $\boldsymbol{\Pi}$. Before proceeding, we introduce additional notation: we denote the Frobenious norm of a matrix \mathbf{X} by $\|\mathbf{X}\|_2 = \sqrt{\text{trace}(\mathbf{X}'\mathbf{X})}$, and the Spectral norm of \mathbf{X} by $\|\mathbf{X}\|_\infty = \max_{i \leq \min(J,R)} \sigma_i$, where the σ_i are the singular values of \mathbf{X} .

In order to derive a bound on the finite-sample error, we need to verify two additional conditions. The first condition is that a form of *restricted strong convexity* (RSC) holds with high probability. For now, we adopt an RSC assumption without explicitly deriving the probability that RSC holds, which we defer for future work:

Assumption 2.8. (Restricted Strong Convexity) *Define the quantity:*

$$D_{\mathbf{N}}(\boldsymbol{\Lambda}^*, \boldsymbol{\Lambda}^0) = \sum_{j=1}^J \sum_{r=1}^R n_{jr} D_{-\log}(\lambda_{jr}^*, \lambda_{jr}^0).$$

For a scalar $\alpha > 0$, the following inequality holds with probability at least $1 - \varepsilon_0$:

$$D_{\mathbf{N}}(\boldsymbol{\Lambda}^*, \boldsymbol{\Lambda}^0) \geq \alpha \|\boldsymbol{\Lambda}^* - \boldsymbol{\Lambda}^0\|_2^2.$$

Denote the set of parameter values where this inequality is satisfied (holding $\boldsymbol{\Lambda}^0$ fixed) by $\mathcal{C}_\alpha(\boldsymbol{\Lambda}^0)$, and let $\sup_{\boldsymbol{\Lambda} \in \mathcal{C}_\alpha(\boldsymbol{\Lambda}^0)} \frac{\|\boldsymbol{\Lambda} - \boldsymbol{\Lambda}^0\|_*}{\|\boldsymbol{\Lambda} - \boldsymbol{\Lambda}^0\|_2} = \beta \leq \sqrt{\min(J, R)}$.

The remainder of this analysis conditions on the event $\boldsymbol{\Lambda}^* \in \mathcal{C}_\alpha(\boldsymbol{\Lambda}^0)$. Assumption 2.8 allows us to place an upper bound on the difference between the EM proxy function $\mathcal{E}^1(\boldsymbol{\Lambda} \mid \boldsymbol{\Lambda}^t)$ evaluated at any point $\boldsymbol{\Lambda} \in \mathcal{C}_\alpha(\boldsymbol{\Lambda}^0)$ and evaluated at the true parameter values (where $\mathcal{E}^1(\cdot \mid \cdot)$ is the $\boldsymbol{\Lambda}$ -specific component of the full EM proxy function). By rearranging terms and adding and subtracting $\langle \nabla \mathcal{E}^1(\boldsymbol{\Lambda}^0 \mid \boldsymbol{\Lambda}^*), \boldsymbol{\Lambda}^* - \boldsymbol{\Lambda}^0 \rangle$, we have:

$$\begin{aligned} \mathcal{E}^1(\boldsymbol{\Lambda}^* \mid \boldsymbol{\Lambda}^*) - \mathcal{E}^1(\boldsymbol{\Lambda}^0 \mid \boldsymbol{\Lambda}^*) &= \langle \nabla \mathcal{E}^1(\boldsymbol{\Lambda}^0 \mid \boldsymbol{\Lambda}^*), \boldsymbol{\Lambda}^* - \boldsymbol{\Lambda}^0 \rangle - D_{\mathbf{N}}(\boldsymbol{\Lambda}^*, \boldsymbol{\Lambda}^0) \\ &\leq \langle \nabla \mathcal{E}^1(\boldsymbol{\Lambda}^0 \mid \boldsymbol{\Lambda}^*), \boldsymbol{\Lambda}^* - \boldsymbol{\Lambda}^0 \rangle - \alpha \|\boldsymbol{\Lambda}^* - \boldsymbol{\Lambda}^0\|_2^2. \end{aligned}$$

The second condition is that a bound on the statistical error of the model, as measured by the spectral norm of the deviation of observed counts from their means. There is a large literature on bounding the spectral norm of random matrices, although most papers focus on matrices with Gaussian entries.

McRae and Davenport (2020) derive the following bound (Lemma 2.4), which we particularize to our setting:

Lemma 2.1. (Tail Bound for Statistical Error) *Let \mathbf{N} be a random matrix where each entry is independently distributed $\text{Poisson}(\lambda_{jr}^0)$, such that $\mathbb{E}[\mathbf{N}] = \mathbf{\Lambda}$. Define the quantity:*

$$A_{\mathbf{\Lambda}^0}(\varepsilon_1) = 2\sigma(\mathbf{\Lambda}^0) + \frac{8\varepsilon_1}{\sqrt{JR}} + C_1 \max \left\{ \max_{jr} \lambda_{jr}^0, 4 \log \left(\frac{2JR}{\varepsilon_1} \right) \right\} \times \sqrt{\log \left(\frac{\max J, R}{\varepsilon_1} \right)},$$

where C_1 is a universal constant and

$$\sigma(\mathbf{\Lambda}^0) = \max_j \sqrt{\sum_{r=1}^R \lambda_{jr}^0} + \max_r \sqrt{\sum_{j=1}^J \lambda_{jr}^0}.$$

Then we have, for $\varepsilon_1 \in (0, 1/2)$:

$$\Pr \left(\|\mathbf{N} - \mathbf{\Lambda}^0\|_\infty \geq A_{\mathbf{\Lambda}^0}(\varepsilon_1) \right) \leq 2\varepsilon_1.$$

We use Lemma 2.1 to construct a bound on the quantity $\langle \nabla \mathcal{E}^1(\mathbf{\Lambda}^0 \mid \mathbf{\Lambda}^*), \mathbf{\Lambda}^* - \mathbf{\Lambda}^0 \rangle = \sum_{j=1}^J \sum_{r=1}^R \frac{\tilde{\pi}_{jr}^*}{\lambda_{jr}^0} (n_{jr} - \lambda_{jr}^0)(\lambda_{jr}^* - \lambda_{jr}^0)$. Let $\lambda_{\min}^0 = \min_{jr} \lambda_{jr}^0$. Applying the Cauchy-Schwartz inequality and 2.1, we can bound this term by

$$\langle \nabla \mathcal{E}^1(\mathbf{\Lambda}^0 \mid \mathbf{\Lambda}^*), \mathbf{\Lambda}^* - \mathbf{\Lambda}^0 \rangle \leq \frac{A_{\mathbf{\Lambda}^0}(\varepsilon_1)}{\lambda_{\min}^0} \|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_*$$

with high probability. We additionally condition on this event for the remaining analysis.

With the RSC and statistical error bounds established, we may now calculate a bound on the discrepancy between our estimates and the true model parameters.

Theorem 2.3. (Upper Error Bound for $\mathbf{\Lambda}$) *Adopt the assumptions and notation of this section, and assume the regularization parameter τ associated with $\|\mathbf{\Lambda}\|_*$ is set such that $\tau \geq \frac{A_{\mathbf{\Lambda}^0}(\varepsilon_1)}{\lambda_{\min}^0}$. Then with high probability, we have that:*

$$\|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_2 \leq 2 \frac{\beta\tau}{\alpha}.$$

Proof: We have already established most of the components of the proof. Because $\mathbf{\Lambda}^* = \arg \max_{\mathbf{\Lambda}} \mathcal{E}^1(\mathbf{\Lambda} \mid \boldsymbol{\theta}^*) - \tau \|\mathbf{\Lambda}\|_*$, we may write:

$$\tau (\|\mathbf{\Lambda}^*\|_* - \|\mathbf{\Lambda}^0\|_*) \leq \mathcal{E}^1(\mathbf{\Lambda}^* \mid \boldsymbol{\theta}^*) - \mathcal{E}^1(\mathbf{\Lambda}^0 \mid \boldsymbol{\theta}^*).$$

Substituting expressions established above on the righthandside of the inequality, and applying the triangle inequality to the lefthand side, we may write:

$$-\tau \|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_* \leq \frac{A_{\mathbf{\Lambda}^0}(\varepsilon_1)}{\lambda_{\min}^0} \|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_* - \alpha \|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_2^2.$$

Rearranging and applying the bound $\|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_* < \beta\|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_2$ gives:

$$\alpha\|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_2^2 \leq \beta \frac{A_{\mathbf{\Lambda}^0}(\varepsilon_1)}{\lambda_{\min}^0} \|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_2 + \beta\tau\|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_2.$$

Finally, applying the assumption that $\tau > \frac{A_{\mathbf{\Lambda}^0}(\varepsilon_1)}{\lambda_{\min}^0}$ and dividing both sides by $\|\mathbf{\Lambda}^* - \mathbf{\Lambda}^0\|_2$ yields the result. \square

It is likely that this error bound could be improved, and such improvements are ripe ground for further work.

2.6 An Algorithm for Sample Selection

Finally, we present an algorithm that uses these parameter estimates from the zero-inflated Poisson assignment model to restrict the data to an appropriate sample in which there is sufficient overlap, and construct propensity scores in that overlap sample. It is well known that if $X_k \stackrel{iid}{\sim} \text{Poisson}(\lambda_k)$ for $k = 1, \dots, K$, then:

$$\left(X_1, \dots, X_K \right) \Big| \sum_{k=1}^K X_k = N \sim \text{Multinomial}(\varrho_1, \dots, \varrho_K; N), \text{ where } \varrho_k = \frac{\lambda_k}{\sum_{k'=1}^K \lambda_{k'}}.$$

If we knew $Z_{jr} = 1$ for all j and r , we could immediately use this transformation to construct propensity scores encoding the model-based probability that individual i in strata r is assigned to unit j :

$$\varrho_{jr}^{\text{full sample}} = \frac{\lambda_{jr}}{\sum_{k=1}^J \lambda_{kr}} = \Pr(D_{ij} = 1 \mid S_{ir} = 1).$$

Because the direct approach is infeasible, we consider decision rules to determine an overlap sample for which a notion of risk is minimized. In this setting, decision rules map evidence – here the matrix of observed counts \mathbf{N} – to a binary decision matrix $\mathbf{\Delta} = \boldsymbol{\delta}^\mu \boldsymbol{\delta}^{\nu'}$, where $\boldsymbol{\delta}^\mu$ is a $J \times 1$ vector of indicators and $\boldsymbol{\delta}^\nu$ is an $R \times 1$ vector of indicators taking the parameters of the model $\boldsymbol{\theta}$ as given. These indicators are set to one when the researcher includes unit j or stratum r , respectively. Note that the decision rule $\mathbf{\Delta}$ is restricted in this way (the outer product of two binary vectors) because the aim of this procedure is to determine the maximal subset of units and randomization strata such that the probability of assignment to each unit is likely to be nonzero in each included stratum (the minimal requirement for overlap).

We frame the sample selection problem as the task of minimizing the risk associated with a particular loss function. To begin, associate each jr cell with a (nonrandom) weight κ_{jr} that reflects the importance the researcher places on cell jr (the matrix \mathbf{W} collects all $J \times R$ weights). The loss from including cell jr when $Z_{jr} = 0$ (a false positive) is proportional to the weight of the cell times a researcher-specified constant $\gamma > 0$, and the loss from omitting cell jr when $Z_{jr} = 1$ (a false negative) is proportional to one times the weight of the cell.

The parameter γ measures the relative cost of false positives and false negatives. The total loss associated with a particular decision rule is given by summing over all cells. We can write the loss from a particular decision rule as:

$$\begin{aligned}\mathcal{L}_\gamma(\Delta \mid \mathbf{W}, \mathbf{Z}) &= \sum_{j=1}^J \sum_{r=1}^R \kappa_{jr} \left[(1 - \delta_j^\mu \delta_r^\nu) Z_{jr} + \gamma \delta_j^\mu \delta_r^\nu (1 - Z_{jr}) \right] \\ &= \sum_{j=1}^J \sum_{r=1}^R \kappa_{jr} \left[Z_{jr} + \delta_j^\mu \delta_r^\nu (\gamma - (1 + \gamma) Z_{jr}) \right] \\ &\propto \sum_{j=1}^J \sum_{r=1}^R \delta_j^\mu \delta_r^\nu \left[\kappa_{jr} \left(\frac{\gamma}{\gamma+1} - Z_{jr} \right) \right] \\ &= \boldsymbol{\delta}^{\mu'} \left(\mathbf{W} \circ \left[\frac{\gamma}{\gamma+1} - \mathbf{Z} \right] \right) \boldsymbol{\delta}^\nu,\end{aligned}$$

where the final two lines omit a constant term that does not depend on Δ . We calculate risk as expected loss conditional on the sample \mathbf{N} and the parameter estimates $\boldsymbol{\theta}$, holding the decision rule Δ fixed:

$$\begin{aligned}\mathcal{R}_\gamma(\Delta \mid \mathbf{N}, \boldsymbol{\theta}) &= \mathbb{E} \left[\mathcal{L}_\gamma(\Delta \mid \mathbf{W}, \mathbf{Z}) \mid \mathbf{N}, \boldsymbol{\theta} \right] \\ &= \boldsymbol{\delta}^{\mu'} \left(\mathbf{W} \circ \left[\frac{\gamma}{\gamma+1} - \widetilde{\boldsymbol{\Pi}} \right] \right) \boldsymbol{\delta}^\nu,\end{aligned}$$

where $\widetilde{\boldsymbol{\Pi}} = \mathbb{E}[\mathbf{Z} \mid \mathbf{N}, \boldsymbol{\theta}] = (\widetilde{\pi}_{jr}^*)_{j=1, r=1}^{J, R}$ is a matrix encoding the posterior expectations of the Z_{jr} variables. The choice of weights is up to the researcher. The simplest choice sets $\kappa_{jr} = 1$ for all cells. Alternately, setting $\kappa_{jr} = \lambda_{jr}$ effectively maximizes the (virtual) number of observations of the overlap sample.

Given a choice of weights, the sample selection procedure is Binary Quadratic Programming (BQP) problem:

$$\left(\widehat{\boldsymbol{\delta}}^\mu, \widehat{\boldsymbol{\delta}}^\nu \right) = \arg \min_{\boldsymbol{\delta}^\mu, \boldsymbol{\delta}^\nu} \begin{bmatrix} \boldsymbol{\delta}^{\mu'} & \boldsymbol{\delta}^{\nu'} \end{bmatrix} \begin{bmatrix} \mathbf{0}_{J \times J} & \mathbf{W} \circ \left[\frac{\gamma}{1+\gamma} - \widetilde{\boldsymbol{\Pi}} \right] \\ \mathbf{0}_{R \times J} & \mathbf{0}_{R \times R} \end{bmatrix} \begin{bmatrix} \boldsymbol{\delta}^\mu \\ \boldsymbol{\delta}^\nu \end{bmatrix} \text{ s.t. } \delta_j^\mu, \delta_r^\nu \in \{0, 1\} \forall j, r.$$

Aside from the constraints that the elements of the decision rules are binary, the BQP is unconstrained. The BQP can be solved using standard integer programming solvers. The output of the BQP determines our estimate of the “largest” feasible overlap set \mathcal{O} :

$$\widehat{\mathcal{O}} = \{jr \text{ s.t. } \widehat{\delta}_j^\mu = 1 \text{ and } \widehat{\delta}_r^\nu = 1\}.$$

Given an estimate of the overlap sample, we form estimated propensity scores as:

$$\widehat{\varrho}_{jr} = \frac{\lambda_{jr}^*}{\sum_{j=1}^J \widehat{\delta}_j^\mu \lambda_{jr}^*} \text{ if } \widehat{\delta}_j^\mu = 1, \widehat{\delta}_r^\nu = 1.$$

Similarly, we construct the estimated weights which enter into the likelihood as:

$$\widehat{w}_{jr} = \widehat{\delta}_j^\mu \widehat{\delta}_r^\nu \times \frac{\sum_{j=1}^J \widehat{\delta}_j^\mu \lambda_{jr}^*}{\lambda_{jr}^*}.$$

2.7 Conclusion

This paper proposes a novel framework for estimating many individual treatment effects, as well as the distribution of those treatment effects. In contrast to much of the literature concerned with estimating treatment effects across many units, which impose both an assumption that assignment of individuals to treatments is unconfounded and functional form or distributional assumptions, the framework of this paper requires only unconfoundedness. In other words, this paper advocates for a *design-based* approach to causal inference, as opposed to a *model-based approach*. As such, estimates produced using the procedures outlined in this paper are robust to functional form mis-specification. The key challenge of this approach is that it relies on obtaining either direct knowledge of, or consistent estimates of, the statistical process that gave rise to the observed assignments of individuals to units. When there are many units, this issue manifests as an explosion in the dimensionality of the propensity score that must be estimated. Further, the high dimensionality of the assignments suggests that typical notions of overlap are not likely to hold. In order to address these issues, we specify a model of assignments that allows for systematic failures of overlap, and suggest an estimator of the assignment probabilities. We then explore the computational and statistical properties of our estimator before suggesting an algorithm for selecting the largest subset of the sample (trimming the sample) such that overlap is likely to hold. The tools developed here have direct application to estimation of many treatment effects (sometimes known as value-added), but other applications are possible. For instance, this procedure can be used in the “judge-IV” context to assess, for example, whether a sample of judges has overlap in the distribution of individuals who may be assigned to their courtrooms. A second possibility is that the procedure could be used to estimate firm- and worker-“effects” in data on the labor market to produce decompositions of the variance in wages that are robust to composition bias.

Chapter 3

Spinning the Wheel: Heterogeneity and Choice in the Provision of Indigent Defense

3.1 Introduction

The United States Supreme Court's landmark 1963 ruling in *Gideon v. Wainwright* established that the Constitution requires states and localities to provide attorneys for criminal defendants too poor to afford legal representation at market rates. In the wake of *Gideon*, the right not just to counsel, but to *effective* counsel, has become a cornerstone of American jurisprudence. In making its ruling, the Court created a constellation of means-tested public assistance programs across federal, state, and local governments that has grown to serve nearly 80% of criminal defendants (Mosteller 2011). Against the backdrop of increased public concern over mass incarceration – the U.S. has the world's highest incarceration rate, at 0.716% (Wamsley 2013) – stark disparities in incarceration rates by race and income (Rehavi and Starr 2014; Carson and Sabol 2012; Lofstrom and Raphael 2016), and a long line of academic research establishing the negative collateral consequences of interaction with the criminal legal system (Mueller-Smith 2015; Aizer and Doyle 2015; Raphael and Smith 2011; Pager 2003, e.g.), some policy makers are considering changes to the administration of these indigent defense systems as part of a larger effort to reform the criminal legal system.

At present, only a handful of academic studies gauge the efficacy of indigent defense provision. A notable exception is Shem-Tov (2020), who uses the plausibly-random assignment of defendants in multiple-defendant cases to public defenders or court-appointed counsel to estimate the relative efficacy of the two systems. He finds that defendants assigned to public defenders are 22% less likely to receive an incarceration sentence relative to those who are assigned to court-appointed counsel. In related research, Agan, Freedman, and Owens (2021) find that controlling for differences in case characteristics eliminates differences in average outcomes between defendants assigned court-appointed counsel and those who privately retain counsel.

This paper explores the extent to which differences in case outcomes across indigent

defendants are attributable to differences in unobserved attorney quality. In contrast to prior studies, this paper does not contrast defendant outcomes between modes of defense; instead, it analyzes the distribution of quality within a single mode of defense. Understanding variability in the quality of indigent defense is of key importance, since such variability is at odds with the normative goals we might wish the criminal punishment system to achieve. Because most indigent defendants cannot choose their attorney, variation in attorney quality may expose some defendants to substantially increased risks of adverse outcomes for reasons entirely unrelated to the facts of their cases. Heterogeneity in attorney quality therefore undermines both horizontal equity (defendants with identical cases should face equal odds and/or expected severity of punishment) and vertical equity (defendants with more severe cases should face higher odds and/or expected severity of punishment) in the criminal legal system. In addition, heterogeneity also undermines the extent to which the criminal legal system produces efficient outcomes, in the sense that the societal benefits of the punishments defendants receive should be weakly greater than the social costs of those punishments.

As in the teacher value-added literature (e.g. Chetty, Friedman, and Rockoff 2014 and Rothstein 2010), a key challenge to obtaining credible estimates of variation in attorney quality is the statistical bias induced by non-random sorting of defendants and cases to attorneys. To overcome this bias, we leverage institutional features of the process by which attorneys are assigned to indigent defendants in Bexar county, Texas. Bexar county is a large, diverse metropolitan county home to San Antonio. In our setting, individuals charged with felonies are assigned quasi-randomly to court-appointed attorneys based on a “wheel” system. After controlling for randomization factors that determine a defendant’s potential pool of attorneys, we find no evidence of systematic sorting of defendants to attorneys of differing quality.

In the first part of this paper, we apply the framework of Chapter 2 to estimate the distribution of treatment effects across attorneys in Bexar county, focusing on both the strategic decision to enter a guilty plea and the ultimate case outcomes (incarceration, probation, or deferral). We find that treatment effects are highly variable across attorneys. For instance, for defendants in felony cases, a one-standard-deviation decrease in attorney quality is associated with a 5.6 percentage-point increase in the probability of incarceration. Variation in treatment effects across attorneys is, in general, not correlated with the observable characteristics of those attorneys. Using estimates of the distribution of attorney treatment effects, we then simulate the effects of policies that either lay off low-predicted-quality attorneys or retain high-predicted-quality attorneys. In both cases, simulations suggest that these personnel policies can meaningfully shift the distribution of case outcomes.

In the second part of this paper, we apply our methodology to gauge the impacts of a first-of-its-kind reform in Comal county, Texas – a suburban county that is part of the San Antonio metropolitan area – that allowed indigent defendants their choice of representation. A near-universal feature of US indigent defense systems is that defendants have little choice over the attorney assigned to represent them. Simple theoretical models that have been applied to the provision of other public services, such as housing and schooling, predict that voucher systems should improve welfare relative to systems with no choice. Models of the provision of indigent defense would suggest that clients, who are better informed

about their potential needs and preferences, should realize increased welfare from a system of client choice than a system of random assignment. If defendants are informed about the quality of attorneys, the demand for attorney services under a client choice program should depend on attorney effectiveness, creating an incentive for attorneys to provide better quality services to clients. However, if clients are ill-informed about the quality of the attorneys they are choosing between, a system of client choice could decrease welfare: in the educational context, Abdulkadiroğlu, Pathak, and Walters (2018) find that school choice reduced student achievement.

Our results suggest that at best, the introduction of client choice had zero impact on the distribution of attorney quality. Rather, we find that clients choose attorneys based on factors largely orthogonal to quality (by, for instance, heavily discriminating against black attorneys). Our results suggest that while there is a role for policies aimed at improving attorney quality, those policies must be carefully designed. Without providing better information to defendants, client choice policies are not likely to produce improvements in the quality of representation.

3.2 Institutional Background

The data used in this paper is drawn from two neighboring counties in Texas that form part of the San Antonio metropolitan area. The first, Bexar county, is home to San Antonio proper. As of 2010, the population of Bexar county was over 1.7 million. The second, Comal county, is home to several suburbs of San Antonio, including New Braunfels. As of 2010, the population of Comal county was slightly less than 110,000.

Indigent Defense in Bexar County

Because the administration of indigent defense is handled independently at the federal, state, and local levels, its provision is relatively heterogeneous. Some jurisdictions maintain public defenders' offices, which employ full-time specialist attorneys who handle the bulk of indigent criminal defense in a particular area. Other jurisdictions rely on "court-appointed counsel" systems, in which judicial officers maintain a list of qualified private attorneys who have agreed to represent indigent clients on a fixed fee schedule. In many cases, attorneys are assigned to clients based on a conditionally-random "wheel" system. In such a system, the court maintains an ordered list of attorneys qualified for particular types of cases and assigns those attorneys to defendants by their position on that list, with the attorney who was least recently assigned to a case assigned first. Bexar county relies almost entirely on the latter system.

When a criminal defendant requests appointed counsel, a determination of indigence is made by Bexar County Pre-Trial Services. According to a 2010 study of Bexar county's indigent defense system that examined cases from 2008-2009, defendants are presumed indigent if their incomes fall below 125% of the federal poverty line, although defendants with incomes higher than this threshold may still be found to be indigent. A defendant's ability to post bail is not considered in the determination of indigence, but counsel must be appointed

sooner if a defendant remains in custody. According to the study, roughly 85% of defendants who apply for appointed counsel automatically qualify, and many of the remainder eventually qualify after a court determination of indigence (Texas Task Force on Indigent Defense 2010).

After the determination of indigence, an attorney is selected from one of five appointment lists (“wheels”) maintained by the court. In order to remain on a felony appointment list, an attorney must complete 10 hours of continuing legal education in criminal law every year. In addition, minimum experience requirements for each of the five lists apply. Importantly, attorneys must accept the cases assigned to them (except in special circumstances) in order to remain on an appointment list. Finally, attorneys must be approved by a majority of judges who preside over relevant cases. By state law, the default system for appointment of counsel is rotational: unless the court makes a “finding of good cause” to deviate, courts “shall appoint attorneys from among the next five names on the appointment list in the order in which the attorney’s names appear[.]”¹ The 2010 study found significant deviations from this procedure for the assignment of counsel in misdemeanor cases, but adherence to this method for felony cases. The report specifies that defendants in the 144th, 175th, 187th, 227th, 290th, 399th and 437th District Courts are always appointed counsel from the applicable wheel, while defendants in the 186th, 187th, and 226th District Courts assign defendants to attorneys who are present in the courtroom. Figure 3.1, which reproduces a flow chart from the 2010 Report, illustrates the attorney assignment process for defendants charged with felonies in Bexar county. Given this, we restrict our main analysis sample in Bexar county to felony cases from 2008-2018 in District Courts that always assign counsel from the wheel system. Attorneys appointed to represent defendants have the choice between flat fee and hourly compensation; the 2010 study found that the “vast majority” of attorneys chose flat fee compensation. Attorney compensation has been updated several times since the beginning of our sample period.

Indigent Defense in Comal County & the Client Choice Program

Like Bexar county, Comal county maintained a rotational court-appointment system for the assignment of counsel to indigent defendants until 2015. Judges in Comal county maintain three appointment lists of attorneys – one for serious felonies, one for less serious felonies, and one for misdemeanors. While less information is readily available about the assignment process in Comal county, a 2017 Justice Management Institute (JMI) report providing an initial evaluation of the Client Choice Program states that exceptions to the rotational system were granted only in extenuating circumstances (Justice Management Institute 2017).

Beginning in the first week of February 2015, Comal county began a pilot of the Client Choice Program. The Client Choice Program is essentially a voucher system for indigent defendants, with some caveats. By Texas statute, indigent criminal defendants in Texas must have counsel appointed by a judge or the designate of a judge (e.g. Pre-Trial Services), and so defendants are not allowed to shop for counsel. Rather, defendants were asked by

¹ Tex. Code Crim. Proc. art. 26.04(a)

magistrate judges if they wished to choose counsel or have counsel appointed for them (using the pre-existing rotational system). According to the JMI Report, in the year after the initial implementation of client choice, 72% of defendants elected to choose their attorney, rather than have an attorney assigned to them by a magistrate. Those who elect to choose are granted leave – usually around fifteen minutes to half an hour – to review a binder of “Lawyer Information Forms,” which contain basic information about each attorney on the appointment list, including: attorney name, law firm, languages spoken, types of cases handled, number of defendants represented in the past 12 months, and an explanation of any disciplinary history. Defendants are not allowed to interview potential counsel. Defendants then provide the magistrate with a ranked list of three choices, from whom an attorney is appointed based on preference rank and availability. According to the JMI Report, the lead magistrate expressed the opinion that the information provided was insufficient to support an informed choice, and that defendants instead chose based on “word of mouth” at the jail (Justice Management Institute 2017). The reform did not affect other aspects of the attorney appointment process, including compensation.

3.3 Data and Summary Statistics

Data Sources

This paper makes use of three separate data sources: 1) administrative data on case outcomes from Bexar county, 2) data compiled from online case records in Comal county, and 3) Texas state bar data on attorney characteristics. Bexar county publicizes detailed administrative records for each case handled by county courts. These records include each defendant’s name, race/ethnicity, gender, date of birth, unique state correctional identification number, detailed information about the offense(s) the defendant was charged with, the identity (name and bar number) of the attorney representing the defendant, whether the attorney was appointed or retained, the courtroom to which the defendant was assigned, detailed information about the disposition of the case, and ultimate sentence (if any). Similar data were constructed for Comal county by accessing individual publicly-available case description pages hosted by the county government. After collecting records on felony and misdemeanor cases for the sample period, we extracted information including defendant characteristics (race, age, gender), attorney identity (name and method of appointment), charged offense(s), courtroom, judge, and final disposition and sentence. We merge our data on criminal cases to records maintained by the Texas State Bar that measure attorney characteristics. These data provide information on attorney gender, race, graduation year, and the law school from which that attorney graduated. In order to provide a rough measure of educational quality, we associate each law school observed in our sample with the most recent US News and World Report ranking for that school.

Summary Statistics

Table 3.1 reports descriptive statistics for our primary analysis samples in Bexar and Comal counties. In our analysis, we focus on felony cases filed between 2008 and 2018. In Bexar county, we drop cases assigned to the 186th, 187th, and 226th District Courts due to the documented deviations in the attorney assignment procedures in those courts from the quasi-random “wheel” system. After applying these restrictions, our initial analysis sample contains 41,574 cases in Bexar county and 3,805 cases in Comal county, respectively. We then apply the propensity score estimation and sample selection procedure outline in Chapter 2 on both the Bexar and Comal samples, with γ set to 0.95. The procedure results in significant reductions in sample size in both counties: after trimming, 30,440 cases remain in the Bexar analysis sample, while just 1,653 remain in the Comal sample. For each of these samples, Table 3.1 reports the means and standard deviations (where informative) of a basic set of defendant characteristics. Despite the decreases in sample size, the trimming procedure leaves these summary statistics basically unchanged. In both counties, just over 75% of defendants are male. 30% of the defendants in Bexar county identify as white, while 18% identify as Black (about twice the share of the population in Bexar county that identifies as Black, which was 8% as of 2010) and 51% identify as Latino (slightly lower than the share of the population in Bexar county that identifies as Latino, which was 56% as of 2010). Comal county’s court records do not differentiate defendants of Latino ethnicity. Just 7% of defendants identify as Black in Comal county. The mean age of defendants in both counties is roughly 33, with defendants in Bexar county slightly younger and defendants in Comal county slightly older.

Table 3.2 provides averages of attorney characteristics in both counties, weighted at the case level. There are 624 and 95 attorneys in the initial Bexar and Comal samples, respectively, which trimming reduces to 404 and 49. Importantly, demographic data is available for most but not all attorneys, and so the averages are taken over differing sub-samples of the data depending on missing-ness. Like defendants, just over three quarters of attorneys in both counties are male. Roughly 60% and 30% of attorneys in both counties are white and Latino, respectively. Attorneys in Comal county graduated slightly earlier, and from marginally better-ranked law schools, than their peers in Bexar county.

Finally, Table 3.3 provides a tabulation of case outcomes in both counties. A staggering 89% of defendants enter a plea of guilty in their cases in Bexar county, while just 63% of defendants in Comal county enter guilty pleas (74% in the trimmed sample). Roughly similar shares of defendants receive incarceration sentences in both counties: 51% of defendants in Bexar county and 48% of defendants in Comal county. Defendants are more likely to be sentenced to probation in Comal county (24%) than in Bexar county (14%), but are far less likely to receive a deferred adjudication (5% in Comal, 27% in Bexar). In both counties, the modal case disposition is a guilty plea followed by a period of incarceration. Defendants in Comal county are more likely to be sentenced to long incarceration spells than their counterparts in Bexar county: about one quarter of all defendants in Comal county are sentenced to incarceration periods greater than four years, while 14% of defendants in Comal county receive incarceration sentences greater than four years.

3.4 Results

Variation in Treatment Effects Across Attorneys

We now turn to our first set of results: estimates variation in treatment effects across attorneys. Following the framework of Chapter 2, we produce estimates \hat{p}_{jy} for each binary outcome of interest y by taking inverse-propensity-score-weighted averages of case outcomes for each attorney across randomization strata. In our setting, randomization strata S_i correspond to unique combinations of felony class, year, and courtroom. We adopt Assumption 2.3, which states that the true parameters are independent draws from a common distribution, such that for each outcome y ,

$$p_{jy} \stackrel{iid}{\sim} G_y(\cdot).$$

Under this assumption moments of $G_y(\cdot)$ of up to order J are identified. The identified moments of $G(\cdot)$ can be used to construct an unbiased estimate of the variance of treatment effects across units:

$$\widehat{\text{Var}}(p_{jy}) = \frac{1}{J-1} \sum_{j=1}^J \left(\hat{p}_{jy} - \frac{1}{J} \sum_{j=1}^J \hat{p}_{jy} \right)^2 - \frac{1}{J} \sum_{j=1}^J \text{SE}(\hat{p}_{jy}),$$

where $\text{SE}(\hat{p}_{jy})$ is an unbiased estimate of the standard error of \hat{p}_{jy} . It can be shown that $\widehat{\text{Var}}(p_{jy})$ is a U-statistic, and therefore:

$$\sqrt{J} \left(\widehat{\text{Var}}(p_{jy}) - \text{Var}(p_{jy}) \right) \rightsquigarrow N(0, V_\sigma/J).$$

We construct an unbiased estimate of V_σ following the method of Wang and Lindsay (2014), and use that estimate to test against the null hypothesis of no heterogeneity.

In order to provide baseline evidence for the validity of the research design, Table 3.4 reports tests for heterogeneous treatment effects of attorneys on *pre-treatment* variables. Under the hypothesis that assignment of cases to attorneys is as good as random conditional on strata S_i , the variance of attorney “effects” on these pre-treatment variables should be zero. To find otherwise would suggest that there is systematic sorting of cases to attorneys on the basis of these characteristics. We discretize each of the pre-treatment controls, and then estimate variance components and conduct tests against the null hypothesis of no heterogeneity using the methods described above. For each pre-treatment control, Table 3.4 reports the estimated variance and standard deviation of attorney effects, the standard error of the estimate of the variance, and a p-value for the test of no-heterogeneity. We find no evidence of systematic sorting of cases to attorneys: each of the estimated variance components is statistically insignificant at the 5% level, and a test of joint significance comfortably fails to reject the null hypothesis that all variance components are identically zero ($p = 0.26$). These results suggest that the observed variation in defendant characteristics across attorneys is not large enough to rule out that that variation is driven by sampling uncertainty alone.

We now turn to estimates of the variation in attorney treatment effects on case outcomes of interest, which are reported in Table 3.5 and were produced using identical methods. In contrast to Table 3.4, nearly all estimates reported in Table 3.5 are highly significant: the p-values associated with tests of no heterogeneity are significant at the 2% level for all but two outcomes. A joint test of significance resoundingly rejects the null hypothesis that all variance components are zero ($p = 0.001$). These results suggest that the observed variation in case outcomes is almost certainly too large to be driven by sampling variation alone.

What do the results suggest about the scale of variation in attorney quality? Across outcomes, the answer is that attorney treatment effects vary substantially. One standard deviation in attorney effects on the probability of entering a guilty plea is 3.4 percentage points. This variation could reflect strategic decisions on the part of defense counsel, or it could reflect varying levels of effort on behalf of defendants. Anecdotal evidence suggests that differences in effort between attorneys plays at least some part. Perhaps most shockingly one standard deviation in attorney effects on the probability of receiving an incarceration sentence is 5.6 percentage points. Variation in attorney effects on the length of incarceration is even larger: one standard deviation in attorney effects on the probability of receiving incarceration sentences of greater than one year, greater than two years, and greater than four years are 7.6, 7.0, and 6.1 percentage points, respectively. Taken on face value, these results suggest that “winning the attorney lottery” can indeed make a huge difference in determining the outcome of one’s case. This necessarily imply that a significant component of the variation in punishments people receive when they interact with the criminal justice system are the product of luck, and not necessarily “deservingness.” The results also suggest that there is significant scope for improving defendants’ outcomes by implementing personnel policies that aim to shift the distribution of attorney quality.

Is estimated attorney quality correlated with observable characteristics of attorneys? Table 3.6 reports regressions of estimated attorney effects on the probabilities of entering a guilty plea, incarceration, probation, and deferred adjudication on a vector of attorney characteristics gathered from the Texas Bar. While there are some significant associations, the overarching story is that attorney quality is difficult to predict on the basis of observables alone. Perhaps counterintuitively, younger (or less experienced) attorneys seem to produce better outcomes for their clients than do older (or more experienced) attorneys. This pattern may be suggestive of a form of dynamic negative selection into indigent defense work: high quality young attorneys who enter into indigent defense work may eventually graduate out of the court-appointed system into better-paying roles, while lower quality attorneys may not receive similar opportunities.

The Full Distribution of Attorney Treatment Effects and Policy Simulations

Next, we consider estimation of the underlying distribution of attorney effects, $G_y(\cdot)$, and again following the framework of Chapter 2. The log integrated likelihood can be written:

$$\ell_y^*(G) = \sum_{j=1}^J \log \left(\int f_{jy}(p | Y_{jy}) dG(p) \right).$$

where $f_{jy}(p | Y_{jy})$ is the (weighted) likelihood of observing attorney j 's vector of case outcomes Y_{jy} for outcome variable y . We produce estimates of the distribution of treatment effects across units by maximizing $\ell_y^*(\cdot)$ with respect to $G_y(\cdot)$, following the methods proposed by Koenker and Mizera (2014) and Efron (2016). Specifically, assume that $G_y(\cdot)$ takes the form:

$$G_y(p) = \sum_{k=0}^K g_{ky} \times \mathbf{1}[p \leq k/K], \text{ with } G_y(1) = 1, g_{ky} \geq 0 \forall k,$$

for some K relatively large. This assumption restricts the support of $G_y(\cdot)$ to a fine grid of points. Given this assumption, the likelihood becomes:

$$\ell^*(G) = \sum_{j=1}^J \log(F_j'G),$$

where $F_j = (f_{jy}(p_k | Y_{jy}))_{k=1}^K$ and $G = (g_k)_{k=1}^K$. Let $F = (F_1, \dots, F_J)$. Estimation of G proceeds via nonlinear convex programming:

$$\min_G -\mathbf{1}'_J \log(F'G) \text{ s.t. } \mathbf{1}'_K G = 1, G \geq 1.$$

The distributions produced by this routine are “spiky,” with estimated $g_k > \epsilon \approx 0$ for approximately $\log(J)$ points only. Efron (2016) proposed an empirical bayes deconvolution estimator that is essentially a smoothed version of Koneker and Mizera’s estimator that imposes additional (exponential family) structure on the g_k . Specifically, Efron sets:

$$g_k = g_k(\alpha) = \exp(Q'_k \alpha - \phi(\alpha)), \text{ with } \phi(\alpha) = \log\left(\sum_{k=0}^K \exp(Q'_k \alpha)\right).$$

Here, α is a p -dimensional parameter vector and $Q = (Q_0 \dots Q_K)$ is a known $p \times K + 1$ design matrix. The full procedure specifies $Q'\alpha$ as a spline in p_k and imposes a penalty function on α .

These methods allow for greatly increased flexibility over standard parametric methods. Gilraine, Gu, and McMillan (2020) show that in the context of teacher value added, making parametric assumptions can lead to overstated predictions of gains from personnel policies. We produce estimates of $G_y(\cdot)$ for each outcome y three ways: 1) without restrictions, 2) imposing exponential family structure, and 3) imposing an approximate restriction that the log-odds of attorney effects are normally distributed (a standard random effects logit).

Figure 3.2 plots the estimated distribution of attorney effects on the probability of entering a guilty plea, while Figure 3.3 plots the estimated distribution of attorney effects on the probability of incarceration. The shape of estimated distributions of treatment effects on guilty pleas vary somewhat by estimation strategy, while there is little variation across estimators in the shape of the estimated distribution of attorney effects on incarceration. Table 3.7 compares the moments of these estimated distributions.

With estimates of the full distribution of attorney effects in hand, we next simulate the effects of two policies aimed at shifting the distribution of attorney quality in the hopes of

improving outcomes for defendants: layoff policies and retention policies. In each of these simulations, we assume that attorneys can be readily replaced with fresh draws from a stable $G(\cdot)$. Given this assumption, the expected reduction in the incarceration rate from laying off attorneys at or above the q -th quantile of $G(\cdot)$ is:

$$\mathbb{E}\left[(p_j - \mathbb{E}[p_j]) \times \mathbf{1}[p_j \geq G^{-1}(q)]\right],$$

while the expected reduction in the incarceration rate from retaining attorneys at or below the q -th quantile of $G(\cdot)$ is:

$$\mathbb{E}\left[(p_j - \mathbb{E}[p_j]) \times \mathbf{1}[p_j < G^{-1}(q)]\right].$$

Since the p_j are not known, we instead simulate 10,000 random draws from the estimated $G(\cdot)$ distributions:

- For each draw p_j^b , take 10 draws from a Bernoulli(p_j^b) (simulated outcomes).
- Construct empirical bayes posterior predictions \hat{p}_j^{EB} given the simulated outcomes.
- Compute averages of *true* p_j above or below quantiles of \hat{p}_j^{EB} .

We plot the results of our simulations of layoff policies in Figure 3.4 and our simulations of retention policies in Figure 3.5. In each of these figures, the x-axis represents the fraction of attorneys either laid off or retained, and the y-axis represents the expected reduction in the incarceration rate from implementing the policy. We simulate the policy using each of the three estimated distributions (unrestricted, exponential family, logit-normal), and find similar results under all three scenarios. Both policies can reduce the incarceration rate, although the layoff policy achieves the same expected reductions in the incarceration rate. The simulations suggest that a policy in which the bottom 5% of attorneys (ranked by EB posteriors) would lead to a nearly 0.7 percentage point reduction in the overall incarceration rate. In Bexar county, over the period of this study, that would amount to laying off the 20 lowest-performing attorneys, and reducing the number of incarceration spells by roughly 200 (0.7% of $\sim 30,000$ cases). These simulations do not account for effects on any of the other outcomes.

Evaluating the Client Choice Program

We now turn to the final part of our analysis, an evaluation of the Client Choice Program. To evaluate the program, we estimate the distribution of attorney quality in both Bexar and Comal counties using the methods outlined above, but restricting to years prior to 2015 (before the program was first piloted). We then construct empirical bayes posterior estimates of attorney quality for each attorney in the sample. If clients made informed choices about their attorneys, then it should be the case that better attorneys are chosen more often in Comal county after the instatement of the Client Choice Program, such that the distribution of attorney quality (as proxied here by EB posterior means) should shift differentially between

Comal county, where defendants were given the ability to choose, and Bexar county, where they were not. We investigate this hypothesis by fitting a basic event study regressions.

Before presenting the event study results, we first explore whether the onset of the Client Choice Program was associated with any changes in the composition of attorneys representing clients in Comal county. As a summary measure of composition, we compute the Herfindahl-Hirschman Index (HHI) of individual attorney's shares of overall representation in both Comal and Bexar counties. Figure 3.6 plots the time series of representation HHI by county, normalizing values by their 2014 levels. While the HHI of attorney representation is near constant in Bexar county throughout the sample period, the HHI of attorney representation sharply and discontinuously increased in Comal county after the onset of the Client Choice Program. This sharp increase suggests that the program indeed changed the composition of attorney representation, in particular that a small number of attorneys were retained by many clients after those clients were given the ability to choose counsel.

Despite the clear change in the distribution of cases to attorneys, there is little evidence that the Client Choice Program actually shifted the distribution of attorney quality. Figure 3.7 and Figure 3.8 plot event study estimates of the effects of the program on the distribution of attorney effects for pleas and incarceration, respectively (controlling for all case characteristics and a linear trend). In neither case is there any evidence of a differential shift in mean predicted attorney quality between Comal and Bexar counties.

The event study evidence suggests that either clients were unaware of differences in quality between attorneys, or did not value those differences in quality. What characteristics of attorneys, if any, are associated with a greater probability of selection after the introduction of choice? In order to gauge how defendants value various attributes, including attorney quality, we conduct a simple discrete choice analysis. We assume clients make their choice of attorney by maximizing the utility of that choice. Specifically, let d_{ij} denote an indicator equal to 1 if client i chooses attorney j . Assume client i assigns utility $v_j + \epsilon_{ij}$ to choice j , where v_j measures the common component of utility across clients for attorney j and $\epsilon_{ij} \stackrel{\text{iid}}{\sim} EV_1$ measures the idiosyncratic client-specific component of utility for attorney j . Given this formulation, the probability of any choice can be written:

$$\Pr(d_{ij} = 1) = \frac{\exp(v_j)}{\sum_{\ell=1}^J \exp(v_\ell)}.$$

We proxy the probability of selection, $\Pr(d_{ij} = 1)$, by the empirical fraction of cases represented by, or market share of, each attorney in Comal county after the introduction of client choice. Denote the market share for attorney j by S_j . We can then write:

$$\log(S_j/S_\ell) = v_j - v_\ell$$

for any pair of attorneys j and ℓ . We further assume that the common component of utility is a linear index of attorney characteristics: $v_j = Z_j' \gamma + \xi_j$, where Z_j includes observable characteristics, and the ξ_j includes all other factors that affect choice (but that we do not observe - such as word of mouth). Fix a reference attorney, $\ell = \text{ref}$. Our assumptions imply:

$$\log(S_j/S_{\text{ref}}) = \dot{Z}_j \gamma + \dot{\xi}_j,$$

where $\dot{Z}_j = Z_j - Z_{\text{ref}}$ and $\dot{\xi}_j = \xi_j - \xi_{\text{ref}}$.

Table 8 reports the results from estimating γ by OLS in the sample of attorneys selected by clients in Comal county after the introduction of choice. Because the number of observations is limited, we first estimate coefficients on each characteristic in isolation, then estimate the equation including all characteristics. The results suggest that, to some extent, client choices are explained by attorney characteristics. In particular, clients seem to 1) weakly prefer Latino attorneys, 2) strongly dis-prefer Black attorneys, and 3) dis-prefer younger attorneys (later graduation years). Of these characteristics, experience was found to *negatively* correlate with predicted attorney quality, while race was not associated with differences in quality. Finally, predicted quality *emph* is positively associated with the probability of choice, but the coefficients are highly insignificant. While the estimates of γ are extremely noisy, the results are suggestive evidence that clients are not well informed about attorney quality, and therefore resort to choosing attorneys on the basis of characteristics that are either uncorrelated with quality, or worse, negatively correlated with quality.

Taken together, the evidence suggests that the Client Choice Program likely did not meaningfully shift the distribution of attorney quality in Comal county. At best, the program may have simply re-shuffled the distribution of cases to attorneys, while keeping average quality constant. At worst, the program may have actually reduced overall attorney quality, since clients appear to make uninformed choices over attorneys.

3.5 Conclusion

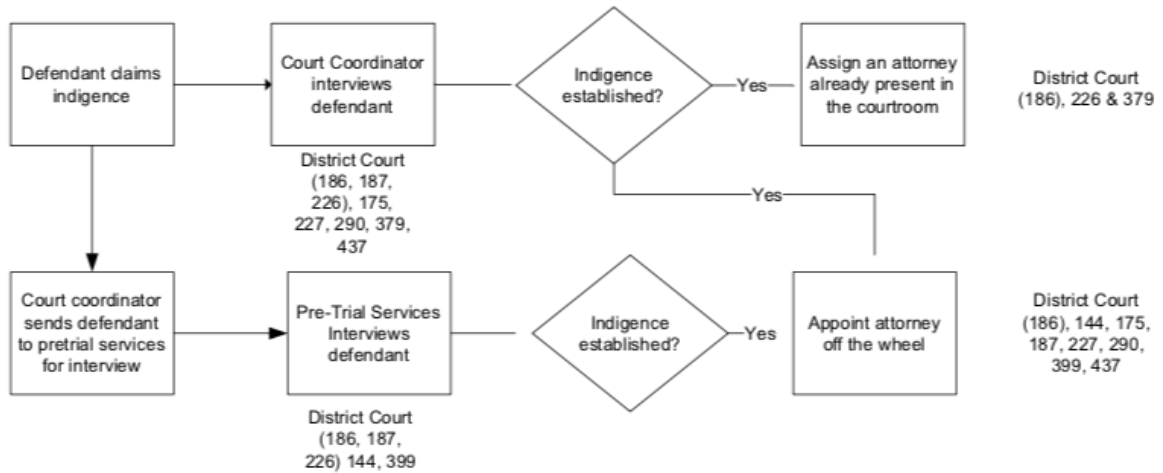
This paper investigated the distribution of attorney quality in the context of indigent defense. We estimate the distribution of attorney quality using tools developed in Chapter 2 and leveraging known institutional features of the assignment process for assigned counsel. We find that there is substantial heterogeneity in treatment effects of attorneys: some attorneys systematically produce better outcomes for clients than others. For instance, a one standard deviation decrease in attorney quality is associated with a 5.6 percentage point increase in the probability that a defendant will be sentenced to a period of incarceration. In simulations, we show that layoff and retention policies can improve outcomes even when decisionmakers are acting on relatively little information.

In future work, it will be crucial to explore whether policies that attempt to align attorney and client incentives, like altering fee structures for compensating attorneys who take on indigent defense work, ultimately produces better outcomes for clients. We test the effects of one such policy: the introduction of client choice in Comal county, Texas. We find that the introduction of choice at best had zero impacts on the distribution of attorney quality, and at worst reduced attorney quality. We found that allowing clients to choose their attorneys had at best zero effect on the distribution of attorney quality. When given the choice, clients discriminate against black attorneys and discriminate in favor of older attorneys. Policies that rely on defendants to ascertain which attorneys are best without providing reliable information on attorney quality are unlikely to produce better outcomes.

Figures

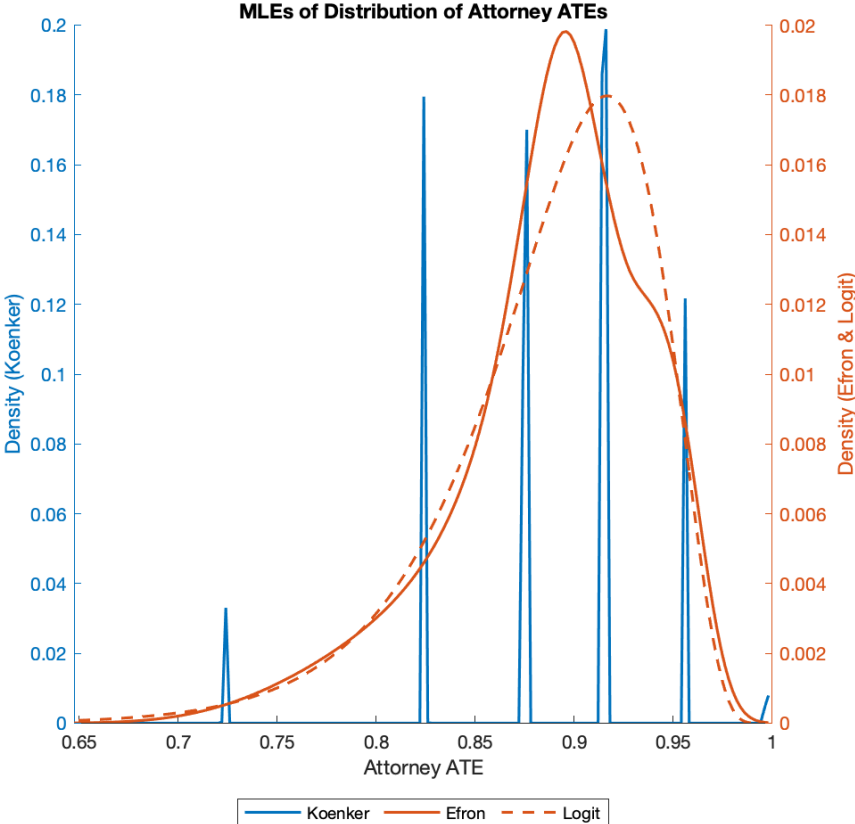
Figure 3.1: Flowchart of Case Assignment Process in Bexar County

Figure 4: District Court in Court Assignments



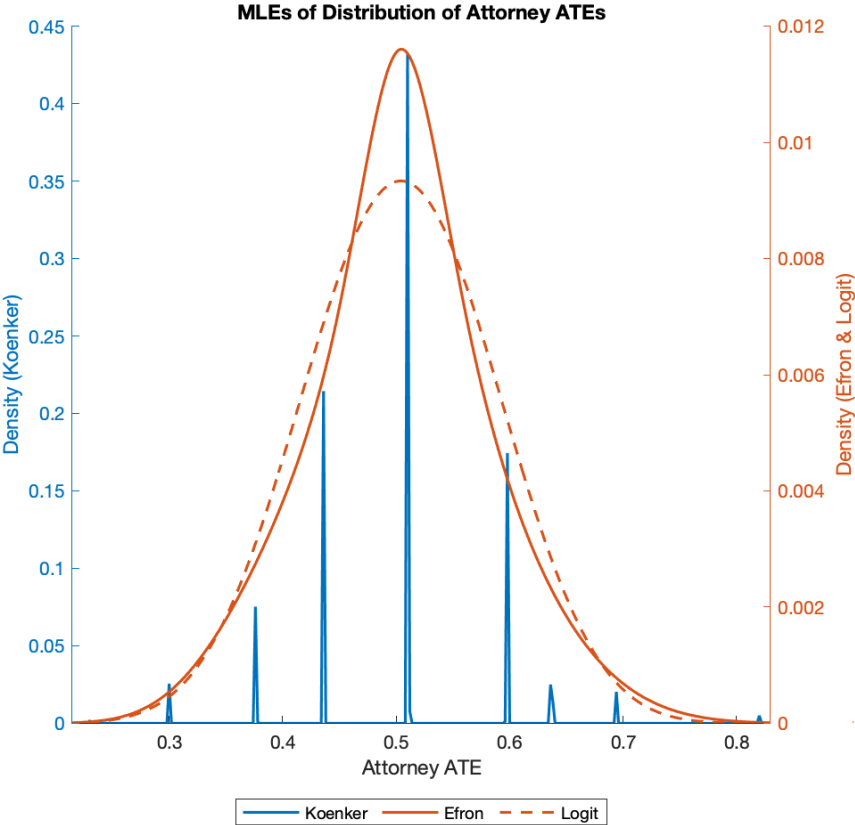
Note: This Figure reproduces a flowchart from the the 2010 Task Force Report illustrating the process by which defendants in felony cases are assigned to attorneys. Because cases in courtrooms 186, 226, and 379 do not assign counsel via the wheel system, we drop those cases from the sample.

Figure 3.2: Maximum Likelihood Estimates of Distribution of Attorney Effects for Pleas

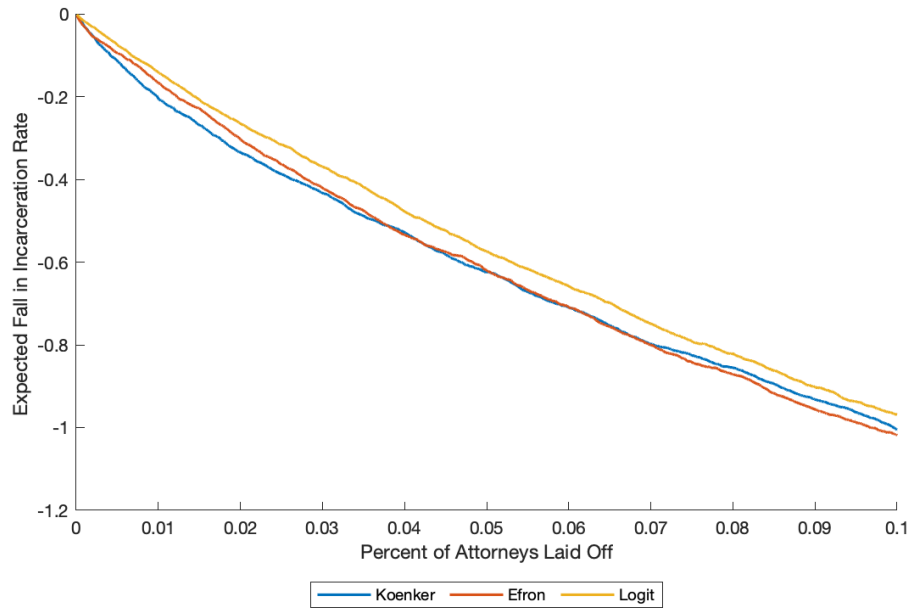


Note: This Figure plots three maximum likelihood estimates of the distribution of attorney effects on the probability that clients enter a guilty plea. The distribution plotted in blue is the unrestricted nonparametric maximum likelihood estimate. The distribution plotted with a solid orange line is exponential family spline estimate. The distribution plotted with a dashed orange line is the logit-normal estimate.

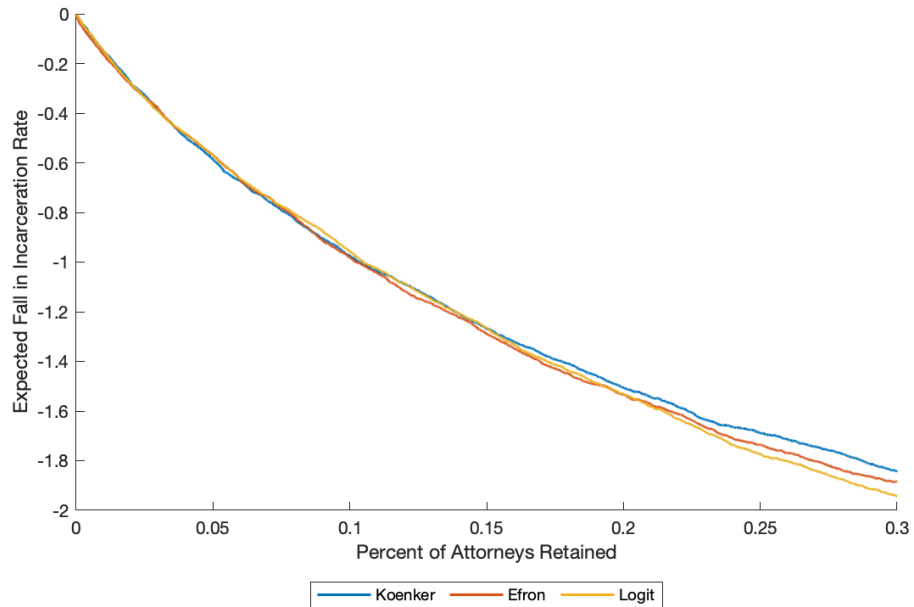
Figure 3.3: Maximum Likelihood Estimates of Distribution of Attorney Effects for Incarceration



Note: This Figure plots three maximum likelihood estimates of the distribution of attorney effects on the probability that clients receive a sentence of incarceration. The distribution plotted in blue is the unrestricted nonparametric maximum likelihood estimate. The distribution plotted with a solid orange line is exponential family spline estimate. The distribution plotted with a dashed orange line is the logit-normal estimate.

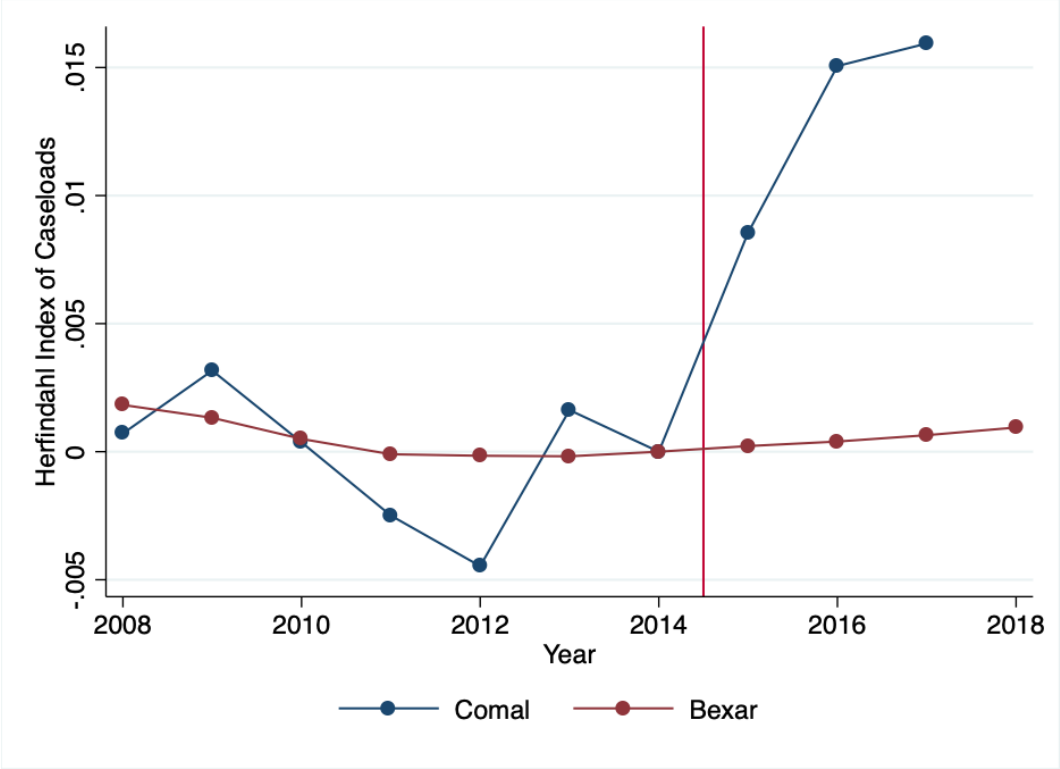
Figure 3.4: Simulated Reduction in Incarceration Rate from a Layoff Policy

Note: This Figure plots the results of simulating layoff policies assuming that true attorney quality is distributed according to each of the three maximum likelihood estimates of $G(\cdot)$. The plotted functions represent the expected decrease in the incarceration rate from laying off the bottom x -percent of attorneys, as ranked by posterior EB means.

Figure 3.5: Simulated Reduction in Incarceration Rate from a Retention Policy

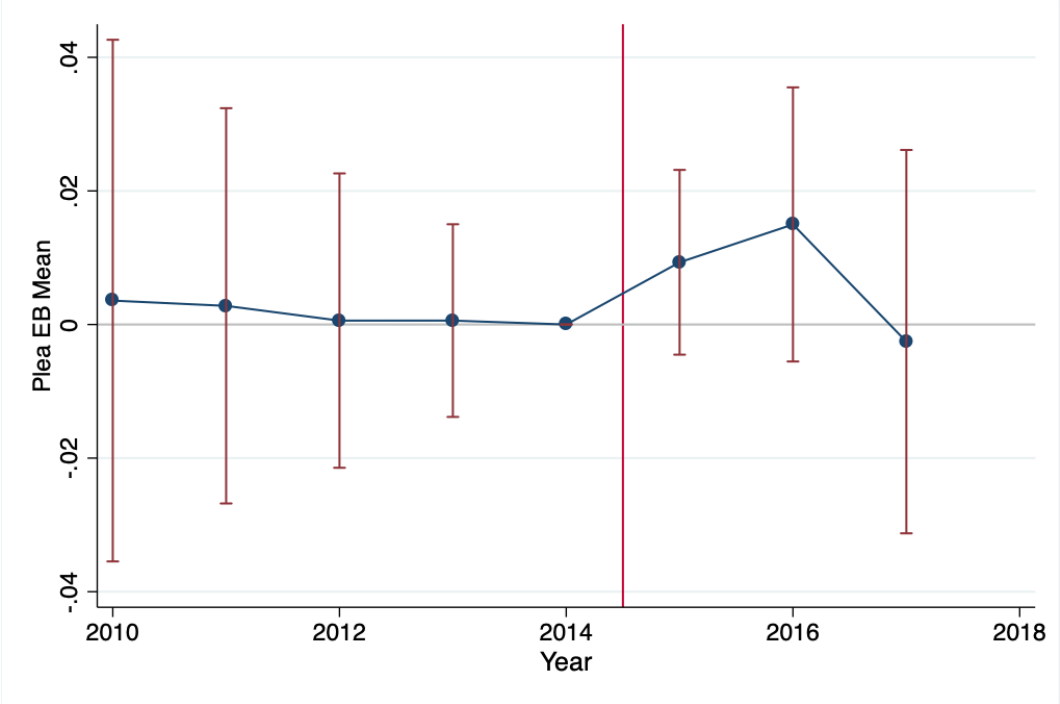
Note: This Figure plots the results of simulating retention policies assuming that true attorney quality is distributed according to each of the three maximum likelihood estimates of $G(\cdot)$. The plotted functions represent the expected decrease in the incarceration rate from retaining only the top x -percent of attorneys, as ranked by posterior EB means.

Figure 3.6: Time series of HHI of Attorney Representation



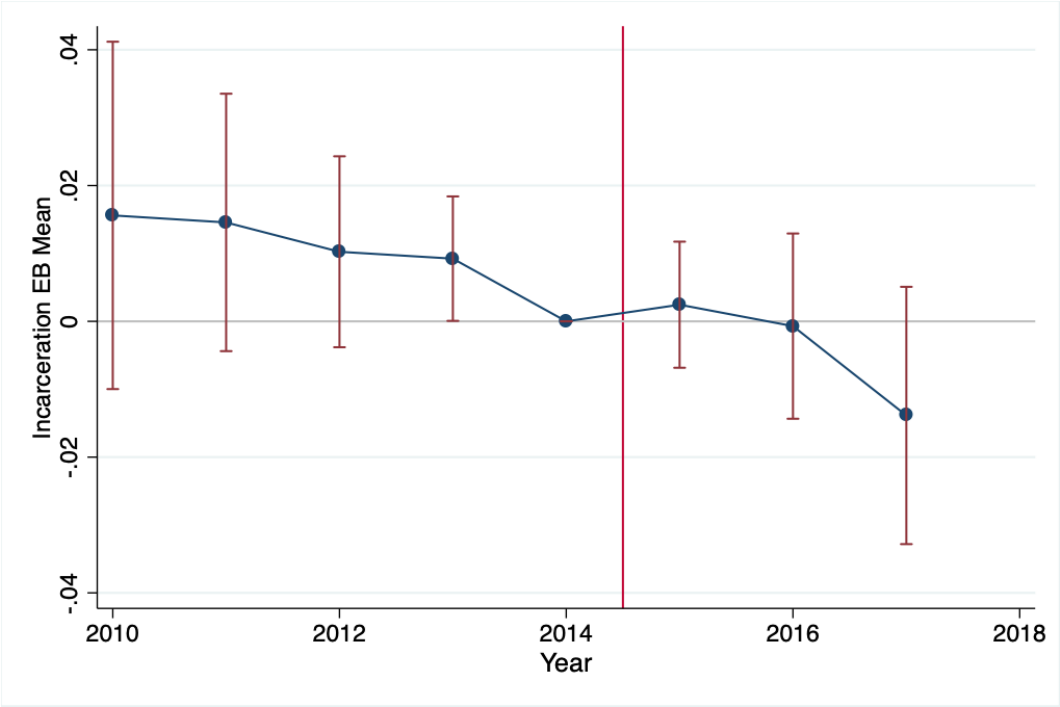
Note: This Figure plots the time series of the Herfindahl-Hirschman Index of attorney representation in Comal and Bexar counties.

Figure 3.7: Event Study: Effect on Distribution of Attorney Effects on Probability to Plea



Note: This Figure plots event study coefficients quantifying the differential change in average posterior predictions of attorney effects on the probability to enter a guilty plea between Comal and Bexar counties.

Figure 3.8: Event Study: Effect on Distribution of Attorney Effects on Probability of Incarceration



Note: This Figure plots event study coefficients quantifying the differential change in average posterior predictions of attorney effects on the probability of incarceration between Comal and Bexar counties.

Tables

Table 3.1: Defendant Summary Statistics

| | Bexar - Full | | Bexar - Trimmed | |
|-------------------|--------------|-----------|-----------------|-----------|
| | Mean | Std. Dev. | Mean | Std. Dev. |
| Male | 0.77 | | 0.78 | |
| White | 0.30 | | 0.31 | |
| Black | 0.18 | | 0.18 | |
| Latino | 0.51 | | 0.51 | |
| Prior Felony | 0.46 | | 0.46 | |
| Prior Misdemeanor | 0.61 | | 0.61 | |
| Age | 33.07 | 10.85 | 33.01 | 10.88 |
| Prior Cases | 8.41 | 17.98 | 8.40 | 18.18 |
| N | 41,574 | | 30,440 | |

| | Comal - Full | | Comal - Trimmed | |
|-------------------|--------------|-----------|-----------------|-----------|
| | Mean | Std. Dev. | Mean | Std. Dev. |
| Male | 0.76 | | 0.77 | |
| White | 0.92 | | 0.91 | |
| Black | 0.07 | | 0.07 | |
| Latino | - | | - | |
| Prior Felony | 0.31 | | 0.34 | |
| Prior Misdemeanor | 0.37 | | 0.36 | |
| Age | 33.76 | 11.17 | 33.32 | 11.31 |
| Prior Cases | 3.44 | 3.65 | 3.47 | 3.45 |
| N | 3,805 | | 1,653 | |

Note: This Table reports summary statistics for defendant characteristics in Bexar and Comal counties. “Full” refers to the initial analysis sample, while “Trimmed” refers to the final analysis sample after conducting the sample selection procedure described in Chapter 2. Data on identification as Latino is not available for defendants in Comal county.

Table 3.2: Attorney Summary Statistics

| | Bexar - Full | | Bexar - Trimmed | |
|--------------------|--------------|-----------|-----------------|-----------|
| | Mean | Std. Dev. | Mean | Std. Dev. |
| Male | 0.74 | | 0.75 | |
| White | 0.57 | | 0.58 | |
| Latino | 0.33 | | 0.33 | |
| Black | 0.05 | | 0.05 | |
| Solo Practitioner | 0.73 | | 0.80 | |
| Disciplinary Hist. | 0.06 | | 0.07 | |
| Graduation Year | 1996 | 11.41 | 1993 | 10.76 |
| Law School Rank | 137 | 54.34 | 56 | 55.55 |
| Experience | 16.21 | 10.73 | 18.79 | 10.23 |
| Prior Cases | 67 | 50.36 | 88 | 48.01 |
| N | 624 | | 404 | |

| | Comal - Full | | Comal - Trimmed | |
|--------------------|--------------|-----------|-----------------|-----------|
| | Mean | Std. Dev. | Mean | Std. Dev. |
| Male | 0.72 | | 0.74 | |
| White | 0.63 | | 0.59 | |
| Latino | 0.29 | | 0.34 | |
| Black | 0.06 | | 0.02 | |
| Solo Practitioner | 0.69 | | 0.83 | |
| Disciplinary Hist. | 0.07 | | 0.13 | |
| Graduation Year | 1995 | 11.12 | 1994 | 10.88 |
| Law School Rank | 118 | 60.39 | 120 | 61.65 |
| Experience | 16.87 | 10.94 | 17.67 | 10.52 |
| Prior Cases | 96 | 73.84 | 138 | 73.89 |
| N | 95 | | 49 | |

Note: This Table reports summary statistics for attorney characteristics in Bexar and Comal counties. “Full” refers to the initial analysis sample, while “Trimmed” refers to the final analysis sample after conducting the sample selection procedure described in Chapter 2.

Table 3.3: Case Summary Statistics

| | Bexar | | Comal | |
|--------------------|--------|---------|-------|---------|
| | Full | Trimmed | Full | Trimmed |
| Plea Guilty | 0.89 | 0.89 | 0.63 | 0.74 |
| Incarceration | 0.52 | 0.51 | 0.43 | 0.48 |
| Probation | 0.13 | 0.14 | 0.19 | 0.24 |
| Deferred Adj. | 0.27 | 0.27 | 0.03 | 0.05 |
| Plea & Incarcerate | 0.51 | 0.51 | 0.41 | 0.46 |
| Incar. > 6m | 0.38 | 0.39 | 0.41 | 0.45 |
| Incar. > 1y | 0.28 | 0.31 | 0.36 | 0.39 |
| Incar. > 2y | 0.19 | 0.22 | 0.28 | 0.30 |
| Incar. > 4y | 0.12 | 0.14 | 0.24 | 0.24 |
| N | 41,574 | 30,440 | 3,805 | 1,653 |

Note: This Table reports summary statistics for final case dispositions in Bexar and Comal counties. “Full” refers to the initial analysis sample, while “Trimmed” refers to the final analysis sample after conducting the sample selection procedure described in Chapter 2.

Table 3.4: Validation Exercise

| | $\text{Var}(p_j)$ | $\text{SD}(p_j)$ | $\text{SE}(\text{Var}(p_j))$ | p -Value |
|-------------------|-------------------|------------------|------------------------------|------------|
| Male | 0.072 | 0.027 | 0.037 | 0.05 |
| White | 0.064 | 0.025 | 0.043 | 0.14 |
| Black | 0.040 | 0.020 | 0.031 | 0.19 |
| Latino | 0.073 | 0.027 | 0.062 | 0.24 |
| Age: 1st Quartile | 0.025 | 0.016 | 0.041 | 0.54 |
| Age: 2nd Quartile | -0.021 | . | 0.034 | 0.54 |
| Age: 3rd Quartile | 0.006 | 0.008 | 0.039 | 0.87 |
| Age: 4th Quartile | -0.029 | . | 0.034 | 0.40 |
| Any Prior Cases | 0.014 | 0.012 | 0.041 | 0.74 |
| Prior Felony | 0.032 | 0.018 | 0.018 | 0.08 |
| Multiple Charges | 0.069 | 0.026 | 0.052 | 0.18 |
| Prior Misdemeanor | 0.058 | 0.024 | 0.072 | 0.42 |

Test of Joint Significance: $p = 0.26$

Note: This Table reports estimates of the variance of attorney effects for pre-treatment variables that are assumed orthogonal to the assignment mechanism. Variances have been multiplied by 100 for readability.

Table 3.5: Estimated Attorney Effect Variances

| | Var(p_j) | SD(p_j) | SE(Var(p_j)) | p -Value |
|--------------------|--------------|-------------|------------------|------------|
| Plea Guilty | 0.116 | 0.034 | 0.059 | 0.05 |
| Incarceration | 0.314 | 0.056 | 0.129 | 0.02 |
| Probation | 0.029 | 0.017 | 0.027 | 0.28 |
| Deferred Adj. | 0.188 | 0.043 | 0.063 | 0.01 |
| Plea & Incarcerate | 0.289 | 0.054 | 0.125 | 0.02 |
| Incar. > 6m | 0.435 | 0.066 | 0.153 | 0.01 |
| Incar. > 1y | 0.580 | 0.076 | 0.181 | 0.01 |
| Incar. > 2y | 0.496 | 0.070 | 0.131 | 0.01 |
| Incar. > 4y | 0.370 | 0.061 | 0.111 | 0.01 |

Test of Joint Significance: $p = 0.001$

Note: This Table reports estimates of the variance of attorney effects for case outcome variables. Variances have been multiplied by 100 for readability.

Table 3.6: Correlation Between Estimated Effects \hat{p}_j and Attorney Characteristics

| | (1) | (2) | (3) | (4) |
|----------------------|---------------------|----------------------|----------------------|---------------------|
| | Plea | Incarceration | Probation | Deferred Adj. |
| Male | 0.0084 (0.0077) | 0.0270* (0.0123) | -0.0078 (0.0071) | -0.0140 (0.0099) |
| Latino | 0.0024 (0.0060) | 0.0188 (0.0102) | -0.0036 (0.0063) | -0.0118 (0.0087) |
| Black | 0.0101 (0.0130) | 0.0323 (0.0197) | -0.00462 (0.0154) | -0.0233 (0.0132) |
| Grad. Year/100 | -0.0483 (0.0296) | -0.1080* (0.0499) | -0.0205 (0.0335) | 0.0624 (0.0438) |
| US News Rank/100 | -0.0024 (0.0050) | -0.0100 (-0.008) | 0.00727 (0.0052) | 0.0041 (0.0069) |
| Solo Practitioner | 0.0199* (0.0094) | 0.0037 (0.0130) | -0.0101 (0.0078) | 0.0162 (0.0124) |
| Disciplinary History | 0.0078 (0.0087) | 0.0099 (0.0174) | 0.00297 (0.0096) | 0.0012 (0.0121) |
| N | 327 | 327 | 327 | 327 |
| Adj. R^2 | 0.0536 | 0.0674 | 0.0188 | 0.0284 |
| F-statistic | 2.408 | 3.123 | 0.873 | 1.801 |
| Joint p -value | 0.0205 | 0.00334 | 0.528 | 0.0864 |

Note: This Table reports regressions of estimated attorney effects on attorney characteristics. A constant is included, but not reported. Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Table 3.7: Comparison of Moments of Estimated $G(\cdot)$ s

| <i>Panel A: Unrestricted G</i> | | | |
|---|-------------|---------------|---------------|
| | Plea Guilty | Incarceration | Deferred Adj. |
| Std. Dev. | 0.052 | 0.083 | 0.070 |
| Skewness | -0.946 | 0.099 | 0.023 |
| Kurtosis | 4.454 | 3.762 | 3.723 |
| <i>Panel B: Exponential Family (Spline) G</i> | | | |
| | Plea Guilty | Incarceration | Deferred Adj. |
| Std. Dev. | 0.058 | 0.084 | 0.075 |
| Skewness | -3.360 | 0.041 | 1.276 |
| Kurtosis | 36.926 | 4.218 | 12.115 |
| <i>Panel C: Logit-Normal G</i> | | | |
| | Plea Guilty | Incarceration | Deferred Adj. |
| Std. Dev. | 0.053 | 0.083 | 0.072 |
| Skewness | -1.154 | -0.007 | 0.443 |
| Kurtosis | 5.020 | 2.803 | 3.136 |

Note: This Table reports moments of the estimated distribution of attorney effects for all three estimators (unrestricted, exponential family, and logit-normal).

Table 3.8: Estimates of Defendant Preferences over Attorney Characteristics

| | log(S_j/S_{ref}) | | | | | | | |
|------------------------|-----------------------------|------------------|---------------------|-------------------|------------------|------------------|------------------|--------------------|
| Female | 0.003 (0.53) | | | | | | | -0.517 (0.93) |
| Latino | | 0.896* (0.45) | | | | | | -0.199 (0.73) |
| Black | | | -1.905*** (0.38) | | | | | -2.509** (0.72) |
| Grad. Year | | | | -3.876* (1.90) | | | | 0.34 (3.07) |
| US News Ranking | | | | | -0.210 (0.36) | | | (0.48) (0.55) |
| EB Mean, Plea Guilty | | | | | | -0.801 (2.20) | | 1.796 (3.81) |
| EB Mean, Incarceration | | | | | | | -4.802 (4.33) | -6.422 (4.49) |
| R^2 | 0.00 | 0.05 | 0.07 | 0.07 | 0.01 | 0.00 | 0.03 | 0.14 |
| N | 60 | 56 | 56 | 60 | 60 | 36 | 36 | 32 |

Note: This Table reports estimates of a simple multinomial model of defendant preferences over attorney characteristics. A constant is included, but not reported. Robust standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Bibliography

- Abdulkadiroğlu, Atila, Parag A. Pathak, and Christopher R. Walters. 2018. “Free to Choose: Can School Choice Reduce Student Achievement?” *American Economic Journal: Applied Economics* 10 (1): 175–206. <https://doi.org/10.1257/app.20160634>. <https://www.aeaweb.org/articles?id=10.1257/app.20160634>.
- Abowd, John, Francis Kramarz, and David Margolis. 1999. “High Wage Workers and High Wage Firms.” *Econometrica* 67 (2): 251–333.
- Agan, Amanda, Matthew Freedman, and Emily Owens. 2021. “Is Your Lawyer a Lemon? Incentives and Selection in the Public Provision of Criminal Defense.” *The Review of Economics and Statistics* 103, no. 2 (May): 294–309. ISSN: 0034-6535. https://doi.org/10.1162/rest_a_00891. eprint: https://direct.mit.edu/rest/article-pdf/103/2/294/1915875/rest_a_00891.pdf. https://doi.org/10.1162/rest_a_00891.
- Aizer, Anna, and Joseph J. Doyle. 2015. “Juvenile Incarceration, Human Capital, and Future Crime: Evidence from Randomly Assigned Judges *.” *The Quarterly Journal of Economics* 130, no. 2 (February): 759–803. <https://doi.org/10.1093/qje/qjv003>. <https://doi.org/10.1093/qje/qjv003>.
- Allison, Paul D., and Nicholas A. Christakis. 1994. “Logit Models for Sets of Ranked Items.” *Sociological Methodology* 24:199–228. ISSN: 00811750, 14679531. <http://www.jstor.org/stable/270983>.
- Angrist, Joshua, Peter Hull, Parag Pathak, and Christopher Walters. 2020. *Simple and Credible Value-Added Estimation Using Centralized School Assignment*. Technical report. December. <https://doi.org/10.3386/w28241>. <https://doi.org/10.3386/w28241>.
- Arnold, David. 2021. *Mergers and Acquisitions, Local Labor Market Concentration, and Worker Outcomes*. Working Paper. UC San Diego.
- Avery, Christopher N., Mark E. Glickman, Caroline M. Hoxby, and Andrew Metrick. 2013. “A Revealed Preference Ranking of U.S. Colleges and Universities.” *The Quarterly Journal of Economics* 128 (1): 425–467. <https://ideas.repec.org/a/oup/qjecon/v128y2013i1p425-467.html>.

- Azar, José, Steven Berry, and Ioana Marinescu. 2019. *Estimating Labor Market Power* [in en]. SSRN Scholarly Paper ID 3456277. Rochester, NY: Social Science Research Network. Accessed October 14, 2019. <https://papers.ssrn.com/abstract=3456277>.
- Azar, José, Ioana Marinescu, Marshall Steinbaum, and Bledi Taska. 2020. “Concentration in US labor markets: Evidence from online vacancy data.” *Labour Economics* 66:101886.
- Backus, Matthew, Christopher Conlon, and Michael Sinkinson. 2021. *Common Ownership and Competition in the Ready-to-Eat Cereal Industry*. Working Paper, Working Paper Series 28350. National Bureau of Economic Research. <https://doi.org/10.3386/w28350>. <http://www.nber.org/papers/w28350>.
- Backus, Matthew, and Gregory Lewis. 2020. *Dynamic Demand Estimation in Auction Markets*. Working Paper, Working Paper Series 22375. National Bureau of Economic Research. <https://doi.org/10.3386/w22375>. <http://www.nber.org/papers/w22375>.
- Bagnoli, Mark, and Ted Bergstrom. 2005. “Log-Concave Probability and Its Applications.” *Economic Theory* 26 (2): 445–469. ISSN: 09382259, 14320479. <http://www.jstor.org/stable/25055959>.
- Bain, Joe S. 1951. “Relation of Profit Rate to Industry Concentration: American Manufacturing, 1936–1940.” *The Quarterly Journal of Economics* 65 (3): 293–324. <https://EconPapers.repec.org/RePEc:oup:qjecon:v:65:y:1951:i:3:p:293-324..>
- Barseghyan, Levon, Maura Coughlin, Francesca Molinari, and Joshua C. Teitelbaum. 2021. “Heterogeneous Choice Sets and Preferences.” *Econometrica* 89 (5): 2015–2048.
- Barth, Erling, Alex Bryson, James C. Davis, and Richard Freeman. 2016. “It’s Where You Work: Increases in the Dispersion of Earnings across Establishments and Individuals in the United States.” *Journal of Labor Economics* 34 (S2): S67–S97.
- Berger, David W., Kyle Herkenhoff, and Simon Mongey. 2017. *Labor Market Power*. Presentation. University of Chicago, Becker Friedman Institute for Economics.
- Berry, Steven T. 2021. *Market Structure and Competition, Redux*. Keynote. FTC Micro Conference.
- Berry, Steven T., and Philip A. Haile. 2014. “Identification in Differentiated Products Markets Using Market Level Data.” *Econometrica* 82 (5): 1749–1797. <https://doi.org/https://doi.org/10.3982/ECTA9027>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.3982/ECTA9027>. <https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA9027>.
- . 2020. *Nonparametric Identification of Differentiated Products Demand Using Micro Data*. Working Paper, Working Paper Series 27704. National Bureau of Economic Research. <https://doi.org/10.3386/w27704>. <http://www.nber.org/papers/w27704>.
- Bhaskar, V., Stephen J. Machin, and Gavin C. Reid. 1991. “Testing a Model of the Kinked Demand Curve.” *Journal of Industrial Economics* 39:241–254.

- Bhaskar, V., Alan Manning, and Ted To. 2002. "Oligopsony and Monopsonistic Competition in Labor Markets." *Journal of Economic Perspectives* 16 (2): 155–174. <https://doi.org/10.1257/0895330027300>. <https://www.aeaweb.org/articles?id=10.1257/0895330027300>.
- Bhaskar, V., and Ted To. 1999. "Minimum Wages for Ronald McDonald Monopsonies: A Theory of Monopsonistic Competition." *Economic Journal* 109 (455): 190–203. <https://EconPapers.repec.org/RePEc:ecj:conjl:v:109:y:1999:i:455:p:190-203>.
- . 2003. "Oligopsony and the Distribution of Wages." *European Economic Review* 47 (2): 371–399.
- Blinder, Alan. 1973. "Wage Discrimination: Reduced Form and Structural Estimates." *Journal of Human Resources* 8 (4): 436–455. <https://EconPapers.repec.org/RePEc:uwp:jhriss:v:8:y:1973:i:4:p:436-455>.
- Boal, William M., and Michael Ransom. 1997. "Monopsony in the Labor Market." *Journal of Economic Literature* 35 (1): 86–112. <https://EconPapers.repec.org/RePEc:aea:jeclit:v:35:y:1997:i:1:p:86-112>.
- Bolotnyy, Valentin, and Natalia Emanuel. 2022. "Why Do Women Earn Less Than Men? Evidence from Bus and Train Operators." *Journal of Labor Economics*, no. forthcoming.
- Bresnahan, Timothy F. 1989. "Chapter 17 Empirical studies of industries with market power," 2:1011–1057. *Handbook of Industrial Organization*. Elsevier. [https://doi.org/https://doi.org/10.1016/S1573-448X\(89\)02005-4](https://doi.org/https://doi.org/10.1016/S1573-448X(89)02005-4). <https://www.sciencedirect.com/science/article/pii/S1573448X89020054>.
- . 1987. "Competition and Collusion in the American Automobile Industry: The 1955 Price War." *Journal of Industrial Economics* 35 (4): 457–82. <https://EconPapers.repec.org/RePEc:bla:jindec:v:35:y:1987:i:4:p:457-82>.
- Burdett, Kenneth, and Dale T. Mortensen. 1998. "Wage Differentials, Employer Size, and Unemployment." *International Economic Review* 39 (2): 257–273. ISSN: 00206598, 14682354. <http://www.jstor.org/stable/2527292>.
- Caldwell, Sydnee, and Nikolaj Harmon. 2019. "Outside Options, Wages, and Bargaining: Evidence from Coworker Networks." *Working Paper*.
- Camerer, Colin, Linda Babcock, George Loewenstein, and Richard Thaler. 1997. "Labor Supply of New York City Cabdrivers: One Day at a Time." *The Quarterly Journal of Economics* 112 (2): 407–441.
- Card, David. 2022. *Design-Based Research in Empirical Microeconomics*. Technical report. January. <http://arks.princeton.edu/ark:/88435/dsp01ft848t765>.
- . 2012. "Model-Based or Design-Based? Competing Approaches in "Empirical Micro"." Woytinsky Lecture, University of Michigan.

- Card, David, Ana Rute Cardoso, Joerg Heining, and Patrick Kline. 2018. "Firms and Labor Market Inequality: Evidence and Some Theory." *Journal of Labor Economics* 36 (S1): 13–70.
- Card, David, Jörg Heining, and Patrick Kline. 2013. "Workplace Heterogeneity and the Rise of West German Wage Inequality." *The Quarterly Journal of Economics* 128, no. 3 (May): 967–1015. ISSN: 0033-5533. <https://doi.org/10.1093/qje/qjt006>.
- Carson, E. Ann, and William J. Sabol. 2012. *Prisoners in 2011*. Technical report NCJ 239808. http://www.prisonstudies.org/sites/prisonstudies.org/files/resources/downloads/wppl_10.pdf.
- Chamberlain, Edward, and Joan Robinson. 1933. *The American Economic Review* 23 (4): 683–685. ISSN: 00028282. <http://www.jstor.org/stable/1807525>.
- Chamberlain, Gary. 1987. "Asymptotic efficiency in estimation with conditional moment restrictions." *Journal of Econometrics* 34 (3): 305–334. <https://EconPapers.repec.org/RePEc:eee:econom:v:34:y:1987:i:3:p:305-334>.
- Chetty, Raj, John N. Friedman, and Jonah E. Rockoff. 2014. "Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-Added Estimates." *American Economic Review* 104 (9): 2593–2632. <https://doi.org/10.1257/aer.104.9.2593>. <https://www.aeaweb.org/articles?id=10.1257/aer.104.9.2593>.
- Crump, Richard K., V. Joseph Hotz, Guido W. Imbens, and Oscar A. Mitnik. 2009. "Dealing with limited overlap in estimation of average treatment effects." *Biometrika* 96, no. 1 (January): 187–199. ISSN: 0006-3444. <https://doi.org/10.1093/biomet/asn055>. eprint: <https://academic.oup.com/biomet/article-pdf/96/1/187/642537/asn055.pdf>. <https://doi.org/10.1093/biomet/asn055>.
- DellaVigna, Stefano, and Matthew Gentzkow. 2019. "Uniform Pricing in U.S. Retail Chains*." *The Quarterly Journal of Economics* 134, no. 4 (June): 2011–2084. ISSN: 0033-5533. <https://doi.org/10.1093/qje/qjz019>.
- Dempster, A. P., N. M. Laird, and D. B. Rubin. 1977. "Maximum Likelihood from Incomplete Data Via the EM Algorithm." *Journal of the Royal Statistical Society: Series B (Methodological)* 39, no. 1 (September): 1–22. <https://doi.org/10.1111/j.2517-6161.1977.tb01600.x>. <https://doi.org/10.1111/j.2517-6161.1977.tb01600.x>.
- Donald, Stephen, and Harry Paarsch. 1996. "Identification, Estimation, and Testing in Parametric Empirical Models of Auctions within the Independent Private Values Paradigm." *Econometric Theory* 12 (3): 517–567. https://EconPapers.repec.org/RePEc:cup:etheor:v:12:y:1996:i:03:p:517-567_00.
- . 1993. "Piecewise Pseudo-Maximum Likelihood Estimation in Empirical Models of Auctions." *International Economic Review* 34 (1): 121–148. ISSN: 00206598, 14682354. <http://www.jstor.org/stable/2526953>.

- Donald, Stephen, and Harry Paarsch. 2002. "Superconsistent estimation and inference in structural econometric models using extreme order statistics." *Journal of Econometrics* 109 (2): 305–340. <https://EconPapers.repec.org/RePEc:eee:econom:v:109:y:2002:i:2:p:305-340>.
- Duarte, Marco, Lorenzo Magnolfi, Mikkel Sølvsten, and Christopher Sullivan. 2021. *Testing Firm Conduct*. Working Paper. Department of Economics, University of Wisconsin.
- Dube, Arindrajit, Jeff Jacobs, Suresh Naidu, and Siddharth Suri. 2020. "Monopsony in Online Labor Markets." *American Economic Review: Insights* 2 (1): 33–46. <https://doi.org/10.1257/aeri.20180150>. <https://www.aeaweb.org/articles?id=10.1257/aeri.20180150>.
- Dube, Arindrajit, Alan Manning, and Suresh Naidu. 2020. *Monopsony and Employer Misoptimization Explain Why Wages Bunch at Round Numbers*. Working Paper, Working Paper Series 24991. National Bureau of Economic Research. <https://doi.org/10.3386/w24991>. <http://www.nber.org/papers/w24991>.
- Dunne, Timothy, Lucia Foster, John Haltiwanger, and Kenneth Troske. 2004. "Wage and Productivity Dispersion in United States Manufacturing: The Role of Computer Investment." *Journal of Labor Economics* 22 (2): 397–430. <https://EconPapers.repec.org/RePEc:ucp:jlabe:v:22:y:2004:i:2:p:397-430>.
- Efron, Bradley. 2016. "Empirical Bayes deconvolution estimates." *Biometrika* 103, no. 1 (February): 1–20. ISSN: 0006-3444. <https://doi.org/10.1093/biomet/asv068>. eprint: <https://academic.oup.com/biomet/article-pdf/103/1/1/24331584/asv068.pdf>. <https://doi.org/10.1093/biomet/asv068>.
- Faggio, Giulia, Kjell G. Salvanes, and John van Reenen. 2010. "The evolution of inequality in productivity and wages: panel data evidence." *Industrial and Corporate Change* 19 (6): 1919–1951. <https://EconPapers.repec.org/RePEc:oup:indcch:v:19:y:2010:i:6:p:1919-1951>.
- Farber, Henry S. 2015. "Why you Can't Find a Taxi in the Rain and Other Labor Supply Lessons from Cab Drivers*." *The Quarterly Journal of Economics* 130, no. 4 (July): 1975–2026. ISSN: 0033-5533. <https://doi.org/10.1093/qje/qjv026>. eprint: <https://academic.oup.com/qje/article-pdf/130/4/1975/30637404/qjv026.pdf>. <https://doi.org/10.1093/qje/qjv026>.
- Flinn, Christopher, and Joseph Mullins. 2021. *Firms' Choices of Wage-Setting Protocols*. Working Paper.
- Freeman, Richard B. 1976. "A Cobweb Model of the Supply and Starting Salary of New Engineers." *Industrial and Labor Relations Review* 29 (2): 236–248. ISSN: 00197939, 2162271X. <http://www.jstor.org/stable/2522143>.

- Gandhi, Amit, and Aviv Nevo. 2021. “Empirical Models of Demand and Supply in Differentiated Products Industries.” In *Handbook of Industrial Organization*, 1st ed., edited by Ali Hortascu, Kate Ho, and Alessandro Lizzeri, vol. 4. Elsevier. <https://EconPapers.repec.org/RePEc:eee:indchp:2-17>.
- Gentry, Matthew, and Tong Li. 2014. “Identification in auctions with selective entry.” *Econometrica* 82 (1): 315–344. ISSN: 00129682, 14680262. <http://www.jstor.org/stable/24029177>.
- Gibbons, Robert, Lawrence Katz, Thomas Lemieux, and Daniel Parent. 2005. “Comparative Advantage, Learning, and Sectoral Wage Determination.” *Journal of Labor Economics* 23 (4): 681–724. <https://EconPapers.repec.org/RePEc:ucp:jlabecon:v:23:y:2005:i:4:p:681-724>.
- Goldin, Claudia, and Lawrence F. Katz. 2016. “A Most Egalitarian Profession: Pharmacy and the Evolution of a Family-Friendly Occupation.” *Journal of Labor Economics* 34 (3): 705–746. <https://doi.org/10.1086/685505>.
- Gourieroux, Christian, Alain Monfort, Eric Renault, and Alain Trognon. 1987. “Generalised residuals.” *Journal of Econometrics* 34 (1-2): 5–32. <https://EconPapers.repec.org/RePEc:eee:econom:v:34:y:1987:i:1-2:p:5-32>.
- Graham, Bryan S. 2020. *Sparse network asymptotics for logistic regression*. Working Paper. UC Berkeley.
- Guerre, Emmanuel, Isabelle Perrigne, and Quang H. Vuong. 2000. “Optimal Nonparametric Estimation of First-Price Auctions.” *Econometrica* 68 (3): 525–574. ISSN: 00129682, 14680262. <http://www.jstor.org/stable/2999600>.
- Guryan, Jonathan, and Kerwin Kofi Charles. 2013. “Taste-based or Statistical Discrimination: The Economics of Discrimination Returns to its Roots.” *The Economic Journal* 123 (572): F417–F432.
- Hall, Robert, and Alan B. Krueger. 2012. “Evidence on the Incidence of Wage Posting, Wage Bargaining, and On-the-Job Search.” *American Economic Journal: Macroeconomics* 4 (4): 56–67.
- Hall, Robert, and Andreas Mueller. 2018. “Wage Dispersion and Search Behavior: The Importance of Nonwage Job Values.” *Journal of Political Economy* 126 (4): 1594–1637. <https://EconPapers.repec.org/RePEc:ucp:jpolecon:doi:10.1086/697739>.
- Hamermesh, Daniel. 1999. “Changing Inequality in Markets for Workplace Amenities.” *The Quarterly Journal of Economics* 114 (4): 1085–1123. <https://EconPapers.repec.org/RePEc:oup:qjecon:v:114:y:1999:i:4:p:1085-1123>.
- Han, Ruijian, Yiming Xu, and Kani Chen. 2020. “A General Pairwise Comparison Model for Extremely Sparse Networks.” *ArXiv* abs/2002.08853.

- Hangartner, Dominik, Daniel Kopp, and Michael Siegenthaler. 2021. "Monitoring hiring discrimination through online recruitment platforms." *Nature* 589:572–576.
- Hirano, Keisuke, Guido W. Imbens, and Geert Ridder. 2003. "Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score." *Econometrica* 71, no. 4 (July): 1161–1189. <https://doi.org/10.1111/1468-0262.00442>. <https://doi.org/10.1111/1468-0262.00442>.
- Horton, John J., Ramesh Johari, and Philipp Kircher. 2021. *Cheap Talk Messages for Market Design: Theory and Evidence from a Labor Market with Directed Search*. Working Paper, Working Paper Series 29445. National Bureau of Economic Research. <https://doi.org/10.3386/w29445>. <http://www.nber.org/papers/w29445>.
- Imbens, Guido W. 2000. "The role of the propensity score in estimating dose-response functions." *Biometrika* 87, no. 3 (September): 706–710. <https://doi.org/10.1093/biomet/87.3.706>. <https://doi.org/10.1093/biomet/87.3.706>.
- Jarosch, Gregor, Jan Sebastian Nimczik, and Isaac Sorkin. 2021. *Granular Search, Market Structure, and Wages*. Working Paper, Working Paper Series 26239. National Bureau of Economic Research. <https://doi.org/10.3386/w26239>. <http://www.nber.org/papers/w26239>.
- Jarosch, Gregor, and Laura Pilossoph. 2018. "Statistical Discrimination and Duration Dependence in the Job Finding Rate." *The Review of Economic Studies* 86, no. 4 (September): 1631–1665. ISSN: 0034-6527. <https://doi.org/10.1093/restud/rdy055>. eprint: <https://academic.oup.com/restud/article-pdf/86/4/1631/28883078/rdy055.pdf>. <https://doi.org/10.1093/restud/rdy055>.
- Justice Management Institute. 2017. *The Power of Choice: The Implications of a System Where Indigent Defendants Choose Their Own Counsel*. Technical report.
- Kline, Patrick, and Christopher Walters. 2021. "Reasonable Doubt: Experimental Detection of Job-Level Employment Discrimination." *Econometrica* 89 (2): 765–792. <https://doi.org/10.3982/ecta17489>. <https://doi.org/10.3982/ecta17489>.
- Koenker, Roger, and Ivan Mizera. 2014. "Convex optimization, shape constraints, compound decisions, and empirical Bayes rules." *Journal of the American Statistical Association* 109 (506): 674–685.
- Kroft, Kory, Fabian Lange, and Matthew J. Notowidigdo. 2013. "Duration Dependence and Labor Market Conditions: Evidence from a Field Experiment*." *The Quarterly Journal of Economics* 128, no. 3 (June): 1123–1167. ISSN: 0033-5533. <https://doi.org/10.1093/qje/qjt015>. eprint: <https://academic.oup.com/qje/article-pdf/128/3/1123/30631486/qjt015.pdf>. <https://doi.org/10.1093/qje/qjt015>.
- Lachowska, Marta, Alexandre Mas, Raffaele Saggio, and Stephen Woodbury. 2021. *Do Workers Bargain Over Wages? A Test Using Dual Jobholders*. Working Paper, Working Paper Series 28409. National Bureau of Economic Research.

- Lagos, Lorenzo. 2021. *Labor Market Institutions and the Composition of Firm Compensation: Evidence from Brazilian Collective Bargaining*. Working Paper. https://www.dropbox.com/s/g1ux5gcl3dh63ep/LMICFC_paper.pdf?dl=0.
- Lamadon, Thibaut, Magne Mogstad, and Bradley Setzler. 2022. “Imperfect Competition, Compensating Differentials and Rent Sharing in the U.S. Labor Market.” *American Economic Review*, no. forthcoming.
- Lee, Daniel, and H. Sebastian Seung. 2000. “Algorithms for Non-negative Matrix Factorization.” In *Advances in Neural Information Processing Systems*, edited by T. Leen, T. Dietterich, and V. Tresp, vol. 13. MIT Press. <https://proceedings.neurips.cc/paper/2000/file/f9d1152547c0bde01830b7e8bd60024c-Paper.pdf>.
- Lindenlaub, Ilse, and Fabien Postel-Vinay. 2021. *The Worker-Job Surplus*. Working Paper, Working Paper Series 28402. National Bureau of Economic Research. <https://doi.org/10.3386/w28402>. <http://www.nber.org/papers/w28402>.
- Lipsky, David B., and Henry S. Farber. 1976. “The Composition of Strike Activity in the Construction Industry.” *ILR Review* 29 (3): 388–404. <https://doi.org/10.1177/001979397602900305>.
- Liu, Ao, Zhibing Zhao, Chao Liao, Pinyan Lu, and Lirong Xia. 2019. “Learning Plackett-Luce Mixtures from Partial Preferences.” *Proceedings of the AAAI Conference on Artificial Intelligence* 33 (01): 4328–4335. <https://doi.org/10.1609/aaai.v33i01.33014328>. <https://ojs.aaai.org/index.php/AAAI/article/view/4342>.
- Lofstrom, Magnus, and Steven Raphael. 2016. “Crime, the Criminal Justice System, and Socioeconomic Inequality.” *Journal of Economic Perspectives* 30 (2): 103–26. <https://doi.org/10.1257/jep.30.2.103>. <https://www.aeaweb.org/articles?id=10.1257/jep.30.2.103>.
- Lu, Haihao, Robert M. Freund, and Yurii Nesterov. 2018. “Relatively Smooth Convex Optimization by First-Order Methods, and Applications.” *SIAM Journal on Optimization* 28 (1): 333–354. <https://doi.org/10.1137/16M1099546>. eprint: <https://doi.org/10.1137/16M1099546>. <https://doi.org/10.1137/16M1099546>.
- Luce, R. Duncan. 1959. *Individual Choice Behavior A Theoretical Analysis*. Riley Wiley.
- Macaluso, Claudia, Brad Hershbein, and Chen Yeh. 2021. *Monopsony in the US Labor Market*. Working Paper. Society for Economic Dynamics.
- Manning, Alan. 2011. “Imperfect Competition in the Labor Market.” Chap. 11, 1st ed., edited by O. Ashenfelter and D. Card, vol. 4B, 973–1041. Elsevier. <https://EconPapers.repec.org/RePEc:eee:labchp:5-11>.
- . 2005. “Monopsony and Labour Demand.” *Brussels Economic Review* 48 (1-2): 95–112. <https://ideas.repec.org/a/bxr/bxrceb/y2005v48i1-2p95-112.html>.

- Mas, Alexandre, and Amanda Pallais. 2017a. "Valuing Alternative Work Arrangements." *American Economic Review* 107 (12): 3722–3759.
- . 2017b. "Valuing Alternative Work Arrangements." *American Economic Review* 107 (12): 3722–59. <https://doi.org/10.1257/aer.20161500>. <https://www.aeaweb.org/articles?id=10.1257/aer.20161500>.
- McRae, Andrew D, and Mark A Davenport. 2020. "Low-rank matrix completion and denoising under Poisson noise." *Information and Inference: A Journal of the IMA* 10, no. 2 (August): 697–720. <https://doi.org/10.1093/imaiai/iaaa020>. <https://doi.org/10.1093/imaiai/iaaa020>.
- Mosteller, Robert P. 2011. "Failures of the American Adversarial System to Protect the Innocent and Conceptual Advantages in the Inquisitorial Design for Investigative Fairness." *North Carolina Journal of International Law and Commercial Regulation* 36:319–364.
- Mueller-Smith, Michael. 2015. *The Criminal and Labor Market Impacts of Incarceration*. Technical report. <https://sites.lsa.umich.edu/mgms/wp-content/uploads/sites/283/2015/09/incar.pdf>.
- Nevo, Aviv. 2001. "Measuring Market Power in the Ready-to-Eat Cereal Industry." *Econometrica* 69 (2): 307–42. <https://EconPapers.repec.org/RePEc:ecm:emetrp:v:69:y:2001:i:2:p:307-42>.
- Neyman, Jerzy, and Elizabeth L. Scott. 1948. "Consistent Estimates Based on Partially Consistent Observations." *Econometrica* 16 (1): 1–32. ISSN: 00129682, 14680262. <http://www.jstor.org/stable/1914288>.
- Oaxaca, Ronald. 1973. "Male-Female Wage Differentials in Urban Labor Markets." *International Economic Review* 14 (3): 693–709. <https://EconPapers.repec.org/RePEc:ier:iecrev:v:14:y:1973:i:3:p:693-709>.
- Pager, Devah. 2003. "The Mark of a Criminal Record." *American Journal of Sociology* 108, no. 5 (March): 937–975. <https://doi.org/10.1086/374403>. <https://doi.org/10.1086/374403>.
- Pesaran, M.H. 1990. "Non-nested Hypotheses." *Eatwell J., Milgate M., Newman P. (eds) Econometrics. The New Palgrave. Palgrave Macmillan, London*.
- Pierce, Brooks. 2001. "Compensation Inequality*." *The Quarterly Journal of Economics* 116, no. 4 (November): 1493–1525. ISSN: 0033-5533. <https://doi.org/10.1162/003355301753265633>. eprint: <https://academic.oup.com/qje/article-pdf/116/4/1493/5372911/116-4-1493.pdf>. <https://doi.org/10.1162/003355301753265633>.
- Plackett, R. L. 1975. "The Analysis of Permutations." *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 24 (2): 193–202. ISSN: 00359254, 14679876. <http://www.jstor.org/stable/2346567>.

- Postel-Vinay, Fabien, and Jean-Marc Robin. 2002. "Equilibrium Wage Dispersion with Worker and Employer Heterogeneity." *Econometrica* 70 (6): 2295–2350. ISSN: 00129682, 14680262. <http://www.jstor.org/stable/3081988>.
- . 2004. "To Match or Not to Match? Optimal Wage Policy With Endogenous Worker Search Intensity." *Review of Economic Dynamics* 7 (2): 297–330. [https://doi.org/10.1016/S1094-2025\(03\)000](https://doi.org/10.1016/S1094-2025(03)000). <https://ideas.repec.org/a/red/issued/v7y2004i2p297-330.html>.
- Raphael, Steven, and Jeffrey Smith. 2011. "11. Improving Employment Prospects for Former Prison Inmates: Challenges and Policy." In *Controlling Crime: Strategies and Tradeoffs*, edited by Philip J. Cook, Jens Ludwig, and Justin McCrary, 521–572. University of Chicago Press. <https://doi.org/doi:10.7208/9780226115139-014>. <https://doi.org/10.7208/9780226115139-014>.
- Rehavi, M. Marit, and Sonja B. Starr. 2014. "Racial Disparity in Federal Criminal Sentences." *Journal of Political Economy* 122, no. 6 (December): 1320–1354. <https://doi.org/10.1086/677255>. <https://doi.org/10.1086/677255>.
- Rivers, Douglas, and Quang Vuong. 2002. "Model selection tests for nonlinear dynamic models." *The Econometrics Journal* 5 (1): 1–39. <https://doi.org/https://doi.org/10.1111/1368-423X.t01-1-00071>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/1368-423X.t01-1-00071>. <https://onlinelibrary.wiley.com/doi/abs/10.1111/1368-423X.t01-1-00071>.
- Robinson, Joan. 1933. *The Economics of Imperfect Competition*. Palgrave Macmillan.
- Rosen, Sherwin. 1986. "The theory of equalizing differences," 1:641–692. *Handbook of Labor Economics*. Elsevier. [https://doi.org/https://doi.org/10.1016/S1573-4463\(86\)01015-5](https://doi.org/https://doi.org/10.1016/S1573-4463(86)01015-5). <https://www.sciencedirect.com/science/article/pii/S1573446386010155>.
- Rosenbaum, Paul R., and Donald B. Rubin. 1983. "The central role of the propensity score in observational studies for causal effects." *Biometrika* 70 (1): 41–55. <https://doi.org/10.1093/biomet/70.1.41>. <https://doi.org/10.1093/biomet/70.1.41>.
- Rothstein, Jesse. 2010. "Teacher Quality in Educational Production: Tracking, Decay, and Student Achievement*." *The Quarterly Journal of Economics* 125, no. 1 (February): 175–214. ISSN: 0033-5533. <https://doi.org/10.1162/qjec.2010.125.1.175>. eprint: <https://academic.oup.com/qje/article-pdf/125/1/175/5314918/125-1-175.pdf>. <https://doi.org/10.1162/qjec.2010.125.1.175>.
- Roussille, Nina. 2021. *The Central Role of the Ask Gap in Gender Pay Inequality*. Working Paper. University of California, Berkeley.
- Rubin, Donald B., Elizabeth A. Stuart, and Elaine L. Zanutto. 2004. "A Potential Outcomes View of Value-Added Assessment in Education." *Journal of Educational and Behavioral Statistics* 29, no. 1 (March): 103–116. <https://doi.org/10.3102/10769986029001103>. <https://doi.org/10.3102/10769986029001103>.

- Schmalensee, Richard. 1989. "Chapter 16 Inter-industry studies of structure and performance," 2:951–1009. *Handbook of Industrial Organization*. Elsevier. [https://doi.org/https://doi.org/10.1016/S1573-448X\(89\)02004-2](https://doi.org/https://doi.org/10.1016/S1573-448X(89)02004-2).
- Schubert, Gregor, Anna Stansbury, and Bledi Taska. 2021. *Employer Concentration and Outside Options*. Working Paper. University of Chicago, Becker Friedman Institute for Economics.
- Shem-Tov, Yotam. 2020. "Make-or-Buy? The Provision of Indigent Defense Services in the U.S." *The Review of Economics and Statistics* (September): 1–27. ISSN: 0034-6535. https://doi.org/10.1162/rest_a_00976. eprint: https://direct.mit.edu/rest/article-pdf/doi/10.1162/rest_a_00976/1891314/rest_a_00976.pdf. https://doi.org/10.1162/rest_a_00976.
- Simons, Gordon, and Yi-Ching Yao. 1999. "Asymptotics when the number of parameters tends to infinity in the Bradley-Terry model for paired comparisons." *The Annals of Statistics* 27 (3): 1041–1060. <https://doi.org/10.1214/aos/1018031267>. <https://doi.org/10.1214/aos/1018031267>.
- Sorkin, Isaac. 2018. "Ranking Firms Using Revealed Preference." *The Quarterly Journal of Economics* 133 (3): 1331–1393.
- . 2017. "The Role of Firms in Gender Earnings Inequality: Evidence from the United States." *American Economic Review* 107 (5): 384–87. <https://doi.org/10.1257/aer.p20171015>. <https://www.aeaweb.org/articles?id=10.1257/aer.p20171015>.
- Soufiani, Hossein Azari, Hansheng Diao, Zhenyu Lai, and David C. Parkes. 2013. "Generalized Random Utility Models with Multiple Types" [in English (US)]. *Advances in Neural Information Processing Systems*, ISSN: 1049-5258.
- Staiger, Douglas O., Joanne Spetz, and Ciaran S. Phibbs. 2010. "Is There Monopsony in the Labor Market? Evidence from a Natural Experiment." *Journal of Labor Economics* 28 (2): 211–236. <https://doi.org/10.1086/652734>.
- Sterba, Sonya K. 2009. "Alternative Model-Based and Design-Based Frameworks for Inference From Samples to Populations: From Polarization to Integration." *Multivariate Behavioral Research* 44, no. 6 (November): 711–740. <https://doi.org/10.1080/00273170903333574>. <https://doi.org/10.1080/00273170903333574>.
- Sweezy, Paul M. 1939. "Demand Under Conditions of Oligopoly." *Journal of Political Economy* 47 (4): 568–573. ISSN: 00223808, 1537534X. <http://www.jstor.org/stable/1824594>.
- Taber, Christopher, and Rune Vejlin. 2020. "Estimation of a Roy/Search/Compensating Differential Model of the Labor Market." *Econometrica* 88 (3): 1031–1069. <https://doi.org/https://doi.org/10.3982/ECTA14441>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.3982/ECTA14441>. <https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA14441>.

- Texas Task Force on Indigent Defense. 2010. *Review of Bexar County's Indigent Defense Systems*. Technical report.
- Vuong, Quang H. 1989. "Likelihood Ratio Tests for Model Selection and Non-Nested Hypotheses." *Econometrica* 57 (2): 307–333. ISSN: 00129682, 14680262. <http://www.jstor.org/stable/1912557>.
- Wamsley, Roy. 2013. *World Prison Population List (Tenth Edition)*. Technical report. http://www.prisonstudies.org/sites/prisonstudies.org/files/resources/downloads/wppl_10.pdf.
- Wang, Qing, and Bruce G. Lindsay. 2014. "Variance estimation of a general u-statistic with application to cross-validation." *Statistica Sinica*, <https://doi.org/10.5705/ss.2012.215>. <https://doi.org/10.5705/ss.2012.215>.
- Wiswall, Matthew, and Basit Zafar. 2018. "Preference for the Workplace, Investment in Human Capital, and Gender." *The Quarterly Journal of Economics* 133 (1): 457–507.
- Yan, Ting, Yaning Yang, and Jinfeng Xu. 2012. "Sparse Paired Comparisons in the Bradley-Terry Model." *Statistica Sinica* 22 (3): 1305–1318. ISSN: 10170405, 19968507. <http://www.jstor.org/stable/24309985>.

Appendix A

Appendix to Chapter 1

A.1 Additional Figures

Figure A.1: Mandatory features of a candidate profile, at the time of the study

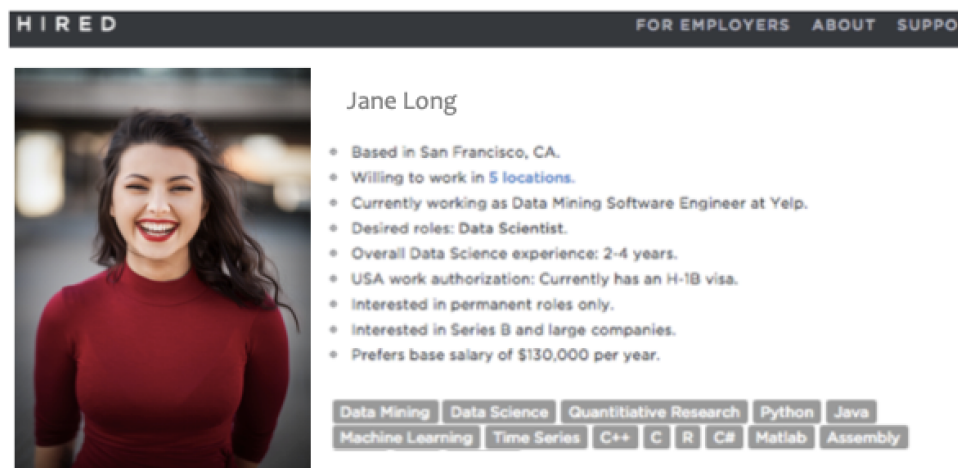


Figure A.2: Typical interview request message sent by a company to a candidate, at the time of the study

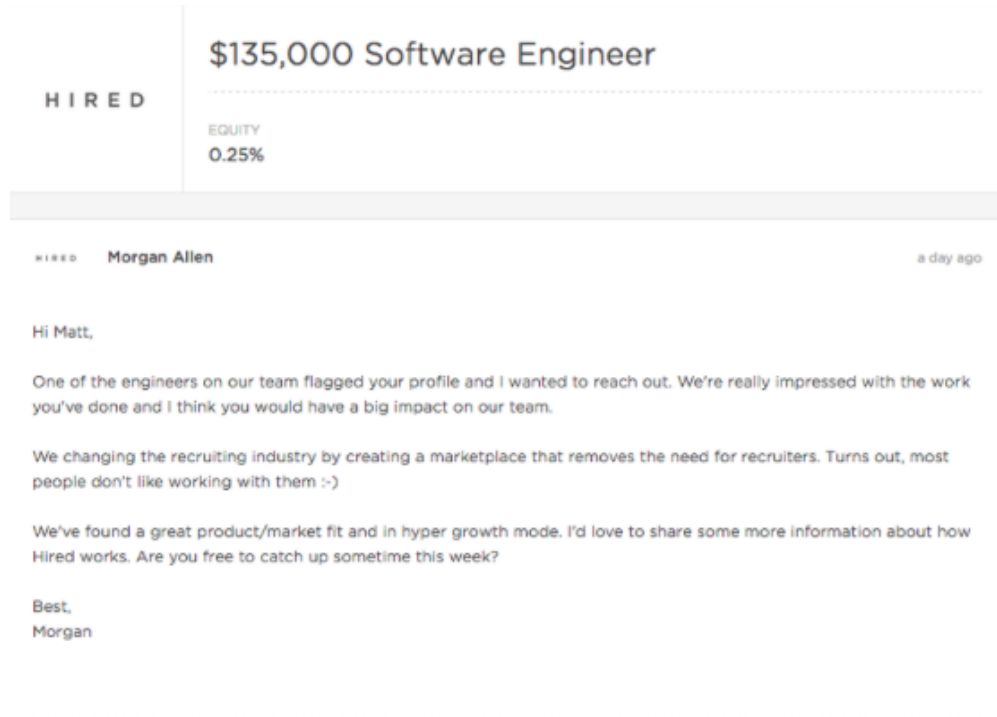
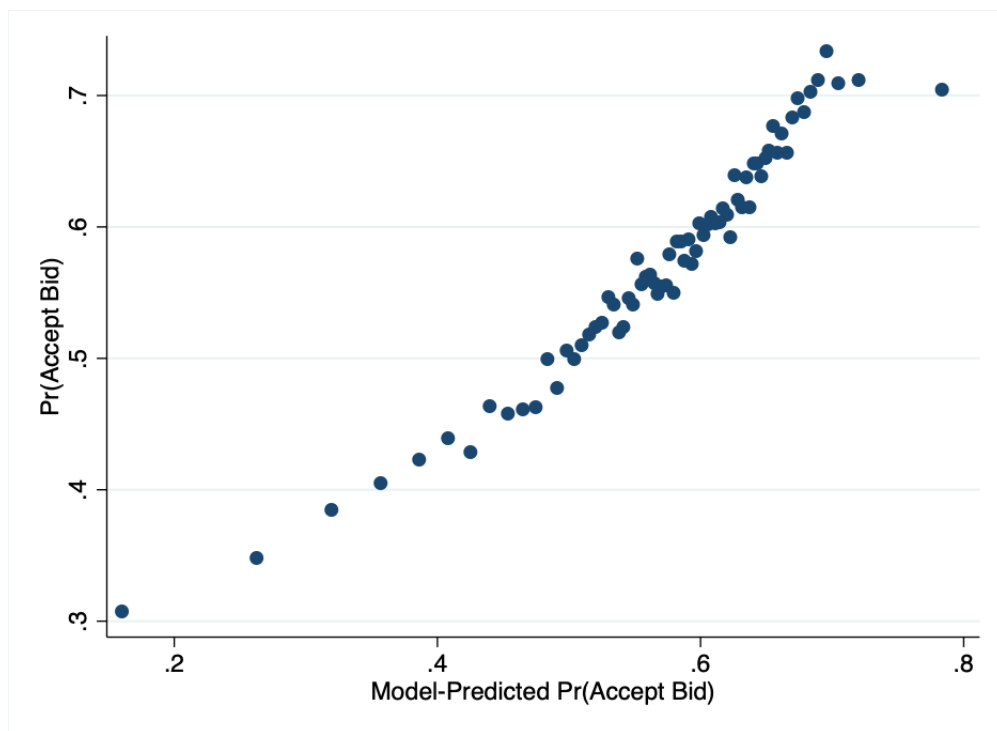
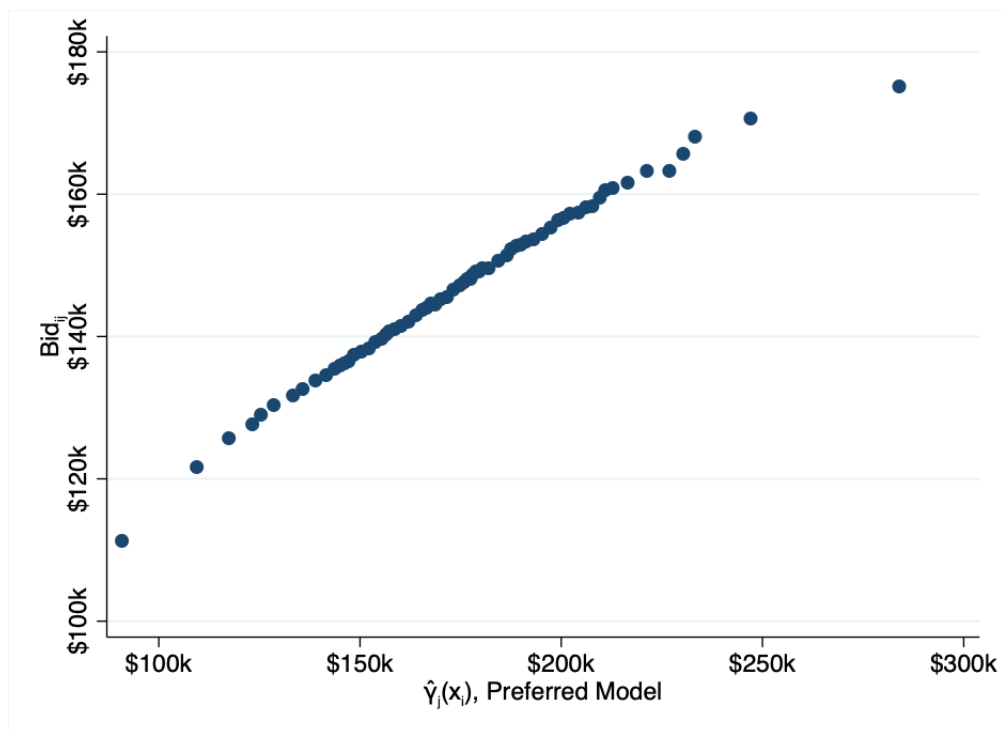


Figure A.3: Model Fit: Labor Supply

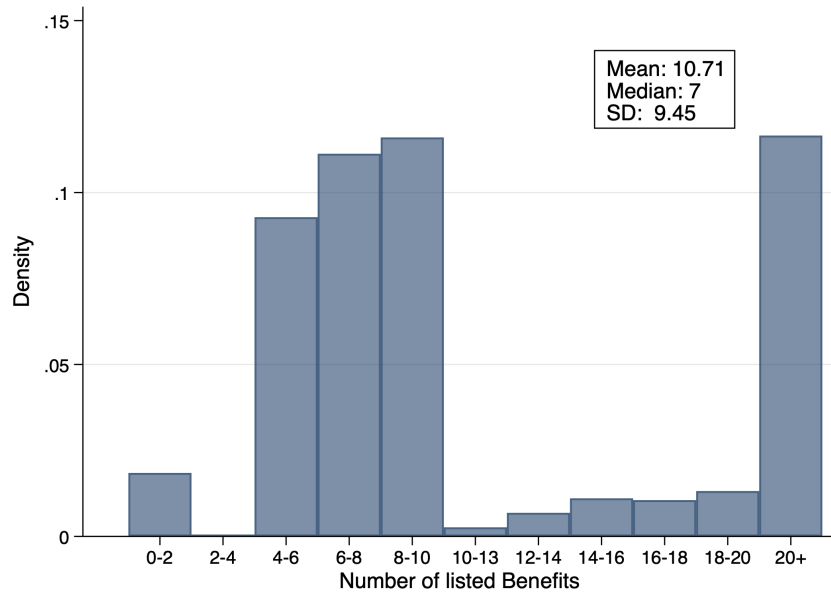
Note: This Figure plots the relationship between the empirical acceptance probability of a bid and the model-implied probabilities that the bid will be accepted.

Figure A.4: Relationship between bids and systematic component of valuations, $\gamma_j(x_i)$ 

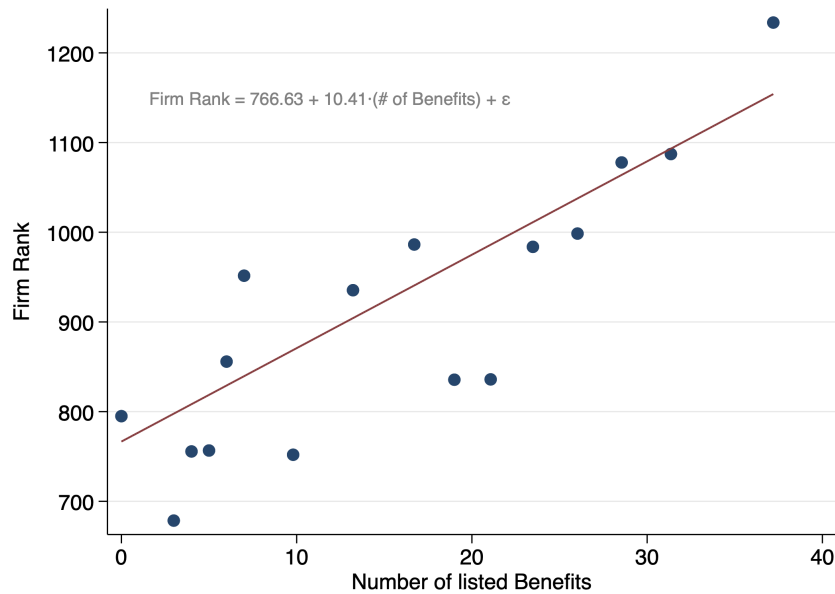
Note: This Figure plots the relationship between observed bids and the systematic component of valuations $\exp(z_j' \Gamma x_i)$ in the preferred model, controlling for the asked salary. Unconditionally, the slope of the relationship between bids and the observed component of valuations is 0.83.

Figure A.5: Summary Statistics of Benefits listed by Firms

(a) Distribution of Number of listed Benefits



(b) Share of listed Benefits



Note: This Figure displays the distribution of benefits listed by firms in the subset of ranked firms. Panel (a) plots the density of the number of listed benefits per firm. The bar “20+” includes numbers of listed benefits greater than 20 up to a maximum of 53. The mean number of benefits is 10.71 (SD 9.45), while the median lies at 7. Panel (b) illustrates the relationship between firm ranking and the number of listed benefits. On average an additional benefit increases the firm’s ranking by 10.41.

A.2 Additional Tables

Table A.1: Comparison of data sources

| Observe... | Admin | Surveys | Experiments | This Paper |
|---------------------------------|---------|---------|-------------|-------------------|
| full choice sets? | No | Depends | Yes | Yes |
| multiple choices per worker? | No | Depends | Depends | Yes |
| info on indiv. characteristics? | Depends | Yes | Yes | Yes |
| high stakes choices? | Yes | Depends | No | Yes |
| exogenous choice sets? | No | No | Yes | No |

Table A.2: Match productivity estimates: $\gamma_j(x_i) = z_j' \Gamma x_i$

| | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) |
|------------------|------------------------|-----------------------|---------------------------|------------------------|------------------------|------------------------|------------------------|---------------------|
| | Soft-Eng | Experience | (Experience) ² | Unemployed | Ivy Plus | CS Degree | FAANG | Previous Jobs |
| Constant | 0.0326*** (0.0029) | 0.0005 (0.0006) | 0.00002 (0.00002) | 0.0009 (0.0010) | -0.0060* (0.0023) | 0.0042 (0.0022) | -0.0012 (0.0028) | -0.0003 (0.0005) |
| 16-50 Employees | -0.0046 (0.0031) | 0.0007 (0.0006) | -0.0000111 (0.00002) | 0.0003 (0.0010) | -0.0035 (0.0025) | -0.0028 (0.0024) | -0.0017 (0.0030) | -0.0008 (0.0005) |
| 51-500 Employees | -0.0144*** (0.0029) | 0.0020*** (0.0006) | -0.00005** (0.0000176) | 0.0002 (0.0010) | 0.0049* (0.0024) | -0.0031 (0.0022) | -0.0020 (0.0028) | -0.0011 (0.0005) |
| 501+ Employees | -0.0167*** (0.0030) | 0.0016** (0.0006) | -0.00005** (0.00002) | -0.0006 (0.0010) | 0.0073** (0.0025) | -0.0020 (0.0023) | 0.0001 (0.0029) | -0.0002 (0.0005) |
| Finance | 0.0084*** (0.0017) | -0.0006 (0.0004) | 0.00001 (0.00001) | 0.0006 (0.0006) | -0.0077*** (0.0015) | -0.0052*** (0.0013) | -0.0047** (0.0017) | 0.0003 (0.0003) |
| Tech | 0.0068*** (0.0014) | -0.0005 (0.0003) | 0.00001 (0.00001) | -0.0008 (0.0005) | -0.0010 (0.0013) | 0.0016 (0.0011) | -0.0022 (0.0014) | -0.0003 (0.0003) |
| Health | 0.0074*** (0.0022) | -0.0004 (0.0005) | 0.00001 (0.00001) | -0.0007 (0.0008) | -0.0027 (0.0021) | -0.0049** (0.0018) | -0.0031 (0.0024) | 0.0004 (0.0004) |
| | (9) | (10) | (11) | (12) | (13) | (14) | (15) | |
| | Fulltime | Sponsorship | Remote | Java | Python | SQL | C | |
| Constant | -0.0035 (0.0022) | -0.0019 (0.0028) | 0.0029 (0.0020) | -0.0002 (0.0021) | 0.0009 (0.0020) | -0.0030 (0.0023) | 0.0077** (0.0026) | |
| 16-50 Employees | 0.0011 (0.0024) | 0.0150** (0.0030) | 0.0032 (0.0022) | -0.0006 (0.0023) | -0.0004 (0.0022) | 0.0065* (0.0025) | -0.0136*** (0.0029) | |
| 51-500 Employees | 0.0039 (0.0022) | 0.0058* (0.0028) | -0.0012 (0.0021) | 0.0042 (0.0022) | -0.0018 (0.0020) | 0.0039 (0.0023) | -0.0076** (0.0027) | |
| 501+ Employees | 0.0034 (0.0023) | 0.0057* (0.0028) | -0.0020 (0.0021) | 0.0064** (0.0022) | -0.0029 (0.0021) | 0.0032 (0.0024) | -0.0087** (0.0027) | |
| Finance | -0.0023 (0.0014) | 0.0025 (0.0016) | 0.0003 (0.0013) | -0.0063*** (0.0013) | 0.0011 (0.0013) | 0.0008 (0.0014) | 0.0012 (0.0016) | |
| Tech | -0.0028* (0.0012) | 0.0004 (0.0013) | 0.0001 (0.0011) | -0.0058*** (0.0011) | 0.0024* (0.0011) | 0.0024 (0.0012) | 0.0021 (0.0013) | |
| Health | 0.0025 (0.0019) | -0.0031 (0.0021) | 0.0027 (0.0017) | 0.0004 (0.0018) | -0.0032 (0.0017) | -0.0003 (0.0019) | 0.0013 (0.0023) | |

Note: This Table presents the remaining set of coefficients corresponding to Table 1.7. The omitted category for the number of employees is “1-15 Employees”. Every cell reports the coefficient on the interaction of the variables specified in the corresponding row and column. Column variables are candidate characteristics (x_i), and row variables are firm characteristics (z_j). Standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

A.3 Illustration of conceptual framework

The following simple model, adapted from Bhaskar, Manning, and To (2002), can be used to illustrate the logic of our conduct testing procedure. In particular, the model illustrates the role of worker preference heterogeneity, the implications of conduct assumptions, and the basic logic of our estimation and testing framework. The basic message is that combinations of assumptions on competition and wage-setting flexibility deliver different wage equations, which can then be used to infer conduct. Our simple model consists of:

- Firms $j = -1, +1$, which are located on either end of a mile-long road;

$$\text{MRPL}_j = \text{ARPL}_j = \gamma_j.$$

- Workers distributed along road with location ξ , which is private information:

$$\xi \sim \text{Unif}[0, 1].$$

- Workers live on either side of the road, given by the variable v , which is public information:

$$v \perp\!\!\!\perp \xi, \quad v = \{-1, +1\} \text{ w.p. } 1/2.$$

- Firms post wages (which may vary by v), and worker utilities are given by:

$$u_{-1}^v(\xi) = w_{-1}^v - \beta(\xi + \alpha v); \quad u_{+1}^v(\xi) = w_{+1}^v - \beta(1 - (\xi + \alpha v)).$$

Under these assumptions, type- v 's labor supply to firm j is:

$$S_j^v(w_j^v; w_{-j}^v) = \frac{1}{2} + \frac{w_j^v - w_{-j}^v}{2\beta} + \alpha v j.$$

Labor demand is determined by profit maximization:

$$\pi_j(\mathbf{w}) = \frac{1}{2} \sum_{v=-1}^{+1} (\gamma_j - w^v) \times S_j^v(w^v; \hat{w}_{-j}^v),$$

where the random variable \hat{w}_{-j}^v encodes j 's knowledge of the competitive environment. Wages are determined by firms' first-order conditions and a market clearing constraint:

$$w_j^v = \frac{1}{2}(\hat{w}_{-j}^v + \gamma_j - \beta) - \alpha\beta v j, \quad S_j^v(w_j^v; \hat{w}_{-j}^v) + S_{-j}^v(w_{-j}^v; \hat{w}_j^v) = 1.$$

We next define what we mean by firm conduct: in this setting, we define conduct as assumptions about the content of \hat{w}_{-j}^v and firms' use of v in wage setting. Applying each conduct assumption, we find that each conduct assumption implies a distinct markdown:

| Conduct | use v ? | Firm's \hat{w}_{-j}^v | Equilibrium Wage(s) w_j^v |
|---------------|-----------|-------------------------|---|
| Perfect Comp. | No | — | γ_j |
| Monopsonistic | No | \bar{w} | $\frac{3}{4}\gamma_j + \frac{1}{4}\gamma_{-j} - \beta$ |
| Monopsonistic | Yes | \bar{w}^v | $\frac{3}{4}\gamma_j + \frac{1}{4}\gamma_{-j} - \beta(1 + \alpha v j)$ |
| Oligopsony | No | w_{-j} | $\frac{2}{3}\gamma_j + \frac{1}{3}\gamma_{-j} - \beta$ |
| Oligopsony | Yes | w_{-j}^v | $\frac{2}{3}\gamma_j + \frac{1}{3}\gamma_{-j} - \beta(1 + \frac{2}{3}\alpha v j)$ |

Next, we consider estimation and model selection. Each model, which we index by m , yields a wage equation of the form:

$$w_j^v = c_{\text{own}}^m \cdot \gamma_j + c_{\text{other}}^m \cdot \gamma_{-j} - c_j^{vm}$$

. A traditional approach in labor economics is to estimate $\hat{\mathbf{c}}$. To do so, one might first construct proxies for firm productivity γ_j and identify instruments that shift γ_j (and/or competitive environment). Then, one would regress w_j^v on γ_j , γ_{-j} , and concentration measures. To conduct inference, we might perform a simple Wald test on the parameter c_j , for instance: $H_0 : c_j \geq 1$, $H_a : c_j < 1$. Our approach (which follows the New Empirical Industrial Organization tradition) is to estimate $\hat{\gamma}$, rather than $\hat{\mathbf{c}}$. A particular conduct assumption m , in combination with labor supply parameters estimated in a prior step, determines the coefficients \mathbf{c}^m . Rather than searching for instruments for productivity, find instruments for markdowns that are excluded from productivity. Then, regress $w_j^v + c_j^{vm}$ on c_{own}^m and c_{other}^m to recover $\hat{\gamma}_j^m$; for example, when firms do not use v in wage setting, we have:

$$\begin{bmatrix} \hat{\gamma}_{-1}^m \\ \hat{\gamma}_{+1}^m \end{bmatrix} = \begin{bmatrix} c_{\text{own}}^m & c_{\text{other}}^m \\ c_{\text{other}}^m & c_{\text{own}}^m \end{bmatrix}^{-1} \begin{bmatrix} w_{-1} + c_{-1}^m \\ w_{+1} + c_{+1}^m \end{bmatrix}$$

Finally, in order to adjudicate between different forms of conduct, we use the Vuong (1989) and Rivers and Vuong (2002) tests, which compare model lack of fit between alternatives.

A.4 Details of EM algorithm

We estimate the parameters of the the preference model via the EM algorithm. Specifically, we use a first-order (or ‘‘Generalized’’) EM (GEM) algorithm, in which we replace full maximization of the surrogate function in the M step with a single gradient ascent update. Our algorithm proceeds as follows:

- **Initialization:** provide an initial guess of parameter values $(\boldsymbol{\beta}^{(0)}, \boldsymbol{\rho}^{(0)})$.

- **E Step:** at iteration t , approximate the average log integrated likelihood at $\boldsymbol{\beta}^{(t)}, \boldsymbol{\rho}^{(t)}$ with the function:

$$\mathcal{E}(\boldsymbol{\beta}, \boldsymbol{\rho} \mid \boldsymbol{\beta}^{(t)}, \boldsymbol{\rho}^{(t)}) = \frac{1}{N} \sum_{i=1}^N \sum_{q=1}^Q \alpha_{iq}^{(t)} \log \left(\alpha_q(x_i \mid \boldsymbol{\beta}) \times \mathcal{P}(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid \boldsymbol{\rho}_q) \right),$$

where the weights $\alpha_{iq}^{(t)}$ are given by:

$$\alpha_{iq}^{(t)} = \frac{\alpha_q(x_i \mid \boldsymbol{\beta}^{(t)}) \times \mathcal{P}(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid \boldsymbol{\rho}_q^{(t)})}{\sum_{r=1}^Q \alpha_q(x_i \mid \boldsymbol{\beta}^{(t)}) \times \mathcal{P}(\mathcal{B}_i^1 \succ \mathcal{B}_i^0 \mid \boldsymbol{\rho}_q^{(t)})}.$$

- **M Step:** Find $\boldsymbol{\beta}^{(t+1)}, \boldsymbol{\rho}^{(t+1)}$ by computing a single gradient ascent update (hence “first-order”).

We initialize our algorithm at 50 random starting values, and report the estimate that yields the highest likelihood.

A.5 Properties of bidding strategies

For clarity, we suppress dependence on m . Under each model m , we may generally write $G_{ij}(b) = \int \tilde{G}_{ij}(b, \lambda) dH(\lambda)$, where either $\tilde{G}_{ij}(b, \lambda) = \exp(u(b, a_i)) / (\exp(u(b, a_i)) + \exp(\lambda))$ under oligopsony or $\tilde{G}_{ij}(b, \lambda) = \exp(u(b, a_i) - \lambda)$ under monopsonistic competition. In the latter case, log concavity of $G_{ij}(b)$ follows directly from the fact that $u(b, a_i)$ is concave (by assumption), since $G_{ij}(b) = \exp(u(b, a_i)) \times \int \exp(-\lambda) dH(\lambda)$. Log concavity in the former case can also be shown via differentiation of $\log(G_{ij}(b))$.

Let the function $G_{ij}^+(b)$ (with derivative $g_{ij}^+(b)$) denote the right-hand side of the $G_{ij}(b)$ function, which replaces $\theta_0 + \theta_1 \cdot \mathbf{1}[b < w_i]$ with θ_0 . We similarly let $G_{ij}^-(b)$ denote the left-hand side function, which replaces $\theta_0 + \theta_1 \cdot \mathbf{1}[b < w_i]$ with $\theta_0 + \theta_1$. Clearly, $G_{ij}(b) = \mathbf{1}[b \geq w_i] \cdot G_{ij}^+(b) + \mathbf{1}[b < w_i] \cdot G_{ij}^-(b)$. Under the assumption that both $G_{ij}^+(b)$ and $G_{ij}^-(b)$ are log-concave, we have that the functions $g_{ij}^+(b)/G_{ij}^+(b)$ and $g_{ij}^-(b)/G_{ij}^-(b)$ are both strictly decreasing functions of b . This implies that both the left-hand and right-hand inverse bidding functions, $\varepsilon_{ij}^-(b) = b + G_{ij}^-(b)/g_{ij}^-(b)$ and $\varepsilon_{ij}^+(b) = b + G_{ij}^+(b)/g_{ij}^+(b)$ are monotone increasing functions of the bid. This in turn implies that the left- and right-hand bidding functions, which we denote by $b_{ij}^-(\varepsilon_{ij})$ and $b_{ij}^+(\varepsilon_{ij})$ are also strictly increasing functions of ε_{ij} . We may also define the left- and right-hand indirect expected profit functions as $\pi_{ij}^{*s}(\varepsilon_{ij}) = G_{ij}^s(b_{ij}^s(\varepsilon_{ij}))^2 / g_{ij}^s(b_{ij}^s(\varepsilon_{ij}))$ for $s \in \{-, +\}$, which are both strictly increasing functions of ε_{ij} . These results establish the monotonicity of firm strategies and payoffs in their unobserved valuations when firms bid on either side of the kink.

A necessary, but not sufficient, condition that the firm bids at the kink is that the derivative of the left-hand expected profit function is positive at the asked wage:

$$g_{ij}^-(w_i)(\varepsilon_{ij} - w_i) - G_{ij}^-(w_i) < 0.$$

We assume that $\varepsilon_{ij} > w_i$, since otherwise the firm would never choose to bid at ask. We additionally assume that both θ_0 and θ_1 are positive. Given these assumptions, we have that

$$g_{ij}^-(w_i)(\varepsilon_{ij} - w_i) - G_{ij}^-(w_i) < 0 \implies g_{ij}^+(w_i)(\varepsilon_{ij} - w_i) - G_{ij}^+(w_i) < 0,$$

since by construction $g_{ij}^+(w_i) < g_{ij}^-(w_i)$ and $G_{ij}^+(w_i) = G_{ij}^-(w_i)$. By the same logic, we can show:

$$g_{ij}^+(w_i)(\varepsilon_{ij} - w_i) - G_{ij}^+(w_i) > 0 \implies g_{ij}^-(w_i)(\varepsilon_{ij} - w_i) - G_{ij}^-(w_i) > 0.$$

These conditions guarantee that the firm's optimal choice of bid is unique, even incorporating the kink. Given these definitions, we can write the condition that firms bid at the kink as:

$$\varepsilon_{ij}^-(w_i) \leq \varepsilon_{ij} \leq \varepsilon_{ij}^+(w_i)$$

Therefore, we may write the firm's optimal bidding function as:

$$b_{ij}(\varepsilon_{ij}) = \begin{cases} b_{ij}^-(\varepsilon_{ij}) & \text{if } \varepsilon_{ij}^-(w_i) \geq \varepsilon_{ij} \\ w_i & \text{if } \varepsilon_{ij}^-(w_i) \leq \varepsilon_{ij} \leq \varepsilon_{ij}^+(w_i) \\ b_{ij}^+(\varepsilon_{ij}) & \text{if } \varepsilon_{ij} \geq \varepsilon_{ij}^+(w_i). \end{cases}$$

We have therefore shown that the firm's optimal strategy is a strictly increasing function of its valuation outside of the interval $[\varepsilon_{ij}^-(w_i), \varepsilon_{ij}^+(w_i)]$, and is flat within that region.

Next, we consider firms' participation decisions. The results established above imply that the firm's indirect expected profit function is a *strictly increasing* function of the firm's valuation:

$$\pi_{ij}^*(\varepsilon_{ij}) = \begin{cases} \pi_{ij}^{*-}(\varepsilon_{ij}) & \text{if } \varepsilon_{ij}^-(w_i) \geq \varepsilon_{ij} \\ G_{ij}(w_i)(\varepsilon_{ij} - w_i) & \text{if } \varepsilon_{ij}^-(w_i) \leq \varepsilon_{ij} \leq \varepsilon_{ij}^+(w_i) \\ \pi_{ij}^{*+}(\varepsilon_{ij}) & \text{if } \varepsilon_{ij} \geq \varepsilon_{ij}^+(w_i). \end{cases}$$

Firms participation decisions are therefore given by the condition:

$$B_{ij} = \mathbf{1} [\pi_{ij}^*(\varepsilon_{ij}) > c_j].$$

Since $\pi_{ij}^*(\varepsilon_{ij})$ is a strictly increasing function of the firm's valuation, an inverse indirect expected profit function exists and is also strictly increasing. Therefore, we may re-write the above equation as:

$$B_{ij} = \mathbf{1} \left[\nu_{ij} > \pi_{ij}^{*-1}(c_j) - \gamma_j(x_i) \right].$$

A.6 Proof of the consistency of \hat{c}_j^m

Our proof of the consistency of \hat{c}_j^m for each firm j (and model m) closely follows the proof of Lemma 1 (ii) of Donald and Paarsch (2002). For clarity, we omit j and m indices. Let n denote the total number of bids, with $n \rightarrow \infty$. A sufficient condition for establishing consistency is the existence of a vector of candidate characteristics $x \in \mathcal{X}$ (including ask

salary a) occurring with positive probability such that there is a positive probability the firm optimally bids below ask for candidates with those characteristics: $\exists x \in \mathcal{X}$ such that $\Pr(a > b_i > 0 \cap x_i = x) > 0$. The vast majority of firms (92%) bid below ask at least once, which suggests that this assumption is reasonable. The vector x need not be the same for all firms. This assumption implies that the distribution of model-implied option value upper bounds $\hat{\pi}_i$ is bounded below by c when $x_i = x$, and that $\Pr(\hat{\pi}_i \in [c, c + \delta] \mid x_i = x) > 0$ for arbitrary $\delta > 0$. Let n_x denote the number of bids made to candidates with characteristics x and let \hat{c}_x^n denote the minimum implied $\hat{\pi}$ among those bids (such that $\hat{c}^n = \min_{x' \in \mathcal{X}} \hat{c}_{x'}^n$). Our sampling assumptions imply $n_x \xrightarrow{\text{a.s.}} \infty$. For an arbitrary $\epsilon > 0$, note that $\Pr(|\hat{\pi}_i - c| > \epsilon \mid x_i = x) = \Pr(\hat{\pi}_i > c + \epsilon \mid x_i = x) = 1 - F_\pi(c + \epsilon \mid x_i = x) < 1$. Let $\bar{F}_{\pi|x}(a) = 1 - F_\pi(a \mid x_i = x)$. We then have that $(\bar{F}_{\pi|x}(c + \epsilon))^{n_x} \xrightarrow{\text{a.s.}} 0$, and therefore $\Pr(|\hat{c}_x^n - c| > \epsilon) = \Pr(\hat{c}_x^n > c + \epsilon) = E\left[(\bar{F}_{\pi|x}(c + \epsilon))^{n_x}\right]$. Since ϵ is arbitrary, $\hat{c}_x^n \xrightarrow{\text{P}} c$, and since $\hat{c}_x^n \geq \hat{c}^n \geq c$, $\hat{c}^n \xrightarrow{\text{P}} c$. Further, $\sup_{m > n} |\hat{c}^m - c| = |\hat{c}^n - c| \xrightarrow{\text{P}} 0$ since \hat{c}^n is non-increasing in n , and so $\hat{c}^n \xrightarrow{\text{a.s.}} c$. \square