

UC Berkeley

UC Berkeley Electronic Theses and Dissertations

Title

Computational modeling of behavior under uncertainty: Commonalities and differences between anxiety and depression

Permalink

<https://escholarship.org/uc/item/561265zh>

Author

Gagne, Christopher R

Publication Date

2019

Peer reviewed|Thesis/dissertation

Computational modeling of behavior under uncertainty: Commonalities and differences
between anxiety and depression

By

Christopher R. Gagne

A dissertation submitted in partial satisfaction of the

requirements for the degree of

Doctor of Philosophy

in

Psychology

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor Sonia Bishop, Chair
Professor Anne Collins
Professor Joni Wallis
Professor Anca Dragan

Fall 2019

Computational modeling of behavior under uncertainty: Commonalities and differences between anxiety and depression

© 2019

By Christopher Gagne

ALL RIGHTS RESERVED.

Abstract

Computational modeling of behavior under uncertainty: Commonalities and differences between anxiety and depression

by

Christopher R. Gagne

Doctor of Philosophy in Psychology

University of California, Berkeley

Professor Sonia Bishop, Chair

Individuals who are prone to experiencing high levels of anxiety and depression often exhibit dysfunctional behavior. For example, anxious individuals often avoid situations that have even the slightest chance of a highly negative outcome (e.g. a plane crash), and depressed individuals often show a reduced pursuit of activities that most people find enjoyable. Progress can be made in understanding dysfunctional behavior by using formal, mathematical frameworks of decision making, which break down behavior into its computational components, and in which we can start to pinpoint the specific abnormalities associated with anxiety and depression. Chapter 1.2 and 1.3 review prior literature, highlighting some of the computations that seem to be altered, such as the overestimation for the probability that rare, extremely negative events will occur. Chapter 2 and Chapter 3 empirically examine behavior in situations that require individuals to accurately estimate the probability that an outcome will (or will not) occur as a result of their actions. In a task where individuals have to estimate action-outcome probabilities by trial-and-error (Chapter 2), individuals with high overall levels of anxiety and depression show a reduced ability to align the rate at which they learn to the rate of change in the environment (i.e. the level of *volatility*). In a task where individuals have to choose between options that have known (*risky*) versus unknown (*ambiguous*) probabilities (Chapter 3), individuals who have high levels of physiological anxiety tend to avoid the ambiguous options more than other individuals, as information is removed about those probabilities. On the other hand, individuals who are prone to experiencing mania are more likely to make the opposite choice, seeking ambiguity, when the outcomes are rewarding. Chapter 4 examines possible sources for dysfunctional beliefs, as opposed to behaviors. In a hypothetical vocational setting where individuals estimate their rank relative to others, individuals with high levels of anhedonia-related symptoms show initial beliefs that are more negative relative to the beliefs of others. Individuals with high levels of anxiety, on the other hand, show negatively biased updating of those beliefs in response to unbiased information. Chapter 5 summarizes the empirical findings and discusses more broadly how anxiety and depression seem to impact behavior (and its underlying computations) in uncertain situations.

Chapter 1: Introduction and Background

Chapter 1.1: Brief introduction

Mood and anxiety disorders are extremely prevalent. Over the course of any twelve-month period, over 9% of adults are diagnosable with an anxiety disorder and 18% with a mood disorder (Kessler et al., 2005). These disorders cause substantial disruption to daily life, interfering with social life, work, and even basic activities. Yet despite the magnitude of this problem, treatments for mood and anxiety disorders remain only partially effective (Ballenger et al., 1999; DeRubeis et al., 2005; Hollon et al., 2005). Better treatments require a better understanding of the source of this disruption.

A large part of the disruption to daily life involves dysfunctional behavior. Anxiety is often associated with avoidance behaviors, for example, avoiding flying for fear of a plane crash and opting to stay at home instead. Depression, on the other hand, is often associated with the reduced pursuit of pleasurable activities, for example, going out to see friends or even enjoying a meal on one's own. As a result of the actions they take, individuals with both anxiety and depression often experience poorer outcomes and less satisfaction in life than others.

One promising approach to understanding this dysfunctional behavior involves the use of formal, mathematical frameworks of decision making, which have been used classically in economics and engineering and have been more recently popular in neuroscience and psychology. These frameworks break down the complex decision making process into simpler, more abstract components (i.e. decision variables and the computations that act on them). The source of aberrant behavior and decisions can then be pinpointed to dysfunction in one or more components.

For example, the decision of whether to fly to a relative's for the Holidays can be thought of in terms of actions (fly or stay at home), potential outcomes (safe flight or plane crash), the probabilities of those outcomes, and the different values of those outcomes should they occur, all of which needs to be estimated and weighed up. Although decisions like this may be easy for some, individuals with anxiety and depression routinely struggle with this type of decision and often make choices (e.g. staying at home) that they later regret. Dissecting decisions using formal frameworks can help us understand why. Continuing the example, anxiety might lead an individual to overestimate the probability of a plane crash, while depression might lead him or her to underestimate of the value of visiting. If anxiety and depression are truly associated with these differences, the optimal tack for treatment would be quite different.

Chapter 1.2 (published in Annual Reviews of Neuroscience) reviews evidence for which computations are likely to be altered in anxiety and depression. It also reviews some of the neural mechanisms that likely underly these alterations. Finally, it repeats the call for the need to disentangle anxiety and depression, which are highly comorbid (Kessler et al., 2005). To that end, we highlight some computations that seem to be altered in both anxiety and depression,

and others that seem uniquely altered in one or the other. This differentiation of shared versus unique aspects of anxiety and depression is a theme for the empirical chapters of this dissertation.

Chapter 1.3 (published in *Current Opinions in Behavioral Science*) is more narrowly focused. In it, analogies are drawn between the more formalized notions of simulation and replay in reinforcement learning and the different types of persistent, negative thoughts that are characteristic of anxiety and post-traumatic stress disorder (PTSD). For anxiety, worry is proposed to be analogous to biased simulation, leading to the overestimation of the probabilities that rare, extremely negative events occur (e.g. plane crashes). This, in turn, results in more avoidance behavior and a further reliance on (biased) simulation. For PTSD, intrusive recollections are proposed to reflect over-generalized and/or unchecked updating of world models (and the values of states and actions contained in them) via simulation and replay.

Both **Chapter 2** and **Chapter 3** empirically investigate how individuals with anxiety and depression make decisions involving uncertainty. Difficulties in handling uncertainty have been associated with both anxiety and depression (Gentes & Ruscio, 2011; Carleton et al., 2012), but most of the previous work has relied on methods that did not differentiate between possible types of uncertainty or look at how these different types might impact decision making. Two main types of uncertainty are typically distinguished. First-order uncertainty refers to the case when outcomes occur probabilistically, rather than with certainty. Second-order uncertainty refers to the case when the probabilities themselves are unknown. **Chapter 2** investigates one particular type of second-order uncertainty, referred to as *volatility*, which reflects change in probabilities over time. **Chapter 3** investigates another type of second-order uncertainty, referred to as *ambiguity*, which reflects missing information about probabilities.

Chapter 2 (to be submitted as a journal article) looks at how individuals, differing in their levels of mood and anxiety symptoms, adjust their rate of learning about probabilities between periods differing in their level of volatility. Individuals with high levels of anxiety have previously been observed to have a difficulty aligning their learning rate to the level of volatility (Browning et al., 2015), either learning too quickly when things are stable or too slowly when they are rapidly changing (volatile). This mismatch can lead to inaccurate estimates and poorer choices. In **Chapter 2**, we show that this deficit extends to clinical levels of both anxious and depressive symptomatology, providing evidence that it acts as a transdiagnostic vulnerability factor to both types of disorders. This chapter also advances the computational modeling of both behavior and symptomatology, by combining state-of-the-art methods from each domain.

Chapter 3 looks at how individuals, differing in their levels of anxiety and mania-related symptomatology, make decisions involving ambiguity. We look at differences in attitudes towards first-order uncertainty (called *risk* in this study) and second-order uncertainty arising due to missing information (i.e. ambiguity). Attitudes towards risk and ambiguity have long been studied in economics, but they only recently have started to be used to understand why individuals with anxiety seem to be avoidant and why individuals prone to mania tend to engage in risky behaviors. For decisions involving potential financial gains or losses, we show that symptoms of mania and symptoms of physiological anxiety are associated with partially opposite attitudes towards ambiguity, but no differences in risk attitudes. If replicated in future work, this would provide a clearer picture of why individuals who have these different

symptoms exhibit different behavior in the real world. We additionally leverage our dataset to weigh in on questions that are still being investigated in economics surrounding typical attitudes towards ambiguity. We show that individuals, on average, tend to be ambiguity averse for decisions involving gains and ambiguity seeking or neutral for decisions involving losses depending on the level of missing information. This finding mirrors a more well-known one for risk attitudes (Kahneman & Tversky, 1979), but has been much less well characterized to date.

Chapter 4 (to be submitted as a journal article) looks at potential sources for negative beliefs and judgements that have been reported by individuals with high levels of anxiety and depression. These beliefs can also be seen as inputs to the decision making processes studied in the earlier chapters. Specifically, we look at whether these individuals tend to bring negative prior beliefs to new situations or whether they update their beliefs using new information in a biased way. An additional aim of this chapter is to examine dysfunction in a more realistic ('ecologically valid') setting. To that end, we used a novel task in which participants imagined themselves in a hypothetical internship, competing against their undergraduate classmates. We show that higher levels of depression are associated with more negative prior beliefs relative to others, but no differences in updating beliefs. In contrast, higher levels of anxiety were associated with overweighting negative versus positive beliefs during belief updating, but not with differences in prior beliefs.

Chapter 1.2: Anxiety, Depression, and Decision Making: A Computational Perspective (A review of the literature)

Previously published as:

Bishop, S.J., Gagne, C. (2018) Anxiety, Depression, and Decision Making: A Computational Perspective. Annu Rev Neurosci. doi: 10.1146/annurev-neuro-080317-062007. Epub 2018. Apr 25. PubMed PMID: 29709209.

The following slight formatting modifications have been made for the sake of coherence as part of this dissertation: section numbers and capitalizations.

1.2.1 Introduction

1.2.1.1 Overview

Anxiety and depressive disorders are highly comorbid (Brown et al. 2001, Kessler et al. 2005) yet also distinctive in their symptomatology. Individuals suffering from both anxiety and depression show difficulties with decision making. In this article, we use the computational decision making literature to outline the component processes that inform our decisions, and consider evidence pertaining to the influence of anxiety and depression on these processes. Although previous authors have considered alterations to decision making in anxiety or depression (Hartley & Phelps, 2012; Huys et al., 2015), there has been little attempt to review both in conjunction. We believe this is an important step if there is to be progress in characterizing unique, versus common, alterations to decision making in anxiety and depression and how these alterations might in turn contribute to the maintenance and distinctive features of these disorders. To this end, we also provide a schematic framework that can be used to integrate and further explore influences of anxiety and depression on both reward- and threat-related decision making. We note that, moving forward, it will be important to identify alterations in decision making specific to different subdimensions of anxiety or depression (examples of such subdimensions can be found in Bijsterbosch et al., 2014). We do not emphasize this finer differentiation in this review, given the current lack of pertinent empirical evidence. We do, however, highlight the few studies that have begun to investigate correlates of specific subdimensions of anxiety or depression.

1.2.1.2 Decision-Making Processes that Guide Reward Seeking and Threat Avoidance

Our actions, and those of other species, can be described in terms of attempts to obtain rewarding, or positive, outcomes and to avoid aversive, or negative, outcomes. Maybe we want a promotion, someone to say yes to a date, to escape being mugged, or to avoid being laid off at work. We can rarely be certain that a given action will achieve our goal. Further, it might be the case that no single action will suffice and we need to consider how alternate action

sequences will play out. Computational approaches to decision making provide us with a framework within which to understand the processes that inform our action selection.

Through our interactions with our environment, we gain information about the value of potential outcomes, the probability that those outcomes will be obtained following different actions, and the effort that different courses of action require. Multiple studies have provided evidence that such inputs are indeed used to inform action selection by humans and other animals (Camerer 1995; Chong et al., 2016; Schultz 2015). In many cases, the information needed to precisely calculate outcome value, probability, and effort costs is not fully available and these parameters must be estimated under varying levels of second-order uncertainty (Bach et al., 2011). Further, people vary in the subjective valuation of outcomes and the relative weighting of outcome probability and outcome value (i.e., risk aversion) (Kahneman & Tversky 1979). This raises the possibility that not only differences in the accuracy of parameter estimation but also differences in parameter weighting might characterize the decision making of anxious or depressed individuals.

Individuals might also vary in their reliance on different methods for estimating the potential value of alternate actions. Here, a distinction has been drawn between model-free and model based decision-making processes (Daw et al., 2005). In model-free learning, the individual is held to use the outcome of past actions (i.e., how often a given action has been followed by a rewarding or aversive outcome and the magnitude of that outcome) to update current estimates of the value of alternate actions. In model-based decision making, the individual is held to form a model of the world that includes the probability of transitions between states and the outcomes linked to each state in question (Daw & Dayan, 2014; Sutton & Barto, 1998). This enables the individual to evaluate actions using the summed long-run value that would be obtained by traversing a series of states. When we think of complex real-life situations (e.g., choosing between jobs or places to live), model-based reasoning provides an intuitively appealing formulization of the decision processes involved. However, full model-based reasoning is intractable; it rapidly becomes unmanageable to evaluate every possible series of states that might follow a given initial choice. It has been argued that off-line simulation of state transitions and associated outcomes might simplify in-the-moment comparison of the long-run value of alternate actions. It has further been proposed that recall, or replay, of actual experiences of state-to-state transitions and state–outcome associations might also inform these long-run value estimates (Foster & Wilson, 2006; Johnson & Redish, 2005; Lin 1992; Sutton 1990). The choice of which state transitions and state–outcome associations to replay or simulate might vary across individuals, as might also the extent to which replay or simulation is engaged. Similarly, in-the-moment comparison of alternate actions might also be influenced by the effects of mood on the accessibility of different states and outcomes.

1.2.1.3 Structure of the Review

In **Section 1.2.2**, we review the evidence for whether anxiety or depression influences the rate of model-free learning. In **Section 1.2.3**, we consider the evidence for whether anxiety or depression influences the accessibility of particular states and outcomes during simulation and recall, given the proposed role of these processes in model-based decision making. We also consider whether the extent of engagement in simulation or recall might be altered in anxiety

or depression. In the remaining sections, we review findings pertaining to the influence of anxiety and depression on the subjective valuation and weighting of parameters that likely inform both model-free and model-based decision making. Specifically, we review the evidence for whether anxiety or depression is associated with altered valuation of rewarding or aversive stimuli (**Section 1.2.4**), and with different willingness to engage in various levels of effort to achieve the desired outcome (**Section 1.2.5**). We also consider studies of risk aversion and the evidence for whether anxiety or depression is associated with differential weighting of outcome magnitude versus probability (**Section 1.2.6**). In **Sections 1.2.7** and **1.2.8**, we conclude with the provision of a schematic framework that can be used to characterize and further investigate the influences of anxiety and depression on the component computational processes that guide our decision making.

1.2.2 Influences of anxiety and depression on the rate of model-free learning

One way that anxiety and depression might affect action values is through model-free learning and the rate at which estimates of action values are updated following an unexpected outcome. Activity in the dopamine system and brain regions innervated by this system, including the striatum and regions within the frontal cortex, signals how much outcomes received diverge from outcomes expected (Schultz 1998; Schultz et al., 1997). This signal is called a prediction error, and changes in both outcome probability and outcome magnitude influence the size of this signal in the context of reward (Rushworth & Behrens, 2008). Similar prediction error signals are generated when aversive outcomes differ from expectation (Li et al., 2011; Mirenowicz & Schultz, 1996; Seymour et al., 2004) (though see Schultz 2016 for caveats regarding the interpretation of dopaminergic prediction errors for aversive outcomes and Dayan & Huys, 2009 for the potential role of serotonergic systems in aversive prediction errors).

The extent to which we use prediction errors to update our estimates of expected value depends on our current rate of learning. For optimal performance, we need to take into account levels of second-order uncertainty (Bach et al., 2011). The more confident we are in our value estimates, the slower we should be to change them. One source of second-order uncertainty is contingency volatility. If action-outcome contingencies are noisy but stable, such as when a given action leads to a given outcome three-quarters of the time, the lower the learning rate that is adopted, the less likely an actor is to suboptimally change behavior following intermittent unexpected outcomes. In contrast, when the probability or magnitude of outcome linked to a given action is rapidly changing, a high learning rate is required for the actor to avoid becoming stuck in a pattern of behavior that is no longer optimal. It might seem a tall order for individuals to be able to differentiate contingency volatility from contingency noise and adjust their behavior by moderating their learning rate accordingly. However, in the cases of both reward and aversive learning, healthy participants are remarkably accurate in their ability to do this (Behrens et al., 2007; Browning et al., 2015). Findings from the human and basic neuroscience literature implicate the anterior cingulate cortex (ACC) and the amygdala in the use of contingency volatility to modulate rate of learning (Behrens et al., 2007; Li et al., 2011; Roesch et al., 2012).

In contrast to low-trait anxious individuals, high-trait anxious individuals struggle to adapt their learning rate to current levels of volatility, especially in the case of aversive outcomes (Browning et al., 2015). There is no difference between high- and low-anxious individuals in mean learning rate, and high-anxious individuals do not show impaired prediction error generation—modulation of both pupil dilation and next-trial reaction time by outcome surprise (i.e., the unsigned prediction error) is unaffected by trait anxiety (Browning et al., 2015). The association between anxiety and impoverished adaptation of learning rate has been shown to also hold in the case of reward loss but not reward gain (Pulcu & Browning, 2017).

Turning to depression, a recent meta-analysis concluded that there was little evidence that learning rate differs between patients with major depressive disorder (MDD) and controls or varies as a function of level of anhedonic depression (Huys et al., 2013). However, the task used in the studies reviewed did not manipulate contingency volatility. Hence, this leaves open the question of whether depression, like anxiety, might be linked to a specific deficit in ability to adjust learning rate to contingency volatility. In the study reported by Browning et al. (2015), both the anxious arousal and the anhedonic depression subscales of the Mood and Anxiety Symptom Questionnaire showed the same inverse relationship to adjustment of learning rate as that reported for the State Trait Anxiety Inventory trait subscale (C. Gagne & S. Bishop, unpublished data). Further analysis revealed that learning rate adjustment was linked primarily to the shared variance of these two subscales. Given the high correlations ($r \geq 0.6$) between scores on the different scales, we need to establish whether this result replicates in a larger sample. For now, this finding provides tentative evidence that impoverished adjustment of learning rate to match environmental volatility might represent a common vulnerability linked to both anxiety and depression.

An inability to adapt learning rate to current levels of volatility is likely to result in individuals being less able to determine the best course of action when faced with unexpected outcomes. That this can affect high-trait anxious individuals' decision making even when contingencies are stable is supported by existing findings (Browning et al., 2015). One potential response may be to simply treat all environments as highly volatile; indeed, there is evidence that anxious participants do sometimes select this approach (Huang et al., 2017). If individuals either incorrectly treat stable environments as volatile or have high levels of uncertainty around their estimates of volatility (high meta-volatility), the breadth of distribution of potential values of a given action derived from model-free learning will be increased. This in turn might reduce subjective confidence in action value estimates.

1.2.3. Evidence for influences of anxiety and depression on simulation and recall processes

Tversky and Kahneman (Kahneman & Tversky, 1982; Tversky & Kahneman, 1974) introduced the idea that individuals might use availability heuristics when judging the probability of future events. Here, the contention is that the more instances we can recall of a given event having happened in the past, or the easier we find it to simulate the given event happening in the future, the higher our subjective judgment of the event's probability will be. When an event type is rare, we are generally less likely to recall or simulate instances of the event's occurrence; hence, we judge the event as less probable. In terms of current theories of

model-based decision making, our judgment of the probability of a given outcome is influenced by the number of state sequences sampled that result in that outcome.

Considerable evidence indicates that the emotional salience of events affects their recall (Cahill et al., 1996; Dolcos et al., 2017). Within the anxiety and depression literatures, it has been argued that mood-congruent biases might affect the relative ease with which individuals recall or simulate negative and positive events and that this in turn might influence judgments of the future probability of events. In the terminology of model-based decision making, anxiety and depression might influence the states and outcomes we consider when we engage in simulation and recall processes and when we use these processes to inform long-run estimates of the probability of various outcomes and the consequent summed value of a given course of action. We break down the empirical evidence in support of this claim below.

First, anxious and depressed individuals do indeed show altered judgments of the future probability of real-world emotional events. Patients with generalized anxiety disorder (GAD) and MDD and individuals with high subclinical levels of anxiety and depression show elevated estimates of the probability that they will experience negative events (Butler & Mathews, 1983; MacLeod et al., 1996; Muris & van der Heiden, 2006). Here, there is some evidence that estimates of the future probability of negative events are more strongly linked to anxiety than to depression when the two are teased apart (Muris & van der Heiden, 2006). Depression is also strongly linked to reduced estimates of the future probability of positive events, with evidence indicating this association is specific to depression rather than shared with anxiety (MacLeod et al. 1996; Muris & van der Heiden, 2006).

Other findings meanwhile link both ease of simulation and ease of recall to judgments of the future probability of real-world events. Macleod et al. (1991) reported that participants' ability to generate reasons why events might or might not occur significantly predicted their estimates of the probability of both future negative and future positive events. Similarly, ease of recall of past negative or past positive events has also been found to predict estimates of the probability of future events (MacLeod & Campbell, 1992). In addition, studies of simulation and recall in anxiety and depression have revealed valence-specific biases. When asked to generate, in a limited time, events that might happen in the future, anxious individuals show increased generation of negative events relative to control participants (MacLeod & Byrne, 1996). Meanwhile, groups characterized by depressed mood show reduced generation of positive events (Bjarehed et al., 2010; MacLeod & Byrne, 1996; MacLeod & Salaminiou, 2001; Moore et al., 2006). Multiple studies have also linked elevated depression levels to heightened recall of negative events and stimuli (Bradley & Mathews, 1983; Clark & Teasdale, 1982; Teasdale et al., 1980), with some additional evidence for depression being linked to reduced recall of positive events and stimuli (Bishop et al., 2004). We note that in the studies reviewed, it is difficult to dissociate the influence of mood-congruent biases on recall or simulation from the influences of differences in life experiences on recall or simulation. However, putting aside the origin of these biases, it does appear that anxiety and depression are associated with differences in the output of the recall and simulation processes central to model based decision making. In relation to negative events, in particular, there is some suggestion that recall biases might be more evident in depression and future-oriented simulation biases might be more evident in anxiety. However, the robustness of this dissociation and the relative influence of recall

processes versus simulation processes on judgments of the future probability of negative events remain to be established.

Additional potential evidence of alterations to simulation and recall processes in anxiety and depression comes from clinical studies of the role of repetitive cognitions, specifically worry and rumination, in anxiety and depressive disorders. Although both rumination and worry tend to be focused on negative outcomes, rumination is largely a past-oriented form of repetitive thought and worry is future-oriented in its focus (Watkins et al., 2005). Multiple cross-sectional and longitudinal studies have linked extent of rumination to levels of both current and future depressive symptomatology; worry has similarly been linked to both current and future anxiety symptomatology (for a comprehensive review of this literature, see Watkins 2008). However, many of these studies focused selectively on depression and rumination or on anxiety and worry. One exception was a study by Hong (2007). Here, the author found that worry predicted both future depressive and anxiety-related symptomatology, whereas rumination was more uniquely associated with future risk for depression. It has since been proposed that rumination and worry might be different facets of a transdiagnostic risk factor (McEvoy et al., 2013). In line with this, Kircanski et al. (2015) reported that rumination and worry were equally elevated across patients with GAD, patients with MDD, and patients comorbid for GAD and MDD, relative to healthy control participants.

Rumination and worry have been interpreted by some as maladaptive attempts at problem solving (Szabo & Lovibond, 2006; Treynor et al. 2003; Watkins 2008). In particular, the worry literature throws light on how simulation processes might become pathological in nature. Here, a focus on negative outcomes has been linked to both elevated reported frequency and uncontrollability of worrying (Szabo & Lovibond, 2006). Difficulty with terminating the simulation process also appears to be of importance. In model-based decision making terms, this difficulty might reflect failure to achieve a given stopping criterion. Szabo & Lovibond (2006) reported that worry characterized as uncontrollable was associated with failure to settle on a good solution to the problem being worried about. In addition, children with clinically significant levels of anxiety were more likely to report an inability to stop worrying until the perceived threat was removed (Szabo & Lovibond, 2004). In adults, anxiety has also been linked to an increased number and duration of periods spent worrying (Verkuil et al. 2007). These findings suggest that anxiety, and possibly depression, given its similar association with worry and rumination (Kircanski et al., 2015), might be characterized by a perceived or actual failure to successfully complete attempts at model-based decision making, leading to prolonged engagement in simulation and replay.

1.2.4. Sensitivity to threat and reward in anxiety and depression

Alterations to model-free or model-based decision making in anxiety and depression might interact with altered subjective valuation of rewarding or aversive outcomes. The main evidence pertaining to whether subjective valuation of aversive and rewarding outcomes is altered in anxiety and depression comes from studies of sensitivity to threat and sensitivity to reward. We review these in turn.

It has long been suggested that anxiety is linked to increased threat sensitivity, potentially as a result of amygdala hyperresponsivity to threat (Etkin et al., 2004; Mathews &

Mackintosh, 1998). However, an increasing number of studies are beginning to challenge this assumption. Blair & Blair (2012) reviewed studies examining both physiological and neural responses to visual threat stimuli in patients with GAD relative to healthy control subjects. Most studies reviewed reported either no difference in threat responsivity between groups or reduced threat responsivity in the GAD group. In other anxiety disorders, amygdala hyperresponsivity has been reported most consistently in response to disorder-related stimuli (e.g., social stimuli in social anxiety disorder, trauma-related cues in post-traumatic stress disorder). In these cases, it is difficult to dissociate elevated threat sensitivity from differential prior Pavlovian learning of conditioned responses to the stimuli in question (Blair & Blair, 2012; Shin & Liberzon, 2010). Within the pain literature, researchers have addressed whether sensitivity to primary aversive stimuli is altered in individuals with preexisting anxiety or depression. However, here, investigations of whether individuals with anxiety and depressive disorders show altered pain sensitivity have also produced inconsistent findings including both hyper- and hyposensitivity to pain (Wiech & Tracey, 2009).

The studies reviewed above examine the response to aversive stimuli upon presentation. The effects of anxiety on threat responsivity have also been studied during expectation of aversive stimuli. Studies of both rodents and humans have reported anxiety-related increases in the magnitude of physiological responses while subjects are waiting for the occurrence of aversive stimuli, especially when the delivery of these stimuli is unpredictable (Davis et al. 2010; Grillon et al. 2008). Further, anxiety has been linked to elevated estimates of the aversiveness, as well as of the probability, of potential future real-life negative events (Butler & Mathews, 1983). One possible explanation for these findings is that anxiety influences expectations about the subjective value of aversive outcomes, as opposed to modulating the immediate response to outcomes of a given magnitude. If this is the case, we would need to address why nonpathological responses to experienced outcomes do not lead to successful updating of expected outcomes. Arguably, this might arise from a combination of deficits in model-free learning (see **Section 1.2.2**) and heightened accessibility of aversive outcomes during model-based simulation (see **Section 1.2.3**). We further explore this possibility in **Section 1.2.7**.

Alterations in reward sensitivity have been studied primarily in relation to depression rather than anxiety. Self-reported anhedonia, the inability to derive pleasure from normally rewarding activities, is a diagnostic feature of MDD (Am. Psychiatr. Assoc. 2013) and a major dimension of depressive symptomatology. However, experimental studies have found little evidence for reductions in primary reward sensitivity in MDD. For example, several studies have failed to find reductions in pleasantness ratings of sucrose, or chocolate, in patients with MDD relative to control participants (Amsterdam et al., 1987; Potts et al., 1997; Scinska et al., 2004). Although we need to be cautious when interpreting null results, one possibility is that depression is associated primarily with altered processing of social rewards. This would still fit with depressed patients' reduced participation in normally rewarding activities as these tend to be social in nature. However, studies have also failed to find an influence of depression on perceived intensity of socially rewarding stimuli (Branco et al., 2017; Schaefer et al., 2010). An alternate possibility is that, in parallel to anxiety and threat, depression might be linked to biases in estimation of future reward value, as opposed to altered responsivity to actual outcomes. In line with this possibility, MacLeod & Salaminiou (2001) found that patients with

MDD gave lower estimates of the pleasure they would experience from various life events than did control participants, whereas Peeters et al. (2003) found that patients with MDD actually reported greater increases in positive affect after experiencing positive events than did control participants. This pair of findings is consistent with depression being linked to a lower expectation of reward value and a positive prediction error upon reward receipt. Here, again, the natural question is why would a positive prediction error not lead to updating of expected reward value over time. In parallel to the argument for anxiety and threat sensitivity outlined above, this could arise as a result of impaired model-free learning in conjunction with reduced sampling of states linked to highly rewarding outcomes during model based simulation and recall. As reviewed in **Section 1.2.2**, there is less evidence for altered model-free updating in the case of depression and rewarding outcomes than in the case of anxiety and aversive outcomes. An alternative possibility is that, in depression, elevated subjective valuation of effort costs might reduce engagement in actions aimed at obtaining reward and decrease opportunities for model-free updating based on actual experiences of reward. In line with this, Peeters et al. (2003) report that MDD patients experience fewer positive events than do controls. We further consider the evidence for altered valuation of effort in depression versus anxiety in the next section.

1.2.5. Valuation of effort: opposing effects of anxiety and depression

In depression, reduced participation in normally rewarding activities might reflect increased subjective valuation of the effort costs involved in pursuing these activities. Findings suggest that depression is associated with reduced preference for high-magnitude rewards that require high effort expenditure over low-magnitude rewards that require low-effort expenditure (Treadway et al., 2009; 2012). However, these findings might reflect either altered subjective effort costs or altered subjective valuation of rewarding outcomes.

In recent work on apathy, Husain and colleagues have used computational modeling to tease apart the influences of reward magnitude and required effort. Findings from these studies suggest that apathy is linked primarily to increased valuation of effort, as opposed to differential sensitivity to the magnitude of reward (Bonnelle et al., 2015, 2016; Chong et al., 2016). Given that individuals with MDD often show high levels of apathy, an important question is how much do apathy levels mediate differences in willingness to exert effort to obtain reward in patients with depression relative to healthy controls. Moving forward, larger-scale studies are required to disentangle the relationship between overall levels of depressive symptomatology, specific levels of apathy and of anhedonia, and effort valuation versus reward sensitivity. Investigation of the influences of anxiety on effort–reward trade-offs is also much needed.

At the neural level, ACC dysfunction is a strong potential candidate for contributing to altered effort–reward trade-offs in depression. ACC lesions, disconnection of the ACC and NAc core, and disconnection of the amygdala and ACC have all been demonstrated to result in a shift in behavior toward lower-effort, lower-reward options (Floresco & Ghods-Sharifi, 2007; Hauber & Sommer, 2009; Walton et al. 2003, 2009). Human neuroimaging findings have also reported altered ACC structural and functional connectivity in individuals with high levels of apathy (Bonnelle et al., 2016).

With regard to aversive outcomes, increased willingness to exert effort to avoid aversive outcomes is commonly observed in animal models of anxiety (Servatius et al. 2008). The forced swim test, more commonly associated with animal models of depression, measures exertion of physical effort when subjects are placed in water without the ability to reach a platform (Porsolt et al., 1977). Reduced time spent immobile in the forced swim test is a predictor of antidepressant effectiveness (Cryan et al., 2005; Porsolt et al. 1977), suggesting that neurochemical changes determining recovery from depression may be linked to alterations in the level of effort that an individual is willing to exert. Recent work has revealed that elevated anxiety leads to above-baseline levels of locomotion in the forced swim test and that this is reduced by anxiolytic agents (Lee et al., 2017). These findings are of interest because they reveal opposing effects of anxiolytics and antidepressants. This suggests that willingness to exert effort to avoid aversive outcomes might be a domain where differential correlates of depression and anxiety will potentially be observed in humans. Research, in human participants, into the trait correlates of willingness to deploy effort to avoid aversive outcomes is in its early stages. Initial findings suggest that negative affect is linked to increased deployment of effort (Nord et al., 2017). Further studies are required to tease apart the effects of anxiety and depression as well as to control for potential inter-subject differences in outcome valuation both prior to action selection and upon receipt.

1.2.6. Increase risk aversion in anxiety

Differences in the relative valuation of outcome magnitude versus outcome probability might also lead to differences in how individuals weigh competing options. One area where this has been studied is in relation to risk aversion and its association with anxiety. Patients with anxiety disorders report fewer risk-taking behaviors than patients with depressive disorders or healthy control participants (Maner et al., 2007). In the context of reward-based decision making, risk aversion has been studied by assessing participants' preferences for higher-probability, lower-value outcomes over lower-probability, higher-value outcomes. Although many individuals show some degree of risk aversion, several studies have reported that risk aversion is more pronounced in individuals with high levels of anxiety (for a review, see Hartley & Phelps, 2012).

When analyzing risk aversion findings, it is important to consider the computations that might give rise to apparent risk aversion. Individual differences in risk aversion might reflect preference for actions with high-probability outcomes, altered subjective valuation of low-magnitude versus high-magnitude outcomes, or a combination of both. In addition, if participants need to estimate outcome probability, individual differences in learning rate, or adaptation of learning rate to second-order uncertainty (both volatility and level of information available with which to estimate outcome probability), might come into play. Finally, if there is the possibility for loss of reward, then individual differences in valuation of reward gain versus reward loss are also likely to be pertinent.

Ragunathan & Pham (1999) investigated risk aversion using a simple paradigm in which information about outcome probability was directly provided. Participants chose between a high-risk (i.e., low-probability), high-magnitude reward option and a low-risk (i.e., high-probability), low-magnitude reward option using a simple gamble matched on expected value

(i.e., the product of outcome probability by outcome magnitude). Raghunathan & Pham (1999) found that induced anxiety led to increased selection of the low-risk, low-reward option, in contrast to induced sadness, which was associated with greater selection of the high-risk, high-reward option. One limitation of this study is that participants were presented with only a single gamble. In contrast, recent studies have typically adopted more complex paradigms with multiple trials.

One task commonly used to investigate risk aversion is the balloon analog risk task. Participants choose how far to pump up a virtual balloon, increasing their financial payout with each pump but losing all their winnings for a given balloon if they reach the unknown point at which the balloon explodes (Lauriola et al., 2014). Individuals with elevated levels of anxiety and worry tend to have earlier stopping points (Maner et al., 2007). This finding is taken as indicative of increased risk aversion. In this task, participants can use their experience with prior balloons to update their estimates of how likely a balloon is to pop at any given point. This popping point is unknown and variable (often drawn from two or more probability distributions); hence, the probability that the balloon will pop at any point needs estimating and there is considerable second-order uncertainty around this estimate. If anxious individuals are less able to choose an appropriate learning rate under such circumstances, as discussed in **Section 1.2.2**, one possible heuristic might be to adopt a safe early stopping point. Hence, risk aversion in this task might reflect either a deficit in learning or a preference for low risk despite uncompromised learning. This is difficult to disentangle. Further, the amount of money already gained and that will be lost if the balloon explodes increases with each pump. It is well established that many individuals weight potential losses more than potential gains (Tversky & Kahneman, 1992; for a review, see Schultz 2015). Individual differences in risk aversion versus loss aversion are also difficult to differentiate within the balloon analog risk task.

In a recent elegant computational study, Charpentier et al. (2017) teased apart the effects of risk aversion from those of loss aversion while investigating the correlates of both anxiety and depression. Patients with GAD showed elevated risk aversion relative to healthy control participants but did not differ from controls in how much they valued losses relative to gains. Patients with GAD and a concurrent diagnosis of MDD showed a level of risk aversion that fell in between that of patients with GAD alone and that of control participants. Analyses using continuous measures of trait anxiety and depression revealed that anxiety levels were positively correlated with risk aversion when controlling for depression, but depression levels showed no significant relationship with risk aversion when controlling for anxiety. The experimental power of these analyses was relatively low, and replication of these results is needed. However, taken together with the other findings reviewed here, these results suggest that anxiety is linked to increased preference for low-risk options. In contrast, there appears to be little evidence of a unique relationship between depression and risk aversion.

1.2.7. A schematic framework of altered computations underlying decision making in anxiety and depression

In the sections above, we have reviewed findings pertaining to the influences of anxiety and depression on the component processes involved in decision making. In this penultimate section of the article, we put forward a schematic framework of altered decision making in

anxiety and depression. Our intention is to both integrate the findings reviewed above and provide a framework of potential value to future computational psychiatry studies. In **Figure 1.2.1**, we illustrate how the altered computations underlying decision making (ACDM) framework can be applied to account for reduced engagement in rewarding activities in depression and increased engagement in avoidance behaviors in anxiety and the vicious circles that might consequently ensue and contribute to maintenance of psychopathology. We expand our discussion of this framework within the rest of this section and in **Section 1.2.8**. We note that the ACDM framework is easily extended to address altered decision making in other forms of psychopathology.

The ACDM framework illustrates how model-free and model-based decision making processes might interact to contribute to altered action selection in both anxiety and depression. Some of the choices we face in everyday life reoccur only intermittently. This influences the potential frequency of model-free updating of action values following actual outcomes. Further, if we choose not to act, we gain no new evidence about outcome value or probability or the effort that would be required and hence do not have the opportunity to update our action values. Similarly, if we act to avoid a given outcome, we gain no new information about how bad or probable that outcome would have been in the absence of the action taken. The potential consequences of such path dependency for anxiety and depression are outlined in **Figure 1.2.1** and further considered below.

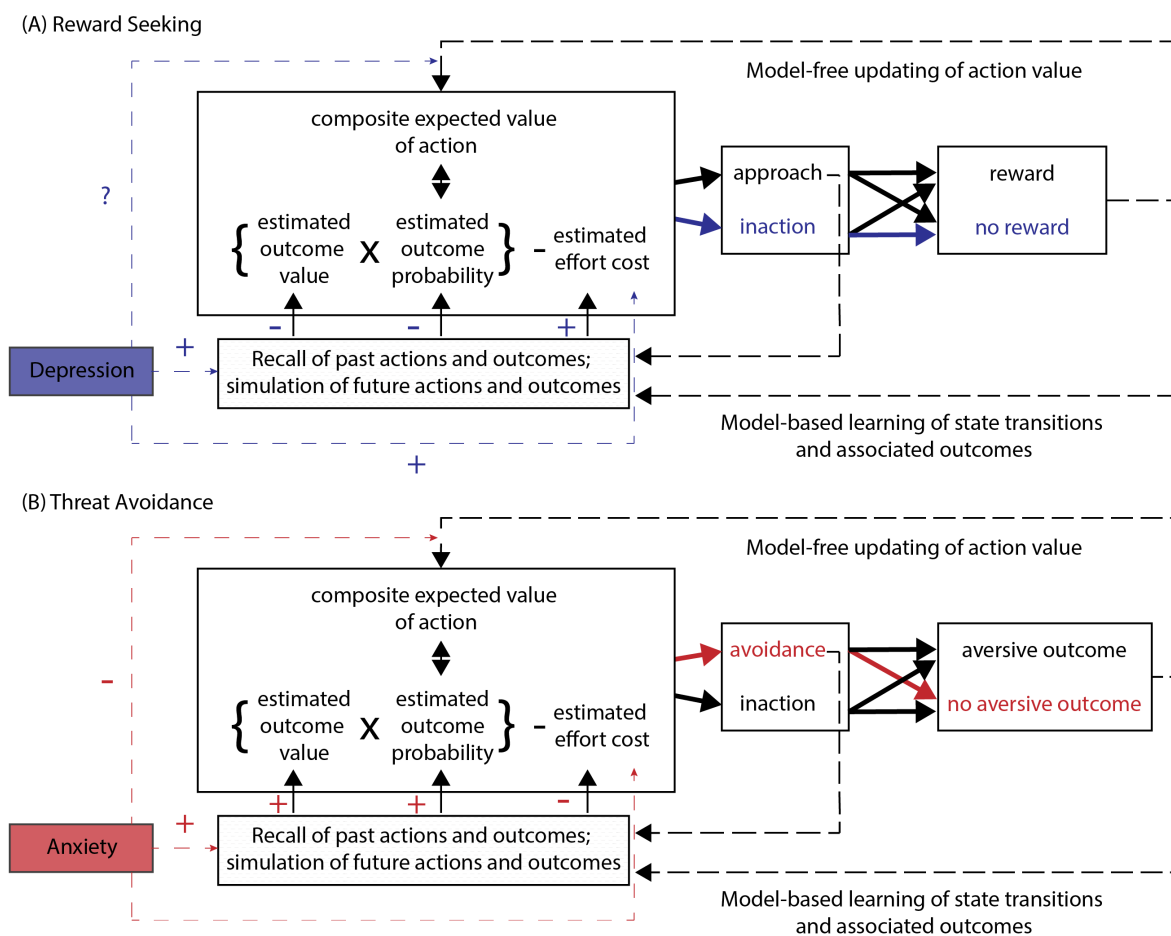


Figure 1.2.1 The altered computations underlying decision making (ACDM) schematic framework: accounting for reduced reward-seeking behaviors in depression and increased threat-avoidance behaviors in anxiety. (a) Engagement in potentially rewarding activities. Model-free estimates of action value are updated over time on the basis of the individual's experiences. These estimates are influenced by the frequency with which the action in question has a positive outcome, how subjectively rewarding the outcome is, and the level of subjective effort required. Model-based recall and simulation mechanisms enable off-line updating of estimates of the probability that a given action will directly, or indirectly, result in one or more outcomes, as well as of the relative subjective value of those outcomes and the level of effort that achieving them, directly or indirectly via the action in question, will entail. These model-based estimates are also updated over time as a result of the individual's experience. However, in addition, current mood affects which state transitions and state–outcome associations are sampled during off-line simulation and recall and in-the-moment evaluation of alternate actions. An interplay between model-free and model-based processes occurs, with the relative influence of these processes on action selection varying across individuals and situations. Depressed mood is predicted to reduce availability of simulated states associated with future positive outcomes, especially those of high value, resulting in decreased estimates of the probability and value of rewarding outcomes. Difficulties recalling past experiences of reward following engagement in similar activities are also expected to affect these estimates. We speculate that depressed mood might similarly increase recall or simulation of state transitions involving high effort costs. On the basis of the forced swim test literature, we also propose a direct effect of depression on maximal effort levels that an individual is willing to exert. Together these influences are expected to lead to reduced engagement in actions that have the potential to result in rewarding outcomes. As a result, experience-based updating of model-free and model-based estimates will decrease, leaving calculations of action value increasingly susceptible to the influences of depressed mood state on simulation and recall processes. For now, we remain agnostic about whether

depression also affects confidence in action value estimates through failure to adjust model-free learning rates to environmental volatility. (b) Avoidance of aversive outcomes. Many of the aversive outcomes that we worry about are relatively rare. Hence, we are likely to have low certainty in our estimates of their probability and severity. This leaves room for simulation processes to have a strong influence on these estimates. Anxiety is hypothesized to modulate this influence through increased simulation of future states associated with feared outcomes, leading to overestimation of aversive outcome probability and magnitude. Anxiety is also predicted to increase the effort an individual is willing to commit to avoidance behaviors. Given the relative infrequency of aversive events, engagement in avoidance behaviors is likely to be reinforced by the nonoccurrence of aversive outcomes, especially if effects of anxiety on simulation processes have upwardly biased estimates of aversive outcome probability and magnitude. In addition, opportunities will be missed for learning an association between inaction and aversive outcome nonoccurrence. Finally, anxiety also affects the adjustment of the rate of learning to match the stability, or volatility, of the current environment, potentially also impairing the ability to learn from actual experience. Both engagement in avoidance behaviors and disrupted learning from actual outcomes are predicted to leave calculations of action value increasingly susceptible to influences of anxiety on simulation processes, enabling a vicious circle to develop.

Experience-based learning also informs model-based estimates of state-to-state transition probabilities and state–outcome associations. However, during both off-line simulation and recall and in-the-moment selection between alternate actions, the relative ease with which different states and outcomes are accessed is held to be susceptible to fluctuations in mood state and affected by mood-congruent biases linked to both anxiety and depression (findings linking anxiety to increased accessibility of negative outcomes when simulating future events and depression to decreased accessibility of positive outcomes when simulating future events are reviewed in **Section 1.2.3**). This might contribute to biases in model-based estimates of both outcome probability and value (see **Sections 1.2.3** and **1.2.4**) and potentially also affect model-based estimates of effort costs. We note that the ACDM framework includes a reciprocal influence between model-free and model-based estimates of action value. This reflects the contention that output from simulation- and recall-based processes influences model-free estimates of action value (Sutton 1990), as well as evidence that model-free action value estimates are integrated into model-based calculations under certain circumstances (Keramati et al., 2016).

In **Figure 1.2.1**, we illustrate how influences of anxiety at various stages of the decision-making process might lead to increased engagement in actions aimed at avoiding aversive outcomes. Learning theorists have long argued that avoidance behaviors play a key role in maintaining anxiety disorders (for a review, see Krypotos et al., 2015). Within the ACDM framework, selection of avoidance related actions reduces opportunities for updating estimates of the value of aversive outcomes, or their probability of occurrence in the absence of the avoidance behavior engaged, on the basis of actual experiences. This maintains high levels of second-order uncertainty around these estimates, which is likely to be compounded by anxious individuals' difficulty in using second-order uncertainty to inform learning rate (Browning et al., 2015). The predicted consequence is that estimates of action value will be highly susceptible to influences of mood-congruent biases on simulation processes. As reviewed in **Section 1.2.3**, anxiety is linked to increased time spent worrying about potential future negative outcomes (Verkuil et al., 2007). In the ACDM framework this translates to heightened engagement in simulation processes, especially ones focused on the occurrence of aversive outcomes. This heightened engagement is expected to lead to an increasing imbalance between model-free

and model-based influences on decision making and to a vicious circle of greater valuation, and selection, of avoidance behaviors, worsening anxiety, and increased time spent worrying, with estimates of the probability and severity of future possible aversive events becoming increasingly reliant on mood-congruent simulations and decreasingly based on actual experience.

The ACDM framework also illustrates how a distinct vicious circle may maintain the association between depression and reduced pursuit of rewarding activities. On the basis of the findings reviewed in **Sections 1.2.3–1.2.5**, depression is proposed to be associated with both increased valuation of effort costs and decreased estimates of the expected value of potentially rewarding outcomes. The latter might initially reflect effects of low positive affect on the ability to simulate experiencing and enjoying future rewarding outcomes (MacLeod & Salaminiou, 2001; MacLeod et al., 1996). Decreased estimates of expected value and increased estimates of effort costs are expected to decrease the choice to pursue reward. This in turn reduces opportunities for updating the action value of engaging in potentially rewarding activities on the basis of actual experience. As a result, action values may become increasingly reliant, across time, on output from simulation processes susceptible to mood-congruent biases.

1.2.8. Caveats and Conclusions

The schematic portrayal of the effects of anxiety and depression on decision-making processes provided in **Figure 1.2.1** is inevitably highly simplified. The studies reviewed in this article support differential influences of anxiety versus depression on threat avoidance versus reward seeking. However, the evidence is far from clear-cut and our ability to draw conclusions is hampered by the relative lack of studies, in humans, addressing influences of anxiety on decision making about rewarding outcomes and influences of depression on decision making about aversive outcomes. It would be of value for future studies to seek to remedy this imbalance.

One area where evidence is mixed concerns whether depression, and not just anxiety, is linked to elevated estimates of the probability and subjective value of real-world aversive outcomes. Here, we note that clinical studies have revealed that patients with MDD show levels of worry similar to those shown by patients with GAD (Kircanski et al., 2015). Further, depression is also highly associated with rumination (Kircanski et al., 2015; Watkins 2008) and increased recall of negative events and stimuli (Bradley & Mathews, 1983; Clark & Teasdale 1982; Teasdale et al. 1980). However, even if elevated levels of rumination or worry exert an influence on depressed individuals' estimates of the expected value of aversive outcomes, an important predicted difference between anxiety and depression pertains to the willingness to exert effort to avoid aversive outcomes. Within the ACDM framework, influences of depression on estimated effort costs are predicted to reduce engagement in the active avoidance behaviors that are thought to play a key role in the maintenance of anxiety disorders. Effectively, in depression, a bias toward inaction may counteract a bias toward overvaluation of aversive future outcomes, facilitating learning about the nonoccurrence of the feared event when avoidance behaviors are not engaged.

In **Figure 1.2.1**, we also do not consider how the ACDM framework can encompass effects of anxiety on reward-related decision making. As reviewed in **Section 1.2.6**, anxious

individuals are more risk averse than healthy control participants when pursuing reward. This effect appears to be unique to anxiety, as opposed to shared with depression (Charpentier et al., 2017). In the context of our ACDM framework, risk-averse behavior is predicted to increase feedback about, and hence decrease uncertainty about, the expected value of low-risk reward outcomes. In other words, over time, low-risk options will also become more information rich (i.e., lower in estimation uncertainty) than high-risk options. Both unpublished empirical data from our laboratory and findings from the clinical literature (Dugas et al., 1998) suggest that anxiety is associated with avoidance of information-poor options. Hence, increases in the relative information level of low-risk versus high-risk options might potentially augment anxious individuals' engagement in risk-averse behaviors.

Inevitably, this review is not exhaustive in its scope. Availability-based heuristics are unlikely to be the only heuristics to influence decision making regarding rewarding and aversive outcomes. In addition, we have not covered the literature on Pavlovian learning or on the influences of Pavlovian-to-instrumental transfer on willingness to engage in approach or avoidance behaviors (Talmi et al., 2008). The role of temporal discounting in decision making is an additional topic that we have not had space to discuss. Finally, we have focused largely on how depression and anxiety influence the component processes supporting decision making, as opposed to discussion of the underlying brain mechanisms. Whereas computational neuroscience studies have advanced our understanding of the neural substrate of decision-making component processes in healthy subjects (for reviews, see Rushworth & Behrens, 2008; Schultz 2015), there is currently limited evidence pertaining to the specific influences of anxiety and depression on this neural circuitry, especially in humans. Over the next few years, the burgeoning field of computational psychiatry will hopefully provide further insight into the neural substrate of anxiety- and depression-related deficits in decision making. In particular, we look forward to being able to draw more concrete conclusions about the relative role of cingulate, striatal, and amygdala dysfunction.

To conclude, our review had two primary objectives. The first was to bring together studies that provide insight into which of the computations supporting decision making are altered in anxiety and depression as well as to highlight areas in which our knowledge is lacking. The second was to provide a schematic framework of how these alterations might lead to decreased engagement in potentially rewarding activities in depression and increased engagement in threat-avoidance behaviors in anxiety. Our review of the literature produced little compelling evidence for altered valuation of primary rewarding or aversive outcomes in anxiety or depression upon outcome receipt. In contrast, anxiety, and possibly depression, appears linked to increased estimates of the future probability and value of aversive outcomes, with depression also being linked to lower estimates of the future probability and value of rewarding outcomes. Findings from studies of rumination and worry and of recall and simulation in anxiety and depression suggest that increased engagement in recall and simulation processes, and mood-congruent influences on state and outcome accessibility, might contribute to these estimate biases. Differential effects of anxiety and depression on willingness to exert effort might additionally contribute to increased engagement in avoidance behaviors in anxiety and reduced engagement in potentially rewarding activities in depression. This in turn might affect opportunities for updating action value estimates on the basis of actual outcomes in a manner that sustains these maladaptive behavioral patterns. Problems with

adjusting learning rate to match levels of second-order uncertainty in anxiety might further impair learning from actual outcomes.

Additional research is needed to establish whether depression also affects use of second-order uncertainty to adjust learning rate for either aversive or rewarding outcomes. It would also be of value to better understand the potential three-way trade-offs of outcome value, outcome probability, and effort costs and whether these trade-offs vary as a function of anxiety or depression levels. Finally, further research is needed to determine whether the effects of anxiety or depression on specific components of decision-making processes are mediated by levels of worry, rumination, apathy, or anhedonia, as well as other dimensions of interest. Results from such research might enable us to better predict patterns of difficulties with decision making and behavioral symptoms across patients with different profiles on these dimensional measures. We hope that the framework we have put forward will be of value in guiding such future studies and in interpreting their findings.

Chapter 1.3: When planning to survive goes wrong: predicting the future and replaying the past in anxiety and PTSD (A review of the literature)

Previously published as:

Gagne, C., Dayan, P., Bishop, S.J. (2018) When planning to survive goes wrong: predicting the future and replaying the past in anxiety and PTSD, Current Opinion in Behavioral Sciences, 24, 89-95. ISSN 2352-1546, <https://doi.org/10.1016/j.cobeha.2018.03.013>.

The following slight formatting modifications have been made for the sake of coherence as part of this dissertation: section numbers and capitalizations.

Introduction

In modern life, aversive events vary both in their frequency and severity. Shootings, terrorist incidents and plane crashes are rare, extremely negative events that might threaten our survival if experienced even just once. Avoiding exposure to such events and handling them appropriately if they occur is critical to our survival and well-being but, we argue, surprisingly hard to integrate smoothly into the course of our day-to-day lives. Here, we lay out this computational problem as a form of approximate Bayesian decision theory (BDT) (Berger 1985) and consider how miscalibrated attempts to solve it might contribute to anxiety and stress disorders.

According to BDT, we should combine a probability distribution over all relevant states of the world with estimates of the benefits or costs of outcomes associated with each state. We must then calculate the course of action that delivers the largest long-run expected value. Individuals can only possibly approximately solve this problem. To do so, they bring to bear different sources of information (e.g. priors, evidence, models of the world) and apply different methods to calculate the expected long-run values of alternate courses of action. Avoiding catastrophic events poses unique difficulties above those in other situations framed in BDT because the rarity of these events renders methods that work well in more typical situations, such as model-free learning, relatively less useful. As a result, model-based processes that more efficiently re-use and extend experience, such as replay and counterfactual simulation, become especially important both before and after these events.

In part 1, we discuss the computations required to take into account the potential future occurrence of yet-to-happen rare, extremely negative events as we plan and navigate our daily lives. We consider how individual differences in these computations might confer vulnerability to anxiety. In part 2, we focus on the computations required to update our models of the world and action policies after the occurrence of rare, extremely negative events. We explore how the re-experiencing symptomatology characteristic of Post Traumatic Stress Disorder (PTSD) might be understood in the context of these computations.

Part 1: anxiety and predicting the future

The survival circuits that are the focus of this issue provide rich, hard-wired (sometimes called Pavlovian) policies that directly determine particular actions in the face of immediate mortal threat. However, waiting until threats materialize is rarely wise; maximizing our chances of survival and well-being requires estimating the probability and cost of extreme negative events and developing strategies for ameliorating or avoiding them ahead of time.

When estimating the expected value of avoidance behaviors, one should weight outcome costs by outcome probabilities. In the case of rare, extremely negative events, the costs are so high that even small differences in probability estimates will have a huge impact on these expected values. High trait anxious children and adults produce higher estimates of the probability that future negative events (e.g. being involved in a road accident) will occur to them than do low trait anxious participants (Butler & Matthews, 1983; Macleod et al., 1996; Muris et al., 2006). Such differences in probability estimates might result in increased selection of avoidance behaviors despite the associated disruption to everyday life activities.

There are a number of potential sources of these anxiety-related differences in probability estimates; these include differences in the method of estimation used, in initial biases in the estimates (priors) and in the precision of estimate calculation. The probability of rare, extremely negative events is hard to calculate precisely because similar events have rarely, if ever, been actually experienced. Thus, probability estimates are likely to be broad, with weak upper bounds. If the world is rapidly changing, that is, volatile, these bounds should be weaker still, as only very recent outcomes are pertinent (Behrens et al., 2007). In anxiety, there is evidence for difficulties in estimating environmental volatility (Browning et al., 2015; Pulcu et al., 2017) and increased adoption of high volatility priors (Huang et al., 2016). Hence, anxious individuals might have even weaker upper bounds than other individuals for probability estimates of rare, extremely negative events. One strategy for robustly avoiding catastrophic failures is to adopt a worse-case scenario (H1 control; Doyle et al., 1989), that is, to rely on the upper bound as opposed to the mean or median of the distribution of outcome probabilities. Consistent with this, clinically anxious individuals are reported to engage in catastrophizing, focusing on worst case outcomes (Sandin et al., 2015). If anxious individuals do indeed show a combination of widened bounds for probability estimates of rare, extremely negative events and reliance on upper bounds during action selection, this might promote more frequent selection of avoidance behaviors.

If we seek to weigh up the benefits of certain behaviors (e. g. going on vacation in London) against the potential probability and cost of rare, extremely negative events (e.g. a plane crash or terrorist incident), we must calculate the long-run values of alternate actions. Long-run values take into account outcomes that might only arise several steps after the initial choice. The methods used for this are often conceived as living on a spectrum between so-called model-free and model-based computations (Doya 1999; Daw & Dayan, 2005). Model-free and model-based methods both aim to produce appropriate policies (which specify the actions that should be taken in different situations), however, they use information from the world differently to do so. Model-free methods such as temporal difference learning (Sutton 1988; Sutton & Barto, 1998) cache information during experience of the environment. They thereby create policies that are computationally straightforward to use to guide subsequent on-the-

spot action selection and have the speed to be well suited to the avoidance or mitigation of extremely negative events. However, use of model-free methods to create multi-step action policies aimed at avoidance of rare, extremely negative events is heavily compromised by the reliance of these methods on past experience, given the typical absence of past experience for such events.

In contrast, model-based methods construct internal representations of alternate states of the world, and of how alternate future courses of action might play out depending on the initial state encountered (Sutton & Barto, 1998; Daw & Dayan, 2015). Construction of such models of the world is informed by direct experience. However, indirect evidence such as vicarious experience or intuitive knowledge about the physical world can also be incorporated (Battaglia et al., 2013). Sampling from the model can be used to play out what might transpire given selection of a particular initial action, even if that action has not been taken in real life. Thus, model-based methods can anticipate states and outcomes that have never been experienced, a characteristic of particular value for working out how to avoid rare, negative events. Such sampling can be used directly for planning (Kocsis et al., 2006). However, it has also been suggested that sampling during off-line periods (such as quiet wakefulness or sleep) can be used to train model-free estimates of action values (Sutton 1990). The putative benefit of this is to create model-free action policies that are fast to use but nevertheless reflect the knowledge contained within the model. However, if the model, or samples drawn from it, is biased, then not only will model-based planning be biased, but the model-free policy trained on the basis of the samples drawn will become biased too.

Biased sampling is likely to impact all of us, to some extent, but might be a particular problem for individuals at risk for anxiety disorders. Given the impossibility of exploring all potential future states, we need strategies for restricting the states we consider. It has been suggested that we focus on states that are easily available (Kahneman et al., 1982). The frequency with which states have been encountered in the past is likely to impact their availability. This may result in low frequency outcomes being overlooked during the estimation of action values. However, relying on state frequency alone might be suboptimal in some situations, and it has been suggested that this might be offset by the oversampling of emotionally salient outcomes, especially those involving extreme (i.e. rare, high value) events (Lieder et al., 2018). In line with this, emotional salience has been shown to facilitate recall of past events (Dolcos et al., 2017), in particular in the case of extreme events such as terrorist attacks (Lichtenstein et al., 2011; Madan et al., 2014). That increased availability of such events might impact action valuation is indicated by findings that availability of positive and negative events during simulation and recall predicts estimates of the future probability of these events (Macleod et al., 1991; 1992). A current example is the reported drop in Southwest bookings the week after an engine broke apart resulting in the death of a passenger on one of their flights.

Anxious individuals potentially oversample negative extreme outcomes and associated antecedents to a greater extent than individuals low in anxiety. In line with this, participants with high anxiety levels selectively generate more negative possible future life events than low anxious participants within a limited period of time (Macleod et al., 1996). If anxiety is linked to oversampling of negative outcomes and their antecedents, the frequency of such simulation might also moderate the extent to which estimates of the values of avoidance behaviors are influenced by sampling biases. Worry, repetitive thinking focused on future potential threat,

imagined catastrophes and their possible prevention (Watkins 2008), is a common form of simulation of the real world. Elevated levels of worry are a defining feature of Generalized Anxiety Disorder and also characterize other anxiety disorders (Am. Psychiat. Assoc. 2013). Anxious adults report more worry episodes and greater overall time engaged in worry (Verkuil et al., 2007), and anxious children report being unable to stop worrying until the focus of worry is removed (Szabo et al., 2006).

Frequent, uncontrollable simulation of negative outcomes and their antecedents might contribute to the maintenance of anxiety disorders by increasing the subjective valuation, and selection, of avoidance behaviors. If anxious participants also show increased reliance on upper bounds of probability estimates for rare negative events (as discussed above), this will have a converging influence upon the overvaluation of avoidance behaviors. These behaviors, in turn, will reduce opportunities for anxious individuals to collect data showing that extreme negative events almost never occur, even if avoidance behaviors are not engaged. Therefore, there will not be the observations needed to correct estimation biases and stabilize a potentially detrimental cycle of increasing miscalibration of action value estimates and selection of avoidance behaviors. Such decision-theoretic path dependencies have been implicated in various other psychiatric contexts (Huys et al., 2015; Dayan et al., 2018).

Part 2. PTSD and replaying the past

Despite our best efforts, extremely negative events do occasionally occur. If such an event is survived, the balance of planning activities should shift toward avoiding the event being experienced again. This is both because the events occurrence might contain information useful for avoiding similar events in the future and because there might be autocorrelation in the occurrence of extremely negative events (e.g. when new predators enter an environment; Travis et al., 2013).

Off-line replay of prior experiences and simulation of counterfactual actions and associated outcomes provide a means to update action values following the occurrence of an extremely negative event (Sutton 1990). It has been argued that previous states should be prioritized for replay based on how much that replay would change value estimates (Mattar & Daw, 2018). One way to accomplish this is by tagging states based on how much their successors value has changed (Moore et al., 1983). If change in value estimates determines priority for replay, the astronomically large discrepancy in outcome value occasioned by the occurrence of a rare, extremely negative event would be expected to result in prioritized replay of that events antecedents (see **Figure 1.3.1**).

By replaying the states that preceded a rare, extremely negative event, we can ascribe more negative values to these states and the actions selected within them that led up to the events occurrence. Equally, actions that were not taken can be simulated, together with the possible outcomes of these actions. Should similar states be encountered again, the model-free system can use the updated action values to choose swiftly a safer course of action. However, this may not be entirely straightforward. Specifically, exposure to a rare, extreme event might increase the salience and availability of similar outcomes and increase the probability, from an otherwise negligible level, of such outcomes being simulated following various actions. This

might result in many courses of action being evaluated more negatively than before the experience of the extreme negative event.

Findings from the traumatic stress literature indicate that most individuals do indeed replay the antecedents of extreme negative events after their occurrence. Following extreme negative (also termed ‘traumatic’) events, such as motor-vehicle accidents, over 50% of individuals report intrusive recollections, flashbacks, or nightmares up to three months following the event (Ehlers et al., 1998; Bryant 2000). These phenomena are collectively referred to as re-experiencing symptomatology. It is also common for people to ruminate repeatedly on their experience, thinking about the events causes and ways that it might have been prevented, for example, ‘I was running late so I cut through town, if I had gone the long way round . . .’ (Ehlers et al., 1998; Ehlers & Clark 2000; El Liethy et al., 2006).

Some degree of re-experiencing, rumination and counterfactual thinking about the past might be functional in the wake of an extremely negative event, or trauma. Indeed, psychological accounts have argued that repetitive thinking in the wake of a traumatic event might be important for resolving the discrepancy between the event and pre-existing core beliefs or assumptions. Horowitz (1985) describes this process as ‘cognitive processing’ and Janoff-Bulman (1992) describes it as ‘integration’, resulting in what Tedeschi and Calhoun (2004) describe as ‘post-traumatic growth’. Critically, whereas reexperiencing symptoms decay rapidly over time for some individuals, for others they remain frequent and cause significant levels of distress and disruption to everyday life. Researchers have struggled to identify aspects of the re-experiencing process or content that predict post-traumatic growth versus disorder (Dekel et al., 2011). Here, we suggest that the re-experiencing symptomatology and negative cognitions (e.g. rumination and counterfactual reasoning) observed following an extreme negative event can be computationally operationalized in terms of replay and simulation. In the rest of this section, we consider how this operationalization might shed light on the potential determinants of healthy versus dis-ordered responses to the experience of an extremely negative event.

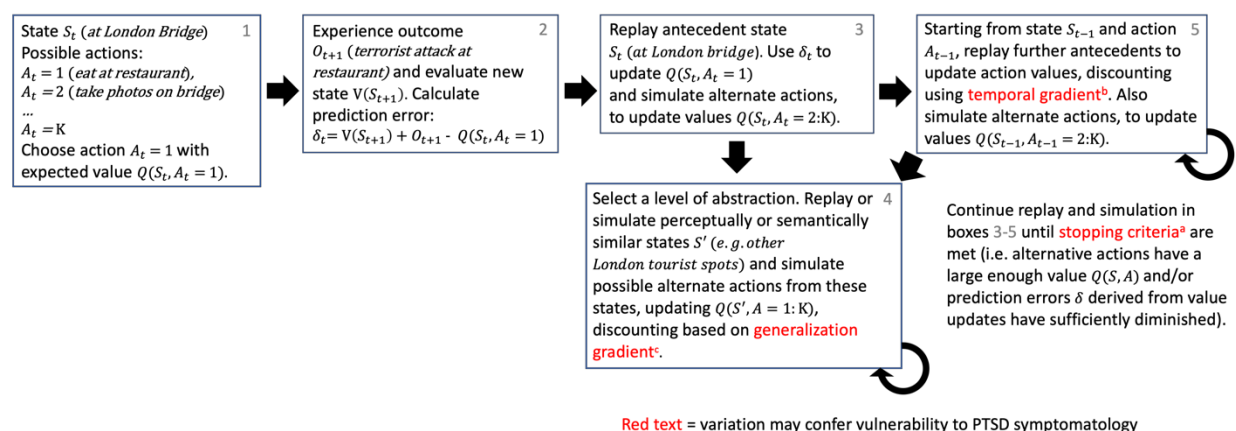


Figure 1.3.1. A replay and simulation account of PTSD. When a traumatic event is experienced, the difference in value between the expected and actual outcome is calculated (box 2). Given the high negative value of the traumatic event and its low prior probability of occurrence, a large negative prediction error (δ) will be experienced. This prediction error triggers the replay of the antecedent state and action (box 3). During this replay,

the antecedent state and action will be ascribed a more negative value. Alternate actions will also be simulated in order to identify a better counterfactual course of action from the antecedent state. The action values employed when the initial decision was made (box 1) will be updated. Pre-trauma, the sampling, and hence consideration during action valuation, of traumatic outcomes is unlikely. Post-trauma, the increased salience and availability of such outcomes will make their sampling more probable. Even very small increases in the estimated probability of actions leading to the same, or other, traumatic outcomes will result in substantial negative revision of action values given the large negative value of such outcomes. Replay and simulation will continue until a counterfactual action is found with a sufficiently large expected value ($Q(S_t, a)$) or until the prediction errors (δ) resulting from the replay and simulations have sufficiently diminished. Increased estimates of the probability of traumatic outcomes and associated downward revision of action values may result in many actions being considered in search for one with an acceptable value. Stopping criteria (thresholds for $Q(S_t, a)$ or δ) may vary across individuals and may potentially be influenced by external factors (e.g. novel stressors). Other states that share perceptual or semantic features with the state antecedent to the traumatic event will likely also be replayed or simulated and their associated actions may be ascribed a more negative value to the extent that shared features increase the availability, during simulation, of traumatic outcomes (box 4). The change in value of the state immediately antecedent to the traumatic event (S_t) will entail that the values of states and actions (starting from S_t, A_t) before that state will also need to be updated (Mattar & Daw, 2018; Moore et al., 1983). This will result in the replay and re-valuation of increasingly earlier states and actions (box 5). Generalization of values to similar states and the simulation and revaluation of alternate courses of action will also occur from these more distant antecedents. This cycle of replay will continue, with gradual discounting as more temporally distant states are revisited, until prediction errors (δ) sufficiently diminish or until counterfactual actions can be found at each temporal point in the antecedent chain and at each level of generalization with sufficiently high expected values. This process of replay of chosen action paths and simulation of alternate action paths might be phenomenologically experienced as intrusive thoughts, dreams, rumination and counterfactual reasoning. Red text is used to signify points where individual differences might confer vulnerability to elevated PTSD symptomatology. (a) Individuals may differ in stopping criteria; for example, individuals vulnerable to PTSD might be reluctant to take actions with even the slightest possibility of future catastrophe or might have a lower tolerance for negative changes in action or state values. (b) A reduced rate of discounting of the prediction error as more temporally distant antecedents are considered and (c) a shallower generalization slope when re-valuing states or actions that resemble those antecedent to the trauma are likely to increase the number of states and actions replayed resulting in higher levels of re-experiencing, avoidance and hyper-vigilance.

Several psychological accounts of PTSD posit that individuals with rigid beliefs about the positive nature of the world are more likely to experience PTSD after a traumatic event due to their assumptions or schema about the world being unable to flexibly accommodate or integrate the traumatic experience (Janoff-Bulman 1992; Park et al., 2012). According to BDT, if an individual has a much lower prior expectation of the occurrence of extremely negative events, this will generate a larger discrepancy between the value of the expected outcome and the value of the actual outcome (the traumatic event), strongly prioritizing the replay of the events antecedent states and actions (Mattar & Daw, 2018; Moore et al., 1983) and likely resulting in greater re-experiencing symptomatology. If rigidity is associated with an unwillingness to update state and action values to make them more negative, then negative prediction errors will stubbornly persist.

Individuals who go on to develop PTSD endorse more negative world views in the immediate aftermath of trauma (Ginzburg 2004). Moreover, elevated pre-trauma levels of anxiety and depression have been found prospectively to predict levels of post-traumatic stress symptomatology (Orr et al., 2012). Since anxiety and depression are linked to negative, not positive, biases in beliefs, interpretations of ambiguous events and judgements about the future (Beck 1976; Teasdale 1983; Mathews & Macleod, 2005), it seems unlikely that the

magnitude of the prediction error occasioned by experiencing a traumatic event is a key predisposing factor, at least for these individuals. Indeed, there is little empirical evidence that pre-trauma possession of optimistic priors confers vulnerability to PTSD.

There are other ways in which individual differences might influence extent of engagement in replay and simulation. One possibility is that individuals with a pre-trauma history of anxiety or depression might be more prone to difficulties with terminating disadvantageous replay or simulation processes. Such stopping difficulties might confer vulnerability to PTSD as well as to anxiety disorders (as described in Part 1) and depressive disorders. In line with this, both anxiety and depression are characterized by elevated levels of repetitive thoughts (worry and rumination) (Kircanski et al., 2015). Further, pre-trauma engagement in repetitive thinking is a significant predictor of post-trauma levels of PTSD symptomatology (Spinhoven et al., 2015). In the anxiety literature, difficulty finding a potential course of action with a sufficiently positive expected value has been associated with increased uncontrollability of worry (Szabo 2006). Post-trauma, inability to identify one or more counterfactual courses of action with high enough expected values to terminate simulation and replay might similarly be linked to elevated PTSD symptomatology (**Figure 1.3.1**).

An increased propensity for repetitive thinking might be further compounded by a disposition to over-generalization (Watkins 2011). Clinical accounts of PTSD describe how everyday sounds, like a balloon popping or a car backfiring, that resemble a gunshot can lead to extreme physiological and emotional responses in individuals for whom the traumatic event involved combat or gun violence. This has led to suggestions that over-generalization might be a key vulnerability factor in PTSD (Ehlers & Clark, 2000). In terms of the replay and simulation framework put forward here (see **Figure 1.3.1**), when an individual uses replay and simulation to update the value of states and actions that preceded a traumatic event, a key issue will be determining how far to go in also updating the value of other states and actions that share features with the antecedent states and actions. Other states may be related to antecedent states at a very concrete level (e.g. a similar looking street-corner to where the accident occurred) or at a very abstract level (e.g. any form of transportation in which you are not in control). Both selection of an abstract level and, at any given level, adoption of less steep (dis)similarity gradients might result in larger numbers of states and actions being re-evaluated both off-line and on-line when the individual encounters a state that shares features with a state antecedent to the traumatic event. Future planning would suffer from a similar problem.

There is strong evidence linking both forms of overgeneralization described above to anxiety, depression and PTSD. Patients with GAD, Panic Disorder, and PTSD have been shown to generalize from a conditioned stimulus across a wider range of perceptually similar shapes than healthy controls (Lissek 2012; Lissek et al., 2014; Kaczurkin et al., 2016). In addition, over-general autobiographical recall has been shown to characterize both patients with depressive disorders and individuals with a history of trauma (Williams et al., 2007). Further, levels of rumination have been shown to increase the influence of over-general memory on both future depressive and posttraumatic stress symptomatology (Hamlat et al., 2015; Kleim & Ehlers 2008). Trauma-analog studies have also reported that participants asked to ruminate abstractly, versus concretely, after viewing a traumatic video show both prolonged negative mood and more negative intrusions (Ehring et al., 2008; 2009).

In addition to re-experiencing, PTSD is also characterized by hyper-arousal and avoidance behaviors, together with other symptoms (DSM-V). Within a replay and simulation account of PTSD, both over-generalization and going further back in the chain of antecedent states as part of the replay process (see **Figure 1.3.1**) will result in more states and associated actions being re-evaluated as potentially dangerous. This, in turn, could lead to an increased sense of current threat, associated physiological reactivity and avoidance of multiple situations.

One important question for future research is whether individuals with a prior history of anxiety or depression are equally likely to show over-generalization between perceptually similar stimuli or states and over-abstract levels of simulation, or if the former might be more associated with anxiety and the latter with depression. In addition, the extent to which individuals vary in stopping criteria for terminating replay and simulation, and the role of this in anxiety, depression and PTSD, remains to be established.

Conclusion

In order to survive and maximize our wellbeing, we need to weigh up actions aimed at avoiding life threatening events against the pursuit of rewarding activities and avoidance of more minor aversive outcomes or losses. Here, we have outlined the computational processes involved in action valuation both in advance of, and subsequent to, the occurrence of rare, extremely negative events, and discussed how both anxiety and PTSD might be understood within this computational framework. Future work would valuably extend our computational analysis to consideration of potential neurobiological markers of disease risk. Given the putative role of hippocampal function in off-line model-based simulation processes (Wilson & McNaughton, 1994; Foster & Wilson, 2006; Ambrose et al., 2016) and evidence for hippocampal dysfunction in both anxiety (Gray & McNaughton, 2003; Khemka et al., 2017; Bach et al., 2014; Oler et al., 2010) and PTSD (Shin et al., 2006; Kheirbek et al., 2012; Stevens et al., 2018), this is a particular structure of interest.

Chapter 2: Decision Making Under Volatility

Introduction

Both mood and anxiety disorders substantially disrupt the daily lives of individuals afflicted by them. Part of this disruption comes from the uncertainty as to what might transpire following different courses of action. Both individuals with high levels of anxiety and those with depression report that uncertain situations evoke distress and disrupt their ability to make decisions (Dugas et al., 2001; Birrell et al., 2011; Gentes & Ruscio, 2011; Carleton et al., 2012). Another part of this disruption seems to come from a differential emphasis of negative versus positive outcomes. For example, both anxious and depressed individuals report that negative outcomes are more likely to occur to them, relative to others, across a variety of different situations (Butlers & Matthews, 1983); healthy individuals, on the other hand, tend to have optimistic expectations when making similar judgments (Sharot et al., 2011). Although a general difficulty in processing uncertainty and a differential processing of positive and negative outcomes has long been recognized in relation to anxiety and depression, it has only recently been investigated more precisely using the tools afforded by a computational approach to psychiatry (Montague et al., 2012; Wang et al., 2014; Huys et al., 2016).

Difficulties in processing uncertainty, in relation to anxiety, have recently been examined by Browning et al. (2015) who dissociated the impact of two different forms of uncertainty—contingency noise and contingency volatility—on the decision making of low and high anxious individuals. Within computational frameworks of decision making (Yu & Dayan, 2005; Behrens et al., 2007; Nassar et al., 2012; Payzan-LeNestour et al., 2013), contingency noise refers to the inherent probability in action-outcome associations (e.g. action 1 leads to outcome 1 with a $p=0.75$), and contingency volatility refers to the rate of change in these associations over time. Both noise and volatility should be taken into account in order to learn, effectively, action-outcome contingencies through experience. When volatility is low (i.e. contingencies are stable), unexpected outcomes should be more often attributed to noise, but when volatility is high (i.e. contingencies are changing frequently), unexpected outcomes should be more often attributed to a change in the contingencies. Whether an individual is making accurate attributions can be measured by the difference in their learning rates between volatile and stable periods, estimated within the context of a reinforcement learning model of their behavior (Behrens et al., 2007; Browning et al., 2015). In Browning et al. (2015), low anxious individuals were observed to have lower learning rates during stable periods and higher learning rates during volatile periods, consistent with an accurate attribution of unpredicted outcomes to noise when conditions were stable and to changes in contingencies when conditions were volatile. High anxious individuals, on the other hand, seemed less able to adjust their attributions in this manner. This deficit in learning rate adjustment to volatility led to more negative outcomes during the experiment, and likely contributes to a self-reported difficulty in handling uncertain situations (Dugas et al., 2001; Birrell et al., 2011).

Other recent computational studies of decision making in depression and anxiety have started to examine differences in processing both the absolute and the relative valence of

outcomes. Absolute valence refers to the type of outcome, either potential rewards or punishments, whereas relative valence refers to whether outcomes are better or worse than the average value within the context of each outcome type. For reward contexts, relatively better outcomes (referred to here as good outcomes) might be the receipt of a monetary reward, whereas for punishment contexts, relatively better outcomes might be the avoidance of an electric shock. Depression-related abnormalities, such as a reduced sensitivity to rewards, have been identified in the context of rewards (Huys et al., 2013), and anxiety-related abnormalities, such as increased punishment learning rates, have been identified in the context of threat-related outcomes (Mkrtchian et al., 2017). Abnormalities in processing relative valence in outcomes have also been observed. For example, Korn et al., (2014) measured differences in updating beliefs following desirable or undesirable information for the probability that adverse life events would occur. They observed that healthy individuals, relative to individuals diagnosed with major depressive disorder, had a bias in updating favoring desirable information, whereas this bias was absent in patients.

In most of these previous computational studies, either symptoms of anxiety or symptoms of depression were investigated, despite the increasing recognition in the clinical literature that their high rate of co-occurrence needs to be taken into account (Kotov et al., 2017). Within the clinical literature, bifactor analysis has become a recently popular tool for dealing with symptom co-occurrence, because it partitions symptom variance into that which is shared between anxiety and depression and that which is unique to one or the other. It represents shared variance using a general factor, which typically contains loadings from many different symptoms, and represents unique variance using one or more specific factors, which typically contain loadings from a single cluster of symptoms. Bifactor analysis has been applied within different groups (e.g. adolescents, students, community members, outpatients, etc.) and to different self-report measures of anxious or depressive symptomatology (e.g. BAI, BDI, IDAS, MFQ etc.) (Clark et al., 1994; Steer et al., 1995; Zinbarg & Barlow, 1996; Steer et al., 1998; Simms et al., 2008; Steer et al., 2008; Brodbeck et al., 2011). It has consistently revealed a substantial amount of shared variance, often termed ‘general distress’ or ‘negative affect’, which has been proposed as a trait vulnerability factor to both mood and anxiety disorders (Clark et al., 1994). In addition, separate specific factors for depression and anxiety have been consistently observed, with the depression-specific factors most often comprised of symptoms of anhedonia (Clark et al., 1994; Steer et al., 1998; Steer et al., 2008) and anxiety-specific factors most often comprised of symptoms of anxious arousal (Clark et al., 1994; Steer et al., 1998; Steer et al., 2008) or worry (Brodbeck et al., 2011).

In the current study, we used bifactor analysis along with a reinforcement learning model of behavior to investigate individual differences in processing volatility, the absolute valence of outcomes, and the relative valence of outcomes, in individuals with pathological levels of anxious and depressive symptoms. Our first research question was whether a deficit in learning rate adjustment to volatility, previously linked to trait anxiety (Browning et al., 2015), is related specifically to anxiety or to both anxiety and depression. We operationalized specificity to anxiety, specificity to depression, and generality to both anxiety and depression by using separate scores for participants on three latent factors from the bifactor analysis: a general factor, a depression-specific factor, and an anxiety-specific factor. Given that maladaptive responses to uncertainty have recently been argued to be a transdiagnostic

marker for both anxiety and depressive disorders (Gentes & Ruscio, 2011; Carleton et al., 2012; Boswell et al., 2013), we predict that learning rate adjustment to volatility will relate to the general factor, rather than the anxiety specific factor.

Our second research question was whether the relationship between learning rate adjustment to volatility and mood and anxiety symptoms is additionally modulated by absolute valence in outcomes (i.e. different learning rate adjustments in the context of rewards or aversive outcomes). In the previous study, the relationship between trait anxiety and learning rate adjustment was observed in the context of aversive outcomes and not rewarding outcomes (Browning et al., 2015). However, in our current experiment, we included individuals who were diagnosed with either major depression (MDD) or generalized anxiety disorder (GAD), and hence had more severe levels of symptoms than in the previous study. Therefore, we might expect that a deficit in learning rate adjustment to volatility also extends to rewarding contexts for individuals with high levels of depressive symptoms or for individuals with high levels of overall symptomatology (i.e. general factor scores).

Our third research question was whether learning rate adjustment to volatility and its potential relationship to mood and anxiety symptoms was additionally modulated by relative outcome valence (i.e. good or bad outcomes) within the context of rewards or aversive outcomes or across both. We did not have a specific prediction for how relative outcome valence may impact learning rate modulation to volatility. However, given the previous work linking depression to a lack of a bias in updating beliefs following desirable versus undesirable information (Korn et al., 2012), we might expect that the depression specific factor or the general factor is related to differences in learning following good versus bad outcomes.

To answer these three questions, we first fit a bifactor model to a set of standard self-report measures of anxious and depressive symptomatology. Using this bifactor model, we calculated a score for each participant on each of the three factors: the general, the depression-specific, and the anxiety-specific factors. We then used a reinforcement learning model to estimate differences in participants' learning rates as a function of volatility level, absolute outcome valence, and relative outcome valence. The relationship between these differences (and their potential interactions) were estimated in a hierarchical Bayesian framework. We first addressed these three research questions in an in-lab participant sample (experiment 1), consisting of healthy controls, patients diagnosed with major depression or generalized anxiety disorder, and a community sample. We then tested the replicability of our findings in an independent online sample (experiment 2).

Methods

Participants

For experiment 1, participants between the ages of 18 and 55 were recruited continuously from the local community until a pre-designated end date of August 2017. By the end date, we had recruited 58 participants in total: 12 participants who met diagnostic criteria for Generalized Anxiety Disorder (GAD), 20 participants who met diagnostic criteria for Major Depressive Disorder (MDD) (three with a secondary diagnosis of GAD), and 26 healthy controls

(i.e. screened to not meet any DSM diagnostic criteria) from the same community. To increase the number of subjects, 30 additional community recruited participants were added from a separate unpublished dataset. The procedure for these subjects was similar. However, they completed the task during an fMRI scanning session and were not screened for psychiatric condition, allowing for the possibility that some of these individuals might have had a diagnosable condition at the time of the experiment. See **Supplemental Table 2.1** for detailed demographics.

For experiment 2, 172 participants from Amazon's Mechanical Turk were recruited to perform an online version of the experiment.

Our third dataset consisted of mood and anxiety questionnaire data from 199 additional participants from UC Berkeley's psychology research pool. This data was used to test the generalizability of the factor structure that was estimated in the in-lab participant sample. This dataset originally contained 325 participants, but only 199 participants had no missing responses, which were required for the confirmatory factor analysis.

Experimental Protocol

Experiment 1 was conducted in lab. Screening for psychiatric diagnosis and taking of informed consent was done during the first experimental session. Diagnoses were determined using the research version of the structured clinical interview for DSM-IV-TR (SCID) administered by trained staff and supervised by an experienced clinical psychologist. We excluded participants if they were currently receiving treatment for psychiatric illness or had been prescribed psychotropic medication within the past 3 months. Participants meeting diagnostic criteria for OCD, PTSD, bipolar disorder, substance abuse, or showing any psychotic symptomatology, were also excluded. Those meeting criteria for inclusion in the study, were invited back for two additional sessions. During the second session, participants completed the aversive-learning version of the task. One control participant dropped out during this session. During the third session, participants completed the reward-learning version of task. The sessions were spaced at least 1 day, but no more than 1-week, apart.

Experiment 2 was conducted online. Participants were directed from Amazon's Mechanical Turk to an externally hosted website. There, participants completed both a reward and a loss version of the experiment within the same session. Participants were required to take two 5-minute breaks, one after filling out the questionnaires and another before completing the second task. The order of the gain and loss task was randomized across participants.

Participant Exclusion

We excluded data from either the reward or aversive task if there was equipment malfunction or if a participant reported after the session that he/she did not understand the task. In experiment 1, data from the aversive-learning task were excluded for 8 participants (4 patients and 4 controls). Data from reward-learning task were excluded for 5 participants (3 patients and 2 controls). Only two participants (both control subjects) had data excluded from both tasks. These exclusions left 86 participants in total for experiment 1.

For experiment 2, 25 participants were excluded from the online dataset for having greater than ten missed responses in each task, leaving 147 participants.

Self-report Measures of Mood and Anxiety Symptoms

Participants in experiment 1 completed several standardized self-report measures of negative affect and anxiety and depression symptomatology. Measures included: the Spielberger State-Trait Anxiety Inventory (STAI form Y; Spielberger, 1983), the Beck Depression Inventory (BDI; Beck et al., 1961), the Mood and Anxiety Symptoms Questionnaire (MASQ; Clark et al., 1995; Watson & Clark, 1991), the Penn State Worry Questionnaire (Meyer, Miller, Metzger, & Borkovec, 1990), the Center for Epidemiologic Studies Depression Scale (CESD; Radloff, 1977), and the Neuroticism subscale from the 80-item Eysenck Personality Questionnaire (EPQ; Eysenck & Eysenck, 1975).

Online participants in experiment 2 completed the BDI, MASQ, and STAI. These participants had similar distributions of scores as the in-lab participant. For STAI, the 25%, 50%, and 75% percentiles (in-lab; online) were (33;33), (48;42), (59;53). For BDI, they were (2;3), (9;7), (23;15). For the MASQ anhedonia subscale, they were (42,49), (68,64), (79,78).

The 199 additional participants from UC Berkeley's psychology research pool completed the same measures as the in-lab participants.

Bifactor Analysis of Mood and Anxiety Symptoms

The goal of the bifactor analysis was to create three sets of orthogonal scores for participants: one representing the overall level of mood and anxiety symptoms (i.e. common variance), a second representing the level of depression-specific symptoms (relative to the overall level of symptoms), and a third representing the level of anxiety-specific symptoms (also relative to the overall level of symptoms). These three sets of scores were calculated from three latent factors, which were estimated as part of a bifactor model. The three factors are referred to as the general factor, the depression-specific factor, and the anxiety-specific factor.

The data used for the bifactor analysis consisted of 128 individual questions from the following questionnaires or subscales within the questionnaires: the MASQ anhedonia subscale, MASQ anxious arousal subscale, the STAI, the BDI, the CESD, the PSWQ, and the EPQ-N neuroticism subscale. Responses were either binary (0-1), quaternary (0-4), or quinary (0-5). Response categories that were endorsed by fewer than 2% of the participants were collapsed into the adjacent category, in order to mitigate the effects of extreme skewness. Positively coded responses, such as "I feel happy", were reversed to ease the interpretation of factor loadings. Polychoric correlations were used to calculate the correlation matrix to adjust for the fact that categorical variables cannot have correlations in the full range of -1 to 1.

Before estimating the bifactor model, we specified its structure to have three factors (one general and two specific). This decision was guided by prior work (Clark & Watson, 1991) together with the results of eigenvalue decomposition of the covariance matrix. Only the first three eigenvalues were significantly greater than chance (as determined by comparison against eigenvalues obtained from a random normal matrix of equivalent size; Humphreys & Montanelli, 1975; Floyd & Widaman, 1995) (**Supplemental Figure 2.14**).

Following model specification, the Schmid-Leiman (SL) procedure was used to estimate the loadings of the individual items onto each factor (Schmid & Leiman 1956). This procedure performs oblique factor analysis followed by a higher-order factor analysis on the lower-order factor correlations to extract a single higher-order factor (i.e. a general factor). Both steps were done using external software: the 'omega' function from the Psych package in R.

Factor scores for each participant were calculated using the Anderson-Rubin method (Anderson & Rubin, 1956), which is a weighted-least squares solution that maintains the orthogonality of the general and specific factor scores. Scores for the participants in the online behavioral validation dataset (experiment 2) were calculated using the questions ($k=80$) from the subset of three measures (MASQ, STAI, BDI) administered to them. To check that factor scores could be reliability estimated on this smaller set of items, we calculated factor scores on the combined in-lab participant sample ($n=86$) and UC Berkeley student sample ($n=199$; which was used for validating the factor structure) separately using the full and using the reduced set of items. We obtained extremely similar scores for the general factor ($r=0.98$), the depression-specific factor ($r=0.96$), and moderately similar scores for the anxiety-specific factor ($r=0.59$) between the full and reduced sets. Prior to the calculation of the scores, each individual question was normalized across both datasets to make the resulting scores commensurate between datasets.

Experimental Tasks

Participants in experiment 1 completed both reward and aversive versions of a probabilistic learning under volatility task (Behrens et al., 2007; Browning et al., 2015). Each task was divided into a stable and volatile block of trials, each 90 trials long. On each trial, participants chose between two shapes with the aim of either accumulating monetary bonus in the reward version of the task or avoiding the delivery of a mildly painful electric shock in the aversive version of the task. Participants were instructed to consider the magnitude of the potential outcome, shown as a number inside each shape (**Figure 2.1a**), as well as the probability that the outcome would occur if the shape was chosen. Outcome magnitudes varied from 1 to 99 independently of the outcome probability and corresponded to different sizes of reward or to different intensities of electric shock. For the aversive task, the magnitudes were mapped different intensities of electric stimulation. The intensity of electric stimulation was calibrated before the task, so that the magnitudes between 1 and 99 mapped onto a subjective pain scaled from 1 (mildly unpleasant) to 7 (very unpleasant). Outcome probabilities could be inferred from the relative number of occurrences that the outcome followed the choice of one shape and not the other. In the stable block (90 trials), one of the two shapes delivered the outcome on 75% of the trials (i.e. with a 75% probability) and the other shape delivered the outcome on the remaining trials. In the volatile block (90 trials), the identity of the shape delivering the outcome with an 80% probability and the identity of the shape delivering the outcome with a 20% probability switched every 20 trials (**Figure 2.1b**). The order of the stable and volatile blocks was randomized across participants and between tasks.

The online versions of the two tasks in experiment 2 were similar to those used in experiment 1. The intra- and inter-trial timings for the online task were shortened slightly. In the reward-gain task, participants started with a total of 0 points. Outcomes with magnitudes

between 1 and 99 either increased this total or kept it the same. In the reward-loss task, participants started with 5000 points, and the magnitudes between -99 and -1 either decreased the total or kept it the same. The summed total number of points that participants had left at the end of both tasks was compared across all the participants. A bonus of \$3 was awarded to participants that scored in the top 5%, \$1 was awarded to those in the top 10%, and \$0.25 was awarded to those the top 50%. Participant performance on the tasks, measured by the percent of rewarded trials (or trials on which electric stimulation or loss was avoided), was similar between the online and in-lab participant samples (see **Supplemental Figure 2.13**).

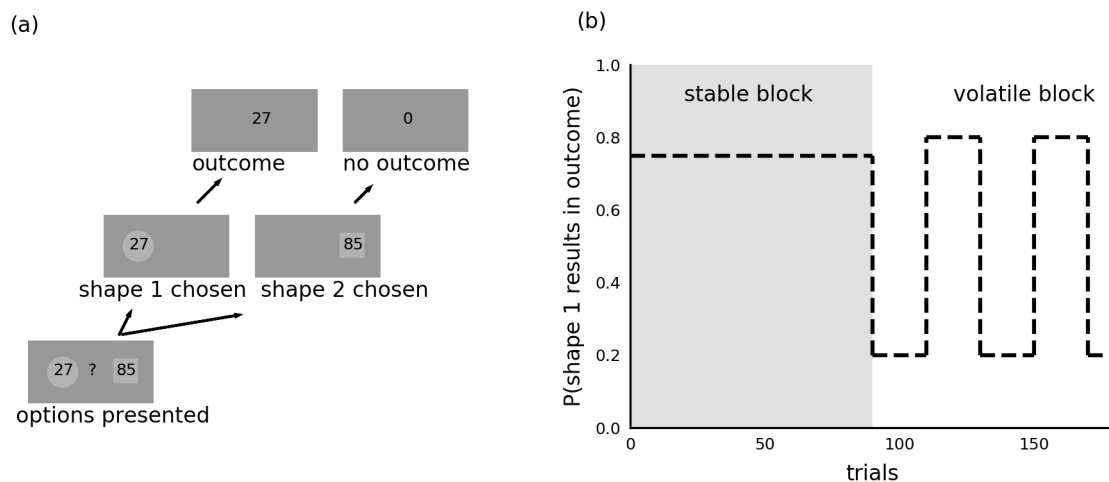


Figure 2.1 Task. (a) On each trial, participants chose between two shapes. Only one of the two shapes led to the outcome on each trial. The magnitude of the potential outcome was shown as a number inside each shape and corresponded to the size of the reward in the reward learning task or intensity of the electric shock in the aversive learning task. (b) Within each task, trials were organized into two 90-trial blocks. During the stable block, the one shape resulted in the outcome on 75% of the trials (i.e. with a 75% probability), while the other shape resulted in the outcome on the remaining trials. During the volatile block, the identity of the shape delivering the outcome with an 80% probability and the identity of the shape delivering the outcome with a 20% probability switched every 20 trials.

Estimating Differences in Learning Rates as a Function of Volatility Level, Absolute Outcome Valence, and Relative Outcome Valence

Our primary goal in modeling participants' choice behavior was to examine how learning rates were modulated by three experimental factors: absolute outcome valence, volatility level, and relative outcome valence. This corresponded to looking at differences in learning rates between the reward and aversive tasks, differences between the volatile and stable blocks within each task, differences between trials that followed good and bad outcomes within each task and block, and the three two-way interactions of these differences. In other words, for each participant, we estimated seven components for learning rate: (1) a baseline learning rate $\alpha_{baseline}$, (2) a difference in learning rates between the volatile and stable blocks $\alpha_{volatile-stable}$, (3) a difference in learning rates between the reward and aversive tasks $\alpha_{reward-aversive}$, (4) a difference in learning rates between trials following good and bad

outcomes $\alpha_{good-bad}$, (5) the interaction of volatility and absolute valence $\alpha_{(volatile-stable) \times (reward-pain)}$, (6) the interaction of volatility and relative valence $\alpha_{(volatile-stable) \times (good-bad)}$, and (7) the interaction of absolute and relative valence $\alpha_{(reward-pain) \times (good-bad)}$. On any given trial, a participant's combined learning rate, denoted simply as α , was calculated as follows:

Eqn 1.

$$\begin{aligned} \alpha = & \text{logistic}(\alpha_{baseline} \\ & + \alpha_{(volatile-stable)}\chi_{(volatile-stable)} \\ & + \alpha_{(reward-pain)}\chi_{(reward-pain)} \\ & + \alpha_{(good-bad)}\chi_{(good-bad)} \\ & + \alpha_{(volatile-stable) \times (reward-pain)}\chi_{(volatile-stable) \times (reward-pain)} \\ & + \alpha_{(volatile-stable) \times (good-bad)}\chi_{(volatile-stable) \times (good-bad)} \\ & + \alpha_{(reward-pain) \times (good-bad)}\chi_{(reward-pain) \times (good-bad)}) \end{aligned}$$

In Eqn 1., $\chi_{(volatile-stable)}$, for example, takes on a value of 1 when the trial is in the volatile block and a value of -1 when the trial is in the stable block.

The division of learning rate into seven components was supported by a model comparison analysis, in which we found that all seven of these differences were necessary for minimizing approximate leave-one-out cross validation error (LOO; Vehtari et al., 2017; see **Supplemental Methods: Model Comparison**).

Full Behavioral Model

Learning rates were estimated along with several other parameters by fitting a reinforcement learning model to participants choices in the two tasks. On each trial, participants were assumed to update an estimate p_t for the probability that the outcome would result from choosing shape 1 and not shape 2. Probability estimates were updated using the delta-rule (given in Eqn. 2), where the combined learning rate α determined how much the estimate was revised by the prediction error (i.e. the difference between the previous estimate p_{t-1} and the most recent outcome O_{t-1}).

Eqn. 2

$$p_t = p_{t-1} + \alpha (O_{t-1} - p_{t-1})$$

The difference between the probability estimate for shape 1 (p_t) and shape 2 ($1 - p_t$) was combined with the difference between the magnitudes of the potential outcome associated with each shape ($M_1 - M_2$) (in Eqn. 3). The mixture weight $\gamma \in [0,1]$ specified whether magnitude or probability was weighted more heavily in calculating the total outcome value for each shape. The mixture weight was also divided into the same seven parameter components; this division was also supported in the model comparison analysis.

The difference in magnitudes was also nonlinearly transformed using r to allow for different impacts of large and small differences (note that the sign for the difference was

temporally removed before exponentiating and then added back again). The nonlinearity parameter, r , was divided into a baseline across all conditions and a difference between reward and aversive learning tasks; we included this parameter, because it also improved LOO (see **Supplemental Results: Model Comparison**).

Eqn. 3

$$v_t = (\gamma)(p_t - (1 - p_t)) + (1 - \gamma)(M_1 - M_2)^r$$

Participants were also assumed to update a choice kernel k_t on each trial using the delta-rule (given by Eqn. 4). The update rate η determined how much to update the kernel using the participant's most recent choice C_{t-1} . The choice kernel keeps track of a participant's tendency to choose one shape over the other, which can influence subsequent choice in addition to outcome value. A single baseline update rate, η , shared across tasks and blocks was estimated for each participant. Adding the choice kernel improved LOO, but further allowing for different update rates for different conditions caused issues with non-convergence in parameter estimation; this indicated that participants' data did not contain enough information to estimate differences in update rates across conditions.

Eqn. 4

$$k_t = k_{t-1} + \eta (C_{t-1} - k_{t-1})$$

Finally, the outcome value v_t and the choice kernel k_t on the current trial were combined using two inverse temperatures, τ and τ_k , to calculate the probability of choosing shape 1 (using Eqn. 5). The outcome value inverse temperature, τ , was divided into the same seven components as learning rate, whereas the choice kernel inverse temperature, τ_k , was only divided into a baseline and a difference between reward and aversive tasks. These parameter divisions were again justified by model comparison.

Eqn. 5

$$P(C_t = 1) = \frac{1}{1 + \exp(-(\tau v_t + \tau_k(k_t - (1 - k_t))))}$$

The full model (Eqn. 2 through Eqn. 5), with parameters taking on different values depending on the condition, was fit to each participant's data across both the reward and aversive learning tasks. Variables were coded to have similar interpretations. The outcome was coded such that $O_t = 1$ if shape 1 was chosen and followed a good outcome (i.e. delivery of reward or absence of electric stimulation) or if shape 2 was chosen and followed by a bad outcome (i.e. absence of reward or delivery of electric stimulation). $O_t = 0$ codes for the opposite two cases in each task.

Hierarchical Bayesian Estimation of Parameters in Behavioral Model

Individual parameter components (e.g. $\alpha_{volatile-stable}$) were estimated using a hierarchical Bayesian approach, which estimates a group-level distribution (with a mean μ and a variance σ^2 across participants) for each parameter component (see Eqn. 6 for an example). These group-level distributions can help to prevent overfitting in individual participants, because parameter components that are not found to be important for explaining behavioral variation in a large number of participants would be estimated to have a group mean and variance near zero; this would effectively be like removing them from the model altogether.

The relationship between symptom factor scores (general, anxiety-specific, and depression-specific) and parameter components were modeled by allowing the mean of each group-level distribution to vary as a function of the three factors. The three factor scores are denoted by (X_g, X_d, X_a) in Eqn 6, where subscripts denote the general factor, depression-specific factor, and anxiety-specific factor, respectively. The strength of the linear relationship between the factor scores and the parameter components are given by the group-level regression coefficients $\{\beta_g, \beta_d, \beta_a\}$.

Eqn. 6

$$\alpha_{volatile-stable} \sim Normal(\mu + \beta_g X_g + \beta_d X_d + \beta_a X_a, \sigma^2)$$

The parameters were transformed to the appropriately constrained space before they were used in the behavioral model (Eqns. 2-5). A logistic transform was used for $\{\alpha, \gamma, \eta\}$ to bound them between $[0,1]$, and a log transform was used for r, τ, τ_k to constrain them to be positive.

The hyperpriors for these group-level parameters $(\mu, \beta_g, \beta_d, \beta_a)$ were uninformative $Normal(0,10)$, which is effectively uniform over the space of reasonable parameter values. The hyperpriors for the population variances, σ^2 , were $Cauchy(2.5)$.

Models were fit using PyMC3 (Salvatier et al., 2016), a Python Bayesian statistical modeling software package, which is similar to STAN. Hamiltonian Monte-Carlo was used to sample from the posterior. Four chains were run with 250 tuning steps and 1000 samples. Visual inspections of the traces as well as Gelman–Rubin statistics (\hat{R}) were used to assess convergence (Gelman & Rubin, 1992). There were no group-level parameters with \hat{R} values greater than 1.1 (most were below 1.01). There were only 8 out of the 2236 participant-level parameters (from two participants) with \hat{R} values greater than 1.1, and these were for η and τ_k , which were not the focus of the main analysis.

The marginal posterior distributions for the group-level parameters $(\mu, \beta_g, \beta_d, \beta_a)$ were used to assess the statistical significance. Group-level parameters with a 95% highest posterior density (HDI) intervals that did not contain zero were deemed statistically significant (see Patzelt et al., 2018 and Aylward et al., 2019 for similar treatment of posterior credible intervals).

Results

A validation of the general, anxiety-specific, and depression-specific factors from the bifactor analysis

Before analyzing participants' behavioral data in the two tasks, we first assessed whether the three factors estimated in our bifactor model were similar to factors that have been estimated in previous bifactor models, whether they captured differences in patient diagnoses, and whether they generalized to an independent participant sample.

In line with previous literature (Clark et al., 1994; Steer et al., 1995; Zinbarg & Barlow, 1996; Steer et al., 1998; Simms et al., 2008; Steer et al., 2008; Brodbeck et al., 2011), the general factor had moderate loadings (>0.2) for almost every symptom item, and it had high loadings (>0.4) both on anxiety-related items and on depression-related items. One specific factor had high loadings (>0.4) on questions related to anhedonia and depressed mood (in agreement with Clark et al., 1994; Steer et al., 1995; Steer et al., 2008). The other specific factor had high loadings (>0.4) on questions related to worry and anxiety (in agreement with Brodbeck et al., 2011).

Scores on the three latent factors (general, depression-specific and anxiety-specific) were calculated for each participant (**Figure 2.2**). It can be seen from **Figure 2.2** that the combined patient group (i.e. participants diagnosed with either MDD or GAD) tended to have higher scores on the general factor dimension (mean=0.85; median=0.76) than the healthy controls (mean=-0.50; median=-0.79) and the community sample (mean=-0.45; median=-0.76). Participants diagnosed with GAD also had a higher average score on the anxiety-specific dimension (mean=0.93; median=0.87) than patients diagnosed with MDD (mean=-0.26; median=-0.12; $t=2.5$, $p=0.02$ for difference). Conversely, participants diagnosed with MDD had a higher average score on the depression-specific factor (mean=0.63; median=0.68) than participants diagnosed with GAD (mean=-0.19; median=-0.19; $t=3.2$, $p=0.003$ for difference).

We tested the generalizability of this factor structure using an independent online dataset ($n=199$). Participants were students at UC Berkeley (120 females, mean age=20±4). This group was fairly distinct from our first sample, differing in nationality, being more homogenous in age and educational status and less homogenous in ethnicity, and not including individuals recruited to meet diagnosis for either GAD or MDD. Replicability of the factor structure across these two datasets is hence a strong test of its generalizability across the population. We first tested how well the factor structure fit in the UC Berkeley student sample. This resulted in a good fit as measured by the comparative fit index (CFI=0.962), the root mean square error of approximation (RMSEA=0.065), and the standardized root mean square residual (SRMR=0.11). Factor scores calculated for these participants using this structure can be seen in **Figure 2.2** to extend across a similar range to the in-lab participants' scores on each of the three dimensions. Next, we re-estimated the factor structure in just the student sample to see if a similar structure would emerge. Indeed, the new loadings were highly congruent with the original loadings (cosine-similarity=0.98 for the general factor loadings, 0.83 for the depression-specific factor loadings, and 0.91 for the anxiety specific factor loadings).

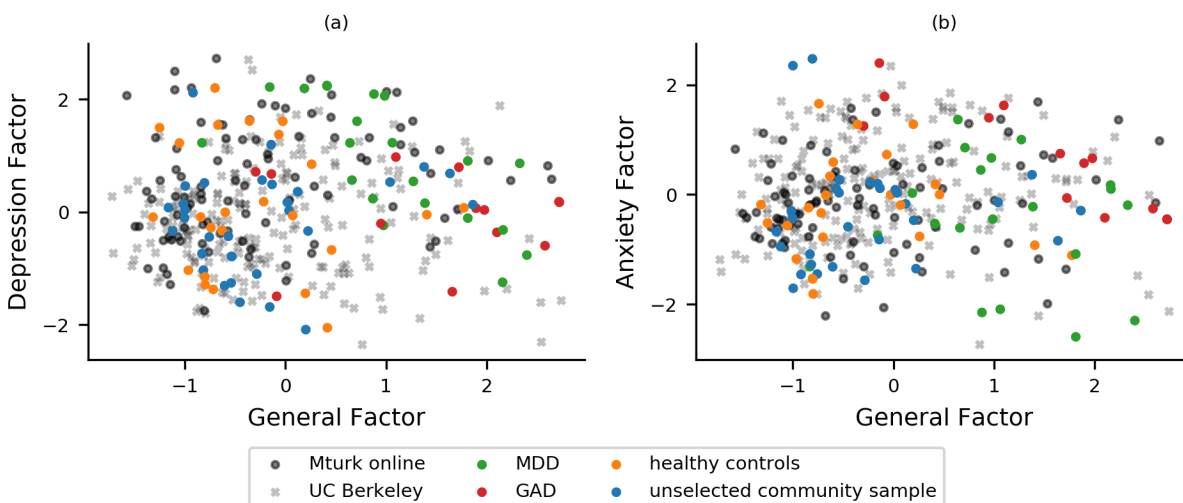


Figure 2.2: Bifactor analysis of self-report anxiety and depression symptoms. Eigenvalue decomposition was applied to the covariance matrix of 128 individual items from self-report standardized measures of anxiety and depressive symptomatology administered to in-lab participants ($n=86$). Eigenvalues from the data were compared to eigenvalues from a random normal matrix of equivalent size, providing evidence that there were three dimensions distinguishable from noise. A bifactor model with three factors, one general and two specific, was estimated from the covariance matrix. The factor structure was confirmed in an independent sample of online participants ($n=199$) from UC Berkeley, resulting in a good fit (Comparative Fit Index = 0.96). Item loadings from the bifactor model were used to calculate scores for individual participants for the general factor (x-axis in a and b), the depression-specific factor (y-axis in a), and the anxiety-specific factor (y-axis in b). These plots show that there is a similar range of scores between the in-lab and UC Berkeley datasets ($n=199$; denoted by x's). We also calculated scores for the online Mturk dataset ($n=147$), which had only a subset of the questionnaires and which was used in experiment 2 to test the replicability of the relationships between symptoms and decision-making. MDD=major depressive disorder; GAD=generalized anxiety disorder.

Learning rates differences by volatility level, absolute outcome valence, and relative outcome valence: Group-level findings (Experiment 1)

After fitting our reinforcement learning model to participants' choice behavior in both tasks, we first looked at whether learning rates differed on average across participants, between blocks (i.e. stable or volatile), tasks (i.e. reward or aversive), or trials categorized by relative outcome valence (i.e. trials following good or bad outcomes). **Figure 2.3** (left panel) shows the posterior means, along with their 95% highest posterior density intervals (HDI's), for the group-level average differences and interactions between conditions. The HDI's significantly excluded zero for differences in learning rates for block ($\alpha_{volatile-stable}$; posterior mean for $\mu=0.16$, 95%-HDI=[0.04,0.3]), task ($\alpha_{reward-aversive}$; $\mu=-0.2$, [-0.38,-0.04]) and relative outcome valence ($\alpha_{good-bad}$; $\mu=0.49$, [0.31,0.65]). This means that, on average, participants had higher learning rates during the volatile block than the stable block, higher learning rates during the aversive task than the reward task, and higher learning rates on trials following good versus bad outcomes. Higher learning rates for volatile than for stable blocks confirmed that, on average, participants were able to infer the level of volatility and correspondingly adjust learning rates in line with previous work (Behrens et al., 2007; Browning et al., 2015). There were no significant

interactions in the differences in learning rate between block, task, or relative outcome valence.

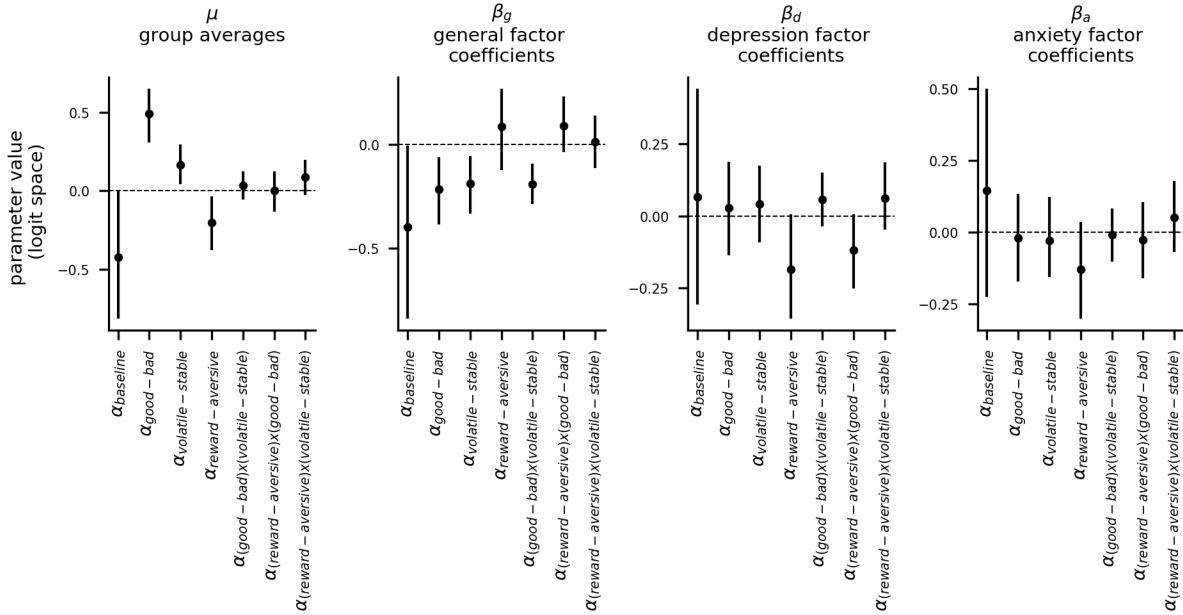


Figure 2.3: Group-level parameters (μ , β_g , β_d , β_a) for learning rate components (in-lab participant sample, $n=86$). For the group-level averages μ (first panel), the 95% highest posterior density intervals (HDI) significantly excluded zero for the learning rates differences for block $\alpha_{volatile-stable}$, task $\alpha_{reward-aversive}$, and relative outcome valence $\alpha_{good-bad}$. This means that participants had higher learning rates during the volatile block than the stable block, higher learning rates during the aversive task then the reward task, and higher learning rates on trials following good versus bad outcomes. For the general factor group-level regression coefficient β_g (second panel), the HDI excluded zero for learning rate differences for block $\alpha_{volatile-stable}$, relative outcome valence $\alpha_{good-bad}$, and the interaction of the two differences $\alpha_{(good-bad) \times (volatile-stable)}$. This meant that participants with low scores on the general factor (i.e. low levels of mood and anxiety symptoms) had higher learning rates in volatile than in stable blocks, following good versus bad outcomes, and highest learning rates following good outcomes in volatile blocks. There were no significant differences in learning rate associated with the depression-specific or anxiety-specific regression coefficients $\{\beta_d, \beta_a\}$ (third and fourth panels).

Learning rates differences by volatility level, absolute outcome valence, and relative outcome valence: Relationship to mood and anxiety symptoms (Experiment 1)

To address our first research question—that is, whether a lack of learning rate adjustment to volatility is related to both symptoms of anxiety and depression or to just symptoms of anxiety—we looked at whether the difference in learning rates between the volatile and stable blocks depended significantly on the general factor scores or on the anxiety-specific factor scores. For the learning rate difference between blocks, $\alpha_{volatile-stable}$, only the HDI for the general factor regression coefficient significantly excluded zero ($\beta_g=-0.19$, [-0.33, -0.05]). Neither the anxiety-specific factor coefficient ($\beta_a=0.03$, [-0.12,0.16]) nor the depression-

specific factor coefficient ($\beta_d = -0.05$, [-0.18, 0.08]) had a significant relationship. This relationship with the general factor scores meant that individuals with low scores (i.e. low overall levels of mood and anxiety symptoms) demonstrated higher learning rates during the volatile than stable block, whereas individuals with high scores did not, with their estimated $\alpha_{volatile-stable}$'s close to zero (see **Supplemental Figure 2.1** for individual participants' parameters). This result suggests that the ability to appropriately adjust learning rates to the level of volatility, previously linked to trait anxiety in Browning et al., (2015), likely underpins both anxiety and depression.

To address our second research question—that is, whether absolute outcome valence (reward or aversive) modulates the relationship between learning rate adjustment to volatility and symptoms of anxiety and depression—we looked at the interaction of differences in learning rates for task and volatility level, $\alpha_{(reward-aversive) \times (volatile-stable)}$. This was not significantly related to any of the three factor scores ($\beta_g = -0.01$ [-0.14, 0.11]; $\beta_a = -0.05$ [-0.18, 0.07]; $\beta_d = -0.06$ [-0.19, 0.05]). Furthermore, there were no significant differences in learning rates, on average, between the reward and aversive tasks related to any of the three factors ($\alpha_{reward-aversive}$; $\beta_a = -0.13$ [-0.3, 0.04]; $\beta_d = -0.19$ [-0.36, 0.01]; $\beta_g = 0.09$ [-0.12, 0.27]).

To address our third research question—that is, whether relative outcome valence (good or bad) modulates the relationship between learning rate adjustment to volatility and symptoms of anxiety and depression—we looked at the interaction of differences in learning rates for relative outcome valence and volatility level, $\alpha_{(good-bad) \times (volatile-stable)}$. This parameter component was significantly related to the general factor ($\beta_g = -0.19$ [-0.29, -0.09]), but not to the anxiety specific factor or to the depression specific factor. This interaction meant that individuals with varying general factor scores showed different degrees of learning rate adjustment to volatility on trials following good versus bad outcomes. We also looked at the difference in learning rates following good versus bad outcome, independently of volatility, i.e. $\alpha_{good-bad}$. Again, the general factor, but not the anxiety-specific or depression-specific factors, was significantly related to the learning rate difference for relative outcome valence ($\alpha_{good-bad}$; $\beta_g = -0.22$ [-0.38, -0.06]). This meant that individuals with low scores on the general factor also updated their action-outcome contingency estimates a greater extent following good outcomes, on average, across blocks.

In order to more easily visualize the combined effects on learning rate associated with the general factor, we calculated the expected mean learning rate for each condition for scores that were ± 1 standard deviation above or below the mean on the general factor. **Figure 2.4** (blue) clearly shows that low general factor scores are associated with higher learning rates during volatile than stable blocks and higher learning rates following good versus bad outcomes. Furthermore, it shows that the difference between the volatile and stable block was more pronounced on trials following good outcomes. In contrast, individuals with high scores on the general factor (**Figure 2.4** red) were associated with a lower baseline learning rate (the effect for $\alpha_{baseline}$ was indeed significant; $\beta_g = -0.4$ [-0.84, -0.01]) and smaller differences in learning rate between volatile and stable blocks and following good and bad outcomes.

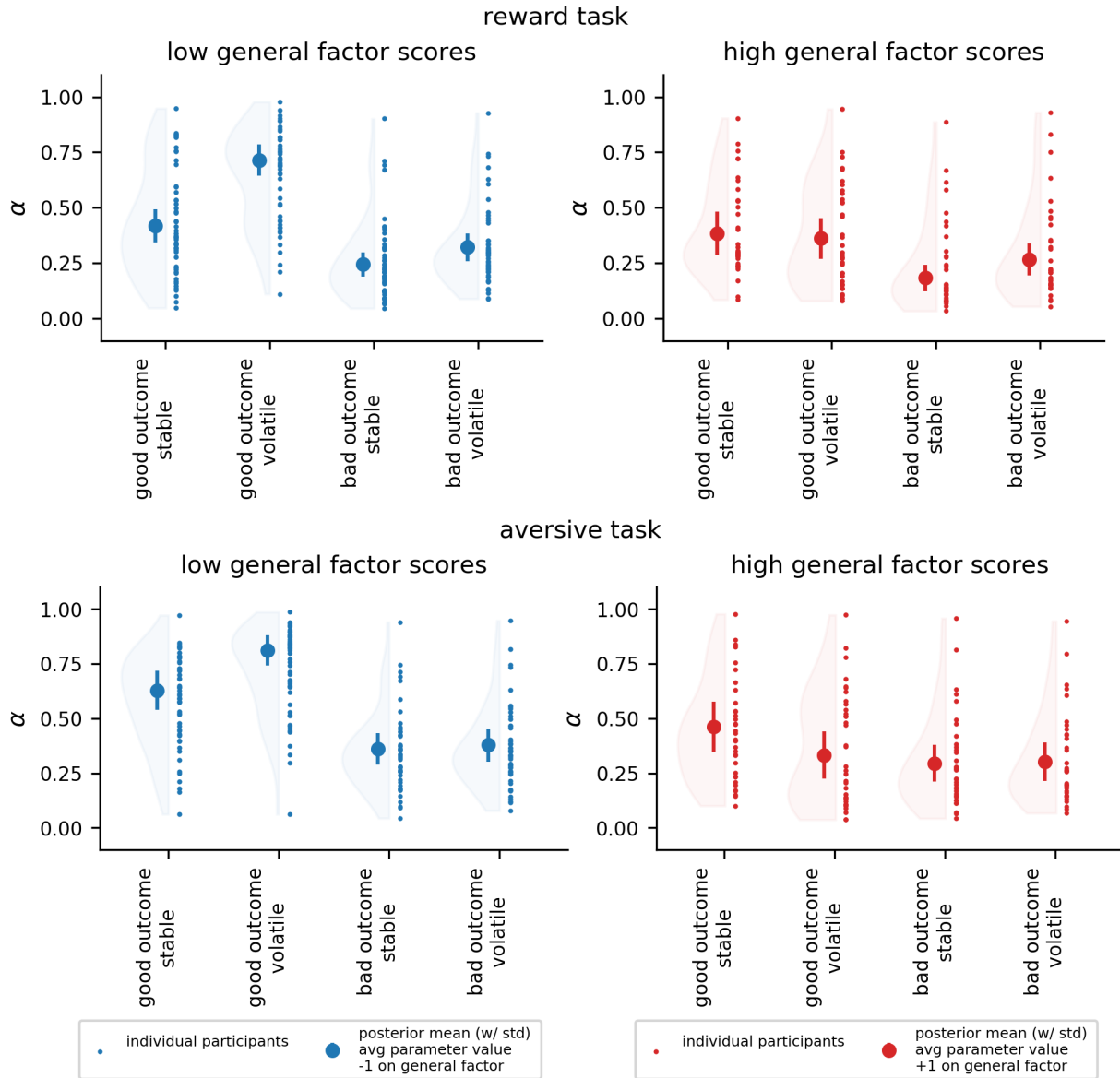


Figure 2.4. Expected mean learning rates per condition for participants with low (-1 std) or high (+1 std) scores on the general factor scores (in-lab participant sample, $n=86$). The expected mean learning rates for participants with ± 1 standard deviation above or below the mean on the general factor were calculated per condition (large data points). These expected means were calculated in the context of the hierarchical Bayesian model using the group-level μ 's and β_g 's relevant to each condition; for example, the negative estimated value for β_g for $\alpha_{volatile-stable}$ leads to a difference in learning rates between volatile and stable blocks observed for individuals with low general factor scores (blue) and also to the lack of difference for individuals with high general factor scores (red). The term 'expected mean' refers to the posterior expectation of the estimate of the group-level mean. Error bars represent the posterior standard deviation for these means (akin to the s.e.m.). Individual parameters for participants who were above or below the mean on the general factor are also plotted (small data points), along with their distributions, in order to visualize individual variability.

Learning rates differences by volatility level, absolute outcome valence, and relative outcome valence: Relationship to mood and anxiety symptoms (Experiment 2)

We fit the behavioral model from experiment 1 to an independent online sample of participants, in order to test whether the findings replicated in an online version of the reward-learning task and whether they generalized to an alternative aversive-learning task, with monetary loss in lieu of primary aversive outcomes. At a group level, learning rates were significantly higher for good versus bad outcomes, $\alpha_{good-bad}$ ($\mu=0.56$ [0.41,0.7]), but were not significantly different between volatile and stable blocks, $\alpha_{volatile-stable}$ ($\mu=-0.03$ [-0.15,0.1]) or between the gain and loss tasks, $\alpha_{gain-loss}$ ($\mu=0.14$ [-0.02,0.28]) (see **Figure 2.5** left panel).

Although there was no effect of volatility at the group level, there was a significant interaction of volatility by participant general factor scores. As in experiment 1, the difference in learning rates between volatile and stable blocks, $\alpha_{volatile-stable}$, had a significant negative relationship with the general factor scores ($\beta_g=-0.14$ [-0.3, -0.01]). In other words, individuals with low general factor scores again showed greater adjustment of learning rates to volatility (learning faster in the volatile block than the stable block), than participants with high general factor scores. Also as observed in experiment 1, neither anxiety nor depression specific scores were related to the adjustment to volatility ($\beta_a=0.03$ [-0.08,0.15]; $\beta_d=0.01$ [-0.13,0.14]).

Similarly to experiment 1 with regards to our second question, none of the three factors were significantly linked to differential learning in the reward gain versus reward loss condition ($\beta_g=-0.03$ [-0.15,0.21]; $\beta_a=-0.13$ [-0.28,0.02]; $\beta_d=0.0$ [-0.13,0.14]).

With regards to our third question, there was partial agreement between experiment 1 and experiment 2. Unlike in experiment 1, there was no association between general factor scores and learning for good versus bad outcomes, $\alpha_{good-bad}$ (i.e. a main effect independent of its interaction with volatility). As can be seen in **Figure 2.6** (red), participants with high scores on the general factor indeed showed higher learning rates following good outcomes. Nonetheless, there was still a significant association between the interaction of volatility and relative outcome valence, $\alpha_{(good-bad) \times (volatile-stable)}$ and scores on the general factor ($\beta_g=-0.13$ [-0.23, -0.01]). This meant that in both experiments, individuals with low general factor scores increased learning rates in volatile relative to stable conditions to a greater extent following good versus bad outcomes, and again, that this pattern was absent in individuals with high general factor scores (see **Figure 2.6**). Neither anxiety-specific nor depression specific scores were associated with this interaction ($\beta_d=-0.07$ [-0.17,0.03]; $\beta_a=-0.0$ [-0.09,0.09]), further supporting the observation from experiment 1 that a lack adjustment to volatility and its modulation by relative outcome valence is a shared feature of both anxiety and depression.

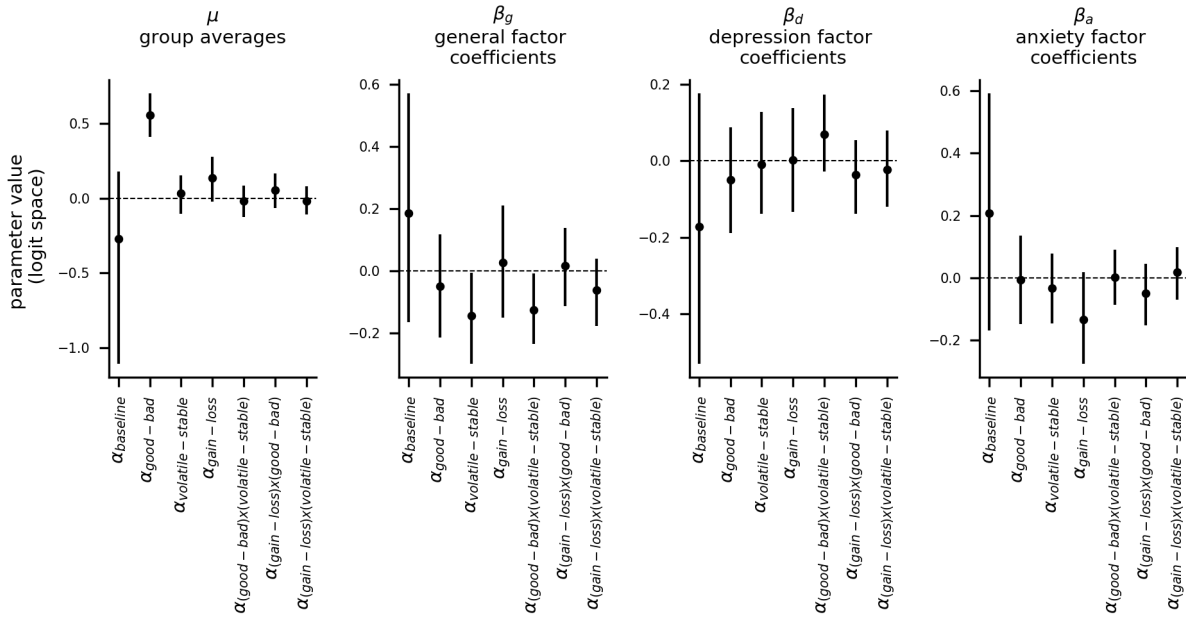


Figure 2.5: Group-level parameters (μ , β_g , β_d , β_a) for learning rate components (online participant sample, $n=147$). For the group-level averages μ (first panel), the 95% highest posterior density intervals (HDI) significantly excluded zero for the learning rates differences for relative outcome valence $\alpha_{good-bad}$, but not block $\alpha_{volatile-stable}$, or task $\alpha_{gain-loss}$. Similarly to experiment 1, the HDI for the general factor regression coefficient β_g (second panel) significantly excluded zero for learning rate differences for block $\alpha_{volatile-stable}$ and the interaction of block and relative outcome valence $\alpha_{(good-bad) \times (volatile-stable)}$. In contrast to experiment 1, the relationship between the general factor and the (main effect) difference between relative outcome valence $\alpha_{good-bad}$ was not significant. Similarly to experiment 1, there were no significant differences in learning rate associated with the depression-specific or anxiety-specific regression coefficients $\{\beta_d, \beta_a\}$ (third and fourth panels).

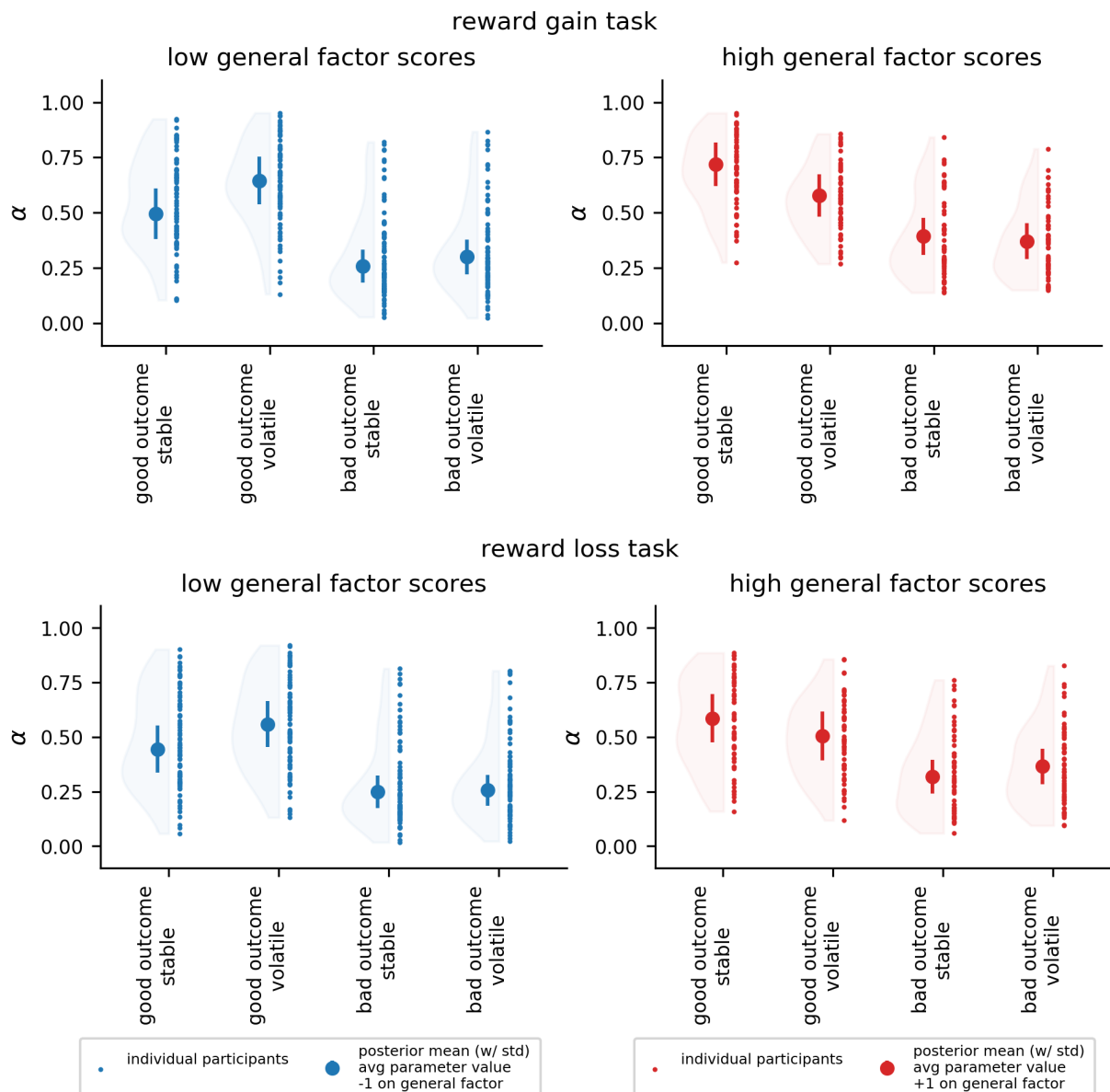


Figure 2.6. Expected mean learning rates per condition for participants with low (-1 std) or high (+1 std) scores on the general factor scores (online participant sample, $n=147$). The expected mean learning rates for participants with ± 1 standard deviation above or below the mean on the general factor were calculated per condition (large data points). Error bars represent the posterior standard deviation for these values (akin to the s.e.m.). Individual parameters for participants who were above or below the mean on the general factor are also plotted (small data points), along with their distributions, in order to visualize individual variability. The pattern of learning rates across conditions for online participants with low scores on the general factor is consistent with that found in experiment 1 for participants with general factor low scores. Online participants with high scores show higher baseline learning rates than participants in experiment 1, but similarly did not exhibit higher learning rates for the volatile versus the stable block or following good outcomes in the volatile block.

Other differences between experiment 1 and 2

As mentioned above, the behavioral model contained five other parameters to account for individual differences in other aspects of choice behavior. None of these other parameters (inverse temperature τ , mixture weight for probability versus magnitude γ , choice kernel update rate η , choice kernel inverse temperature τ_k , or the subjective magnitude parameter r) showed consistent across-dataset relationships to any of the three factors. However, there are a few significant relationships to parameters within one or the other dataset. These results are detailed in the supplemental materials along with other group-level findings (see **Supplemental Results**).

Discussion

In this study, we examined individual differences in learning across individuals with varying levels of anxiety and depression. Specifically, we examined individual differences in learning rates as a function of volatility level (stable or volatile blocks), absolute outcome valence (reward or aversive tasks), and relative outcome valence (good or bad outcomes). We used three latent factors to differentiate whether individual differences in learning rates were specifically associated with anxiety, specifically associated with depression, or associated with both.

Consistently across both experiment 1 (in-lab participants, including patients) and experiment 2 (online participants), we observed that individuals with low levels of overall mood and anxiety symptoms (i.e. low general factor scores) adjusted their learning rates to be higher in volatile than stable conditions, whereas individuals with high levels of overall symptoms did not. Neither the anxiety-specific nor the depression-specific factors were related to this adjustment, suggesting that the deficit is a general feature of both anxiety and depression, rather than being specific to anxiety, as might be inferred from Browning et al. (2015). We also observed, consistently across both experiments, that individuals with low levels of overall symptoms had larger learning rate adjustments to volatility following good outcomes. We discuss the implications of this below. Contrary to our expectations, we did not observe any baseline differences in learning rates, in individuals with high levels of anxious or depressive symptoms, related to absolute outcome valence. Moreover, we did not observe any interaction of absolute outcome valence with the modulation of learning rate by volatility or by relative outcome valence, meaning that the learning rate alterations in individuals with high levels of overall symptoms were consistent across reward and aversive contexts.

That a deficit in learning rate adjustment to volatility is associated with both anxiety and depression broadens our understanding of the range of those who may have difficulty making decisions in noisy and volatile real-world situations. Learning rate adjustment to volatility can more accurately align an individual's internal estimate of action-outcome contingencies with the true underlying reality, which can lead to better decisions and outcomes in the long run. This may be relevant to a broad range of real-world situations, such as interpersonal settings, in which you might need to decide whether an unexpected outcome reflects a one-off occurrence or a fundamental change in your relationship. Volatility is also likely to occur in a vocational

setting, in which your recent performance might sometimes reflect luck (good or bad) and sometimes reflect a true change in your effectiveness. Given the ubiquity of volatility in the real-world, targeting behavior in this type of setting might be useful as part of a transdiagnostic treatment for both anxiety and depression. In addition, experimentally measuring learning rate adjustment to volatility might be a useful complement to the self-report measures, such as the intolerance to uncertainty scale (IUS), which are being investigated for use as transdiagnostic markers for treatment success (Boswell et al., 2013).

As observed in our experiment, a lack of learning rate adjustment to volatility may be associated with different average learning rates across individuals. In experiment 1, we observed that individuals with high levels of mood and anxiety symptoms (individuals with high general factor scores mainly consisting of patients) had significantly lower baseline learning rates across all conditions than the controls, whereas in experiment 2, similarly high-symptomatic individuals did not significantly differ in baseline learning rates relative to low-symptomatic individuals. This difference between experiments could potentially reflect different assumptions that individuals resort to when they are unable to accurately infer the level of volatility. High-symptomatic individuals in experiment 1, on average, may have erroneously assumed that things were always stable, which could be related to individual differences in attributional style (i.e. whether someone believes in more global and stable causes for negative outcomes; Abramson et al., 1978). On the other hand, high-symptomatic participants online may have been more heterogeneous in their assumptions, with some participants erroneously assuming things were always stable and others erroneously assuming things were always volatile. Evidence that anxiety is associated with treating noisy contexts as more volatile has been previously observed (Huang et al., 2017). It would be interesting for future work to explore the factors or other individual differences that could lead to one default assumption over the other.

In both experiment 1 and experiment 2, we also observed that individuals with low levels of mood and anxiety symptoms learned more following good outcomes than bad outcomes, specifically in volatile contexts. In our experimental task, this behavior does not confer the same advantage as learning rate adjustment to volatility, because bad outcomes were equally informative as to whether contingencies had changed or not. However, this asymmetric learning for good and bad outcomes may reflect a more general attitude adopted by individuals with low levels of mood and anxiety symptoms in situations characterized by other types of second order uncertainty (of which volatility is one type). For example, higher learning rates for positive versus negative prediction errors has been previously used to explain differences in risk preferences, because higher positive than negative learning rates tend to lead to the overestimation of an action's value when outcomes are noisy versus certain, leading to risk seeking behavior (Mihatsch & Neuneier, 2002; Niv et al., 2012). Similarly in our model, a higher learning rate for good outcomes would lead to an exaggerated difference in outcome probability between the two shape's (e.g. estimates of $p=0.9$ and $p=0.1$ versus the true $p=0.8$ and $p=0.2$ for shape 1 and shape 2). This might lead to something like an optimism for the current course of action. Whether advantageous or not, it would be interesting for future work to explore the relationship between asymmetric learning and volatility, and whether a lack of asymmetry is similarly associated with mood and anxiety symptoms in other learning contexts containing second order uncertainty.

Regarding our second research question, we did not observe different learning rates or a change in the learning rate adjustment to volatility associated with the absolute outcome valence (i.e. rewards, primary aversive outcomes, or financial losses) – i.e., implying that we found no evidence that depression is more associated with the absence of rewards and anxiety, the presence of punishments. This null result is contrary to some proposals that depression is associated with biased learning from feedback (Elliot et al., 1997; Steele et al., 2007), however, it is in line with a recent review that cites mixed evidence for the relationship between anhedonia and differences in reward versus punishment learning (Robinson & Chase, 2017).

This null result reemphasizes the need to identify factors that truly differentiate anxiety and depression. We propose that one differentiating factor may be related to decisions involving effort (Bishop & Gagne 2017; **Chapter 1.2**). Specifically, a reduction in the willingness to exert effort to pursue rewards may be more strongly associated with depression, whereas an increase in willingness to exert effort to avoid punishment may be more strongly associated with anxiety.

Two previous studies that have linked trait anxiety to learning rate adjustment to volatility (Browning et al., 2015; Pulcu et al., 2017) observed a significant correlation in the context of aversive outcomes (electric stimulation or financial loss, respectively) and a non-significant trend in the context of rewards. Using a larger sample size in the current study, we confirmed these trends, providing evidence that the link between the learning rate deficit and anxiety (and now also depression) extends to reward contexts. Nonetheless, it would be useful to further investigate whether a deficit in learning rate adjustment to volatility is amplified in certain contexts and what additional factors might drive that amplification.

In summary, we showed that a deficit in learning rate adjustment to volatility is related more broadly to internalizing symptomology (i.e. anxiety and depression), rather than being specifically related to anxious symptomology. This has different implications for how and in whom this deficit might impact decision making under volatility. We also observed that individuals, specifically those with low levels of internalizing symptoms, adjust their learning rate to volatility more following good rather than bad outcomes within both the context of rewards and aversive outcomes. Further work is needed to determine whether this represents an adaptive strategy, like learning rate adjustment to volatility on its own can be seen to be, or whether it carries over from other individual differences that covary with mood and anxiety symptoms.

Supplementary Results

Model Comparison: Alternative Behavioral Models

We compared fourteen alternative models, including the one from the main text (the full list of models can be found in **Supplemental Table 2.2**). Since the model space was too large to search exhaustively, we performed a series of pairwise comparisons between similar models. For each comparison, we chose the model with the best approximate leave-one-out cross-validation (LOO) (Vehtari et al., 2017).

We started with a Model #1 that was similar to the one used in Browning et al. (2015), except that it was fit using the Bayesian Hierarchical framework. Model #1 calculated the estimated outcome probability in the same way as the main model (Eqn. 2 in main text). In Model #1, the outcome probability estimate was then adjusted to account for differences in risk preference (using Eqn S1). Following that, the expected value for each shape was calculated, multiplying the outcome probability and outcome magnitude, before calculating the difference in expected value between shapes (all using Eqn. S2).

Eqn. S1

$$p'_t = \min(\max((p_t - 0.5)^\lambda + 0.5, 0), 1)$$

Eqn. S2

$$v_t = p'_t M1 - (1 - p'_t) M2$$

We first compared Model #1 with Model #2, which calculated the differences in magnitude and probability separately for each shape and then combined them as a mixture (using Eqn. 3 in the main text; excluding r). Each model was composed of three parameters, each of which was divided into four components: a baseline, a difference between volatile and stable blocks, a difference between the reward and aversive task, and the interaction of block and task. The λ parameter in both models determined the relative weight given to outcome magnitude and outcome probability, yet in different ways. LOO was substantially improved for the Model #2 (Model #1-Model #2; dLOO=2055).

Next, Model #2 was compared to Model #3, which additionally divided learning rate α to allow for differences on trials following good versus bad outcomes (i.e. outcome valence) and for the interaction of that difference with task and block. Model #3 substantially improved LOO (Model #2-Model #3; dLOO=617). We also compared Model #4, which broke down inverse temperature τ and the mixture weight γ by outcome valence instead of learning rate. Model #4 also improved LOO, albeit to a lesser extent than Model #3 (Model #2-Model #4; dLOO=122). Combining the two models together, that is, breaking down all three parameters by outcome valence, achieved even better LOO (Model #2-Model #5; dLOO=704). Adding the triple interaction between block, task, and outcome valence, however, slightly worsened the fit (Model #5-Model #6; dLOO=-27).

Next, we tried Model #7 that non-linearly transformed the difference in magnitudes between the shapes (using Eqn. 3 in the main text; including r). r was divided into a baseline

and a difference between task. LOO was substantially improved by adding this non-linearity (Model #5 vs Model #7; dLOO=295).

Next, we tried Model #11 included a choice kernel to account for people's tendency to repeat past choices regardless of outcomes (Lau & Glimcher, 2005; Ito & Doya, 2009; Akaishi et al., 2014; calculated using Eqn. 4 in the main text). The choice kernel update rate η was composed of a single baseline, while the choice kernel inverse temperature τ_k was divided into a baseline and a difference between tasks. This model improved LOO over a model that was equivalent in all respects except for the inclusion of the choice kernel (Model #7 vs Model #11; dLOO=103).

Next, we tried including a lapse parameter ϵ , which can account for mistakes that participants make that do not reflect the value differences between the shapes. ϵ could be added to any model by transforming the choice probability using Eqn S3. Neither adding ϵ to Model #7 (without choice kernel) and Model #11 (with choice kernel), improved fit (Model #11#-Model #12; dLOO=-11) or the model without it (Model #7-Model #8; dLOO=-8).

Eqn. S3

$$P'(C_t) = (1 - \epsilon) P(C_t) + \epsilon/2$$

For the next few comparisons, we tried alternative forms for the estimation of outcome probabilities (Eqn. 2 in the main text). The primary behavioral model assumes that participants have a single estimate for the probability that shape 1 and not shape 2 is followed by the outcome. Alternatively, participants could have two separate stimulus-specific probability estimates, one for shape 1 and one for shape 2. Using two separate estimates would be akin to using Q-values (Li & Daw, 2011; Mkrtchian et al., 2017; Aylward et al., 2019). We tried two versions of a model that update stimulus-specific probability estimates.

In the first version, the outcome probabilities for each shape were updated using Eqn S4 substituted for Eqn 2 in the main text. After the update, both probability estimates were decayed towards 0.5 using Eqn S5. The decay parameter $\delta \in [0,1]$ was composed of a baseline and a difference between task. This version of the model did not improve LOO either with a choice kernel (Model #11-Model #13; dLOO=-193), or without a choice kernel (Model #7-Model #9; dLOO=-220).

Eqn. S4

$$\begin{aligned} p_t &= p_{t-1} + \alpha(O_{t-1} - p_{t-1}) \text{ if participant chose shape 1} \\ q_t &= q_{t-1} + \alpha(O_{t-1} - q_{t-1}) \text{ if participant chose shape 2} \end{aligned}$$

Eqn. S5

$$\begin{aligned} p'_t &= (1 - \delta)p_t + (\delta) * 0.5 \\ q'_t &= (1 - \delta)q_t + (\delta) * 0.5 \end{aligned}$$

The second version used a Beta-Bernoulli Bayesian model to calculate the outcome probability estimates for each shape. Outcome probabilities were estimated by keeping track of the counts of good or bad outcomes that occurred following the choice of each shape (counts were updated using Eqn S6). The outcome probabilities for each shape were calculated as the

mean of the Beta distribution using Eqn S7. These probabilities were also decayed towards 50% using Eqn S5. This second version did not improve LOO with (Model #7-Model #10; dLOO=-165) or without the choice kernel (Model #11-Model #14; dLOO=-149).

Eqn. S6

$$\begin{aligned} a_t &= \delta a_{t-1} + \alpha && \text{If participant chose shape 1 and received good outcome} \\ b_t &= \delta b_{t-1} + \alpha && \text{If participant chose shape 1 and received bad outcome} \\ c_t &= \delta c_{t-1} + \alpha && \text{If participant chose shape 2 and received good outcome} \\ d_t &= \delta d_{t-1} + \alpha && \text{If participant chose shape 2 and received bad outcome} \end{aligned}$$

Eqn S7.

$$p_t = \frac{a_t + 1}{a_t + b_t + 2} \quad q_t = \frac{c_t + 1}{c_t + d_t + 2}$$

All fourteen alternative models were estimated hierarchically with the three symptom factors entered into the group-level distribution for each parameter component. We also fit the winning model (i.e. the one described in the main text) without symptom factors (μ only; LOO=24,647), with just the general factor (μ and β_g ; LOO=24,668), and with just the significant learning rate effects for the general factor (μ for all parameters and β_g for the four significant learning rate effects; LOO=24,640). We compared these alternative models to the one used in the main text (LOO=24,681). The best LOO was obtained by the model that just included the significant learning rate effects for general factor (dLOO=-7, compared to second best model).

Model Parameter Recovery

Parameter recovery was performed to check that the parameters were both identifiable and estimable (i.e. identifiable given the level of noise in our dataset). The posterior means for each participant's parameter components (e.g. $\alpha_{baseline}$, $\alpha_{volatile-stable}$ etc.) were used to simulate new choice data from the main behavioral model. Each simulated dataset had 86 new participants. The simulated datasets were fit with the same model to estimate new parameters. The original parameters from the actual dataset (referred to as the 'generative' parameters) were correlated with the newly estimated parameters (referred to as 'recovered' parameters) for each simulated dataset. An example for one simulated dataset can be seen in **Supplemental Figure 2.3** for learning rates and in **Supplemental Figure 2.5** for the other parameters. This procedure was repeated with 10 simulated datasets. The mean correlation across simulated datasets and parameters was $r=0.76$ (std=0.15). The average correlation was slightly higher for learning rate parameters ($r=0.88$, std=0.13). We also looked at our ability to identify and estimate the four types of group-level parameters (μ , β_g , β_a , β_d). In **Supplemental Figure 2.4** and in **Supplemental Figure 2.6**, it can be seen that the means and ranges of the recovered parameters were very similar to the original estimates. This provides good evidence that the symptom-parameter relationships can also be reliably estimated.

Model Reproduction of the Number of Switch Trials

The number of trials on which a participant stays with the same choice or switches to the other choice is a basic qualitative feature of the data that our model should be able to reproduce (even though it was not optimized to do so). In each of the simulated datasets described in the previous section, we counted the number of trials on which each simulated participant switched from one choice to the other. The number of switch trials made by the simulated participants and those made by the actual participants were extremely correlated ($r's > 0.88$ for all conditions and datasets), demonstrating that the model can indeed reproduce basic features of the data. The correlation between switch trials for each participant (simulated and real) in one of the simulated datasets can be seen in **Supplemental Figure 2.7**.

Additional Group Average Differences between Tasks

In addition to learning rate differences discussed in the main text, there were a few other notable differences in participant behavior between tasks. On average, participants in experiment 1 tended to use probability more than magnitude ($\gamma_{reward-aversive}; \mu=0.27$ [0.03,0.51]) and make choices that were more deterministically based on outcome value ($\tau_{reward-aversive}; \mu=0.26$ [0.05,0.49]) during the reward task compared to the aversive task. Participants, on average, also tended not to repeat their previous choice as much in the aversive task compared to the reward task (indexed by the choice kernel inverse temperature, $\tau_{k\ reward-aversive}; \mu=-0.52$ [-0.93,-0.17]). These differences suggest that participants may have had an easier time estimating and employing outcome probability when learning about rewards rather than aversive outcomes.

In experiment 2, online participants, on average, also tended to use probability more than magnitude during the gain task relative to the loss task ($\gamma_{gain-loss} \mu=0.31$ [0.06,0.59]). However, they made choices that were less deterministic as a function of total outcome value during the gain task ($\tau_{gain-loss} \mu=-0.16$ [-0.25,-0.09]). Therefore, it was less clear in experiment 2, whether participants had an easier time learning and utilizing outcome probabilities in the context of gains or losses. None of these parameter differences varied as a function of the general, anxiety-specific or depression-specific factor scores.

Additional Relationships between the General Factor and Model Parameters

In experiment 1, there was a significant dependence on the general factor for the difference in mixture weight on trials following good versus bad outcomes $\gamma_{good-bad}$ ($\beta_g=-0.39$ [0.78,-0.0]), meaning that participants with low scores on the general factor tended to use probability slightly less following good outcomes (see **Supplemental Figure 2.8**). However, this effect, which counteracts the effect on learning rate $\alpha_{good-bad}$ (i.e. learning more following good outcomes), only slightly reduced model fit when removed (dLOO=-30, se=26), compared to removing the same effect for learning rate (dLOO=-112, se=40). This effect was also not replicated in experiment 2 (also shown in **Supplemental Figure 2.8**).

In experiment 1, there was also a significant dependence on the general factor for the interaction between task and outcome valence for inverse temperature ($\tau_{(reward-aversive) \times (good-bad)}$; $\beta_g = -0.36$ [-0.71, -0.03]). This was because participants with high scores on the general factor had particularly low inverse temperatures on the trials following bad outcomes in the pain task (see **Supplemental Figure 2.9**). However, this effect was not observed in the loss task in experiment 2 and removing outcome valence modulations on inverse temperature had a very small change on overall model fit (dLOO=-7, se=22).

In experiment 1, the general factor was associated with a lower baseline mixture weight, $\gamma_{baseline}$, and also a lower baseline inverse temperature, $\tau_{baseline}$ ($\beta_g = -0.36$ [0.71, -0.03]; $\beta_g = -0.39$ [-0.78, -0.0]). A lower baseline inverse temperature can indicate a deficient model fit, which along with a lower mixture weight (i.e. relying mostly on outcome magnitude instead of outcome probability), can make it difficult to estimate a participant's learning rates. We performed two checks to reduce concerns that these effects were resulting in the low learning rates observed in individuals with high scores on the general factor (and therefore driving the observed associations between learning rate adaption to volatility and the general factor scores). First, we showed that model fit (which was easier to interpret than inverse temperature), mixture weight, and baseline learning rate were not strongly correlated ($r=0.24$, $r=0.31$, $r=0.33$). Moreover, as can be seen in **Supplemental Figure 2.12**, there were a number of participants with high scores on the general factor who seemed to have genuinely lower baseline learning rates—i.e., they also had good model fit and high values for the mixture weight. Secondly, we refit the model after removing the 11 participants who fell below 66.5% correctly predicted choices (a threshold at a clear break in the histogram; chance=50%) and the 11 more participants that had a baseline mixture weight less than -2 (a threshold corresponding to less than 12% reliance on probability versus magnitude and that was at another clear break in the histogram). In the model fit on this subset of participants, the learning rate difference between volatile and stable blocks $\alpha_{volatile-stable}$ and the interaction of that difference with outcome valence $\alpha_{(good-bad) \times (volatile-stable)}$ still had significant correlations with the general factor.

In experiment 2, the general factor was associated with the interaction of block and outcome valence $\gamma_{(good-bad) \times (volatile-stable)}$ ($\beta_g = 0.11$ [0.0, 0.21]) and the interaction of task and outcome valence for the mixture weight $\gamma_{(reward-aversive) \times (good-bad)}$ ($\beta_g = 0.12$ [-0.23, -0.0]) (these can be seen in **Supplemental Figure 2.8**). These effects were not observed in experiment 1 and removing modulations of mixture weight by outcome valence had a smaller change on overall model fit (dLOO=-60, se=30) than removing the same modulations from learning rate (dLOO=-126, se=31).

Additional Relationships between the Anxiety-Specific and Depression-Specific Factors and Model Parameters

In experiment 1, the depression specific factor was associated with an interaction between task and block for the inverse temperature parameter $\tau_{(reward-aversive) \times (volatile-stable)}$ ($\beta_d = -0.12$ [-0.22, -0.03]) (see **Supplemental Figure 2.10**).

However, this effect was not observed in experiment 2. Moreover, removing this interaction had a very little change on overall model fit (dLOO=-15, se=20).

In both experiment 1 and experiment 2, the anxiety specific factor scores were associated with a negative difference between the volatile and stable blocks for the mixture weight parameter $\gamma_{(volatile-stable)}$ ($\beta_a=-0.12$ [-0.23,-0.01]; $\beta_a-0.14$ [-0.28,-0.0]). However, the difference was driven by qualitatively different relationships in each experiment. In experiment 1, high anxiety scoring participants had a higher reliance on probability during all blocks, except the volatile blocks in the aversive-learning task, and low scoring participants had relatively flat mixture weights across all conditions within each task (see **Supplemental Figure 2.11**). In experiment 2, high anxiety scoring participants had relatively flat mixture weights across conditions and low anxiety scoring participants had higher mixture weights for the volatile blocks (also see **Supplemental Figure 2.11**).

Supplemental Table 2.1: Demographics

Participant Group	<i>MDD (in-lab)</i>	<i>GAD (in-lab)</i>	<i>Healthy Controls (in-lab)</i>	<i>Unselected Community Sample (in-lab)</i>	<i>Mturk Online (experiment 2)</i>	<i>UC Berkeley online (bifactor model validation)</i>
Female (total N)	10 (20)	11 (12)	16 (26)	14 (30)	65 (147)	120 (199)
Age Mean±SD [Min, Max]	31±10 [20,51]	32±9 [19,45]	27±6 [20,46]	27±5 [18,40]	Not recorded	21±4 [18,53]
STAI	59±6 [48,73]	58±9 [40,74]	40±12 [20,63]	36±12 [20,64]	43±13 [20,76]	44±9 [23,73]
BDI	24±9 [3,42]	20±11 [5,42]	7±7 [0,31]	6±8 [0,28]	11±11 [0,44]	7±6 [0,29]
MASQ-AD	80±10 [55,96]	74±16 [35,93]	55±18 [26,84]	53±19 [27,91]	64±18 [22,103]	54±15 [24,97]
MASQ-AS	28±7 [17,43]	33±10 [17,53]	21±4 [17,37]	22±6 [17,37]	22±7 [17,50]	24±8 [4,57]
General Factor	1.1±0.8 [-0.8,2.4]	1.3±1.0 [-0.3,2.7]	-0.3±0.8 [-1.3,1.8]	-0.3±0.8 [-1.2,1.9]	-0.1±1.0 [-1.6,3.2]	-0.1±0.9 [-1.7,2.7]
Depression Specific Factor	0.8±1.0 [-1.2,2.2]	-0.1±0.8 [-1.5,1.0]	0.1±1.1 [-2.1,2.2]	-0.2±0.9 [-2.1,2.1]	0.4±1.0 [-2.0,2.7]	-0.2±0.9 [-2.4,2.7]
Anxiety Specific Factor	-0.5±1.1 [-2.6,1.4]	0.8±0.9 [-0.4,2.4]	-0.2±0.9 [-1.8,1.7]	-0.4±1.0 [-1.7,2.5]	-0.1±1.0 [-2.7,2.9]	0.1±1.0 [-2.7,2.3]

Supplemental Table 2.2: Model Comparison

14 total models were fit to the behavioral data. Models were fit hierarchically and compared using approximate leave-one-out cross validation (LOO). The model with the lowest LOO was selected for the main analyses.

Model (Number) and Name	Parameters	# of Parameter Components	LOO
(#1) Expected value (most similar to Browning 2015)	α, r, τ^1	12	27,838
(#2) Probability difference and magnitude difference	α, γ, τ	12	25,783
(#3) Probability difference and magnitude difference, good/bad learning rates	$\alpha^{gb}, \gamma, \tau$	15	25,166
(#4) Probability difference and magnitude difference, good/bad inverse temperature and mixture weight	$\alpha, \gamma^{gb}, \tau^{gb}$	18	25,661
(#5) Probability difference and magnitude difference, good/bad all three	$\alpha^{gb}, \gamma^{gb}, \tau^{gb}$	21	25,079
(#6) Probability difference and magnitude difference, good/bad all three and triple interaction	$\alpha^{gb}, \gamma^{gb}, \tau^{gb}$	24	25,106
(#7) Probability difference and nonlinear magnitude difference, good/bad all three	$\alpha^{gb}, \gamma^{gb}, \tau^{gb},$ $r^{rp \text{ only}}$	23	24,784
(#8) Probability difference and nonlinear magnitude difference, good/bad all three, lapse parameter	$\alpha^{gb}, \gamma^{gb}, \tau^{gb},$ $r^{rp \text{ only}}, \epsilon^{rp \text{ only}}$	25	24,792
(#9) Stimulus-specific probability difference and nonlinear magnitude difference, good/bad all three	$\alpha^{gb}, \gamma^{gb}, \tau^{gb},$ $r^{rp \text{ only}}, \delta$	27	25,004

(#10) Stimulus-specific Bayesian probability difference and nonlinear magnitude difference, good/bad all three	$\alpha^{gb}, \gamma^{gb}, \tau^{gb},$ $r^{rp \text{ only}}, \delta$	27	24,949
(#11) Probability difference and nonlinear magnitude difference, good/bad all three, choice kernel (primary model used in main text)	$\alpha^{gb}, \gamma^{gb}, \tau_v^{gb},$ $r^{rp \text{ only}}, \tau_k^{rp \text{ only}}, \eta^{ba}$	26	24,681
(#12) Probability difference and nonlinear magnitude difference, good/bad all three, choice kernel, lapse parameter	$\alpha^{gb}, \gamma^{gb}, \tau_v^{gb},$ $r^{rp \text{ only}}, \tau_k^{rp \text{ only}}, \eta^{ba}$ $, \epsilon^{rp \text{ only}}$	28	24,692
(#13) Stimulus-specific probability difference and nonlinear magnitude difference, good/bad all three, choice kernel	$\alpha^{gb}, \gamma^{gb}, \tau_v^{gb},$ $r^{rp \text{ only}}, \tau_k^{rp \text{ only}}, \eta^{ba}$ δ	32	24,874
(#14) Stimulus-specific Bayesian probability and magnitude difference, good/bad all three, choice kernel	$\alpha^{gb}, \gamma^{gb}, \tau_v^{gb},$ $r^{rp \text{ only}}, \tau_k^{rp \text{ only}}, \eta^{ba}$ δ	32	24,830

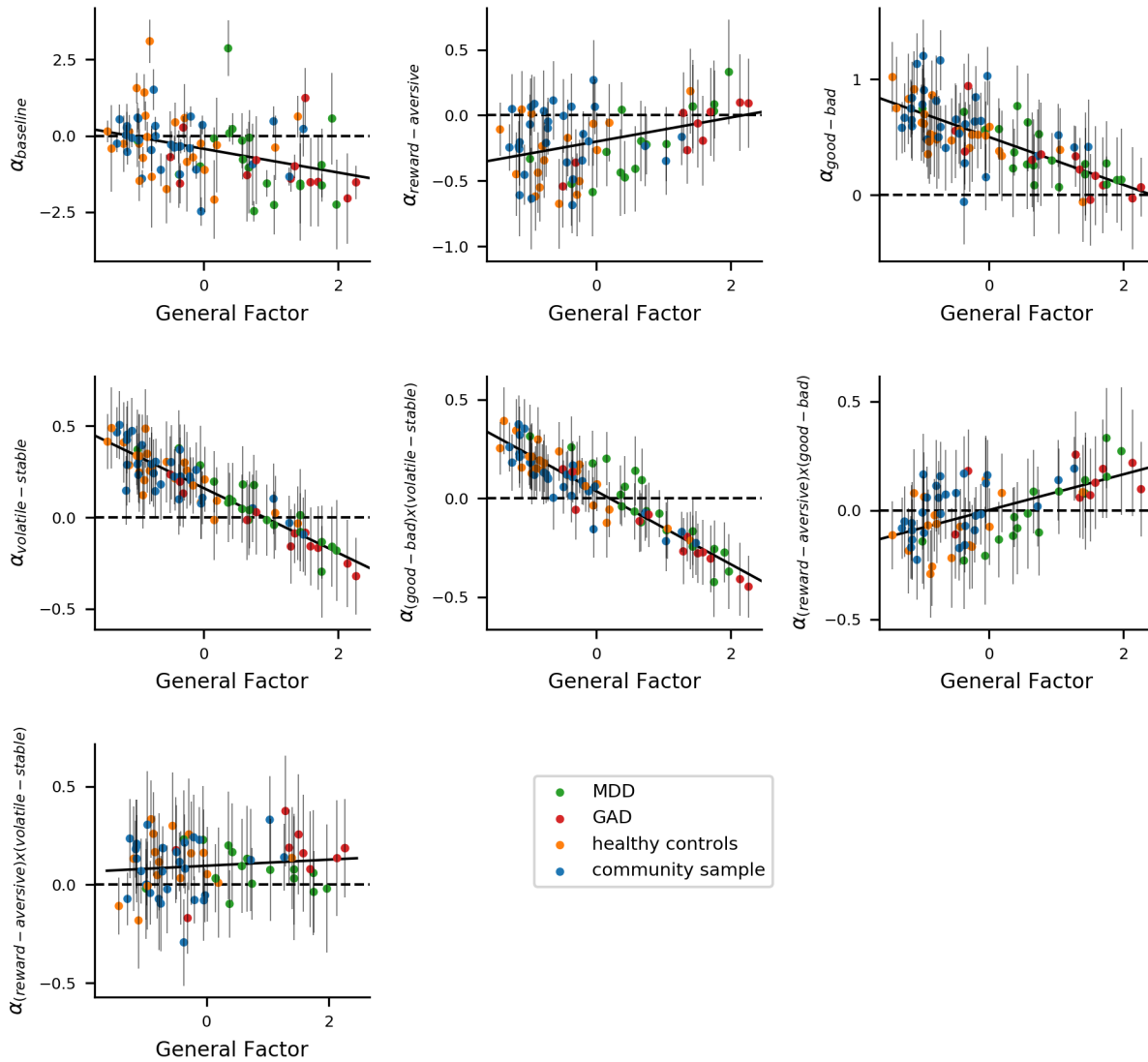
Note.

¹: Unless otherwise stated, each parameter is divided into 4 parameter components: a shared baseline parameter across blocks and tasks, and differences in the parameter between stable and volatile blocks (volatile-stable), between the reward and aversive tasks (reward-aversive) and an interaction of those differences (reward-aversive)x(volatile-stable).

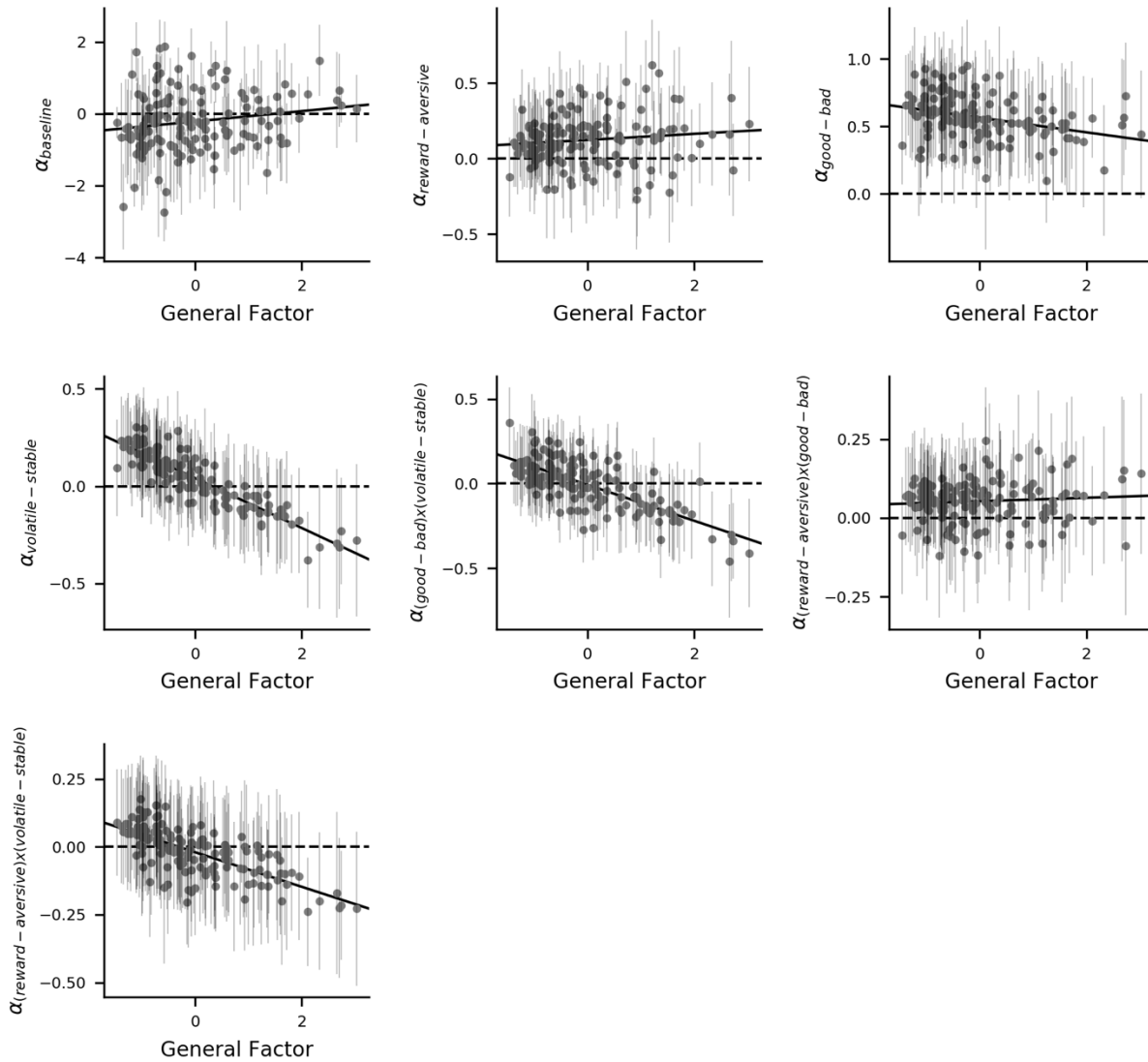
^{gb}: For each parameter with this superscript, three additional parameter components were added for difference in the parameter on trials following good versus bad outcomes (good-bad) and the interactions of this difference with block type (volatile-stable)x(good-bad), and task (reward-aversive)x(good-bad).

^{rp only}: For each parameter with this superscript, only differences in the parameters between the reward and aversive task versions were included (reward-aversive).

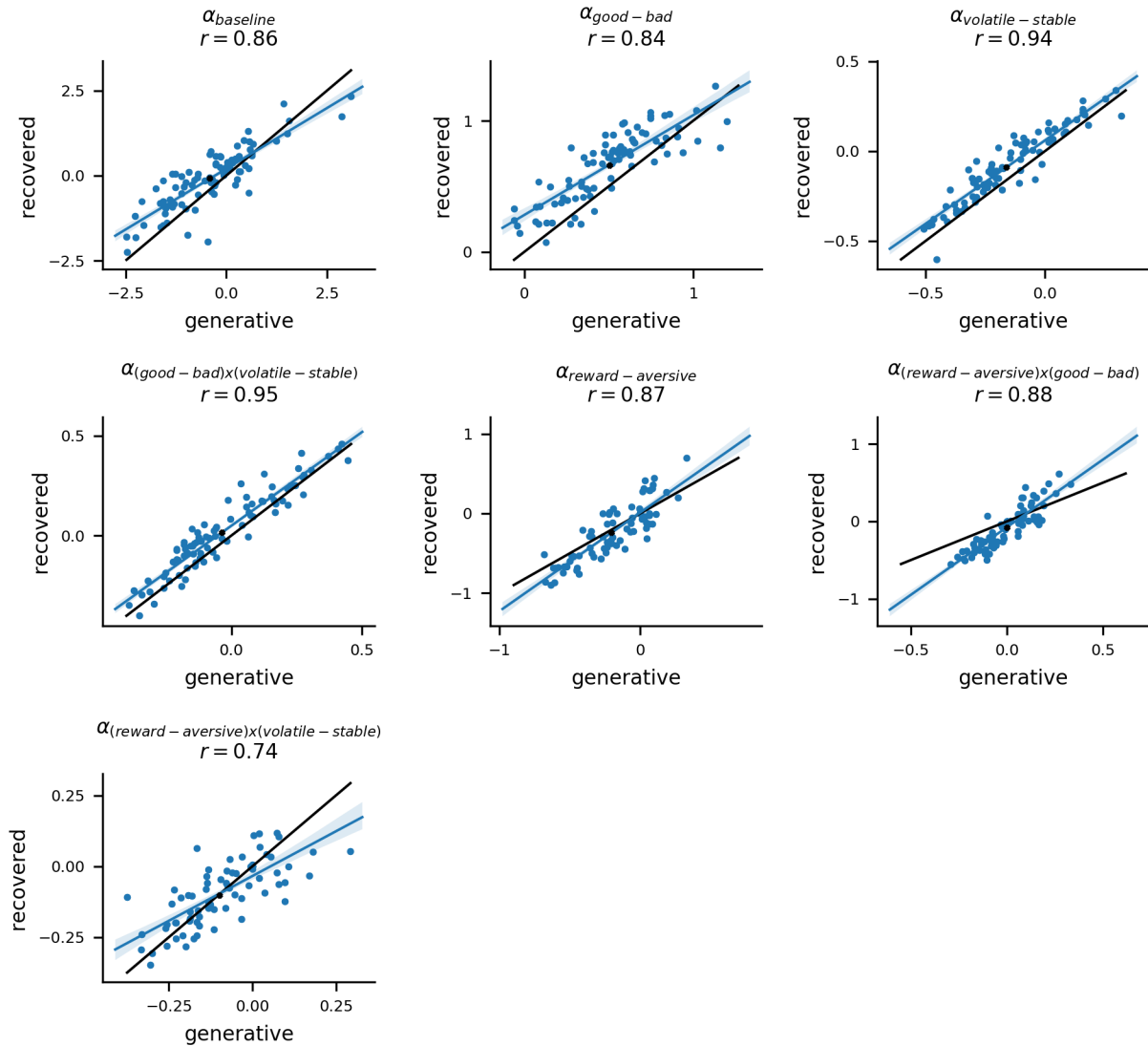
^{baseline}: For each parameter with this superscript, only one single baseline parameter was used, across both task versions and volatile and stable blocks.



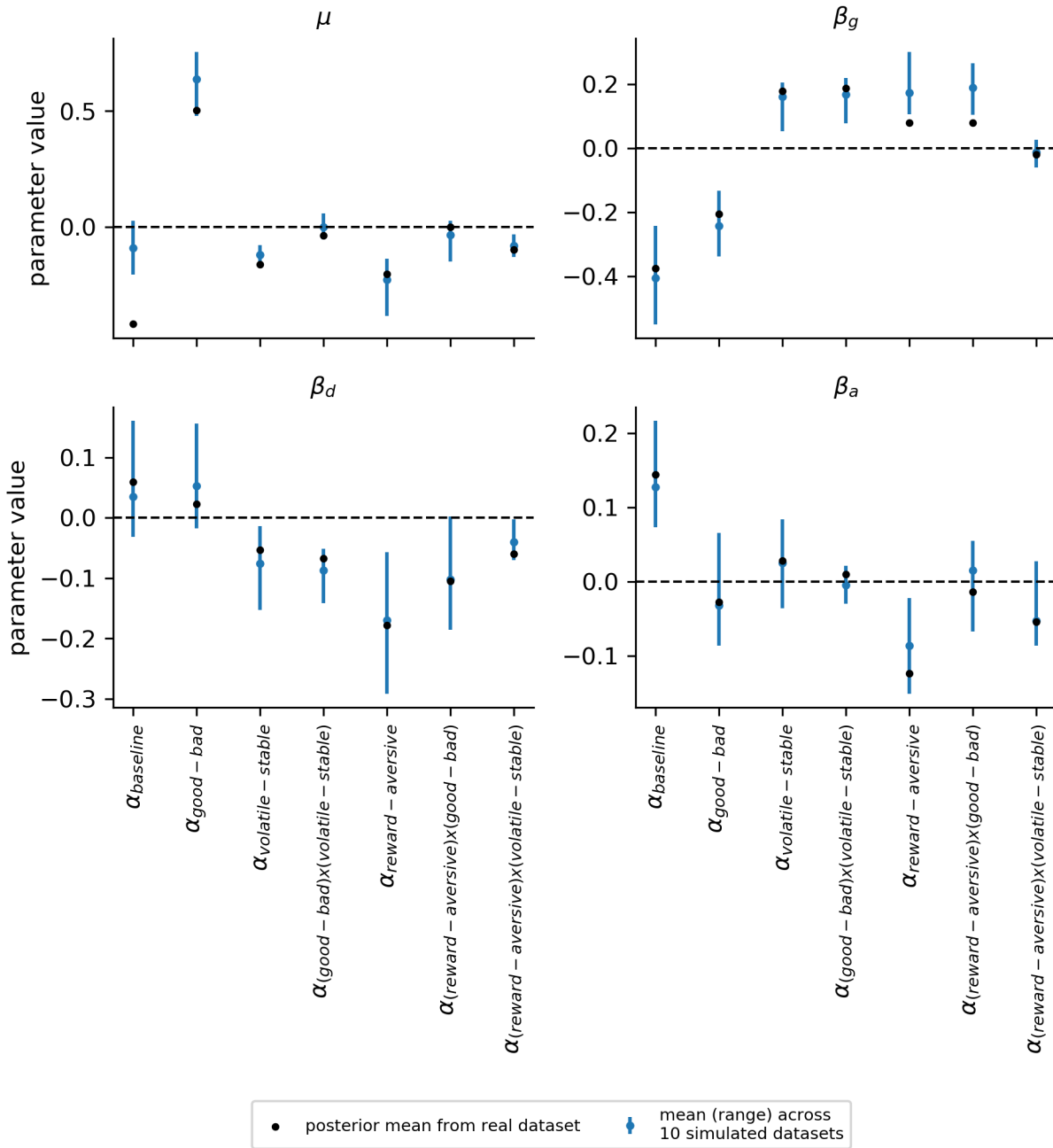
Supplemental Figure 2.1: Learning rate parameter components for individual participants (in-lab participants, $n=86$). Note that parameters exhibit strong linear relationships, because they were estimated using the general factor in the group-level distribution. Hence, statistical tests were not done on these parameter estimates, but instead were done using the posterior distribution for the group-level parameters (i.e. β_g).



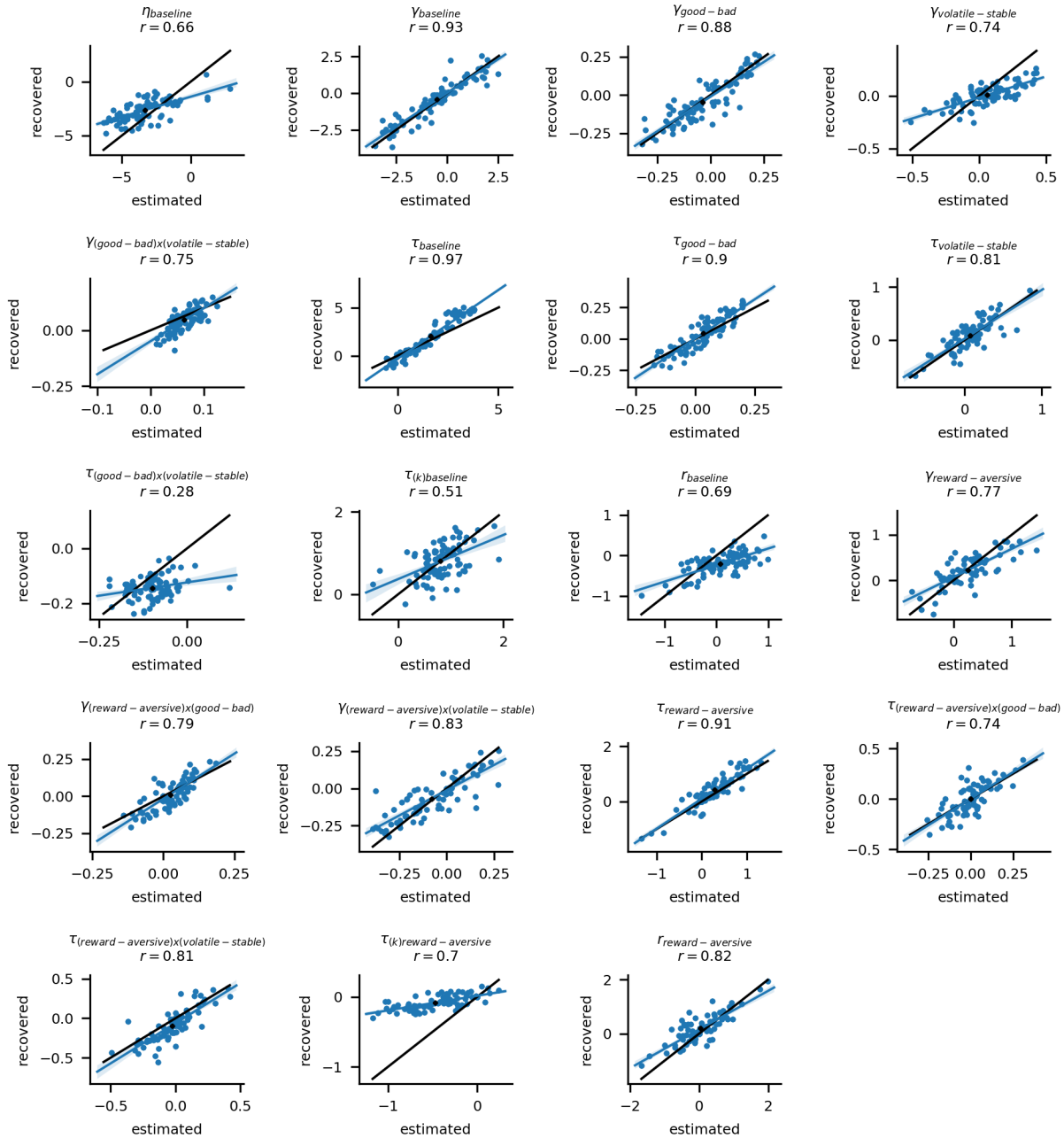
Supplemental Figure 2.2: Learning rate parameter components for individual participants (online Mturk dataset, n=147). Note similarity to in-lab participant sample (previous figure) for $\alpha_{\text{volatile} - \text{stable}}$ and $\alpha_{(\text{good} - \text{bad}) \times (\text{volatile} - \text{stable})}$.



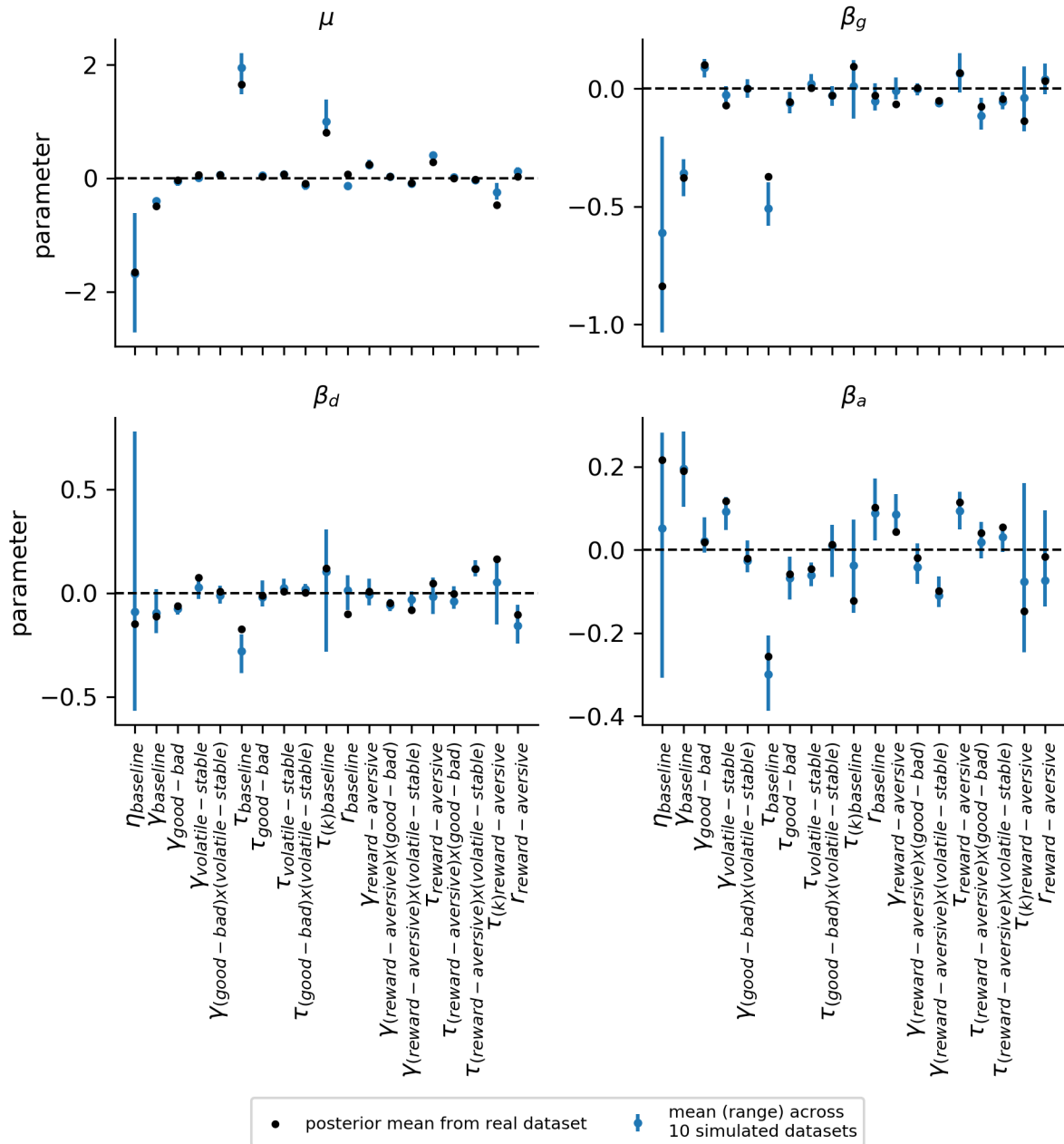
Supplemental Figure 2.3: Recovery of individual-level learning rate parameters. Parameter components estimated from the actual dataset for each participant (i.e. generative parameters) were used to simulate new behavioral data. Ten datasets with 86 participants each were simulated. Within each simulated dataset, newly estimated parameters (i.e. recovered parameters) were correlated with the generative parameters. A single example dataset is shown here. The average correlation between generative and recovered parameters across the 10 datasets for learning rate components was $r=0.88$ ($std=0.13$), confirming that parameters were recoverable.



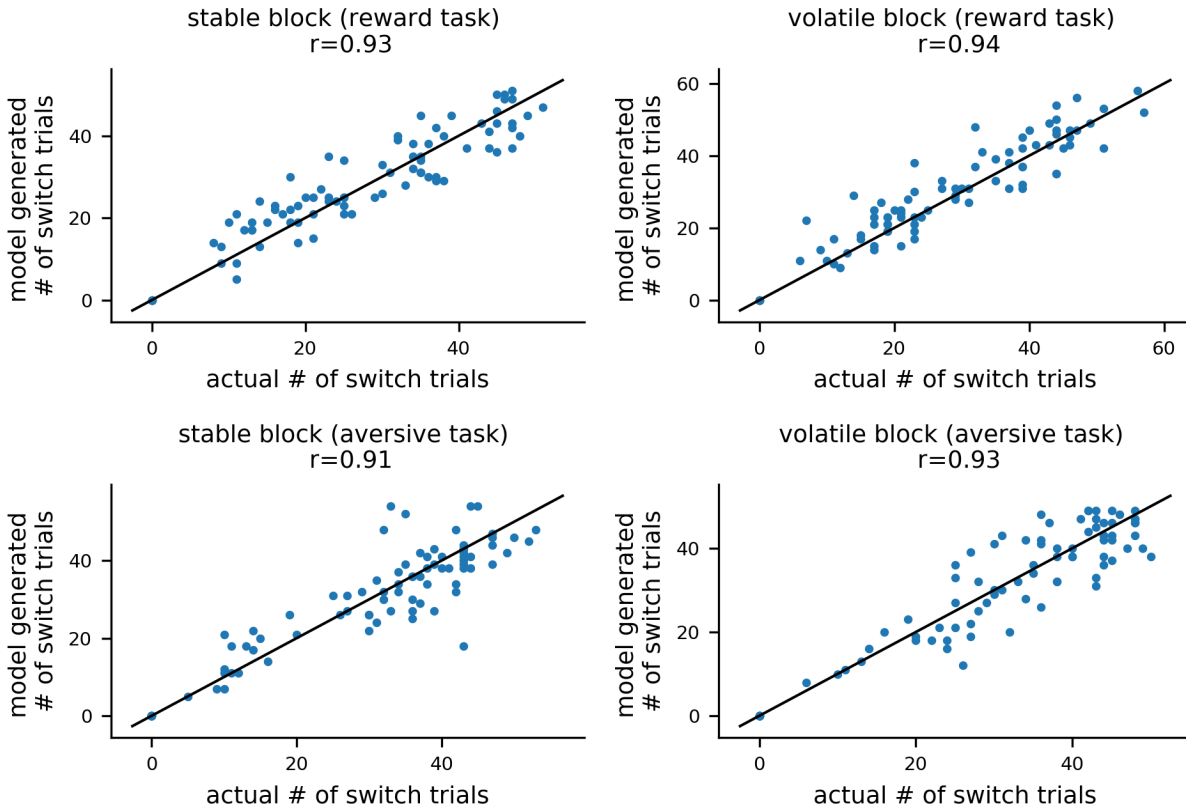
Supplemental Figure 2.4: Recovery of group-level learning rate parameters. Parameters estimated from the actual dataset for each participant were used to simulate new behavioral data. Ten simulated datasets with 86 participants each were generated. Recovered group-level averages μ and regression coefficients $\{\beta_g, \beta_d, \beta_a\}$ (i.e. their mean and range across 10 datasets shown as blue error bars) were very similar to the actual estimates from the real dataset (black data points).



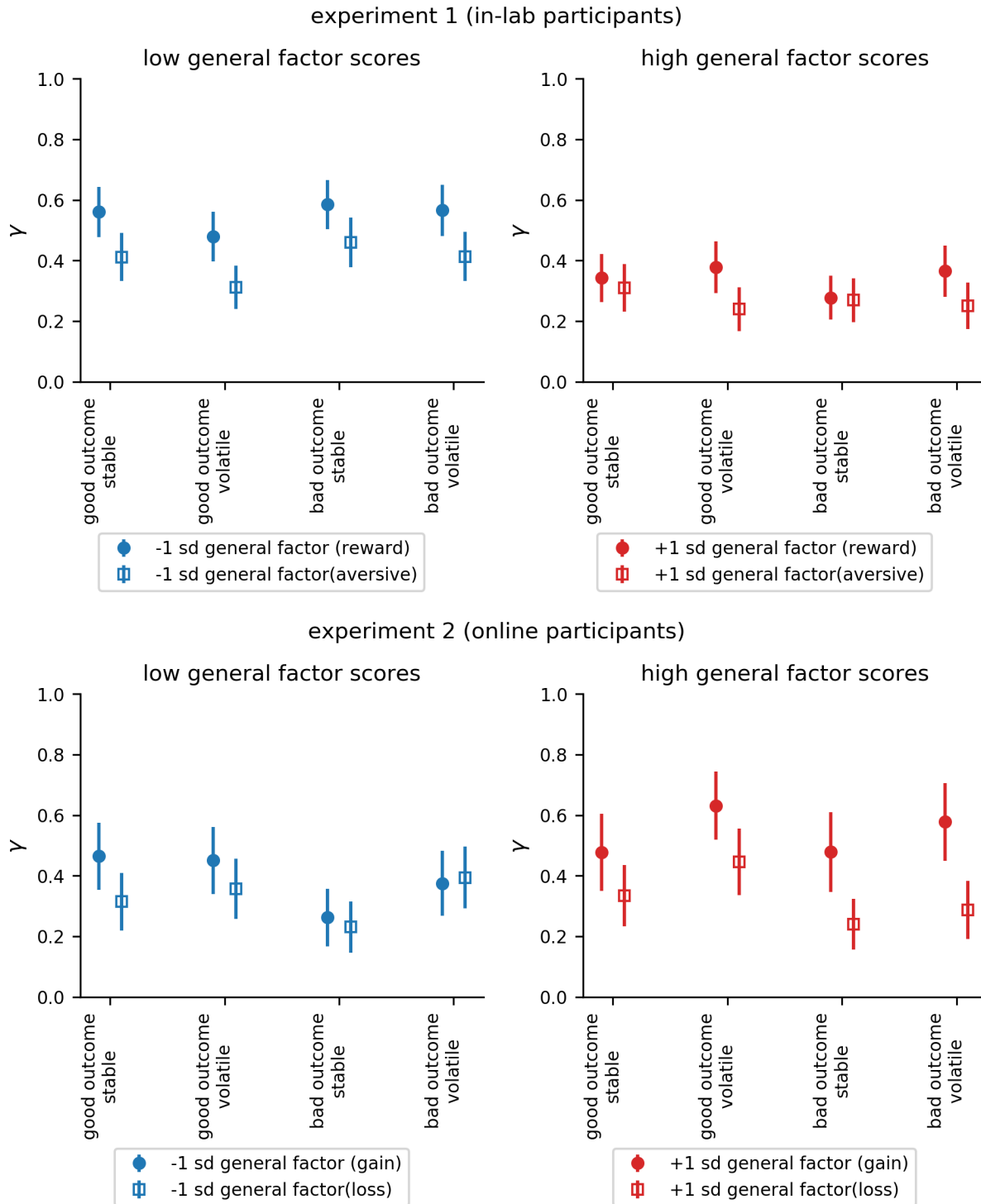
Supplemental Figure 2.5: Recovery of other (non-learning rate) individual-level parameters. Parameter components estimated from the actual dataset for each participant (i.e. generative parameters) were used to simulate new behavioral data. Ten datasets with 86 participants each were simulated. Within each simulated dataset, newly estimated parameters (i.e. recovered parameters) were correlated with the generative parameters. A single example dataset is shown here. The average correlation between generative and recovered parameters across the 10 datasets for all parameters was $r=0.76$ ($\text{std}=0.15$), confirming that parameters were largely recoverable.



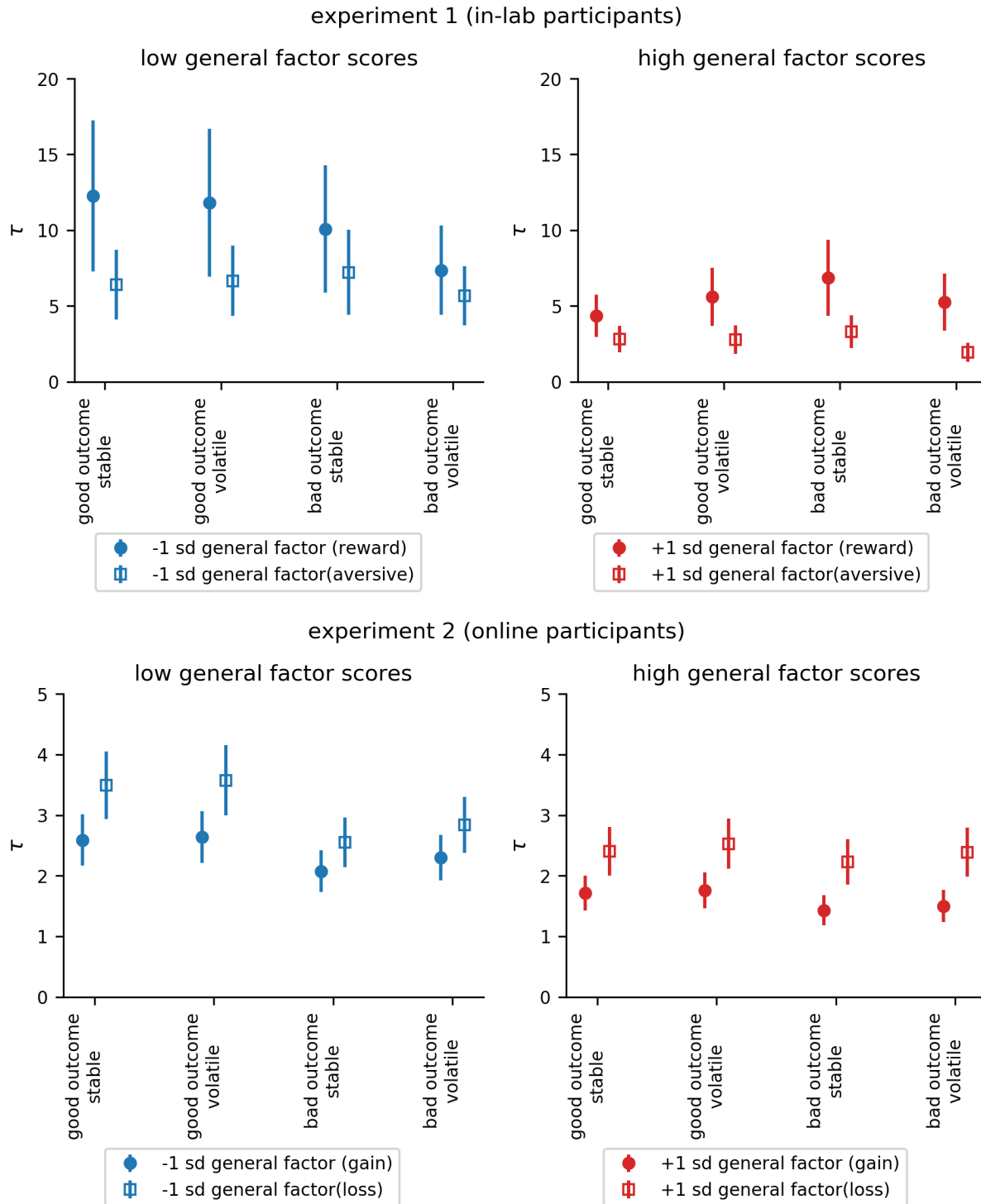
Supplemental Figure 2.6: Recovery of group-level learning rate parameters. Parameters estimated from the actual dataset for each participant were used to simulate new behavioral data. Ten simulated datasets with 86 participants each were generated. Recovered group-level averages μ and regression coefficients $\{\beta_g, \beta_d, \beta_a\}$ (i.e. their mean and range across 10 datasets shown as blue error bars) were very similar to the actual estimates from the real dataset (black data points).



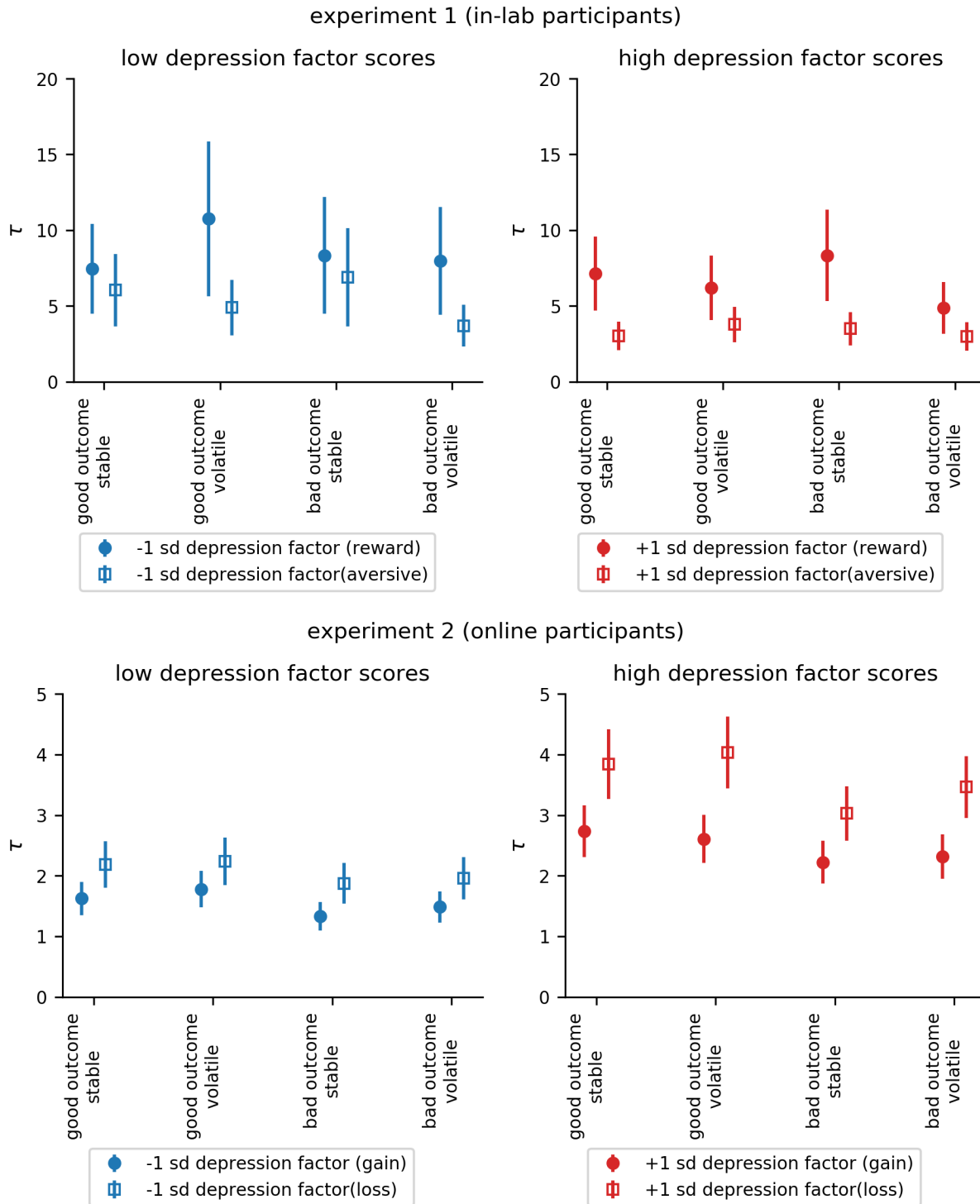
Supplemental Figure 2.7: Comparing actual and model generated numbers of switch trials. Parameters estimated from the actual dataset for each participant were used to simulate new behavioral data. Ten simulated datasets, with 86 participants each, were generated. The number of switch trials for each simulated participant was correlated with the number of switch trial for each actual participant. One of the ten simulated datasets is shown here as an example. The correlations between actual and generative switch trials were high (r 's > 0.88 across the four conditions shown above and across all datasets), demonstrating that the model can reproduce basic qualitative features of the data.



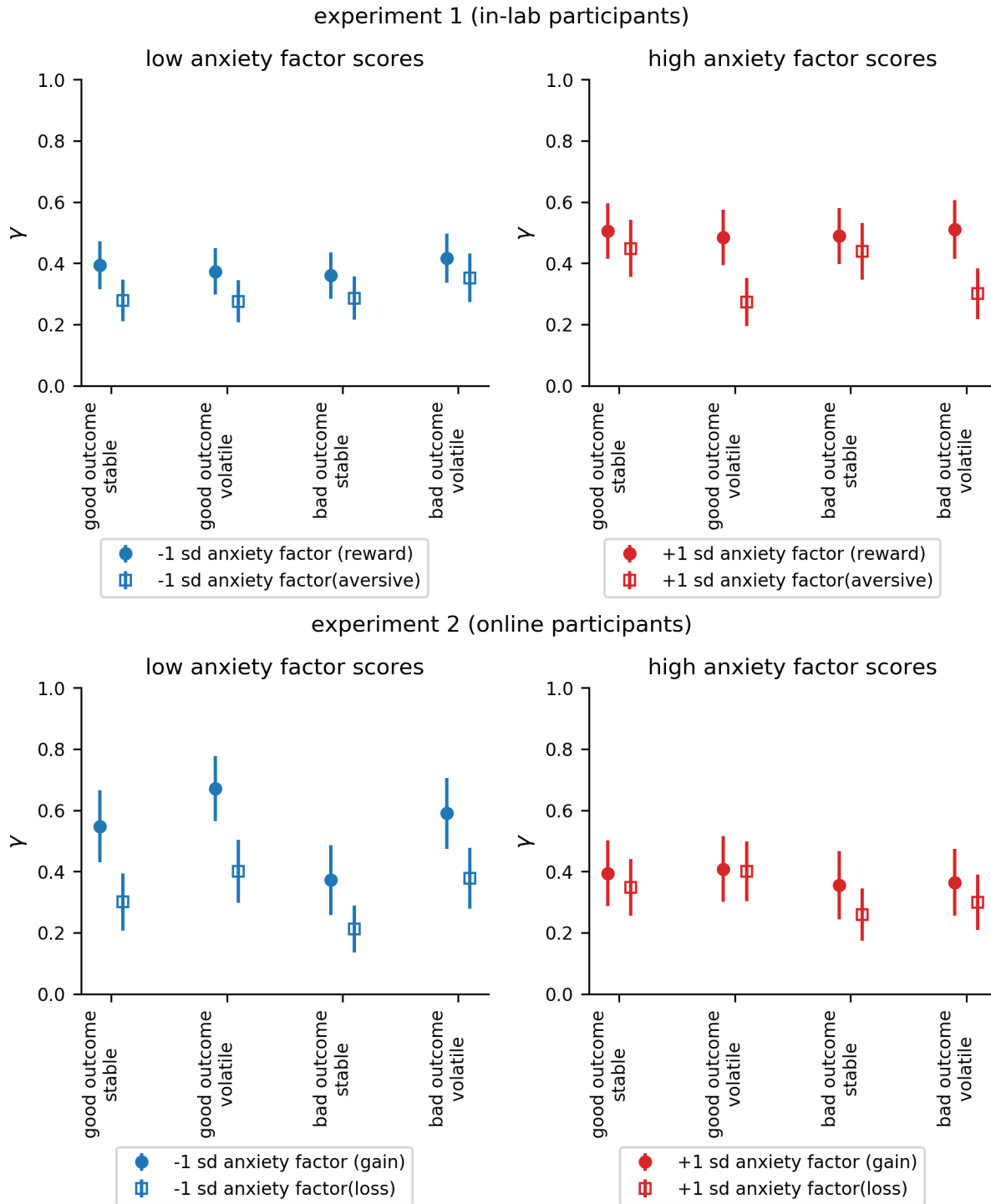
Supplemental Figure 2.8: Expected mean for the mixture weight (usage of probability versus magnitude) varying as a function of the general factor in the in-lab and online participant samples. The expected mean for the parameter for each condition was calculated in the same way as **Figure 2.4**.



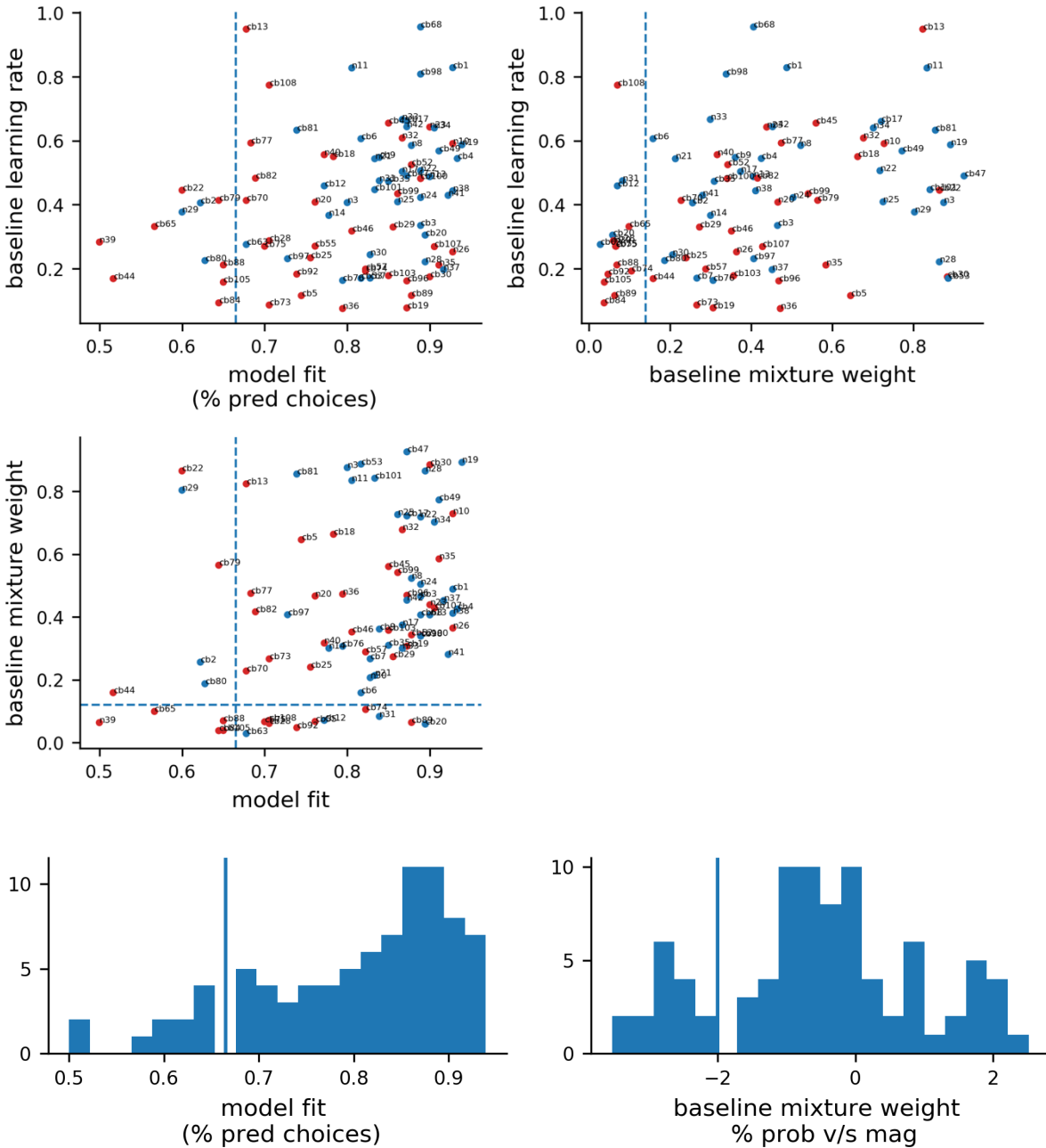
Supplemental Figure 2.9: Expected mean for the inverse temperature varying as a function of the general factor in the in-lab and online participant samples. The expected mean for the parameter for each condition was calculated in the same way as **Figure 2.4**.



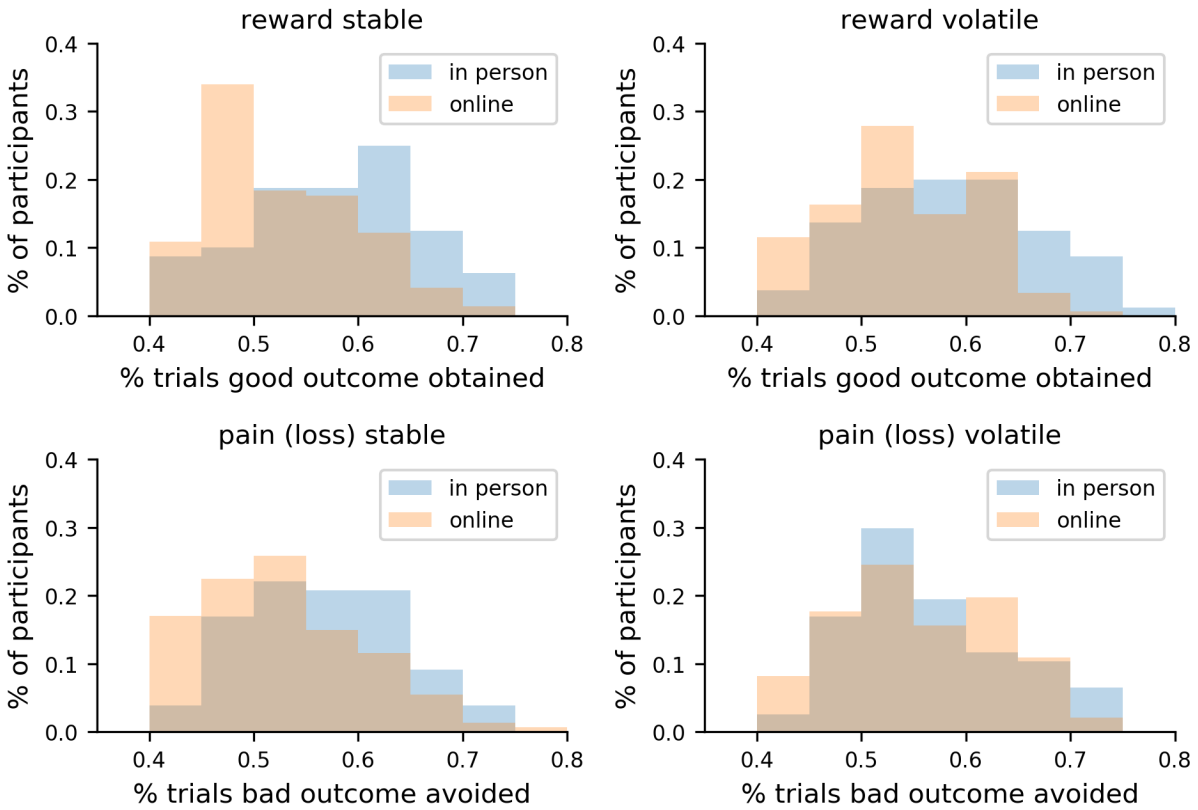
Supplemental Figure 2.10: Expected mean for the inverse temperature varying as a function of the depression-specific factor in the in-lab and online participant samples. The expected mean for the parameter for each condition was calculated in the same way as Figure 2.4.



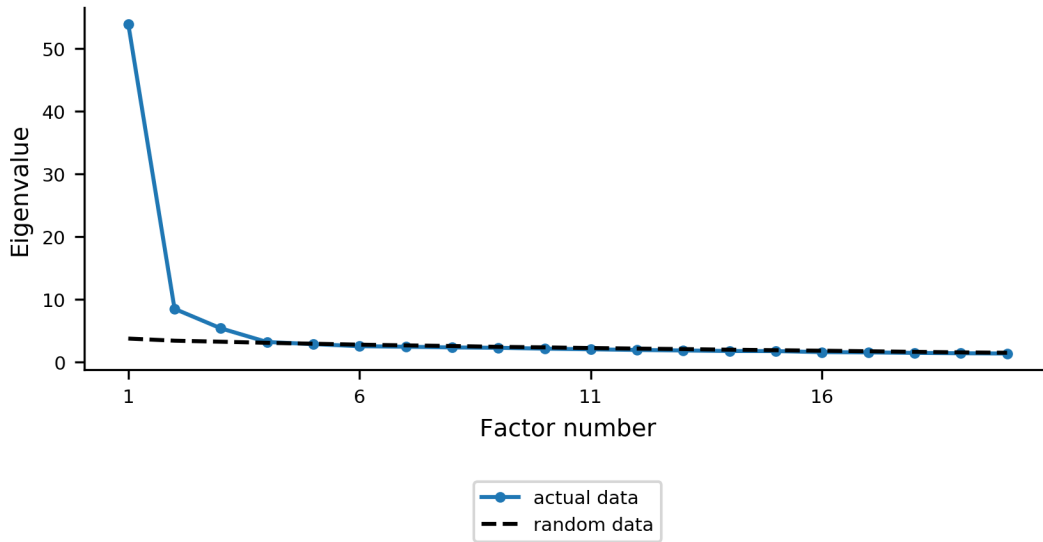
Supplemental Figure 2.11: Expected mean for the mixture weight (usage of probability versus magnitude) varying as a function of the anxiety-specific factor in the in-lab and online participant samples. The expected mean for the parameter for each condition was calculated in the same way as **Figure 2.4**.



Supplemental Figure 2.12: Baseline learning rate, baseline mixture weight, and model fit. Red data points denote individual participants who have scores on the general factor above the mean; Blue data points denote individuals with scores below the mean on general factor. High scoring (and low scoring) participants on the general factor occupy each part of the parameter space, which is important for ensuring accurate estimation of the relationship between learning rates and factor scores. Dotted lines correspond to the breaks in the histograms at the bottom of the plot. Participants below the horizontal line (or to the left of the vertical line) were excluded and the model was re-fit to ensure that the learning rate effects associated with the general factor were not driven by these poorer fitting participants.



Supplemental Figure 2.13: Performance on each task and block for the in-lab (n=86) and online datasets (n=147). The percent of rewarded trials (or trials on which no shock or a loss was delivered) is a basic measure of performance. There is a large range for the percentages, reflecting different learning and choice behavior across participants. However, the distributions are similar for in-lab and online participants.



Supplemental Figure 2.14: Scree plot for the eigenvalue decomposition of the covariance matrix of individual items from self-report symptom measures (in-lab participant sample, experiment 1). This suggests, along with the parallel analysis (described in the main text), that there were three dimensions of symptom variation that were distinguishable from noise.

Chapter 3: Decision Making Under Ambiguity

Introduction

We often have to make decisions under uncertainty. Two types of uncertainty have been classically distinguished from one another in economics. The first type of uncertainty, called risk, arises when outcomes occur probabilistically rather than deterministically. The second type of uncertainty, called ambiguity, arises when the outcome probabilities are themselves unknown. In general, people tend to be averse to both types of uncertainty when making decisions, preferring certain or highly probable outcomes to low probability outcomes and preferring outcomes with known probabilities to outcomes with unknown probabilities. These two types of aversion, known as risk aversion and ambiguity aversion, have been used to explain many behavioral phenomena from finance to medicine (Dow & Werlang, 1991; 1992; Curley et al., 1984; see Camerer & Weber, 1992 for review).

There is also substantial variability in individual attitudes towards risk and ambiguity. One source of variability is the framing of the decision, that is, whether financial outcomes are construed as gains or losses. It is well established that most individuals tend to be risk averse when the decision involves potential gains, but many people switch to being risk seeking when the decisions involve potential losses (Kahneman & Tversky, 1979; 1992; reviewed in Fox et al 2015). A gain-loss framing effect has also been observed for ambiguity (Cohen, Jaffray, & Said, 1985; Einhorn & Hogarth, 1986), but it is less well established and tends to be ignored in theoretical economic models (Kocher et al., 2018). Moreover, it is unclear how the gain-loss framing effect depends on level of ambiguity—most studies that have reported ambiguity seeking (or neutrality) in the case of loss contrasted known (risky) with completely unknown probabilities (completely ambiguous), without varying the level of possible information about the probabilities (partial ambiguity). Hence, the first aim of the current study was to investigate how individual attitudes towards ambiguity (i.e. aversion or seeking) under reward gain or reward loss vary as a function of ambiguity level (i.e. the level of missing information).

Differences in attitudes towards risk and ambiguity are also implicated in various forms of psychopathology. Anxiety disorders have been classically associated with risk avoidant behavior (Raghunathan et al., 1999; Maner et. al, 2007; Giorgetta et. al, 2012) and an intolerance for uncertainty (Birrell et al., 2011). However, most of the previous work on anxiety did not use direct behavioral measures of risk aversion or ambiguity aversion. One previously used task to study risk aversion is the balloon analog risk task (BART), in which participants choose how much to pump up a virtual balloon, increasing financial payout with each pump but losing everything if the balloon is inflated too much and explodes (Lauriola et al., 2014). In this task, both individuals with high levels of trait anxiety and individuals diagnosed with an anxiety disorder choose to pump fewer times than individuals with lower levels of anxiety, receiving a lower payout as a result of their cautious (risk averse) behavior (Maner et al., 2007). However, in the BART, the probability that the balloon explodes is always ambiguous, so ambiguity aversion and risk aversion cannot be teased apart. Moreover, each pump presents both the possibility of a gain in reward (if the balloon inflates successfully) and a loss in reward (if the

balloon explodes), so any potential impact of gain-loss framing for ambiguity (or risk) cannot be examined in relation to anxiety.

Two recent studies have used computational approaches to directly examine anxiety-related differences in risk aversion and ambiguity avoidance, although in separate studies. One study modeled risk aversion using a variant of a prospect theory model (Kahneman & Tversky, 1979) and observed that individuals with pathological anxiety were more risk averse than healthy controls, preferring certain smaller rewards to the risky larger rewards (Charpentier et al., 2017). Another study modeled two separate effects of ambiguity: an average (or categorical) aversion and an aversion that depends on the level of missing information (Lawrance et al., in review). In this study, elevated levels of trait anxiety were associated specifically with the information level dependent ambiguity aversion and not the categorical ambiguity aversion, for decisions involving primary aversive outcomes.

In light of these two more recent computational studies and the earlier studies investigating risk avoidance (e.g. Maner et al., 2007), there are a number of open questions regarding anxiety and its relation to ambiguity and risk. For one, do individuals with high levels of anxiety show elevated ambiguity aversion for decisions involving loss? This might be expected given the association between ambiguity aversion and anxiety for decisions involving primary aversive outcomes (Lawrance et al., in review) and the association between risk avoidant behavior and anxiety in the BART task, which mixes gain with loss and mixes risk with ambiguity. If so, does this potential relationship between anxiety and ambiguity aversion vary as a function of missing information, like it did in Lawrence et al. (in review)? Third, do individuals with high levels of anxiety show risk aversion for decisions involving only losses, given that they show risk aversion for decisions involving only gains and decisions involving a mixture of gains and losses (Charpentier et al., 2017). The second aim of the current study addresses these three related questions.

In contrast to anxiety, individuals experiencing symptoms of mania or hypomania (a less severe form of mania; Am. Psychiatr. Assoc. 2013) often show risk-taking behaviors (e.g. buying sprees or sexual indiscretions). Engagement of risky behavior even constitutes part of the diagnosis for manic episodes (see the DSM-V; Am. Psychiatr. Assoc. 2013). Experimentally, this behavior has been studied using Iowa gambling task (IGT), in which participants choose between two 'risky' decks of cards (both containing large potential gains and losses, but a negative expected value on average) and two 'safe' decks of cards (both containing small gains and losses, but a positive expected value on average). Bipolar patients have been shown to select the risky decks more often than healthy controls (Adida et al., 2008; 2011), but the evidence for this effect is inconsistent across studies, especially for patients who are not in the midst of experiencing a manic episode (i.e. euthymic patients; see Edge et al., 2013 for a meta-analysis). The BART task has also been used to examine differences associated with hypomania (Devlin et al., 2015), but the results from this study were somewhat contradictory across measures; individuals with higher levels of hypomania made fewer balloon pumps (i.e. were risk averse) in the BART, even though these same individuals reported that they were more likely to engage in risky real-world behaviors. Real-world risk-taking behaviors and risk-related behavioral differences in the IGT and BART tasks are typically thought to reflect altered reward processing (Johnson et al., 2012), which broadly includes a number of individual differences, such as in the willingness to expend effort in pursuit of reward, the elevation in mood following

the receipt of reward, etc. However, it is unknown whether ambiguity and risk sensitivity might also be involved, potentially contributing to risk-taking behavior alongside alterations in reward processing. The third aim of the current study is to test whether individuals with elevated mania-related symptoms show differences in attitudes towards risk and ambiguity. Given the association between mania and altered reward processing, we might expect individual differences specifically in the context of financial gains, rather than losses. This contrasts with our expectations for anxiety, which are that individual differences will occur in the context of losses, given the association between anxiety and enhanced threat-sensitivity (Barlow 2002).

Methods

Participants

Participants (N=1400; 565 females) participants were recruited from Amazon's Mechanical Turk (Mturk) platform and completed the experiment on an externally hosted website, within a single session. 1150 participants identified as white for their race/ethnicity, 85 as black or African American, 72 as Asian, 76 as more than one race, and 17 chose not to respond. Participants' ages were not recorded.

Self-Report Measures of Anxiety, Depression, Mania, and Schizophrenia-related Symptoms

Participants completed a number of standardized self-report symptom questionnaires: the Spielberger State-Trait Anxiety Inventory (STAI form Y; Spielberger, 1983), the Mood and Anxiety Symptoms Questionnaire (MASQ; Clark et al., 1995; Watson & Clark, 1991), the Penn State Worry Questionnaire (Meyer, Miller, Metzger, & Borkovec, 1990), the Center for Epidemiologic Studies Depression Scale (CESD; Radloff, 1977), the Hypomanic Personality Scale (HPS; Eckblad & Chapman, 1986), and the Oxford-Liverpool Inventory of Feelings and Experiences (OLIFE; Mason 1995).

Trait anxiety scores (i.e. the STAI-trait) had a mean of 40 (SD=12) across participants, which is slightly higher than has previously been observed in a large community sample (mean=33, SD=7.8; Knight 1983), but is in line with the observation that Mturk participants report higher levels of anxiety (Shapiro et al., 2013). 17% of our participants scored higher than 55 on the STAI, which was the mean score for pathologically anxious participants in Charpentier et al. (2017), suggesting that our sample likely includes individuals diagnosed (or diagnosable) with a mood or anxiety disorder.

Hypomanic personality scale scores (i.e. the HPS) had a mean of 13 (SD=8.1), which is slightly lower than has been observed previously (M=20, SD=12; Eckblad & Chapman, 1986). Fewer than 1% of our participants had HPS scores greater 36, which has been previously used as cutoff for 'high' levels of hypomania, and above which a participant is likely to have had at least one hypomanic episode (Eckblad & Chapman, 1986). Therefore, participants in our sample who have medium to high scores on the HPS should be considered at-risk rather than diagnosable with bipolar disorder.

Using Bifactor analyses to estimate symptom factor scores

In order to determine whether any potential relationship to risk aversion or ambiguity aversion was specific to anxiety or shared with depression given the substantial comorbidity between the two (Kessler et al., 2005), we applied bifactor analysis to the questionnaire subscales related to mood and anxiety disorders (8 subscales: STAI anxiety, STAI depression, MASQ anxious arousal, MASQ anhedonia, MASQ anxious symptoms, MASQ depressive symptoms, CESD, and PSWQ). In this bifactor analysis, we estimated a general factor to represent shared variance, which we refer to as the *negative affect general factor*. The general factor explained 77% of the total variance captured by the bifactor model. In previous work, it is often the case that only two specific factors are estimated in addition to the general factor, one for anxiety and one for depression (see **Chapter 2**). In the current study, however, we were able to estimate three specific factors: a *physiological anxiety factor*, a *cognitive anxiety factor*, and a *depression factor*. Constraining the model to estimate only two specific factors from our dataset did not result in a single anxiety factor containing both physiological anxiety and worry, like some previous work, but instead resulted in the PSWQ (i.e. worry) questionnaire loading onto the depression specific factor. Hence, we felt that the construct validity was improved by retaining three specific factors.

Scores were estimated for each participant on each of these four factors. **Supplemental Figure 3.2** shows the correlations between participant factor scores and their scores on questionnaire subscales. The *negative affect general* factor scores correlated broadly with all the subscales related to mood and anxiety symptoms. The *depression factor* scores correlated most strongly the anhedonia-related items from the MASQ (**Supplemental Figure 3.2d**; $r=0.64$), in line with the tripartite model (Watson & Clark, 1991). The two anxiety-related factor scores (*physiological anxiety* and *cognitive anxiety*) correlated most strongly with the anxious arousal subscale from the MASQ ($r=0.57$) and the PSWQ (i.e. worry) questionnaire. Similar specific factors have been observed separately in previous bifactor analyses (Clark & Watson, 1994; Steer et al., 2008 for anxious arousal; and Brodbeck et al., 2011 for worry).

We applied a second bifactor analysis to the mania- and schizophrenia-related questionnaire subscales (5 subscales: HPS, OLIFE unusual experiences (e.g. delusions and hallucinations), OLIFE introverted anhedonia (i.e. negative or depressive symptoms), OLIFE impulsive nonconformity (e.g. antisocial or impulsive behavior), and OLIFE cognitive disorganization (e.g. poor attention and concentration)). Both mania and symptoms of schizophrenia relate more broadly to a dimension of psychopathology that is referred to as thought disorder symptomatology (Kotov et al., 2011). For this bifactor analysis, we estimated a general factor and two specific factors from these five subscales. Trying to estimate more factors resulted in a factor containing loadings all equal to zero, meaning that we could not extract more than three factors from the five subscales. Scores were estimated for each participant for each of three factors. Scores on the *thought-disorder general* factor were correlated positively with scores on four out of the five subscales (with the exception of OLIFE anhedonia). Scores on the *mania-related* factor correlated most strongly with the hypomanic personality scale (HPS; $r=0.86$), but also correlated with the OLIFE impulsive nonconformity subscale ($r=0.38$) and weakly correlated with the OLIFE unusual experiences subscale ($r=0.2$). The *mania-related* factor can therefore be thought of as a broadened measure of mania-related

symptoms, which includes some variance related impulsivity, and which is orthogonal to the shared (thought-disorder related) variance. It is also orthogonal to the other specific factor, the *negative and cognitive symptoms* factor, which correlated with the OLIFE anhedonia subscale ($r=0.8$) and the OLIFE cognitive disorganization subscale ($r=0.69$).

In contrast to the *negative affect general* factor, which explained 77% of the total variance captured by the bifactor model, the *thought disorder general* factor explained only 41%. Given the apparently lower level of shared variance between mania-related and schizophrenia-related subscales, we also used the hypomanic personality scale (HPS) to relate to behavior; however, the results using these two measures were expected to be very similar, given high correlation ($r=0.86$) between them.

Task: Decision Making under Risk and Ambiguity

The experimental task was adapted from the classic Ellsberg urn task (Ellsberg 1961) and consisted of 300 trials. On half the trials, participants could gain points, and in the other half of trials, participants could lose points. At the end of the experiment, points were converted to a monetary bonus between \$0-\$5. Participants with point totals in the top 5% received \$5, those with point totals in the top 25% received \$2.50, and those with point totals the top 50% received \$1. Gains or loss trials were grouped into blocks of 5 trials. The order of these blocks was pseudorandomized with the condition that there could be no more than 4 blocks of the same type in a row. The type of trial was indicated by the background color of the screen (blue for gain, red for loss).

On each trial, participants were asked to make a choice between two “urns”, each of which was filled with 50 tokens (X and O symbols) (see **Figure 3.1** for an example). When one of the two urns was chosen by the participant, a token was drawn uniformly at random from that urn. If an ‘O’ was drawn, the participant’s point total remained unchanged. If an ‘X’ was drawn, the participant’s point total was either increased if the trial was a gain trial or decreased if the trial was a loss trial. The amount of potential increase or decrease was displayed as a number above each urn and referred to as outcome magnitude. Outcome magnitudes ranged from 0 to 150. Participants were instructed that the number of X’s divided by the total number of tokens determined the probability that they would receive an outcome (a gain or loss in points versus no change in points). Participants were encouraged to consider both the outcome probability and the magnitude of the potential outcome when making their decision.

On half of the trials, the outcome probability for one of the two urns was rendered ambiguous by hiding some of the tokens using an ‘=’ sign. Participants were instructed that there was a still a true but unknown probability of receiving an X, but that they couldn’t observe this probability directly. Instead, they had to infer the underlying probability from the tokens whose identities were shown. There were 6 different levels of ‘missing information’ (*MI*). This was defined as: $MI = 1 - \sqrt{\frac{n}{50}}$, where 50 is the total number of tokens and n could be: 1, 3, 5, 10, 20 or 40 tokens revealed. This definition was used to match Lawrance et al. (in review). Outcome probability and magnitude were manipulated orthogonally to missing

information level.

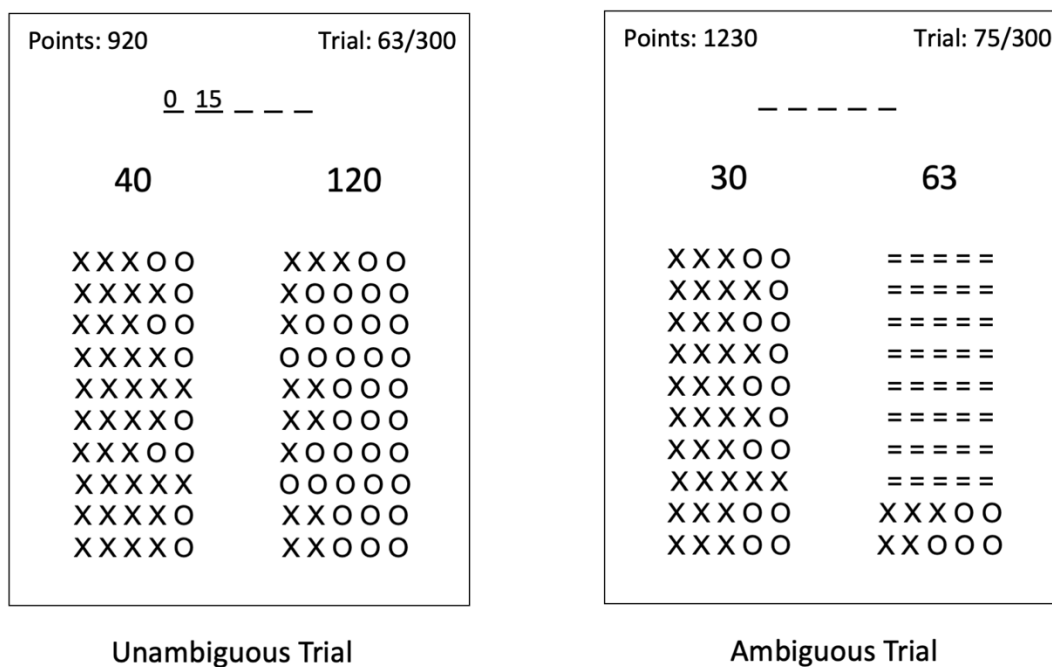


Figure 3.1: Two example trials from the task. On unambiguous trials (shown on the left), participants choose between the two ‘urns’, each containing a different number of tokens (X’s and O’s). After the participant chooses an urn, a token is drawn uniformly at random. If an X is drawn, the number shown above the urn (e.g. 40 or 120) is added to the participant’s point total, otherwise the total remains the same. To maximize their point total, participants are instructed to consider both the probability that an X will be drawn and the amount of points that would be delivered if it is. Points are added to the total on gain trials (designated by a blue background; not shown) and points are subtracted on loss trials (designated by a red background and negative signs in front of the outcome magnitudes; not shown). On ambiguous trials, the probability for one urn (on the far right) is made to be ambiguous by covering some of the tokens with ‘=’. Participants are told that the tokens that they can see are a random sample from the urn and can be used to infer the true underlying probability (ratio of X’s to the total number of tokens). The number of tokens revealed could be: 1, 3, 5, 10, 20 or 40, which correspond to different levels of missing information.

Participant Exclusion

Data from participants was excluded if they failed to answer all four ‘catch questions’ correctly on the self-report measures, failed to make a choice between the two urns before 8 seconds expired on more than 5% of the trials, or failed to achieve better than 85% accuracy on ‘no brainer’ trials in both the first and second half of the experiment. ‘Catch questions’ instructed the participants to select one of the possible answers. For example, on the STAI, one question instructed participants to select “almost always”, which was one of the four possible answers. Three other similar catch questions were included, two questions in the MASQ and one question in the PSWQ. ‘No brainer’ trials were unambiguous trials that had one urn that was clearly better than the other in that it was associated with a better outcome and a higher probability of obtaining that outcome. There were 34 ‘no brainer’ trials in total. After applying these exclusion criteria, roughly one quarter of participants (n=359) remained for the analyses.

Modeling Choice Behavior

To measure attitudes (aversion or seeking) towards risk and ambiguity, we fit participant's choice behavior using a modified expected utility model. In the model, we assume that participants first calculate the expected utility (EU) of each urn by multiplying the outcome magnitude (M) by the outcome probability (P) (Eqn. 1a-b; urn identity is denoted by subscript). The magnitudes were exponentiated using a risk parameter λ , converting nominal values into subjective values (i.e. utilities). We estimated separate risk parameters for gain and loss trials (denoted as λ_{gain} and λ_{loss} in the text). On gain trials, a risk parameter less than one ($\lambda_{gain} < 1$) creates a concave utility function, corresponding to risk aversion, and a risk parameter greater than one ($\lambda_{gain} > 1$) creates a convex utility function, corresponding to risk seeking. On loss trials, the correspondence between the risk parameter and the risk attitude is reversed since EU is multiplied by -1 ($\lambda_{loss} < 1$ implies risk seeking and $\lambda_{loss} > 1$ implies risk aversion).

Eqn. 1a-b (for unambiguous trials)

$$\begin{aligned} EU_1 &= M_1^\lambda P_1 \quad (\text{for urn \#1}) \\ EU_2 &= M_2^\lambda P_2 \quad (\text{for urn \#2}) \end{aligned}$$

When choosing between urns, we assume that participants take the difference between the expected utilities (EU) of the two urns and choose probabilistically according to logistic choice function (Eqn. 2). The inverse temperature parameter β (again separated for gain and loss trials: β_{gain} and β_{loss}), determines the degree to which choices were driven by the difference in EU between the two urns.

Eqn. 2 (for unambiguous trials)

$$P(\text{choice}_{urn \#1}) = \frac{1}{1 + \exp(-\beta(EU_1 - EU_2))}$$

On ambiguous trials, the outcome probability for one of the two urns is only partially observable, and therefore needs to be estimated from the information given. As an estimate for this probability, we use the posterior mean (denoted as P_a in Eqn. 3b) from a Bayesian (beta-binomial) model: $beta(K + 1, N - K + 1)$, where K is the number of revealed X 's and N is the total number of revealed tokens. The effect of using this Bayesian estimate is that the observed probability (ratio of X 's to total number of tokens) is adjusted towards 50% by an amount determined by the level of missing information. Although participants may not be using this approach exactly to estimate the probability, we use it so that we can define ambiguity aversion (or seeking) as an effect above and beyond what is normative (i.e. rational). Indeed, ambiguity aversion is typically defined as a 'preference' or 'attitude' that cannot be explained by subjective expected utilities (SEU; Savage 1954), which combines subjective (Bayesian) probabilities and subjective values (utilities) to create a rational model for choice. The Bayesian adjusted expected utility is denoted by SEU in Eqn. 3a-b, in reference to subjective expected utility.

Eqn. 3a-b. (for ambiguous trials)

$$\begin{aligned} EU_u &= M_u^\lambda P_u \quad (\text{for the unambiguous urn}) \\ SEU_a &= M_a^\lambda P_a \quad (\text{for the ambiguous urn}) \end{aligned}$$

Having incorporated ambiguity into choice using a simple, yet rational Bayesian approach, ambiguity aversion (or seeking) was then modeled using two additional parameters (see Eqn. 4). β_0 represents a tendency to choose or avoid choosing the ambiguous urn at an average level of ambiguity, above and beyond differences in the subjective expected utility. β_3 represents an additional dependence on the level of missing information (MI), which is z-scored. These two parameters are also separated for gain and loss trials ($\beta_{0_{loss}}$ and $\beta_{0_{gain}}$, $\beta_{3_{loss}}$ and $\beta_{3_{gain}}$).

Eqn. 4 (for ambiguous trials)

$$P(\text{choice}_{unambiguous\ urn}) = \frac{1}{1 + \exp(-1 * (\beta(EU_u - SEU_a) + \beta_0 + \beta_3 * MI))}$$

All trials were concatenated for parameter estimation. Weak independent priors, either a Normal(mean=0, sd=10) or a HalfNormal(mean=0, sd=10) distribution, were used to prevent extreme outliers. Parameters were estimated by maximizing the likelihood multiplied by these priors (i.e. using Maximum A Posteriori). All but one of the parameters had distributions of fitted values that deviated from normality (Shapiro-Wilk test, $p < 0.05$). Therefore, Spearman's rank correlation was used.

Confirming that participants' behavior was fit reasonably well by our model, the average number of choices predicted correctly by the model was $77\% \pm 4\%$ SD (SE=0.1%), and the average pseudo R^2 was 0.30 ± 0.04 SD (SE=0.0002).

We note that the main model had a better penalized fit (BIC) than the model used previously (Lawrance et al., in review) for the current dataset (but not the previous dataset). Including ambiguity effects (β_0 and β_3) also provided a better penalized fit than a model that omitted the effects of ambiguity (see **Supplemental Figure 3.1**).

Results

Risk aversion versus risk seeking under conditions of reward gain and reward loss: group level findings

Before addressing our three aims, we first tested whether we could replicate the previously observed finding of a gain-loss framing effect for risk attitude. We measured this in our model using the risk parameters, λ_{gain} and λ_{loss} . We observed that participants' risk parameters, on average, were significantly less than one for both gain ($\lambda_{gain} < 1$, $t(358) = -9.2$, $p < 0.001$) and loss trials ($\lambda_{loss} < 1$, $t(358) = -4.2$, $p < 0.001$) (**Figure 3.3a**). In line with prior work (Kahneman & Tversky, 1979; Tversky & Kahneman, 1992), this corresponded to a concave utility

function on gain trials and a convex utility function on loss trials, on average across participants (**Figure 3.3b**). This meant that participants tended to be risk averse in the context of gains, preferring high probability small gains to low probability large gains, and that participants tended to be risk seeking in the context losses, preferring the opposite.

We also observed this gain-loss reversal at the individual participant level (**Figure 3.3c**). Most participants (64.9%, SE=2.5%; z-test for proportions against 50%, $z(358) = 5.9$, $p < 0.001$) had risk parameters that were both less than one ($\lambda_{gain} < 1$ and $\lambda_{loss} < 1$), meaning that most individuals switched from being risk averse to risk seeking. Eleven percent (11%) of participants showed the reverse pattern, nineteen percent (19%) of participants were risk averse in both cases, and six percent (6%) of participants were risk seeking in both cases.

We validated the reversal results reported above using a model-agnostic exploration of the performance on trials where the two urns had very similar expected values (there were no trials with exactly equivalent expected value). This represented 25% of unambiguous trials (see **Supplemental Table 3.1** for individual trials). On gain trials where the two urns had very similar expected values, participants tended to choose the urn with the higher probability outcome (mean proportion of times higher probability urn chosen=66.7%, between participant SD=15%, and between participant SE=0.7%; one sample t-test against 50%, $t(358) = 21$, $p < 0.001$). Conversely, on loss trials, on which the two urns had very similar expected values, participants only chose the urn with the higher probability 35% of the time (SD=18%, SE=0.9%; one sample t-test against 50%, $t(358) = -14$, $p < 0.001$). The difference between these percentages was highly significant (paired sample t-test, $t(358) = 21$, $p < 0.001$). Moreover, this effect was present even on the trials on which the expected value difference worked against the effect. On gain trials, on which the higher probability outcome had slightly lower expected value, participants still chose it 73% of the time. On loss trials, on which the higher probability outcome had slightly higher expected value, participants still only chose it 36% of the time.

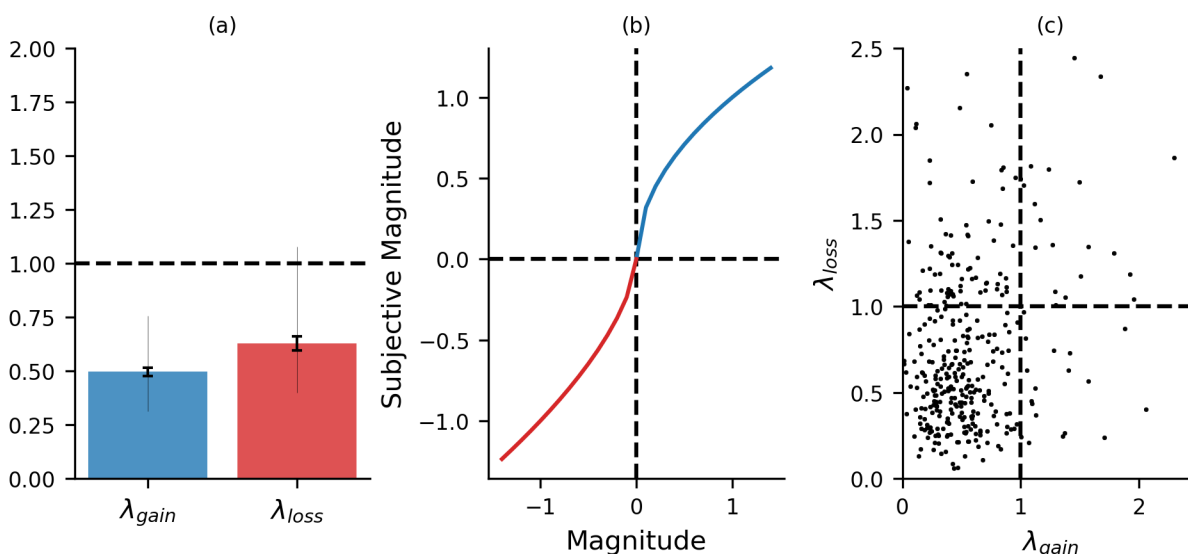


Figure 3.3: Risk avoidance versus risk seeking on gains and loss trials. A utility function transforms nominal values for outcome magnitudes into subjective values. The curvature of this function, which is determined by the risk parameters (λ_{gain} , λ_{loss}), is used to model risk averse and risk seeking behavior (as concave and convex utility functions, respectively). (a) Median risk parameters for both gain and loss trials were significantly less than one ($\lambda_{gain} < 1$, $p < 0.001$; $\lambda_{loss} < 1$, $p < 0.001$). Thick error bars represent bootstrapped standard error, while thin error

bars represent 25th and 75th percentiles across participants. (b) The utility function implied by the median group risk parameters is shown (magnitudes 0 to 1.5 in the model correspond to 0-150 points in the experiment). This median utility function is concave for gain trials and convex for loss trials, corresponding to risk aversion for gains and risk seeking for losses, respectively. (c) Individual parameter estimates show that the majority (68%) of participants have risk parameters that are both less than one ($\lambda_{gain} < 1$ and $\lambda_{loss} < 1$), which means that most people switch from being risk averse in gains to risk seeking in loss.

Ambiguity aversion versus ambiguity seeking under conditions of reward gain and reward loss across different levels of missing information: Group level findings

The first aim of the current study was to explore the effects of gain-loss framing on ambiguity aversion (or seeking), and the potential modulation of that effect by missing information level. To that aim, we estimated separate ambiguity parameters for average ambiguity aversion (versus seeking) (β_0), and ambiguity aversion (versus seeking) as a function of level of missing information (β_3). Separate parameters were estimated for reward gain trials and reward loss trials.

On reward gain trials where one of the two urns had missing information (i.e. ambiguous trials), $\beta_{0_{gain}}$ was significantly greater than zero on average across participants ($\beta_{0_{gain}} > 0$; one sample t-test, $t(358) = 26$; $p < 0.001$; **Figure 3.4a**). Moreover, $\beta_{0_{gain}}$ was greater than zero for the vast majority of participants (94%, $SE = 1.2\%$; z-test for proportions against 50%, $z(358) = 36$, $p < 0.001$), implying that most participants would choose an unambiguous urn over an ambiguous urn containing the average level of missing information (around 10 tokens revealed), when the two urns were equated for subjective expected utility. As the level of missing information in the ambiguous urn increased, participants on average tended to choose the unambiguous urn to an even greater extent ($\beta_{3_{gain}} > 0$; $t(358) = 7.5$, $p < 0.001$; **Figure 3.4b**). The $\beta_{3_{gain}}$ parameter was also greater than zero in the majority of participants (65%, $SE = 2.5\%$; z-test for proportions against 50%, $z(358) = 5.7$, $p < 0.001$). The combined attitude toward ambiguity as a function of missing information ($\beta_0 + \beta_3 * MI$) can be seen in **Figure 3.4c** for the group mean (thick line) and for individual participants (blue thin lines). As shown in this figure, for gains, most participants exhibited some degree of ambiguity aversion at all levels of missing information (i.e. $(\beta_0 + \beta_3 * MI) > 0$).

On reward loss trials where one of the two urns had missing information (i.e. ambiguous trials), we observed the reverse pattern for the categorical effect of ambiguity (β_0), which was below zero, on average ($\beta_{0_{loss}} < 0$; $t(358) = -6.2$, $p < 0.001$; **Figure 3.4a**), and for the majority of participants (63%, $SE = 2.5\%$; z-test for proportions against 50%, $z(358) = 5.3$, $p < 0.001$). This implies that most participants would choose an ambiguous urn containing an average level of missing information over an unambiguous urn of equivalent (negative) subjective expected value. This reversal from categorical ambiguity avoidance for gains to categorical ambiguity avoidance for losses mirrors that the loss gain reversal observed in relation to risk. As the level of missing information increased on loss trials, however, participants tended to become more ambiguity averse ($\beta_{3_{loss}} > 0$; $t(358) = 14.8$, $p < 0.001$; 79% of participants, $SE = 2.1\%$; **Figure 3.4b**), similar to their behavior on gain trials. **Figure 3.4d** shows that participants, on average (thick line), were ambiguity seeking ($(\beta_0 + \beta_3 * MI) < 0$) at low levels of missing information

(40 to 20 tokens revealed) and approximately ambiguity neutral ($(\beta_0 + \beta_3 * MI) \cong 0$) at higher levels of missing information.

To validate these observations, we examined the percentage of trials where participants chose the unambiguous over the ambiguous urn at each level of missing information. As shown in **Supplemental Figure 3.3**, on gain trials, the middle 50% of participants (i.e. those between the 1st and 3rd quartiles) chose the unambiguous more often than the ambiguous urn at each level of missing information, except at n=10 tokens revealed. On loss trials, however, the middle 50% of participants chose the ambiguous urn more often than the unambiguous urn for the lowest level of missing information (i.e. 40 out of 50 revealed tokens). For higher levels of missing information (20 or fewer tokens revealed), the middle 50% of participants did not all choose the ambiguous over the unambiguous urn, on average across trials. In general, these results support those involving β_0 and β_3 .

We also observed a high degree of correlation among the four ambiguity parameters across participants. On gain trials and loss trials separately, the categorical effect of ambiguity was positively correlated with the information-level dependent effect ($\{\beta_{0_{gain}}, \beta_{3_{gain}}\}$, $rs(358) = 0.49$, $p < 0.001$); $\{\beta_{0_{loss}}, \beta_{3_{loss}}\}$ $rs(358) = 0.46$, $p < 0.001$; rs denotes rank correlation), meaning that participants' aversion to average levels of ambiguity was related to their information-level dependent aversion (**Figure 3.4e-f**). Across gain and loss trials, the four ambiguity parameters were all positively correlated as well (rs 's > 0.2 , p 's < 0.001 ; only $\beta_{0_{gain}}$ to $\beta_{0_{loss}}$ and $\beta_{3_{gain}}$ to $\beta_{3_{loss}}$ are shown in **Figure 3.4g-h**), suggesting that individuals have a generalized tendency to be more or less averse across the domains.

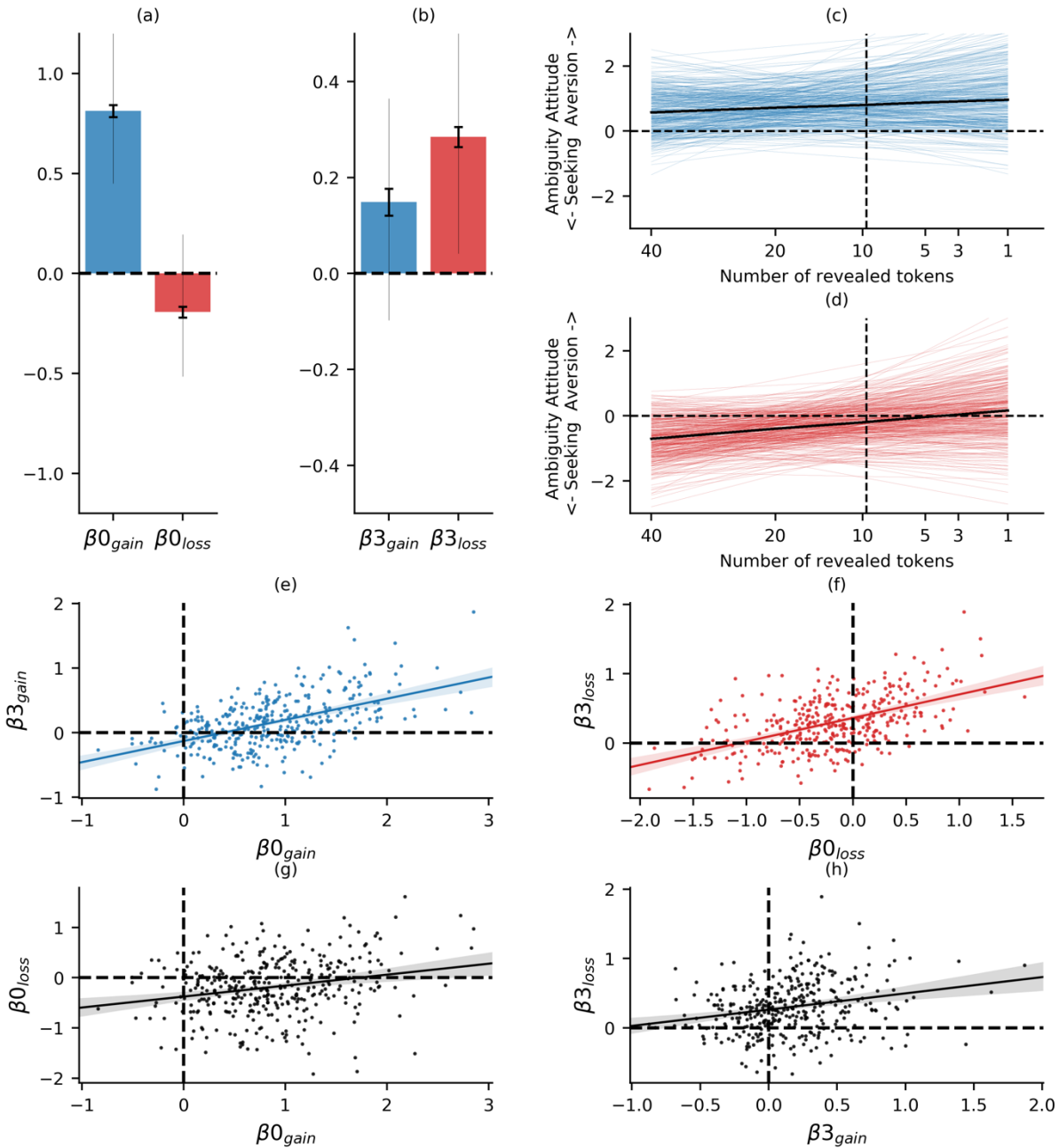


Figure 3.5: Ambiguity parameters on gain and loss trials (a-b) Median parameters for average ambiguity parameter (β_0) and missing information-level dependent ambiguity parameter (β_3), for both gain and loss trials. Thick error bars are bootstrapped standard error; thin error bars are 25th and 75th percentiles. (c-d) Combined attitude towards ambiguity (aversion or seeking) as a function of missing information ($\beta_0 + \beta_3 * MI$) is plotted for group-level median parameters (thick lines), and for individual participants (thin lines). On gain trials, most participants show ambiguity aversion ($(\beta_0 + \beta_3 * MI) > 0$) at all levels of missing information. On loss trials, most participants show ambiguity seeking ($(\beta_0 + \beta_3 * MI) < 0$) at low levels of missing information, but have a more neutral attitude at higher levels of ambiguity ($\beta_0 + \beta_3 * MI \cong 0$). (e-h) Individual ambiguity parameters, within each domain (e-f) and across domains (g-h) are significantly correlated.

Correlation between attitudes towards risk and attitudes towards ambiguity

To explore the relationship between attitudes towards risk and attitudes towards ambiguity, we examined the correlations of risk and ambiguity parameters within each outcome domain (gains or losses). On gain trials, the risk parameter (λ_{gain}) was positively correlated with both ambiguity parameters ($\beta_{0_{gain}}$, $rs(358) = -0.12$, $p=0.02$; $\beta_{3_{gain}}$, $rs(358) = -0.19$, $p<0.001$). On loss trials, the risk parameter (λ_{loss}) was only weakly related to the average ambiguity parameter ($\beta_{0_{loss}}$, $rs(358) = -0.1$, $p=0.05$), and not at all related to the information level dependent parameter ($\beta_{3_{loss}}$, $rs(358) = 0.05$, $p=0.25$). These results suggest that the relationship between attitudes towards risk and attitudes towards ambiguity differ depending on whether decisions are framed as gains or as losses.

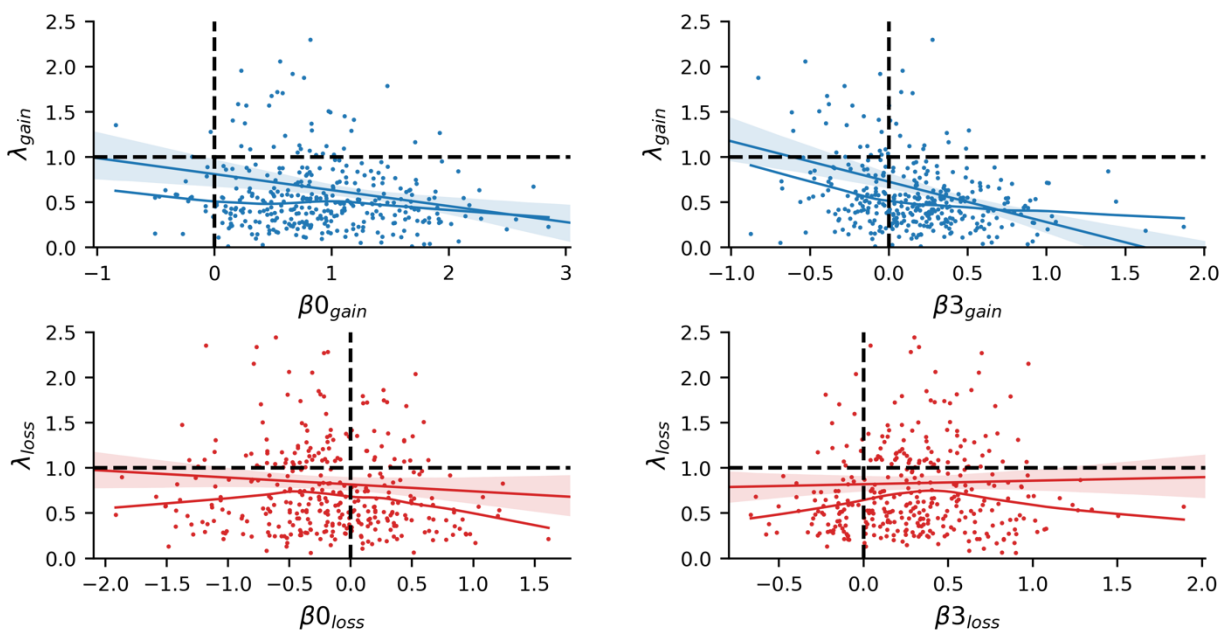


Figure 3.6: The relationship between attitudes towards risk and attitudes towards ambiguity. For gain trials, the risk parameter (λ_{gain}) was significantly, but moderately, correlated with both ambiguity parameters ($\beta_{0_{gain}}$, $rs=-0.12$, $p=0.02$; $\beta_{3_{gain}}$, $rs=-0.19$, $p<0.001$). For loss trials, the risk parameter (λ_{loss}) was only marginally correlated with the categorical ambiguity parameter ($\beta_{0_{loss}}$, $rs=-0.1$, $p=0.05$) and uncorrelated with the information-level dependent ambiguity parameter ($\beta_{3_{loss}}$, $rs=0.05$, $p=0.25$).

Anxiety-related symptoms and attitudes towards risk and ambiguity

Under the second aim, our first question was whether individuals reporting higher levels of anxiety-related symptoms would show elevated ambiguity aversion for decisions involving potential loss. Our second question was, if so, whether the relationship between anxiety and ambiguity aversion would be modulated by the level of missing information. We observed no correlation between the categorical ambiguity parameter ($\beta_{0_{gain}}$, $\beta_{0_{loss}}$) and either the *physiological anxiety* or the *cognitive anxiety* factor scores (all p 's > 0.17). In other words, participants with higher levels of anxiety did not show elevated aversion to ambiguity at the average level of missing information. However, there was a modest but significant correlation

between scores on the *physiological anxiety* factor and the information-level dependent ambiguity parameter, for losses ($\beta_{3_{loss}}$, $rs(358) = 0.12$, $p=0.029$) as well as for gains ($\beta_{3_{gain}}$, $rs(358)= 0.12$, $p=0.024$; p -values are uncorrected; see **Figure 3.7**). That is, individuals who reported higher levels of *physiological anxiety*, relative those with lower levels, showed a faster rate of increase in their aversion to ambiguity as levels of missing information increased. Scores on the *cognitive anxiety* factor, as well as the *depression* factor and the *negative affect general* factor, were not correlated with the information-level dependent ambiguity parameter for either gain or loss trials (p 's >0.19), suggesting specificity to physiological anxiety related symptoms.

Our third question, under the second aim, was whether anxiety would be related to risk aversion for decisions involving only losses, given the previous association between risk aversion and anxiety for gain-only and mixed gain-loss decisions (Charpentier et al., 2017). We did not observe any significant correlations of the *physiological anxiety* or *cognitive anxiety* factor scores (or the depression or general factor scores) with the risk parameters (λ_{gain} , λ_{loss} ; p 's >0.16 , uncorrected). These correlations as well as the ones for other parameters and factors can be found in **Supplemental Figure 3.3**.

Mania-related symptoms and attitudes towards risk and ambiguity

Under the third aim, we tested whether mania-related symptoms were related to differences in attitudes towards risk or ambiguity for decisions involving reward gain or loss, given the previous associations between risk-taking behavior and bipolar disorder. Analogously to our procedure for mood and anxiety symptoms, we measured mania-related symptoms in the context of other closely related symptoms (i.e. those related to schizophrenia).

There was no significant correlation of scores on the *mania-related* factor with risk parameters for either gain (λ_{gain}) or loss (λ_{loss}) trials (p 's >0.2 ; see **Supplemental Figure 3.3**). Similarly, the scores on the HPS were not significantly correlated with the risk parameter on loss trials (λ_{loss} , $rs(358) = -0.01$, $p = 0.87$). However, correlation on gain trials was trend-level significant (λ_{gain} , $rs(358)=0.1$, $p=0.06$), which, if replicated, would imply that higher scores on the hypomanic personality scale are associated modestly with less risk aversion.

For decisions involving ambiguity, there were no significant correlations between the *mania-related* factor and the categorical ambiguity parameter ($\beta_{0_{gain}}$, $\beta_{0_{loss}}$, p 's >0.07). However, we observed that scores on the *mania-related* factor were negatively correlated with the information-level dependent ambiguity parameter ($\beta_{3_{gain}}$) on gain trials ($rs(358) = -0.12$, $p=0.05$), but not on loss trials ($\beta_{3_{loss}}$, $rs(358) = 0.01$, $p=0.81$). As can be seen in **Figure 3.7**, participants with higher *mania-related* factor scores had $\beta_{3_{gain}}$ values that were closer to zero. This meant that, in the context of potential gains, these individuals showed a smaller rate of increase in ambiguity aversion as the level of missing information increased. A weaker trend-level relationship between scores on the HPS and $\beta_{3_{gain}}$ was also observed ($rs(358) = -0.1$, $p=0.06$), suggesting that the *mania-related* correlation was largely driven by scores on the HPS. This relationship also survives controlling for any of the mood or anxiety factors (physiological anxiety, cognitive anxiety, depression, or general factor) by including them individually as regressors alongside the *mania-related* factor (p 's <0.05); this was important to test, because

concomitant mood and anxiety symptoms likely modulate reward- or threat-related processing biases in bipolar disorder (Johnson et al., 2012).

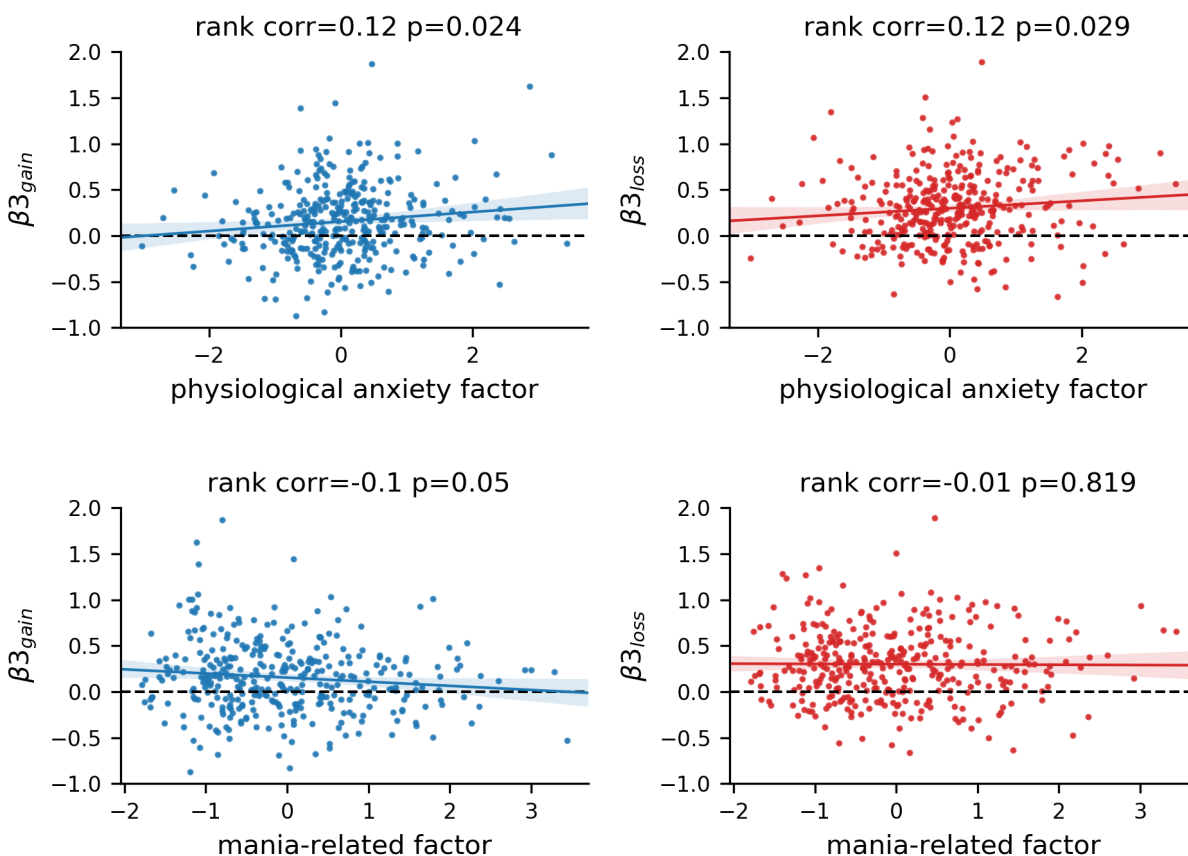


Figure 3.7: Information-level dependent ambiguity parameter and scores on the *physiological* and *mania-related* factors. Scores on the *physiological anxiety* factor were significantly correlated with the information-level dependent ambiguity parameter for both gain (top left; blue) and loss trials (top right; red). This meant that individuals with higher scores exhibited ambiguity aversion that increased at a faster rate to increases in the level of missing information, relative to individuals with lower scores. In contrast, scores on the *mania-related* factor showed the opposite pattern for gain trials (bottom left; blue). Participants who had high scores on this factor showed less sensitivity to missing information—meaning that they did not increase their aversion to ambiguity as quickly as the level of missing information increased (less sensitivity is indicated by $\beta_{3_{gain}}$ parameter values near zero).

Discussion

The current study had three aims. Under the first aim, we looked at how a gain-loss framing effect for ambiguity, observed in a few previous studies, might change based on the level of ambiguity (i.e. missing information). Under the second aim, we explored the relationship of anxiety with both ambiguity aversion and risk aversion. We asked whether individuals with higher anxiety showed elevated ambiguity aversion for decisions involving loss, and if so, whether the aversion depended on the level of missing information. We also

examined whether anxiety was related to risk aversion for decisions involving loss, given previous associations risk aversion and anxiety for decisions involving gains or mixed gains and losses (Charpentier et al., 2017). Under the third aim, we tested whether mania-related symptoms were related to differences in attitudes towards risk or ambiguity, especially for decisions involving financial gains, given previous associations between risk-taking behavior, altered reward processing, and bipolar disorder. Under both aims two and three, we used factor scores derived from bifactor analyses to test whether anxiety- or mania-related symptoms were associated uniquely, or together with co-occurring symptoms (i.e. depression or schizophrenia-related symptoms), to differences in attitudes towards risk and ambiguity.

In line with prior theory (Kahneman & Tversky, 1979), most participants in our study exhibited behavior consistent with a concave utility function for gains and convex utility function for losses. In other words, they exhibited risk aversion for gains and risk seeking for losses. In addition to replicating previous results, our study supports this theory in a relatively under-utilized experimental paradigm. Most previous studies that explored gain-loss reversals in risk attitudes asked participants to choose between a risky option and certain alternative with equal or greater expected value (Fishburn & Kochenberger, 1979; Cohen, Jaffray, & Said, 1985, Tversky & Kahneman, 1992). In contrast, we asked participants to choose between two risky options, which complicates the decision and the measurement of risk attitude; nonetheless, we showed that a risk attitude reversal is robustly observable for this type of decision.

Regarding attitudes towards ambiguity, we extended previous work (Cohen, Jaffray, & Said, 1985; Einhorn & Hogarth, 1986; Kocher et al., 2018) by investigating how gain-loss reversals are potentially modulated by the level of missing information; most previous studies have investigated these reversals using complete ambiguity and did not include decisions involving partial ambiguity. For decisions involving gains, we observed that participants tended to be averse to ambiguity across all levels, with the aversion increasing at higher levels of missing information. For decisions involving loss, on the other hand, we observed that participants were ambiguity seeking at low levels of missing information, but they became more averse to ambiguity (being approximately neutral towards ambiguity) at higher levels. In other words, the attitudes towards ambiguity did not seem to reverse entirely under conditions of gain and loss; for losses relative to gains, average ambiguity attitude moved from aversion towards seeking, but the effect of missing information had the same direction of effect. This suggests potentially two different mechanisms—for example an instinctual attitude towards ambiguity and a graded attitude that might reflect further calculations involving the amount of missing information. In either case, the influence of ambiguity on choice was not captured by a simple, yet rational, adjustment of expected utilities using Bayesian (subjective) probabilities.

In Lawrance et al. (in review), we previously observed that individuals with higher trait anxiety, relative to those with lower trait anxiety, showed a higher rate of increase in ambiguity aversion as the level of missing information increased for decisions involving primary aversive outcomes. This effect was captured by a parameter that corresponded to the β_3 parameter in our model. In the current study, we recruited a large enough sample size to examine whether ambiguity aversion related to different subdimensions of anxiety (i.e. physiological or cognitive) or to negative affect more broadly (i.e. the general factor). We observed that the same β_3 parameter was weakly, but significantly related to scores on the *physiological anxiety* factor

under conditions of both gain and loss. If replicated, this finding has a number of interesting implications.

The fact that we observed anxiety-linked ambiguity aversion for decisions involving both gains and losses, in addition to decisions involving threat (Lawrance et al., in review), suggests that this attitude may contribute to avoidance behavior in anxious individuals in a variety of different situations, both experimental and real-world. For example, in the balloon analog risk task (BART), which involves ambiguity, the avoidance behavior (Maner et al., 2007) may be driven by ambiguity aversion rather than (or in addition to) risk aversion, loss or threat sensitivity, or differences in learning about the probabilities across trials. For real-world decisions, ambiguity aversion may also contribute to avoidance, for example, of unfamiliar social settings or during travel, alongside the overestimation of the probability that negative events will occur (proposed in **Chapter 1.2** and **1.3**).

That the physiological anxiety factor was primarily involved and not the cognitive anxiety factor was an unexpected result—we might have predicted that cognitive anxiety would be related to ambiguity aversion, given our proposal that worry relates to biased (model-based) simulations and to the overestimation of the probability that uncertain negative outcomes will occur (**Chapter 1.2** and **1.3**). We speculate, however, that physiological anxiety may be relevant in the context of the current experiment, because decisions were made on a very short time scale (typically under 2 seconds). A historically influential theory of decision making (the somatic marker hypothesis: Bechara & Damasio, 2005) argued that decisions are driven to a large extent by feedback involving physiological (somatic) information. A process like this, which is likely to be computationally faster than a model-based simulation, might be more heavily relied on in our task, therefore implicating physiological anxiety rather than cognitive anxiety like we observed.

In contrast to Charpentier et al. (2017), we did not observe a significant association of risk attitude with either the anxiety (or the depression) factors. There are a few possible reasons that we did not observe this association. One could be that our participant sample did not contain as many participants with clinical levels of anxiety; however, 17% of our participants reported STAI trait anxiety scores that were above the mean score in the pathologically anxious group from Charpentier et al. (2017), so this is unlikely to be the main reason. Another possible reason is that risk attitudes may change when participants are also dealing with ambiguity. Differences in probability between a risky and certain alternative, which may normally evoke aversion in anxious individuals, may be overshadowed by the difference between known probabilities and unknown probabilities, which occurred on half of the trials in our task.

In our study, we also observed a weak, but significant relationship between attitudes towards ambiguity and scores on our mania-related factor, which were predominantly related to variance in the hypomanic personality scale (HPS; correlation $r=0.86$ between the two). We observed that high scores on the mania-related factor were associated with a slower rate of increase in ambiguity aversion as the level of missing information increased for decisions involving gains but not losses. To our knowledge, no previous study has investigated ambiguity attitude at various levels of missing information and its relation to symptoms of mania. Our finding, if replicated, would suggest another potential driver, among those related to altered reward processing (Johnson et al., 2012), for risk-taking behavior during episodes of mania.

Risk-related behavior, in relation to bipolar I disorder, has been previously studied experimentally using the Iowa Gambling Task (IGT; Rubinsztein et al., 2006; Adida et al., 2008; 2011). In the IGT, both patients experiencing acute mania and those who were euthymic (not experiencing an episode of mania), tended to select cards from risky decks over safer decks more often than healthy controls (Adida et al., 2008; 2011). This effect has been observed at a meta-analytic level (Edge et al., 2013), however the authors conclude that the IGT is likely to be an insufficiently sensitive measure of differences between euthymic patients and healthy controls. Many of these previous studies suggest that risky choices (both in the IGT task and in the real-world) arise from altered reward processing (Edge et al., 2013; reviewed in Johnson et al., 2012). Our result suggests that in addition to other potential individual differences in reward processing, a vulnerability to experience mania may also be associated with a decreased sensitivity (i.e. aversion) to ambiguity. This aversion, which may normally contribute in part to the selection of the safer choices in healthy individuals, may not be there to put on the breaks for these individuals.

Limitations

One caveat regarding our findings on risk attitudes is that risk attitudes typically follow a fourfold pattern (Fox et al, 2015): risk aversion for gains and risk seeking for losses for moderate to high probabilities ($p > 0.5$), and risk seeking for gains and risk aversion for losses for low probabilities ($p < 0.1$). In our study, we only had 3 unambiguous trials with probability less than 0.1, so we would be unable to estimate whether participants showed different risk attitudes for low probabilities.

To model attitudes towards ambiguity, we chose a particular form for how ambiguity impacts choice, namely that ambiguity and missing information contribute additively to subjective expected value. There are a dozen or more alternative theoretical models for the effects of ambiguity that we could have based our model on, each with their advantages and disadvantages (see Camerer & Weber, 1992 for a review). However, the aim of the current study was not to test alternative specific accounts for the way ambiguity influences choice—doing that effectively would require a task designed to amplify subtle differences between models. Instead, our aim was to capture the essential differences in attitudes towards ambiguity between the domains of gains and loss, and how those differences related to symptoms of psychopathology. To that aim, we chose a simple, interpretable model that measured risk attitude and ambiguity attitude separately, and that measured ambiguity attitude in a way that was translatable to our previous work (Lawrance et al., in prep).

Another caveat is that we did not include diagnostic measures, such as the General Behavior Inventory for bipolar disorder (Depue et al., 1981). It is therefore difficult to know whether attitudes towards ambiguity differ between the group of individuals who meet criteria for a diagnosis and the group of individuals who have elevated symptoms but do not meet that same criteria.

Finally, we note that the statistical significance for the symptom-to-parameter correlations do not survive correction for multiple comparisons. Replication would be required before strong conclusions can be drawn regarding anxiety and mania and attitudes towards ambiguity.

Supplemental Results

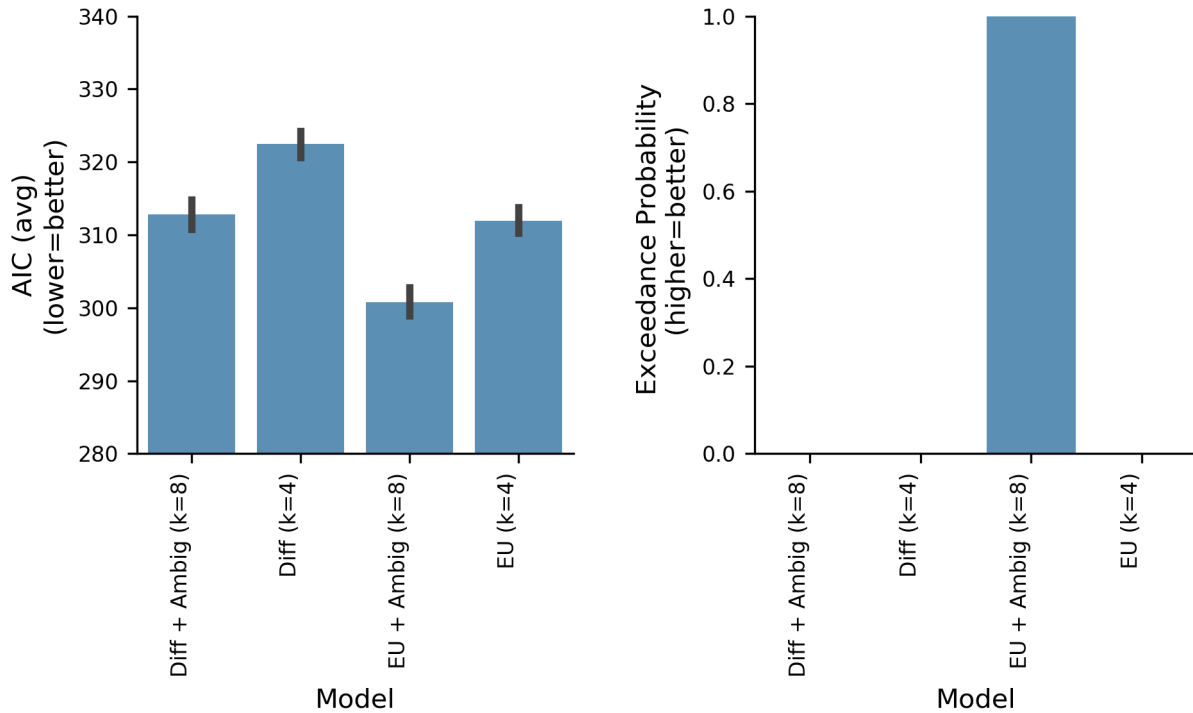
(a) Gain trials with $|EV \text{ difference} < 2|$

	% chose P1	P1	M1	P2	M2	EV1-EV2
trial #						
223	0.56	0.44	68.0	0.24	130.0	-1.28
75	0.65	0.64	44.0	0.28	106.0	-1.52
208	0.86	0.70	22.0	0.26	62.0	-0.72
119	0.85	0.82	7.0	0.46	14.0	-0.70
87	0.86	0.84	18.0	0.34	47.0	-0.86
4	0.86	0.88	7.0	0.38	20.0	-1.44
135	0.52	0.92	5.0	0.84	6.0	-0.44
269	0.71	0.92	5.0	0.04	125.0	-0.40
255	0.31	0.20	78.0	0.18	85.0	0.30
268	0.87	0.48	11.0	0.26	14.0	1.64
222	0.44	0.58	17.0	0.58	15.0	1.16
19	0.74	0.70	23.0	0.66	22.0	1.58
68	0.82	0.76	42.0	0.50	61.0	1.42
270	0.46	0.80	38.0	0.72	42.0	0.16
242	0.88	0.82	7.0	0.78	5.0	1.84
29	0.64	0.88	34.0	0.74	39.0	1.06
86	0.47	0.92	29.0	0.88	30.0	0.28
125	0.64	0.96	5.0	0.68	7.0	0.04

(b) Loss trials with $|EV \text{ difference} < 2|$

	% chose P1	P1	M1	P2	M2	EV1-EV2
trial #						
184	0.81	0.20	-78.0	0.18	-85.0	-0.30
114	0.23	0.48	-11.0	0.26	-14.0	-1.64
290	0.22	0.58	-17.0	0.58	-15.0	-1.16
233	0.15	0.70	-23.0	0.66	-22.0	-1.58
32	0.29	0.76	-42.0	0.50	-61.0	-1.42
236	0.62	0.80	-38.0	0.72	-42.0	-0.16
176	0.40	0.88	-34.0	0.74	-39.0	-1.06
83	0.48	0.92	-29.0	0.88	-30.0	-0.28
103	0.26	0.96	-5.0	0.68	-7.0	-0.04
230	0.37	0.44	-68.0	0.24	-130.0	1.28
227	0.33	0.64	-44.0	0.28	-106.0	1.52
175	0.42	0.70	-22.0	0.26	-62.0	0.72
113	0.33	0.82	-7.0	0.46	-14.0	0.70
262	0.29	0.84	-18.0	0.34	-47.0	0.86
84	0.25	0.88	-7.0	0.38	-20.0	1.44
160	0.51	0.92	-5.0	0.84	-6.0	0.44
264	0.35	0.92	-5.0	0.04	-125.0	0.40

Supplemental Table 1: Model agnostic analysis of risk aversion for gains and risk seeking for losses. Risk aversion for gains implies that when deciding between two options (both with unambiguous probabilities) of equivalent or near equivalent expected value, participants will tend to choose the urn with the higher probability (denoted as P1 above) but lower magnitude (denoted as M1 above) outcome. Risk seeking for losses implies that participants will tend to choose the urn with the lower probability (denoted as P2 above) but larger absolute magnitude (denoted as M2 above) outcome. The table above lists the 17 gain trials and 17 loss trials, constituting 25% of the unambiguous trials, which had an absolute expected value difference less than 2 ($|EV \text{ difference}| < 2$) between the urns; ($|EV \text{ difference}|$ ranged from 0.04 to 85 across all trials). (a) On 14 out of the 17 gains trials, more than 50% of participants selected the urn with the higher probability (denoted as % chose P1). Averaged across trials and participants, the urn with higher probability was selected 66.7% of the time (between participant $SE=1.1\%$). (b) On 14 out of the 17 loss trials, fewer than 50% of participants chose the urn with the higher probability. Averaged across trials and participants, the urn with higher probability was chosen only 34.6% of the time (between participant $SE=1.3\%$). The difference between the gain and loss trials for the average percentages was highly significant (binomial test; $p<0.001$). Importantly, the tendency to choose the high probability (P1) urn on gains and the low probability (P2) urn on loss trials was present even on trials where the expected value difference worked against this effect (in the table above, these trials were those where $EV1-EV2 < 0$ for gains, and those where $EV1-EV2 > 0$ for losses).



Supplemental Figure 1: Model comparison. Four models were compared using penalized fit averaged across participants (AIC; left plot; lower = better) and exceedance probability (Stephan 2009; right plot; higher = better), which better takes into account heterogeneity in fit across participants. The model used in the main text is labeled as "EU + Ambig". The model used in Lawrance et al is labeled as "Diff + Ambig" and assumes that participants separately compare the difference in magnitude and the difference in probability between the two urns (see Eqn. S1 and S2 below), rather than use the expected utility difference. We also compared a version of each model that did not have effects of ambiguity (no β_0 or β_3); these two models are labeled as "EU" and "Diff". The main effect for AIC of using expected utility ("EU", "EU + Ambig") versus probability and magnitude differences ("Diff", "Diff + Ambig") and the main effect of using ambiguity parameters v/s not were both significant ($F=34$, $p<0.001$; $F=29$, $p<0.001$; left panel). In the panel on the right, we see that the exceedance probability, that the "EU + Ambig" model is more prevalent than any of the other three models, is $p>0.99$. (k=4, for example, denotes the number of parameters).

Eqn. S1 (Diff + Ambig Model, for unambiguous trials)

$$p(\text{choice}_{\text{urn \#1}}) = \frac{1}{1 + \exp(-1 * (\beta_1(M_1 - M_2) + \beta_2(|\log|(P_1 - P_2)))}$$

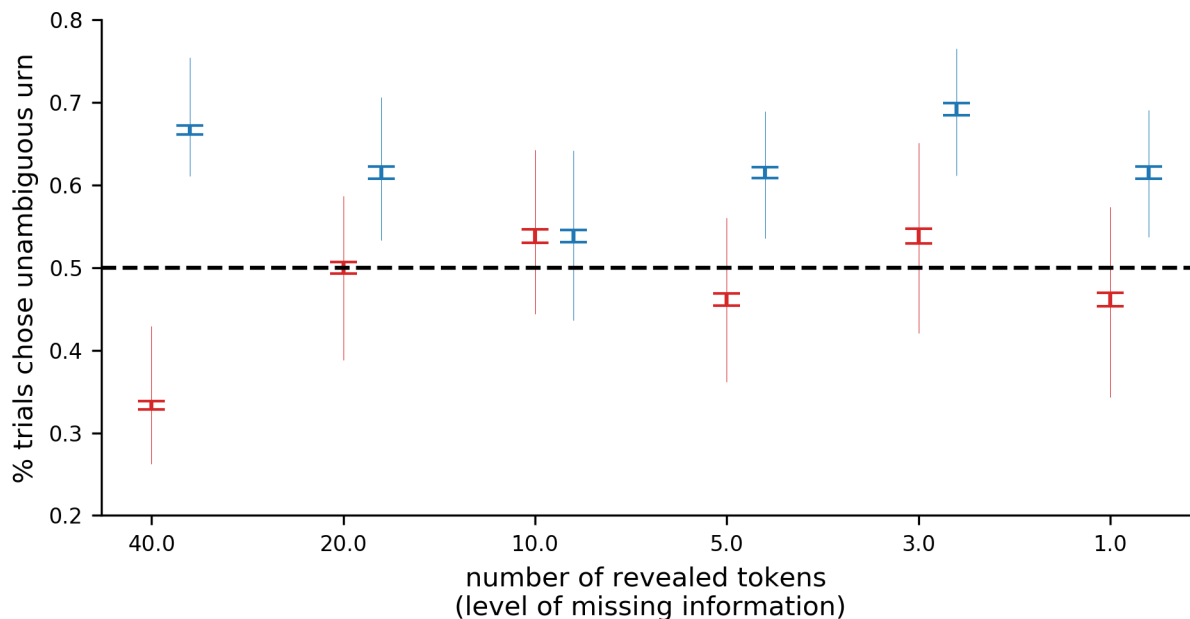
Eqn. S2 (Diff + Ambig Model, for ambiguous trials)

$$p(\text{choice}_{\text{unambiguous urn}}) = \frac{1}{1 + \exp(-1 * (\beta_1(M_a - M_u) + \beta_2(|\log|P_a - P_u) + \beta_0 + \beta_3 * MI))}$$

Note: The log modulus function ($|\log|(P_a - P_u))$) compresses larger probability differences in a way that is symmetric about zero but also signed. This was used in Lawrance et al. Removing it does not appreciably impact model fit.

negative affect (general)	0.95	0.92	0.9	0.64	0.75	0.92	0.84	0.73	0.89	0.25	0.47	0.65	0.35	0.54
depression (specific)	-0.11	-0.02	0.19	0.04	0.64	0.2	0.01	0.07	0.29	-0.21	-0.06	0.08	0.19	0.05
physiological anxiety (specific)	-0.07	-0.0	-0.11	0.57	0.04	0.02	0.52	0.02	0.09	0.17	0.17	0.1	0.0	0.11
cognitive anxiety (specific)	-0.14	0.28	0.04	0.08	0.12	-0.1	0.03	0.59	-0.04	-0.11	0.04	0.12	0.15	0.03
thought disorder (general)	0.49	0.51	0.45	0.52	0.37	0.47	0.52	0.42	0.55	0.47	0.9	0.71	0.03	0.61
mania-related (specific)	-0.08	0.02	-0.14	0.06	-0.13	0.02	0.06	-0.17	-0.0	0.86	0.2	-0.08	0.03	0.38
neg/cog symptoms (specific)	0.51	0.49	0.49	0.3	0.51	0.44	0.4	0.43	0.48	0.04	0.03	0.69	0.8	0.46
	stai_trait	stai_trait_anx	stai_trait_dep	masq_aa	masq_ad	masq_ds	masq_as	pswq	cesd	hps	olife_ue	olife_cd	olife_ia	olife_in

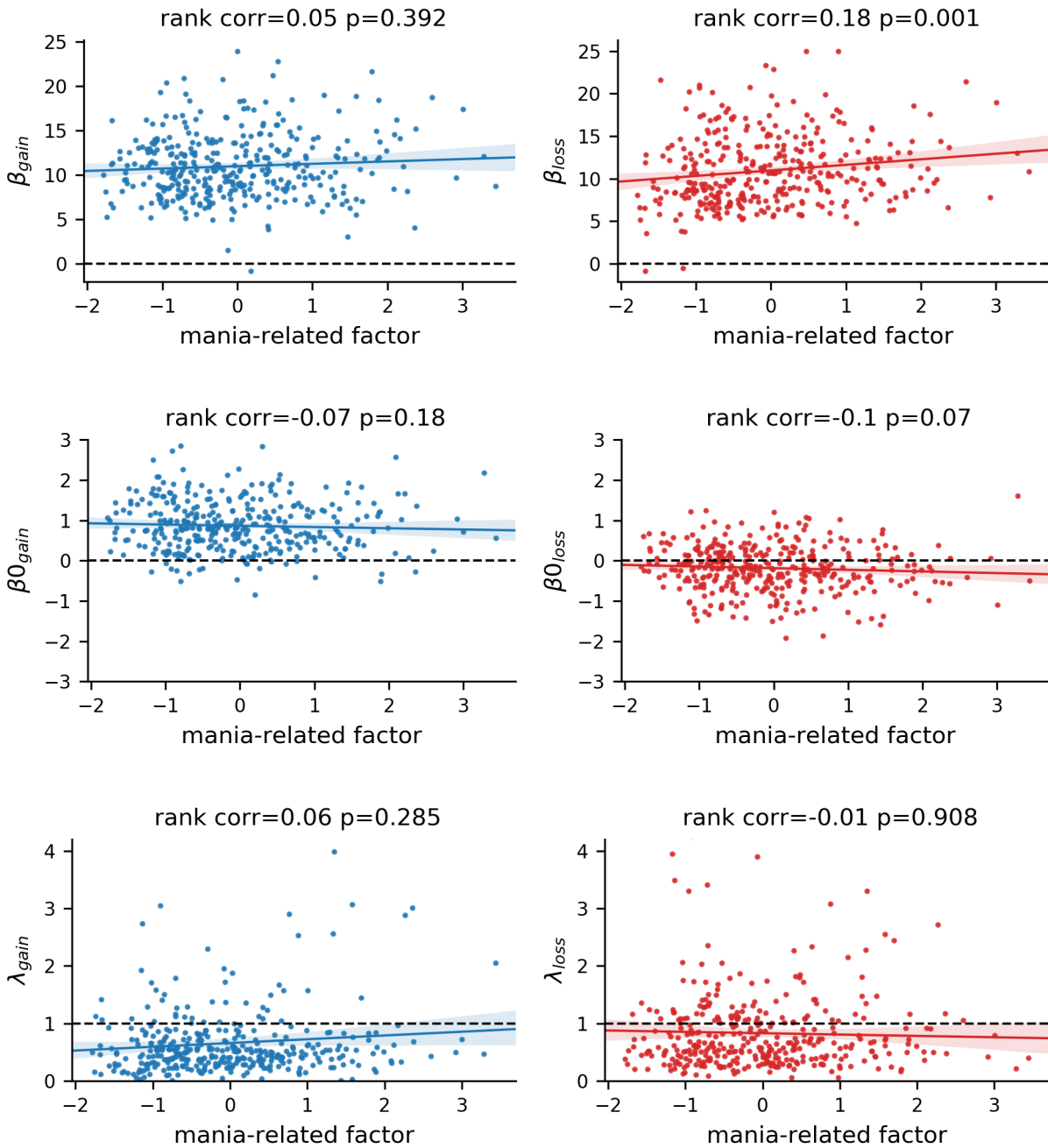
Supplemental Figure 2: Correlation between factor scores and questionnaire subscales. Two bifactor different analyses were applied to the data to separate common from unique variance (represented by general factors and specific factors, respectively). In the first bifactor analysis, participant scores for a general factor and three specific factors were estimated from several mood and anxiety related questionnaire subscales: STAI anxiety (anx), STAI depression (dep), MASQ anxious arousal (aa), MASQ anhedonia (ad), MASQ anxious symptoms (as), MASQ depressive symptoms (ds), CESD, and PSWQ. The correlations between the questionnaire subscale scores and the scores on the four mood and anxiety factors and are shown in the top four rows. A second bifactor analysis was applied to the hypomanic personality scale (HPS) and four subscales of the OLIFE (which measures symptoms related to schizophrenia): OLIFE unusual experiences (ue), OLIFE introverted anhedonia (ia), OLIFE impulsive nonconformity (in), and OLIFE cognitive disorganization (cd). A general factor and two specific factors were estimated for this second bifactor analysis. The correlation between the questionnaire subscale scores and the scores on these three additional factors are shown in the bottom three rows.



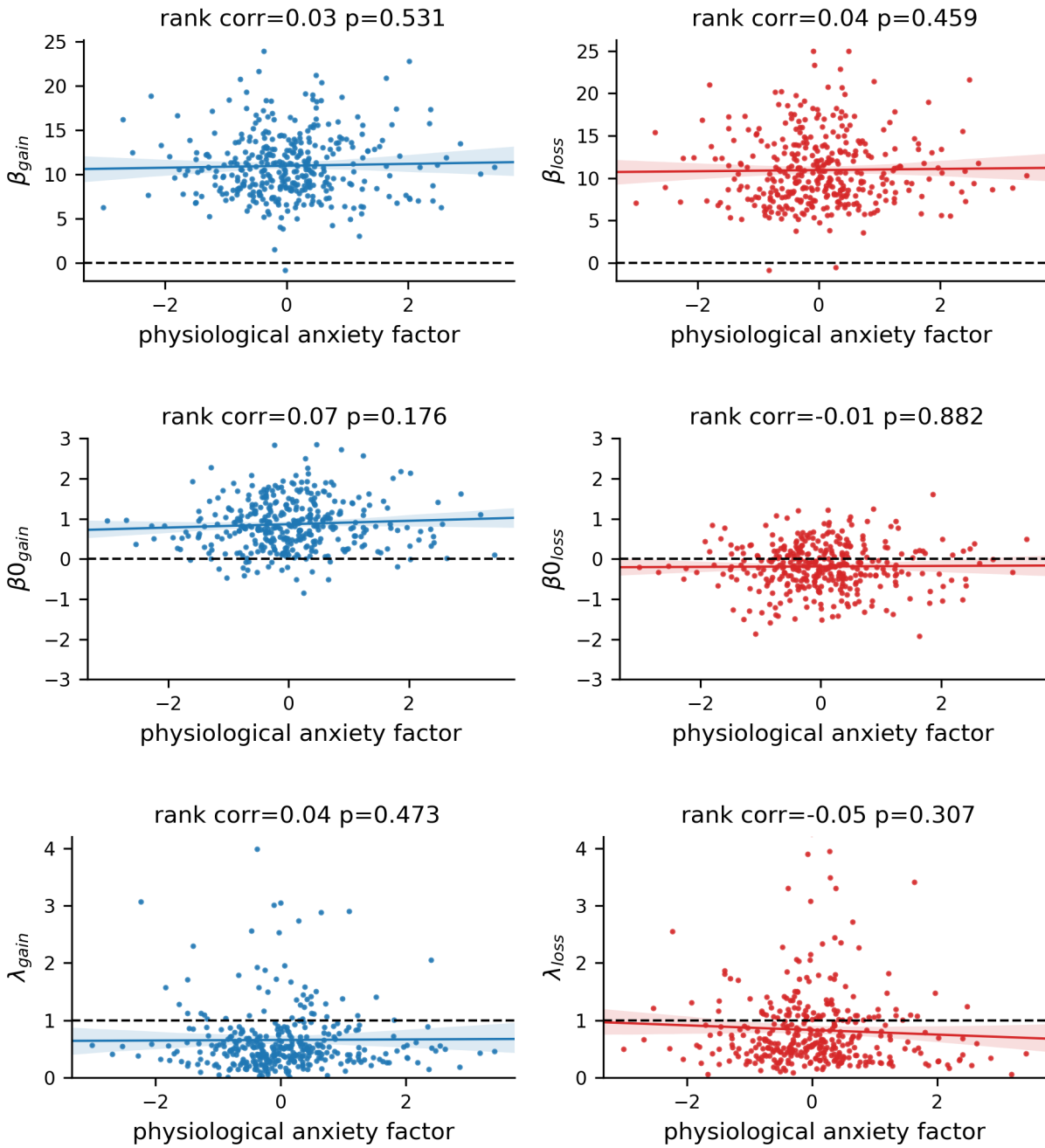
Supplemental Figure 3: Percentage of trials on which the unambiguous urn was chosen over the ambiguous urn (y-axis). This percentage provides a rough measure of ambiguity aversion (if above 50%) or ambiguity seeking (if below 50%). Percentages for each participant were averaged over trials at each level of missing information. The 1st and 3rd quartiles, which bracket the average percentages for the middle 50% of participants, are plotted as thin error bars. The standard errors for the means are plotted as thick error bars. On gain trials (in blue), the middle 50% of participants chose the unambiguous urn more often than the ambiguous urn (above dotted line) at all levels of missing information except at $n=10$ tokens revealed. On loss trials (in red), the middle 50% of participants chose the ambiguous urn more often than the unambiguous urn (below dotted line) at the lowest level of missing information (corresponding to 40 tokens revealed). However, the middle 50% of participants became more ambiguity neutral at higher levels of missing information, with some participants in the middle 50% choosing the unambiguous or and others the ambiguous urn more often. Statistical significance for differences can be inferred from the standard errors; all differences would be significant, except at $n=10$ tokens.

β_{gain}	-0.03 (0.629)	-0.01 (0.894)	-0.03 (0.531)	0.04 (0.431)	0.07 (0.189)	-0.05 (0.392)	-0.04 (0.436)
β_{loss}	-0.01 (0.921)	-0.01 (0.867)	0.04 (0.459)	-0.01 (0.816)	-0.09 (0.073)	0.18 (0.001)	0.05 (0.32)
λ_{gain}	0.03 (0.615)	-0.05 (0.36)	0.04 (0.473)	-0.01 (0.921)	0.07 (0.166)	0.06 (0.285)	-0.02 (0.648)
λ_{loss}	-0.0 (0.97)	-0.06 (0.281)	-0.05 (0.307)	-0.0 (0.938)	-0.0 (0.938)	-0.01 (0.908)	-0.09 (0.079)
β_0_{gain}	0.01 (0.839)	-0.04 (0.455)	0.07 (0.176)	-0.03 (0.618)	0.02 (0.67)	-0.07 (0.18)	-0.04 (0.399)
β_0_{loss}	-0.03 (0.531)	0.01 (0.88)	-0.01 (0.882)	-0.05 (0.339)	-0.01 (0.804)	-0.1 (0.07)	-0.01 (0.834)
β_3_{gain}	0.02 (0.659)	0.07 (0.197)	0.12 (0.024)	-0.01 (0.808)	-0.04 (0.428)	-0.1 (0.05)	0.02 (0.642)
β_3_{loss}	-0.04 (0.503)	0.01 (0.892)	0.12 (0.029)	-0.03 (0.591)	0.03 (0.586)	-0.01 (0.819)	-0.01 (0.834)
	negative affect (general)	depression (specific)	physiological anxiety (specific)	cognitive anxiety (specific)	thought disorder (general)	mania-related (specific)	neg/cog symptoms (specific)

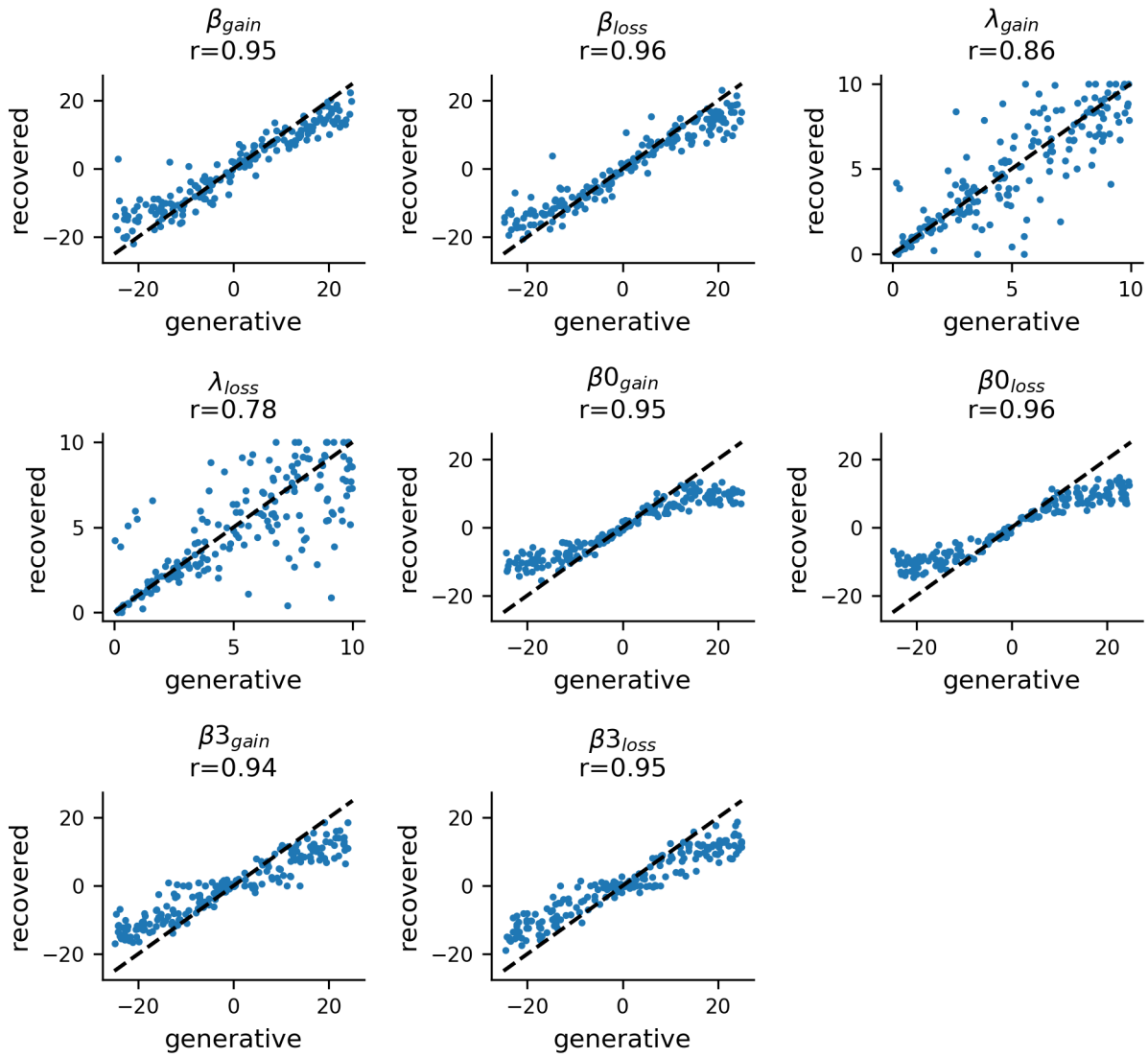
Supplemental Figure 3: Rank Correlations between model parameters and symptom factors. The x-axis contains the parameters from the model. ($\beta_{gain}, \beta_{loss}$) are the inverse temperature parameters, ($\lambda_{gain}, \lambda_{loss}$) are the risk parameters, ($\beta_0_{gain}, \beta_0_{loss}$) are the categorical (or average) ambiguity parameters, and ($\beta_3_{gain}, \beta_3_{loss}$) are the information-level dependent ambiguity parameters. The y-axis contains the symptom factors. The rank correlations and p-values (in parentheses) between the symptom factor scores and model parameters are displayed in each cell. Significant rank correlations (uncorrected for multiple comparisons) are shown in red or blue, depending on whether the correlations are positive or negative. The effects related to ($\beta_3_{gain}, \beta_3_{loss}$) are discussed in the main text. The significant positive correlation between the *mania-related* factor and β_{loss} mean that higher scores are associated with choices that are better predicted by the model on loss trials.



Supplemental Figure 4: The relationships between scores on the *mania-related* factor and model parameters. These plots were included to show the range for parameters that were not shown in the main text.



Supplemental Figure 5: The relationships between scores on the *physiological anxiety* factor and model parameters. These plots were included to show the range for parameters that were not shown in the main text.



Supplemental Figure 6: Parameter recovery for main model. A range of different parameter values were selected at random and used to generate new choice data from the model. The same procedure that was used to estimate parameters from the actual participants' data was used to estimate parameters from the simulated data. The parameters values used to simulate data (referred to as generative on x-axes above) were strongly correlated with the parameter values that were estimated from the simulated data (referred to as recovered on y-axes above).

Chapter 4: Prior Beliefs and Belief Updating

Introduction

In both clinical and subclinical populations, individuals with depressed mood show reduced estimates of the likelihood that positive things will happen to them in the future, and increased estimates of the likelihood that negative things will happen to them in the future, relative to non-depressed individuals (Butlers & Matthews, 1983). Individuals with high levels of anxiety have also been reported to show increased estimates for the likelihood of future negative events happening to them, with there being more mixed evidence for decreased estimates regarding future positive events (Butlers & Matthews, 1983).

But why should adults with high levels of anxiety or depression show systematically altered, self-referential judgements? One possibility is that they developed and then stabilized dysfunctional schema, likely in childhood, and so have inflexible negative self-referential beliefs that are applied automatically and globally to many different situations (Beck 1976; Abramson et al., 1989). These would act as priors. For example, before a new job interview, an individual with depression might reflexively think “I will never be a success”, negatively biasing their expectations. These negative priors may or may not reflect actual differences in the experience of adverse life events. Another possibility is that they have biased updating in the light of unbiased experience, for instance weighing the news of not getting a job more heavily than the news of getting it.

There has been some investigation of biased weighting of positive versus negative information during belief updating with regards to potential future events (Sharot et. al, 2011), however, various criticisms (Harris et al, 2011; 2017; Shah et. al, 2016) make it desirable to investigate beliefs for which there is a clear ground truth and experimental control over the information concerned. As a step in this direction, Eil & Rao (2010) asked healthy participants to estimate their true rank relative to other participants for performance on an IQ test and how highly others rated them in terms of physical attractiveness. They then gave participants objective feedback, whether or not they ranked higher or lower than another randomly selected participant, and then asked them to update their beliefs. Participants incorporated negative feedback into their beliefs to a lesser extent and less reliably than positive feedback. In a similar design, Mobius et al. (2010) gave participants probabilistic feedback about whether they were in the top half of performers on an IQ test, and similarly observed that participants asymmetrically update more for positive than for negative information and are conservative in updating relative to a Bayesian reference point (Mobius et al., 2010).

Our current study was informed by the work of Eil & Rao (2010) and Mobius et al. (2010) but extended this previous work to investigate the basis of negative self-referential beliefs in individuals with high levels of anxiety and depression. We employed a novel experimental paradigm, which we hoped would be highly ecologically valid for our student participants by virtue of investigating beliefs that were likely to have direct importance, namely beliefs about how they compared against their peers in a competitive hypothetical internship. Ecological validity is critical since it has also been argued that negative biases associated with depression

are more strongly evident in studies with high ecological validity (Dobson & Franche, 1989). Indeed, in both Eil & Rao (2010) and Mobius et al. (2010), the positive bias observed in belief updating was more prominent for beliefs about the self than for neutral beliefs, such as the performance of a robot on the same IQ test.

We used this new paradigm in conjunction with a refined, bifactor analysis that better disentangles the substantial comorbidity of affective disorders to investigate whether elevated levels of depression and/or anxiety were associated with negative initial beliefs, prior to the receipt of any feedback, and whether these prior beliefs were subject to revision in response to concrete information. We further investigated whether belief updating, if observed, was biased—that is whether anxiety and/or depression were linked to greater belief change following negative versus positive feedback.

Results

During the experiment, participants were asked to imagine themselves in a hypothetical internship, vying against one another to be selected as partners. In the first session, participants made profiles about themselves, which consisted of their actual grades, SAT scores, and a brief description of why they would be a good candidate for the internship. In the second session, participants were shown pairs of other participants' profiles and were asked to choose with whom they would rather work of each pair. In the third and final session (depicted in **Figure 4.1**), they were shown the results (i.e. feedback), one pair at a time, of whether other participants had chosen to work with them or not. Beliefs about the likelihood of being in the top (or bottom) half of participants, in terms of how often they were selected to work with, were elicited from participants before any feedback (i.e. prior beliefs) and following each piece of feedback (i.e. belief updating).

Participants were also administered several questionnaires, which were selected to measure a wide range of mood and anxiety symptomatology (see **Methods** for a list). We analyzed participants' item-level responses on these questionnaires using a variant of the bifactor model that we had previously used (**Chapter 2**; see **Methods** for details about this model and the statistically minor variations associated with the use of a slightly different collection of questionnaires). This bifactor model, consistent with others (see Clark et al., 1994 for anhedonia and Brodbeck et. al, 2011 for worry), separates out a general factor, representing shared variance, and two specific factors, one for depression-specific variance (i.e. anhedonia) and one for anxiety-specific variance (i.e. worry). We calculated scores for participants on each of these three factors to test whether differences in initial beliefs or belief updating were associated uniquely with depression, uniquely with anxiety, or with both. The correlations between the factor scores and the full set of questionnaire subscales can be found in **Supplemental Figure 4.1**.

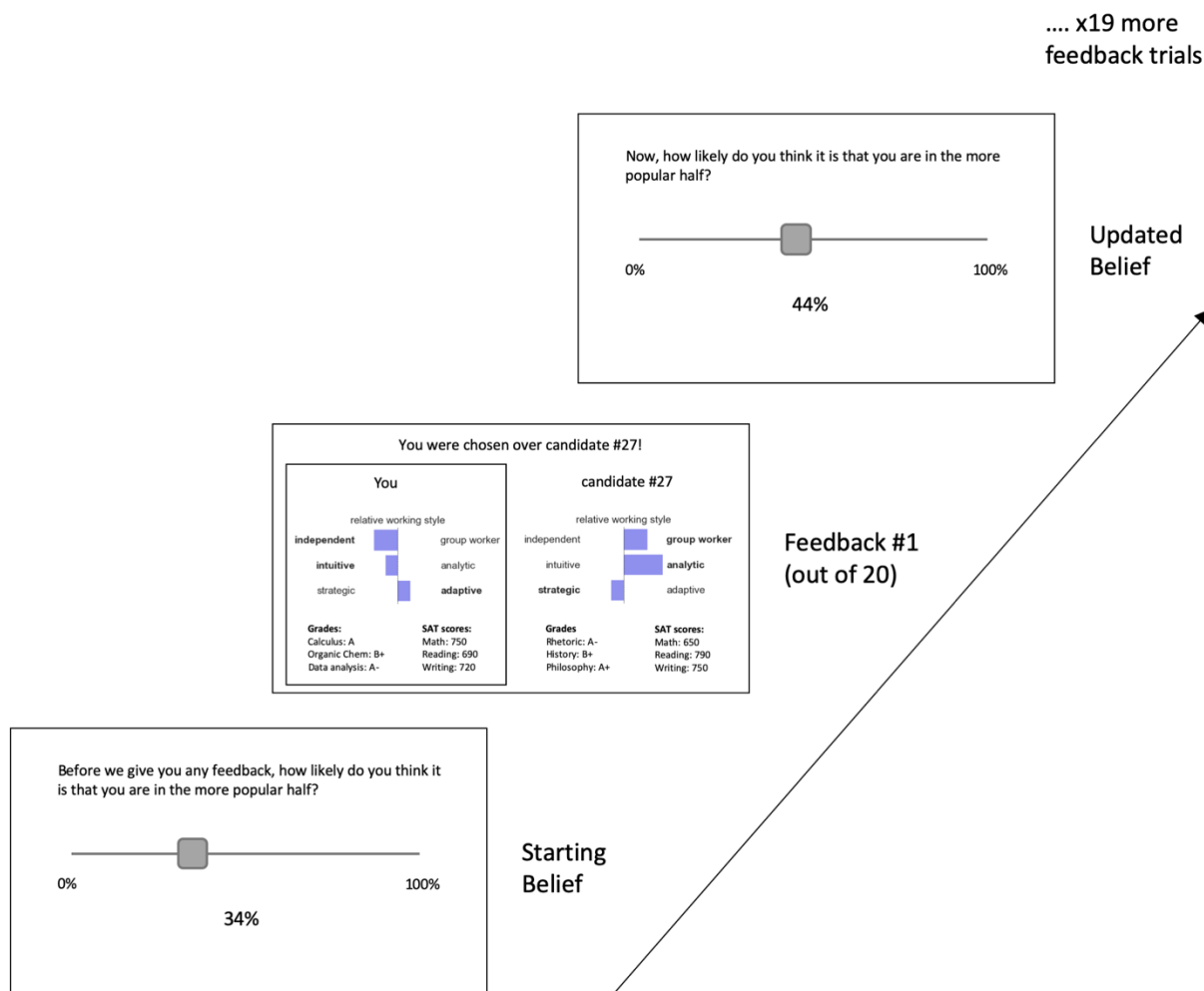


Figure 4.1: Experimental Session 3. Participants were first asked to estimate their *likelihood of being in the more popular half of students*; these estimates were considered their starting (i.e. prior) beliefs. Participants reported their beliefs using a “slider” that when moved presented values from 0% to 100% on the screen. Participants were then shown a pair of profiles, containing their profile and another participant’s profile, and whether or not they had been chosen by a third participant during session 2 (depicted by a black outline, as shown above). After the receiving this feedback, participants were again asked to estimate their likelihood of being in the more popular half of students. Participants were given twenty pieces of feedback in total and asked to update their belief each time.

Model Agnostic Analyses of Beliefs and Belief Updating

Starting Beliefs

Before participants were shown any feedback, we asked them to estimate *their likelihood of being in the more popular half of students* (i.e. their starting or prior belief). Scores on the depression-specific factor were negatively correlated with starting belief ($r=-0.36$, $p=0.003$). This meant that individuals with high depression factor scores were more likely initially to endorse the belief that they were among the 50% of participants least selected as project partners. Neither the general factor nor the anxiety-specific factor scores correlated

with starting belief (**Figure 4.2**). For the interested reader and to aid comparison across studies, we report correlations with individual questionnaire scales in the Supplementary Materials (see **Supplemental Figure 4.2**). These analyses, together with the correlations between factor scores and questionnaire measure (**Supplemental Figure 4.1**) suggest that more negative starting beliefs relate predominantly to anhedonic depressive symptomatology.

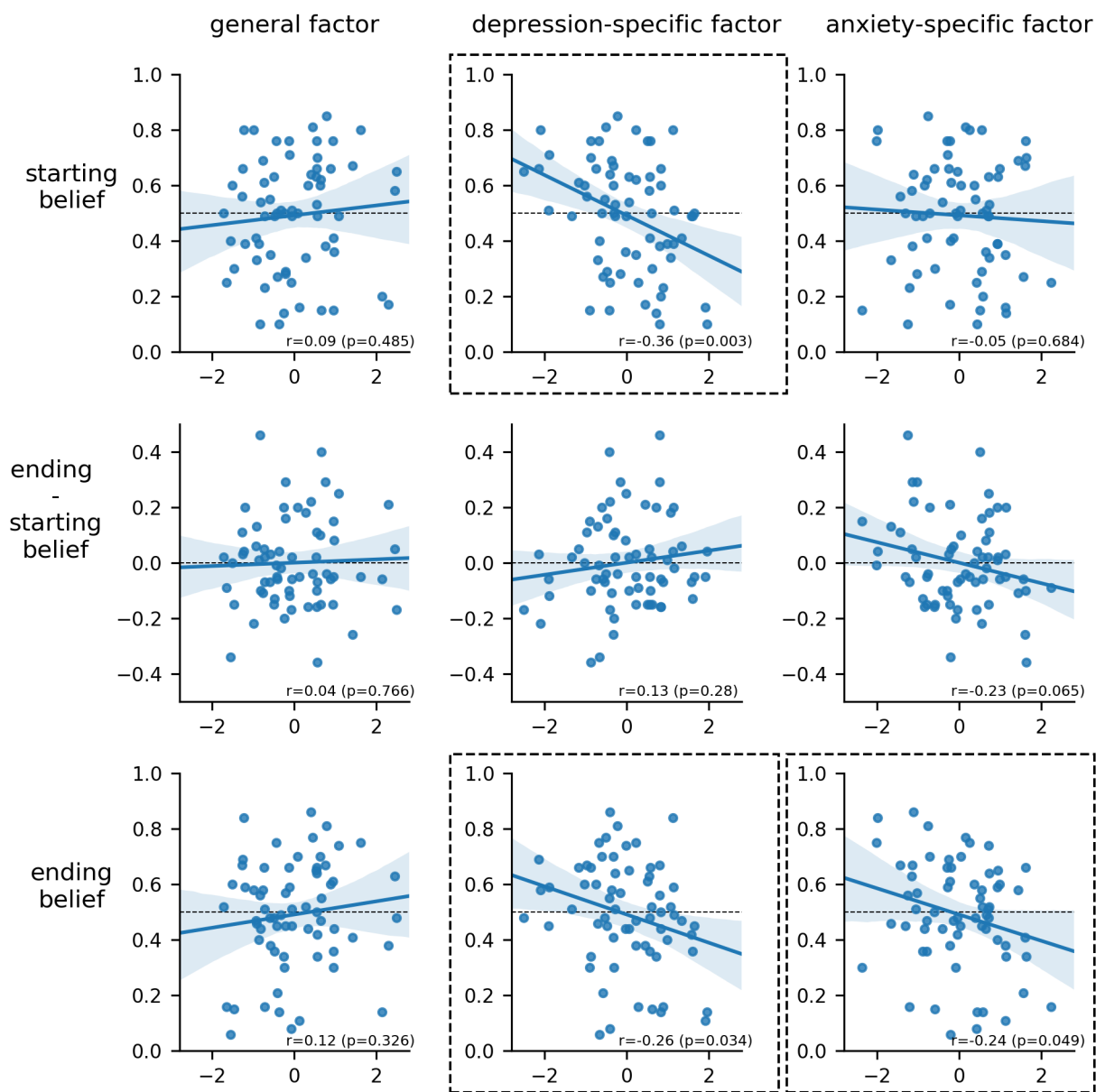


Figure 4.2: Participants' standardized scores (x-axes) for the general factor, anxiety-specific factor and depression-specific factor (columns) are plotted against their reported beliefs at the start of the feedback period, the end of the feedback period, and the difference between beliefs at the start and the end (y-axes and rows).

Starting Beliefs and Ground Truth: Depressive Realism?

Depressive realism is the hypothesis that the negative beliefs held by individuals with depression are actually more accurate, as opposed to being negatively biased, than the positively biased belief held by healthy individuals (Alloy & Abramson, 1979). We can examine this by measuring the number of times that a participant was actually chosen as a potential partner in session 2. This gives us an index of actual 'profile competitiveness'. Scores on the depression-specific factor were weakly correlated with profile competitiveness ($r=-0.21$, $p=0.1$). However, profile competitiveness did not mediate the relationship between scores on the depression-specific factor and the negative starting beliefs ($p=0.4$). The lack of mediation can be attributed to the lack of correlation between profile competitiveness and starting belief across participants ($r=-0.07$, $p=0.52$). There was also no relationship between anhedonia and accuracy at predicting one's own profile competitiveness (see **Supplemental Figure 4.3**).

Updating of beliefs following feedback

After reporting their initial belief about the *likelihood of being in the more popular half of students*, participants saw a balanced set of ten instances of positive feedback and ten instances of negative feedback, selected from a larger set of possible feedback available from session 2 and delivered in one of two possible pre-randomized sequences. Participants varied in amount that they updated their beliefs in response to this feedback. The average absolute belief change, following each instance of feedback, ranged between 1-10% (mean=4.86%) across participants. Six participants fell outside this range (four participants updated between 10-20% and two participants updated between 0-1% on each trial). None of the three factor scores correlated with the average absolute value of the trial-to-trial updating ($p's > 0.19$).

We next looked at participants' change in belief from start to end. On average, participants changed their beliefs in one direction or the other by an overall amount of +/- 12% (std=9%) by the end. Because the feedback was balanced, we expected participants to move their estimates towards 50%. Indeed, by the end, participants reported beliefs somewhere between their starting value and 50% (see **Figure 4.3**). Neither scores on the depression specific factor nor scores on the general factor were significantly correlated with the difference between ending belief and starting belief ($r=0.04$, $p=0.77$; $r=0.13$, $p=0.28$; **Figure 4.2**). There was a trend towards a negative relationship between scores on the anxiety-specific factor and change in beliefs from start to end of the feedback period, but this did not reach significance ($r=-0.23$, $p=0.065$).

By the last feedback presentation, anxiety-specific scores were significantly negatively correlated with the reported beliefs ($r=-0.24$, $p=0.049$) and depression-specific scores were still significantly negatively correlated with the reported beliefs ($r=-0.26$, $p=0.034$). In the case of anxiety, this relationship was absent ($r=-0.05$, $p=0.68$) at the start of the feedback period and hence was driven at least to some extent by the marginally significant change in belief as a result of feedback. In the case of depression, the persistence of the relationship likely reflects the fact that most participants did not substantially move their beliefs from start to end; indeed, ($n=43$) participants ended with a belief on the same side of 50% as their starting belief (these participants can be seen in bottom-left and top-right quadrants in **Figure 4.3**). General factor scores were not correlated with ending belief (**Figure 4.2**).

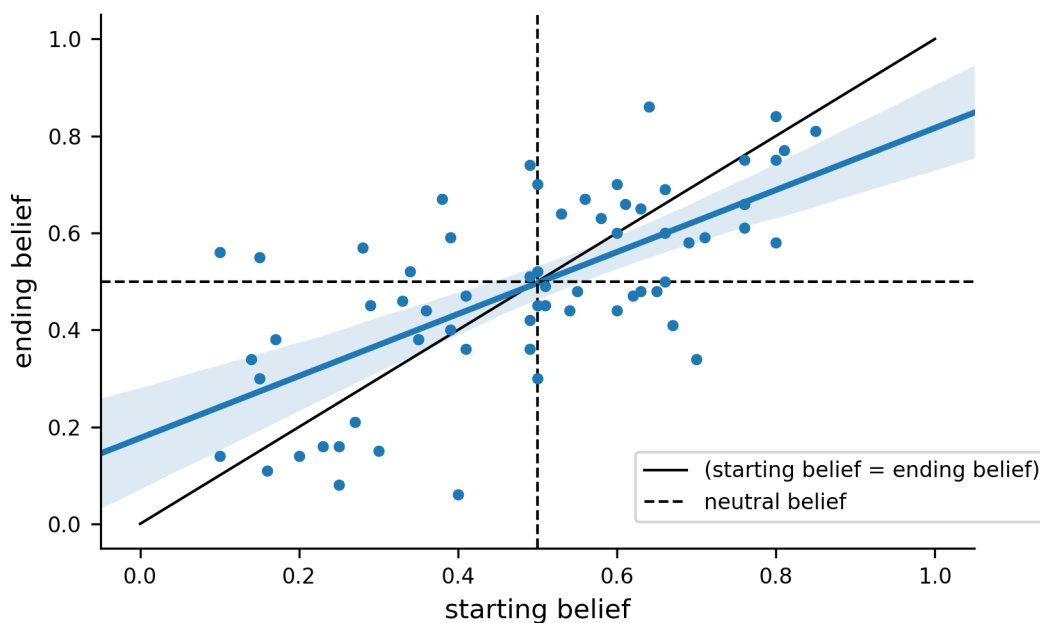


Figure 4.3: Regressing participants' starting beliefs against their ending beliefs shows that participants update in the direction of but not completely to a neutral belief (i.e. 50%) after receiving twenty instances of balanced feedback (10 positive and 10 negative instances). The slope of the regression line was significantly greater than 0 and significantly less than 1, indicating partial updating towards 50%.

Effects of Feedback Order

There were two distinct sequences of feedback that different participants saw. The positive-first feedback sequence started with positive feedback for the first two trials and had a total of six positive feedbacks in the first ten trials. The negative-first feedback sequence was exactly the opposite. There was a significant difference in the ending minus starting belief between the groups of participants receiving different feedback sequences ($t=3.36$, $p=0.001$). Participants who received positive feedback first shifted their beliefs more in the positive direction from start to end. Across the entire participant sample, however, these biases seemed to cancel out; there was no significant difference in the ending minus starting belief on average ($t=-0.01$, $p=0.98$). Moreover, feedback order did not interact significantly with any of the symptom factors in their influence on aggregate bias ($p's > 0.6$).

Other-related beliefs

After engaging in the task described above, participants underwent the same procedure for updating beliefs about another randomly chosen participant. They were first shown the anonymized profile and asked to estimate the *likelihood that the other participant is in the more popular half of students*. They were then given ten instances of positive and ten instances of negative feedback (i.e. profile pairings where the other individual had or had not been

selected) in one of two new possible sequences. None of the three factors (general, anxiety or depression) were correlated with other-related starting beliefs, ending beliefs or starting minus ending beliefs (see **Supplemental Table 4.1**).

Modeling Starting Beliefs and Bias in Belief Updating

The model-agnostic measures in the previous section characterized individual differences in belief updating using two point-estimates given by the subjects: one at the start and one at the end of the feedback period. In these, participants reported the likelihood that they were in the top or bottom half of potential teammates. In this section, we analyze the trial-to-trial updating of beliefs using models that have to make a few additional assumptions about the nature of peoples' beliefs and the way that they update them, but can thereby draw conclusions based on the entire sequence of participant reports.

To model belief updating, we assume that participants report an estimate for the probability that they were in the top half of participants based on a distribution of possible values that they consider plausible. We then assume that participants update this belief distribution in response to feedback using one of four models, based on either a more exact or a more approximate Bayesian inference, and with or without a bias. The more exact Bayesian inference (which we just call 'Bayesian') integrates prior expectations with the information on the current trial (i.e. the likelihood) to generate a posterior, which then becomes the prior for the next trial. This produces distributions of belief that become narrower as information is progressively incorporated, and so the updates become smaller over trials. The more approximate form of inference uses the Rescorla-Wagner (RW) rule, which can be derived as a form of Bayesian inference but with a fixed update size. This is associated with distributions of belief that have a constant width.

We adjusted both Bayesian and RW updating to account for any additional potential bias in responding to positive or negative feedback. For Bayesian updating, the bias was parameterized as ω and $1/\omega$ for positive and negative feedback (a value of $\omega = 1$ corresponds to a lack of bias). For the RW updating, the bias was parameterized by changing the feedback values from $\{0,1\}$ for negative and positive feedback to $\{0, r\}$ ($r = 1$ also corresponds to unbiased updating).

We compared the resulting four models according to their exceedance probabilities (Stephan et al., 2009). The biased RW model provided a conclusively better fit to the data than the other three models (exceedance probability > 0.99). It also provided a better qualitative fit to the data; simulating data from both models showed that the biased RW model better matched the distribution of average (per participant) belief updates following positive and negative feedback (see **Supplemental Figure 4.5**). It also had a higher exceedance probability than a variety of other alternative models (**Supplemental Figure 4.7**), which we nevertheless excluded from our main analyses, because model recovery simulations showed that they were not suitably distinguishable from one another in our task.

In the biased RW model, the parameter for the mean of the belief distribution prior to feedback, μ_0 , was significantly correlated with both the starting belief reported by participants ($r=0.97$, $p<0.001$) and the depression-specific factor scores ($r=-0.36$, $p=0.003$; **Figure 4.4**),

matching the model-agnostic results of the previous section. The bias parameter r significantly correlated with the model-agnostic measure of belief updating (ending-starting beliefs; $r=0.51$, $p<0.001$). The bias parameter r was also significantly negatively correlated with the anxiety-specific factor scores ($r=-0.32$, $p=0.008$), meaning that high anxiety scores were associated with negatively biased updating, an effect that we observed at a trend-level of significance when using the model-agnostic measures of the previous section. The relationships between the factor scores and the two other parameters of the biased RW model were non-significant (see **Supplemental Figure 4.4**).

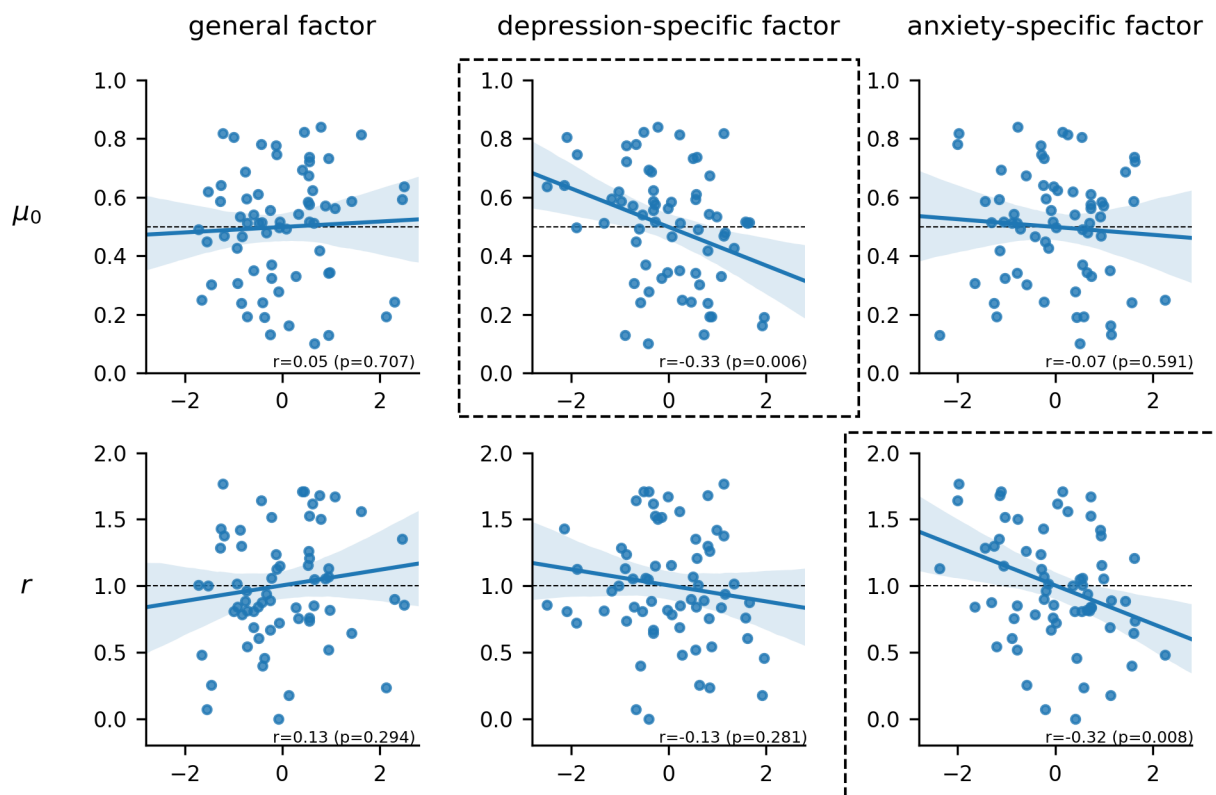


Figure 4.4: Participants' standardized scores (x-axes) for the general factor, anxiety-specific factor and depression-specific factor (columns) are plotted against the estimated parameters in the biased RW model. μ_0 is the estimated mean of the participant's belief distribution at the start of the feedback period and r represents a bias in updating in response to feedback ($r > 1$ corresponds to a positive bias, whereas $r < 1$ corresponds to a negative bias).

Discussion

In this study, we looked at whether participants with high levels of anxious and depressive symptoms reported more negative beliefs prior to receiving any feedback, whether they updated those beliefs in response to concrete feedback, and if so, whether they exhibited a bias in updating following negative or positive feedback. We devised a novel experiment task resembling a scenario that we expected our undergraduate participants to find to be important, namely a hypothetical internship in which they were compared to their peers. We observed that prior to receiving feedback, individuals with high scores on a depression-specific factor reported beliefs about themselves that were more negative than the beliefs reported by others. This relationship continued to persist until the end of the feedback period, with no differential updating to positive or negative feedback (i.e. no bias) associated with depression-specific scores. Individuals with high scores on the anxiety-specific factor, on the other hand, were associated with biased updating in our computational model, updating to a greater extent following negative rather than positive feedback. By the end of the experiment, this bias resulted in significantly more negative beliefs being reported by individuals with high scores on the anxiety-specific factor.

The association between depression-specific factor scores and the negative beliefs, reported prior to receiving feedback, raises the question of where these beliefs come from. Although our current study cannot fully address the origins of these beliefs, classic cognitive theories (Beck 1976) would argue that they are generalized from more global beliefs that participants might have about themselves. Depression has been associated with negative, non-specific beliefs, such as “no one really likes me” or “I will never be a success”, which may be being relied heavily upon prior to feedback. Depression has also been associated with negatively biased retrieval from memory (Mathews & MacLeod, 2005), which may also have impacted beliefs at the start of the experiment.

That negative initial beliefs were specific to depression and not related to anxiety-specific scores or general factor scores is noteworthy and may be because the content of the beliefs in this experiment were strongly tied to conceptions of self-worth and competence. Although negative beliefs and schemas are thought to underly both anxiety and depression, the content of those schemas is thought to differ, with failure being one of the central themes in depressive schemas and vulnerability to threat being a central theme in anxious schemas (Beck 2005; Clark & Beck, 2010). It would be interesting to conduct a similar experiment that targeted anxious schemas, perhaps by shifting the focus to something more akin to threat, or from past to future rejection by peers, to see if the relationship with initial beliefs switches to anxiety.

During belief updating, anxiety-specific factor scores, rather than depression-specific scores, were associated with biases in responding to negative versus positive feedback, with individuals with high levels of anxious-specific symptomology exhibiting larger updates following negative than positive feedback. This finding is generally fits with proposals that anxiety is characterized by hypervigilance to external threats (Barlow 2000), and the finding that children who are vulnerable to developing anxiety disorders are more likely to exhibit ‘monitoring’ behavior, seeking (as opposed to avoiding) information about potential danger (Muris et al., 2000). It would be interesting for future work to investigate if those who exhibit monitoring behavior are also those who exhibit biased updating.

It would be also interesting in future work to investigate how beliefs change over the course of longer time scales. Individuals with levels of depression symptoms, might for example, have their beliefs drift back to their initial negative state, even if they have been updated to a more neutral value during a particular situation (like our experiment). This might be hypothesized to occur given the biases in memory retrieval associated more strongly with depression than anxiety, as mentioned above (Mathews & MacLeod, 2005).

In our previous work, we observed that scores on the general factor, but not the anxiety-specific factor nor the depression-specific factor, were related to a deficit in learning rate adjustment to volatility (**Chapter 2**). In contrast, we observed no association between general factor scores and behavior in the current study. Instead, we observed that the depression-specific scores were related to biases in initial belief and the anxiety-specific scores were related to biases in belief updating. Both studies together provide evidence for a triple dissociation between different aspects of symptomology and different biases. Investigating all three in the same study, such as one in which a person's popularity in the hypothetical internship is volatile, would be interesting for future work. Furthermore, this dissociation was likely made more easily detectable by the use of an extremely similar bifactor model (modified only slightly to accommodate differences in the number of items) in both studies. This suggests that in addition to being a useful tool for understanding which aspects of symptomatology is shared versus unique (Clark & Watson, 1994; Steer et al., 1995; Brodbeck et al., 2011), bifactor models may also be useful as a general tool for categorizing potential biases and deficits along similar lines.

Eil & Rao (2010) and Mobius et al. (2010) did not investigate how symptoms of anxiety and depression impact initial beliefs or belief updating. However, they did look at belief updating in undergraduate participants in an experiment that was similar to ours. In contrast to our study, both previous studies observed that participants, on average, showed a positivity bias during belief updating. However, these two previous experiments differed from ours in a number of key respects, which may have given rise to differences in average biases. Eil & Rao (2010) asked participants to assign individual probabilities, required to sum to one, to each possible rank they could be in (rank #1 through #10), whereas we asked participants for a point estimate for the probability that they were in the top (or bottom) half. Asking participants to visualize and manually edit their entire belief distribution may have led participants to update their beliefs very differently. Mobius et al. (2010) also asked participants for a point estimate, similarly to our study, but delivered a probabilistic signal as feedback (that was correct 75% of the time) for whether a participant was in the top or bottom half. Another key difference between our study and the two previous studies is the content of the beliefs that participants were updating. Both Eil & Rao (2010) and Mobius et al. (2010) looked at beliefs about IQ, and Eil & Rao (2010) additionally looked at beliefs about beauty. Our study looked at beliefs about competitiveness in a hypothetical internship. Although undergraduates may find IQ, beauty, and competitiveness in a workplace-like setting to be similarly important, they may not have had much experience with the latter. A lack of experience may have reduced a potential optimism bias on average in our study. Finally, there may also have been differences in the overall levels of anxiety or depression symptomatology between participants in our study and theirs, which we could not compare directly.

Methods

Participant Recruitment

Participants were UC Berkeley students recruited during three separate semesters, Spring 2018, Fall 2018, and Spring 2019. Each session of the experiment had different numbers of participants, because only students passing quality assurance checks and who desired to continue in the experiment were advanced to subsequent sessions. In Spring 2018, 110 participants completed session 1. 73 participants wanted to come back for session 2, but only 64 were invited (9 excluded because of failed multiple catch questions). 41 participants came back to complete session 2, all of whom were invited back for session 3. 26 participants came back to complete session 3. In Fall 2018, 118 participants completed session 1. 74 participants wanted to come back, but only 66 were invited (so 8 excluded because of failed multiple catch questions). 41 participants came back to complete session 2, all of whom were invited back for session 3. 17 participants came back to complete session 3. In Spring 2019, 121 participants completed session 1. 92 participants wanted to come back, and all were invited (no one failed multiple catch questions). 63 participants came back to complete session 2, and all were invited back for session 3. 32 participants came back to complete session 3.

Experimental Task

Session 1: Participants construct a personal profile

UC Berkeley students enrolled in our study by using the Psychology department's Research Participation Program. In session 1 of the experiment, participants visited the website that we used for the experiment. On the website, participants were asked to read about and consent to the experiment, following which they were presented with a backstory about the 'hypothetical internship' that they would be a part of for the experiment. Participants were told that they would be selecting amongst one another to choose teammates for the internship project and that participants would use personal profiles to select amongst each other. In this session, participants answered questions about their working style, their grades and SAT scores, and wrote a brief three sentence explanation about why they should be chosen to work with. We told participants that we would use this information to construct a profile for them that would be shown to the other participants. Participants also filled out the Penn-State Worry Questionnaire (PSWQ) and Dysfunctional Attitudes Scale (DAS). They were told (truthfully) that the answers to these questionnaires would not be shown to other participants.

Session 2: Participants choose who to work with

Participants who wanted to continue the study and who passed quality assurance checks, returned to the website 1-3 weeks later. They viewed 26 pairs of profiles from the other participants and while viewing each pair chose one of the two people that they'd prefer to work with in the hypothetical internship. They also filled out several more mood and anxiety symptom questionnaires, the MASQ, the STAI, and the CESD.

Session 3: Participants receive feedback and update their beliefs

Participants who wanted to continue the study and who passed additional quality assurance checks, were invited into the lab to complete the final part of the experiment 2-6 weeks after session 2. They were told that their profile along with a competitor's profile was shown to the other participants in the experiment and each time either they or the competitor had been selected. They were told that their goal is to estimate *their likelihood of being in the more popular half of students* (in terms of how many times they were chosen to work with in session 2). Some participants were asked about the being in less popular half instead of the more popular half. Participants were first asked to estimate this probability just based on how confident they were about their profile, not knowing anything about how many times they had been chosen to work with but having seen profiles of 52 other participants in the second part of the experiment. Then, they were given 20 instances of positive or negative feedback and asked to update their estimate each time. For each instance of feedback, they were shown their profile and the competitor's profile and which profile was selected by a third participant, along with a statement such as *Classmate #29 compared you and classmate #14. You were chosen to work with!* Participants were shown actual pairs of profiles and choices made by other participants. However, the sequence of pairs shown were chosen to be balanced (there were 10 positive and 10 negative instances of feedback), and they were presented in only one of two predetermined sequences. Participants all repeated the same procedure for another participant's profile after doing it for their own profile, judging the probability that this other participant is in the top half. During session 3, participants also filled out the STAI and CESD for a second time.

Participant Exclusion

Seventy-five (n=75) UC Berkeley students voluntarily completed all three experimental sessions. Prior to conducting analyses, we excluded participants who had poor data quality. To check whether participants were following the instructions about reporting their beliefs, we looked at the percentage of trials on which the participant updated his/her belief in the correct direction or kept his/her belief the same. 48 participants made a directionally wrong update on 2 trials or fewer. Another 25 participants made between 3 and 8 directional errors. Three participants made 9 or more directional errors and were excluded. One participant was additionally excluded because he reported misunderstanding the task in the post-experiment interview. To check whether participants were paying attention to the self-report questionnaires, we looked at 'catch questions'. Catch questions instructed the participant to select a particular response (e.g. 'sometimes') from the set of available responses (e.g. 'never', 'sometimes', 'almost always', 'always'). Four participants did not answer at least 2 out of the 4 of the "catch questions" questions correctly. One participant answered all 2's for the MASQ. This left a total of 66 participants whose data we analyzed in detail.

Mood and Anxiety Symptom and Trait Measures

The full set of questionnaires includes: the Penn-State Worry Questionnaire (PSWQ) (administered in session 1), the Dysfunctional Attitudes Scale (DAS) (administered session 1), the Mood and Anxiety Symptom Questionnaire (MASQ) (administered in session 2), the Spielberger Trait-State Anxiety Inventory (STAI) (administered in session 2 and session 3), and the Center for Epidemiologic Studies Depression Scale (CESD) (administered in session 2 and session 3). There were 227 questions in total (40 of these were repeats from STAI and CESD given again in session 3).

Calculating Factor Scores from the Previous Established Factor Loadings

To calculate factor scores for participants in the current study from the previously estimated bifactor model (**Chapter 2**), we regressed the current participant's responses onto the previous factor loadings. We used the Anderson-Rubin method (Anderson & Rubin, 1956) to preserve orthogonality. The previous bifactor model used 128 questions coming from the STAI, the CESD, the BDI, the MASQ anxious arousal subscale, the MASQ anhedonic depression subscales, the PSWQ, and neuroticism items from the EPQ. The BDI or the EPQ were not administered in the current study, so their loadings could not be used for calculating the factor scores in this dataset. Nonetheless, scores calculated in the previous dataset from the full set of items ($n=128$) and the reduced set of items ($n=95$) were extremely correlated ($r \geq 0.98$ for each of the three factors). Participants in the current study completed the STAI and the CESD in two different sessions, so the responses in each session were concatenated in order to calculate the factor scores (averaging the scores across session yielded extremely similar factor scores; $r > 0.98$ for each factors). Therefore, the score calculation in this dataset relied on 135 total (95 unique) questions.

Administering STAI and CESD in two different sessions also allowed us to calculate test-retest reliability. The test-retest reliability for STAI was very good ($r=0.91$), whereas it was more moderate for CESD ($r=0.53$). The CESD test-retest reliability varied widely across individual items, ranging from $r=0.69$ for "I could not get going" to $r=0.2$ for "my appetite was poor", partially motivating the use of individual items as opposed to subscales in the factor analysis.

Computational Models of Belief Updating

In all of the computational models, we assume that participants distribute their belief over different possible values for the probability $q \in [0,1]$ that they are in the top half of participants. This belief distribution is parameterized as a Beta distribution with two parameters, α_t and β_t . When reporting their beliefs on a trial t , participants are assumed to sample a probability estimate \hat{q}_t from the belief distribution (sampling is denoted by Eqn. 1).

Eqn. 1

$$\hat{q}_t \sim \text{Beta}(\alpha_t, \beta_t)$$

The mean μ_t and precision v_t (i.e. belief certainty) of this belief distribution are related to α_t and β_t through Eqn. 2a-b.

Eqn. 2a-b

$$\mu_t = \frac{\alpha_t}{\alpha_t + \beta_t}$$

$$v_t = \alpha_t + \beta_t$$

The Bayesian models, discussed next, update α_t and β_t in response to feedback, which indirectly change both the mean and the precision. The RW models, discussed after that, update μ_t directly in response to feedback and estimate a separate, constant value for the precision.

Bayesian Updating Model

The two Bayesian models assume that participants update their belief distribution using Bayes rule. After observing either positive or negative feedback X_t on trial t ($X_t = 1$ denotes positive feedback; $X_t = 0$ denotes negative feedback), participants update α_t and β_t according to Eqn. 3a-b. (Note that Bayes rule can be reduced to this specific form, because we are using a Beta-Bernoulli model).

Eqn. 3a-b

$$\alpha_t = \alpha_{t-1} + X_t$$

$$\beta_t = \beta_{t-1} + (1 - X_t)$$

Following each update, participants sample a new probability estimate, again according to Eqn. 1. The only two free parameters in this model, which are used to characterize individual differences, are the belief distribution parameters before the participant has received any feedback (on trial 0): $\alpha_0 \in [0,500]$ and $\beta_0 \in [0,500]$.

Biased Bayesian Updating Model

To assess whether participants show biased updating (i.e. have different rates of updating following positive or negative feedback), we include a bias parameter ω . This is achieved by replacing Eqn. 3a-b with Eqn. 4a-b below ($\omega = 1$ represents unbiased updating, whereas negative and positive bias are given by $\omega < 1$ and $\omega > 1$, respectively). The free parameters for this model are $\alpha_0 \in [0,500]$, $\beta_0 \in [0,500]$, and $\omega \in [0.1,10]$.

Eqn. 4a-b

$$\alpha_t = \alpha_{t-1} + \omega X_t$$

$$\beta_t = \beta_{t-1} + \frac{1}{\omega} (1 - X_t)$$

Rescorla-Wagner (RW) Model

In the Rescorla-Wagner models, participants update the mean of their belief distribution based on a difference between the feedback $X_t \in [0,1]$ and the mean of the belief distribution on the previous trial. This difference, called the prediction error, is scaled by a learning rate parameter η .

Eqn 5.

$$\mu_t = \mu_{t-1} + \eta(X_t - \mu_{t-1})$$

Participants are assumed to use a fixed precision for their belief distribution v , which is a parameter of the model. Then, after each update, participants sample a new belief from the distribution, using Eqn. 1, with $\alpha_t = \mu_t * v$, and $\beta_t = v - \alpha_t$. The free parameters for this model are the precision of the belief distribution $v \in [2,1000]$, the mean of the distribution on the trial before any feedback $\mu_0 \in [0,1]$, and the learning rate $\eta \in [0,1]$.

Biased Rescorla-Wagner (RW) Model

To incorporate a bias into the Rescorla-Wagner model, we scaled the feedback X_t by $r \in [0,5]$ (replacing Eqn. 5 with Eqn. 6); a value of $r < 1$ will bias beliefs in the negative direction and $r > 1$ will bias beliefs in the positive direction.

Eqn 6.

$$\mu_t = \mu_{t-1} + \eta(rX_t - \mu_{t-1})$$

The free parameters for this model are the precision of the belief distribution $v \in [2,1000]$, the mean of the distribution on the trial before any feedback $\mu_0 \in [0,1]$, the learning rate $\eta \in [0,1]$, and the feedback bias $r \in [0,5]$.

Supplemental Results

Supplemental Table 4.1: Beliefs about the Other Participant

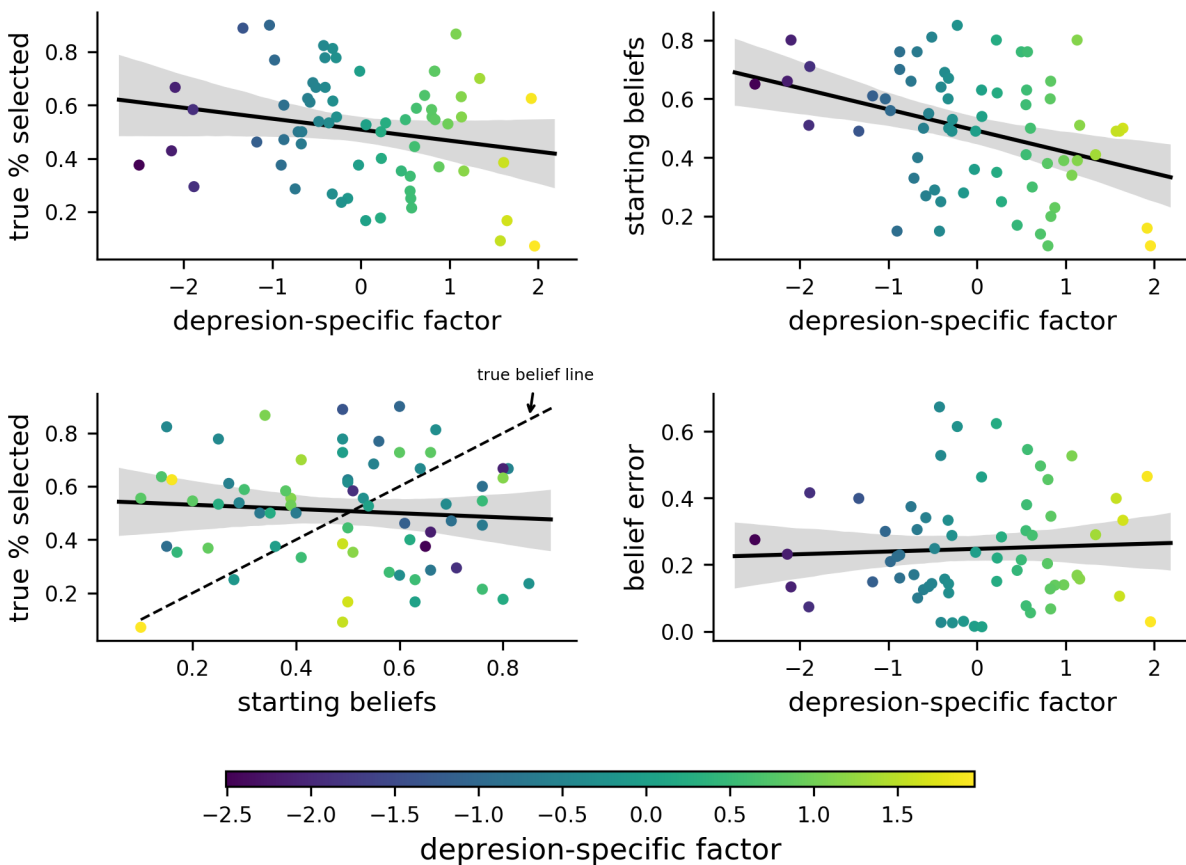
	General Factor	Depression-Specific	Anxiety-Specific
Starting Belief	$r=-0.08, p=0.531$	$r=-0.14, p=0.25$	$r=0.06, p=0.619$
Ending – Starting Belief	$r=-0.11, p=0.378$	$r=0.06, p=0.648$	$r=-0.09, p=0.486$
Ending Belief	$r=-0.15, p=0.226$	$r=-0.06, p=0.605$	$r=-0.02, p=0.857$

general factor	0.36 (0.0)	0.71 (0.0)	0.75 (0.0)	0.35 (0.0)	0.4 (0.0)	0.53 (0.0)	0.59 (0.0)	0.27 (0.03)	0.5 (0.0)	0.64 (0.0)	0.77 (0.0)	0.55 (0.0)	0.31 (0.01)	0.63 (0.0)	0.66 (0.0)	0.5 (0.0)
depression-specific factor	0.3 (0.01)	0.26 (0.04)	0.13 (0.29)	0.54 (0.0)	0.56 (0.0)	0.12 (0.34)	0.05 (0.69)	0.66 (0.0)	0.49 (0.0)	0.24 (0.05)	0.13 (0.32)	0.17 (0.17)	0.74 (0.0)	0.14 (0.27)	0.28 (0.03)	0.2 (0.11)
anxiety-specific factor	0.09 (0.47)	0.43 (0.0)	0.34 (0.01)	0.43 (0.0)	0.37 (0.0)	-0.17 (0.18)	-0.06 (0.6)	0.1 (0.4)	0.03 (0.8)	0.16 (0.19)	0.11 (0.36)	0.16 (0.21)	-0.01 (0.94)	0.25 (0.05)	-0.1 (0.41)	0.78 (0.0)
	STAI state (sess3)	STAI anx (sess2)	STAI anx (sess3)	STAI dep (sess2)	STAI dep (sess3)	CESD dep (sess2)	CESD dep (sess3)	CESD anh (sess2)	CESD anh (sess3)	CESD som (sess2)	CESD som (sess3)	MASQ anxars (sess2)	MASQ anh (sess2)	MASQ anxgen (sess2)	MASQ depden (sess2)	PSWQ (sess1)

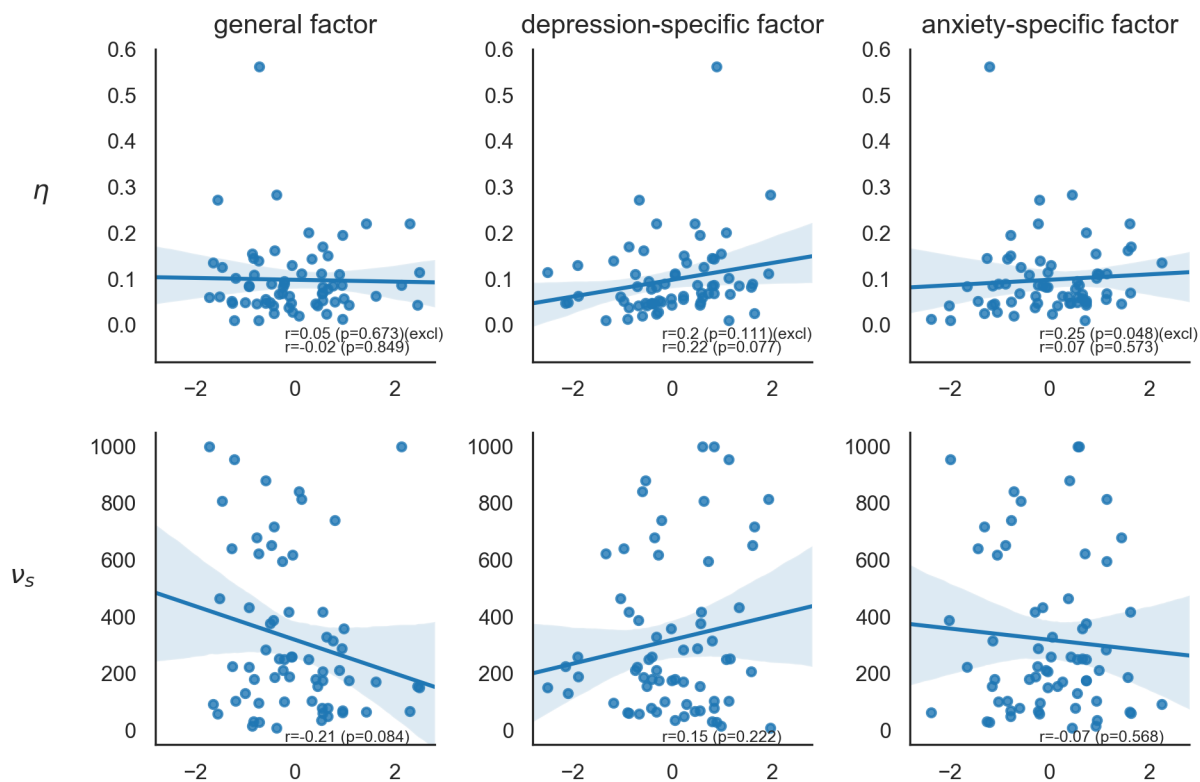
Supplemental Figure 4.1: Correlation of questionnaire subscales with factor scores calculated using loadings from a concurrent study (**Chapter 2**). The general factor scores correlate with most of the subscales. The depression-specific factor scores correlate most strongly with the anhedonia-related subscales (e.g. MASQ anhedonia). The anxiety-specific factor scores correlate most strongly with a worry-related questionnaire (i.e. the PSWQ).

starting belief	0.09 (0.48)	-0.36 (0.0)	-0.05 (0.68)	0.02 (0.87)	-0.23 (0.06)	-0.04 (0.76)	-0.28 (0.02)	-0.27 (0.03)	-0.07 (0.6)	-0.11 (0.37)	-0.4 (0.0)	-0.21 (0.09)	0.01 (0.94)	-0.03 (0.83)	0.11 (0.4)	-0.3 (0.01)	0.06 (0.63)	-0.1 (0.41)	-0.07 (0.59)
end-start	0.04 (0.77)	0.13 (0.28)	-0.23 (0.07)	-0.03 (0.8)	0.05 (0.7)	-0.04 (0.73)	-0.04 (0.77)	-0.04 (0.74)	0.18 (0.14)	0.14 (0.28)	0.27 (0.03)	0.08 (0.54)	0.08 (0.53)	0.01 (0.93)	-0.04 (0.74)	0.17 (0.18)	-0.02 (0.86)	0.11 (0.36)	-0.14 (0.26)
ending belief	0.12 (0.33)	-0.26 (0.03)	-0.24 (0.05)	-0.01 (0.97)	-0.2 (0.11)	-0.08 (0.54)	-0.33 (0.01)	-0.32 (0.01)	0.08 (0.5)	-0.01 (0.96)	-0.19 (0.12)	-0.16 (0.21)	0.07 (0.55)	-0.02 (0.87)	0.08 (0.54)	-0.18 (0.16)	0.05 (0.72)	-0.01 (0.91)	-0.19 (0.13)
	general factor	depression-specific factor	anxiety-specific factor	STAI state (sess3)	STAI anx (sess2)	STAI anx (sess3)	STAI dep (sess2)	STAI dep (sess3)	CESD dep (sess2)	CESD dep (sess3)	CESD anh (sess2)	CESD anh (sess3)	CESD som (sess2)	CESD som (sess3)	MASQ anxars (sess2)	MASQ anh (sess2)	MASQ anxgen (sess2)	MASQ depden (sess2)	PSWQ (sess1)

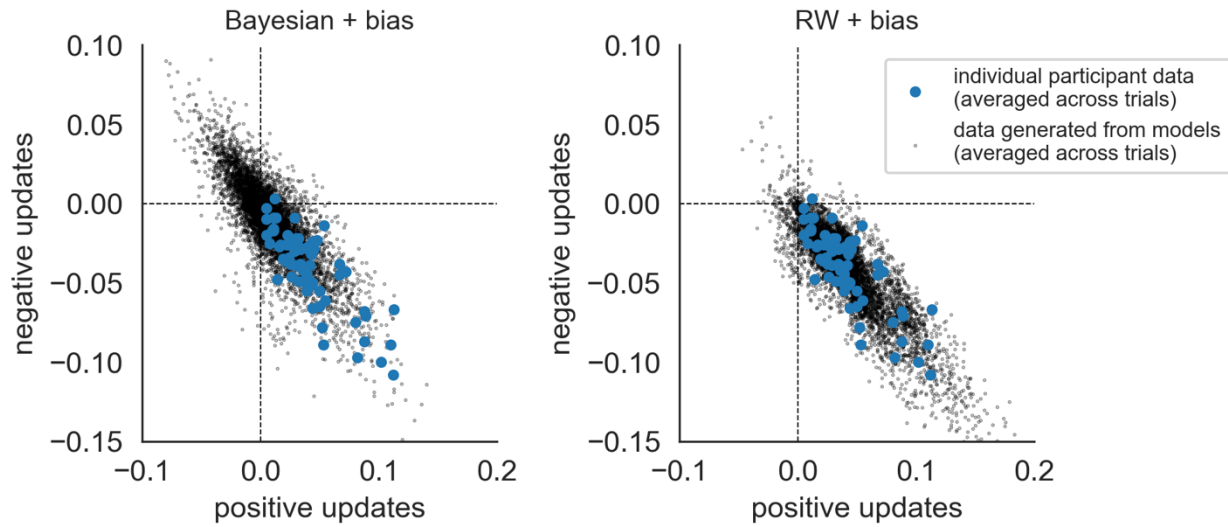
Supplemental Figure 4.2: Correlations of starting belief, ending belief minus starting belief, and ending belief with individual subscales and factor scores. P-values are shown in parentheses underneath correlation values. In line with the factor score results, all of the anhedonia related measures were negatively correlated with starting belief (CESD anhedonia, MASQ anhedonia, STAI depression). The subscales that measure more cognitive or somatic symptoms of depression (CESD depression, CESD somatic and MASQ-DS) were only weakly related to starting belief. The STAI state anxiety measure from session 3 was uncorrelated with starting belief, suggesting that negative starting beliefs are associated with traits rather than temporally depressed or anxious mood.



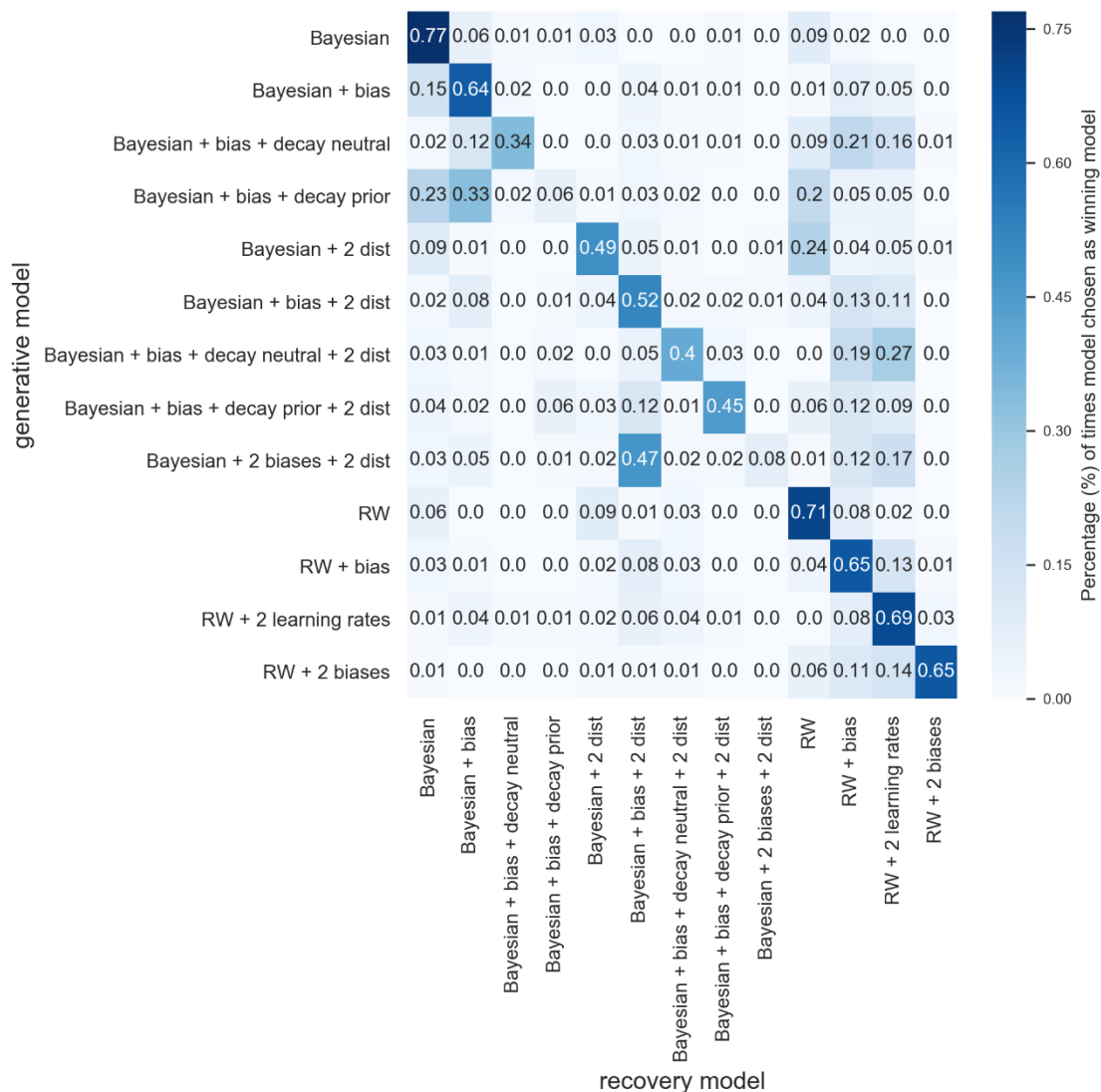
Supplemental Figure 4.3: Testing for Depressive Realism. (a) Scores on the depression-specific factor were weakly correlated with 'profile competitiveness' ($r=-0.21$, $p=0.1$), which is defined as the percentage of times that a participant was actually chosen as a potential partner in session 2 (denoted as 'true % selected' y-axis). (b) The relationship between scores on the depression-specific factor and the negative starting beliefs was not, however, mediated by profile competitiveness ($p=0.4$). (c) The lack of correlation between profile competitiveness and starting belief across participants ($r=-0.07$, $p=0.52$) is responsible for the lack of mediation. (d) The depression-specific factor is also not related to belief error, calculated as the root squared error between profile competitiveness and the starting beliefs.



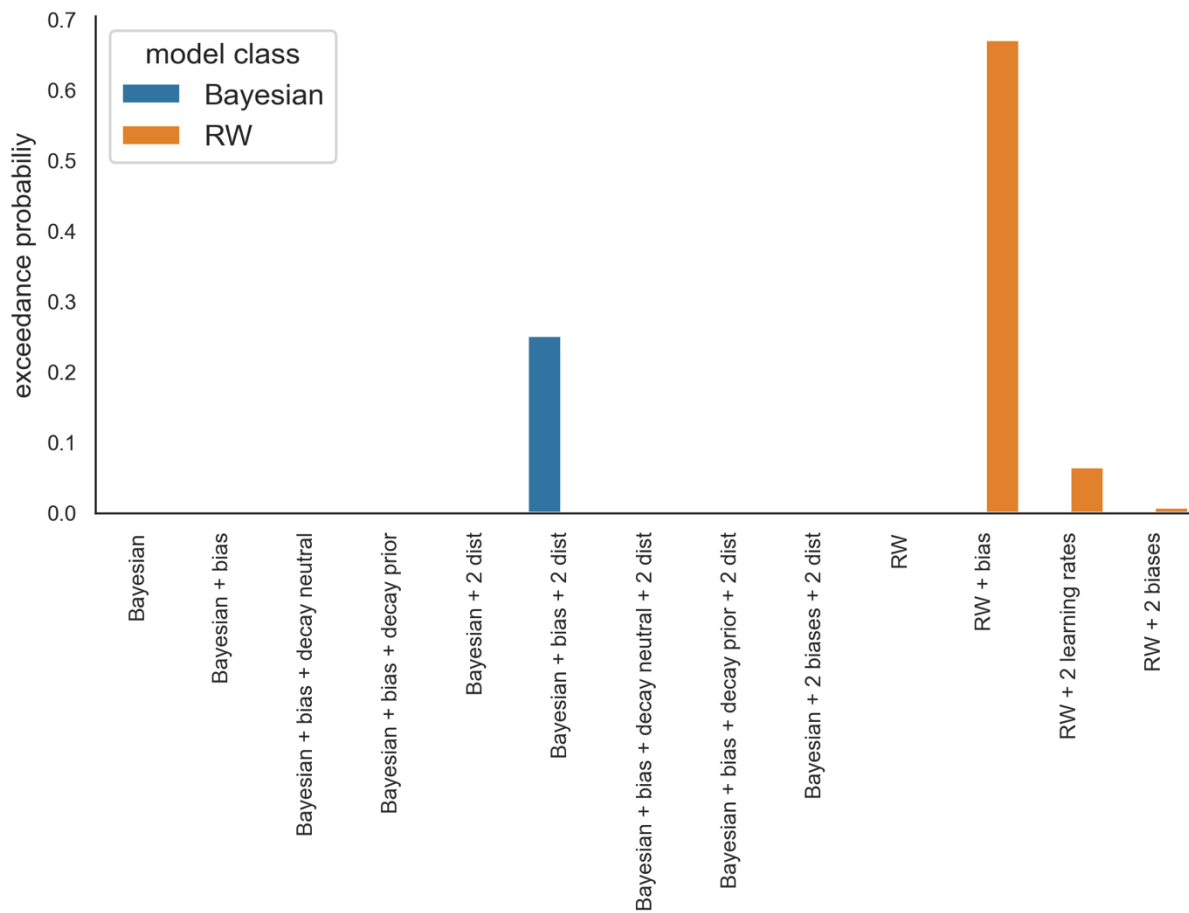
Supplemental Figure 4.4: Participants' standardized scores for the general factor, anxiety-specific factor and depression-specific factor (columns and x-axes) are plotted against the two other parameters in biased RW model. η is the estimated learning rate, which measures the absolute amount that beliefs are updated in response to feedback, and ν_s is the estimated precision of the reporting belief distribution, which measures how closely participants report beliefs to the mean of their belief distribution.



Supplemental Figure 4.5: Qualitative model comparison. (Left subplot) Simulated data from the biased Bayesian model (black points) fails to match the distribution of positive and negative updates averaged across trials per participant (blue points). (Right subplot) Simulated data from the biased RW model provides a better match to the distribution. In each plot, average updates in the all but the bottom right quadrant are in the wrong direction (e.g. updating negatively to positive feedback). For each participant, 100 new beliefs were simulated for each trial from the participant's estimated parameters (hence the larger number of black data points than blue datapoints).



Supplemental Figure 4.6: Confusion matrix for model recovery simulations. For each model on the y-axis, data from 100 new participants was simulated; (new parameter values were chosen for each simulated participant by sampling from distributions informed by the distribution of actual parameter estimates). These models, listed again on the x-axis, were each fit to the data of each simulated participant, and the percentage of participants for which the model had the best penalized fit (BIC) compared to all other models is shown in the cells of the matrix. The percentages add up to 100% for each simulated dataset (i.e. across the rows in the matrix). A high percentage along the diagonal indicates high model recoverability—i.e., that the model that generated the data was chosen as the best fitting model against all others. The first two Bayesian models and the first two RW models were used in the main text.



Supplemental Figure 4.7: Quantitative model comparison. The y-axis is the exceedance probability that a particular model is more prevalent in the population than the others (see Stephan et al. 2009 for more details).

Supplemental Materials

Alternative Bayesian Models

Bayesian + Bias + Decay Neutral

In this model, the belief distribution decays (using Eqn. 7) towards a neutral distribution $Beta(1,1)$ after updating each feedback (using Eqn. 4). The parameter $\gamma \in [0,1]$ determines the degree to which α_t and β_t persist versus decay from one trial to the next.

Eqn. 7

$$\begin{aligned}\alpha_t &= \gamma\alpha_t + (1 - \gamma) * 1 \\ \beta_t &= \gamma\beta_t + (1 - \gamma) * 1\end{aligned}$$

Bayesian + Bias + Decay Prior

In this model, the belief distribution decays (using Eqn. 8) towards the prior distribution estimated for each participant $Beta(\alpha_0, \beta_0)$ after updating each feedback (using Eqn. 4). The parameter $\gamma \in [0,1]$ determines the size of the decay.

Eqn. 8

$$\begin{aligned}\alpha_t &= \gamma\alpha_t + (1 - \gamma)\alpha_0 \\ \beta_t &= \gamma\beta_t + (1 - \gamma)\beta_0\end{aligned}$$

Bayesian + 2 Distributions (separate reporting and updating distributions)

This model divides the belief distribution that participants used for both reporting beliefs (Eqn. 1) and updating beliefs (Eqn. 3a-b) into two separate distributions, one for reporting and one for updating. The rationale for this separation is to capture behavior like the following: a participant may update his beliefs substantially to new information (i.e. have a wide updating distribution), but he may report beliefs close to the mean of the distribution (i.e. have a narrow reporting distribution).

The reporting distribution is still given by Eqn. 1. The new distribution for updating is denoted by $Beta(\alpha_t^u, \beta_t^u)$, where the superscript u denotes ‘updating’ and is used to differentiate these parameters from the α_t and β_t of the reporting distribution. The parameters of the updating distribution are updated according to Eqn. 9a-b, which replaces Eqn. 3a-b.

Eqn. 9a-b

$$\alpha_t^u = \alpha_{t-1}^u + X_t$$

$$\beta_t^u = \beta_{t-1}^u + (1 - X_t)$$

The α_t and β_t from the reporting distribution are no longer updated. Instead, they are calculated using the updating distribution parameters (using Eqn 10 and 11a-b), where v controls the width of the reporting distribution.

Eqn. 10

$$\mu_t = \frac{\alpha_t^u}{\alpha_t^u + \beta_t^u}$$

Eqn. 11a-b

$$\alpha_t = v * \mu_t$$

$$\beta_t = v - \alpha$$

The free parameters for this model are $\alpha_0^u \in [0,100]$, $\beta_0^u \in [0,100]$, and $v \in [1,1000]$.

Bayesian + Bias + 2 Distributions (separate reporting and updating distributions)

This model is the same as the previous model except that it uses a bias ω in Eqn 9a-b. The free parameters for this model are $\alpha_0^u \in [0,100]$, $\beta_0^u \in [0,100]$, $\omega \in [0.1,10]$, $v \in [1,1000]$.

Bayesian + Bias + Decay Neutral + 2 Distributions (separate reporting and updating distributions)

This model decays the updating distribution back to a neutral distribution centered around 50% using Eqn 7. The free parameters for this model are $\alpha_0^u \in [0,100]$, $\beta_0^u \in [0,100]$, $\omega \in [0.1,10]$, $\gamma \in [0.2,1]$, $v \in [1,1000]$.

Bayesian + Bias + Decay Prior + 2 Distributions (separate reporting and updating distributions)

This model decays the updating distribution parameters α^u, β^u back to their starting values using Eqn 8. The free parameters for this model are $\alpha_0^u \in [0,100]$, $\beta_0^u \in [0,100]$, $\omega \in [0.1,10]$, $\gamma \in [0.2,1]$, $v \in [1,1000]$.

Bayesian + 2 Biases + Decay Neutral + 2 Distributions (separate reporting and updating distributions)

This model has two bias parameters, one for positive and one for negative feedback. The free parameters for this model are $\alpha_0^u \in [0,100]$, $\beta_0^u \in [0,100]$, $\omega_{pos} \in [0.1,10]$, $\omega_{neg} \in [0.1,10]$, $v \in [1,1000]$.

Alternative Rescorla-Wagner Models

Rescorla-Wagner (RW) Model + 2 Learning Rates

This model is similar to the Rescorla-Wagner model, but has two separate learning rate parameters η_{pos} and η_{neg} , for updates following positive and negative feedback respectively.

Rescorla-Wagner (RW) Model + 2 Biases

This model is similar to the biased Rescorla-Wagner model, however the values for positive and negative feedback have been changed from $\{0, r\}$ to $\{r_{pos}, r_{neg}\}$.

Eqn. 12a-b

$$\begin{aligned}\mu_t &= \mu_{t-1} + \eta(r_{pos}X_t - \mu_{t-1}) \\ \mu_t &= \mu_{t-1} + \eta(r_{neg}(1 - X_t) - \mu_{t-1})\end{aligned}$$

Chapter 5: Brief General Discussion

The broad aim of the empirical Chapters (**2**, **3**, and **4**) was to examine dysfunctional behavior associated with anxiety and depression using computational frameworks of decision making. All three chapters investigated behavior (or beliefs) under second order uncertainty (i.e. uncertainty regarding probabilities). **Chapter 2** investigated volatility, which arises when probabilities change over time. **Chapter 3** investigated ambiguity, which arises when information is missing about those probabilities. **Chapter 4** also investigated ambiguous probabilities but focused on how individuals may bring different prior beliefs to new situations and update their beliefs differently as information is incrementally provided. Although the results from all three studies (especially in **Chapter 3**) are preliminary, a few words can be said in summary about how dysfunctional behavior and beliefs associated with anxiety and depression seem to be impacted by second order uncertainty.

Before that, it is worth pointing out that another core aim of the empirical studies was to disentangle anxiety and depression. To that aim, each of the three studies used bifactor factor analysis to calculate scores for participants on a general factor and two or more specific factors. Scores on the general factors represent the broad elevation of both mood and anxiety symptoms, and scores on the specific factors represent the elevation of a smaller set of more closely related symptoms (e.g. those related to anhedonia, worry, or physiological anxiety, etc.). By using general and specific factor scores, rather than the pre-existing anxiety and depression measures, we hoped that we might identify processes that may differentially confer vulnerability to common or unique aspects of symptomology. Indeed, we implicated a general factor in processing volatility in **Chapter 2**, an anxiety-specific factor in processing ambiguity in **Chapter 3**, another anxiety-specific factor in belief updating in **Chapter 4**, and a depression-specific factor in holding more negative prior beliefs also in **Chapter 4**.

In **Chapter 2**, higher scores on the general factor, but not the anxiety-specific factor nor depression-specific factor, were associated with a lack of learning rate adjustment to volatility (i.e. not having higher learning rates in the volatile versus stable block), during the course of learning about action-outcome probabilities. This means that individuals with high levels of anxious symptoms, or depressive symptoms, or both, might have difficulties making optimal (or even good) decisions in real-world situations containing volatility (e.g. personal relationships, jobs, etc.). Worse life experiences as a result of the poorer choices made could, in turn, lead to the development of (or the exacerbation of pre-existing) mood and anxiety disorders.

In **Chapter 3**, higher scores on an anxiety-specific factor (i.e. physiological anxiety), but not the general factor, were linked to ambiguity aversion. Both the ambiguity task and the volatility task involve processing second order uncertainty, in some sense. This raises the question of why did we not see any general factor effects in the ambiguity task, or see any anxiety-specific effects in the volatility task?

First, in addressing this question, it is important to clarify what processing second order uncertainty (SOU) might mean in both tasks, and how that might be reflected in our primary behavioral measurements (i.e. learning rate adjustment and ambiguity aversion).

One aspect of processing SOU could be the creation a distribution around one's (first order) probability point estimates. Difficulties here might come in the form of having too much or too

little uncertainty in this second-order distribution (i.e. its width), which would be reflected in the two tasks in different ways.

In the volatility task (**Chapter 2**), differences in the amount of uncertainty in a second-order distribution should translate into differences in learning rates, with more uncertainty in one's existing estimates leading to higher learning rates, because new observations are given more relative weight (Kalman 1960). A lack of learning rate adjustment would then be observed if someone has consistently too much second-order uncertainty, consistently too little, or an amount that fluctuates but does not track the actual level of volatility, changing throughout in the experiment. Interestingly, in our in-lab experiment, higher scores on the general factor were associated with lower average learning rates, but in the online experiment, no differences in average learning rates were associated with the general factor. This suggests that some individuals, who had high general factor scores, had too much SOU and others too little. Either way, their miscalibration manifested as a lack of an appropriate adjustment in learning rate to volatility level.

In contrast, in the ambiguity task (**Chapter 3**), neither having too much nor too little SOU would necessarily lead to ambiguity aversion—it would depend on the probabilities in the decision. For example, if the ambiguous probability (e.g. $p=0.9$) is slightly higher than the unambiguous probability (e.g. $p=0.8$) for a rewarding outcome, high amounts of SOU can lead to ambiguity aversion if the participant adjusts the probability estimate towards 50% on the basis of that uncertainty (e.g. from $p=0.9$ to $p=0.7$); on the other hand, if the ambiguous probability (e.g. $p=0.2$) is slightly lower than the unambiguous probability (e.g. $p=0.25$), high amounts of SOU would lead to ambiguity seeking as the ambiguous probability is adjusted towards 50% (e.g. from $p=0.2$ to $p=0.3$). The reverse argument can be made for low amounts of SOU. Therefore, given the indirect relationship between the amount of SOU and ambiguity aversion, it is perhaps not surprising after all that there was no relationship between general factor scores and ambiguity aversion.

In response to the same question, it is also important to note that second order uncertainty entered into the two experimental situations very differently. For the ambiguity task, the SOU (i.e. the amount of missing information) is presented clearly, and all at once, to the participant. For the volatility task, the SOU (i.e. coming from potential change in action-outcome contingencies) must be estimated from incremental experience with the environment. For first order probabilities, different modes of delivery lead people to treat probabilities very differently; this is often referred to as the 'description-experience' gap (reviewed in Fox et al., 2015). Thus, another reason for the apparent dissociation between the physiological anxiety factor and the general factor may be that the general factor is more closely related to difficulties inferring and using SOU from incremental experience, and physiological anxiety is more related to an attitude towards SOU when it is presented all up front.

Finally, it is also worth pointing out that some people may choose to avoid options that involve second order uncertainty when given the chance regardless of their ability to estimate or use it correctly, because it is computationally intensive to process. The ambiguity task involves a choice between a risky and an ambiguous urn, so some participants could occasionally have chosen the unambiguous option to avoid spending the effort deliberating between the two urns; this would be picked up as ambiguity aversion, on average, because the two urns were balanced for expected value across trials. In contrast, in the volatility task, both

options were subject to the same volatility since their probabilities were p and $1-p$, so participants could not avoid processing SOU. If individual differences in a willingness to expend effort doing computations relates to the general factor or to the physiological factor, this could additionally contribute a difference in effects between the two tasks.

In **Chapter 4**, higher scores on the depression-specific factor (i.e. anhedonia symptoms), but not scores on the anxiety factor nor on the general factor, were associated with negatively biased prior beliefs relative to other participants. That the depression factor was solely implicated may be due to the fact that this experiment, unlike the others, involved beliefs about self-worth or competence relative to others. A key feature that has been proposed to distinguish depression from anxiety is the content of the dysfunctional beliefs (Beck 1979), which are thought to revolve around hopelessness and worthlessness for depression and threat and vulnerability for anxiety (Clark & Beck, 2010).

In **Chapter 4**, we also observed that the scores on the anxiety-specific factor were associated with negatively biased updating in response to objective feedback. One might expect that the anxiety-specific factor would have therefore been related to increased learning rates following bad outcomes as opposed to good outcomes in the volatility task. However, the presence of volatility and the interaction of that with asymmetric learning rates for good versus bad outcomes, may have obscured any potential relationship between anxiety-specific variance and asymmetric learning in the volatility task.

In closing, the work in this dissertation hopefully makes a small contribution to the understanding of why individuals who suffer from anxiety or depression make decisions, hold beliefs, and behave in dysfunctional ways. It will also hopefully inspire future work in examining dysfunction related to second order uncertainty or using similar computational methods.

References

- Abramson, L. Y., Metalsky, G. I., & Alloy, L. B. (1989). Hopelessness depression: A theory-based subtype of depression. *Psychological review*, 96(2), 358.
- Adida, M., Clark, L., Pomietto, P., Kaladjian, A., Besnier, N., Azorin, J. M., ... & Goodwin, G. M. (2008). Lack of insight may predict impaired decision making in manic patients. *Bipolar disorders*, 10(7), 829-837.
- Adida, M., Jollant, F., Clark, L., Besnier, N., Guillaume, S., Kaladjian, A., ... & Courtet, P. (2011). Trait-related decision-making impairment in the three phases of bipolar disorder. *Biological psychiatry*, 70(4), 357-365.
- Akaishi, R., Umeda, K., Nagase, A., & Sakai, K. (2014). Autonomous mechanism of internal choice estimate underlies decision inertia. *Neuron*, 81(1), 195-206.
- Alloy, L. B., & Abramson, L. Y. (1979). Judgment of contingency in depressed and nondepressed students: Sadder but wiser?. *Journal of experimental psychology: General*, 108(4), 441.
- Am. Psychiatr. Assoc. (2013) *Diagnostic and Statistical Manual of Mental Disorders*. Arlington, VA: American Psychiatric Publishing. 5th ed.
- Ambrose, R. E., Pfeiffer, B. E., & Foster, D. J. (2016). Reverse replay of hippocampal place cells is uniquely modulated by changing reward. *Neuron*, 91(5), 1124-1136.
- Amsterdam, J. D., Settle, R. G., Doty, R. L., Abelman, E., & Winokur, A. (1987). Taste and smell perception in depression. *Biological psychiatry*, 22(12), 1481-1485.
- Anderson, T. W., & Rubin, H. (1956). Statistical inference in factor analysis. In *Proceedings of the third Berkeley symposium on mathematical statistics and probability* (Vol. 5, pp. 111-150).
- Aylward, J., Valton, V., Ahn, W. Y., Bond, R. L., Dayan, P., Roiser, J. P., & Robinson, O. J. (2019). Altered learning under uncertainty in unmedicated mood and anxiety disorders. *Nature human behaviour*, 1.
- Bach, D. R., Hulme, O., Penny, W. D., & Dolan, R. J. (2011). The known unknowns: neural representation of second-order uncertainty, and ambiguity. *Journal of Neuroscience*, 31(13), 4811-4820.
- Bach, D. R., Guitart-Masip, M., Packard, P. A., Miró, J., Falip, M., Fuentemilla, L., & Dolan, R. J. (2014). Human hippocampus arbitrates approach-avoidance conflict. *Current Biology*, 24(5), 541-547.

- Ballenger, J. C. (1999). Current treatments of the anxiety disorders in adults. *Biological psychiatry*, 46(11), 1579-1594.
- Barlow, D. H. (2000). Unraveling the mysteries of anxiety and its disorders from the perspective of emotion theory. *American psychologist*, 55(11), 1247.
- Battaglia, P. W., Hamrick, J. B., & Tenenbaum, J. B. (2013). Simulation as an engine of physical scene understanding. *Proceedings of the National Academy of Sciences*, 110(45), 18327-18332.
- Bechara, A., & Damasio, A. R. (2005). The somatic marker hypothesis: A neural theory of economic decision. *Games and economic behavior*, 52(2), 336-372.
- Beck, A. T. (1976). *Cognitive therapy and the emotional disorders*. New York: International University Press.
- Beck, A. T. (2005). The current state of cognitive therapy: a 40-year retrospective. *Archives of general psychiatry*, 62(9), 953-959.
- Beck, A.T., & Emery, G. (with Greenberg, R. L.) (1985) *Anxiety disorders and phobias: A cognitive perspective*, Basic Books
- Beck, A. T., Ward, C. H., Mendelson, M., Mock, J., & Erbaugh, J. (1961). An inventory for measuring depression. *Archives of general psychiatry*, 4(6), 561-571.
- Behrens, T. E., Woolrich, M. W., Walton, M. E., & Rushworth, M. F. (2007). Learning the value of information in an uncertain world. *Nature neuroscience*, 10(9), 1214.
- Berger, J. O. (2013). *Statistical decision theory and Bayesian analysis*. Springer Science & Business Media.
- Bijsterbosch, J., Smith, S., Forster, S., John, O. P., & Bishop, S. J. (2014). Resting state correlates of subdimensions of anxious affect. *Journal of Cognitive Neuroscience*, 26(4), 914-926.
- Birrell, J., Meares, K., Wilkinson, A., & Freeston, M. (2011). Toward a definition of intolerance of uncertainty: A review of factor analytical studies of the Intolerance of Uncertainty Scale. *Clinical psychology review*, 31(7), 1198-1208.
- Bishop, S., Dalgleish, T., & Yule, W. (2004). Memory for emotional stories in high and low depressed children. *Memory*, 12(2), 214-230.
- Bjärehed, J., Sarkohi, A., & Andersson, G. (2010). Less positive or more negative? Future-directed thinking in mild to moderate depression. *Cognitive behaviour therapy*, 39(1), 37-45.

- Blair KS, Blair RJR. (2012). A cognitive neuroscience approach to generalized anxiety disorder and social phobia. *Emot. Rev.* 4:133–38
- Bonnelle, V., Manohar, S., Behrens, T., & Husain, M. (2015). Individual differences in premotor brain systems underlie behavioral apathy. *Cerebral cortex*, 26(2), 807-819.
- Bonnelle, V., Veromann, K. R., Heyes, S. B., Sterzo, E. L., Manohar, S., & Husain, M. (2015). Characterization of reward and effort mechanisms in apathy. *Journal of Physiology-Paris*, 109(1-3), 16-26.
- Boswell, J. F., Thompson-Hollands, J., Farchione, T. J., & Barlow, D. H. (2013). Intolerance of uncertainty: A common factor in the treatment of emotional disorders. *Journal of clinical psychology*, 69(6), 630-645.
- Bradley, B., & Mathews, A. (1983). Negative self-schemata in clinical depression. *British Journal of Clinical Psychology*, 22(3), 173-181.
- Branco, L. D., Cotrena, C., Ponsoni, A., Salvador-Silva, R., Vasconcellos, S. J. L., & Fonseca, R. P. (2017). Identification and perceived intensity of facial expressions of emotion in bipolar disorder and major depression. *Archives of Clinical Neuropsychology*, 33(4), 491-501.
- Brodbeck, J., Abbott, R. A., Goodyer, I. M., & Croudace, T. J. (2011). General and specific components of depression and anxiety in an adolescent population. *BMC psychiatry*, 11(1), 191.
- Browning, M., Behrens, T. E., Jocham, G., O'reilly, J. X., & Bishop, S. J. (2015). Anxious individuals have difficulty learning the causal statistics of aversive environments. *Nature neuroscience*, 18(4), 590.
- Bryant R. (2000) Acute stress disorder. *PTSD Res Q* 11.
- Butler, G., & Mathews, A. (1983). Cognitive processes in anxiety. *Advances in behaviour research and therapy*, 5(1), 51-62.
- Cahill, L., Haier, R. J., Fallon, J., Alkire, M. T., Tang, C., Keator, D., ... & Mcgaugh, J. L. (1996). Amygdala activity at encoding correlated with long-term, free recall of emotional information. *Proceedings of the National Academy of Sciences*, 93(15), 8016-8021.
- Camerer, C., & Weber, M. (1992). Recent developments in modeling preferences: Uncertainty and ambiguity. *Journal of risk and uncertainty*, 5(4), 325-370.
- Camerer, C. (1995). Individual Decision Making. *Handbook of Experimental Economics*. J. Kagel and AE Roth.

- Carleton, R. N., Mulvogue, M. K., Thibodeau, M. A., McCabe, R. E., Antony, M. M., & Asmundson, G. J. (2012). Increasingly certain about uncertainty: Intolerance of uncertainty across anxiety and depression. *Journal of Anxiety Disorders*, 26(3), 468-479.
- Charpentier, C. J., Aylward, J., Roiser, J. P., & Robinson, O. J. (2017). Enhanced risk aversion, but not loss aversion, in unmedicated pathological anxiety. *Biological psychiatry*, 81(12), 1014-1022.
- Chong, T. J., Bonnelle, V., & Husain, M. (2016). Quantifying motivation with effort-based decision-making paradigms in health and disease. In *Progress in brain research* (Vol. 229, pp. 71-100). Elsevier.
- Clark, D. A., & Beck, A. T. (2010). Cognitive theory and therapy of anxiety and depression: Convergence with neurobiological findings. *Trends in cognitive sciences*, 14(9), 418-424.
- Clark, D. A., Steer, R. A., & Beck, A. T. (1994). Common and specific dimensions of self-reported anxiety and depression: implications for the cognitive and tripartite models. *Journal of abnormal psychology*, 103(4), 645.
- Clark, D. A., & Teasdale, J. D. (1982). Diurnal variation in clinical depression and accessibility of memories of positive and negative experiences. *Journal of abnormal psychology*, 91(2), 87.
- Clark, L. A., & Watson, D. (1991). Tripartite model of anxiety and depression: psychometric evidence and taxonomic implications. *Journal of abnormal psychology*, 100(3), 316.
- Clark, L. A., & Watson, D. (1995). The mini mood and anxiety symptom questionnaire (Mini-MASQ). *Unpublished manuscript, University of Iowa*.
- Cohen, M., Jaffray, J. Y., & Said, T. (1985). Individual behavior under risk and under uncertainty: An experimental study. *Theory and Decision*, 18(2), 203-228.
- Cryan, J. F., Valentino, R. J., & Lucki, I. (2005). Assessing substrates underlying the behavioral effects of antidepressants using the modified rat forced swimming test. *Neuroscience & Biobehavioral Reviews*, 29(4-5), 547-569.
- Curley, S. P., Eraker, S. A., & Yates, J. F. (1984). An investigation of patient's reactions to therapeutic uncertainty. *Medical Decision Making*, 4(4), 501-511.
- Davis, M., Walker, D. L., Miles, L., & Grillon, C. (2010). Phasic vs sustained fear in rats and humans: role of the extended amygdala in fear vs anxiety. *Neuropsychopharmacology*, 35(1), 105.
- Daw, N. D., & Dayan, P. (2014). The algorithmic anatomy of model-based evaluation. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 369(1655), 20130478.

- Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nature neuroscience*, 8(12), 1704.
- Dayan, P., & Huys, Q. J. (2009). Serotonin in affective control. *Annual review of neuroscience*, 32, 95-126.
- Dayan, P., Roiser, J. P., & Viding, E. (2018). The first steps on long marches: the costs of active observation. In *Rethinking Biopsychosocial Psychiatry*. Oxford University Press.
- Dekel, S., Mandl, C., & Solomon, Z. (2011). Shared and unique predictors of post-traumatic growth and distress. *Journal of clinical psychology*, 67(3), 241-252.
- Depue, R. A., Slater, J. F., Wolfstetter-Kausch, H., Klein, D., Goplerud, E., & Farr, D. (1981). A behavioral paradigm for identifying persons at risk for bipolar depressive disorder: a conceptual framework and five validation studies. *Journal of abnormal psychology*, 90(5), 381.
- DeRubeis, R. J., Hollon, S. D., Amsterdam, J. D., Shelton, R. C., Young, P. R., Salomon, R. M., ... & Gallop, R. (2005). Cognitive therapy vs medications in the treatment of moderate to severe depression. *Archives of general psychiatry*, 62(4), 409-416.
- Devlin, H. C., Johnson, S. L., & Gruber, J. (2015). Feeling good and taking a chance? Associations of hypomania risk with cognitive and behavioral risk taking. *Cognitive therapy and research*, 39(4), 473-479.
- Dobson, K., & Franche, R. L. (1989). A conceptual and empirical review of the depressive realism hypothesis. *Canadian Journal of Behavioural Science/Revue canadienne des sciences du comportement*, 21(4), 419.
- Dolcos, F., Katsumi, Y., Weymar, M., Moore, M., Tsukiura, T., & Dolcos, S. (2017). Emerging directions in emotional episodic memory. *Frontiers in psychology*, 8, 1867.
- Dow, J., & Werlang, S. R. D. C. (1991). Excess volatility of stock prices and Knightian uncertainty.
- Dow, J., & da Costa Werlang, S. R. (1992). Uncertainty aversion, risk aversion, and the optimal choice of portfolio. *Econometrica: Journal of the Econometric Society*, 197-204.
- Doya, K. (1999). What are the computations of the cerebellum, the basal ganglia and the cerebral cortex?. *Neural networks*, 12(7-8), 961-974.
- Doyle, J., Glover, K., Khargonekar, P., & Francis, B. (1988, June). State-space solutions to standard H₂ and H_∞ control problems. In *1988 American Control Conference* (pp. 1691-1696). IEEE.

- Dugas, M. J., Gagnon, F., Ladouceur, R., & Freeston, M. H. (1998). Generalized anxiety disorder: A preliminary test of a conceptual model. *Behaviour research and therapy*, *36*(2), 215-226.
- Dugas, M. J., Gosselin, P., & Ladouceur, R. (2001). Intolerance of uncertainty and worry: Investigating specificity in a nonclinical sample. *Cognitive therapy and Research*, *25*(5), 551-558.
- Eckblad, M., & Chapman, L. J. (1986). Development and validation of a scale for hypomanic personality. *Journal of abnormal psychology*, *95*(3), 214.
- Edge, M. D., Johnson, S. L., Ng, T., & Carver, C. S. (2013). Iowa gambling task performance in euthymic bipolar I disorder: A meta-analysis and empirical study. *Journal of affective disorders*, *150*(1), 115-122.
- Ehlers, A., & Clark, D. M. (2000). A cognitive model of posttraumatic stress disorder. *Behaviour research and therapy*, *38*(4), 319-345.
- Ehlers, A., Mayou, R. A., & Bryant, B. (1998). Psychological predictors of chronic posttraumatic stress disorder after motor vehicle accidents. *Journal of abnormal psychology*, *107*(3), 508.
- Ehring, T., Frank, S., & Ehlers, A. (2008). The role of rumination and reduced concreteness in the maintenance of posttraumatic stress disorder and depression following trauma. *Cognitive therapy and research*, *32*(4), 488-506.
- Ehring, T., Fuchs, N., & Kläsener, I. (2009). The effects of experimentally induced rumination versus distraction on analogue posttraumatic stress symptoms. *Behavior Therapy*, *40*(4), 403-413.
- Eil, D., & Rao, J. M. (2011). The good news-bad news effect: asymmetric processing of objective information about yourself. *American Economic Journal: Microeconomics*, *3*(2), 114-38.
- Einhorn, H. J., & Hogarth, R. M. (1986). Decision making under ambiguity. *Journal of Business*, S225-S250.
- El Leithy, S., Brown, G. P., & Robbins, I. (2006). Counterfactual thinking and posttraumatic stress reactions. *Journal of abnormal psychology*, *115*(3), 629.
- Elliott, R., Sahakian, B. J., Herrod, J. J., Robbins, T. W., & Paykel, E. S. (1997). Abnormal response to negative feedback in unipolar depression: evidence for a diagnosis specific impairment. *Journal of Neurology, Neurosurgery & Psychiatry*, *63*(1), 74-82.
- Ellsberg, D. (1961). Risk, ambiguity, and the Savage axioms. *The quarterly journal of economics*, 643-669.

- Etkin, A., Klemenhagen, K. C., Dudman, J. T., Rogan, M. T., Hen, R., Kandel, E. R., & Hirsch, J. (2004). Individual differences in trait anxiety predict the response of the basolateral amygdala to unconsciously processed fearful faces. *Neuron*, *44*(6), 1043-1055.
- Eysenck, H. J., & Eysenck, S. B. G. (1975). *Manual of the Eysenck Personality Questionnaire (junior and adult)*. Hodder and Stoughton.
- Fishburn, P. C., & Kochenberger, G. A. (1979). Two-piece von Neumann-Morgenstern utility functions. *Decision Sciences*, *10*(4), 503-518.
- Floresco, S. B., & Ghods-Sharifi, S. (2006). Amygdala-prefrontal cortical circuitry regulates effort-based decision making. *Cerebral cortex*, *17*(2), 251-260.
- Floyd, F. J., & Widaman, K. F. (1995). Factor analysis in the development and refinement of clinical assessment instruments. *Psychological assessment*, *7*(3), 286.
- Foster, D. J., & Wilson, M. A. (2006). Reverse replay of behavioural sequences in hippocampal place cells during the awake state. *Nature*, *440*(7084), 680.
- Fox, C. R., Erner, C., & Walters, D. J. (2015). Decision under risk: From the field to the laboratory and back. *The Wiley Blackwell handbook of judgment and decision making*, 43-88.
- Freeston, M. H., Rhéaume, J., Letarte, H., Dugas, M. J., & Ladouceur, R. (1994). Why do people worry?. *Personality and individual differences*, *17*(6), 791-802.
- Gelman, A., & Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical science*, *7*(4), 457-472.
- Gentes, E. L., & Ruscio, A. M. (2011). A meta-analysis of the relation of intolerance of uncertainty to symptoms of generalized anxiety disorder, major depressive disorder, and obsessive-compulsive disorder. *Clinical psychology review*, *31*(6), 923-933.
- Ginzburg, K. (2004). PTSD and world assumptions following myocardial infarction: A longitudinal study. *American Journal of Orthopsychiatry*, *74*(3), 286-292.
- Giorgetta, C., Grecucci, A., Zuanon, S., Perini, L., Balestrieri, M., Bonini, N., ... & Brambilla, P. (2012). Reduced risk-taking behavior as a trait feature of anxiety. *Emotion*, *12*(6), 1373.
- Gray, J. A., & McNaughton, N. (2003). Fundamentals of the septo-hippocampal system. *The Neuropsychology of Anxiety: An Enquiry into the Functions of Septo-hippocampal System*, 2nd ed. Oxford University Press, Oxford, 204-232.

- Grillon, C., Lissek, S., Rabin, S., McDowell, D., Dvir, S., & Pine, D. S. (2008). Increased anxiety during anticipation of unpredictable but not predictable aversive stimuli as a psychophysiological marker of panic disorder. *American Journal of Psychiatry*, *165*(7), 898-904.
- Hamlat, E. J., Connolly, S. L., Hamilton, J. L., Stange, J. P., Abramson, L. Y., & Alloy, L. B. (2015). Rumination and overgeneral autobiographical memory in adolescents: An integration of cognitive vulnerabilities to depression. *Journal of youth and adolescence*, *44*(4), 806-818.
- Hartley, C. A., & Phelps, E. A. (2012). Anxiety and decision-making. *Biological psychiatry*, *72*(2), 113-118.
- Harris, A. J., & Hahn, U. (2011). Unrealistic optimism about future life events: A cautionary note. *Psychological review*, *118*(1), 135.
- Harris, A. J., de Molière, L., Soh, M., & Hahn, U. (2017). Unrealistic comparative optimism: An unsuccessful search for evidence of a genuinely motivational bias. *PloS one*, *12*(3), e0173136.
- Hauber, W., & Sommer, S. (2009). Prefrontostriatal circuitry regulates effort-related decision making. *Cerebral Cortex*, *19*(10), 2240-2247.
- Hollon, S. D., DeRubeis, R. J., Shelton, R. C., Amsterdam, J. D., Salomon, R. M., O'Reardon, J. P., ... & Gallop, R. (2005). Prevention of relapse following cognitive therapy vs medications in moderate to severe depression. *Archives of general psychiatry*, *62*(4), 417-422.
- Hong, R. Y. (2007). Worry and rumination: Differential associations with anxious and depressive symptoms and coping behavior. *Behaviour research and therapy*, *45*(2), 277-290.
- Horn, J. L. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, *30*(2), 179-185.
- Horowitz, M. J. (1985). Disasters and psychological responses to stress. *Psychiatric Annals*, *15*(3), 161-167.
- Huang, H., Thompson, W., & Paulus, M. P. (2017). Computational dysfunctions in anxiety: Failure to differentiate signal from noise. *Biological psychiatry*, *82*(6), 440-446.
- Humphreys, L. G., & Montanelli Jr, R. G. (1975). An investigation of the parallel analysis criterion for determining the number of common factors. *Multivariate Behavioral Research*, *10*(2), 193-205.
- Huys, Q. J., Daw, N. D., & Dayan, P. (2015). Depression: a decision-theoretic analysis. *Annual review of neuroscience*, *38*, 1-23.

Huys, Q. J., Maia, T. V., & Frank, M. J. (2016). Computational psychiatry as a bridge from neuroscience to clinical applications. *Nature neuroscience*, 19(3), 404.

Huys, Q. J., Pizzagalli, D. A., Bogdan, R., & Dayan, P. (2013). Mapping anhedonia onto reinforcement learning: a behavioural meta-analysis. *Biology of mood & anxiety disorders*, 3(1), 12.

Ito, M., & Doya, K. (2009). Validation of decision-making models and analysis of decision variables in the rat basal ganglia. *Journal of Neuroscience*, 29(31), 9861-9874.

Janoff-Bulman R: Shattered Assumptions: Towards a New Psychology of Trauma. Free Press; 1992.

Jennrich, R. I., & Bentler, P. M. (2011). Exploratory bi-factor analysis. *Psychometrika*, 76(4), 537–549. doi:10.1007/bf02294706

Johnson, A., & Redish, A. D. (2005). Hippocampal replay contributes to within session learning in a temporal difference reinforcement learning model. *Neural Networks*, 18(9), 1163-1171.

Johnson, S. L., Edge, M. D., Holmes, M. K., & Carver, C. S. (2012). The behavioral activation system and mania. *Annual review of clinical psychology*, 8, 243-267.

Kaczurkin, A. N., Burton, P. C., Chazin, S. M., Manbeck, A. B., Espensen-Sturges, T., Cooper, S. E., ... & Lissek, S. (2016). Neural substrates of overgeneralized conditioned fear in PTSD. *American Journal of Psychiatry*, 174(2), 125-134.

Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1), 35-45.

Kahneman, D., Slovic, S. P., Slovic, P., & Tversky, A. (Eds.). (1982). *Judgment under uncertainty: Heuristics and biases*. Cambridge university press.

Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2), 363-391.

Keramati, M., Smittenaar, P., Dolan, R. J., & Dayan, P. (2016). Adaptive integration of habits into depth-limited planning defines a habitual-goal-directed spectrum. *Proceedings of the National Academy of Sciences*, 113(45), 12868-12873.

Kessler, R. C., Chiu, W. T., Demler, O., & Walters, E. E. (2005). Prevalence, severity, and comorbidity of 12-month DSM-IV disorders in the National Comorbidity Survey Replication. *Archives of general psychiatry*, 62(6), 617-627.

- Kheirbek, M. A., Klemenhagen, K. C., Sahay, A., & Hen, R. (2012). Neurogenesis and generalization: a new approach to stratify and treat anxiety disorders. *Nature neuroscience*, *15*(12), 1613.
- Khemka, S., Barnes, G., Dolan, R. J., & Bach, D. R. (2017). Dissecting the function of hippocampal oscillations in a human anxiety model. *Journal of Neuroscience*, *37*(29), 6869-6876.
- Kircanski, K., Thompson, R. J., Sorenson, J. E., Sherdell, L., & Gotlib, I. H. (2015). Rumination and worry in daily life: Examining the naturalistic validity of theoretical constructs. *Clinical Psychological Science*, *3*(6), 926-939.
- Kleim, B., & Ehlers, A. (2008). Reduced autobiographical memory specificity predicts depression and posttraumatic stress disorder after recent trauma. *Journal of consulting and clinical psychology*, *76*(2), 231.
- Knight, R. G., Waal-Manning, H. J., & Spears, G. F. (1983). Some norms and reliability data for the State-Trait Anxiety Inventory and the Zung Self-Rating Depression scale. *British Journal of Clinical Psychology*, *22*(4), 245-249.
- Kocher, M. G., Lahno, A. M., & Trautmann, S. T. (2018). Ambiguity aversion is not universal. *European Economic Review*, *101*, 268-283.
- Kocsis, L., & Szepesvári, C. (2006, September). Bandit based monte-carlo planning. In *European conference on machine learning* (pp. 282-293). Springer, Berlin, Heidelberg.
- Korn, C. W., Sharot, T., Walter, H., Heekeren, H. R., & Dolan, R. J. (2014). Depression is related to an absence of optimistically biased belief updating about future life events. *Psychological Medicine*, *44*(3), 579-592.
- Kotov, R., Ruggero, C. J., Krueger, R. F., Watson, D., Yuan, Q., & Zimmerman, M. (2011). New dimensions in the quantitative classification of mental illness. *Archives of general psychiatry*, *68*(10), 1003-1011.
- Kryptos, A. M., Eftting, M., Kindt, M., & Beckers, T. (2015). Avoidance learning: a review of theoretical models and recent developments. *Frontiers in Behavioral Neuroscience*, *9*, 189.
- Kwapil, T. R., Miller, M. B., Zinser, M. C., Chapman, L. J., Chapman, J., & Eckblad, M. (2000). A longitudinal study of high scorers on the Hypomanic Personality Scale. *Journal of Abnormal Psychology*, *109*(2), 222.
- Lau, B., & Glimcher, P. W. (2005). Dynamic response-by-response models of matching behavior in rhesus monkeys. *Journal of the experimental analysis of behavior*, *84*(3), 555-579.

- Lauriola, M., Panno, A., Levin, I. P., & Lejuez, C. W. (2014). Individual differences in risky decision making: A meta-analysis of sensation seeking and impulsivity with the balloon analogue risk task. *Journal of Behavioral Decision Making, 27*(1), 20-36.
- Lawrance, E., O'Reilly, J.X., Bjisterbosch, J., Gagne, C., Bishop, S.J. (in review). The computational and neural substrate of ambiguity avoidance in anxiety.
- Lee, K. M., Coelho, M. A., Sern, K. R., Class, M. A., Bocz, M. D., & Szumlinski, K. K. (2017). Anxiolytic effects of buspirone and MTEP in the porsolt forced swim test. *Chronic Stress, 1*, 2470547017712985.
- Li, J., & Daw, N. D. (2011). Signals in human striatum are appropriate for policy update rather than value prediction. *Journal of Neuroscience, 31*(14), 5504-5511.
- Li, J., Schiller, D., Schoenbaum, G., Phelps, E. A., & Daw, N. D. (2011). Differential roles of human striatum and amygdala in associative learning. *Nature neuroscience, 14*(10), 1250.
- Lichtenstein, S., Slovic, P., Fischhoff, B., Layman, M., & Combs, B. (1978). Judged frequency of lethal events. *Journal of experimental psychology: Human learning and memory, 4*(6), 551.
- Lieder, F., Griffiths, T. L., & Hsu, M. (2018). Overrepresentation of extreme events in decision making reflects rational use of cognitive resources. *Psychological review, 125*(1), 1.
- Lin, L. J. (1992). Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning, 8*(3-4), 293-321.
- Lissek S. (2012) Toward an account of clinical anxiety predicated on basic, neurally mapped mechanisms of Pavlovian fear learning: the case for conditioned overgeneralization. *Depress Anxiety, 29*:257-263.
- Lissek, S., Kaczkurkin, A. N., Rabin, S., Geraci, M., Pine, D. S., & Grillon, C. (2014). Generalized anxiety disorder is associated with overgeneralization of classically conditioned fear. *Biological psychiatry, 75*(11), 909-915.
- MacLeod, A. K. (1996). Affect, emotional disorder, and future-directed thinking. *Cognition & Emotion, 10*(1), 69-86.
- MacLeod, A. K., & Byrne, A. (1996). Anxiety, depression, and the anticipation of future positive and negative experiences. *Journal of abnormal psychology, 105*(2), 286.
- MacLeod, C., & Campbell, L. (1992). Memory accessibility and probability judgments: An experimental evaluation of the availability heuristic. *Journal of Personality and Social Psychology, 63*(6), 890.

- MacLeod, A. K., & Salaminiou, E. (2001). Reduced positive future-thinking in depression: Cognitive and affective factors. *Cognition & Emotion, 15*(1), 99-107.
- MacLeod, A. K., Williams, J. M., & Bekerian, D. A. (1991). Worry is reasonable: The role of explanations in pessimism about future personal events. *Journal of Abnormal psychology, 100*(4), 478.
- Madan, C. R., Ludvig, E. A., & Spetch, M. L. (2014). Remembering the best and worst of times: Memories for extreme outcomes bias risky decisions. *Psychonomic bulletin & review, 21*(3), 629-636.
- Maner, J. K., Richey, J. A., Cromer, K., Mallott, M., Lejuez, C. W., Joiner, T. E., & Schmidt, N. B. (2007). Dispositional anxiety and risk-avoidant decision-making. *Personality and Individual Differences, 42*(4), 665-675.
- Mason, O., Claridge, G., & Jackson, M. (1995). New scales for the assessment of schizotypy. *Personality and Individual differences, 18*(1), 7-13.
- Mathews, A., & Mackintosh, B. (1998). A cognitive model of selective processing in anxiety. *Cognitive therapy and research, 22*(6), 539-560.
- Mathews, A., & MacLeod, C. (2005). Cognitive vulnerability to emotional disorders. *Annu. Rev. Clin. Psychol., 1*, 167-195.
- Mattar, M. G., & Daw, N. D. (2018). Prioritized memory access explains planning and hippocampal replay. *Nature neuroscience, 21*(11), 1609.
- McEvoy, P. M., Watson, H., Watkins, E. R., & Nathan, P. (2013). The relationship between worry, rumination, and comorbidity: Evidence for repetitive negative thinking as a transdiagnostic construct. *Journal of affective disorders, 151*(1), 313-320.
- Meyer, T. J., Miller, M. L., Metzger, R. L., & Borkovec, T. D. (1990). Development and validation of the penn state worry questionnaire. *Behaviour research and therapy, 28*(6), 487-495.
- Mihatsch, O., & Neuneier, R. (2002). Risk-sensitive reinforcement learning. *Machine learning, 49*(2-3), 267-290.
- Mirenowicz, J., & Schultz, W. (1996). Preferential activation of midbrain dopamine neurons by appetitive rather than aversive stimuli. *Nature, 379*(6564), 449.
- Mkrtchian, A., Aylward, J., Dayan, P., Roiser, J. P., & Robinson, O. J. (2017). Modeling avoidance in mood and anxiety disorders using reinforcement learning. *Biological psychiatry, 82*(7), 532-539.

- Mobius (2011). Managing Self-Confidence: Theory and Experimental Evidence. NBER Working Paper No. 17014
- Montague, P. R., Dolan, R. J., Friston, K. J., & Dayan, P. (2012). Computational psychiatry. *Trends in cognitive sciences*, 16(1), 72-80.
- Moore, A. C., MacLeod, A. K., Barnes, D., & Langdon, D. W. (2006). Future-directed thinking and depression in relapsing-remitting multiple sclerosis. *British journal of health psychology*, 11(4), 663-675.
- Moore, A. W., & Atkeson, C. G. (1993). Prioritized sweeping: Reinforcement learning with less data and less time. *Machine learning*, 13(1), 103-130.
- Muris, P., Merckelbach, H., Gadet, B., & Meesters, C. (2000). Monitoring and anxiety disorders symptoms in children. *Personality and individual differences*, 29(4), 775-781
- Muris, P., & van der Heiden, S. (2006). Anxiety, depression, and judgments about the probability of future negative and positive events in children. *Journal of anxiety disorders*, 20(2), 252-261.
- Nassar, M. R., Rumsey, K. M., Wilson, R. C., Parikh, K., Heasley, B., & Gold, J. I. (2012). Rational regulation of learning dynamics by pupil-linked arousal systems. *Nature neuroscience*, 15(7), 1040.
- Niv, Y., Edlund, J. A., Dayan, P., & O'Doherty, J. P. (2012). Neural prediction errors reveal a risk-sensitive reinforcement-learning process in the human brain. *Journal of Neuroscience*, 32(2), 551-562.
- Nord, C. L., Prabhu, G., Nolte, T., Fonagy, P., Dolan, R., & Moutoussis, M. (2017). Vigour in active avoidance. *Scientific reports*, 7(1), 60.
- Oler, J. A., Fox, A. S., Shelton, S. E., Rogers, J., Dyer, T. D., Davidson, R. J., ... & Kalin, N. H. (2010). Amygdalar and hippocampal substrates of anxious temperament differ in their heritability. *Nature*, 466(7308), 864.
- Orr, S. P., Lasko, N. B., Macklin, M. L., Pineles, S. L., Chang, Y., & Pitman, R. K. (2012). Predicting post-trauma stress symptoms from pre-trauma psychophysiologic reactivity, personality traits and measures of psychopathology. *Biology of mood & anxiety disorders*, 2(1), 8.
- Park, C. L., Mills, M. A., & Edmondson, D. (2012). PTSD as meaning violation: Testing a cognitive worldview perspective. *Psychological Trauma: Theory, Research, Practice, and Policy*, 4(1), 66.
- Patzelt, E. H., Kool, W., Millner, A. J., & Gershman, S. J. (2018). Incentives Boost Model-Based Control Across a Range of Severity on Several Psychiatric Constructs. *Biological psychiatry*

Payzan-LeNestour, E., Dunne, S., Bossaerts, P., & O'Doherty, J. P. (2013). The neural representation of unexpected uncertainty during value-based decision making. *Neuron*, *79*(1), 191-201.

Peeters, F., Nicolson, N. A., Berkhof, J., Delespaul, P., & deVries, M. (2003). Effects of daily events on mood states in major depressive disorder. *Journal of abnormal psychology*, *112*(2), 203.

Porsolt, R. D., Le Pichon, M., & Jalfre, M. L. (1977). Depression: a new animal model sensitive to antidepressant treatments. *Nature*, *266*(5604), 730.

Potts, A. J., Bennett, P. J., Kennedy, S. H., & Vaccarino, F. J. (1997). Depressive symptoms and alterations in sucrose taste perception: cognitive bias or a true change in sensitivity?. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale*, *51*(1), 57.

Pulcu, E., & Browning, M. (2017). Affective bias as a rational response to the statistics of rewards and punishments. *Elife*, *6*, e27879.

Radloff, L. S. (1977). The CES-D scale: A self-report depression scale for research in the general population. *Applied psychological measurement*, *1*(3), 385-401.

Raghunathan, R., & Pham, M. T. (1999). All negative moods are not equal: Motivational influences of anxiety and sadness on decision making. *Organizational behavior and human decision processes*, *79*(1), 56-77.

Robinson, O. J., & Chase, H. W. (2017). Learning and Choice in Mood Disorders: Searching for the Computational Parameters of Anhedonia. *Computational Psychiatry*, *1*, 208-233.

Roesch, M. R., Esber, G. R., Li, J., Daw, N. D., & Schoenbaum, G. (2012). Surprise! Neural correlates of Pearce–Hall and Rescorla–Wagner coexist within the brain. *European Journal of Neuroscience*, *35*(7), 1190-1200.

Rubinsztein, J. S., Michael, A., Underwood, B. R., Tempest, M., & Sahakian, B. J. (2006). Impaired cognition and decision-making in bipolar depression but no 'affective bias' evident. *Psychological Medicine*, *36*(5), 629-639.

Rushworth, M. F., & Behrens, T. E. (2008). Choice, uncertainty and value in prefrontal and cingulate cortex. *Nature neuroscience*, *11*(4), 389.

Sandin, B., Sánchez-Arribas, C., Chorot, P., & Valiente, R. M. (2015). Anxiety sensitivity, catastrophic misinterpretations and panic self-efficacy in the prediction of panic disorder severity: Towards a tripartite cognitive model of panic disorder. *Behaviour Research and Therapy*, *67*, 30-40.

Salvatier J., Wiecki T.V., Fonnesbeck C. (2016) Probabilistic programming in Python using PyMC3. *PeerJ Computer Science* 2:e55 DOI: 10.7717/peerj-cs.55.

Savage, L. J. (1954). The foundations of. *Statistics*, 11-34.

Schaefer, K. L., Baumann, J., Rich, B. A., Luckenbaugh, D. A., & Zarate Jr, C. A. (2010). Perception of facial emotion in adults with bipolar or unipolar depression and controls. *Journal of Psychiatric Research*, 44(16), 1229-1235.

Schmid, J., & Leiman, J. M. (1957). The development of hierarchical factor solutions. *Psychometrika*, 22(1), 53-61.

Schultz, W. (1998). Predictive reward signal of dopamine neurons. *Journal of neurophysiology*, 80(1), 1-27.

Schultz, W. (2015). Neuronal reward and decision signals: from theories to data. *Physiological reviews*, 95(3), 853-951.

Schultz, W. (2016). Reward functions of the basal ganglia. *Journal of Neural Transmission*, 123(7), 679-693.

Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and reward. *Science*, 275(5306), 1593-1599.

Scinska, A., Sienkiewicz-Jarosz, H., Kuran, W., Ryglewicz, D., Rogowski, A., Wrobel, E., ... & Bienkowski, P. (2004). Depressive symptoms and taste reactivity in humans. *Physiology & behavior*, 82(5), 899-904.

Servatius, R. J., Jiao, X., Beck, K. D., Pang, K. C. H., & Minor, T. R. (2008). Rapid avoidance acquisition in Wistar–Kyoto rats. *Behavioural Brain Research*, 192(2), 191-197.

Seymour, B., O'Doherty, J. P., Dayan, P., Koltzenburg, M., Jones, A. K., Dolan, R. J., ... & Frackowiak, R. S. (2004). Temporal difference models describe higher-order learning in humans. *Nature*, 429(6992), 664.

Shah, P., Harris, A. J., Bird, G., Catmur, C., & Hahn, U. (2016). A pessimistic view of optimistic belief updating. *Cognitive psychology*, 90, 71-127.

Shapiro, D. N., Chandler, J., & Mueller, P. A. (2013). Using Mechanical Turk to study clinical populations. *Clinical psychological science*, 1(2), 213-220.

Sharot, T., Korn, C. W., & Dolan, R. J. (2011). How unrealistic optimism is maintained in the face of reality. *Nature neuroscience*, 14(11), 1475.

Shin, L. M., Rauch, S. L., & Pitman, R. K. (2006). Amygdala, medial prefrontal cortex, and hippocampal function in PTSD. *Annals of the New York Academy of Sciences*, 1071(1), 67-79.

Shin, L. M., & Liberzon, I. (2010). The neurocircuitry of fear, stress, and anxiety disorders. *Neuropsychopharmacology*, 35(1), 169.

Simms, L. J., Grös, D. F., Watson, D., & O'hara, M. W. (2008). Parsing the general and specific components of depression and anxiety with bifactor modeling. *Depression and anxiety*, 25(7), E34-E46.

Speilberger, C. D., Gorsuch, R. L., Lushene, R., Vagg, P. R., & Jacobs, G. A. (1983). State-trait anxiety inventory for adults. *Redwood City: Mind Garden Inc.*

Spinhoven, P., Penninx, B. W., Krempeniou, A., van Hemert, A. M., & Elzinga, B. (2015). Trait rumination predicts onset of Post-Traumatic Stress Disorder through trauma-related cognitive appraisals: A 4-year longitudinal study. *Behaviour Research and Therapy*, 71, 101-109.

Steele, J. D., Kumar, P., & Ebmeier, K. P. (2007). Blunted response to feedback information in depressive illness. *Brain*, 130(9), 2367-2374.

Steer, R. A., Clark, D. A., Beck, A. T., & Ranieri, W. F. (1995). Common and specific dimensions of self-reported anxiety and depression: A replication. *Journal of Abnormal Psychology*, 104(3), 542.

Steer, R. A., Clark, D. A., Beck, A. T., & Ranieri, W. F. (1998). Common and specific dimensions of self-reported anxiety and depression: the BDI-II versus the BDI-IA. *Behaviour research and therapy*, 37(2), 183-190.

Steer, R. A., Clark, D. A., Kumar, G., & Beck, A. T. (2008). Common and specific dimensions of self-reported anxiety and depression in adolescent outpatients. *Journal of Psychopathology and Behavioral Assessment*, 30(3), 163-170.

Stephan, K. E., Penny, W. D., Daunizeau, J., Moran, R. J., & Friston, K. J. (2009). Bayesian model selection for group studies. *Neuroimage*, 46(4), 1004-1017.

Stevens, J. S., Reddy, R., Kim, Y. J., van Rooij, S. J., Ely, T. D., Hamann, S., ... & Jovanovic, T. (2018). Episodic memory after trauma exposure: Medial temporal lobe function is positively related to re-experiencing and inversely related to negative affect symptoms. *NeuroImage: Clinical*, 17, 650-658.

Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine learning*, 3(1), 9-44.

Sutton, R. S. (1990). Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Machine Learning Proceedings 1990* (pp. 216-224). Morgan Kaufmann.

Sutton, R. S., & Barto, A. G. (1998). Reinforcement learning: an introduction MIT Press. Cambridge, MA.

Szabó, M., & Lovibond, P. F. (2004). The cognitive content of thought-listed worry episodes in clinic-referred anxious and nonreferred children. *Journal of Clinical Child and Adolescent Psychology, 33*(3), 613-622.

Szabó, M., & Lovibond, P. F. (2006). Worry episodes and perceived problem solving: A diary-based approach. *Anxiety, stress, and coping, 19*(2), 175-187.

Talmi, D., Seymour, B., Dayan, P., & Dolan, R. J. (2008). Human Pavlovian–instrumental transfer. *Journal of Neuroscience, 28*(2), 360-368.

Teasdale, J. D., Taylor, R., & Fogarty, S. J. (1980). Effects of induced elation-depression on the accessibility of memories of happy and unhappy experiences. *Behaviour research and therapy, 18*(4), 339-346.

Teasdale, J. D. (1983). Negative thinking in depression: Cause, effect, or reciprocal relationship?. *Advances in Behaviour Research and Therapy, 5*(1), 3-25.

Tedeschi, R. G., & Calhoun, L. G. (2004). " Posttraumatic growth: Conceptual foundations and empirical evidence". *Psychological inquiry, 15*(1), 1-18.

Travis, J. M., Palmer, S. C., Coyne, S., Millon, A., & Lambin, X. (2013). Evolution of predator dispersal in relation to spatio-temporal prey dynamics: how not to get stuck in the wrong place!. *PLoS One, 8*(2), e54453.

Treadway, M. T., Bossaller, N. A., Shelton, R. C., & Zald, D. H. (2012). Effort-based decision-making in major depressive disorder: a translational model of motivational anhedonia. *Journal of abnormal psychology, 121*(3), 553.

Treadway, M. T., Buckholtz, J. W., Schwartzman, A. N., Lambert, W. E., & Zald, D. H. (2009). Worth the 'EEfRT'? The effort expenditure for rewards task as an objective measure of motivation and anhedonia. *PloS one, 4*(8), e6598.

Treynor, W., Gonzalez, R., & Nolen-Hoeksema, S. (2003). Rumination reconsidered: A psychometric analysis. *Cognitive therapy and research, 27*(3), 247-259.

- Tversky, A., & Kahneman, D. (1974). Judgment under uncertainty: Heuristics and biases. *science*, *185*(4157), 1124-1131.
- Tversky, A., & Kahneman, D. (1992). Advances in prospect theory: Cumulative representation of uncertainty. *Journal of Risk and uncertainty*, *5*(4), 297-323.
- Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using leave-one-out cross-validation and dLOO. *Statistics and Computing*, *27*(5), 1413-1432.
- Verkuil, B., Brosschot, J. F., & Thayer, J. F. (2007). Capturing worry in daily life: Are trait questionnaires sufficient?. *Behaviour Research and Therapy*, *45*(8), 1835-1844.
- Wang, X. J., & Krystal, J. H. (2014). Computational psychiatry. *Neuron*, *84*(3), 638-654.
- Walton, M. E., Bannerman, D. M., Alterescu, K., & Rushworth, M. F. (2003). Functional specialization within medial frontal cortex of the anterior cingulate for evaluating effort-related decisions. *Journal of Neuroscience*, *23*(16), 6475-6479.
- Walton, M. E., Groves, J., Jennings, K. A., Croxson, P. L., Sharp, T., Rushworth, M. F., & Bannerman, D. M. (2009). Comparing the role of the anterior cingulate cortex and 6-hydroxydopamine nucleus accumbens lesions on operant effort-based decision making. *European Journal of Neuroscience*, *29*(8), 1678-1691.
- Watkins, E. R. (2008). Constructive and unconstructive repetitive thought. *Psychological bulletin*, *134*(2), 163.
- Watkins, E. D., Moulds, M., & Mackintosh, B. (2005). Comparisons between rumination and worry in a non-clinical population. *Behaviour research and therapy*, *43*(12), 1577-1585.
- Watkins, E. (2011). Dysregulation in level of goal and action identification across psychological disorders. *Clinical psychology review*, *31*(2), 260-278.
- Watson, D., & Clark, L. A. (1991). The mood and anxiety symptom questionnaire (MASQ). *Unpublished manuscript, University of Iowa, Iowa City*.
- Wiech, K., & Tracey, I. (2009). The influence of negative emotions on pain: behavioral effects and neural mechanisms. *Neuroimage*, *47*(3), 987-994.
- Williams, J. M. G., Barnhofer, T., Crane, C., Herman, D., Raes, F., Watkins, E., & Dalgleish, T. (2007). Autobiographical memory specificity and emotional disorder. *Psychological bulletin*, *133*(1), 122.
- Wilson, M. A., & McNaughton, B. L. (1994). Reactivation of hippocampal ensemble memories during sleep. *Science*, *265*(5172), 676-679.

Yu, A. J., & Dayan, P. (2005). Uncertainty, neuromodulation, and attention. *Neuron*, *46*(4), 681-692.

Zinbarg, R. E., & Barlow, D. H. (1996). Structure of anxiety and the anxiety disorders: a hierarchical model. *Journal of Abnormal Psychology*, *105*(2), 181.