

UC San Diego

UC San Diego Previously Published Works

Title

An adaptive nearest neighbor rule for classification

Permalink

<https://escholarship.org/uc/item/64622969>

Authors

Balsubramani, Akshay
Dasgupta, Sanjoy
Freund, Yoav
[et al.](#)

Publication Date

2019

Peer reviewed

An adaptive nearest neighbor rule for classification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce a variant of the k -nearest neighbor classifier in which k is chosen
2 adaptively for each query, rather than supplied as a parameter. The choice of
3 k depends on properties of each neighborhood, and therefore may significantly
4 vary between different points. (For example, the algorithm will use larger k for
5 predicting the labels of points in noisy regions.)

6 We provide theory and experiments that demonstrate that the algorithm performs
7 comparably to, and sometimes better than, k -NN with an optimal choice of k . In
8 particular, we derive bounds on the convergence rates of our classifier that depend
9 on a local quantity we call the “advantage” which is significantly weaker than the
10 Lipschitz conditions used in previous convergence rate proofs. These generalization
11 bounds hinge on a variant of the seminal Uniform Convergence Theorem due to
12 Vapnik and Chervonenkis; this variant concerns conditional probabilities and may
13 be of independent interest.

14 1 Introduction

15 We introduce an adaptive nearest neighbor classification rule. Given a training set with labels $\{\pm 1\}$,
16 its prediction at a query point x is based on the training points closest to x , rather like the k -nearest
17 neighbor rule. However, the value of k that it uses can vary from query to query. Specifically, if there
18 are n training points, then for any query x , the smallest k is sought for which the k points closest to x
19 have labels whose average is either greater than $+\Delta(n, k)$, in which case the prediction is $+1$, or less
20 than $-\Delta(n, k)$, in which case the prediction is -1 ; and if no such k exists, then “?” (“don’t know”)
21 is returned. Here, $\Delta(n, k) \sim \sqrt{(\log n)/k}$ corresponds to a confidence interval for the average label
22 in the region around the query.

23 We study this rule in the standard statistical framework in which all data are i.i.d. draws from some
24 unknown underlying distribution P on $\mathcal{X} \times \mathcal{Y}$, where \mathcal{X} is the data space and \mathcal{Y} is the label space. We
25 take \mathcal{X} to be a separable metric space, with distance function $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, and we take $\mathcal{Y} = \{\pm 1\}$.
26 We can decompose P into the marginal distribution μ on \mathcal{X} and the conditional expectation of the
27 label at each point x : if (X, Y) represents a random draw from P , define $\eta(x) = \mathbb{E}(Y|X = x)$. In
28 this terminology, the Bayes-optimal classifier is the rule $g^* : \mathcal{X} \rightarrow \{\pm 1\}$ given by

$$g^*(x) = \begin{cases} \text{sign}(\eta(x)) & \text{if } \eta(x) \neq 0 \\ \text{either } -1 \text{ or } +1 & \text{if } \eta(x) = 0 \end{cases} \quad (1)$$

29 and its error rate is the Bayes risk, $R^* = \frac{1}{2} \mathbb{E}_{X \sim \mu} [1 - |\eta(X)|]$. A variety of nonparametric classi-
30 fication schemes are known to have error rates that converge asymptotically to R^* . These include
31 k -nearest neighbor (henceforth, k -NN) rules [FH51] in which k grows with the number of training
32 points n according to a suitable schedule (k_n) , under certain technical conditions on the metric
33 measure space (\mathcal{X}, d, μ) .

34 In this paper, we are interested in consistency as well as rates of convergence. In particular, we find
35 that the adaptive nearest neighbor rule is also asymptotically consistent (under the same technical

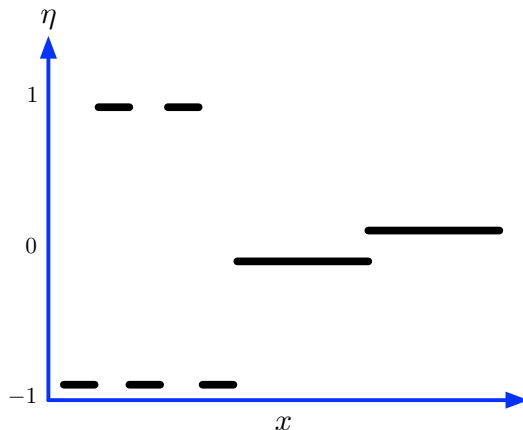


Figure 1: For values of x on the left half of the shown interval, the pointwise bias $\eta(x)$ is close to -1 or 1 , and thus a small value of k will yield an accurate prediction. Larger k will not do as well, because they may run into neighboring regions with different labels. For values of x on the right half of the interval, $\eta(x)$ is close to 0 , and thus large k is essential for accurate prediction.

36 conditions) while converging at a rate that is about as good as, and sometimes significantly better
 37 than, that of k -NN under any schedule (k_n) .

38 Intuitively, one of the advantages of k -NN over nonparametric classifiers that use a fixed bandwidth
 39 or radius, such as Parzen window or kernel density estimators, is that k -NN automatically adapts to
 40 variation in the marginal distribution μ : in regions with large μ , the k nearest neighbors lie close to
 41 the query point, while in regions with small μ , the k nearest neighbors can be further afield. The
 42 adaptive NN rule that we propose goes further: it also adapts to variation in η . In certain regions of
 43 the input space, where η is close to 0 , an accurate prediction would need large k . In other regions,
 44 where η is near -1 or 1 , a small k would suffice, and in fact, a larger k might be detrimental because
 45 neighboring regions might be labeled differently. See Figure 1 for one such example. A k -NN
 46 classifier is forced to pick a single value of k that trades off between these two contingencies. Our
 47 adaptive NN rule, however, can pick the right k in each neighborhood separately.

48 Our estimator allows us to give rates of convergence that are tighter and more transparent than those
 49 customarily obtained in nonparametric statistics. Specifically, for any point x in the instance space
 50 \mathcal{X} , we define a notion of the *advantage at x* , denoted $\text{adv}(x)$, which is rather like a local margin.
 51 We show that the prediction at x is very likely to be correct once the number of training points
 52 exceeds $\tilde{O}(1/\text{adv}(x))$. Universal consistency follows by establishing that almost all points have
 53 positive advantage.

54 1.1 Relation to other work in nonparametric estimation

55 For linear separators and many other *parametric* families of classifiers, it is possible to give rates
 56 of convergence that hold without any assumptions on the input distribution μ or the conditional
 57 expectation function η . This is not true of nonparametric estimation: although any target function can
 58 in principle be captured, the number of samples needed to achieve a specific level of accuracy will
 59 inevitably depend upon aspects of this function such as how fast it changes [DGL96, chapter 7]. As a
 60 result, nonparametric statistical theory has focused on (1) asymptotic consistency, ideally without
 61 assumptions, and (2) rates of convergence under a variety of smoothness assumptions.

62 Asymptotic consistency has been studied in great detail for the k -NN classifier, when k is allowed
 63 to grow with the number of data points n . The risk of the classifier, denoted R_n , is its error rate
 64 on the underlying distribution P ; this is a random variable that depends upon the set of training
 65 points seen. Cover and Hart [CH67] showed that in general metric spaces, under the assumption
 66 that every x in the support of μ is either a continuity point of η or has $\mu(\{x\}) > 0$, the expected
 67 risk $\mathbb{E}R_n$ converges to the Bayes-optimal risk R^* , as long as $k \rightarrow \infty$ and $k/n \rightarrow 0$. For points
 68 in finite-dimensional Euclidean space, a series of results starting with Stone [Sto77] established
 69 consistency without any assumptions on μ or η , and showed that $R_n \rightarrow R^*$ almost surely [DGKL94].

70 More recent work has extended these *universal consistency* results—that is, consistency without
 71 assumptions on η —to arbitrary metric measure spaces (\mathcal{X}, d, μ) that satisfy a certain differentiation
 72 condition [CG06, CD14].

73 Rates of convergence have been obtained for k -nearest neighbor classification under various smooth-
 74 ness conditions including Holder conditions on η [KP95, Gyö81] and “Tsybakov margin” condi-
 75 tions [MT99, AT07, CD14]. Such assumptions have become customary in nonparametric statistics,
 76 but they leave a lot to be desired. First, they are uncheckable: it is not possible to empirically
 77 determine the smoothness given samples. Second, they view the underlying distribution P through
 78 the tiny window of two or three parameters, obscuring almost all the remaining structure of the
 79 distribution that also influences the rate of convergence. Finally, because nonparametric estimation is
 80 *local*, there is the intriguing possibility of getting different rates of convergence in different regions
 81 of the input space: a possibility that is immediately defeated by reducing the entire space to two
 82 smoothness constants.

83 The first two of these issues are partially addressed by recent work of [CD14], who analyze the finite
 84 sample risk of k -NN classification without any assumptions on P . Their bounds involve terms that
 85 measure the probability mass of the input space in a carefully defined region around the decision
 86 boundary, and are shown to be “instance-optimal”: that is, optimal for the specific distribution P ,
 87 rather than minimax-optimal for some very large class containing P . However, the expressions for
 88 the risk are somewhat hard to parse, in large part because of the interaction between n and k .

89 In the present paper, we obtain finite-sample rates of convergence that are “instance-optimal” not
 90 just for the specific distribution P but also for the specific query point. This is achieved by defining
 91 a *margin*, or *advantage*, at every point in the input space, and giving bounds (Theorem 1) entirely
 92 in terms of this quantity. For parametric classification, it has become common to define a notion
 93 of margin that controls generalization. In the nonparametric setting, it makes sense that the margin
 94 would in fact be a function $\mathcal{X} \rightarrow \mathbb{R}$, and would yield different generalization error bounds in different
 95 regions of space. Our adaptive nearest neighbor classifier allows us to realize this vision in a fairly
 96 elementary manner.

97 **Organization.** Due to space limitations, all proofs are relegated to appendices.

98 We begin by formally defining the setup and notation in Section 2. Then, a formal description of
 99 the adaptive k -NN algorithm is given in Section 3. In Sections 4 and 5 and appendix A, we state
 100 and prove consistency and generalization bounds for this classifier, and compare them with previous
 101 bounds in the k -NN literature. These bounds exploit a general VC-based uniform convergence
 102 statement which is presented and proved in a self-contained manner in Appendix B.

103 2 Setup

104 Take the instance space to be a separable metric space (\mathcal{X}, d) and the label space to be $\mathcal{Y} = \{\pm 1\}$.
 105 All data are assumed to be drawn i.i.d. from a fixed unknown distribution P over $\mathcal{X} \times \mathcal{Y}$.

Let μ denote the marginal distribution on \mathcal{X} : if (X, Y) is a random draw from P , then

$$\mu(S) = \Pr(X \in S)$$

for any measurable set $S \subseteq \mathcal{X}$. For any $x \in \mathcal{X}$, the conditional expectation, or *bias*, of Y given x , is

$$\eta(x) = \mathbb{E}(Y|X = x) \in [-1, 1].$$

Similarly, for any measurable set S with $\mu(S) > 0$, the conditional expectation of Y given $X \in S$ is

$$\eta(S) = \mathbb{E}(Y|X \in S) = \frac{1}{\mu(S)} \int_S \eta(x) d\mu(x).$$

106 The risk of a classifier $g : \mathcal{X} \rightarrow \{-1, +1, ?\}$ is the probability that it is incorrect on pairs $(X, Y) \sim P$,
 107

$$R(g) = P(\{(x, y) : g(x) \neq y\}). \quad (2)$$

108 The Bayes-optimal classifier g^* , as given in (1), depends only on η , but its risk R^* depends on μ . For
 109 a classifier g_n based on n training points from P , we will be interested in whether $R(g_n)$ converges
 110 to R^* , and the rate at which this convergence occurs.

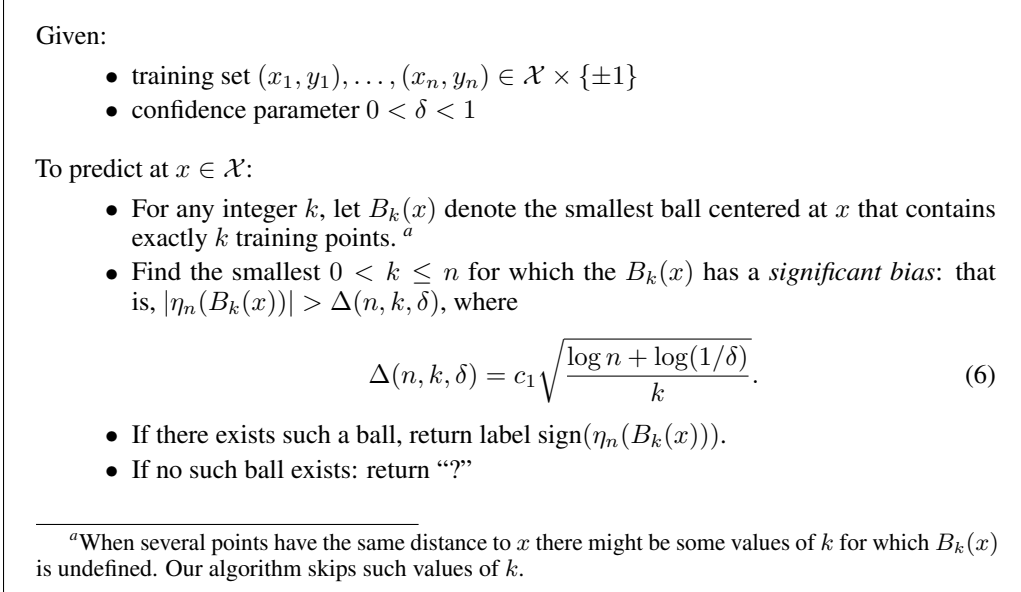


Figure 2: The adaptive k -NN (AKNN) classifier. The absolute constant c_1 is from Lemma 6.

The algorithm and analysis in this paper depend heavily on the probability masses and biases of balls in \mathcal{X} . For $x \in \mathcal{X}$ and $r \geq 0$, let $B(x, r)$ denote the closed ball of radius r centered at x ,

$$B(x, r) = \{z \in \mathcal{X} : d(x, z) \leq r\}.$$

111 For $0 \leq p \leq 1$, let $r_p(x)$ be the smallest radius r such that $B(x, r)$ has probability mass at least p ,
 112 that is,

$$r_p(x) = \inf\{r \geq 0 : \mu(B(x, r)) \geq p\}. \quad (3)$$

113 It follows that $\mu(B(x, r_p(x))) \geq p$.

The *support* of the marginal distribution μ plays an important role in convergence proofs and is formally defined as

$$\text{supp}(\mu) = \{x \in \mathcal{X} : \mu(B(x, r)) > 0 \text{ for all } r > 0\}.$$

114 It is a well-known consequence of the separability of \mathcal{X} that $\mu(\text{supp}(\mu)) = 1$ [CH67].

115 3 The adaptive k -nearest neighbor algorithm

116 The algorithm is given a labeled training set $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$. Based on these points,
 117 it is able to compute empirical estimates of the probabilities and biases of different balls.

118 For any set $S \subseteq \mathcal{X}$, we define its empirical count and probability mass as

$$\begin{aligned} \#_n(S) &= |\{i : x_i \in S\}| \\ \mu_n(S) &= \frac{\#_n(S)}{n}. \end{aligned} \quad (4)$$

119 If this is non-zero, we take the empirical bias to be

$$\eta_n(S) = \frac{\sum_{i: x_i \in S} y_i}{\#_n(S)}. \quad (5)$$

120 The adaptive k -NN algorithm (AKNN) is shown in Figure 2. It makes a prediction at x by growing a
 121 ball around x until the ball has significant bias, and then choosing the corresponding label. In some
 122 cases, a ball of sufficient bias may never be obtained, in which event “?” is returned. In what follows,
 123 let $g_n : \mathcal{X} \rightarrow \{-1, +1, ?\}$ denote the AKNN classifier.

124 Later, we will also discuss a variant of this algorithm in which a modified confidence interval,

$$\Delta(n, k, \delta) = c_1 \sqrt{\frac{d_0 \log n + \log(1/\delta)}{k}}, \quad (7)$$

125 is used, where d_0 is the VC dimension of the family of balls in (\mathcal{X}, d) .

126 4 Pointwise advantage and rates of convergence

127 We now provide finite-sample rates of convergence for the adaptive nearest neighbor rule. For
128 simplicity, we give convergence rates that are specific to any query point x and that depend on a
129 suitable notion of the “margin” of distribution P around x .

130 Pick any $p, \gamma > 0$. Recalling definition (3), we say a point $x \in \mathcal{X}$ is (p, γ) -salient if the following
131 holds for either $s = +1$ or $s = -1$:

- 132 • $s\eta(x) > 0$, and $s\eta(B(x, r)) > 0$ for all $r \in [0, r_p(x))$, and $s\eta(B(x, r_p(x))) \geq \gamma$.

133 In words, this means that $g^*(x) = s$ (recall that g^* is the Bayes classifier), that the biases of all balls
134 of radius $\leq r_p(x)$ around x have the same sign as s , and that the bias of the ball of radius $r_p(x)$ has
135 margin at least γ . A point x can satisfy this definition for a variety of pairs (p, γ) . The *advantage*
136 of x is taken to be the largest value of $p\gamma^2$ over all such pairs:

$$\text{adv}(x) = \begin{cases} \sup\{p\gamma^2 : x \text{ is } (p, \gamma)\text{-salient}\} & \text{if } \eta(x) \neq 0 \\ 0 & \text{if } \eta(x) = 0 \end{cases} \quad (8)$$

137 We will see (Lemma 3) that under a mild condition on the underlying metric measure space, almost
138 all x with $\eta(x) \neq 0$ have a positive advantage.

139 4.1 Advantage-based finite-sample bounds

140 The following theorem shows that for every point x , if the sample size n satisfies $n \gtrsim 1/\text{adv}(x)$,
141 then the label of x is likely to be $g^*(x)$, where g^* is the Bayes optimal classifier. This provides a
142 pointwise convergence of $g(x)$ to $g^*(x)$ with a rate which is sensitive to the “local geometry” of x .

143 **Theorem 1** (Pointwise convergence rate). *There is an absolute constant $C > 0$ for which the*
144 *following holds. Let $0 < \delta < 1$ denote the confidence parameter in the AKNN algorithm (Figure 2),*
145 *and suppose the algorithm is used to define a classifier g_n based on n training points chosen i.i.d.*
146 *from P . Then, for every point $x \in \text{supp}(\mu)$, if*

$$n \geq \frac{C}{\text{adv}(x)} \max\left(\log \frac{1}{\text{adv}(x)}, \log \frac{1}{\delta}\right)$$

147 *then with probability at least $1 - \delta$ we have that $g_n(x) = g^*(x)$.*

148 If we further assume that the family of all balls in the space has finite VC dimension d_0 then we
149 can strengthen Theorem 1 so that the guarantee holds with high probability simultaneously
150 for all $x \in \text{supp}(\mu)$. This is achieved by a modified version of the algorithm that uses confidence interval
151 (7) instead of (6).

152 **Theorem 2** (Uniform convergence rate). *Suppose that the set of balls in (\mathcal{X}, d) has finite VC*
153 *dimension d_0 , and that the algorithm of Figure 2 is used with confidence interval (7) instead of (6).*
154 *Then, with probability at least $1 - \delta$, the resulting classifier g_n satisfies the following: for every*
155 *point $x \in \text{supp}(\mu)$, if*

$$n \geq \frac{C}{\text{adv}(x)} \max\left(\log \frac{1}{\text{adv}(x)}, \log \frac{1}{\delta}\right)$$

156 *then $g_n(x) = g^*(x)$.*

157 A key step towards proving Theorems 1 and 2 is to identify the subset of \mathcal{X} that is likely to be
158 correctly classified for a given number of training points n . This follows the rough outline of [CD14],
159 which gave rates of convergence for k -nearest neighbor, but there are two notable differences. First,
160 we will see that the likely-correct sets obtained in that earlier work (for k -NN) are subsets of those
161 we obtain for the new adaptive nearest neighbor procedure. Second, the proof for our setting is
162 considerably more streamlined; for instance, there is no need to devise tie-breaking strategies for
163 deciding the identities of the k nearest neighbors.

164 **4.2 A comparison with k -nearest neighbor**

165 For $a \geq 0$, let \mathcal{X}_a denote all points with advantage greater than a :

$$\mathcal{X}_a = \{x \in \text{supp}(\mu) : \text{adv}(x) > a\}. \quad (9)$$

166 In particular, \mathcal{X}_0 consists of all points with positive advantage.

167 By Theorem 1, points in \mathcal{X}_a are likely to be correctly classified when the number of training points
 168 is $\tilde{\Omega}(1/a)$, where the $\tilde{\Omega}(\cdot)$ notation ignores logarithmic terms. In contrast, the work of [CD14]
 169 showed that with n training points, the k -NN classifier is likely to correctly classify the following set
 170 of points:

$$\begin{aligned} \mathcal{X}'_{n,k} = & \{x \in \text{supp}(\mu) : \eta(x) > 0, \eta(B(x, r)) \geq k^{-1/2} \text{ for all } 0 \leq r \leq r_{k/n}(x)\} \\ & \cup \{x \in \text{supp}(\mu) : \eta(x) < 0, \eta(B(x, r)) \leq -k^{-1/2} \text{ for all } 0 \leq r \leq r_{k/n}(x)\}. \end{aligned}$$

Such points are $(k/n, k^{-1/2})$ -salient and thus have advantage at least $1/n$. In fact,

$$\bigcup_{1 \leq k \leq n} \mathcal{X}'_{n,k} \subseteq \mathcal{X}_{1/n}.$$

171 In this sense, the adaptive nearest neighbor procedure is able to perform roughly as well as all choices
 172 of k simultaneously (logarithmic factors prevent this from being a precise statement).

173 **5 Universal consistency**

174 In this section we study the convergence of $R(g_n)$ to the Bayes risk R^* as the number of points n
 175 grows. An estimator is described as universally consistent in a metric measure space (\mathcal{X}, d, μ) if it
 176 has this desired limiting behavior for all conditional expectation functions η .

177 Earlier work [CD14] has established the universal consistency of k -nearest neighbor (for $k/n \rightarrow 0$
 178 and $k/(\log n) \rightarrow \infty$) in any metric measure space that satisfies the Lebesgue differentiation condition:
 179 that is, for any bounded measurable $f : \mathcal{X} \rightarrow \mathbb{R}$ and for almost all (μ -a.e.) $x \in \mathcal{X}$,

$$\lim_{r \downarrow 0} \frac{1}{\mu(B(x, r))} \int_{B(x, r)} f d\mu = f(x). \quad (10)$$

180 This is known to hold, for instance, in any finite-dimensional normed space or any doubling metric
 181 space [Hei01, Chapter 1].

182 We will now see that this same condition implies the universal consistency of the adaptive nearest
 183 neighbor rule. To begin with, it implies that almost every point has a positive advantage.

Lemma 3. *Suppose metric measure space (\mathcal{X}, d, μ) satisfies condition (10). Then, for any conditional
 expectation η , the set of points*

$$\{x \in \mathcal{X} : \eta(x) \neq 0, \text{adv}(x) = 0\}$$

184 *has zero μ -measure.*

185 *Proof.* Let $\mathcal{X}' \subseteq \mathcal{X}$ consist of all points $x \in \text{supp}(\mu)$ for which condition (10) holds true with $f = \eta$,
 186 that is, $\lim_{r \downarrow 0} \eta(B(x, r)) = \eta(x)$. Since $\mu(\text{supp}(\mu)) = 1$, it follows that $\mu(\mathcal{X}') = 1$.

Pick any $x \in \mathcal{X}'$ with $\eta(x) \neq 0$; without loss of generality, $\eta(x) > 0$. By (10), there exists $r_o > 0$
 such that

$$\eta(B(x, r)) \geq \eta(x)/2 \text{ for all } 0 \leq r \leq r_o.$$

187 Thus x is (p, γ) -salient for $p = \mu(B(x, r_o)) > 0$ and $\gamma = \eta(x)/2$, and has positive advantage. \square

188 Universal consistency follows as a consequence; the proof details are deferred to Section A.

189 **Theorem 4** (Universal consistency). *Suppose the metric measure space (\mathcal{X}, d, μ) satisfies condi-
 190 tion (10). Let (δ_n) be a sequence in $[0, 1]$ with (1) $\sum_n \delta_n < \infty$ and (2) $\lim_{n \rightarrow \infty} (\log(1/\delta_n))/n = 0$.
 191 Let the classifier $g_{n, \delta_n} : \mathcal{X} \rightarrow \{-1, +1, ?\}$ be the result of applying the AKNN procedure (Figure 2)
 192 with n points chosen i.i.d. from P and with confidence parameter δ_n . Letting $R_n = R(g_{n, \delta_n})$ denote
 193 the risk of g_{n, δ_n} , we have $R_n \rightarrow R^*$ almost surely.*

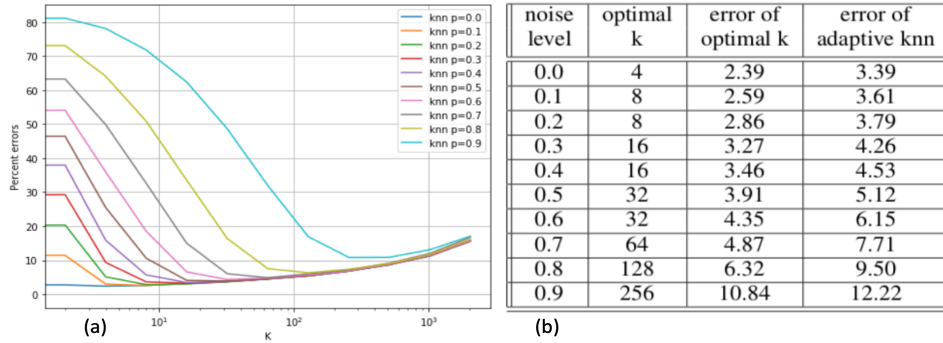
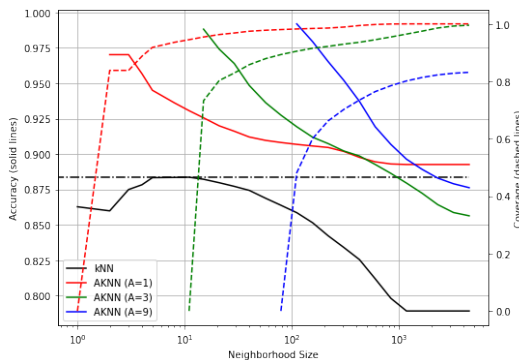


Figure 3: Effect of label noise on k -NN and AKNN. Performance on MNIST for different levels of random label noise p and for different values of k . Each line in the figure on the left (a) represents the performance of k -NN as a function of k for a given level of noise. The optimal choice of k increases with the noise level, and that the performance degrades severely for too-small k . The table (b) shows that AKNN, with a fixed value of A , performs almost as well as k -NN with the optimal choice of k .



At left: performance of AKNN on notMNIST for different settings of the confidence parameter ($A = 1, 3, 9$), as a function of the neighborhood size. For each confidence level we show two graphs: an accuracy graph (solid lines) and a coverage line (dashed line). For each value of k we plot the accuracy and the coverage of AKNN which is restricted to using a neighborhood size of at most k . Increasing A generally causes an increase in the accuracy and a decrease in coverage. Larger values of A cause AKNN to have coverage zero for values of k that are too small. For comparison, we plot the performance of k -NN as a function of k . The highest accuracy (≈ 0.88) is achieved for $k = 10$ (dotted horizontal line), and is surpassed by AKNN with high coverage (100% for $A = 1$).

Figure 4: Performance of AKNN on notMNIST. See also Figure 5.

194 6 Experiments

195 We performed a few experiments using real-world datasets from computer vision and genomics (see
 196 Section C). These were conducted with some practical alterations to the algorithm of Fig. 2.

197 **Multiclass extension:** Suppose the set of possible labels is \mathcal{Y} . We replace the binary rule “find
 198 the smallest k such that $|\eta_n(B_k(x))| > \Delta(n, k, \delta)$ ” with the rule: “find the smallest k such that
 199 $\eta_n^y(B_k(x)) - \frac{1}{|\mathcal{Y}|} > \Delta(n, k, \delta)$ for some $y \in \mathcal{Y}$, where $\eta_n^y(S) \doteq \frac{\#\{x_i \in S \text{ and } y_i = y\}}{\#_n(S)}$.”

200 **Parametrization:** We replace Equation (6) with $\Delta = \frac{A}{\sqrt{k}}$, where A is a confidence parameter.

201 **Resolving multilabel predictions:** Our algorithm can output answers that are not a single label. The
 202 output can be “?”, which indicates that no label has sufficient evidence. It can also be a subset of \mathcal{Y}
 203 that contains more than one element, indicating that more than one label has significant evidence. In
 204 some situations, using subsets of the labels is more informative. However, when we want to compare
 205 head-to-head with k -NN, we need to output a single label. We use a heuristic to predict with a single
 206 label $y \in \mathcal{Y}$ on any x : the label for which $\max_k \eta_n^y(B_k(x)) / \sqrt{k}$ is largest.

207 We briefly discuss our main conclusions from the experiments, with further details deferred to
 208 Appendix C.

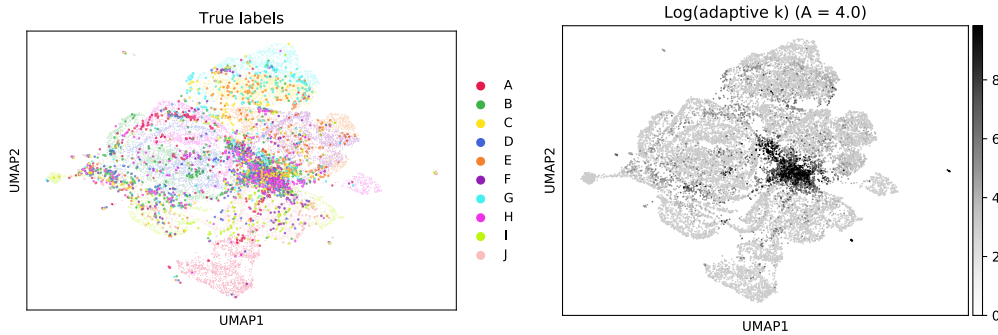


Figure 5: A visualization of the performance of AKNN on notMNIST. **(a)** The correct labels, with prediction errors of AKNN ($A = 4$) highlighted. **(b)** The value of k chosen by the algorithm when predicting each datapoint. Zooming in reveals more details. An interactive explorer for our experiments is available at <http://35.239.251.24/aknn/>.

209 **AKNN is comparable to the best k -NN rule.** In Section 4.2 we prove that AKNN compares
 210 favorably to k -NN with any fixed k . We demonstrate this in practice in different situations. With
 211 simulated independent label noise on the MNIST dataset (Fig. 3), a small value of k is optimal for
 212 noiseless data, but performs very poorly when the noise level is high. On the other hand, AKNN
 213 adapts to the local noise level automatically, as demonstrated without adding noise on the more
 214 challenging notMNIST and single-cell genomics data (Fig. 4, 5, 6).

215 **Varying the confidence parameter A controls abstaining.** The parameter A controls how conser-
 216 vative the algorithm is deciding to abstain, instead of incurring error by predicting. $A \rightarrow 0$ represents
 217 the most aggressive setting, in which the algorithm never abstains, essentially predicting according to
 218 a 1-NN rule. Higher settings of A cause the algorithm to abstain on some of these predicted points,
 219 for which there is no sufficiently small neighborhood with a sufficiently significant label bias (Fig. 7).

220 **Adaptively chosen neighborhood sizes reflect local confidence.** The number of neighbors chosen
 221 by AKNN is a local quantity that gives a practical pointwise measure of the confidence associated with
 222 label predictions. Small neighborhoods are chosen when one label is measured as significant nearly
 223 as soon as statistically possible; by definition of the AKNN stopping rule, this is not true where large
 224 neighborhoods are necessary. In our experiments, performance on points with significantly higher
 225 neighborhood sizes dropped monotonically, with the majority of the dataset having performance
 226 significantly exceeding the best k -NN rule over a range of settings of A (Fig. 4, 6; Appendix C).

227 References

- 228 [AT07] J.-Y. Audibert and A.B. Tsybakov. Fast learning rates for plug-in classifiers. *Annals of*
 229 *Statistics*, 35(2):608–633, 2007.
- 230 [BBL05] Stéphane Boucheron, Olivier Bousquet, and Gábor Lugosi. Theory of classification: A
 231 survey of some recent advances. *ESAIM: probability and statistics*, 9:323–375, 2005.
- 232 [C⁺18] Tabula Muris Consortium et al. Single-cell transcriptomics of 20 mouse organs creates a
 233 tabula muris. *Nature*, 562(7727):367, 2018.
- 234 [CD10] K. Chaudhuri and S. Dasgupta. Rates of convergence for the cluster tree. In *Advances in*
 235 *Neural Information Processing Systems*, pages 343–351, 2010.
- 236 [CD14] K. Chaudhuri and S. Dasgupta. Rates of convergence for nearest neighbor classification.
 237 In *Advances in Neural Information Processing Systems*, pages 3437–3445. 2014.
- 238 [CG06] F. Cerou and A. Guyader. Nearest neighbor classification in infinite dimension. *ESAIM:*
 239 *Probability and Statistics*, 10:340–355, 2006.
- 240 [CH67] T. Cover and P.E. Hart. Nearest neighbor pattern classification. *IEEE Transactions on*
 241 *Information Theory*, 13:21–27, 1967.

- 242 [DCL11] Wei Dong, Moses Charikar, and Kai Li. Efficient k-nearest neighbor graph construction
243 for generic similarity measures. In *Proceedings of the 20th international conference on*
244 *World wide web*, pages 577–586. ACM, 2011.
- 245 [DGKL94] L. Devroye, L. Györfi, A. Krzyzak, and G. Lugosi. On the strong universal consistency
246 of nearest neighbor regression function estimates. *Annals of Statistics*, 22:1371–1385,
247 1994.
- 248 [DGL96] L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*.
249 Springer, 1996.
- 250 [Dud79] R.M. Dudley. Balls in \mathbb{R}^k do not cut all subsets of $k + 2$ points. *Advances in Mathematics*,
251 31(3):306–308, 1979.
- 252 [FH51] E. Fix and J. Hodges. Discriminatory analysis, nonparametric discrimination. *USAF*
253 *School of Aviation Medicine, Randolph Field, Texas, Project 21-49-004, Report 4, Con-*
254 *tract AD41(128)-31*, 1951.
- 255 [Gyö81] L. Györfi. The rate of convergence of k_n -nn regression estimates and classification rules.
256 *IEEE Transactions on Information Theory*, 27(3):362–364, 1981.
- 257 [Hei01] J. Heinonen. *Lectures on Analysis on Metric Spaces*. Springer, 2001.
- 258 [KP95] S. Kulkarni and S. Posner. Rates of convergence of nearest neighbor estimation under
259 arbitrary sampling. *IEEE Transactions on Information Theory*, 41(4):1028–1039, 1995.
- 260 [MNI]
- 261 [Mou18] Mouse cell atlas dataset. [ftp://ngs.sanger.ac.uk/production/teichmann/](ftp://ngs.sanger.ac.uk/production/teichmann/BKNN/MouseAtlas.zip)
262 [BKNN/MouseAtlas.zip](ftp://ngs.sanger.ac.uk/production/teichmann/BKNN/MouseAtlas.zip), 2018. Accessed: 2019-05-02.
- 263 [MT99] E. Mammen and A.B. Tsybakov. Smooth discrimination analysis. *The Annals of*
264 *Statistics*, 27(6):1808–1829, 1999.
- 265 [not11] notmnist dataset. <http://yaroslavvb.com/upload/notMNIST/>, 2011. Accessed:
266 2019-05-02.
- 267 [RS98] Martin Raab and Angelika Steger. "balls into bins" - A simple and tight analysis. In *Ran-*
268 *domization and Approximation Techniques in Computer Science, Second International*
269 *Workshop, RANDOM'98, Barcelona, Spain, October 8-10, 1998, Proceedings*, pages
270 159–170, 1998.
- 271 [Sto77] C. Stone. Consistent nonparametric regression. *Annals of Statistics*, 5:595–645, 1977.
- 272 [VC71] Vladimir N Vapnik and A Ya Chervonenkis. On the uniform convergence of relative
273 frequencies of events to their probabilities. *Theory of Probability & Its Applications*,
274 16(2):264–280, 1971.

275 **A Analysis and proofs**

276 The first step in establishing advantage-dependent rates of convergence is to bound the accuracy
 277 of empirical estimates of probability mass and bias. This is achieved by a careful choice of large
 278 deviation bounds.

279 **A.1 Large deviation bounds**

280 Suppose we draw n points $(x_1, y_1), \dots, (x_n, y_n)$ from P . If n is reasonably large, we would expect
 281 the empirical mass $\mu_n(S)$ of any set $S \subset \mathcal{X}$, as defined in (4), to be close to its probability mass
 282 under μ . The following lemma, from [CD10], quantifies one particular aspect of this.

Lemma 5 ([CD10], Lemma 7). *There is a universal constant c_o such that the following holds. Let \mathcal{B}
 be any class of measurable subsets of \mathcal{X} of VC dimension d_0 . Pick any $0 < \delta < 1$. Then with
 probability at least $1 - \delta^2/2$ over the choice of $(x_1, y_1), \dots, (x_n, y_n)$, for all $B \in \mathcal{B}$ and for any
 integer k , we have*

$$\mu(B) \geq \frac{k}{n} + \frac{c_o}{n} \max\left(k, d_0 \log \frac{n}{\delta}\right) \implies \mu_n(B) \geq \frac{k}{n}.$$

283

284 Likewise, we would expect the empirical bias $\eta_n(S)$ of a set $S \subset \mathcal{X}$, as defined in (5), to be close to
 285 its true bias $\eta(S)$. The latter is defined whenever $\mu(S) > 0$.

Lemma 6. *There is a universal constant c_1 for which the following holds. Let \mathcal{C} be a class of subsets
 of \mathcal{X} with VC dimension d_0 . Pick any $0 < \delta < 1$. Then with probability at least $1 - \delta^2/2$ over the
 choice of $(x_1, y_1), \dots, (x_n, y_n)$, for all $C \in \mathcal{C}$,*

$$|\eta_n(C) - \eta(C)| \leq \Delta(n, \#_n(C), \delta)$$

286 where $\#_n(C) = |\{i : x_i \in C\}|$ is the number of points in C and

$$\Delta(n, k, \delta) = c_1 \sqrt{\frac{d_0 \log n + \log(1/\delta)}{k}}. \quad (11)$$

287

288 Lemma 6 is a special case¹ of a uniform convergence bound for conditional probabilities (Theorem 8)
 289 that we present and prove in Appendix B.

290 **A.2 Proof of Theorem 1**

291 **Theorem** (Theorem 1 restatement). *There is an absolute constant $C > 0$ for which the following
 292 holds. Let $0 < \delta < 1$ denote the confidence parameter in the AKNN algorithm (Figure 2), and
 293 suppose the algorithm is used to define a classifier g_n based on n training points chosen i.i.d. from P .
 294 Then, for every point $x \in \text{supp}(\mu)$, if*

$$n \geq \frac{C}{\text{adv}(x)} \max\left(\log \frac{1}{\text{adv}(x)}, \log \frac{1}{\delta}\right)$$

295 then with probability at least $1 - \delta$ we have that $g_n(x) = g^*(x)$.

296 *Proof.* Define $c_2 = \max(c_1, 1/2)\sqrt{1 + c_o}$, where c_o and c_1 are the constants from Lemmas 5 and 6,
 297 and take $c_3 = 16c_2^2$.

298 Suppose $\eta(x) > 0$; the negative case is symmetric. The set \mathcal{B} of all balls centered at x is easily seen
 299 to have VC dimension $d_0 = 1$. By Lemmas 5 and 6, we have that with probability at least $1 - \delta^2$, the
 300 following two properties hold for all $B \in \mathcal{B}$:

- 301 1. For any integer k , we have $\#_n(B) \geq k$ whenever $n\mu(B) \geq k + c_o \max(k, \log(n/\delta))$.
- 302 2. $|\eta_n(B) - \eta(B)| \leq \Delta(n, \#_n(B), \delta)$.

¹Indeed, Lemma 6 follows from Theorem 8 by plugging in it $\mathcal{A} = \{\mathcal{X} \times \{+1\}\}$, $\mathcal{B} = \{C \times \{\pm 1\} : C \in \mathcal{C}\}$.

303 Assume henceforth that these hold.

304 By the definition of advantage, point x is (p, γ) -salient for some $p, \gamma > 0$ with $\text{adv}(x) = p\gamma^2$. The
 305 lower bound on n in the theorem statement implies that

$$\gamma \geq 2c_2 \sqrt{\frac{\log n + \log(1/\delta)}{np}}, \quad (12)$$

306 or equivalently that $n \cdot \text{adv}(x) \geq 4c_2^2(\log n + \log(1/\delta))$.

307 Set $k = np/(1 + w)$. By (12) we have $np \geq 4c_2^2 \log(n/\delta)$ and thus $k \geq \log(n/\delta)$. As a result,
 308 $np \geq k + w \max(k, \log(n/\delta))$, and by property 1, the ball $B = B(x, r_p(x))$ has $\#_n(B) \geq k$. This
 309 means, in turn, that by property 2,

$$\begin{aligned} \eta_n(B) &\geq \eta(B) - \Delta(n, k, \delta) = \gamma - c_1 \sqrt{\frac{\log(n/\delta)}{k}} \\ &\geq 2c_2 \sqrt{\frac{\log(n/\delta)}{np}} - c_1 \sqrt{\frac{\log(n/\delta)}{k}} \geq 2c_1 \sqrt{\frac{\log(n/\delta)}{k}} - c_1 \sqrt{\frac{\log(n/\delta)}{k}} \\ &= c_1 \sqrt{\frac{\log(n/\delta)}{k}} \geq \Delta(n, \#_n(B), \delta). \end{aligned}$$

310 Thus ball B would trigger a prediction of $+1$.

At the same time, for any ball $B' = B(x, r)$ with $r < r_p(x)$,

$$\eta_n(B') \geq \eta(B') - \Delta(n, \#_n(B'), \delta) > -\Delta(n, \#_n(B'), \delta)$$

311 and thus no such ball will trigger a prediction of -1 . Therefore, the prediction at x must be $+1$. \square

312 A.3 Proof of Theorem 2

313 This proof follows much the same outline as that of Theorem 1. A crucial difference is that uniform
 314 large deviation bounds (Lemmas 5 and 6) are applied to the class of all balls in \mathcal{X} , which is assumed²
 315 to have finite VC dimension d_0 . In contrast, the proof of Theorem 1 only applies these bounds to the
 316 class of balls centered at a specific point, which has VC dimension at most 1 in any metric space.

317 A.4 Proof of Theorem 4

318 Recall from (9) that \mathcal{X}_a denotes the set of points with advantage $> a$.

Lemma 7. *Pick any $0 < \delta < 1$ as a confidence parameter for the AKNN estimator of Figure 2. Fix any $a > 0$. If the number of training points n satisfies*

$$n \geq \frac{c_3}{a} \max\left(\log \frac{c_3}{a}, \log \frac{1}{\delta}\right),$$

then with probability at least $1 - \delta$, the resulting classifier g_n has risk

$$R(g_n) - R^* \leq \delta + \mu(\mathcal{X}_0 \setminus \mathcal{X}_a).$$

319

Proof. From Theorem 1, we have that for any $x \in \mathcal{X}_a$,

$$\Pr_n(g_n(x) \neq g^*(x)) \leq \delta^2,$$

where \Pr_n denotes probability over the choice of training points. Thus, for $X \sim \mu$,

$$\mathbb{E}_n \mathbb{E}_X \mathbf{1}(g_n(X) \neq g^*(X) | X \in \mathcal{X}_a) \leq \delta^2,$$

and by Markov's inequality,

$$\Pr_n[\Pr_X(g_n(X) \neq g^*(X) | X \in \mathcal{X}_a) \geq \delta] \leq \delta.$$

²This is motivated by finite-dimensional Euclidean space \mathbb{R}^D , where it holds with $d_0 = D + 1$ ([Dud79]).

Thus, with probability at least $1 - \delta$ over the training set,

$$\Pr_X(g_n(X) \neq g^*(X) | X \in \mathcal{X}_a) \leq \delta.$$

320 On points with $\eta(x) = 0$, both g_n and the Bayes-optimal g^* incur the same risk. Thus

$$\begin{aligned} R(g_n) - R^* &\leq \Pr_X(g_n(X) \neq g^*(X) | X \in \mathcal{X}_a) + \Pr_X(X \notin \mathcal{X}_a, \eta(X) \neq 0) \\ &\leq \delta + \Pr_X(X \in \mathcal{X}_0 \setminus \mathcal{X}_a) + \Pr_X(\text{adv}(X) = 0, \eta(X) \neq 0) \\ &\leq \delta + \mu(\mathcal{X}_0 \setminus \mathcal{X}_a), \end{aligned}$$

321 where we invoke Lemma 3 for the last step. \square

We now complete the proof of Theorem 4. Given the sequence of confidence parameters (δ_n) , define a sequence of advantage values (a_n) by

$$a_n = \frac{c_3}{n} \max \left(2 \log n, \log \frac{1}{\delta_n} \right).$$

322 The conditions on (δ_n) imply $a_n \rightarrow 0$.

Pick any $\epsilon > 0$. By the conditions on (δ_n) , we can pick N so that $\sum_{n \geq N} \delta_n \leq \epsilon$. Let ω denote a realization of an infinite training sequence $(X_1, Y_1), (X_2, Y_2), \dots$ from P . By Lemma 7, for any positive integer N ,

$$\Pr(\omega : \exists n \geq N \text{ s.t. } R(g_n(\omega)) - R^* > \delta_n + \mu(\mathcal{X}_0 \setminus \mathcal{X}_{a_n})) \leq \sum_{n \geq N} \delta_n \leq \epsilon.$$

Thus, with probability at least $1 - \epsilon$ over the training sequence ω , we have that for all $n \geq N$,

$$R(g_n(\omega)) - R^* \leq \delta_n + \mu(\mathcal{X}_0 \setminus \mathcal{X}_{a_n}),$$

323 whereupon $R(g_n(\omega)) \rightarrow R^*$ (since $\delta_n, a_n \rightarrow 0$ and $\lim_{a \downarrow 0} \mu(\mathcal{X}_0 \setminus \mathcal{X}_a) = 0$). Since this holds for
324 any $\epsilon > 0$, the theorem follows.

325 **B Uniform Convergence of Empirical Conditional Measures**

326 **B.1 Formal Statement**

327 Let P be a distribution over X , and let \mathcal{A}, \mathcal{B} be two collections of events. Consider n independent
328 samples from P , denoted by x_1, \dots, x_n . We would like to estimate $P(A|B)$ simultaneously for
329 all $A \in \mathcal{A}, B \in \mathcal{B}$. It is natural to consider the empirical estimates:

$$P_n(A|B) = \frac{\sum_i 1_{[x_i \in A \cap B]}}{\sum_i 1_{[x_i \in B]}}.$$

330 We study when (and to what extent) these estimates provide a good approximation. Note that the
331 case where $\mathcal{B} = \{X\}$ (i.e., in which one estimates $P(A)$ using $P_n(A)$ simultaneously for all $A \in \mathcal{A}$)
332 is handled by the classical VC theory. Throughout this section we assume that both \mathcal{A}, \mathcal{B} have a finite
333 VC-dimension, and we let d_0 denote an upper bound on both $\text{VC}(\mathcal{A})$ and $\text{VC}(\mathcal{B})$.

334 To demonstrate the kinds of statements we would like to derive, consider the case where each of \mathcal{A}, \mathcal{B}
335 contains only one event: $\mathcal{A} = \{A\}$, and $\mathcal{B} = \{B\}$, and set $\#_n(B) = \sum_i 1_{[x_i \in B]}$. A Chernoff
336 bound implies that conditioned on the event that $\#_n(B) > 0$, the following holds with probability at
337 least $1 - \delta$:

$$|P(A|B) - P_n(A|B)| \leq \sqrt{\frac{2 \log(1/\delta)}{\#_n(B)}}. \quad (13)$$

338 To derive it, use that conditioned on $x_i \in B$, the event $x_i \in A$ has probability $P(A|B)$, and therefore
339 the random variable “ $\#_n(B) \cdot p_n(A|B)$ ” has a binomial distribution with parameters $\#_n(B)$ and
340 $P(A|B)$.

341 Note that the bound on the error in Equation (13) depends on $\#_n(B)$ and therefore is data-dependent.
342 We stress that this is the type of statement we want: the more samples belong to B , the more certain
343 we are with the empirical estimate. Thus, we would want to prove a statement as follows:

344 With probability at least $1 - \delta$,

$$(\forall A \in \mathcal{A}) (\forall B \in \mathcal{B}) : |P(A|B) - P_n(A|B)| \leq O\left(\sqrt{\frac{d_0 \log(1/\delta)}{\#_n(B)}}\right),$$

345 where $\#_n(B) = \sum_{i=1}^n 1[x_i \in B]$.

346 The above statement is, unfortunately, false. As an example, consider the probability space defined
 347 by drawing $x \sim [n]$ uniformly, and then coloring x by $c_x \in \{\pm 1\}$ uniformly. For each i let B_i
 348 denote the event that i was drawn, and let A denote the event that the drawn color was $+1$. (formally,
 349 $B_i = \{i\} \times \{\pm 1\}$, and $A = [n] \times \{+1\}$). One can verify that the VC dimension of $\mathcal{B} = \{B_i : i \leq n\}$
 350 and of $\mathcal{A} = \{A\}$ is at most 1. The above statement fails in this setting: indeed, one can verify that if
 351 we draw n samples from this space then with a constant probability there will be some j such that:

- 352 (i) j always gets the same color (say $+1$), and
- 353 (ii) j is sampled at least $\Omega(\log n / \log \log n)$ times³.

354 Therefore, with constant probability we get that

$$P_n(A|B_i) = 1, P(A|B_i) = 1/2,$$

355 and so the difference between the error is clearly $1 - (1/2) = 1/2$, which is clearly not upper bounded
 356 by $O(\sqrt{\log \log n / \log n})$.

357 We prove the following (slightly weaker) variant:

358 **Theorem 8** (UCECM). *Let P be a probability distribution over X , and let \mathcal{A}, \mathcal{B} be two families*
 359 *of measurable subsets of X such that $\text{VC}(\mathcal{A}), \text{VC}(\mathcal{B}) \leq d_0$. Let $n \in \mathbb{N}$, and let $x_1 \dots x_n$ be n i.i.d*
 360 *samples from P . The, the following event occurs with probability at least $1 - \delta$:*

$$(\forall A \in \mathcal{A}) (\forall B \in \mathcal{B}) : |P(A|B) - P_n(A|B)| \leq \sqrt{\frac{k_o}{\#_n(B)}},$$

361 where $k_o = 1000(d_0 \log(8n) + \log(4/\delta))$, and⁴ $\#_n(B) = \sum_{i=1}^n 1[x_i \in B]$.

362 **Discussion.** Theorem 8 can be combined with Lemma 5 to yield a bound on the minimal n for
 363 which $P_n(A|B)$ is a non-trivial approximation of $P(A|B)$. Indeed, Lemma 5 implies that if n is large
 364 enough so that $P(B) = \Omega\left(\frac{d_0 \log n}{n}\right)$, then the empirical estimate $P_n(A|B)$ is a decent approximation.
 365 In the context of the adaptive nearest neighbor classifier, this means that the empirical biases provide
 366 meaningful estimates of the true biases for balls whose measure is $\tilde{\Omega}\left(\frac{d_0}{n}\right)$. This resembles the
 367 learning rate in realizable settings.

368 We remark that a weaker statement than Theorem 8 can be derived as a corollary of the classical
 369 uniform convergence result [VC71]. Indeed, since the VC dimension of $\{B \cap A : i \in \mathcal{I}\}$ is at most
 370 d_0 , it follows that

$$P_n(A|B) \approx \frac{P(A \cap B) \pm \sqrt{d_0/n}}{P(B) \pm \sqrt{d_0/n}}.$$

371 However, this bound guarantees non-trivial estimates only once $P(B)$ is roughly $\sqrt{d_0/n}$. This is
 372 similar to the learning rate in agnostic (i.e., non-realizable) settings.

373 Another major advantage of the uniform convergence bound in Theorem 8 is that it is data-dependent:
 374 if many points from the sample belong to $B \in \mathcal{B}$ (i.e. $\#_n(B)$ is large), then we get better guarantees
 375 on the approximation of $P(A|B)$ by $P_n(A|B)$ for all $A \in \mathcal{A}$.

³This follows from analyzing the maximal bin in a uniform assignment of $\Theta(n)$ balls into n bins [RS98]

⁴Note that the above inequality makes sense also when $k(B) = 0$, by identifying $\frac{\cdot}{0}$ as ∞ , and using the convention that $\infty - \infty = \infty$ and that $\infty \leq \infty$.

376 **B.2 Proof of Theorem 8**

377 As noted above, the standard uniform convergence bound for VC classes can not yield the bound
 378 in Theorem 8. Instead, we use a variant of it due to [BBL05] which concerns *relative deviations*
 379 (see [BBL05]: Theorem 5.1 and the discussion before Corollary 5.2). In order to state the theorem,
 380 we need the following notation: Let \mathcal{C} be a family of subsets of \mathcal{X} . We denote by $\mathbb{S}_{\mathcal{C}} : \mathbb{N} \rightarrow \mathbb{N}$ the
 381 *growth function* of \mathcal{C} , which is defined by:

$$\mathbb{S}_{\mathcal{C}}(n) = \max\{|\mathcal{C}|_R : R \subseteq X, |R| = n\},$$

382 where $\mathcal{C}|_R = \{C \cap R : C \in \mathcal{C}\}$ is the projection of \mathcal{C} to R .

383 **Theorem 9** ([BBL05]). *Let \mathcal{C} be a family of subsets of \mathcal{X} and let P be a distribution over \mathcal{X} . Then,*
 384 *the following holds with probability $1 - \delta$:*

$$(\forall C \in \mathcal{C}) : |P(C) - P_n(C)| \leq 2\sqrt{P_n(C) \frac{\log \mathbb{S}_{\mathcal{C}}(2n) + \log(4/\delta)}{n}} + 4 \frac{\log \mathbb{S}_{\mathcal{C}}(2n) + \log(4/\delta)}{n}.$$

385 Set $\mathcal{C} = \mathcal{B} \cup \{A \cap B : A \in \mathcal{A}, B \in \mathcal{B}\}$. We prove Theorem 8 by applying Theorem 9 on \mathcal{C} ; to this end
 386 we first upper bound $\mathbb{S}_{\mathcal{C}}(n)$. Let $\mathcal{D} = \{A \cap B : A \in \mathcal{A}, B \in \mathcal{B}\}$, so that $\mathcal{C} = \mathcal{B} \cup \mathcal{D}$. Then:

$$\mathbb{S}_{\mathcal{C}}(n) \leq \mathbb{S}_{\mathcal{B}}(n) + \mathbb{S}_{\mathcal{D}}(n) \leq \mathbb{S}_{\mathcal{B}}(n) + \mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{B}}(n) \leq 2\mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{B}}(n) \leq 2 \binom{n}{\leq d_0}^2 \leq 2(2n)^{2d_0},$$

387 where the second inequality follows since $\mathbb{S}_{\mathcal{D}}(n) \leq \mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{B}}(n)$, the second to last inequality
 388 follows from the Sauer-Shelah-Perles Lemma, and the last inequality follows since $\binom{a}{\leq b} \leq (2a)^b$.
 389 Therefore, applying Theorem 9 on \mathcal{C} yields that with probability $1 - \delta$ the following event holds:

$$(\forall C \in \mathcal{C}) : |P(C) - P_n(C)| \leq 4\sqrt{P_n(C) \frac{d_0 \log 8n + \log(4/\delta)}{n}} + 8 \frac{d_0 \log 8n + \log(4/\delta)}{n}. \quad (14)$$

390 For the remainder of the proof we assume that the event in Equation (14) holds and argue that it
 391 implies the conclusion in Theorem 8. Let $A \in \mathcal{A}, B \in \mathcal{B}$, and let $k = n \cdot P_n(B) = \#_n(B)$ denote
 392 the number of data points in B . We want to show that

$$|P(A|B) - P_n(A|B)| \leq \sqrt{\frac{k_o}{k}}, \quad (15)$$

393 where $k_o = 1000(d_0 \log(8n) + \log(4/\delta))$. Let $j = k \cdot P_n(A|B) = \#_n(A \cap B)$ denote the number
 394 of data points in $A \cap B$. We establish Equation (15) by showing that

$$P(A|B) \leq P_n(A|B) + \sqrt{\frac{k_o}{k}} \quad \text{and} \quad P(A|B) \geq P_n(A|B) - \sqrt{\frac{k_o}{k}}.$$

395 In the following calculation it will be convenient to denote $D := d_0 \log(8n) + \log(4/\delta)$. By
 396 Equation (14) we get:

$$\begin{aligned} P(A|B) &= \frac{P(A \cap B)}{P(B)} \\ &\leq \frac{P_n(A \cap B) + 4\sqrt{P_n(A \cap B) \frac{D}{n}} + 8 \frac{D}{n}}{P_n(B) - 4\sqrt{P_n(B) \frac{D}{n}} - 8 \frac{D}{n}} \\ &= \frac{\frac{P_n(A \cap B)}{P_n(B)} + 4\sqrt{\frac{P_n(A \cap B)}{P_n(B)} \frac{D}{nP_n(B)}} + 8 \frac{D}{nP_n(B)}}{1 - 4\sqrt{\frac{D}{nP_n(B)}} - 8 \frac{D}{nP_n(B)}} s = P_n(A|B) \frac{1 + 4\sqrt{\frac{D}{j}} + 8 \frac{D}{j}}{1 - 4\sqrt{\frac{D}{k}} - 8 \frac{D}{k}}, \end{aligned}$$

397 where the first inequality follows from Equation (14) and the following equalities are trivial. Thus,

$$P(A|B) \leq \frac{j}{k} \left(\frac{1 + 4\sqrt{\frac{D}{j}} + 8 \frac{D}{j}}{1 - 4\sqrt{\frac{D}{k}} - 8 \frac{D}{k}} \right). \quad (16)$$

398 Next, note that we may assume that $k \geq k_o = 1000D$, as otherwise Equation (15) trivially holds.
 399 Therefore,

$$\frac{1}{1 - 4\sqrt{\frac{D}{k}} - 8\frac{D}{k}} \leq 1 + 8\sqrt{\frac{D}{k}} + 16\frac{D}{k}. \quad ((\forall x < \frac{1}{2}) : \frac{1}{1-x} \leq 1 + 2x)$$

400 Plugging this in Equation (16), and using first that $j \leq k$ and then that $1000D \leq k$, yields:

$$\begin{aligned} P(A|B) &\leq \frac{j}{k} \left(1 + 4\sqrt{\frac{D}{j}} + 8\frac{D}{j}\right) \left(1 + 8\sqrt{\frac{D}{k}} + 16\frac{D}{k}\right) \\ &= \frac{j}{k} + 8\frac{j}{k}\sqrt{\frac{D}{k}} \left(1 + 2\sqrt{\frac{D}{k}}\right) + \left(\frac{4\sqrt{jD} + 8D}{k}\right) \left(1 + 4\sqrt{\frac{D}{k}}\right)^2 \\ &\leq \frac{j}{k} + 8\sqrt{\frac{D}{k}} \left(1 + 2\sqrt{\frac{D}{k}}\right) + \left(4\sqrt{\frac{D}{k}} + \frac{8D}{k}\right) \left(1 + 4\sqrt{\frac{D}{k}}\right)^2 \\ &\leq \frac{j}{k} + 30\sqrt{\frac{D}{k}} = P_n(A|B) + \sqrt{\frac{k_o}{k}}, \end{aligned}$$

401 and so

$$P(A|B) \leq P_n(A|B) + \sqrt{\frac{k_o}{k}}.$$

402 A symmetric argument yields similarly to Equation (16) that:

$$P(A|B) \geq \frac{j}{k} \left(\frac{1 - 4\sqrt{\frac{D}{j}} - 8\frac{D}{j}}{1 + 4\sqrt{\frac{D}{k}} + 8\frac{D}{k}}\right).$$

403 Then, a similar calculation (using the relation $(\forall x > 0) : \frac{1}{1+x} \geq 1 - 2x$) implies that

$$P(A|B) \geq P_n(A|B) - \sqrt{\frac{k_o}{k}},$$

404 which finishes the proof. □

405 C Experimental Results

406 C.1 Datasets

407 We test AKNN on the notMNIST dataset ([not11]), consisting of extracted glyphs of the letters
 408 A-J from publicly available fonts. We use the 18724 labeled examples from this set, preprocessed
 409 feature-wise to be in $[-\frac{1}{2}, \frac{1}{2}]$ using $x \mapsto \frac{x}{255} - \frac{1}{2}$.

410 We also test on the MNIST dataset ([MNI]).

411 We use AKNN on a challenging binary classification task of independent and continuing interest,
 412 involving gene expression data on a population single cells from different mouse organs collected
 413 by the Tabula Muris consortium ([C⁺18], as processed in [Mou18]). This constitutes 45291 cells
 414 (training examples). Each cell has its data collected using one of two approaches. The task is to
 415 classify between them.

416 The data are collected using representative protocols of the two currently dominant approaches
 417 to isolate and measure single cells: a “plate”-based approach using microwells on a chip, and a
 418 “droplet”-based approach manipulating cells within droplets using microfluidic technologies. Each
 419 approach has its own set of technical biases, about which much remains to be understood. Identifying
 420 and characterizing these biases to discriminate between such approaches is currently of great interest.

421 Both approaches measure effectively the same cells for our purposes, so there is a large decision
 422 boundary in the binary classification problem.

423 **C.2 A note on efficient implementation**

424 In this paper, we computed the nearest neighbors of data exactly when running AKNN, to faithfully
 425 demonstrate its behavior. In practice, this would be done using approximate nearest-neighbor search
 426 to build a k -NN graph using a small fixed k (say 10), and then using pairwise distances on this
 427 graph to compute neighborhoods as needed by AKNN. We tried this (using the nearest-neighbor
 428 method of [DCL11]) on notMNIST without substantive differences in the results, and will release
 429 this implementation upon publication.

430 **C.3 Supplemental Figures**

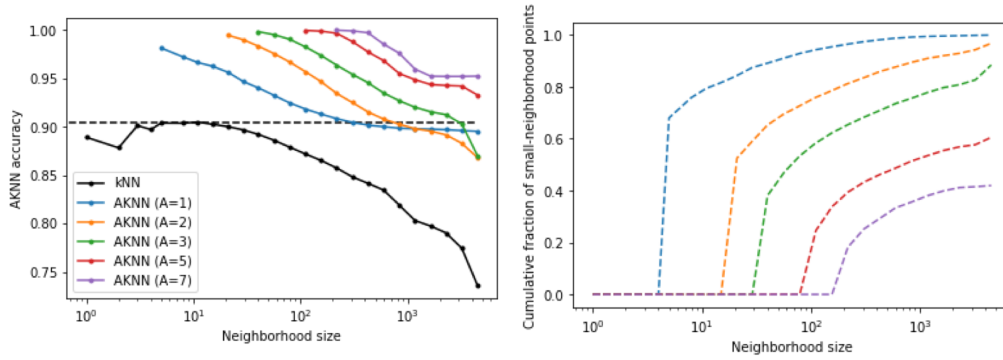


Figure 6: As Fig. 4, on single-cell mouse data.

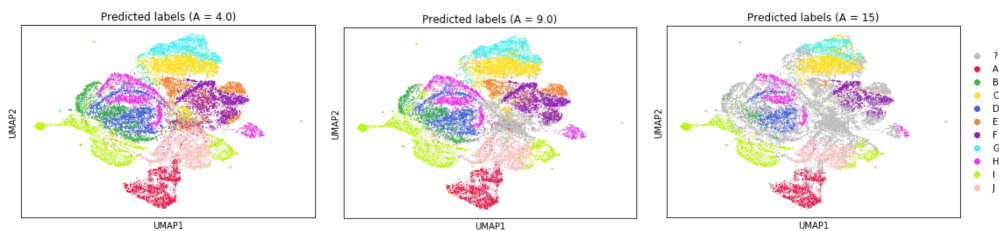


Figure 7: AKNN predictions on notMNIST, for different settings of A .

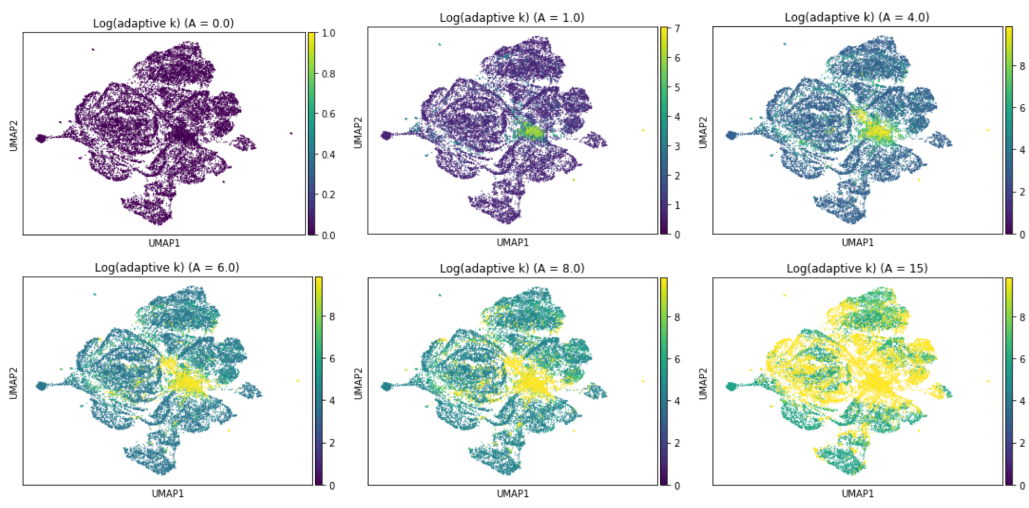


Figure 8: AKNN neighborhood sizes on notMNIST, in increasing order of A , plotted on a log scale. Top left figure ($A = 0$) represents a 1-NN classifier. Bottom right figure ($A = 15$) shows that many of the points' neighborhoods are maximally large, which can be compared to the right panel of Fig. 7.

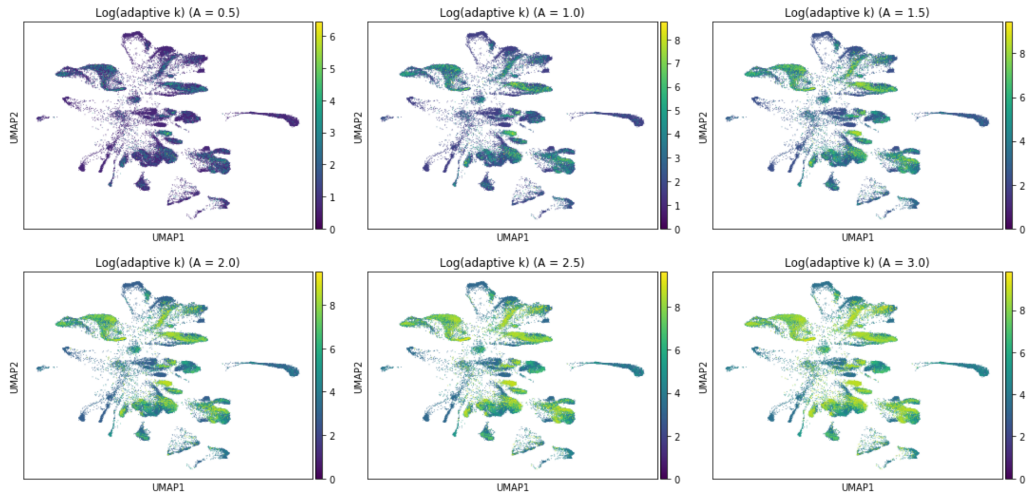


Figure 9: As Fig. 8, on single-cell mouse data.