

A first model of learning

Let's restrict our attention to binary classification

- our labels belong to $\mathcal{Y} = \{1, 0\}$ (or $\mathcal{Y} = \{+1, -1\}$)

We observe the data

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

where each $\mathbf{x}_i \in \mathbb{R}^d$

Suppose we are given an ensemble of possible hypotheses / classifiers

$$\mathcal{H} = \{h_1, \dots, h_m\}$$

From the **training data** \mathcal{D} , we would like to select the best possible classifier from \mathcal{H}

Empirical risk minimization (ERM)

Recall the definitions of risk/empirical risk

$$R(h_j) := \mathbb{P}[h_j(X) \neq Y]$$

$$\widehat{R}_n(h_j) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{i: h_j(\mathbf{x}_i) \neq y_i\}}(i)$$

Ideally, we would like to choose

$$h^\# = \arg \min_{h_j \in \mathcal{H}} R(h_j)$$

Since $\widehat{R}_n(h_j)$ is supposed to be a good estimate of $R(h_j)$, an incredibly natural (and common) strategy is to pick

$$h^* = \arg \min_{h_j \in \mathcal{H}} \widehat{R}_n(h_j)$$

The excess risk

Note that by definition, we must have $R(h^\sharp) \leq R(h^*)$

We would like to guarantee that $R(h^*) - R(h^\sharp)$ is small

Last time we took an extended detour through the exciting world of concentration inequalities to show (using Hoeffding) that for any fixed h_j

$$\mathbb{P} \left[|\widehat{R}_n(h_j) - R(h_j)| \geq \epsilon \right] \leq 2e^{-2\epsilon^2 n}$$

or equivalently, that with probability at least $1 - \delta$

$$|\widehat{R}_n(h_j) - R(h_j)| \leq \sqrt{\frac{1}{2n} \log(2/\delta)}$$

Getting a uniform guarantee

We would like to get a bound on $|\widehat{R}_n(h^*) - R(h^*)|$

Since the choice of h^* depends on the empirical risk of all the classifiers in \mathcal{H} , one way to do this is to ensure that $|\widehat{R}_n(h_j) - R(h_j)|$ is small for **every** $h_j \in \mathcal{H}$

That is, we want to show that with probability at least $1 - \delta$

$$\max_{h_j \in \mathcal{H}} |\widehat{R}_n(h_j) - R(h_j)| \leq \epsilon$$

for some suitable choice of δ, ϵ

We do this using the **union bound**

Result of the union bound

$$\begin{aligned} \mathbb{P} \left[\max_{h_j \in \mathcal{H}} \left| \widehat{R}_n(h_j) - R(h_j) \right| \geq \epsilon \right] \\ \leq \sum_{j=1}^m \mathbb{P} \left[\left| \widehat{R}_n(h_j) - R(h_j) \right| \geq \epsilon \right] \\ \leq \sum_{j=1}^m 2e^{-2\epsilon^2 n} \\ = 2me^{-2\epsilon^2 n} \end{aligned}$$

Bounding the excess risk

Note that when we have a bound such as

$$\max_{h_j \in \mathcal{H}} |\widehat{R}_n(h_j) - R(h_j)| \leq \epsilon$$

we can not only make guarantees about $|\widehat{R}_n(h^*) - R(h^*)|$, but can also relate the performance of h^* to h^\sharp :

$$\begin{aligned} R(h^*) - R(h^\sharp) &= R(h^*) - \widehat{R}_n(h^*) + \widehat{R}_n(h^*) - R(h^\sharp) \\ &\leq |R(h^*) - \widehat{R}_n(h^*)| + |\widehat{R}_n(h^*) - R(h^\sharp)| \end{aligned}$$

Note that $R(h^\sharp) \leq R(h^*) \leq \widehat{R}_n(h^*) + \epsilon$

$$\widehat{R}_n(h^*) \leq \widehat{R}_n(h^\sharp) \leq R(h^\sharp) + \epsilon$$

Thus $R(h^*) - R(h^\sharp) \leq 2\epsilon$

The upshot

As long as m isn't too big ($m \lesssim e^n$) then we can be reasonably confident that $\widehat{R}_n(h^*) \approx R(h^*)$ and furthermore that $R(h^*)$ isn't too much larger than $R(h^\#)$

Of course, the trick in a doing a good job of learning is ensure that $R(h^*)$ is actually *small*

To achieve this, we need a rich set of possible hypotheses...

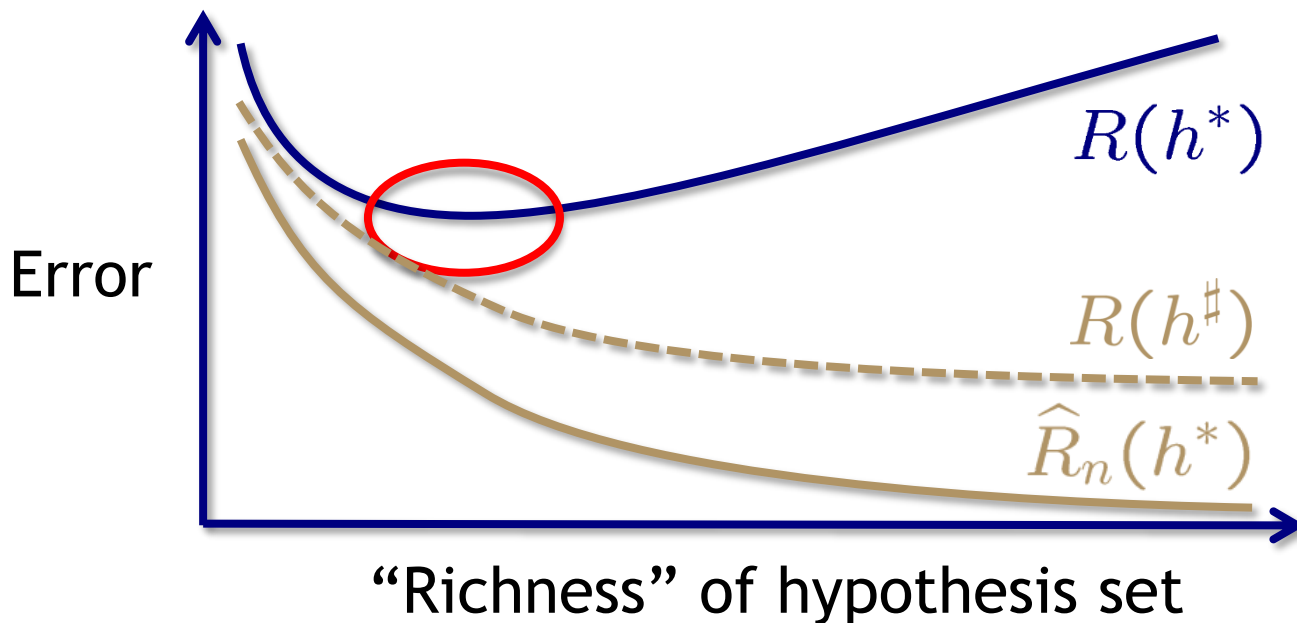
unfortunately...

Fundamental tradeoff

More hypotheses ultimately sacrifices our guarantee that

$$\widehat{R}_n(h^*) \approx R(h^*)$$

Richer set of hypotheses \rightarrow $\left\{ \begin{array}{l} \widehat{R}_n(h^*) \downarrow \quad R(h^\#) \downarrow \\ \widehat{R}_n(h^*) - R(h^*) \uparrow \end{array} \right.$



What is a good hypothesis?

Ideally, we would like to have a small number of hypotheses, so that $\widehat{R}_n(h^*) \approx R(h^*)$, while also being lucky (or smart) enough to have $R(h^*) \approx R(h^\#) \approx 0$

In general, this may not be possible, since there may not be **any** function h with $R(h) \approx 0$

Why not?

Noise: $Y = h(X) + N$

Suppose we knew the joint distribution of our data

- what is the optimal classification rule h^* ?
- what are the fundamental limits on how small $R(h^*)$ can be?

Known distribution case

Consider (X, Y) where

- X is a random vector in \mathbb{R}^d
- $Y \in \{0, \dots, K - 1\}$ is a random variable (depending on X)

Let $h : \mathbb{R}^d \rightarrow \{0, \dots, K - 1\}$ be a **classifier** with **probability of error/risk** given by

$$R(h) := \mathbb{P}[h(X) \neq Y]$$

Denote the **a posteriori class probabilities** by

$$\eta_k(\mathbf{x}) := \mathbb{P}[Y = k | X = \mathbf{x}]$$

for $k = 0, \dots, K - 1$

The Bayes classifier

Theorem

The classifier $h^*(\mathbf{x}) := \arg \max_k \eta_k(\mathbf{x})$ satisfies

$$R^* = R(h^*) \leq R(h)$$

for any possible classifier h

Terminology:

- h^* is called a **Bayes classifier**
- R^* is called the **Bayes risk**

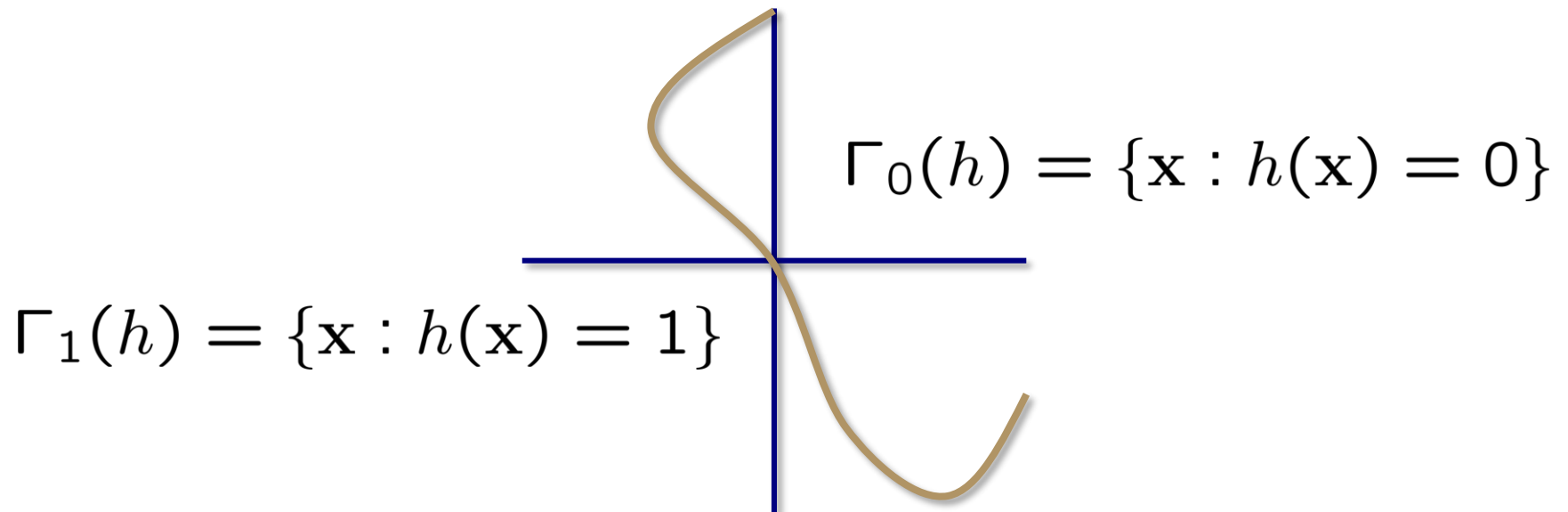
Proof

For convenience, assume $X|Y = k$ is a continuous random variable with density $f_{X|Y}(\mathbf{x}|k)$

Let $\pi_k = \mathbb{P}[Y = k]$ denote the *a priori class probabilities*

Consider an arbitrary classifier h . Denote the decision regions

$$\Gamma_k(h) := \{\mathbf{x} : h(\mathbf{x}) = k\}$$



Proof (Part 2)

We can write $1 - R(h) = \mathbb{P}[h(X) = Y]$

$$= \sum_{k=0}^{K-1} \pi_k \cdot \mathbb{P}[h(X) = k | Y = k]$$
$$= \sum_{k=0}^{K-1} \pi_k \cdot \int_{\Gamma_k(h)} f_{X|Y}(\mathbf{x}|k) d\mathbf{x}$$

We want to **maximize** this expression, we should design our classifier h such that

$$\mathbf{x} \in \Gamma_k(h) \quad \longleftrightarrow \quad \pi_k f_{X|Y}(\mathbf{x}|k) \text{ is maximal}$$

Proof (Part 3)

Therefore, the optimal h has

$$\begin{aligned}h^*(\mathbf{x}) &= \arg \max_k \pi_k f_{X|Y}(\mathbf{x}|k) \\&= \arg \max_k \frac{\pi_k f_{X|Y}(\mathbf{x}|k)}{\sum_{\ell=0}^{K-1} \pi_\ell f_{X|Y}(\mathbf{x}|\ell)} \\&= \arg \max_k \mathbb{P}[Y = k | X = \mathbf{x}]\end{aligned}$$

Bayes rule!

Note that in addition to our rigorous derivation, this classifier also coincides with “common sense”

Variations

Different ways of expressing the Bayes classifier

- $h^*(\mathbf{x}) = \arg \max_k \eta_k(\mathbf{x})$
- $h^*(\mathbf{x}) = \arg \max_k \pi_k f_{X|Y}(\mathbf{x}|k)$
- When $K = 2$

$$\frac{f_{X|Y}(\mathbf{x}|1)}{f_{X|Y}(\mathbf{x}|0)} \underset{1}{\overset{0}{\leq}} \frac{\pi_0}{\pi_1} \quad \text{likelihood ratio test}$$

- When $\pi_0 = \pi_1 = \dots = \pi_{K-1}$

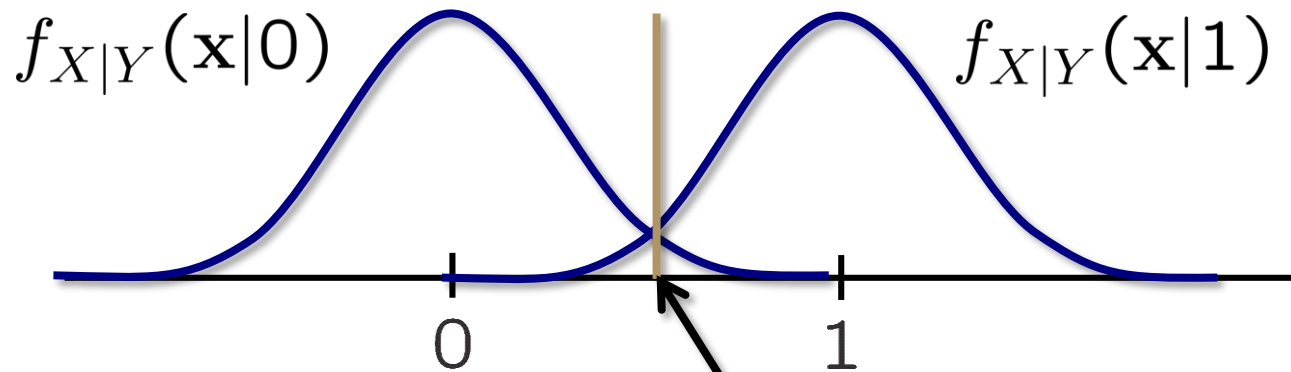
$$h^*(\mathbf{x}) = \arg \max_k f_{X|Y}(\mathbf{x}|k) \quad \text{maximum likelihood classifier/detector}$$

Example

Suppose that $K = 2$ and that

$$X|Y = 0 \sim \mathcal{N}(0, 1)$$

$$X|Y = 1 \sim \mathcal{N}(1, 1)$$



$$\frac{f_{X|Y}(x|1)}{f_{X|Y}(x|0)} \underset{1}{\overset{0}{\gtrless}} \frac{\pi_0}{\pi_1}$$

If $\pi_0 = \pi_1$

Example

How do we calculate the *Bayes risk*?

$$\begin{aligned} R(h^*) &= \mathbb{P} [h^*(X) \neq Y] \\ &= \mathbb{P} [\text{declare } 0 | Y = 1] \cdot \pi_1 \\ &\quad + \mathbb{P} [\text{declare } 1 | Y = 0] \cdot \pi_0 \end{aligned}$$

In the case where $\pi_0 = \pi_1 = \frac{1}{2}$, our test reduced to declaring 1 iff $x \geq \frac{1}{2}$, thus

$$\begin{aligned} R(h^*) &= \frac{1}{2} \mathbb{P} \left[X < \frac{1}{2} | Y = 1 \right] + \frac{1}{2} \mathbb{P} \left[X > \frac{1}{2} | Y = 0 \right] \\ &= \frac{1}{2} \int_{-\infty}^{\frac{1}{2}} f_{X|Y}(x|1) dx + \frac{1}{2} \int_{\frac{1}{2}}^{\infty} f_{X|Y}(x|0) dx \\ &= \Phi \left(-\frac{1}{2} \right) \end{aligned}$$

Alternative cost/loss functions

So far we have focused on minimizing the risk $\mathbb{P}[h(X) \neq Y]$

There are many situations where this is not appropriate

- cost-sensitive classification
 - type I/type II errors or misses/false alarms may have very different costs, in which case a loss function of the form
$$C_0\mathbb{P}[h(X) \neq Y|Y = 0] + C_1\mathbb{P}[h(X) \neq Y|Y = 1]$$
 - alternatively, it may be better to focus on them directly a la Neyman-Pearson classification
- unbalanced datasets
 - when one class dominates the other, the probability of error will place less emphasis on the smaller class
 - the class proportions in our dataset may not be representative of the “wild”

What about learning?

We have just seen that when we know the true distribution underlying our dataset, solving the classification problem is straightforward

In practical learning problems, all we have is the data...

One natural approach is to use the data to estimate the distribution, and then just plug this into the formula for the Bayes classifier

Plugin methods

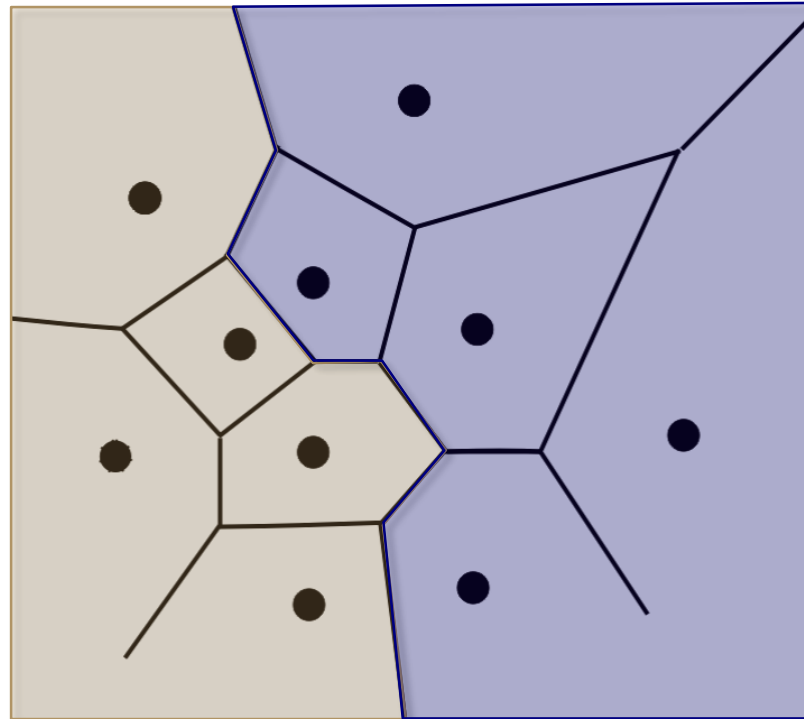
Before we get to these, we will first talk about what is quite possibly the absolute simplest learning algorithm there is...

Nearest neighbor classifier

The *nearest neighbor classifier* is easiest to state in words:

Assign \mathbf{x} the same label as the closest training point \mathbf{x}_i to \mathbf{x}

The nearest neighbor rule defines a *Voronoi partition* of the input space



Risk of the nearest neighbor classifier

Recall that the Bayes classifier simply chooses the k that maximizes $\eta_k(\mathbf{x}) := \mathbb{P}[Y = k|X = \mathbf{x}]$

If $k^* = h^*(\mathbf{x})$, then the risk of the Bayes classifier at \mathbf{x} is

$$R^*(\mathbf{x}) := \mathbb{P}[Y \neq k^*|X = \mathbf{x}] = 1 - \eta_{k^*}(\mathbf{x})$$

If \mathbf{x}_{NN} denotes the nearest neighbor to \mathbf{x} and has label y_{NN} , then the risk of the nearest neighbor classifier at \mathbf{x} is given by

$$\begin{aligned} R_{\text{NN}}(\mathbf{x}) &:= \mathbb{P}[Y \neq y_{\text{NN}}|X = \mathbf{x}] \\ &= \sum_k \mathbb{P}[Y_{\text{NN}} = k|X = \mathbf{x}] \cdot \mathbb{P}[Y \neq k|X = \mathbf{x}] \\ &= \sum_k \eta_k(\mathbf{x}_{\text{NN}})(1 - \eta_k(\mathbf{x})) \end{aligned}$$

Intuition from asymptotics

In the limit as $n \rightarrow \infty$, we can assume that $\mathbf{x}_{\text{NN}} \rightarrow \mathbf{x}$

Thus, as $n \rightarrow \infty$ we have

$$R_{\text{NN}}(\mathbf{x}) \approx \sum_k \eta_k(\mathbf{x})(1 - \eta_k(\mathbf{x}))$$

Consider the binary classification case and suppose (without loss of generality) that $\eta_1(\mathbf{x}) \geq \eta_0(\mathbf{x})$

Bayes risk: $R^*(\mathbf{x}) = 1 - \eta_1(\mathbf{x})$

NN risk: $R_{\text{NN}}(\mathbf{x}) = \eta_0(\mathbf{x})(1 - \eta_0(\mathbf{x})) + \eta_1(\mathbf{x})(1 - \eta_1(\mathbf{x}))$
 $= 2\eta_1(\mathbf{x})(1 - \eta_1(\mathbf{x}))$
 $\leq 2(1 - \eta_1(\mathbf{x}))$

Asymptotically, the risk of the nearest neighbor classifier is *at most twice the Bayes risk*

k -nearest neighbors

We can drive the factor of 2 in this result down to 1 by generalizing the nearest neighbor rule to the **k -nearest neighbor** rule as follows:

Assign a label to \mathbf{x} by taking a majority vote over the k training points \mathbf{x}_i closest to \mathbf{x}

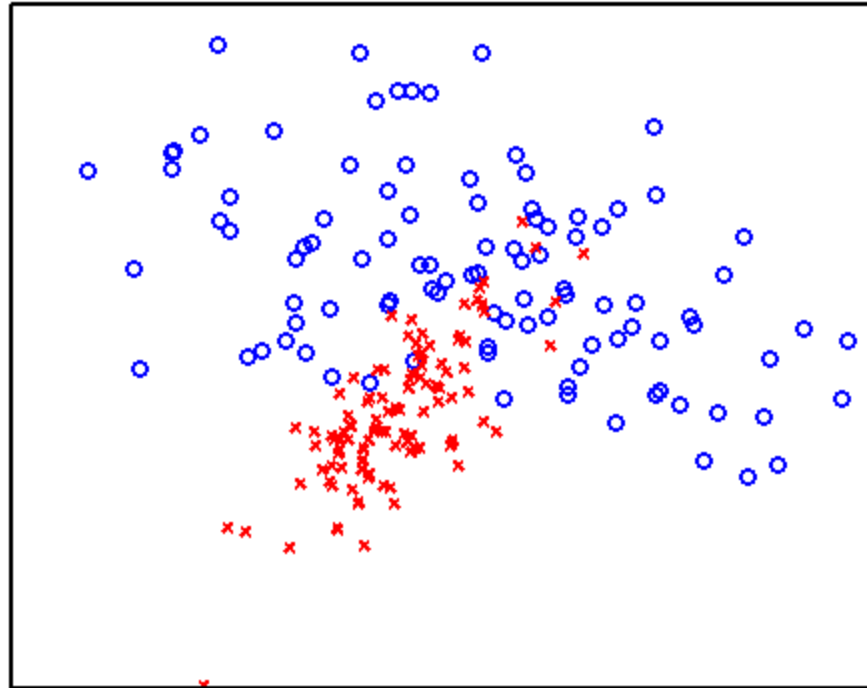
How do we define this more mathematically?

$I_k(\mathbf{x}) :=$ indices of the k training points closest to \mathbf{x}

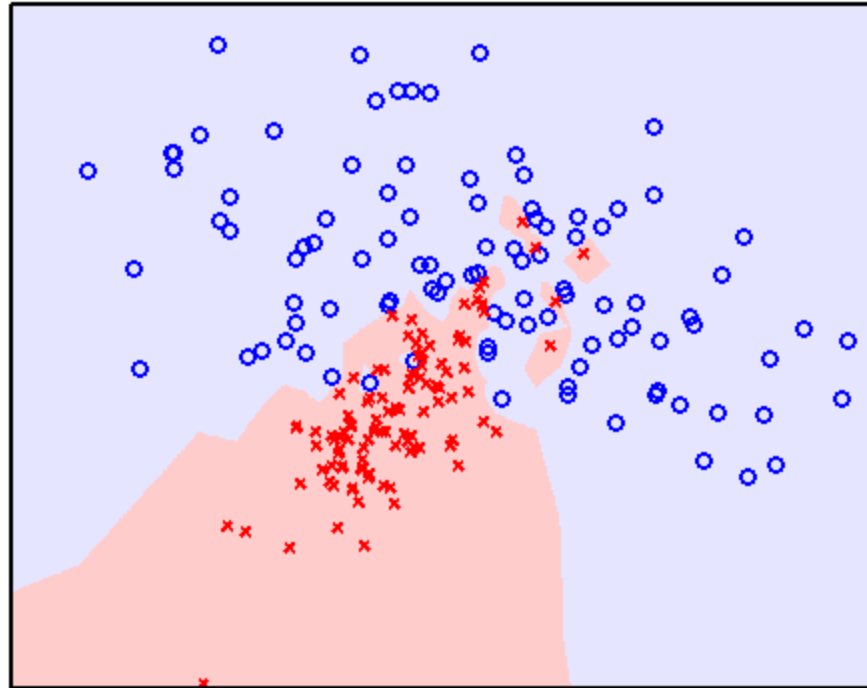
If $y_i = \pm 1$, then we can write the k -nearest neighbor classifier as

$$\hat{h}_k(\mathbf{x}) := \text{sign} \left(\sum_{i \in I_k(\mathbf{x})} y_i \right)$$

Example

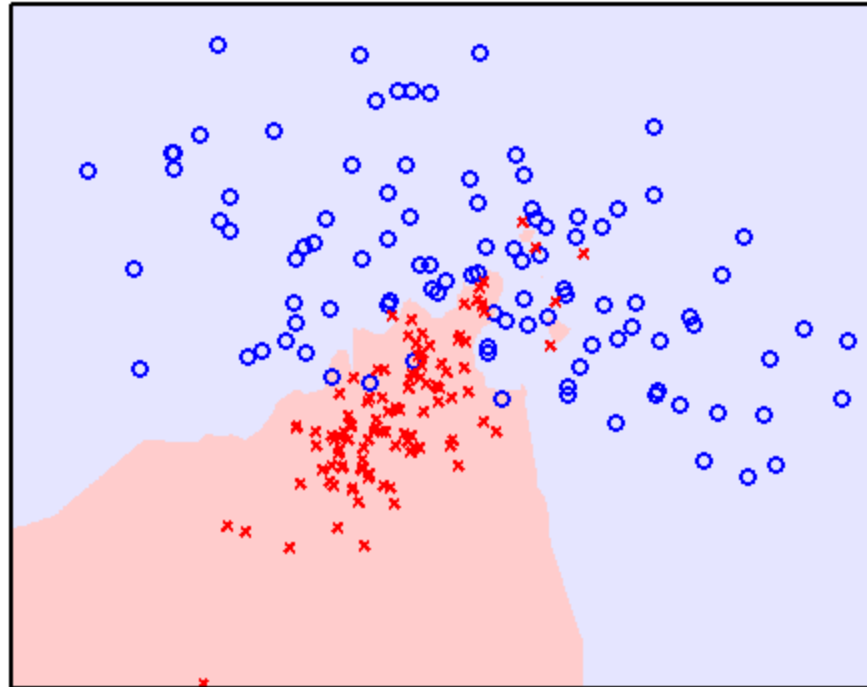


Example



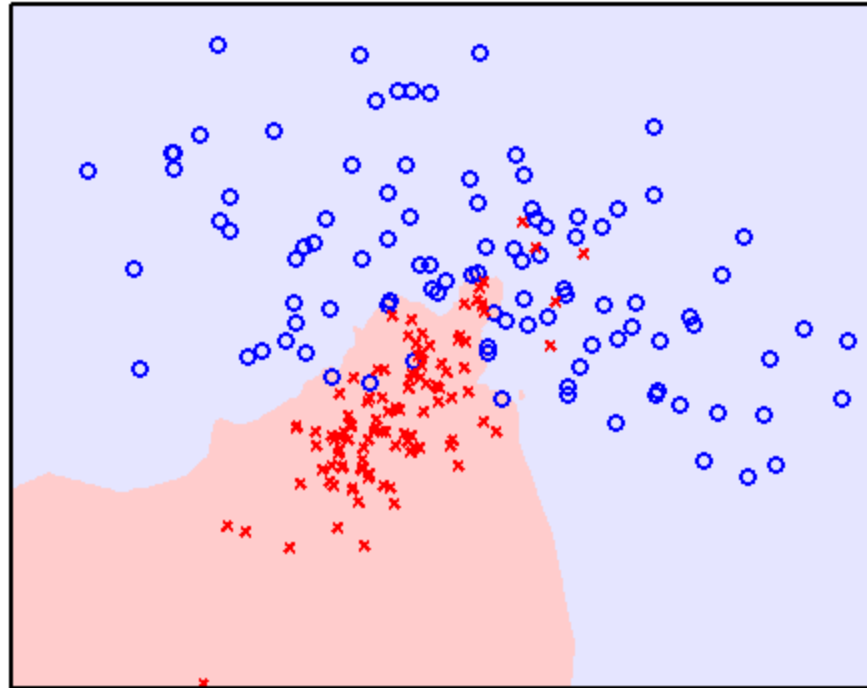
$$k = 1$$

Example



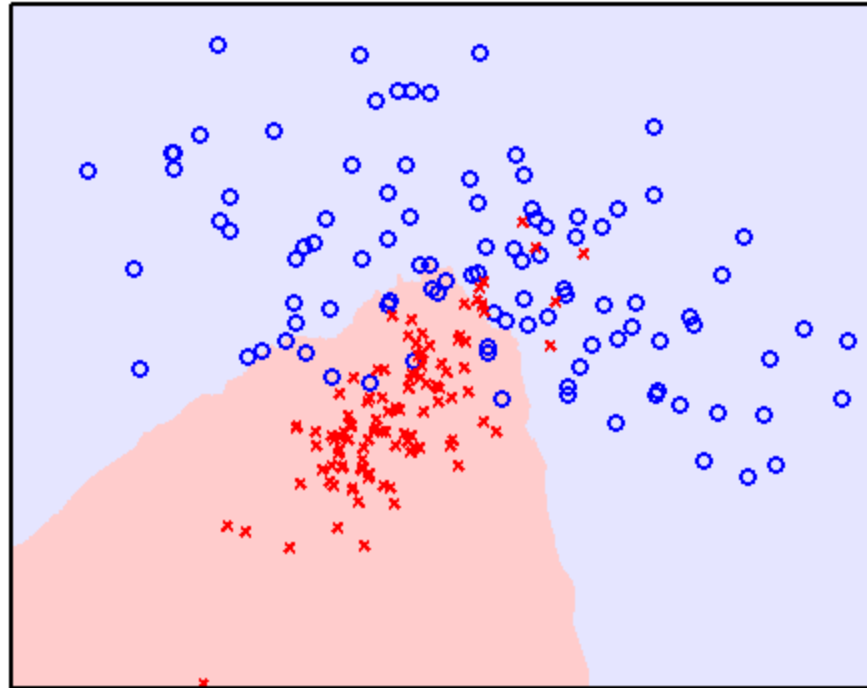
$$k = 3$$

Example



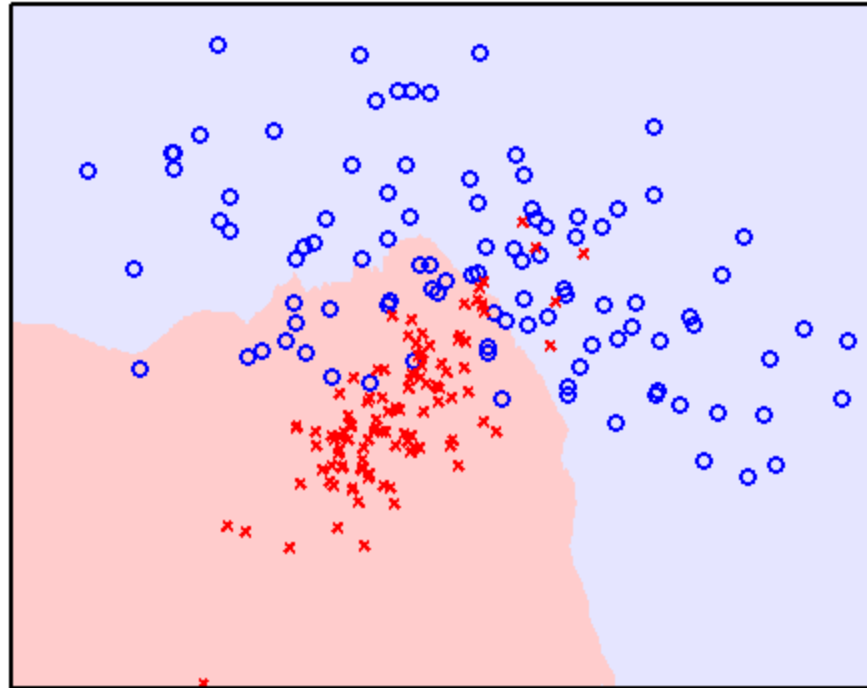
$$k = 5$$

Example



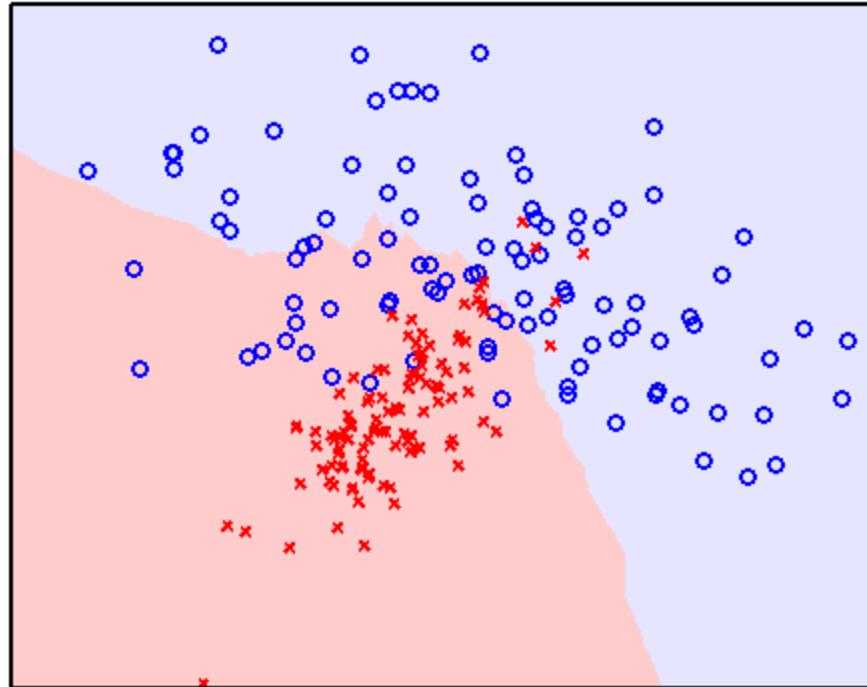
$$k = 25$$

Example



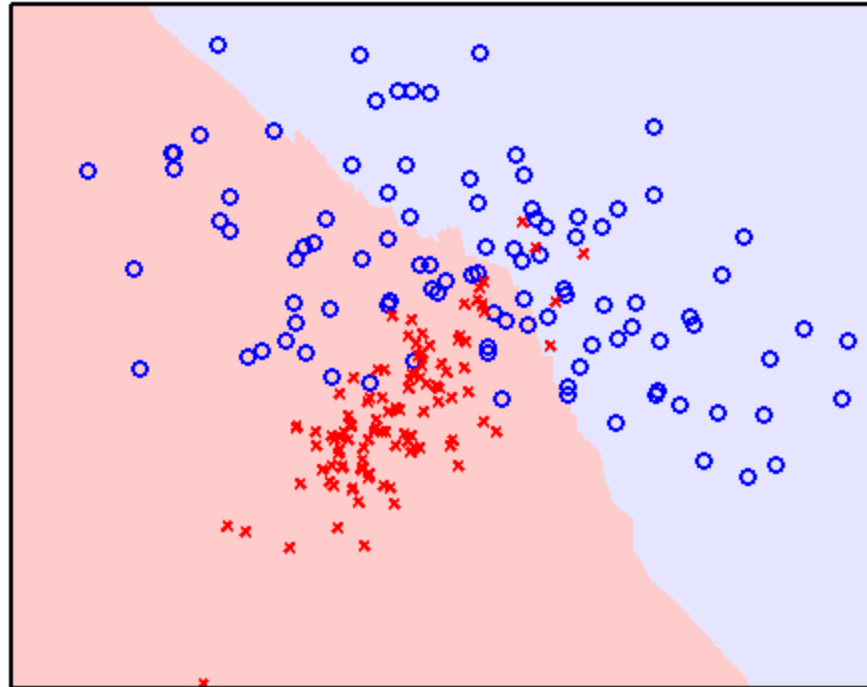
$$k = 51$$

Example



$$k = 75$$

Example



$$k = 101$$

Choosing the size of the neighborhood

Setting the parameter k is a problem of *model selection*

Setting k by trying to minimize the training error is a particularly bad idea

$$\hat{R}_n(\hat{h}_k) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\hat{h}_k(\mathbf{x}_i) \neq y_i\}}$$

What is $\hat{R}_n(\hat{h}_1)$?

No matter what, we always have $\hat{R}_n(\hat{h}_1) = 0$

Not much practical guidance from the theory, so we typically must rely on estimates based on holdout sets or more sophisticated model selection techniques

Consistency

We say that a classification algorithm is **consistent** if when the size of the training set $n \rightarrow \infty$, we have that

$$\mathbb{E} \left[R(\hat{h}_n) \right] \rightarrow R^*$$

where \hat{h}_n is a classifier learned from a training set of size n and R^* is the Bayes risk

Theorem

Let $\hat{h}_{k,n}$ denote the k -nearest neighbor rule with a training set of size n . If $n \rightarrow \infty$, $k \rightarrow \infty$, and $k/n \rightarrow 0$, then $\hat{h}_{k,n}$ is consistent, i.e.,

$$\mathbb{E} \left[R(\hat{h}_{k,n}) \right] \rightarrow R^*$$

Summary

Given enough data, the k -nearest neighbor classifier will do just as well as pretty much any other method

Catch

- The amount of required data can be huge, especially if our feature space is high-dimensional
- The parameter k can matter a lot, so model selection will can be very important
- Finding the nearest neighbors out of a set of millions of examples is still pretty hard
 - can be sped up using k -d trees, but can still be relatively expensive to apply
 - in contrast, many of the other algorithms we will study have an expensive “training” phase, but application is cheap