

Empirical risk minimization (ERM)

Recall the definitions of risk/empirical risk

$$R(h) := \mathbb{P}[h(X) \neq Y] \quad \hat{R}_n(h) := \frac{|\{i : h(\mathbf{x}_i) \neq y_i\}|}{n}$$

Ideally, we would like to choose $h^\# = \arg \min_{h_j \in \mathcal{H}} R(h_j)$

Since we cannot compute $R(h_j)$, instead we choose $h^* = \arg \min_{h_j \in \mathcal{H}} \hat{R}_n(h_j)$

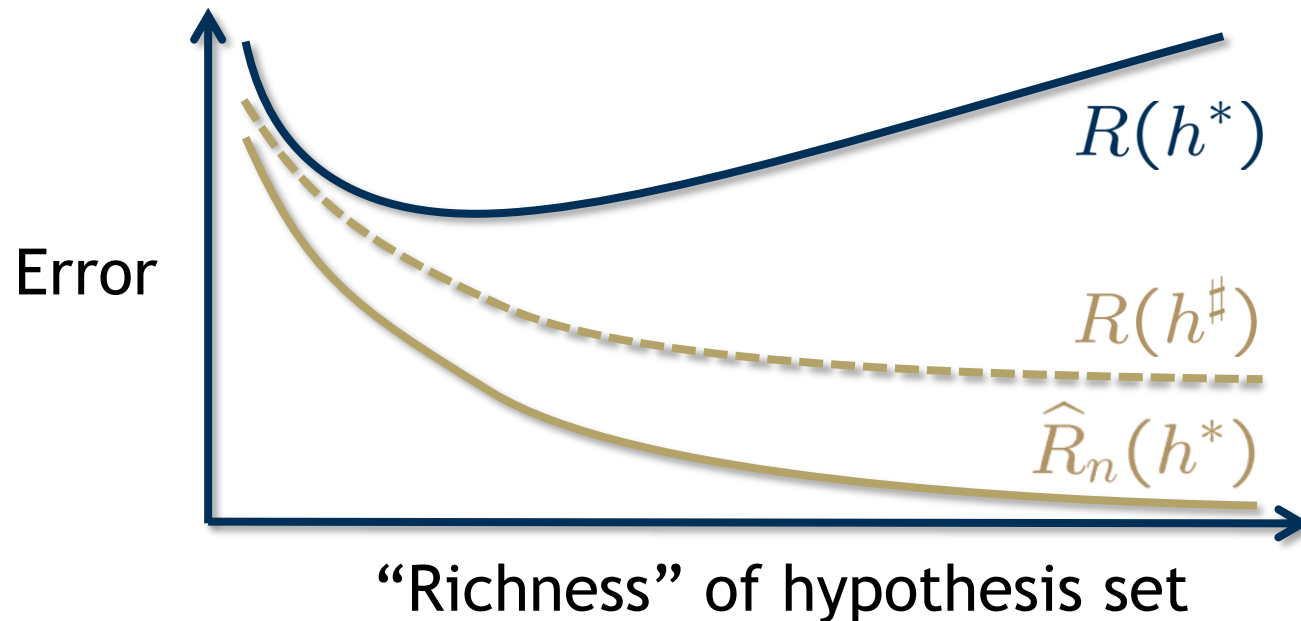
This makes sense if $m = |\mathcal{H}|$ is not too large so that $R(h_j) \approx \hat{R}_n(h_j)$ for all h_j

Unfortunately, we also want m to be large so that $R(h^\#)$ can be as small as possible...

Fundamental tradeoff

More hypotheses ultimately sacrifices our guarantee that $\hat{R}_n(h^*) \approx R(h^*)$

Richer set of hypotheses \rightarrow $\begin{cases} \hat{R}_n(h^*) \downarrow & R(h^\#) \downarrow \\ \hat{R}_n(h^*) - R(h^*) \uparrow \end{cases}$



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 enough to have $R(h^*) \approx R(h^\#) \approx 0$

In general, this may not be possible

There may not be **any** function h with $R(h) \approx 0$

Why not?

$$\text{Noise: } Y = h(X) + N$$

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

Our goal is to formulate a simple rule for minimizing $R(h)$ when the joint distribution of (X, Y) is known

We will let $f_{X,Y}(\mathbf{x}, y)$ denote this joint distribution of (X, Y)

The joint distribution

For any $\mathcal{X} \subset \mathbb{R}^d$ and any $\mathcal{Y} \subset \{0, \dots, K - 1\}$, $f_{X,Y}(\mathbf{x}, y)$ gives us a way to compute the probability that a randomly drawn (X, Y) will satisfy $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$

$$\mathbb{P}[X \in \mathcal{X}, Y \in \mathcal{Y}] = \sum_{y \in \mathcal{Y}} \int_{\mathbf{x} \in \mathcal{X}} f_{X,Y}(\mathbf{x}, y) d\mathbf{x}$$

Conditioning on $X = \mathbf{x}$ results in a conditional distribution on the class labels known as the ***a posteriori distribution***:

$$\eta_y(\mathbf{x}) := p_{Y|X}(y|\mathbf{x}) = \mathbb{P}[Y = y|X = \mathbf{x}]$$

Conditioning on $Y = y$ results in the ***class conditional distribution***:

$$f_{X|Y}(\mathbf{x}|y)$$

Factoring the joint distribution

It is often useful to think about the joint distribution in terms of these conditional distributions

For any fixed (\mathbf{x}, y) we can write

$$f_{X,Y}(\mathbf{x}, y) = \mathbb{P}[Y = y] f_{X|Y}(\mathbf{x}|y)$$

or

$$f_{X,Y}(\mathbf{x}, y) = \eta_y(\mathbf{x}) f_X(\mathbf{x})$$

Both ways of thinking will be useful!

The Bayes classifier

Theorem

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

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

for any possible classifier h

Note: h^* is not restricted to any particular set \mathcal{H} , and hence we will have $R(h^*) \leq R(h^\#) \leq R(h^*)$

Terminology:

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

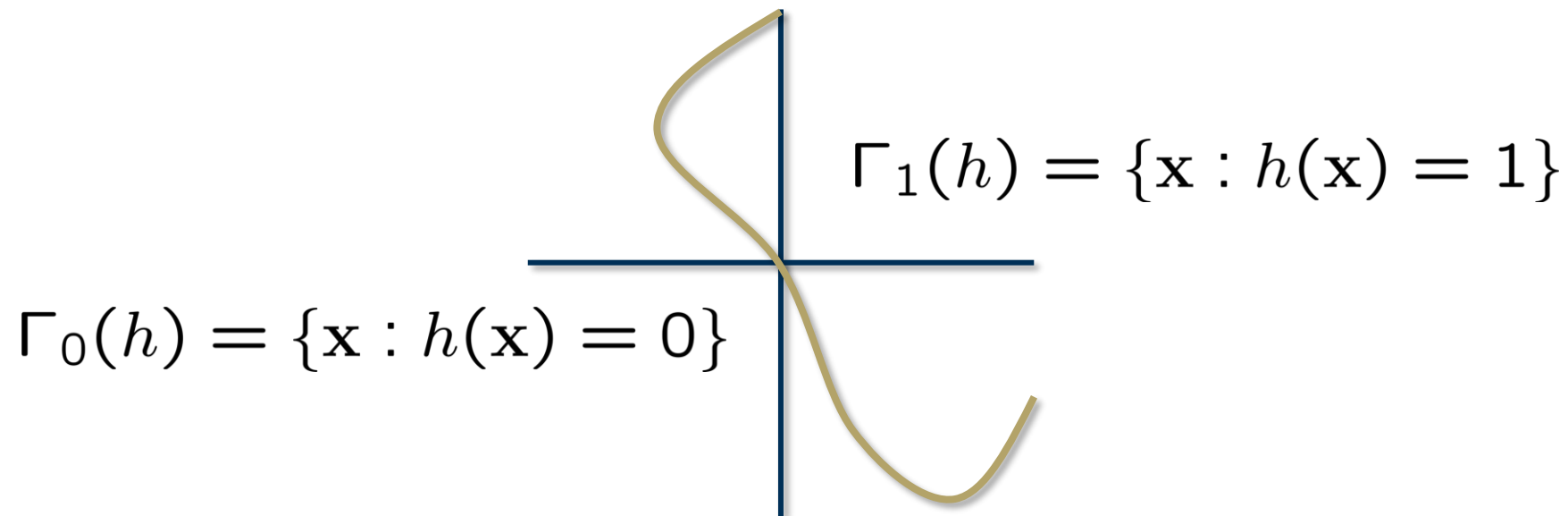
Proof

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

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

Consider an arbitrary classifier h . Denote the decision regions

$$\Gamma_y(h) := \{\mathbf{x} : h(\mathbf{x}) = y\}$$



Proof (Part 2)

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

$$\begin{aligned} &= \sum_{y=0}^{K-1} \pi_y \cdot \mathbb{P}[h(X) = y | Y = y] \\ &= \sum_{y=0}^{K-1} \pi_y \cdot \int_{\Gamma_y(h)} f_{X|Y}(\mathbf{x}|y) d\mathbf{x} \\ &= \sum_{y=0}^{K-1} \int_{\Gamma_y(h)} \pi_y f_{X|Y}(\mathbf{x}|y) d\mathbf{x} \end{aligned}$$

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


$$\mathbf{x} \in \Gamma_y(h) \quad \longleftrightarrow \quad \pi_y f_{X|Y}(\mathbf{x}|y) \text{ is maximal}$$

Proof (Part 3)

Therefore, the optimal h has

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

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

- $h^*(\mathbf{x}) = \arg \max_y \eta_y(\mathbf{x})$
- $h^*(\mathbf{x}) = \arg \max_y \pi_y f_{X|Y}(\mathbf{x}|y)$
- When $K = 2$

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

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

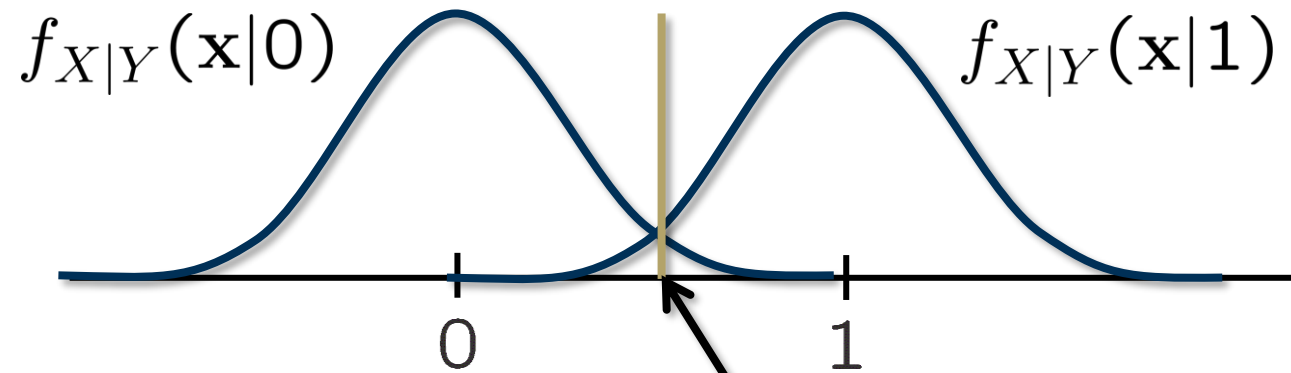
$$h^*(\mathbf{x}) = \arg \max_y f_{X|Y}(\mathbf{x}|y) \quad \text{maximum likelihood classifier/detector}$$

Example

Suppose that $K = 2$ and that

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

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



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

If $\pi_0 = \pi_1$

Example

How do we calculate the *Bayes risk*?

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

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

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

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 it may be desirable to instead minimize

$$C_0\mathbb{P}[h(X) \neq Y|Y = 0] + C_1\mathbb{P}[h(X) \neq Y|Y = 1]$$

- alternatively, it may be better to focus on them directly a la Neyman-Pearson classification

$$\underset{h}{\text{minimize}} \mathbb{P}[h(X) \neq Y|Y = 1]$$

$$\text{subject to } \mathbb{P}[h(X) \neq Y|Y = 0] \leq \alpha$$

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

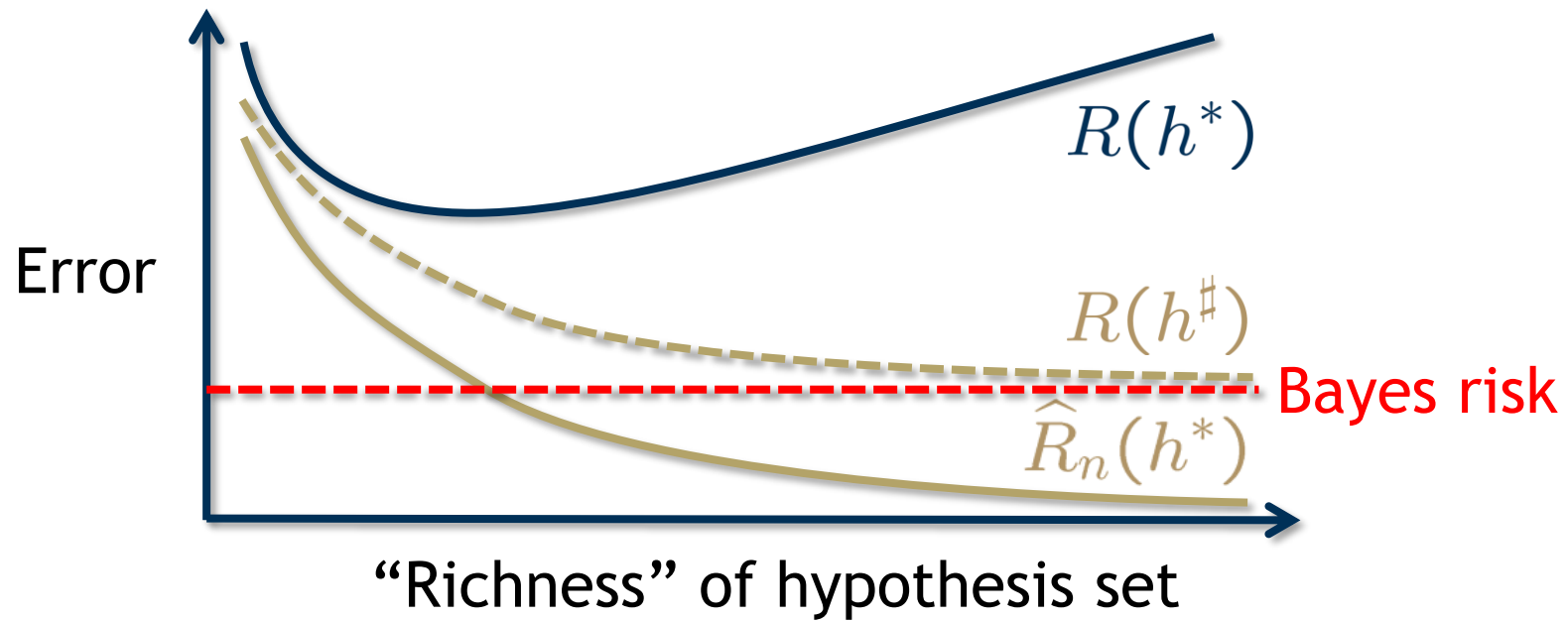
- unbalanced datasets
 - when one class dominates the other, the probability of error will place less emphasis on the smaller class
 - the class proportions in our dataset may not be representative of the “wild”
 - one can use the same ideas as before or alternatively simply minimize something like

$$\mathbb{P}[h(X) \neq Y|Y = 0] + \mathbb{P}[h(X) \neq Y|Y = 1]$$

or

$$\max(\mathbb{P}[h(X) \neq Y|Y = 0], \mathbb{P}[h(X) \neq Y|Y = 1])$$

Fundamental tradeoff



What about learning?

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

Can we get close when 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

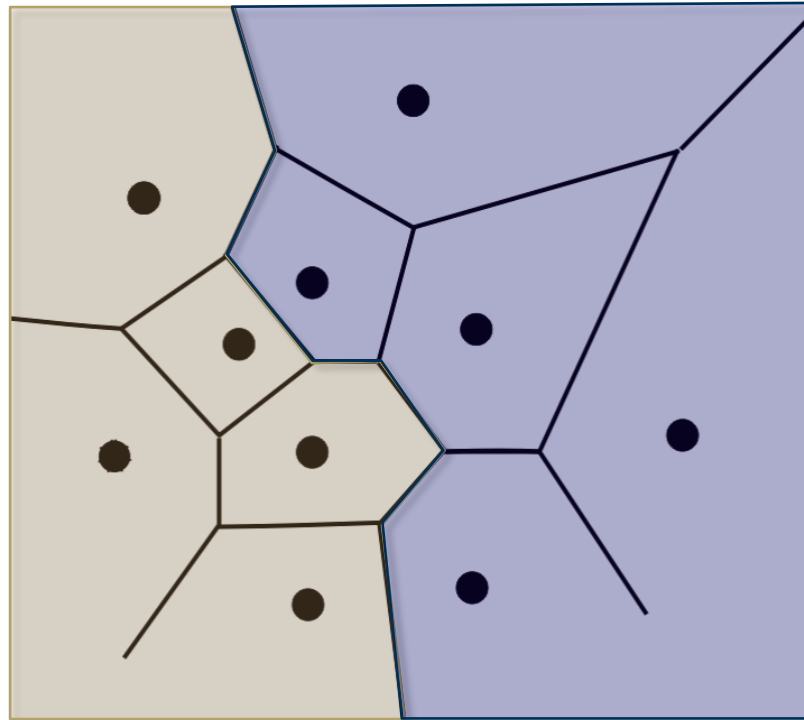
Before we get to these, we will first talk about what is quite possibly the absolute simplest learning algorithm there is...

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

We will begin by restricting our attention to the binary case where $Y \in \{0, 1\}$

Consider the Bayes risk conditioned on $X = \mathbf{x}$:

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

Note that if $h^*(\mathbf{x}) = 0$, then we must have $R^*(\mathbf{x}) = \eta_1(\mathbf{x})$

Similarly, if $h^*(\mathbf{x}) = 1$, then $R^*(\mathbf{x}) = \eta_0(\mathbf{x})$

Since $h^*(\mathbf{x})$ selects the label that maximizes $\eta_y(\mathbf{x})$, we thus have that

$$R^*(\mathbf{x}) = \min(\eta_0(\mathbf{x}), \eta_1(\mathbf{x}))$$

Risk of the nearest neighbor classifier

Now consider the risk of the nearest neighbor classifier conditioned on $X = \mathbf{x}$

$$R^{\text{NN}}(\mathbf{x}) := \mathbb{P}[Y \neq h^{\text{NN}}(\mathbf{x}) | X = \mathbf{x}]$$

Note that for a fixed \mathbf{x} , we are treating Y as random

Here we will further treat $h^{\text{NN}}(\mathbf{x})$ as being random since *it depends on the training dataset*

Thus we have that

$$\begin{aligned} R^{\text{NN}}(\mathbf{x}) &= \mathbb{P}[Y = 0 | X = \mathbf{x}] \mathbb{P}[h^{\text{NN}}(\mathbf{x}) = 1 | X = \mathbf{x}] \\ &\quad + \mathbb{P}[Y = 1 | X = \mathbf{x}] \mathbb{P}[h^{\text{NN}}(\mathbf{x}) = 0 | X = \mathbf{x}] \end{aligned}$$

Risk of the nearest neighbor classifier

$$R^{\text{NN}}(\mathbf{x}) = \mathbb{P}[Y = 0|X = \mathbf{x}] \mathbb{P}[h^{\text{NN}}(\mathbf{x}) = 1|X = \mathbf{x}] \\ + \mathbb{P}[Y = 1|X = \mathbf{x}] \mathbb{P}[h^{\text{NN}}(\mathbf{x}) = 0|X = \mathbf{x}]$$

Note that if \mathbf{x}_{NN} is the nearest neighbor to \mathbf{x} , then

$$\mathbb{P}[h^{\text{NN}}(\mathbf{x}) = y|X = \mathbf{x}] = \mathbb{P}[Y = y|X = \mathbf{x}_{\text{NN}}] \\ = \eta_y(\mathbf{x}_{\text{NN}})$$

Thus, we can write $R^{\text{NN}}(\mathbf{x}) = \eta_0(\mathbf{x})\eta_1(\mathbf{x}_{\text{NN}}) + \eta_1(\mathbf{x})\eta_0(\mathbf{x}_{\text{NN}})$

Intuition from asymptotics

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

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

$$\begin{aligned} R^{NN}(\mathbf{x}) &\rightarrow \eta_0(\mathbf{x})\eta_1(\mathbf{x}) + \eta_1(\mathbf{x})\eta_0(\mathbf{x}) \\ &= 2\eta_0(\mathbf{x})\eta_1(\mathbf{x}) \end{aligned}$$

It is easy to see that

$$\begin{aligned} 2\eta_0(\mathbf{x})\eta_1(\mathbf{x}) &\leq 2 \min(\eta_0(\mathbf{x}), \eta_1(\mathbf{x})) \\ &= 2R^*(\mathbf{x}) \end{aligned}$$

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}

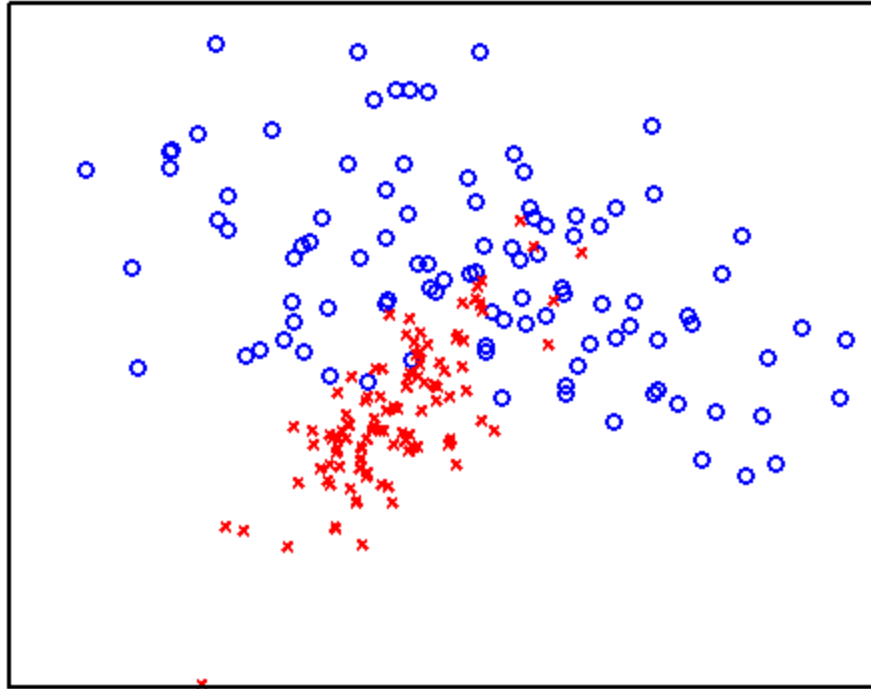
How do we define this more mathematically?

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

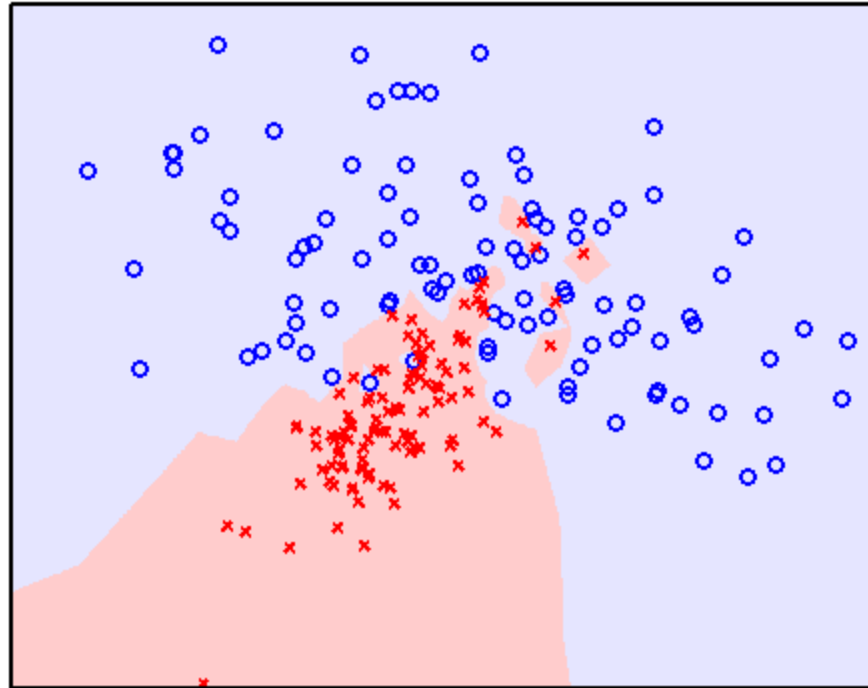
If $y_i = \pm 1$, then we can write the k -nearest neighbor classifier as

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

Example

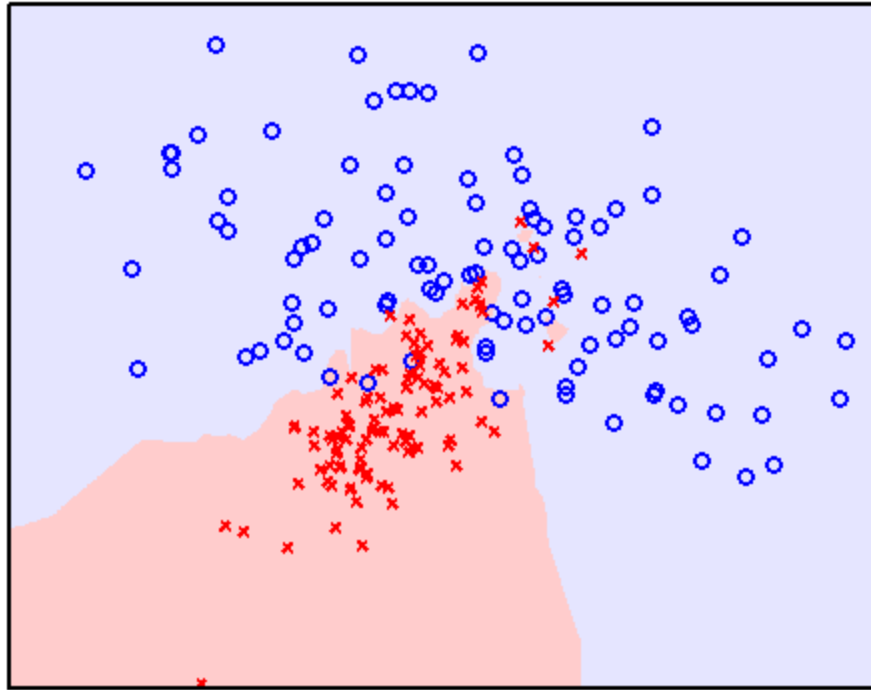


Example



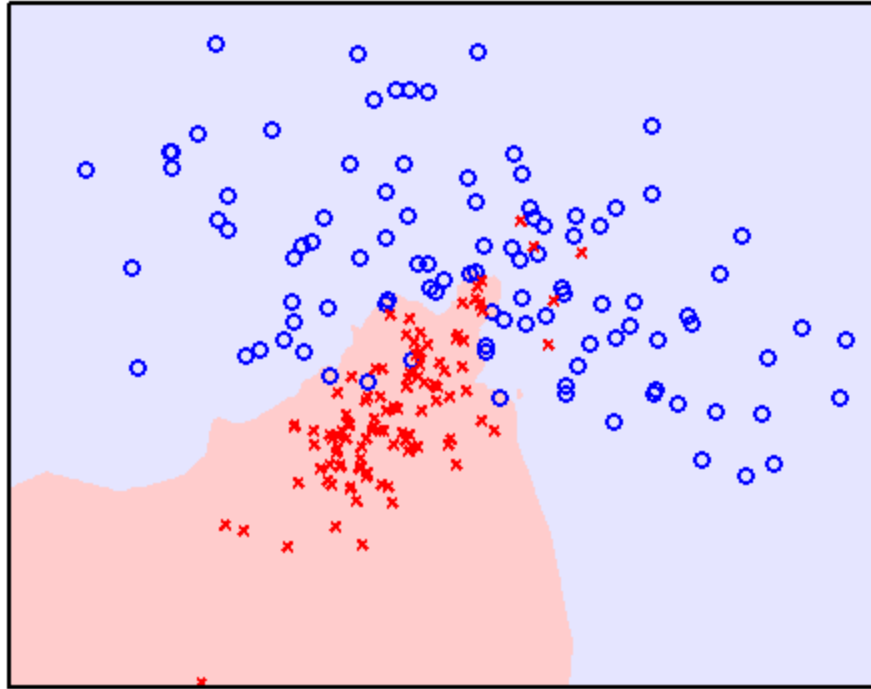
$$k = 1$$

Example



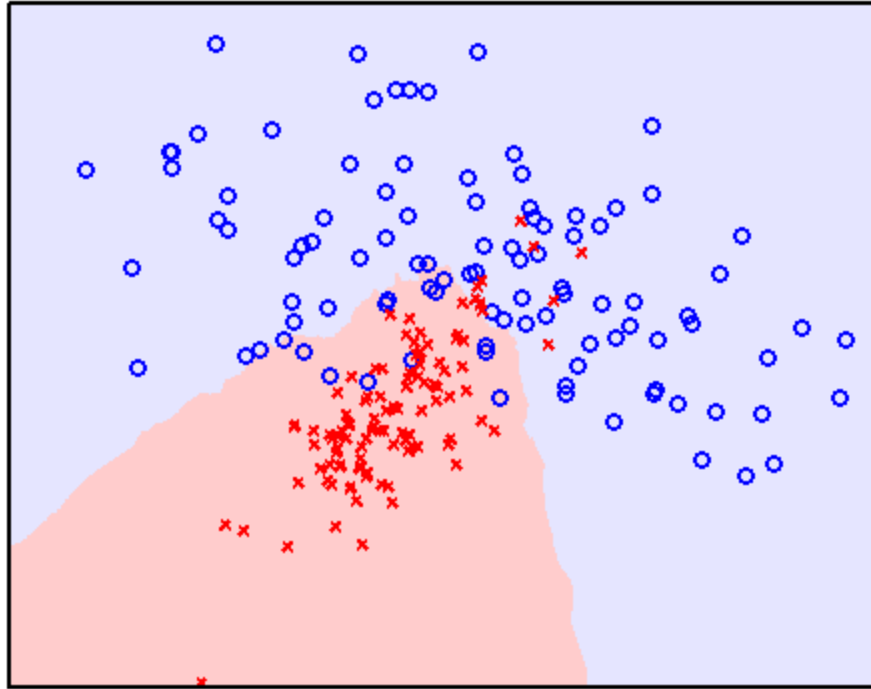
$$k = 3$$

Example



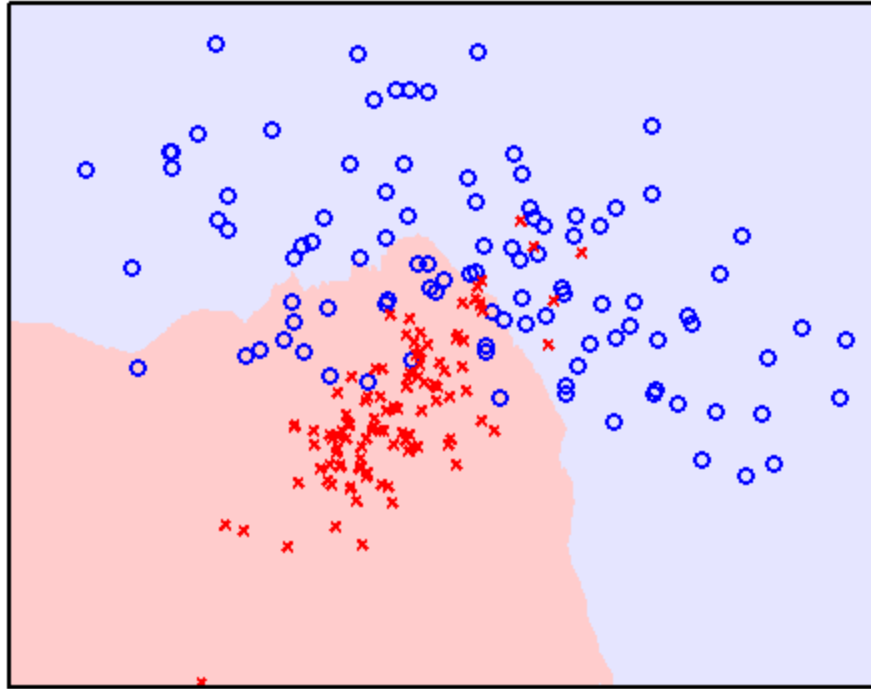
$$k = 5$$

Example



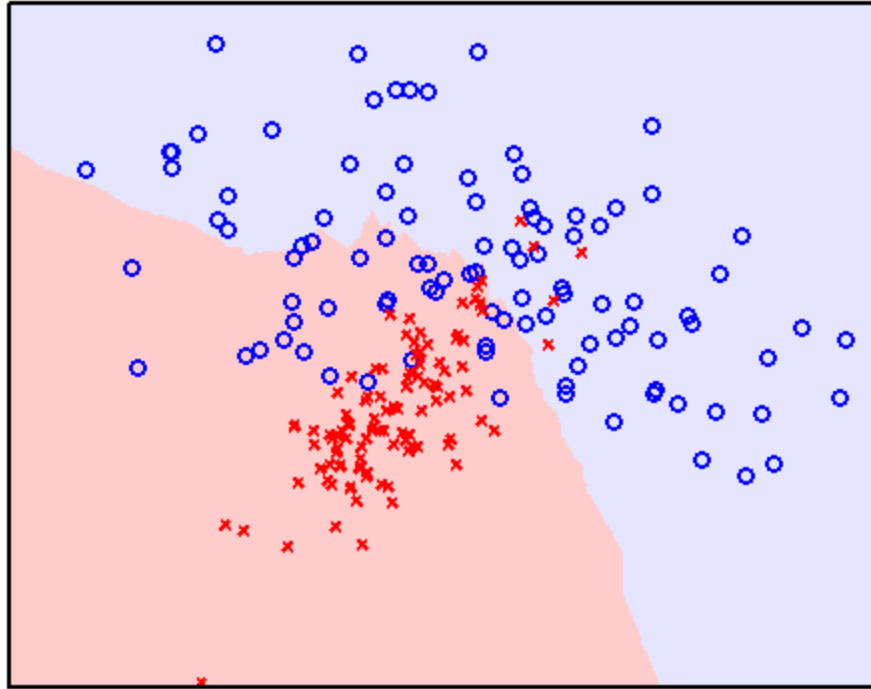
$$k = 25$$

Example



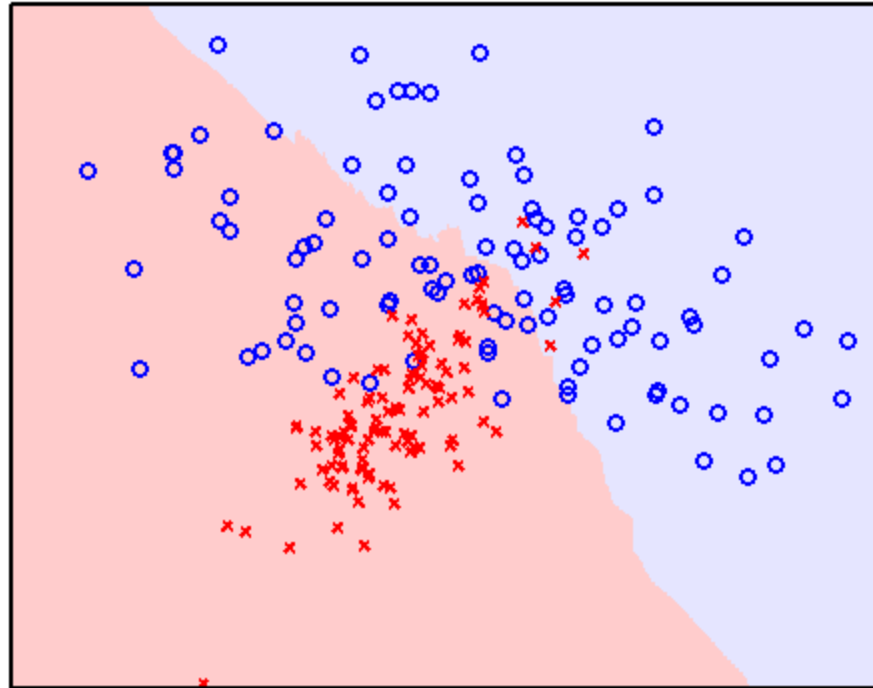
$$k = 51$$

Example



$$k = 75$$

Example



$$k = 101$$

Choosing k : Practice

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

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

$$\hat{R}_n(\hat{h}_k) = \frac{|\{i : \hat{h}_k(\mathbf{x}_i) \neq y_i\}|}{n}$$

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

Choosing k : Theory

Using a similar argument as before, one can show that

$$\lim_{n \rightarrow \infty} R^{\text{kNN}}(\mathbf{x}) \leq \left(1 + \sqrt{2/k}\right) R^*(\mathbf{x})$$

Thus, by letting $k \rightarrow \infty$ and $k/n \rightarrow 0$ as $n \rightarrow \infty$, we can (asymptotically) expect to perform arbitrarily close to the Bayes risk

This is known as ***universal consistency***: given enough data, the algorithm will eventually converge to a classifier that matches the Bayes risk

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