

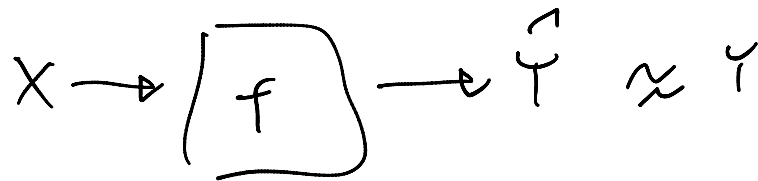
Supervised learning:

Given examples $X_1, Y_1, \dots, X_n, Y_n$

Want to learn a prediction rule f

s.t. for any new unseen example

$$Y \approx f(x) = \hat{Y}$$



The output Y takes values in \mathcal{C}

We classify the cases as:

- $\mathcal{C} = \{a, b\}$ binary classification

where a, b can be anything

e.g.:

- $\mathcal{C} = \{0, 1\}$

- $\mathcal{C} = \{\text{red}, \text{blue}\}$

- $\mathcal{C} = \{\text{high risk}, \text{low risk}\}$

- $\mathcal{C} = \{a, b, c, \dots\}$ multi-class classification

- $\mathcal{C} = \mathbb{R}$ regression

Our first classification algo: KNN

Given a new query point x :

- find k "closest" x -values in

- Given a new x
- Find the k "closest" x -values in the examples x_1, \dots, x_n
 - Note their indices $i_1, \dots, i_k \in [n] = \{1, \dots, n\}$
 - Classify the example as the most common label in the corresponding set of training labels

$$\hat{y} = \text{mode}(\{y_{i_1}, \dots, y_{i_k}\})$$

What does "closest" mean?

Need a distance measure.

Suppose X is a vector of features

$$X_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^p$$

Can use Euclidean distance

$$\|x - x'\|_2 = \sqrt{\sum_j (x_j - x'_j)^2}$$

OR can use any other ^{vector} distance

$$\text{e.g. } \|x - x'\|_\infty = \max_j |x_j - x'_j|$$

Constructing feature vectors

1. If using these vector distances, mean

1. If using these vector distances, may want to first standardize the data:

$$X_{ij} \leftarrow \frac{(X_{ij} - \hat{\mu}_j)}{\hat{\sigma}_j}$$

mean & std dev
of X_{i1}, \dots, X_{ij}

2. If some of the features are categorical

e.g. Which of {Yale, Columbia, Cornell} did you attend

often we encode this using

one-hot encoding:

use three ^{binary} variables to encode

Yale as $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$ Columbia $\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$ Cornell $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

3. Later in class, we'll also talk about more abstract ways to measure similarity.

E.g. $X =$ Content of an email
in words

$Y =$ Whether spam

Out-of-sample evaluation

KNN tries to make \hat{y}_i equal Y

KNN tries to make \hat{y} equal y
 To evaluate how well it — or
 any other supervised learning
 algo — does is to ask how
 often is that true.

If X, Y are r.v.s representing
 new random examples drawn from
 the population of examples, then
 we want low risk

$$R(f) = \mathbb{E}[l(Y, f(X))]$$

where l is a loss fn

$$\text{For now: } l(y, \hat{y}) = \mathbb{I}[y \neq \hat{y}] = \begin{cases} 1 & y \neq \hat{y} \\ 0 & y = \hat{y} \end{cases}$$

To estimate $R(f)$ we can take
 a test set of examples $x_1^{\text{test}}, y_1^{\text{test}}, \dots, x_{n_{\text{test}}}^{\text{test}}, y_{n_{\text{test}}}^{\text{test}}$
 drawn at random from the population
 of examples & compute an empirical avg

$$\hat{R}_{n_{\text{test}}}(f) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} l(y_i^{\text{test}}, f(x_i^{\text{test}}))$$

This is a consistent & unbiased

estimate of $R(f)$.

Often, we get the test set by splitting off from the training data.

Bayes classifier & Bayes (Error) Rate

So "NN" is just one classifier.

But: What is the BEST classifier?

The one that minimizes $R(f)$

$$\begin{aligned}R(f) &= \mathbb{E}[l(Y, f(x))] \\&= \mathbb{E}[\mathbb{I}[Y \neq f(x)]] \\&= \mathbb{E}[1 - \mathbb{I}[Y = f(x)]] \\&= \mathbb{E}[\mathbb{E}[1 - \mathbb{I}[Y = f(x)] | X]] \quad (\text{total expectation}) \\&= \mathbb{E}\left[\sum_{y \in \mathcal{G}_Y} P(Y=y|X) (1 - \mathbb{I}[y=f(x)])\right] \\&= 1 - \mathbb{E}\left[\sum_{y \in \mathcal{G}_Y} P(Y=y|X) \mathbb{I}[y=f(x)]\right]\end{aligned}$$

For each x , we can choose one y [^] so that $f(x) = y$

Which y to choose to minimize $R(f)$?

Maximize the 2nd term

for each x , choose $f(x)$ to maximize

$$\sum_{y \in \mathcal{C}} \mathbb{P}(Y=y|X=x) \mathbb{I}[y=f(x)]$$

\Rightarrow choose $f(x)$ to be the maximizer
of $\mathbb{P}(Y=y|X=x)$

$$f^*(x) = \operatorname{argmax}_{y \in \mathcal{C}} \mathbb{P}(Y=y|X=x) = \operatorname{mode}(\mathbb{P}(Y|X=x))$$

f^* is known as the Bayes classifier

$R(f^*)$ is known as the Bayes (error) rate

\mathcal{R} the best possible risk

f^* gets at the fundamental limit
of the predictability of Y from X

kNN as a probabilistic classifier

Can modify the output of kNN as
a probability estimate:

$$\hat{\mathbb{P}}(Y=y|X) = \frac{1}{k} \sum_{j=1}^k \mathbb{I}[Y_{i_k} = y]$$

fraction of k nearest
neighbors that have

fraction of
neighbors that have
label y

Our k NN prediction can be written

$$\hat{Y} = \underset{y \in \mathcal{C}}{\operatorname{argmax}} \hat{P}(Y=y|X)$$

Mimics Bayes classifier w/ an estimated
conditional prob

For $\mathcal{C} = \{0, 1\}$ (binary classification),

this translates to

$$\hat{Y} = \begin{cases} 0 & \hat{P}(Y=1|X) \leq \frac{1}{2} \\ 1 & \hat{P}(Y=1|X) > \frac{1}{2} \end{cases}$$

$$= \mathbb{I} \left[\hat{P}(Y=1|X) > \frac{1}{2} \right]$$

Can modify this and use other thresholds

$$\hat{Y} = \mathbb{I} \left[\hat{P}(Y=1|X) > \theta \right] \quad \theta \in \mathbb{R}$$

What this does is change the

classification rates.

Classification rates: binary case

Classification rates: binary case

Exhaustively categorize all cases of prediction:

-	$Y=1, \hat{Y}=1$	TP	} T = true F = false P = positive N = negative
-	$Y=0, \hat{Y}=1$	FP	
-	$Y=0, \hat{Y}=0$	TN	
-	$Y=1, \hat{Y}=0$	FN	

$$\# TP = \sum_{i=1}^{n_{\text{test}}} \mathbb{I} [Y_i^{\text{test}} = 1, f(x_i^{\text{test}}) = 1]$$

$$\# FP = \dots$$

Confusion matrix:

$\hat{Y} \backslash Y$	1	0
1	#TP	#FP
0	#FN	#TN

Sometimes row numbers

Sometimes divide everything by n_{test}

Rates:

$$\text{True positive rate} = \text{TPR} = \frac{\#TP}{\#P} = \frac{\#TP}{\#TP + \#FN}$$

aka Recall

$$\text{TNID} = \frac{\#TN}{\dots} = \frac{\#TN}{\dots}$$

$$TNR = \frac{\#TN}{\#FN} = \frac{\#TN}{\#TN + \#FP}$$

$$\begin{aligned} FPR &= 1 - TNR \\ &= \frac{\#FP}{\#TN + \#FP} \end{aligned}$$

$$FNR = (1 - TPR) = \frac{\#FN}{\#TP + \#FN}$$

Accuracy $Acc = 1 - (\text{Any 0-1 loss})$

$$= \frac{\#TP + \#TN}{\#TP + \#TN + \#FP + \#FN}$$