

Recap: Supervised Learning

$$x_1, y_1, \dots, x_n, y_n \xrightarrow{\text{Algo}} \boxed{f} \rightarrow \hat{y} \stackrel{?}{\approx} y$$

$\uparrow$   
 $x$

kNN as a prob classifier

$$\hat{P}(Y=y|x) = \frac{1}{K} \sum_{j=1}^K \mathbb{I}[Y_j = y]$$

For  $\mathcal{G} = \{0, 1\}$

$$\hat{y} = \underset{j \in \{0, 1\}}{\text{argmax}} \hat{P}(Y=j|x) = \mathbb{I}[\hat{P}(Y=1|x) > \frac{1}{2}]$$

let's use  $\mathbb{I}[\hat{P}(Y=1|x) > \theta]$  instead

this changes the classification rates.

Classification rates

#TP = no. of test examples  
w/  $Y=1$   $\hat{y}=1$

#TN, #FN, #FP def similarly

Confusion matrix

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

Rates:

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

$$\text{TNR} = \frac{\#TN}{\#N} = \frac{\#TN}{\#TN + \#FP}$$

$$1 - \text{FNR} = \frac{\#P - \#FP}{\#P} = \frac{\#TP + \#FN}{\#P}$$

$$\text{TNR} = \frac{\#TN}{\#N} = \frac{\#TN}{\#TN + \#FP}$$

$$\text{FNR} = 1 - \text{TPR} \quad \text{FPR} = 1 - \text{TNR}$$

$$\text{Acc} = 1 - (\text{Avg 0-1 loss}) = \frac{\#TP + \#TN}{\#TP + \#TN + \#FP + \#FN}$$

$$\text{Recall} = \text{TPR}$$

$$\text{Precision} = \frac{\#TP}{\#TP + \#FP}$$

$$\text{F-score} = \frac{2}{\frac{1}{\text{prec}} + \frac{1}{\text{recall}}}$$

## Receiver Operating Characteristic (ROC) curve

The ROC curve is the plot of FPR (x-axis) vs TPR (y-axis) as we vary the classification threshold,  $\theta$ .

ROC tells us how good a prob classification is, regardless of the specific threshold.

Can quantify how good an ROC curve looks in terms of its Area Under the curve (AUC)

$$\text{AUC} = \frac{1}{2} \rightarrow \text{corresponds to}$$

$AUC = \frac{1}{2} \rightarrow$  corresponds to  
random guessing

$AUC = 1 \rightarrow$  corresponds to  
perfect prediction

Fun & useful fact:

AUC exactly equal to the prob  
that given <sup>two</sup> examples, one pos & one neg,  
that our score sorts them in the  
right order.

## Regression

Recall regression is the setting  
when  $Y$  can be real-valued

$$Y \in \mathcal{G} = \mathbb{R}$$

## kNN as a regression method

$$\hat{Y} = \frac{1}{K} \sum_{j=1}^K Y_{i_j}$$

## Lazy vs Eager training

kNN is an example of "lazy" training  
or "memory-based"

- wait until we get a query  $x$   
in order to compute  $f(x)$

" " " " " "

in order to compute  $f'(x)$   
"Eager" training tries to build our  
predictor  $f'$  ahead of time  
& summarize it into a single  
model.

## The Linear Model

Posit a model parametrized  
by  $\beta \in \mathbb{R}^{p+1}$

$$\begin{aligned} f_{\beta}(x) &= \beta_0 + \sum_{i=1}^p \beta_i \cdot x_i \\ &= \beta_0 + \beta_{1:p}^T x \end{aligned}$$

|  
intercept/bias

coefficients/  
slopes/weights

To make things simple,  
we'll just assume that one of  
the features in  $x$  is the constant  
feature:  $x_{i+1} = 1 \quad \forall i$

Then we write  $\beta \in \mathbb{R}^p$

$$f_{\beta}(x) = \beta^T x$$

The question  $\Rightarrow$  how how to  
choose  $f_{\beta}$  in this model  
class — i.e. how to choose  $\beta$

# Least Squares

To evaluate the fit of any  $f$  to regression data we'll use the squared error loss:

$$l(y, \hat{y}) = (y - \hat{y})^2$$

This leads to the avg squared residuals (empirical) risk:

$$\begin{aligned}\hat{R}_n(f) &= \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)) \\ &= \frac{1}{n} \sum_{i=1}^n \underbrace{(y_i - f(x_i))}_{\text{residual}}^2\end{aligned}$$

The linear model defines a particular class of  $f$ 's:

$$f \in \{f_\beta : \beta \in \mathbb{R}^p\}$$

Seek the  $f$  in this class that does the best on the training data.

$$\begin{aligned}\hat{R}_n(f_\beta) &= \frac{1}{n} \sum_{i=1}^n (y_i - f_\beta(x_i))^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \beta^T x_i)^2 \\ &= \frac{1}{n} \|Y - X\beta\|_2^2\end{aligned}$$

Where we define  $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$

design matrix  $\rightarrow \bar{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \in \mathbb{R}^{n \times p}$

$$\bar{X}\beta = \begin{pmatrix} x_1^T \beta \\ \vdots \\ x_n^T \beta \end{pmatrix} = \begin{pmatrix} f_\beta(x_1) \\ \vdots \\ f_\beta(x_n) \end{pmatrix} \in \mathbb{R}^n$$

We want

$$\min_{\beta \in \mathbb{R}^p} \hat{R}_n(f_\beta)$$

Call the optimizer  $\hat{\beta}$   $\nearrow$  differentiable & convex (cup-like)  $f_n$  so min occurs at critical pts

$$\begin{aligned} \frac{\partial \hat{R}_n(f_\beta)}{\partial \beta_j} &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \beta_j} (\beta^T x_i - y)^2 \\ &= \frac{1}{n} \sum_i 2 (\beta^T x_i - y) x_{ij} \\ &= \frac{2}{n} \sum_i (\beta^T x_i - y) \bar{X}_{ij} \end{aligned}$$

$$\nabla_{\beta} \hat{R}_n(f_\beta) = \begin{pmatrix} \partial \hat{R}_n(f_\beta) / \partial \beta_1 \\ \vdots \\ \partial \hat{R}_n(f_\beta) / \partial \beta_p \end{pmatrix} = \frac{2}{n} \bar{X}^T (\bar{X}\beta - \bar{Y})$$

Want to solve for  $\beta$  in

$$0 = \nabla_{\beta} \hat{R}_n(f_\beta) = \frac{2}{n} \bar{X}^T (\bar{X}\beta - \bar{Y})$$