

# STATISTICAL DECISION THEORY

23/02/2021

SUPERVISED LEARNING : Given observations  $(x_i, y_i), i=1, \dots, n$

find a rule

$\hat{g}$

regression function (quantitative response  $y$ )

classifier (qualitative response  $y$ )

that allows to predict  $y$  from  $x$ .

So  $\hat{g}(x)$  will be some "approximation" of a true  $g(x)$ , describing the "true" relationship between  $y$  and  $x$ .

What "true" is though depends on the way we set up our statistical learning problem (how would we like to describe/capture this relationship?)

Statistically, this boils down to choosing a loss function that decides in which way  $g(x)$  should describe the data.

We are going to explore this in the regression and classification setting.

## REGRESSION SETTING

In regression, a popular choice is the squared error loss:

SQUARED  
ERROR  
LOSS

$$L(Y, f(X)) = (Y - f(X))^2$$

Leading to the criterion: find  $f$  that minimizes

EXPECTED  
PREDICTION  
ERROR

$$EPE(f) = E_{x,y} [(Y - f(X))^2]$$

We will see that setting the problem in this way is equivalent to considering  $f(x) = E(Y | X=x)$ , the conditional mean of  $Y$  given  $x$ , also called regression function.

Theorem  $E[(y - f(x))^2]$  has a minimum when  
 $f(x) = E[y | X=x] \quad \forall x$

Proof

$$\begin{aligned} E[(y - f(x))^2] &= \iint (y - f(x))^2 \underbrace{g_{x,y}(x,y)}_{\text{joint density of } x \text{ and } y} dx dy \\ &= \iint (y - f(x))^2 g_{Y|X}(y|x) g_X(x) dy dx \\ &= \int \left[ \int (y - f(x))^2 g_{Y|X}(y|x) dy \right] g_X(x) dx \\ &\quad \xrightarrow{\text{E}_{Y|X}[(y - f(x))^2 | X=x]} \end{aligned}$$

$$g_{x,y}(x,y) = g_{Y|X}(y|x) g_X(x)$$

Iterated rule of expectation:

$$E_{X,Y}(h(x,y)) = E_X[E_{Y|X}(h(x,y)|X)]$$

Note:  $Z \sim h(z)$   
 $\min_a E[(z-a)^2] = E[(z - E(z))^2]$

$$\begin{aligned} \frac{d}{da} E(Z^2 - 2az + a^2) &= \\ &= \frac{d}{da} (E(Z^2) - 2a E(Z) + a^2) = -2E(Z) + 2a \\ \rightarrow \hat{a} &= E(Z) ! \end{aligned}$$

Therefore  $\int (y - f(x))^2 g_{Y|X}(y|x) dy \geq \int (y - E(Y|X=x))^2 g_{Y|X}(y|x) dy$ .

But this is true for all  $x$ , so  $E[(y - f(x))^2] \geq E[(y - E(Y|X))^2]$ .

## SQUARED ERROR LOSS $\leftrightarrow$ $f(x) = E(Y|X=x)$

- ① Different losses would lead to different functions,  
eg  $L(Y, f(X)) = |Y - f(X)| \rightarrow f(x) = \text{median}(Y|X=x)$
- ② Statistical learning under a squared error loss consists in different methods for estimating  $E(Y|X=x)$ , which is unknown.

QUANTILE  
REGRESSION

For example :

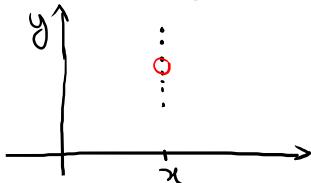
- LINEAR REGRESSION (parametric)

$$E(Y|X=x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

with  $\beta_0, \beta_1, \dots, \beta_p$  estimated from training data (least squares).

## K-NEAREST NEIGHBOUR (non-parametric)

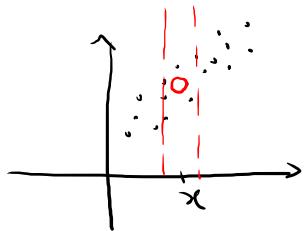
At each point  $x$ , calculate the average (sample mean) of all training points  $y_i$  with  $x_i = x$ .



Typically only few observations for each  $x$ , often only one, so

$$\hat{f}(x) = \text{Average } (y_i \mid x_i \in N_k(x))$$

neighbourhood:  
k points closest to x



Two approximations:

① Expectation approximated by sample mean

② Conditioning on a point relaxed by conditioning on a region around that point.

Generally good only for very large sample sizes (hence small neighbourhood) and also small p. This simple method however inspired more complex methods, called kernel methods (which use weights that decrease smoothly as we go away from the target point).

## CLASSIFICATION SETTING

What is a reasonable loss function here?

Response  $y$  is categorical, with the typical case BINARY (yes/no).

A classifier is therefore going to decide which class to assign an object  $x$ .

Our  $\hat{y}$  in this case is therefore a predicted class and the "error" that we commit is to assign  $x$  to the wrong class.

Since predictions are classes, we can represent the loss as a  $K \times K$  table, with  $K$  the number of classes.

In the binary case:

ZERO-ONE LOSS

		PREDICTED	
		0	1
EXPECTED	0	0	1
	1	1	0

This leads to the criterion: find  $f(x)$  that minimizes

EXPECTED  
0-1 LOSS

$$\mathbb{E}_{x,y} L(y, f(x))$$

## MINIMIZING EXPECTED 0-1 LOSS

Let us consider a fixed  $x$ . If we predict  $x$  to class 0 ( $\hat{y}=0$ ) then

$$E[L(y, 0)] = 0 \cdot p(0|x) + 1 \cdot p(1|x) = p(1|x)$$

Whereas, if we predict  $x$  to class 1 ( $\hat{y}=1$ ),

$$E[L(y, 1)] = 0 \cdot p(1|x) + 1 \cdot p(0|x) = p(0|x)$$

So according to 0-1 loss, we should classify  $x$  to class 1 if

$$p(1|x) > p(0|x), \text{ and } 0 \text{ otherwise.}$$

In general, with more than 2 classes, the loss is minimized by assigning each observation to the most likely class given the predictor values:

ASSIGN  $x$  to CLASS  $j = \operatorname{argmax}_{c \in \text{classes}} P(y=c|x)$

BAYES  
CLASSIFIER

## BAYES CLASSIFIER

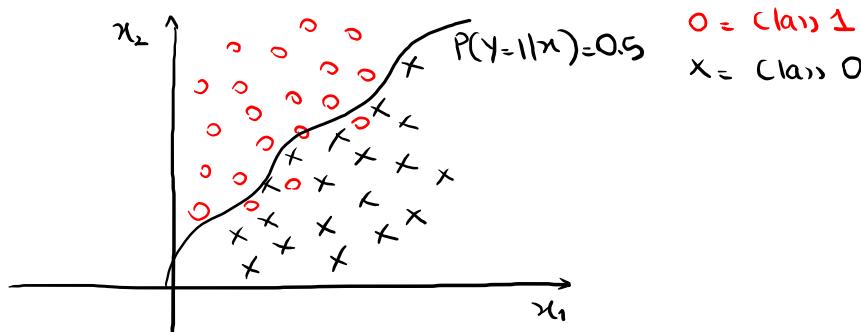
Going back to the two classes, the Bayes classifier results in:

ASSIGN  $x$  to CLASS  $\begin{cases} 1 & \text{if } P(1|x) > 0.5 \\ 0 & \text{if } P(1|x) \leq 0.5 \end{cases}$

$$\begin{aligned} P(1|x) &> P(0|x) \\ \Rightarrow P(1|x) &> 1 - P(1|x) \\ \Rightarrow 2P(1|x) &> 1 \\ \Rightarrow P(1|x) &> 0.5 \end{aligned}$$

## BAYES DECISION BOUNDARY

: All  $x$  such that  $P(Y=1|x) = 0.5$

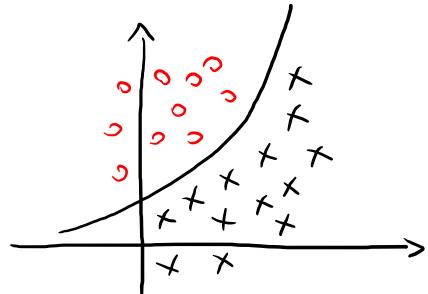


## BAYES ERROR RATE

The overall error that one commits can be defined by

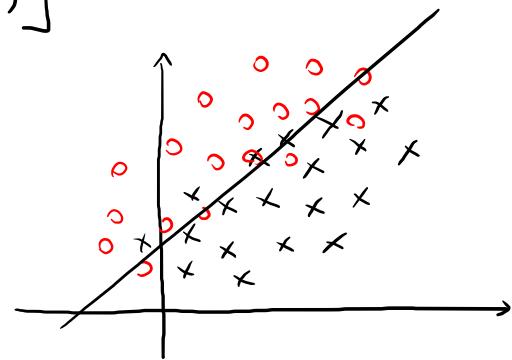
BAYES  
ERROR  
RATE

$$\text{BER} = 1 - \mathbb{E}_x \left[ \max_c P(y=c | X) \right]$$



NO OVERLAP,  $\text{BER} = 0$

(For each  $x$ ,  $\max_c P(y=c | x)$  is 1 for one class and zero for all the others)



OVERLAP,  $\text{BER} > 0$

(most cases, analogous to irreducible error in regression problems)

## MORE GENERAL LOSS FUNCTION

A 0-1 loss assumes equal costs of misclassification.

This can be relaxed by choosing a more general loss function:

0	1
1	0

MISCLASSIFICATION LOSS

		PREDICTED	
		0	1
TRUE	0	0	$c_0$
	1	$c_1$	0

In this case:

$$E[L(Y, 0)] = c_1 p(1|x)$$

$$E[L(Y, 1)] = c_0 p(0|x)$$

=> ASSIGN  $x$  to class 1 if

$$c_1 p(1|x) > c_0 p(0|x) \Rightarrow$$

$$p(0|x) = 1 - p(1|x)$$

Useful in many contexts:

- ① One misclassification more costly than another (wrong medical diagnosis, giving mortgages to bad customers, ...)
- ② Unbalanced class sizes  
(only 5% in one class so 99% if we classify everything to the large class)

$p(1|x) > \frac{c_0}{c_1 + c_0}$

NEW THRESHOLDS DEFINING THE DECISION BOUNDARY (0.5 if  $c_0 = c_1$ )

## REGRESSION vs CLASSIFICATION SETTING

REGRESSION Under a squared error loss, we are interested in

$$f(x) = E(Y|X=x)$$

(Different methods will provide different estimates.)

CLASSIFICATION Assign  $x$  to the class 1 according to  
(binary)

$$f(x) = P(Y=1|X=x)$$

( $P(Y=1|X=x)$  is unknown. Different methods will provide different ways of estimating this quantity.)

## ASSESSING MODEL ACCURACY

In practice, we use our chosen method to estimate  $\hat{g}(x)$  from a training data  $(x_i, y_i)$ ,  $i=1, \dots, n$  and then test its performance on some test data  $(x_i^{(t)}, y_i^{(t)})$ ,  $i=1, \dots, m$  (unseen data, not used to build  $\hat{g}(x)$ ).

### REGRESSION

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i^{(t)} - \hat{g}(x_i^{(t)}))^2$$

MEAN-SQUARED  
ERROR

### CLASSIFICATION ( $c_0 = 1$ , binary)

Counting misclassification errors

CONFUSION  
MATRIX

		PREDICTED	
		0	1
TRUE	0	TN	FP
	1	FN	TP

TN : True Negative

TP : True Positive

FP : False Positive

FN : False Negative

$$m = TN + TP + FP + FN$$

## ERROR RATE (ON TEST DATA)

		PREDICTED
		0      1
TRUE	0	TN      FP
	1	FN      TP

$\Rightarrow$

$$TPR = \frac{TP}{TP + FN} \quad \text{SENSITIVITY}$$

$$FPR = \frac{FP}{TN + FP} = 1 - \frac{TN}{TN + FP} \quad 1 - \text{SPECIFICITY}$$

ERROR RATE

$$\frac{FN + FP}{TN + TP + FN + FP} = \frac{FN + FP}{3} = \frac{1}{3} \sum_{i=1}^3 (y_i^{(t)} \neq \hat{y}_i^{(t)})$$

If  $c_0 \neq c_1$ :

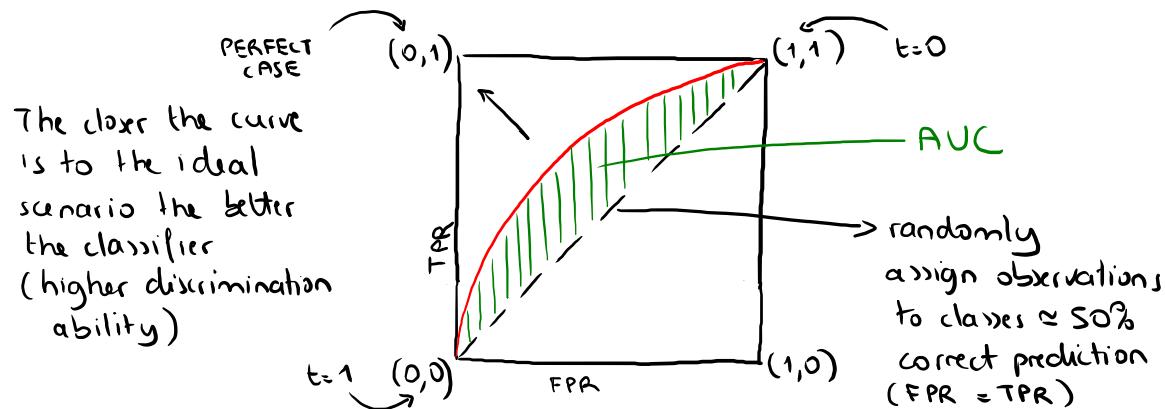
MISCLASSIFICATION COST

$$: \frac{c_0 FP + c_1 FN}{3} \quad (\text{weighted error rate})$$

## RECEIVER OPERATING CHARACTERISTIC (ROC) CURVE

Since misclassification costs may not be known exactly, it is also common to evaluate a classifier across the full range of thresholds.

ROC curve : Plots TP vs FP for  $t \in [0, 1]$



It is common to summarise the curve with the Area Under the Curve (AUC) as a single measure of accuracy. Different models can be compared based on their AUCs, though careful with crossing curves!