

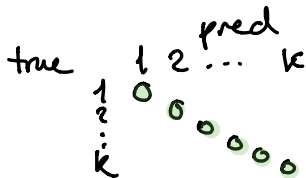
CLASSIFICATION

Bayes classifier: $P(Y=k | X=x)$
 $k=1, \dots, K$ ↑ covariates
↑ explanatory variables

training data x, y → fit method to produce $\hat{P}(Y=k | X=x)$
 • LDA
 • logistic regression
↑ test data

Test data: true class and $\hat{P}(Y=k | X=x)$
 x, y

confusion matrix



how to use this to classify?
 choose k with $\hat{P}(Y=k | X=x)$ max over $1, \dots, k$

misclassification rate

to evaluate the goodness of our method

Special case of $k=2$ →

÷ nondisease $Y=0$

+ disease $Y=1$

		predict		total
		-	+	
true	-	TN	FP	N ← truth
	+	FN	TP	P ← truth
total		N*	P*	

↙ ↘
 predictions

$$\text{sensitivity} = \frac{TP}{P}$$

$$\text{specificity} = \frac{TN}{N}$$

With $k=2$ then the max class will have $\hat{P}(Y=k | X=x) > 0.5$
 ↑
 Bayes class. rule

1) Let $p(x) = \hat{P}(Y=1 | X=x)$ then we classify as disease (1)
 if $p(x) > 0.5$.

→ example (Ath. Heart disease)

$$\text{sens} = 0.625$$

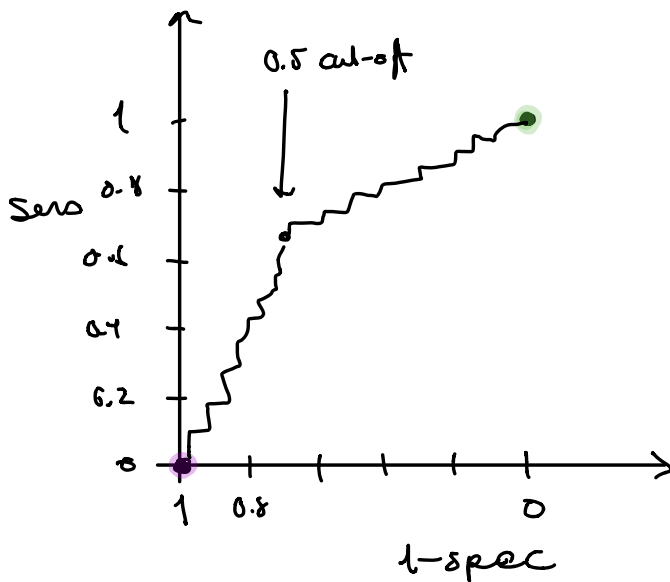
$$\text{spec} = 0.778$$

2) But - is the cost of misclassification the same for both mistakes?
 say + if - Type I, false positive
 say - if + Type II, false negative
 → want to investigate different cut-offs

on $p(x)$ for classification.

$$\begin{array}{l} p(x) > 0.1 \rightarrow \text{sens, spec} \\ > 0.2 \rightarrow \quad \vdots \quad \vdots \\ \vdots \end{array}$$

⇒ plot (sens, 1-spec) = ROC curve



$$\begin{array}{l} p(x) \geq 0 \Rightarrow \text{all } y=1 \\ \text{sens} = \frac{TP}{P} = \frac{P}{P} = 1 \\ \text{spec} = \frac{TN}{N} = \frac{0}{N} = 0 \end{array}$$

$$\begin{array}{l} p(x) \geq 1 \Rightarrow \text{all } y=0 \\ \text{sens} = \frac{TP}{P} = \frac{0}{P} = 0 \\ \text{spec} = \frac{TN}{N} = \frac{N}{N} = 1 \end{array}$$

Good if sens & spec

both high

⇒ curve hng upper left corner

AUC = area under (ROC) curve

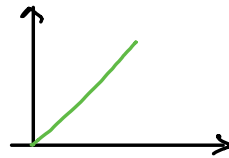
What if we just do "random guessing"?

* if 0.1 is used as cut-off, for each individual we draw $\text{bin}(1, 0.1) \rightarrow \frac{TP}{P} \approx 0.1, \frac{TN}{N} \approx 0.9$

* 0.2 $\text{bin}(1, 0.2) \rightarrow \frac{TP}{P} \approx 0.2, \frac{TN}{N} \approx 0.8$

This is similar to assign uniformly drawn p's [0,1] to each observation

\Rightarrow this will give a ROC curve with $AUC = 0.5$



This is often used for comparison. An $AUC \approx 0.5$ is thus not good.