

Binary predictions

- Challenge 5 (logistic regression in Python) released

Types of research questions

So far, we have learned how to answer the following questions:

- + what is the probability for observation i in the data?
- + What is the relationship between the explanatory variable(s) and the response? (fitting & interpreting model)
- + What is a "reasonable range" for a parameter in this relationship? (confidence intervals)
- + Do we have strong evidence for a relationship between these variables? (hypothesis testing)

What other kinds of research questions might we ask?

- . How well will I predict on new observations? /
How well do I predict the response?
- . what model should I use to predict the response? /
what variables are important?

Making predictions with the Titanic data

- + For each passenger, we calculate \hat{p}_i (estimated probability of survival)
- + But, we want to predict *which* passengers actually survive

How do we turn \hat{p}_i into a binary prediction of survival / no survival?

$$\hat{y}_i = \begin{cases} 1 & \hat{p}_i \geq 0.5 \leftarrow \text{threshold} \\ 0 & \hat{p}_i < 0.5 \end{cases}$$

Confusion matrix

		Actual	
		$Y = 0$	$Y = 1$
Predicted	$\hat{Y} = 0$	344	70
	$\hat{Y} = 1$	80	220

Annotations:

- Cell (344, $\hat{Y} = 0$, $Y = 0$) is circled and labeled "Correct (True negatives)".
- Cell (220, $\hat{Y} = 1$, $Y = 1$) is circled and labeled "Correct (True Positive)".
- Cell (70, $\hat{Y} = 0$, $Y = 1$) is labeled "Incorrect (False Negatives)".
- Cell (80, $\hat{Y} = 1$, $Y = 0$) is labeled "Incorrect (False positives)".

Did we do a good job predicting survival?

$$\text{Accuracy} = \frac{\# \text{ correct predictions}}{\# \text{ observations}} = \frac{TP + TN}{n} = \frac{220 + 344}{714} = 0.79$$

If I randomly select an observation, what is the probability my prediction is correct?

$1 - \text{Accuracy} = \text{classification error}$

Another confusion matrix (Dengue) :

Precision: $P(Y=1 | \hat{Y}=1)$
(PPV)

		$Y=0$	$Y=1$
		3957	1631
$\hat{Y}=0$	$\hat{Y}=0$	66	66
	$\hat{Y}=1$		

Accuracy:

$$\frac{3957 + 66}{5720} = 0.703$$

$\approx 70\%$ of the patients don't have Dengue

Problem: Accuracy is misleading with imbalanced data

How well did we do within group?

$$\underbrace{P(\hat{Y}=1 | Y=1)}_{\text{sensitivity} \text{ (aka recall, or TPR)}} = \frac{66}{1631 + 66} = \frac{TP}{TP+FN} = 0.039$$

$$\underbrace{P(\hat{Y}=0 | Y=0)}_{\text{specificity} \text{ (aka TNR)}} = \frac{3957}{3957 + 66} = \frac{TN}{TN+FP} = 0.984$$

Why a threshold of 0.5?

Consider data (x, y) $x \in \mathbb{R}^d$ and $y \in \{0, 1\}$

Fit a model to estimate $p(x) = P(Y=1 | X=x)$

We want binary predictions: $\hat{y} = \begin{cases} 1 & p(x) \geq h \\ 0 & p(x) < h \end{cases}$

Classification error: $P(Y \neq \hat{y})$
(1-accuracy)

Claim: $h = 0.5$ minimizes classification error
among all binary classifiers

Proof: Let $c(x)$ be an arbitrary classification function. $c(x) \in \{0, 1\}$, $\hat{y} = c(x)$

$$P(Y \neq \hat{Y}) = P(Y \neq C(X)) = E[1\{\hat{C}(X) \neq Y\}]$$

$$\mathbb{1}_A = \begin{cases} 1 & x \in A \\ 0 & x \notin A \end{cases}$$

$$P(X \in A) = \int f(x) dx =$$

$$\int_{-\infty}^{\infty} f(x) \cdot \mathbb{1}_A(x) dx = E[\mathbb{1}_A(X)]$$

$$E[1\{\hat{C}(X) \neq Y\}] = E\left[E\left[1\{\hat{C}(X) \neq Y\} | X\right]\right]$$

(iterated expectations)

$$= \begin{cases} E[1\{Y=0\}] & C(X)=1 \\ E[1\{Y=1\}] & C(X)=0 \end{cases}$$

$$= E[1\{Y=0\}] C(X) + E[1\{Y=1\}] (1 - C(X))$$

$$= P(Y=0) C(X) + P(Y=1) (1 - C(X))$$

$$= (1 - p(X)) C(X) + p(X) (1 - C(X))$$

$$\Rightarrow P(Y \neq \hat{Y}) = E[(1 - p(X)) C(X) + p(X) (1 - C(X))]$$

$$\Rightarrow P(Y \neq \hat{Y}) = \mathbb{E} \left[(1 - p(x)) c(x) + p(x) (1 - c(x)) \right]$$

$$= \int_x \left[(1 - p(x)) c(x) + p(x) (1 - c(x)) \right] f(x) dx$$

minimize $P(Y \neq \hat{Y})$: make $(1 - p(x)) c(x) + p(x) (1 - c(x))$ small!

$$(f \quad 1 - p(x) < p(x) : c(x) = 1$$

$$1 - p(x) > p(x) : c(x) = 0$$

$$\Rightarrow c(x) = \begin{cases} 1 & 1 - p(x) \leq p(x) \\ 0 & 1 - p(x) > p(x) \end{cases} = \begin{cases} 1 & p(x) \geq 0.5 \\ 0 & p(x) < 0.5 \end{cases}$$

Bayes classifier

\Rightarrow to minimize classification error, choose $h = 0.5$

if

Tradeoff: increase threshold \Rightarrow sensitivity: \downarrow
 specificity \uparrow

Changing the threshold

$$\text{Accuracy} = \text{sensitivity } P(Y=1) + \text{spec. } P(Y=0)$$

Using a threshold of 0.7:

		Actual	
		$Y = 0$	$Y = 1$
Predicted	$\hat{Y} = 0$	412	136
	$\hat{Y} = 1$	12	154

sensitivity: $\frac{154}{290}$

specificity: $\frac{412}{424}$

Using a threshold of 0.3:

		Actual	
		$Y = 0$	$Y = 1$
Predicted	$\hat{Y} = 0$	309	49
	$\hat{Y} = 1$	115	241

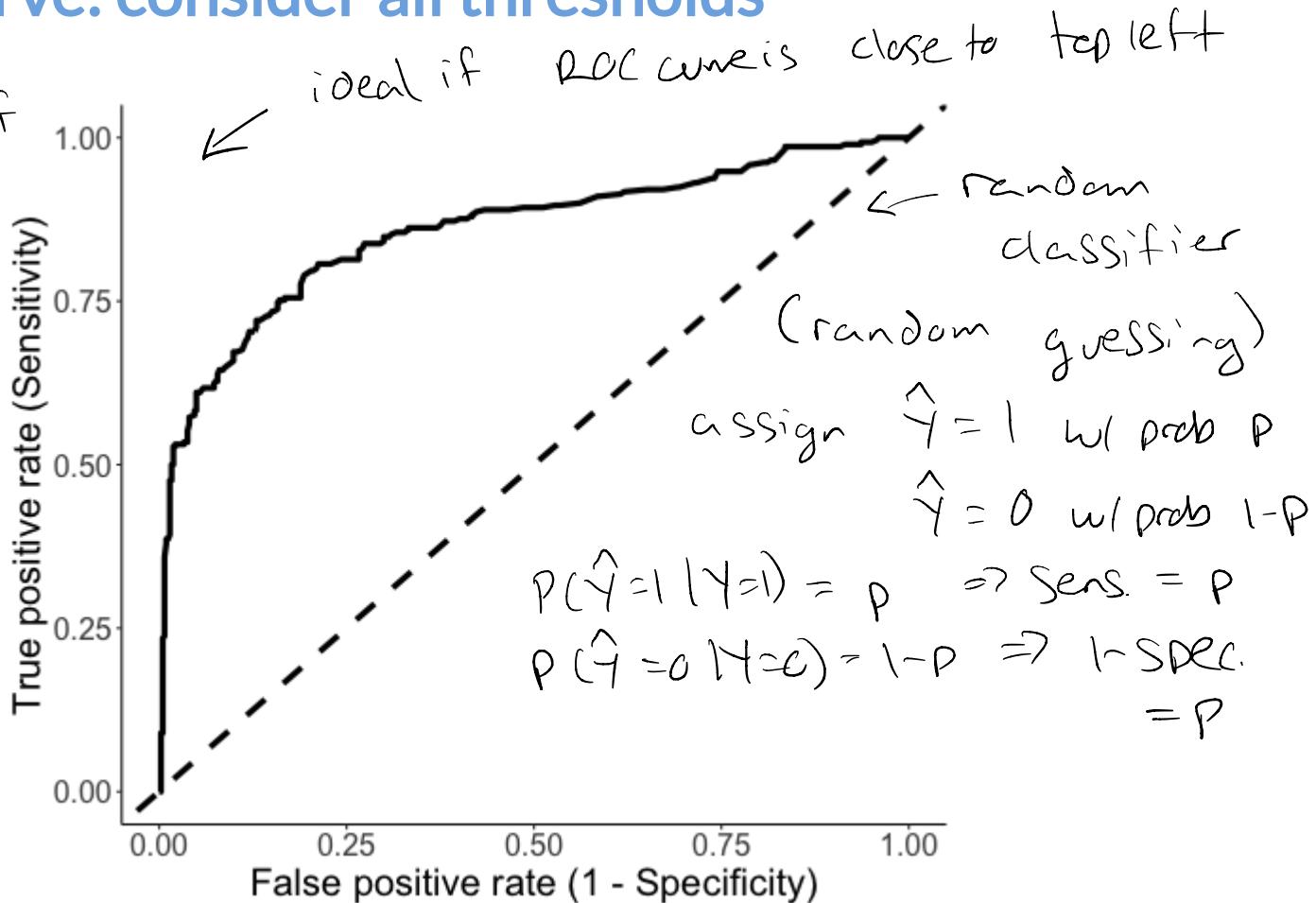
sensitivity: $\frac{241}{290}$

specificity: $\frac{309}{424}$

\nwarrow receiver operating characteristic

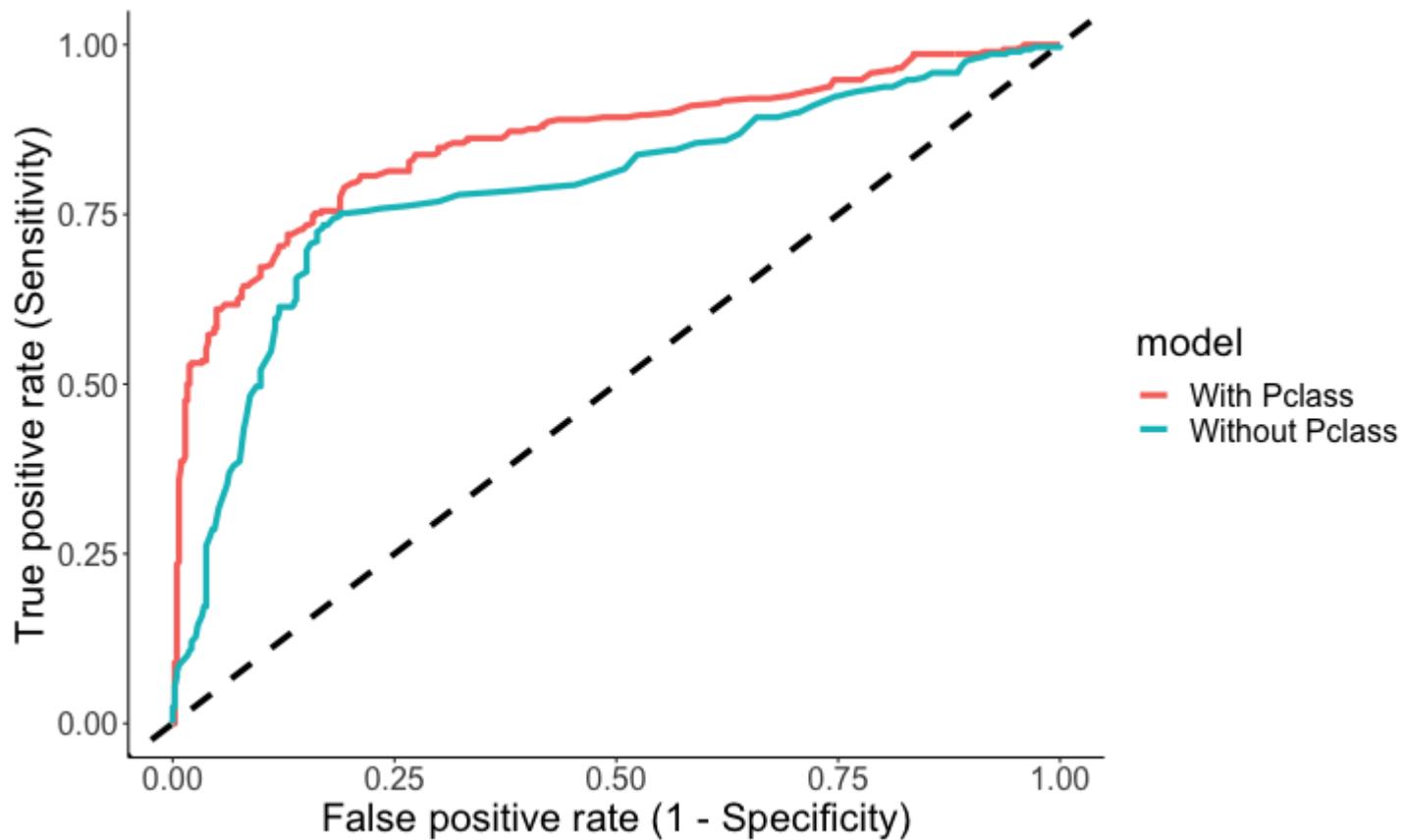
ROC curve: consider all thresholds

plots trade-off
between
sensitivity
specificity



Summarizing : Area under curve (AUC)
 $0.5 \leq \text{AUC} \leq 1$, closer to 1 is better

Comparing models with ROC curves



Problem: reusing data...

It is generally a bad idea to assess performance of a model on the same data we used to train it. This can lead to overfitting.

What can we do instead?