

Math 340 Lec 26 5/18/18

11

Logistic regression gives \hat{p} 's. How to evaluate?

"The" \hat{p} is $f_{pr}(\tilde{x})$ but you don't get to see it.

You only get y 's which are $\in \mathcal{Y} = \{0, 1\}$.

You need a scoring rule $S(\hat{p}, y)$.

A proper scoring rule is:

$$\forall i \quad f_{pr}(\tilde{x}_i) = \arg\max \{S(\hat{p}, y)\} \quad \text{i.e. maximal if you see the truth.}$$

Two proper rules:

- Log scoring rule:

$$\text{let } s_i = y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i)$$

$$S = \frac{1}{n} \sum_{i=1}^n s_i$$

- Brier score (1950)

$$\text{let } s_i = -(y_i - \hat{p}_i)^2, \quad S = \frac{1}{n} \sum s_i$$

These scores are used as an " R^2 ", a means to evaluate model but don't have an interpretation like "1% variance explained."

Using prob. classifiers to classify. Need a way to turn

$$\hat{p} \rightarrow \hat{y}. \quad \hat{p} = 90\% \Rightarrow \hat{y} = 1, \quad \hat{p} = 2\% \Rightarrow \hat{y} = 0$$

Naive:

$$\hat{y} = \mathbb{1}_{\hat{p} \geq 50\%} \quad \text{HARD threshold at 50\%}$$

What if ...

$$\hat{y} = \mathbb{1}_{\hat{p} \geq 90\%} \quad ? \quad \text{Makes it difficult to predict } \hat{y} = 1.$$

You want this if $CFP > CFN$!

this shifts confusion table ...

		\hat{y}	
		0	1
y	0		FP
	1	FN	
		FN	PP

which reduces FP at the expense of FN.

Generally ...

$$\hat{y} = \mathbb{1}_{\hat{p} \geq p_{th}} \quad \text{varying parameter based on asymmetric costs. Each value gives a different classifier. } p_{th} \in (0, 1)$$

Which one to choose?

Not relevant...

Look at all of them and choose!

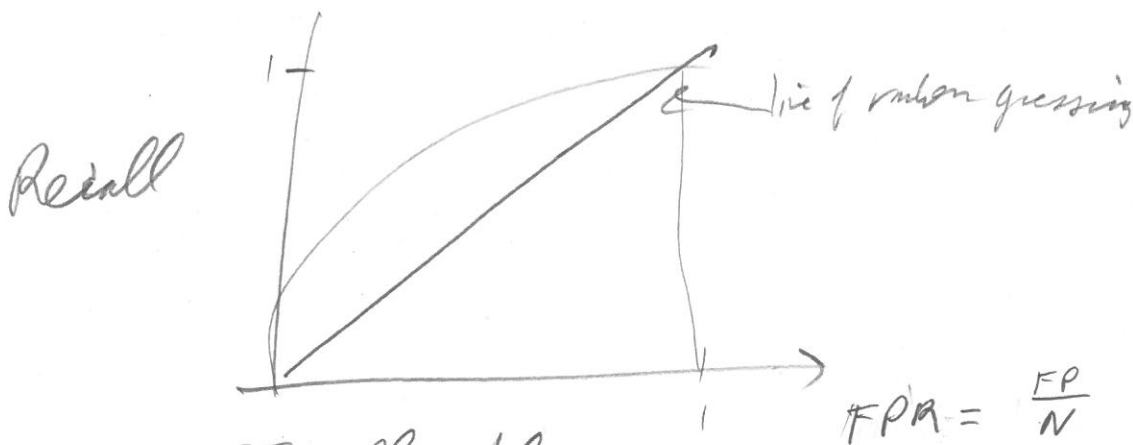
$$C_{FP}FP + C_{FN}FN$$

3

p_{cs}	TP	TN	FP	FN	precision	recall	FDR	FOR	etc etc
0.01									
0.02									
0.03									
...									
0.98									
0.99									

all can be derived

the following curve is used



The null model...

why? $\hat{y} = \mathbb{1}_{u \geq p_{cs}}$ where $u \sim U(0,1)$

AUC = 50% of 'null prob. classifier'

Ident: AUC = 1

AUC > 1 \Rightarrow better than chance

	\hat{y}	0	1	
0	$n(-p_{cs})$ $(1-p_{cs})$	$n(-p_{cs})$ $(1-p_{cs})$	$n(-p_{cs})$ $(1-p_{cs})$	$n(-p_{cs})$ $(1-p_{cs})$
1	$n(p_{cs})$ p_{cs}	$n(p_{cs})$ p_{cs}	$n(p_{cs})$ p_{cs}	$n(p_{cs})$ p_{cs}
	n	n	n	n

$$\text{recall} = \frac{n(p_{cs})}{n(p_{cs})} = 1$$

$$\text{FPR} = \frac{n(-p_{cs})}{n(-p_{cs})} = 1$$