

Lec 21 7/30/10 Math 390.3

The model  $y = t(\bar{x})$

But there's stuff you don't know, so it "appears" as:

$$Y \sim \text{Bern}(f_{\text{pr}}(\bar{x}))$$

i.e. Given the same  $\bar{x}$ , sometimes  $y=1$  and sometimes  $y=0$ .

So if we care about  $P(Y=1|x)$  our target is  $f_{\text{pr}}(\bar{x})$ .

But  $f_{\text{pr}}(\bar{x})$  can be arbitrarily complicated so we try to approximate it using a function in class  $\mathcal{H}$ . we pick

$$\mathcal{H}_{\text{pr}} = \left\{ \frac{e^{\vec{w} \cdot \bar{x}}}{1 + e^{\vec{w} \cdot \bar{x}}} : \vec{w} \in \mathbb{R}^{p+1} \right\}$$

Our A was to maximize likelihood assuming each  $y_i$  ind. and we found:

$$\vec{b} = \underset{\vec{w}}{\text{argmax}} \left\{ \sum_{i=1}^n \ln(1 + e^{-z_i \vec{w} \cdot \bar{x}_i}) \right\} \quad \text{where } z_i = 2y_i - 1$$

convenience

why does this make sense?

looks like  $\sum e_i$  error

$$\text{where } e_i = \ln(1 + e^{-z_i \vec{w} \cdot \bar{x}_i}) > 0$$

why?

$$\text{if } y_i = 1 \quad e_i = \ln(1 + e^{-\vec{w} \cdot \bar{x}_i})$$

$$\text{if } y_i = 0 \quad e_i = \ln(1 + e^{\vec{w} \cdot \bar{x}_i})$$

if  $\vec{w} \cdot \bar{x}_i$  is large  $\Rightarrow$  error small  
if ... neg & large  $\Rightarrow$  error small

$\frac{d}{d\vec{w}} [\sum e_i] \stackrel{??}{=} 0$  does not have any gradient sol.

~~Best no analytic solution exists...~~ must use numerical methods just like the SVM algorithm (which min. sum hinge loss)

How to use?

Okay ...  $\hat{p}_i = g(\vec{x}_i) = \left(1 + e^{-\vec{b} \cdot \vec{x}_i}\right)^{-1} = P(Y_i=1 | \vec{x}_i)$

↑  
prob estimate for row case,  $\vec{x}_i$

The linear term  $\vec{b} \cdot \vec{x}$  is buried inside. Does it have meaning?

$$\Rightarrow 1 - \hat{p} = \left(1 + e^{\vec{b} \cdot \vec{x}_i}\right)^{-1} = P(Y_i=0 | \vec{x}_i)$$

$$\Rightarrow \frac{\hat{p}}{1 - \hat{p}} = \frac{1 + e^{\vec{b} \cdot \vec{x}}}{1 + e^{-\vec{b} \cdot \vec{x}}} \cdot \frac{e^{\vec{b} \cdot \vec{x}}}{e^{\vec{b} \cdot \vec{x}}} = e^{\vec{b} \cdot \vec{x}} \frac{(1 + e^{\vec{b} \cdot \vec{x}})}{(1 + e^{-\vec{b} \cdot \vec{x}})}$$

<sup>$P(Y=1 | \vec{x})$</sup>   
 <sup>$(1 - P(Y=1 | \vec{x}))$</sup>   
Odds( $Y | \vec{x}$ )

Interpretation?  $\vec{b}$  means in  $X$ , because log odds of  $Y=1$  by  $\vec{b}_1$ !

$$\Rightarrow \ln\left(\frac{\hat{p}}{1 - \hat{p}}\right) = \vec{b} \cdot \vec{x} = b_0 + b_1 x_1 + \dots + b_p x_p$$

DEMO

"log-Odds( $Y | \vec{x}$ )" if you know how to read them, they're cool

very negative  $\Rightarrow$  low prob.

very positive  $\Rightarrow$  high prob.

near zero  $\Rightarrow$  50% prob.

Validating prob. models. What is the best prob. model you can come up with? fpr. So we should validate against

fpr. Can we? No! we have to validate  $\hat{p}$  vs.  $y$ !

This is very awkward! Enter the theory of "Scoring Functions" or "Scoring rules".

Let  $S(\hat{p}_i, y_i)$  be the scoring rule for obs.  $i$

A "proper scoring rule" has the following property

$$\forall i \quad \underset{\text{expected}}{f_{pr}(\vec{x}_i)} = \underset{\hat{p}}{\operatorname{argmax}} \{ S(\hat{p}_i, y_i) \}$$

the scoring rule is maximised if you use the true prob.

Popular scoring rules are the • log scoring rule:

$$S_i = y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i)$$

evaluate or  
 $\Rightarrow \text{avg. score} = \frac{1}{n} \sum_{i=1}^n S_i$

If  $y_i = 1$ , this is when  $\hat{p} \approx 1$

If  $y_i = 0$ , ... when  $\hat{p} \approx 0$

• Brier score (1950)

$$S_i = -(y_i - \hat{p}_i)^2$$

Does the avg. score mean anything in the problem? No... just as a means to validate and compare models.

To validate, you need to think in Brier scores e.g. Very opaque!

What's another way to validate? Use prob. est. to do

Classification! Then validate the classification!



Why not look at all of them or really all of them.

$p_{th}$	TP	TN	FP	FN	Precision	Recall	FDR	FOR	FPR
0.01									
0.02									
0.03									
$\vdots$									
0.99									

all can be derived

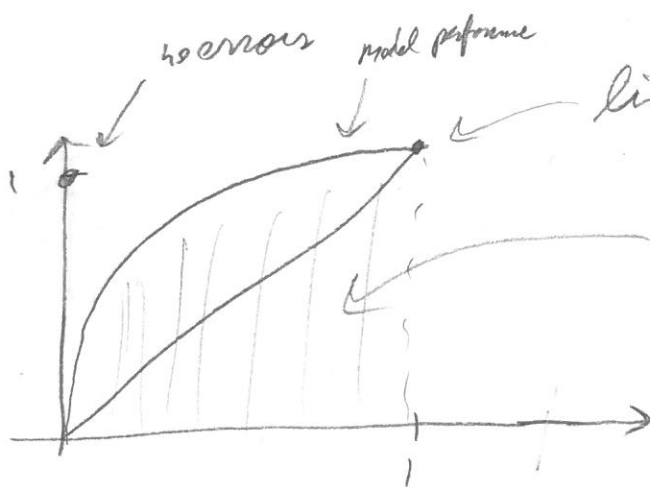
curve

Now let's plot a performance of all possible models

Recall-Opposing Characteristic Curve

(ROC)

Recall



line of random guessing

Area Under the Curve (AUC)

If  $> 0.5 \Rightarrow$  majority of models perform better than chance.

False Positive Rate (FPR)  $\leftarrow$  another metric to discuss accuracy

Before we go further... let's consider the case

If you use the null models, your classification is just prob.  $p_{th}$

become the null models can be  $Y \sim \text{Bern}(p_{th})$  where  $p_{th} \in (0, 1)$ .

If  $p_Y$  is  $P(Y=1)$  the marginal base rate, then the following is the conf table:

predicted (Y)

	0	1	
true (Y)	$(1-p_{th})(1-p_Y)$ TN	$p_{th}(1-p_Y)$ FP	$1-p_Y$ #N
	$(1-p_{th})p_Y$ FN	$p_{th}p_Y$ TP	$p_Y$ #P
	$1-p_{th}$	$p_{th}$	1

multiply by

$$= \frac{FP}{\#N}$$

$$FPR = \frac{p_{th}(1-p_Y)}{(1-p_Y)} = p_{th}$$

$$\text{recall} = \frac{p_{th}p_Y}{p_Y} = p_{th}$$