

Previously in Binary Classification, we presented error rate:

$$MAE = \frac{1}{n} \sum \mathbb{1}_{y_i \neq \hat{y}_i} \quad (\text{misclassification error/rate})$$

There are two types of errors:

\*  $y_i = 0$  &  $\hat{y}_i = 1$  (called false positive)

\*  $y_i = 1$  &  $\hat{y}_i = 0$  (called false omission).

How can you visualize these type errors?

Predict = ( $\hat{y}$ )

	0	1	
truth (y)	0	1	
0	TN	FP	#N
1	FN	TP	#P
	#PN	#PP	n

FP = False Positive

FN = False Negative

TN = True Negative

TP = True Positive

#PN : Number of predictive Neg

#PP : Number of predictive Pos

Called CONFUSION TABLE

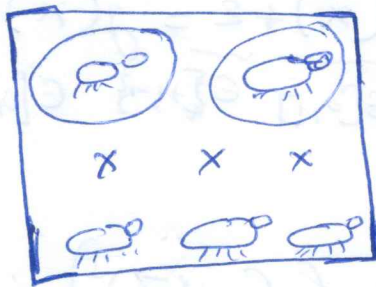
$$\text{Misclassification error} = \frac{FP + FN}{n}$$

$$\text{or Accuracy} = 1 - \frac{FP + FN}{n}$$

$$\text{or} = \frac{TP + TN}{n}$$

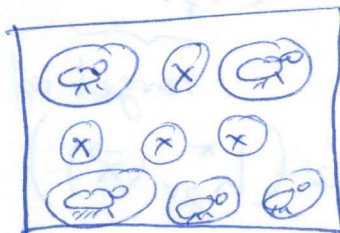
$$\text{Precision} = \frac{TP}{\#PP}$$

$$\text{Recall or Sensitivity} = \frac{TP}{\#P}$$



$$\text{precision} = \frac{2}{2} = 100\%$$

$$\text{recall} = \frac{2}{5} = 40\%$$



$$\text{precision} = \frac{5}{9} = 56\%$$

$$\text{recall} = \frac{5}{5} = 100\%$$

$$F_1 := \frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$$

$$F_1 \propto > 2\%$$

False discovery rate =  $1 - \text{precision} = \frac{FP}{\#PP}$   
(FDR)

False omission rate =  $\frac{FN}{\#PN}$

$$y = [0, 1]$$

TN	FP
FN	TP

## NEW UNIT:

$$y = \{0, 1\}$$

Targets is  $P(y=1/\vec{x}) \in [0, 1]$

Prob. Estimate

Probability Classification.

$$Y = t(\bar{x}) = f(\bar{x}) + \varepsilon = \underbrace{h^*(\bar{x})}_{\substack{\in [0,1] \\ \downarrow \\ \in \{+1,-1\}}} + \underbrace{\varepsilon}_{\substack{\in [0,1] \\ \in \{+1,-1\}}} = g(\bar{x}) + e \Rightarrow f, h^*, g, \text{ do not rescale prob's}$$

Consider:

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

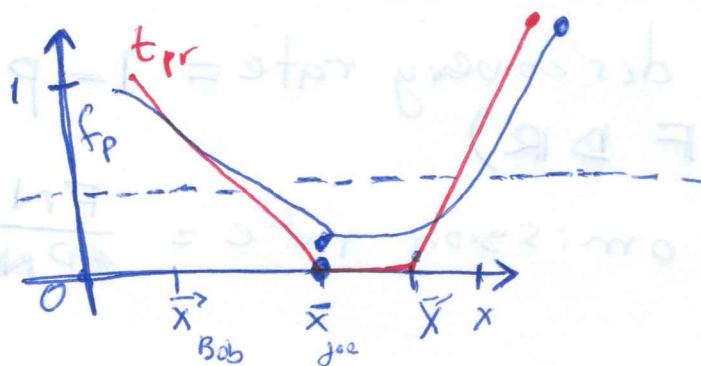
target of estimate

$$Y \sim \text{Bernoulli}(h_{pr}^*(\bar{x}))$$

$$Y \sim \text{Bernoulli}(g_{pr}(\bar{x}))$$

$$t(\bar{x}) = 1 \quad g_{pr}(\bar{x}) = 0.2$$

$$t_{pr}, f_{pr}, h_{pr}^*, g_{pr}$$



$$g_{pr} = \mathcal{A}(\mathcal{X}, \mathbb{D})$$

What is  $\mathcal{X}, \mathcal{A}$ ?

observe that:

$$Y_1 \sim \text{Bern}(f_{pr}(\bar{x}_1)) = f_{pr}(\bar{x}_1)^{Y_1} (1 - f_{pr}(\bar{x}_1))^{1-Y_1}$$

$$Y_2 \sim \text{Bern}(f_{pr}(\bar{x}_2))$$

$$\vdots$$

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

$$\text{What is } p(Y_1, Y_2, \dots, Y_n)?$$



Unknown ... unless you know / assume (42)

Big Assumpt<sup>n</sup>:  $Y_1, \dots, Y_n$  are Indpt.

How to check this assumpt<sup>n</sup>: Ask how  $\mathbb{D}$  was collected

► Indep<sup>t</sup>: If indpt

$$\underbrace{P(Y_1, Y_2, \dots, Y_n | x)}_{\text{likelihood}} = \prod_{i=1}^n P(Y_i | \vec{x}_i) = \prod_{i=1}^n f_p(\vec{x})^{y_i} (1 - f_p(\vec{x})^{1-y_i})$$

likelihood

We want to pick  $f_{pr}$  such that:

$P(Y_1, \dots, Y_n | x)$  is MAX.

But  $f_{pr}$  is arbitrary complicated so we construct our candidate set to  $\mathcal{H}$ .

How about ...  $\mathcal{H} = \{ \vec{w} \cdot \vec{x} : \vec{w} \in \mathbb{R}^{p+1} \}$ ?

Illegal since  $\vec{w} \cdot \vec{x} \in \mathbb{R}$  and  $\forall b(\vec{x}) \in (0, 1)$ .

How about  $\mathcal{H} = \{ \mathbb{1}_{\vec{w} \cdot \vec{x} \geq 0} : \vec{w} \in \mathbb{R}^{p+1} \}$ ? Illegal since

$\mathbb{1}_{\vec{w} \cdot \vec{x}} \in \{0, 1\}$  and we said we need probs!

But, ... I still like  $\vec{w} \cdot \vec{x}$ ! How can I use it?

We need a "like fun<sup>t</sup>",  $\phi$ , such that  $\phi(\vec{w} \cdot \vec{x}) \in (0, 1)$ .

Because we can't be too sure and say prob = 0 or 1.

$\phi: \mathbb{R} \rightarrow (0, 1)$  and is strictly increasing.

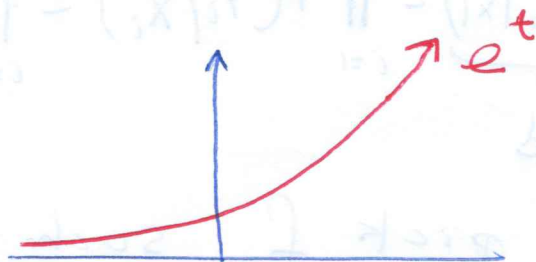
There are many ~~possibilities~~ possible  $\phi$ 's. Popular

Popular choices include:

a)  $\phi(x) = \frac{e^n}{1+e^n} = \frac{1}{1+e^{-n}}$   
(sigma)

b)  $\phi(x) = \Phi^{-1}(c_i)$   
(probit)

c)  $\phi(u) = 1 - e^{-u}$



This is called complementary log-log

d) Hyperbolic tangent:

$$\phi(u) = \tanh(u) = \left( \frac{e^u - e^{-u}}{e^u + e^{-u}} \right) \cdot \frac{1}{2}$$

▷ Logistic Regression

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

Called Generalized Linear Model (GLM)

$$\text{Let } \vec{b} = \arg \max \left\{ \prod \left( \frac{e^{\vec{w} \cdot \vec{x}_i}}{1 + e^{\vec{w} \cdot \vec{x}_i}} \right)^{y_i} \cdot \left( 1 - \frac{e^{\vec{w} \cdot \vec{x}_i}}{1 + e^{\vec{w} \cdot \vec{x}_i}} \right)^{1-y_i} \right\}$$

$\frac{1}{1+e^{-\vec{w} \cdot \vec{x}_i}} \quad \frac{1}{1+e^{\vec{w} \cdot \vec{x}_i}}$

$(1+e^{-\vec{w} \cdot \vec{x}})^{-1} \quad (1+e^{\vec{w} \cdot \vec{x}_i})^{-1}$

$$\Rightarrow \vec{b} = \operatorname{argmax} \left\{ \prod_{i=1}^n \left( (1 + e^{-\vec{w} \cdot \vec{x}_i})^{-1} \right)^{y_i} \left( (1 + e^{\vec{w} \cdot \vec{x}_i})^{-1} \right)^{1-y_i} \right\} \quad (43)$$

$$= \begin{cases} (1 + e^{-\vec{w} \cdot \vec{x}_i})^{-1} & \text{if } y_i = 1 \\ (1 + e^{\vec{w} \cdot \vec{x}_i})^{-1} & \text{if } y_i = 0 \end{cases}$$

$$= \prod_{i=1}^n \left( 1 + e^{(1-2y_i) \vec{w} \cdot \vec{x}_i} \right)^{-1}$$

$$= \prod_{i=1}^n \left( 1 + e^{-z_i \vec{w} \cdot \vec{x}_i} \right)^{-1}$$

$$\vec{b} = \left\{ -\sum \ln \left( 1 + e^{-z_i \vec{w} \cdot \vec{x}_i} \right) \right\}$$

$$= \operatorname{argmin}_{\vec{w}} \left\{ \sum \ln \left( 1 + e^{-z_i \vec{w} \cdot \vec{x}_i} \right) \right\}$$

Take  $\frac{d}{d\vec{x}} \left[ \right] \stackrel{\text{Set}}{=} 0$  but no closed form analytical solution

Need: numerical optimization:

Typically "Gradient"