

# logistic regression

used for classification

training data :  $\{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$

▢ feature vector  $\vec{x}_i \in \mathbb{R}^d \leftarrow$  how many features

▢ label  $y_i \in \{0, 1\}$

in logistic regression, we learn from the data a probabilistic model for

$$\Pr(Y=y | \vec{x}) \quad \text{for } Y=0 \text{ and } Y=1.$$

label ← feature

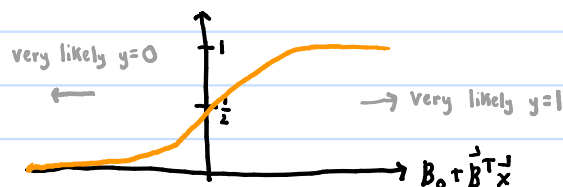
given a data pt. w/ feature vector  $\vec{x}$ , what is the probability that its label is  $y=1$ ?  $\Pr(Y=1 | \vec{x})$

the logistic model: a linear model for the log odds

$$\log\left(\frac{\Pr(Y=1 | \vec{x})}{1 - \Pr(Y=1 | \vec{x})}\right) = \beta_0 + \vec{\beta}^T \vec{x}$$

odds that  $Y=1$  given  $\vec{x}$       intercept  $\beta_0 \in \mathbb{R}$       coefficients  $\vec{\beta} \in \mathbb{R}^d$

$$\Pr(Y=1 | \vec{x}) = \frac{e^{\beta_0 + \vec{\beta}^T \vec{x}}}{1 + e^{\beta_0 + \vec{\beta}^T \vec{x}}} \in (0, 1) \quad \text{"}$$



$$\text{and } \Pr(Y=0 | \vec{x}) = 1 - \Pr(Y=1 | \vec{x}).$$

how to learn  $\beta_0, \vec{\beta}$  from data?

under this probabilistic model, write the likelihood, the probability of seeing the data given the probabilistic model & its parameters, then choose the parameters that maximize the likelihood.

$$\mathcal{L}(\beta_0, \vec{\beta}) = \prod_{i=1}^n \Pr(Y=y_i | \vec{x}_i) = \prod_{i=1}^n \Pr(Y=1 | \vec{x}_i)^{y_i} [1 - \Pr(Y=1 | \vec{x}_i)]^{1-y_i}$$

↑  
likelihood

assumes independent,  
identically distributed data

$$\begin{aligned} \log \mathcal{L} &= \sum_{i=1}^n y_i \log[\Pr(Y=1 | \vec{x}_i)] + (1-y_i) \log[1 - \Pr(Y=1 | \vec{x}_i)] \\ &= \sum_{i=1}^n y_i \log\left[\frac{e^{\beta_0 + \vec{\beta}^T \vec{x}_i}}{1 + e^{\beta_0 + \vec{\beta}^T \vec{x}_i}}\right] + (1-y_i) \log\left[\frac{1}{1 + e^{\beta_0 + \vec{\beta}^T \vec{x}_i}}\right] \\ &= \sum_{i=1}^n y_i (\beta_0 + \vec{\beta}^T \vec{x}_i) - \log(1 + e^{\beta_0 + \vec{\beta}^T \vec{x}_i}) \end{aligned}$$

maximizing  $\log \mathcal{L}$  gives same  $\beta_0, \vec{\beta}$  as maximizing  $\mathcal{L}$ .  
 $\log \mathcal{L}$  is monotonic...

$$\rightarrow \vec{\nabla}_{\vec{\beta}} \log \mathcal{L} = \vec{0} \quad \text{to find } \vec{\beta} \text{ that fits the data ("learn", "train")}$$

$$\sum_{i=1}^n \vec{x}_i (y_i - \Pr(Y=1 | \vec{x}_i)) = \vec{0}$$

... non-linear in  $\vec{\beta}$

one option to solve for  $\vec{\beta}$ : the Newton-Raphson algorithm.

## confusion matrix

once a cutoff  $P(Y=1|\vec{x})$  is decided upon for a decision boundary, the confusion matrix summarizes classification performance.

TRUTH

PREDICTION		TRUTH	
		$y=1$	$y=0$
$y=1$	TP	FN	
$y=0$	FP	TN	

TP : true positive

# data whose true and predicted label is  $y=1$

TN : true negative

# data whose true and predicted label is  $y=0$

FP : false positive

# data whose predicted label is  $y=1$ , but true label is  $y=0$

FN : false negative

# data whose predicted label is  $y=0$ , but true label is  $y=1$

## receiver operator characteristic (ROC) curve

LR gives us  $\Pr(Y=1|\vec{x}) \in (0,1)$ .

to classify (map  $\vec{x}$  to 0 or 1), we must choose a threshold  $p^*$  for the classification rule:

$$y(\vec{x}) = \begin{cases} 0 & \text{if } \Pr(Y=1|\vec{x}) \leq p^* \\ 1 & \text{if } \Pr(Y=1|\vec{x}) > p^* \end{cases}$$

choice of  $p^*$  should reflect the costs of FP/FN's

$p^* \rightarrow 0 \Rightarrow$  fewer FN, more FP

$p^* \rightarrow 1 \Rightarrow$  fewer FP, more FN

in the limit, all data is labeled by model as  $y=1$   
 $y=0$

the ROC curve scans all possible thresholds,  $p^*$ , and plots TPR vs. FPR

TPR: true positive rate  $\equiv$  recall  $\equiv$  sensitivity

$$= \frac{TP}{P} = \frac{TP}{TP+FN}$$

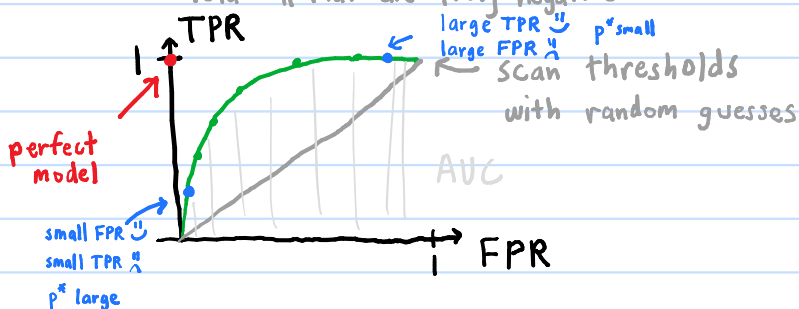
total # that are truly positive

FPR: false positive rate  $\equiv$  fall out  $\equiv$  1-selectivity

$$= \frac{FP}{N} = \frac{FP}{TN+FP}$$

total # that are truly negative

commonly used to  
evaluate LR model w/o  
imposing a  $p^*$



## area under ROC curve (AUC):

- 1 randomly select a positive instance  $\vec{x}_+$
- 2 randomly select a negative instance  $\vec{x}_-$
- 3  $AUC = \Pr[\Pr(Y=1|\vec{x}_+) > \Pr(Y=1|\vec{x}_-)]$

i.e. AUC is probability LR properly  
ranks the two instances.