

TAM 598 Lecture 16 :

## Classification

---

Announcements:

- Hw 4 covers lectures 13-16; due on Fri Apr 4

↑  
updated!

## I. Classification / logistic regression

observations  $\underline{x}_{1:n} = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n)$

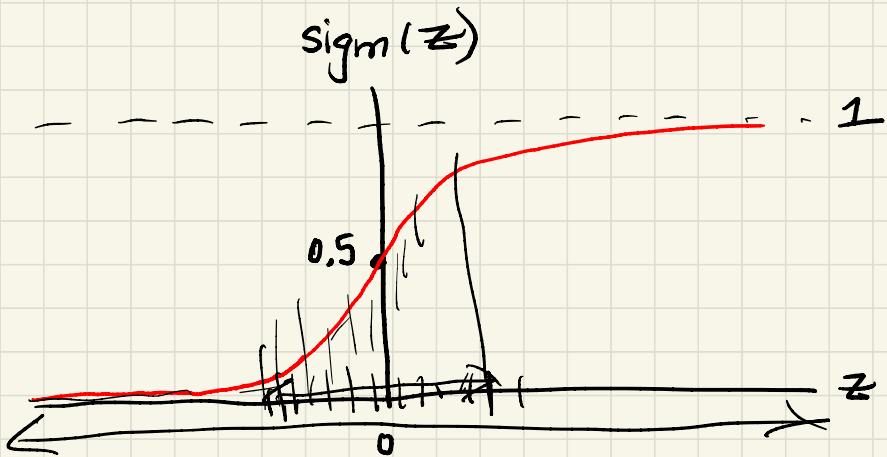
targets  $\underline{y}_{1:n} = (y_1, \dots, y_n) \quad \leftarrow \text{discrete labels}$

consider binary classification, where  $y=0$  or  $y=1$ .  
probability that  $y=1$  conditioned on  $x$ :

$$P(y=1 \mid \underline{x}, \underline{w}) = \text{sigm}\left(\sum_{j=1}^m w_j \phi_j(\underline{x})\right) = \text{sigm}\left(\underline{w}^T \underline{\phi}(\underline{x})\right)$$

where:  $w_j$ 's are weights,  $\phi_j(x)$  are basis functions, and

"sigm" is the sigmoid function  $\text{sigm}(z) = \frac{1}{1 + e^{-z}}$



logistic regression → a generalized linear model pushed through a sigmoid function so it is mapped to  $[0, 1]$

$$\text{then } p(y=1) = \text{sigm}(\underline{w}^T \underline{\phi}(x))$$

$$p(y=0) = 1 - \text{sigm}(\underline{w}^T \underline{\phi}(x))$$

or

$$\underline{p(y|x, w)} = [\text{sigm}(\underline{w}^T \underline{\phi})]^y [1 - \text{sigm}(\underline{w}^T \underline{\phi})]^{1-y}$$

## II. Likelihood of all observed data:

$$P(y_{1:n} | \underline{x}_{1:n}, \underline{w}) = \prod_{i=1}^n P(y_i | \underline{x}_i, \underline{w})$$

$$= \prod_{i=1}^n [\text{sigm}(\underline{w}^T \underline{\phi}(\underline{x}_i))]^{y_i} [1 - \text{sigm}(\underline{w}^T \underline{\phi}(\underline{x}_i))]^{1-y_i}$$

Then we obtain the best weight vector  $\underline{w}$  using MLE :

$$\underline{w}^* = \max_{\underline{w}} \log P(y_{1:n} | \underline{x}_{1:n}, \underline{w})$$

$$= \max_{\underline{w}} \sum_{i=1}^n \left\{ y_i \text{sigm}(\underline{w}^T \underline{\phi}(\underline{x}_i)) + (1-y_i) [1 - \text{sigm}(\underline{w}^T \underline{\phi}(\underline{x}_i))] \right\}$$

Finding  $\max \underline{w}$  is equivalent to minimizing the loss

$$L(\underline{w}) = -\sum_{i=1}^n \{y_i \operatorname{sigm}(\underline{w}^\top \underline{\phi}(x_i)) + (1-y_i) [1 - \operatorname{sigm}(\underline{w}^\top \underline{\phi}(x_i))] \}$$

this is called a **cross-entropy loss function**.  
used also, e.g., for DNNs that classify images.

III Making Decisions : say we have a point estimate  $\underline{w}^*$

↳  $p(y|x, \underline{w}=\underline{w}^*)$  is the prob. that  $y=1$

↳ but you need to make a decision,  $y=0$  or  $y=1$

As before with decision making, need a loss function

$l(\hat{y}, y) = \text{cost of picking } \hat{y} \text{ if the true value is } y$

choice of cost function is subjective.

For binary classification,  $l(\hat{y}, y)$  is a  $2 \times 2$  matrix

		$y$ (true)	
		0	1
$\hat{y}$	0	0	1, 0
	1	1, 0	0

Given  $l(\hat{y}, y)$ , your decision is the  $\hat{y}$  that minimizes expected loss

=

decision is

$$\min_{\hat{y}}$$

$$\sum_{y=0,1} l(\hat{y}, y) p(y | \underline{x}, \underline{w} = \underline{w}^*)$$

a different optimization problem  
for each  $\underline{x}$

## IV. Diagnostics for Classification

- split dataset into training & validation
- two important metrics are

balanced accuracy:

average of  
recall (TP rate) 2) confusion matrix

for each class:

$$\left. \begin{array}{l} \frac{TP}{TP+FN} \\ \frac{TN}{TN+FP} \end{array} \right\} \text{avg}$$

		predicted	
		0	1
actual	0	TP	FN
	1	FP	TN

$$\text{accuracy} = \frac{1}{n_{\text{samples}}} \sum_{i=1}^{n_{\text{samples}}} \mathbb{1} (\hat{y}_i = y)$$

would like the numbers on the diagonal to be large

## V. Multi-class classification

observations  $\underline{x}_{1:n} = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n)$

targets  $\underline{y}_{1:n} = (\underline{y}_1, \dots, \underline{y}_n)$  ← discrete labels

Now we have K possible values for labels so