

# 8/31/2023: Classification

Goal: Predict "label" or "class" for each input from a discrete set of options

- Tumor: Benign or malignant?
- Email: Spam or not spam?
- Handwritten digits: 0, 1, 2, ..., 9
- Image: Bird, snake, dog, ...

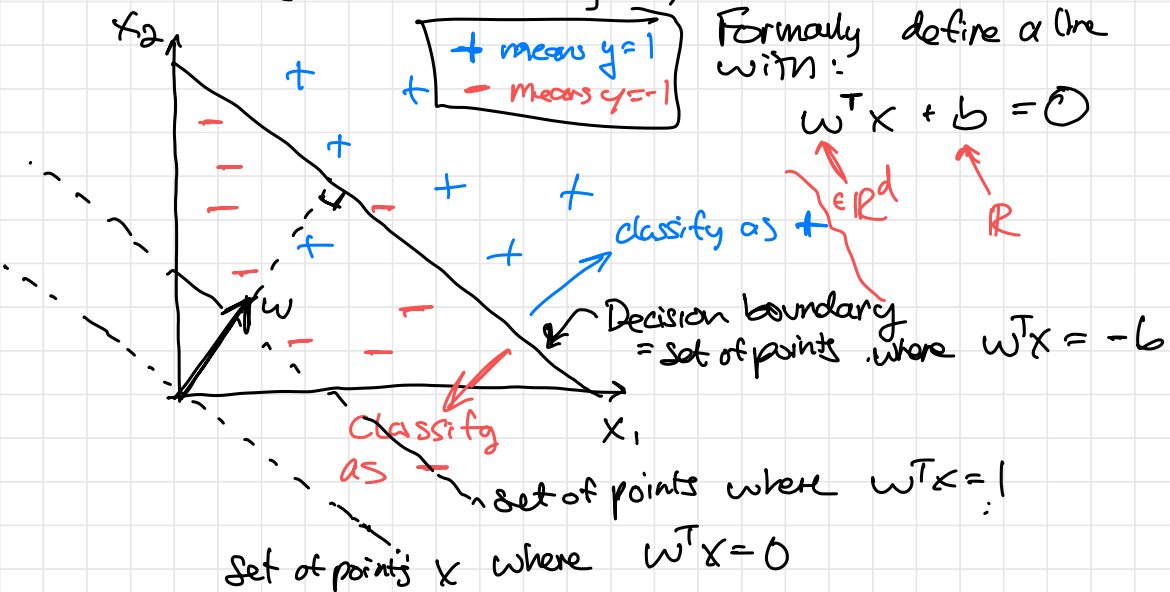
Binary classification  
2 possible labels

"10-way classification problem"

multi-class classification  
> 2 possible labels

## Binary classification

- One label is  $y = 1$  "positive"
  - Other label is  $y = -1$  "negative"
- (can also use  $y = 0$ )



Machine learning goal: Learn  $w \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$  such that the decision boundary correctly separates '+'s and '-'s

Model predictions:

If  $w^T x + b > 0$  predict  $y = 1$   
If  $w^T x + b < 0$  predict  $y = -1$

## Maximum Likelihood Estimation

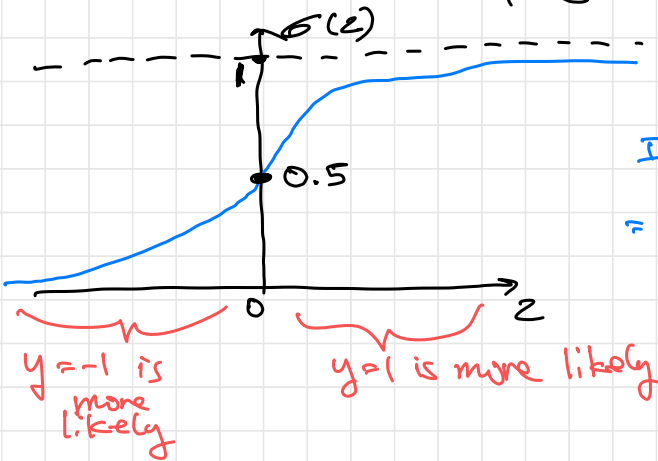
Same idea as linear regression,  
just need a different probabilistic story

$$P(y=1 | x; w) = \frac{1}{1 + \exp(-w^T x)} = \sigma(w^T x)$$

omit  $b$  for same reason as before

"sigmoid" or "logistic" function

where  $\sigma(z) = \frac{1}{1 + e^{-z}}$



Convenient fact:

$$P(y | x; w) = \sigma(y w^T x)$$

If  $y = 1$ , clearly true

If  $y = -1$ ,  $\sigma(-w^T x)$

$$= \frac{1}{1 + e^{w^T x}} = \frac{e^{-w^T x}}{e^{-w^T x} + 1}$$
$$= 1 - \frac{1}{1 + e^{-w^T x}}$$

maximize log-likelihood:

$$\log \prod_{i=1}^n P(y^{(i)} | x^{(i)}; w)$$
$$= \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}; w)$$
$$= \sum_{i=1}^n \log \sigma(y^{(i)} w^T x^{(i)})$$

Last step: multiply by  $-\frac{1}{n}$ , swap maximization to minimization

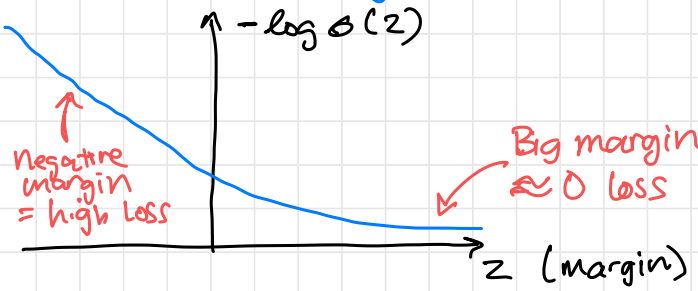
Final Loss function:  $L(w) = \frac{1}{n} \sum_{i=1}^n -\log \sigma(y^{(i)} w^T x^{(i)})$

Logistic Regression

"margin"  
large margin is good

$-\log \sigma(z)$  is a function that measures "badness" of our margin

If  $y^{(i)} = 1$ , want  $w^T x^{(i)} > 0$   
If  $y^{(i)} = -1$ , want  $w^T x^{(i)} < 0$   
margin  $> 0 \iff$  prediction is correct



Minimize  $L(w)$  with gradient descent  
- Fact: this  $L(w)$  is also convex

Gradient  $\nabla_w \frac{1}{n} \sum_{i=1}^n -\log \sigma(y^{(i)} w^T x^{(i)})$   
 $= \frac{1}{n} \sum_{i=1}^n -\sigma(-y^{(i)} w^T x^{(i)}) \cdot \underbrace{y^{(i)} x^{(i)}}_{\text{"constants"}}$

Fact:  $\frac{d}{dz} -\log \sigma(z) = -\sigma(-z)$

$= \frac{1}{n} \sum_{i=1}^n \underbrace{-\sigma(-y^{(i)} w^T x^{(i)})}_{\text{positive \#}} \cdot \underbrace{y^{(i)}}_{\pm 1} \underbrace{x^{(i)}}_{\text{vector}}$

If  $y^{(i)} = 1$ : gradient has [negative #]  $\cdot x^{(i)}$   
add multiple of  $x^{(i)}$  to  $w$   
increase  $w^T x^{(i)}$ , increase  $p(y=1 | x; w)$

If  $y^{(i)} = -1$ : gradient has (positive #)  $\cdot x^{(i)}$   
 subtract multiple of  $x^{(i)}$  from  $w$   
 decreases  $w^T x^{(i)}$ , increase  $p(y^{(i)} = -1 | x^{(i)}, w)$

What about  $\sigma(-y^{(i)} w^T x^{(i)})$ ?  
 $= \sigma(-\text{margin})$

If margin large:  $\approx 0$

If margin small:  $\approx 1$

we're already doing well on this example  $\Rightarrow$  don't need to update

we're getting this example wrong need to update  $w$

## Softmax Regression

AKA "Multinomial Logistic Regression"

Similar to logistic regression but for  $\geq 2$  classes

We now have  $C$  classes,  $x^{(i)} \in \mathbb{R}^d$

Model will have  $C \times d$  parameters  $w^{(1)}, \dots, w^{(C)} \in \mathbb{R}^d$

$w^{(j)T} x$  measure how much  $x$  looks like class  $j$

"dog"

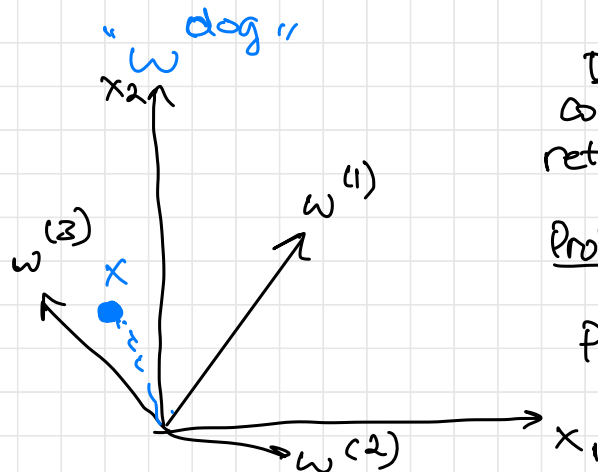
"dog"

### Decision Rule:

compute  $w^{(1)T} x, \dots, w^{(C)T} x$   
 return  $j$  with largest  $w^{(j)T} x$

### Probabilistic Story:

$$p(y=j | x; w) = \frac{\exp(w^{(j)T} x)}{\sum_{k=1}^C \exp(w^{(k)T} x)}$$



$$\begin{aligned}
 w^{(1)T} x &= 1 \xrightarrow{\text{exp}} \approx 2.7 \\
 w^{(2)T} x &= -3 \rightarrow \approx 0.1 \\
 w^{(3)T} x &= 2 \rightarrow \approx 7.4
 \end{aligned}
 \left. \vphantom{\begin{aligned} w^{(1)T} x \\ w^{(2)T} x \\ w^{(3)T} x \end{aligned}} \right\} \text{Normalize} \begin{aligned} p(y=1|x;w) &\approx 0.27 \\ p(y=2|x;w) &\approx 0.01 \\ p(y=3|x;w) &\approx 0.72 \end{aligned}$$

Predict  $y=3$  10.2 1.0

## Maximum Likelihood Estimation:

minimize negative log-likelihood ( $\times \frac{1}{n}$ )

$$L(w) = \frac{1}{n} \sum_{i=1}^n -\log p(y^{(i)} | x^{(i)}; w)$$

$$= -\frac{1}{n} \sum_{i=1}^n w^{(y^{(i)})T} x^{(i)} = -\log \sum_{k=1}^c \exp(w^{(k)T} x^{(i)})$$

(Finishing on 9/5) Now: apply gradient descent w.r.t.  $w^{(j)}$

$$\nabla_{w^{(j)}} L(w) = -\frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{I}[y^{(i)}=j]}_{\substack{\text{1 if } y^{(i)}=j \\ 0 \text{ else}}} x^{(i)}$$

$$= \frac{1}{\sum_{k=1}^c \exp(w^{(k)T} x^{(i)})} \cdot \exp(w^{(j)T} x^{(i)}) \cdot x^{(i)}$$

$= p(y^{(i)}=j | x^{(i)}; w)$

$$= \frac{1}{n} \sum_{i=1}^n \left( p(y^{(i)}=j | x^{(i)}; w) - \mathbb{I}[y^{(i)}=j] \right) \cdot x^{(i)}$$

If  $y^{(i)} \neq j$ :  $\nabla$  has positive  $\cdot x^{(i)} \Rightarrow$  GD subtracts multiple of  $x^{(i)}$  from  $w^{(j)}$   
 $\Rightarrow$  makes  $w^{(j)T} x^{(i)}$  smaller  
 $\rightarrow$  reduces  $p(y^{(i)}=j | x^{(i)}; w)$

If  $y^{(i)} = j$ :  $\nabla$  has negative  $\cdot x^{(i)} \Rightarrow$  GD adds multiple of  $x^{(i)}$  to  $w^{(j)}$