

## Logistic regression

- Logistic model for distributions over  $\mathbb{R}^d \times \{+1, -1\}$  (or  $\mathbb{R}^d \times \{0, 1\}$ ).

Definition: We say that a distribution  $D$  over  $\mathbb{R}^d \times \{+1, -1\}$  satisfies the logistic model if there exists  $w \in \mathbb{R}^d$  such that if  $(x, y) \sim D$  then

$$\Pr[Y = y \mid X = x] = \frac{1}{1 + e^{-y w^T x}}$$

for any  $(x, y) \in \mathbb{R}^d \times \{+1, -1\}$ .

---

Note: Logistic model does not place any assumptions on the

place any assumptions  
marginal distribution of  $X$ :

$$\Pr[X=x, Y=y] = \Pr[X=x] \cdot \Pr[Y=y | X=x]$$

↑  
this can be anything

↑  
this has to satisfy logistic model assumption

---

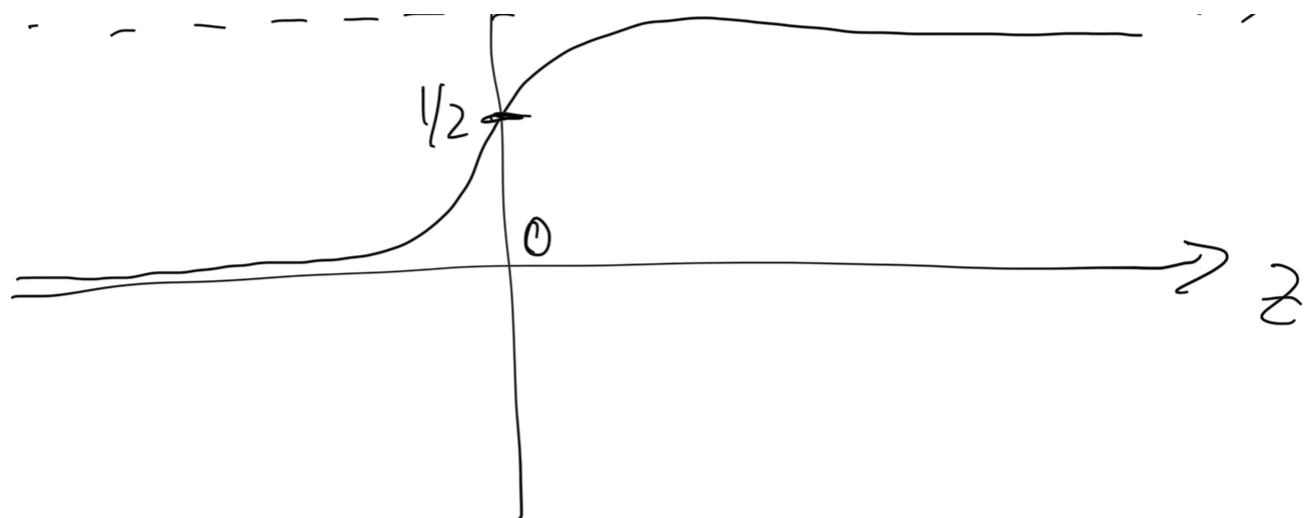
Why logistic regression?

1) Define  $\sigma: \mathbb{R} \rightarrow (0, 1)$

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

↑

$\sigma(z)$



$\sigma(z)$  is called

a) the sigmoid function

because of its S-shape

b) logistic function

(Etymology of the word  
"logistic" is bit fuzzy.)

$$\Pr[Y=y | X=x] = \sigma(yw^T x)$$

or

$$\Pr[Y=1 \mid X=x] = \sigma(w^T x)$$



Generalized linear  
model

2) "Odds" of probability event  
A are

$$\text{"Pr}[A] : \text{Pr}[A^c]\text{"}$$

• Mathematically, odds are  $\frac{\text{Pr}[A]}{\text{Pr}[A^c]}$

• The conditional odds for

logistic regression are

$$\frac{\Pr[Y=1 | X=x]}{\Pr[Y=-1 | X=x]}$$

$$\stackrel{||}{=} \frac{\sigma(w^T x)}{\sigma(-w^T x)}$$

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

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

$$= \frac{e^{w^T x} + (e^{w^T x})^2}{e^{w^T x} + 1}$$

$$= e^{w^T x}$$

- Log odds (logarithm of odds) are a linear function of  $x$ :

$$\ln \left( \frac{\Pr[Y=1 | X=x]}{\Pr[Y=-1 | X=x]} \right) = w^T x$$

- Logistic function  $\sigma: \mathbb{R} \rightarrow (0,1)$  maps log-odds to probabilities:

$$\sigma(w^T x) = \sigma \left( \ln \left( \frac{\Pr[Y=1|X=x]}{\Pr[Y=-1|X=x]} \right) \right)$$

$$= \Pr[Y=1|X=x]$$

Probability

log odds

---

## Training logistic regression model

- Let  $S = ((x_1, y_1), \dots, (x_m, y_m))$  be a labeled data set where  $x_1, x_2, \dots, x_m \in \mathbb{R}$  and  $y_1, y_2, \dots, y_m \in \{+1, -1\}$ .
- We want find model

parameters  $w \in \mathbb{R}^d$  that  
"best fit the data."

• This is (usually) done  
with maximum likelihood  
principle.

• Likelihood function  $L: \mathbb{R}^d \rightarrow \mathbb{R}$

$$L(w) = \Pr_w [Y_1, Y_2, \dots, Y_m \mid X_1, \dots, X_m]$$

where subscript  $w$  indicates  
logistic model with parameter  
 $w$ .

• Likelihood  $\neq$  Probability



1) Likelihood is a function of  $w$

2) Probability is a measure i.e. a function of events.

• Suppose  $(X_1, Y_1) \dots (X_m, Y_m)$  is an i.i.d. sample. Then,

$$L(w) = \prod_{i=1}^m \Pr_w [Y_i | X_i]$$

$$= \prod_{i=1}^m \sigma(-Y_i w^T X_i)$$

by independence

• We want to find

$$\hat{w} = \operatorname{argmax}_{w \in \mathbb{R}^d} L(w)$$

• This is equivalent to

$$\hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \underbrace{-\ln(L(w))}_{\text{Negative log-likelihood}}$$

$$-\ln(L(w)) = -\ln \prod_{i=1}^m \sigma(y_i w^T x_i)$$

$$= -\sum_{i=1}^m \ln \sigma(y_i w^T x_i)$$

$$= -\sum_{i=1}^m \ln \frac{1}{1 + e^{-y_i w^T x_i}}$$

$$\sum_{i=1}^m (1 + e^{-\gamma_i w^T x_i})$$

$$= \sum_{i=1}^m \ln (1 + e^{-\gamma_i w^T x_i})$$

• We define log-loss or logistic loss

$$l(w, x, \gamma) = \ln (1 + e^{-\gamma w^T x})$$

• There is no "closed form" expression for

$$\hat{w} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^m \ln (1 + e^{-\gamma_i w^T x_i})$$

... L L

- The minimum needs to be computed by "numerical" methods.

- One can show that the function

$$f: \mathbb{R}^d \rightarrow \mathbb{R},$$

$$f(w) = \sum_{i=1}^m \ln(1 + e^{-\gamma_i w^T x_i})$$

is convex.

- Convexity makes finding minimum of  $f(w)$  easy to find.