

Lecture 21

- why does the cross entropy loss make sense?

Imagine this scenario

$$\text{coin, } P_H = ? = \theta$$

$$P_T = 1 - P_H = 1 - ? = 1 - \theta$$

Sample size = 10

{0, 0, 0, 0, 0, 0, 1, 1, 1, 1}

$$\begin{aligned} \hookrightarrow n &= 10 \\ n_H &= 6 \\ n_T &= 4 \end{aligned}$$

- Our best guess is $\hat{\theta} = 0.6$



likelihood

Probability that I see this result given an estimate for θ

$$P(Y_{1-6} = 0, Y_{7-10} = 1 \mid \theta = 0.6) \quad \text{High}$$

$$P(Y_{1-6} = 0, Y_{7-10} = 1 \mid \theta = 0.3) \quad \text{Low}$$

Ok... Cool but what does this
have to do w/ logistic regression.



Our data is a sample; the coin
is the world model we try to
predict.

$$f_{\theta}(x) = \sigma(x^T \theta) = p(Y=1 | x) = p$$

Our goal maximize likelihood
for $\sigma(x^T \theta) = f_{\theta}(x) = p$

Back to the coin

$$p(Y_{1-6} = 0, Y_{7-10} = 1 | \theta = p)$$

o
o

$$\text{Bernoulli}(p) \\ p^y (1-p)^{(1-y)}$$

$$\prod_{i=1}^{10} p^{y_i} (1-p)^{(1-y_i)}$$



generalize

for n samples of data

and our model $f_{\theta}(x) = \sigma(x^T \theta)$

$$\prod_{i=1}^n p^{y_i} (1-p)^{(1-y_i)}$$

" p "

y_i : observed data

p : predicted probability,

as output by

$$\sigma(x^T \theta)$$



$$\prod_{i=1}^n (\sigma(x_i^T \theta))^{y_i} (1 - \sigma(x_i^T \theta))^{(1-y_i)}$$

our goal for $\hat{\theta}$:

find best $\hat{\theta}$ parameters that

maximize the likelihood of observing

our data samples $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n \sigma(x_i^T \theta)^{y_i} (1 - \sigma(x_i^T \theta))^{(1-y_i)}$$

Trivial

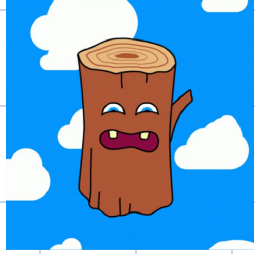
log is monotonically increasing;

therefore

$$\underset{\theta}{\operatorname{argmax}} f(x) = \underset{\theta}{\operatorname{argmax}} \log(f(x))$$

iff $f(x) > 0$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log \left(\prod_{i=1}^n \sigma(x^T \theta)^{y_i} (1 - \sigma(x^T \theta))^{(1-y_i)} \right)$$



$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log(\sigma(x^T \theta)^{y_i}) + \log((1 - \sigma(x^T \theta))^{(1-y_i)})$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n y_i \log(\sigma(x^T \theta)) + (1 - y_i) \log(1 - \sigma(x^T \theta))$$

$$= \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^n y_i \log(\sigma(x^T \theta)) + (1 - y_i) \log(1 - \sigma(x^T \theta))$$



CROSS entropy loss