

# Logistic regression & cross entropy Loss

①

Setup: Logistic regression

Given data  $\{(x_i, y_i)\}_i$ ,  $x_i \in \mathbb{R}^n$ ,  $y_i \in \{0, 1\}$

We will assume  $y_i$  was generated according to distribution

$$f(x) = P(y=1|x)$$

Idea:  $x$  is a data point &  $y$  is corresponding label.  $f(x)$  is going to be a smooth function denoting the probability of a positive label. We want to do regression in the sense that we want to fit some function  $f_\theta(x)$  to  $f(x)$ . However, we don't have access to  $f$ , only noisy samples  $y \sim f(x)$ .

In logistic regression, we assume the following parameter form for  $f_\theta$ .

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

Note: No bias b/c this can actually be added by augmenting  $x$  with a 1 at the end.  
We'll try to fit  $f_\theta$  to  $D = \{(x_i, y_i)\}_i$  by maximizing the likelihood of the data.

~~Terminology: If  $f_\theta$  is a parametric class of probability distributions, then the likelihood~~

Terminology: Let  $\underbrace{P(y|x; \theta)}_{:= f_{\theta}(x)}$  denote the probability (density) of  $y$  given  $x$

under the density induced by  $\theta$ . In an abuse of notation, the likelihood function given data  $D = \{(x_i, y_i)\}$  is

$$\begin{aligned} L(\theta) &= P(\{y_i\} | \{x_i\}; \theta) \\ &= \prod_{i=1}^n P(y_i | x_i; \theta) = \prod_{i=1}^n f_{\theta}(x_i, y_i) \end{aligned}$$

Note

b/c here we'll always assume independence

$$P(y_i = 1 | x_i; \theta) = \sigma(\theta^T x_i) = \frac{1}{1 + e^{-\theta^T x_i}}$$

$$P(y_i = y_i | x_i; \theta) = \sigma(\theta^T x_i)^{y_i} (1 - \sigma(\theta^T x_i))^{1-y_i}$$

Let

$$LL(\theta) := \log L(\theta) \quad \text{denote the log-likelihood.}$$

Have

$$LL(\theta) = \log \prod_{i=1}^n P(y_i | x_i; \theta)$$

$$= \sum_{i=1}^n \ln P(y_i | x_i; \theta)$$

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

Comments:

• You don't need to use the  $\sigma(\theta^T x_i) (1 - \sigma(\cdot))^{1-y_i}$  trick.  
Just split over positive & negative labels.

• In book they use  $y_i \in \{-1, 1\}$  and use the loss

$$\sum_i \ln(1 - e^{-y_i \theta^T x_i})$$

This is the same thing.

Lets generalize: multiple classes & General model

$$\text{Let } f_n(x; \theta) = P(y_i = n | x; \theta)$$

Where "n" denotes a class number

Let  $N_n = \#\{x: \text{occurs in } D\}$   $P(y = n | x) :=$  Empirical distribution of  $y$  given  $x$

For a single  $x$  (for which there may be  $N > 1$  indep  $\checkmark$   $\text{label} = \# \{y = n \text{ when } x \text{ was sampled } y\}$  feature vectors)

$$L(\theta | x) = \prod_n f_n(x; \theta)^{N P(y=n|x)}$$

$$\frac{1}{N} LL(\theta | x) = \frac{1}{N} \ln \prod f_n(x; \theta)^{N P(y=n)}$$

$$= \sum_n \ln f_n(x; \theta)^{P(y=n)}$$

$$= \sum_n P(y=n) \ln f_n(x; \theta)$$

$$= H(P, f(x; \theta)) \text{ where } f = (f_n)_n$$

For many  $x$  (but typically only 1 label  $y$  per  $x$ )

(4)

$$LL(\theta) = \log \prod_i L(\theta | x_i)$$

$$= \sum_i LL(\theta | x_i)$$

$$= \sum_i H(p, f(x; \theta))$$

$\uparrow$   
Typically an indicator on  $y$