

## Alternative to logistic regression (12 pts)

In class, we discussed the logistic regression model for binary classification problem. Here, we consider an alternative model. We have a training set  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$  where  $\mathbf{x}_n \in \mathbb{R}^{D+1}$  and  $y_n \in \{0, 1\}$ . Like in logistic regression, we will construct a probabilistic model for the probability that  $y_n$  belongs to class 0 or 1, given  $\mathbf{x}_n$  and the model parameters,  $\theta_0$  and  $\theta_1$  ( $\theta_0, \theta_1 \in \mathbb{R}^{D+1}$ ). More specifically, we model the target  $y_n$  as:

$$\begin{aligned} p(y_n = 0 | \mathbf{x}_n; \theta_0, \theta_1) &= Ce^{\theta_0^\top \mathbf{x}} \\ p(y_n = 1 | \mathbf{x}_n; \theta_0, \theta_1) &= Ce^{\theta_1^\top \mathbf{x}} \end{aligned} \quad (1)$$

- (a) (2 pts) Find the value of  $C$  that makes Equation 1 a valid probability distribution for  $y_n$ .

$$\begin{aligned} ① \quad 1 &= \int_{-\infty}^{\infty} f(x) dx = \int_{-\infty}^{\infty} C e^{\theta_0^\top x} dx = C \int_{-\infty}^{\infty} e^u \frac{du}{\theta_0^\top} = \left( \frac{e^{\theta_0^\top x}}{\theta_0^\top} \right) \Big|_{-\infty}^{\infty} \\ \text{let } u &= \theta_0^\top x \Rightarrow du = \theta_0^\top dx \Rightarrow dx = \frac{du}{\theta_0^\top} \\ \Rightarrow 1 &= C \left( \frac{e^{\theta_0^\top (1)}}{\theta_0^\top} - \frac{e^{\theta_0^\top (-\infty)}}{\theta_0^\top} \right) = C \end{aligned}$$

$$1 = C * \left( \frac{e^{\theta_0^\top} - 1}{\theta_0^\top} \right) \Rightarrow C = \frac{\theta_0^\top}{e^{\theta_0^\top} - 1} = \frac{\theta_1^\top}{e^{\theta_1^\top} - 1}$$

- (b) (5 pts) Write the log likelihood of the parameters for a single instance  $(\mathbf{x}_n, y_n)$ :  
$$l(\boldsymbol{\theta}_0, \boldsymbol{\theta}_1) = \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}_0, \boldsymbol{\theta}_1).$$
 Express your answer in terms of  $y_n, \mathbf{x}_n, \boldsymbol{\theta}_0, \boldsymbol{\theta}_1$ .  
(Hint: use the notation in class where for a Bernoulli random variable  $Y_n$  that takes values 1 with probability  $\theta$  and 0 with probability  $1 - \theta$ , we can write its probability mass function:  $p(y_n) = \theta^{y_n}(1 - \theta)^{1-y_n}$ ).

(c) (2 pts) Write the log likelihood of the parameters  $\mathcal{LL}(\theta_0, \theta_1)$  for the full training data.

(d) (3 pts) Recall that in logistic regression, we model the target  $y_n$  as :

$$p(y_n = 0 | \mathbf{x}_n; \boldsymbol{\theta}) = 1 - \sigma(\boldsymbol{\theta}^\top \mathbf{x})$$

$$p(y_n = 1 | \mathbf{x}_n; \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \mathbf{x})$$

Here  $\sigma(x) = \frac{1}{1+e^{-x}}$  is the sigmoid function. Show that our new model (Equation 1) is exactly the same as the logistic regression model with  $\boldsymbol{\theta} = \boldsymbol{\theta}_1 - \boldsymbol{\theta}_0$ .

## Identities

### Probability density/mass functions for some distributions

$$\text{Normal} : P(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\text{Multinomial} : P(\mathbf{x}; \boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{x_k}$$

$\mathbf{x}$  is a length  $K$  vector with exactly one entry equal to 1  
and all other entries equal to 0

$$\text{Poisson} : P(x; \lambda) = \frac{\lambda^x \exp(-\lambda)}{x!}$$

## Matrix calculus

Here  $\mathbf{x} \in \mathbb{R}^n, \mathbf{b} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{n \times n}$ .  $\mathbf{A}$  is symmetric.

$$\nabla \mathbf{x}^T \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}, \quad \nabla \mathbf{b}^T \mathbf{x} = \mathbf{b}$$

## Entropy

The entropy  $H(X)$  of a Bernoulli random variable  $X \sim \text{Bernoulli}(p)$  for different values of  $p$ :

$p$	$H(X)$
$\frac{1}{2}$	1
$\frac{1}{3}$	0.92
$\frac{1}{4}$	0.81
$\frac{1}{5}$	0.73
$\frac{2}{5}$	0.97