

### 8. Linear Regression

$$\{(y_n, x_n)\}_{n=1}^N, \quad x_n \in \mathbb{R}^d, \quad y_n \in \mathbb{R}$$

What we train?

$$\sum_{n=1}^N (y_n - \langle x_n, \beta \rangle)^2 \xrightarrow{\text{min } \beta}$$

2. Statistical model?

$$T_n | X_n \sim \mathcal{N}(\mu(X_n), \sigma_n^2)$$

$$\rightarrow \mu(X_n) = \mathbb{E} T_n | X_n = \langle x_n, \beta \rangle$$

Maximum Likelihood Estimation (MLE)

$$\sum_{n=1}^N \log \mathcal{N}(y_n | \langle x_n, \beta \rangle, \sigma_n^2) \rightarrow \max_{\beta}$$

$$\sum_{n=1}^N -\frac{1}{2} \cdot \frac{1}{\sigma_n^2} (y_n - \langle x_n, \beta \rangle)^2 \rightarrow \max_{\beta}$$

$$\sum_{n=1}^N \frac{1}{\sigma_n^2} (y_n - \langle x_n, \beta \rangle)^2 \rightarrow \min_{\beta}$$

$$T_n \in \mathbb{R}$$



link function

$$\mathcal{N}(y_n | \mu, \sigma^2)$$

$$\rightarrow T_n \in \mathbb{N}_+$$

$$T_n \in \{0, 1\}$$

$$g(\mathbb{E} T_n | X_n) = \bar{x}_n^\top \beta \in \mathbb{R}$$

## 2. Example

$\mu_n$  is expected number of dearsises on the day  
 $t_n$

- model  $\mu_n = \gamma \exp(\delta t_n)$

$$\log \mu_n = \underbrace{\log \gamma + \delta t_n}_{\in \mathbb{R}}$$

linear model

because

### 3. Exponential Family of Distributions

$$p_\theta(y) = \exp(\langle \theta, \tau(y) \rangle - B(\theta)) h(y)$$

$$B(\theta) = \log \int \exp(\langle \theta, \tau(y) \rangle) h(y) dy$$

$$\log p_\theta(y) = \langle \theta, \tau(y) \rangle - B(\theta) + \log h(y)$$

$$\nabla_\theta \log p_\theta(y) = \tau(y) - \nabla_\theta B(\theta)$$

$$\nabla_\theta B(\theta) = \mathbb{E}_\theta \log \int \exp(\langle \theta, \tau(y) \rangle) h(y) dy =$$

$$= \frac{1}{\int \exp(\langle \theta, \tau(y) \rangle) h(y) dy} \int \exp(\langle \theta, \tau(y) \rangle) \tau(y) h(y) dy = \mathbb{E} \tau(y).$$

$$\nabla_\theta \log p_\theta(y) = \tau(y) - (\mathbb{E} \tau(y))$$

$$(\mathbb{E} \nabla_\theta \log p_\theta(y) = (\mathbb{E} \tau(y)) - (\mathbb{E} \tau(y)) = 0.$$

$$\sum_{n=1}^N \log p_\theta(y_n) \rightarrow \max_{\theta}; \sum_{n=1}^N \langle \theta, \tau(y_n) \rangle - N B(\theta)$$

Take the gradient  $\rightarrow$

first we can

$$\sum_{n=1}^N T(y_n) - N \mathbb{E}_\theta B(\theta) = 0$$

↑ because problem  
is come x

$$\frac{1}{N} \sum_{n=1}^N T(y_n) = \mathbb{E}_\theta B(\theta) = \mathbb{E} T(y_n)$$

$\underbrace{\quad}_{\text{empirical avg}}$

$\underbrace{\quad}_{\text{true average}}$

$$g(\mathbb{E} T_n | X_n) = x_n^\top \beta ? \in \mathbb{R}$$

$$\Downarrow \mathbb{E}_\theta B(\theta) = \mathbb{E} T(y_n)$$

$$x_n^\top \beta = \theta = (\mathbb{E}_\theta B)^{-1} (\mathbb{E} T(y_n)) .$$

mean  
statistics  
only

we solve for 1-D ,

we could restrict our desire  $T(y) = y$

$$p_\theta(y) = \exp \left\{ \frac{y \cdot \theta - b(\theta)}{\varphi} + c(y, \varphi) \right\}$$

$$\log p_\theta(y) = \underbrace{y \cdot \theta - b(\theta)}_{\text{"mean"} \downarrow} \underbrace{\varphi}_{\text{"variance"} \downarrow} + c(y, \varphi)$$

$$\nabla_\theta: [y - \mathbb{E}_\theta b(\theta)] \frac{1}{\varphi}$$

$$\mathbb{E} y = \mathbb{E}_\theta b(\theta), \quad \theta = (\mathbb{E}_\theta b)^{-1} (\mathbb{E} y)$$

# Poisson Regression

$$y_n \in \mathbb{N}^+$$

$$p(y_n) = \frac{\mu^{y_n}}{y_n!} \exp(-\mu)$$

$$\log p(y_n) = y_n \cdot \log \mu - \mu - \log(y_n!)$$

$$\theta = \log \mu, \mu = \exp(\theta)$$

$$(x_n, y_n) \quad y_n \in \mathbb{N}^+ \quad \text{MLE}$$

$$\sum_n \log p(y_n) \rightarrow \max$$

$$\sum_n y_n \cdot \log \mu - \mu \rightarrow \max$$

$$\sum_{n=1}^N y_n \cdot \theta - \exp(\theta) \rightarrow \max_{\theta}$$

$$\theta = x_n^\top \beta \quad \sum_{n=1}^N y_n \cdot (x_n^\top \beta) - \exp(x_n^\top \beta) \rightarrow \max_{\beta}$$

$$\nabla_{\beta} : \sum_{n=1}^N y_n \cdot x_n - \exp(x_n^\top \beta) x_n$$

$$\nabla_{\beta} : \underbrace{\sum_{n=1}^N [y_n - \exp(x_n^\top \beta)]}_{y_n - \mathbb{E}[y_n | x_n]} x_n = \Theta \quad \downarrow \text{features}$$

$$\beta^{k+1} = \beta^k + \alpha \nabla_{\beta} = \beta^k + \sum_{n=1}^N (y_n - \mathbb{E} y_n | x_n) x_n$$

## Classification Problems

Bernoulli Distribution

$$y_n \in \{0, 1\}$$

$$p(y_n) = p^y_n (1-p)^{1-y_n}, p \in (0, 1)$$

probability  
of success  $P(Y=1)$

$$\begin{aligned} \log p(y_n) &= y_n \log p + (1-y_n) \log (1-p) = \\ &= y_n \cdot \underbrace{\log \frac{p}{1-p}}_{\theta} + \underbrace{\log(1-p)}_{?} \end{aligned}$$

$$\theta = \log \frac{p}{1-p}, \exp(\theta) = \frac{p}{1-p} = \frac{1}{\frac{1}{p} - 1}$$

$$\exp(-\theta) + 1 = \frac{1}{p}$$

Sigmoid function

$$\left\{ \hat{p} = \frac{1}{\exp(-\theta) + 1} \right.$$

$$\begin{array}{l} x \in 1 \times 5 \quad x = 2 \\ 0 \times 5, 1 \times 10 \end{array}$$

$$y \cdot \theta - \underbrace{\log(1 + \exp(\theta))}_{D\theta}$$

$$\mathbb{E} y = \frac{\exp(\theta)}{1 + \exp(\theta)} = \frac{1}{\exp(-\theta) + 1}$$

$$\begin{array}{l} P(y=1) = \frac{10}{15} = \\ \frac{2}{3}. \end{array}$$

$$\left( \frac{2}{3} \right) = \frac{1}{\exp(-\theta) + 1}$$

$$\log \frac{E_y}{1-E_y} = \theta$$

↙ ↗

logit function inverse of sigmoid

$E_y|x=1 = 0$   
 $E_y|x=2 = 1$

$$\log \frac{E_y}{1-E_y} = x^T \beta$$

likelihood of  
Bernoulli distib.

MLL:

$$\sum_{n=1}^N y_n \theta_n - \log(1 + \exp(\theta_n))$$

$\uparrow x_n^T \beta$

$$\sum_{n=1}^N y_n x_n^T \beta - \log(1 + \exp(x_n^T \beta))$$

$$\hat{\beta}_{\text{ML}}: \sum_{n=1}^N y_n x_n - \frac{\exp(x_n^T \beta)}{1 + \exp(x_n^T \beta)} x_n = 0$$

$$\sum_{n=1}^N (y_n - \frac{\exp(x_n^T \beta)}{1 + \exp(x_n^T \beta)}) x_n$$

$y_n$  - observed  
 $E y_n | x_n$  - conditional exp.

$$\frac{\exp(x_n^T \beta)}{1 + \exp(x_n^T \beta)} = \frac{1}{1 + \exp(-x_n^T \beta)}$$

$E y|X_n$

$$p(y) \sim N(y | \mu, \sigma^2) \quad |E y = \mu$$

$$p(y|x) \sim N(\tilde{y} | x^\top \beta, \sigma^2) \quad |E \tilde{y} = x^\top \beta$$

$$|E y|X = x^\top \beta$$

$$|E y = |E_x |E y|X - |E_x x^\top \beta = (|E_x x|)^\top \beta.$$