

Exercise Sheet 4

Exercise 1 - Hessian of logistic regression

$$L(w, \lambda, y) = -\sum_{i=1}^N \left[y_i \log p(y_i = 1 | x_i, w) + (1-y_i) \log (1 - p(y_i = 1 | x_i, w)) \right]$$

$$L_w = -\sum_{i=1}^N \left[y_i \log p_i + (1-y_i) \log (1-p_i) \right]$$

Gradient:

$$\nabla_w L_i = \left(y_i \frac{\nabla_w p_i}{p_i} + (1-y_i) \frac{-\nabla_w p_i}{1-p_i} \right) = -\nabla_w p_i \left(\frac{y_i}{p_i} - \frac{1-y_i}{1-p_i} \right)$$

$$= -\nabla_w p_i \left(\frac{p_i - y_i}{p_i(1-p_i)} \right) \quad \text{---} \xrightarrow{*} \delta(w^T x_i) (1 - \delta(w^T x_i))$$

$$\nabla_w p_i = \nabla_w \delta(w^T x_i) = \cancel{\delta'(w^T x_i)} \cdot \cancel{x_i} = p_i(1-p_i)x_i \quad \text{---} \times \times$$

* and **

$$\Rightarrow \nabla_w L_i = (p_i - y_i) x_i \Rightarrow \text{First gradient} \quad \left. \begin{matrix} \nabla_w \left(\begin{matrix} p_1 - y_1 \\ \vdots \\ p_n - y_n \end{matrix} \right) = \left[\begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} \right] \end{matrix} \right\}_{m \times 1}$$

Sum for all samples $\Rightarrow \nabla_w L = \sum_{i=1}^N (p_i - y_i) x_i = X^T (P - Y)$

Hessian: for entry in (j, k) th position in Hessian: 2

$$\left[\nabla_w^2 L \right]_{jk} = \frac{\partial}{\partial w_j} \left(\sum_{i=1}^N (p_i - y_i) x_{i,j} \right) = \sum_{i=1}^N \frac{\partial p_i}{\partial w_j} \cdot x_{ik} \quad \left. \begin{matrix} \nabla_w \left(\begin{matrix} p_1 - y_1 \\ \vdots \\ p_n - y_n \end{matrix} \right) = \left[\begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} \right] \end{matrix} \right\} \Rightarrow$$

$$\frac{\partial p_i}{\partial w_j} = \frac{\partial}{\partial w_j} (\delta(w^T x_i)) = \delta'(w^T x_i) \cdot x_{ij} = p_i(1-p_i) x_{ij}$$

$$\Rightarrow \left[\nabla_w^2 L \right]_{jk} = \sum_{i=1}^N p_i(1-p_i) x_{ij} \cdot x_{ik} \quad \text{--- element-wise results} \quad \left. \begin{matrix} \nabla_w \left(\begin{matrix} p_1 - y_1 \\ \vdots \\ p_n - y_n \end{matrix} \right) = \left[\begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} \right] \end{matrix} \right\} \in \mathbb{R}^{N \times N}$$

Diagonal matrix $S := \text{diag}(p_1(1-p_1), p_2(1-p_2), \dots, p_N(1-p_N))$

$$\Rightarrow \sum_{i=1}^N x_{ij} p_i(1-p_i) x_{ik} \quad \text{for all entries} \quad \Rightarrow \nabla_w^2 L(w) = X^T S X$$

Exercise 2 - Kullback-Leibler Divergence:

$$D_{KL}(P||\theta) := \sum_{x \in X} P(x) \cdot \log \frac{P(x)}{\theta(x)}$$

① $D_{KL}(P||\theta) \geq 0$ and $D_{KL}(P||\theta) = 0 \iff P = \theta$

$$D_{KL}(P||\theta) = \sum_x P(x) \cdot \log \left(\frac{\theta(x)}{P(x)} \right)^{-1}, \quad \text{~~(derivative of KL divergence)~~}$$

$$\leftarrow \cancel{E(P(x) \log \theta(x))} - \sum_x P(x) \log \frac{\theta(x)}{P(x)}$$

property of $\log \Rightarrow$ for $\forall x > 0 \Rightarrow \log x \leq -1 \Rightarrow -\log x \geq 1 - x$

$$\Rightarrow -\cancel{E(P(x) \log \frac{\theta(x)}{P(x)})} \geq 0 \quad \Rightarrow -\log \frac{\theta(x)}{P(x)} \geq 1 - \frac{\theta(x)}{P(x)}$$

multiply

$$\underbrace{\text{both sides}}_{\text{by } P(x) \geq 0} \Rightarrow -P(x) \log \frac{\theta(x)}{P(x)} \geq P(x) - \theta(x) \xrightarrow{\substack{\text{Sum over all } x \in X \\ (\text{sum of elements for any probability distribution})}} -\sum_x P(x) \log \frac{\theta(x)}{P(x)} \geq \sum_x P(x) - \sum_x \theta(x) = 1 - 1 = 0$$

$$\Rightarrow \sum_{x \in X} P(x) \log \frac{P(x)}{\theta(x)} \geq 0$$

$$D_{KL}(P||\theta) = 0 \iff P = \theta$$

$$\therefore \Rightarrow D_{KL}(P||\theta) = 0 \Rightarrow \sum_{x \in X} P(x) \cdot \log \frac{P(x)}{\theta(x)} = 0 \Rightarrow 1 \cdot \log \frac{P(x)}{\theta(x)} = 0 \Rightarrow \frac{P(x)}{\theta(x)} = 1$$

$$\text{for each element, } P(x) \log \frac{P(x)}{\theta(x)} = 0 \xrightarrow{\text{divide by } P(x)} \log \frac{P(x)}{\theta(x)} = 0 \Rightarrow P(x) = \theta(x) \Rightarrow P(x) = \theta(x)$$

$$\therefore P = \theta \Rightarrow \log \frac{P(x)}{\theta(x)} = \log 1 = 0 \Rightarrow D_{KL}(P||\theta) = \sum_x (P(x) - 1) \cdot \log \frac{P(x)}{\theta(x)} = 0$$

Exercise 2 - Kullback-Leibler Divergence

2- Show that $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$ with an explicit example

Let $X = \{a, b, c\}$

$$\text{dist } P \Rightarrow P = \left(\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\right)$$

$$\text{dist } Q \Rightarrow Q = \left(\frac{1}{8}, \frac{1}{8}, \frac{3}{4}\right)$$

$$D_{KL}(P \parallel Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)} = \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{8}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{8}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{3}{4}}$$

$$= \underbrace{\frac{1}{2} \log 4}_{0.602} + \underbrace{\frac{1}{4} \log 2}_{0.301} + \underbrace{\frac{1}{4} \log \frac{1}{3}}_{-0.477} \approx 0.257 *$$

$$0.301 \quad 0.075 \quad -0.119$$

$$D_{KL}(Q \parallel P) = \sum_{x \in X} Q(x) \log \frac{Q(x)}{P(x)} = \frac{1}{8} \log \frac{\frac{1}{8}}{\frac{1}{2}} + \frac{1}{8} \log \frac{\frac{1}{8}}{\frac{1}{4}} + \frac{3}{4} \log \frac{\frac{3}{4}}{\frac{1}{4}}$$

$$= \frac{1}{8} \log \frac{1}{2} + \frac{1}{8} \log \frac{1}{2} + \frac{3}{4} \log 3 \approx 0.245 **$$

$$0.257 \neq 0.245 \Rightarrow D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$$