

# Problem 1: Gradient

$$g(A, \beta, x) = \sum_{i=1}^p \beta_{(i)} \left( \sum_{j=1}^k A_{(ij)} x_{(j)} \right)^3$$

$$\beta \in \mathbb{R}^p, A \in \mathbb{R}^{p \times k}, x \in \mathbb{R}^k, \text{ compute } \nabla_A g, \nabla_x g, \nabla_\beta g$$

$$g \in \mathbb{R} \quad g(A, \beta, x) = \beta^T (Ax)^3$$

$$\begin{cases} \nabla_A g = 3 [\beta \circ (Ax)^2] x^T \\ \nabla_x g = 3 A^T [\beta \circ (Ax)^2] \\ \nabla_\beta g = (Ax)^3 \end{cases}$$

Proof:

$$\begin{aligned} & \begin{array}{c} g(x) \\ \uparrow \\ \boxed{\beta^T a^{(2)}} \end{array} \rightarrow \frac{\partial g}{\partial \beta} = a^{(2)} = (Ax)^3 \\ & \begin{array}{c} \uparrow \\ a^{(2)} : \frac{\partial g}{\partial a^{(2)}} = \beta \end{array} \quad \frac{\partial g}{\partial a^{(2)}_{(j)}} = \beta_{(j)} \\ & \begin{array}{c} \uparrow \\ \boxed{(z^{(2)})^3} \end{array} \\ & \begin{array}{c} \uparrow \\ z^{(2)} : \frac{\partial g}{\partial z^{(2)}} = \sum_j \frac{\partial g}{\partial a^{(2)}_{(j)}} \frac{\partial a^{(2)}_{(j)}}{\partial z^{(2)}} = \sum_j \beta_{(j)} \frac{\partial e_j^T (z^{(2)})^3}{\partial z^{(2)}} = \sum_j \beta_{(j)} 3(z^{(2)})^2 e_j \\ = 3 \beta \circ (Ax)^2 \end{array} \\ & \begin{array}{c} \frac{\partial g}{\partial A} = \sum_j \frac{\partial g}{\partial z^{(2)}_{(j)}} \frac{\partial z^{(2)}_{(j)}}{\partial A} \leftarrow \begin{array}{c} \uparrow \\ \boxed{Ax} \end{array} \end{array} \\ & \begin{array}{c} = \sum_j 3 [\beta \circ (Ax)^2]_{(j)} e_j x^T \\ = 3 [\beta \circ (Ax)^2] x^T \end{array} \quad \begin{array}{c} x : \frac{\partial g}{\partial x} = \sum_j \frac{\partial g}{\partial z^{(2)}_{(j)}} \frac{\partial z^{(2)}_{(j)}}{\partial x} = \sum_j 3 [\beta \circ (Ax)^2]_{(j)} \frac{\partial e_j^T Ax}{\partial x} \\ = 3 \sum_j [\beta \circ (Ax)^2]_{(j)} (e_j^T A)^T \\ = 3 \sum_j [\beta \circ (Ax)^2]_{(j)} A^T e_j \\ = 3 A^T [\beta \circ (Ax)^2] \end{array} \\ & \begin{array}{c} \boxed{\frac{\partial z^{(2)}_{(j)}}{\partial A}} = \frac{\partial e_j^T Ax}{\partial A} = \frac{\partial \text{tr}(e_j^T Ax)}{\partial A} \\ = \frac{\partial \text{tr}(A x e_j^T)}{\partial A} = \frac{\partial \text{tr}[A (e_j x^T)^T]}{\partial A} \\ = e_j x^T \end{array} \end{aligned}$$

## Problem 2: Weighted Logistic Regression

suppose  $x_i \in \mathbb{R}^d$   $\theta \in \mathbb{R}^d$

Since  $L(\theta)$  is a function of a vector  $\theta \in \mathbb{R}^d$

the gradient is also a vector  $\in \mathbb{R}^d$

the  $j$ th entry of gradient is

$$\begin{aligned} [\nabla L(\theta)]_{(j)} &= \frac{\partial L(\theta)}{\partial \theta_{(j)}} \\ &= \frac{\partial \sum_{i=1}^n d_i [\log(1 + \exp(x_i^T \theta)) - y_i x_i^T \theta]}{\partial \theta_{(j)}} \\ &= \sum_{i=1}^n d_i \frac{\partial [\log(1 + \exp(x_i^T \theta)) - y_i x_i^T \theta]}{\partial \theta_{(j)}} \quad \text{linearity of derivative} \\ &= \sum_{i=1}^n d_i \left[ \frac{\partial \log(1 + \exp(x_i^T \theta))}{\partial \theta_{(j)}} - \frac{\partial y_i x_i^T \theta}{\partial \theta_{(j)}} \right] \\ &= \sum_{i=1}^n d_i \left[ \frac{\exp(x_i^T \theta) [x_i]_{(j)}}{1 + \exp(x_i^T \theta)} - y_i [x_i]_{(j)} \right] \\ &= \sum_{i=1}^n \left[ d_i \left( \frac{\exp(x_i^T \theta)}{1 + \exp(x_i^T \theta)} - y_i \right) [x_i]_{(j)} \right] \end{aligned}$$

$$\text{Then } \nabla L(\theta) = \sum_{i=1}^n \left[ d_i \left( \frac{\exp(x_i^T \theta)}{1 + \exp(x_i^T \theta)} - y_i \right) x_i \right]$$

### Problem 3 Weighted Linear Regression

$$y_i = x_i^T \beta^* + w_i \quad w_i \sim N(0, \sigma_i^2)$$

$$y_i \sim N(x_i^T \beta^*, \sigma_i^2)$$

$$L(y; x, \beta^*, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i^2} \exp\left[-\frac{(y_i - x_i^T \beta^*)^2}{2\sigma_i^2}\right]$$

$$NLL = -\log L = -\log\left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i^2} \exp\left[-\frac{(y_i - x_i^T \beta^*)^2}{2\sigma_i^2}\right]\right)$$

$$= -\sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi}\sigma_i^2} \exp\left[-\frac{(y_i - x_i^T \beta^*)^2}{2\sigma_i^2}\right]\right)$$

$$= -\sum_{i=1}^n \left[ \log \frac{1}{\sqrt{2\pi}\sigma_i^2} - \frac{(y_i - x_i^T \beta^*)^2}{2\sigma_i^2} \right]$$

$$= \frac{n}{2} \log 2\pi + \sum_{i=1}^n \log \sigma_i + \frac{1}{2\sigma_i^2} \sum_{i=1}^n (y_i - x_i^T \beta^*)^2$$

The effect of adding a variable  $w_i \sim N(0, \sigma_i^2)$  to a Gaussian compared to adding a constant  $w \sim N(0, \sigma^2)$  to a Gaussian in standard linear regression is that when we try to find the optimal  $\beta$  by solving  $\arg\min_{\beta} NLL(\beta)$

by setting  $\nabla_{\beta} NLL(\beta) = 0$ , every sample will contribute a  $\frac{1}{2\sigma_i^2}$  ratio which means every sample is weighted



P4: Regularization

$$L(\theta) = f(\theta) + \frac{\lambda}{2} \theta^T \theta$$

$$\nabla_{\theta} L(\theta) = \nabla f(\theta) + \lambda \theta$$

$$\theta_{k+1} = \theta_k - \eta_k \nabla_{\theta} L(\theta_k)$$

$$\boxed{\theta_{k+1} = \theta_k - \eta_k (\nabla f(\theta_k) + \lambda \theta_k)}$$

where  $\eta_k$  is the learning rate

## P5 Concepts

a)  $X \in \mathbb{R}^{100 \times 4}$ ,  $x_i \in \mathbb{R}^4$  is the  $i^{\text{th}}$  flower sample in the training set.

$$y_i = \begin{cases} 1 & \text{species A} \\ 0 & \text{species B} \end{cases}$$

For a new test case  $v \in \mathbb{R}^4$ , we calculate its distance between each flower sample in the set, so we have 100 distances, then we select the top 5 flower sample data points with shortest distances.

Then we look at the species of these 5 data points, we choose the most common species to be the label of the new test case.

b)  $k=1$  will lead to an overfitted model, label of a new test case will be the species of the 1st nearest data point, without capturing comprehensive features of the data set. Test error will be very large.

c) Because the data set only has 100 data points. Set  $k=100$ . then the label of a new test case will always be the species of majority of the data set.

# P6. Gradient Descent

① if  $|y - y'| < 1$   $l(y, y') = \frac{1}{2}(y - y')^2$

$$\begin{aligned} \nabla_{\theta} l(x_i^T \theta, y_i) &= l'(x_i^T \theta, y_i) \nabla_{\theta}(x_i^T \theta) \\ &= \left[ (y - y') \middle|_{(x_i^T \theta, y_i)} \right] \cdot x_i \\ &= (y_i - x_i^T \theta) x_i \end{aligned}$$

② if  $|y - y'| \geq 1$   $l(y, y') = |y - y'| - \frac{1}{2}$

$$\begin{aligned} \nabla_{\theta} l(x_i^T \theta, y_i) &= l'(x_i^T \theta, y_i) \nabla_{\theta}(x_i^T \theta) \\ &= \left[ \text{sign}(y - y') \middle|_{(x_i^T \theta, y_i)} \right] \cdot x_i \\ &= \text{sign}(y_i - x_i^T \theta) x_i \end{aligned}$$

$$\Rightarrow \nabla_{\theta} l(x_i^T \theta, y_i) = \begin{cases} (y_i - x_i^T \theta) x_i & \text{if } |y - y'| < 1 \\ \text{sign}(y_i - x_i^T \theta) x_i & \text{otherwise} \end{cases}$$

P7 MLE

$$\text{Likelihood } L(\theta^*) = \prod_{i=1}^n f_{\theta}(y_i) = \begin{cases} \prod_{i=1}^n \exp[-(y_i - x_i^T \theta^*)] & y_i - x_i^T \theta^* \geq 0 \\ 0 & y_i - x_i^T \theta^* < 0 \end{cases}$$

$$\begin{aligned} \log L(\theta^*) &= \log L(\theta^*) = \log \prod_{i=1}^n \exp[-(y_i - x_i^T \theta^*)] \\ &= \sum_{i=1}^n (x_i^T \theta^* - y_i) \end{aligned}$$

$$\Rightarrow \boxed{\log L(\theta^*) = \sum_{i=1}^n (x_i^T \theta^* - y_i) \quad \text{for } y_i - x_i^T \theta^* \geq 0.}$$



# Problem 8 MLE

$$x_i \sim N(0, \sigma^2) \quad \text{i.i.d}$$

$$p(x_i | \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x_i^2}{2\sigma^2}\right)$$

$$L(x; \sigma^2) = \prod_{i=1}^n \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x_i^2}{2\sigma^2}\right) \right]$$

$$NLL(\sigma^2) = -\log L(x; \sigma^2) = -\sum_{i=1}^n \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x_i^2}{2\sigma^2}\right) \right]$$

$$= \frac{n}{2} \log 2\pi + \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2$$

$$\frac{\partial NLL(\sigma^2)}{\partial \sigma^2} = \frac{\partial \left( \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2 \right)}{\partial \sigma^2}$$

$$= \frac{n}{2} \frac{1}{\sigma^2} - \frac{\frac{\sum_{i=1}^n x_i^2}{2}}{(\sigma^2)^2}$$

$$\text{set } \frac{\partial NLL(\sigma^2)}{\partial \sigma^2} = 0$$

$$\Rightarrow \boxed{\sigma^2 = \frac{\sum_{i=1}^n x_i^2}{n}}$$

The MLE estimate for  $\sigma^2$  is  $\frac{\sum_{i=1}^n x_i^2}{n}$



Problem 9 MLE

$$\hat{\theta}_j = \log \frac{N_j}{n} \quad N_j \text{ is the number of } [y_i]_{(j)} = 1 \text{ in } n \text{ samples}$$

$$y_i \in \mathbb{R}^k$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \rightarrow 2^{\text{th}} \text{ row } [y_i]_{(2)} = 1 \Rightarrow Z_{(1)} = 2$$

$$\text{vector } Z = \begin{bmatrix} 2 \\ \vdots \\ j \\ \vdots \end{bmatrix} \rightarrow 1^{\text{th}} \text{ row} \\ \rightarrow j^{\text{th}} \text{ row}$$

$$y_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \rightarrow j^{\text{th}} \text{ row} \quad [y_i]_{(j)} = 1 \Rightarrow Z_{(i)} = j$$

For convenience, denote vector  $Z \in \mathbb{R}^n$   $Z_{(i)} = k$  if  $[y_i]_{(k)} = 1$   
 define  $N_j =$  number of  $Z_{(i)} = j$   $N \in \mathbb{R}^k$   $\sum_{j=1}^k N_{(j)} = n$

$$L(y; \theta) = L(Z; \theta) = \prod_{i=1}^n P(y) = \prod_{j=1}^k (p_j)^{N_j}$$

$$- \text{NLL} = -\log \prod_{j=1}^k (p_j)^{N_j}$$

$$= -\sum_{j=1}^k \log(p_j)^{N_j}$$

$$= -\sum_{j=1}^k N_{(j)} \log(p_j)$$

$$= -\left[ \sum_{j=1}^k N_{(j)} \log(p_j) + \lambda \left( 1 - \sum_{j=1}^k p_j \right) \right]$$

$$\sum_{j=1}^k p_j = 1$$

$$\begin{cases} \frac{\partial \text{NLL}(p_j, \lambda)}{\partial p_j} = \frac{N_{(j)}}{p_j} - \lambda = 0 \end{cases}$$

$$\frac{\partial \text{NLL}(p_j, \lambda)}{\partial \lambda} = 1 - \sum_{j=1}^k p_j = 0$$

$$\Rightarrow \begin{cases} \hat{p}_j = \frac{N_{(j)}}{\lambda} = \frac{N_{(j)}}{n} = \exp(\hat{\theta}_j) \\ \lambda = n \end{cases}$$

$$\hat{\theta}_j = \log \frac{N_{(j)}}{n}$$

Also we could verify:

$$\sum_{j=1}^k \exp(\hat{\theta}_j) = \sum_{j=1}^k \exp\left(\log \frac{N_{(j)}}{n}\right) = \sum_{j=1}^k \frac{N_{(j)}}{n} = \frac{1}{n} \cdot n = 1$$

### Problem 10 Concept

First I will standardize the new data point by the same way training set has done

$x' = \frac{x-m}{s}$  (where  $m$  is mean,  $s$  is standard deviation of training set). Then apply function  $h$  to  $x'$  to predict its label.

Because the nearness of NN is based on distance, without standardize columns,

if one variable is on the scale of billions while another is on tens. then scale of billions will contribute more to error. Thus the prediction will bias towards

samples that are close in billions scale feature and ignore the other. Since

the training set was standardized, the test set should also be standardized

to the same scale.