

### Question 3

The logistic model for  $q$  classes can be defined as:

$$P(Y_i = k | X_i) = \frac{e^{x_i B_k}}{\sum_{k'=1}^q e^{x_i B_{k'}}} = P_{ik}$$

where  $x_i = i^{\text{th}}$  row of the feature matrix,  $X$   
 $B_k = k^{\text{th}}$  column of the parameter matrix,  $B$

As calculated in class, we want to maximize  $L(B)$  using log-likelihood

$$L(B) = \prod_{i=1}^n \prod_{k=1}^q P(Y[i]=k | X[i,:])^{11[i,k]}$$

lets take log of  $L(B)$

$$\log L(B) = \log \left( \prod_{i=1}^n \prod_{k=1}^q P(Y[i]=k | X[i,:])^{11[i,k]} \right)$$

$$= \sum_{i=1}^n \sum_{k=1}^q 11[i,k] \cdot \log P(Y[i]=k | X[i,:])$$

Substituting the  $P_{ik}$  function

$$= \sum_{i=1}^n \sum_{k=1}^q 1[i,k] \cdot \log \left( \frac{e^{X[i,:]\cdot B[:,k]}}{\sum_{k'=1}^q e^{X[i,:]\cdot B[:,k']}} \right)$$

simplifying wrt log

$$= \sum_{i=1}^n \sum_{k=1}^q 1[i,k] \cdot \left( X[i,:]\cdot B[:,k] - \log \left( \sum_{k'=1}^q e^{X[i,:]\cdot B[:,k']} \right) \right)$$

Let's calculate the gradient for this  $\log L(B)$  function.  
We will do it partially for each term

$$\frac{d}{dB} (X[i,:]\cdot B[:,k]) = X[i,j] \quad \text{--- (1)}$$

$$\begin{aligned} \frac{d}{dB} \left( \log \left( \sum_{k'=1}^q e^{X[i,:]\cdot B[:,k']} \right) \right) &= \frac{X[i,j] \cdot e^{X[i,:]\cdot B[:,k]}}{\sum_{k'=1}^q e^{X[i,:]\cdot B[:,k']}} \\ &= X[i,j] \cdot P_{ik} \quad \text{--- (2)} \end{aligned}$$

Let's substitute these for the entire derivative  
Furthermore  $1[i,k]$  is treated as a constant.

$$\frac{d}{dB} \log L(B) = \sum_{i=1}^n (1[i,k] \cdot X[i,j] - X[i,j] P_{ik})$$

taking  $X[i,j]$  as common

$$= \sum_{i=1}^n X[i,j] (11[i,k] - P_{ik})$$

Converting to matrix form,  
we get

$$\frac{d \log L(\beta)}{d \beta} = X^T (I - P)$$

Question 1A

Given

$$X = [x_1, x_2, \dots, x_n]^T$$

$$y = [y_1, y_2, \dots, y_n]^T$$

the model can be defined as

$$y = XB \quad \text{where } B \text{ is the parameter matrix to be estimated.}$$

Then Mean Square Error  $L(B)$  is given by

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

Since, the goal is to minimize MSE wrt  $B$

We can iteratively update  $B$  using

$$B^{(k+1)} = B^{(k)} - \alpha \frac{dL(B)}{dB}$$

Assume that the iterative update rule for the parameter,  $B$  behaves like a fixed point iteration

$B_{t+1} = B_t - \alpha l(B)$  where  $l(B)$  is the derived loss function for MSE.

$$\begin{aligned} \text{and since } \frac{dL(B)}{dB} &= -2 x^T e \\ &= -2 \sum_{i=1}^n x_i (y_i - x_i B) \end{aligned}$$

$$MSE(B + 2\alpha(x^T e)) < MSE(B)$$

the MSE after updating  $B$  is given by

$$\begin{aligned} MSE(B + 2\alpha(x^T e)) &= \frac{1}{n} \sum_{i=1}^n (y_i - x_i(B + 2\alpha x^T e))^2 \\ &= \frac{1}{n} \sum_{i=1}^n y_i^2 + x_i^2(B + 2\alpha x^T e)^2 - 2 y_i x_i (B + 2\alpha x^T e) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{i=1}^n y_i^2 + x_i^2 (B^2 + (2\alpha x_i^T e)^2 + 4\alpha x_i^T e) - 2y_i x_i (B + 2\alpha x_i^T e) \\
&= \frac{1}{n} \sum_{i=1}^n y_i^2 + x_i^2 B^2 + x_i^2 (2\alpha x_i^T e)^2 + 4\alpha x_i^T e - 2y_i x_i (B + 2\alpha x_i^T e)
\end{aligned}$$

MSE for B

$$\begin{aligned}
\text{MSE}(B) &= \frac{1}{n} \sum_{i=1}^n (y_i - x_i B)^2 \\
&= \frac{1}{n} \sum_{i=1}^n y_i^2 + (x_i B)^2 - 2x_i B y_i
\end{aligned}$$

$$\leq x_i^2 (2\alpha x_i^T e)^2 + 4\alpha x_i^T e - 2y_i x_i (B + 2\alpha x_i^T e) \leq -2x_i B y_i$$

$$\leq x_i^2 (2\alpha x_i^T e)^2 + 4\alpha x_i^T e - 2y_i x_i B - 4\alpha x_i^T e \leq -2x_i B y_i$$

Which we can solve for  $\alpha$

$$\alpha < \sum_{i=1}^n \frac{x_i y_i}{x_i^2}$$

## Question 1 B

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - xB)^2$$

This is a quadratic function, therefore the plot will be parabolic.

Since the LSE of B minimizes the MSE, we set the derivative of MSE wrt B to 0

$$\frac{d}{dB} l(B) = \frac{2}{n} \left( B \sum_{i=1}^n x_i^2 - \sum_{i=1}^n x_i y_i \right) = 0$$

$$B \sum x_i^2 - \sum x_i y_i = 0$$

$$B \sum x_i^2 = \sum x_i y_i$$

$$B_{LSE} = \frac{\sum x_i y_i}{\sum x_i^2}$$

Case 1:  $\alpha < B_{LSE}$  at this the parameter B will slowly approach  $B_{LSE}$

Case 2:  $\alpha > B_{LSE}$

