

1) Linear Regression

$$L(g) = \sum_i (g(x_i) - y_i)^2, L(h) = \sum_i (h(x_i) - y_i)^2$$

$$g(x) = ax + b, h(x) = cx^2 + dx + f$$

Answer :

(a)

The squared loss is :

$$J(g) = \sum_i (g(x_i) - y_i)^2, J(h) = \sum_i (h(x_i) - y_i)^2$$

Assume that we put the parameters in a vector w , then the loss function $J(g)$ becomes

$$J(g) = \sum_i (w^T x_i - y_i)^2, w = \begin{bmatrix} b \\ a \end{bmatrix}$$

Then, we collect x_i and y_i in a vector X and we will have

$$J(g) = (w^T X - Y)^T (w^T X - Y), X = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 \end{bmatrix}, Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix}$$

then, taking partial derivatives with respect to each parameter, we will have

$$\frac{\partial J(g)}{\partial w} = \begin{bmatrix} \frac{\partial J}{\partial w_0} \\ \vdots \\ \frac{\partial J}{\partial w_n} \end{bmatrix}$$

Setting the partial derivatives to zero $\frac{\partial J}{\partial w} = 0$,

$$\frac{\partial}{\partial w} ((w^T X - Y)^T (w^T X - Y)) = \frac{\partial}{\partial w} (y^2 + X^T X w^T w - X^T w^T Y - X w^T Y)$$

$$0 = 2 X^T X w - 2 X^T Y$$

$$X^T X w = X^T Y$$

Multiplying each side with $(X^T X)^{-1}$ gives

$$\cancel{(X^T X)^{-1}} (X^T X w) = (X^T X)^{-1} X^T y$$

$$\hat{w} = (X^T X)^{-1} X^T y$$

Then, inserting the values that we have from the problems, we will have

$$w = \begin{bmatrix} b \\ a \end{bmatrix}, X = \begin{bmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix}, y = \begin{bmatrix} 3 \\ 5.2 \\ 8 \\ 10.3 \\ 11.6 \end{bmatrix}, X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \end{bmatrix}$$

$$\hat{w} = (X^T X)^{-1} X^T y$$

$$= \left(\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ 5.2 \\ 8 \\ 10.3 \\ 11.6 \end{bmatrix}$$

$$= \begin{pmatrix} 5 & 0 \\ 0 & 10 \end{pmatrix}^{-1} \begin{bmatrix} 38.6 \\ 10.2 \end{bmatrix} = \frac{1}{50} \begin{bmatrix} 10 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 38.6 \\ 21.8 \end{bmatrix}$$

$$= \begin{bmatrix} 0.2 & 0 \\ 0 & 0.1 \end{bmatrix} \begin{bmatrix} 38.6 \\ 21.8 \end{bmatrix}$$

$$\hat{w} = \begin{bmatrix} 7.72 \\ 2.18 \end{bmatrix} = \begin{bmatrix} b \\ a \end{bmatrix}$$

Hence, the model that we have is

$$g(x) = \hat{w}^T x = ax + b =$$

$$2.18x + 7.72$$

For the second loss, we can construct the loss function in matrix/vector form as below

$$L(h) = \sum_i (ex_i^2 + dx_i + f - y_i)^2$$

Similar with the first loss function, we construct the loss function in a matrix/vector form?

$$w = \begin{bmatrix} f \\ d \\ c \end{bmatrix}, X = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & \vdots & \vdots \\ 1 & \vdots & \vdots \\ 1 & x_N & x_N^2 \end{bmatrix} = \begin{bmatrix} 1 & -2 & 4 \\ 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix}, y = \begin{bmatrix} 3 \\ 5.7 \\ 8 \\ 10.3 \\ 11.6 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \\ 4 & 1 & 0 & 1 & 4 \end{bmatrix}$$

Estimating the model parameters,

$$\hat{w} = (X^T X)^{-1} X^T y$$

$$= \left(\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \\ 4 & 1 & 0 & 1 & 4 \end{bmatrix} \begin{bmatrix} 1 & -2 & 4 \\ 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \\ 4 & 1 & 0 & 1 & 4 \end{bmatrix} \begin{bmatrix} 3 \\ 5.7 \\ 8 \\ 10.3 \\ 11.6 \end{bmatrix}$$

$$= \left(\begin{bmatrix} 5 & 0 & 10 \\ 0 & 10 & 0 \\ 10 & 0 & 34 \end{bmatrix} \right)^{-1} \begin{bmatrix} 38.6 \\ 21.8 \\ 74.4 \end{bmatrix}$$

$$\hat{w} = \begin{bmatrix} 0.986 & 0 & -0.143 \\ 0 & 0.1 & 0 \\ -0.143 & 0 & 0.071 \end{bmatrix} \begin{bmatrix} 38.6 \\ 21.8 \\ 74.4 \end{bmatrix} = \begin{bmatrix} 8.12 \\ 2.18 \\ -0.24 \end{bmatrix} = \begin{bmatrix} f \\ d \\ c \end{bmatrix}$$

Hence, the model that we have is

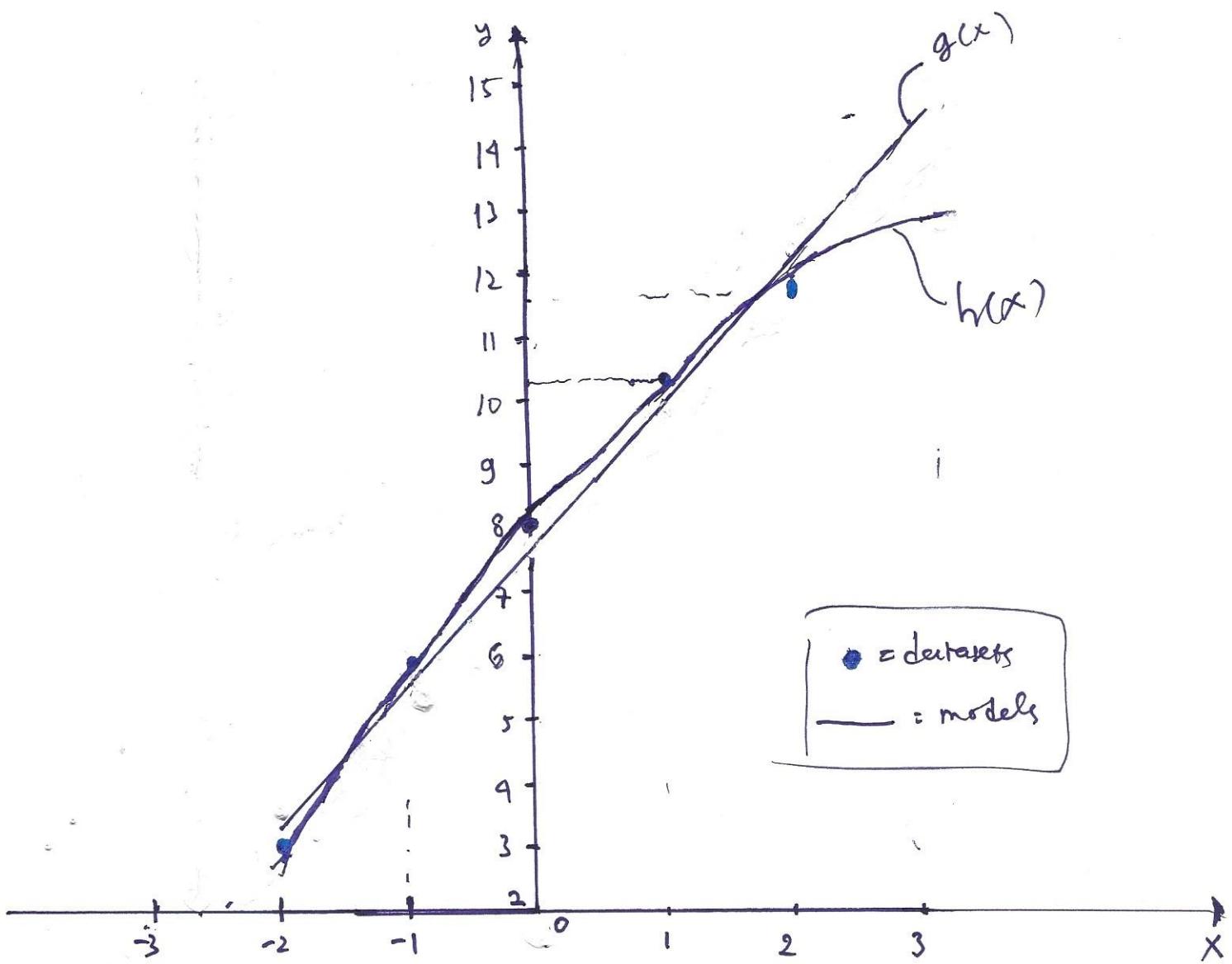
$$h(x) = ex^2 + dx + f = -0.24x^2 + 2.18x + 8.12$$

Finding the points to plot according to the models

table 1. $\alpha=0$

$g(x)$	$h(x)$	x_i	x_i^2	y	$L(g(x,y))$	$L(h(x,y))$
7.18	-0.22	-3	9	NA	NA	NA
3.36	2.96	-2	4	3	0.1296	0.0016
5.54	5.74	-1	1	5.7	0.0256	0.0016
2.72	8.12	0	0	8	0.0784	0.0144
9.90	10.1	1	1	10.3	0.16	0.040
12.08	11.68	2	4	11.6	0.2304	0.0064
14.26	12.86	3	9	NA	NA	NA

Hence, plotting both model in a single curve together with the dataset gives



b) Let $\alpha=6$, then re-calculating the parameters estimation, for $g(x)$

$$w = \begin{bmatrix} b \\ a \end{bmatrix}, X = \begin{bmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix}, y = \begin{bmatrix} 3 \\ 5.7 \\ 14 \\ 10.3 \\ 11.6 \end{bmatrix}, X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \end{bmatrix}$$

$$\hat{w} = \underbrace{(X^T X)^{-1}}_{\text{fixed term}} \underbrace{X^T Y}_{\text{changing terms}}$$

$$= \frac{1}{50} \begin{bmatrix} 10 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ 5.7 \\ 14 \\ 10.3 \\ 11.6 \end{bmatrix}$$

$$= \begin{bmatrix} 0.2 & 0 \\ 0 & 0.1 \end{bmatrix} \begin{bmatrix} 44.6 \\ 21.8 \end{bmatrix} = \begin{bmatrix} 8.92 \\ 2.18 \end{bmatrix} = \begin{bmatrix} b \\ a \end{bmatrix}$$

Hence, the model that we have become

$$g(x) = \hat{w}x = ax+b = 2.18x + 8.92$$

for the second model $h(x)$, re-calculating the parameters estimation

$$w = \begin{bmatrix} f \\ d \\ c \end{bmatrix}, X = \begin{bmatrix} 1 & -2 & 4 \\ 1 & -1 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \end{bmatrix}, y = \begin{bmatrix} 3 \\ 5.7 \\ 14 \\ 10.3 \\ 11.6 \end{bmatrix}$$

$$\hat{w} = \underbrace{(X^T X)^{-1}}_{\text{fixed term}} \underbrace{X^T Y}_{\text{changing term}}$$

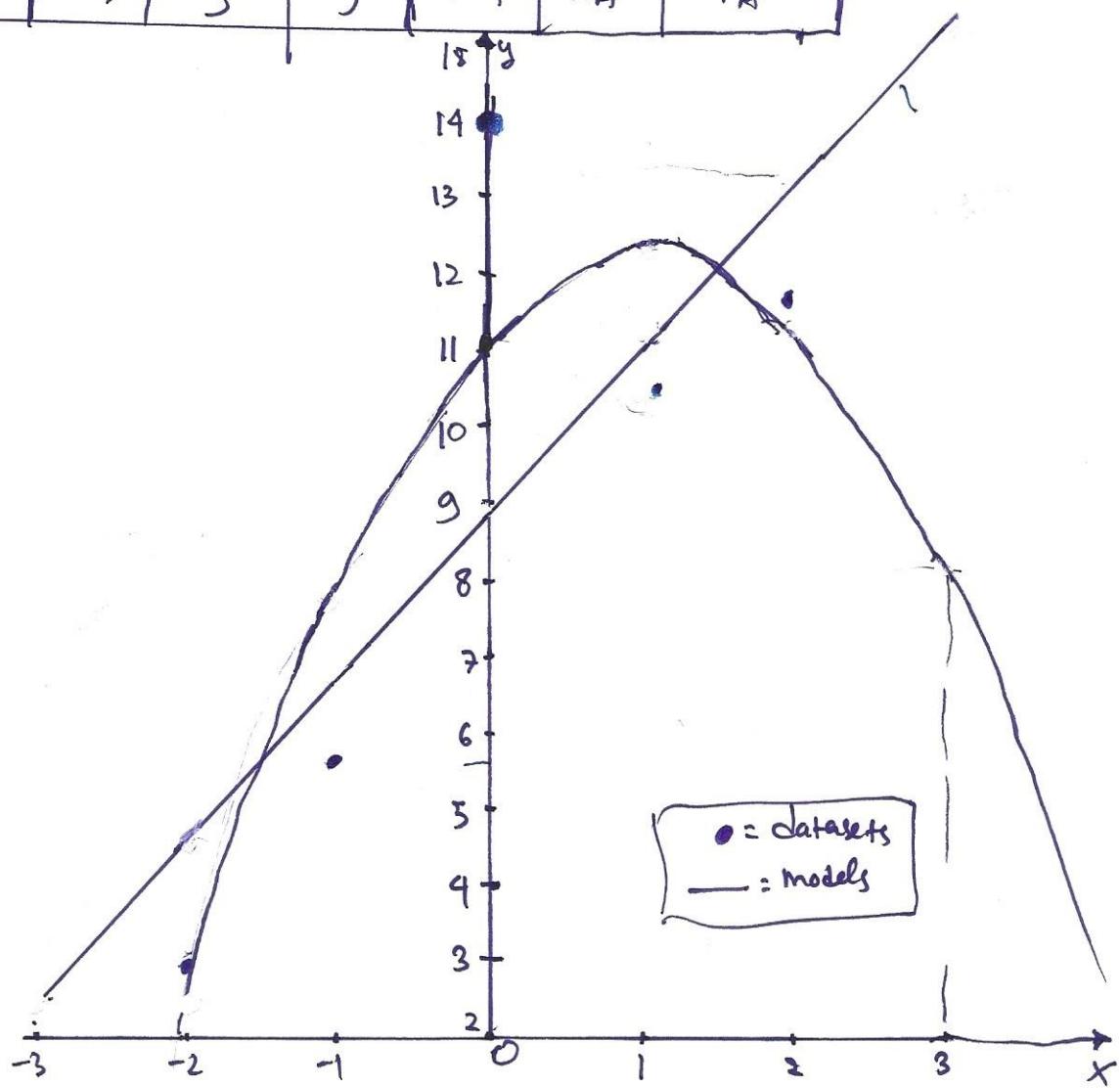
$$= \left(\begin{bmatrix} 5 & 0 & 10 \\ 0 & 10 & 0 \\ 10 & 0 & 34 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -2 & -1 & 0 & 1 & 2 \\ 4 & 0 & 1 & 4 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 5.7 \\ 14 \\ 10.3 \\ 11.6 \end{bmatrix}$$

$$= \begin{bmatrix} 0.486 & 0 & -0.143 \\ 0 & 0.1 & 0 \\ -0.143 & 0 & 0.071 \end{bmatrix} \begin{bmatrix} 44.6 \\ 21.8 \\ 74.4 \end{bmatrix} = \begin{bmatrix} 11.03 \\ 2.18 \\ -1.06 \end{bmatrix} = \begin{bmatrix} f \\ d \\ c \end{bmatrix}$$

The model $h(x)$ becomes $h(x) = cx^2 + dx + f = 11.03x^2 + 2.18x - 1.06$

Finding the points to plot according to the model
table 2. $\alpha=6$

$g(x)$	$h(x)$	x_i	x^2	y	$\frac{L(g, x_i, y)}{L(h, x_i, y)}$	$L(h, x_i, y)$
2.38	-5.05	-3	9	NA	NA	NA
4.56	2.43	-2	4	3	2.43	0.325
6.74	7.79	-1	1	5.7	1.08	4.87
8.92	11.03	0	0	14	25.81	8.82
11.1	12.15	1	1	10.3	0.64	3.92
13.28	11.15	2	4	11.6	2.82	0.203
15.46	8.03	3	9	NA	NA	NA



We could observe that both models becomes underfit. and the loss becomes considerable larger. In previous question, we could easily determine that the quadratic model is the best model. However, here as we add $\alpha=6$, it becomes harder to determine the best model.

This is due to the fact that at $x=0$, data is omitted from the other data. Hence, fitting becomes hard.

c) To show that we always have $\min L(g_i, x_i, y) \geq \min L(h, x_i, y)$,

Answer: In the first model with $\alpha = 0$,

we could observe that the minimum loss of the linear model is

$$\min L(g_i, x_i, y) = -0.0256$$

for $\alpha = 6$,

$$\min L(g_i, x_i, y) = 0.64$$

In the second model with $\alpha = 0$,

$$\min L(h, x_i, y) = 0.0016$$

For $\alpha = 6$,

$$\min L(h, x_i, y) = 0.203$$

conclude

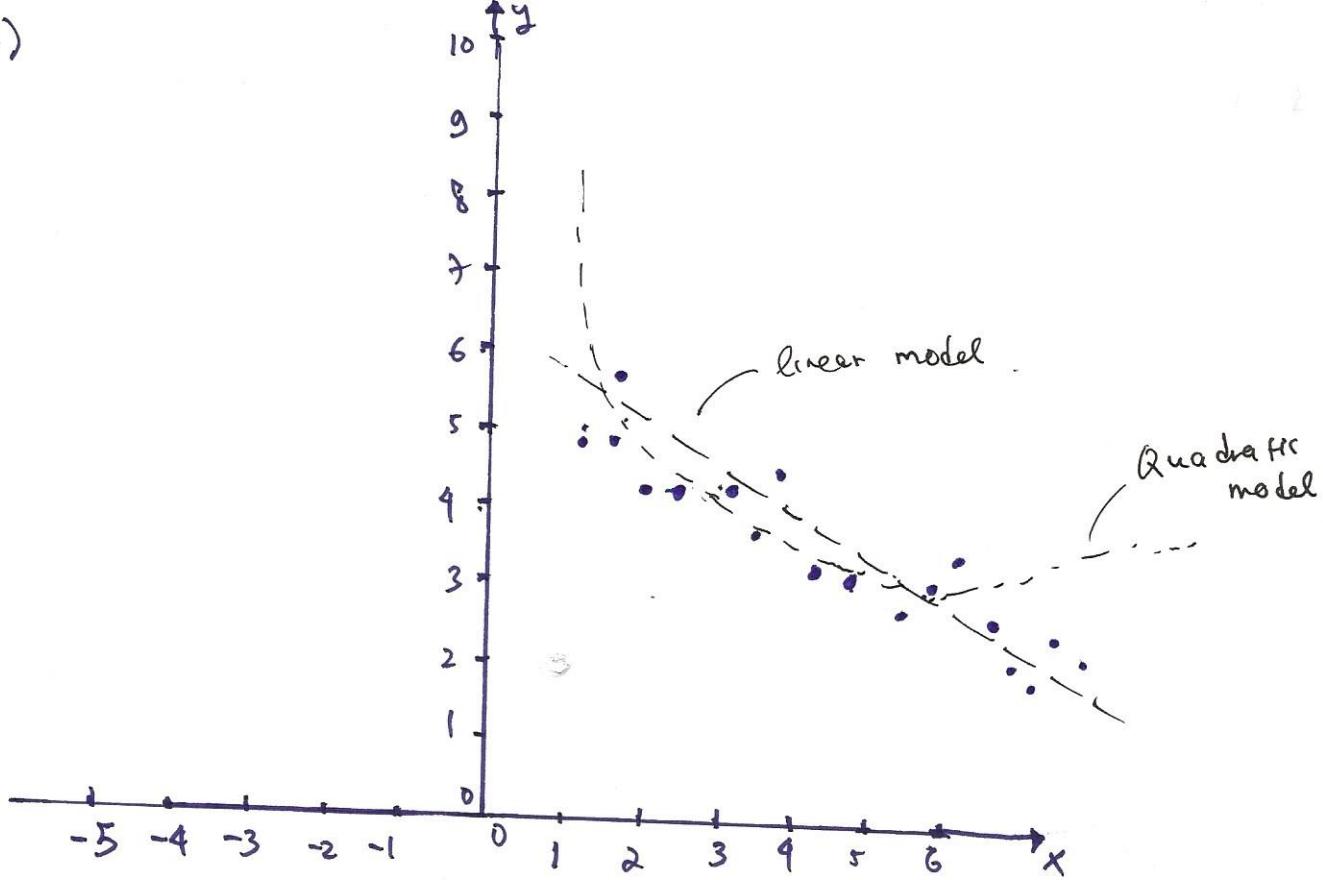
From these two datasets, we could conclude that we always have $\min L(g_i, x_i, y) \geq \min L(h, x_i, y)$

Theoretically, we could also say that since loss always decreases as the model is made more complex, we will always have

$$\min L(g_i, x_i, y) \geq \min L(h, x_i, y)$$

due to the fact that model $h(x)$ has higher complexity than model $g(x)$

e)



In this kind of data, where the data has a trend almost perfectly linearly decreasing combined with the fact that the

Variance is small / where data is more compact to each other, is better fitted using linear regression model.

It is shown that using quadratic model, the data is underfitted whereas using linear model, the data was fitted well.

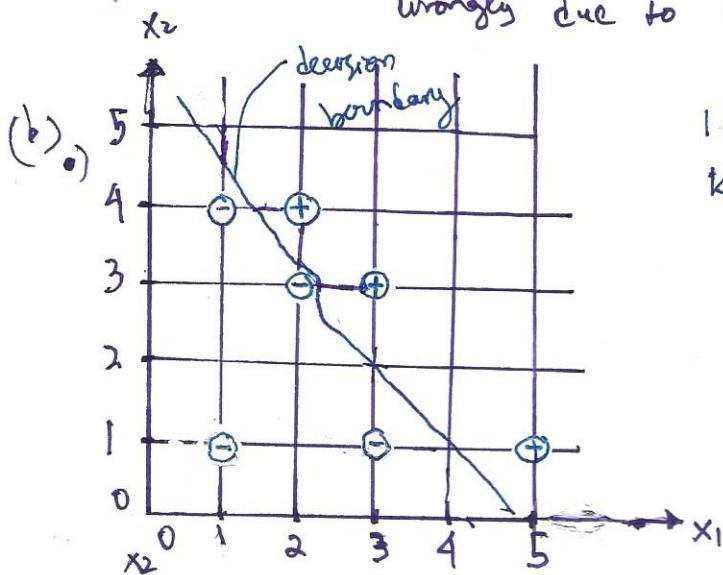
~~scribble~~

3.) Classification

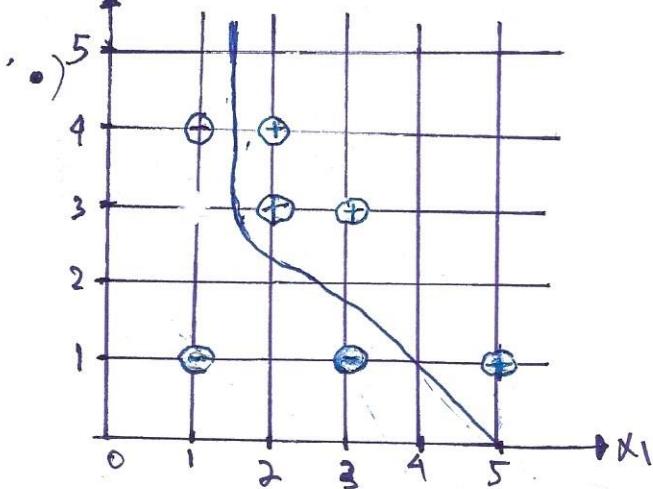
Answer:

- (a) The figure with decision boundary 1
 is the one \rightarrow belongs to 1-nearest neighbor classifier.
 thus it is due to the fact that we have an island
 as a result of overfitting. We could also
 say that for 1-nearest neighbor, it is
 heavily influenced by noise exemplified by the island.
 $(k \geq 2)$

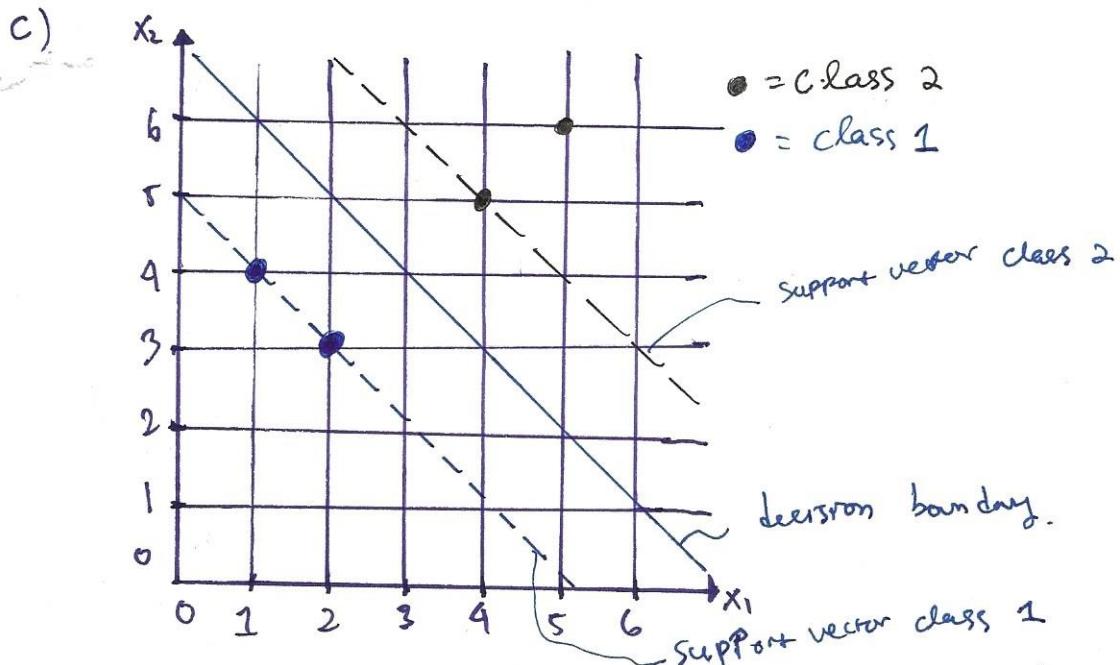
The figure with decision boundary 2 \rightarrow k -nearest neighbor,
 the figure will be smoother than
 with 1-nearest neighbor. However,
 there is one outlier which should have been
 classified as class (+), however classified
 wrongly due to the decision boundary.



1-nearest neighbor \rightarrow ? is nearest neighbor
 $k=1 \rightarrow$? would be
 classified as (-)



3 nearest neighbor
 \rightarrow {2 (+) class,
 {1 (-) class
 $k=3 \rightarrow$? would be
 classified as (+)



The plot for the 4 training points could be illustrated as above. From the plot, we could say the following:

- the decision boundary could be expressed mathematically by linear equation as: $x_2 = -x_1 + 7$
- the support vector for class 1, could be expressed mathematically as: $x_2 = -x_1 + 5$
- The support vector for class 2 could be expressed as: $x_2 = -x_1 + 9$

Then, we know that the general equation for decision boundary is:

$$w_1x_1 + w_2x_2 + b = 0$$

~~$$x_2 = -x_1 + 7 \leftrightarrow x_1 + x_2 - 7 = 0$$~~

Hence, the weights and biases are

$$w_1 = 1, w_2 = 1, b = -7$$

d) Limitations of linear hard SVM

- 1) In linear hard SVM, it can only work well when data is completely linearly separable without any errors in the form of noise or outliers. If there is an error, since a single outlier can determine the boundary, the hard margin SVM would not work well.

Hence, it possibly overfits to a particular dataset and therefore cannot generalize well.

- 2) Linear hard SVM is not very suited for large datasets.
- 3) In case of when we have more features for each data point than the number of training data samples, we will have bad performance of SVM.

(Q). Assume $p_i = p_j \forall i, j$, what is the probability of the document "turn left"?

Answer:

A multinomial distribution is given by:

$$P(x_1, \dots, x_k | p_1, \dots, p_k) = d! \frac{(p_1)^{x_1} \cdots (p_k)^{x_k}}{(x_1)! \cdots (x_k)!} = \frac{d!}{\prod_{t=1}^d x_t!} \prod_{t=1}^d p_t^{x_t}$$

Assuming random variables x_1, \dots, x_k such that $\sum_{k=1}^K x_k = d$.

Consider the vocabulary

$$V = \{\text{left, motorway, the, leave, turn, and, Gothenburg}\}$$

$$|V| = d = 7.$$

Let d^M is the multinomial feature vector, then we would have

$$d^M = [1, 0, 0, 0, 1, 0, 0]^T$$

Assuming $p_{\text{turn}} = 1$. Since $p_i = p_j$, Then, the probability of the document "turn left" is

$$\text{By Multinomial Naive-Bayes } \rightarrow P(x_1, x_2 | c) = \underbrace{p(c)}_{\substack{\text{we have only} \\ \text{single class}}} \cdot \prod P(x_i | c) = 1 \cdot \frac{1}{7} \cdot \frac{1}{7} = \frac{1}{49}$$

b) A multinomial distribution is given by

general

$$P(x_1, \dots, x_k | p_1, \dots, p_k) = d! \frac{(p_1)^{x_1} \times \dots \times (p_k)^{x_k}}{(x_1)! \times \dots \times (x_k)!}$$

$$= \frac{d!}{\prod_{i=1}^k x_i!} \prod_{i=1}^k p_i^{x_i}$$

$$= \binom{d}{x_1, x_2, \dots, x_k} \cdot p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$

Changing the variable x_i with fdk_i ,

$$= \binom{d}{\sum_{i=1}^k fdk_i} \cdot \prod_{i=1}^k p_i^{fdk_i}$$

then, taking the log-likelihood

$$\ell(p_1, p_2, \dots, p_k) = \log d! + \left(\sum_{i=1}^k (fdk_i \times \log p_i) \right) - \sum_{i=1}^k \log fdk_i!$$

The arrow shows the relationship before and after log-likelihood calculation.

4.) Deep Learning

a) The convolution operation of the weights (2×2 filters) on the image ($2 \times n$) size will result in output of size ($1 \times n-1$) as shown below

$$\text{Image} \odot \text{Filter} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n} \end{bmatrix} \odot \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix}$$

$$= [z_1 \ z_2 \ \dots \ z_{n-1}]$$

Let's say that the output has an index m , therefore we could generalize the operation with the following equation

$$z_m = \sum_{i=1}^2 \sum_{j=1}^2 w_{i,j} x_{i,j+m-1} + b$$

where index $m = 1, \dots, n-1$, $b \rightarrow \text{bias}$

b) Given $\frac{dE}{dz_i} \rightarrow$ gradients of loss function E wrt outputs z_i

$$\frac{dE}{dw_{i,j}} =$$

- c) For $k > 2$ (multi-class classification), to design the network, we need to take several common measures:
1. Define the output label as a one-hot encoded vector such that when the output is predicted, the value of node corresponding to output should be 1 whereas the other nodes should be 0.
 2. Make sure to use softmax and cross-entropy function, softmax is used due to the fact that it takes a vector as input as well as it produces an output between 0 and 1. Hence, it is beneficial to use it in multi-classification problem. Regarding cross entropy loss function, it is showing a very good performance for multiclass classification. Hence, it is common to use it.
 3. Define the feed forward phase. Mind the dimension of the nodes.
 4. Design the Backward propagation surely.
 5. Choose the best initialization method. (e.g. He initialization)
 6. Interpreting result of multiclass classification is better by using quantitative metrics (e.g. confusion matrix) rather than visual metric since it is easier to evaluate. (e.g. ROC)
 7. Tune the Network architecture by increasing network complexity, checkpoints/early stopping or learning rate scheduler, dropout, batch normalization etc.
 8. choose the best optimizer which performs best for multiclass classification.

d) $\hat{p}_{j_i}(x) \rightarrow$ probability computed by model for input x in class j

$y_i \in \{1, \dots, k\} \rightarrow$ corresponding class labels

$i = 1, \dots, N$

$x_i \rightarrow$ training dataset consisting of N samples

Answer:

By using Maximum Likelihood principle, the likelihood of input x_i 's in class y_i could be written as:

$$\text{negative } P(\hat{y}_1, \dots, \hat{y}_N | x_i) = \prod_{i=1}^N \hat{p}_{y_i}(x_i)$$

The log-likelihood could be expressed as

$$\log(P(\hat{y}_1, \dots, \hat{y}_N | x_i)) = -\log \prod_{i=1}^N \hat{p}_{y_i}(x_i)$$

The maximization of the negative log-likelihood wrt y_i will result in

$$L(\hat{y}_i, x_i) = -\sum_{i=1}^N \log \hat{p}_{y_i}(x_i)$$

Averaging ^{loss} over the whole dataset will give us the Cross Entropy Loss as

$$CE = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}_{y_i}(x_i))$$

S.) Mixture Models and EM algorithm

a) ... $\Theta = \{\pi_1; \lambda_1; \lambda_2\}$ be the set of parameters

, exponential distribution with parameter λ has the density function

$$f(x; \lambda) = \lambda e^{-\lambda x}$$

The distribution of inaction times is a mixture of 2 exponential distributions:

$$p(x=x) = \pi_1 \lambda_1 e^{-\lambda_1 x} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x}$$

z_k → cluster variable

Inactivity times samples x_n

3	5	2
---	---	---

Answer:

$$p(x_n | z_1=1, \Theta) = ?$$

Since we have latent variable, thus only one element of z is equal to 1 and all other elements are zero.

$$z_k \in \{0, 1\}, \sum_k z_k = 1$$

Then, we know that $z_1=1$,

Hence $\boxed{z_2=0}$

The marginal (prior) distribution over z is specified in terms of the mixing coefficients π_{ik} .

$$P(z_k=1) = \pi_{ik}$$

$$\downarrow$$

$$P(z_1=1) = \pi_1$$

for $p(x)$ to be a proper distribution
 $\sum_{k=1}^2 \pi_{ik} = 1$ and $\pi_{ik} \geq 0, 1 \leq k \leq 2$

Therefore, the conditional distribution of the sample x_n given a particular value of z is

$$P(x_n | z_1=1, \pi_1, \lambda_1, \lambda_2) = \pi_1 \lambda_1 e^{-\lambda_1 x_n} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_n}$$

$$P(x_1 | z_1=1, \pi_1, \lambda_1, \lambda_2) = \pi_1 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_1}$$

$$P(x_2 | z_1=1, \pi_1, \lambda_1, \lambda_2) = \pi_1 \lambda_1 e^{-\lambda_1 x_2} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_2}$$

$$P(x_1, x_2 | z_1=1, \pi_1, \lambda_1, \lambda_2) = \pi_1 \lambda_1 e^{-\lambda_1 x_1} (1-\pi_1) \lambda_2 e^{-\lambda_2 x_2}$$

Find full data log-likelihood

- b) $P(x_1, x_2, x_3 | \Theta)$ in terms of different distributions and parameters?
 $\Sigma = (x_1, x_2, x_3), N=3,$

Answer:

The full data log likelihood could be expressed as.

$$\ln P(\Sigma | \pi_1, \pi_2, \lambda_1, \lambda_2) = \sum_{n=1}^N \ln \left\{ \pi_1 (\lambda_1 e^{-\lambda_1 x_n}) + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_n}) \right\}$$
$$= \ln \left[\underbrace{\pi_1 (\lambda_1 e^{-\lambda_1 x_1}) + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_1})}_{a} \right] + \ln \left[\underbrace{\pi_1 (\lambda_1 e^{-\lambda_1 x_2}) + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_2})}_{b} \right]$$
$$+ \ln \left[\underbrace{\pi_1 (\lambda_1 e^{-\lambda_1 x_3}) + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_3})}_{c} \right]$$

By the identity $\ln(a) + \ln(b) + \ln(c) = \ln(abc)$, we get

$$= \ln \left[(\pi_1 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_1})) (\pi_1 (\lambda_1 e^{-\lambda_1 x_1} + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_1})) \right. \\ \left. (\pi_1 (\lambda_1 e^{-\lambda_1 x_3}) + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_3})) \right]$$
$$= \ln \left[(\pi_1 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_1})) \left\{ \pi_1^2 (\lambda_1^2 e^{-\lambda_1(x_1+x_3)} + (1-\pi_1)^2 (\lambda_2^2 e^{-\lambda_2(x_1+x_3)}) \right. \right. \\ \left. \left. + (1-\pi_1) \pi_1 (\lambda_1 \lambda_2 e^{-\lambda_2 x_1 - \lambda_1 x_3}) + (1-\pi_1) \pi_1 (\lambda_1 \lambda_2 e^{-\lambda_1 x_1 - \lambda_2 x_3}) \right\} \right]$$
$$= \ln \left[(\pi_1 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_1) (\lambda_2 e^{-\lambda_2 x_1})) \left\{ \pi_1^2 (\lambda_1^2 e^{-\lambda_1(x_1+x_3)} - \lambda_2^2 e^{-\lambda_2(x_1+x_3)}) + (1-\pi_1) \pi_1 \lambda_1 \lambda_2 \right. \right. \\ \left. \left. e^{-\lambda_2(x_1+x_3) - \lambda_1(x_1+x_3)} + \lambda_2^2 e^{-\lambda_2(x_1+x_3)} \right\} \right]$$

Inserting $(x_1, x_2, x_3) = (3, 5, 2)$

$$= \ln \left[(\pi_1 \lambda_1 e^{-3\lambda_1} + (1-\pi_1) (\lambda_2 e^{-3\lambda_2})) \left(\pi_1^2 (\lambda_1^2 e^{-8\lambda_1} - \lambda_2^2 e^{-8\lambda_2}) + (\pi_1 - \pi_1^2) \lambda_1 \lambda_2 e^{-5(\lambda_1 + \lambda_2)} \right) \right]$$
$$= \ln \left[(\pi_1^3 \lambda_1^3 e^{-11\lambda_1} - \pi_1^2 \lambda_1^2 \lambda_2 e^{-3\lambda_1 - 5\lambda_2}) + (\pi_1^2 - \pi_1^3) \lambda_1^2 \lambda_2 e^{-8\lambda_1 - 5\lambda_2} + \lambda_2^2 e^{-5\lambda_2} \right. \\ \left. + (1-\pi_1) \pi_1^2 (\lambda_1^2 e^{-8\lambda_1} - \lambda_2^2 e^{-8\lambda_2}) (\lambda_2 e^{-3\lambda_2}) + (1-\pi_1) \lambda_2^3 e^{-8\lambda_2} \right. \\ \left. + (1-\pi_1) (\pi_1 - \pi_1^2) \lambda_1^2 \lambda_2 e^{2 - 5\lambda_1 - 5\lambda_2} \right]$$

= - - - - -

$$\gamma(z_{22}) = \frac{\pi_2 \lambda_1 e^{-\lambda_1 z_2} + (1-\pi_2) \lambda_2 e^{-\lambda_2 z_2}}{0.0926}$$

$$= \frac{((0.6)(2)e^{-10} + (0.4)(3)e^{-65})}{0.0926}$$

$$= 7.72 \times 10^+$$

$$\gamma(z_{32}) = \frac{\pi_2 \lambda_1 e^{-\lambda_1 z_3} + (1-\pi_2) \lambda_2 e^{-\lambda_2 z_3}}{0.0926}$$

$$= \frac{((0.6)(2)e^{-9} + (0.4)(3)e^{-6})}{0.0926}$$

$$= 0.35$$

d) $N_k = \sum_{n=1}^3 \gamma(z_{nk})$

$$\uparrow \downarrow$$

$$N_1 = \gamma(z_{11}) + \gamma(z_{21}) + \gamma(z_{31})$$

$$= 0.052 + 0.019 + 0.18 = 0.249$$

Re-estimating the mixing coefficient as

$$\pi_1 = \frac{N_1}{N} = \frac{0.249}{3} = 0.083$$

$$N_2 = \gamma(z_{12}) + \gamma(z_{22}) + \gamma(z_{32})$$

$$= 0.044 + 7.7 \times 10^+ + 0.35 = 0.399$$

Re-estimating the mixing coeffs.

$$\pi_2 = \frac{N_2}{N} = \frac{0.399}{3} = 0.133$$

The new parameters are

$$\textcircled{2}_1 = \left\{ \pi_1 = 0.083, \lambda_1 = 2, \lambda_2 = 3, \pi_2 = 0.133 \right\}$$

$$c) \quad \Theta_0 = \{ \pi_1, \pi_2, \pi_3, \pi_2 \pi_3 = 3 \}, \quad x_1 = 3, x_2 = 5, x_3 = 2$$

calculate responsibilities

$$\gamma(z_{nk}) = p(z_k=1 | x_n, \Theta_0) \text{ for all } n \in \{1, 2, 3\}, k = \{1, 2\}$$

$$\sum_{k=1}^2 \pi_k = 1, \quad \pi_1 = 0.4 \rightarrow \pi_2 = 1 - 0.4 = 0.6$$

Answer:

$$\gamma(z_{nk}) = p(z_k=1 | x_n, \Theta_0) = \frac{p(z_{nk}=1) p(x_n | z_{nk}=1)}{\sum_{j=1}^k p(z_{nj}=1) p(x_n | z_{nj}=1)}$$

$$\gamma(z_{11}) = \frac{(\pi_1 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_1})}{\sum_{j=1}^k} (0.4)$$

$$\left\{ \begin{array}{l} (\pi_1 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_1}) \\ + (\pi_1 \lambda_1 e^{-\lambda_1 x_2} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_2}) \\ + (\pi_1 \lambda_1 e^{-\lambda_1 x_3} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_3}) \end{array} \right\} \times 0.4 \quad \left\{ \begin{array}{l} (\pi_2 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_2) \lambda_2 e^{-\lambda_2 x_1}) \\ + (\pi_2 \lambda_1 e^{-\lambda_1 x_2} + (1-\pi_2) \lambda_2 e^{-\lambda_2 x_2}) \\ + (\pi_2 \lambda_1 e^{-\lambda_1 x_3} + (1-\pi_2) \lambda_2 e^{-\lambda_2 x_3}) \end{array} \right\} \times 0.6$$

$$= (0.4)(2)e^{-6} + (0.6)3e^{-9}$$

$$\left\{ \begin{array}{l} ((0.4)(2)e^{-6} + (0.6)3e^{-9}) + ((0.4)(2)e^{-15} + (0.6)(3)e^{-15}) \\ + ((0.6)(2)e^{-6} + (0.4)3e^{-15}) + ((0.6)(2)e^{-10} + (0.4)(3)e^{-10}) \\ + ((0.6)(2)e^{-4} + (0.4)(3)e^{-6}) \end{array} \right\}$$

$$= 2.2 \times 10^{-3}$$

$$(0.04) + (0.002) + (0.019) + (0.003) + (0.003) + (0.025)$$

$$= \frac{0.0022}{0.024 + 0.0186} = \frac{0.0022}{0.0426} = 0.052$$

$$\gamma(z_{21}) = \frac{0.4(\pi_1 \lambda_1 e^{-\lambda_1 x_2} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_2})}{0.0426} = \frac{((0.4)(2)e^{-15} + (0.6)(3)e^{-15})}{0.0426} \times 0.4$$

$$\gamma(z_{31}) = \frac{\pi_1 \lambda_1 e^{-\lambda_1 x_3} + (1-\pi_1) \lambda_2 e^{-\lambda_2 x_3}}{0.0426} = \frac{(0.4)(2)e^{-4} + (0.6)(3)e^{-6}}{0.0426} \times 0.4$$

$$\gamma(z_{12}) = \frac{\pi_2 \lambda_1 e^{-\lambda_1 x_1} + (1-\pi_2) \lambda_2 e^{-\lambda_2 x_1}}{0.0426} = \cancel{0.019} \quad 0.18$$

$$= \frac{((0.6)(2)e^{-6} + (1-0.6)(3)e^{-9})}{0.0426} \times 0.6 = 0.099$$