Logistic Regression

Q1)

a) Given,

$$\omega = \{1, 5, 100\}$$

To plot the sigmoid function $1/(1+e^{-\omega x})$ vs $x \in R$ or $\omega$

For $\omega = 1$



Observations from the above graphs :-

o Like you can see that as the weight value is increasing from 1 to 100 the line in between become almost vertical or at 90° to the x-axis.

→ As # with higher weight values, when we change the input values slightly then also overall change will be high or more.

→ Thus high weights leads to overfitting as # because of high variation in output

→ With higher weights, the output probabilities are either 0 or 1 generally. (0 denotes one class and 1 denotes other class).

b) Given,

$P(w_0, \ldots, w_d)$ is a prior on weights.

→ Standard Gaussian prior $N(0, I)$ where I is a identity matrix.

→ To prevent overfitting in, we are replacing MLE (Maximum likelihood Estimation) with MAP (Maximum a Posteri)

→ The derivation for the second part of the question proceeds as follows,

Here, $\omega = [\omega_0, \ldots, \omega_d]^T$

Log conditional posterior is

$$L(\omega) = \log\left(P(\omega) \prod_{j=1} P(y^j | x^j, \omega)\right) \quad —①$$

$$P(\omega) = \prod_{i=0}^{d} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\omega_i^2}{2}\right) \quad —②$$

Now putting this $P(\omega)$ value of equation ② in the equation ①

the MAP estimate is

$$\omega^* = \arg\max_\omega L(\omega) = \arg\max_\omega \left[\sum_{j=1} \log(P(y^j|x^j,\omega)) - \sum_i \frac{\omega_i^2}{2}\right]$$

Now applying gradient ascent update rule

$$\omega_i^{t+1} \leftarrow \omega_i^t + \eta \frac{\partial L(\omega)}{\partial \omega_i} \quad —③$$

where

weight at time $t+1$ $\leftarrow$ weight at time $t$

$\eta \rightarrow$ step size

Here,

In equation ③

$$\frac{\partial L(\omega)}{\partial \omega_i} = \frac{\partial \log P(\omega)}{\partial \omega_i} + \frac{\partial \log \left(\prod_{j=1}^{n} P(y^j|x^j,\omega)\right)}{\partial \omega_i}$$

Now putting $P(\omega)$ value from equation ② into equation ③

$$\frac{\partial \log P(\omega)}{\partial \omega_i} = \frac{\partial}{\partial \omega_i} \log\left(\prod_{i=0}^{d} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\omega_i^2}{2}\right)\right)$$

$$= \frac{d}{d\omega_i} \left( \sum_{i=0} \frac{1}{\sqrt{2\pi}} + \frac{d}{d\omega_i} \log \exp\left(\frac{-\omega_i^2}{2}\right) \right)$$

$$= \frac{d}{d\omega_i} \left( \sum_{i=0} \frac{1}{\sqrt{2\pi}} \right)^0 + \frac{d}{d\omega_i} \left(\frac{-\omega_i^2}{2}\right)$$

$$= 0 + -2 \times \frac{1}{2} \times \omega_i$$

$$= -\omega_i$$

$$\Rightarrow \frac{d}{d\omega_i} \log(P(\omega)) = -\omega_i$$

And, $\dfrac{d}{d\omega_i} \log \left( \prod_{j=1} P(y^j | x^j, \omega) \right) = \sum_{j=1} x_i^j \left( y^j - P(y=1 | x^j, \omega) \right)$

Final update rule :-

$$\omega_i^{(t+1)} \leftarrow \omega_i^{(t)} + \eta \left( -\omega_i^{(b)} + \sum_{j=1} x_i^j \left( y^j - P(y=1 | x^j, \omega^{(t)}) \right) \right)$$

c) Extending logistic regression to multi-class (k class labels) :-

Since sum of all probabilities must sum to 1, we should have

$$P(Y=y_k | x) = 1 - \sum_{k=1}^{k-1} P(Y=y_k | x)$$

Defining logistic regression as we defined in Page-5 binary classification,

$$P(Y = y_k | x) = \frac{1}{1 + \sum_{k=1}^{k-1} \exp -(w_{k0} + \sum_{i=1}^{d} w_{ki} x_i)}$$
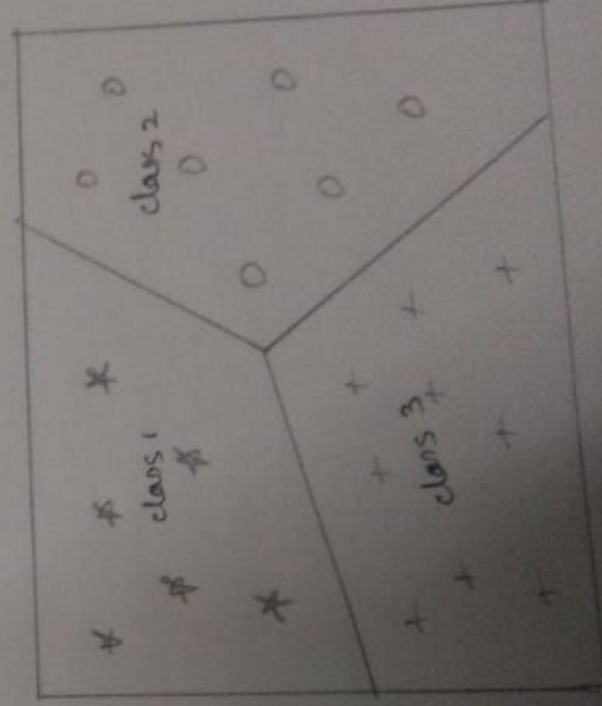
and for $k = 1, \dots k-1$

$$P(Y = y_k | x) = \frac{\exp(w_{k0} + \sum_{i=1}^{d} w_{ki} x_i)}{1 + \sum_{k=1}^{k-1} \exp(w_{k0} + \sum_{i=1}^{d} w_{ki} x_i)}$$

Classification Rule :-

$$y = y_k$$

where $k = \arg\max_{k \in \{1, \dots k\}} P(Y = y_k | x)$

thus, it picks the class which has highest probability.

d) 3 - class logistic regression Decision boundary

In multi-class logistic regression, the decision boundary is decided like we do in one vs all SVM, here also we find like that only.

→ Decision boundary between each pair of class will be linear, Hence the overall decision boundary will be piece - wise linear.