

## 機器學習理論 HW2

108064535 陳文遠

### Step 1. 使用 Maximum Likelihood & Least Squares 法來訓練 Training\_set.csv 檔中的資料

□ 假設有一線性方程式，而題目給定此自變數為一個 feature vector 如下：

$$\phi(x) = [\phi_1(x), \phi_2(x), \dots, \phi_P(x), \phi_{P+1}, \phi_{P+2}]^T \text{ for } P = O_1 \times O_2$$

其中，

$$\phi_k(x) = e^{-\frac{(x_1 - \mu_i)^2}{2s_1} - \frac{(x_2 - \mu_j)^2}{2s_2}} \text{ for } 1 \leq i \leq O_1, 1 \leq j \leq O_2$$

$$k = O_2 \times (i - 1) + j$$

$$\mu_i = s_1 \times (i - 1) + x_{1\_min}, \quad \mu_j = s_2 \times (j - 1) + x_{2\_min}$$

$$s_1 = \frac{x_{1\_max} - x_{1\_min}}{O_1 - 1}, \quad s_2 = \frac{x_{2\_max} - x_{2\_min}}{O_2 - 1}$$

最後，

$$\phi_{P+1}(x) = x_3(\text{Research Experience}) \text{ and } \phi_{P+2}(x) = 1(\text{bias})$$

### Step 2. 將 Testing\_set.csv 中的測試資料送進 Maximum Likelihood & Least Squares 的 Model 中來預測 chance of admit 以及計算其 squared error $(y(x) - t(x))^2$

□ 假設線性函數為( $\epsilon$  為 Added Gaussian Noise)：

$$t = y(x, w) + \epsilon = \sum_{j=1}^{P+2} w_j \phi_j(x) + \epsilon \text{ where } p(\epsilon|\beta) = N(\epsilon|0, \beta^{-1})$$

其機率分布可以寫成：

$$p(t|x, w, \beta) = N(t|y(x, w), \beta^{-1})$$

其中， $X = \{x_1, \dots, x_n\}$  是我們輸入的資料(分數)， $t = [t_1, t_2, \dots, t_n]^T$  為 Targets (chance of admit)。

則我們可以獲得下方的 likelihood function：

$$p(t|X, w, \beta) = \prod_{n=1}^N N(t_n|w^T \phi(x_n), \beta^{-1})$$

再將 likelihood function 取對數：

$$\begin{aligned} \ln[p(t|X, w, \beta)] &= \sum_{n=1}^N N(t_n|w^T \phi(x_n), \beta^{-1}) \\ &= \frac{N}{2} \ln(\beta) - \frac{N}{2} \ln(2\pi) - \beta \times \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2 \end{aligned}$$

最後將上述式子對  $w$  微分令為 0 來求極值：

$$\frac{d}{dw} \ln[p(t|X, w, \beta)] = \beta \sum_{n=1}^N \{t_n - w^T \phi(x_n)\} \phi(x_n)^T = 0$$

$$\Rightarrow w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t ; \text{ proved}$$

□ 接下來，將 Testing\_csv 檔中的測試資料代入  $\phi_k(x)$  來構成 testing data 的 feature vector ( $\phi(x)$ ):

$$\phi_k(x) = e^{-\frac{(x_1 - \mu_i)^2}{2s_1} - \frac{(x_2 - \mu_j)^2}{2s_2}} \text{ for } 1 \leq i \leq O_1, 1 \leq j \leq O_2$$

同樣地，

$$\phi_{p+1}(x) = x_3(\text{Research Experience}) \text{ and } \phi_{p+2}(x) = 1(\text{bias})$$

□ 最後求出 predict 的值：

$$\hat{y} = \Phi w_{ML} = \phi(x) w_{ML}$$

Squared Error 為：

$$\text{Squared Error} = (y(x) - t(x))^2 = (\hat{y} - y)^2 = (\hat{y} - x_4)^2$$

### Step 3. 使用 Bayesian Linear Regression 法來訓練 Training\_csv 檔中的資料

□ 根據貝式定理我們得知

$$p(w|t) \propto p(t|w)p(w)$$

假設事前機率 (priori probability) 為

$$p(w) \sim N(w|m_0, S_0) = N(w|0, \alpha^{-1}I)$$

$$m_0 = 0, \quad S_0 = \alpha^{-1}I$$

假設後驗機率 (posteriori probability) 為

$$p(w|t) \sim N(w|m_N, S_N)$$

經推導後可得  $m_N = \beta S_N \Phi^T t$  以及  $S_N^{-1} = \alpha I + \beta \Phi^T \Phi$ ，其中  $\Phi$  與 **Step 1.** 的求法一樣

□ 在此實驗中，我僅先將  $\alpha$  以及  $\beta$  設為 1

### Step 4. 將 Testing\_set.csv 中的測試資料送進 Bayesian Linear Regression 的 Model 中來預測 w, chance of admit 以及計算其 squared error $(y(x) - t(x))^2$

□ 權重值為  $w = w_{MAP} = m_N$ ，此題計算出的 w 分別為：

-0.18731298      0.07518969      0.03911943      0.49466673      0.35914929

□ 接下來，將 Testing\_csv 檔中的測試資料代入  $\phi_k(x)$  來構成 testing data 的 feature vector ( $\phi(x)$ ):

$$\phi_k(x) = e^{-\frac{(x_1 - \mu_i)^2}{2s_1} - \frac{(x_2 - \mu_j)^2}{2s_2}} \text{ for } 1 \leq i \leq O_1, 1 \leq j \leq O_2$$

同樣地，

$$\phi_{p+1}(x) = x_3(\text{Research Experience}) \text{ and } \phi_{p+2}(x) = 1(\text{bias})$$

□ 最後求出 predict 的值：

$$\hat{y} = \Phi w_{\text{MAP}} = \phi(x) w_{\text{MAP}}$$

Squared Error 為：

$$\text{Squared Error} = (y(x) - t(x))^2 = (\hat{y} - y)^2 = (\hat{y} - x_4)^2$$

## Step 5. 討論

### 1. Maximum likelihood and least squares 以及 Bayesian linear regression 的差異：

已知貝式定理如下，其中  $\theta$  為欲估計的數， $D$  為輸入的資料集

$$p(\theta|D) = \frac{p(D|\theta) \times p(\theta)}{p(D)} \Rightarrow \text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

Maximum likelihood and least squares 的目的是為了找出一個  $\theta$  來使得 likelihood  $p(D|\theta)$ 。換句話說在此方法中，他將  $\frac{p(\theta)}{p(D)}$  視為一個常數而非隨機變數，並且不需要借助 prior probability 來估計  $\theta$

而在 Bayesian linear regression 中會完全計算出 posterior probability  $p(\theta|D)$ ，並且在此方法中會將  $\theta$  視為隨機變數，倘若  $p(\theta|D)$  的變異數足夠小，就將他的期望值視為估計值。Bayesian 法與 Maximum likelihood 法的關鍵區別就是他允許採用 prior information

### 2. $O_1, O_2$ 對結果的影響

為了比較其影響，我分別將 Maximum likelihood 法以及 Bayesian 法對 100 的樣本估計的 squared error 進行加總，接著修改  $O_1, O_2$  的值來觀察何時會得到最小 squared error

從觀察的結果中會發現  $O_1, O_2$  對 Bayesian 法的結果影響不大，對 Maximum likelihood 法的影響較大，而當  $O_1 = 2, O_2 = 4$  時我可以得到最小的總 squared error

當  $O_1 = 2, O_2 = 4$  時，Maximum likelihood 法之 100 個樣本的 squared error 加總為 0.4363359203046397，而 Bayesian 法之 100 個樣本的 squared error 加總為 0.44922167547056285