

學號：R10922111 系級：資工所碩一 姓名：黃冠瑋

1. (1%) 請比較說明**generative model**、**logistic regression**兩者的異同為何？再分別列出本次使用的資料中五個分得正確/不正確的**sample**，並說明為什麼如此？

generative model 是根據數學理論推出可能值，logistic regression則是一直調整直到尋找到適合的參數。

2. (1%) 請實作兩種**feature scaling**的方法 (**feature normalization**, **feature standardization**)，並說明哪種方法適合用在本次作業？

	normalization	standardization
logistic regression	1.69399	1.68803

normalization 較適合，因為他會把資料轉換成0~1，且因為搜集資料的feature只有前面幾個需要做調整，後面都屬於binary feature。因此轉換後，這時所有的feature都會在同一個範圍內，在進行梯度下降時便可以減少迭代次數，增加精準度。

3. (1%) 請說明你實作的**best model**及其背後「原理」為何？你覺得這次作業的**dataset**比較適合哪個**model**？為什麼？

我使用decision tree 來做分類。他運行的方式就像一棵樹，從root往下到leaf，每個節點都是一個feature，根據feature再細分成不同的節點。

我覺得這次作業適合decision tree這個模型，因為我們的feature有很大一部分都屬於離散的，很容易做劃分。

4. (3%) Refer to math problem

<https://hackmd.io/@GfOkB4kgS66YhhM7j6TJew/BJ-wGv8HY>

$$1. P(X_n, t_n) = \prod_{k=1}^K (P(X_n | C_k) P(C_k))^{t_{nk}} = \prod_{k=1}^K (P(X_n | C_k) \pi_k)^{t_{nk}}$$

$$\Rightarrow \text{Max likelihood function } L(\theta) = \prod_{n=1}^N \prod_{k=1}^K (P(X_n | C_k) \pi_k)^{t_{nk}}$$

$$\log L(\theta) = \ell(\theta) = \sum \sum t_{nk} [\log P(X_n | C_k) + \log \pi_k]$$

$$\text{Since } \sum_{k=1}^K \pi_k = 1, \therefore \ell(\theta, \lambda) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} [\log P(X_n | C_k) + \log \pi_k] + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial \ell(\theta)}{\partial \pi_k} = \frac{1}{\pi_k} \sum_{n=1}^N t_{nk} + \lambda = 0 \quad \Rightarrow \quad \pi_k = -\frac{\sum_{n=1}^N t_{nk}}{\lambda} = \frac{-N_k}{\lambda}$$

$$\uparrow$$

$$\text{To maximize } \ell(\theta), \frac{\partial \ell(\theta, \lambda)}{\partial \pi_k} = 0 = \frac{\partial \ell(\theta, \lambda)}{\partial \lambda}$$

$$\frac{\partial \ell(\theta, \lambda)}{\partial \lambda} = \sum_{k=1}^K \pi_k - 1 = 0 \quad \Rightarrow \quad \sum_{k=1}^K \pi_k = 1 = \sum_{k=1}^K \frac{-N_k}{\lambda} = \frac{-N}{\lambda}$$

$$\Rightarrow \lambda = -N$$

$$\text{So, } \pi_k = \frac{-N_k}{-N} = \frac{N_k}{N}$$

$$2. \text{ Assume } A \in \mathbb{R}^{m \times m}, \text{ and it's non-singular. } |A| = \sum_{j=1}^m (-1)^{i+j} a_{ij} |A_{ij}| = \sum_{i=1}^m (-1)^{i+j} a_{ij} |A_{ij}|$$

$$\text{then } \frac{\partial |A|}{\partial a_{ij}} = (-1)^{i+j} |A_{ij}|$$

$$\text{And since } A \text{ is non-singular, } \therefore \exists x \in \mathbb{R}^{m \times 1}, \text{ st } Ax = e_i^T \Rightarrow x = A^{-1} e_i^T$$

$$\Rightarrow x^{(j)} = e_j A^{-1} e_i^T$$

$$\text{By Cramer's rule, } x^{(j)} = \frac{(-1)^{i+j} |A_{ij}|}{|A|} = e_j A^{-1} e_i^T$$

$$\text{By chain rule, } \frac{\partial \log |A|}{\partial a_{ij}} = \frac{1}{|A|} \frac{\partial |A|}{\partial a_{ij}} = \frac{(-1)^{i+j}}{|A|} |A_{ij}| = e_j A^{-1} e_i^T$$

$$\text{let } A = \Sigma, \text{ we get } \frac{\partial \log |\Sigma|}{\partial \sigma_{ij}} = e_j \Sigma^{-1} e_i^T$$

3. By (1), maximum likelihood function = $\prod_{n=1}^N (z_n | N(x_n | \mu_1, \Sigma))^{t_{n1}} (z_n | N(x_n | \mu_2, \Sigma))^{t_{n2}} \dots (z_n | N(x_n | \mu_K, \Sigma))^{t_{nK}}$

$$\log L = \sum_{n=1}^N t_{n1} \log N(x_n | \mu_1, \Sigma) + t_{n2} \log N(x_n | \mu_2, \Sigma) + \dots + t_{nK} \log N(x_n | \mu_K, \Sigma)$$

Then, to maximize likelihood for mean, $\frac{\partial \log L}{\partial \mu_k} = 0$

$$\frac{\partial \log L}{\partial \mu_k} = \sum_{n=1}^N t_{nk} \frac{\partial}{\partial \mu_k} \left(-\frac{k}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - \mu_k)^T \Sigma^{-1} (x_n - \mu_k) \right)$$

$$\rightarrow = \sum_{n=1}^N t_{nk} \frac{\partial}{\partial (x_n - \mu_k)} \left(-\frac{k}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - \mu_k)^T \Sigma^{-1} (x_n - \mu_k) \right) \cdot \frac{\partial (x_n - \mu_k)}{\partial \mu_k}$$

$$\because \frac{\partial N(A)}{\partial A} = 2A$$

If k independent from A and $A^T = A$

$$= \sum_{n=1}^N t_{nk} \left(-\frac{1}{2} \cdot 2 \Sigma^{-1} (x_n - \mu_k) \right) \cdot (-1)$$

$$= \sum_{n=1}^N t_{nk} \Sigma^{-1} (x_n - \mu_k) = 0$$

$$\Rightarrow 0 = \sum_n t_{nk} x_n - \sum_n t_{nk} \mu_k$$

$$\Rightarrow \sum_n t_{nk} x_n = N_k \mu_k \Rightarrow \mu_k = \frac{\sum_{n=1}^N t_{nk} x_n}{N_k}$$

$$\log L(\mu, \Sigma | x_n) = \sum_{n=1}^N C_n + t_{n1} \log N(x_n | \mu_1, \Sigma) + \dots + t_{nK} \log N(x_n | \mu_K, \Sigma)$$

$$\frac{\partial \log L(\mu, \Sigma | x_n)}{\partial \Sigma^{-1}} = \sum_{n=1}^N C_n + t_{n1} \left(C'_1 + \frac{1}{2} \log |\Sigma^{-1}| - \frac{1}{2} \times \text{tr}[(x_n - \mu_1)(x_n - \mu_1)^T \Sigma^{-1}] \right) + \dots + t_{nK} \left(C'_K + \frac{1}{2} \log |\Sigma^{-1}| - \frac{1}{2} \times \text{tr}[(x_n - \mu_K)(x_n - \mu_K)^T \Sigma^{-1}] \right)$$

$$= \sum_{n=1}^N \left(t_{n1} \left(\frac{1}{2} \Sigma - \frac{1}{2} (x_n - \mu_1)(x_n - \mu_1)^T \right) + \dots + t_{nK} \left(\frac{1}{2} \Sigma - \frac{1}{2} (x_n - \mu_K)(x_n - \mu_K)^T \right) \right)$$

$$= 0$$

$$\Rightarrow \sum_{n=1}^N \frac{1}{2} \Sigma (t_{n1} + \dots + t_{nK}) = \frac{1}{2} \sum_{n=1}^N t_{n1} (x_n - \mu_1)(x_n - \mu_1)^T + \dots + t_{nK} (x_n - \mu_K)(x_n - \mu_K)^T$$

$$t_{n1} + \dots + t_{nK} = 1$$

$$\sum_{k=1}^K \frac{N_k}{N} S_k = \#$$

$$\therefore N \cdot \Sigma = \sum_{n=1}^N \sum_{k=1}^K t_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

$$\Rightarrow \Sigma = \frac{1}{N} \cdot \sum_{k=1}^K \sum_{n=1}^N t_{nk} (x_n - \mu_k)(x_n - \mu_k)^T = \sum_{k=1}^K \frac{N_k}{N} \cdot \frac{1}{N_k} \sum_{n=1}^N t_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$