

學號：b05602052 系級：電機四 姓名：舒泓諭

1. (0.5%) 請比較你實作的 generative model、logistic regression 的準確率，何者較佳？

	Public	Private
generative model	0.84606	0.84399
logistic regression	0.85380	0.84964

從以上可以看出 logistic regression 不論在 public score，還是 private score，都來的比 generative model 來的更高。符合老師在上課提到大部分的狀況下 logistic regression 都會來的比 generative model 更好。然而這取決於 training data，是否符合當初 generative 的機率分佈。如果 generative 的 training data 符合當初假設的機率分布，那麼在這個條件下，generative 效果通常會比 logistic regression 來得更好。

2. (0.5%) 請實作特徵標準化(feature normalization)並討論其對於你的模型準確率的影響

logistic regression:

	Public	Private
Non-feature normalization	0.84815	0.84277
feature normalization	0.85380	0.84964

Generative model:

	Public	Private
Non-feature normalization	0.84373	0.84289
feature normalization	0.84606	0.84399

有做 feature normalization 的話不論是 logistic regression 還是 Generative model，public score 和 private score 都會上升。然而如果不做 feature normalization 的話 Private score 連 simple baseline 都過不了，可見 feature normalization 的重要性。

	Change of Public	Change of Private
logistic regression	5.65*e-3	6.87*e-3
Generative model	2.33*e-3	1.1*e-3

由上表可以發現 feature normalization，對於 logistic regression 的影響會比 Generative model 來的更高，主要原因在於 logistic regression，再找 w, b 值時有進行 gradient descent，然而 Generative model 只是純粹從機率分佈去推估 w, b 。所以可以看出符合老師在上課提到的 feature normalization，可以有效解決當 feature range 不一樣時，gradient 會前後 oscillate 的問題。

3. (1%) 請說明你實作的 best model，其訓練方式和準確率為何？

Data preprocessing:

因為經過 one hot encoding 之後，大部分的 feature，都變為 0 和 1，只有 age,fnlwgt,capital_gain,capital_loss,hours_per_week，的數值不是離散的，所以只有在這幾個 feature 進行 normalization。

Model training:

本次使用的 model 只疊了一層 sklearn GradientBoostingClassifier，訓練參數如下：

	parameter
Loss function	logistic regression
learning_rate	0.1
n_estimators(epoch)	100
max_depth	3
min_impurity_split (early stopping)	1*e-7

Predict outcome:

	Training accuracy	Public	Private
generative model	0.841927	0.84606	0.84399
logistic regression	0.849973	0.85380	0.84964
gradient boosting	0.866689	0.86977	0.86549

原本的兩個 model 在同樣的 data preprocessing 的方法下，不管怎麼樣都過不了 strong baseline，然而使用 GradientBoosting 的方式後，一下子就過了 strong baseline。

此外，如果從這三種方法進行比較我們可以發現，從 Training accuracy 來看，gradient boosting，在 training data fit 的比起另外兩種方式來得更好。然而從 public score 和 private score 來看，gradient boosting 也表現的比另外兩種方式來的更好，所以沒有發現 overfitting 的現象。

手寫作業：

(1)

No. Date

$$P(X, Y) = \prod_{k=1}^K P(X|C_k) \pi_k$$

$$L(X, Y) = \prod_{n=1}^N \prod_{k=1}^K P(X_n|C_k) \pi_k^{t_{n,k}}$$

$$\ell(\theta) = \log L(X, Y) = \sum_{n=1}^N \sum_{k=1}^K t_{n,k} (\log P(X_n|C_k) + \log \pi_k)$$

$$\ell(\pi, \lambda) = \sum_{n=1}^N \sum_{k=1}^K t_{n,k} (\log P(X_n|C_k) + \log \pi_k) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial \ell(\pi, \lambda)}{\partial \pi_k} = \frac{1}{\pi_k} \sum_{n=1}^N t_{n,k} + \lambda = 0 \Rightarrow \pi_k = -\frac{1}{\lambda} \sum_{n=1}^N t_{n,k}$$

$$= -\frac{1}{\lambda} N_k$$

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

$$\Rightarrow \lambda = -\sum_{k=1}^K N_k$$

$$\Rightarrow \pi_k = -\frac{N_k}{\lambda} = \frac{N_k}{\sum_{k=1}^K N_k} = \frac{N_k}{N} \quad \text{(Q.E.D.)}$$

$|N_{\text{total}}|$ number of data #

(2)

$$\text{adj}(\Sigma) = \begin{pmatrix} M_{11} & M_{12} & \dots \\ M_{21} & M_{22} & \dots \\ \vdots & \vdots & \ddots \\ M_{m1} & M_{m2} & \dots \end{pmatrix}$$

M_{ij} : cofactor matrix (removing i th row and j th column of Σ)

$$\begin{aligned} \frac{\partial \log(\det \Sigma)}{\partial \sigma_{ij}} &= \frac{1}{\det \Sigma} \frac{\partial (\det \Sigma)}{\partial \sigma_{ij}} \quad \left. \begin{array}{l} \text{cofactor} \\ \text{expansion} \end{array} \right. \\ &= \frac{1}{\det \Sigma} \frac{\partial (\Sigma (-1)^{ij} \sigma_{ij} M_{ij})}{\partial \sigma_{ij}} \\ &= \frac{1}{\det \Sigma} (-1)^{ij} M_{ij} \quad \downarrow j \text{ th row.} \end{aligned}$$

$$\text{右式} = e_j \Sigma^{-1} e_j^T = e_j \frac{\text{adj}(\Sigma)}{\det \Sigma} e_j^T =$$

$$\begin{aligned} \left(\begin{array}{c} \text{adj}(\Sigma) \\ \vdots \\ M_{ij} \\ \vdots \\ M_{mm} \end{array} \right) &= \frac{(-1)^{ij}}{\det \Sigma} \left(\begin{array}{c} e_j^T \\ \vdots \\ M_{ij} \\ \vdots \\ e_j^T \end{array} \right) \quad (M_{ij}: i \text{ th row of adj}(\Sigma)) \\ &= \frac{1}{\det \Sigma} (-1)^{ij} M_{ij} \end{aligned}$$

$$\Sigma^{-1} = \frac{\text{adj}(\Sigma)}{\det(\Sigma)} \quad \text{左式} = \text{右式} \quad Q.E.D. \quad \#$$

$$\begin{aligned} \Sigma \sum_{k=1}^K N_{ik} &= \sum_{k=1}^K \sum_{n=1}^N t_{nk} (x_n - M_{ik}) (x_n - M_{ik})^T \\ \therefore \Sigma &= \frac{1}{N} \sum_{k=1}^K N_k \sum_{n=1}^N t_{nk} (x_n - M_{ik}) (x_n - M_{ik})^T \end{aligned}$$

(3)

$$\begin{aligned}
 f(x_n | c_k) &= \frac{1}{(2\pi)^{\frac{N}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x_n - M)^T \Sigma^{-1} (x_n - M)\right) \\
 P(\Sigma | M) &= \prod_{k=1}^K \left(P(x_n | c_k) \pi_k \right)^{t_{n,k}} \\
 L(\Sigma | M) &= \prod_{n=1}^N \prod_{k=1}^K \left(P(x_n | c_k) \pi_k \right)^{t_{n,k}} \\
 l(\theta) &= \log L(\Sigma, M) = \sum_{n=1}^N \sum_{k=1}^K t_{n,k} \left[\log f(x_n | c_k) + \log \pi_k \right] \\
 &= \sum_{n=1}^N \sum_{k=1}^K t_{n,k} \left(-\frac{1}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - M)^T \Sigma^{-1} (x_n - M) + \log \pi_k \right) \\
 \frac{\partial l(\theta)}{\partial M_k} &= 0 = \sum_{n=1}^N t_{n,k} \left(-\frac{1}{2} \Sigma^{-1} (x_n - M) \right) \stackrel{(1)}{=} x - 1 \\
 \Rightarrow \sum_{n=1}^N t_{n,k} \left(\Sigma^{-1} (x_n - M) \right) &= 0 \\
 \Rightarrow \sum_{n=1}^N (x_n - M) t_{n,k} - M t_{n,k} &= 0 \\
 \Rightarrow \sum_{n=1}^N x_n t_{n,k} &= M t_{n,k} \stackrel{(2)}{=} M_k \cdot N_k \\
 \Rightarrow M_k &= \frac{1}{N_k} \sum_{n=1}^N x_n t_{n,k} \quad (Q.E.D.) \\
 \left(\frac{\log \det \Sigma^{-1}}{\partial \Sigma^{-1}} = \frac{1}{2} \Sigma = \Sigma^{-1} \right) \quad \left(\because (x_n - M)^T \Sigma^{-1} (x_n - M) \text{ is scalar,} \right. \\
 &\quad \left. \therefore \text{Trace}((x_n - M)^T \Sigma^{-1} (x_n - M)) \right) \\
 l(\theta) &= \sum_{n=1}^N \sum_{k=1}^K t_{n,k} \left(-\frac{1}{2} \log 2\pi + \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_n - M)^T \Sigma^{-1} (x_n - M) \right) \\
 \frac{\partial l(\theta)}{\partial \Sigma^{-1}} &= \sum_{n=1}^N \sum_{k=1}^K \frac{t_{n,k}}{2} \Sigma - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K t_{n,k} \text{Trace}((x_n - M)^T \Sigma^{-1} (x_n - M)) \cdot t_{n,k} \\
 &= \frac{1}{2} \sum_{k=1}^K \left(\sum_{n=1}^N t_{n,k} \right) \Sigma - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (x_n - M)^T (x_n - M) t_{n,k} = 0 \\
 &= \sum_{k=1}^K N_k \frac{\Sigma}{2} - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (x_n - M)^T (x_n - M) t_{n,k} = 0
 \end{aligned}$$

$$\begin{aligned}
 \bar{\Sigma} \sum_{k=1}^K N_k &= \sum_{k=1}^K \sum_{n=1}^N t_{nk} (x_n - M_k) (x_n - M_k)^T \\
 \Rightarrow \bar{N} \cdot \sum_{k=1}^K N_k &= \sum_{k=1}^K N_k \sum_{n=1}^N \frac{t_{nk}}{N_k} (x_n - M_k) (x_n - M_k)^T \\
 \Rightarrow \bar{\Sigma} &= \sum_{k=1}^K \frac{N_k}{N} \underbrace{\frac{1}{N_k} \sum_{n=1}^N (x_n - M_k) (x_n - M_k)^T}_{= \frac{1}{N} \sum_{k=1}^K N_k \cdot S_k} \\
 &= \sum_{k=1}^K \frac{N_k}{N} \cdot S_k. \quad \text{Q.E.D.}
 \end{aligned}$$

NO. _____ Date _____

$$\left. \begin{aligned}
 \frac{\partial W^T A W}{\partial w_i} &= 2 A W. \\
 \frac{\partial (\frac{1}{2} (x_n - M_k)^T \bar{\Sigma}^{-1} (x_n - M_k))}{\partial M_k} &= -2 \bar{\Sigma}^{-1} (x_n - M_k)
 \end{aligned} \right) \quad \textcircled{1}$$

$\therefore x^T A x$ is scalar

$$\left. \begin{aligned}
 \therefore \frac{\partial}{\partial n} x^T A x &= \frac{\partial}{\partial A} \text{tr}(x^T A x) = (x x^T)^T = x^T x^T \\
 \frac{\partial}{\partial \bar{\Sigma}^{-1}} \text{tr} \text{re}[(x_n - M_k)^T \bar{\Sigma}^{-1} (x_n - M_k)] &= (x_n - M_k) (x_n - M_k)^T
 \end{aligned} \right) \quad \textcircled{2}$$