

Machine Learning HW1

土木五 B08501011 何子勤

Programming part:

1. 解釋什麼樣的 data preprocessing 可以 improve 你 training/testing accuracy。請提供數據(例如 kaggle public score RMSE)以佐證你的想法。

挑選可能與 PM2.5 相關的數據，從數值組成的長條圖挑選和 PM2.5 相似的數據組，再挑出粒子類的項目，最後由繪製折線圖觀察與未來 PM2.5 相似的項目。

	my_sol.csv Complete · 21s ago · with 2, 3, 6, 14	3.49722	
	my_sol.csv Complete · 15m ago · all factor in second order	5.00581	

2. 請實作 2nd-order polynomial regression model (不用考慮交互項)。(1%)

貼上 polynomial regression 版本的 Gradient descent code 內容

```
for num in range(epoch):
    for b in range(int(x.shape[0]/batch_size)):
        t+=1
        x_batch = x[b*batch_size:(b+1)*batch_size]
        y_batch = y[b*batch_size:(b+1)*batch_size].reshape(-1,1)
        # Prediction of linear regression
        pred = np.dot(np.square(x_batch), w2) + np.dot(x_batch, w1) + bias
        # loss
        loss = y_batch - pred
        # Compute gradient
        g_t1 = np.dot(x_batch.transpose(), loss) * (-2)
        g_t2 = np.dot(np.square(x_batch).transpose(), loss) * (-2)
        g_t_b = loss.sum(axis=0) * (-2)
        m_t1 = beta_1*m_t1 + (1-beta_1)*g_t1
        m_t2 = beta_1*m_t2 + (1-beta_1)*g_t2
        v_t1 = beta_2*v_t1 + (1-beta_2)*np.multiply(g_t1, g_t1)
        v_t2 = beta_2*v_t2 + (1-beta_2)*np.multiply(g_t2, g_t2)
        m_cap1 = m_t1/(1-(beta_1**t))
        m_cap2 = m_t2/(1-(beta_1**t))
        v_cap1 = v_t1/(1-(beta_2**t))
```

```

v_cap2 = v_t2/(1-(beta_2**t))
m_t_b = 0.9*m_t_b + (1-0.9)*g_t_b
v_t_b = 0.99*v_t_b + (1-0.99)*(g_t_b*g_t_b)
m_cap_b = m_t_b/(1-(0.9**t))
v_cap_b = v_t_b/(1-(0.99**t))

# Update weight & bias
w1 -= ((lr*m_cap1)/(np.sqrt(v_cap1)+epsilon)).reshape(-1, 1)
w2 -= ((lr*m_cap2)/(np.sqrt(v_cap2)+epsilon)).reshape(-1, 1)
bias -= (lr*m_cap_b)/(math.sqrt(v_cap_b)+epsilon)

```

(b) 在只使用 NO 數值作為 feature 的情況下，紀錄該 model 所訓練出的 parameter 數值(w2, w1, b)以及 kaggle public score.



my_sol.csv

Complete · 8d ago · polynomial with only NO (corrected)

5.42307



```

w2:
[[0.00932009]
 [0.00618544]
 [0.00431545]
 [0.00440361]
 [0.00486678]
 [0.00545708]
 [0.00825261]
 [0.0139166 ]]
w1:
[[0.13713366]
 [0.13115987]
 [0.12713296]
 [0.12231566]
 [0.12191603]
 [0.12310089]
 [0.12555388]
 [0.13238816]]
bias:
[0.36294667]

```

1. (a)

$$\nabla_W W^T A W = \frac{(W+\Delta W)^T A (W+\Delta W) - W^T A W}{(\Delta W)^T} = \frac{\Delta W^T A W + W^T A \Delta W + (\Delta W)^T A \Delta W}{(\Delta W)^T}$$

$$\therefore f(W+\Delta W) - f(W) \approx (\Delta W)^T \nabla f(W)$$

$$\Rightarrow (\Delta W)^T A W + (\Delta W)^T (W^T A)^T \Rightarrow \nabla_W W^T A W = A W + (W^T A)^T = A W + A^T W$$

$$\text{If } A \text{ is symmetric } A^T = A \Rightarrow \nabla_W W^T A W = 2 A W$$

(b)

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mm} \end{bmatrix} \quad B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1m} \\ b_{21} & b_{22} & \dots & b_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & \dots & b_{mm} \end{bmatrix} \quad \text{tr}(AB) = \sum_{i=1}^m \sum_{j=1}^m a_{ij} b_{ji} \quad \frac{\partial \text{tr}(AB)}{\partial a_{ij}} = b_{ji}$$

$$AB = \begin{bmatrix} (a_{11}b_{11} + a_{12}b_{21} + \dots) & \vdots \\ (a_{21}b_{11} + a_{22}b_{21} + \dots) & \vdots \\ \vdots & \vdots \end{bmatrix} \quad \text{a}_{ij} \text{ 係數即 } b_{ji}$$

(c)

$$\det(A) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(A_{ij}) \quad \frac{\partial}{\partial a_{ij}} \det(A) = (-1)^{i+j} \det(A_{ij})$$

又根據 Cramer's rule

$$A x = e_i \Rightarrow x_j = e_j^T x = e_j^T A^{-1} e_i = \frac{(-1)^{i+j} \det(A_{ij})}{\det(A)} = \frac{\partial \log(\det(A))}{\partial a_{ij}}$$

2. (a)(i)

$$L(\theta) = \prod_{i=1}^N p(x_i | \theta) = \prod_{i=1}^N [\pi_1 f_{\mu_1, \Sigma_1}(x_i)]^{\hat{y}_i} [\pi_2 f_{\mu_2, \Sigma_2}(x_i)]^{1-\hat{y}_i} \quad \text{Assume } \hat{y}_i = \begin{cases} 1 & \text{if } y_i = 1 \\ 0 & \text{if } y_i = 2 \end{cases}$$

(ii)

$$\ln L(\theta) = \sum_{i=1}^N \ln p(x_i | \theta) = \sum_{i=1}^N \hat{y}_i \ln \pi_1 f_{\mu_1, \Sigma_1}(x_i) + (1-\hat{y}_i) \ln \pi_2 f_{\mu_2, \Sigma_2}(x_i)$$

$$= \sum_{i=1}^N \left\{ \hat{y}_i \left[\ln \pi_1 - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} (x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1) \right] + (1-\hat{y}_i) \left[\ln \pi_2 - \frac{1}{2} \ln |\Sigma_2| - \frac{1}{2} (x_i - \mu_2)^T \Sigma_2^{-1} (x_i - \mu_2) \right] \right\}$$

 $\ln L(\theta)$ depend on μ

$$\Rightarrow -\frac{1}{2} \sum_{i=1}^N \hat{y}_i (x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1)$$

 Σ^{-1} is symmetric

$$\Rightarrow \sum_{i=1}^N \hat{y}_i \Sigma_1^{-1} (x_i - \mu_1) = 0$$

取 Σ_1 及利用 1.(b)(c) Σ^{-1} is symmetric

$$\Rightarrow -\frac{N_1}{2} (\Sigma_1^{-1}) = -\frac{N_1}{2} \left(\sum_{i=1}^N \hat{y}_i (x_i - \mu_1)(x_i - \mu_1)^T \right)^{-1}$$

$$\Rightarrow N_1 (\Sigma_1) = \sum_{i=1}^N \hat{y}_i (x_i - \mu_1)(x_i - \mu_1)^T$$

 $\ln L(\theta)$ depend on Σ_1

$$\Rightarrow -\frac{1}{2} \sum_{i=1}^N \hat{y}_i \ln |\Sigma_1| - \frac{1}{2} \sum_{i=1}^N \hat{y}_i (x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1)$$

 $\ln L(\theta)$ depend on π_1

$$\Rightarrow \sum_{i=1}^N \hat{y}_i \ln \pi_1 + (1-\hat{y}_i) \ln \pi_2 \Rightarrow \nabla_{\pi_1} \sum_{i=1}^N \hat{y}_i \ln \pi_1 + (1-\hat{y}_i) \ln \pi_2 = 0 \Rightarrow \frac{N_1}{\pi_1} - \frac{N_2}{1-\pi_1} = 0$$

$$\Rightarrow \pi_1 + \pi_2 = 1$$

$$\Rightarrow N_1 - N_1 \pi_1 - N_2 \pi_1 = 0$$

$$\Rightarrow \begin{cases} \pi_1^* = \frac{N_1}{N} \\ \pi_2^* = \frac{N_2}{N} \end{cases}$$

2 (a)(iii)

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} = \frac{\pi_1 f_{\mu_1, \Sigma_1}(x)}{\pi_1 f_{\mu_1, \Sigma_1}(x) + \pi_2 f_{\mu_2, \Sigma_2}(x)} \rightarrow \text{給定 } x, \text{ 求 } x \text{ 為 Class 1 機率}$$

$$P(x|C_1) = f_{\mu_1, \Sigma_1}(x) \Rightarrow x \text{ 在 Class 1 的 gaussian distribution 中對應的機率}$$

$$(iv) P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)} = \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}} \xRightarrow{\text{令 } z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}} \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} z &= \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} = \ln \frac{\frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{|\Sigma_1|^{\frac{1}{2}}} \exp\{-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1)\} P(C_1)}{\frac{1}{(2\pi)^{\frac{1}{2}}} \frac{1}{|\Sigma_2|^{\frac{1}{2}}} \exp\{-\frac{1}{2}(x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)\} P(C_2)} \\ &= \ln \frac{|\Sigma_2|^{\frac{1}{2}}}{|\Sigma_1|^{\frac{1}{2}}} - \frac{1}{2} [(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)] + \ln \frac{P(C_1)}{P(C_2)} \\ &= \ln \frac{|\Sigma_2|^{\frac{1}{2}}}{|\Sigma_1|^{\frac{1}{2}}} - \frac{1}{2} [(x-\mu_1)^T \Sigma_1^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)] + \ln \left(\frac{\pi_1}{\pi_2} \right) \end{aligned}$$

3. (b)(i)

$$\text{similar to (a)(i)} \quad L(\theta) = \prod_{i=1}^N [\pi_1 f_{\mu_1, \Sigma_1}(x_i)]^{\hat{y}_i} [\pi_2 f_{\mu_2, \Sigma_2}(x_i)]^{1-\hat{y}_i}$$

(ii)

μ_1^* and μ_2^* same with (a)(ii)

$$\ln L(\theta) \text{ depend on } \Sigma \Rightarrow -\frac{1}{2} \sum_{i=1}^N \hat{y}_i \ln |\Sigma_1| - \frac{1}{2} \sum_{i=1}^N \hat{y}_i (x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1) - \frac{1}{2} \sum_{i=1}^N (1 - \hat{y}_i) \ln |\Sigma_2| - \frac{1}{2} \sum_{i=1}^N (1 - \hat{y}_i) (x_i - \mu_2)^T \Sigma_2^{-1} (x_i - \mu_2)$$

$$\Rightarrow -\sum_{i=1}^N \hat{y}_i \ln |\Sigma_1| - \frac{1}{2} \sum_{i=1}^N \hat{y}_i (x_i - \mu_1)^T \Sigma_1^{-1} (x_i - \mu_1) - \sum_{i=1}^N (1 - \hat{y}_i) \ln |\Sigma_2| - \frac{1}{2} \sum_{i=1}^N (1 - \hat{y}_i) (x_i - \mu_2)^T \Sigma_2^{-1} (x_i - \mu_2)$$

(Σ_1 and Σ_2 are calculated in (a)(ii))

π_1^* and π_2^* same with (a)(ii)

(iii)

$$P(C_1|x) = \frac{\pi_1 f_{\mu_1, \Sigma}(x)}{\pi_1 f_{\mu_1, \Sigma}(x) + \pi_2 f_{\mu_2, \Sigma}(x)} \Rightarrow \text{給定 } x, \text{ 求 } x \text{ 為 Class 1 的機率}$$

$$P(x|C_1) = f_{\mu_1, \Sigma}(x) \Rightarrow x \text{ 在 Class 1 的 gaussian distribution 中對應的機率}$$

(iv)

$$\Sigma_1 = \Sigma_2 = \Sigma$$

$$\begin{aligned} z &= -\frac{1}{2} [(x-\mu_1)^T \Sigma^{-1} (x-\mu_1) - (x-\mu_2)^T \Sigma^{-1} (x-\mu_2)] + \ln \left(\frac{\pi_1}{\pi_2} \right) \\ &= (\mu_1 - \mu_2)^T \Sigma^{-1} x - \frac{1}{2} (\mu_1)^T \Sigma^{-1} \mu_1 + \frac{1}{2} (\mu_2)^T \Sigma^{-1} \mu_2 + \ln \left(\frac{\pi_1}{\pi_2} \right) \end{aligned}$$

reference:

Pattern recognition and machine learning
Christopher M Bishop

3.

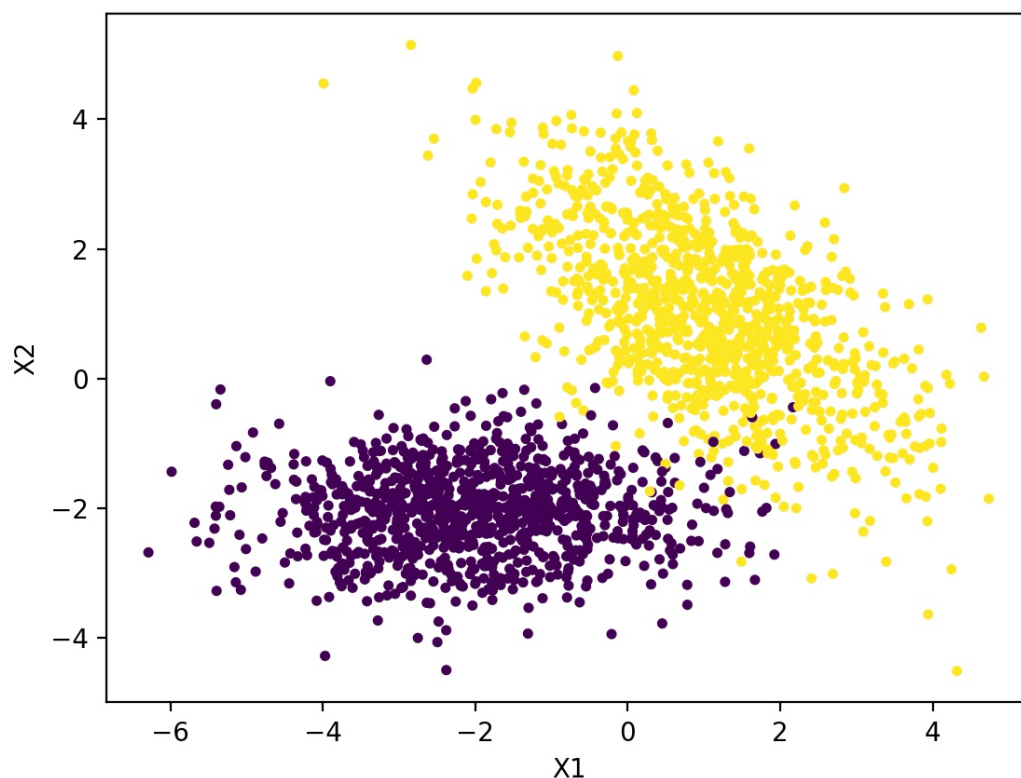
(a)

```
N = 2000 N1 = 1000 N2 = 1000
pi1 = 0.5 pi2 = 0.5
C1avg =
[[-2.02571697]
 [-2.04619501]]
C2avg =
[[1.01143637]
 [1.00493194]]
Cvar =
[[ 1.85889712 -0.51610136]
 [-0.51610136  1.14373928]]
```

(b)

```
N = 2000 N1 = 1000 N2 = 1000
pi1 = 0.5 pi2 = 0.5
C1avg =
[[-2.02571697]
 [-2.04619501]]
C2avg =
[[1.01143637]
 [1.00493194]]
C1Var =
[[2.01130388  0.03386452]
 [0.03386452  0.45977354]]
C2Var =
[[ 1.70649036 -1.06606724]
 [-1.06606724  1.82770502]]
```

(c)



我認為(b)的表現會比較好，因為 Class 1 和 Class2 的分佈走向顯示兩個種類的 covariance matrix 並不相同，Class 1 的兩變數 covariance 應接近於零，但 Class 兩變數應有一定程度相關，故兩者的 covariance matrix 不應假設相同。

Code:

```
import numpy as np
import matplotlib.pyplot as plt

# datashape = (2000, 3)
data = np.load("./data.npy")

# class count
N = 0
N1 = 0
N2 = 0
for i in range(2000):
    N += 1
    if data[i, 2] == 0:
        N1 += 1
    else:
        N2 += 1
print("N = " + str(N), "N1 = " + str(N1), "N2 = " + str(N2))
print("pi1 = " + str(N1 / N), "pi2 = " + str(N2 / N))

# calculate average
X1 = np.array([data[:, 0] * (1 - data[:, 2]), data[:, 1] * (1 - data[:, 2])])
X2 = np.array([data[:, 0] * data[:, 2], data[:, 1] * data[:, 2]])
C1avg = (np.sum(X1, axis = 1) / N1).reshape(2, 1)
C2avg = (np.sum(X2, axis = 1) / N2).reshape(2, 1)
print("C1avg = ")
print(C1avg)
print("C2avg = ")
print(C2avg)

# calculate variance
X1 = np.array([(data[:, 0] - C1avg[0]) * (1 - data[:, 2]), (data[:, 1] - C1avg[1]) * (1 - data[:, 2])])
X2 = np.array([(data[:, 0] - C2avg[0]) * data[:, 2], (data[:, 1] - C2avg[1]) * data[:, 2]])
C1Var = np.matmul(X1, np.transpose(X1)) / N1
C2Var = np.matmul(X2, np.transpose(X2)) / N2
```

```
'''  
print("C1Var = ")  
print(C1Var)  
print("C2Var = ")  
print(C2Var)  
'''  
  
CVar = N1 / N * C1Var + N2 / N * C2Var  
print("Cvar = ")  
print(CVar)  
# plot data  
plt.scatter(data[:, 0], data[:, 1], marker = '.', c = data[:, 2])  
plt.xlabel("X1")  
plt.ylabel("X2")  
plt.show()
```


$$4. (a) \sum_i k_i (y_i - x_i^T \theta)^2 + \lambda \sum_j \omega_j^2 = (y - X\theta)^T K (y - X\theta) + \lambda \theta^T \theta$$

$$= y^T K y - y^T K X \theta - \theta^T X^T K y + \theta^T X^T K X \theta + \lambda \theta^T \theta$$

$$\nabla_{\theta} \theta = \theta + \Delta \theta \text{ 找 } (\Delta \theta)^T \text{ 係數 } \Rightarrow -(y^T K X)^T - X^T K y + 2 X^T K X \theta + \lambda 2 \theta = 0$$

$$K^T = K \Rightarrow -2 X^T K y + 2 X^T K X \theta + \lambda 2 \theta = 0 \Rightarrow (X^T K X + \lambda I) \theta = X^T K y$$

$$\Rightarrow \theta^* = (X^T K X + \lambda I)^{-1} X^T K y$$

5.

$$\tilde{L}_{SS}(w, b) = E \left[\frac{1}{2N} \sum_{i=1}^N (f_{w,b}(x_i + \eta_i) - y_i)^2 \right]$$

$$= \frac{1}{2N} E \left[\sum_{i=1}^N (w^T (x_i + \eta_i) + b - y_i)^2 \right] = \frac{1}{2N} E \left[\sum_{i=1}^N (w^T x_i + w^T \eta_i + b - y_i)(w^T x_i + w^T \eta_i + b - y_i) \right]$$

$$= \frac{1}{2N} E \left[\sum_{i=1}^N (f_{w,b}(x_i) - y_i)^T (f_{w,b}(x_i) - y_i) + (w^T \eta_i)^T (f_{w,b}(x_i) - y_i) + (f_{w,b}(x_i) - y_i)^T w^T \eta_i + (w^T \eta_i)^T w^T \eta_i \right]$$

$$E(\eta_i) = 0, \text{ 中間兩項為 } 0 \quad \eta_i^T \cdot f_{w,b}(x_i) = 0$$

$$\Rightarrow = \frac{1}{2N} \sum_{i=1}^N E[(f_{w,b}(x_i) - y_i)^2] + \frac{1}{2N} \sum_{i=1}^N E(\eta_i^T w w^T \eta_i)$$

$$\eta_i^T w, w^T \eta_i \text{ 均 } 1 \times 1 \text{ matrix}$$

$$\Rightarrow = \frac{1}{2N} \sum_{i=1}^N E(f_{w,b}(x_i) - y_i)^2 + \frac{1}{2N} \sum_{i=1}^N E(w^T \eta_i)(\eta_i^T w) = \frac{1}{2N} \sum_{i=1}^N E(f_{w,b}(x_i) - y_i)^2 + \frac{1}{2N} \sum_{i=1}^N E(w^T \eta_i \eta_i^T w)$$

$$\sum_{i=1}^N E(\eta_i \eta_i^T) = N \sigma^2 I \quad w^T w = \|w\|^2$$

$$\Rightarrow = \frac{1}{2N} \sum_{i=1}^N E(f_{w,b}(x_i) - y_i)^2 + \frac{N \sigma^2}{2N} \|w\|^2 = \frac{1}{2N} \sum_{i=1}^N E(f_{w,b}(x_i) - y_i)^2 + \frac{\sigma^2}{2} \|w\|^2$$