

# Machine Learning

## Homework 1

系級：醫工三    學號：B12508025    姓名：盧雅筠

## Problem 1: cross-validation, feature selection, and classification

### I. Reports

#### 1. The multivariate Gaussian Distribution models and Parameter Estimation

- For each class  $y \in \{0,1\}$ ,  $\mathbf{x} \in \mathbb{R}^d$  assume the multivariate Gaussian distribution:

$$\mathbf{x} \mid y = c \sim \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c), \quad c \in \{0,1\}$$

$$p(\mathbf{x} \mid y = c) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_c|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \right]$$

where  $\boldsymbol{\mu}_c$  and  $\boldsymbol{\Sigma}_c$  denote the mean vector and covariance matrix of class  $c$ .

- Parameter Estimation: maximum likelihood estimation (MLE) use

`np.mean(X_c, axis=0)` and `cov = np.cov(X_c, rowvar=False)`

to compute:

$$\ell(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) = \sum_{i:y_i=c} \log p(\mathbf{x}_i \mid \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$$

$$\hat{\boldsymbol{\mu}}_c = \frac{1}{N_c} \sum_{i:y_i=c} \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}}_c = \frac{1}{N_c - 1} \sum_{i:y_i=c} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)^T$$

- Posterior Probability: use

`Multivariate_Gaussian_Distribution_likelihood()` and

`Bayesian_decision_classifier()` to compute:

$$P(y = 1 \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid y = 1)p(y = 1)}{p(\mathbf{x} \mid y = 1)p(y = 1) + p(\mathbf{x} \mid y = 0)p(y = 0)}$$

- Discriminant Function: use

`Multivariate_Gaussian_Distribution_likelihood()` and

`Bayesian_decision_classifier()` to compute:

$$g_c(\mathbf{x}) = \ln P(\mathbf{x} \mid y = c) + \ln p(y = c)$$

$$\Delta(\mathbf{x}) = g_1(\mathbf{x}) - g_0(\mathbf{x}), \text{ the decision boundary is } \Delta(\mathbf{x}) = 0$$

## 2. Forward Feature Selection and its cost Function

- Feature selection:

use `Forward_Selection()` and its cost function `BIC_score()`

Starting from an empty feature set, at each iteration the algorithm tests all remaining features and adds the one that minimizes the Bayesian Information Criterion (BIC):

$$\text{BIC} = k \ln n - 2 \ln(\hat{L})$$

where  $n$  is the number of training samples,  $k$  is the number of parameters, and  $\hat{L}$  is the maximized likelihood.

- Determine the number of features to be selected

A minimum of 3 features is always kept to ensure sufficient discriminative information. The maximum features are determined by

$$\text{max\_features} = \min(\max(1, \min(n_0, n_1) - 1), d)$$

where  $n_0, n_1$  is the number of two classes samples,  $d$  is  $X_{\text{train}}. \text{shape}[1]$

## 3. Selected Features per Fold

- Each of the 103 data samples serves as a test point once in leave-one-out cross-validation, yielding 103 different feature subsets. The selected features for each fold are stored in `selected_features`. The list of selected features per fold is shown in

Figure 1

[illegible][illegible][illegible]

Figure 1

#### 4. The Test Performance

- Accuracy, Sensitivity, Specificity, AUC are shown in Figure 2

```
Test Performance
Confusion Matrix:
[[56  6]
 [24 17]]
Accuracy=0.709, Sensitivity=0.415, Specificity=0.903, AUC=0.751
Top-2 features: [('LipD', 103), ('UVermilionH', 103)]
```

Figure 2

- The ROC curve is shown in Figure 3

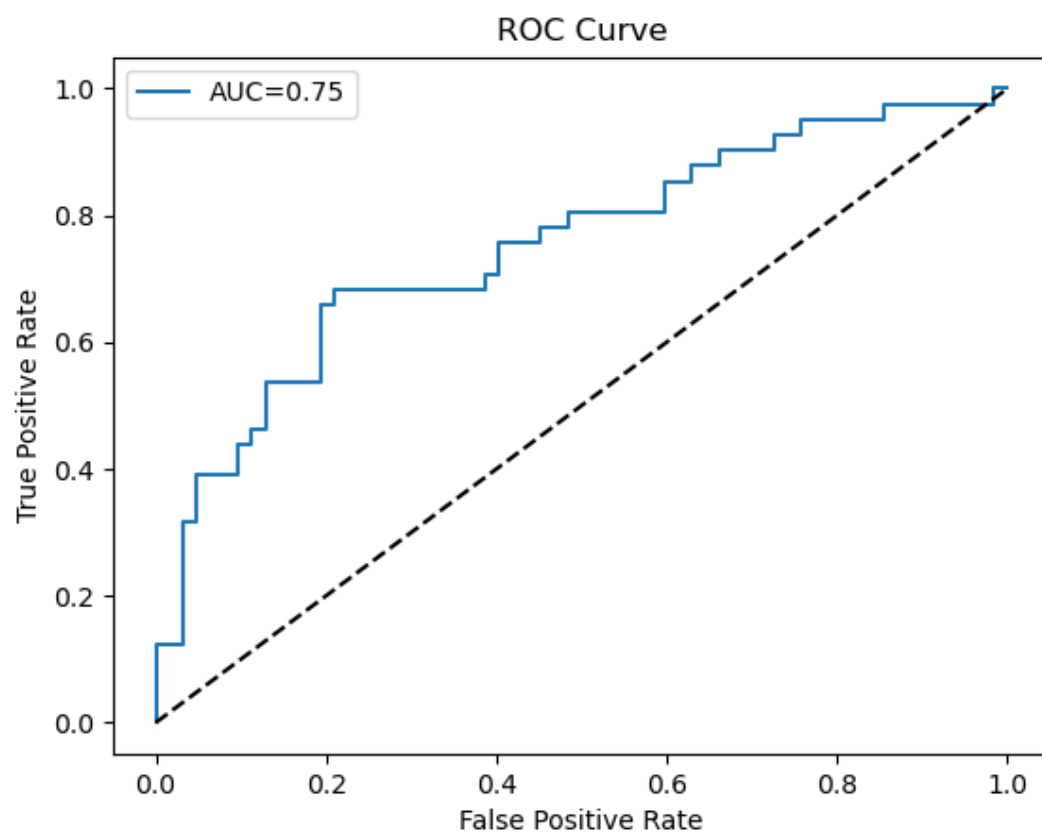


Figure 3

## 5. Top-2 Most Frequently Selected Features and Bivariate Gaussian Bayesian Model Visualization

- Top-2 Most Frequently Selected Features use

`counts = Counter()` and `counts.most_common(2)` to find ,  
shown in Figure 4

```
Top-2 features: [('LipD', 103), ('UVermilionH', 103)]
```

Figure 4

- Bivariate Gaussian Bayesian Model Visualization is shown in Figure 5
  - Scatter plot: class 0 as circles o(gray), class 1 as plus signs +(green)
  - Contour maps: equal-probability contours of each class distribution
  - Decision boundary: the curve where  $\Delta(\mathbf{x}) = 0$

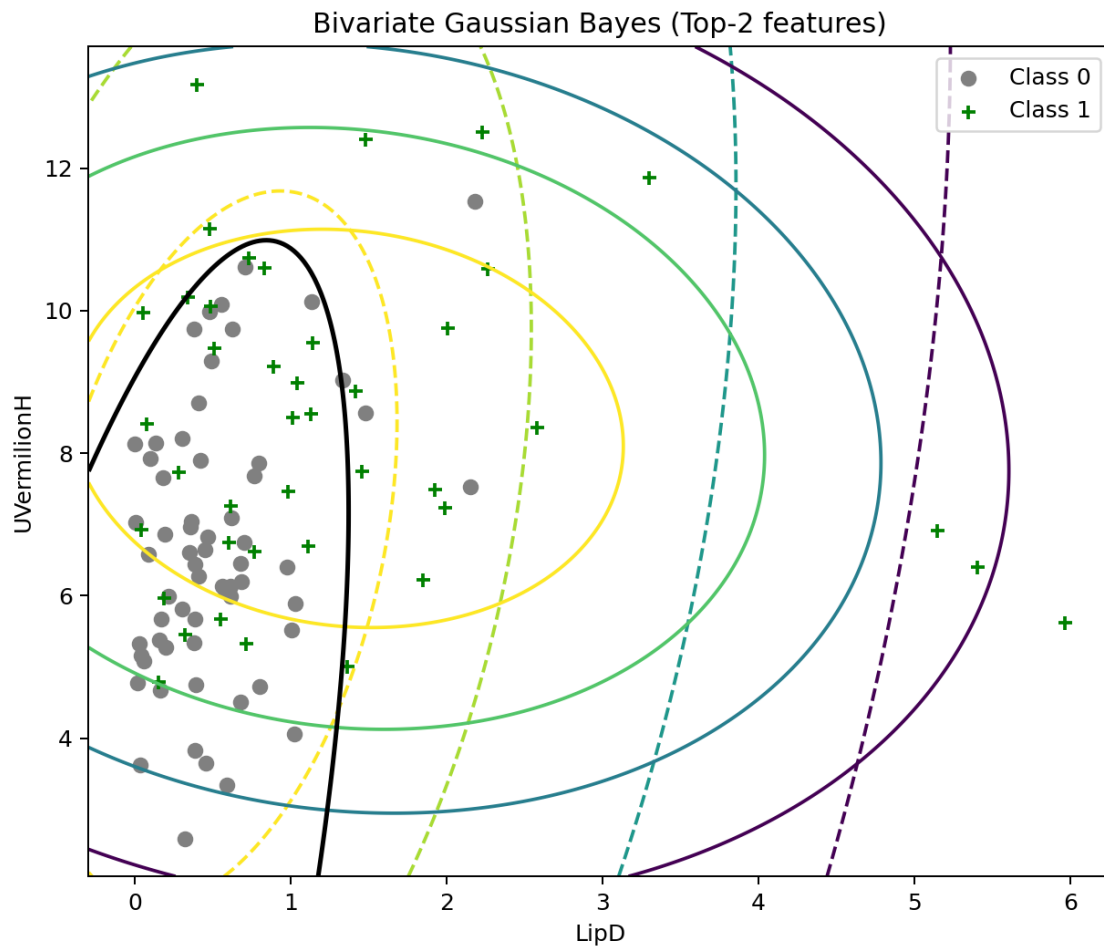


Figure 5

## II. Code

Set up the environment `environment.yml` and run `1.py` shown in Figure 7, making sure to remove unnecessary information from `AcromegalyFeatureSet.xlsx` (eg. Figure 6)

98	97	1	2	153.062	154.967
99	98	1	2	137.798	135.533
100	99	1	2	151.329	150.528
101	100	1	2	153.679	152.093
102	101	1	1	175.044	166.512
103	102	1	1	161.504	160.874
104	103	1	2	156.574	148.620
105					
106		0: control	1: male		
107		1: patient	2: female		
108					
109					
110					
111					
112					
113					

Figure 6. The blue rectangle must be deleted

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.metrics import roc_curve, auc, confusion_matrix
4 import matplotlib.pyplot as plt
5 from collections import Counter
6
7 ''' Multivariate Gaussian Bayesian classifier '''
8 def Multivariate_Gaussian_Distribution_likelihood(x, mu, cov, ridge=1e-6):
9     d = len(x)
10    cov = cov + np.eye(d) * ridge
11    diff = x - mu
12    likelihood = np.exp(-0.5 * diff @ np.linalg.inv(cov) @ diff) / np.sqrt((2 * np.pi)**d * np.linalg.det(cov))
13    return likelihood
14
15 def Bayesian_decision_classifier(x, mu0, mu1, cov0, cov1, p0, p1):
16    l0 = Multivariate_Gaussian_Distribution_likelihood(x, mu0, cov0, ridge=1e-6)
17    l1 = Multivariate_Gaussian_Distribution_likelihood(x, mu1, cov1, ridge=1e-6)
18    evidence = (l0 * p0) + (l1 * p1)
19    posterior_0 = (l0 * p0) / evidence
20    posterior_1 = (l1 * p1) / evidence
21
22    # Discriminant functions
23    g1_x = np.log(l1) + np.log(p1)
24    g0_x = np.log(l0) + np.log(p0)
25    delta = g1_x - g0_x
26
27    return delta, posterior_1
28
29 ''' BIC score (Bayesian Information Criterion)'''
30 def BIC_score(X_train, y_train):
31    classes = np.unique(y_train)
32    log_like = 0.0
33    n = len(y_train)
34    k = 0
35    for c in classes:
36        X_c = X_train[y_train == c]
37        mu = np.mean(X_c, axis=0)
38        cov = np.atleast_2d(np.cov(X_c, rowvar=False))
39        for x in X_c:
40            l = Multivariate_Gaussian_Distribution_likelihood(x, mu, cov, ridge=1e-6)
41            log_like = log_like + np.log(l)
42        k = k + len(mu) + (len(mu)*(len(mu)+1))/2
43    bic = np.log(n) * k - 2 * log_like # smaller is better
44    return bic
45
46 ''' Forward Feature Selection based on BIC '''
47 def Forward_Selection(X_train, y_train, max_features, min_features):
48    remaining = list(range(X_train.shape[1]))
49    selected = []
50    best_score = np.inf
51
52    while remaining and len(selected) < max_features:
53        scores = []
54        for f in remaining:
55            trial = selected + [f]
56            score = BIC_score(X_train[:, trial], y_train)
57            scores.append((score, f))
58            #print(f" Trial features: {[feature_names[j] for j in trial]}, BIC score: {score:.2f}")
59
60        new_score, new_f = min(scores, key=lambda x: x[0]) # 注意這裡要取最小的分數 (BIC越小越好)
61        if new_score < best_score or len(selected) < min_features:
62            best_score = new_score
63            selected.append(new_f)
64            remaining.remove(new_f)
65        else:
66            break
67    return selected
```

```

69 ''' Load Data '''
70 data = pd.read_excel('AcromegalyFeatureSet.xlsx')
71 data.rename(columns=lambda s: s.strip() if isinstance(s, str) else s, inplace=True)
72 X = data.drop(columns=['SeqNum', 'Gender', 'GroundTruth']).values
73 y = data['GroundTruth'].values
74 feature_names = data.drop(columns=['SeqNum', 'Gender', 'GroundTruth']).columns
75 n = len(y)
76
77 ''' Leave-One-Out Cross Validation '''
78 postsl = []
79 selected_features = []
80 deltas = []
81 print(f"Leave-One-Out Cross Validation")
82 for i in range(n):
83     X_train = np.delete(X, i, axis=0)
84     y_train = np.delete(y, i)
85     X_test = X[i]
86
87     n0 = int(np.sum(y_train == 0))
88     n1 = int(np.sum(y_train == 1))
89     print(f"-----")
90     print(f"Fold: {i+1}/{n}, Class 0 samples: {n0}, Class 1 samples: {n1}")
91
92     max_features = min(max(1, min(n0, n1) - 1), X_train.shape[1])
93
94     selected = Forward_Selection(X_train, y_train, max_features, 3)
95     selected_features.append([feature_names[j] for j in selected])
96     print(f"Selected features: {[feature_names[j] for j in selected]}")
97     X_train_sel = X_train[:, selected]
98     X_test_sel = X_test[selected]
99
100     # mu0, mu1 = mean vectors for class 0 and 1
101     mu0 = X_train_sel[y_train == 0].mean(axis=0)
102     mu1 = X_train_sel[y_train == 1].mean(axis=0)
103
104     # cov0, cov1 = covariance matrices for class 0 and 1
105     cov0 = np.atleast_2d(np.cov(X_train_sel[y_train == 0], rowvar=False))
106     cov1 = np.atleast_2d(np.cov(X_train_sel[y_train == 1], rowvar=False))
107
108     # p0 = prior for class 0, p1 = prior for class 1
109     p0, p1 = np.mean(y_train == 0), np.mean(y_train == 1)
110
111     delta, posteriorl = Bayesian_decision_classifier(X_test_sel, mu0, mu1, cov0, cov1, p0, p1)
112     postsl.append(posteriorl)
113     deltas.append(delta)
114 postsl = np.array(postsl)
115 deltas = np.array(deltas)
116
117 ''' Performance '''
118 fpr, tpr, _ = roc_curve(y, postsl)
119 roc_auc = auc(fpr, tpr)
120 preds1 = (deltas >= 0).astype(int)
121 cm = confusion_matrix(y, preds1)
122 TN, FP, FN, TP = cm.ravel()
123 acc = (TP + TN) / np.sum(cm)
124 sen = TP / (TP + FN)
125 spe = TN / (TN + FP)
126 print(f"-----")
127 print("Test Performance")
128 print(f"Confusion Matrix:\n{cm}")
129 print(f"Accuracy={acc:.3f}, Sensitivity={sen:.3f}, Specificity={spe:.3f}, AUC={roc_auc:.3f}")
130
131 # ROC curve plot
132 plt.plot(fpr, tpr, label=f"AUC={roc_auc:.2f}")
133 plt.plot([0, 1], [0, 1], 'k--')
134 plt.xlabel('False Positive Rate')
135 plt.ylabel('True Positive Rate')
136 plt.legend()
137 plt.title('ROC Curve')
138 plt.savefig('roc_curve.png')
139
140 ''' Bivariate Gaussian Bayes classifier with top-2 features '''
141 counts = Counter([f for fs in selected_features for f in fs])
142 print('Top-2 features:', counts.most_common(2))
143
144 Top2_features = [f for f, _ in counts.most_common(2)]
145
146 Top2_X = data[Top2_features].values.astype(float)
147 Top2_y = data['GroundTruth'].values.astype(int)
148
149 Top2_X0, Top2_X1 = Top2_X[Top2_y==0], Top2_X[Top2_y==1]
150 Top2_mu0, Top2_mu1 = Top2_X0.mean(0), Top2_X1.mean(0)
151 Top2_cov0, Top2_cov1 = np.cov(Top2_X0, rowvar=False), np.cov(Top2_X1, rowvar=False)
152 Top2_p0, Top2_p1 = len(Top2_X0)/len(Top2_X), len(Top2_X1)/len(Top2_X)
153
154 # grid points
155 m, M = Top2_X.min(0), Top2_X.max(0); pad = 0.05*(M-m)
156 xs = np.linspace(m[0]-pad[0], M[0]+pad[0], 300)
157 ys = np.linspace(m[1]-pad[1], M[1]+pad[1], 300)
158 xx, yy = np.meshgrid(xs, ys)
159 pts = np.c_[xx.ravel(), yy.ravel()]
160
161 # log-likelihoods and decision boundary
162 LL0 = np.array(np.log(Multivariate_Gaussian_Distribution_likelihood(p, Top2_mu0, Top2_cov0) for p in pts)).reshape(xx.shape)
163 LL1 = np.array(np.log(Multivariate_Gaussian_Distribution_likelihood(p, Top2_mu1, Top2_cov1) for p in pts)).reshape(xx.shape)
164 g0 = LL0 + np.log(Top2_p0)
165 g1 = LL1 + np.log(Top2_p1)
166 Top2_delta = g1 - g0 # decision boundary: delta=0
167
168 # plot
169 plt.figure(figsize=(7,6))
170 plt.scatter(Top2_X0[:,0], Top2_X0[:,1], marker='o', label='Class 0', color='gray')
171 plt.scatter(Top2_X1[:,0], Top2_X1[:,1], marker='v', label='Class 1', color='green')
172 levels0 = np.unique(np.sort(np.percentile(LL0, [20, 40, 60, 80])))
173 levels1 = np.unique(np.sort(np.percentile(LL1, [20, 40, 60, 80])))
174
175 plt.contour(xx, yy, LL0, levels=levels0, linestyle='dashed')
176 plt.contour(xx, yy, LL1, levels=levels1, linestyle='solid')
177
178 # The decision boundary
179 plt.contour(xx, yy, Top2_delta, levels=[0.0], linewidths=2, colors='black')
180
181 plt.xlabel(Top2_features[0]); plt.ylabel(Top2_features[1])
182 plt.title('Bivariate Gaussian Bayes (Top-2 features)')
183 plt.legend(); plt.tight_layout()
184 plt.savefig('bivariate_bayes_gda.png', dpi=180)

```

Figure 7



## Problem 2. Proof of Bayesian estimator.

Suppose  $x^t \sim N(\theta, \sigma^2)$  and  $\theta \sim N(u_0, \sigma_0^2)$ , where  $u_0, \sigma_0^2, \sigma^2$  are known. That is

$$p(X|\theta) = \frac{1}{(2\pi)^{N/2} \sigma^N} \exp\left[-\frac{\sum_t (x^t - \theta)^2}{2\sigma^2}\right]$$

$$p(\theta) = \frac{1}{\sqrt{2\pi} \sigma_0} \exp\left[-\frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right]$$

Please show that

$$E[\theta|X] = \frac{N/\sigma^2}{N/\sigma^2 + 1/\sigma_0^2} m + \frac{1/\sigma_0^2}{N/\sigma^2 + 1/\sigma_0^2} \mu_0$$

where  $m$  is the maximum likelihood estimator of the sample mean.

$$\left. \begin{array}{l} \text{suppose Data: } x^t \sim N(\theta, \sigma^2) \quad t \in \{1, \dots, N\} \\ \text{Prior distribution } \theta \sim N(\mu_0, \sigma_0^2) \end{array} \right\} \Rightarrow p(X|\theta) = \frac{1}{(2\pi)^{N/2} \sigma^N} \exp\left[-\frac{\sum_t (x^t - \theta)^2}{2\sigma^2}\right]$$

$$p(\theta) = \frac{1}{\sqrt{2\pi} \sigma_0} \exp\left[-\frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right]$$

由 Bayes' Rule 得  $p(\theta|X) \propto p(X|\theta) p(\theta)$

$$\Rightarrow p(\theta|X) \propto \exp\left[-\frac{1}{2\sigma^2} \sum_t (x^t - \theta)^2 - \frac{1}{2\sigma_0^2} (\theta - \mu_0)^2\right] \because \mu_0, \sigma_0^2, \sigma^2 \text{ are known}$$

$$\sum_t (x^t - \theta)^2 = \sum_t (x^t)^2 - 2\theta \sum_t x^t + N\theta^2, \quad (\theta - \mu_0)^2 = \theta^2 - 2\theta\mu_0 + \mu_0^2$$

$$\Rightarrow -\frac{1}{2\sigma^2} \sum_t (x^t - \theta)^2 - \frac{1}{2\sigma_0^2} (\theta - \mu_0)^2 = -\frac{1}{2} \left[ \frac{N}{\sigma^2} \theta^2 - 2 \frac{\sum_t x^t}{\sigma^2} \theta + \frac{1}{\sigma_0^2} \theta^2 - 2 \frac{\mu_0}{\sigma_0^2} \theta \right] + \underbrace{\left( -\frac{\mu_0^2}{2\sigma_0^2} - \frac{1}{2\sigma^2} \sum_t (x^t)^2 \right)}_{\text{constant}}$$

$$\text{set } A = \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}, \quad B = \frac{\sum_t x^t}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}$$

$$\Rightarrow p(\theta|X) \propto \exp\left[-\frac{1}{2}(A\theta^2 - 2B\theta)\right] = \exp\left[-\frac{A}{2}(\theta^2 - 2\frac{B}{A}\theta)\right] \propto \exp\left[-\frac{A}{2}(\theta^2 - 2\frac{B}{A}\theta + (\frac{B}{A})^2) + \underbrace{\frac{B^2}{2A}}_{\text{constant}}\right]$$

$$= \exp\left[-\frac{A}{2}(\theta - \frac{B}{A})^2\right]$$

$$\Rightarrow p(\theta|X) = N(\text{mean} = \frac{B}{A}, \text{variance} = \frac{1}{A})$$

$$\Rightarrow E(\theta|X) = \frac{B}{A} = \frac{\frac{\sum_t x^t}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}}{\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}}$$

' $\because m = \frac{1}{N} \sum_t x^t$  the maximum likelihood of sample mean  $\Rightarrow \sum_t x^t = Nm$

$$\Rightarrow E(\theta|X) = \frac{\frac{Nm}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}}{\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}} = \frac{N/\sigma^2}{N/\sigma^2 + 1/\sigma_0^2} m + \frac{1/\sigma_0^2}{N/\sigma^2 + 1/\sigma_0^2} \mu_0$$

得證 #