# EECE 5644 Assignment 1    Jingcheng Wang

*Repository:*    *https://github.com/tomwang777/2025-Fall-EECE-5644-Machine-Learning/tree/main/Assignment%201*

## Question 1

Generated N0 = 6526, N1 = 3474

Part A



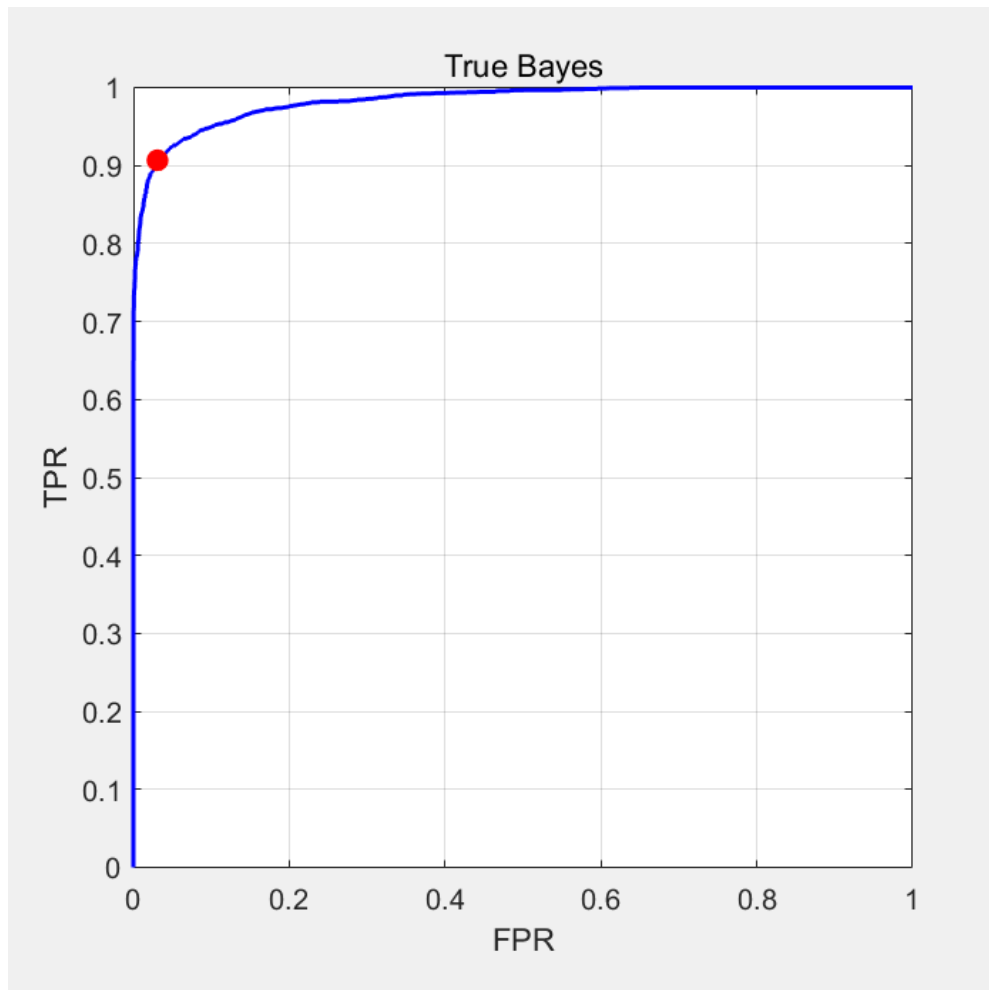The empirical γ value I chose is much smaller than the theoretical optimal threshold, which means that the actual sample situation will deviate greatly from the theoretical result and requires a comprehensive analysis instead of calculating only one of them.
True Bayes:
  min empirical error = 0.052762, threshold = 1.416269
  TPR = 0.9067, FPR = 0.0310
  theoretical gamma = 1.8571, empirical error at gamma = 0.053613

**True Bayes**

## Part B

By incorrectly assuming that the covariance matrix is equal to the identity matrix and recalculating the likelihood ratio corresponding to each threshold, and drawing the ROC curve, we can find that the model mismatch leads to a significant decrease in the threshold results when designing the naive Bayes classifier. Only the values of C0 and C1 are changed, and the others remain unchanged. The minimum error probability of the model increases, which has a negative impact on the model performance. The reason is that the ROC curve moves downward, the minimum error probability becomes larger.

Naive Bayes (I):
min empirical error = 0.063794, threshold = 0.914395

Naive Bayes

Part C



Part C.

Fisher LDA classifier.

No. $L=0$.  $N_1$.  $L=1$.

$\hat{m}_j = \frac{1}{N_j} \sum_{x \in j} x$.

$\hat{C}_j = \frac{1}{N_j - 1} \sum_{x \in j} (x - \hat{m}_j)(x - \hat{m}_j)^T$.

$S_w = \hat{C}_0 + \hat{C}_1$.

$w_{LDA} = S_w^{-1}(\hat{m}_1 - \hat{m}_0)$.

Projection. for every $x$.  $y = w_{LDA}^T x$.

Decision Rule.  $y > \tau \Rightarrow D = 1$.  $y \leq \tau. D = 0$.

For every $\tau$.  $FPR = \dfrac{\text{Amount of } (w_{LDA}^T x > \tau \text{ and } L=0)}{\text{Amount of } (L=0)}$.  $TPR = \dfrac{\text{Amount of } (w_{LDA}^T x > \tau \text{ and } L=1)}{\text{Amount of } (L=1)}$.

The performance of Fisher LDA is between True Bayes and Naive Bayes, mainly because it effectively utilizes the correlation information, but the projection causes some loss.

Fisher LDA:

min empirical error = 0.060396, threshold = 0.554445

The red dot in each figure indicates the minimum error point.

Overall

## Question 2

Part A



Question 2.

Part A.

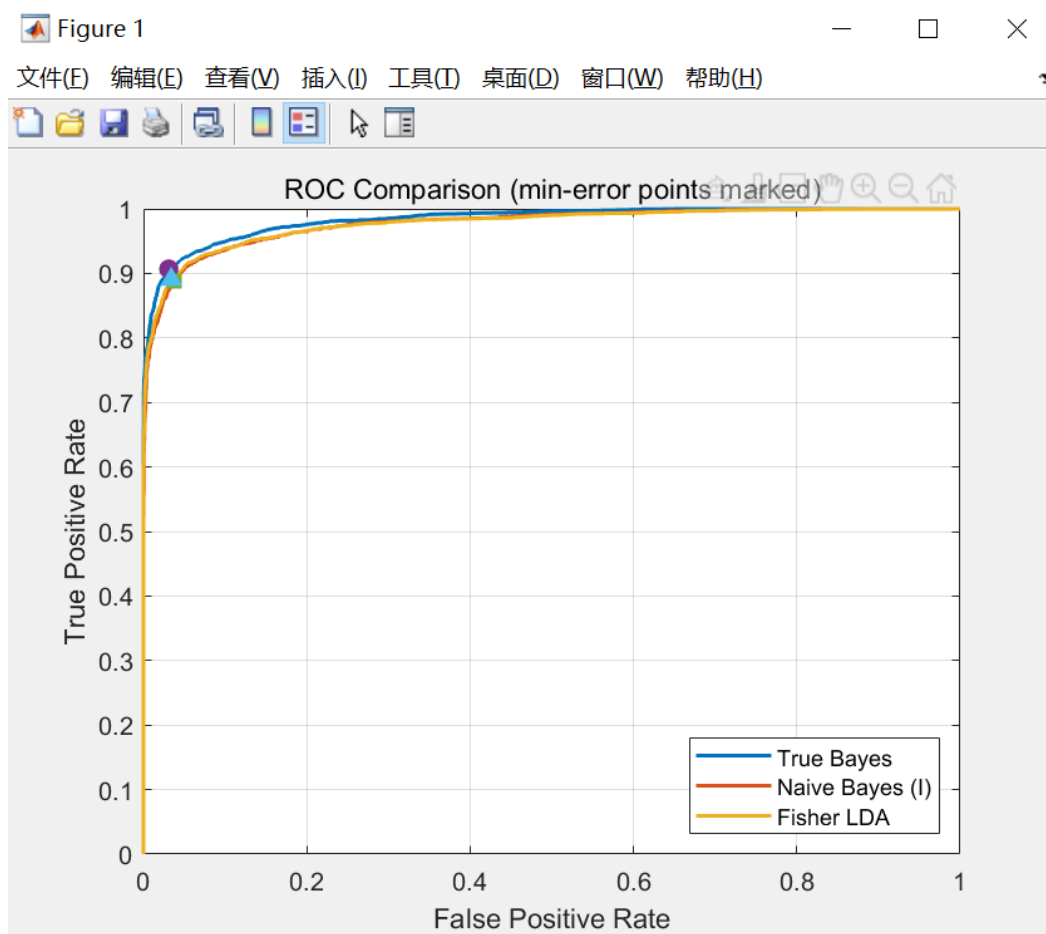$m_1 = [0 \ 0]$ $\quad C_1 = \begin{bmatrix} 0.8 & 0.2 \\ 0.2 & 0.6 \end{bmatrix}$

$m_2 = [3.5 \ 0.5]$ $\quad C_2 = \begin{bmatrix} 0.6 & -0.15 \\ -0.15 & 0.5 \end{bmatrix}$

$m_3 = [0.5 \ 3.0]$ $\quad C_3 = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$

$m_4 = [3.0 \ 3.0]$ $\quad C_4 = \begin{bmatrix} 1.0 & 0.3 \\ 0.3 & 0.9 \end{bmatrix}$

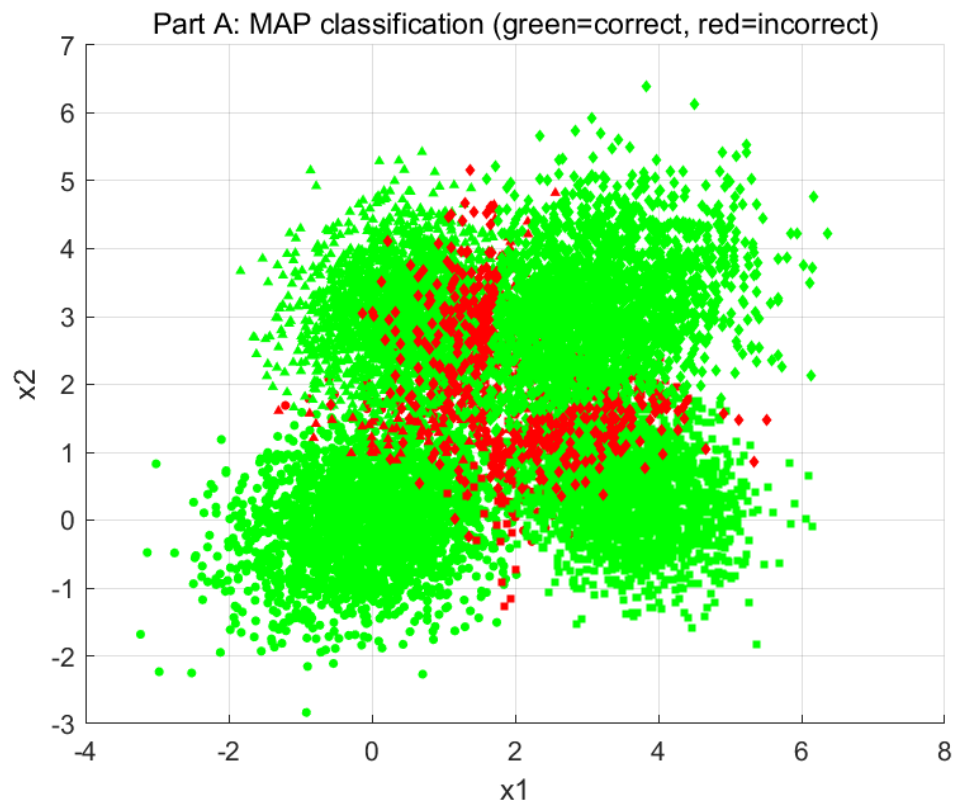$P(L=j) = 0.25.$ $\quad j = 1,2,3,4.$

1. $N = 1000.$

   For every sample. randomly choose a true label $L$ based on $P(L=j) = 0.25$.
   generate 2-dimensional vector $x$ from corresponding $p(x|L)$.

2. $p(L=j \mid x) = \dfrac{p(x \mid L=j) \, P(L=j)}{p(x)}.$

   MAP Decision Rule. $D = \arg\max\limits_{j \in \{1,2,3,4\}} p(L=j \mid x) = \arg\max\limits_{j \in \{1,2,3,4\}} p(x \mid L=j).$

   $= \arg\max\limits_{j \in \{1,2,3,4\}} \dfrac{1}{\sqrt{(2\pi)^2 |C_j|}} e^{-\frac{1}{2}(x-m_j)^T C_j^{-1}(x-m_j)}$

   for confusion matrix $M$.

   $P(D=i \mid L=j) = \dfrac{\text{Amount of } (x \text{ classified as } D=i \ \& \ \text{true label } L=j)}{\text{Amount of } (\text{True label } L=j).}$

   $= \dfrac{}{N_j}.$

Total samples: 10000 | incorrect: 813



Part A: MAP classification (green=correct, red=incorrect)

Confusion Matrix P(D=i|L=j) - MAP

| Decision D=i | | | | |
|---|---|---|---|---|
| 1 | 0.957 | 0.012 | 0.018 | 0.013 |
| 2 | 0.014 | 0.934 | 0.001 | 0.061 |
| 3 | 0.020 | 0.001 | 0.933 | 0.077 |
| 4 | 0.009 | 0.053 | 0.048 | 0.849 |
| | 1 | 2 | 3 | 4 |

True Label L=j

Part B



Part B.

$$R(D=i|x) = \sum_{j=1}^{4} \lambda_{ij} P(L=j|x).$$

ERM Decision Rule. $D = \arg\min_{i \in \{1,2,3,4\}} R(D=i|x).$

$$= \arg\min_{i \in \{1,2,3,4\}} \sum_{j=1}^{4} \lambda_{ij} P(L=j|x).$$

$$= \arg\min_{i \in \{1,2,3,4\}} \sum_{j=1}^{4} \lambda_{ij} \frac{P(x|L=j) \, P(L=j)}{P(x)}$$

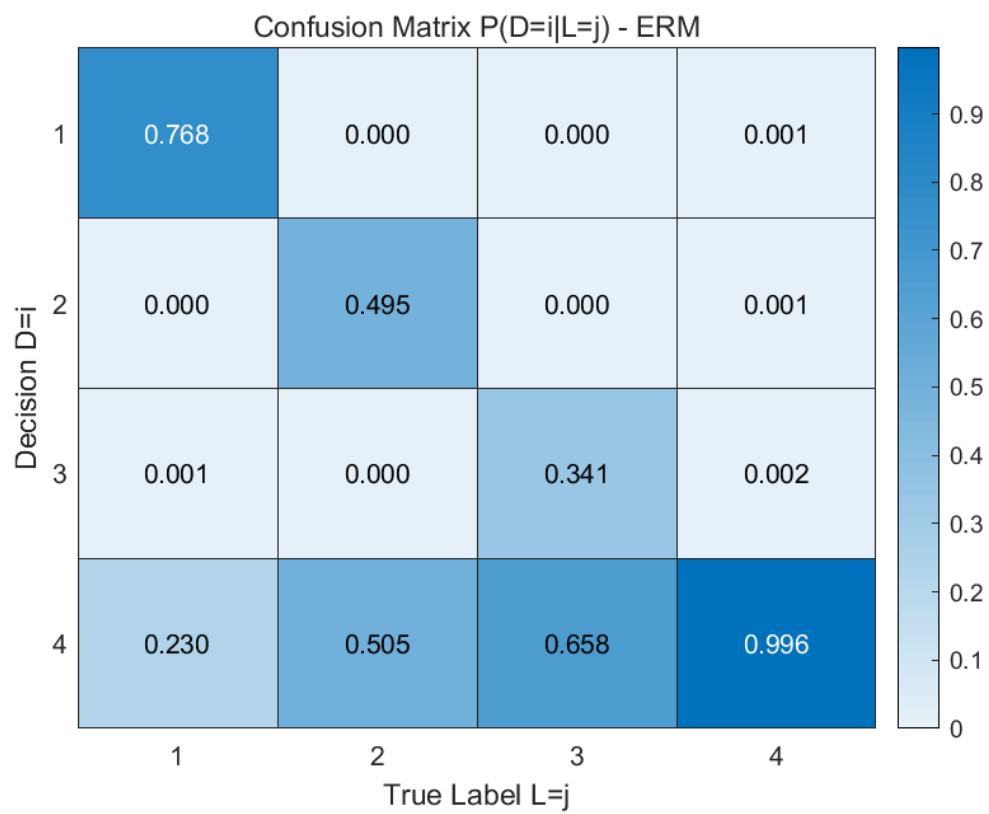$$= \arg\min_{i \in \{1,2,3,4\}} \sum_{j=1}^{4} \lambda_{ij} \, P(x|L=j). \qquad 0.4514.$$
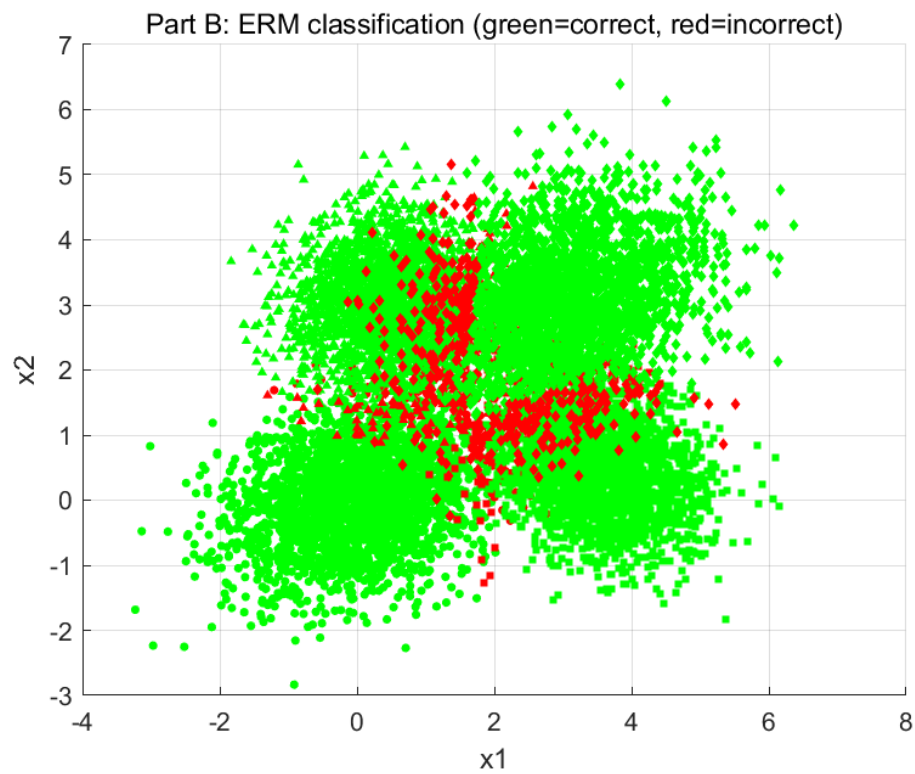
Minimum expected risk. $R = E[loss] = \sum_{j=1}^{4} \sum_{i=1}^{4} \lambda_{ij} P(D=i, L=j)$

$$P(D=i, L=j) = \frac{\text{Amount of } (D=i \ \& \ L=j)}{N}.$$

$$\hat{R}_{min} = \frac{1}{N} \sum_{x} loss(x)$$

$$= \frac{1}{N} \sum_{j=1}^{4} \sum_{i=1}^{4} \lambda_{ij} \frac{\text{Amount of } (D=i \ \& \ L=j)}{N} \qquad 3.8343.$$

Total samples: 10000 | incorrect: 813
MAP empirical risk (Lambda) = 3.8343
ERM empirical risk = 0.4514

Part B: ERM classification (green=correct, red=incorrect)



Confusion Matrix P(D=i|L=j) - ERM

| Decision D=i \ True Label L=j | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0.768 | 0.000 | 0.000 | 0.001 |
| 2 | 0.000 | 0.495 | 0.000 | 0.001 |
| 3 | 0.001 | 0.000 | 0.341 | 0.002 |
| 4 | 0.230 | 0.505 | 0.658 | 0.996 |

## Question 3

Question 3.

Class Priors. $\hat{P}(L=j) = \dfrac{\text{Amount of (samples in class } j)}{\text{Sample amounts.}}$

Mean Vectors. $\hat{m}_j = \dfrac{1}{N_j} \sum\limits_{x \in \text{Class } j} x$ .

Covariance Matrix. $C_{\text{Sample Average}, j} = \dfrac{1}{N_j - 1} \sum\limits_{x \in \text{Class } j} (x - \hat{m}_j)(x - \hat{m}_j)^T$ .

If ill-conditioned for sample Covariance Matrix $C$ .

Regularization .

$$C_{\text{Regularized}, j} = C_{\text{sampleAverage}, j} + \lambda I$$
$$= \dfrac{1}{N_j - 1} \sum\limits_{x \in \text{Class } j} (x - \hat{m}_j)(x - \hat{m}_j)^T + \lambda I . \quad \lambda > 0. \qquad \lambda = \alpha \cdot \dfrac{\text{trace } (C_{\text{sample Average}})}{\text{rank } (C_{\text{Sample Average}})} .$$

MAP Rule.

$$P(L=j \mid x) = \dfrac{P(x \mid L=j) \, P(L=j)}{P(x)} .$$

$$D = \underset{j}{\text{argmax}} \; P(L=j \mid x)$$

$$= \underset{j}{\text{argmax}} \left\{ P(x \mid L=j) \, P(L=j) \right\} .$$

$$= \underset{j}{\text{argmax}} \left\{ \hat{P}(L=j) \cdot \dfrac{1}{\sqrt{(2\pi)^d \, |C_{\text{Reg}, j}|}} \; e^{-\frac{1}{2}(x - \hat{m}_j)^T \, C_{\text{Reg}, j}^{-1} \, (x - \hat{m}_j)} \right\} .$$

$$P(\text{error}) = \dfrac{\sum\limits_{i \neq j} \text{Amount of } (D=i \,\&\, L=j)}{N} .$$

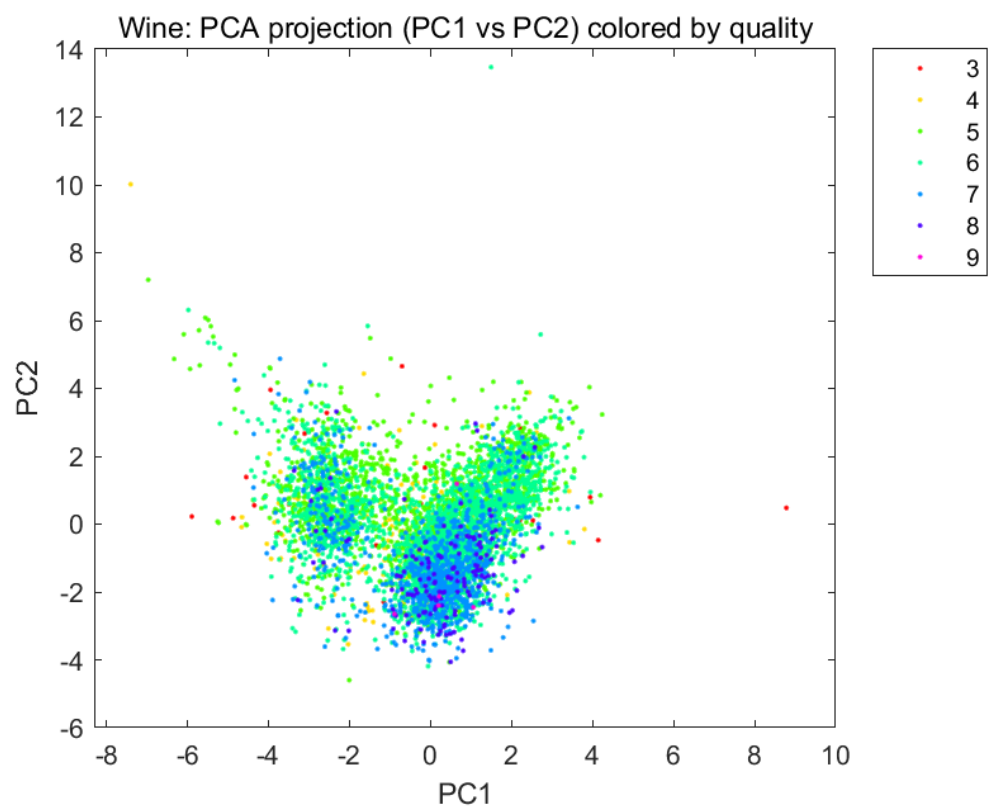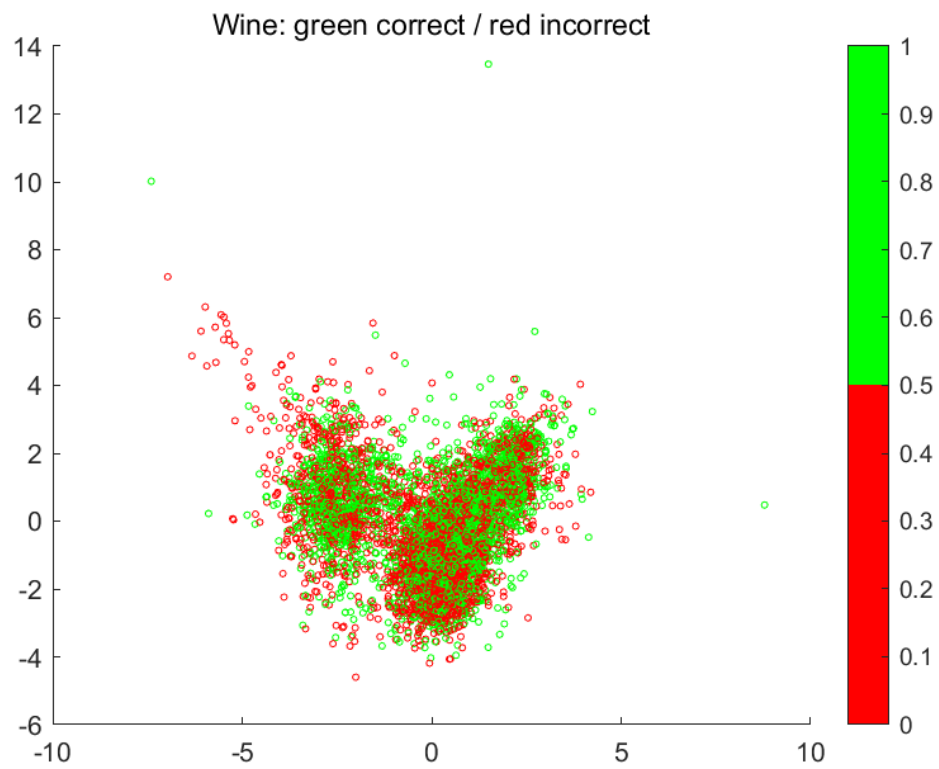Confusion Matrix.

$$M_{ij} = \text{Amount of } (D=i \,\&\, L=j) .$$
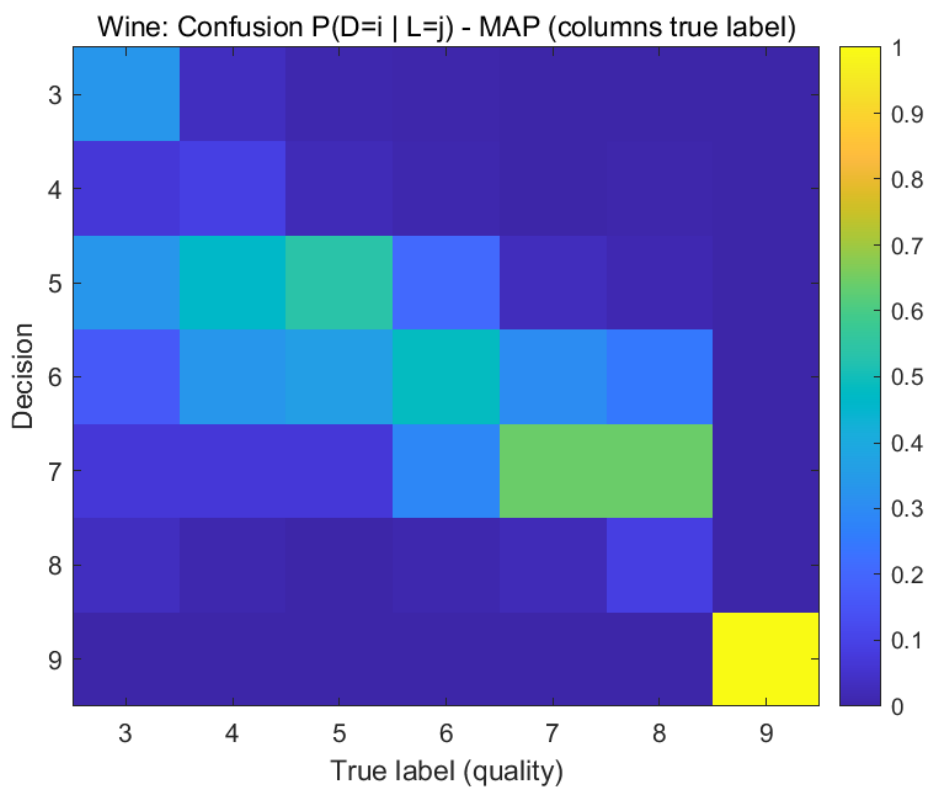
Wine dataset (all data) results: N = 6497
d = 11
classes = 7
Training MAP error = 0.4982

Wine: green correct / red incorrect

Wine: PCA projection (PC1 vs PC2) colored by quality

Wine: Confusion P(D=i | L=j) - MAP (columns true label)

HAR dataset results: N = 10299
d = 561
classes = 6
Training MAP error = 0.0149



HAR: PCA (PC1 vs PC2) colored by activity label

HAR: Confusion P(D=i|L=j) - MAP

The Gaussian class conditional model may not be appropriate for these datasets, as the PCA projections show significant aliasing, suggesting an oversimplified assumption of a Gaussian distribution for the same class. The Gaussian model choice may be too simplistic for this dataset, leading to model bias, higher training errors, and increased confusion. GMMs and mixture models may be needed to improve training quality and discriminative performance.

Citation

1. Course recording

2. Course notes

3. Course codes provided on Canvas

4. Discussion with classmates

5. Generative AI models