**Problem1**

1. The goal of the project
   This project is aimed at explore the quality rating for many well-known US hospitals. Since health care is an important issue for current society, an increasing number of US hospitals have emerged among different states, which support service in a diverse range of professional fields. This project hopes to help people decide which hospital is the best choice depending on many factors, such as charges, distance, patient reviews, doctor qualifications, etc. Through weighting these factors in an analysis model for different kinds of patients, quality for hospitals based on their needs is possible to be ranked.

2. Data collection
   The dataset is offered by Centers for Medicare & Medicaid Services (CMS). It contains information of 90 hospitals, including locations, hospital type, mortality records, safety of care national comparison, patient experience national comparison, etc. More information can be retrieved from its website.

3. Current/ Future tools
   PCA, LDA, linear/non-linear regression, decision tree

   (Note: This is only a preliminary proposal plan, it may be changed or modifies if I find something more interesting afterwards)

**Problem 2**

P2.

① MAP Decision. Rule.

Two class $c \in \{0,1\}$ $(y_i, x_i)$ are iid. drawn from.

$P_x(x|y=c) \sim N(\bar{x}^c, \Sigma)$ $P_{xy}(x, y=c) = P_x(x|y=c) P(y=c)$

↑
same for both classes.

$$P(y=c|x) = \frac{P_{xy}(x, y=c)}{P_x(x)} = \frac{P_x(x|y=c) P(y=c)}{P_x(x)}.$$

MAP decision rule : $\hat{C}(x) = \begin{cases} 1 & P(y=1|x) \geqslant P(y=0|x) \\ 0 & o/w. \end{cases}$

$$P(y=1|x) \underset{<}{\geqslant} P(y=0|x) \qquad (\hat{C}(x) = 1)$$
$$\qquad\qquad (\hat{C}(x) = 0)$$

$$\frac{P(x|y=1) P(y=1)}{P_x(x)} \underset{\substack{< \\ \hat{c}_x=0}}{\overset{\hat{c}_x=1}{\gtrless}} \frac{P(x|y=0) P(y=0)}{P_x(x)}$$

$P(y=1) = \tau_1$ $P(y=0) = \tau_0$.

$P(x|y=1) \sim N(\bar{x}^1, \Sigma) \Rightarrow \frac{1}{(2\pi)^{\frac{d}{2}}|\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x-\bar{x}^1)^T \Sigma^{-1}(x-\bar{x}^1))$

$P(x|y=0) \sim N(\bar{x}^0, \Sigma) \Rightarrow \frac{1}{(2\pi)^{\frac{d}{2}}|\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x-\bar{x}^0)^T \Sigma^{-1}(x-\bar{x}^0))$

$$\Rightarrow \frac{\tau_1}{(2\pi)^{\frac{d}{2}}|\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x-\bar{x}^1)^T \Sigma^{-1}(x-\bar{x}^1)) \underset{\substack{< \\ C(x)=0}}{\overset{C(x)=1}{\gtrless}} \frac{1}{(2\pi)^{\frac{d}{2}}|\Sigma|^{\frac{1}{2}}} \tau_0 \exp(-\frac{1}{2}(x-\bar{x}^0)^T \Sigma^{-1}(x-\bar{x}^0))$$

$$\frac{\exp(-\frac{1}{2}(x-\bar{x}^1)^T \Sigma^{-1}(x-\bar{x}^1))}{\exp(-\frac{1}{2}(x-\bar{x}^0)^T \Sigma^{-1}(x-\bar{x}^0))} \frac{\tau_1}{\tau_0} \underset{\substack{< \\ C(x)=0}}{\overset{C(x)=1}{\gtrless}} 1$$

$$\log \Rightarrow -\frac{1}{2}(x-\bar{x}^1)^T \Sigma^{-1}(x-\bar{x}^1) + \frac{1}{2}(x-\bar{x}^0)^T \Sigma^{-1}(x-\bar{x}^0) + \ln\frac{\tau_1}{\tau_0} \underset{\substack{< \\ c(x)=0}}{\overset{C(x)=1}{\gtrless}} 0$$
①

$$-\frac{1}{2}x^T\Sigma^{-1}x + \frac{1}{2}\bar{x}^{1T}\Sigma^{-1}x + \frac{1}{2}x\Sigma^{-1}\bar{x}^1 - \frac{1}{2}\bar{x}^{1T}\Sigma^{-1}\bar{x}^1 +$$

$$\frac{1}{2}x^T\Sigma^{-1}x - \frac{1}{2}\bar{x}^{0T}\Sigma^{-1}x - \frac{1}{2}x^T\Sigma^{-1}\bar{x}^0 + \frac{1}{2}\bar{x}^0\Sigma^{-1}\bar{x}^0 + \ln\frac{z_1}{z_0} \underset{c(x)=0}{\overset{c(x)=1}{\gtrless}} 0$$

$$\frac{1}{2}x^T\Sigma^{-1}(\bar{x}^1 - \bar{x}^0) + \frac{1}{2}(\bar{x}^{1T} - \bar{x}^{0T})\Sigma^{-1}x + \frac{1}{2}\bar{x}^0\Sigma^{-1}\bar{x}^0$$

$$-\frac{1}{2}\bar{x}^{1T}\Sigma^{-1}\bar{x}^1 + \ln\frac{z_1}{z_0} \underset{c(x)=0}{\overset{c(x)=1}{\gtrless}} 0$$

$\circledast$ 
$$\underbrace{x^T\Sigma^{-1}(\bar{x}^1 - \bar{x}^0)}_{\theta} + \underbrace{\frac{1}{2}\bar{x}^{0T}\Sigma^{-1}\bar{x}_0 - \frac{1}{2}\bar{x}^{1T}\Sigma^{-1}\bar{x}^1 + \log\frac{z_1}{z_0}}_{\theta_0} \underset{c(x)=0}{\overset{c(x)=1}{\gtrless}} 0$$

$$\Rightarrow \quad \theta^T x + \theta_0 \underset{c(x)=0}{\overset{c(x)=1}{\gtrless}} 0 \qquad \theta = \Sigma^{-1}(\bar{x}^1 - \bar{x}^0)$$

$$\theta_0 = \frac{1}{2}\bar{x}^{0T}\Sigma^{-1}\bar{x}^0 - \frac{1}{2}\bar{x}^{1T}\Sigma^{-1}\bar{x}^1 + \log\frac{z_1}{z_0}$$

---

② Minimum - distance decision rule.
(Euclidean)

Assign x to class $C_i$ : $i = \text{argmin}_j (x - \bar{x}^j)^T(x - \bar{x}^j)$ $\qquad j = 1,2$

Normal $\quad C_0 \sim N(\bar{x}_0, \Sigma) \qquad C_1 \sim (\bar{x}_1, \Sigma) \qquad \Sigma = \sigma^2 I$

if. $\quad (x - \bar{x}^0)^T(x - \bar{x}^0) > (x - \bar{x}^1)^T(x - \bar{x}^1)$

assign x to class 1

o/w assign x to class 0.

Compare it with Bayesian rule,

as show in ① $g = -\frac{1}{2}(x - \bar{x}^1)^T\Sigma^{-1}(x - \bar{x}^1) + \frac{1}{2}(x - \bar{x}^0)^T\Sigma^{-1}(x - \bar{x}^0) + \ln\frac{z_1}{z_0}$

$z_1, z_0$ are the same in this case.

$$(X - \bar{X}^1)^T \Sigma^{-1} (X - \bar{X}^1) \underset{\text{class 1}}{\overset{\text{class 0}}{\gtrless}} (X - \bar{X}^0)^T \Sigma^{-1} (X - \bar{X}^0)$$

this is called Minimum Mahalanobis distance classifier.

Unlike Euclidean decision rule. $(\Sigma = \sigma^2 I)$, Mahalanobis distance decision rule is equivalent to bayesian decision rule, which is more general as a weighted form of Euclidean decision rule.

## Problem 3

1,2,3. See codes

4. For the first case that two class has the same covariance, the classification error using Bayesian classification rule is 24, while the error for logistic regression is 32. They look similar as the data can be easily separated by a linear line. Here, Bayesian rule can be simplified as a linear classification and logistic regression is a generalized linear model, which have similar performance for this scenario.

5. After changing the covariance matrix, the points from two classes are overlapped in the centre part. In this case, Bayesian rule works no longer as a linear model, instead, considering the quadratic term, it is equivalent to Mahalanobis distance rule, which can classify the overlapping points based on the distance between them and mean for each class. On the other hand, logistic regression still works as a linear model and cannot distinguish the overlapping parts very well. Therefore, as the test errors shows, error for Bayesian rule is 275, while error for logistic regression is 489. This result validates our analysis.