

1. The pdf for two jointly Gaussian random variables  $X$  and  $Y$  is of the following form parameterized by the scalars  $m_1, m_2, \sigma_1, \sigma_2$  and  $\rho_{XY}$ :

$$f_{X,Y}(x,y) = \frac{\exp\left\{-\frac{1}{2(1-\rho_{XY}^2)}\left[(\frac{x-m_1}{\sigma_1})^2 - 2\rho_{XY}\left(\frac{x-m_1}{\sigma_1}\right)\left(\frac{y-m_2}{\sigma_2}\right) + (\frac{y-m_2}{\sigma_2})^2\right]\right\}}{2\pi\sigma_1\sigma_2\sqrt{1-\rho_{XY}^2}}. \quad (1)$$

The pdf for multivariate jointly Gaussian random variable  $Z \in \mathbb{R}^k$  is of the following form parameterized by  $\mu \in \mathbb{R}^k$  and  $\Sigma \in \mathbb{R}^{k \times k}$ .

$$f_Z(z) = \frac{\exp\left\{-\frac{1}{2}(z-\mu)^T\Sigma^{-1}(z-\mu)\right\}}{\sqrt{(2\pi)^k|\Sigma|}}. \quad (2)$$

Suppose  $Z = [X, Y]^T$ , i.e.,  $z = [x, y]^T$ .

- (a) Find  $\mu, \Sigma^{-1}$  and  $\Sigma$  in terms of  $m_1, m_2, \sigma_1, \sigma_2$  and  $\rho_{XY}$ .  
 (b) Suppose  $\rho_{XY} = 0$ , what is  $\Sigma$  in this case? Can you write  $f_{X,Y}(x,y)$  as the product of two single variate Gaussian distributions? Are  $X$  and  $Y$  independent?

a) By Comparing the two equations above, we see that when Multivariate jointly Gaussian Density is applied to two variables:

$$\begin{aligned} \mu &= \begin{bmatrix} m_1 \\ m_2 \end{bmatrix} \\ \Sigma^{-1} &= \frac{1}{\sigma_1^2 \sigma_2^2 (1-\rho_{XY}^2)} \begin{bmatrix} \sigma_2^2 & -\rho_{XY} \sigma_1 \sigma_2 \\ -\rho_{XY} \sigma_1 \sigma_2 & \sigma_1^2 \end{bmatrix} \\ \Sigma &= \begin{bmatrix} \sigma_1^2 & \rho_{XY} \sigma_1 \sigma_2 \\ \rho_{XY} \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}. \end{aligned}$$

We double check by plugging  $\mu, \Sigma^{-1}, \Sigma$  in to equation (2):

$$\begin{aligned} f_Z([x, y]) &= \frac{\exp\left\{-\frac{1}{2}([x, y] - [\begin{smallmatrix} m_1 \\ m_2 \end{smallmatrix}])^T \frac{1}{\sigma_1^2 \sigma_2^2 (1-\rho_{XY}^2)} \begin{bmatrix} \sigma_2^2 & -\rho_{XY} \sigma_1 \sigma_2 \\ -\rho_{XY} \sigma_1 \sigma_2 & \sigma_1^2 \end{bmatrix} \cdot ([x, y] - [\begin{smallmatrix} m_1 \\ m_2 \end{smallmatrix}])\right\}}{\sqrt{(2\pi)^2 (\sigma_1^2 \sigma_2^2 (1-\rho_{XY}^2))}} \\ &= \frac{\exp\left\{-\frac{1}{2 \sigma_1^2 \sigma_2^2 (1-\rho_{XY}^2)} [(x-m_1)\sigma_2^2 - (y-m_2)\rho_{XY} \sigma_1 \sigma_2, -(x-m_1)\rho_{XY} \sigma_1 \sigma_2 + (y-m_2)\sigma_1^2] \begin{bmatrix} x-m_1 \\ y-m_2 \end{bmatrix}\right\}}{2\pi\sigma_1\sigma_2\sqrt{1-\rho_{XY}^2}} \\ &= \frac{\exp\left\{-\frac{1}{\sigma_1^2 \sigma_2^2 (1-\rho_{XY}^2)} [(x-m_1)^2 \sigma_2^2 - (y-m_2)^2 \sigma_1^2 + (y-m_2)(x-m_1)\rho_{XY} \sigma_1 \sigma_2 + (-xm_1)(ym_2)\rho_{XY} \sigma_1 \sigma_2 + (y-m_2)^2 \sigma_1^2]\right\}}{2\pi\sigma_1\sigma_2\sqrt{1-\rho_{XY}^2}} \\ &= \frac{\exp\left\{-\frac{1}{(1-\rho_{XY}^2)^2} \left[ \left(\frac{x-m_1}{\sigma_1}\right)^2 + \left(\frac{y-m_2}{\sigma_2}\right)^2 - 2\rho_{XY} \frac{(x-m_1)(y-m_2)}{\sigma_1 \sigma_2} \right]\right\}}{2\pi\sigma_1\sigma_2\sqrt{1-\rho_{XY}^2}} \end{aligned}$$

b) If  $\rho_{XY}=0$ ,  $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$

$$\begin{aligned} f_{X,Y}(x,y) &= \frac{\exp\left\{-\frac{1}{2} \left[ \left(\frac{x-m_1}{\sigma_1}\right)^2 + \left(\frac{y-m_2}{\sigma_2}\right)^2 \right]\right\}}{2\pi\sigma_1\sigma_2} \\ &= \frac{1}{2\pi\sigma_1\sigma_2} e^{-\frac{1}{2} \left(\frac{x-m_1}{\sigma_1}\right)^2} \cdot e^{-\frac{1}{2} \left(\frac{y-m_2}{\sigma_2}\right)^2} \\ &= f(x) \cdot f(y) \end{aligned}$$

Yes,  $X$  and  $Y$  are independent as  $f(x,y) = f(x) \cdot f(y)$ .

2. The Gaussian Discriminant Analysis (GDA) models the class conditional distribution as multivariate Gaussian, i.e.,  $P(X|Y) \sim \mathcal{N}(\mu_Y, \Sigma)$ . Suppose we want to enforce the **Naive Bayes (NB) assumption**, i.e.  $P(X_i|Y, X_j) = P(X_i|Y), \forall j \neq i$ , to GDA. Show that all off diagonal elements of  $\Sigma$  equal to 0:  $\Sigma_{i,j} = 0, \forall i \neq j$  with the **NB assumption**.

$$\text{Cov}(x_i|y) = E((x - E(x))(y - E(y))) \text{ by definition.}$$

$$* E(E(x_j|y)) \\ = E(x_j|y)$$

$$\text{Cov}(x_i|y, x_j|y) = E((x_i|y - E(x_i|y))(x_j|y - E(x_j|y)))$$

$$= E((x_i|y)(x_j|y) - (x_i|y)E(x_j|y) - (x_j|y)E(x_i|y) + E(x_i|y)E(x_j|y))$$

$$= 2 E(x_i|y)E(x_j|y) - 2 E(x_i|y)E(x_j|y)$$

$$= 0$$

3. Consider the classification problem for two classes,  $C_0$  and  $C_1$ . In the generative approach, we model the class-conditional distribution  $P(x|C_0)$  and  $P(x|C_1)$ , as well as the class priors  $P(C_0)$  and  $P(C_1)$ . The posterior probability for class  $C_0$  can be written as

$$P(C_0|x) = \frac{P(x|C_0)P(C_0)}{P(x|C_0)P(C_0) + P(x|C_1)P(C_1)}.$$

(a) Show that  $P(C_0|x) = \sigma(a)$  where  $\sigma(a)$  is the *sigmoid* function defined by

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

Find  $a$  in terms of  $P(x|C_0)$ ,  $P(x|C_1)$ ,  $P(C_0)$  and  $P(C_1)$ .

$$P(C_0|x) = \frac{1}{1 + \frac{P(x|C_1)P(C_1)}{P(x|C_0)P(C_0)}} \\ \exp(-a) = \frac{P(x|C_1)P(C_1)}{P(x|C_0)P(C_0)}$$

$$a = -\ln\left(\frac{P(x|C_1)P(C_1)}{P(x|C_0)P(C_0)}\right) = \ln\left(\frac{P(x|C_0)P(C_0)}{P(x|C_1)P(C_1)}\right)$$

(b) In the GDA model, we have the class conditional distribution as follows

$$P(x|C_0) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right),$$

$$P(x|C_1) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right).$$

Suppose we are able to find the maximum likelihood estimation of  $\mu_0, \mu_1, \Sigma, P(C_0)$ , and  $P(C_1)$ . Show that  $a = w^T x + b$  for some  $w$  and  $b$ . Find  $w$  and  $b$  in terms of  $\mu_0, \mu_1, \Sigma, P(C_0)$ , and  $P(C_1)$ . This shows that the decision boundary is linear.

$$a = \ln\left(\frac{\frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right) \cdot P(C_0)}{\frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right) \cdot P(C_1)}\right) \\ = \ln\left(\exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right) \cdot P(C_0)\right) - \ln\left(\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right) \cdot P(C_1)\right) \\ = \left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0) + \ln(P(C_0)) - \left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1) + \ln(P(C_1))\right)\right) \\ = -\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0) + \frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1) + \ln\left(\frac{P(C_0)}{P(C_1)}\right) \\ = -\frac{1}{2}(x^T - \mu_0^T) \Sigma^{-1} (x - \mu_0) + \frac{1}{2}(x^T - \mu_1^T) \Sigma^{-1} (x - \mu_1) + \ln\left(\frac{P(C_0)}{P(C_1)}\right) \\ = -\frac{1}{2}(x^T \Sigma^{-1} - \mu_0^T \Sigma^{-1}) (x - \mu_0) + \frac{1}{2}(x^T \Sigma^{-1} - \mu_1^T \Sigma^{-1}) (x - \mu_1) + \ln\left(\frac{P(C_0)}{P(C_1)}\right) \\ = -\frac{1}{2}(x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x + \mu_0^T \Sigma^{-1} \mu_0) + \frac{1}{2}(x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1) + \ln\left(\frac{P(C_0)}{P(C_1)}\right) \\ = -\frac{1}{2}(-2x^T \Sigma^{-1} \mu_0 + \mu_0^T \Sigma^{-1} \mu_0) + \frac{1}{2}(-2x^T \Sigma^{-1} \mu_1 + \mu_1^T \Sigma^{-1} \mu_1) + \ln\left(\frac{P(C_0)}{P(C_1)}\right) \\ = x^T \Sigma^{-1} \mu_0 - \frac{1}{2}\mu_0^T \Sigma^{-1} \mu_0 - x^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \ln\left(\frac{P(C_0)}{P(C_1)}\right) \\ = x^T \Sigma^{-1} (\mu_0 - \mu_1) - \frac{1}{2}\mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \ln\left(\frac{P(C_0)}{P(C_1)}\right)$$

Constant C

$$\begin{aligned}
& \text{transpose } X^T \Sigma^{-1} (\mu_0 - \mu_1) \text{ as a number transpose is still a number} \\
& = (\mu_0 - \mu_1)^T (\Sigma^{-1})^T X + C \quad * \Sigma \text{ is symmetric square matrix,} \\
& = (\mu_0 - \mu_1)^T \Sigma^{-1} X + C \quad (\Sigma^{-1})^T = \Sigma^{-1} \\
& W^T = (\mu_0 - \mu_1)^T \Sigma^{-1}
\end{aligned}$$

Inverse of symmetric matrix,  
is still symmetric.

$$W = \Sigma^{-1} (\mu_0 - \mu_1)$$

$$= \Sigma^{-1} \mu_0 - \Sigma^{-1} \mu_1$$

$$b = -\frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \ln \left( \frac{P(C_0)}{P(C_1)} \right)$$

- (c) In (b), we modeled the class conditional distribution with same covariance matrix  $\Sigma$ . Now let us consider two classes that have different covariance matrix as follows

$$P(x|C_0) = \frac{1}{(2\pi)^{n/2} |\Sigma_0|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) \right),$$

$$P(x|C_1) = \frac{1}{(2\pi)^{n/2} |\Sigma_1|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right).$$

Suppose we are able to find the maximum likelihood estimation of  $\mu_0, \mu_1, \Sigma_0, \Sigma_1, P(C_0)$ , and  $P(C_1)$ . Show that  $a = x^T A x + w^T x + b$  for some  $A, w$  and  $b$ . Find  $w$  and  $b$  in terms of  $\mu_0, \mu_1, \Sigma_0, \Sigma_1, P(C_0)$ , and  $P(C_1)$ . This shows that the decision boundary is quadratic.

$$\begin{aligned}
a &= \ln \left( \frac{\frac{1}{(2\pi)^{n/2} |\Sigma_0|^{1/2}} \cdot \exp \left( -\frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) \right) \cdot P(C_0)}{\frac{1}{(2\pi)^{n/2} |\Sigma_1|^{1/2}} \cdot \exp \left( -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right) \cdot P(C_1)} \right) \\
&= \ln \left( \frac{1}{|\Sigma_0|^{1/2}} \cdot \exp \left( -\frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) \right) - \ln \left( \frac{1}{|\Sigma_1|^{1/2}} \cdot \exp \left( -\frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) \right) \right) + \ln \left( \frac{P(C_0)}{P(C_1)} \right) \right) \\
&= \ln \left( \frac{1}{|\Sigma_0|^{1/2}} \right) - \frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) - \ln \left( \frac{1}{|\Sigma_1|^{1/2}} \right) + \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \ln \left( \frac{P(C_0)}{P(C_1)} \right) \\
&= -\frac{1}{2} (x^T - \mu_0^T) \Sigma_0^{-1} (x - \mu_0) + \frac{1}{2} (x^T - \mu_1^T) \Sigma_1^{-1} (x - \mu_1) + \ln \left( \frac{1}{|\Sigma_0|^{1/2}} \right) - \ln \left( \frac{1}{|\Sigma_1|^{1/2}} \right) + \ln \left( \frac{P(C_0)}{P(C_1)} \right) \\
&= -\frac{1}{2} (x^T \Sigma_0^{-1} - \mu_0^T \Sigma_0^{-1})(x - \mu_0) + \frac{1}{2} (x^T \Sigma_1^{-1} - \mu_1^T \Sigma_1^{-1})(x - \mu_1) + C \\
&= -\frac{1}{2} (x^T \Sigma_0^{-1} x - x^T \Sigma_0^{-1} \mu_0 - \mu_0^T \Sigma_0^{-1} x + \mu_0^T \Sigma_0^{-1} \mu_0) + \frac{1}{2} (x^T \Sigma_1^{-1} x - x^T \Sigma_1^{-1} \mu_1 - \mu_1^T \Sigma_1^{-1} x + \mu_1^T \Sigma_1^{-1} \mu_1) + C \\
&= -\frac{1}{2} (x^T \Sigma_0^{-1} x - 2x^T \Sigma_0^{-1} \mu_0 + \mu_0^T \Sigma_0^{-1} \mu_0) + \frac{1}{2} (x^T \Sigma_1^{-1} x - 2x^T \Sigma_1^{-1} \mu_1 + \mu_1^T \Sigma_1^{-1} \mu_1) + C \\
&= \frac{1}{2} x^T (\Sigma_1^{-1} - \Sigma_0^{-1}) x + x^T \Sigma_0^{-1} \mu_0 - x^T \Sigma_1^{-1} \mu_1 - \frac{1}{2} \mu_0^T \Sigma_0^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 + C
\end{aligned}$$

← constants

We transpose  $x^T \Sigma_0^{-1} \mu_0$  and  $x^T \Sigma_1^{-1} \mu_1$  as they are just numbers.

$$= x^T \left( \frac{1}{2} \Sigma_1^{-1} - \frac{1}{2} \Sigma_0^{-1} \right) x + \mu_0^T (\Sigma_0^{-1})^T x - \mu_1^T (\Sigma_1^{-1})^T x + C,$$

$$= x^T \left( \frac{1}{2} \Sigma_1^{-1} - \frac{1}{2} \Sigma_0^{-1} \right) x + (\mu_0^T \Sigma_0^{-1} - \mu_1^T \Sigma_1^{-1}) x + C,$$

$$A = \left( \frac{1}{2} \Sigma_1^{-1} - \frac{1}{2} \Sigma_0^{-1} \right)$$

$$W^T = \mu_0^T \Sigma_0^{-1} - \mu_1^T \Sigma_1^{-1}$$

$$W = \Sigma_0^{-1} \mu_0 - \Sigma_1^{-1} \mu_1$$

$$\begin{aligned}
b &= -\frac{1}{2} \mu_0^T \Sigma_0^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 \\
&\quad + \ln \left( \frac{|\Sigma_1|^{1/2}}{|\Sigma_0|^{1/2}} \right) + \ln \left( \frac{P(C_0)}{P(C_1)} \right)
\end{aligned}$$

4. We are given a training set  $\{(x^{(i)}, y^{(i)}); i = \{1, \dots, m\}\}$ , where  $x^{(i)} \in R^n$  and  $y^{(i)} \in \{0, 1\}$ . We consider the Gaussian Discriminant Analysis (GDA) model, which models  $P(x|y)$  using multivariate Gaussian. Writing out the model, we have:

$$P(y = 1) = \phi = 1 - P(y = 0)$$

$$P(x|y=0) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right)$$

$$P(x|y=1) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right)$$

The log-likelihood of the data is given by:

$$L(\phi, \mu_0, \mu_1, \Sigma) = \ln P(x^{(1)}, \dots, x^{(m)}, y^{(1)}, \dots, y^{(m)}) = \ln \prod_{i=1}^m P(x^{(i)}|y^{(i)})P(y^{(i)}).$$

In this exercise, we want to maximize  $L(\phi, \mu_0, \mu_1, \Sigma)$  with respect to  $\phi, \mu_0$ . The maximization over  $\Sigma$  is left for discussion.

- (a) Write down the explicit expression for  $P(x^{(1)}, \dots, x^{(m)}, y^{(1)}, \dots, y^{(m)})$  and  $L(\phi, \mu_0, \mu_1, \Sigma)$ .

$$\begin{aligned} P(x^1, \dots, x^m, y^1, \dots, y^m) &= \prod_{i=1}^m [P(x_i|y_i=0)P(y_i=0)]^{1-y_i} \cdot [P(x_i|y_i=1)P(y_i=1)]^{y_i} \\ &= \prod_{i=1}^m \left[ \frac{1-\phi}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu_0)^T \Sigma^{-1}(x_i - \mu_0)\right) \right]^{1-y_i} \cdot \left[ \frac{\phi}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu_1)^T \Sigma^{-1}(x_i - \mu_1)\right) \right]^{y_i} \\ L(\phi, \mu_0, \mu_1, \Sigma) &= \ln P(x^1, \dots, x^m, y^1, \dots, y^m) \\ &= \sum_{i=1}^m (1-y_i) \ln(P(x_i|y_i=0)P(y_i=0)) + y_i \ln(P(x_i|y_i=1)P(y_i=1)) \\ &= \sum_{i=1}^m (1-y_i) \ln\left(\frac{1-\phi}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu_0)^T \Sigma^{-1}(x_i - \mu_0)\right)\right) + y_i \ln\left(\frac{\phi}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu_1)^T \Sigma^{-1}(x_i - \mu_1)\right)\right) \\ &= \sum_{i=1}^m (1-y_i) \cdot \left[ \ln\left(\frac{1-\phi}{(2\pi)^{n/2}|\Sigma|^{1/2}}\right) - \frac{1}{2}(x_i - \mu_0)^T \Sigma^{-1}(x_i - \mu_0) \right] + y_i \left[ \ln\left(\frac{\phi}{(2\pi)^{n/2}|\Sigma|^{1/2}}\right) - \frac{1}{2}(x_i - \mu_1)^T \Sigma^{-1}(x_i - \mu_1) \right] \\ &= \sum_{i=1}^m (1-y_i) \cdot \left[ \ln(1-\phi) - \frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma|) - \frac{1}{2}(x_i - \mu_0)^T \Sigma^{-1}(x_i - \mu_0) \right] \\ &\quad + y_i \left[ \ln(\phi) - \frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma|) - \frac{1}{2}(x_i - \mu_1)^T \Sigma^{-1}(x_i - \mu_1) \right] \end{aligned}$$

- (b) Find the maximum likelihood estimate for  $\phi$ . How do you know such  $\phi$  is the “best” but not the “worst”? Hint: Show that the second derivative of  $L(\phi, \mu_0, \mu_1, \Sigma)$  with respect to  $\phi$  is negative.

$$\begin{aligned} L(\phi, \mu_0, \mu_1, \Sigma) &= \sum_{i=1}^m (1-y_i) \ln(1-\phi) + y_i \ln(\phi) + \dots \\ &= \ln(1-\phi) \cdot \sum_{i=1}^m (1-y_i) + \ln(\phi) \cdot \sum_{i=1}^m (y_i) + \dots \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial \phi} &= \frac{-1}{1-\phi} \cdot \sum_{i=1}^m (1-y_i) + \frac{1}{\phi} \cdot \sum_{i=1}^m (y_i) = 0 \\ -\frac{1}{1-\phi} \cdot \sum_{i=1}^m y_i &= \frac{1}{\phi} \cdot \sum_{i=1}^m (1-y_i) \quad \text{let } S_1 = \sum_{i=1}^m y_i \quad S_0 = \sum_{i=1}^m (1-y_i) \\ \frac{S_1}{\phi} &= \frac{S_0}{1-\phi} \quad S_1 - S_0 \phi = S_0 \phi \quad \frac{S_1}{S_0 + S_1} = \phi \end{aligned}$$

$$\frac{\partial^2 L}{\partial \phi^2} = -(1-\phi)^{-2} \cdot S_0 - \phi^{-2} \cdot S_1 = \frac{-1}{(1-\phi)^2} \cdot S_0 - \frac{1}{\phi^2} \cdot S_1 \leq 0$$

Since second derivative is always 0 or negative, we know that original function is always concave down, so when 1st derivative equals 0, must be a max.

- (c) Find the maximum likelihood estimate for  $\mu_0$ . How do you know such  $\mu_0$  is the "best" but not the "worst"? Hint: Show that the Hessian Matrix of  $L(\phi, \mu_0, \mu_1, \Sigma)$  with respect to  $\mu_0$  is negative definite. You may use the following: if  $A$  is positive definite, then  $A^{-1}$  is also positive definite. Also  $B$  is negative definite if  $-B$  is positive definite.

$$\begin{aligned}
 L(\phi, \mu_0, \mu_1, \Sigma) &= \sum_{i=1}^m (1-y_i) \left[ -\frac{1}{2} (\mathbf{x}_i - \mu_0)^T \Sigma^{-1} (\mathbf{x}_i - \mu_0) \right] + \dots \\
 &= \sum_{i=1}^m -\frac{1}{2} (1-y_i) \left[ (\mathbf{x}_i^T - \mu_0^T) \Sigma^{-1} (\mathbf{x}_i - \mu_0) \right] + \dots \\
 &= \sum_{i=1}^m -\frac{1}{2} (1-y_i) (\mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i - \mathbf{x}_i^T \Sigma^{-1} \mu_0 + \mu_0^T \Sigma^{-1} \mu_0) + \dots \\
 &= \sum_{i=1}^m (1-y_i) (\mu_0^T \Sigma^{-1} \mathbf{x}_i - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0) + \dots
 \end{aligned}$$

\* Rules for Matrix Calculus:

$\nabla_w (w^T b) = b$

$\nabla_w (w^T A w) = (A + A^T)w$

\*  $(\Sigma^{-1})^T = \Sigma^{-1}$  because  $\Sigma$  is symmetric.

$$\begin{aligned}
 \nabla_{\mu_0} L &= \sum_{i=1}^m \{(1-y_i)(\Sigma^{-1} \mathbf{x}_i)\} - \sum_{i=1}^m \{(1-y_i) \frac{1}{2} (\Sigma^{-1} + (\Sigma^{-1})^T) \mu_0\} \\
 &= \sum_{i=1}^m \{(1-y_i)(\Sigma^{-1} \mathbf{x}_i)\} - \sum_{i=1}^m \{(1-y_i) \Sigma^{-1} \mu_0\} \\
 &= \sum_{i=1}^m \{(1-y_i)(\Sigma^{-1} \mathbf{x}_i - \Sigma^{-1} \mu_0)\} = 0
 \end{aligned}$$

$$\sum_{i=1}^m \{(1-y_i)(\Sigma^{-1} \mathbf{x}_i)\} = \sum_{i=1}^m \{(1-y_i) \Sigma^{-1} \mu_0\}$$

$$\frac{\sum_{i=1}^m \{(1-y_i)(\Sigma^{-1} \mathbf{x}_i)\}}{\sum_{i=1}^m \{(1-y_i) \Sigma^{-1}\}} = \mu_0 \quad \mu_0 = \frac{\sum_{i=1}^m (1-y_i) \mathbf{x}_i}{\sum_{i=1}^m (1-y_i)} = \frac{\sum_{i=1}^m (1-y_i) \mathbf{x}_i}{S_0}$$

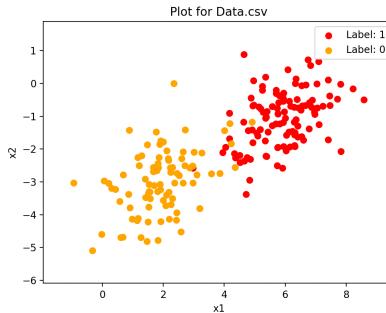
$$H_{\mu_0} = \nabla_{\mu_0}^2 L = \sum_{i=1}^m (1-y_i)(-\Sigma^{-1}) = -S_0 \Sigma^{-1}$$

From class, we know that  $\Sigma$  is a positive definite matrix, so  $\Sigma^{-1}$  is also positive definite.  $H_{\mu_0}$  is negative definite since  $-H_{\mu_0}$  is a positive definite matrix.

Since  $H_{\mu_0}$  is negative definite, we know  $\mu_0$  found by setting gradient to zero must be a MAX.

5. In this exercise, you will implement a binary classifier using the Gaussian Discriminant Analysis (GDA) model in MATLAB. The data is given in *data.csv*. The first two columns are the feature values and the last column contains the class labels.

- (a) Visualization. Plot the data from different classes in different colors. Is the data linearly separable?



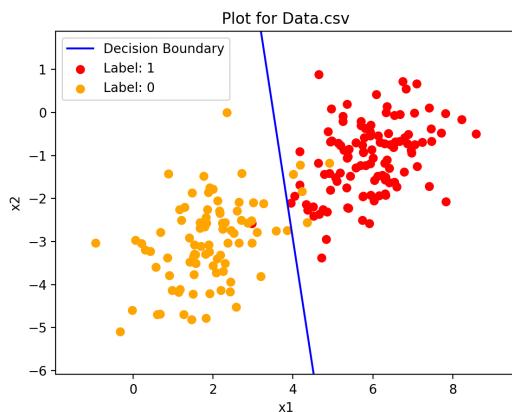
Data is not linearly separable.

- (b) In the GDA model, we assume the class label follows a Bernoulli distribution and we model the class conditional distribution as multivariate Gaussian with same covariance matrix ( $\Sigma$ ) and different means ( $\mu_0$  and  $\mu_1$ ). Find the maximum likelihood estimate of the parameters  $P(y = 0)$  (parameter for the Bernoulli distribution),  $\mu_0, \mu_1$  and  $\Sigma$  given this data set.

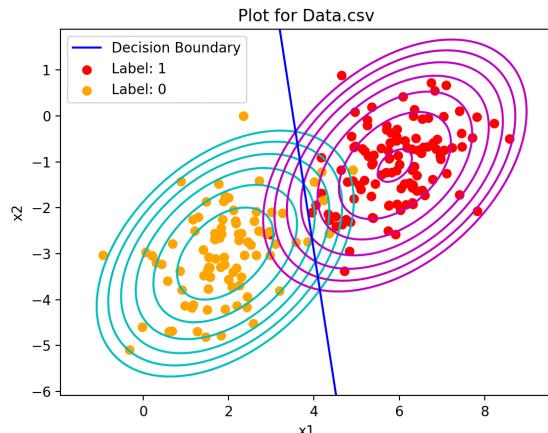
```
p(y = 0): 0.445
mu 0:
[[ 1.9195151]
 [-2.9972116]]
mu 1:
[[ 5.89818829]
 [-1.07926031]]
sigma:
[[[1.0180746  0.38866005]
 [0.38866005  0.80363858]]]
```

- (c) Using the result you find in Question 3 and your ML estimate of model parameters, find the decision boundary parameterized by  $w^T x + b = 0$ . Report  $w, b$  and plot the decision boundary on the same plot.

```
w:
[[ -3.67554649]
 [-0.60899668]]
b: 12.90499335183926
```



- (d) Visualize your results by plotting the contour of the two distributions  $P(x, y = 0)$  and  $P(x, y = 1)$ . For consistency, set 'LevelList' ('level' for python) to `logspace(-3,-1,7)`. Does your decision boundary pass through the points where the two distributions have equal probabilities ? Explain why.



Yes, decision boundary passes through points where two distributions have equal probabilities because decision boundary is defined by  
 $a = \ln\left(\frac{P(x|y=0)P(y=0)}{P(x|y=1)P(y=1)}\right) = \ln\left(\frac{P(x,y=0)}{P(x,y=1)}\right) = w^T x + b = 0$   
 When  $P(x,y=0) = P(x,y=1)$ , it's when two distributions have equal probability which means  $a = \ln(1) = 0$ , so whenever  $P(x,y=0) = P(x,y=1)$  the points will be on the decision boundary.