

Pavan Rathnakar Shetty

EECE 5644

HW3

Due date: November 16, 11:59 pm

GitHub: Files in the hw3 folder

<https://github.com/Pavan-r-shetty/5644.git>



1) Question 1

Given:  $p(x) = p(x|L=0)p(L=0) + p(x|L=1)p(L=1)$   
 $L$  is two class label

$$p(L=0) = 0.6$$

$$p(L=1) = 0.4$$

$$p(x|L=0) = w_0 g(x|w_0, c_0)$$

$$p(x|L=1) = w_1 g(x|w_1, c_1)$$

$$g(x|w, c) = \text{Multivariate Gaussian PDF}$$

$$w_1 = w_2 = 0.5$$

Class 0:  $p(L=0) = 0.6$

$$w_0 = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad w_0 = \begin{bmatrix} 0 \\ 4 \end{bmatrix}$$

$$c_0 = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} \quad c_0 = \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix}$$

Class 1:  $p(L=1) = 0.4$

$$w_1 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$c_1 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

4 datasets were generated with the given means and covariances

$D_{\text{train}}^{100}$  - 100 samples and their labels for training

$D_{\text{train}}^{1000}$  - 1000 samples and their labels for training

$D_{\text{train}}^{10000}$  - 10,000 samples and their labels for training

$D_{\text{validate}}^{20000}$  - 20,000 samples and their labels for validation.

Part 1: Theoretical optimal classifier that achieves minimum probability of error using the known parameters and the pdf.

Maximum expected risk classification wk in the form of likelihood ratio is given by :

$$(D=1) \quad \frac{p(X|L_1)}{p(X|L_0)} \geq \frac{\pi_{10} - \pi_{00}}{\pi_{01} - \pi_{11}} \times \frac{p(L_0)}{p(L_1)} = \gamma (D=0)$$

Minimize probability of misclassification by setting Penalty.

Penalty for wrong classification = 1

Penalty for correct classification = 0

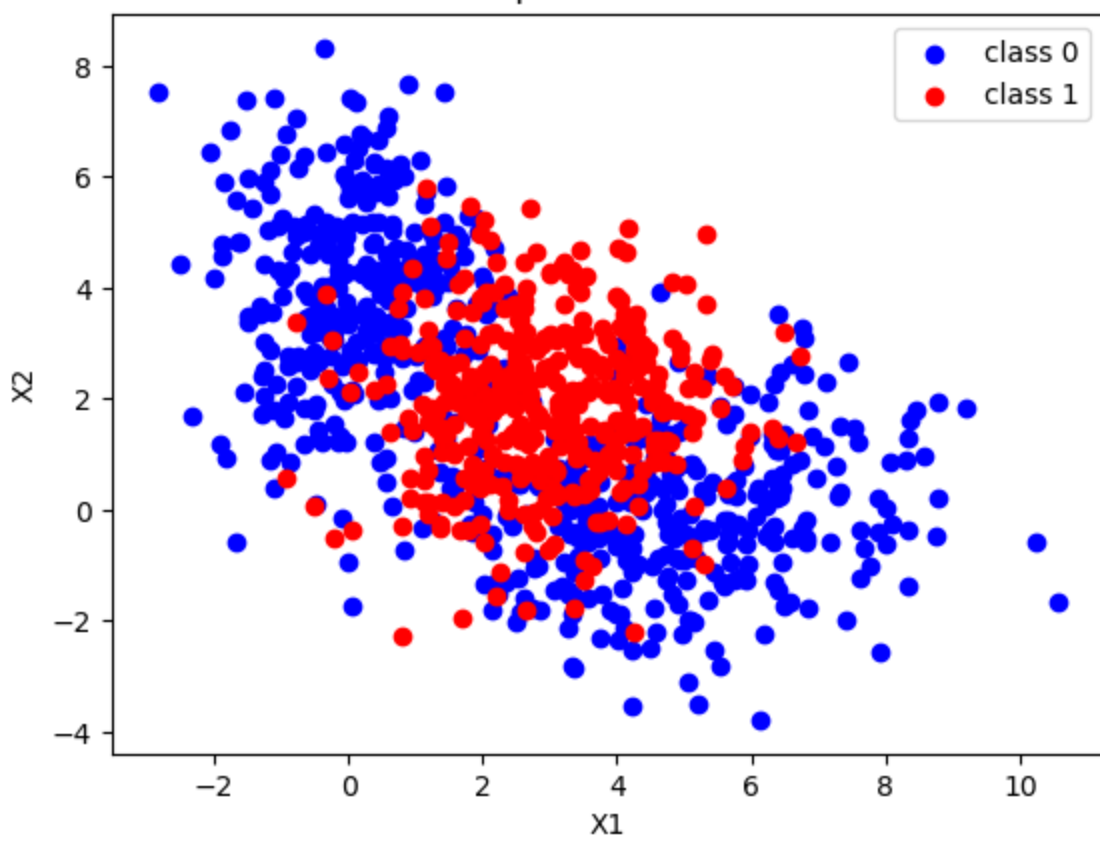
$$(D=1) \quad \frac{p(X|L_1)}{p(X|L_0)} \geq \frac{1.0}{1.0} \times \frac{0.6}{0.4} (D=0)$$

$$= 1.5 = \gamma_{\text{theoretical}}$$

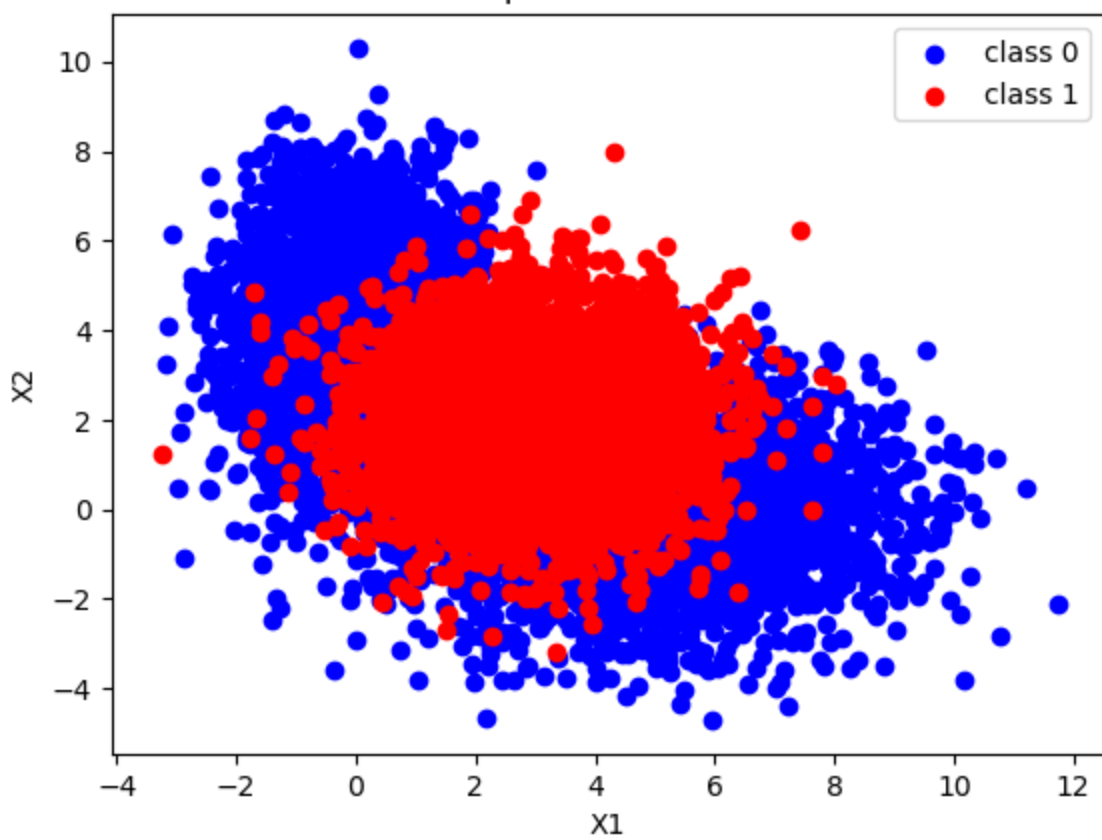
The classifier is implemented on the validation set and ROC curve is plotted as described :

|                     | $\gamma$ | Min Pen |
|---------------------|----------|---------|
| Theoretical         | 1.5      | 0.1740  |
| Estimated from data | 1.64     | 0.168   |

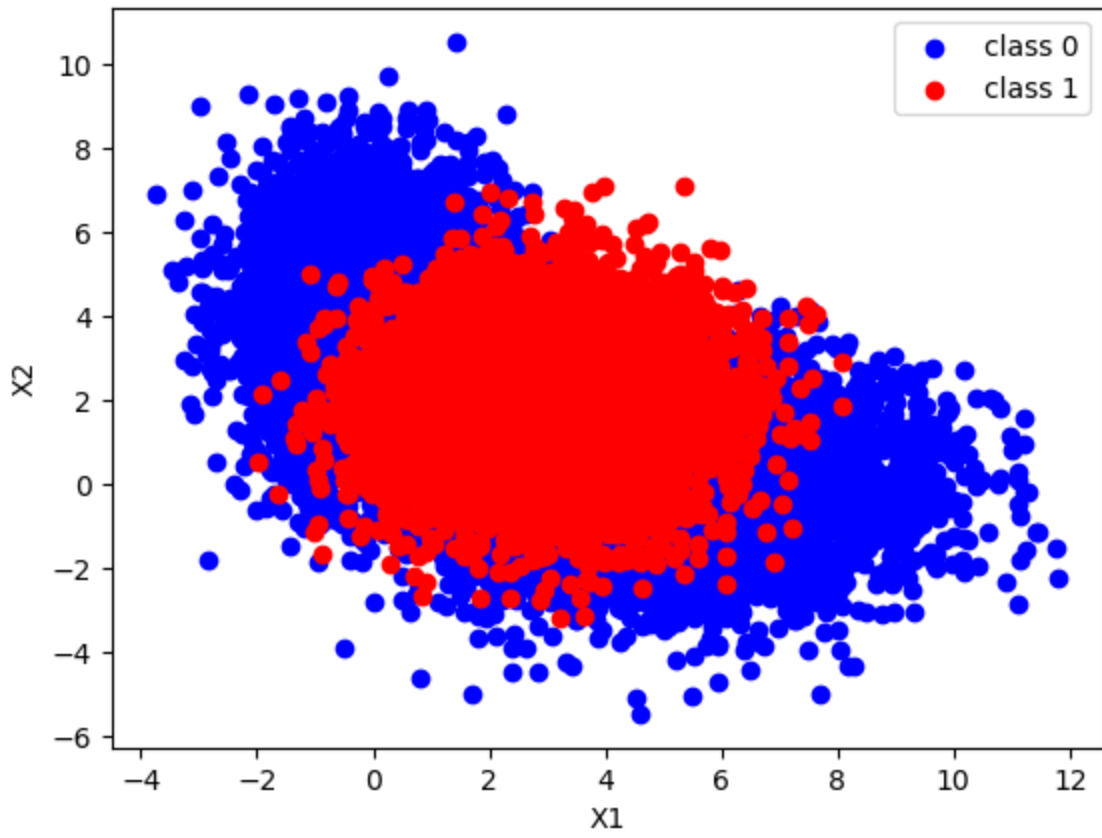
1000 datapoints from 2 classes

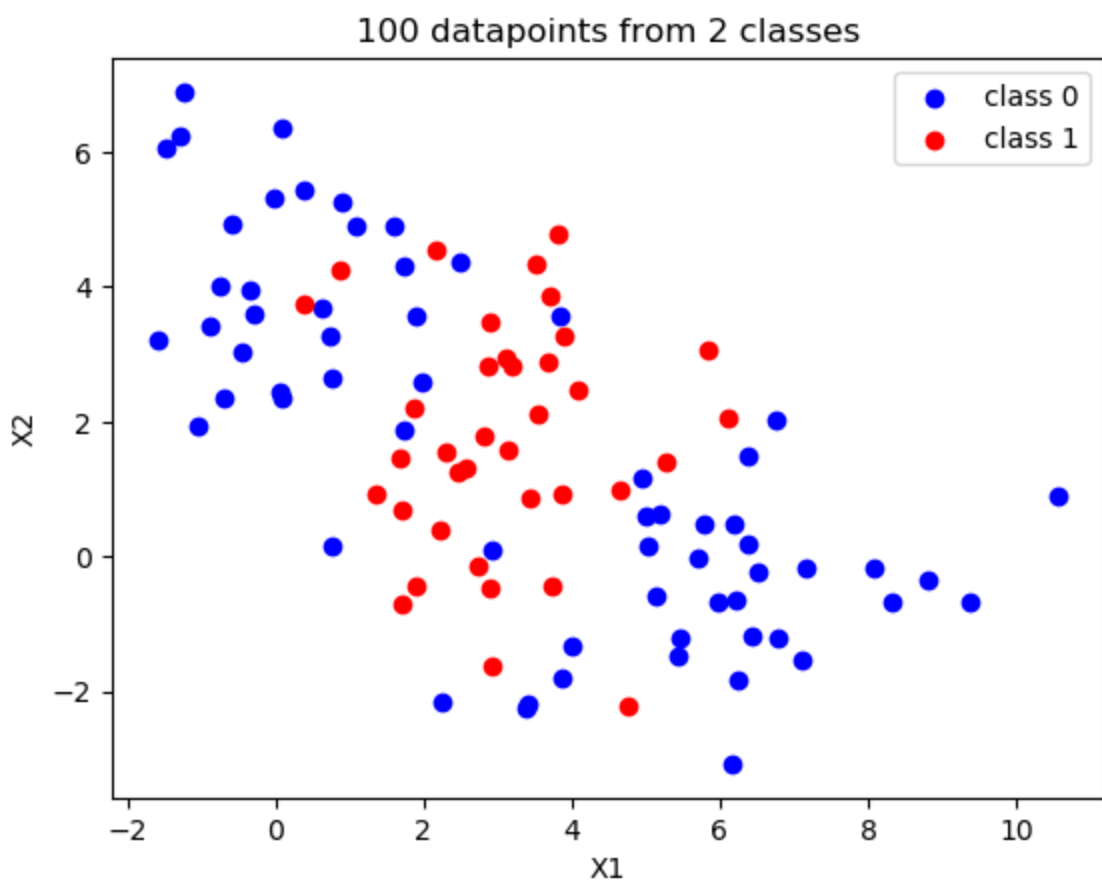


10k datapoints from 2 classes

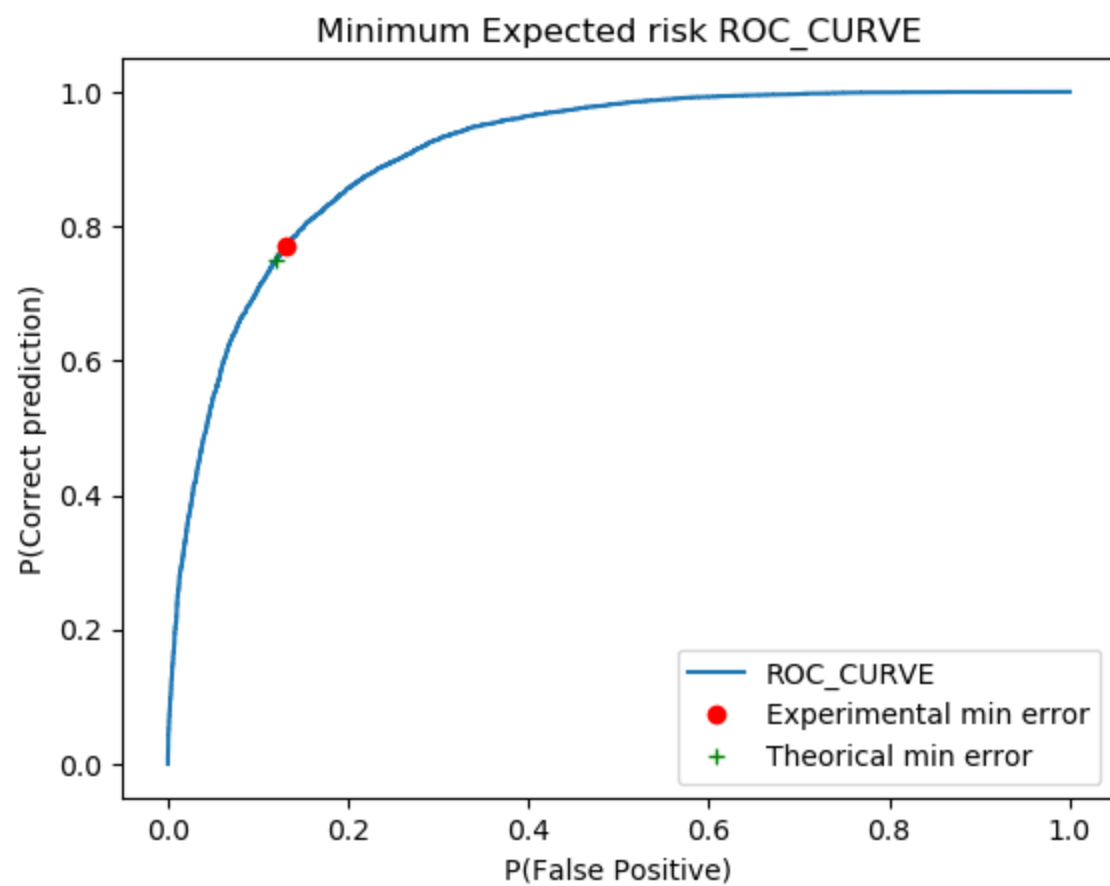


20k datapoints from 2 classes









## Part B:

ML estimation is performed on the three training sets.

Parameter Estimation for class 0 and class 1 was done using sklearn.mixture.  
Gaussian Mixture.

EM Estimation was based on the Maximum likelihood Estimation of the class prior as well as the mean and covariance of the conditional pdfs.

class 0:

$$\arg \max_{\theta} Q(\theta, \theta^0) = \sum_{c=1}^M \sum_{i=1}^N \ln(L_i) p(L_i | x_i; \theta^0) \\ + \sum_{c=1}^M \sum_{i=1}^N p(L_i | x_i; \theta) \ln(p_c(x_i | \theta_c))$$

Maximizing with respect to  $\mu$ ,  $\Sigma$ ,  $\pi$  we get:

$$\hat{\mu}_c = \frac{1}{N} \sum_{i=1}^N p_c(x_i | \mu_c, \Sigma_c, \pi_c)$$

$$\hat{\Sigma}_c = \frac{\sum_{i=1}^N x_i p_c(x_i | \mu_c, \Sigma_c, \pi_c)}{\sum_{i=1}^N p_c(x_i | \mu_c, \Sigma_c, \pi_c)}$$

$$\hat{\pi}_c = \frac{\sum_{i=1}^N p(L_i | x_i; \theta^0) (x_i - \mu_c) (x_i - \mu_c)^T}{\sum_{i=1}^N p(L_i | x_i; \theta^0)}$$

class 1: since it has 1 gaussian component, the maximum likelihood estimates are sample average and covariance

$$\hat{\mu} = \frac{1}{N} \sum x_i \quad \hat{\Sigma} = \frac{1}{N} \sum (x_i - \hat{\mu}) (x_i - \hat{\mu})^T$$

|          | $D^{trial}_{100}$                             | <u>Expanded Means</u><br>$D^{trial}_{1000}$    | $D^{trial}_{10000}$                           |
|----------|---|--|---|
| $w_{01}$ | $\begin{bmatrix} -0.13 \\ 3.91 \end{bmatrix}$ | $\begin{bmatrix} 5.20 \\ 0.13 \end{bmatrix}$   | $\begin{bmatrix} 5.11 \\ -0.09 \end{bmatrix}$ |
| $w_{02}$ | $\begin{bmatrix} 4.94 \\ -0.08 \end{bmatrix}$ | $\begin{bmatrix} -0.017 \\ 3.89 \end{bmatrix}$ | $\begin{bmatrix} 0.026 \\ 3.95 \end{bmatrix}$ |
| $w_{03}$ | $\begin{bmatrix} 2.58 \\ 2.081 \end{bmatrix}$ | $\begin{bmatrix} 2.91 \\ 2.06 \end{bmatrix}$   | $\begin{bmatrix} 2.95 \\ 1.98 \end{bmatrix}$  |

|            | $D^{trial}_{100}$  | $D^{trial}_{1000}$   | $D^{trial}_{10000}$  |
|------------|--|--|--|
| $cov_{01}$ | $\begin{bmatrix} 0.86 & -0.35 \\ -0.35 & 2.66 \end{bmatrix}$ | $\begin{bmatrix} 3.019 & 0.130 \\ 0.130 & 1.887 \end{bmatrix}$ | $\begin{bmatrix} 3.65 & 0.049 \\ 0.049 & 1.95 \end{bmatrix}$ |
| $cov_{02}$ | $\begin{bmatrix} 4.94 & 0.137 \\ 0.137 & 1.79 \end{bmatrix}$ | $\begin{bmatrix} 0.93 & 0.11 \\ 0.11 & 3.15 \end{bmatrix}$     | $\begin{bmatrix} 1.03 & -0.04 \\ -0.04 & 3.10 \end{bmatrix}$ |
| $cov_{03}$ | $\begin{bmatrix} 2.16 & -0.66 \\ -0.66 & 1.91 \end{bmatrix}$ | $\begin{bmatrix} 1.90 & -0.02 \\ -0.02 & 2.02 \end{bmatrix}$   | $\begin{bmatrix} 2.03 & 0.04 \\ 0.04 & 2.01 \end{bmatrix}$   |

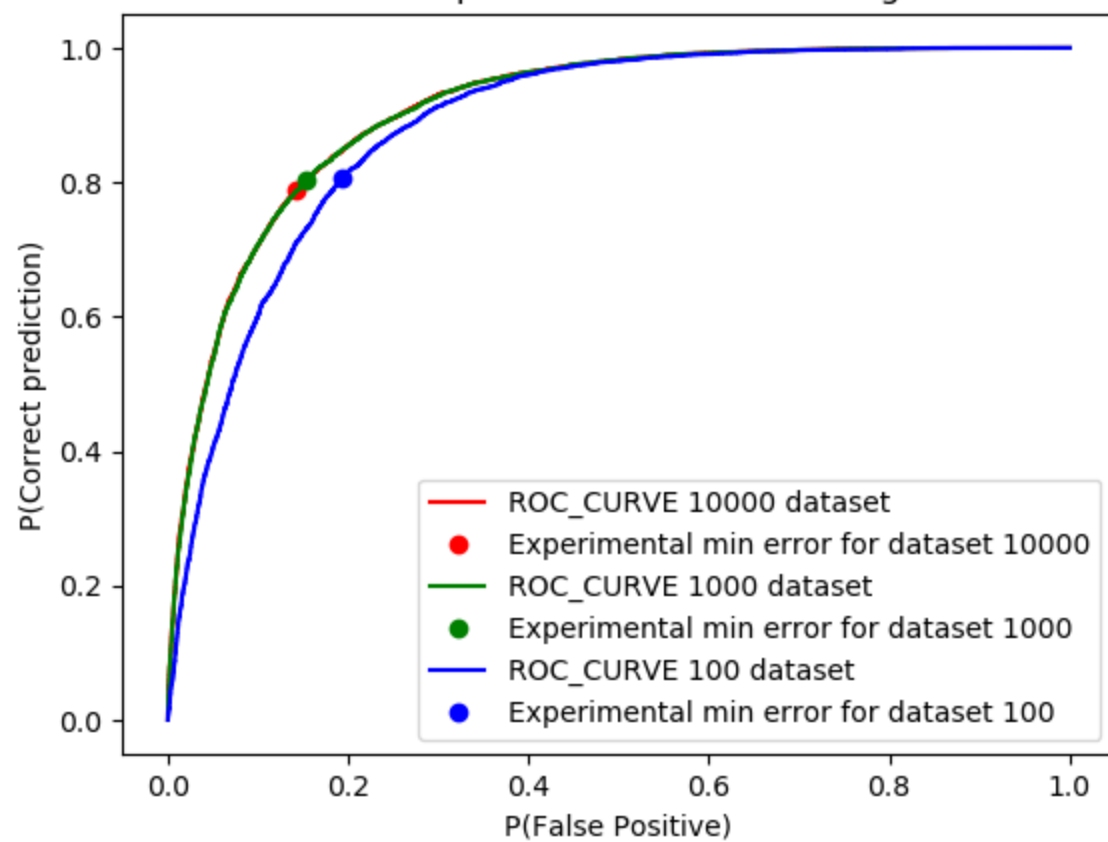
|          | <u>Estimated Alphas</u>  |                           |                            |
|----------|--------------------------|---------------------------|----------------------------|
|          | $D^{\text{train}}_{100}$ | $D^{\text{train}}_{1000}$ | $D^{\text{train}}_{10000}$ |
| $L_{01}$ | 0.46                     | 0.515                     | 0.49                       |
| $L_{02}$ | 0.55                     | 0.484                     | 0.50                       |

|          | <u>Prior Probabilities</u> |                           |                            |
|----------|----------------------------|---------------------------|----------------------------|
|          | $D^{\text{train}}_{100}$   | $D^{\text{train}}_{1000}$ | $D^{\text{train}}_{10000}$ |
| $P(L=0)$ | 0.590                      | 0.58                      | 0.59                       |
| $P(L=1)$ | 0.402                      | 0.41                      | 0.40                       |

Minimum Probability Error for the datasets

| Training Samples           | $D^{\text{train}}$ | Min Probability Error |
|----------------------------|--------------------|-----------------------|
| $D^{\text{train}}_{100}$   | 1.49               | 0.1823                |
| $D^{\text{train}}_{1000}$  | 1.54               | 0.179                 |
| $D^{\text{train}}_{10000}$ | 1.53               | 0.176                 |

Minimum Expected risk roc for training data





### Part C: classifier using logistic function

Maximum likelihood estimation techniques were used to train logistic linear and logistic quadratic based approximation of class label posterior functions given sample.

Logistic function:

$$h(x, w) = \frac{1}{1 + e^{-w^T z(x)}}$$

for linear function  $z(x) = [1, x_1, x_2]^T$   
 for quadratic  $z(x) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2]^T$

The vectors are estimated using numerical optimization techniques with the cost function

$$\hat{w}_{MLE} = -\frac{1}{N} \sum_{i=1}^N \ln \left( h(x_i, 0) \right)^{y_i} (1 - h(x_i, 0))^{(1 - y_i)}$$

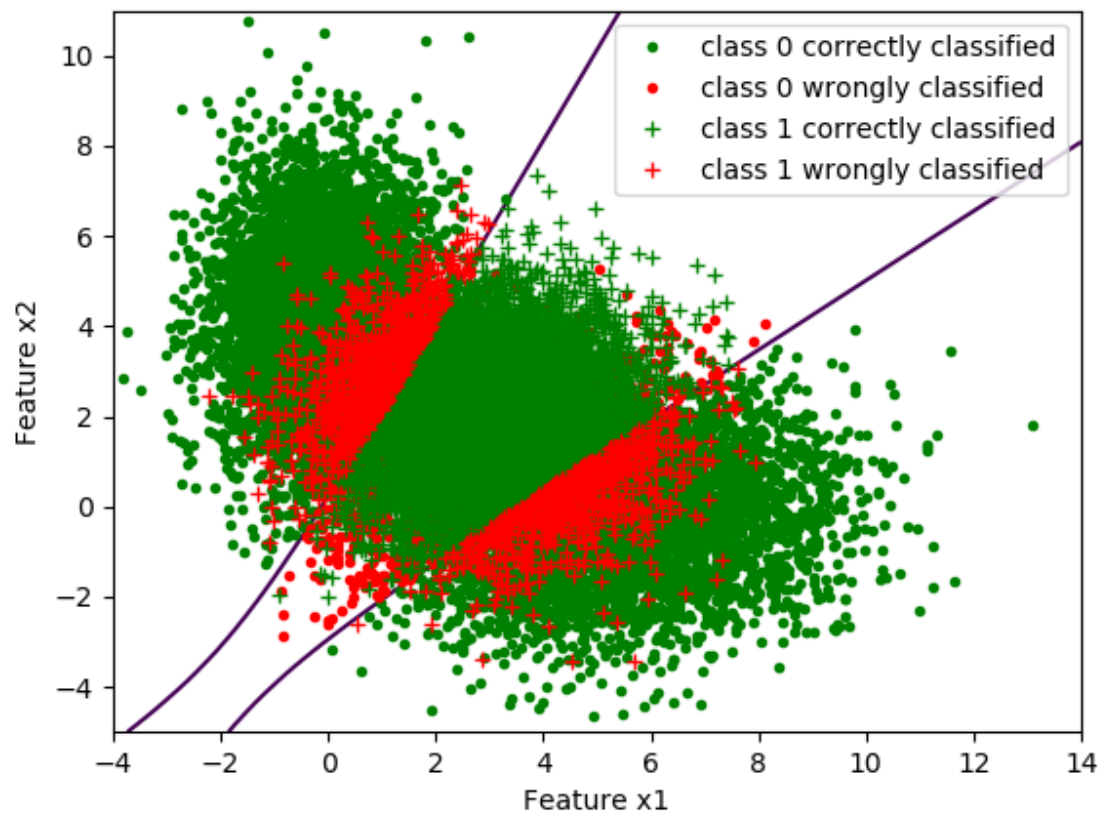
The minimum expected risk classification criteria is

$$(y_i = 1) \quad w^T z(x) \geq 0 \quad (y_i = 0)$$

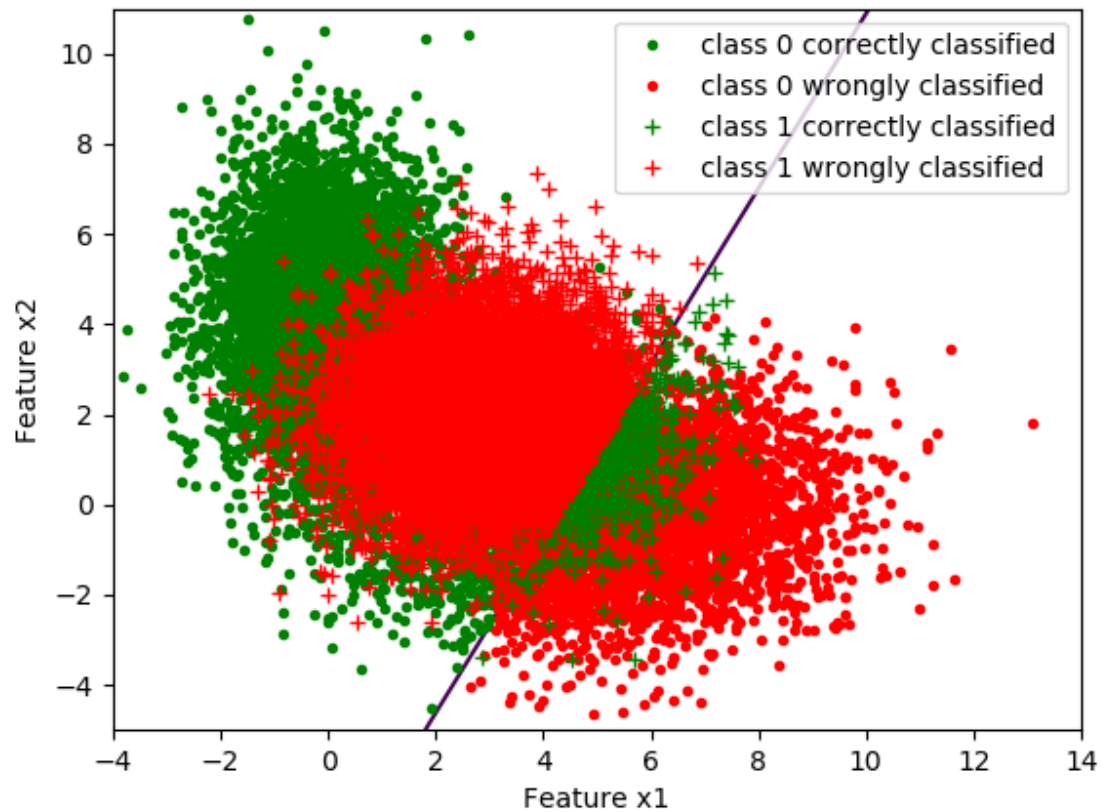
logistic function Probability error

|                    | Pen <sub>error</sub> - Linear | Pen <sub>error</sub> - Quadratic |
|--------------------|-------------------------------|----------------------------------|
| $D_{train, 100}$   | 0.46                          | 0.18                             |
| $D_{train, 1000}$  | 0.12                          | 0.17                             |
| $D_{train, 10000}$ | 0.04                          | 0.17                             |

Distribution after classification overlapped by decision boundaries

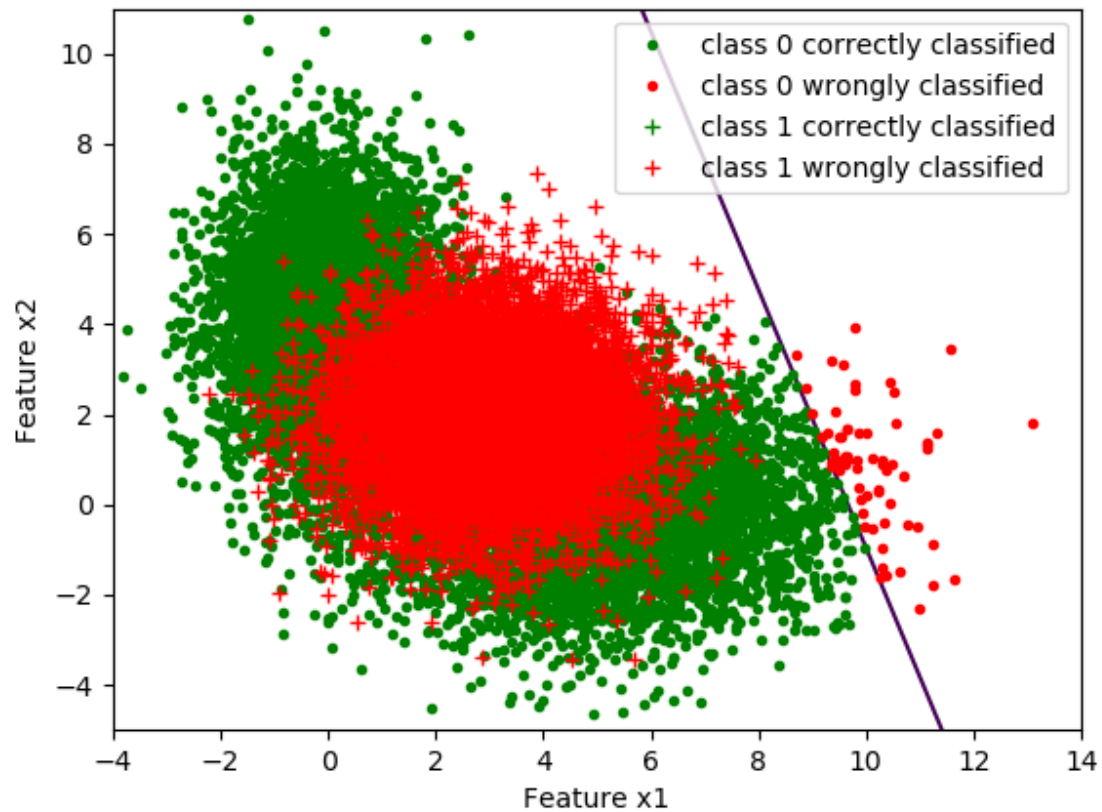


Distribution after classification overlapped by decision boundaries

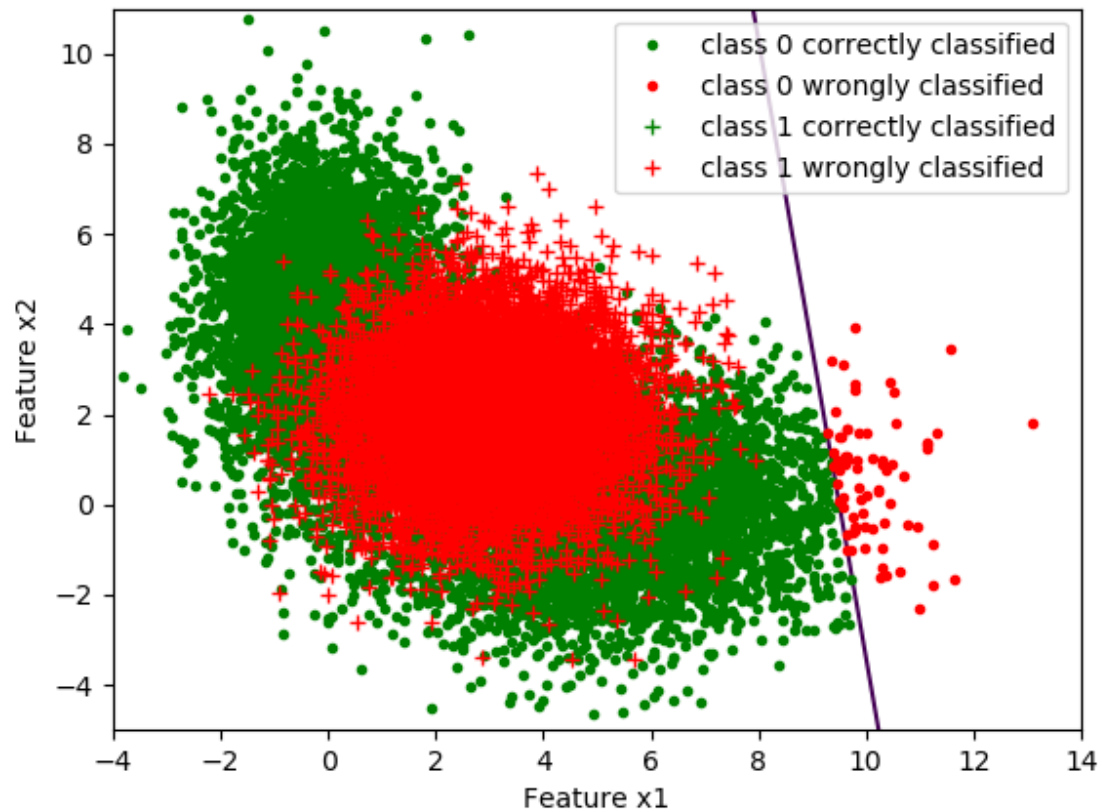




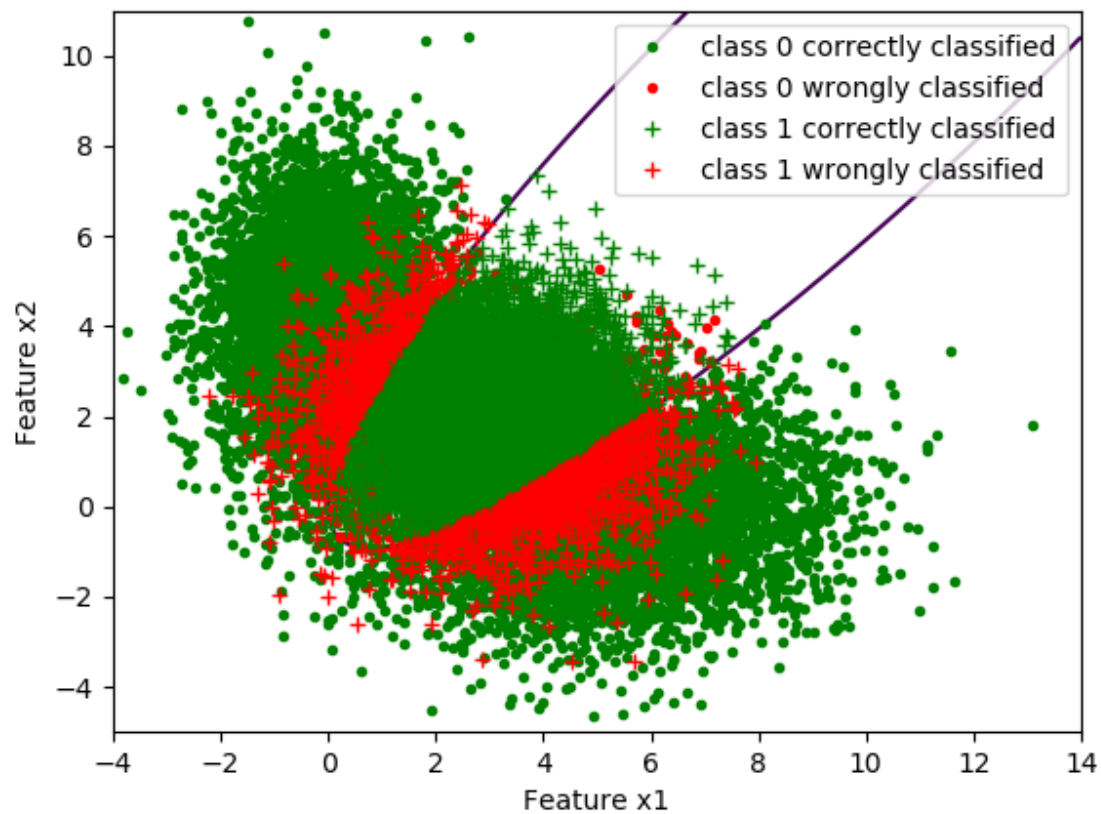
Distribution after classification overlapped by decision boundaries



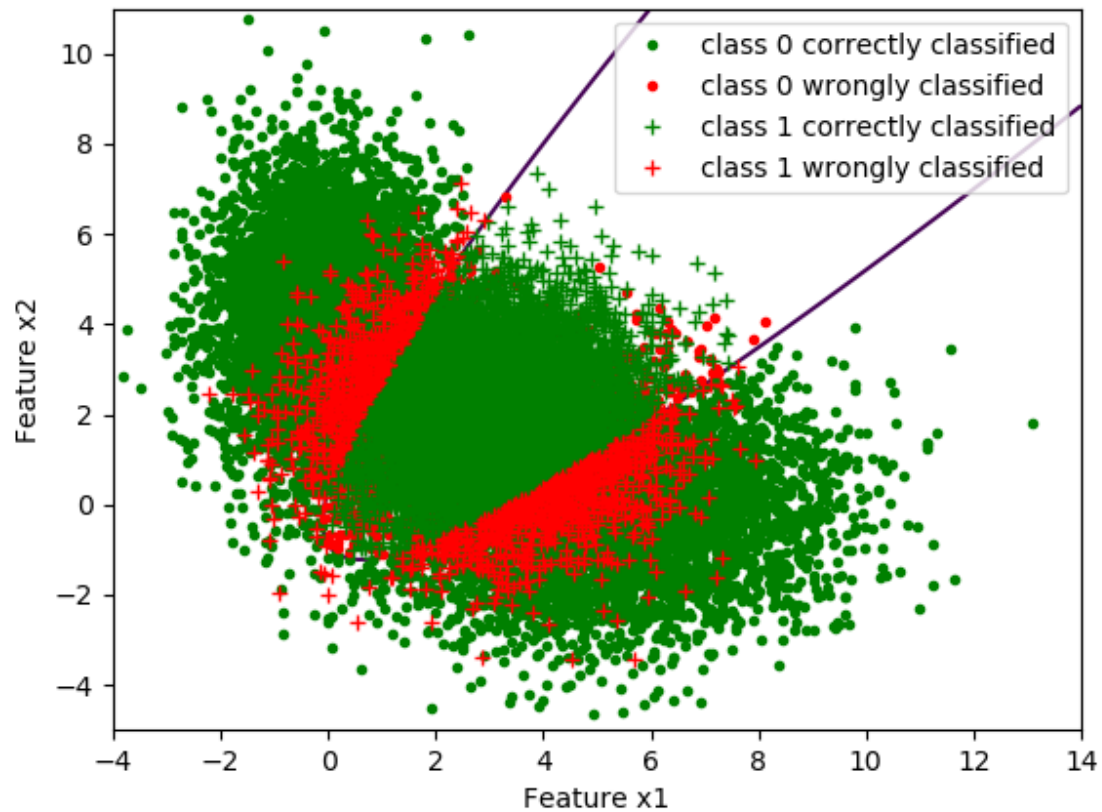
Distribution after classification overlapped by decision boundaries



Distribution after classification overlapped by decision boundaries



Distribution after classification overlapped by decision boundaries



Question 2:

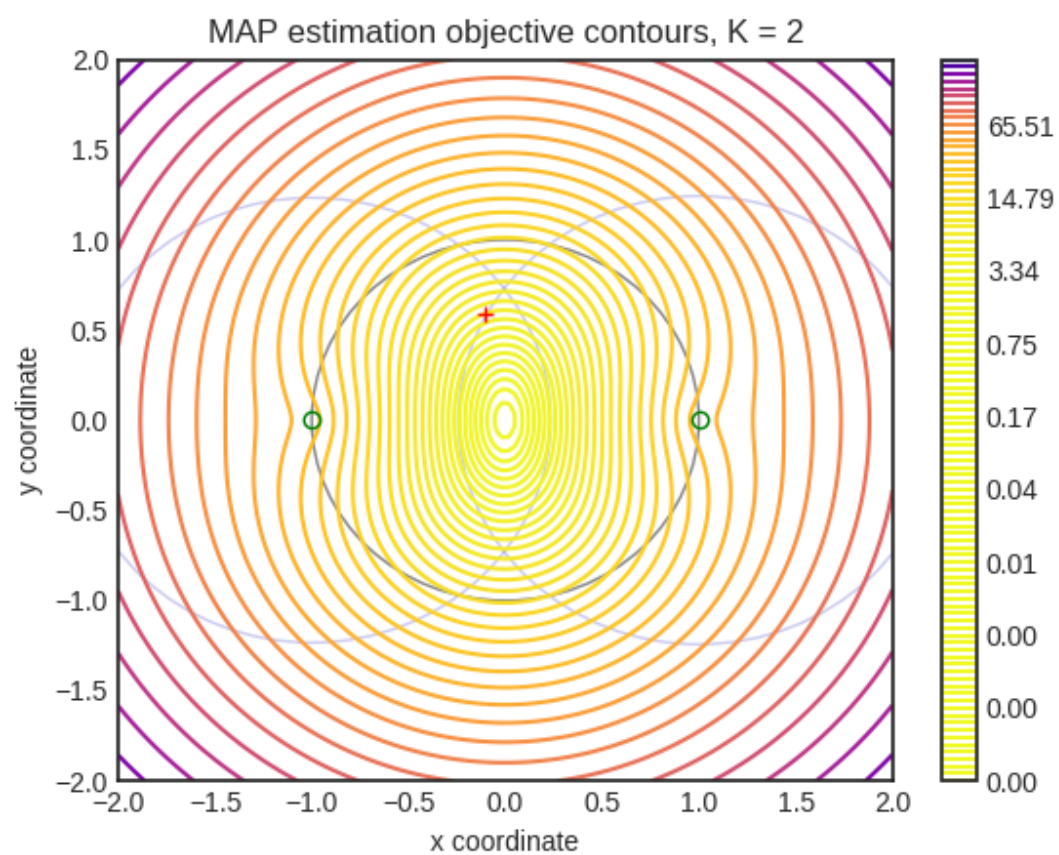
To find:  $[x, y]^T$  coordinate position with highest probability.

given: Prior distribution of range of measurements for each of the  $k$  reference coordinates.

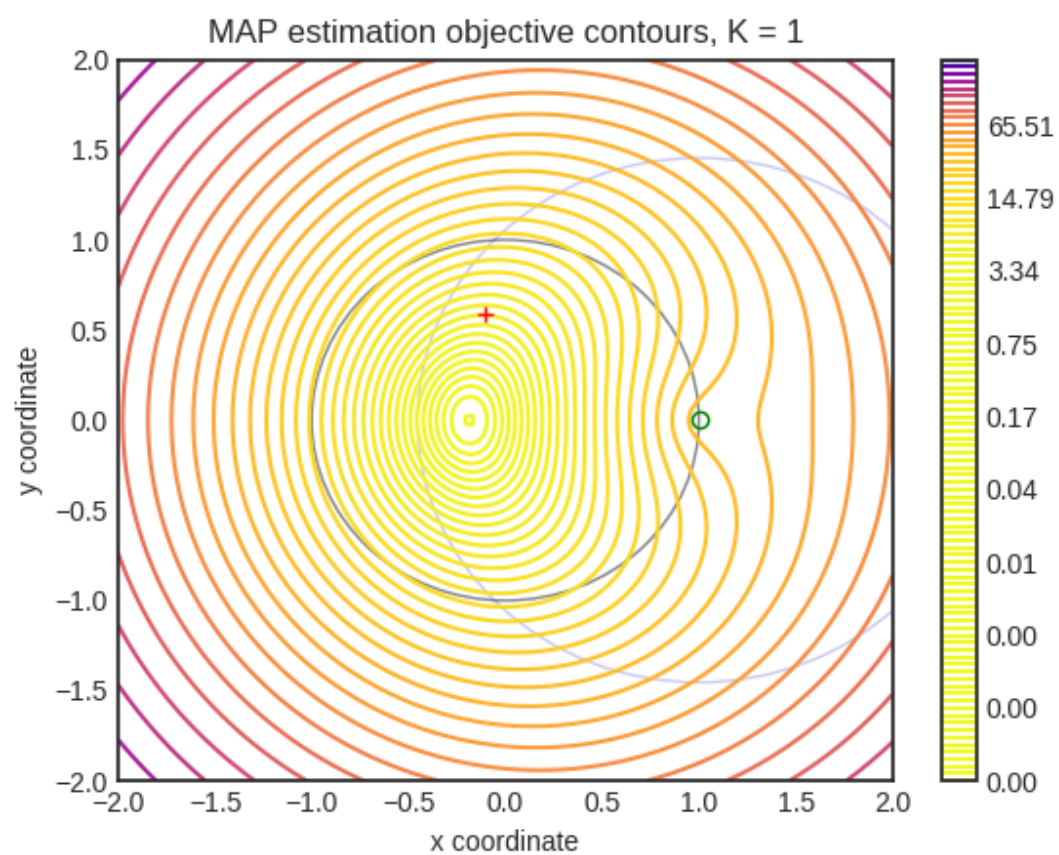
$$\begin{aligned} \begin{bmatrix} x_{map} \\ y_{map} \end{bmatrix} &= \underset{\begin{bmatrix} x \\ y \end{bmatrix}}{\operatorname{argmax}} p[\hat{y}] \{r_1, \dots, r_k\} \\ &= \underset{\begin{bmatrix} x \\ y \end{bmatrix}}{\operatorname{argmax}} \left( \frac{1}{\sqrt{2\pi} \sigma_x \sigma_y} e^{-\frac{1}{2} [x, y] \begin{bmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{bmatrix}^{-1} \begin{bmatrix} x \\ y \end{bmatrix}} \right) \prod_{i=1}^k p[\hat{y}] | r_i \end{aligned}$$

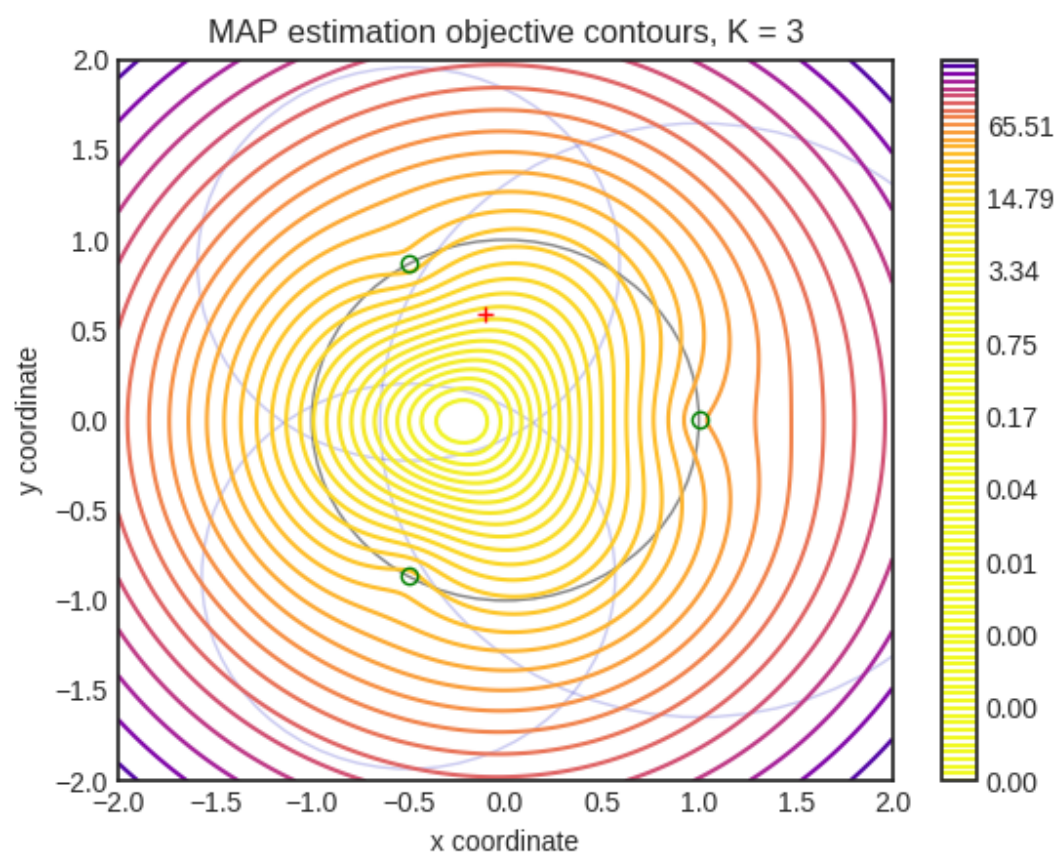
$$= \underset{\begin{bmatrix} x \\ y \end{bmatrix}}{\operatorname{argmax}} -\frac{1}{2} [x, y] \begin{bmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{bmatrix}^{-1} \begin{bmatrix} x \\ y \end{bmatrix} + \sum_{i=1}^k -\frac{(r_i d_i)^2}{2\sigma_i^2}$$

$$\text{where } d_i = \left\| \begin{bmatrix} x \\ y \end{bmatrix} - \begin{bmatrix} x_i \\ y_i \end{bmatrix} \right\|^2$$

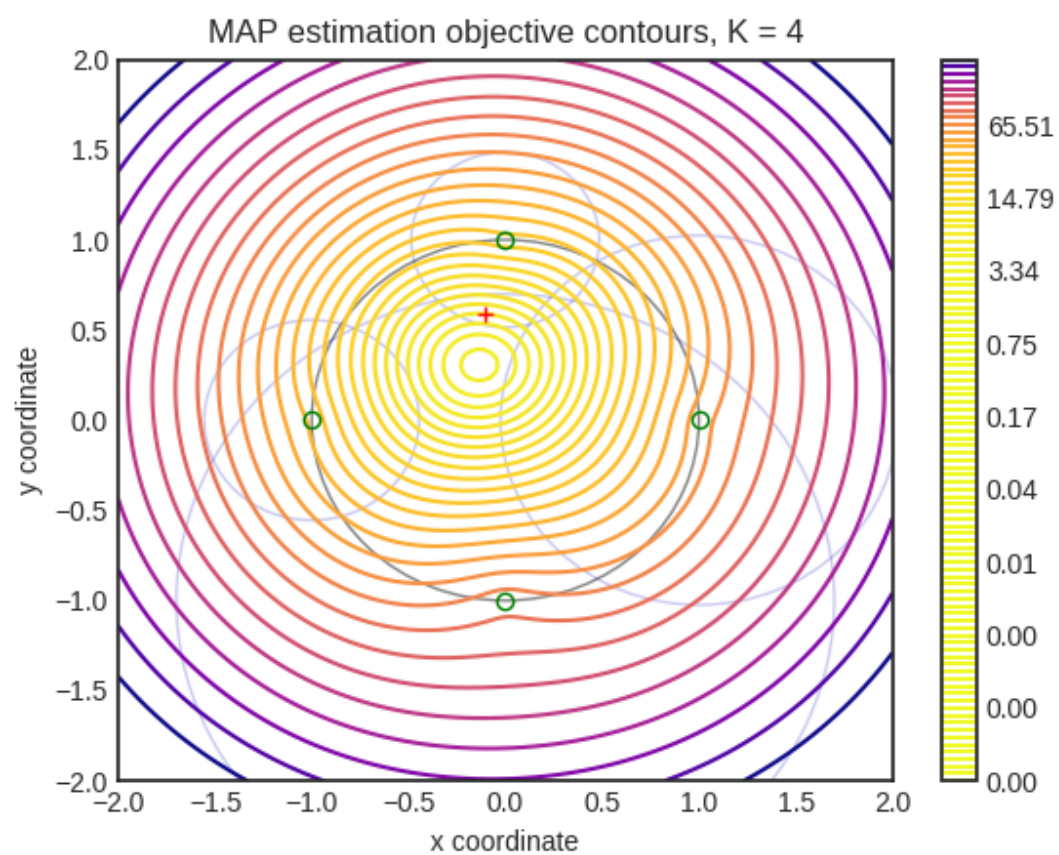












### Explanation of working of the code:

- A random coordinate is generated at the origin as true coordinate.
- For every  $k$ ,  $k$  evenly distributed landmarks around the circle are used to calculate distance from the true vehicle position.
- These measurements have Gaussian white noise.
- The MAP estimation objective function is calculated for all points on  $128 \times 128$  mesh grid, the values are then plotted as equilevel contour.
- Additionally, the unit circle (shown in gray) true location (red +), landmark location (in green) and the range reported by each landmark are plotted. (shown in faint blue circle)



### Behaviour of MAP estimate of $\mu$

- Map estimate for  $k < 3$  is not accurate, estimates are symmetric around  $x$ , since the prior bias have  $y$  coordinate  $= 0$ .
- The estimate is much accurate of  $k=3$  &  $k=4$ .
- In general as  $k$  increases, estimate gets better, however it is not always true for  $k: 1 \rightarrow 2$ .
- In general as  $k$  increases, the estimator increases. This can be visualised on the contour graph by a shrinkage of area of location with a high probability.





### Question 3

$$\lambda(x_i | w_j) = \begin{cases} 0 & i \neq j \\ \lambda_r & i = j \\ \lambda_s & \text{otherwise} \end{cases} \quad i, j = 1, \dots, c$$

to choose class with minimum risk, we have

$$\text{decision} = \underset{j=1, \dots, c}{\operatorname{argmin}} \left\{ \lambda_{ij} P(w_i | x) \right\}$$
$$j = 1, \dots, c$$

we know that  $\lambda_{ii} = 0$ ,

$$\therefore \operatorname{argmin} \lambda_{i1} P(w_1 | x) \dots 0, \dots, \lambda_{ic} P(w_c | x)$$

$\therefore$  minimum possible value is 0 when  $j = i$ .



But  $j=i$ , will be chosen only in  $P(w_i|x)$  is the highest.

$\therefore$  to choose  $j=i$ ,

$$P(w_i|x) \geq P(w_j|x) \text{ for } \forall j \in \{1, \dots, c\}$$

The average risk of choosing class  $w_i$  is given by

$$R(D|x) = \sum_{j=1}^c \lambda_{ij} \frac{P(x|w_j)P(w_j)}{P(x)}$$

$$R(D=w_i|x) = \sum_{j=1}^c \lambda_{ij} P(w_j|x)$$

$$= \sum_{j=1}^{i-1} \lambda_{ij} P(w_j|x) + \lambda_{ii} P(w_i|x) + \sum_{j=i+1}^c \lambda_{ij} P(w_j|x)$$

$$\text{w.k.t } \lambda_{ii} = 0$$

$$\therefore R(D=w_i|x) = \sum_{j=1, j \neq i}^c \lambda_{ij} P(w_j|x) + 0$$

$$= \lambda_{i\cdot} \sum_{j=1, j \neq i}^c P(w_j|x) = \lambda_{i\cdot} (1 - P(w_i|x))$$

$$\therefore P(w_i|x) \uparrow \text{ then risk } \downarrow$$

for  $i=c+1$ , the risk is given by

$$R(D=c+1|x) = \lambda_{r\cdot}$$



Hence for achieving minimum risk (choose  $w_0$ )

we have

$$P(D=0|X) \leq P(D=1|X)$$

$$P_D(1 - P(w_1|X)) \leq 1$$

$$1 - P(w_1|X) \leq \frac{P_r}{P_s}$$

$$\therefore P(w_1|X) = 1 - \frac{P_r}{P_s}$$

Case 1:  $P_r = 0$

$$\Rightarrow P(w_1|X) \geq 1 \Rightarrow P(w_1|X) = 0 \text{ at } P_r = 0$$

$\Rightarrow$  reject always as cost of rejecting is 0

Case 2:  $P_r \geq P_s$

$$\Rightarrow \frac{P_r}{P_s} \geq 1$$

$$\therefore P(w_1|X) \geq 1 - \frac{P_r}{P_s}$$

$$\Rightarrow P(w_1|X) \leq \frac{P_r}{P_s} - 1$$

Implies cost of rejection is higher than cost of choosing any other value, then we never reject.