

1. Universal Function Approximation

a) $\mathcal{G}_M = \{x_i = i/(M-1)\}_{i=0}^{M-1}$

How many parameters are in the hidden layer? How many total parameters in the ANN? Specify how each of the weights and biases in both layers should be chosen to obtain the approximation.

For the hidden layer, there would be $2(M-1)$ weights.

For the second layer, there would be M weights. So, for the ANN there would be $3M-2$ parameters.

The approximation function is:

$$\hat{f}(x) = w_{10}^{(2)} v_0(x) + \sum_{m=1}^{M-1} w_{im}^{(2)} v_m(x)$$

$$\text{In the problem's case, } v_m(x) = h_{\text{ReLU}}(w_{m1}^{(1)} x + w_{mo}^{(1)})$$

$$\text{Let ReLU parameter } b=1, \text{ let } w_{m1}^{(1)} = 1, \forall m \Rightarrow v_m(x) = \max\{0, x + w_{mo}^{(1)}\}$$

$w_{mo}^{(1)}$ is the bias in layer 1.

Order the grid so that $x_0 < x_1 < \dots < x_{N-1}$

Let $\hat{f}_{n-1}(x)$ be the function approximation based on data points x_0, x_1, \dots, x_{N-1}

$$\text{Set } \hat{f}_0(x) = y_0, w_{10}^{(2)} = y_0, v_0(x) = 1.$$

For $\hat{f}_1(x)$, use $\hat{f}_0(x)$ and (x_1, y_1)

$$\therefore \hat{f}_1(x) = \hat{f}_0(x) + w_{11}^{(2)} \max\{0, x + w_{10}^{(1)}\}$$

$$\text{Choose } w_{10}^{(1)} = -x_0$$

$$\therefore \hat{f}_1(x) = \hat{f}_0(x) + w_{11}^{(2)} \underbrace{v_1(x)}_{\max\{0, x_1 - x_0\}} = x_1 - x_0$$

$$\text{To make } \hat{f}_1(x) = y_1, w_{11}^{(2)} = \frac{y_1 - y_0}{x_1 - x_0}$$

Iterate until all the grid point, and get the final approximation

$$\hat{f}_n(x) = \hat{f}_{n-1}(x) + w_{1n}^{(2)} v_n(x), v_n(x) = \{0, x - x_{n-1}\}$$

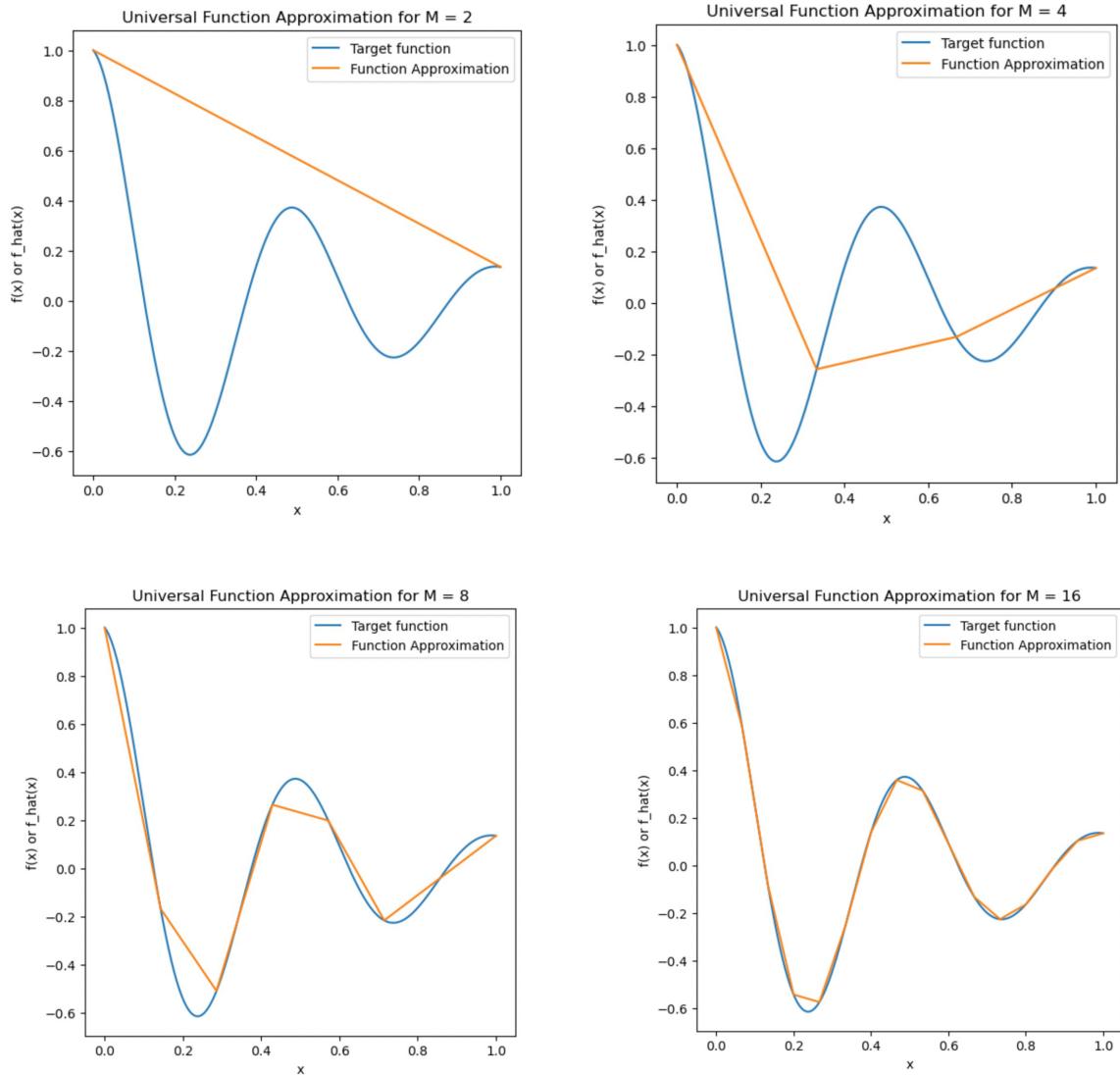
$$\text{To make } \hat{f}_n(x) = y_n, w_{1n}^{(2)} = \frac{y_n - \hat{f}_{n-1}(x_n)}{x_n - x_{n-1}}$$

To conclude:

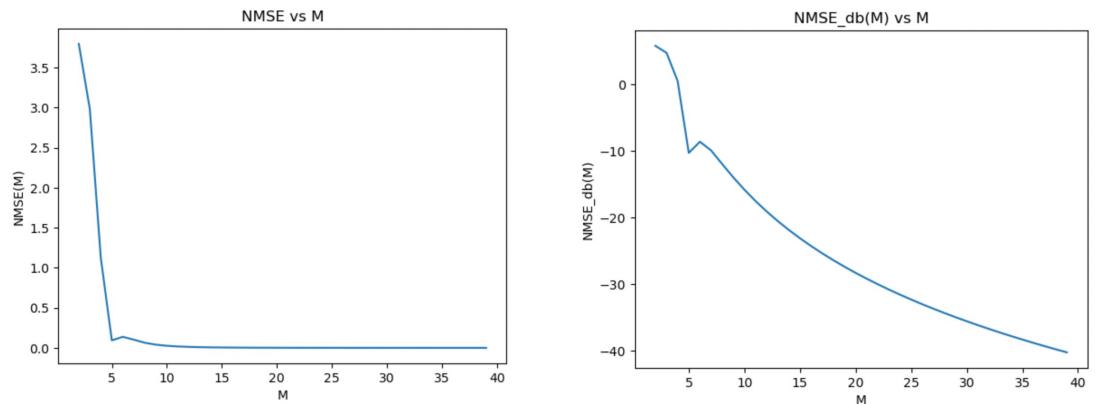
$$\text{Layer 1 weights} = 1, \text{ bias } w_{mo}^{(1)} = -x_{m-1} \quad \forall m \geq 1$$

$$\text{Output layer weights} = w_{1n}^{(2)} = \frac{y_n - \hat{f}_{n-1}(x_n)}{x_n - x_{n-1}}, n \geq 2, \text{ bias } w_{10}^{(2)} = y_0$$

(b)



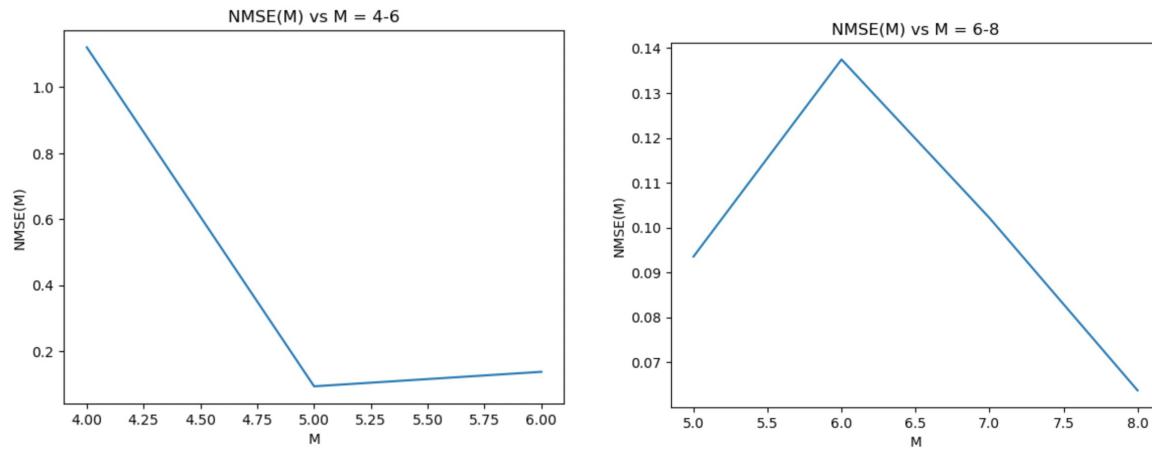
(c)



(d) Interpretation of the NMSE results:

(i) Yes. This means that the approximation function is higher or lower than the original function. This is caused by an approximation function with a constant output regardless the input value.

(ii) No.



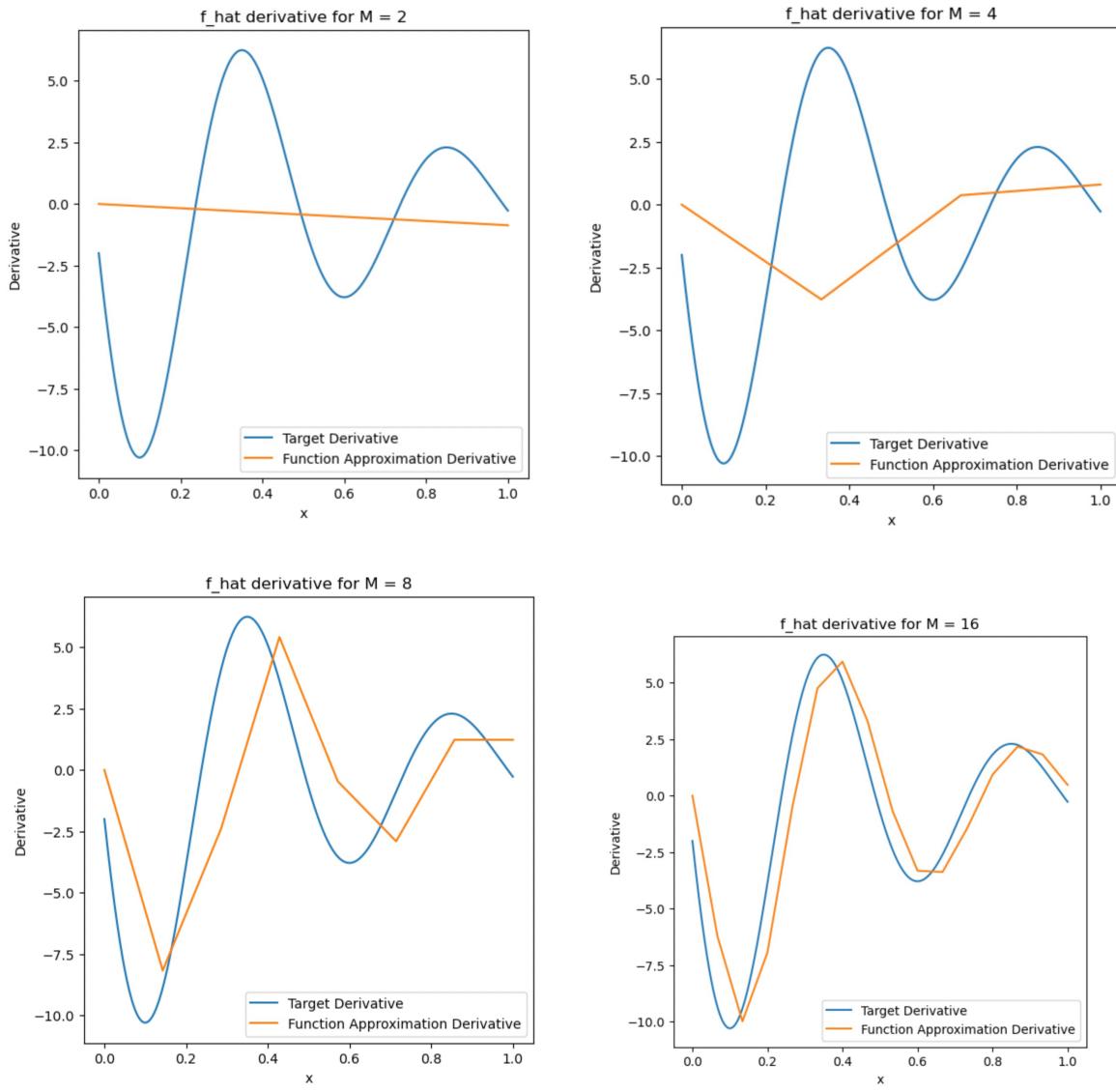
This is because the approximation function fit the original function too well but there exists a bump in the original function. When the approximation function increases after training the left hand side points of the bump, the original function decreases immediately and increases the NMSE. We are using the mean squared error (MSE) and normalized mean squared error (NMSE).

(iii)

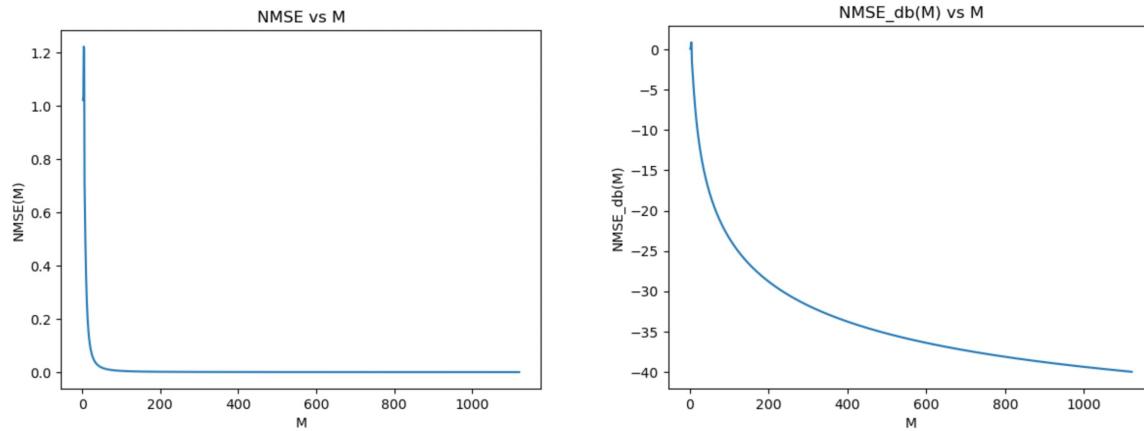
$M = 39$

Total = $3 * M - 2 = 72$ parameters.

(e)



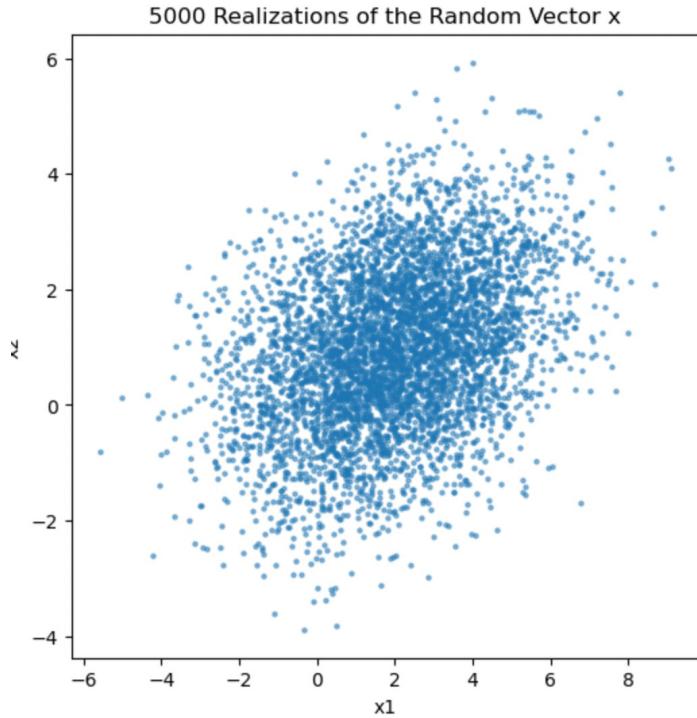
(f)



(g) It's not. Because the two points in the middle are on the two sides of a bump. The left one should increase but it doesn't because it's lower than the initial value. And the right hand side point should decrease. However, it is lower than the last points, the derivative increased instead.

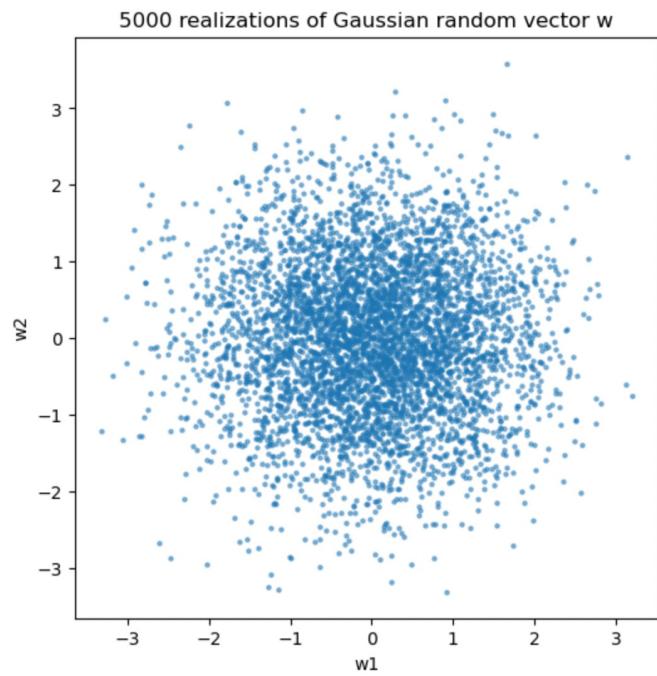
2. Whitening and Simulation of Gaussian Random Vectors

(a)

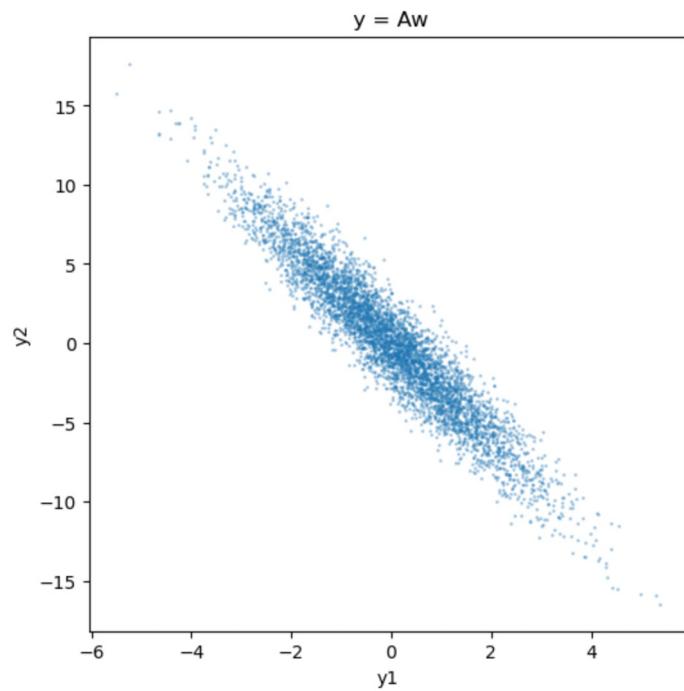


Use `np.random.normal(0,1,size=(5000, 2))` to generate 5000 points, and then calculate eigenvalues and eigenvectors of the covariance matrix. Transform the realizations and add the mean to each realization to get the final points.

(b) (i)



(ii)



$$2. \underline{\underline{y}} = \underline{\underline{A}} \underline{\underline{w}}$$

$$E[\underline{\underline{y}}] = E[\underline{\underline{A}} \underline{\underline{w}}] = \underline{\underline{A}} E[\underline{\underline{w}}]$$

$$\text{As we already know } E[\underline{\underline{w}}] = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\therefore E[\underline{\underline{y}}] = \underline{\underline{A}} E[\underline{\underline{w}}] = \begin{bmatrix} -1 & -1 \\ 2 & 4 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{Var}[\underline{\underline{y}}] = E[(\underline{\underline{y}} - E[\underline{\underline{y}}])(\underline{\underline{y}} - E[\underline{\underline{y}}])^T]$$

$$= E[\underline{\underline{A}} \underline{\underline{w}} \underline{\underline{w}}^T \underline{\underline{A}}^T]$$

$$= \underline{\underline{A}} E[\underline{\underline{w}} \underline{\underline{w}}^T] \underline{\underline{A}}^T$$

$$= \underline{\underline{A}} \underline{\underline{I}} \underline{\underline{A}}^T$$

$$= \begin{bmatrix} -1 & -1 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} -1 & 2 \\ -1 & 4 \end{bmatrix} = \begin{bmatrix} 2 & -6 \\ -6 & 20 \end{bmatrix}$$

$$\therefore \text{Var}[y_1] = 2 \quad \text{Var}[y_2] = 20 \quad 6y_1 y_2 = -6$$

$$\rho_{y_1 y_2} = \frac{6y_1 y_2}{\sqrt{2} \sqrt{20}} = \frac{-6}{\sqrt{2} \sqrt{20}} = -0.94868$$

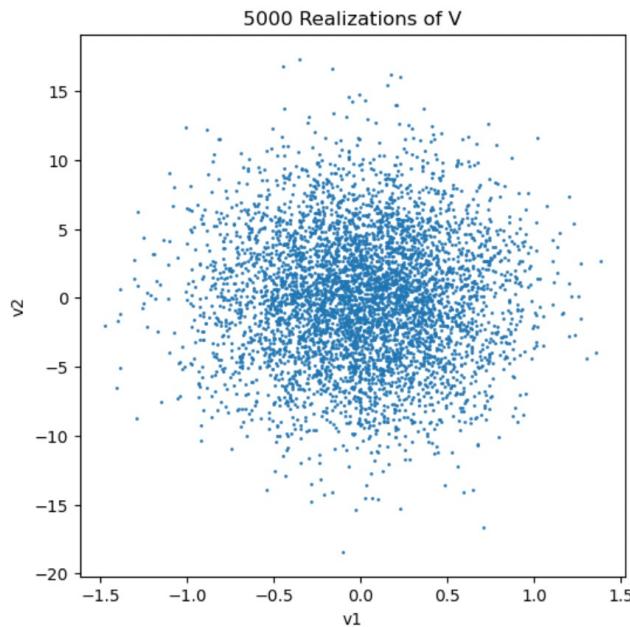
$\therefore y_1$ and y_2 are negatively correlated.

The sample covariance matrix is [[2.00033385 -5.99667272]
[-5.99667272 19.96204433]]

Comparing this to the covariance matrix of y as derived : [[2 -6],[-6 20]], we can find that they have very approximate value.

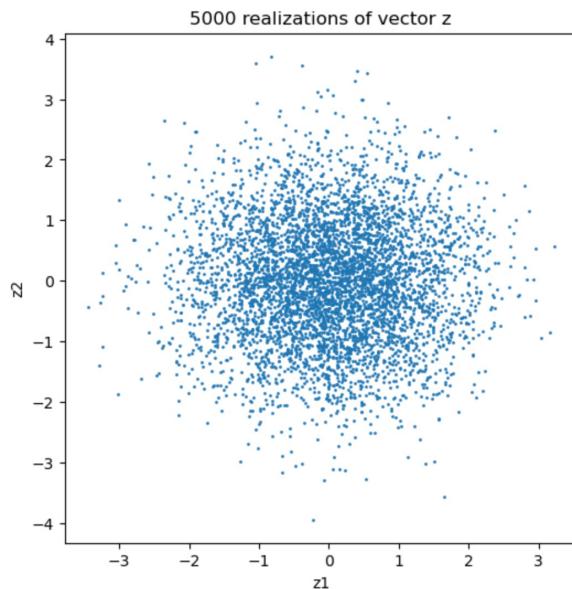
(iii)

The covariance of v is [[1.82310308e-01 -7.85113202e-03]
[-7.85113202e-03 2.17800679e+01]]



(iv)

The covariance of z is [[0.99435022 -0.00392557]
[-0.00392557 0.99832303]]



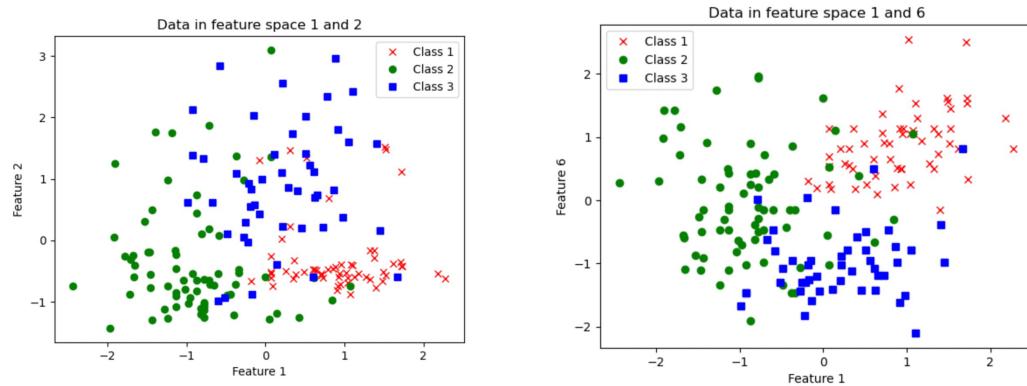
(v) What is the pdf of w and the pdf of z? Is $z = w$ -- i.e., since we colored w to get y, then we whitened y to get z? If not explain why.

The pdf of w is a Gaussian distribution and same as the pdf of z. As w and z have the same pdf, $z = w$ because we colored w to get y, the mean of x equals to zero and the covariance of z is also an identity matrix.

3. Comparison of PCA and MDA on wine dataset.

(a) Baseline for comparison.

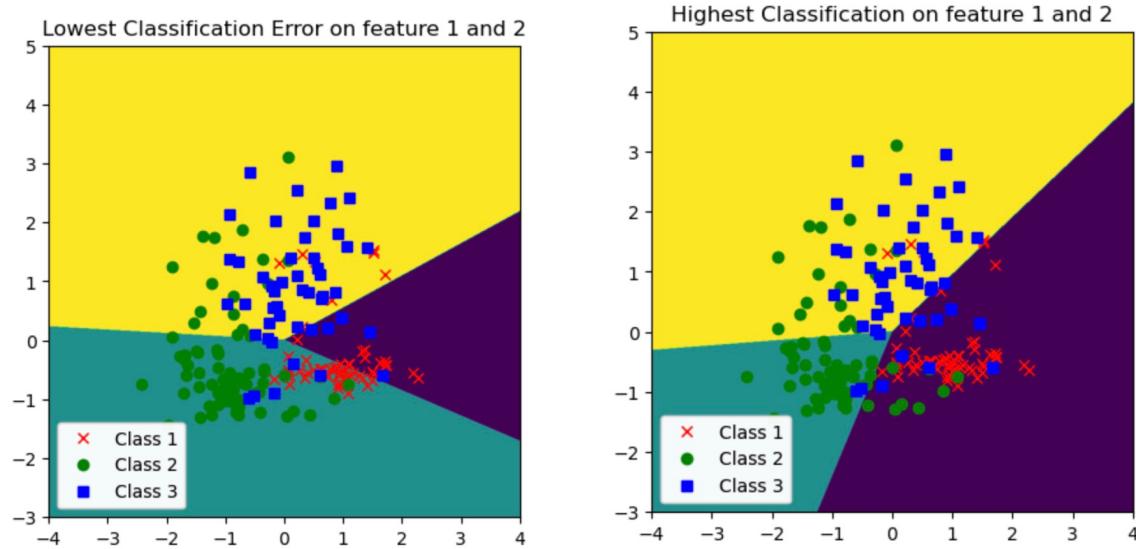
(i)



(ii) Feature 1 and 2:

```

The mean classification error rate in run 1 = 0.28125000000000006
The mean classification error rate in run 2 = 0.2597222222222224
The mean classification error rate in run 3 = 0.2368055555555557
The mean classification error rate in run 4 = 0.225
The mean classification error rate in run 5 = 0.3152777777777778
The average and standard deviation of the mean classification error over the 5 runs is 0.2636111111111113 and 0.03
226469901407941
  
```

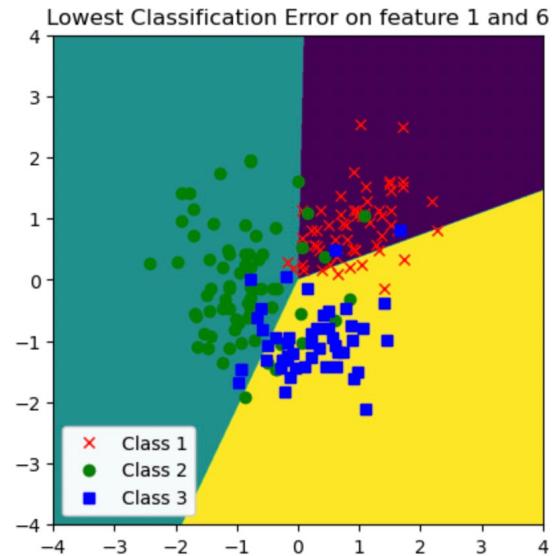
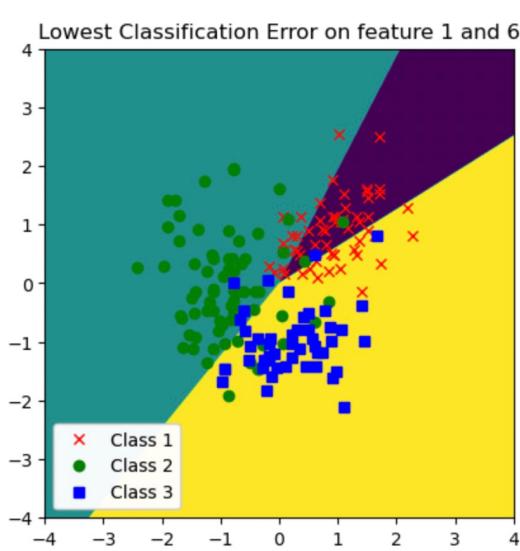


(iii) Feature 1 and 6:

```

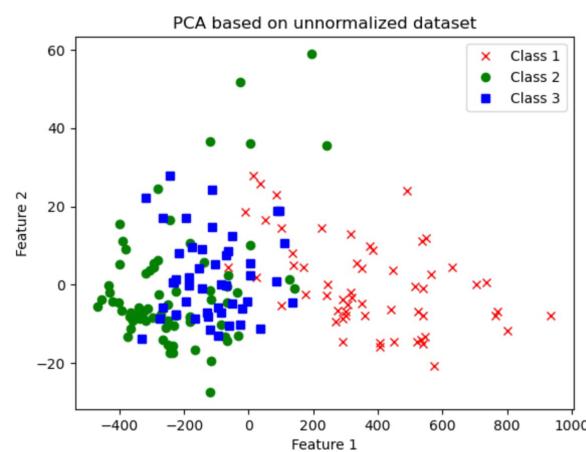
The mean classification error rate in run 1 = 0.24930555555555559
The mean classification error rate in run 2 = 0.23819444444444446
The mean classification error rate in run 3 = 0.2208333333333335
The mean classification error rate in run 4 = 0.2534722222222222
The mean classification error rate in run 5 = 0.22916666666666666
The average and standard deviation of the mean classification error over the 5 runs is 0.2381944444444446 and 0.01
2163685581006285

```



(b) PCA based on unnormalized dataset.

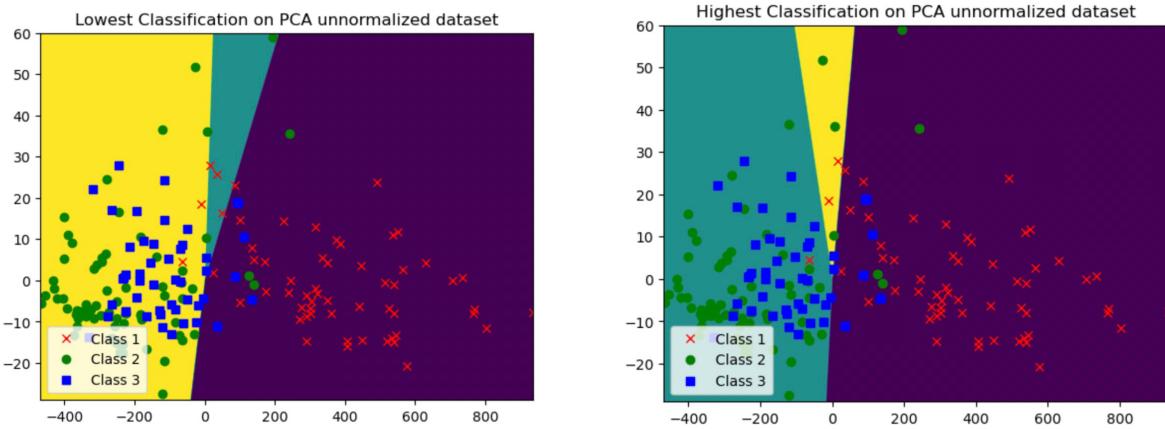
(i)



Class 2 and Class 3 almost overlap, I don't think there would be a better classification result with PCA.

(ii)

The mean classification error rate in run 1 = 0.34375000000000006
 The mean classification error rate in run 2 = 0.36944444444444446
 The mean classification error rate in run 3 = 0.4694444444444444
 The mean classification error rate in run 4 = 0.4583333333333333
 The mean classification error rate in run 5 = 0.4513888888888895
 The average and standard deviation of the mean classification error over the 5 runs is 0.4184722222222225 and 0.051493137502744396

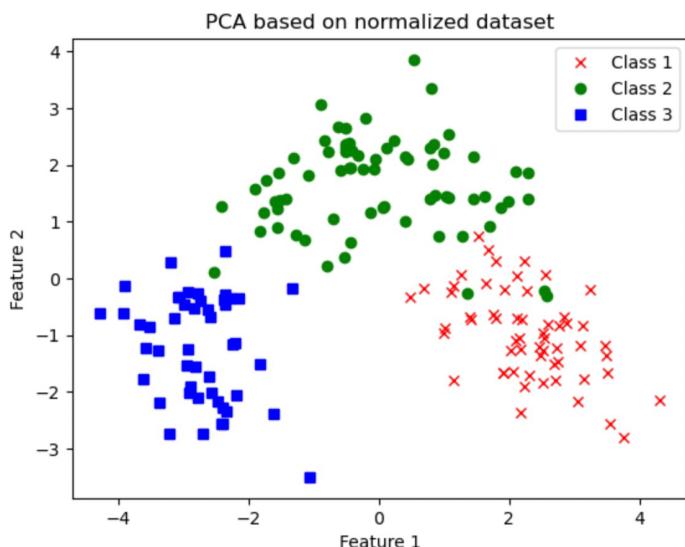


(iii) How does it compare with the baselines in (a)(ii)?

The shape of decision regions seem to like the lowest error result from baselines.

(c) PCA based on standardized dataset.

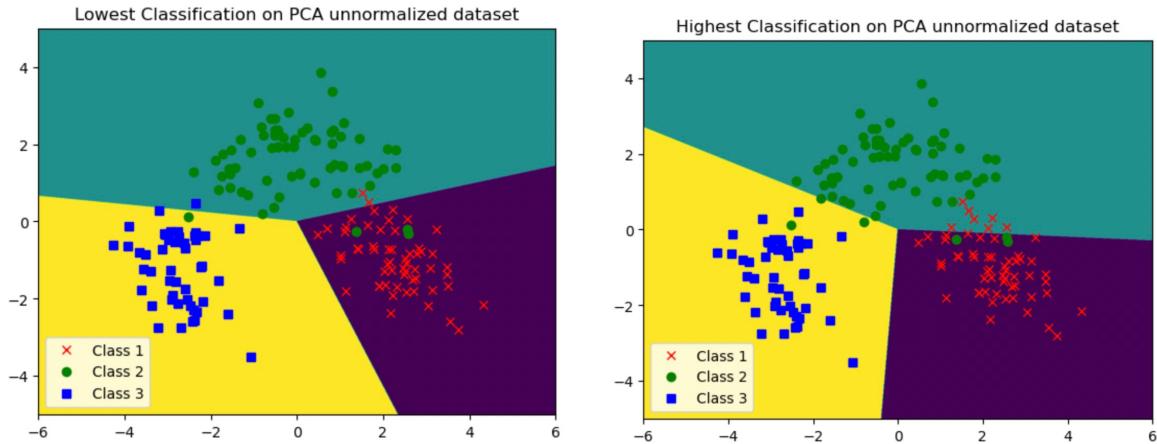
(i)



It will have a better classification result.

(ii)

The mean classification error rate in run 1 = 0.0520833333333335
The mean classification error rate in run 2 = 0.0520833333333336
The mean classification error rate in run 3 = 0.0388888888888889
The mean classification error rate in run 4 = 0.04444444444444445
The mean classification error rate in run 5 = 0.10277777777777779
The average and standard deviation of the mean classification error over the 5 runs is 0.0580555555555557 and 0.022908668637994484

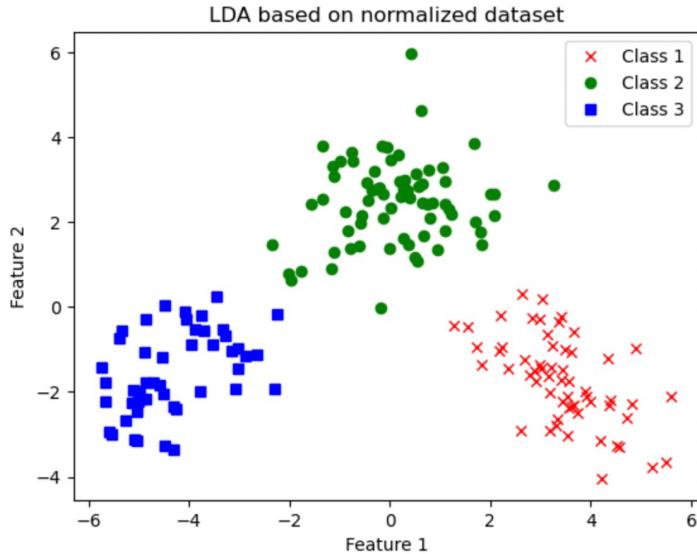


(iii)

The PCA on normalized data make a better performance on classification.

(d) MDA (using LDA as an approximation to MDA).

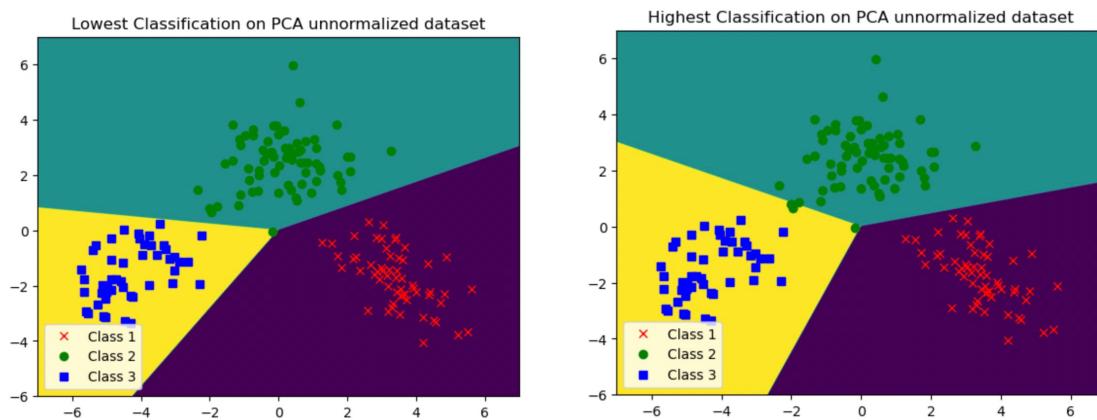
(i)



I expect that the classification result from those data would have low error rate.

(ii)

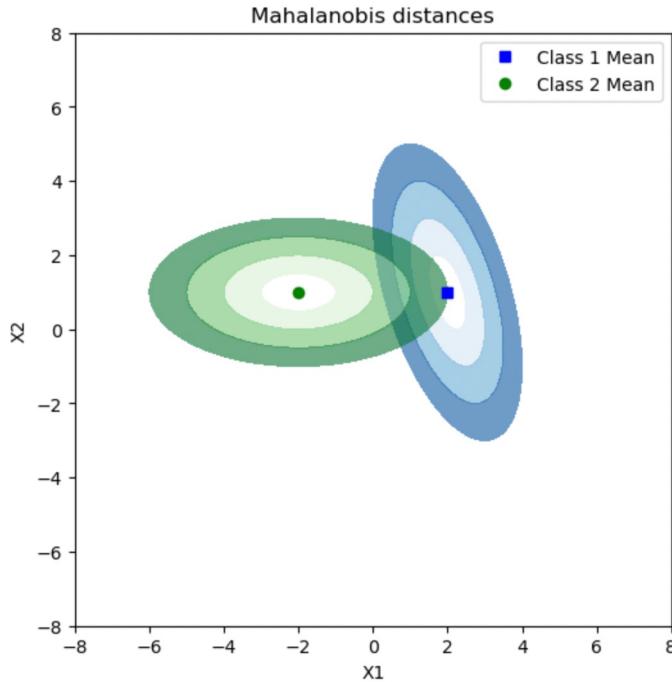
The mean classification error rate in run 1 = 0.0166666666666667
 The mean classification error rate in run 2 = 0.01111111111111117
 The mean classification error rate in run 3 = 0.005555555555555558
 The mean classification error rate in run 4 = 0.016666666666666673
 The mean classification error rate in run 5 = 0.01111111111111117
 The average and standard deviation of the mean classification error over the 5 runs is 0.01222222222222226 and 0.004157397096415491



(iii) According to the mean classification error, the LDA perform best compare to baseline and PCA because the centers of each classes much farer after LDA.

4. Mahalanobis distance and Bayes classification for minimum error

(a)



(b)

$$(b) P(X=x|S_i)P(S_i) \Rightarrow P(\underline{x}|S_i) = N(\underline{x}|\underline{m}_i, \underline{\Sigma}_i), i=1,2$$

$$\ln \frac{P(X=x|S_1)P(S_1)}{P(X=x|S_2)P(S_2)} = \ln \frac{N(\underline{x}, \underline{m}_1, \underline{\Sigma}_1) P(S_1)}{N(\underline{x}, \underline{m}_2, \underline{\Sigma}_2) P(S_2)}$$

$$= \ln \frac{\exp\{-\frac{1}{2}[(\underline{x}-\underline{m}_1)^T \underline{\Sigma}_1^{-1} (\underline{x}-\underline{m}_1)]\} P(S_1)}{(2\pi)^n |\underline{\Sigma}_1|^{\frac{1}{2}}} - \ln \frac{\exp\{-\frac{1}{2}[(\underline{x}-\underline{m}_2)^T \underline{\Sigma}_2^{-1} (\underline{x}-\underline{m}_2)]\} P(S_2)}{(2\pi)^n |\underline{\Sigma}_2|^{\frac{1}{2}}}$$

$$= \frac{1}{2} \ln |\underline{\Sigma}_1| - \frac{1}{2} (\underline{x}-\underline{m}_1)^T \underline{\Sigma}_1^{-1} (\underline{x}-\underline{m}_1) + \ln P(S_1) - \frac{1}{2} \ln |\underline{\Sigma}_2| + \frac{1}{2} (\underline{x}-\underline{m}_2)^T \underline{\Sigma}_2^{-1} (\underline{x}-\underline{m}_2) - \ln P(S_2)$$

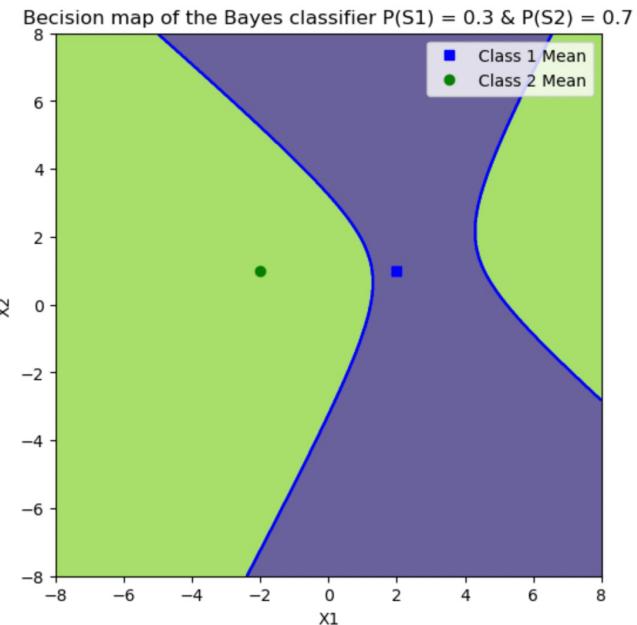
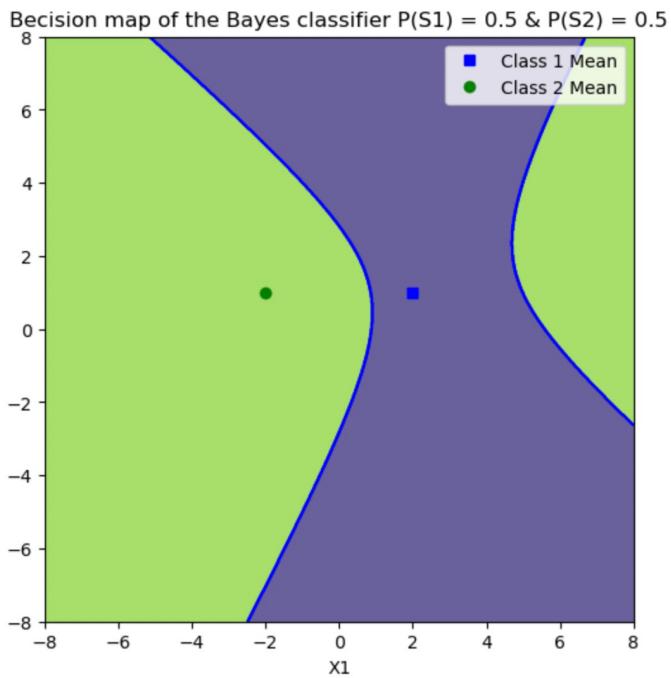
$$= \frac{1}{2} \ln |\underline{\Sigma}_1| - \frac{1}{2} d_m^2(\underline{x}, \underline{m}_1) + \ln P(S_1) - \frac{1}{2} \ln |\underline{\Sigma}_2| + \frac{1}{2} d_m^2(\underline{x}, \underline{m}_2) - \ln P(S_2)$$

$$\ln P(S_1) - \ln P(S_2) + \frac{1}{2} \ln |\underline{\Sigma}_1| - \frac{1}{2} \ln |\underline{\Sigma}_2| = \frac{1}{2} d_m^2(\underline{x}, \underline{m}_1) - \frac{1}{2} d_m^2(\underline{x}, \underline{m}_2)$$

\therefore The decision boundary is:

$$d_m^2(\underline{x}, \underline{m}_1) - d_m^2(\underline{x}, \underline{m}_2) = 2 \ln \frac{P(S_1)}{P(S_2)} + \ln \frac{|\underline{\Sigma}_1|}{|\underline{\Sigma}_2|}$$

(c)



Decision map of the Bayes classifier $P(S1) = 0.1$ & $P(S2) = 0.9$

