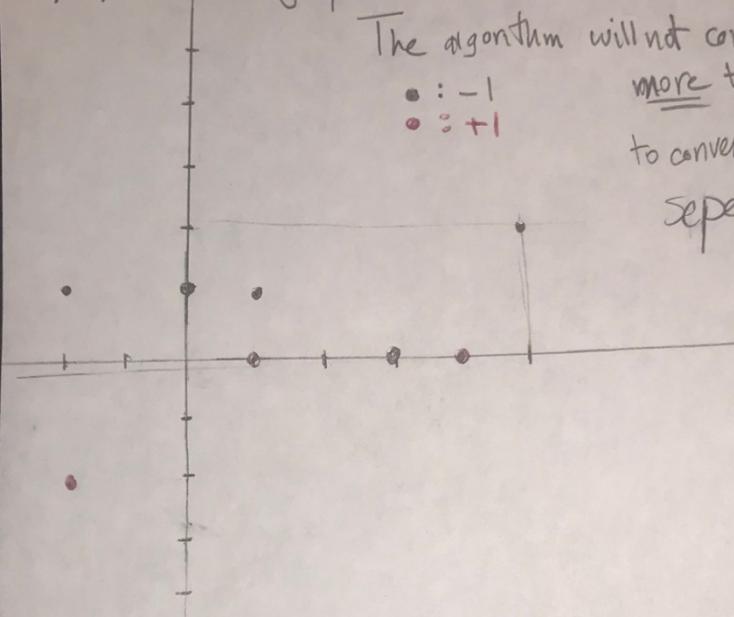


Milad Nourian
EECE M146, HW2

1) a) From the graph it is clear that data is no linearly separable.



The algorithm will not converge if we run the algorithm more times. For the perceptron algorithm to converge, the data needs to be linearly separable.

b) Algorithm for perceptron training:

$$w \leftarrow 0, d \in \{0, 1\}$$

for iter in MaxIter:

for (x, y) in data:

$$a \leftarrow w^T x$$

If $(ay_n < 0)$:

$$w^{t+1} \leftarrow w^t + y_n x_n$$

else:

do nothing

$$(1) \quad w = [4] \text{ # initialize}$$

$$1) \quad a = 1 \Rightarrow ay_1 > 0 \Rightarrow \text{do nothing}$$

$$(2) \quad 2) \quad a = 4 \Rightarrow ay_2 < 0 \Rightarrow w = [4] + [-1][1] = [-1]$$

$$3) \quad a = -1 \Rightarrow ay_3 > 0 \Rightarrow \text{# do nothing}$$

$$4) \quad a = -4 \Rightarrow ay_4 < 0 \Rightarrow w = [-4] + (1)(-2) = [1]$$

$$5) \quad a = -5 \Rightarrow ay_5 > 0 \Rightarrow \text{# nothing}$$

$$6) \quad a = 1 \Rightarrow ay_6 > 0 \Rightarrow \text{# nothing}$$

$$7) \quad a = -1 \Rightarrow ay_7 > 0 \Rightarrow \text{# nothing}$$

$$(8) \quad 8) \quad a = 3 \Rightarrow ay_8 < 0 \Rightarrow w = [3] + (-1)[0] = [-3]$$

DC) perceptron closed form: $\underline{w} = \begin{bmatrix} -2 \\ -3 \end{bmatrix} \rightarrow y = \text{sign}(\underline{w}^T \underline{x})$

The algorithm for weighted perceptron:

Initialize $\underline{w} \leftarrow \underline{0}$, $u \leftarrow \underline{0}$ \hookrightarrow includes bias

For $t = 1$ to M do:

For $(x, y) \in D$:

If $y_n(\underline{w}^T \underline{x}) < 0$: # update

$$\underline{w} \leftarrow \underline{w} + y_n \underline{x}$$

$$u \leftarrow u + y_n$$

and P
what?

We know: If $w_1, \dots, w_k \Rightarrow w_1, \dots, w_k$

$y_1, \dots, y_n \quad c_1, \dots, c_k$

$$w_{\text{avg}} = \sum c_k w_k$$

$$w_{\text{avg}} = 1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 3 \\ -1 \end{bmatrix} + 4 \begin{bmatrix} 1 \\ -3 \end{bmatrix} + \begin{bmatrix} -2 \\ -3 \end{bmatrix}$$

using part a results

$$w_{\text{avg}} = \begin{bmatrix} 12 \\ -17 \end{bmatrix} \rightarrow \hat{y} = \text{sign}(w_{\text{avg}}^T \underline{x})$$

majority vote perceptron: take majority from w_1, w_2, \dots, w_k

$$\hat{y} = \text{sign}\left(\sum_{k=1}^K c_k \text{sign}(w_k^T \underline{x})\right)$$

$$\hat{y} = \text{Sign}\left(1 \text{sign}(\underline{1}^T \underline{x}) + 2 \text{sign}(\underline{3}^T \underline{x}) + 4 \text{sign}(\underline{1}^T \underline{x}) + \text{sign}(\underline{-2}^T \underline{x})\right)$$

1) d) Find training error?

$$W = \begin{bmatrix} -2 \\ -3 \end{bmatrix}$$

$$W_{avg} = \begin{bmatrix} 12 \\ -17 \end{bmatrix}$$

$$W_{majority} = \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

$$a + b = 10$$

	Perception	Average perception	Majority perception
1)	I	C	C
2)	C	C	C
3)	C	C	C
4)	C	C	C
5)	I	C	C
6)	I	C	C
7)	C	I	C
8)	C	I	I

$$\text{Perception Accuracy} = \frac{5}{8} = 0.625$$

$$\text{Average perception Accuracy} = \frac{0.75}{8} = 0.375$$

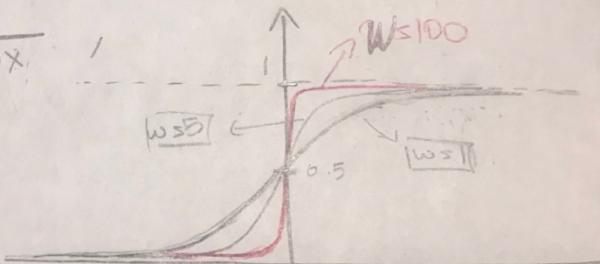
$$\text{Majority Accuracy} = \frac{0.875}{8} = 0.125$$

$$\text{Perception error} = 0.375$$

$$\text{Average perception error} = 0.25$$

$$\text{Majority perception error} = 0.125$$

$$2) a) J(wx) = \frac{1}{1+e^{-wx}}$$



For large values of w , the curve is steep. In this situation, a small change in the input results in a large change in the output which is an indication of overfit. To avoid overfitting, we would want the curves to be smooth as possible, meaning a small change in the input would result in a small change in the output.

$$\begin{aligned}
 2) b) J(\underline{\omega}) &= -\sum_{n=0}^N y_n \log h_{\underline{\omega}}(\underline{x}_n) + (1-y_n) \log(1-h_{\underline{\omega}}(\underline{x}_n)) + \frac{1}{2} \sum_i \omega_i^2 \\
 &= -\sum_{n=0}^N y_n \underbrace{\log(\underline{\omega}^T \underline{x}_n)}_{-\log(1+e^{-\underline{\omega}^T \underline{x}_n})} + (1-y_n) \log(1-\underline{\omega}^T \underline{x}_n) + \frac{1}{2} \sum_i \omega_i^2 \\
 \frac{\partial J}{\partial \omega_j} &= -\sum_{n=0}^N y_n \frac{\partial}{\partial \omega_j} \log(\underline{\omega}^T \underline{x}_n) + (1-y_n) \frac{\partial}{\partial \omega_j} \log(1-\underline{\omega}^T \underline{x}_n) + \frac{1}{2} \sum_i \frac{\partial}{\partial \omega_i} \omega_i^2 \\
 &= -\sum_{n=0}^N y_n \left[-\frac{1}{1+e^{-\underline{\omega}^T \underline{x}_n}} \right] [\underline{x}_n^j] + (1-y_n) \left[\frac{1}{1-(1-\underline{\omega}^T \underline{x}_n)} \right] [\underline{x}_n^j] + \underline{\omega}_j^2 \\
 &= (-\underline{\omega}_j) \sum_{n=0}^N \left[\frac{-y_n}{1+e^{-\underline{\omega}^T \underline{x}_n}} + \frac{(1-y_n)}{1-(1-\underline{\omega}^T \underline{x}_n)} \right] + \underline{\omega}_j^2
 \end{aligned}$$

$$\boxed{\nabla G(\underline{\omega}) = \sum_n [\underline{\omega}^T \underline{x}_n - y_n] \underline{x}_n + \underline{\omega}^2} \quad \checkmark$$

This is the gradient descent step for linear regression now

$$\underline{\omega}^{t+1} \leftarrow \underline{\omega}^t - \eta \left[\sum_n (\underline{\omega}^T \underline{x}_n - y_n) \underline{x}_n + \underline{\omega}^2 \right] \Rightarrow \text{So now we add } \underline{\omega}$$

2(c) We know $L(D) = \mathbb{P}(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n)$
we know every (x_i, y_i) pair is IID, so

$$\log L(\omega) = \log \left[\prod_{i=1}^n P(y_i | x_i, \theta) f(w_0 + w_1 x_i) \right]$$

$$= \sum_{i=1}^n \log P(y_i | x_i, \theta) + \log f(w_0 + w_1 x_i)$$

$$= \sum_{i=1}^n \log \left(\frac{1}{1 + e^{-(w_0 + w_1 x_i)}} \right) \left[\frac{y_i}{1 - h_\theta(x_i)} \right] + \log \frac{1}{\sqrt{2\pi}} e^{-\frac{(w_0 + w_1 x_i)^2}{2}}$$

$$= \sum_{i=1}^n y_i \log h_\theta(x_i) + (1-y_i) \log [1-h_\theta(x_i)] - \frac{1}{2} \log 2\pi - \frac{\sum w_i^2}{2}$$

$$\log L(\omega) = \sum_{i=1}^n y_i \log h_\theta(x_i) + (1-y_i) \log [1-h_\theta(x_i)] - \log 2\pi - \frac{\sum w_i^2}{2}$$

Now we found this in class:
Comparing $\log L(\omega)$ with $J(\omega)$ given in the problem, we see that

$$\log L(\omega) = -J(\omega)$$

So the update rule for gradient descent is

And in part a, we found $\frac{\partial J(\omega)}{\partial \omega}$. So, we can conclude that

Minimizing the cost function is same as maximizing likelihood $L(\omega)$.

We can see that minimizing the MLEAP is equivalent of minimizing logistic regression objective function ($J(\omega) = -\log L(\omega)$)

objective
function

likelihood
function

3) a)

$$J(w_0, w_1) = \sum_n \alpha_n [w_0 + w_1 x_{n,1} - y_n]^2$$

$$\frac{\partial J}{\partial w_0} = 2 \sum_n [\alpha_n] [2] [w_0 + w_1 x_{n,1} - y_n] [1] = 2 \sum_n [\alpha_n] [w_0 + w_1 x_{n,1} - y_n]$$

$$\frac{\partial J}{\partial w_1} = 2 \sum_n [\alpha_n] [w_0 + w_1 x_{n,1} - y_n] [x_{n,1}]$$

so now find gradient $\nabla J = \begin{bmatrix} \frac{\partial J}{\partial w_0} \\ \frac{\partial J}{\partial w_1} \end{bmatrix} = 2 \sum_n [\alpha_n] [w^T x - y_n] x_n$

* So the only change is that now we have a coefficient α_n (as weight).

b) We can first calculate Hessian Matrix:

$$H = \begin{bmatrix} \frac{\partial^2 J}{\partial w_0^2} & \frac{\partial^2 J}{\partial w_0 \partial w_1} \\ \frac{\partial^2 J}{\partial w_1 \partial w_0} & \frac{\partial^2 J}{\partial w_1^2} \end{bmatrix} = \begin{bmatrix} 2 \sum_n \alpha_n & 2 \sum_n \alpha_n x_{n,1} \\ 2 \sum_n [\alpha_n] (x_{n,1}) & 2 \sum_n [\alpha_n] (x_{n,1})^2 \end{bmatrix}_{2 \times 2}$$

Now, we know If $z^T H z > 0 \Rightarrow$ then Convex $z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_{2 \times 1}$

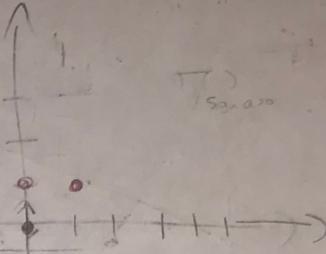
$$2 \begin{bmatrix} \sum_n \alpha_n z_1 + \sum_n \alpha_n x_{n,1} z_2 & \sum_n \alpha_n x_{n,1} z_1 + \sum_n \alpha_n (x_{n,1})^2 z_2 \end{bmatrix}_{1 \times 2} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$= \left\{ \sum_n \alpha_n z_1^2 + \sum_n \alpha_n x_{n,1} z_1 z_2 + \sum_n \alpha_n x_{n,1} z_2 z_1 + \sum_n \alpha_n (x_{n,1})^2 z_2^2 \right\}$$

$$\Rightarrow \left[\sum_n \alpha_n \left(z_1^2 + 2x_{n,1} z_1 z_2 + x_{n,1}^2 z_2^2 \right) \right] > 0 \Rightarrow \text{so } J(\omega) \text{ is convex and}$$

therefore, there exists a global minimum ($\underset{\min}{\text{global}} = \text{local min}$)

$$\begin{array}{ll} \text{4) a)} & \begin{array}{l|l} x_1 = -1 & \rightarrow y = 1 \\ x_2 = 0 & \rightarrow y = 1 \\ x_3 = 1 & \rightarrow y = 1 \\ x_4 = -3 & \rightarrow y = 1 \\ x_5 = -2 & \rightarrow y = 1 \\ x_6 = 3 & \rightarrow y = -1 \end{array} \end{array}$$



One optimal solution

$$\underline{w}_{\text{opt}} = \begin{bmatrix} 1 \\ 1.5 \end{bmatrix}, \text{ Accuracy} = \frac{5}{6}$$

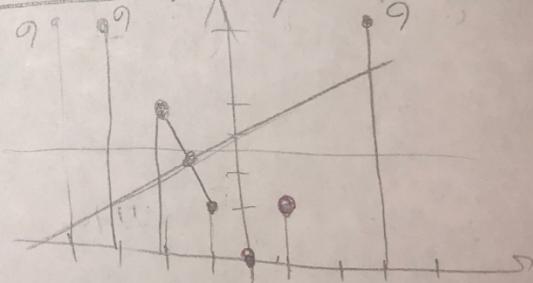
$$\text{b) } K(x_1 z) = x_2(1+xz)$$

$$= x_2 f(x^2)$$

\Rightarrow Therefore feature map $\phi(x) = (x_1, x^2)$ transforms from $R^2 \rightarrow R^3$

Now calculating $\phi(x)$ for all points

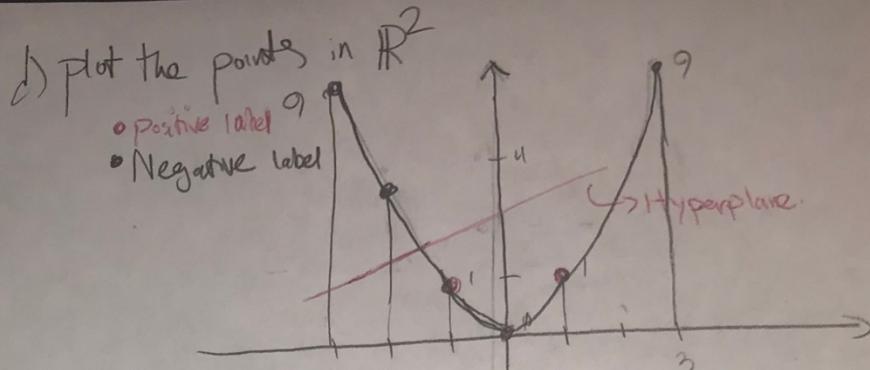
$$\begin{bmatrix} -1 \\ 1 \\ 0 \\ 0 \\ 1 \\ -3 \\ 9 \\ -2 \\ 4 \\ 3 \\ 9 \end{bmatrix}$$



Answer: We can see that the line must pass between $v_1 = [-2, 4]$ and $v_2 = [1, 1]$ to separate the data.

The separating line can also go through middle of v_1, v_2 . So we can get $-x_1 + 3x_2 - 9 = 0$

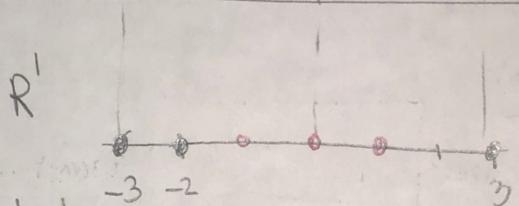
$$\underline{w} = \begin{bmatrix} -1 \\ 3 \\ -9 \end{bmatrix} \Rightarrow \text{is one optimal solution.}$$



This graph agrees with our intuition, as $\phi(x) = (x, x^2)$ was our feature map and we expected a parabola (polynomial).

Optimal \underline{w} (w_{opt}) has accuracy 1.0 & correctly classifies all the points. $\Rightarrow \boxed{\text{Accuracy} = 1.0}$

e)



We can see that the data is not linearly separable in $\underline{\mathbb{R}^1}$.

ECE M146, HW 2

5)a)

- We can see from the data the scatter plot, it seems to follow a cubic behavior. From this behavior we can expect that using linear regression to model this data would most likely **underfit**, as we need to increase the order of polynomial (M) in our model to better fit the given data that we see from the scatter plot.

5)d)

For this table I iterated through all the input points and updated the theta vector using the following formula (used 2 in the gradient descent formula):

$$w_j \leftarrow w_j - 2\alpha \sum_{n=1}^N (h_w(\mathbf{x}_n) - y_n) x_{n,j} \quad (\text{simultaneously update } w_j \text{ for all } j)$$

:	Coefficients	:	Number of iterations	:	Objective Function
[2.27040138 -2.46075209]	:	10000	:	4.08633386426	
[2.44566862 -2.81750141]	:	10000	:	3.91261070545	
[2.43797706 -2.82790485]	:	1233	:	3.91658245895	
[2.39884098 -2.86600362]	:	332	:	4.0160981729	

For this table I used the matrix formula given to us in the class as shown below (did not use the 2 coefficient):

Algorithm 2 Gradient Descent (J)

```

1:  $t \leftarrow 0$ 
2: Initialize  $\theta^{(0)}$ 
3: repeat
4:    $\nabla J(\theta^{(t)}) = \mathbf{X}^T \mathbf{X} \theta^{(t)} - \mathbf{X}^T \mathbf{y} = \sum_n (\mathbf{x}_n^T \theta^{(t)} - y_n) \mathbf{x}_n$ 
5:    $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla J(\theta^{(t)})$ 

```

For each of the step values $10^{-4}, 10^{-3}, 10^{-2}, 0.0407$

:	Coefficients	:	Number of iterations	:	Objective Function
[1.91573585 -1.74358989]	:	10000	:	5.49356558874	
[2.4463815 -2.81630184]	:	10000	:	3.91257640947	
[2.44640698 -2.81635335]	:	1480	:	3.91257640579	
[2.44640705 -2.81635349]	:	377	:	3.91257640579	

- The coefficients for the case with small step size (10^{-4}) are smaller than cases with larger step size. Basically, we see a trend that increasing the step size will make the algorithm for linear regression to converge faster. Also, increasing the step size will give larger weight values (coefficients), so smaller step sizes have smaller coefficients and larger step sizes have larger coefficient values. For step size of 10^{-4} and 10^{-3} , we never converge suggesting that the movement has been too slow, but for the step size of 10^{-2} , the algorithm for the gradient descent converges after 1480 iterations and for step size of 0.0407, we converge after 377 iterations. In addition, it can be seen that the

cost function (J , or objective function) for small step size is large, and it gets smaller as the step size get increases.

5)e)

- The coefficients for theta using the **close form solution** is found to be:

: Coefficients	:	Objective Function
: [2.44640709 -2.81635359]	:	3.91257640579

Comparing with the solution we obtained from the Gradient descent algorithm, we can see that the coefficients and the cost function from Gradient descent (for step size 0.0407 and 10^{-2}) are identical to the values for coefficients and cost found in the closed form. Therefore, we can see that when the gradient descent converges, it reaches the same answer as the closed form solution for the optimal theta, and our expectation holds true.

The closed-form answer found takes longer than the gradient descent approach. The reason for this is that finding the inverse of matrix $(X^T X)^{-1}$ takes $O(D^3)$, which computationally heavy due to its high order complexity whereas updating theta in the gradient descent takes $O(ND)$ on each iteration and therefore, has less complexity.

5)f)

by changing the value of the step size on each iteration, the following result was obtained:

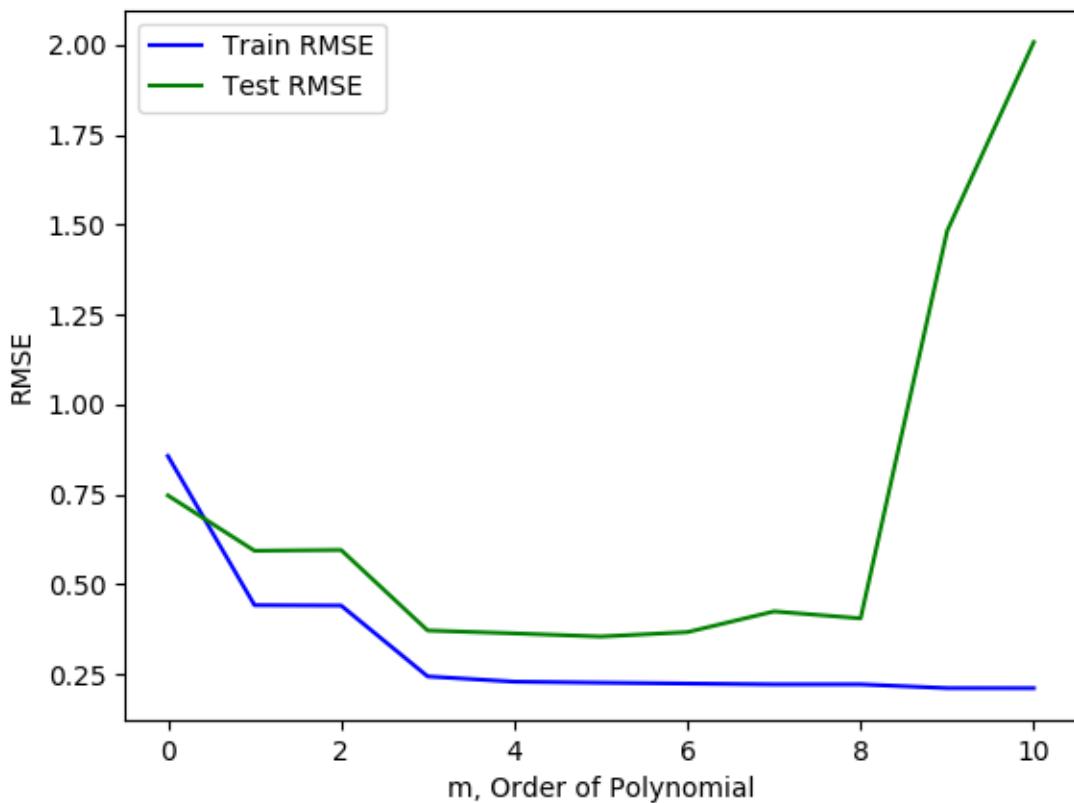
: Coefficients	:	Number of iterations	:	Objective Function
[2.44635846 -2.81652261]	:	10000	:	3.91257677039

It can be seen from the above result that by changing the value of step size, the algorithm does converge but it takes all 10000 iterations. This can be explained because the algorithm uses step = $1 / (1+k)$, so it starts with big step sizes and toward the end the step size becomes very small, taking very long for the algorithm to converge.

5)h)

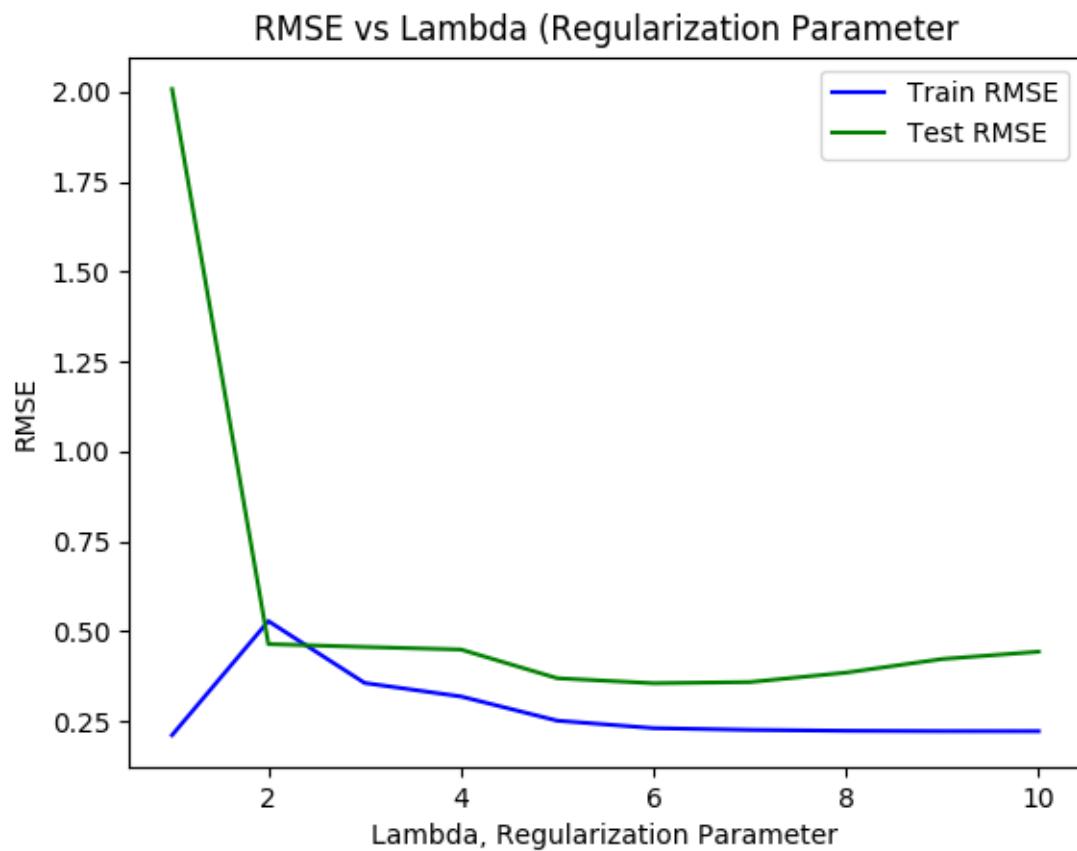
- One reason why MRS is used instead of regular $J(w)$ (cost function) is that MRS provides a normalized value, where the number of samples N is already being taken into the account. This can eliminate the performance measurement dependency on the input data. Using MRS instead of $J(w)$ is a better representative of performance measurement.

5)i)



- From the diagram, we can see that $m = 3$ best fits the data due to the following reasons:
 - For $m = 3$, the test error is at minimum
 - When visualizing data in 5)a) we also anticipated the polynomial fit to be a cubic
 - $m = 3$ is relatively a low degree polynomial, and to avoid overfitting, we tend to use the “simplest models” that fit the data pretty well, so in this case $m = 3$
- There can be seen a significant amount of overfitting for large values of m . From the MSRE graph we see that by increasing m for $m >= 3$, the training error reduces meaning that our model fits the training data very well, however, as m increases further, the test error increases very rapidly, showing a case of overfitting.

5)k)



- From the graph, we can see that the best value for λ is 5 because this value minimizes the RMSE.