

# Project 1 : Curve Fitting

Pramod Kumar

February 5, 2020

## Abstract

Curve fitting is process of finding a approx function which can mimic actual function behind the dataset. In other words, its about finding function behind the given dataset. if we know about actual/approx function, we can calculate/predict value of a unknown input. To achieve this we will implement 2 methods i.e Linear regression and Bayesian approach and we will see what are the short coming in each approach and how to rectify them with regularization.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Approach</b>	<b>2</b>
2.1	Data . . . . .	2
2.2	Methods . . . . .	2
2.2.1	Linear regression using error minimization . . . . .	2
2.2.2	Linear regression using error minimization with regularization . . . . .	3
2.2.3	Maximum Likelihood Estimation . . . . .	4
2.2.4	Maximum Posterior Estimation . . . . .	5
<b>3</b>	<b>Results</b>	<b>7</b>
3.1	Linear regression using error minimization . . . . .	7
3.2	Linear regression using error minimization with regularization . . . . .	8
3.3	Maximum Likelihood Estimation . . . . .	10
3.4	Maximum Posterior Estimation . . . . .	12
<b>4</b>	<b>Conclusion</b>	<b>13</b>
<b>5</b>	<b>References</b>	<b>13</b>

# 1 Introduction

Basic idea revolve around finding the parameter values in linear equation (w.r.t  $W$ ) by using randomly generated data-set, which is called training phase. Later use those calculated  $W^*$  values to predict new input. In practice, instead of regression mostly such problem are solved using gradient descent approach as dimension of data in real world is very high. But for the scope of this project, we ill look at 4 different approach, namely:

- 1) Linear regression using error minimisation
- 2) Linear regression using error minimisation with regularisation
- 3) Maximum Likelihood Estimator
- 4) Maximum posterior estimation

## 2 Approach

### 2.1 Data

In this project we will be using 2 set of data with size 50 and 100. This data is generated using *generateData.m* script in the project directory. This script generate data by calculating  $\sin x$  of random data and add noise to resemble it with real data.

### 2.2 Methods

#### 2.2.1 Linear regression using error minimization

Error Minimization (non-regularized)

A linear equation (wrt.  $W$ ) is defined as

$$y(x, w) = \sum_{i=1}^M w_i x^i \quad (1)$$

$$W = [w_0 \ w_1 \ \dots \ w_M] \quad \text{and} \quad X = \begin{bmatrix} x_1^0 & x_1 & \dots \\ \vdots & \ddots & \\ x_N^0 & & x_N^M \end{bmatrix}_{N \times M}$$

$$\text{and} \quad T = [t_0 \ t_1 \ \dots \ t_N]$$

Error associated with this equation (least square error) can be expressed as follows:

$$E(\vec{w}) = \left(\frac{1}{2}\right) \sum (y(\vec{x}_n) - \vec{t}_n)^2.$$

A best solution will have LSE as small as possible. In geometric interpretation, we need to determine value of  $W$  such that sum of projection of all points on curved formed by  $W$  is minimized.

To minimize error function, we can take differentiate and equate to 0.  
Differentiating with respect to  $w$ , we get

*Matrix representation of  $y$  :  $Y = XW$*

*Error function in matrix form :*

$$E(W) = \frac{1}{2}(XW - T)^2$$

$$E(W) = \frac{1}{2}(XW - T)^T(XW - T)$$

$$E(W) = \frac{1}{2}(W^T X^T - T^T)(XW - T)$$

$$E(w) = \frac{1}{2}(W^T X^T XW - W^T X^T T - T^T XW + T^T T)$$

$$(T^T XW)^T = (W^T X^T T) \quad \text{and both are } 1 \times 1 \text{ matrix}$$

$$E(w) = \frac{1}{2}(W^T X^T XW - 2W^T X^T T + T^T T)$$

$$\begin{aligned} \frac{\partial}{\partial W} E(W) &= \frac{\partial}{\partial W} \frac{1}{2}(W^T X^T XW - 2W^T X^T T + T^T T) \\ &= \frac{1}{2}(2X^T XW - 2X^T T + 0) \end{aligned}$$

$$= X^T XW - X^T T \tag{2}$$

To minimize equation 2, we need to equate to 0

$$X^T XW - X^T T = 0$$

$$X^T XW = X^T T$$

$$W^* = (X^T X)^{-1} X^T T$$

### 2.2.2 Linear regression using error minimization with regularization

Here we will use same sum square error function but we will have another parameter  $\lambda$ , which will help in solving problem of over fitting. A linear equation (wrt.  $W$ ) is defined as

$$E(\vec{w}) = \left(\frac{1}{2}\right) \sum (y(x_n) - \vec{t}_n)^2 + \frac{\lambda}{2} \vec{W}^2.$$

Similar to minimizing non-regularized error function, we will minimize this function too.

*Matrix representation of  $y$  :*  $Y = XW$

*Error function in matrix form :*

$$\begin{aligned}
 E(W) &= \frac{1}{2}(X\vec{W} - \vec{T})^2 + \frac{\lambda}{2}\vec{W}^2 \\
 E(W) &= \frac{1}{2}(X\vec{W} - \vec{T})^T(X\vec{W} - \vec{T}) + \frac{\lambda}{2}W^T W \\
 E(W) &= \frac{1}{2}(\vec{W}^T \vec{X}^T - \vec{T}^T)(X\vec{W} - \vec{T}) + \frac{\lambda}{2}W^T W \\
 E(w) &= \frac{1}{2}(\vec{W}^T \vec{X}^T \vec{X} \vec{W} - \vec{W}^T \vec{X}^T \vec{T} - \vec{T}^T X \vec{W} + \vec{T}^T \vec{T}) + \frac{\lambda}{2}W^T W \\
 (T^T X W)^T &= (W^T X^T T) \quad \text{and both are } 1 \times 1 \text{ matrix} \\
 E(w) &= \frac{1}{2}(W^T X^T X W - 2W^T X^T T + T^T T) + \frac{\lambda}{2}W^T W \\
 \frac{\partial}{\partial W} E(W) &= \frac{\partial}{\partial W} \frac{1}{2}(W^T X^T X W - 2W^T X^T T + T^T T) + \frac{\partial}{\partial W} \frac{\lambda}{2}W^T W \\
 &= \frac{1}{2}(2X^T X W + \lambda 2IW - 2X^T T + 0) \\
 &= X^T X W + \lambda IW - X^T T
 \end{aligned} \tag{3}$$

To minimize equation 2, we need to equate to 0

$$\begin{aligned}
 X^T X W + \lambda IW - X^T T &= 0 \\
 X^T X + \lambda IW &= X^T T \\
 \vec{W}^* &= (\vec{X}^T \vec{X} + \lambda \vec{I})^{-1} \vec{X}^T \vec{T}
 \end{aligned}$$

So minimization function becomes :

$$E(\vec{W}^*) = \frac{1}{2}(X\vec{W}^* - \vec{T})^2 + \frac{\lambda}{2}\vec{W}^{*2} \tag{4}$$

### 2.2.3 Maximum Likelihood Estimation

It is defined as

$$p(W | X, t, \beta) \propto p(t | X, W, \beta)$$

$p(t | X, W, \beta)$  is called likelihood and left side of above equation is call postprior.  
Likelihood defined on N data points as:

$$p(t | X, W, \beta) = \prod_{n=1}^N N(t_n | y(x_n, W), \beta^{-1}) \tag{5}$$

Where

$$N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} e^{\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}}$$

where  $\beta^{-1} = \sigma^2$ . Taking log on equation 5, we get

$$\ln(p(t|x, w, \beta)) = \left(\frac{-\beta}{2}\right) \sum (y(\vec{x}_n) - \vec{t}_n)^2 + \left(\frac{N}{2}\right) \ln \beta - \left(\frac{N}{2}\right) \ln 2\pi. \quad (6)$$

To estimate W, Differentiate equation 6 w.r.t W

$$\frac{1}{p(t|x, w, \beta)} = \left(\frac{-\beta}{2}\right) \frac{\partial (E(W))}{\partial W} \quad (7)$$

Calculating maximization of likelihood, is equal to calculating minimization of E(W), which is equal to minimization of Task 1.

$$W^* = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{T}$$

Now optimize for  $\beta$ . Differentiate equation 6 W.r.t  $\beta$  and equate it to 0 for maximizing.

$$\begin{aligned} \frac{1}{p(t|x, w, \beta)} &= \frac{1}{2} \sum_{n=1}^N (y(x_n, W) - t_n)^2 + \frac{N}{2\beta} \\ \frac{N}{2\beta} &= \frac{1}{2} \sum_{n=1}^N (y(x_n, W) - t_n)^2 \\ \beta^{-1} &= \frac{1}{N} \sum_{n=1}^N (y(x_n, W) - t_n)^2 \end{aligned} \quad (8)$$

#### 2.2.4 Maximum Posterior Estimation

$$p(W | X, t, \alpha, \beta) \propto p(t | X, W, \beta) p(W | \alpha) \quad (9)$$

[1] In posterior, we need to find optimum (best fit) **W** for the given data. where

$$p(W | \alpha) = \frac{\alpha^{M+1}}{2\pi} e^{\left(\frac{-\alpha}{2} w^T W\right)} \quad (10)$$

[1] and

$$p(t | X, W, \beta) = \prod_{n=1}^N N(t_n | y(x_n, W), \beta^{-1})$$

[1] Take log :

$$\ln p(W | X, t, \alpha, \beta) = \ln(t | X, W, \beta) + \ln p(W | \alpha)$$

To find minimum, differentiate both side.

$$\ln(p(t|x, w, \beta)) = \left(\frac{-\beta}{2}\right) \sum (y(\vec{x}_n) - \vec{t}_n)^2 + \left(\frac{N}{2}\right) \ln \beta - \left(\frac{N}{2}\right) \ln 2\pi + \frac{M+1}{2\pi} \ln \alpha + \left(\frac{-\alpha}{2} w^T W\right) \ln e.$$

In above equation,  $\left(\frac{N}{2}\right) \ln \beta$ ,  $\left(\frac{N}{2}\right) \ln 2\pi$  and  $\frac{M+1}{2\pi} \ln \alpha$  are constants W.r.t  $W$ . Hence it wont affect our maximization. The remaining terms are as follows

$$\ln(p(t|x, w, \beta)) = \left(\frac{-\beta}{2}\right) \sum (y(x_n) - t_n)^2 + \left(\frac{-\alpha}{2} w^T W\right) \ln e.$$

Differentiate W.r.t  $\vec{W}$ . and maximize  $p(t|x, w, \beta)$ , which is equal to minimizing  $\frac{1}{p(t|x, w, \beta)}$ . Hence, Minimizing equation becomes:

$$E(w) = \left(\frac{-\beta}{2}\right) \sum (y(x_n) - t_n)^2 + \left(\frac{-\alpha}{2} w^T W\right) \quad (11)$$

After solving above equation 11 similar to task 1, we get the following result

$$W_{MAP}^* = \left(\vec{X}^T \vec{X} + \frac{\alpha}{\beta} \vec{I}\right)^{-1} \vec{X}^T \vec{T} \quad (12)$$

Now to solve for  $\beta^*$ , since  $p(W | \alpha)$  is independent of  $\beta$ . we are left with only likelihood function, Hence maximization of  $\beta^{-1}$  posterior is same as Maximization of Likelihood in Task3.

$$\beta^{-1} = \frac{1}{N} \sum \left(\vec{X}^T \vec{W}_{MAP}^* - \vec{T}\right)^2 \quad (13)$$

### 3 Results

#### 3.1 Linear regression using error minimization

We have 2 parameter to vary,  $N$  (Number of observation) and  $M$  (Degree of polynomial).  $M$  will give more flexibility to fix the data point. with  $M = 0$ , Square sum error from all data points will be high. As we increase degree of freedom, curve start fitting better.

At  $M = 6$ , curve is almost equal to ground truth. when we increase  $M$  to 9, it start over-fitting. Follow below diagram and table for more intuition.

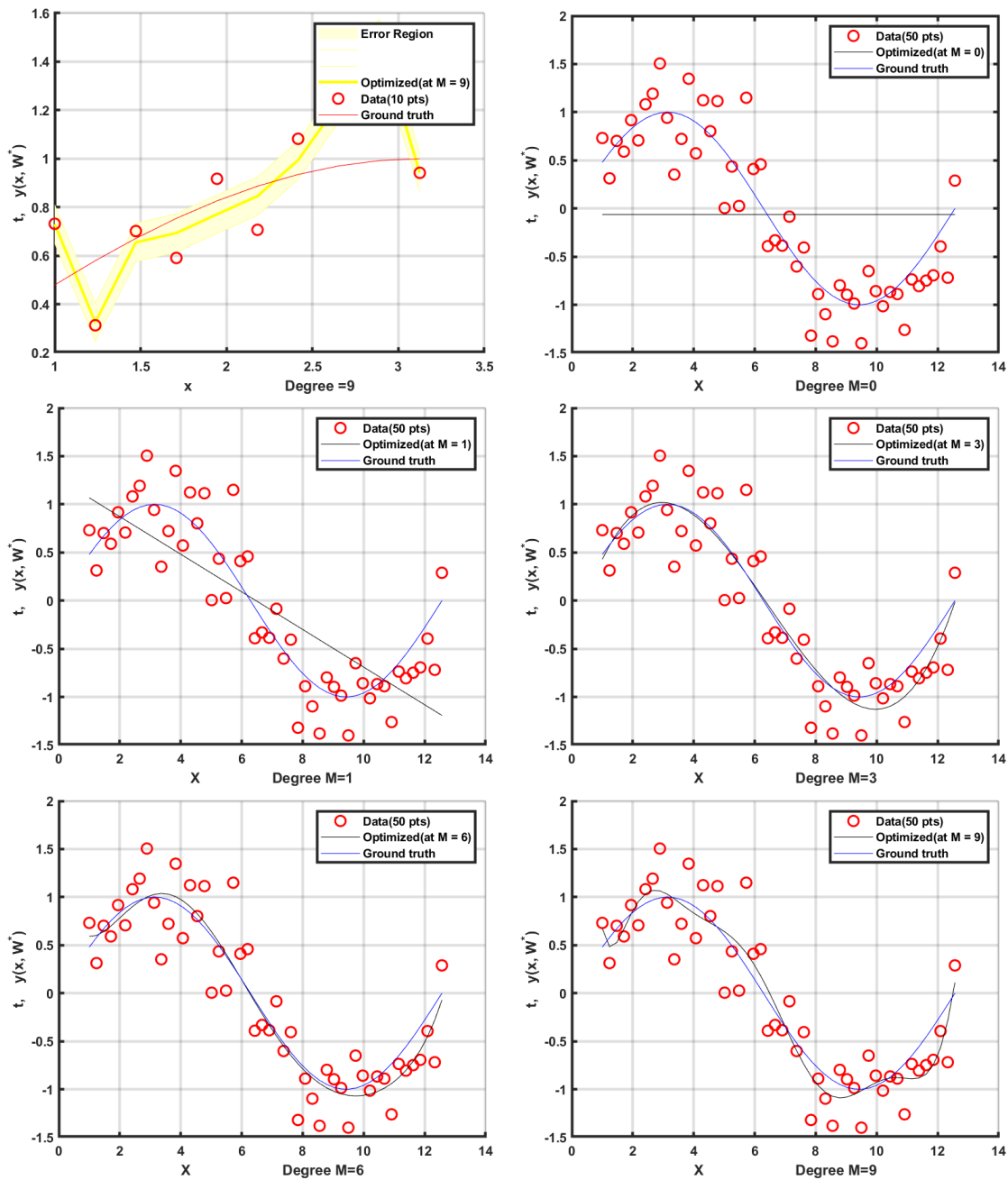


Figure 1:  $M = [0, 1, 3, 6, 9]$  with 50 data points

•

$W^*$ value Comparison					
	M = 0	M = 1	M = 3	M = 6	M = 9
$w_0^*$	-0.0621	1.2625	-0.4407	1.2186	11.0194
$w_1^*$		-0.1953	1.0991	-1.4575	-26.046
$w_2^*$			-0.2415	1.1059	25.1370
$w_3^*$			0.0125	-0.3127	-12.4559
$w_4^*$				0.0393	3.5895
$w_5^*$				-0.0023	-0.6346
$w_6^*$				0.0001	0.0696
$w_7^*$					-0.0046
$w_8^*$					0.0002
$w_9^*$					-0.0000

As we can see in below table, error is reducing along with increase in degree of polynomial

Erms				
M = 0	M = 1	M = 3	M = 6	M = 9
0.8421	0.5163	0.3177	0.3112	0.2929

On increasing N from 10 to 50, Error start getting low. Which means, W gets better and better fit when N value increases.

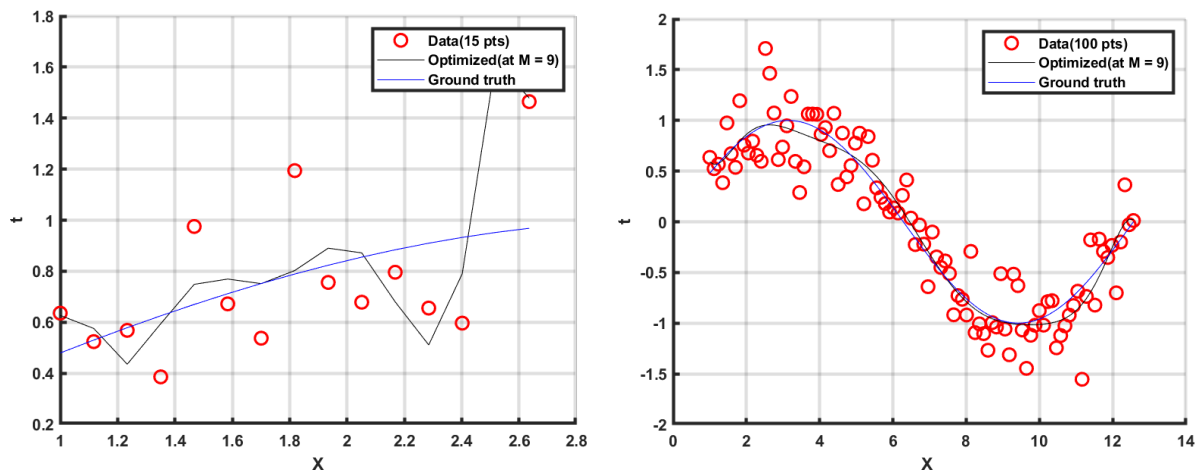


Figure 2:  $N = [15, 100]$  with  $M = 9$

•

### 3.2 Linear regression using error minimization with regularization

Without regularization large weights tries try to overfit when order of equation is 9. By adding regularization, we introduce penalty on large weight, which makes the curve more



smooth.

$\lambda$  here is called hyper-parameter. Changing its value will change the smoothness in the curve fitting.

Below images describe the smoothness in the curve with different values of lambda.

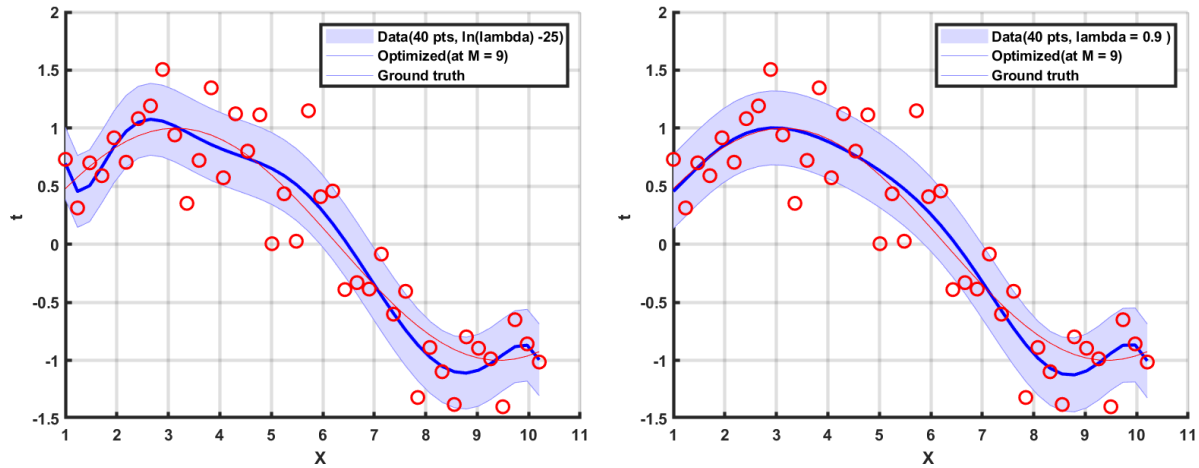


Figure 3:  $\lambda = [\exp^{-25}, 0.9]$  with  $M = 9$

- Similarly, we can see error rate varying with varying  $\lambda$

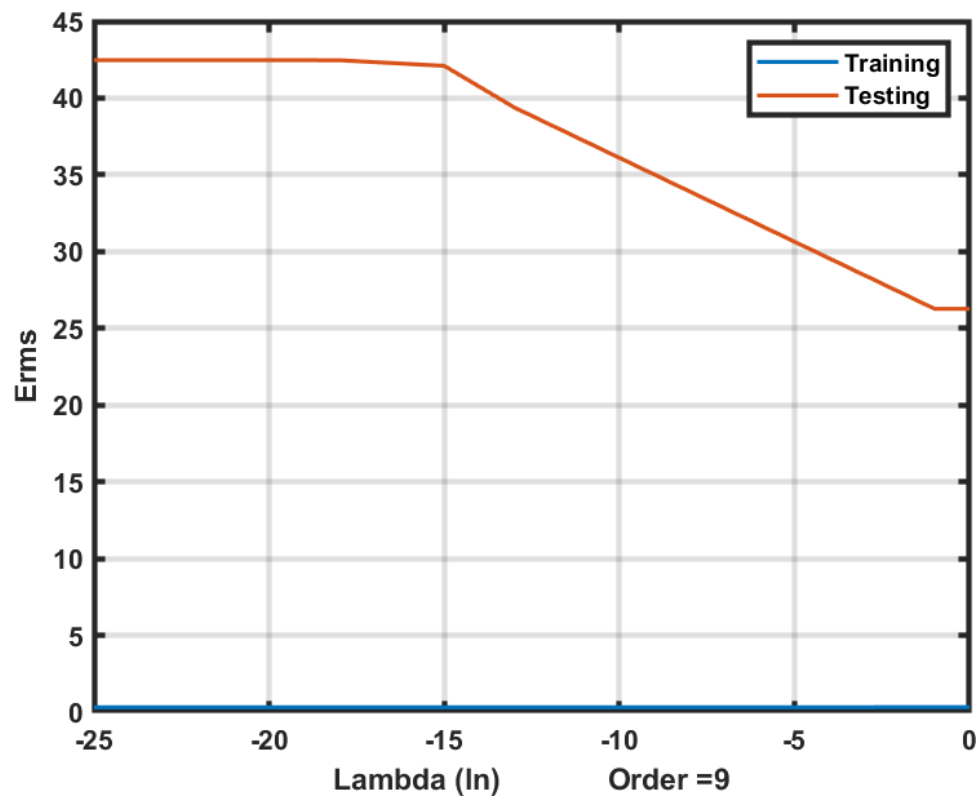


Figure 4:  $\lambda = [\lambda = e^{-25}, \lambda = e^{-18}, \lambda = e^{-15}, \lambda = e^{-13}, \lambda = e^{-1}, \lambda = 1]$  with  $M = 9$

	Error w.r.t $\lambda$					
	$\lambda = e^{-25}$	$\lambda = e^{-18}$	$\lambda = e^{-15}$	$\lambda = e^{-13}$	$\lambda = e^{-1}$	$\lambda = 1$
Train error	0.3093	0.3093	0.3093	0.3096	0.3165	0.3186
Test error	42.4763	42.4650	42.1000	39.3892	26.2731	26.2731

### 3.3 Maximum Likelihood Estimation

MLE has same results as Task 1, Hence avoiding writing duplication.

Below is figure with 10 data points and we can see over-fitting with degree 9

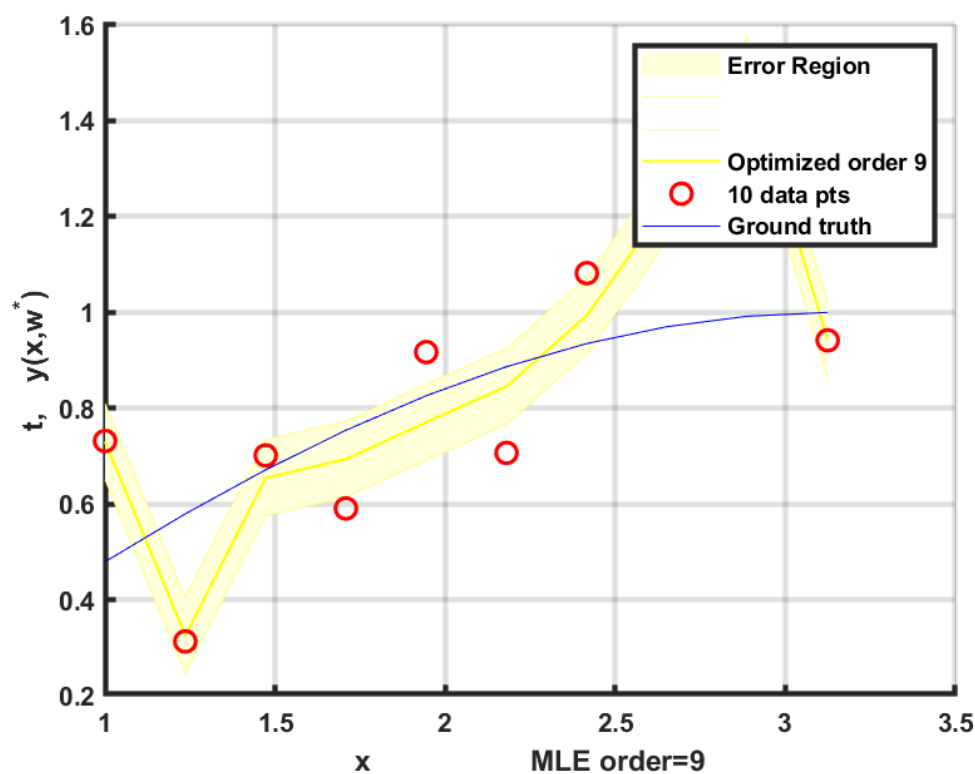


Figure 5: Data: 10pts,  $M = 9$

Comparison between Error minimization and MLE, both looks exactly same.

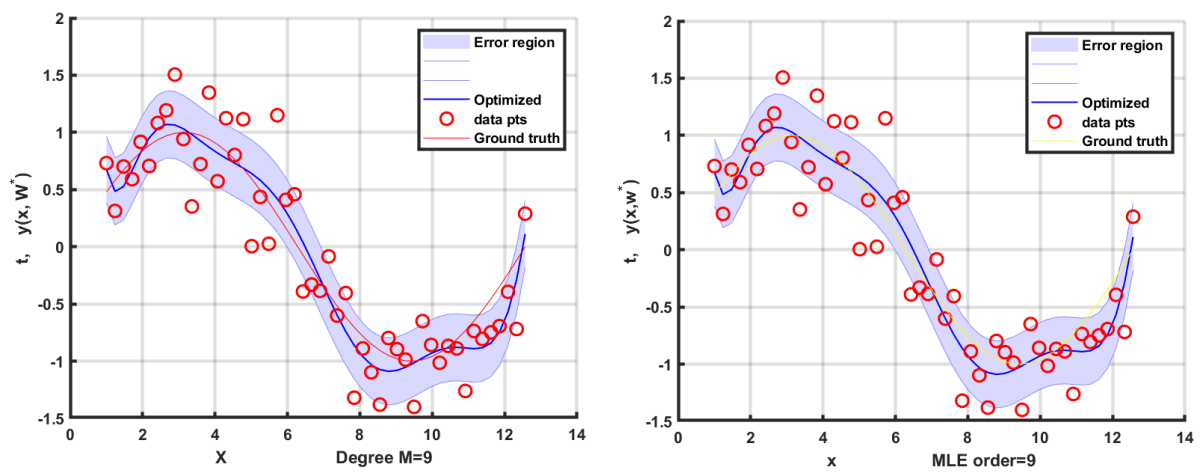


Figure 6: Comparison between EM and MLE with  $M = 9$  , data 50pts

### 3.4 Maximum Posterior Estimation

Maximum posterior draws its advantages from prior information, which is absent in maximum likelihood. In Posterior, we don't have to add  $\lambda$ , it comes internally in form of  $(\alpha$  and  $\beta)$ . Most of observation related to variation w.r.t  $N$  and  $M$  as same as previous 3 tasks, Hence omitted repetition.

Below figure is with 10 data points and error in training data set is about 0.1088, which is little better than what we got from various  $\lambda$  values in task 2 (though comparison is not valid as  $\lambda$  and  $\alpha/\beta$  is not equal, but still it is good improvement compared to moving  $\lambda$  from  $e^{-25}$  to 1).

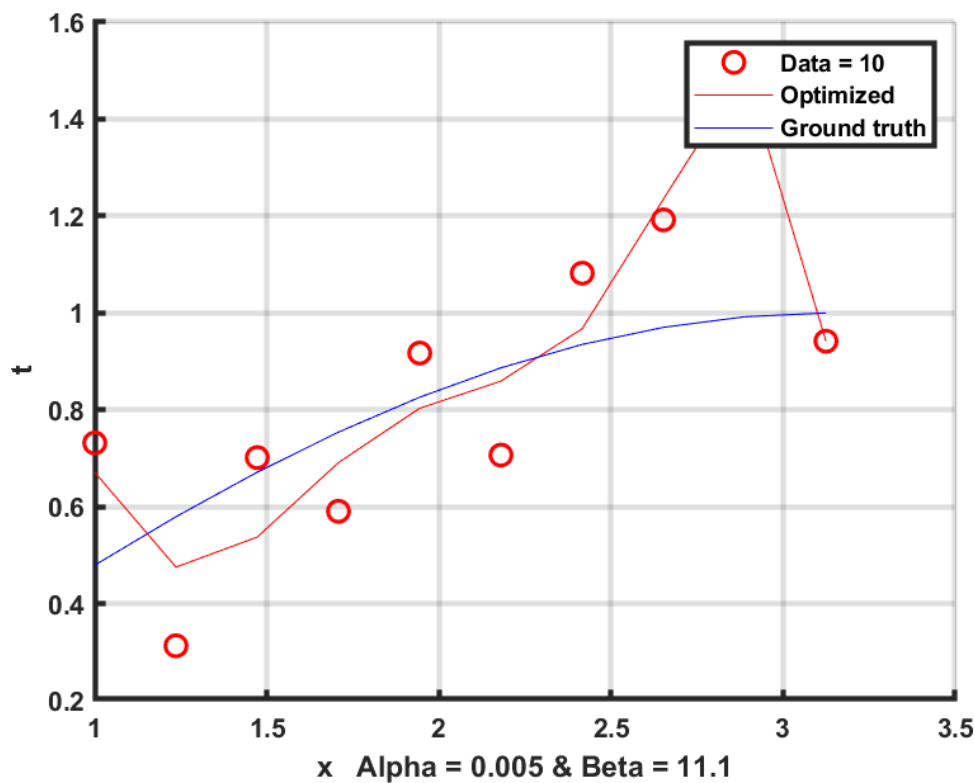


Figure 7: Data: 10pts,  $\alpha = 0.005$ ,  $\beta = 11.1$ ,  $M = 9$

Varying  $\alpha$  value. In below graph we can see that on increasing  $\alpha$  variance is also getting increased. Note: Difference in image is not so conspicuous, for instance it can be seen at leftmost 2 data point.

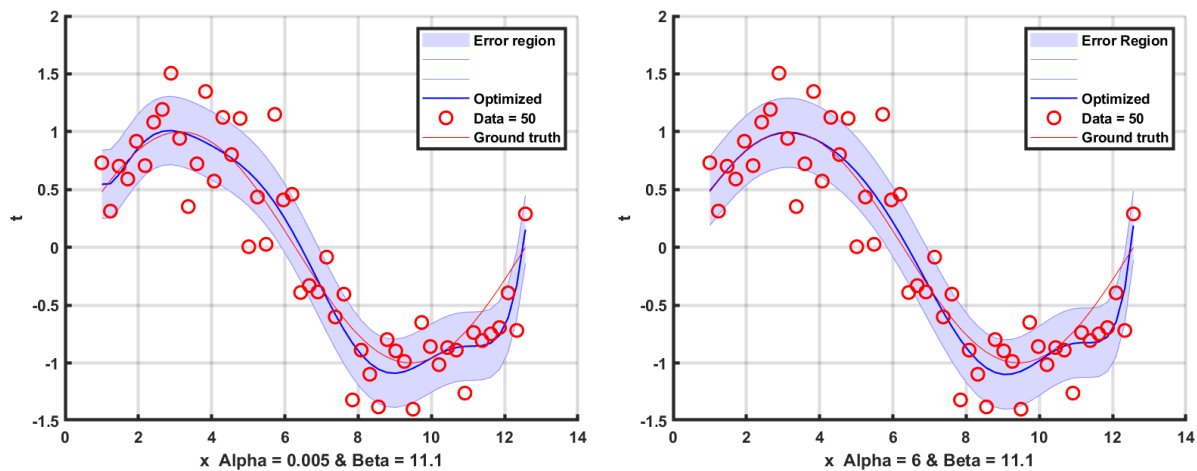


Figure 8: a)  $\alpha = 0.005$ , b)  $\alpha = 6$  with  $\beta = 11.1$ , data 50pts

## 4 Conclusion

Increasing number of data points ( $N$ ) in the dataset will reduce the error training error (Given that degree of polynomial sufficient to handle data distribution). Increasing number degree of polynomial certainly reduce the training error but may increase testing error. To avoid  $W^*$  values bigger we can introduce penalty, which is called regularization. We need to make balance between  $N$  and  $M$  such that we can avoid testing error. Another way of solving this curve fitting is using probability and decision theory. From the observation, it can be concluded that Maximum posterior estimation performs better then linear regression with regularization, but its performance depends on prior. Prior can depend on subject knowledge, hence it may take different time to converge.

**Improvement** Since most of the problem statement are open ended. 85% of the time was spent in writing project report, which actually turned starting fun into pain. It would have been great if it was just programming assignment.

## 5 References

### References

- [1] Bishop, Christopher M. Pattern Recognition and Machine Learning.