

## Question 2

1) We can use a categorical/multinoulli distribution to describe this scenario.

We will have six parameters:

$\mu_1, \mu_2, \mu_3, \mu_4, \mu_5, \mu_6$

the probabilities for side 1, 2, 3, 4, 5, 6 comes up.

2) If we have a fair dice then

$$\mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5 = \mu_6 = \frac{1}{6}$$

3) If the die always rolls two then

$$\mu_2 = 1, \quad \mu_1 = \mu_3 = \mu_4 = \mu_5 = \mu_6 = 0.$$

4) the domain of parameters is  $[0, 1]$

$$\mu_1 \in [0, 1] \quad \mu_2 \in [0, 1] \quad \mu_3 \in [0, 1]$$

$$\mu_4 \in [0, 1] \quad \mu_5 \in [0, 1] \quad \mu_6 \in [0, 1]$$

## Question 2

$$E_B(w) = \frac{1}{2} \sum_{n=1}^N \alpha_n \{t_n - w^T \phi(x_n)\}^2$$

$$\frac{\partial}{\partial w} E_B(w) = \frac{1}{2} \sum_{n=1}^N 2(t_n - w^T \phi(x_n)) (-\phi(x_n)) \cdot \alpha_n$$

$$= \sum_{n=1}^N (t_n - w^T \phi(x_n)) (-\phi(x_n)) \cdot \alpha_n$$

$$\nabla E_B(w) = \sum_{n=1}^N (t_n \alpha_n - w^T \alpha_n \phi(x_n)) (-\phi(x_n))^T$$

$$\phi(x_n) = \begin{bmatrix} \phi_0(x_n) \\ \phi_1(x_n) \\ \vdots \\ \phi_M(x_n) \end{bmatrix}$$

$$0^T = [0, 0, 0, \dots, 0]$$

$$\nabla E_B(w) = \left[ \frac{\partial}{\partial w_0} \ln(\cdot), \frac{\partial}{\partial w_1} \ln(\cdot), \dots, \frac{\partial}{\partial w_M} \ln(\cdot) \right]$$

Set the gradient to 0.

$$0^T = \nabla E_B(w) = \sum_{n=1}^N (t_n \alpha_n - w^T \alpha_n \phi(x_n)) (-\phi(x_n))^T$$

$$0^T = \sum_{n=1}^N -t_n \alpha_n \phi(x_n) + w^T \underbrace{\sum_{n=1}^N \alpha_n \phi(x_n) \phi(x_n)^T}_{\text{}}.$$

$$0^T = t^T \alpha^T \Phi - w^T \Phi^T \alpha^T \Phi$$

Why?

Using this part as an example:

$$\underbrace{\phi(x_n)}_{M \times 1} \underbrace{\phi(x_n)^T}_{1 \times M} = \begin{bmatrix} \phi_0(x_n) \\ \phi_1(x_n) \\ \vdots \\ \phi_M(x_n) \end{bmatrix} \begin{bmatrix} \phi_0(x_n) & \phi_1(x_n) & \dots & \phi_M(x_n) \end{bmatrix}$$

$$= \begin{pmatrix} \phi_0(x_n) \phi_0(x_n) & \phi_0(x_n) \phi_1(x_n) & \dots & \phi_0(x_n) \phi_M(x_n) \\ \phi_1(x_n) \phi_0(x_n) & \phi_1(x_n) \phi_1(x_n) & \dots & \phi_1(x_n) \phi_M(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_M(x_n) \phi_0(x_n) & \phi_M(x_n) \phi_1(x_n) & \dots & \phi_M(x_n) \phi_M(x_n) \end{pmatrix}$$

Treat  $\alpha_n$  as a scalar then  $\alpha_n \phi(x_n) \phi(x_n)^T$  become

$$\begin{pmatrix} \alpha_n \phi_0(x_n) \phi_0(x_n) & \dots \\ \alpha_n \phi_1(x_n) \phi_0(x_n) & \dots \\ \vdots & \ddots \end{pmatrix}$$

$$\sum_{n=1}^N \alpha_n \phi(x_n) \phi(x_n)^T ?$$

$$\alpha_1 \phi_0(x_1) \phi_0(x_1) + \alpha_2 \phi_0(x_2) \phi_0(x_2) + \dots + \alpha_N \phi_0(x_N) \phi_0(x_N)$$

$$\alpha_1 \phi_1(x_1) \phi_0(x_1) + \alpha_2 \phi_1(x_2) \phi_0(x_2) + \dots + \alpha_N \phi_1(x_N) \phi_0(x_N)$$

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{pmatrix} \quad \underline{N \times M}$$

$$\Phi^T = \begin{pmatrix} \phi_0(x_1) & \phi_0(x_2) & \dots & \phi_0(x_N) \\ \phi_1(x_1) & \phi_1(x_2) & \dots & \phi_1(x_N) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{M-1}(x_1) & \phi_{M-1}(x_2) & \dots & \phi_{M-1}(x_N) \end{pmatrix} \quad M \times N$$

$$\alpha = \begin{pmatrix} \alpha_1 & 0 & 0 & \dots & 0 \\ 0 & \alpha_2 & 0 & \dots & 0 \\ 0 & 0 & \alpha_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \alpha_N \end{pmatrix} \quad N \times N$$



$$\frac{\Phi^T \alpha}{M \times N} = \begin{pmatrix} \alpha_1 \phi_0(x_1) & \alpha_2 \phi_0(x_2) & \dots & \alpha_N \phi_0(x_N) \\ \alpha_1 \phi_1(x_1) & \alpha_2 \phi_1(x_2) & \dots & \alpha_N \phi_1(x_N) \\ \vdots & \vdots & & \vdots \\ \alpha_1 \phi_{M-1}(x_1) & \alpha_2 \phi_{M-1}(x_2) & \dots & \alpha_N \phi_{M-1}(x_N) \end{pmatrix}$$

$$\frac{\Phi^T \alpha \Phi}{M \times M} = \begin{pmatrix} \text{ } \\ \text{ } \end{pmatrix}$$

$$\rightarrow \alpha_1 \phi_0(x_1) \phi_0(x_1) + \alpha_2 \phi_0(x_2) \phi_0(x_2) + \dots$$

$$\rightarrow \alpha_1 \phi_1(x_1) \phi_0(x_1) + \alpha_2 \phi_1(x_2) \phi_0(x_2) + \dots$$

Same as  $\begin{pmatrix} \alpha_N \phi_0(x_N) \phi_0(x_N) & \dots \\ \alpha_N \phi_1(x_N) \phi_0(x_N) & \dots \\ \vdots & \ddots \end{pmatrix}$

Therefore,  $\Phi^T = t^T \alpha^T \Phi - w^T \Phi^T \alpha^T \Phi$

$$w^T \Phi^T \alpha^T \Phi = t^T \alpha^T \Phi$$

$$(\Phi^T \alpha \Phi) w = \Phi^T \alpha t$$

$$w = (\Phi^T \alpha \Phi)^{-1} (\Phi^T \alpha t)$$

$\left. \begin{aligned} (ABC)^T &= C^T B^T A^T \\ (AB)^T &= B^T A^T \end{aligned} \right\}$ 
 $(\Phi^T \alpha \Phi)^T = t^T \alpha^T \Phi$   
 $(\Phi^T \alpha \Phi) w = w^T \Phi^T \alpha^T \Phi$

### Question 3

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2$$

$$E(\tilde{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \frac{\lambda}{2} \|w\|^2$$

$$RMS: E_{RMS} = \sqrt{2E(w^*)/N}$$

1) No, The training set and validation set are both randomly distributed data, from the dataset. There is no guarantee on the relationship between training error and validation error.

If the model is overfit the validation error is probably higher than the training error.

Good fit: Validation error low, slightly higher than the training error.

Unknown fit: Validation error low, training error high.

Under fit: Validation error and training error both high.

Generally speaking, training error will almost always underestimate the validation error.

But it is possible for the validation error to be less than the training error.

2) Yes, Degree 10 polynomial contains

Degree 9 polynomial. The unregularized regression gives us the optimal solution which means

Degree 10 polynomial mostly fits the data better.

In the worst case, the training error



for them two are equal.

But if we change training error to testing error for this question then the answer should be "No".

3) No.

In most cases the testing error for regularized regression is lower than unregularized regression since the degree=20 is very high and is highly likely to cause overfitting.

But this is not guaranteed. If we got a weak model then the regularized will underfitting slash the predictive power even more and make the testing error larger compared with the unregularized one.

Question 4.

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{t_n - w^T \phi(x_n)\}^2 + \frac{1}{2} \sum_{j \in J_1} \lambda_j |w_j| + \frac{1}{2} \sum_{j \in J_2} \lambda_j |w_j|^2.$$

Here,  $\lambda_{j \in J_1}$  and  $\lambda_{j \in J_2}$  are two subsets/subvector of  $\lambda_{n \in N}$ ,

$$\begin{aligned} \frac{\partial E(w)}{\partial w} &= \frac{1}{2} \sum_{n=1}^N 2(t_n - w^T \phi(x_n))(-\phi(x_n)) + \frac{1}{2} \sum_{j \in J_1} \frac{w_j}{|w_j|} \lambda_j \\ &\quad + \frac{1}{2} \sum_{j \in J_2} \lambda_j 2w_j \\ &= \sum_{n=1}^N (t_n - w^T \phi(x_n))(-\phi(x_n)) + \frac{1}{2} \sum_{j \in J_1} \frac{w_j}{|w_j|} \lambda_j \\ &\quad + \sum_{j \in J_2} \lambda_j w_j. \end{aligned}$$

Define matrices:

$$\varphi = \begin{cases} \frac{w_j}{|w_j|} \lambda_j & \text{for } j \in J_1 \\ 0 & \text{otherwise} \end{cases}$$

$$r = \begin{cases} \lambda_j & \text{for } j \in J_2 \\ 0 & \text{otherwise} \end{cases}$$

$$I = \text{Identity matrix} = [1, 1, 1, 1, \dots]^T.$$

$$\nabla E(w) = -t^T \Phi + W^T \Phi^T \Phi + \frac{1}{2} \varphi I^T + r \cdot W^T$$

Here  $\Phi$  is the same design matrix we declared in class



Question 5.

5.1) Niger.  $313.7/1000 = 31.37\%$

Sierra Leone.  $185.3/1000 = 18.53\%$

"na. values" = "\_" will set the missing features to NA values.

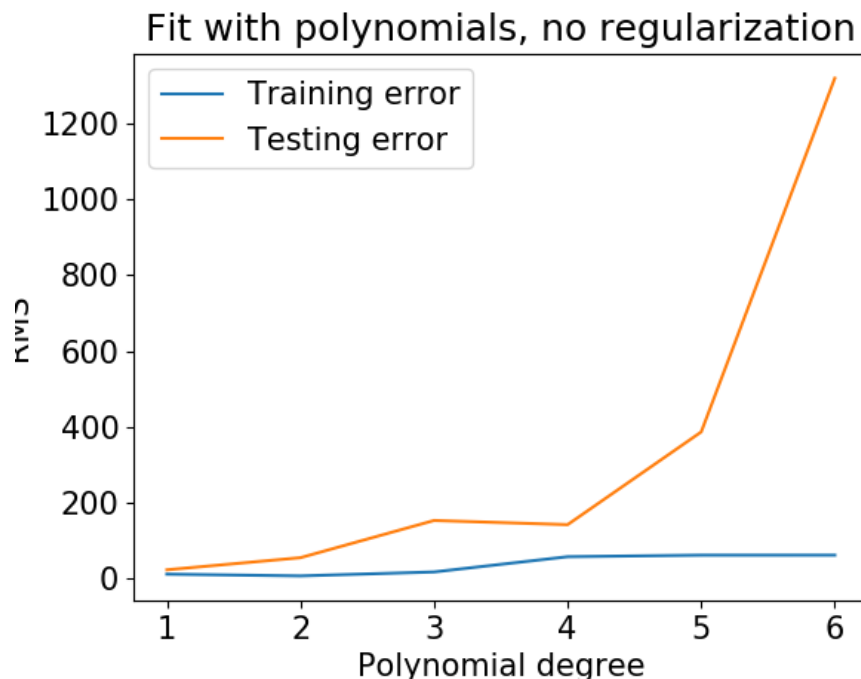
Then we will use `nanmean()` to find the average value for each column/feature, `np.where(np.isnan)` will find the coordinates/indices for the NA values, and we assign the average value to the missing parts according to their indices we found.

NB: `np.where(np.isnan)` will return the row and column indices. But `np.take(mean_vals, inds[1])` → we only need the column indices since the average is a  $1 \times 40$  matrix/vector.  
 $1 \times N$



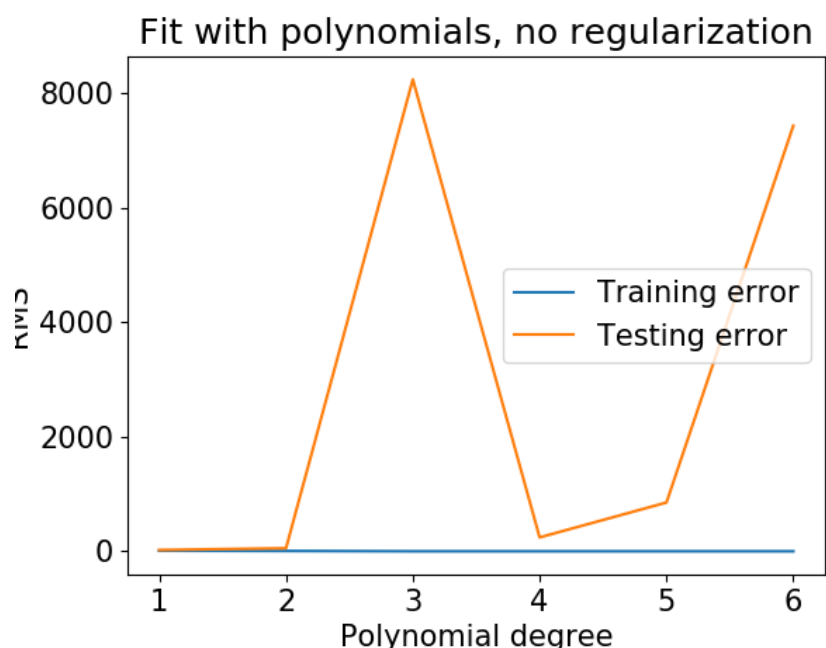
## 5.2

1)



There are two things wrong here. The first one is usually as polynomial degree increases the training error will decrease since higher degree usually fits the data better. After we apply the normalization on input features this problem is solved. The reason is that we may have some features whose orders of magnitudes are larger than others, and they may dominate the objective function and make the estimator unable to learn from other features correctly as expected.

Another problem is as orders of degree increase, the model is being more over-fitted. This problem can be solved with regularized regression which we will do in 5.4.

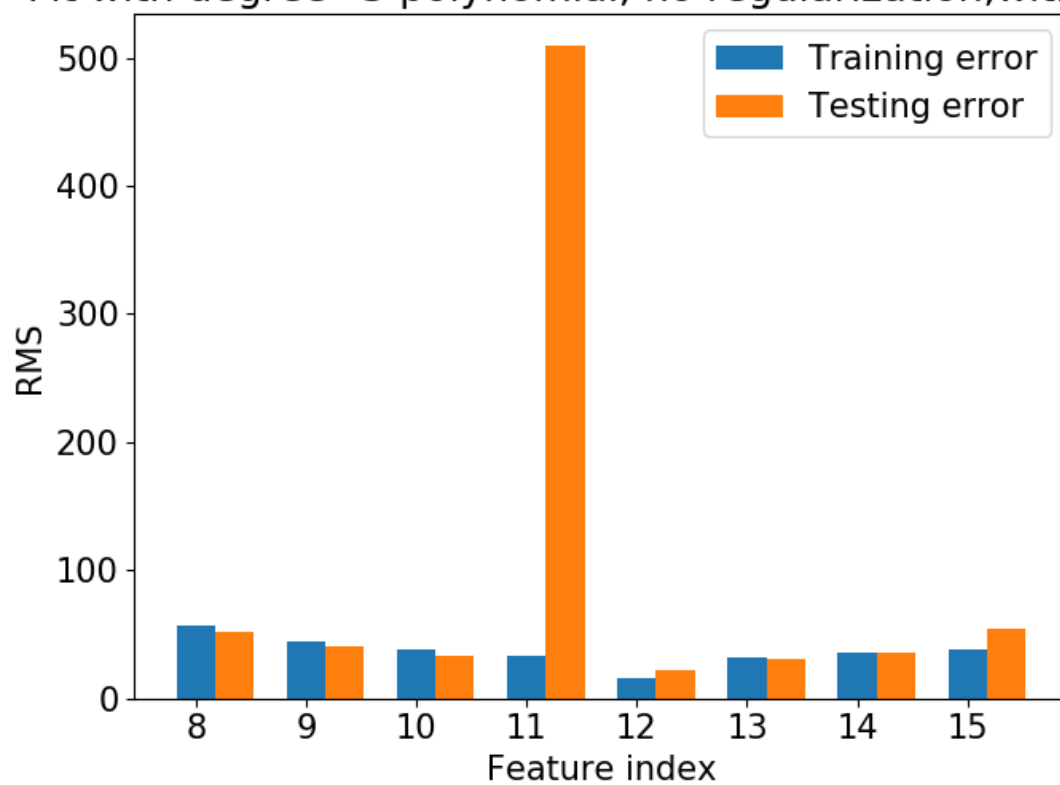


2)

Fit with degree=3 polynomial, no regularization, without bias

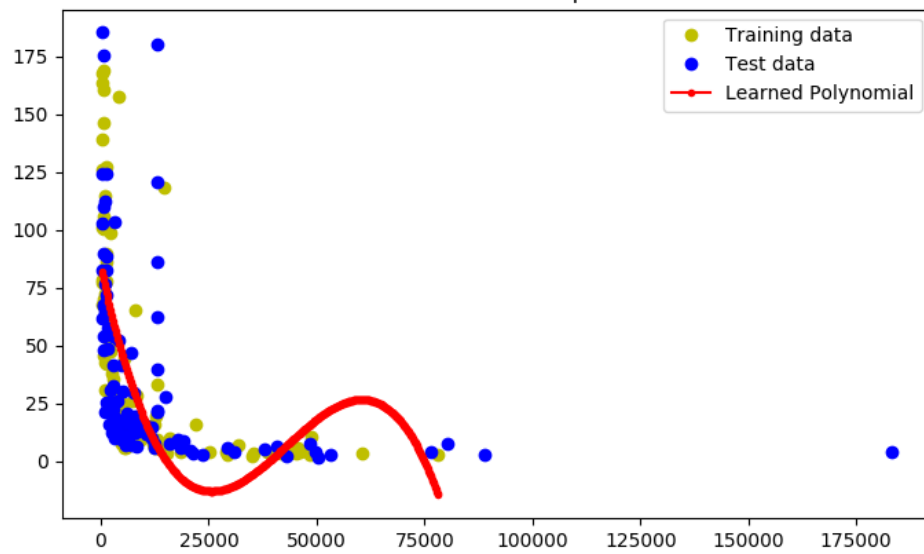


Fit with degree=3 polynomial, no regularization, with bias



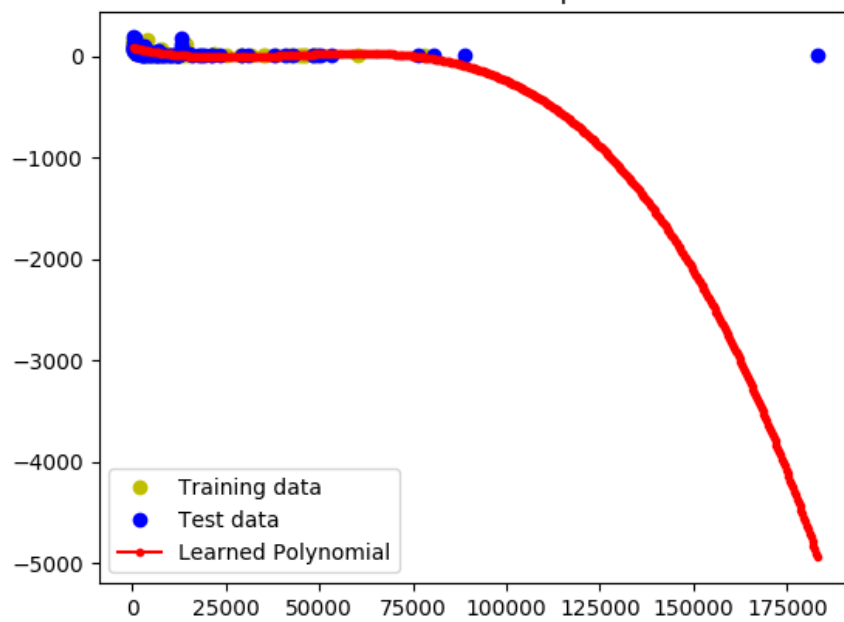


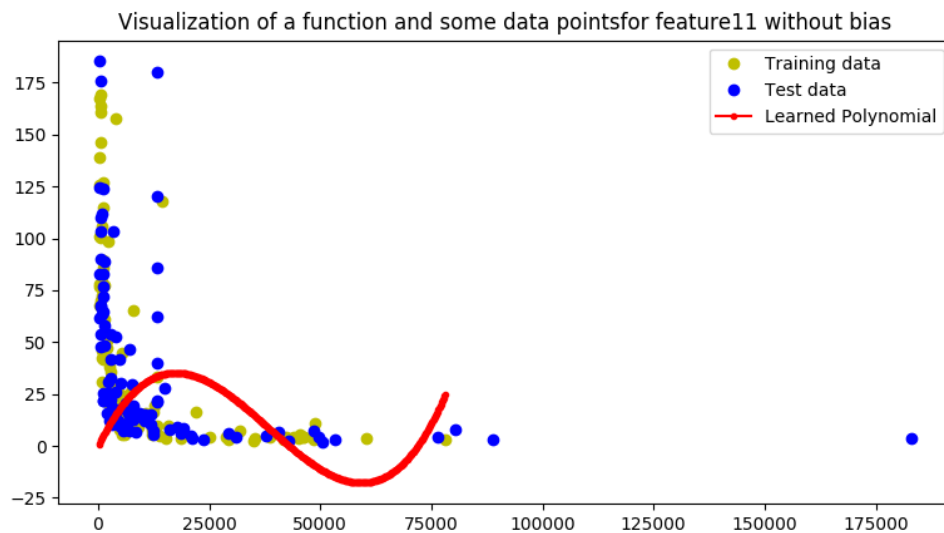
Visualization of a function and some data points for feature11 with bias



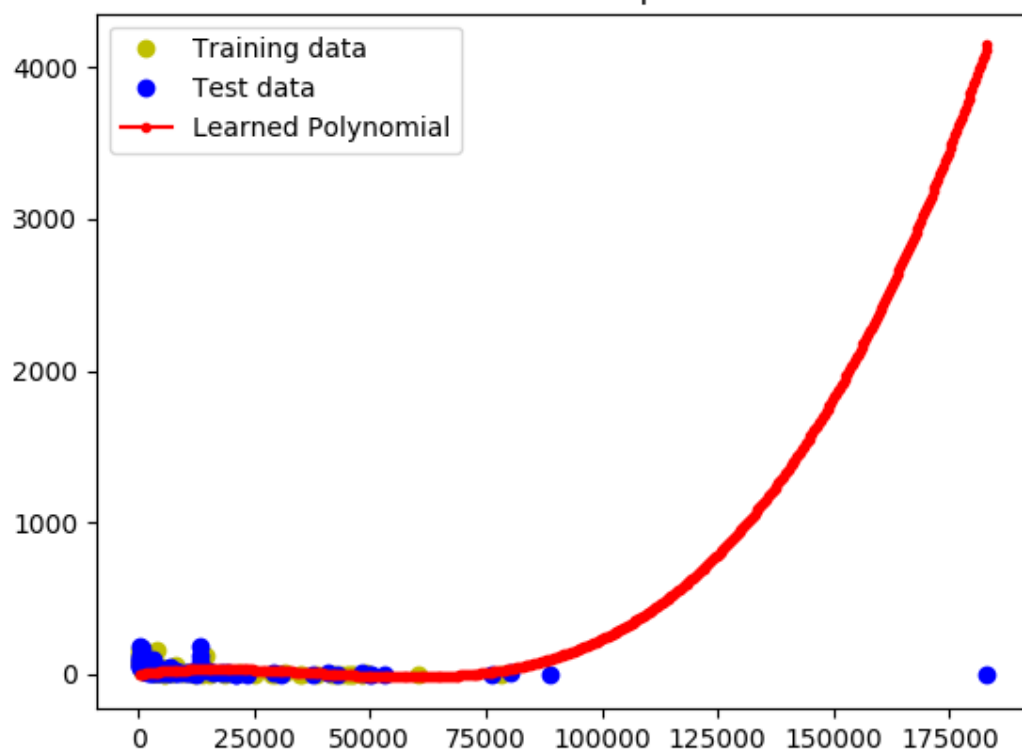
Since there is an outlier so we need to use a larger scale.

Visualization of a function and some data points for feature11 with bias



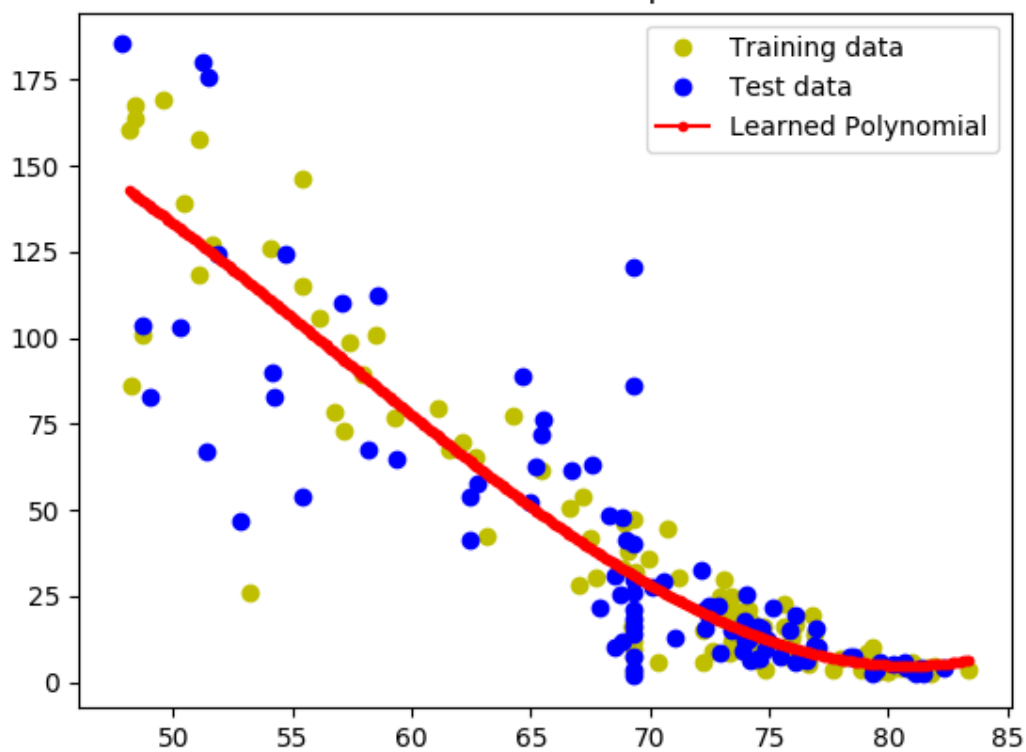


Visualization of a function and some data points for feature11 without bias

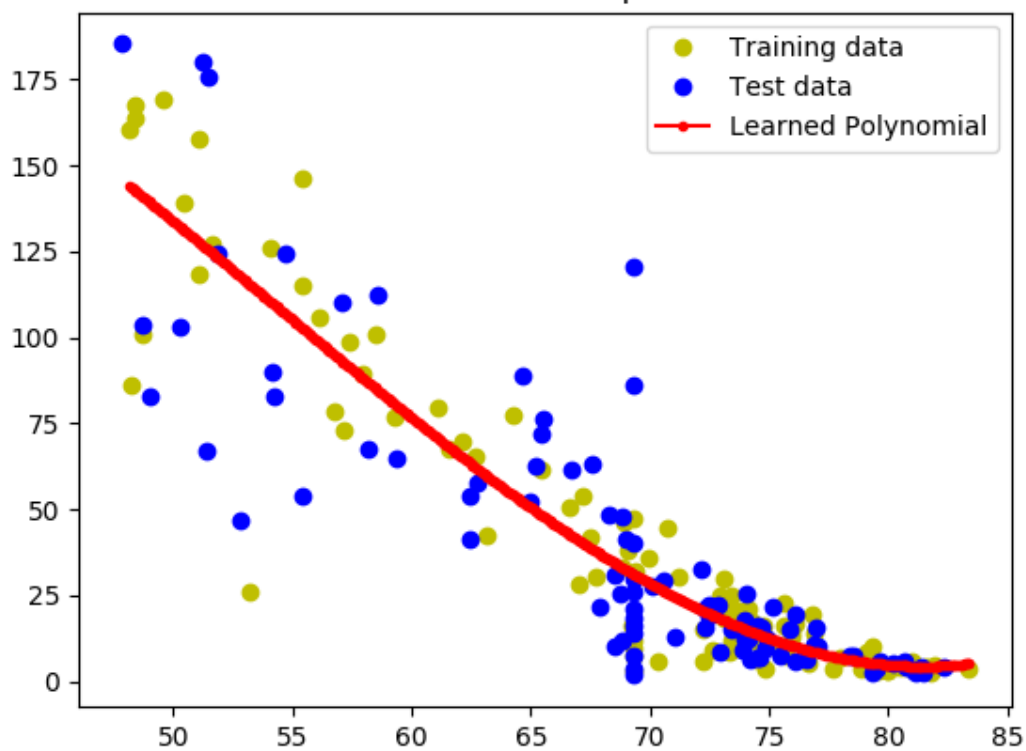




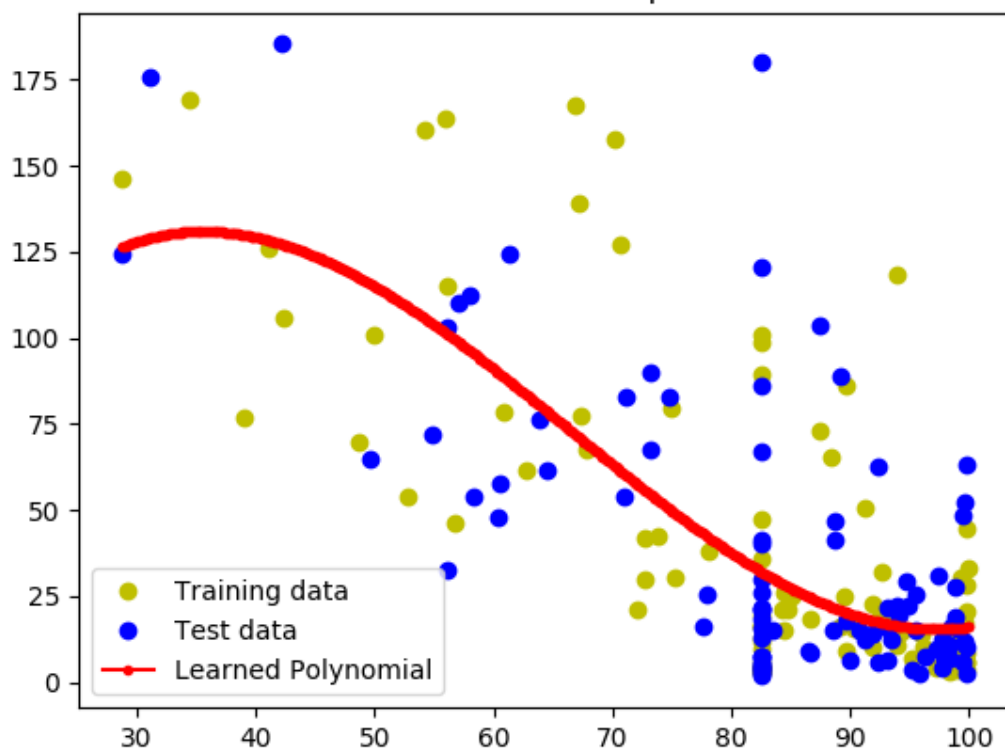
Visualization of a function and some data points for feature12 with bias



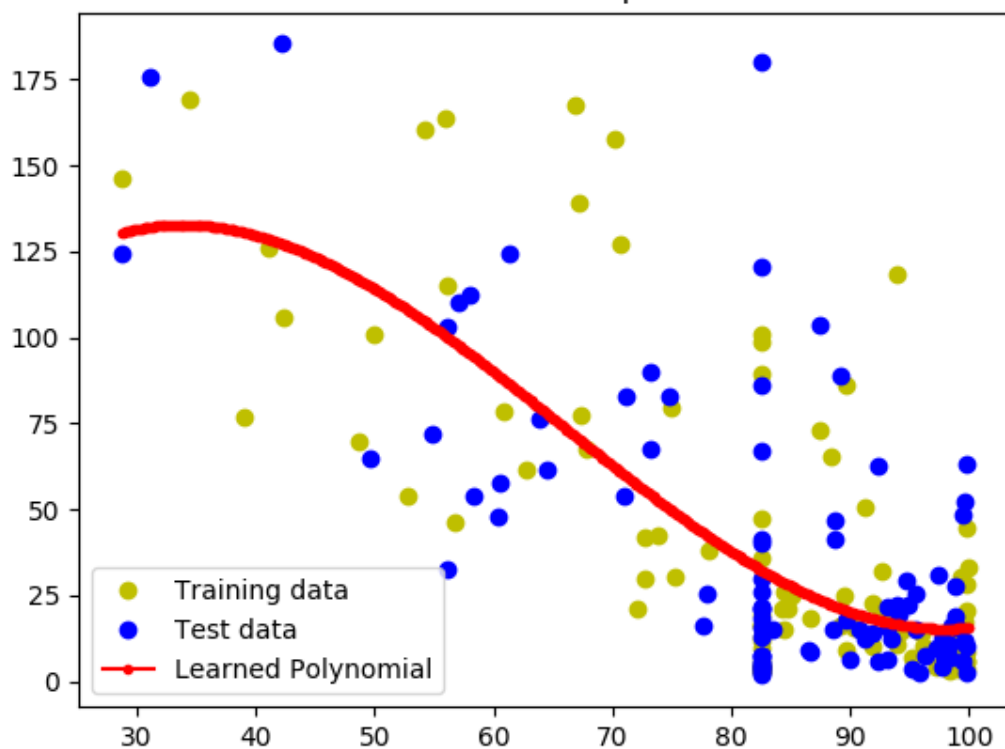
Visualization of a function and some data points for feature12 without bias



Visualization of a function and some data points for feature13 with bias

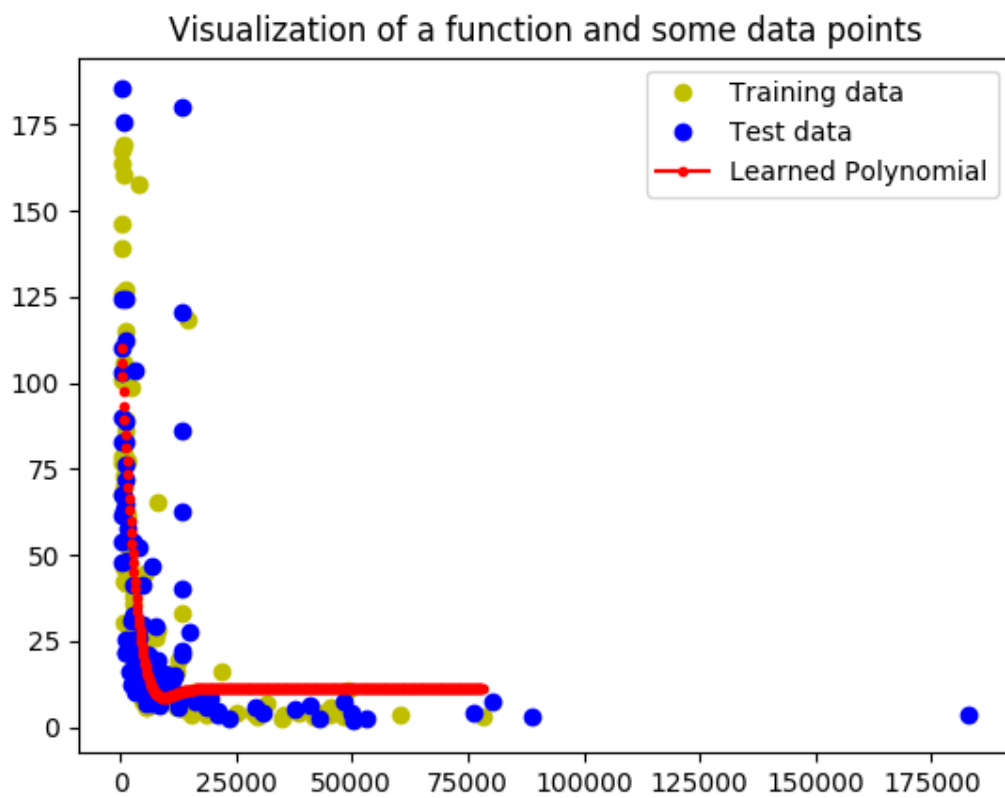


Visualization of a function and some data points for feature13 without bias





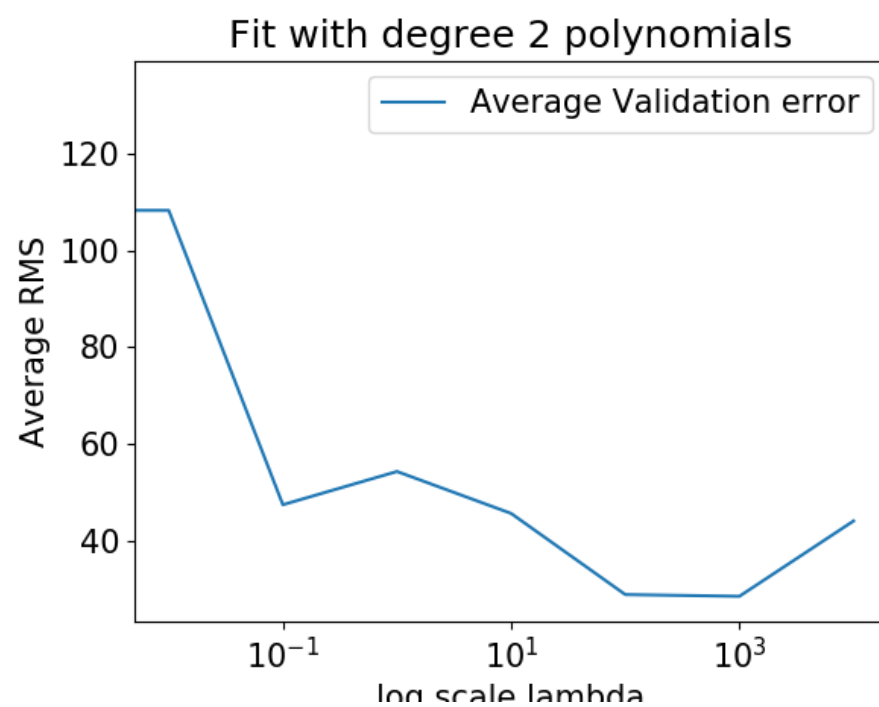
5.3



training error: `[[28.45793776]]`

testing error: `[[33.8067249]]`

5.4



the errors for lambda 0 to 1000 are :

lambda=0	134.08724800120225
lambda=0.01	108.33938137863713
lambda=0.1	47.42014676680511
lambda=1	54.30153018685256
lambda=10	45.61747891270658
lambda=100	28.827211650792474
lambda=1000	28.46922327978877
lambda=10000	44.06076704806573

Lambda=1000 yields the lowest average validation error so we will choose lambda=1000