

Name: Bharat Kashyap Karri
ID: 2020AAPS0319H

FODS Assignment 1

1A

1A) Given $\text{prior} = \beta(u|a,b) \rightarrow \text{Beta distribution}$

NAME: BHARAT KASHYAP KARRE ID: 2020AAPS01194H

$$\text{prior } \beta(u|2,2) = \frac{\Gamma(2+2)}{\Gamma(2)\Gamma(2)} \cdot u^{2-1} \cdot (1-u)^{2-1} \quad \Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$$

$$\text{In general } \beta(u|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot u^{a-1} \cdot (1-u)^{b-1} = P(u)$$

$u \rightarrow \text{true mean}$ $D \rightarrow \text{dataset}$ $u = P(\text{people liking the update})$

$$P(u|D) = \frac{P(D|u) \cdot P(u)}{P(D)} \quad P(D) \rightarrow \text{constant}$$

$$P(u|D) \propto P(D|u) \cdot P(u)$$

$$\propto P(\{ \text{liked} \} | u) \cdot P(\{ \text{liked} \} | u) \dots P(\{ \text{liked} \} | u) \cdot P(\{ \text{disliked} \} | u) \dots P(\{ \text{disliked} \} | u) \cdot P(u)$$

$$\propto u^m \cdot (1-u)^n \cdot P(u) \quad \text{where } m \text{ people in the sample liked the update and } n \text{ people did not like the update.}$$

$$\propto u^m (1-u)^n \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot u^{a-1} \cdot (1-u)^{b-1}$$

$$\propto u^{m+a-1} \cdot (1-u)^{n+b-1}$$

$$P(u|D) = \frac{u^{m+a-1} \cdot (1-u)^{n+b-1}}{\int_0^1 u^{m+a-1} \cdot (1-u)^{n+b-1} du} = \frac{\Gamma(a+b+m+n-1)}{\Gamma(m+a-1)\Gamma(n+b-1)} \cdot u^{m+a-1} \cdot (1-u)^{n+b-1}$$

$$P(u|D) = \beta(u|a+m, b+n)$$

$$\text{Prior: } \beta(u|2,2) \quad \text{Posterior 1: } \beta(u|2+40, 2+10) = \beta(u|42, 12) \rightarrow \text{Prior to next posterior}$$

$$\text{Posterior 2: } \beta(u|42+15, 12+17) = \beta(u|57, 29)$$

$$\text{Posterior 3: } \beta(u|57+175, 29+59) \rightarrow \text{Posterior of 0}$$

Name: WHAT KARYN KARRI

ID: 2020AAP50198

If likelihood θ is a binomial distribution then the posterior θ given n trials distribution

Now the θ Beta distribution is given by $\beta(a/105, 55)$.

Using the prior probability $\mu : \mu = 0.41$

Likelihood of $\theta = P(\text{label the update}) = 0.41$

Probability for m people in a sample size of N to have liked the update

$$= \frac{N!}{m!(N-m)!} \cdot (0.41)^m \cdot (0.59)^{N-m}$$

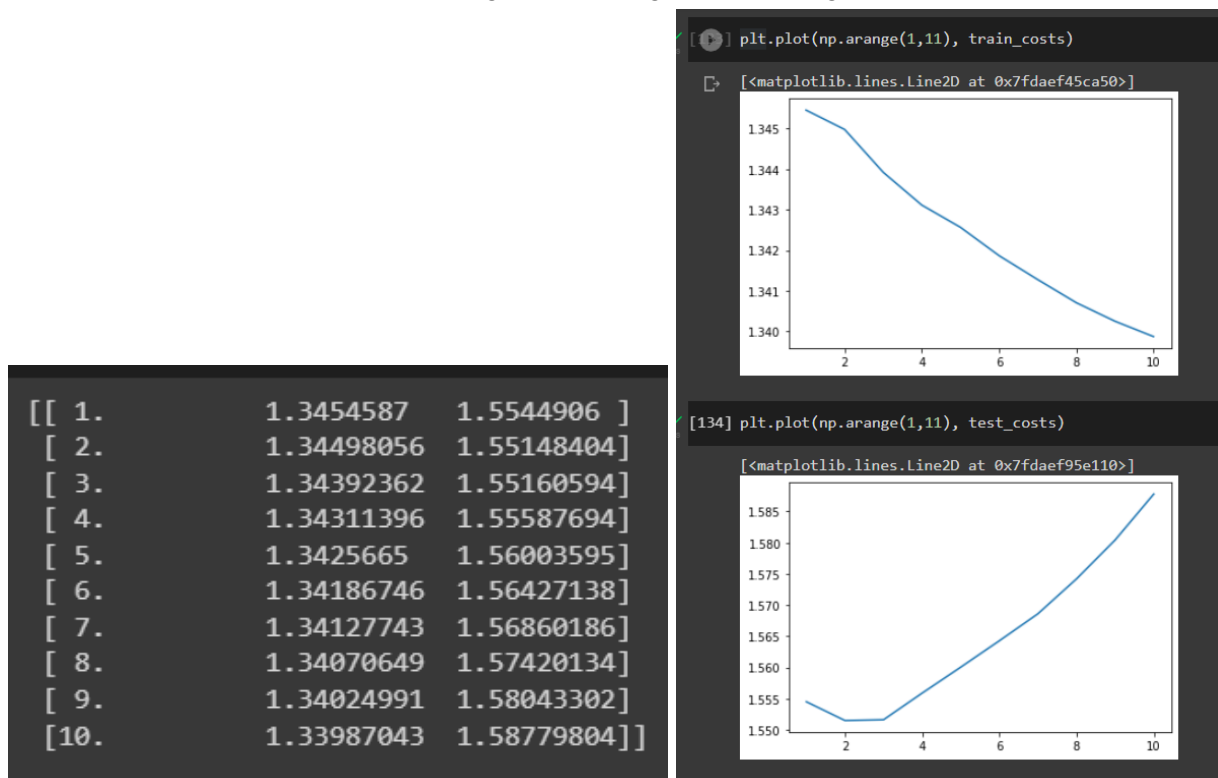
1B

i) The model is a regression model. In the first few implementations it is a polynomial regression model where each column is gaussian normalized to prevent overflow. Each column corresponds to the transformed input data.

In batch gradient descent all of the sample data is used to train the model in each epoch. This results in a slower training process with more number of iterations. With stochastic gradient descent training is done with only one sample for each epoch. Thus training is faster with approximately the same accuracy.

Regularization is taken as the sum of absolute values of each weight to the power of the required norm. It is also multiplied by a hyper-parameter called the regularization constant.

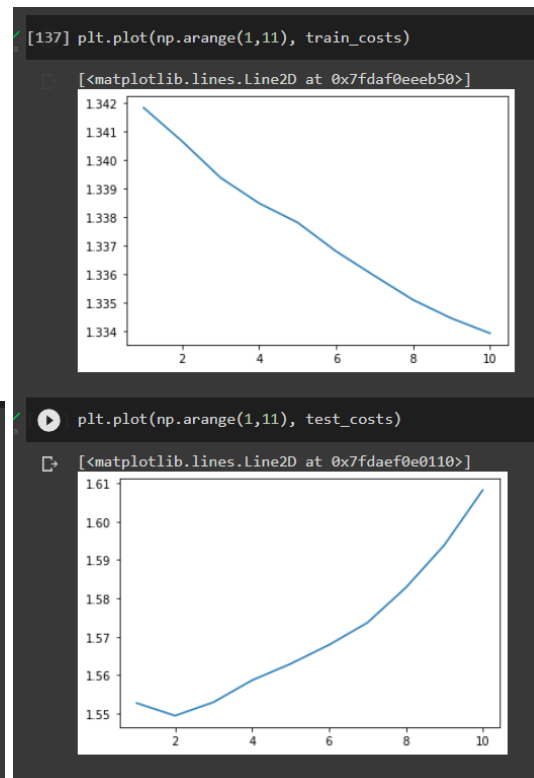
ii) Batch Gradient Descent errors: Degree, Training error, Testing error



The first graph shows the training error and shows that the higher the degree of the polynomial lesser is the training error. However it can be seen that the testing error rises after degree 2 thus showing that the model starts to overfit after degree 2.

Stochastic Gradient Descent errors: Degree, Training error, Testing error

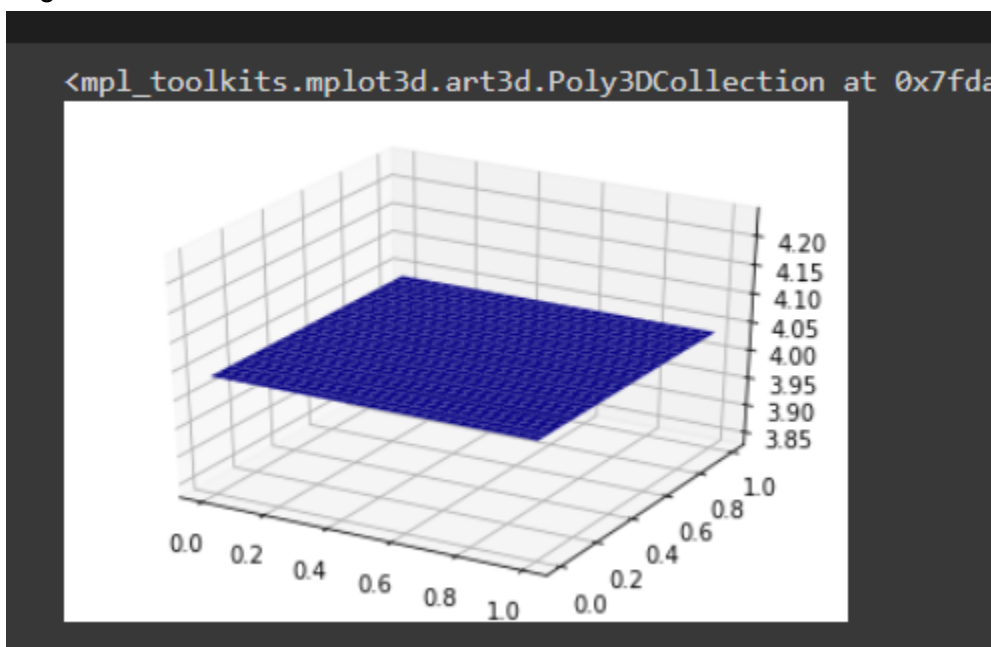
[[1.	1.34181892	1.55266317]
[2.	1.34064605	1.5493665]
[3.	1.33937258	1.55289108]
[4.	1.33848216	1.55859585]
[5.	1.33781033	1.56287786]
[6.	1.33680346	1.56782025]
[7.	1.33594635	1.57361976]
[8.	1.33510799	1.58276294]
[9.	1.33446126	1.59386312]
[10.	1.33393898	1.60820828]]



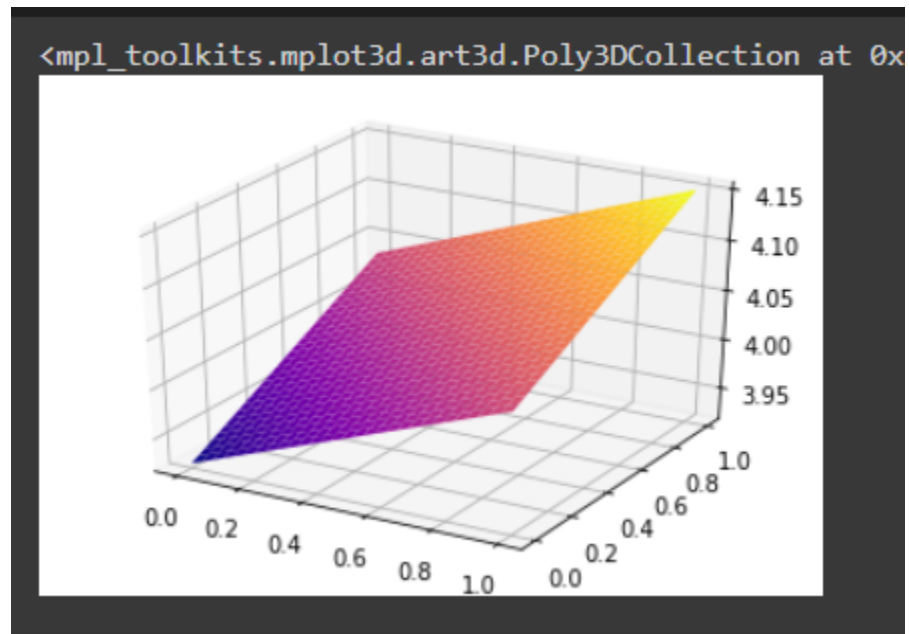
Similar to batch gradient descent, stochastic gradient descent also shows that the model overfits after degree 2.

iii) Surface plots for different degrees

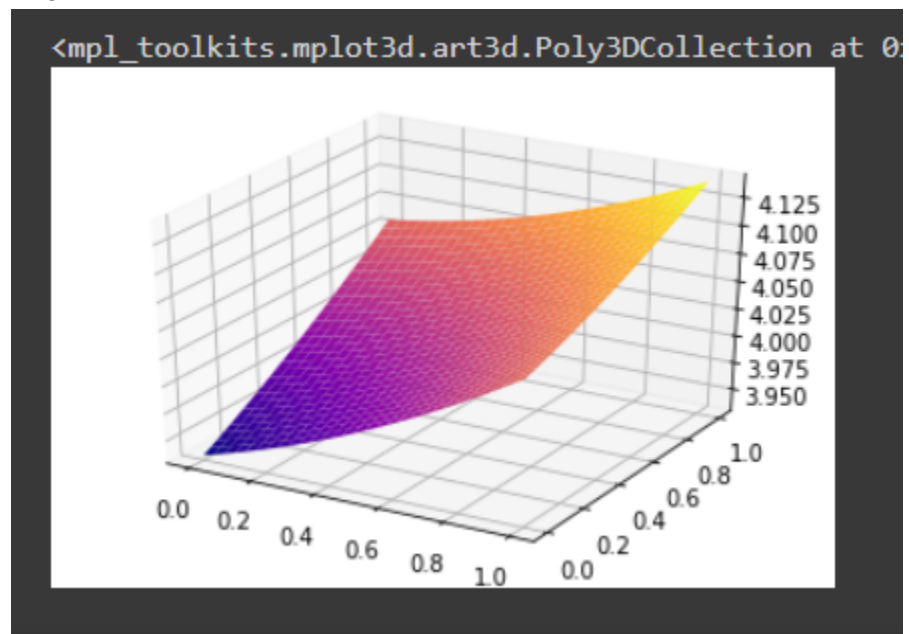
Degree 0:



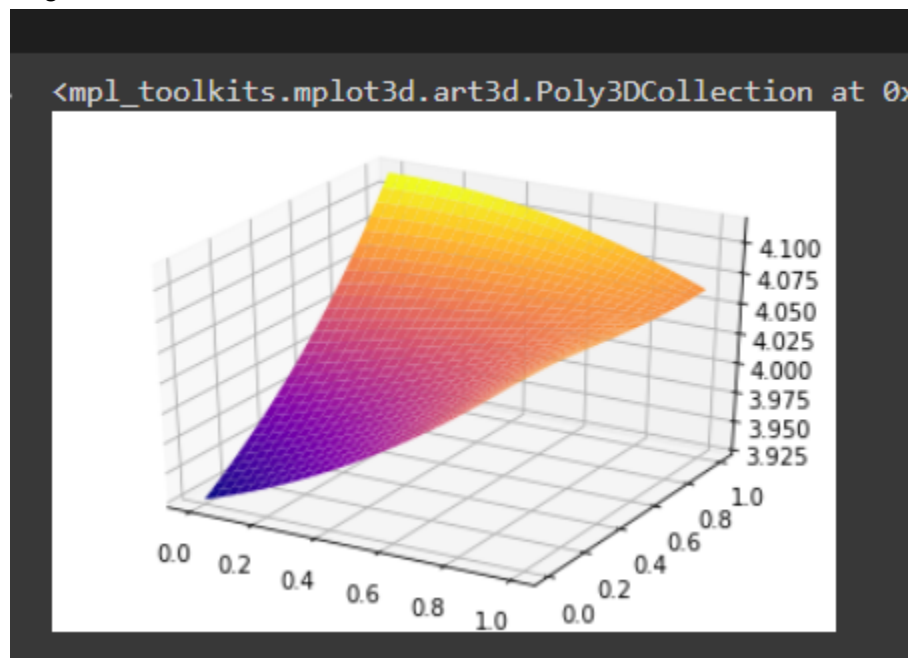
Degree 1:



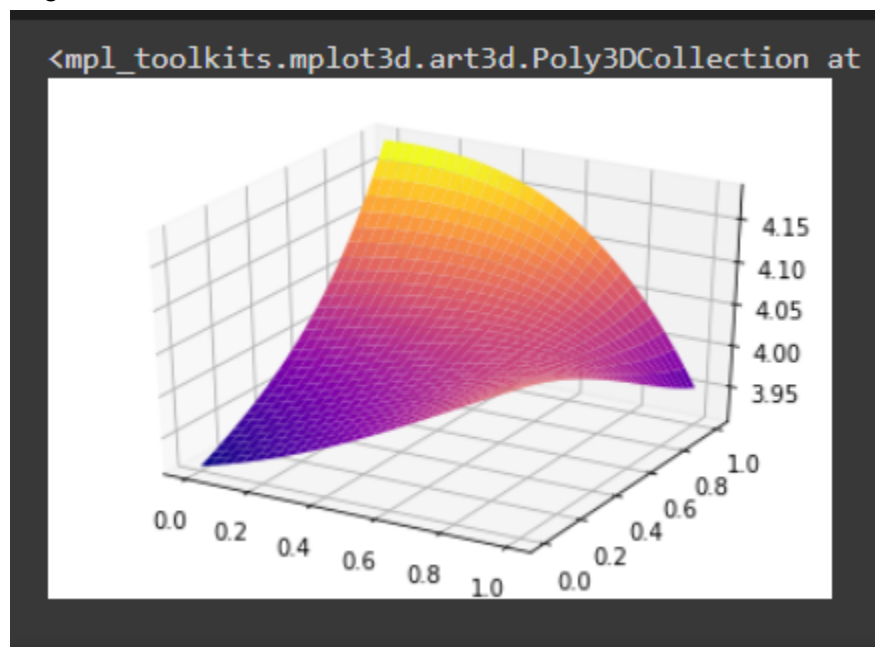
Degree 2:



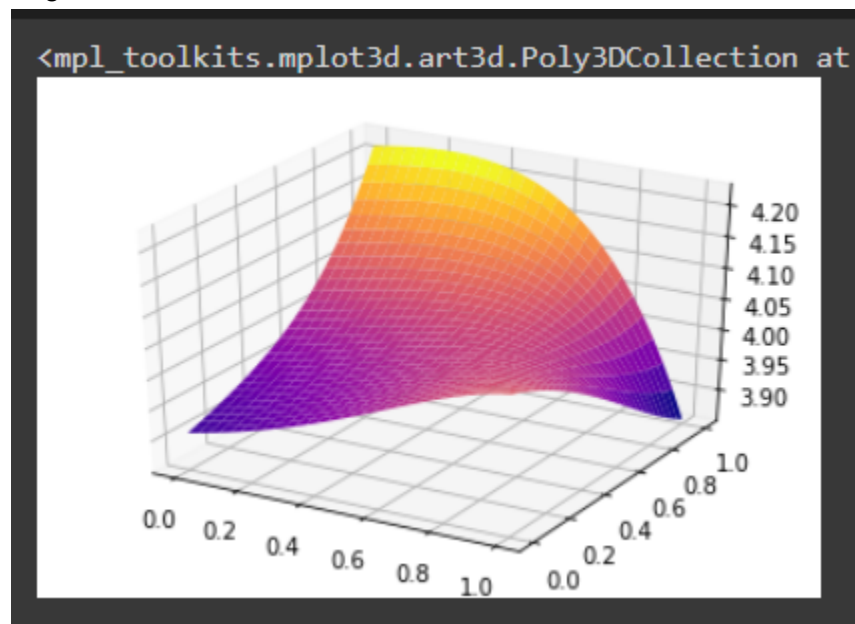
Degree 3:



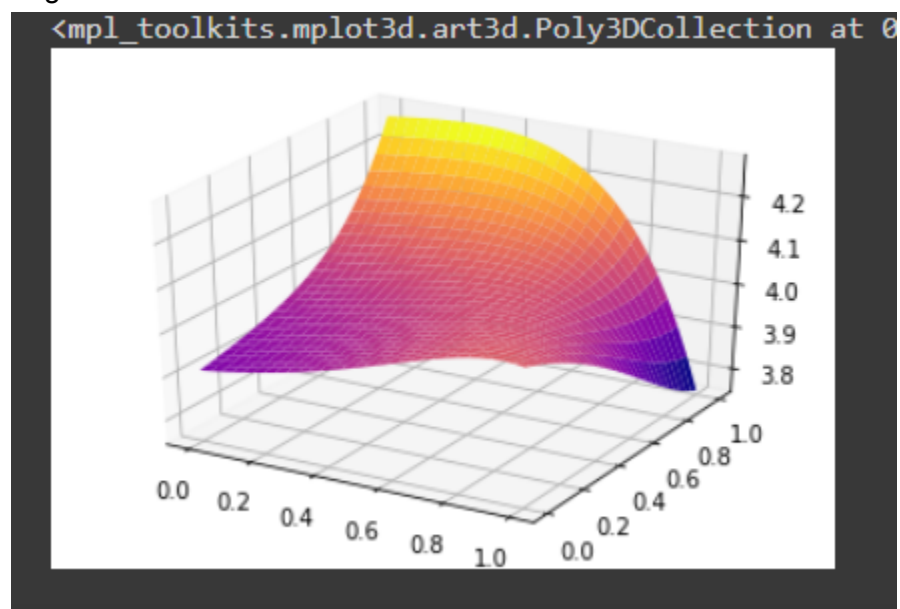
Degree 4:



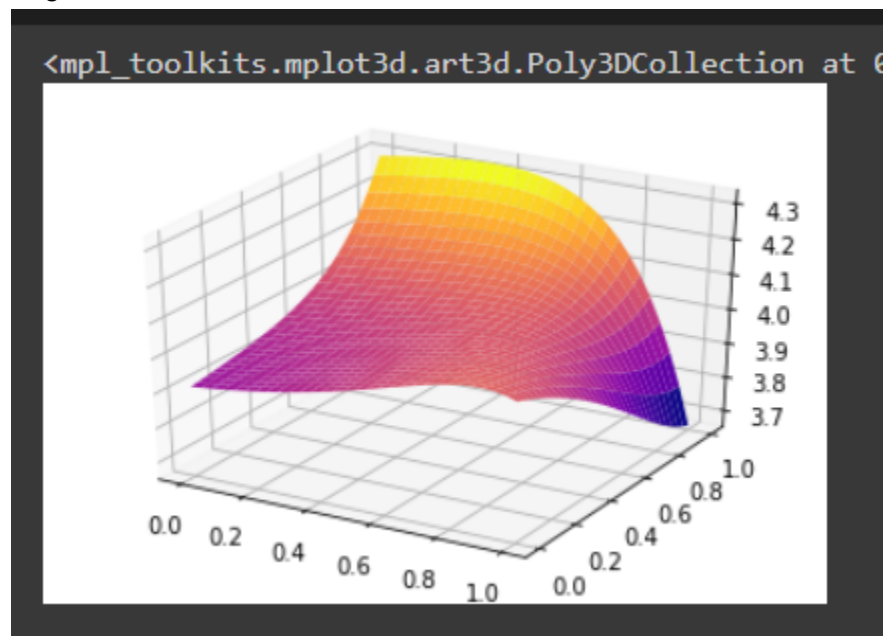
Degree 5:



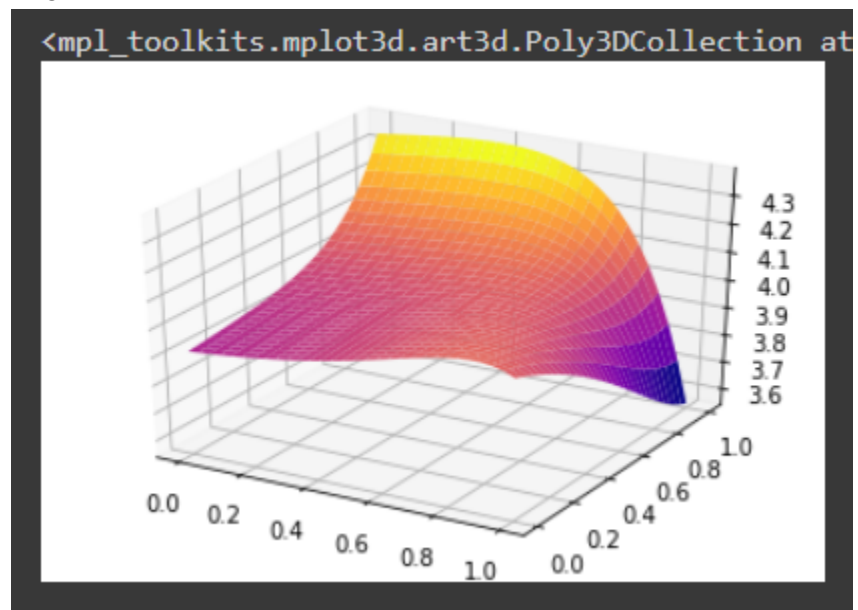
Degree 6:



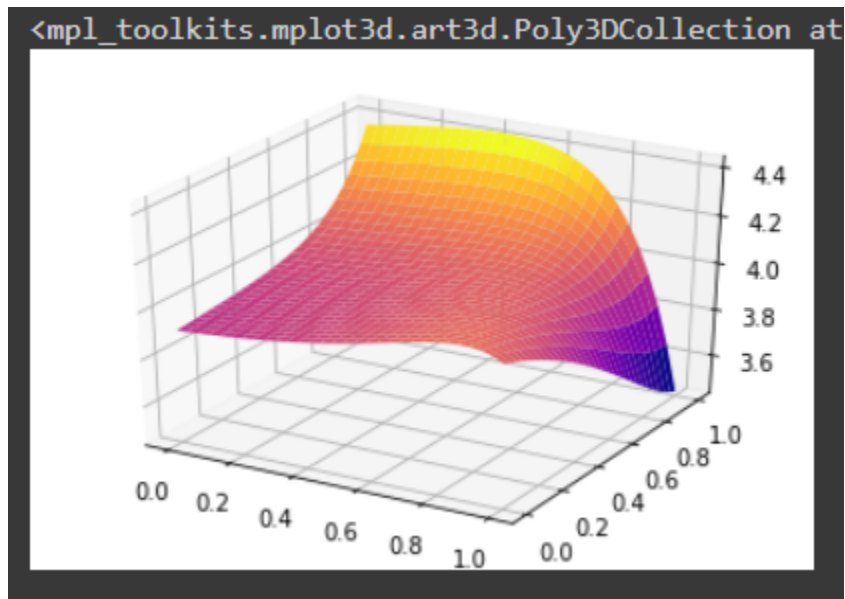
Degree 7:



Degree 8:



Degree 9:



iv) After using grid search to find the optimum regularization rate for each of the values of q . Given below are the regularization rates corresponding to each q and the minimum error rates.

$q = 0.5$, $\text{reg_const} = 10\text{e-}5$, $\text{loss} = 1.05514$

$q = 1$, $\text{reg_const} = 10\text{e-}4$, $\text{loss} = 1.05407$

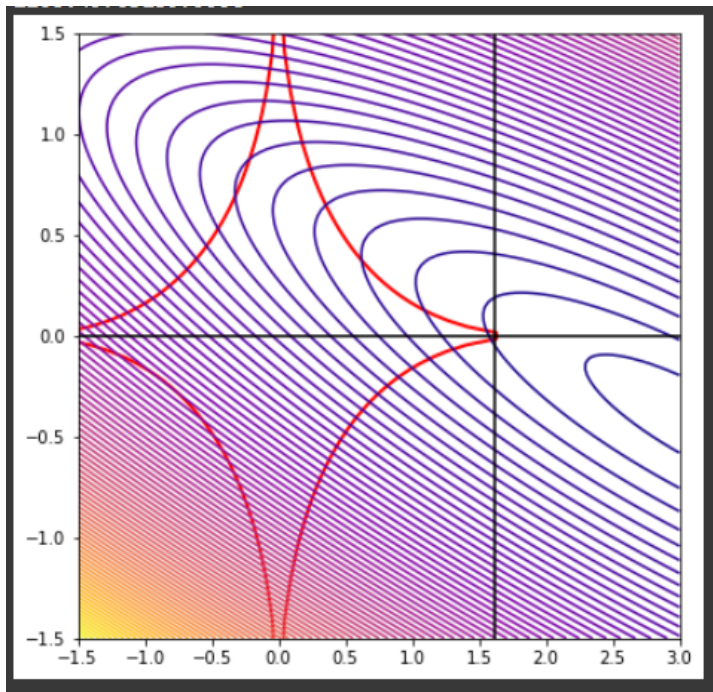
$q = 2$, $\text{reg_const} = 10\text{e-}5$, $\text{loss} = 1.04617$

$q = 4$, $\text{reg_const} = 10\text{e-}7$, $\text{loss} = 1.05677$

Best model with respect to testing error without regularization is degree 2 with error rate 1.5585

1C

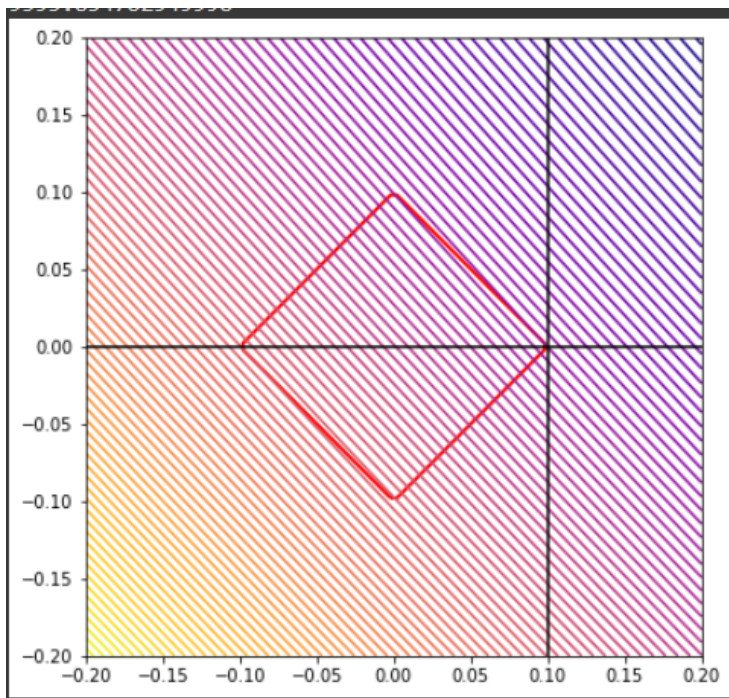
i) $q = 0.5$



Point of intersection - $w_1, w_2 = 0, 1.62$

Error at the point of intersection = 2108.5

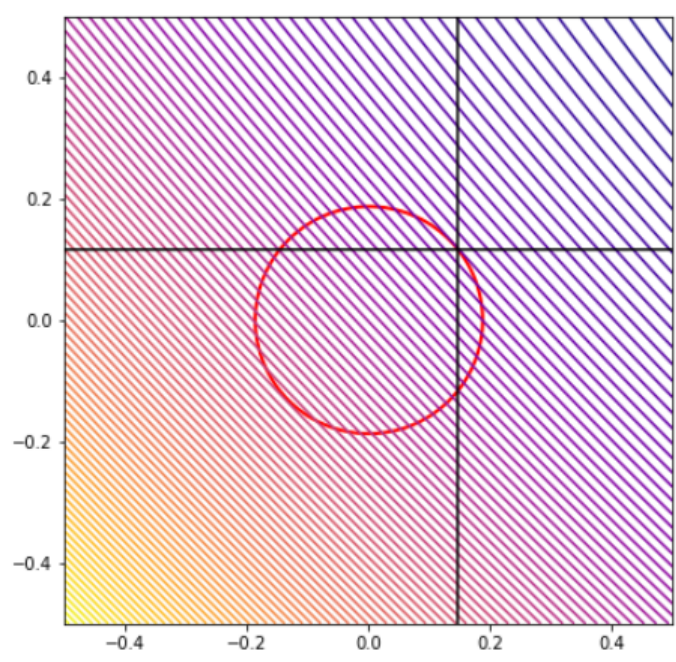
$q = 1$



Point of intersection - $w_1, w_2 = 0, 0.1$

Error at the point of intersection = 9595.85

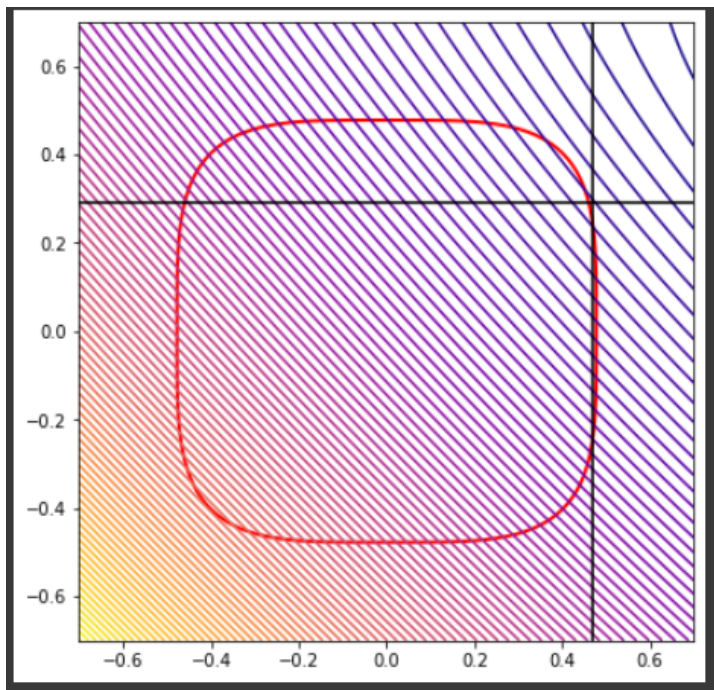
$q = 2$



Point of intersection - $w_1, w_2 = 0.115, 0.145$

Error at the point of intersection = 8236.35

$q = 4$



Point of intersection - $w_1, w_2 = 0.29, 0.47$

Error at the point of intersection = 5042.88