# Foundations of DataScience

## Assignment 1

| 1 | ANIRUDH A | 2018B4A70936H |
|---|---|---|
| 2 | VISHAL KUMAR NK | 2019A4PS0693H |
| 3 | SHIVANI THIRUNAGARI | 2019A4PS0754H |

# Table of Contents

## Objectives of the Assignment:
1. To implement Gradient-Descent algorithm and Stochastic Descent algorithm for polynomials of degree 0-9.
2. Implementation of regularisation using ridge regression and lasso regression for a polynomial of degree 9.
3. Visualization of surface plots for all the polynomial degrees.

## Data-Preprocessing:
For pre-processing, the data was shuffled initially, and then normalised using min-max normalisation. Then the data was split into training and testing portions of 70:30 proportion.

## Description of the Model:
Our method is about developing 10 different models of degrees from 0 to 9. Each model of degree d contains all polynomial combinations of the features with degree less than or equal to d as independent variables. The dataset contains 2 features, Strength and Temperature which are used to predict the pressure. The assumption of Strength as $x1$ and temperature as $x2$ has been considered. For a model of degree 1, we have three different terms in our feature matrix, for a model of degree 2, we have 6 different terms. The resulting model would be:

For a polynomial of degree 1: $\qquad y = w0 + w1x1 + w2x2$

For a polynomial of degree 2:
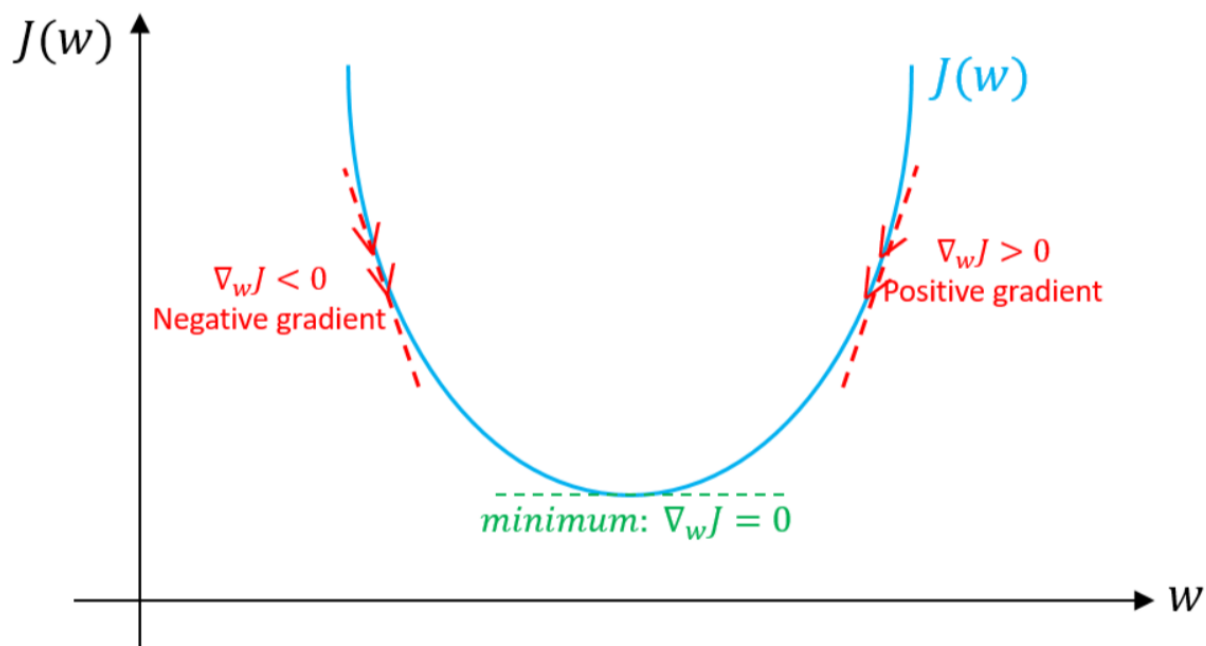$$y = w0 + w1x1 + w2x2 + w3x1^2 + w4x1x2 + w5x2^2$$

For a polynomial of degree n:
$$y = w0 + \Sigma\Sigma Wk(x1)^i(x2)^j \ where \ i + j <= \ n, \text{i,j belongs to (0,n)} ,$$

To generate the feature matrix, we have defined a function named featureVector with input parameters degree and the training data. This processes the featureVector taking the degree as input by using a nested-for loop using the principle that the sum of powers of the features x1, x2 will be less than or equal to that of the degree d. This helped us generate the featureVector as required for each degree of the polynomial.

## Gradient Descent Algorithm:

Gradient Descent is an iterative approach used in optimization to predict the values of the coefficients of the features in order to minimize the cost function. By taking steps of sizes proportional to the negative gradient, the gradient descent iteratively finds the optimal value. This in fact is the global minimum since a convex function is being considered.



For the implementation of this algorithm, we define a function, gradientDescent(deg,eta), where deg stands for degree of the polynomial, eta is the learning rate, which has a while loop that runs till there convergence of the error values. The error values are computed using the function compute_cost(). We initialised the W, cost from previous iteration, error from this iteration. At each iteration of the loop, we calculate the predicted y by taking the dot product of the featureVector with respect to the degree given and current model. We then calculate the derivative by finding the dot product of the featureVector and the difference between the predicted y and the actual. The equations are vectorised to save computation time. Then we update our model(W) by subtracting the derivative found and the learning rate in each iteration of the current model. This is in accordance with the gradient descent algorithm.

After every 50 iterations, we make a note of the cost corresponding to the 50th iteration, and print them. We utilised these arrays to plot the required graphs after every 50ths iteration.

In order to obtain the models of degrees 0-9, we added a for-loop which gives the input of degree deg to the above mentioned gradientDescent() function. It plots the graphs for no. of iterations v/s cost for degrees 0-9.
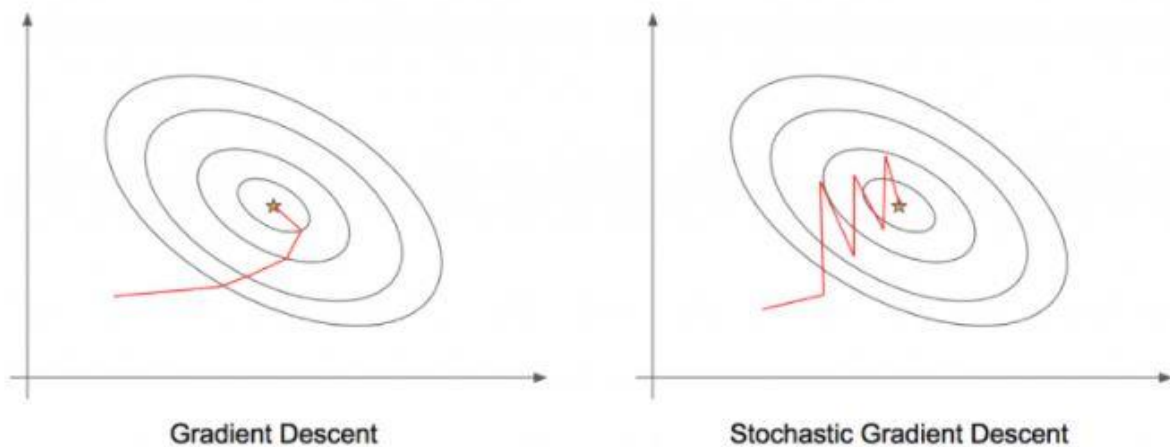
**Stochastic Gradient Descent Algorithm:**

As the name suggests, there is some stochasticity involved in this algorithm and is almost the same as how gradient descent works. Here, we take any random set of values in each iteration to evaluate the gradient and find the global optimum. This algorithm shows a faster convergence than the normal Gradient Descent algorithm, thus helping in saving time to train the data.

In our model, we have defined the function Stochastic_GD(). We set the number of iterations to 2000 and inside the while-loop, we generate a random value of "i" in the range (number of rows) which selects a training data sample randomly. Here, we use the derivation value based on only this one observation and calculate the cost of each iteration. Then at a step of 50 iterations, we print the error to obtain its minimum value and plot the required graphs.

In addition to this, for the ridge and lasso forms of the Stochastic Gradient Descent algorithm, we have calculated the RMSE separately for training and testing data and plotted it against log lambda.

Finally, we incorporated a loop that gives the models for polynomials from degrees 0 to 9, along with their training and testing errors, RMSE, and minimum errors.



Gradient Descent          Stochastic Gradient Descent

**Fig: Comparison of Gradient Descent and Stochastic Gradient Descent algorithms.**

**REGULARIZATION:**

Regularization is a technique to obtain different models of a set of data by changing the values of a parameter called "lambda".

## Ridge regularization:

In the ridge regularization technique, the cost function gets modified by adding a square of the absolute value of the coefficients whole multiplied by lambda. Different values of lambda give us different models. The cost function **without regularisation** is:

$$\sum_{i=1}^{M} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{M} \left( y_i - \sum_{j=0}^{p} w_j \times x_{ij} \right)^2$$

The **cost function with ridge regularisation** is as follows:

$$\sum_{i=1}^{M} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{M} \left( y_i - \sum_{j=0}^{p} w_j \times x_{ij} \right)^2 + \lambda \sum_{j=0}^{p} w_j^2$$

Thus, we have defined a new cost function called cost_function_ridge() that evaluates this cost and functions called gradientDescent_ridge() and Stochastic_GD_ridge() for the model. The regularization parameter lambda has been incorporated into the function to modify it accordingly (as $d/dx(lambda * (w ** 2) = 2 * lambda * w)$.
Again, after every 50 iterations, we print out the error values and plot the respective graphs. A plot of log($lambda$) vs RMSE has also been shown.

## Lasso Regularization:

When the term added to the error function is lambda times the sum of absolute values of the coefficients, the regularization is known as lasso regularization. We use different values of lambda to get different models. The cost function using lasso regularization is:

$$\sum_{i=1}^{M} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{M} \left( y_i - \sum_{j=0}^{p} w_j \times x_{ij} \right)^2 + \lambda \sum_{j=0}^{p} |w_j|$$

A function canned cost_function_lasso() has been defined to evaluate this cost and functions called gradientDescent_lasso() and Stochastic_GD_lasso() have been defined for the models.
The models have also been updated by adding lambda to the derivative(as $d/dx(lambda * abs(w)) = lambda)$.20 different values of lambda have been used to obtain the different models and RMSE for training and testing has been calculated

accordingly. At a step of 50 iterations, the relevant error values are printed and graphs have been plotted accordingly. A plot of log($lambda$) vs RMSE has also been shown.

## Errors : Without Regularisation

### (a) Gradient Descent

| Degree of the polynomial | Minimum Training Error | Minimum Testing Error |
|:---:|---|---|
| 0 | 66.123709285827 | 30.16854042476644 |
| 1 | 13.934476949748674 | 6.603623992313004 |
| 2 | 15.923729583766763 | 7.497449964102895 |
| 3 | 15.8131670343265 | 7.325443698852643 |
| 4 | 15.004249176538965 | 7.020358370894085 |
| 5 | 14.788130681679464 | 6.970498463852234 |
| 6 | 14.838305223907017 | 7.027752353628535 |
| 7 | 15.026731673159844 | 7.130305854345865 |
| 8 | 15.269975209673264 | 7.239002932405782 |
| 9 | 15.488068604155732 | 7.318209202503303 |

**Conclusion:** For polynomial regression using gradient descent, the minimum training and testing error can be observed for a polynomial of degree 1. Hence, a linear polynomial best fits the model considering a learning rate of 0.1.

### (b) Stochastic Gradient Descent

| Degree of the polynomial | Minimum Training Error | Minimum Testing Error |
|---|---|---|

| | | |
|---|---|---|
| 0 | 172.49577394231972 | 68.87068590149865 |
| 1 | 505.5911585695828 | 202.02435222053836 |
| 2 | 1014.2494176270469 | 405.27239778460444 |
| 3 | 1744.750570541063 | 696.8266506430223 |
| 4 | 2656.345692944137 | 1060.8236235596203 |
| 5 | 3840.921899039334 | 1533.536072399309 |
| 6 | 5209.762445168131 | 2080.200727712149 |
| 7 | 6821.695159972138 | 2723.8676409864975 |
| 8 | 8617.624574971656 | 3440.208864367946 |
| 9 | 10719.874652655957 | 4280.705715149436 |

**Conclusion:** ==For polynomial regression using stochastic gradient descent, the minimum training and testing error can be observed for a polynomial of degree 0. However, this is a case of underfitting the model and hence regularization is required.==

Hence, minimum training and testing errors have been tabulated for polynomials of degree 0-9 using gradient descent and stochastic gradient descent.

## Errors : With Regularization

*In this section, we tabulate the minimum training and testing errors along with the RMSE values for ==20 different values of lambda== taken from the range 0 to 0.1. A polynomial of degree 9 has been considered.*
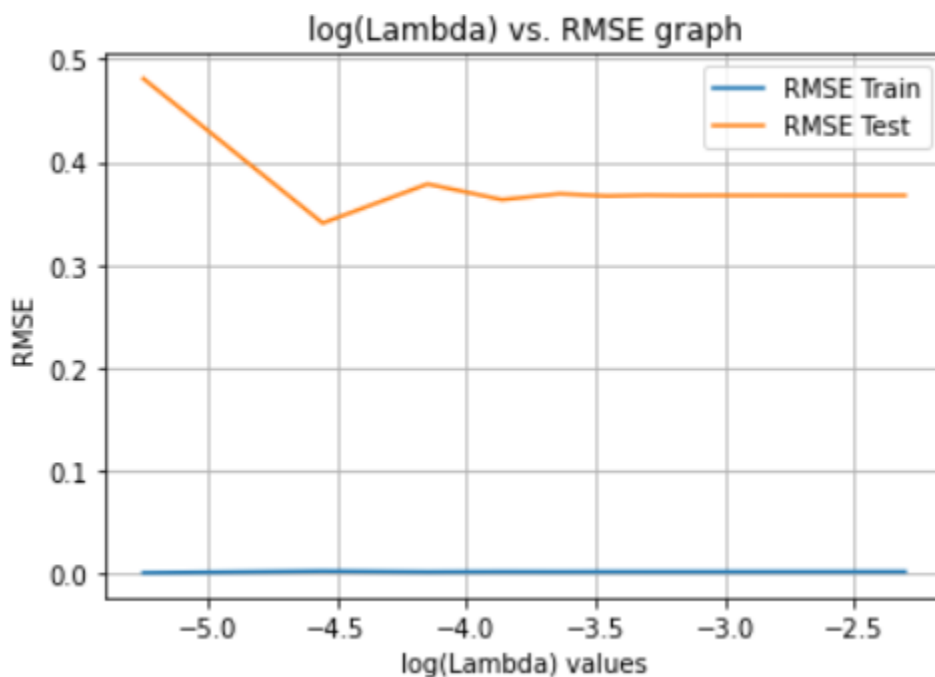
### (a) Gradient Descent with Ridge Regularisation

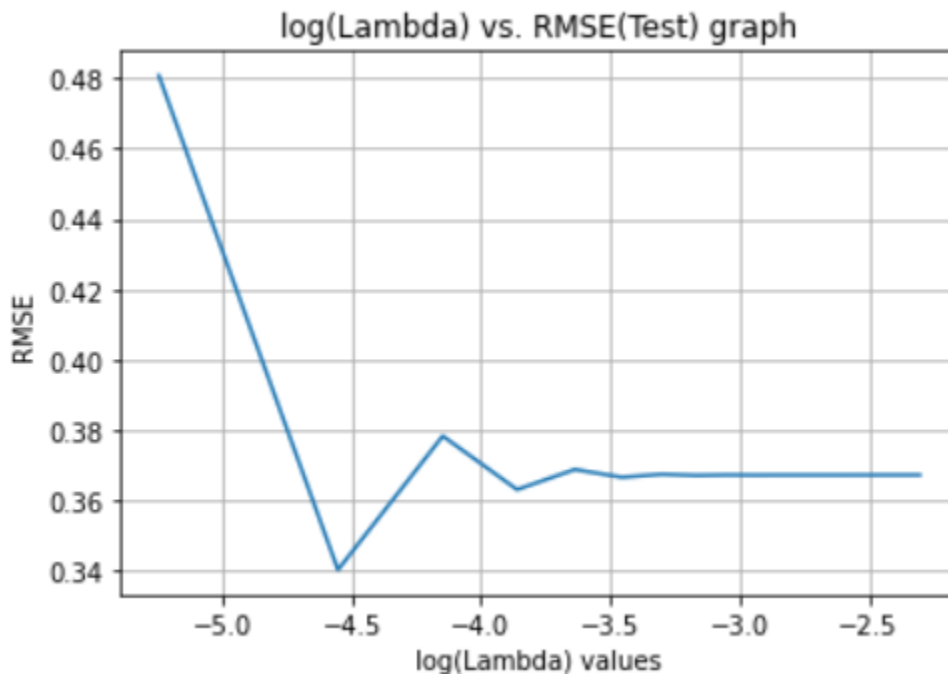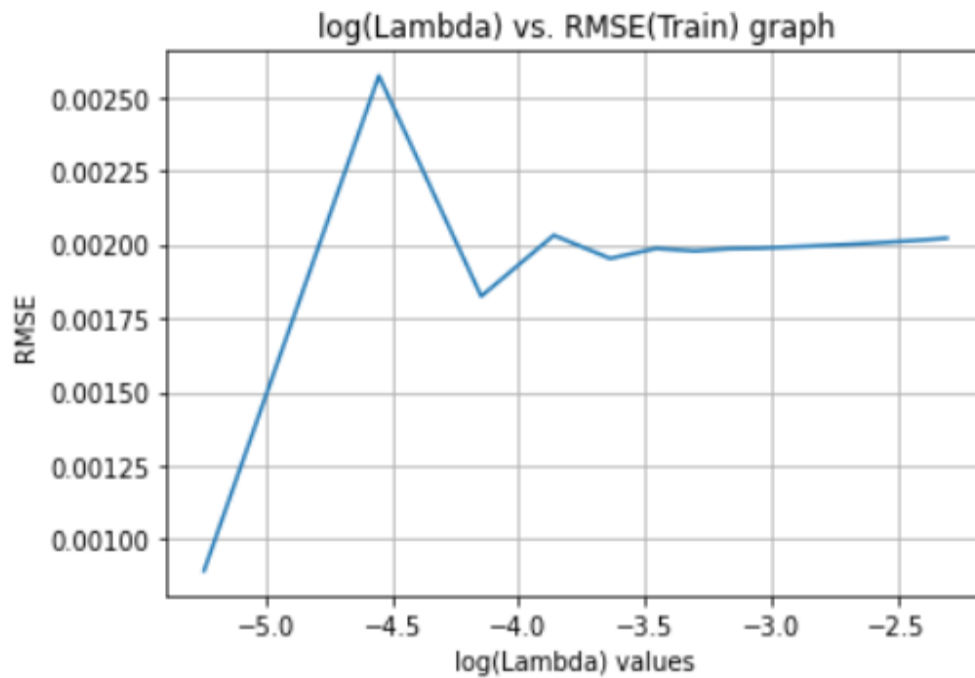| Penalty (Lambda) | RMSE training | RMSE testing | Minimum Training Error | Minimum Testing Error |
|---|---|---|---|---|

| | | | | |
|---|---|---|---|---|
| 0. | 0.0029598229 41240777 | 0.3576048186 5054716 | 0.2956569659 605091 | 0.1260036542 3793083 |
| 0.00526316 | 0.0019093735 02703186 | 0.3684403634 332047 | 0.0001506000 1484298665 | 1.0401124771 370023 |
| 0.01052632 | 0.0019703225 973550977 | 0.3650944312 9325453 | 0.0024166652 773281398 | 0.5237960626 876018 |
| 0.01578947 | 0.0019555148 08122437 | 0.3661007692 4320234 | 0.0018000993 817877764 | 0.6460885226 40739 |
| 0.02105263 | 0.0019722150 063530485 | 0.3657992721 107745 | 0.0029511220 958116867 | 0.5964863922 354798 |
| 0.02631579 | 0.0019655294 92154176 | 0.3658934322 3881976 | 0.0033878592 948978726 | 0.6156981691 8168 |
| 0.03157895 | 0.0019691840 98542376, | 0.3658681359 30915 | 0.0041832490 17209606 | 0.6092338384 659439 |
| 0.03684211 | 0.0019722150 063530485 | 0.3658789246 4305294 | 0.0048117521 48001734 | 0.6126445409 642288 |
| 0.04210526 | 0.0019754587 93790997 | 0.3658788130 8386744 | 0.0055133328 46471726 | 0.6122990868 465671 |
| 0.04736842 | 0.0019786621 926774876 | 0.3658819921 5254397 | 0.0061805865 869152595 | 0.6133883123 980713 |
| 0.05263158 | 0.0019819016 686619404 | 0.3658841757 4806503 | 0.0068607134 51772674 | 0.6139302338 734418 |
| 0.05789474 | 0.0019851539 377191297 | 0.3658866583 597918 | 0.0075335430 48285398 | 0.6146808109 633879 |
| 0.06315789 | 0.0019884258 92378603 | 0.3658890490 3751834 | 0.0081905255 76160615 | 0.6153516029 894314 |
| 0.06842105 | 0.0019917152 96886712 | 0.3658914658 514681 | 0.0088379160 37052584 | 0.6160526456 461533 |
| 0.07368421 | 0.0019950226 772641414 | 0.3658938731 467393 | 0.0094835269 78979806 | 0.6167419612 337254 |
| 0.07894737 | 0.0019983477 24542667 | 0.3658962816 9242516 | 0.0101268976 45738252 | 0.6174355633 067768 |
| 0.08421053 | 0.0020016903 82781046 | 0.3658986882 382593 | 0.0107684029 671625 | 0.6181273435 303892 |
| 0.08947368 | 0.0020050505 199608873 | 0.3659010937 674944 | 0.0114080291 126534 | 0.6188196321 533176 |

| 0.09473684 | 0.0020084280 27543965 | 0.3659034979 84992 | 0.0120458945 47349233 | 0.6195115405 700207 |
|---|---|---|---|---|
| 0.1 | 0.0020118227 903673353 | 0.3659059009 816345 | 0.0126820463 541897 | 0.6202034079 488877 |

**Conclusion:** <mark>For ridge regression using gradient descent, the minimum RMSE training error can be observed for a lambda value of 0.00526316. Hence, among the 20 chosen values, this best fits the model considering a learning rate of 0.1.</mark>
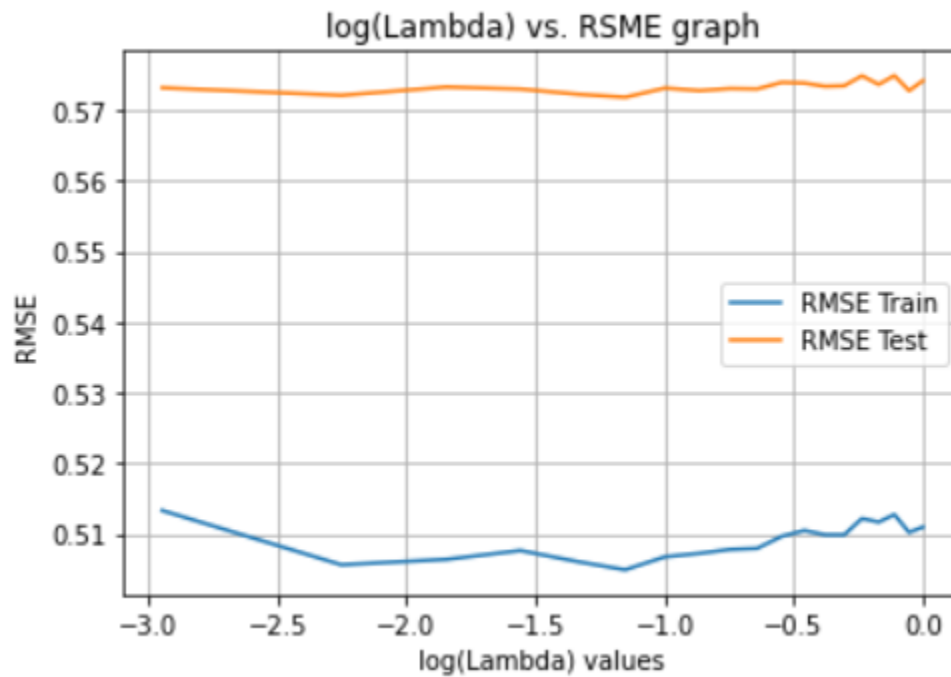
log(Lambda) vs. RMSE(Train) graph



log(Lambda) vs. RMSE(Test) graph

**(b) Stochastic Gradient Descent with Ridge Regularisation**

| Penalty (Lambda) | RMSE training | RMSE testing | Minimum Training Error | Minimum Testing Error |
|---|---|---|---|---|
| | | | | |

| | | | | |
|---|---|---|---|---|
| 0. | 0.3232768691 952058 | 0.5801618925 605911 | 60.352782907 21532 | 83.305485840 92301 |
| 0.05263158 | 0.5133023860 872065 | 0.5732109503 410885 | 152.18456478 337836 | 81.351462000 13483 |
| 0.10526316 | 0.5056077355 131638 | 0.5721498239 61949 | 147.68611240 596823 | 81.077859416 21476 |
| 0.15789474 | 0.5063686973 676841 | 0.5733167911 728955 | 148.16170303 757212 | 81.437358078 46046 |
| 0.21052632 | 0.5076226090 082823 | 0.5730342057 181997 | 148.91832215 845494 | 81.380118620 204 |
| 0.26315789 | 0.5060418874 375047 | 0.5722979712 831239 | 148.02392630 097566 | 81.201084576 47581 |
| 0.31578947 | 0.5049006014 115437 | 0.5718727135 462962 | 147.38206377 600548 | 81.105101405 39554 |
| 0.36842105 | 0.5067399172 613061 | 0.5731709484 527802 | 148.47485513 676614 | 81.495249197 31253 |
| 0.42105263 | 0.5072243742 868127 | 0.5728080092 852189 | 148.79343827 509572 | 81.423202821 36674 |
| 0.47368421 | 0.5077792145 670067 | 0.5731058607 398886 | 149.13266978 55303 | 81.524279577 21714 |
| 0.52631579 | 0.5079321256 243928 | 0.5730464811 495236 | 149.25196862 393378 | 81.534442290 1432 |
| 0.57894737 | 0.5096470545 723275 | 0.5739806918 233357 | 150.27569721 241758 | 81.815601837 17255 |
| 0.63157895 | 0.5104660022 976055 | 0.5738804368 60377 | 150.77903553 550217 | 81.810097318 29869 |
| 0.68421053 | 0.5098592990 69625 | 0.5734221587 205475 | 150.42945063 15116 | 81.698724831 54745 |
| 0.73684211 | 0.5098501851 38831 | 0.5734914730 858356 | 150.44701769 433834 | 81.731005914 86116 |
| 0.78947368 | 0.5121651060 743913 | 0.5749189128 282666 | 151.82301822 626334 | 82.156891179 9266 |
| 0.84210526 | 0.5116635764 025407 | 0.5736843882 813348 | 151.54761883 86171 | 81.830023770 77828 |
| 0.89473684 | 0.5127273896 03675 | 0.5749415800 875519 | 152.16667018 61574 | 82.195113419 56066 |

| 0.94736842 | 0.5102028772773332 | 0.5727993065560243 | 150.75497420038852 | 81.63220935453127 |
| 1. | 0.5109688868212708 | 0.574237567728231 | 151.18856307683987 | 82.03179798151098 |

**Conclusion:** For ridge regression using stochastic gradient descent, the minimum RMSE training error can be observed for a lambda value of 0.31578947. Hence, among the 20 chosen values, this best fits the model considering a learning rate of 0.1.
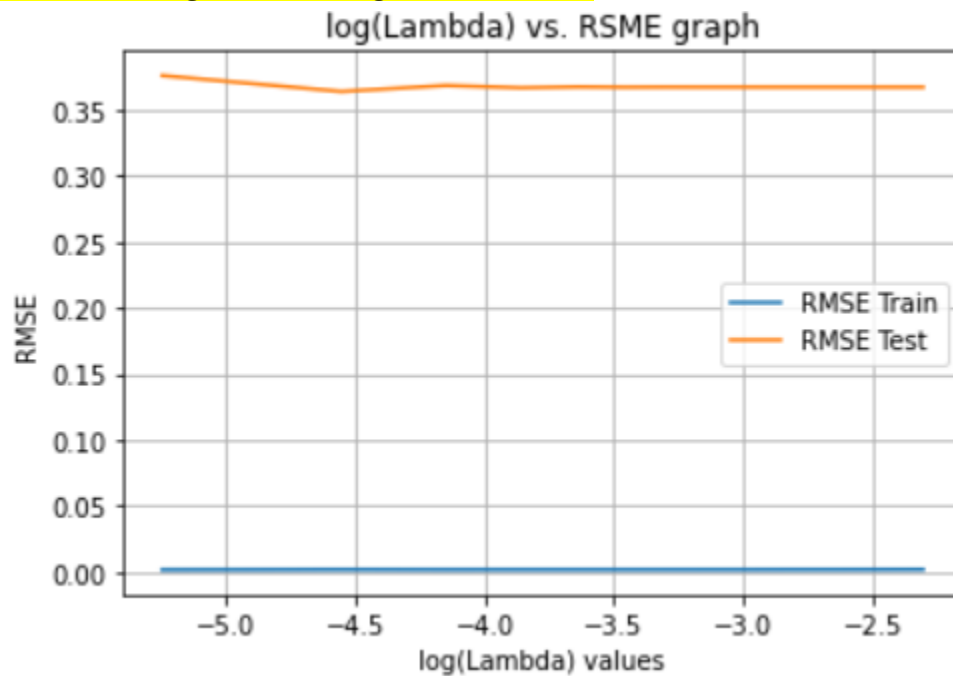
log(Lambda) vs. RSME(Train) graph



log(Lambda) vs. RSME(Test) graph

**(c) Gradient Descent with Lasso Regularisation**

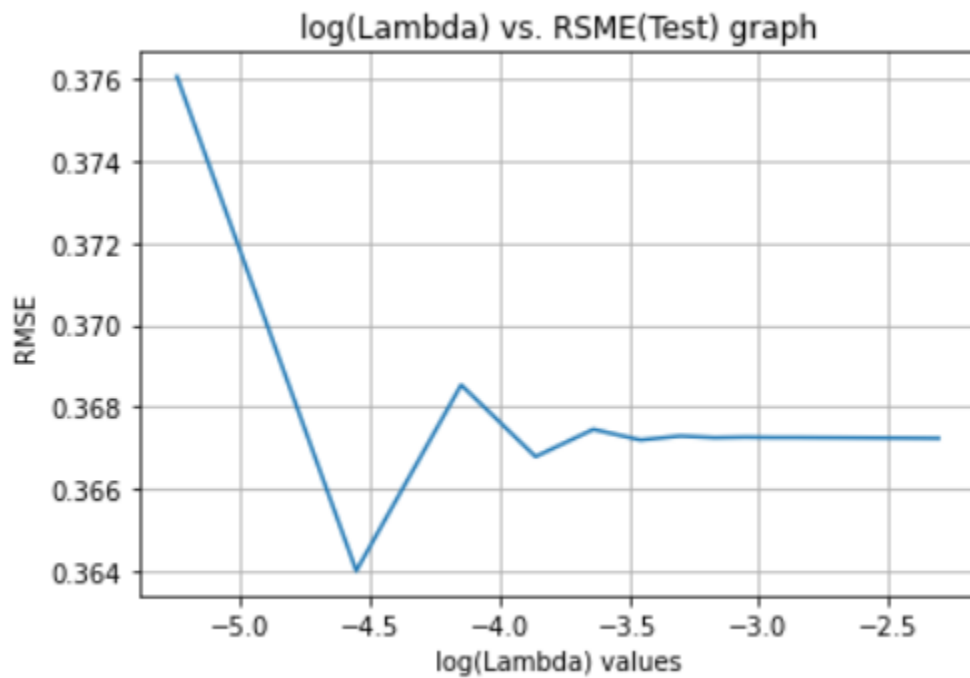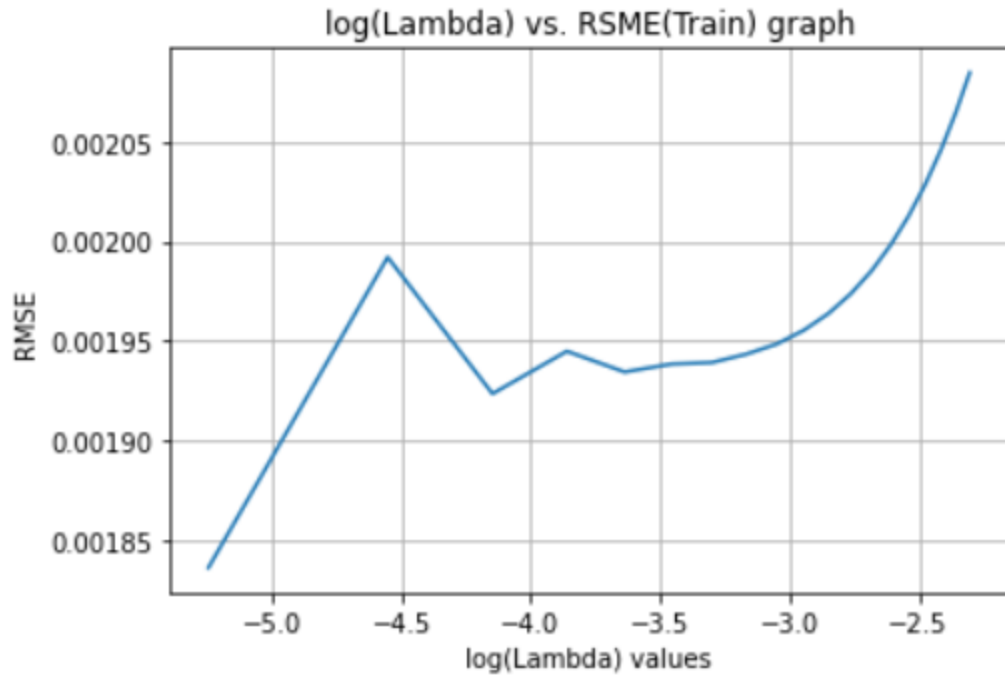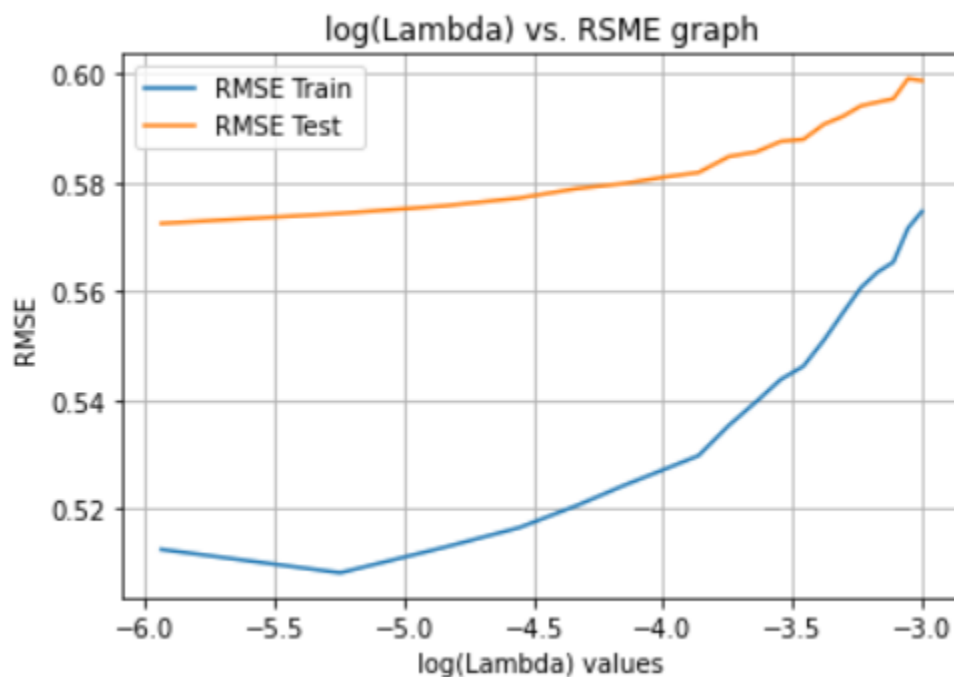| Penalty (Lambda) | RMSE training | RMSE testing | Minimum Training | Minimum Testing Error |
|---|---|---|---|---|
| | | | | |

| | | | Error | |
|---|---|---|---|---|
| 0. | 0.0030626439 77942353 | 0.3457596849 1936665 | 9.8487775424 07888e-05 | 0.5379739187 199286 |
| 0.00526316 | 0.0018366181 083689377 | 0.3760434190 2425015 | 0.0025832194 81021006 | 0.6388867396 987381 |
| 0.01052632 | 0.0019920770 49222626 | 0.3640118303 4140117 | 0.0055670366 465625405 | 0.6017961255 796137 |
| 0.01578947 | 0.0019237002 315991177 | 0.3685367756 8395765 | 0.0080330520 45218604 | 0.6191812931 499923 |
| 0.02105263 | 0.0019449803 282015718 | 0.3667907306 991781 | 0.0107985341 47236865 | 0.6161682937 58994 |
| 0.02631579 | 0.0019345989 732409054 | 0.3674544695 968518 | 0.0133809508 86727451 | 0.6209441953 384277 |
| 0.03157895 | 0.0019385686 68056153 | 0.3671969170 4425356 | 0.0160280599 41361732 | 0.6227396919 229335 |
| 0.03684211 | 0.0019392957 886821467 | 0.3672923418 2214276 | 0.0186202494 8476774 | 0.6256472499 945015 |
| 0.04210526 | 0.0019434752 073149173 | 0.3672526395 8019383 | 0.0212155047 09511602 | 0.6281111009 565452 |
| 0.04736842 | 0.0019485449 437108041 | 0.3672647019 5362226 | 0.0238135837 76360534 | 0.6307488429 435586 |
| 0.05263158 | 0.0019554636 527468155 | 0.3672569512 785462 | 0.0264230611 63295913 | 0.6333324180 441312 |
| 0.05789474 | 0.0019638452 919350886 | 0.3672568011 6480713 | 0.0290459898 1339794 | 0.6359545055 940562 |
| 0.06315789 | 0.0019738102 829525445 | 0.3672537514 0648203 | 0.0316667882 6007826 | 0.6385648116 758342 |
| 0.06842105 | 0.0019852820 636579284 | 0.3672518241 3378665 | 0.0342943133 7952899 | 0.6411854897 415419 |
| 0.07368421 | 0.0019982546 56945496 | 0.3672494785 423012 | 0.0369472898 1596327 | 0.6438301707 915508 |
| 0.07894737 | 0.0020126914 79534199 | 0.3672473049 267678 | 0.0396157689 3489343 | 0.6464908575 933642 |
| 0.08421053 | 0.0020285641 312360234 | 0.3672450771 2345367 | 0.0422885696 0197177 | 0.6491556213 627605 |

| | | | | |
|---|---|---|---|---|
| 0.08947368 | 0.0020458381 097653096 | 0.3672428817 895564 | 0.0449863136 2023679 | 0.6518453703 706788 |
| 0.09473684 | 0.0020644786 482884895 | 0.3672406857 1834345 | 0.0476914389 7295541 | 0.6545424328 269744 |
| 0.1 | 0.0020844489 29944153 | 0.3672385016 347565 | 0.0504033688 5151298 | 0.6572462739 876614 |

**Conclusion:** ==For lasso regression using gradient descent, the minimum RMSE training error can be observed for a lambda value of 0.00526316. Hence, among the 20 chosen values, this best fits the model considering a learning rate of 0.1.==

log(Lambda) vs. RSME(Train) graph



log(Lambda) vs. RSME(Test) graph

**(d) Stochastic Gradient with Lasso Regularisation**

| Penalty (Lambda) | RMSE training | RMSE testing | Minimum Training Error | Minimum Testing Error |
|---|---|---|---|---|
|  |  |  |  |  |

| | | | | |
|---|---|---|---|---|
| 0. | 0.32355923926517616 | 0.5803662019220462 | 60.45881070875387 | 83.36416976252178 |
| 0.00263158 | 0.5126226503275542 | 0.5724393996735007 | 151.8350520372166 | 81.18733491462798 |
| 0.00526316 | 0.5083369546323548 | 0.5742787645689261 | 149.38749288898376 | 81.78635340962273 |
| 0.00789474 | 0.5130456540204075 | 0.5756781051740151 | 152.23726567062573 | 82.25292326729101 |
| 0.01052632 | 0.5166281490541821 | 0.5771072609331792 | 154.4359340606776 | 82.72906760137305 |
| 0.01315789 | 0.5206975855628655 | 0.5788364899324561 | 156.933840979512 | 83.2858796886147 |
| 0.01578947 | 0.5243869239295973 | 0.5797738050902412 | 159.21166859157248 | 83.60879360062474 |
| 0.01842105 | 0.5272919615968767 | 0.5809425061700967 | 161.0264283631437 | 83.99813549870642 |
| 0.02105263 | 0.5298319976487207 | 0.5817851714035112 | 162.6385239933143 | 84.2970967834184 |
| 0.02368421 | 0.5353958071816015 | 0.5847247445283491 | 166.08311806380135 | 85.17359067217849 |
| 0.02631579 | 0.5397340608160104 | 0.5855632400812566 | 168.8270881988225 | 85.46660289887699 |
| 0.02894737 | 0.5438205394533024 | 0.5875041399208681 | 171.33562468079967 | 86.05712128988652 |
| 0.03157895 | 0.5462263556633338 | 0.5878547426013558 | 172.96920954330903 | 86.22075068071929 |
| 0.03421053 | 0.5510436297913769 | 0.5906514322700394 | 176.03542972738667 | 87.05243890752412 |
| 0.03684211 | 0.5560987053700288 | 0.5921169035045885 | 179.28252203953446 | 87.52083901744456 |
| 0.03947368 | 0.5606629329049276 | 0.5940231556331397 | 182.27418284211404 | 88.09938498806912 |
| 0.04210526 | 0.5635170907325913 | 0.5947093607485944 | 184.0031004471282 | 88.37200078283331 |
| 0.04473684 | 0.5653364922766713 | 0.5952999646097221 | 185.4176176163358 | 88.59131481008075 |

| 0.04736842 | 0.5716043803 44485 | 0.5990165218 52475 | 189.57467840 04326 | 89.695344474 11218 |
|---|---|---|---|---|
| 0.05 | 0.5746684597 675228 | 0.5986349064 360172 | 191.53741533 107453 | 89.648228340 73384 |

**Conclusion:** <mark>For lasso regression using stochastic gradient descent, the minimum RMSE training error can be observed for a lambda value of close to 0. Hence, among the 20 chosen values, this best fits the model considering a learning rate of 0.1.</mark>
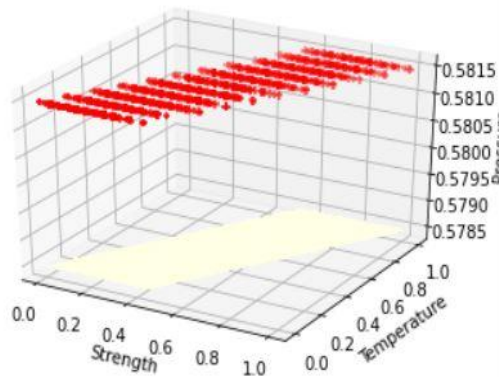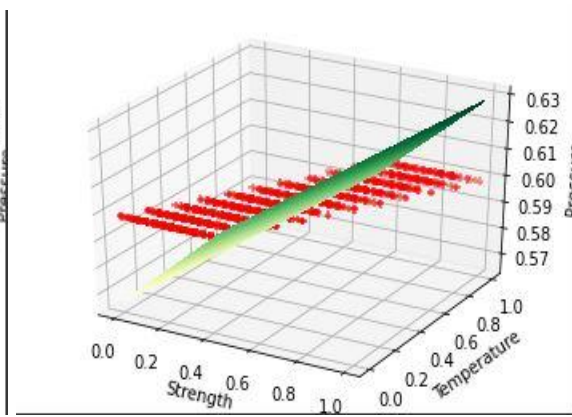
log(Lambda) vs. RSME(Train) graph



log(Lambda) vs. RSME(Test) graph

# Surface Plots

In this section, surface plots have been drawn for polynomial regression with degrees 0-9 based on the gradient descent algorithm.
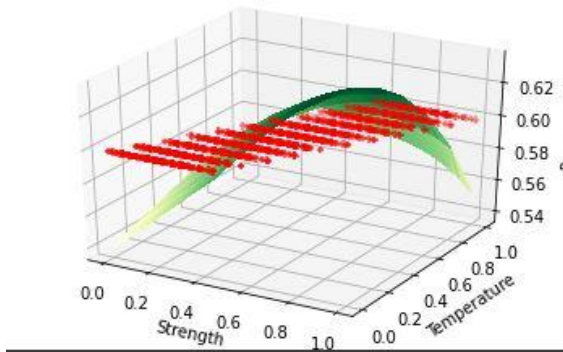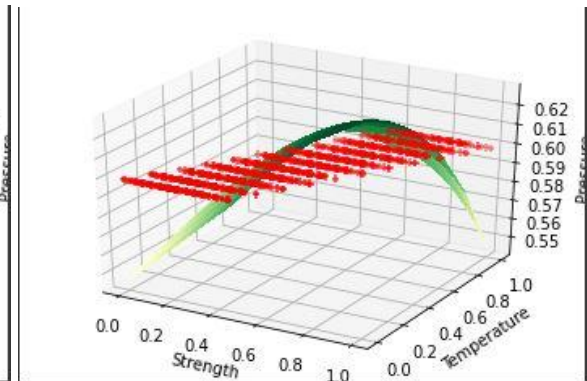
**Degree 0**
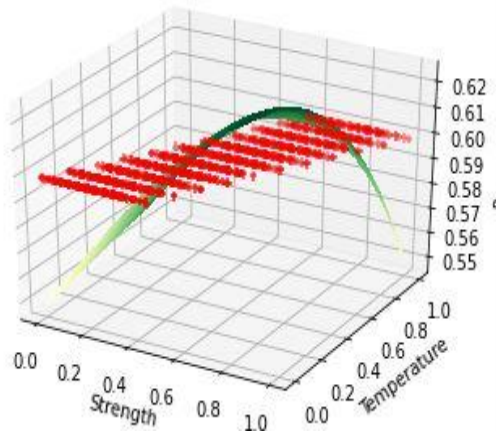


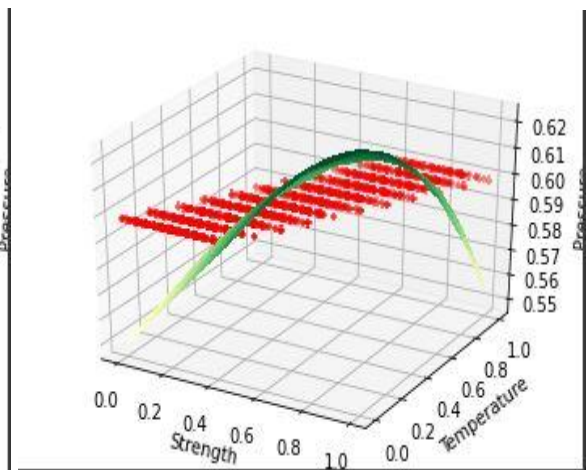**Degree 1**



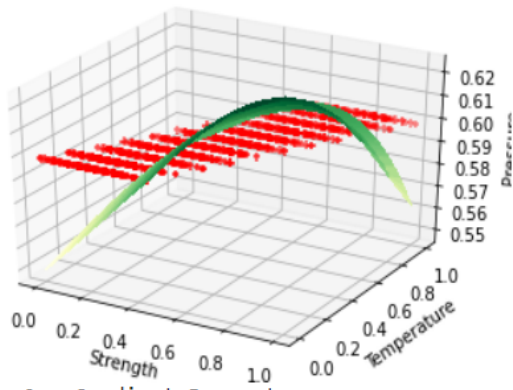**Degree 2**
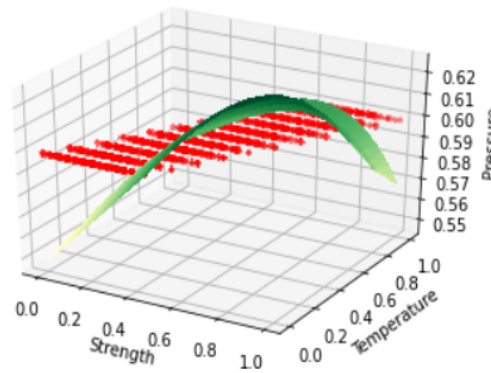


**Degree 3**

**Degree 4**

**Degree 5**

**Degree 6**

**Degree 7**

**Degree 8**

**Degree 9**

**How does overfitting work?**

For the model described above, we can clearly see that the model covers more data points as the number of degrees of the polynomial increases. For the zeroth-degree polynomial, the model looks flat but as the number of the degrees increases, the model tries to fit more data points in order to reduce the training error but this results in an increase in testing error. This is called overfitting and mainly happens when we have a large number of features but less number of data points.

**Compare between the best model obtained in part a) and the best model obtained in part b).**

- **Gradient Descent: -**

  Without regularization, the testing error values for gradient descent algorithm are in the range of 6.5 to 7.5. However, after regularization, the testing error values (in both ridge and lasso regressions) can be seen to be in the range of 0.5-0.7. Hence, the model based on regularization proves to have a better predictive ability.

- **Stochastic Gradient Descent: -**

  Without regularization, the testing error values for stochastic gradient descent algorithm are above 1000. However, after regularization, the testing error values (in both ridge and lasso regressions) can be seen to be around 80. Hence, the model based on regularization proves to have a better predictive ability.

# Conclusion: -

The models built using regularization better predict the values of pressure given the values of strength and temperature applied to a certain piece of plastic. Regularization decreases the training and testing error, thereby building a more reliable model.