

Machine Learning Assignment 1

Sayan Sinha
16CS10048

INTRODUCTION

Supervised machine Learning generally boils down to approximating a continuous function. One method of doing so is by making use of Linear Regression. In any method whatsoever, there always remains a set of hyperparameters which need to be tuned intuitively. In this assignment, we experiment with various values of such hyperparameters and try to draw a conclusion out of them.

EXPERIMENTATION

Part 1

A synthetic dataset is generated and is divided into train and test in a 4:1 ratio randomly. This is done by permuting the data in a random fashion and choosing the first 80% of the dataset to be the training set and the rest to be the test set. This part is kept as a separate function to make sure new data is generated only when required, and not all the time before training. The data is saved in the present working directory.

The dataset is then trained using Linear regression. The weights are initialised to random values between zero and one. At the end of each iteration, the loss up to that iteration is sent to STDOUT. Gradient descent is performed as per the following rule:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Here θ is the weight vector, with θ_j marking each element of the vector. α is the learning rate, m is the number of training examples and h_{θ} is the hypothesis function.

The training function returns the training error, testing error and the RMSE error. For the error, we use squared loss, i.e.:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (W^T \Phi_n(x) - y)^2$$

The training takes place separately for degrees 1 to 9. The degrees denote the polynomial degree of the function we are trying to approximate. Since bias is also introduced, there are a total of $degree + 1$ parameters to be tuned while training. At the end of a certain number of iterations, these values are saved in the local directory.

The parameter values learned have been provided in Table 1. The first row denotes the degree.

1	2	3	4	5	6	7	8	9
0.797531 9619	0.823028 0506	0.818393 5063	1.166862 8	1.442892 006	1.594227 469	1.293419 923	1.216725 803	1.056935 601
-1.441639 135	-1.628034 357	-1.265472 884	-1.637685 999	-1.876832 283	-2.252152 497	-1.266530 609	-1.278813 704	-0.940861 88
	0.197333 2989	-1.059530 691	-1.869621 587	-2.140028 412	-1.565307 6	-1.752213 111	-1.222106 24	-1.410866 293
		0.969801 3993	0.454302 5178	-0.177134 7154	-0.651195 9638	-1.541795 004	-1.481039 603	-0.929530 0843
			1.535320 926	0.413476 0332	-0.039120 77531	-0.036397 05869	-0.428483 5354	-0.410347 7217
				2.233017 764	1.128940 674	0.862934 4587	0.357140 9422	-0.016717 29166
					1.847131 128	0.816115 8519	0.716485 5357	0.103302 1222
						1.831891 609	0.845241 2498	0.523226 7054
							1.552610 11	0.886524 7169
								1.451567 311

Table 1: Weights for 10 datapoints

Part 2

Part 2 takes up the job of looking into how good the training went about using visual plots. First, it asks the user if there is a desire to generate a new dataset, or if at all the training needs to be performed again. After such confirmations, it proceeds to plot the data. First, it reads the saved data and weights after which it recomputes the test and train errors. This is done as the test and train errors, though computed, were not saved in Part 1. For plotting the approximated function using the given weights, the output of the function is obtained for inputs between 0 and 1, separated by 0.01. Thus, when a line plot is drawn on this, it gives the illusion of a curve.

Fig 1: Curve fitting against 10 datapoints

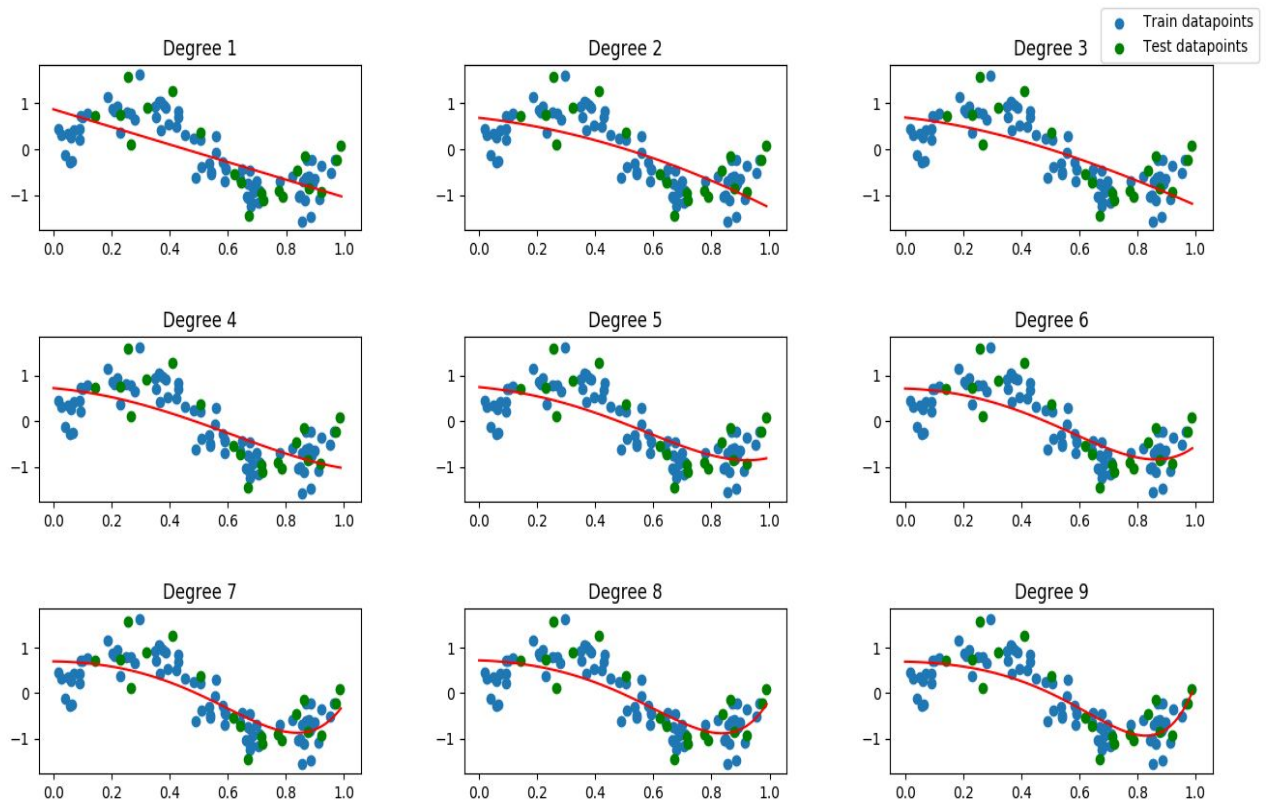
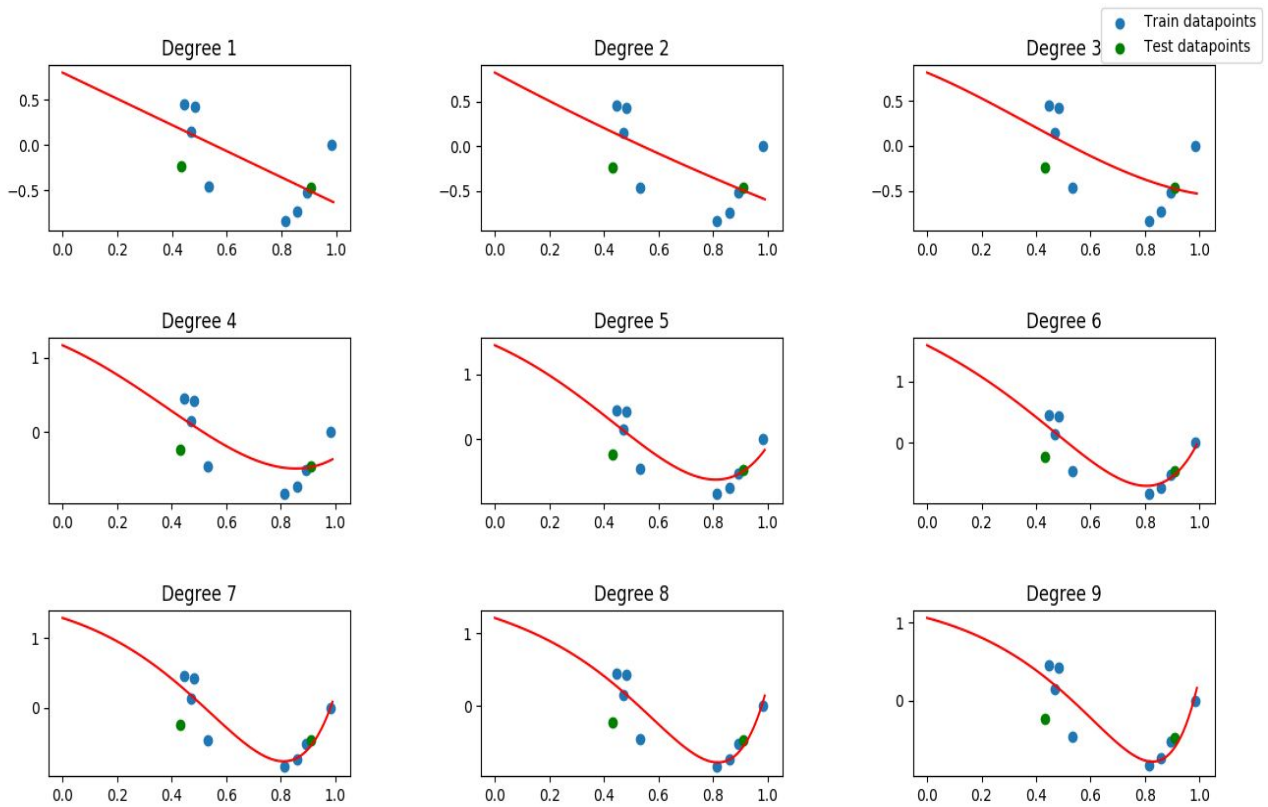


Fig 2: Curve fitting against 100 datapoints

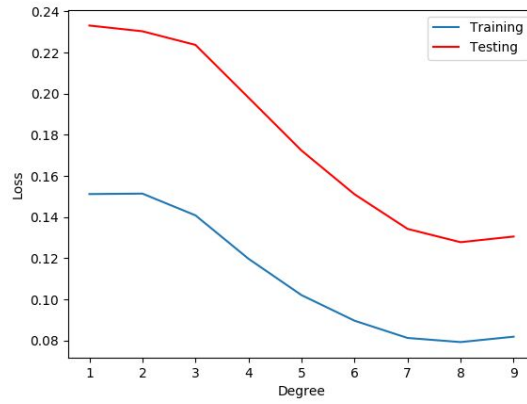


Fig 3: Loss v/s degree for 100 datapoints

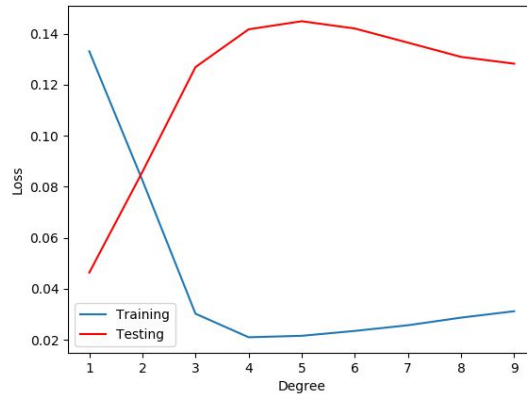


Fig 4: Loss v/s degree for 10 datapoints

Next, the error values obtained in the previous part is plotted against the degree of the polynomial function which is being approximated. Lastly, the degree for which the minimum value of test error was obtained, is returned.

In the experiment, it was found that if the number of datapoints is restricted to 10, the plots obtained were a bit arbitrary. When increased to 100, they looked like they had been anticipated.

Part 3

Part 3 makes a call to Part 2 to obtain the best value of n . Accordingly, it then instructs Part 1 to generate data with various sizes, namely 10, 100, 1,000, and 10,000. The weights obtained have been presented in Table 2, 3 and 4. Upon testing a number of times, the value of the best n is found according to the following distribution (Table 5). The plot against no of datapoints is given in Fig 5.

Table 5: Distribution of degrees

Degree	7	8	9
%	6.67	13.33	80.00

Table 2: Weights for 100 datapoints

1	2	3	4	5	6	7	8	9
0.871026 2544	0.829211 183	0.875575 157	0.908950 669	0.906696 2384	0.884664 807	0.844219 5304	0.821845 2823	0.783383 0117
-1.894498 614	-1.656683 609	-1.625247 068	-1.363524 754	-1.085309 647	-0.879548 715	-0.623205 9524	-0.529345 3787	-0.390148 1505
	-0.226027 0969	-1.296194 565	-1.916144 381	-2.230179 573	-2.199814 806	-2.001010 521	-1.938706 417	-1.637137 163
		1.160183 327	-0.460245 4562	-0.882485 2599	-0.876267 4572	-1.675845 952	-1.206534 577	-1.623533 107
			2.240392 562	1.066217 635	-0.038717 92457	0.057640 72756	-0.917430 5237	-0.875185 3222
				1.875091 938	0.908160 266	0.605138 7781	0.439045 7082	0.381673 158
					2.073271 886	1.295730 354	1.074630 988	0.304503 3555
						1.563315 841	1.146619 257	0.708973 1967
							1.346124 277	1.513295 952
								1.273084 84

Table 3: Weights for 1000 datapoints

1	2	3	4	5	6	7	8	9
0.795294 9234	0.795661 4451	0.671008 6979	0.674772 2721	0.676215 1944	0.741085 8763	0.680414 0972	0.797758 8119	0.701040 2116
-1.606634 293	-1.160043 219	-0.635297 9473	-0.771208 8956	-0.597161 0846	-1.063300 577	-0.659153 834	-1.048135 427	-0.675897 5782
	-0.643251 0116	-0.593078 5415	-0.377599 4579	-1.008405 992	-0.543273 5428	-1.003218 382	-0.616417 6367	-0.770524 1842
		-0.603542 5272	-0.709381 2278	-0.529990 1328	-0.156100 4362	-0.325332 3324	-1.015222 288	-0.702825 1569
			0.082220 34768	0.234952 8062	-0.111481 7494	-0.363662 6983	-0.038812 81459	-0.778274 6306
				0.267195 1549	-0.173285 7535	0.475429 1399	0.417499 5428	0.278040 7732
					0.421244 2658	0.043486 57395	0.413433 1526	0.087793 55429
						0.387009 0428	-0.177975 0626	0.338968 5656

							0.767720 4145	0.901786 2445
								0.197253 533

Table 4: Weights for 10,000 datapoints

1	2	3	4	5	6	7	8	9
0.781714 3592	0.709914 3363	0.682912 0663	0.690717 5004	0.682811 1931	0.752605 908	0.731046 685	0.709935 9958	0.853432 7123
-1.572001 909	-0.728133 1692	-0.597355 4095	-0.741943 3226	-0.936088 8669	-0.907285 8963	-0.647606 6363	-0.827493 4965	-1.121219 198
	-1.033013 008	-1.085452 251	-0.708731 8272	-0.319905 1655	-1.038105 144	-1.134167 163	-0.696408 7892	-0.978014 8136
		-0.077089 84932	-0.557043 3458	-0.274101 1187	-0.078618 73709	-0.830478 5904	-0.628142 2084	-0.171826 8267
			0.292768 2211	-0.307388 3862	-0.296437 6787	0.032630 71575	-0.196983 0868	-0.722522 0256
				0.135145 2717	0.453599 0327	-0.065458 67159	0.078122 1145	0.038532 14125
					0.355022 3538	0.560727 4715	0.171170 8567	0.550561 6662
						0.817546 4423	0.118368 9984	-0.161832 919
							0.691812 0453	0.653741 6599
								0.828182 0676

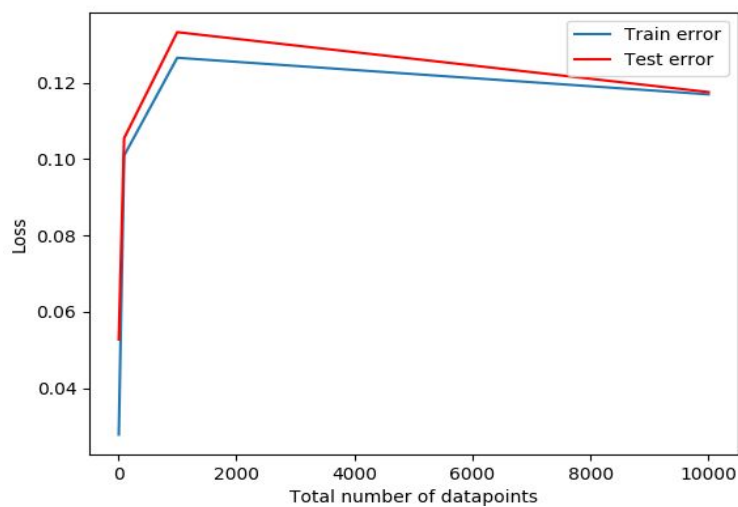


Fig 5: Loss v/s no of datapoints

Part 4

Two other cost functions were experimented with, namely the mean absolute error and the fourth power error. For the mean absolute error, the following formula was used for computing cost

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m |W^T \Phi_n(x) - y|$$

and the following for gradient.

$$\frac{dMAE}{dy_{pred}} = \begin{cases} +1, & y_{pred} > y_{true} \\ -1, & y_{pred} < y_{true} \end{cases}$$

For the fourth power error, the cost function was

$$\frac{1}{2m} \sum_{i=1}^m (W^T \Phi_n(x) - y)^4$$

Accordingly, the gradient comes out to be

$$\frac{2}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Here, the weights have been initialised randomly between 0 and 0.001 to prevent exploding of values.

The RMSE error is trained and computed for each degree. For each instance of the learning rate (0.025, 0.05, 0.1, 0.2 and 0.5) and for every cost function, the degree having the minimum value of RMSE is chosen and plotted across the learning rate values. The weights have been given in Table 6, 7 and 8. The provided weights are for that degree for which the minimum value of RMSE was obtained. The first row of the table mentions the learning rate.

Table 6: Weights for Squared error loss

0.025	0.05	0.1	0.2	0.5
0.7293215757	0.7465717494	0.9291357642	1.042771469	0.9614016857
-0.7186918233	-0.8293871128	-0.9888155837	-1.183347055	-0.5719301969
-0.845257154	-0.9279066871	-1.329807527	-1.330827547	-1.885688898
-0.8664141269	-0.2307502415	-0.8151101417	-0.9837553255	-1.699878445
0.07024936875	-0.1309270971	-0.3318826047	-0.6089554393	-0.173358147
0.1790967975	-0.1041283308	0.246028049	-0.1855961717	-0.2156919053
-0.1095506467	0.1267188376	0.4721587248	0.9466496538	0.669988285
0.3165117195	-0.1687833015	0.7737945748	0.4701290394	0.8501394341
0.2331473181	0.7872001551	0.8259655405	1.190339921	1.096247705
0.5834135234	0.242288971		0.6975961542	1.155297139

Table 7: Weights for Mean Absolute Error

0.025	0.05	0.1	0.2	0.5
0.5897437497	0.7458329931	0.990000795	1.225462428	1.269331133
-0.5366108515	-0.7737905044	-1.352076637	-1.661870657	-1.43829365
-0.5703084878	-0.5708972985	-0.8887757569	-1.341224788	-1.558887571
-0.5431326869	-1.033364742	-0.3741909892	-0.4470348273	-1.569393699
-0.3934620506	0.06160323211	-0.5582285952	-0.1297491422	-0.2512687934
0.263836122	0.266539504	0.1825713155	-0.3301895356	0.3245040276
0.2240216534	0.06892601893	0.4764327226	0.2615325467	0.9356641203
0.02726837111	0.2371244755	0.3778189537	0.3538152723	1.070499122
-0.1784790498	0.5112638983	0.3292100977	1.07589597	1.193997085
0.5924522313		0.5246380648	0.9210709979	

Table 8: Weights for Fourth Power Error

0.025	0.05	0.1	0.2	0.5
0.6245150067	0.8007088396	0.9432650773	0.989307045	0.9290307589
-0.670324757	-0.9542326425	-1.168411553	-1.129097379	-0.6344867589
-0.6319956396	-0.883328324	-1.152236085	-1.400483947	-1.812908766
-0.3972964747	-0.5198589783	-0.6711720606	-0.870859511	-1.294542826
-0.1873818257	-0.1955922497	-0.2225846011	-0.3046772061	-0.5224040152
-0.0342709929				
7	0.04018260693	0.1083178591	0.1296636942	0.1047213058
0.0724038637	0.201124684	0.3334643396	0.4263289758	0.5362645949
0.1431212527	0.3063118265	0.4794926944	0.618876233	0.8057831371
0.1888810998	0.3731522693	0.5704175346	0.7354502344	0.9603547974
0.2180880341	0.4133952611	0.6234962827	0.8010226496	1.034530432

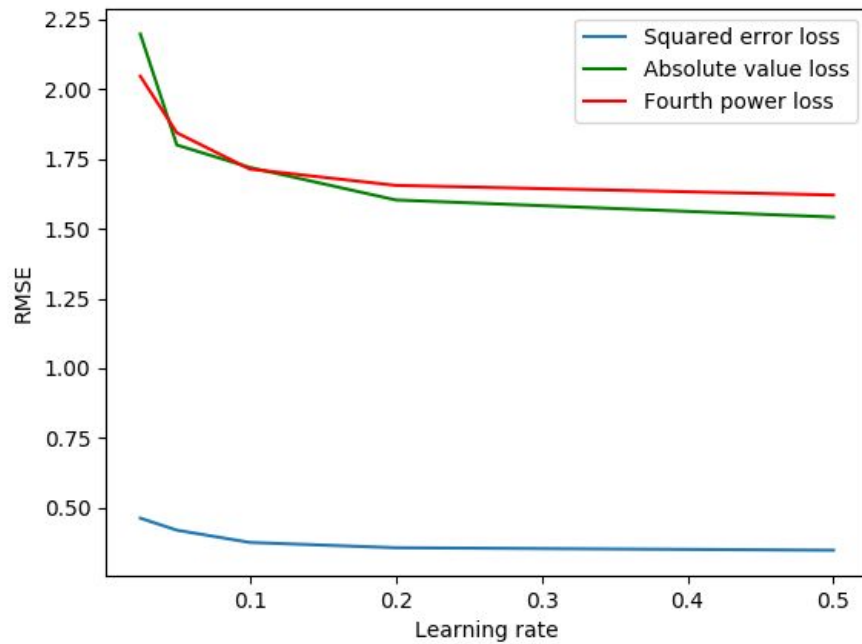


Fig 6: Test RMSE v/s Learning rate

CONCLUSION

We conclude that the best degree for fitting the curve has been 9 (as evident from Table 5), as seen in Part 4. The best cost function is Squared Error Loss (as evident from Fig 6). In reality, higher powers are prone to noises, and since noises were manually added to the synthetic dataset generated, perhaps that's why it did not perform well. The mean absolute error is not differentiable and hence, not used in general. The best learning rate has been experimentally found to be 0.5. Along with that, we can infer that we should have sufficient number of datapoints to draw valid conclusions about our model. For instance, in Part 2, having 10 datapoints gave us a wrong notion about our model (ref to Fig 4), whereas having a 100 datapoints could fit the data better (Fig 2). Having an even higher no of datapoints decreases the loss, as evident in Part 3 (Fig 5).

REFERENCES

- Seber, George AF, and Alan J. Lee. Linear regression analysis. Vol. 329. John Wiley & Sons, 2012.
- Needell, Deanna, Rachel Ward, and Nati Srebro. "Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm." Advances in Neural Information Processing Systems. 2014.
- Robert, Christian. "Machine learning, a probabilistic perspective." (2014): 62-63.