
Learning to Rank using Linear Regression

Anunay Rao

anunayra@buffalo.edu

1 Introduction

The project is to train linear regression model on LeToR dataset by using closed-form solution and then by using stochastic gradient descent approach. For the closed form solution we will be using Moore-Penrose pseudo-inverse method to apply matrix inversion. On the other hand, for stochastic gradient descent we will first randomly initialize the weights and then iteratively update them to get the optimum weights for which our model performs better. Then we will study the effect of various hyperparameters in both the closed form approach and the stochastic gradient descent approach.

2 Preprocessing of the Dataset and Basis Function

The LeToR dataset consists of input which is derived from query-document pair and the target variable which takes the value 0,1 or 2, But here we will obtain real value using linear regression which is more useful for ranking tasks as it avoids collision into only three possible values.

For this linear regression task we need to find the linear function which is a function of some functions called basis functions. For both the closed form solution and stochastic gradient descent solution we are using Gaussian function as the basis function which is given by:

$$e^{-(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

where x is the input, μ is the mean around which gaussian function is spread i.e it denotes where the function is located and Σ is the covariance matrix which gives the variance of a particular feature with respect to others. As it requires the inverted covariance matrix and it is known that if a matrix has any column with zero variance it is non-invertible. In the LeToR dataset we have column number 6,7,8,9 and 10 denoting IDF of body, IDF of anchor, IDF of title, IDF of URL and IDF of whole document respectively with zero variance, infact these features have zero value for all the samples therefore we need to filter out these columns in order to proceed with the linear regression task. Here, we are taking 10 Gaussian basis functions which will require 10 means so we form 10 clusters of the dataset and for each cluster we will have μ vector of dimension 1x41. Gaussian basis function is used because its value depends on the origin or certain point. It is symmetrical about central maximum.

40 **3 Performance Metric**

41 We will evaluate the solution obtained by using Root Mean Square (RMS)

42 error defined as $E_{RMS} = \sqrt{2E(w^*)/N_V}$

43 where w^* is the solution and N_V is the size of dataset. Accuracy is not a good
44 performance metric for this linear regression tasks and hence we will study
45 the effect of various hyper-parameters on E_{RMS} and not accuracy.

46 **4 Effect of various Hyper-parameters**

47 This section presents the experiments conducted by varying the values of
48 Hyper-parameters and how they affect the performance of the model.

49

50 **4.1 CLOSED FORM SOLUTION**

51 Here, we are using Moore-Penrose pseudo-inverse method to find the inverse
52 of a matrix because we want the inverse of ϕ which is not a square matrix. We
53 will be using regularization term λ to take care of overfitting.

$$W_{10 \times 1} = (\underbrace{\phi_{10 \times 58000}^T \phi_{58000 \times 10}}_{10 \times 10})^{-1} \underbrace{\phi_{10 \times 58000}^T}_{10 \times 58000} \underbrace{t}_{58000 \times 1}$$

54
55

Figure 1: Dimensionality for Training

$$t_{1 \times 6000} = W_{1 \times 10}^T \underbrace{\phi_{val}(x)_{10 \times 6000}}_{1 \times 6000}$$

56
57

Figure 2: Getting target value from Validation Design Matrix

58 Default configuration:

- 59
- 60 • Regularization term, $\lambda = 0.03$
 - 61 • Training Percent = 80
 - 62 • Validation Percent = 10
 - 63 • Testing Percent = 10
 - 64 • Number of Basis Function, $M = 10$

65 **4.1.1 Effect of Number of Basis Function (M)**

	10	20	30	40	50	60	70	80	90
Erms_Training	0.549469	0.546123	0.543265	0.541604	0.540736	0.539609	0.539366	0.539054	0.538997
Erms_Validation	0.538428	0.53836	0.53718	0.536806	0.536946	0.536299	0.535902	0.53599	0.536013
Erms_Testing	0.627979	0.62619	0.622963	0.620712	0.619718	0.619219	0.618967	0.618747	0.619126

66
67

Table 1: E_{RMS} values for different values of M

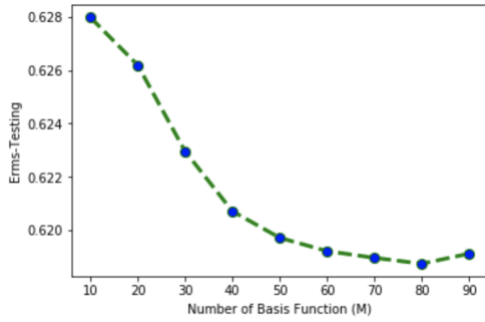


Figure 3: Variation of E_{RMS} -Testing with M

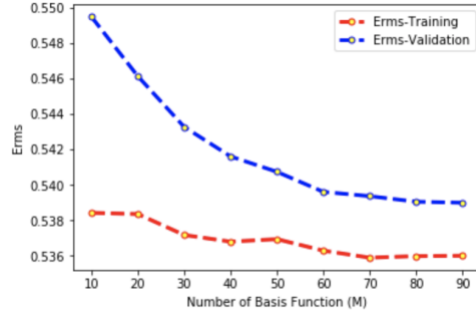


Figure 4: Variation of E_{RMS} -Training and Validation with M

68 **Conclusion**

69 As the number of basis function increases E_{RMS} decreases but if we will use
70 too many basis functions then our model will overfit.

71
72

4.1.2 Effect of Regularization term (λ)

	0.03	0.3	3	30	300
Erms_Training	0.549469	0.549618	0.550344	0.551271	0.554998
Erms_Validation	0.538428	0.538764	0.539797	0.540525	0.543709
Erms_Testing	0.627979	0.628213	0.628197	0.62847	0.631591

73
74

Table 2: Respective E_{RMS} values for different λ

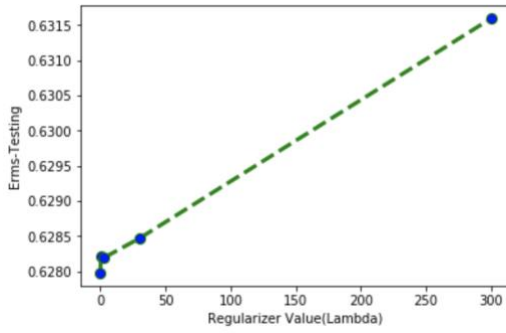


Figure 5: Variation of Erms-testing with regularizer

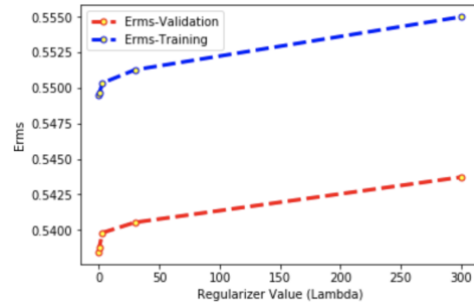


Figure 6: Variation of Erms Training and Validation with regularizer

75

76

77 Conclusion

78 As the the value of λ increases the Erms also increases although for smaller
 79 values of λ there is not any drastic change in Erms. It is used to avoid overfitting
 80 so that our model does not strictly fit to the training data points.

81

82 4.1.3 Effect of Data Split (Training: Validation: Testing)

83

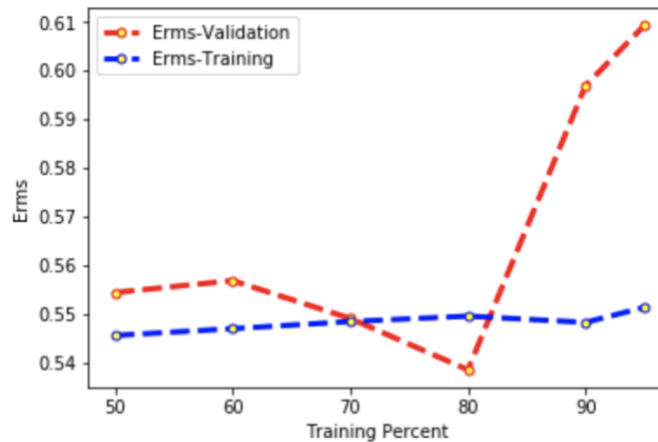
	50:25:25	60:20:20	70:15:15	80:10:10	90:5:5	95:2.5:2.5
Erms_Training	0.545498	0.546892	0.548359	0.549469	0.548179	0.551264
Erms_Validation	0.55432	0.55678	0.548985	0.538428	0.596895	0.609278
Erms_Testing	0.580895	0.584775	0.604691	0.627979	0.657624	0.704038

84

85

Table 3: Erms values for different Data Split

86



87

88 **Conclusion**

89 As the training data set increases the model is validated and tested on less number
90 of data sets. As seen from the data collected by varying the data split we have that
91 for 70:15:15 the Erms training and validation are almost same but on further
92 increasing the training percent thereby reducing the validation and testing data set
93 the Erms of Validation and Testing increased by big margin thereby indicating the
94 overfitting of the model. The weights we have obtained are not optimum weights.

95

96 **4.2 STOCHASTIC GRADIENT DESCENT**

97 The gradient descent approach can be break down in three steps:

- 98 • Randomly initialize the weights.
- 99 • Update the weights iteratively.

100

$$w_{new} = w_{old} + \Delta w$$

101

$$\Delta w_{old} = -\eta_{old} \nabla E$$

102

103

$$\nabla E = \nabla E_D + \lambda \nabla E_W$$

104

105

$$\nabla E_D = -(t_n - w_{old}^T \phi(x_n)) \phi(x_n)$$

106

107

$$\nabla E_W = w_{old}$$

108

109

- Calculate the E_{RMS} value.

110 Here η is a learning rate which determines the step size in which weight gets
111 updated. The lower the value of learning rate the model will converge slowly
112 as it will take small step towards the minima whereas if the value of learning
113 rate is high it will take bigger steps but then it is not guaranteed that it will
114 converge to local minima. Therefore we need an optimum value of learning
115 rate.

116 **Learning Process and why training over only 400 sample:**

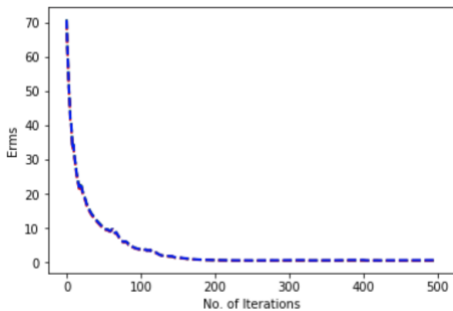


Figure 7: Variation of Erms with No. of Iterations

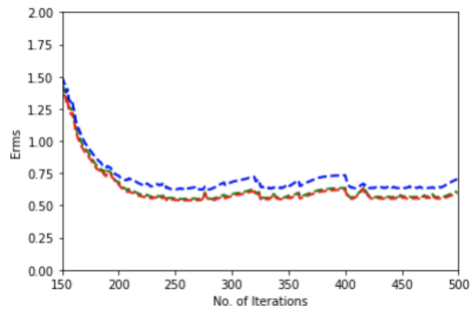


Figure 8: Magnified view of Figure 7

117 **Conclusion**

118 As it is quite evident from the above two graphs that after 400 iterations(samples)
119 the Erms values becomes quite stable and so there is not any significant change in

120 weights and so we stop at 400.

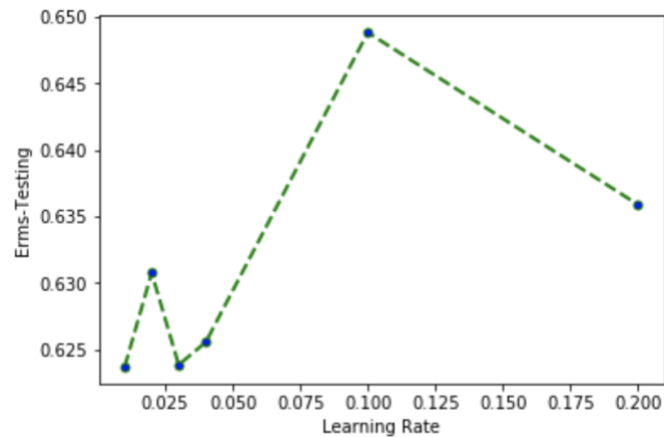
121

122 4.2.1 Effect of learning rate

	0.0001	0.001	0.01	0.02	0.03	0.04	0.1	0.2
Erms_Training	53.6082	13.1223	0.54964	0.55869	0.54955	0.55315	0.58417	0.56518
Erms_Validation	53.0331	12.7047	0.53846	0.54741	0.53842	0.5416	0.57105	0.55413
Erms_Testing	53.476	13.0067	0.62372	0.63076	0.62386	0.62556	0.64879	0.63588

123
124

Table 4: Erms values for different values of learning rate



125

126 Conclusion

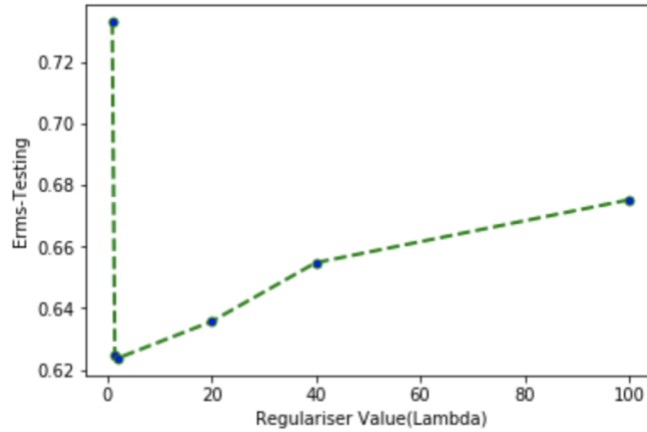
127 As we have considered only 400 samples(iterations) then it is quite evident that
128 for $\eta = 0.0001$ we have higher Erms because it will converge slower as it is
129 taking smaller steps towards the minima. The plot is not shown for $\eta = 0.0001$
130 and 0.001 to get the better understanding from the plot. As the learning rate
131 increases it will take bigger steps towards minima but it may or may not
132 converge. From the plot we have minimum Erms for $\eta = 0.01$. Therefore good
133 learning rate is not a smaller value nor a higher value but it is a optimum value.

134

135 4.2.1 Effect of Regularization term (λ)

	0.2	1	1.5	2	20	40	100
Erms_Testing	12.9095	0.73303	0.62483	0.62372	0.63566	0.65481	0.67524

136



137
138

Conclusion

139 As the value of λ increases Erms also increases. There is an optimum value of
140 λ for which the Erms is minimum. For a very low value of λ we have higher
141 Erms which means the weight we have are not the optimum weights.

142

5 Other Variations

5.1 Deleting Column with no information:

145 In the LeToR dataset the column number 46 namely, Number of child page has
146 the value 0 for all the samples except one. Therefore deleting this column and
147 checking out its effect:

148

149

Closed Form:	
Before deleting:	After Deleting:
E_rms Training = 0.5488846063481925	E_rms Training = 0.5494694067137487
E_rms Validation = 0.5373461870705957	E_rms Validation = 0.5384281741389716
E_rms Testing = 0.6273650154146978	E_rms Testing = 0.6279788453854206
Gradient Descent:	
Before deleting	After deleting:
E_rms Training = 0.54986	E_rms Training = 0.55044
E_rms Validation = 0.53777	E_rms Validation = 0.53899
E_rms Testing = 0.6247	E_rms Testing = 0.62518

150

Conclusion

152 Negligible change in Erms and thus column 46 (Number of child page) was
153 useless.

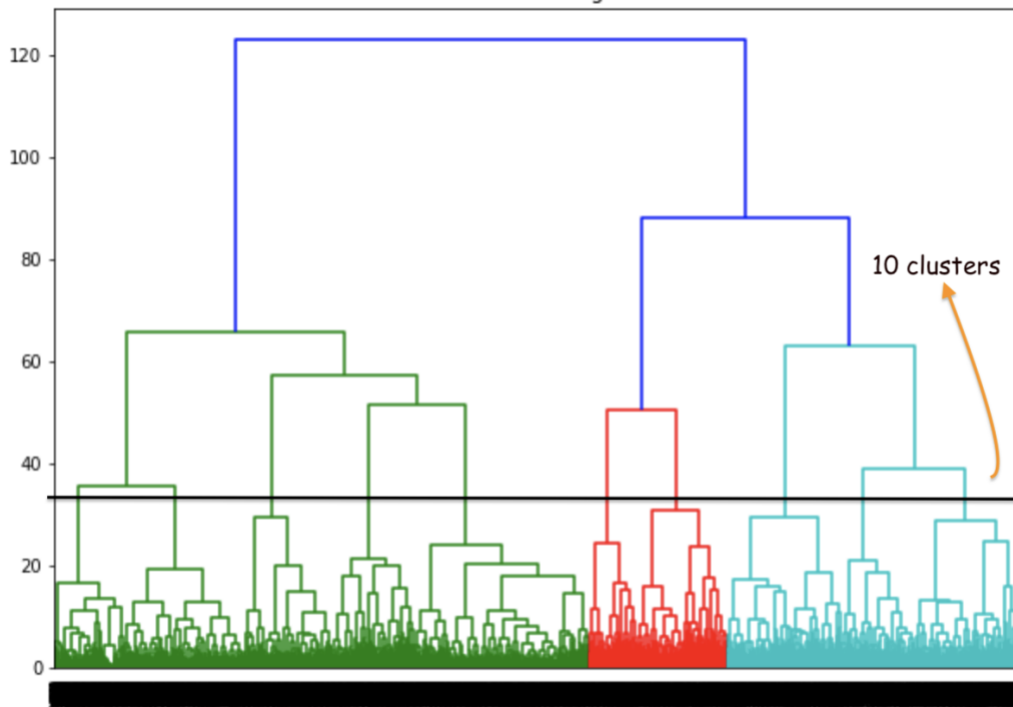
154

155

5.2 Hierarchical Clustering

157 An attempt to find the optimum number of clusters but with limited computational
158 power and time I could only get the dendrogram for 30% of the training dataset.

159 Therefore, determining clusters for it so the longest vertical distance without any
160 horizontal line passing through it is selected and a horizontal line is drawn through it. The number
161 of vertical lines this newly created horizontal line passes is equal to number of clusters
162



163
164

Figure 9: Dendrogram for 30% Validation Data

165

166 Conclusion

167 So after cutting the dendrogram we have 10 clusters.

168

169 References

170

171 [1] Stack Abuse. (2018). *Hierarchical Clustering with Python and Scikit-Learn*. [online] Available at:
172 <https://stackabuse.com/hierarchical-clustering-with-python-and-scikit-learn/> [Accessed 5 Oct. 2018].

173 [2] Towards Data Science. (2018). *Machine Learning – Towards Data Science*. [online] Available at:
174 <https://towardsdatascience.com/machine-learning/home> [Accessed 10 Sep. 2018].

175 [3] Brownlee, J. (2018). *Machine Learning Mastery*. [online] Machine Learning Mastery. Available at:
176 <https://machinelearningmastery.com/>