
CSE 574: Introduction to Machine Learning

Atrayee Nag
#50288651
atrayeen@buffalo.edu

Abstract

In this project, we are going to use Machine learning to solve the Learning to rank problem that arises in Information Retrieval process. We will solve the problem by formulating it into linear regression where we map an input vector x to a real-valued scalar target $y(x, w)$. We have used the Closed form Solution to find out the weights and Stochastic Gradient Descent to train our model.

1 General Concepts

1.1 What is Learning to Rank Problem ?

Learning to rank is a sub-area of machine learning, studying methodologies and theories for automatically constructing a model from data for a ranking problem suppose $[0, 1, 2]$. Learning to rank uses supervised machine learning to discover the best order for a list of items, using features extracted from each item to give it a ranking. It does not provide a definite solution but instead assigns a relative order to each item over another. For web results some of the features are as follows:

1. Whether the page title matches the query
2. How many hits have been there on the page
3. How long is the content of the page
4. Age or popularity of the page
5. Whether the page can be viewed on a phone browser

Learning to rank ties machine learning to search engines is the process of identifying the most important features of document set to the users of those documents, and using those features to tune the search engine to return the best fit documents to each user on each search. Tokens are generated at search time which are looked up from the query at the inverted index, ranks the matching documents and retrieves and returns text associated with it to the user. The diagram below shows the LTR associated with webpages which corresponds to our training , validation and testing data.

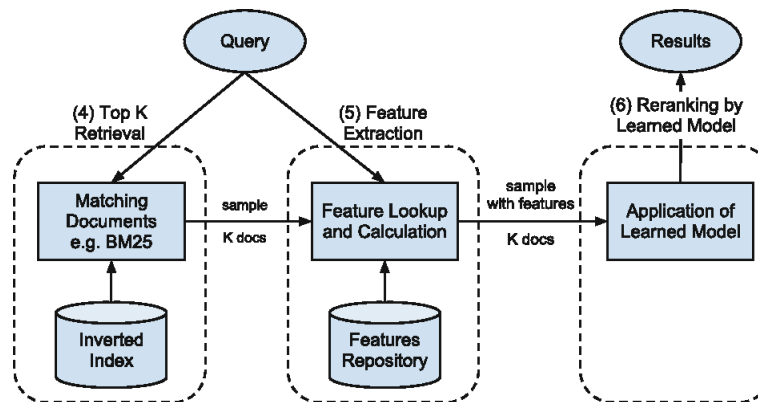


Figure 1: LTR on Search Engines

24 1.2 What is Linear Regression?

25 Simple Linear Regression is a type of regression where there is only one independent variable and
 26 the dependent variable y is linear to x. The target function is defined as below:

$$27 \quad y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 * x_1$$

28 Multivariate linear regression can be defined as a linear combination of input variables, where there
 29 are n number of independent variables. The goal of linear regression is find a polynomial function to
 30 approximate the target function. The target function for multivariate linear regression is as follows:

$$31 \quad y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 * x_1 + w_2 * x_2 + \dots + w_n * x_n$$

32 The generalized equation is given by :

$$33 \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x})$$

34 where $\mathbf{w} = (w_0, w_1, \dots, w_{M-1})$ is a weight vector to be computed from training samples and $\phi = (\phi_0, \dots,$
 35 $\phi_{M-1})$ is a vector of M basis functions.

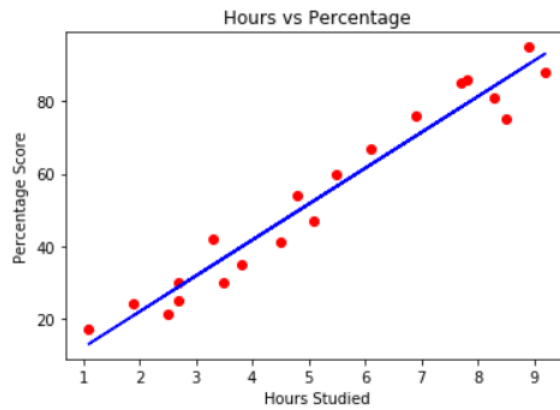


Figure 2: Linear Regression

36 2 What is Basis Function?

37 The plain form of Linear Regression is polynomial basis functions, one of big limitation is that they
 38 are global functions of the input variable, so that changes in one region of input space affect all
 39 other regions. This can be resolved by dividing the input space up into regions and fitting a different
 40 polynomial in each region, leading to spline functions. A common type of basis function for such
 41 models is the Gaussian basis function. This type of model uses the kernel of the normal (or Gaussian)
 42 probability density function (PDF) as the basis function. The equation for a Gaussian radial Basis
 43 Function is given as follows:

$$44 \quad \phi(\mathbf{x}) = \exp(-0.5(\mathbf{x} - \mu_j)^\top \Sigma_j^{-1}(\mathbf{x} - \mu_j))$$

45 where μ_j is the centroid of the basis function and Σ_j decides how broadly the basis function spreads.

46 The below diagram displays a graph for Gaussian basis function:

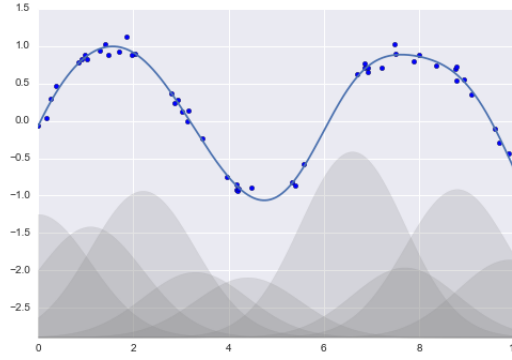


Figure 3: Basis Function

3 Steps involved in the Learning to Rank Problem

There are two steps involved in our Learning to rank problem.

1. We calculate the Root Mean Square Errors using the Closed form solution.
2. We train our model using the Stochastic Gradient Descent

3.1 Closed Form Solution

An equation is said to be a closed-form solution if it solves a given problem in terms of functions and mathematical operations from a given generally-accepted set. For example, an infinite sum would generally not be considered closed-form. However, the choice of what to call closed-form and what not is rather arbitrary since a new closed-form function could simply be defined in terms of the infinite sum.

The closed form solution for calculating the weights without regularization has the form:

$$w_{ML} = (\phi^T \phi)^{-1} \phi^T t$$

where $t = \{t_1, \dots, t_n\}$ is the vector of outputs in the training data and ϕ is the design matrix

The closed-form solution with least-squared regularization has the below form:

$$w^* = (\lambda I + \phi^T \phi)^{-1} \phi^T t$$

where λ is the regularization factor.

```

1  kmeans = KMeans(n_clusters=M, random_state=0).fit(np.transpose(TrainingData))
2  Mu = kmeans.cluster_centers_
3
4  BigSigma      = GenerateBigSigma(RawData, Mu, TrainingPercent, IsSynthetic)
5  TRAINING_PHI = GetPhiMatrix(RawData, Mu, BigSigma, TrainingPercent)
6  W            = GetWeightsClosedForm(TRAINING_PHI, TrainingTarget, (C_Lambda))
7  TEST_PHI     = GetPhiMatrix(TestData, Mu, BigSigma, 100)
8  VAL_PHI      = GetPhiMatrix(ValData, Mu, BigSigma, 100)
9
10 TR_TEST_OUT  = GetValTest(TRAINING_PHI, W)
11 VAL_TEST_OUT = GetValTest(VAL_PHI, W)
12 TEST_OUT     = GetValTest(TEST_PHI, W)
13

```

```

14 TrainingAccuracy = str(GetErms(TR_TEST_OUT, TrainingTarget))
15 ValidationAccuracy = str(GetErms(VAL_TEST_OUT, ValDataAct))
16 TestAccuracy = str(GetErms(TEST_OUT, TestDataAct))

```

65 The steps involved in the Closed form solution to solve the Learning to rank problems are as follows:

- 66 a **Creating the training Dataset:** The training data is divided and a percentage of 80% of it
67 is stored in the numpy array. While creating the dataset we are deleting 5 features which has
68 no change across the dataset. It eliminates values which will give the covariance as 0, thus
69 the inverse cannot be calculated. We further can eliminate features in the dataset which does
70 not add anything to get the data.
- 71 b **Clustering the DataSet using K-means clustering:** This is a clustering method which
72 takes in 2parameters: the number of basis functions we want to generate and the random
73 state. We can input an array of centroids if we have previous knowledge of the centroids, in
74 our case we do not have prior knowledge so just initialize it to 0 or None and let K means
75 calculate it using the Euclidean distance.
- 76 c **Calculating BigSigma:**
77 1. Initialize the BigSigma into a 2-d numpy array of size of the features
78 2. For each row across the dataset it calculates the variance and stores it in a diagonal
79 matrix Bigsigma which represents the feature vector with respect to itself. We use a
80 diagonal matrix because we are calculating variance for a feature and since 2 features are
81 unrelated the variance will be 0.
82
- 83 d **Get PHI Matrix:** We formulate the Gaussian Radial Basis function which converts each
84 row of the vectors - Data and Mu multiplied by inverse of BigSigma into a single scalar
85 value and store it in the PHI matrix. The equation for calculating the PHI matrix is given by
86 : $\phi_j(\mathbf{x}) = \exp(-0.5(\mathbf{x} - \mu_j)^T \sum_j^{-1}(\mathbf{x} - \mu_j))$
87
- 88 e **Get Weights using Closed Form:** We use the hyper parameter lambda to bound the
89 weights within a certain range , so that even if one of the weight is spiked upto a large value,
90 the lambda brings down the weight, so that no particular feature dominates and result in
91 over fitting. We calculate the weights with the formula Moore-Penrose pseudo-inverse,i.e.
92 the sum of squared inverse with regularization of the PHI matrix. The equation is given by:
93 $w^* = (\lambda \mathbf{I} + \phi^T \phi)^{-1} \phi^T \mathbf{t}$
94
- 95 f **Calculate Erms for the Dataset** This method returns the root mean square errors when
96 comparing the training / validation dataset with the target value. We first calculate the
97 residuals i.e. the difference between the actual values and the predicted value and then
98 square it for every dataset and divide it by the total number of datasets and then take a square
99 root which gives us the corresponding name of the square calculation. The equation to find
100 the same is given by:
101 $E_{RMS} = \sqrt{2E(w^*)/N_v}$
102 where w^* is the solution and N_v is the size of the test dataset.

103 3.2 Stochastic Gradient Descent

104 Stochastic gradient descent (often shortened to SGD), also known as incremental gradient descent,
105 is an iterative method for optimizing the whole dataset at every step. It is a very efficient approach
106 to discriminative learning of linear classifiers under convex loss functions such as (linear) Support
107 Vector Machines and Logistic Regression. As the algorithm sweeps through the training set, it
108 performs the above update for each training example. Several passes can be made over the training
109 set until the algorithm converges. If this is done, the data can be shuffled for each pass to prevent
110 cycles.

```

1  W_Now      = np.dot(220, W)
2  La         = 2
3  learningRate = 0.01
4  L_Erms_Val  = []
5  L_Erms_TR   = []
6  L_Erms_Test = []
7  W_Mat       = []
8
9  for i in range(0,400):
10     Delta_E_D = -np.dot((TrainingTarget[i] - np.dot(np.transpose(W_Now)
11     , TRAINING_PHI[i])), TRAINING_PHI[i])
12     La_Delta_E_W = np.dot(La, W_Now)
13     Delta_E      = np.add(Delta_E_D, La_Delta_E_W)
14     Delta_W      = -np.dot.learningRate, Delta_E)
15     W_T_Next     = W_Now + Delta_W
16     W_Now        = W_T_Next
17
18     TR_TEST_OUT  = GetValTest(TRAINING_PHI, W_T_Next)
19     Erms_TR      = GetErms(TR_TEST_OUT, TrainingTarget)
20     L_Erms_TR.append(float(Erms_TR.split(',')[1]))
21
22     VAL_TEST_OUT = GetValTest(VAL_PHI, W_T_Next)
23     Erms_Val     = GetErms(VAL_TEST_OUT, ValDataAct)
24     L_Erms_Val.append(float(Erms_Val.split(',')[1]))
25
26     TEST_OUT     = GetValTest(TEST_PHI, W_T_Next)
27     Erms_Test    = GetErms(TEST_OUT, TestDataAct)
28     L_Erms_Test.append(float(Erms_Test.split(',')[1]))

```

The steps involved in the Stochastic Gradient Descent solution to solve the Learning to rank problems are as follows:

- a **Randomize the Weights:** We randomized the weights by multiplying a scalar value with the weights obtained from the closed form weight solution. However the very first step in gradient descent is randomizing the weights and we could have taken new random weights instead of using the weights from the previous solution.
- b **Update the Weights:** We process the entire training set in one go and update all the weights by calculating the error multiplied by the learning rate. We train the model in a loop of 400 instead of the total training data because it has been observed after 400 datasets the weights are stagnant and do not change. This method is called Early Stopping.
- c **Displaying the Erms values:** We are displaying training, validation and testing data on the fly during iteration. In a real world scenario, training a model could take multiple days, so it's the norm to printing out changes during the iteration as they can monitored in real time.

We are optimizing the weights for next iteration using this formula:

$$\mathbf{w}^{r+1} = \mathbf{w}^r + \Delta \mathbf{w}^r$$

where weight update $\Delta \mathbf{w}^r = \eta^r \nabla E$, η^r is learning rate

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

in which $\nabla E_D = -(t_n - \mathbf{w}^{r\top} \phi(x_n))\phi(x_n)$

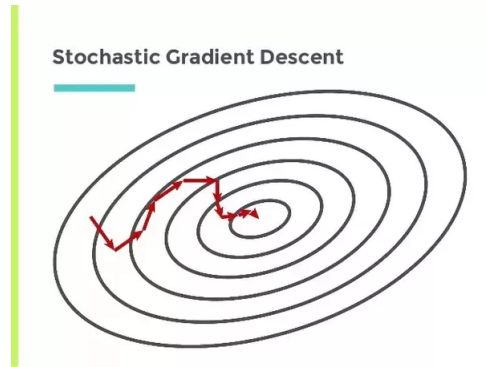


Figure 4: Stochastic Gradient Descent

129 4 Results on changing the HyperParameters

130 1. Case 1: Effect of changing Number of Clusters/Basis Functions

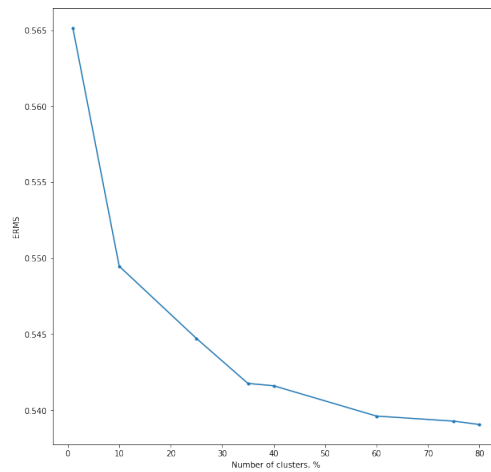


Figure 5: Ermc vs Cluster

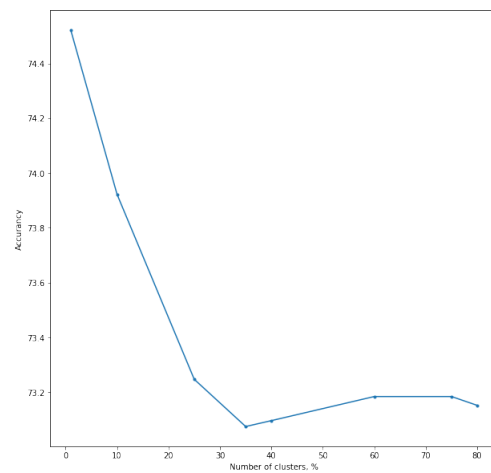


Figure 6: Accuracy Vs Cluster

2. Case 2: Effect of changing Lambda/Regularization

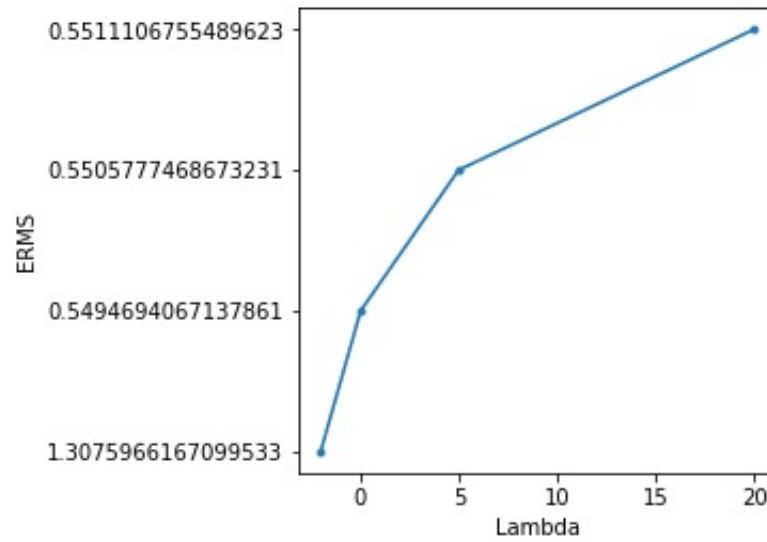


Figure 7: Ermc vs Lambda

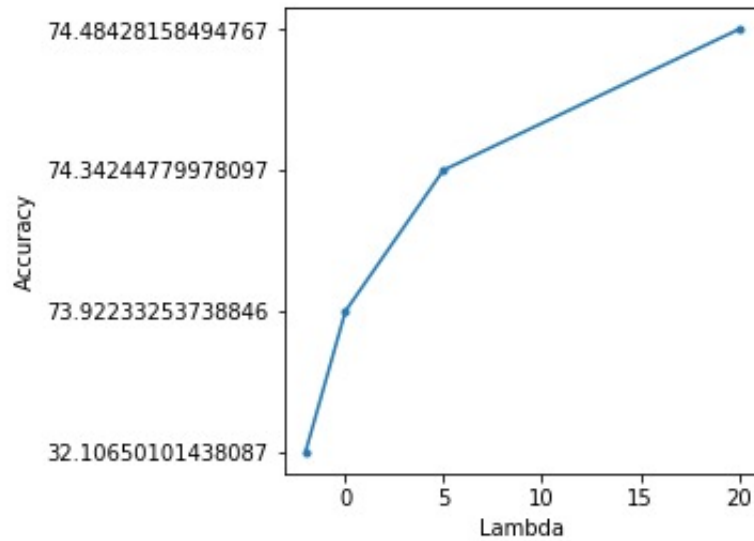


Figure 8: Accuracy Vs Lambda

5 Conclusion

1. The change in the number of clusters is not affecting the Erms or Accuracy by a big margin. However we can observe that there is slight decrease in both Erms and accuracy.
2. We can see from the graph that the increase in lambda decreases both the Erms and accuracy. However the values are stabilizing after a certain threshold.

137 **6 References**

138 The resources that has helped me understand the concepts of the Linear Regression are:

- 139 • Linear Regression : <http://gwansiu.com/2018/06/17/Linear-Regression/>
- 140 • Modeling Data with Linear Combinations of Basis Functions :
141 <http://www.utstat.utoronto.ca/radford/sta414.S11/week1b.pdf>
- 142 • Learning to Rank : [https://opensourceconnections.com/blog/2017/02/24/what-is-learning-](https://opensourceconnections.com/blog/2017/02/24/what-is-learning-to-rank/)
143 [to-rank/](https://opensourceconnections.com/blog/2017/02/24/what-is-learning-to-rank/)
- 144 • Andrew Ng Coursera : <https://www.coursera.org/learn/machine-learning>