

# Learning to Rank using Linear Regression

Haoquan Zhao

50189594

[haoquanz@buffalo.edu](mailto:haoquanz@buffalo.edu)

## 1. Abstract

The main goal of this project is to use machine learning linear regression techniques to solve the Learning to Rank (LeToR) problems on a real dataset and a synthetic dataset. We train our model using the training set which comprises of 80% of the given data set and validate it using the 10 % validation data set and finally test it using the 10% testing data set.

## 2. Introduction

In this project, we apply linear regression on a synthetic dataset from Microsoft LETOR 4.0, train our model on part of it and evaluate the performance on another part, then by tuning Hyper-Parameters, we get to know with which value we can get the minimum error. And apply this into test part of dataset. And finally, by using Root Mean Square error to evaluate the solution on a test part of the dataset, we could find out how the solution works.

The linear regression function is

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

where the  $\mathbf{w}$  is a weight vector to be learn from training samples and  $\boldsymbol{\phi}$  is a vector of  $M$  basis functions, and the Gaussian radial basis function is shown below:

$$\phi_j(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1}(\mathbf{x} - \boldsymbol{\mu}_j)\right)$$

where  $\boldsymbol{\mu}_j$  is the center of the basis function and  $\boldsymbol{\Sigma}_j$  decides how broadly the basis function spreads.

Closed form solution is calculated using basis function, Sigma, lambda and  $\mathbf{M}$ . Its equation is:

$$\mathbf{w}_{ML} = (\lambda \mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{t}$$

For the Stochastic Gradient Descent Solution, We start with an initial random value of solution and then modify this solution for each data sample.

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \Delta \mathbf{w}^{(\tau)}$$

Where

$$\Delta \mathbf{w}^{(\tau)} = -\eta^{(\tau)} \nabla E$$

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

In which

$$\nabla E_D = -(t_n - \mathbf{w}^{(\tau)\top} \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)$$

$$\nabla E_W = \mathbf{w}^{(\tau)}$$

Then we can calculate the error by using:

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N_V}$$

### 3. Process Data

In the code below, we do Partition to decide the whole dataset into three parts. The first one is used to train our model, so that we could generate the optimal parameter set as expected. The second part is reserved as test dataset, with which we could evaluate the efficiency of our model. 80% of the dataset is used as train data, 10% of the dataset is used as valid data, and the rest 10% is used for test. In the dataset, there are 69623 query-document pairs each consisting of 46 features. The last part is called valid part, which is used as a confirmation that the model we build by training the train part dataset is reliable.

```

def GenerateTrainingTarget(rawTraining, TrainingPercent = 80):
    TrainingLen = int(math.ceil(len(rawTraining)*(TrainingPercent*0.01)))
    t = rawTraining[:TrainingLen]
    #print(str(TrainingPercent) + "% Training Target Generated..")
    return t

#In this function , we split 80% of my input data for training as Training data and save it into d2
def GenerateTrainingDataMatrix(rawData, TrainingPercent = 80):
    T_len = int(math.ceil(len(rawData[0])*0.01*TrainingPercent))
    d2 = rawData[:,0:T_len]
    #print(str(TrainingPercent) + "% Training Data Generated..")
    return d2

#In this function, we split 10% of my input data for validation as validation data and save it into dataMatrix
def GenerateValData(rawData, ValPercent, TrainingCount):
    valSize = int(math.ceil(len(rawData[0])*ValPercent*0.01))
    V_End = TrainingCount + valSize
    dataMatrix = rawData[:,TrainingCount+1:V_End]
    #print(str(ValPercent) + "% Val Data Generated..")
    return dataMatrix

```

In code below, we reading the dataset we need to use to do LeToR. “Querylevelnorm\_x” is the input dataset with 46 features. “Querylevelnorm\_t” is the target dataset which has the expected result of the input dataset.

```

In [4]: #Reading the target data and output data
RawTarget = GetTargetVector('Querylevelnorm_t.csv')
RawData = GenerateRawData('Querylevelnorm_X.csv',IsSynthetic)

```

## 4. Closed Form Solution

Closed form solution is calculated using basis function, Sigma, lambda and M.

$$\mathbf{w}_{ML} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

In the code below, first we pick up M (which is 10 in this case) basis functions, and set centers and spreads of them. This process could also be accomplished by K-means function.

```

In [8]: ErmsArr = []
AccuracyArr = []
#Create M number of cluster.
kmeans = KMeans(n_clusters=M, random_state=0).fit(np.transpose(TrainingData))
#Set Mu as the center of the cluster
Mu = kmeans.cluster_centers_

```

Then implement the closed form solution

```
def GetWeightsClosedForm(PHI, T, Lambda):
    Lambda_I = np.identity(len(PHI[0]))
    for i in range(0,len(PHI[0])):
        Lambda_I[i][i] = Lambda
    PHI_T = np.transpose(PHI)
    PHI_SQR = np.dot(PHI_T,PHI)
    PHI_SQR_LI = np.add(Lambda_I,PHI_SQR)
    PHI_SQR_INV = np.linalg.inv(PHI_SQR_LI)
    INTER = np.dot(PHI_SQR_INV, PHI_T)
    W = np.dot(INTER, T)
    ##print ("Training Weights Generated..")
    return W
```

Then we are able to calculate the accuracy for the training, testing and validation

accuracy by applying the Erms function  $ERMS = \sqrt{2E(\mathbf{w}^*)/N_V}$

```
def GetErms(VAL_TEST_OUT,ValDataAct):
    sum = 0.0
    t=0
    accuracy = 0.0
    counter = 0
    val = 0.0
    for i in range (0,len(VAL_TEST_OUT)):
        sum = sum + math.pow((ValDataAct[i] - VAL_TEST_OUT[i]),2)
        if(int(np.around(VAL_TEST_OUT[i], 0)) == ValDataAct[i]):
            counter+=1
    accuracy = (float((counter*100))/float(len(VAL_TEST_OUT)))
    ##print ("Accuracy Generated..")
    ##print ("Validation E_RMS : " + str(math.sqrt(sum/len(VAL_TEST_OUT))))
    return (str(accuracy) + ',' + str(math.sqrt(sum/len(VAL_TEST_OUT))))
```

```
In [10]: TR_TEST_OUT = GetValTest(TRAINING_PHI,W)
VAL_TEST_OUT = GetValTest(VAL_PHI,W)
TEST_OUT = GetValTest(TEST_PHI,W)

TrainingAccuracy = str(GetErms(TR_TEST_OUT,TrainingTarget))
ValidationAccuracy = str(GetErms(VAL_TEST_OUT,ValDataAct))
TestAccuracy = str(GetErms(TEST_OUT,TestDataAct))
```

## 5. Stochastic Gradient Descent Solution

Then we define how dose SGD works in the dataset, firstly, we define the Delta\_E\_D, implement equation  $\Delta \mathbf{w}^{(\tau)} = -\eta^{(\tau)} \nabla E$ . In which, we also get  $\Delta \mathbf{w}^{(\tau)} = -\eta^{(\tau)} \nabla E$ , which is called weight update, it goes against with direction of gradient of the error. As well as delta\_E equation, then we are able to update the value of w.

```

#print ('-----Iteration: ' + str(i) + '-----')
#Implement the equation 11
Delta_E_D = -np.dot((TrainingTarget[i] - np.dot(np.transpose(W_Now),TRAINING_PHI[i])),TRAINING_PHI[i])
#Implement  $\Delta E = \nabla E$ 
La_Delta_E_W = np.dot(La,W_Now)
#Implement  $\nabla E = \nabla E_D + \lambda \nabla E$ 
Delta_E = np.add(Delta_E_D,La_Delta_E_W)
#Implement  $\Delta w(\tau) = -\eta(\tau)\nabla$ 
Delta_W = -np.dot(learningRate,Delta_E)

#Implement equation  $w(\tau+1) = w(\tau) + \Delta w(\tau)$ 
W_T_Next = W_Now + Delta_W
W_Now = W_T_Next

```

Then, we can print out the Gradient Descent result.

```

#-----TrainingData Accuracy-----#
TR_TEST_OUT = GetValTest(TRAINING_PHI,W_T_Next)
Erms_TR = GetErms(TR_TEST_OUT,TrainingTarget)
L_Erms_TR.append(float(Erms_TR.split(',')[1]))

#-----ValidationData Accuracy-----#
VAL_TEST_OUT = GetValTest(VAL_PHI,W_T_Next)
Erms_Val = GetErms(VAL_TEST_OUT,ValDataAct)
L_Erms_Val.append(float(Erms_Val.split(',')[1]))

#-----TestingData Accuracy-----#
TEST_OUT = GetValTest(TEST_PHI,W_T_Next)
Erms_Test = GetErms(TEST_OUT,TestDataAct)
L_Erms_Test.append(float(Erms_Test.split(',')[1]))

: print ('-----Gradient Descent Solution-----')
print ("M = 15 \nLambda = 0.0001\neta=0.01")
print ("E_rms Training = " + str(np.around(min(L_Erms_TR),5)))
print ("E_rms Validation = " + str(np.around(min(L_Erms_Val),5)))
print ("E_rms Testing = " + str(np.around(min(L_Erms_Test),5)))

```

## 6. Hyper parameter Tuning

Comparison of different  $M$  and  $\lambda$ :

Validation root mean square for LeToR:

M	0.01	0.05	0.1	0.15	0.25	0.35	0.5
2	0.5529503322 26	0.5529519154 91	0.5529538887 16	0.5529558554 76	0.5529597697. 21	0.5529636587 54	0.55296544454 71
4	0.5529503322 26	0.5529519154 91	0.5529538887 16	0.5529558554 76	0.55295976978	0.5529636587 54	0.55296944547 1
6	0.5529503322 26	0.5529519154 91	0.5529538887 16	0.5529558554 76	0.55295976978	0.5529636587 54	0.55296944547 1
8	0.5529503322 26	0.5529519156 45	0.5529531287 78	0.5529558554 76	0.55295976978	0.5529636587 54	0.55296944547 1
10	0.5529503322 26	0.5529519154 91	0.5529538887 16	0.5529558554 76	0.55295976978	0.5529636587 54	0.55296944547 1

1 2	0.5529503322 26	0.5529519154 91	0.5529538887 16	0.5529558554 75	0.55295976978	0.5529636857 54	0.55296944531 1
1 4	0.5529503322 26	0.5529519154 91	0.5529538887 16	0.5529558554 76	0.55295976978	0.5529636587 54	0.55296944547 1
1 6	0.5529503322 26	0.5529519154 91	0.5529538887 16	0.5529558554 76	0.55295144697 8	0.5529636587 54	0.55296944547 1

*SGD validation root mean square error for LeToR:*

M λ	0.01	0.05	0.1	0.15	0.25	0.35	0.5
2	0.5668693121 51	0.5655659527 83	0.5647398346 12	0.564406149085	0.5647814258 73	0.566101503304	0.569100092 71
4	0.5673462313 97	0.5655590382 3	0.5647401975 52	0.564435132955	0.5647814258 67	0.566101503304	0.569100092 71
6	0.5666399269 89	0.5659266141 86	0.5647396514	0.564406097985	0.5647814258 68	0.566101503304	0.569100092 71
8	0.5667967872 82	0.5655722871 26	0.5647403044 95	0.564406141412	0.5647814258 37	0.566101503304	0.569100092 71
10	0.5670654660 21	0.5655610197 34	0.5647404389 79	0.564406102586 54	0.5647814258 28	0.566101503304	0.569100092 71
12	0.5670243441 47	0.5655641570 65	0.5647398470 65	0.564406154582	0.5647814258 48	0.566101503306 54	0.569100092 71
14	0.5666152820 05	0.5655530230 72	0.5647394050 78	0.564406106722 25	0.5647814258 7	0.566101503304	0.569100092 71
16	0.5667542992 09	0.5655477306 39	0.5647402942 15	0.564406178662 5	0.5647814258 37	0.566101503304	0.569100092 71

According to the result, the best M and  $\lambda$  setting are M= 12,  $\lambda=0.05$

*Comparison of different learning rate setting if M and lambda is fixed:*

$\eta$	0.01	0.02	0.05	0.1	0.2	0.3
Erms	0.56901431 24	0.56516499 51	0.57305260 18	0.60680180 19	0.63672282 18	0.64789114 16

Hence, when learning rate is 0.01, we will have the minimum root mean square value.

