# CSE474/574: Introduction to Machine Learning(Fall 2018)
## Project 1.2: Learning to Rank using Linear Regression

## Objective

The goal of this project is to use machine learning to solve a problem that arises in Information Retrieval, one known as the Learning to Rank (LeToR) problem. We formulate this as a problem of linear regression where we map an input vector x to a real-valued scalar target y(x;w).
There are two tasks:
1. Train a linear regression model on LeToR dataset using a closed-form solution.
2. Train a linear regression model on the LeToR dataset using stochastic gradient descent (SGD).

## Introduction

Notations included: $x^i$ and $y^i$ to denote (i) th input variable and target variable respectively. So, the pair $(x^i, y^i)$ denotes the (i) th training example. Each training example consists of n features.

When the target variable we have to predict a continuous real value, the learning problem is called Regression. It is an example of supervised learning problem. In this project, we will develop linear model for Learning to Rank problem and use Gradient Descent to learn parameters of the model and also derive a closed form solution for the parameters.

## Linear Regression

In linear regression, a linear hypothesis function $h\theta(x)$ which approximates target variable y.

$$y \approx h_\theta(x) = \theta_0 + \theta_1 x_1 + ... + \theta_n x_n$$

Here θ's are the parameters of the model. When the context of parameters is clear, we can drop θ in $h\theta(x)$.

To simplify the notation, we define $x_0 = 1$, so

$$h(x) = \sum_{i=1}^{n} \theta_i x_i = \theta^T x$$

Job of our learning problem is to learn these parameters. Obvious method is to choose parameters which make h(x) as close to y as possible for the training examples provided to us.

Cost function:

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

Arooshi Avasthy                    arooshia@buffalo.edu

**Gradient Descent**

We want to choose θ which minimizes J(θ). A general strategy would be to start with some random θ and keep changing θ to reduce our objective function J(θ).

More specifically, Gradient descent starts with some initialization of θ and repeatedly performs the following update until convergence:

$$\theta_i := \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i} \qquad \forall i \in [0, 1, 2, 3, .., n]$$

Here, α is the learning rate which is a hyperparameter and we will have to tune it. In order to implement this algorithm, we will first have to calculate the partial derivative.

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

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

This method looks at all training examples for every iteration of the loop.

Closed Form Solution

Another method to find parameters is to set derivative of objective function to zero and obtain required parameters without resorting to iterative algorithm. Before doing this, we quickly introduce notations for matrix.

Given a training set, we define the design matrix X to be m × n matrix as

$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

Similarly, let Y be the m dimensional vector containing the target labels as

$$Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

Writing J(θ) in matrix representation, we get

$$J(\theta) = \frac{1}{2}(X\theta - Y)^T(X\theta - Y)$$

Finally, we take derivative of J(θ) with respect to θ and set it to zero.

$$\nabla_\theta J(\theta) = X^T(X\theta - Y) = 0$$

Arooshi Avasthy                                                         arooshia@buffalo.edu

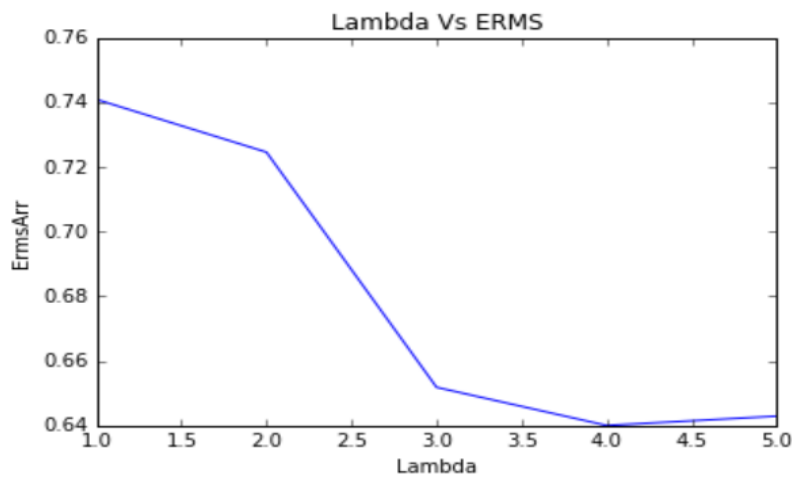Thus, the value of θ which minimizes J(θ) is given in closed form by

$$\theta = (X^T X)^{-1} X^T Y$$

**Experiments with the parameters for Closed form Solution:**

**1.**

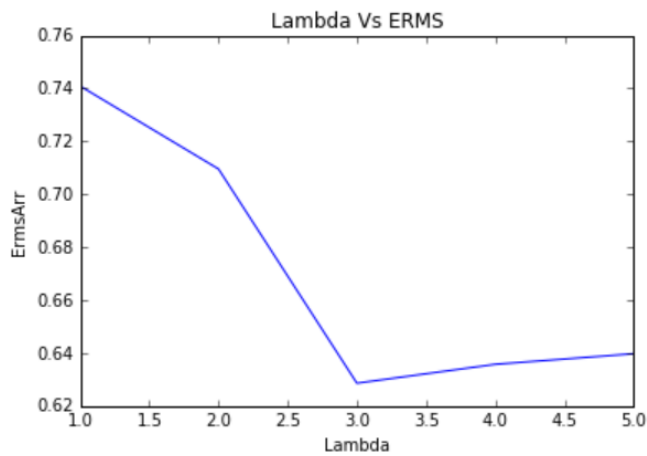| No. of Clusters (M) | Lambda | ERMS Training | ERMS Testing | ERMS Validation | Accuracy Testing |
|---|---|---|---|---|---|
| 1 | 0.1 | 0.64275290353 00575 | 0.74088824104 61395 | 0.62826281331 69708 | 70.234161758 36805 |
| 1 | 2 | 0.62917269896 97519 | 0.72456738898 95873 | 0.61452473788 28577 | 69.716994684 67175 |
| 1 | 30 | 0.57394187635 80454 | 0.65183394906 40811 | 0.56265467089 19288 | 70.234161758 36805 |
| 1 | 400 | 0.56470411541 15904 | 0.64003193013 10835 | 0.55394450441 77406 | 70.234161758 36805 |
| 1 | 5000 | 0.56498437307 40433 | 0.64291194219 89584 | 0.55382704613 6877 | 70.234161758 36805 |

ERMS Vs Lambda Graph depicting the above readings:

**2.**

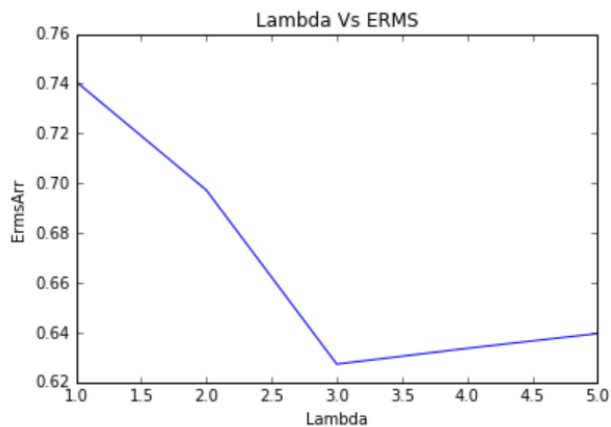| No. of Clusters (M) | Lambda | ERMS Training | ERMS Testing | ERMS Validation | Accuracy Testing |
|---|---|---|---|---|---|
| 11 | 0.1 | 0.64275289843 52141 | 0.74088823868 22628 | 0.62826281254 86431 | 70.234161758 36805 |
| 11 | 2 | 0.61491663260 86989 | 0.70955426186 21802 | 0.60159415098 41954 | 68.840683809 79744 |
| 11 | 30 | 0.55218914663 86234 | 0.62861122105 4567 | 0.54247018908 60169 | 70.119235741 99109 |
| 11 | 400 | 0.55996308184 61021 | 0.63575508886 41107 | 0.54880958250 78251 | 70.234161758 36805 |
| 11 | 5000 | 0.56440054051 1381 | 0.63974518367 61452 | 0.55366646570 38907 | 70.234161758 36805 |

ERMS Vs Lambda Graph depicting the above readings:



**3.**

| No. of Clusters (M) | Lambda | ERMS Training | ERMS Testing | ERMS Validation | Accuracy Testing |
|---|---|---|---|---|---|
| 21 | 0.1 | 0.64275285630 8148 | 0.74088820234 25261 | 0.62826274301 25443 | 70.234161758 36805 |
| 21 | 2 | 0.60663062470 51635 | 0.69722471512 20713 | 0.59519286417 6806 | 68.826318057 75033 |
| 21 | 30 | 0.54944650373 68979 | 0.62731791567 00758 | 0.54049847842 63781 | 69.616434420 3419 |
| 21 | 400 | 0.55725110939 963 | 0.63364096095 02813 | 0.54595261283 92471 | 70.234161758 36805 |

Arooshi Avasthy                                                                                    arooshia@buffalo.edu

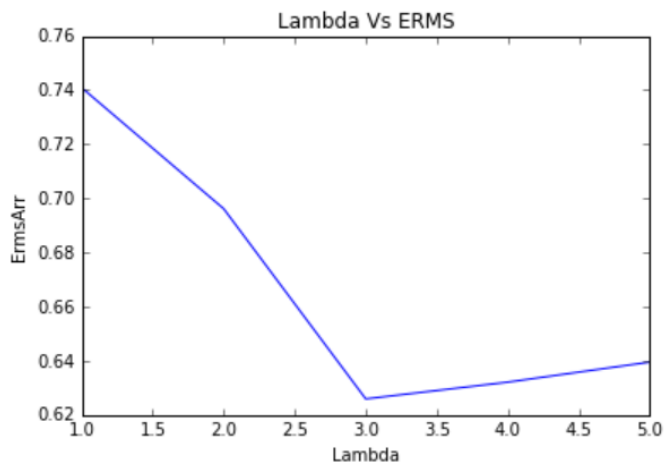| 11 | 5000 | 0.56438799793 52604 | 0.63959611429 63466 | 0.55367535515 56489 | 70.234161758 36805 |
|---|---|---|---|---|---|

ERMS Vs Lambda Graph depicting the above readings:



**4.**

| No. of Clust ers (M) | Lamb da | ERMS Training | ERMS Testing | ERMS Validation | Accuracy Testing |
|---|---|---|---|---|---|
| 31 | 0.1 | 0.64275264564 27277 | 0.74088770733 28477 | 0.62826264814 8756 | 70.234161758 36805 |
| 31 | 2 | 0.60459790580 74382 | 0.69617269520 65399 | 0.59306651147 95357 | 68.768855049 56185 |
| 31 | 30 | 0.54688934995 53777 | 0.62594655777 63952 | 0.53937411486 5649 | 69.587702916 24766 |
| 31 | 400 | 0.55553169530 32732 | 0.63954494373 46774 | 0.55368285593 6783 | 70.234161758 36805 |
| 31 | 5000 | 0.56438508074 73036 | 0.63213221838 8237 | 0.54432967732 99638 | 70.234161758 36805 |

ERMS Vs Lambda Graph depicting the above readings:



Graph Depicting ERMS Vs No.of clusters:

Arooshi Avasthy                                                                                              arooshia@buffalo.edu

**Experiments with the parameters for Gradient Descent:**

.

| No. of Data points | Lambda | Learning Rate | ERMS Training | ERMS Validation | ERMS Testing |
|---|---|---|---|---|---|
| 100 | 1 | 0.1 | 0.56255 | 0.55183 | 0.63338 |
| 200 | 2 | 0.02 | 0.56255 | 0.55183 | 0.63338 |
| 300 | 3 | 0.003 | 0.56255 | 0.55183 | 0.63338 |
| 400 | 4 | 0.0004 | 0.56255 | 0.55183 | 0.63338 |
| 500 | 5 | 0.00005 | 0.56255 | 0.55183 | 0.63338 |

Arooshi Avasthy                                                                            arooshia@buffalo.edu