

## TIM245 Lecture 8 (4/26/17)

Agenda

- 1) Roadmap and schedule for the quarter
- 2) Review Lecture 7
- 3) Multiple Linear Regression
- 4) Model training
- 5) Extension to linear regression: ridge, lasso, elastic nets

# ① Roadmap for the course

Schedule for the next seven weeks

Week	Lecture	Homework	Project	Exams
Week 4 4/24	Supervised Prediction	HW 1 Due	Phase II Due	
Week 5 5/1/17	Supervised Classification	HW 2 Assigned	Phase II Assigned	
Week 6 5/8	Unsupervised Clustering	HW 2 due	Phase III due	Midterm Assigned
Week 7 5/15	Unsupervised Association			Midterm Due
Week 8 5/22	Text Mining	HW 3 Assigned	Phase IV Assigned	
Week 9 5/29	Time Series Analysis	HW 3 Due	Phase IV due	
Week 10 6/5	Graph Mining		Phase V assigned	Final assigned
Finals 6/12			Phase V due	Final Due

### ③ Multiple Linear Regression

Multiple linear regression: linear regression with multiple independent variables or attributes and a single dependent variable or target

(Multivariate linear regression refers to multiple dependent variables)

$$\hat{y} = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im}$$

The values  $\beta_0, \beta_1, \beta_2 \dots \beta_m$  are the linear regression model.

The general process or learning algorithm for creating the model

1) Model Training: estimate coefficients that minimize squared error (least squared error) on the training data

2) Model Testing: evaluate coefficients on the test data set

### ④ Model Training

How do we find the best fit line for the training data set

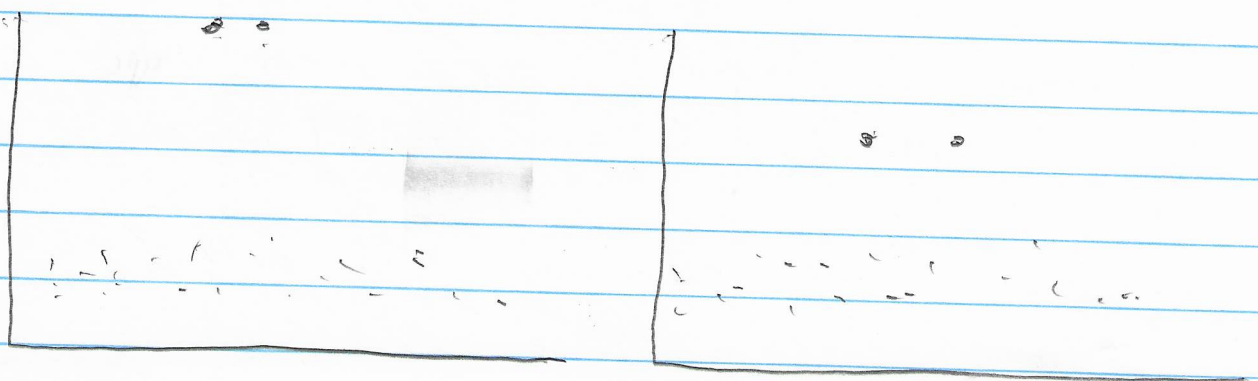


$A_1$



Squared Error  $(y - \hat{y})^2$

Absolute Error  $|y - \hat{y}|$



Two general ways to find the best fit line:

1) Minimize squared error (L2 Norm)

$$RSS(\beta) = \sum_{i=1}^n (y_i - \beta X_i)^2$$

2) Minimize absolute error (L1 Norm)

$$AE(\beta) = \sum_{i=1}^n |y_i - \beta X_i|$$

Where

$$\beta = [\beta_1, \beta_2, \dots, \beta_m]$$

$$X_i = [X_{i1}, X_{i2}, \dots, X_{im}]$$

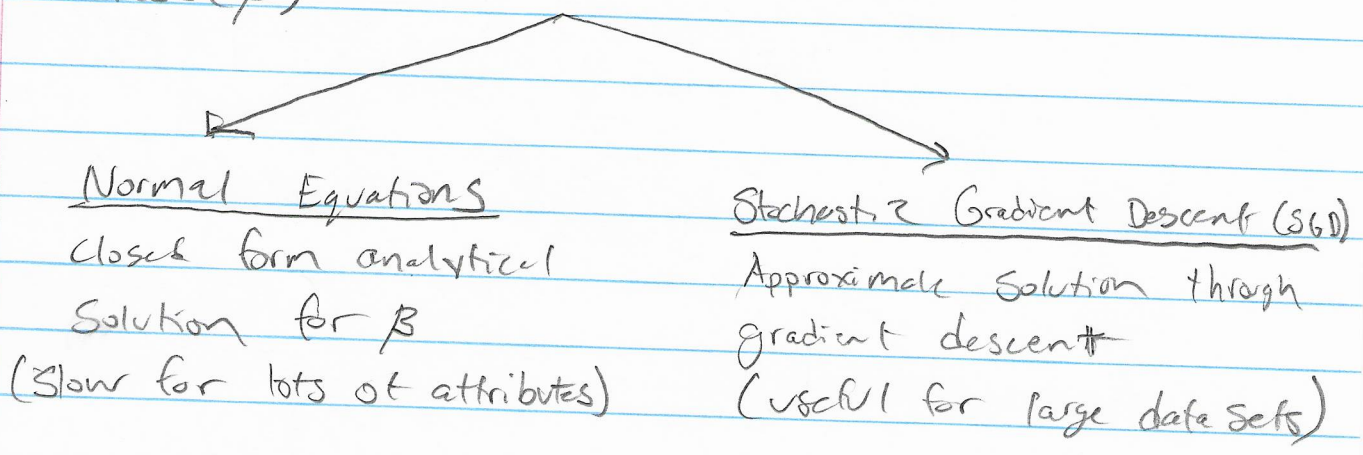
$$\beta X_i = \hat{y}_i$$

L2 Norm loss function: not robust, stable, one solution

L1 Norm loss function: robust, not stable, multiple solution

L2 Norm loss function typically used because of computational efficiency

Two general approaches to minimize  $RSS(\beta)$



Normal Equations:

$$\begin{aligned} \text{RSS}(\beta) &= \sum_{i=1}^n (y_i - \beta x_i)^2 \\ &= (y - \beta X)^T (y - \beta X) \end{aligned}$$

where

Training Data Set

$A_1$	$A_2$	...	$Y$
	X		Y

$\beta_0$  placeholder

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1m} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix}$$

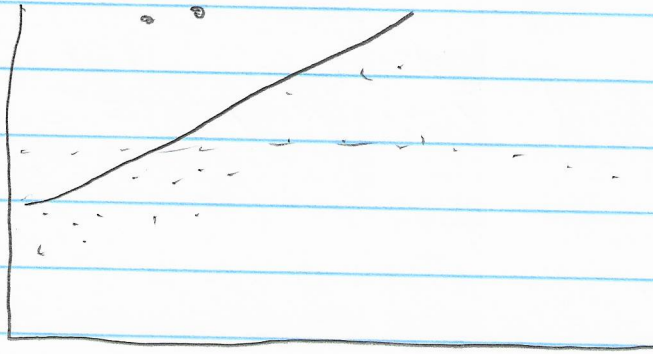
$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Differentiate with respect to  $\beta$ ,  
Set equal to 0, and solve for  
 $\beta$

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



## ⑤ Regularization: Ridge, Lasso, and Elastic Net



Least Squares  
is an unbiased  
estimator,  
Sometimes adding  
bias can help

Regularization is an important tool  
that allows us to prevent overfitting

Three main types of regularization  
used with linear regression models

- 1) Ridge Regression: Coefficient Shrinkage
- 2) Lasso Regression: Coefficient Shrinkage + Feature Selection
- 3) Elastic Net: Combination of ridge and lasso

Ridge regression Controls the variance by imposing a constraint on the coefficients

Minimize  $RSS(\beta)$

$$\text{s.t.}, \sum_{j=1}^m \beta_j^2 < t \quad (\text{L2 Norm})$$

Where  $t$  is a user defined value called the ridge parameter

This leads to a penalized loss function

$$RSS(\beta) + \lambda \sum_{j=1}^m \beta_j^2$$

Which has the closed form solution

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y$$

$\lambda$  controls the regularization, i.e. the size of the coefficients

$\lambda \rightarrow 0$  : Least Squares

$\lambda \rightarrow \infty$  :  $\beta_0$  (intercept only)

Least Absolute Shrinkage and Selection Operator (Lasso) allows  $\beta_j = 0$ , i.e. feature selection

Minimize  $RSS(\beta)$

s.t.  $\sum_{j=1}^m |\beta_j| < t$  (L1 Norm)

this leads to the penalized loss function

$$RSS(\beta) + \lambda \sum_{j=1}^m |\beta_j|$$

No closed form solution

Large value for  $\lambda$  will set some coefficients to 0, i.e. feature selection