

**CSE 626:  
DATA MINING INTRO:  
Part 2**

# Knowledge Discovery in Databases-- Steps

- Selecting Target Data
- Preprocessing the Data
- Transforming them
- Data Mining to extract Patterns and Relationships
- Interpreting/Assessing Discovered Structures

Course  
Focus:  
Algorithms  
for Modeling  
the data

# Process of Data Mining

1. Representation
2. Score function
3. Process to optimize score
4. Implementation:  
data management,  
efficiency

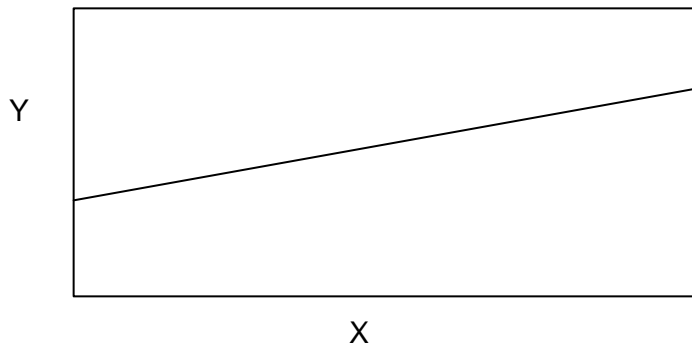
## EXAMPLE

1. Regression:  
 $Y = a + bX$   
Predictor variable = X  
(income)  
Response variable = Y  
(credit card spending)
2. Score: sum of squared errors

# Linear Regression Process

Linear regression with one variable

Data of the form  $(x_i, y_i)$ ,  $i = 1, \dots, n$  samples



Database

<b>Y</b>	<b>X</b>
1	3
8	9
11	11
4	5
3	2

# Score: sum of squared errors

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - y_i')^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

Where  $y_i'$  is the response value obtained from the model

**We wish to minimize SSE.**

# Minimizing SSE for Regression

Differentiating SSE with respect to  $a$  and  $b$  we have

$$\frac{\partial(SSE)}{\partial a} = -2 \sum_{i=1}^n (y_i - a - bx_i)$$

$$\frac{\partial(SSE)}{\partial b} = -2 \sum_{i=1}^n (y_i - a - bx_i)x_i$$

Setting partial derivatives equal to zero and rearranging terms we have equations

$$na + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$

$$a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i$$

Which we solve for  $a$  and  $b$ ,  
the *regression coefficients*

# Regression Coefficients

$$b = \frac{\sum_{i=1}^n (x_i - \text{mean}_x) \cdot (y_i - \text{mean}_y)}{\sum_{i=1}^n (x_i - \text{mean}_x)^2}$$

$$a = \text{mean}_y - b \cdot \text{mean}_x$$

# Application to database

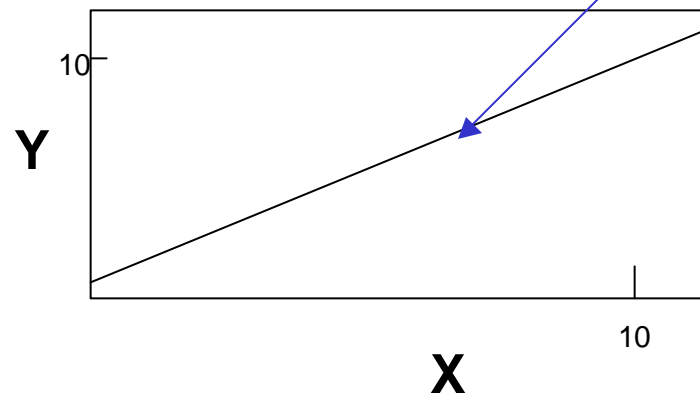
Y	X
1	3
8	9
11	11
4	5
3	2

$$\text{mean}_y = 5$$

$$\text{mean}_x = 6$$

$$a = 0.8, b = 1.04$$

Optimal regression line is  
 $Y = 0.8 + 1.04X$





# Multiple Regression

$p$  predictor variables

$y$	$x_1$	$x_2$	.....	$x_p$
$y(1)$	$x_1(1)$			
$y(n)$	$x_1(n)$			

Solution:

$$a = (X^t X)^{-1} \cdot X^t y$$

$$y = a_0 + \sum_{i=1}^p a_j x_j$$

$n$  objects

$X = n$  by  $p+1$  matrix

Where a column of 1's are added to incorporate  $a_0$  in model

$$y = Xa + e$$

$y$  is a column vector,  $a = (a_0, \dots, a_p)$

$e$  is a  $n$  by  $1$  vector containing residuals

# Implementation of Regression

$$y = a_0 + \sum_{i=1}^p a_j x_j$$

Solution:

$$a = (X^t X)^{-1} \cdot X^t y$$

Simple summaries of the data; sums, sums of squares and sums of products of X and Y are sufficient to compute estimates of a and b

Implies single pass through the data will yield estimates