

# DS 4400

## Machine Learning and Data Mining I

Alina Oprea  
Associate Professor, CCIS  
Northeastern University

January 17 2019

# Logistics

- HW 1 is on Piazza and Gradescope
- Deadline: Friday, Jan. 25, 2019
- Office hours
  - Alina: Thu 4:30-6:00pm (ISEC 625)
  - Ewen: Mon 5:30-6:30pm (ISEC 605)
- How to submit HW
  - Create a PDF and submit on Gradescope before 11:59pm the day assignment is due
  - Submit zip of code in Google form
  - Should include ReadMe file on how to run code
  - Preferred: Use Jupyter notebook in R or Python

# Collaboration policy

- What is allowed
  - You can discuss the homework with your colleagues
  - You can post questions on Piazza and come to office hours
  - You can search for online resources to better understand class concepts
- What is not allowed
  - Sharing your written answers with colleagues
  - Sharing your code or receiving code from colleague
  - Do not use directly code from the Internet!

# Outline

- Terminology for supervised learning
- Multiple linear regression
  - Derivation of optimal model in matrix form
- Practical issues
  - Feature scaling and normalization
  - Outliers
  - Categorical variables
- Lab

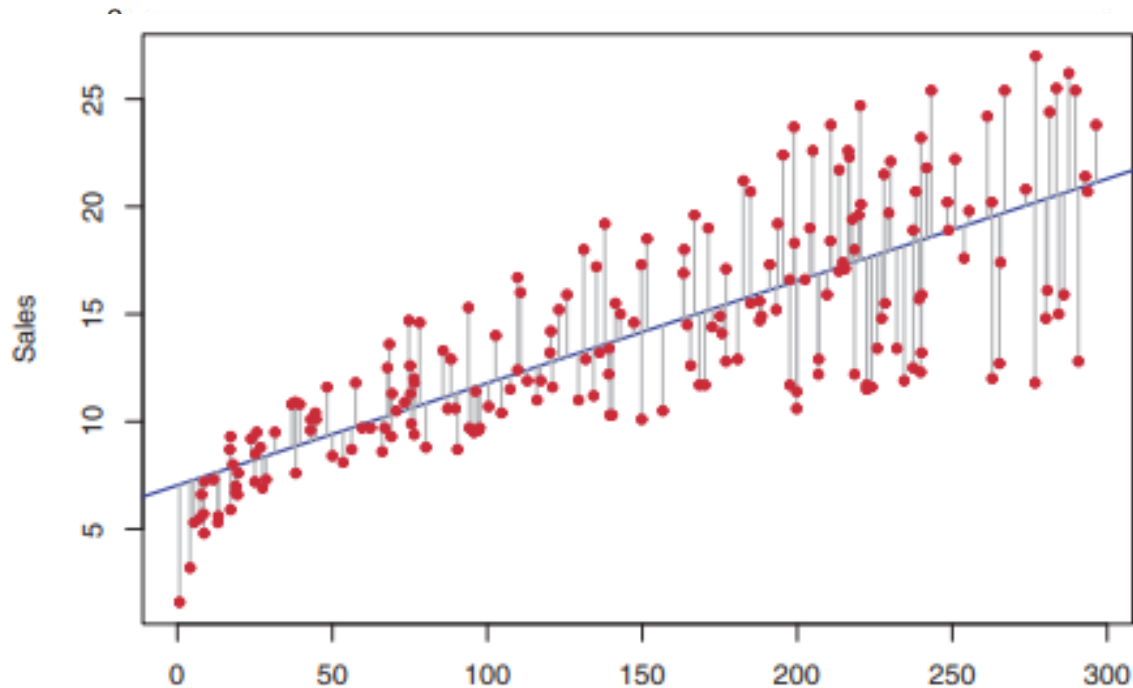
# Terminology

- **Hypothesis space**  $H = \{f: X \rightarrow Y\}$
- **Training data**  $D = (x_i, y_i) \in X \times Y$
- **Features**:  $x_i \in X$
- **Labels**  $y_i \in Y$ 
  - Classification: discrete  $y_i \in \{0,1\}$
  - Regression:  $y_i \in \mathbb{R}$
- **Loss function**:  $L(f, D)$ 
  - Measures how well  $f$  fits training data
- **Training algorithm**: Find hypothesis  $\hat{f}: X \rightarrow Y$ 
  - $\hat{f} = \operatorname{argmin}_{f \in H} L(f, D)$

# Linear regression

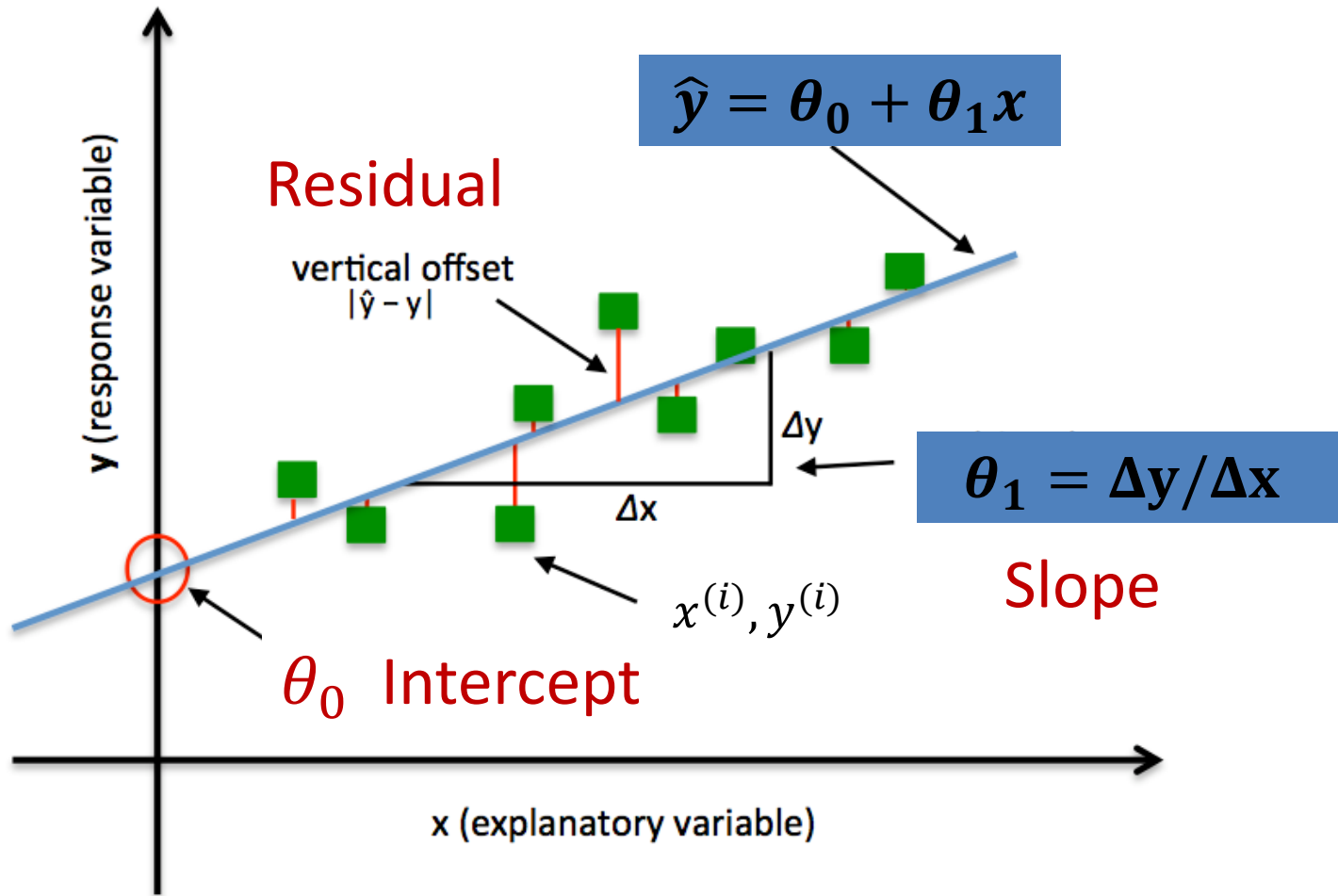
Given:

- Data  $\mathbf{X} = \{x^{(1)}, \dots, x^{(n)}\}$  where  $x^{(i)} \in \mathbb{R}^d$  **Features**
- Corresponding labels  $\mathbf{y} = \{y^{(1)}, \dots, y^{(n)}\}$  where  $y^{(i)} \in \mathbb{R}$  **Response variables**



Simple Linear Regression: 1 predictor

# Interpretation

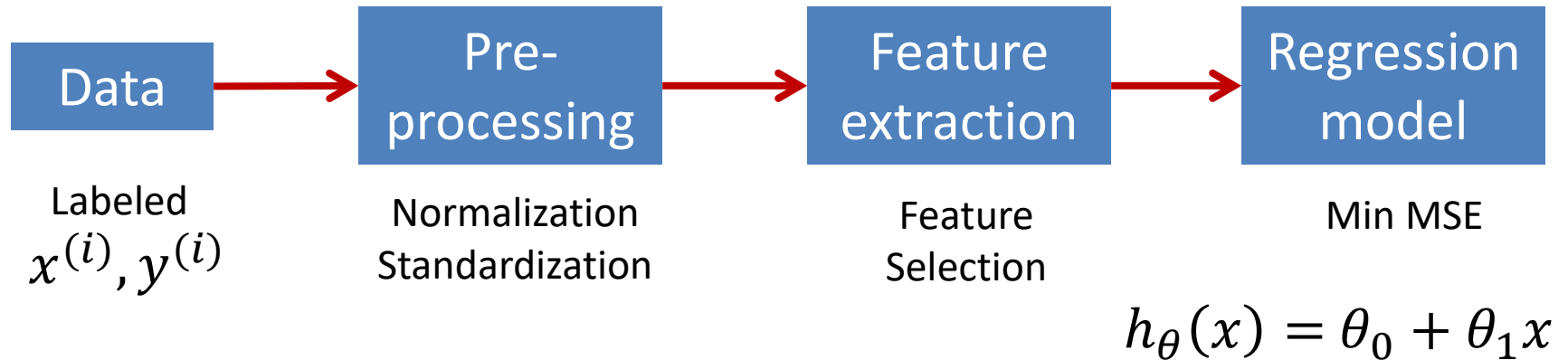


**Hypothesis:**  $h_{\theta}(x) = \theta_0 + \theta_1 x$

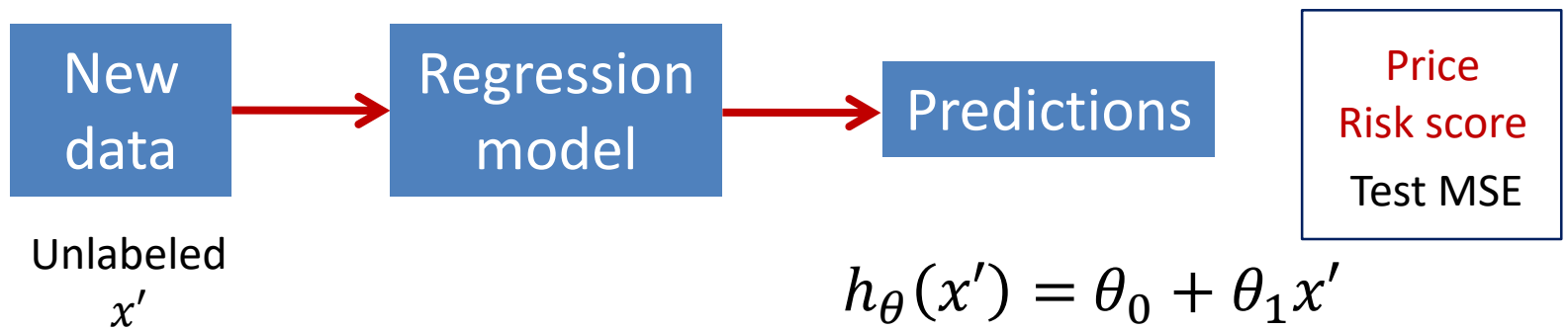
**Loss: MSE**  $= \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$

# Regression Learning

## Training



## Testing





# Simple Linear Regression

- Dataset  $x^{(i)} \in R, y^{(i)} \in R, h_{\theta}(x) = \theta_0 + \theta_1 x$

- $J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$

**MSE / Loss**

- Solution of min loss

$$-\theta_0 = \bar{y} - \theta_1 \bar{x}$$
$$-\theta_1 = \frac{\sum (x^{(i)} - \bar{x})(y^{(i)} - \bar{y})}{\sum (x^{(i)} - \bar{x})^2}$$

Variance of x

Co-variance of x and y

$$\bar{x} = \frac{\sum_{i=1}^n x^{(i)}}{n}$$
$$\bar{y} = \frac{\sum_{i=1}^n y^{(i)}}{n}$$

# How Well Does the Model Fit?

- Correlation between feature and response
  - Pearson's correlation coefficient

$$\text{Cor}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Co-variance of x and y

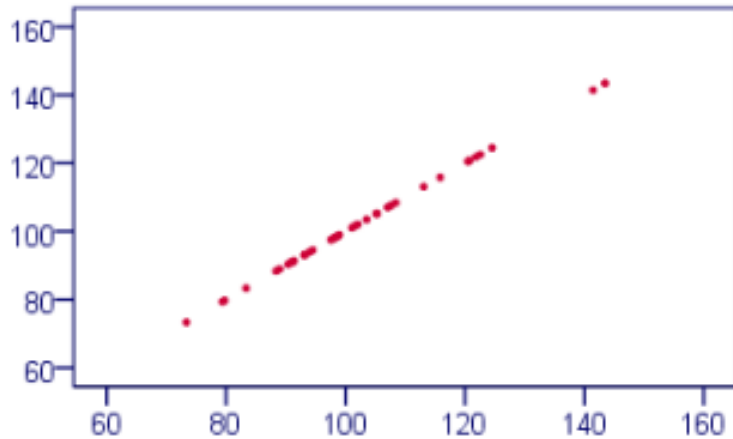
Standard deviation x

Standard deviation y

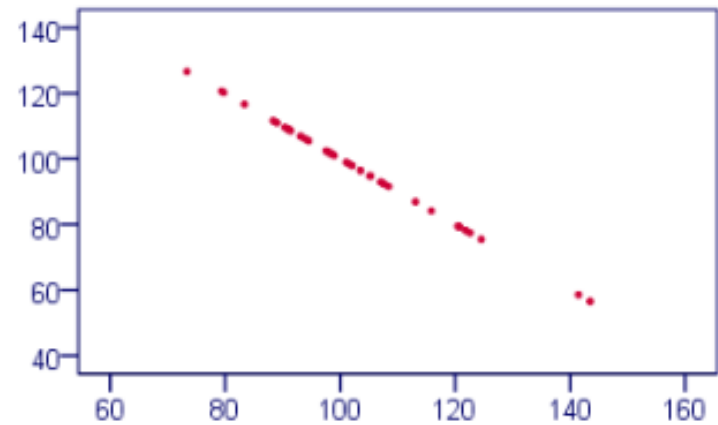
- Measures linear dependence between x and y
- Positive coefficient implies positive correlation
  - The closer to 1 the coefficient is, the stronger the correlation
- Negative coefficient implies negative correlation
  - The closer to -1 the coefficient is, the stronger the correlation

# Correlation Coefficient

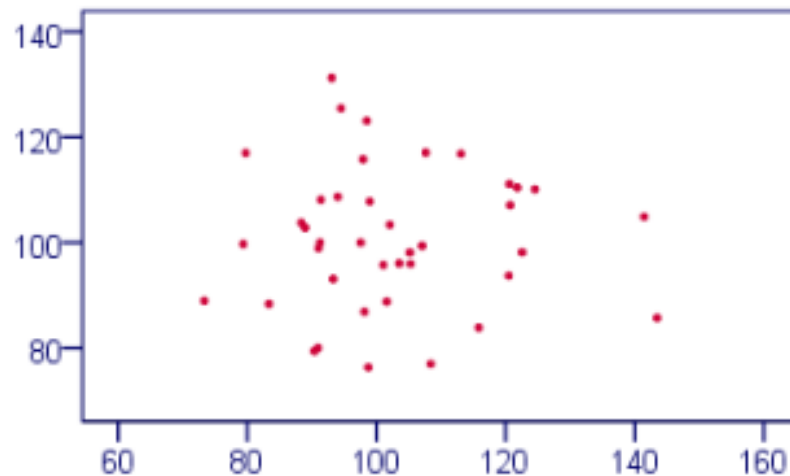
**Correlation Coefficient = 1**



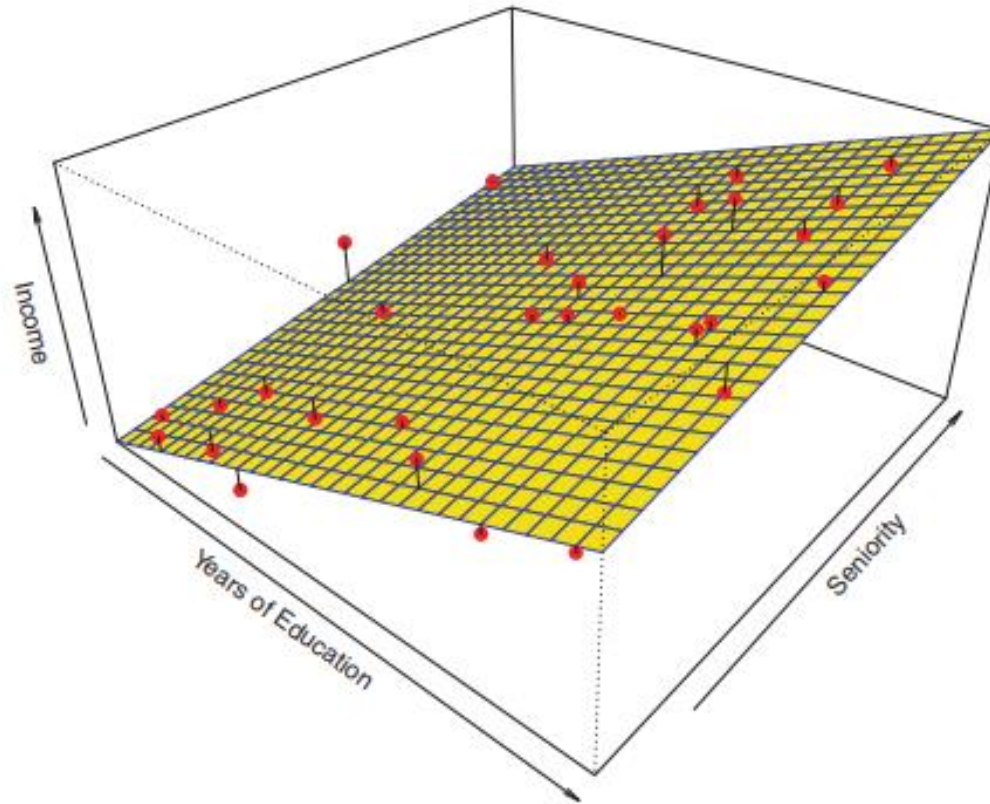
**Correlation Coefficient = -1**



**Correlation Coefficient = 0**

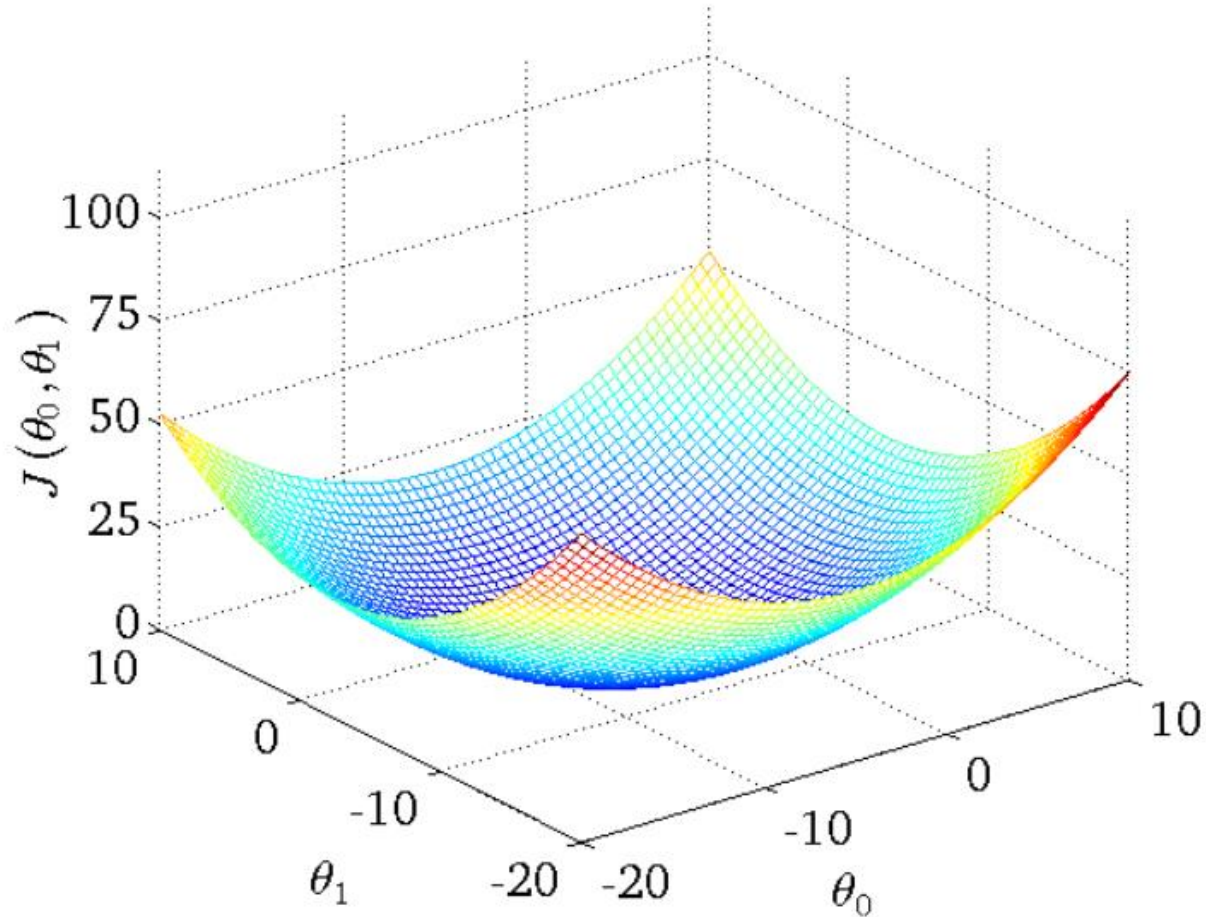


# Multiple Linear Regression



- Linear Regression with 2 predictors
- Dataset:  $x^{(i)} \in R^d, y^{(i)} \in R$

# MSE function



Convex function implies unique minimum

# Vector Norms

**Vector norms:** A norm of a vector  $\|x\|$  is informally a measure of the “length” of the vector.

$$\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$

– Common norms:  $L_1$ ,  $L_2$  (Euclidean)

$$\|x\|_1 = \sum_{i=1}^n |x_i| \quad \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

–  $L_{\text{infinity}}$

$$\|x\|_{\infty} = \max_i |x_i|$$

# Vector products

We will use lower case letters for vectors

The elements are referred by  $x_i$ .

- **Vector dot (inner) product:**

$$x^T y \in \mathbb{R} = [x_1 \quad x_2 \quad \cdots \quad x_n] \begin{bmatrix} y_1 \\ x_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i.$$

- **Vector outer product:**

$$xy^T \in \mathbb{R}^{m \times n} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} [y_1 \quad y_2 \quad \cdots \quad y_n] = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_m y_1 & x_m y_2 & \cdots & x_m y_n \end{bmatrix}$$

# Hypothesis Multiple LR

- Linear Model

- Consider our model:

$$h(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

- Can write the model in vectorized form as  $h(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$

Vector inner product



# Training data

	Feature 1		Feature d		
$\mathbf{X} =$	1	$x_1^{(1)}$	$\dots$	$x_d^{(1)}$	Training example $i$
	$\vdots$	$\vdots$	$\ddots$	$\vdots$	
	1	$x_1^{(i)}$	$\dots$	$x_d^{(i)}$	
	$\vdots$	$\vdots$	$\ddots$	$\vdots$	
	1	$x_1^{(n)}$	$\dots$	$x_d^{(n)}$	

$\mathbb{R}^{n \times (d+1)}$

- Total number of training example:  $n$
- Dimension of training data point (number of features):  $d$

# Use Vectorization

- Consider our model for  $n$  instances:

$$h(\mathbf{x}^{(i)}) = \sum_{j=0}^d \theta_j x_j^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)}$$

- Let

Model  
parameter

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$$\mathbb{R}^{(d+1) \times 1}$$

$\mathbf{X} =$

$$\begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}$$

$$\mathbb{R}^{n \times (d+1)}$$

Training  
data

- Can write the model in vectorized form as  $h_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{X}\boldsymbol{\theta}$

Model prediction vector  $\hat{\mathbf{y}}$

# Loss function MSE

- For the linear regression cost function:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left( h_{\theta} \left( \mathbf{x}^{(i)} \right) - y^{(i)} \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^n \left( \hat{y}^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{n} \|\hat{\mathbf{y}} - \mathbf{y}\|^2$$
$$= \frac{1}{n} \|\mathbf{X}\theta - \mathbf{y}\|^2$$

Let:

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(n)} \end{bmatrix}$$

# Matrix and vector gradients

If  $y = f(x)$ ,  $y \in R$  scalar,  $x \in R^n$  vector

$$\frac{\partial y}{\partial \mathbf{x}} = \left[ \frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \cdots \quad \frac{\partial y}{\partial x_n} \right]$$

Vector gradient  
(row vector)

If  $\mathbf{y} = f(\mathbf{x})$ ,  $\mathbf{y} \in R^m$ ,  $\mathbf{x} \in R^n$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

Jacobian  
matrix  
(Matrix  
gradient)

# Properties

- If  $w, x$  are  $(d \times 1)$  vectors,  $\frac{\partial w^T x}{\partial x} = w^T$
- If  $A: (n \times d)$   $x: (d \times 1)$ ,  $\frac{\partial Ax}{\partial x} = A$
- If  $A: (d \times d)$   $x: (d \times 1)$ ,  $\frac{\partial x^T Ax}{\partial x} = (A + A^T)x$
- If  $A$  symmetric:  $\frac{\partial x^T Ax}{\partial x} = 2Ax$
- If  $x: (d \times 1)$ ,  $\frac{\partial ||x||^2}{\partial x} = 2x^T$

# Min loss function

- Notice that the solution is when  $\frac{\partial}{\partial \theta} J(\theta) = 0$

$$J(\theta) = \frac{1}{n} \|X\theta - y\|^2$$

Using chain rule

$$f(\theta) = h(g(\theta)), \frac{\partial f(\theta)}{\partial \theta} = \frac{\partial h(g(\theta))}{\partial \theta} \frac{\partial g(\theta)}{\partial \theta}$$

$$h(x) = \|x\|^2, g(\theta) = X\theta - y$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{2}{n} [(X\theta - y)^T X] = 0 \Rightarrow X^T (X\theta - y) = 0$$

$$(X^T X)\theta = X^T y$$

Closed Form Solution:

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

# Vectorization

- Two options for operations on training data
  - Matrix operations
  - For loops to update individual entries
- Most software packages are highly optimized for matrix operations
  - Python numpy
  - Preferred method!
- Matrix operations are much faster than loops!

# Closed-form solution

- Can obtain  $\theta$  by simply plugging  $X$  and  $y$  into

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

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- If  $X^T X$  is not invertible (i.e., singular), may need to:
  - Use pseudo-inverse instead of the inverse
    - In python, `numpy.linalg.pinv(a)`
  - Remove redundant (not linearly independent) features
  - Remove extra features to ensure that  $d \leq n$

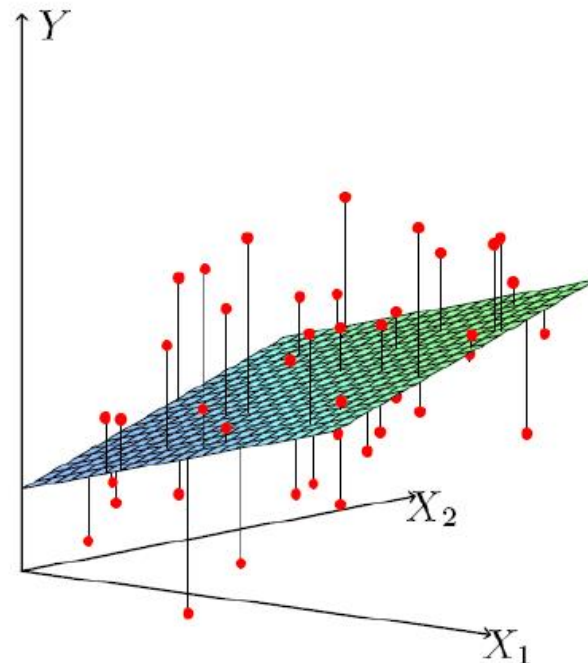
$$AGA = A$$



# Multiple Linear Regression

- Dataset:  $x^{(i)} \in R^d, y^{(i)} \in R$
- Hypothesis  $h_{\theta}(x) = \theta^T x$
- $\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2$  **Loss / cost**

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



# Feature Standardization

- Rescales features to have zero mean and unit variance

– Let  $\mu_j$  be the mean of feature j:  $\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$

– Replace each value with:

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j} \quad \text{for } j = 1 \dots d \quad \text{(not } x_0 \text{!)}$$

- $s_j$  is the standard deviation of feature j

- Must apply the same transformation to instances for both training and prediction
- **Mean 0 and Standard Deviation 1**

# Other feature normalization

- **Min-Max rescaling**

- $x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \min_j}{\max_j - \min_j} \in [0,1]$

- $\min_j$  and  $\max_j$ : min and max value of feature  $j$

- **Mean normalization**

- $x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\max_j - \min_j}$

- Mean 0

# Feature standardization/normalization

- Goal is to have individual features on the same scale
- Is a pre-processing step in most learning algorithms
- Necessary for linear models and Gradient Descent
- Different options:
  - Feature standardization
  - Feature min-max rescaling
  - Mean normalization

# Review

- Solution for multiple linear regression can be computed in closed form
  - Matrix inversion is computationally intense
  - We will discuss an efficient training algorithms (gradient descent)
- In practice several techniques can help generate more robust models
  - Outlier removal
  - Feature scaling

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
- Thanks!