# DS 4400

# Machine Learning and Data Mining I

Alina Oprea

Associate Professor, CCIS

Northeastern University

September 18 2018

# Logistics

- HW 1 is on Piazza and Gradescope

- Deadline: Friday, Sept. 28, 2018

- Office hours
  - Alina: Thu 4:30-6:00pm (ISEC 625)
  - Anand: Tue 2-3pm (ISEC 605)

- How to submit HW
  - Create a PDF and submit on Gradescope before 11:59pm the day assignment is due
  - Include link to code and ReadMe file
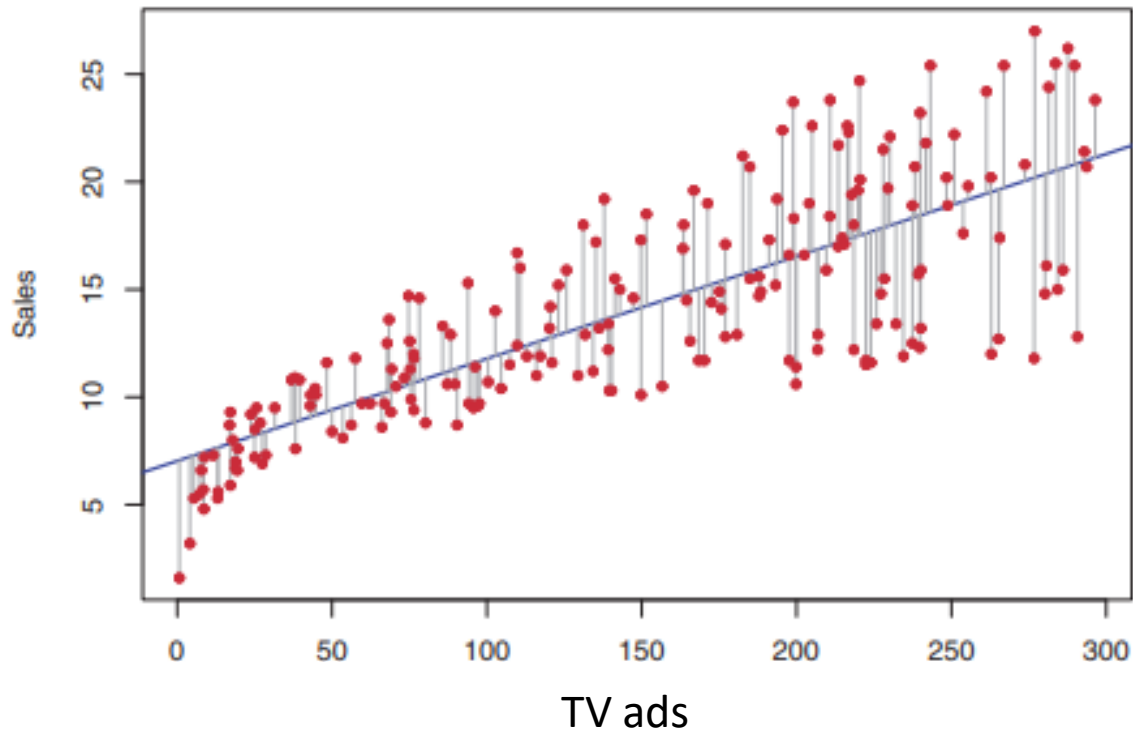  - Use Jupyter notebook in R or Python

# 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 code from the Internet!

# Linear regression

Given:

- Data $X = \left\{ x^{(1)}, \ldots, x^{(n)} \right\}$ where $x^{(i)} \in \mathbb{R}^d$    Features

- Corresponding labels $y = \left\{ y^{(1)}, \ldots, y^{(n)} \right\}$ where $y^{(i)} \in \mathbb{R}$

Response variables



Sales

TV ads

# Simple linear regression

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

- $J(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left(\theta_0 + \theta_1 x^{(i)} - y^{(i)}\right)^2$  <span style="color:red">loss</span>

  $\frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^{n} \left(\theta_0 + \theta_1 x^{(i)} - y^{(i)}\right) \quad = 0$

  $\frac{\partial J(\theta)}{\partial \theta_1} = \frac{2}{n} \sum_{i=1}^{n} x^{(i)} \left(\theta_0 + \theta_1 x^{(i)} - y^{(i)}\right) \quad = 0$
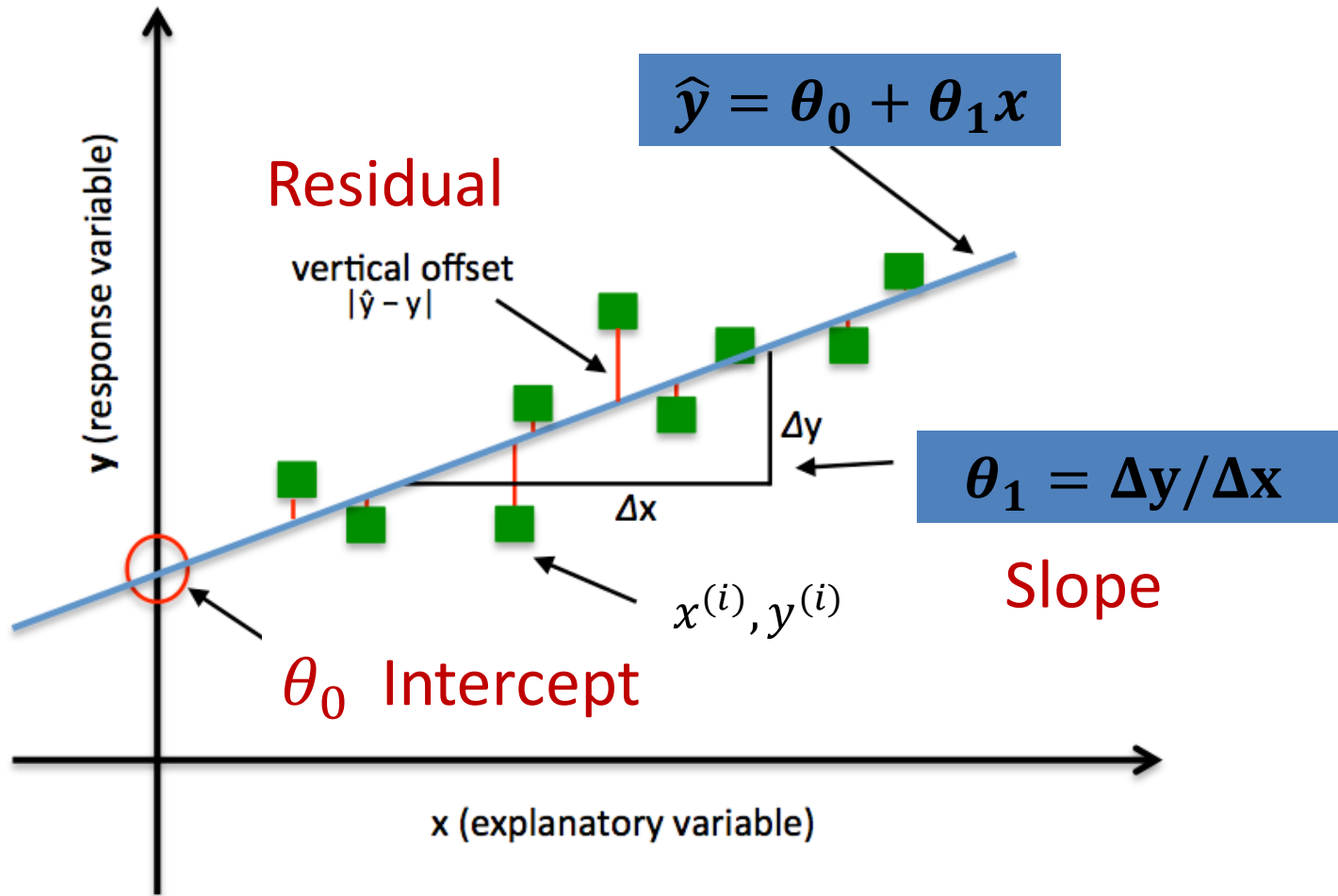
- Solution of min loss

$$- \theta_0 = \bar{y} - \theta_1 \bar{x}$$

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

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

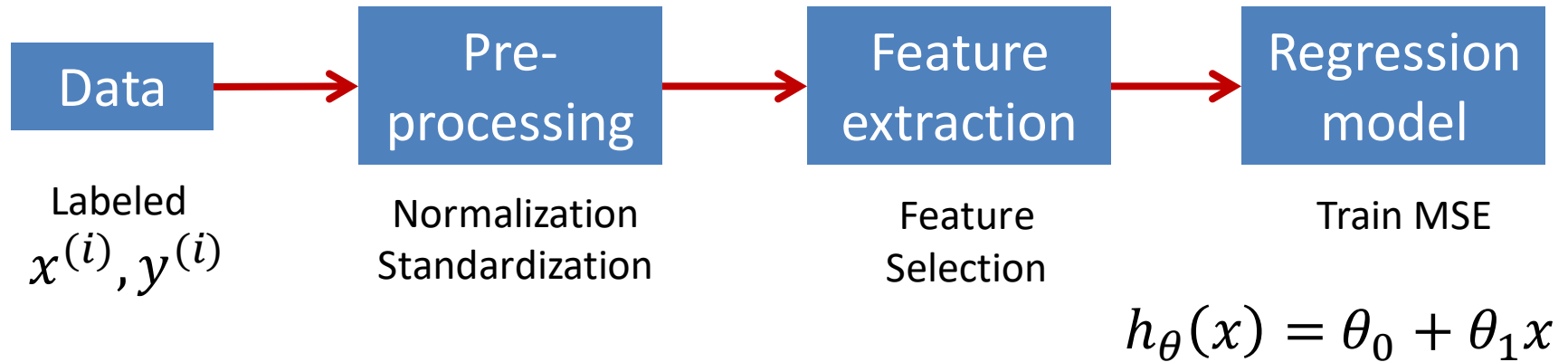$$\bar{y} = \frac{\sum_{i=1}^{n} y^{(i)}}{n}$$

# Interpretation



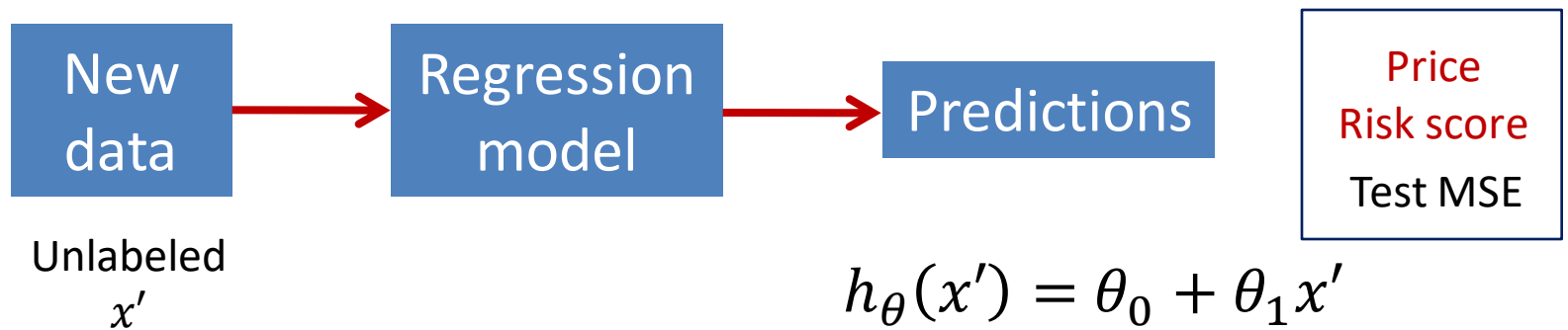$$h_\theta(x) = \theta_0 + \theta_1 x$$

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

# Regression Learning

**Training**

Data $\rightarrow$ Pre-processing $\rightarrow$ Feature extraction $\rightarrow$ Regression model

Labeled
$x^{(i)}, y^{(i)}$

Normalization
Standardization

Feature
Selection

Train MSE

$$h_\theta(x) = \theta_0 + \theta_1 x$$

**Testing**

New data $\rightarrow$ Regression model $\rightarrow$ Predictions

Price
Risk score
Test MSE

Unlabeled
$x'$

$$h_\theta(x') = \theta_0 + \theta_1 x'$$
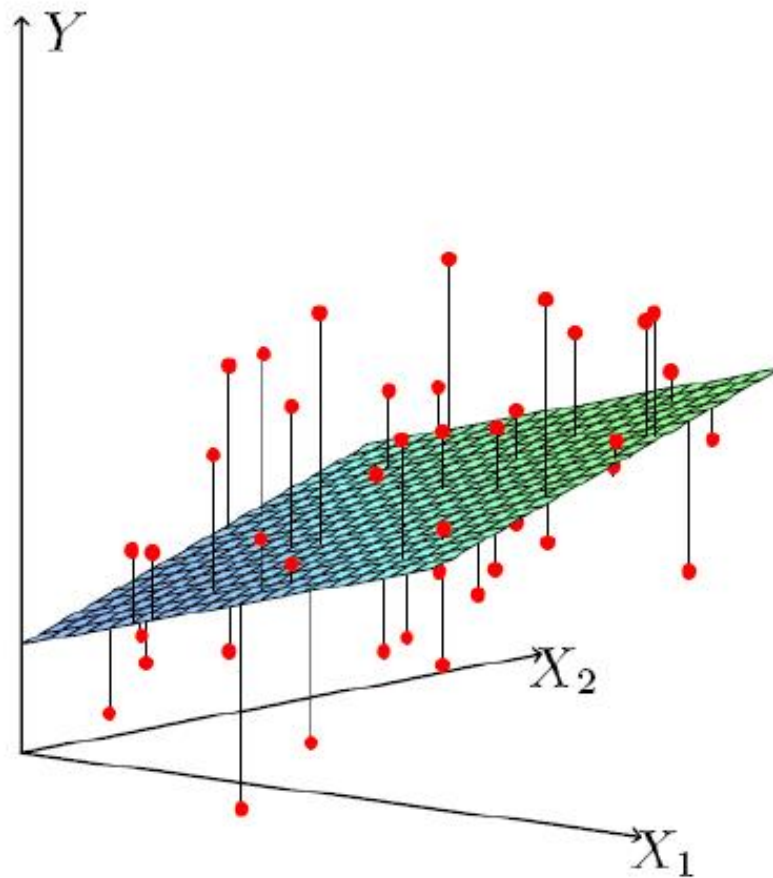
# Outline

- Multiple linear regression
  - Derivation in matrix form
- Practical issues
  - Feature scaling and normalization
  - Outliers
  - Categorical variables
- Gradient descent
  - Efficient algorithm for optimizing loss function
  - Training LR with Gradient Descent

# Multiple Linear Regression

- Dataset: $x^{(i)} \in R^d, y^{(i)} \in R$

# Use Vectorization

- Benefits of vectorization
  - More compact equations
  - Faster code (using optimized matrix libraries)
- Consider our model:

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

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \qquad \boldsymbol{x}^\mathsf{T} = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

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

# Use Vectorization

- Consider our model for n instances:

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

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \qquad \boldsymbol{X} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \cdots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$$

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

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

# Loss function

- For the linear regression cost function:

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

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{n} \left\| \mathsf{X}\theta - y \right\|^2$$

Let:

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

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

# Matrix and vector gradients

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

$$\frac{\partial y}{\partial x} = \left[ \begin{array}{cccc} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} & \cdots & \frac{\partial y}{\partial x_n} \end{array} \right]$$

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

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

Jacobian matrix

# Properties

- If $w, x$ are $(d \times 1)$ vectors, $\frac{\partial w^T x}{\partial x} = w$

- 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 $\frac{\partial}{\partial\boldsymbol{\theta}}J(\boldsymbol{\theta})=0$

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

Using chain rule

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

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

$$h'(x)=2x^T,g'(\theta)=X$$

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

Closed Form Solution: $\boldsymbol{\theta}=(\boldsymbol{X}^\mathsf{T}\boldsymbol{X})^{-1}\boldsymbol{X}^\mathsf{T}\boldsymbol{y}$

# Closed-form solution

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

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

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

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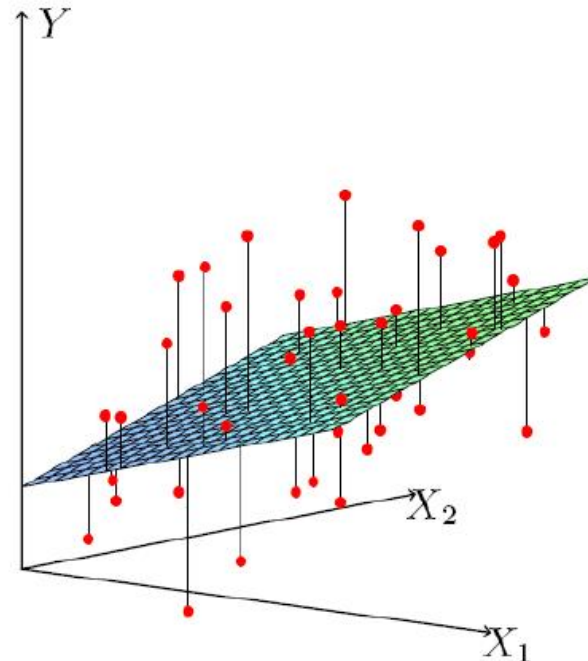    - In python, `numpy.linalg.pinv(a)`
  - Remove redundant (not linearly independent) features
  - Remove extra features to ensure that $d \leq n$
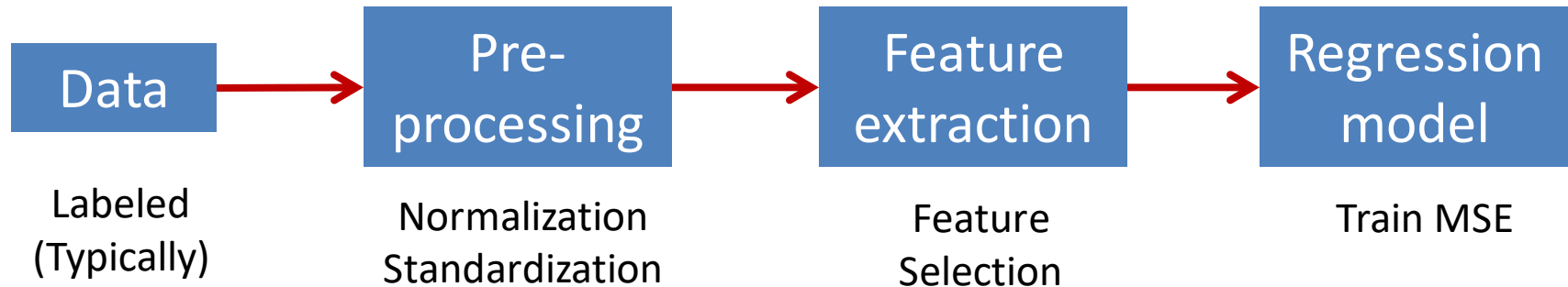
$$AGA = A$$

# Multiple Linear Regression

- Dataset: $x^{(i)} \in R^d, y^{(i)} \in R$
- Hypothesis $h_\theta(x) = \theta^T x$
- MSE $= \frac{1}{n} \sum_{i=1}^{n} (\theta^T x^{(i)} - y^{(i)})^2$ <span style="color:red">loss / cost</span>
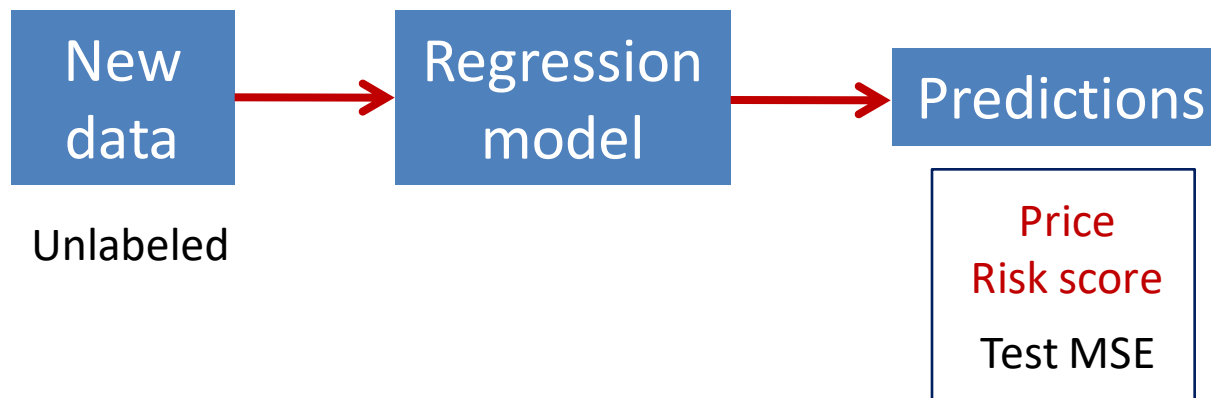
$$\theta = (X^\intercal X)^{-1} X^\intercal y$$

# Regression Learning

**Training**



| Data | → | Pre-processing | → | Feature extraction | → | Regression model |
|------|---|----------------|---|--------------------|---|------------------|

Labeled (Typically)

Normalization Standardization

Feature Selection

Train MSE

**Testing**

| New data | → | Regression model | → | Predictions |
|----------|---|------------------|---|-------------|

Unlabeled

Price
Risk score

Test MSE

# Feature Standardization

- Rescales features to have zero mean and unit variance

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

  - Replace each value with:

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

    - $s_j$ is the standard deviation of feature j

- Must apply the same transformation to instances for both training and prediction
- Outliers can cause problems

# Other feature normalization

- **Re-scaling**

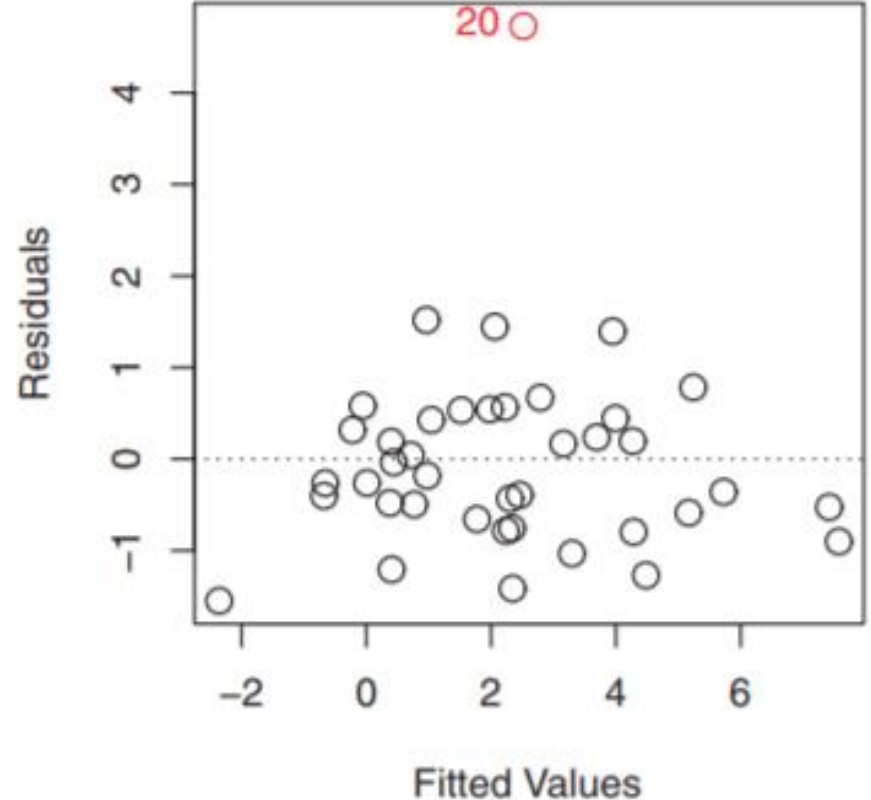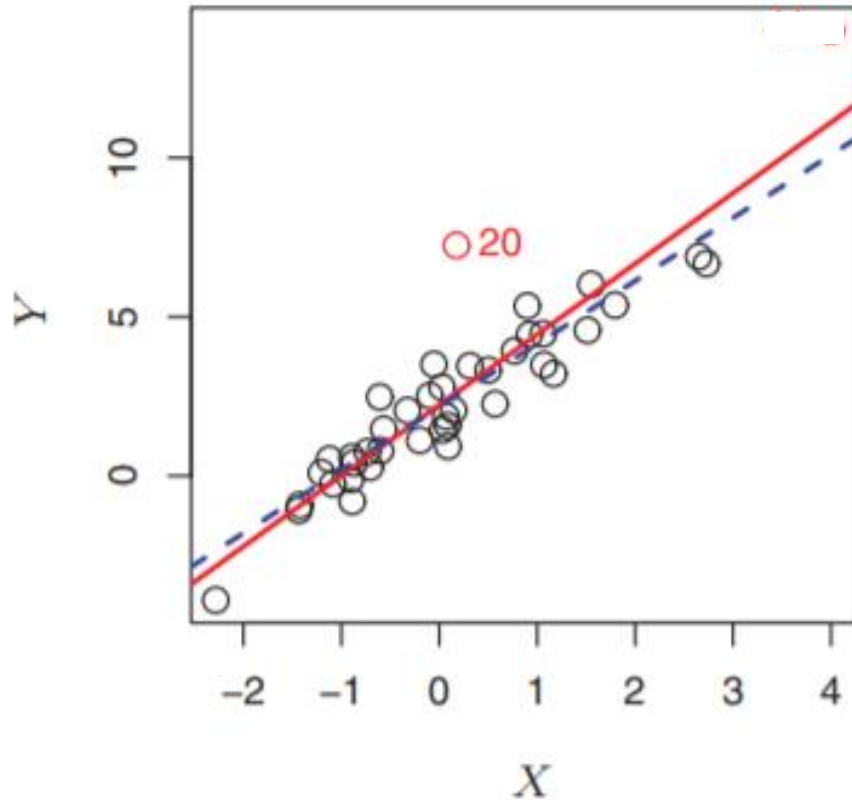$$- x_j^{(i)} \leftarrow \frac{x_j^{(i)} - min_j}{max_j - min_j} \in [0,1]$$

  - $min_j$ 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

# Outliers



- Dashed model is without outlier point
- Linear regression is not resilient to outliers!
- Outliers can be eliminated based on residual value
  - Other techniques for outlier detection

# Categorical variables

- Predict credit card balance
  - Age
  - Income
  - Number of cards
  - Credit limit
  - Credit rating
- Categorical variables
  - Student (Yes/No)
  - State (50 different levels)

# Indicator Variables

- Binary (two-level) variable
  - Add new feature $x_j = 1$ if student and 0 otherwise
- Multi-level variable
  - State: 50 values
  - $x_{MA} = 1$ if State $=$ MA and 0, otherwise
  - $x_{NY} = 1$ if State $=$ NY and 0, otherwise
  - …
  - How many indicator variables are needed?
- Disadvantages: data becomes too sparse for large number of levels

# Comparison with ANOVA

- ANOVA
  - General statistical method for comparing populations
  - Example 1: Is the income of MA and NY residents similar?
  - Example 2: Is there any difference between patients with certain treatment or no treatment?
- Linear regression
  - Learning algorithm used for predicting responses on new data
  - Example 1: Predict the income of US residents
  - Example 2: Predict survival of patients
  - Hypothesis testing for coefficient equal to zero is similar to ANOVA

# Outline

- Multiple linear regression
  - Derivation in matrix form
- Practical issues
  - Feature scaling and normalization
  - Outliers
  - Categorical variables
- Gradient descent
  - Efficient algorithm for optimizing loss function
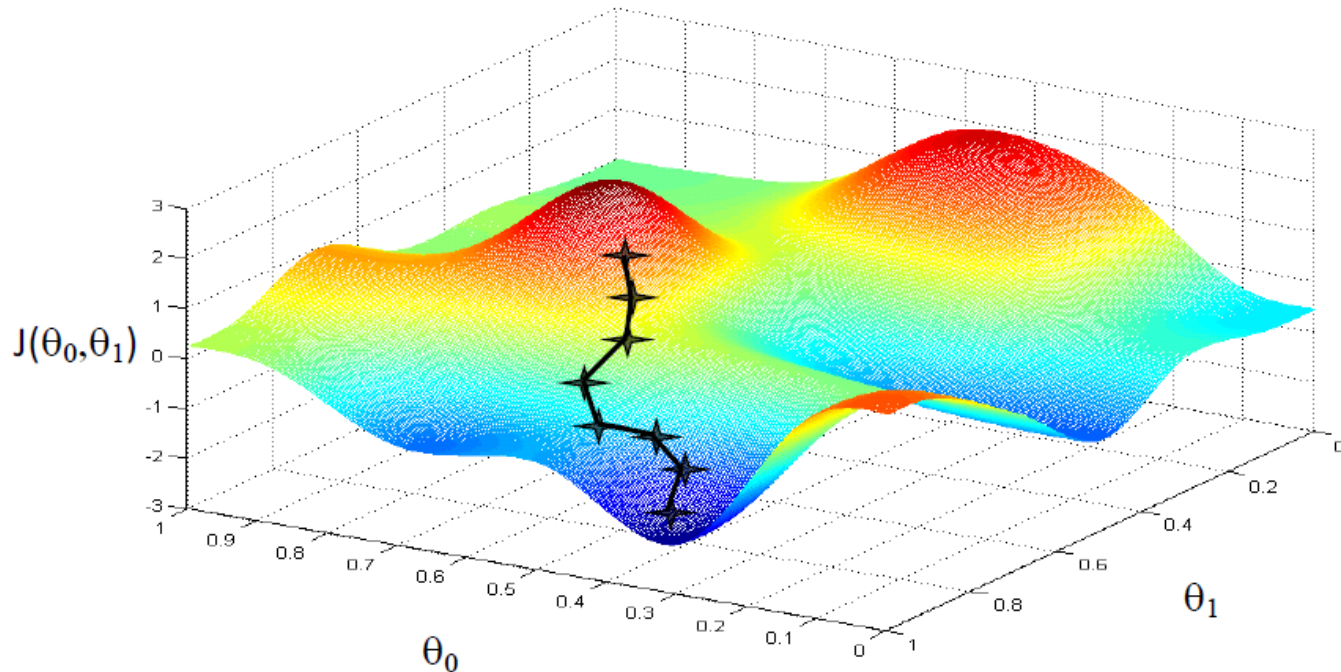  - Training LR with Gradient Descent

# What Strategy to Use?

# Follow the Slope
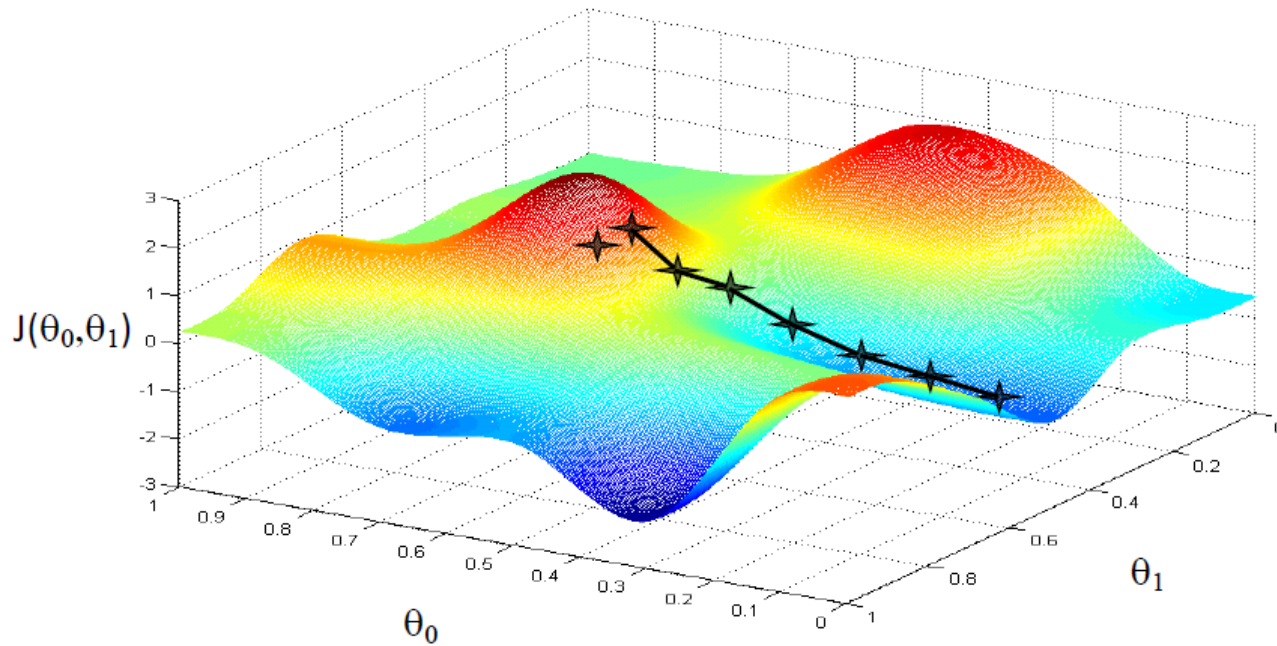


Follow the direction of steepest descent!

# How to optimize $J(\theta)$?

- Choose initial value for $\theta$
- Until we reach a minimum:
  - Choose a new value for $\boldsymbol{\theta}$ to reduce $J(\boldsymbol{\theta})$

# How to optimize $J(\theta)$?

- Choose initial value for $\theta$
- Until we reach a minimum:
  - Choose a new value for $\boldsymbol{\theta}$ to reduce $J(\boldsymbol{\theta})$
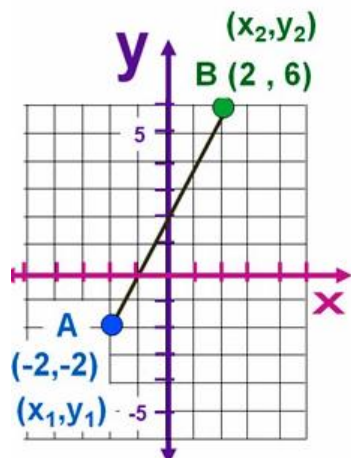


Different starting point

# Gradient Descent

- Initialize $\theta$

- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update
for j = 0 … d
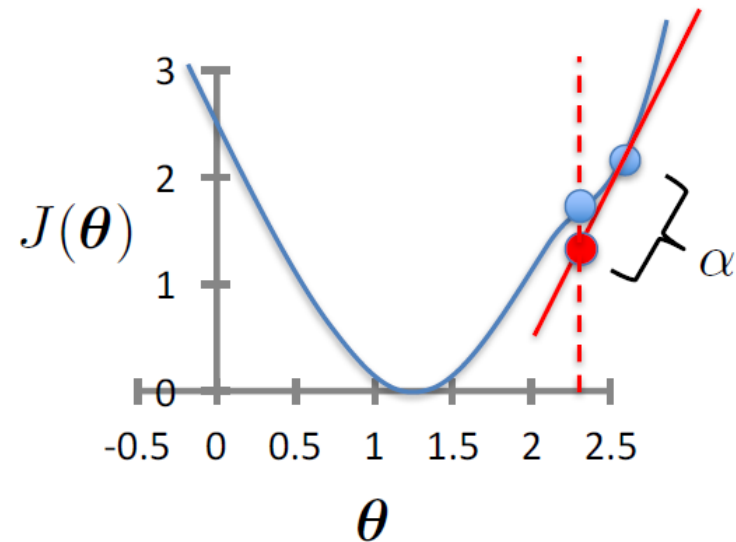
learning rate (small)
e.g., α = 0.05



The Gradient "m" is:

$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta Y}{\Delta X}$$

$$m = \frac{6 - ^-2}{2 - ^-2}$$

$$m = 8 / 4 = 2 \checkmark$$



Gradient = slope of line tangent
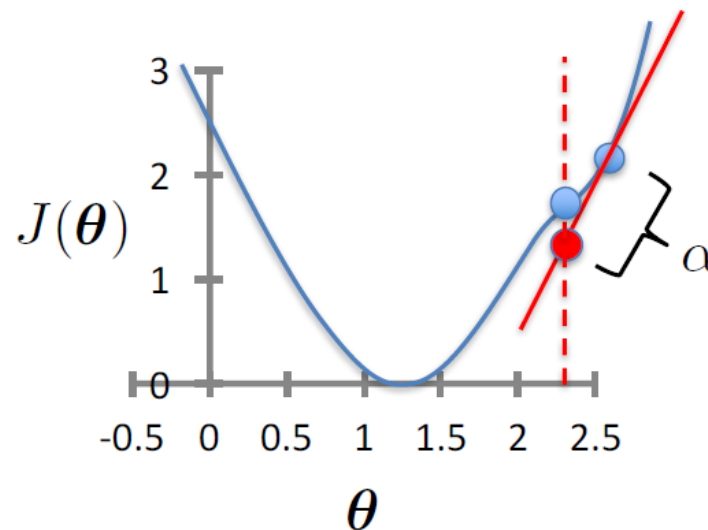to curve at the same point

# Gradient Descent

- Initialize $\boldsymbol{\theta}$

- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update
for j = 0 ... d

learning rate (small)
e.g., α = 0.05



$J(\boldsymbol{\theta})$

$\alpha$

$\boldsymbol{\theta}$

- What happens when $\theta$ reaches a local minimum?
- The slope is 0, and gradient descent converges!

# Review

- Solution for multiple linear regression can be computed in closed form
  - Matrix inversion is computationally intense
- Gradient descent is an efficient algorithm for optimization and training LR
  - The most widely used algorithm in ML!
  - Many variants (SGD, Coordinate descent, etc.)
  - Converges if objective is convex
- 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!