

DS 4400

Machine Learning and Data Mining I

Alina Oprea
Associate Professor
Khoury College of Computer Science
Northeastern University

September 22 2020

Announcements

- HW 1 PIAZZA: PDF) LATEX ; GRADESCOPE - PDF
CODE - GOOGLE FORM
ZIP FILE
– Is due on Monday, Sept. 28
- Python tutorials
 - Panda data frames tutorial by Alex Wang
 - Wed, Sept. 23, 5-6pm
 - Same Zoom link as office hours
 - Recording of first tutorial is available on Canvas under “Lecture Recording”

CLASS TIMES: [TUE. 11:45 AM - 1:25 PM
THU: 2:50 PM - 4:30 PM

Outline

MODULE 2

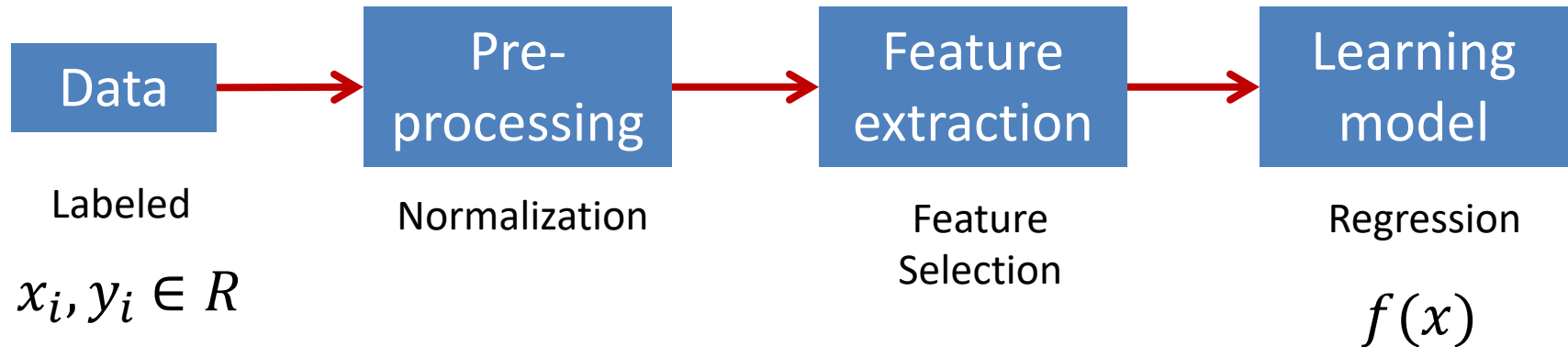
- GRADIENT DESCENT
- REGULARIZATION

- Linear regression
- Simple linear regression
 - MSE as loss function
 - Derivation of optimal solution
 - Correlation coefficient, covariance, and connection to regression
 - Example of linear regression fit
 - Lab in Python
- Multiple linear regression

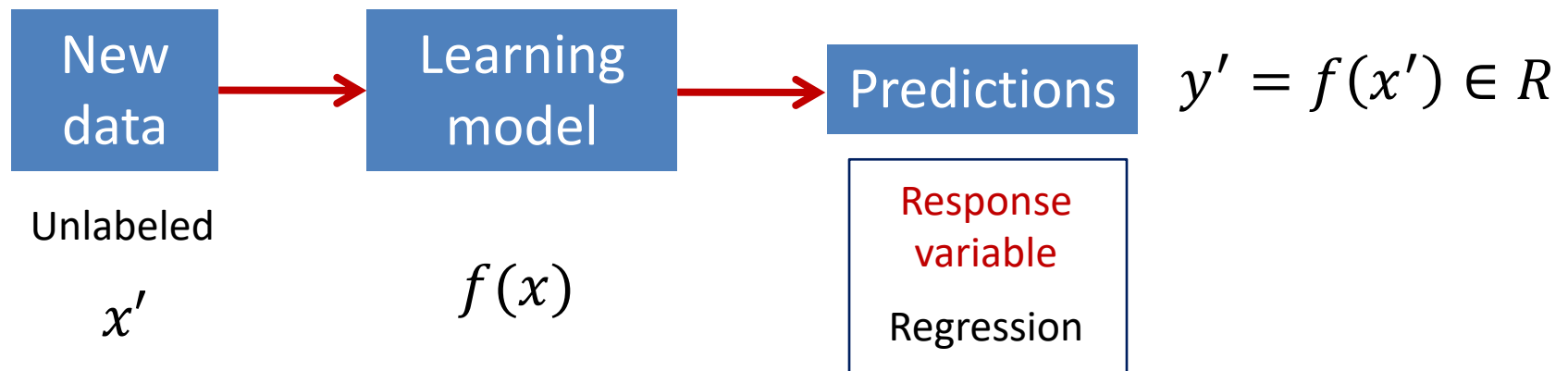
Linear regression

Supervised Learning: Regression

Training



Testing



Steps to Learning Process

- Define problem space
 - Collect data
 - Extract feature
 - Pick a model (hypothesis)
 - Develop a learning algorithm
 - Train and learn model parameters
 - Make predictions on new data
 - Testing phase
 - In practice, usually re-train when new data is available and use feedback from deployment
- Handwritten red annotations:*
- A bracket groups the first four steps (Define problem space, Collect data, Extract feature, Pick a model (hypothesis)) with the word **HUMERICAL**.
 - A bracket groups the fifth step (Develop a learning algorithm) and its sub-step (Train and learn model parameters) with the phrase **OPTIMIZATION**.
 - A bracket groups the sixth step (Make predictions on new data) and its sub-step (Testing phase) with the phrase **FIT MODEL TO DATA**.

Linear regression

- One of the most widely used techniques
- Fundamental to many complex models
 - Generalized Linear Models
 - Logistic regression
 - Neural networks
 - Deep learning
- Easy to understand and interpret
- Efficient to solve in closed form
- Efficient practical algorithm (gradient descent)

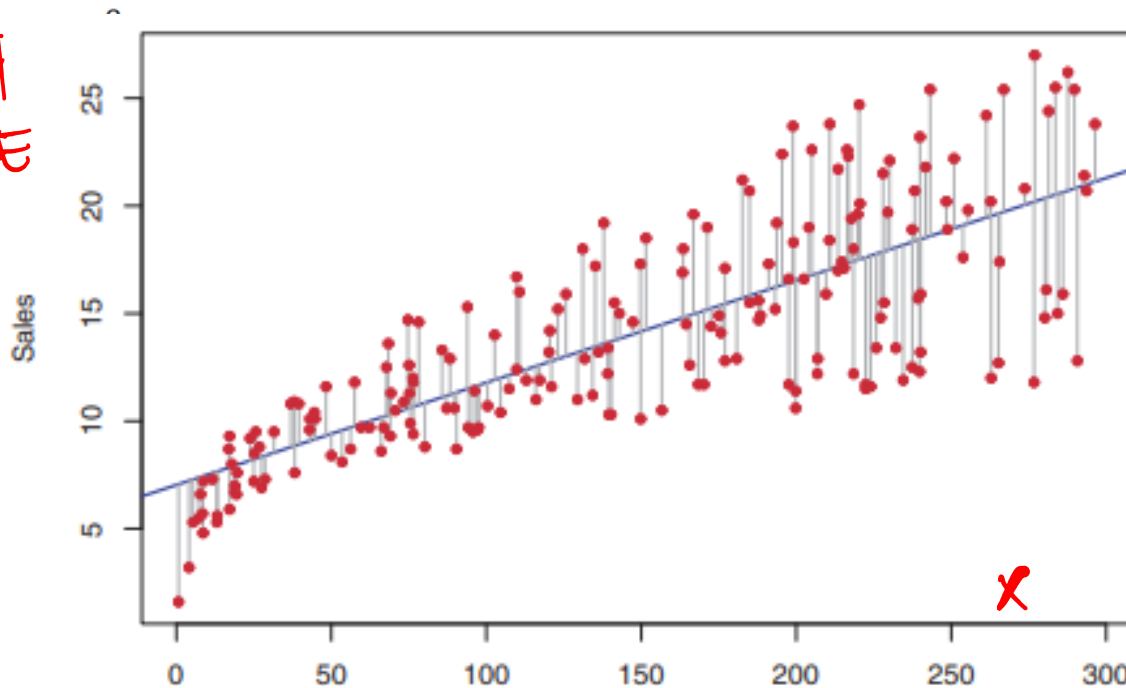
Linear regression

Given:

- Data $X = \{x_1, \dots, x_N\}$, where $x_i \in \mathbb{R}^d$
- Corresponding labels $Y = \{y_1, \dots, y_N\}$, where $y_i \in \mathbb{R}$

TRAINING
DATA

Y
RESPONSE



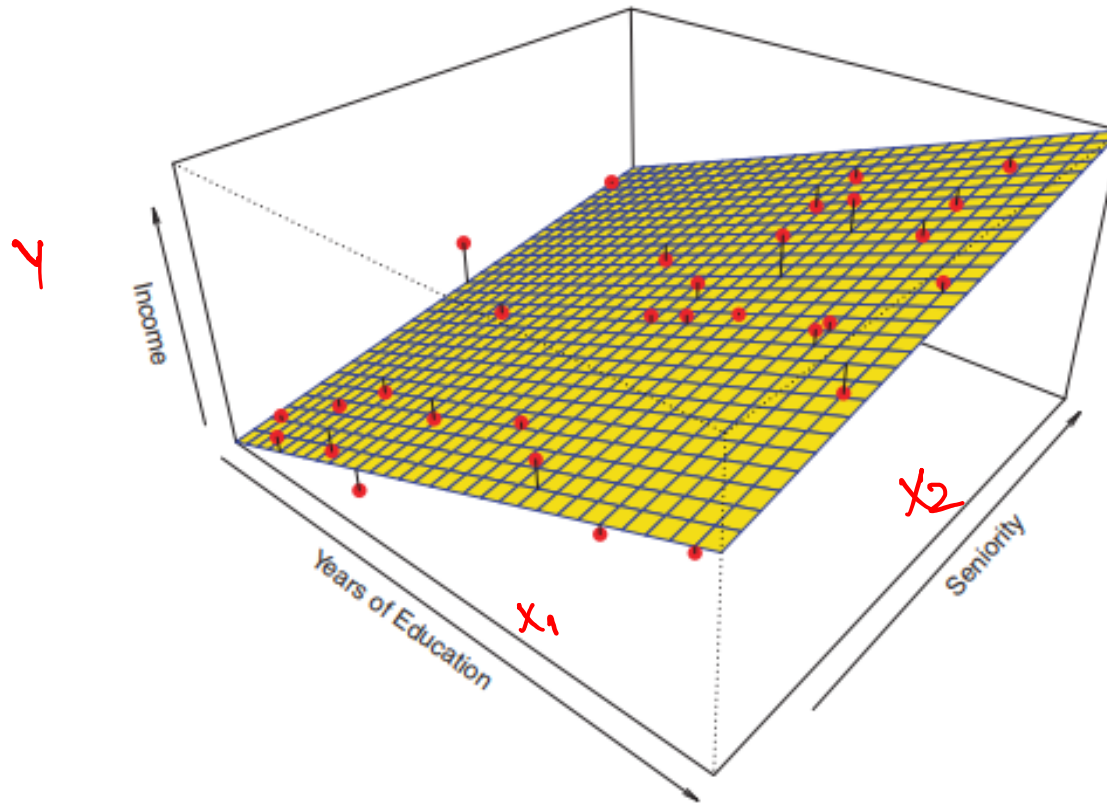
X

PREDICTOR

SIMPLE LINEAR REGR
1 PREDICTOR

ORDINARY LEAST SQUARES
(OLS)

Income Prediction



MULTIPLE LINEAR REGRESSION

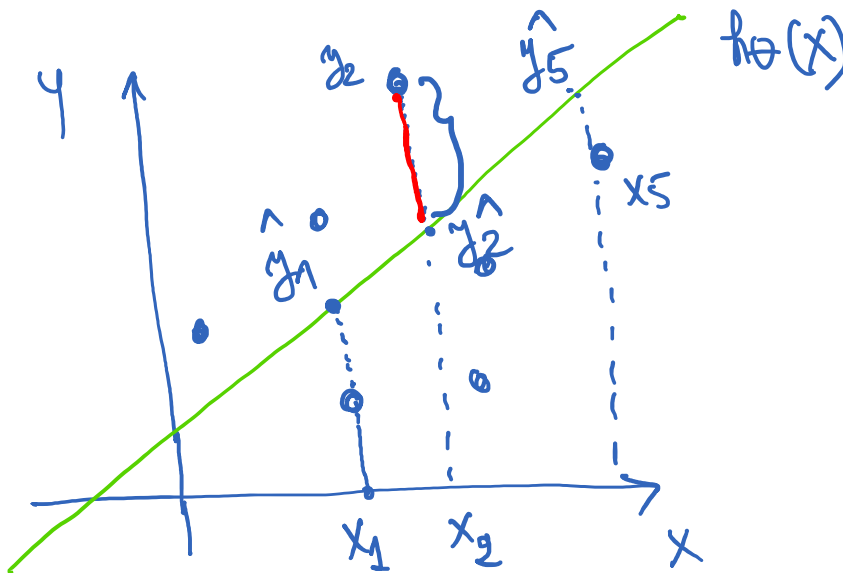
Hypothesis: Linear Model

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

(x_i, y_i)

TRAINING DATA

Simple linear regression: line with 2 parameters: θ_0, θ_1



$$\hat{y}_1 = h_{\theta}(x_1) \quad \text{PREDICTED VALUE}$$

$$y_1 = \text{TRUE VALUE}$$

$$\hat{y}_1 - y_1 : \text{ERRORS}$$

$$\hat{y}_2 - y_2 \quad \text{RESIDUALS}$$

Least-Squares Linear Regression

- Cost Function

LOSS FUNCTION

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N [h_{\theta}(x_i) - y_i]^2$$

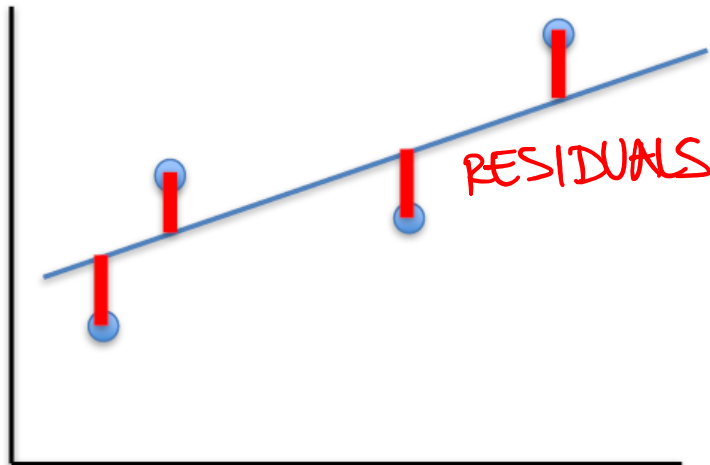
TRAINING DATA POINTS

RESIDUAL

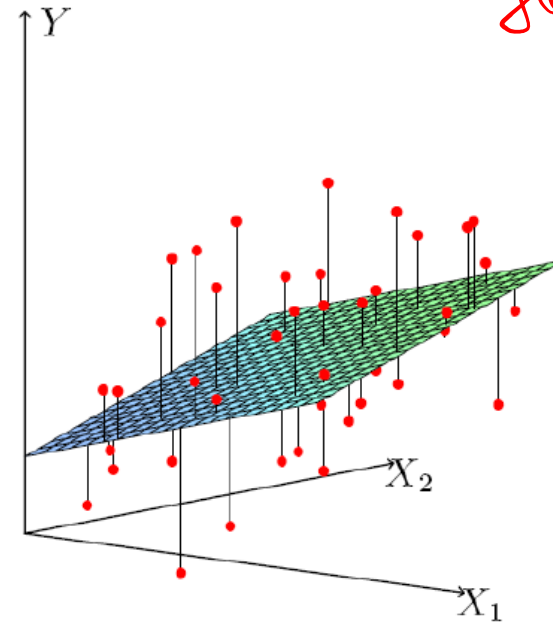
Mean Square Error (MSE)

$J(\theta)$: CONVEX

- Fit by solving $\min_{\theta} J(\theta)$



SIMPLE LR



MULTIPLE LR

Terminology and Metrics

- **Residuals**

- Difference between predicted values and actual values

- Predicted value for example i is: $\hat{y}_i = h_{\theta}(x_i)$

- $R_i = |y_i - \hat{y}_i| = |y_i - (\theta_0 + \theta_1 x_i)|$

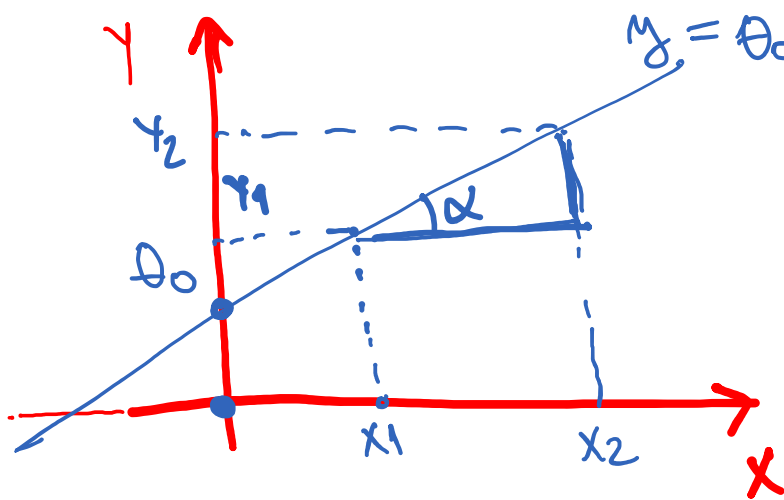
- **Residual Sum of Squares (RSS)**

- $RSS = \sum R_i^2 = \sum [y_i - (\theta_0 + \theta_1 x_i)]^2$

- • **Mean Square Error (MSE)**

- $MSE = \frac{1}{N} \sum R_i^2 = \frac{1}{N} \sum [y_i - (\theta_0 + \theta_1 x_i)]^2$

Interpretation



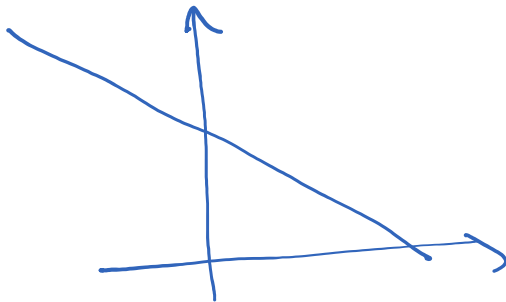
$$\theta_1 = \frac{y_2 - y_1}{x_2 - x_1} \quad \text{SLOPE}$$

$$\tan(\alpha) = \theta_1$$

$$\text{If } \theta_1 = 1 \Rightarrow \alpha = 45^\circ$$

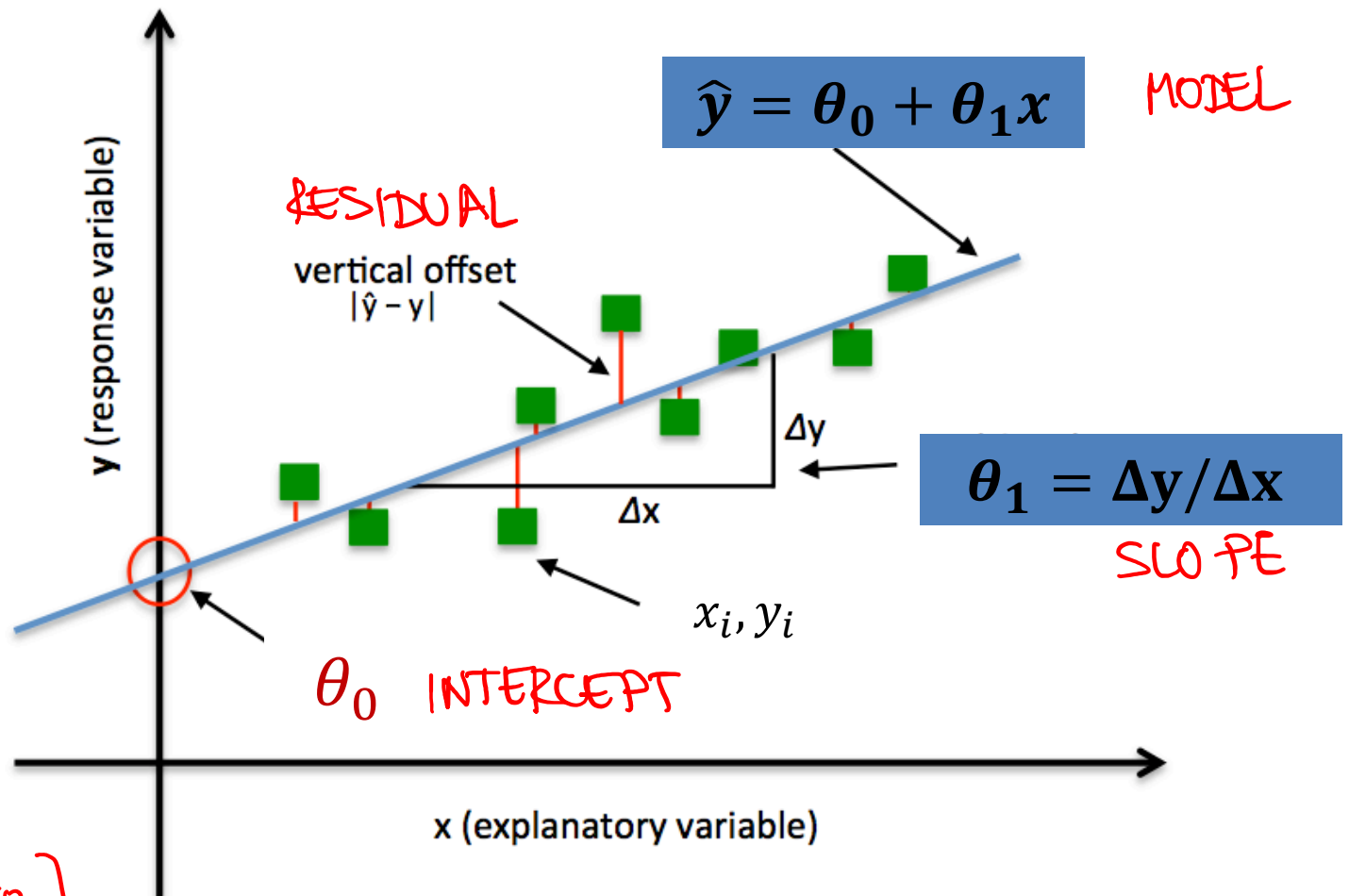
$$\theta_1 = -1 \Rightarrow \alpha = 135^\circ$$

$$x=0 \Rightarrow y=\theta_0 \quad \text{INTERCEPT}$$



$$\theta_1 = -1$$

Interpretation



$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

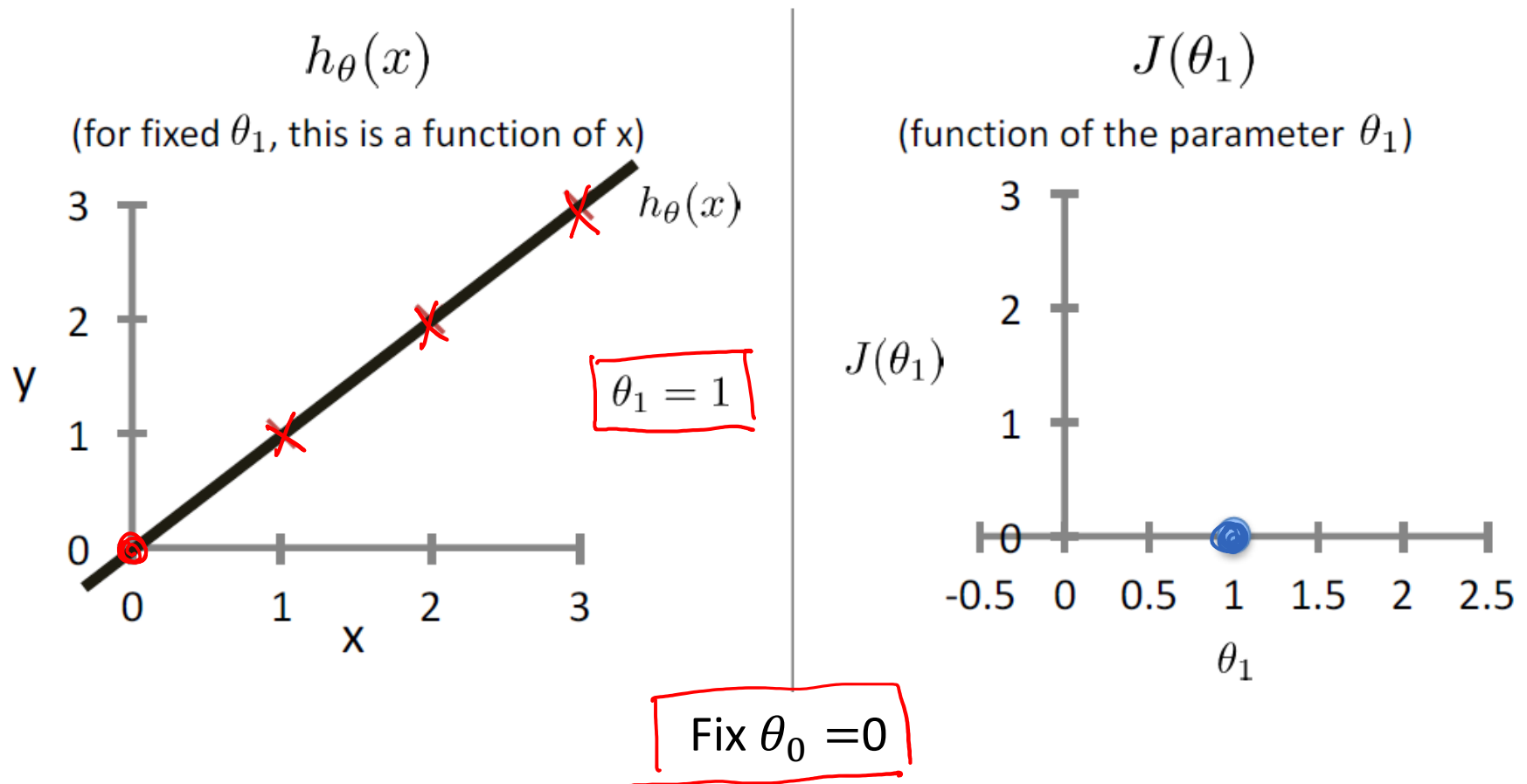
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$
 MODEL

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N [h_{\theta}(x_i) - y_i]^2$$
 ERROR METRIC

Intuition on MSE

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N [h_{\theta}(x_i) - y_i]^2$$

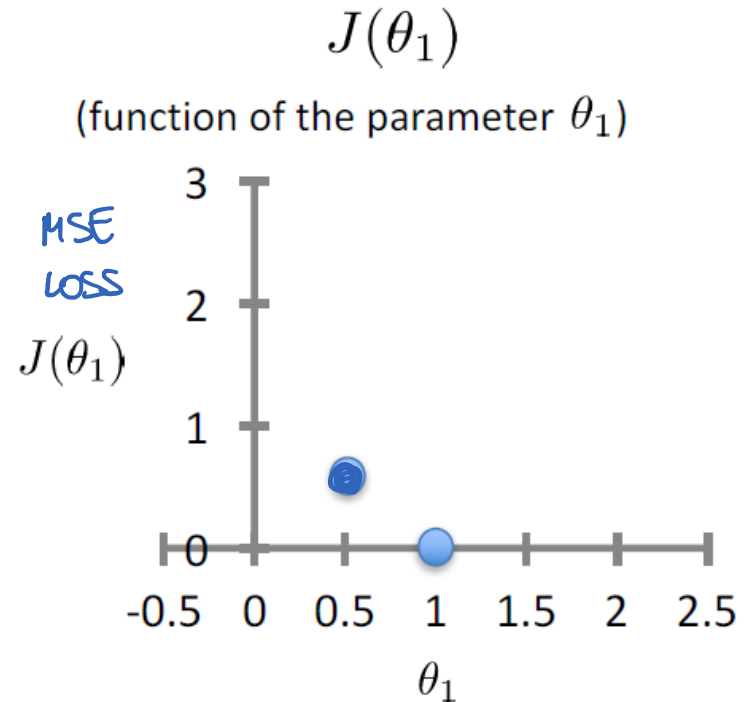
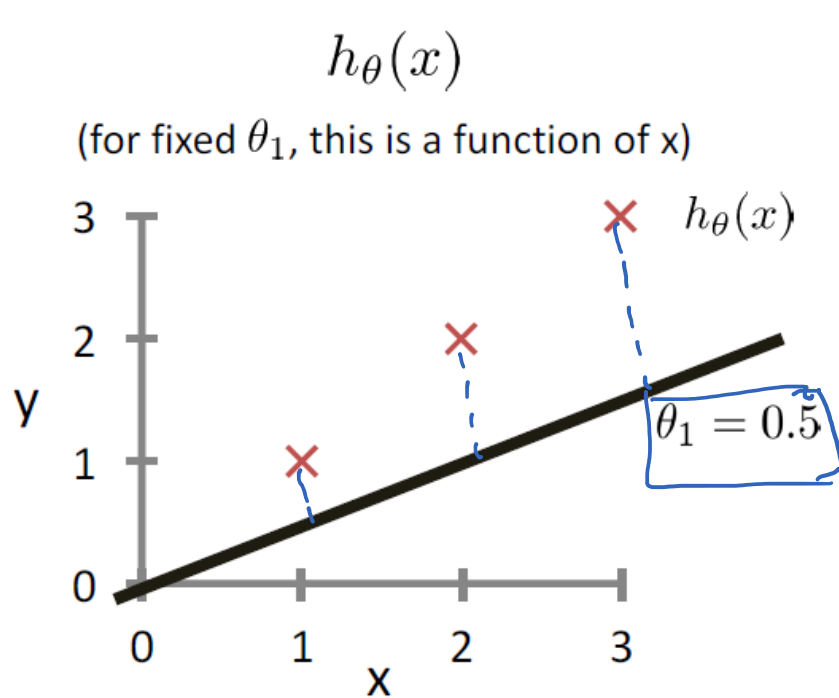
For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$



Intuition on MSE

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N [h_{\theta}(x_i) - y_i]^2$$

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$



Based on example
by Andrew Ng

$$J([0, 0.5]) = \frac{1}{3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \approx 0.58$$

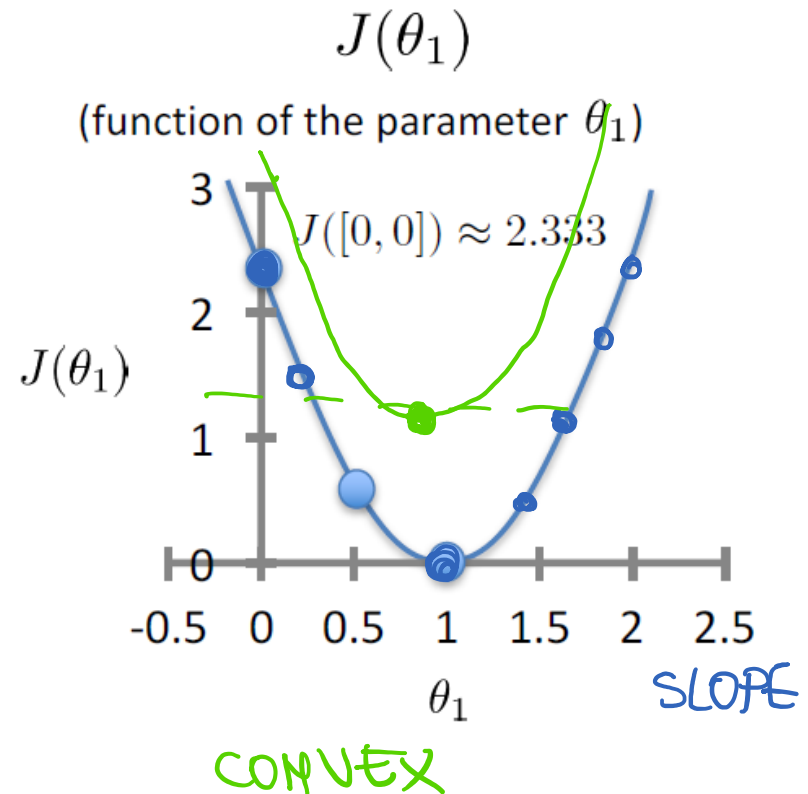
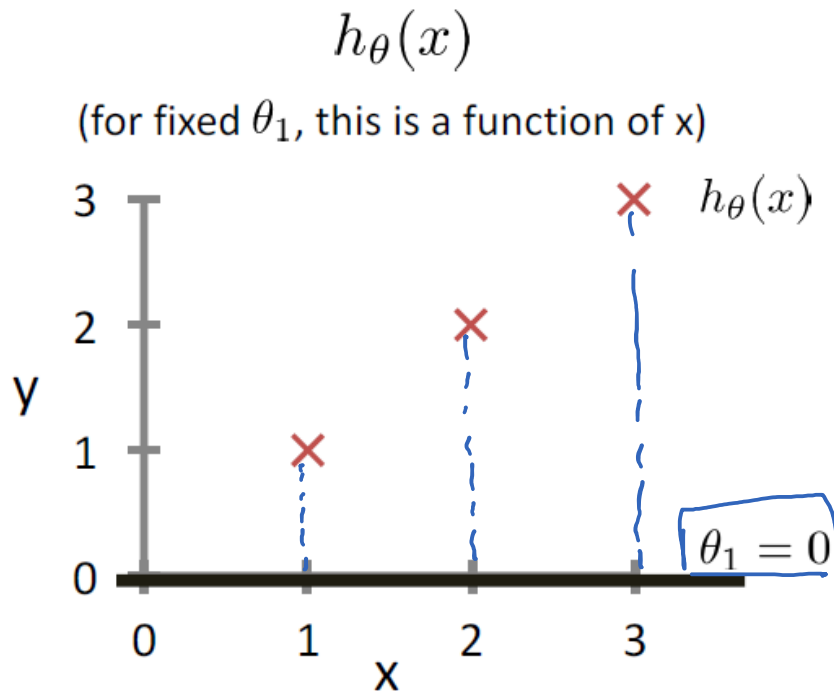
MSE

Intuition on MSE

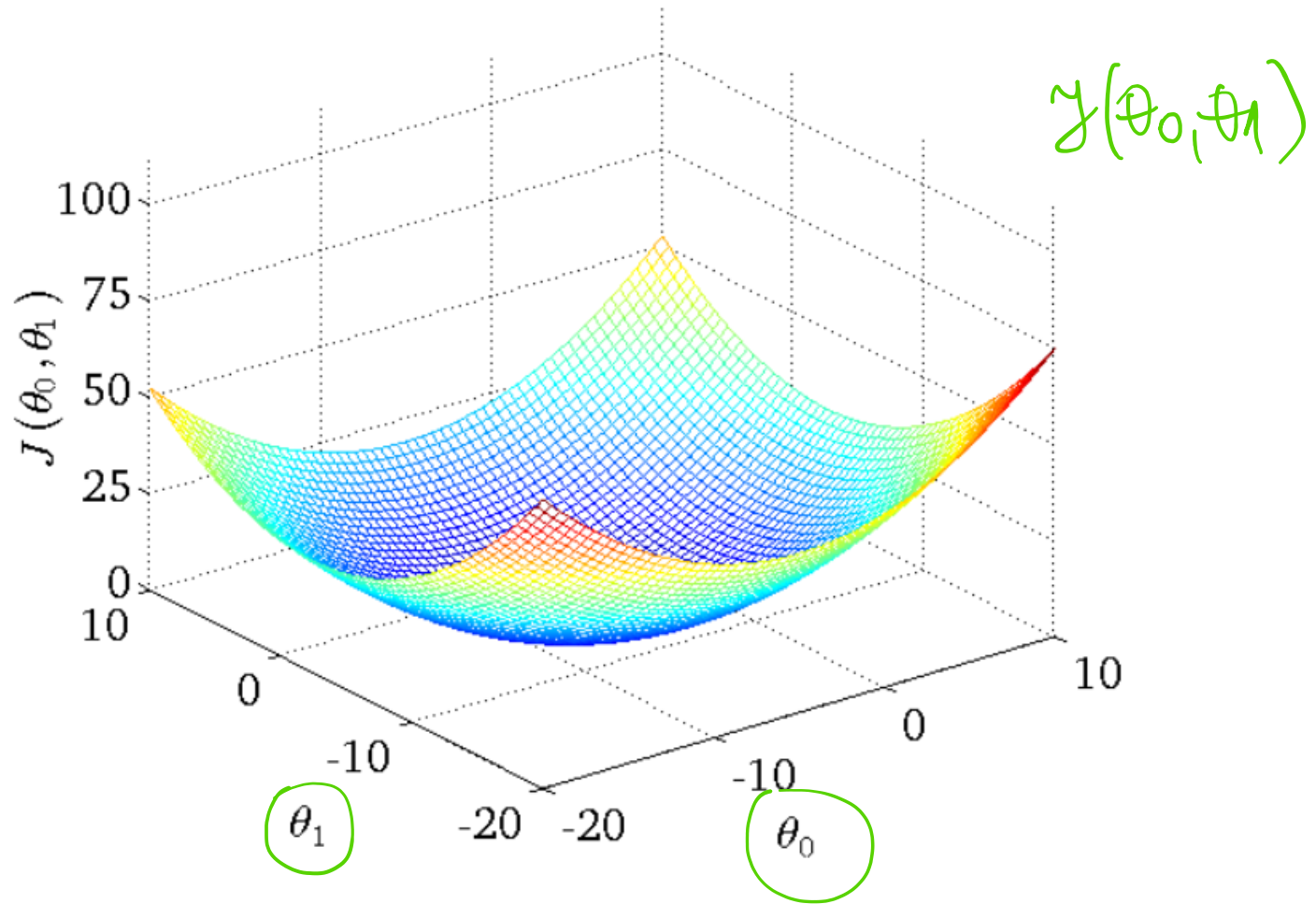
INPUT:
 (x_i, y_i)

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N [h_{\theta}(\overset{\text{INPUT}}{x_i}) - \overset{\text{OUTPUT}}{y_i}]^2 \Rightarrow \theta_0, \theta_1: \min J(\theta)$$

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$



MSE function



Convex function, unique minimum

Solution for simple linear regression

- Dataset $\overbrace{x_i \in R, y_i \in R}^{N \text{ points}}, h_{\theta}(x) = \theta_0 + \theta_1 x$
- $J(\theta) = \frac{1}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_i - y_i)^2$ **MSE / Loss**

Find θ_0 and θ_1
for which $\min J(\theta_0, \theta_1)$

Relationship between Two Random Variables

- Model X (feature / predictor) and Y (response) as two random variables
- Fit of simple linear regression depends on dependence between X and Y
- Covariance
 - Measures the strength of relationship between two random variables
- Pearson correlation
 - Normalized between $[-1,1]$
 - Proportional to covariance

Covariance

- X and Y are random variables
- $Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$
- Properties

$$(1) \quad Cov(X, X) = E[(X - E(X))(X - E(X))] = E[(X - E(X))^2] = Var(X)$$

$$(2) \quad Cov(X, Y) = Cov(Y, X)$$

$$(3) \quad Cov(aX, Y) = a \cdot Cov(X, Y)$$

Covariance

- X and Y are random variables
- $Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$ DEF

$$\begin{aligned} Cov(X, Y) &= E[XY - XE(Y) - YE(X) + E(X)E(Y)] \\ &= E[XY] - E(X)E(Y) - E(X)E(Y) + E(X)E(Y) \\ &= E[XY] - E(X)E(Y) \end{aligned}$$

If X and Y are indep. $\Rightarrow E[XY] = E[X]E[Y]$

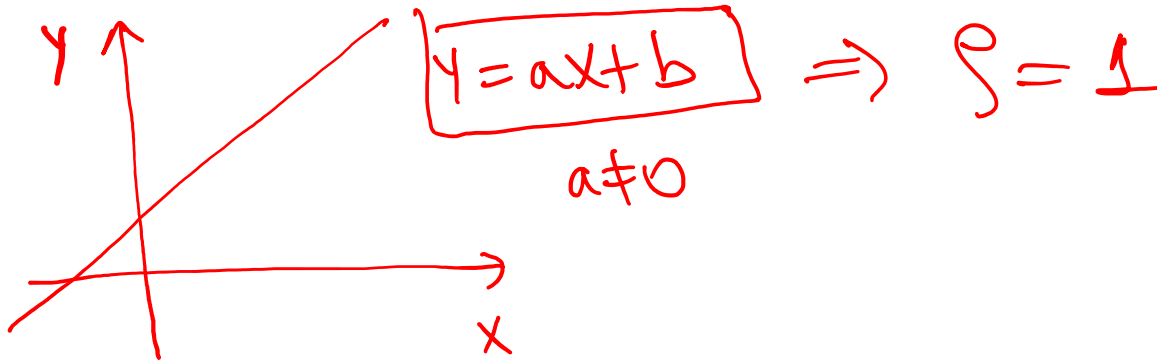
A, B ; $\underbrace{P(A \cap B) = P(A)P(B)}_{DEF} \Rightarrow \boxed{Cov(X, Y) = 0}$

Pearson Correlation

$$\rho = \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \in [-1, 1]$$

Standard deviation

$$\sigma_X = \sqrt{\text{Var}(X)}$$



$$\text{Cov}(X, Y) = \text{Cov}(X, aX + b) = \text{Cov}(X, aX) = a \cdot \text{Cov}(X, X) = a \cdot \text{Var}(X)$$

$$\sigma_Y = \sqrt{\text{Var}(Y)} = \sqrt{a^2 \text{Var}(X)} \Rightarrow \sigma_Y = a \sigma_X$$

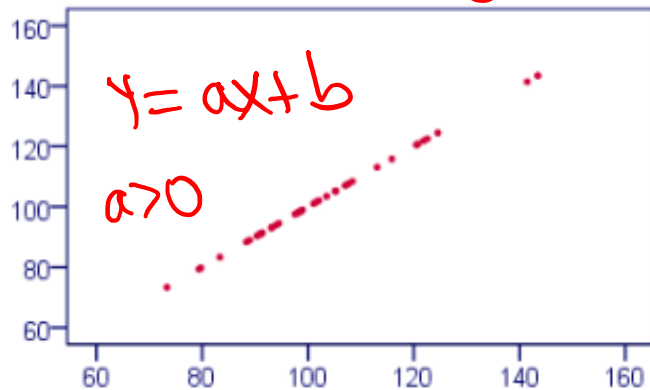
$$\rho = \frac{a \cdot \text{Var}(X)}{\sigma_X \cdot a \cdot \sigma_X} = 1$$

Pearson Correlation

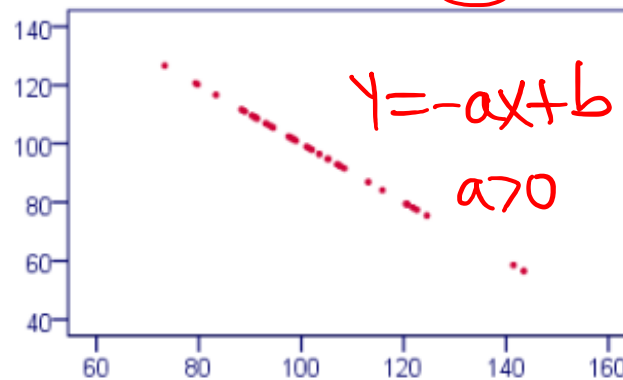
$$\rho = \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \in [-1, 1]$$

Standard deviation
 $\sigma_X = \sqrt{\text{Var}(X)}$

Correlation Coefficient = 1



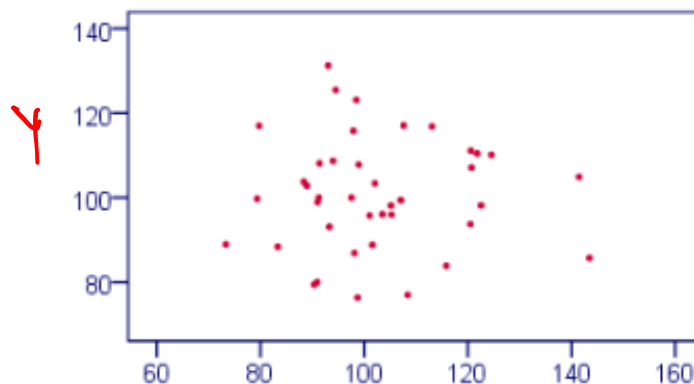
Correlation Coefficient = -1



POSITIVE CORR

NEGATIVE CORR

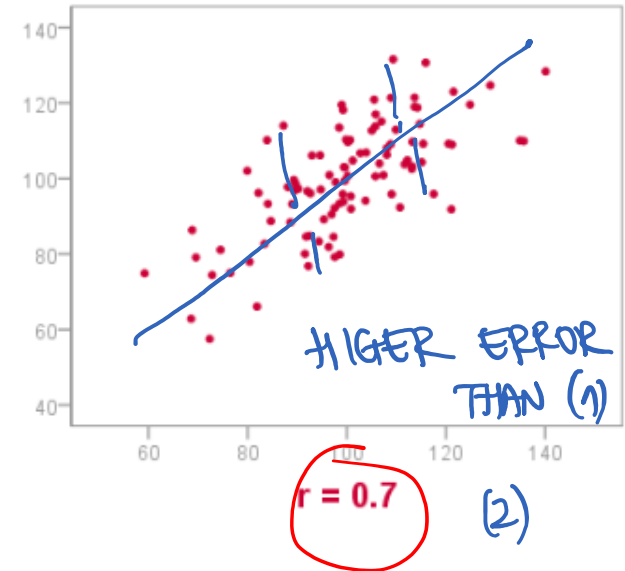
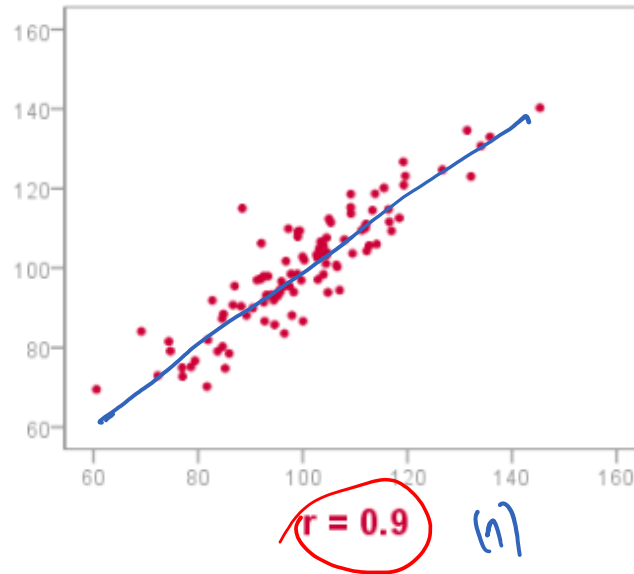
Correlation Coefficient = 0



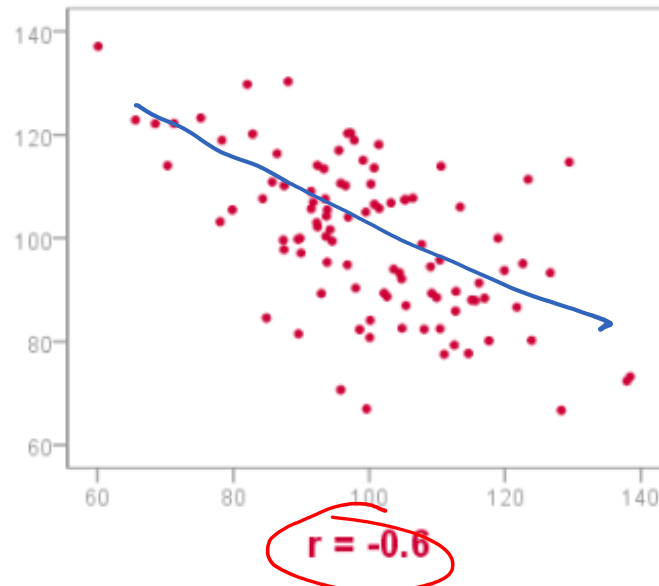
X

Positive/Negative Correlation

Positive
Correlation



Negative
Correlation



HIGHER ERROR
THAN (1) & (2)

How Well Does the Model Fit?

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

$$\rho = \text{Corr}(X, Y) = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\theta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \quad \text{SLOPE}$$

$$\text{If } \sigma_X = \sigma_Y \Rightarrow \boxed{\rho = \theta_1}$$

Regression vs Correlation

- Correlation

- Find a numerical value expressing the relationship between variables

- Regression

- Estimate values of response variable on the basis of the values of predictor variable
- The slope of linear regression is related to correlation coefficient
- Regression scales to more than 2 variables, but correlation does not