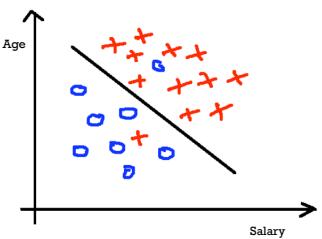




CHULA ENGINEERING  
Foundation toward Innovation

COMPUTER



## Linear Regression

2110574: AI for Engineers

Peerapon Vateekul, Ph.D.

Department of Computer Engineering,  
Faculty of Engineering, Chulalongkorn University

[Peerapon.v@chula.ac.th](mailto:Peerapon.v@chula.ac.th)



# Outlines

- Introduction
- Simple Linear Regression
- Multiple Linear Regression
- Other topics
- Demo



# Introduction

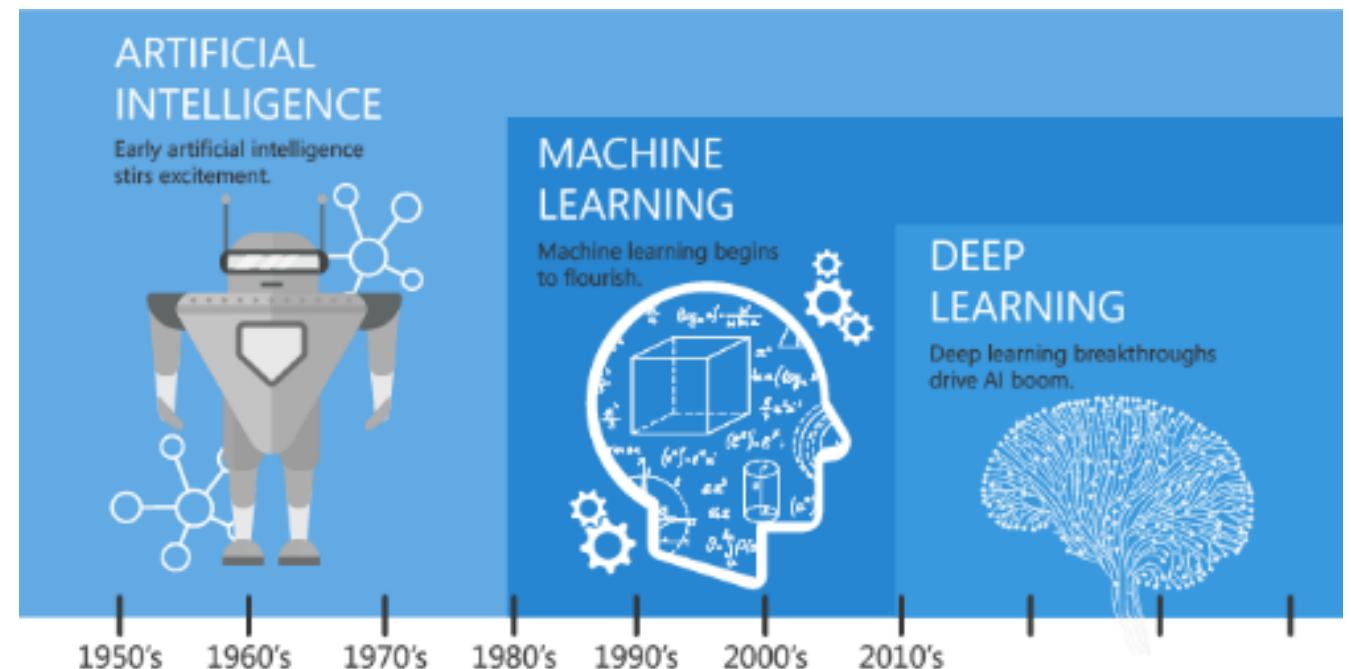
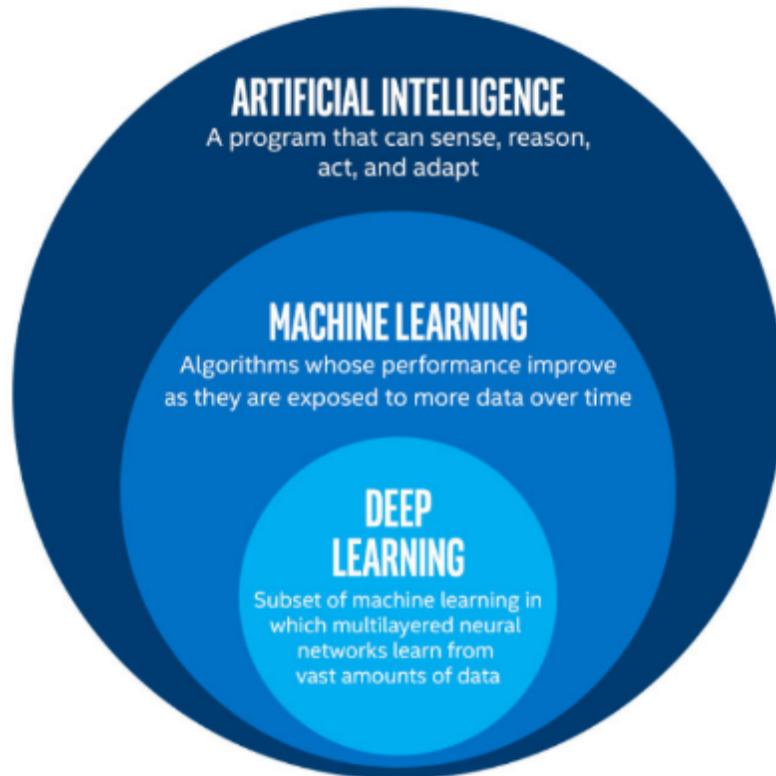


# AI, Machine Learning, and Deep Learning

- Machine Learning (ML) is a subfield in AI focusing on making to learn by itself without human intervention.

*“Machine learning is the science of getting computers to act without being explicitly programmed.” — Stanford University*

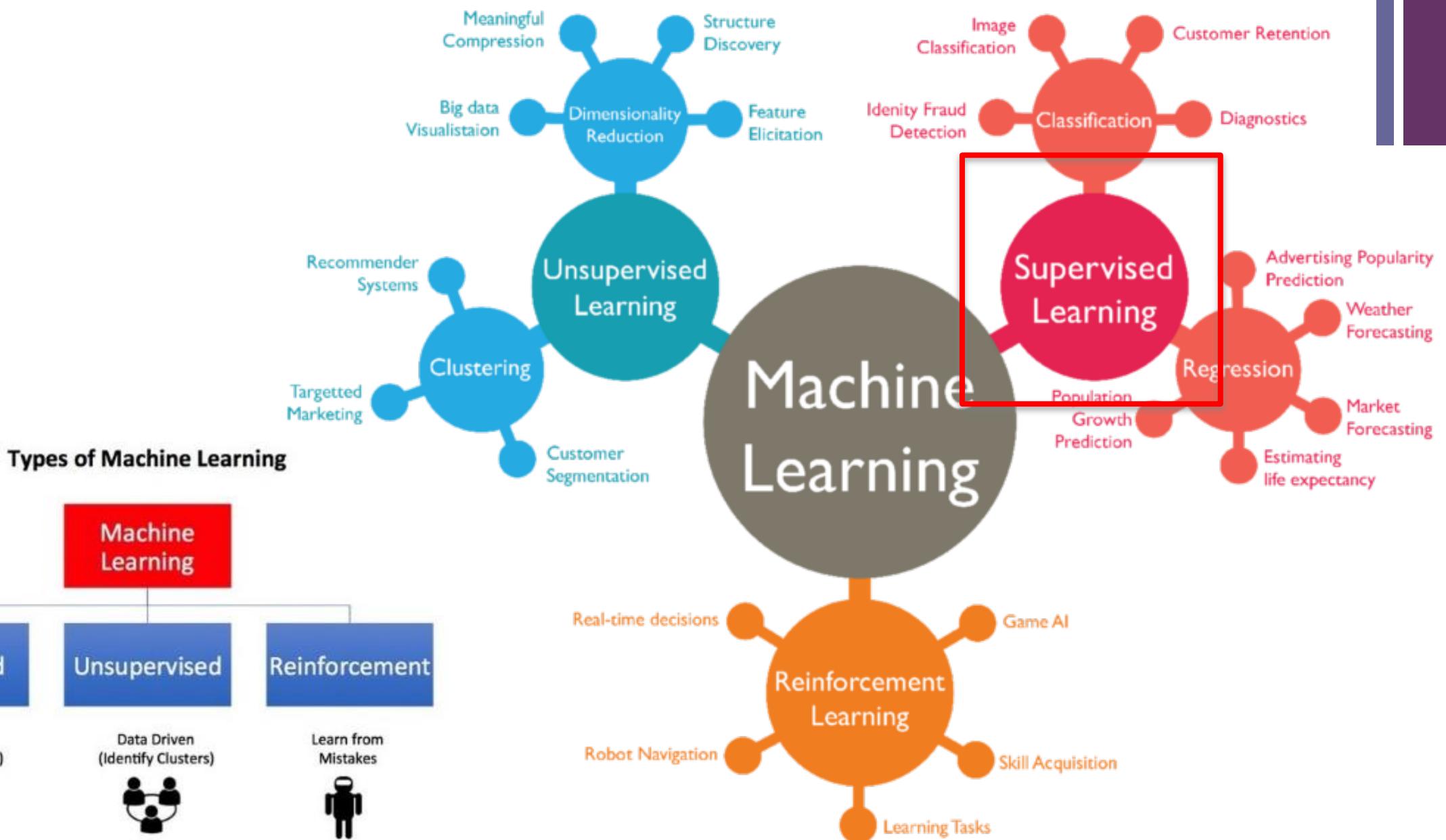
<https://towardsdatascience.com/cousins-of-artificial-intelligence-dda4edc27b55>



Since an early flush of optimism in the 1950's, smaller subsets of artificial intelligence - first machine learning, then deep learning, a subset of machine learning - have created ever larger disruptions.

# + Machine Learning (cont.)

5





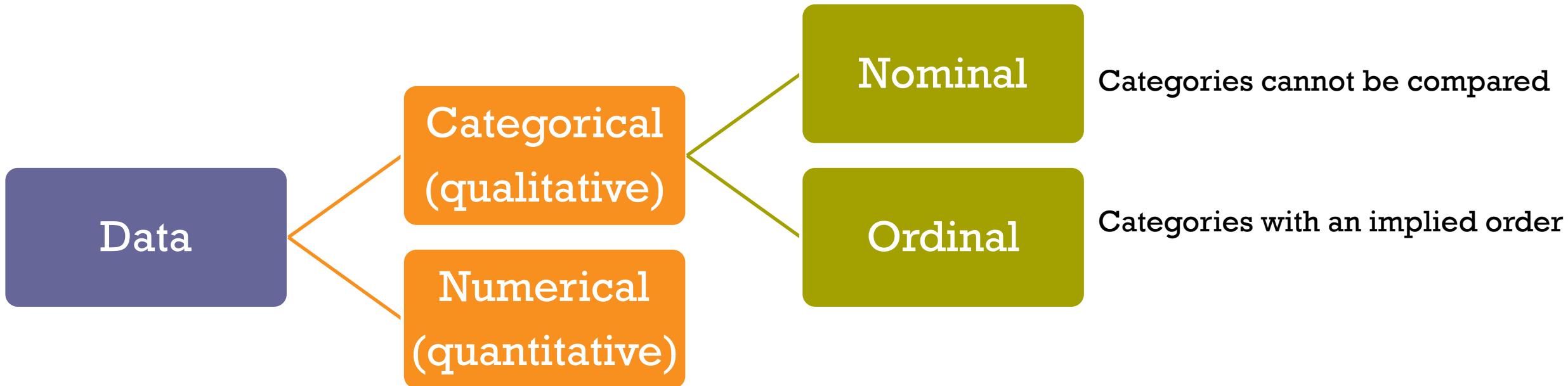
# Terminology: Data table

inputs				target
Age	Income	Gender	Province	Purchase
25	25,000	Female	Bangkok	Yes
35	50,000	Female	Nontaburi	Yes
32	35,000	Male	Bangkok	No

- Row
  - Example, instance, case, observation, subject
- Column
  - Feature, variable, attribute
- Input
  - Predictor, independent, explanatory variable
- Target
  - Output, outcome, response, dependent variable



# Terminology: Kinds of data

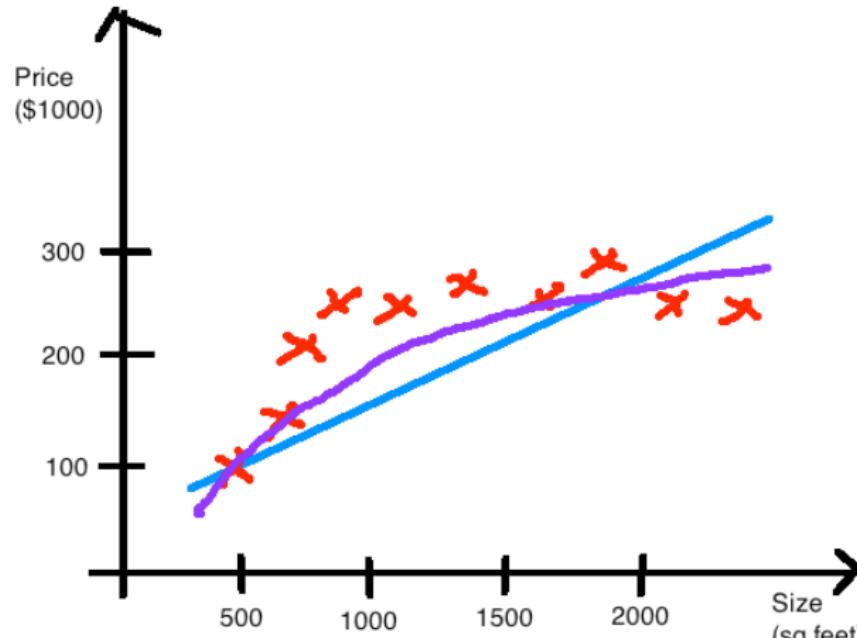






# Regression: predict a continuous value

## Linear Regression



Predict a sale price of each house

### ■ Some techniques:

- Linear Regression / GLM
- Decision Trees
- Support vector regression
- Neural Network
- Ensembles

### ■ Sample Applications

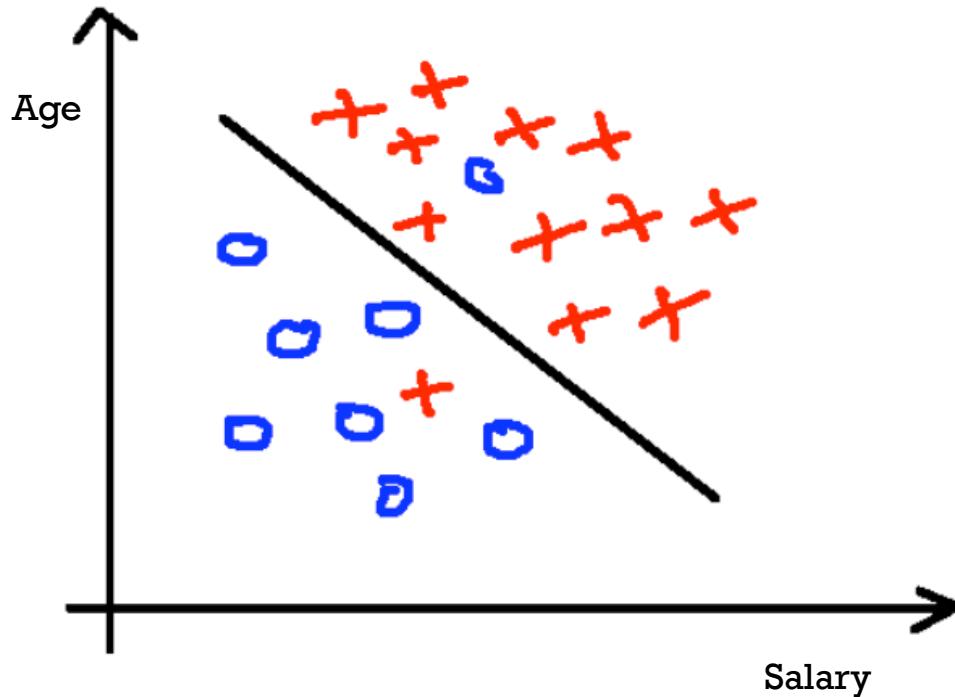
- Financial risk management
- Revenue forecasting





# Classification: predicting a category

## Logistic Regression



Predict targeted customers who  
tend to buy our product (yes/no)

- **Some techniques:**
  - Naïve Bayes
  - Decision Tree
  - Logistic Regression
  - Support Vector Machines
  - Neural Network
  - Ensembles
  
- **Sample Applications**
  - Database marketing
  - Fraud detection
  - Pattern detection
  - Churn customer detection

# + Prediction algorithms

- Decision Tree
- (Logistic) Regression
- kNN
- Support Vector Machine
- Neural Networks (NN)
- Deep Learning

**BASIC REGRESSION**

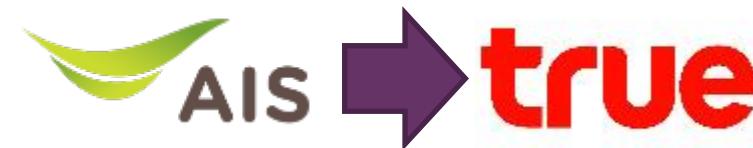
- LINEAR linear\_model.LinearRegression()  
Lots of numerical data
- LOGISTIC linear\_model.LogisticRegression()  
Target variable is categorical or

**CLASSIFICATION**

- NEURAL NET neural\_network.MLPClassifier()  
Complex relationships. Prone to overfitting  
Basically magic.
- K-NN neighbors.KNeighborsClassifier()  
Group membership based on proximity
- DECISION TREE tree.DecisionTreeClassifier()  
If/then/else. Non-contiguous data  
Can also be regression
- RANDOM FOREST ensemble.RandomForestClassifier()  
Find best split randomly  
Can also be regression
- SVM svm.SVC() svm.LinearSVC()  
Maximum margin classifier. Fundamental Data Science algorithm
- NAIVE BAYES GaussianNB() MultinomialNB() BernoulliNB()  
Updating knowledge step by step with new info



# Supervised learning (recap)



Training Data



inputs				target
Age	Income	Gender	Province	Churn
25	25,000	Female	Bangkok	Yes
35	50,000	Female	Nontaburi	Yes
32	35,000	Male	Bangkok	No

Testing Data



Age	Income	Gender	Province	Churn
25	25,000	Female	Bangkok	?

Application: Direct Target Customer

+

# Simple Linear Regression



# Problem: 1 input (predictor) & 1 output

- Collect data of 7 patients
- Systolic Blood Pressure (y) & Cholesterol (x)

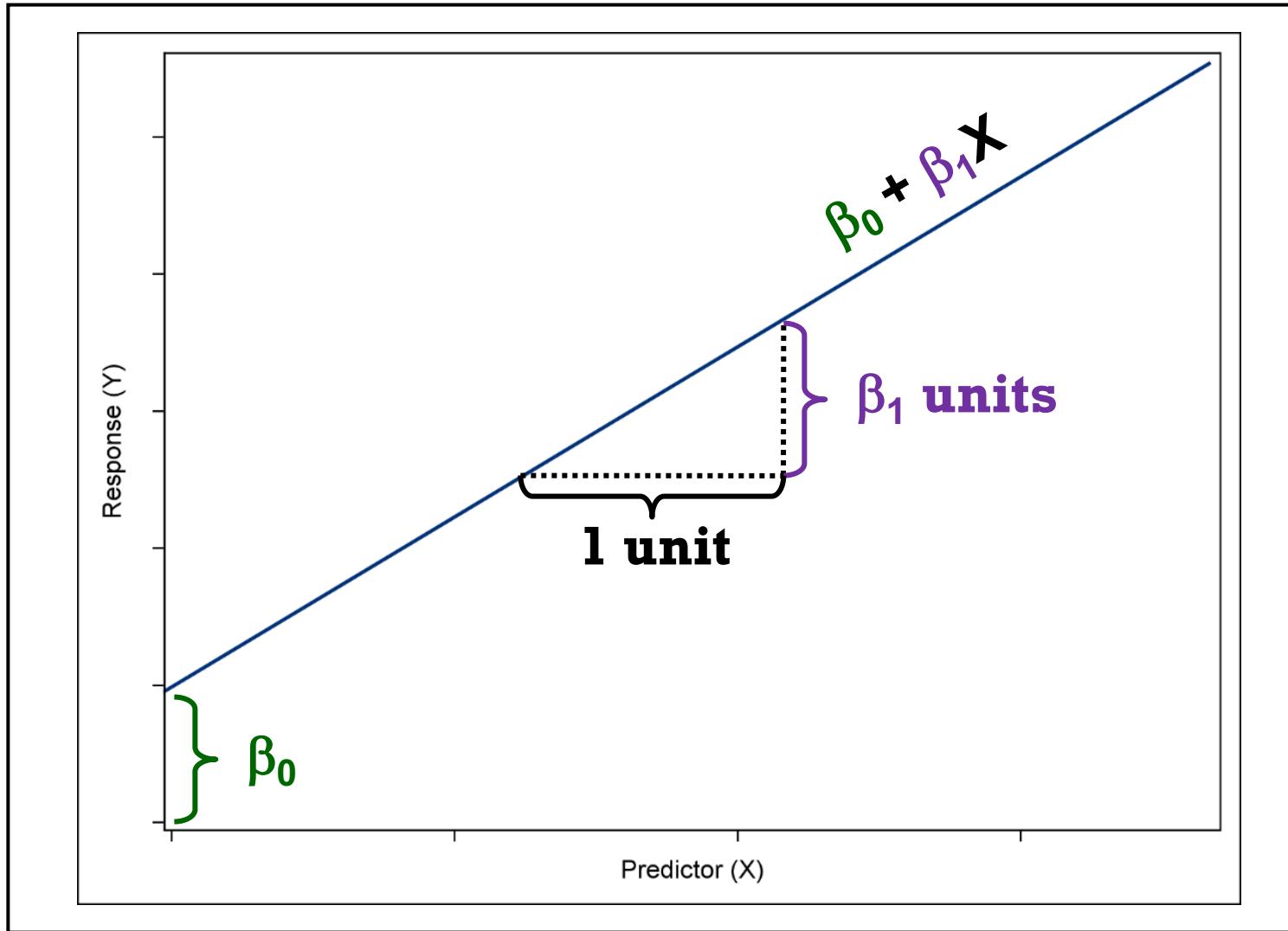
<b>idno</b>	<b>chol (x)</b>	<b>sysbp (y)</b>
1	437	194
2	264	121
3	249	131
4	297	159
5	243	123
6	272	161
7	161	115
รวม	1923	1004



$$\hat{y} = \beta_0 + \beta_1 x$$

$$\widehat{bp} = \beta_0 + \beta_1 chol$$

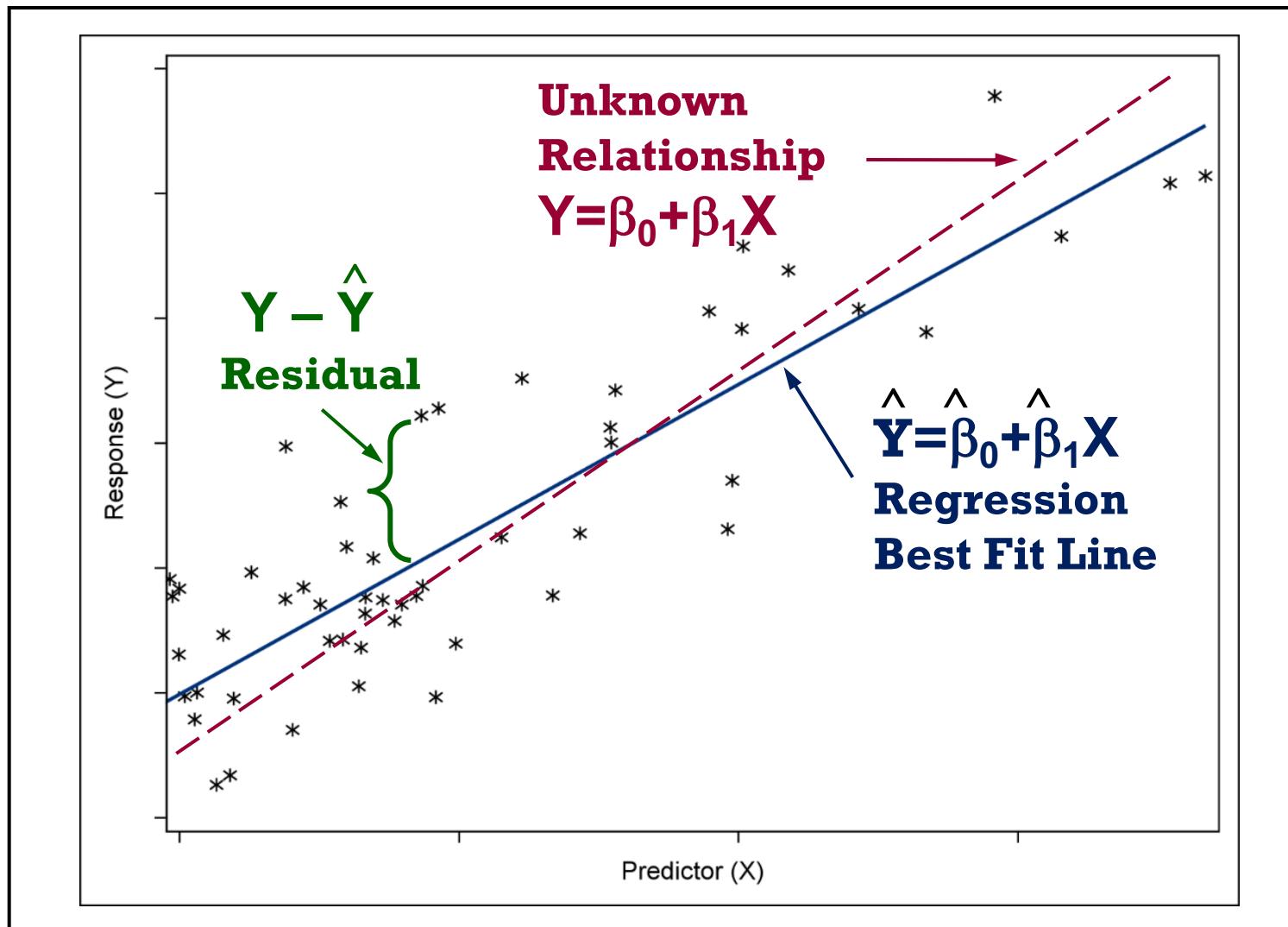
# Simple Linear Regression Model



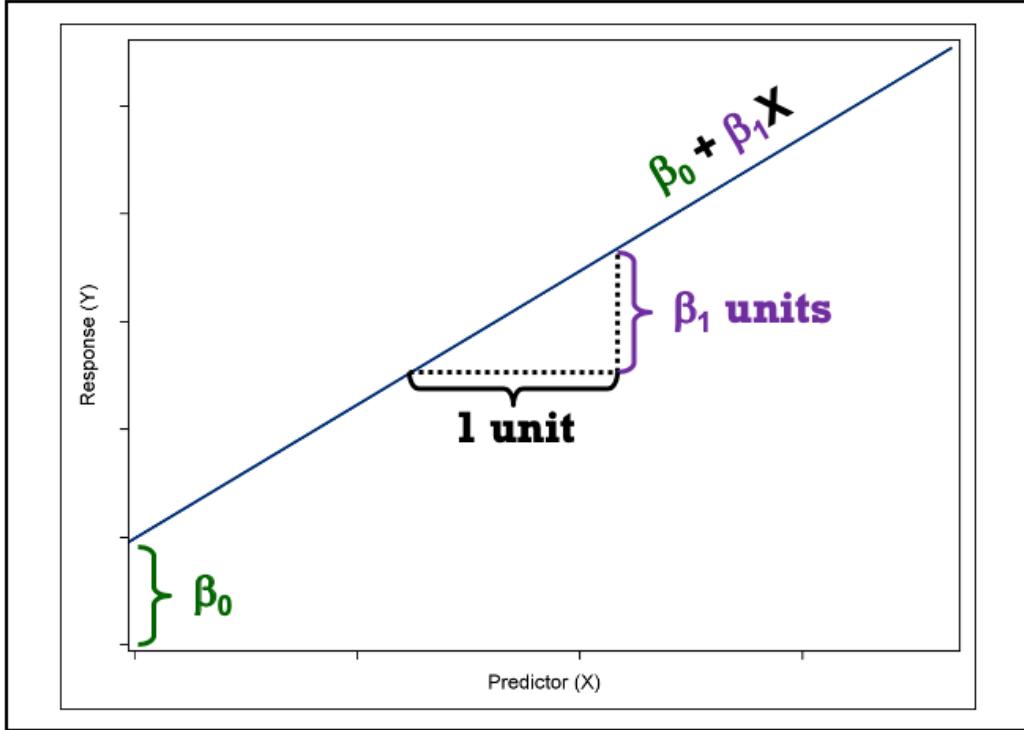
$$\hat{y} = \beta_0 + \beta_1 x$$

$$\widehat{bp} = \beta_0 + \beta_1 chol$$

# Ordinary Least Squares (OLS) Regression



# How to estimate parameters



$$\hat{y} = \beta_0 + \beta_1 x$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_1 = \frac{s_{xy}}{s_{xx}} = \frac{\left( \sum x_i y_i - \frac{\sum x_i \sum y_i}{n} \right)}{\left( \sum x_i^2 - \frac{(\sum x_i)^2}{n} \right)}$$

$$[Y] = [X][\beta]$$

$$[\beta] = [X]^{-1}[Y]$$

# Example

$$\hat{y} = \beta_0 + \beta_1 x$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_1 = \frac{s_{xy}}{s_{xx}} = \frac{\left( \sum x_i y_i - \frac{\sum x_i \sum y_i}{n} \right)}{\left( \sum x_i^2 - \frac{(\sum x_i)^2}{n} \right)}$$

- Systolic Blood Pressure (y)
- Cholesterol (x)

<b>idno</b>	<b>chol (x)</b>	<b>sysbp (y)</b>	<b>x<sup>2</sup></b>	<b>xy</b>	<b>y<sup>2</sup></b>
1	437	194	190969	84778	37636
2	264	121	69696	31944	14641
3	249	131	62001	32619	17161
4	297	159	88209	47223	25281
5	243	123	49049	29889	15129
6	272	161	73984	43792	25921
7	161	115	25921	18515	13225
总数	1923	1004	569829	288760	148994

$$\bar{x} = \frac{1923}{7} = 247.7143, \bar{y} = \frac{1004}{7} = 143.4286$$

$$\beta_1 = \frac{s_{xy}}{s_{xx}} = \frac{\left( 288760 - \frac{1923 \times 1004}{7} \right)}{\left( 569829 - \frac{(1923)^2}{7} \right)} = 0.3116$$

$$\beta_0 = 143.4286 - (0.3116)(247.7143) = 57.8355$$

$$\hat{y} = 57.8355 + 0.3116x$$

$$\widehat{bp} = 57.8355 + 0.3116 \times chol$$

How to read an equation

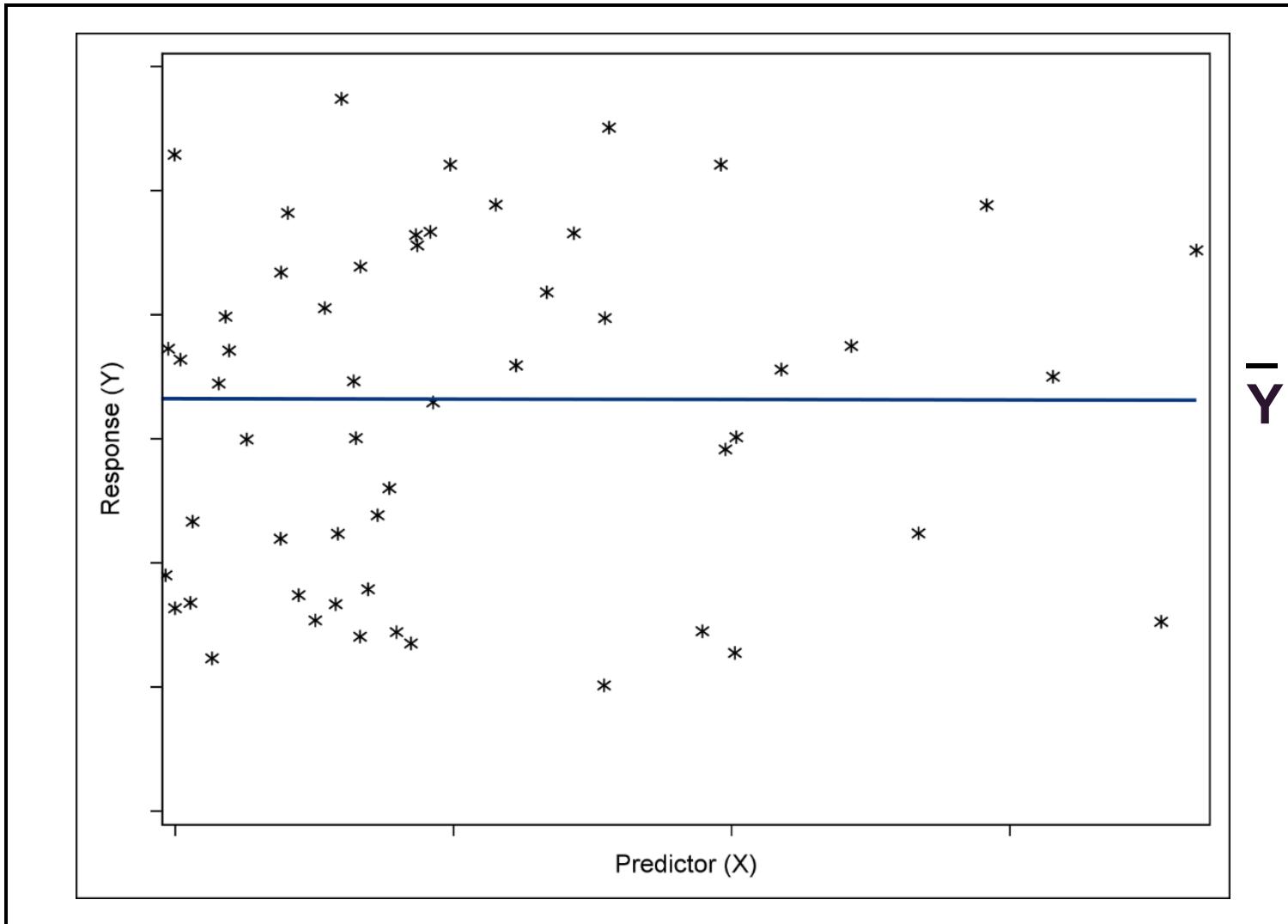
# Example: Prediction

- Systolic Blood Pressure (y)
- Cholesterol (x)

$$\widehat{bp} = 57.8355 + 0.3116 \times chol$$

<b>idno</b>	<b>chol(x)</b>	<b>sysbp(y)</b>	<b>predict</b>
1	437	194	196.1897
2	264	121	141.4179
3	249	131	136.6689
4	297	159	151.8657
5	243	123	134.7693
6	272	161	143.9507
7	161	115	108.8081
รวม	1923	1004	

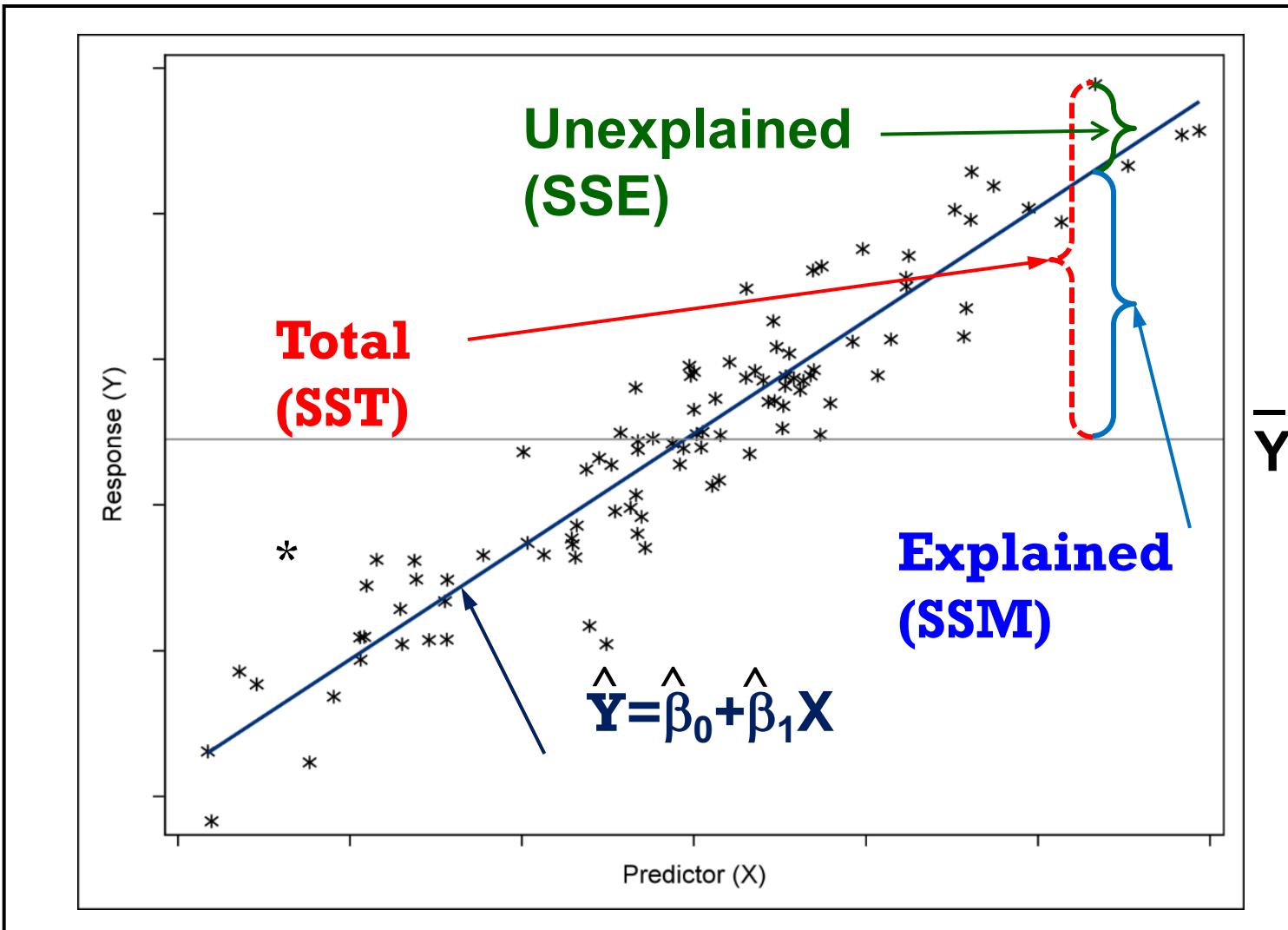
# The Baseline Model (Null Hypothesis)



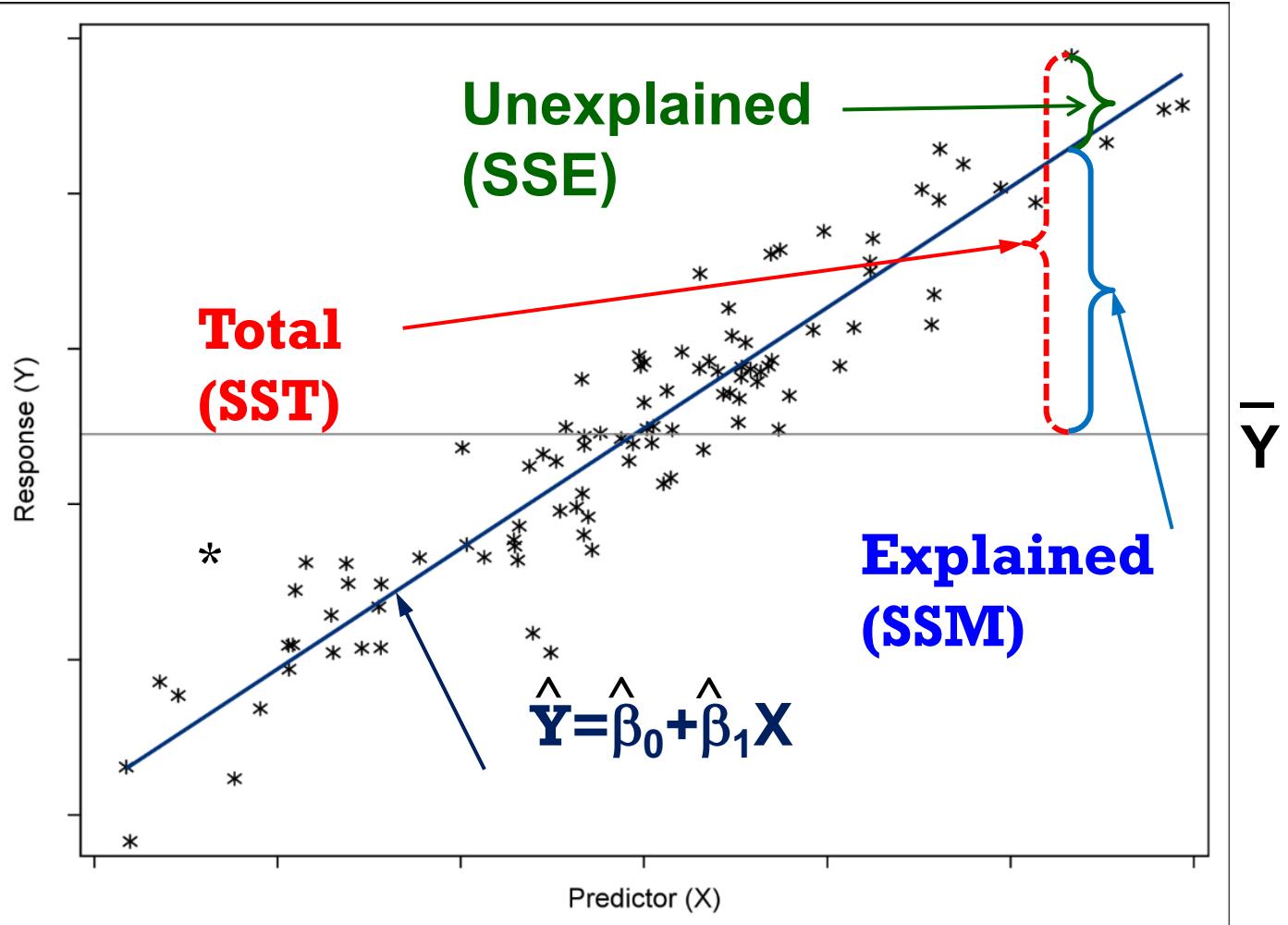
# Explained versus Unexplained Variability

$$SST = SSM + SSE$$

21



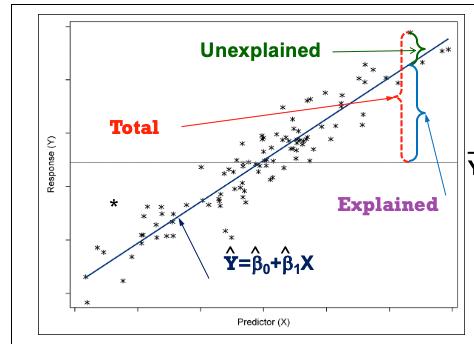
# Coefficient of Determination



$$R^2 = \frac{SSM}{SST} = 1 - \frac{SSE}{SST}$$

- “Proportion of variance accounted for by the model”

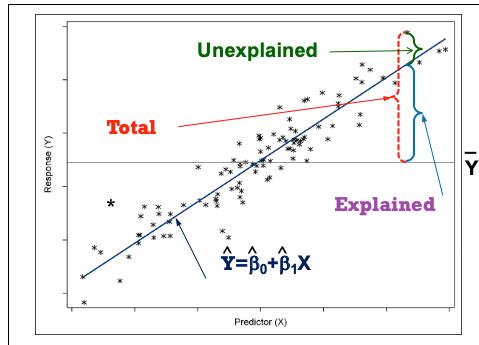
# Coefficient of Determination (cont.)



$$R^2 = \frac{SSM}{SST} = 1 - \frac{SSE}{SST}$$

id	chol (x)	bp (y)	predict	error	squared error (SE)	guess	(y - y_bar )	squared total (ST)
1	437	194	196.1897	(2.1897)	4.7948	143.4286	50.5714	2,557.4694
2	264	121	141.4179	(20.4179)	416.8906	143.4286	(22.4286)	503.0408
3	249	131	136.6689	(5.6689)	32.1364	143.4286	(12.4286)	154.4694
4	297	159	151.8657	7.1343	50.8982	143.4286	15.5714	242.4694
5	243	123	134.7693	(11.7693)	138.5164	143.4286	(20.4286)	417.3265
6	272	161	143.9507	17.0493	290.6786	143.4286	17.5714	308.7551
7	161	115	108.8081	6.1919	38.3396	143.4286	(28.4286)	808.1837
average	274.7143	143.4286		SSE	972.2548	SST		4,991.7143
				MSE	138.8935			
				RMSE	<b>11.7853</b>			
	<b>R^2</b>	<b>1 - (SSE/SST)</b>	<b>0.8052</b>					

# Coefficient of Determination (cont.)



- Train:  $R^2$ , RMSE
- Test:  $R^2$ , RMSE (honest estimate)

Training Data



Testing Data



id	chol (x)	bp (y)	predict	error	squared error (SE)	guess	(y - y_bar )	squared total (ST)
1	437	194	196.1897	(2.1897)	4.7948	143.4286	50.5714	2,557.4694
2	264	121	141.4179	(20.4179)	416.8906	143.4286	(22.4286)	503.0408
3	249	131	136.6689	(5.6689)	32.1364	143.4286	(12.4286)	154.4694
4	297	159	151.8657	7.1343	50.8982	143.4286	15.5714	242.4694
5	243	123	134.7693	(11.7693)	138.5164	143.4286	(20.4286)	417.3265
6	272	161	143.9507	17.0493	290.6786	143.4286	17.5714	308.7551
7	161	115	108.8081	6.1919	38.3396	143.4286	(28.4286)	808.1837
average	274.7143	143.4286		SSE	972.2548	SST		4,991.7143
				MSE	138.8935			
				RMSE	<b>11.7853</b>			
	<b>R^2</b>	<b>1 - (SSE/SST)</b>	<b>0.8052</b>					

# Model Hypothesis Test

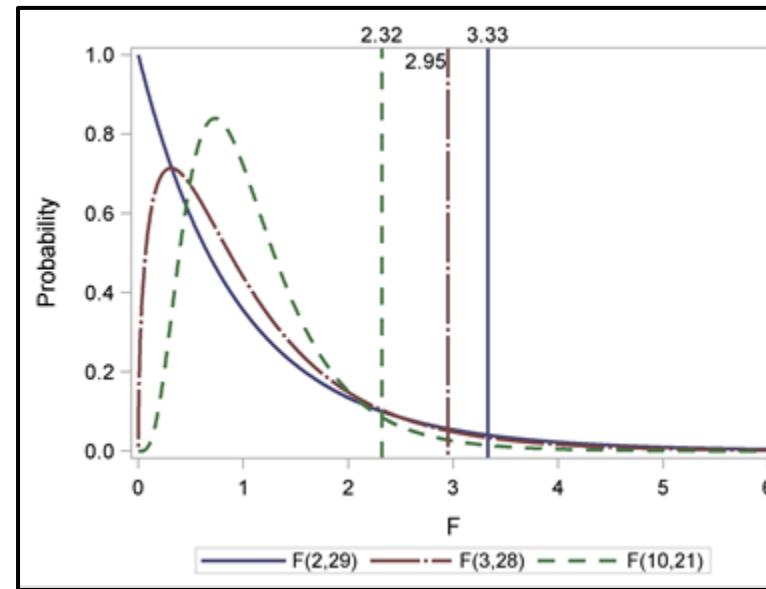
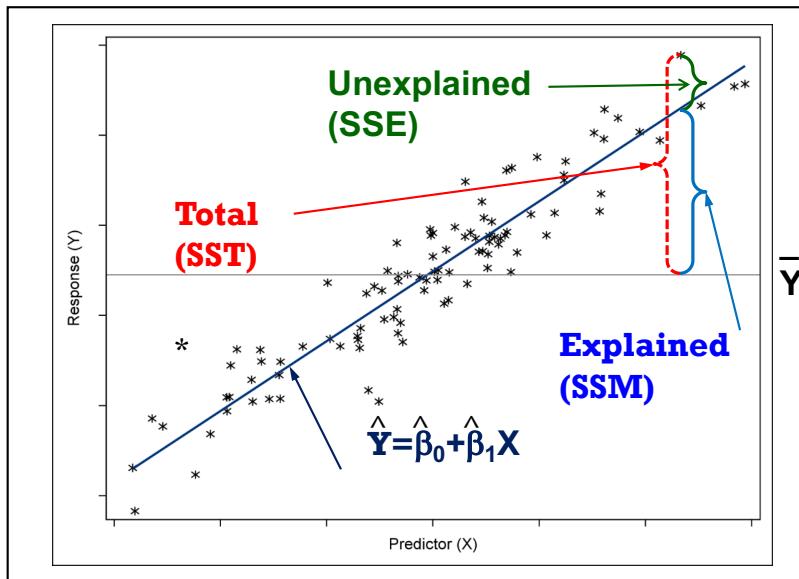
## F Statistic and Critical Values at $\alpha=0.05$

### ■ Null Hypothesis:

- The simple linear regression model does not fit the data better than the baseline model.
- $\beta_1=0$

### ■ Alternative Hypothesis:

- The simple linear regression model does fit the data better than the baseline model.
- $\beta_1 \neq 0$



$$F(\text{Model df, Error df}) = MS_M / MS_E$$

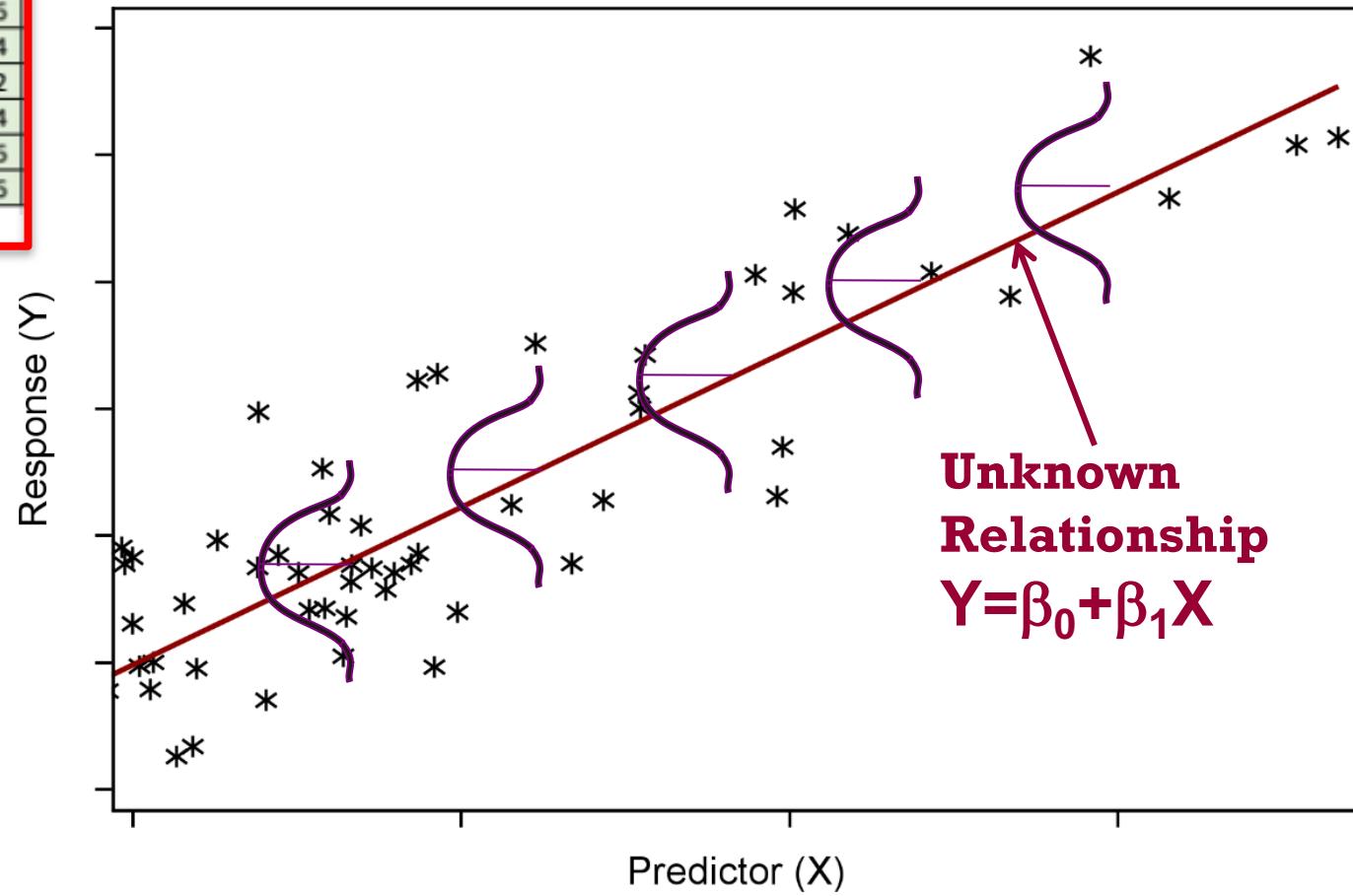
Model df = p-1  
Error df = n-p

Overall df = n-1

# Assumptions of Simple Linear Regression

id	chol (x)	bp (y)	predict	error	squared error (SE)
1	437	194	196.1897	(2.1897)	4.7948
2	264	121	141.4179	(20.4179)	416.8906
3	249	131	136.6689	(5.6689)	32.1364
4	297	159	151.8657	7.1343	50.8982
5	243	123	134.7693	(11.7693)	138.5164
6	272	161	143.9507	17.0493	290.6786
7	161	115	108.8081	6.1919	38.3396

- The mean of the Ys is accurately modeled by a linear function of the X.
- The random error term,  $\varepsilon$ , is assumed to have a normal distribution with a mean of zero.
- The random error term,  $\varepsilon$ , is assumed to have a constant variance,  $\sigma^2$ .
- Not skew**
- The errors are independent.



+

## Multiple Linear Regression

# Multiple Linear Regression with Two Variables

28

- Consider the two-variable model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

- where

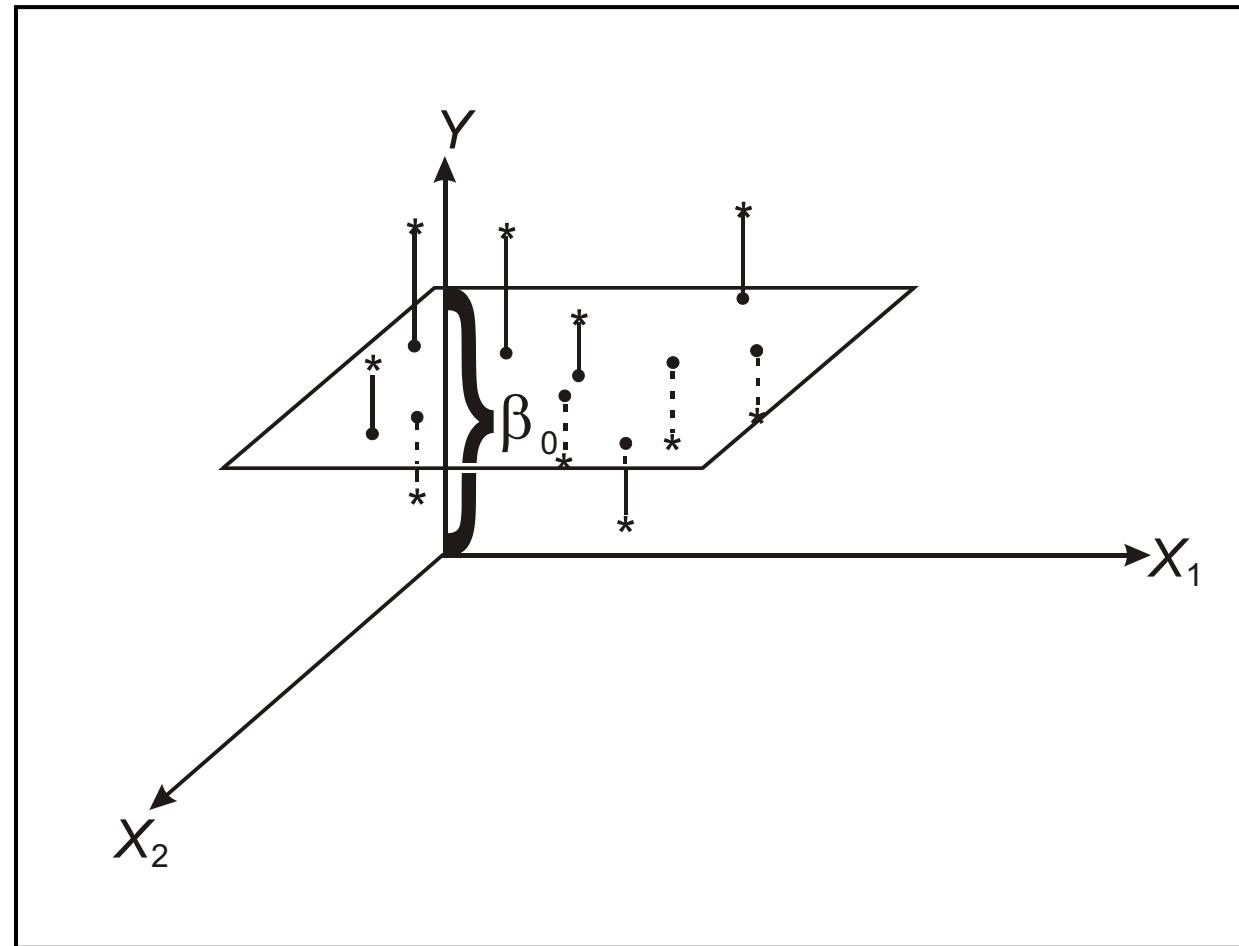
$Y$  is the dependent variable.

$X_1$  and  $X_2$  are the independent or predictor variables.

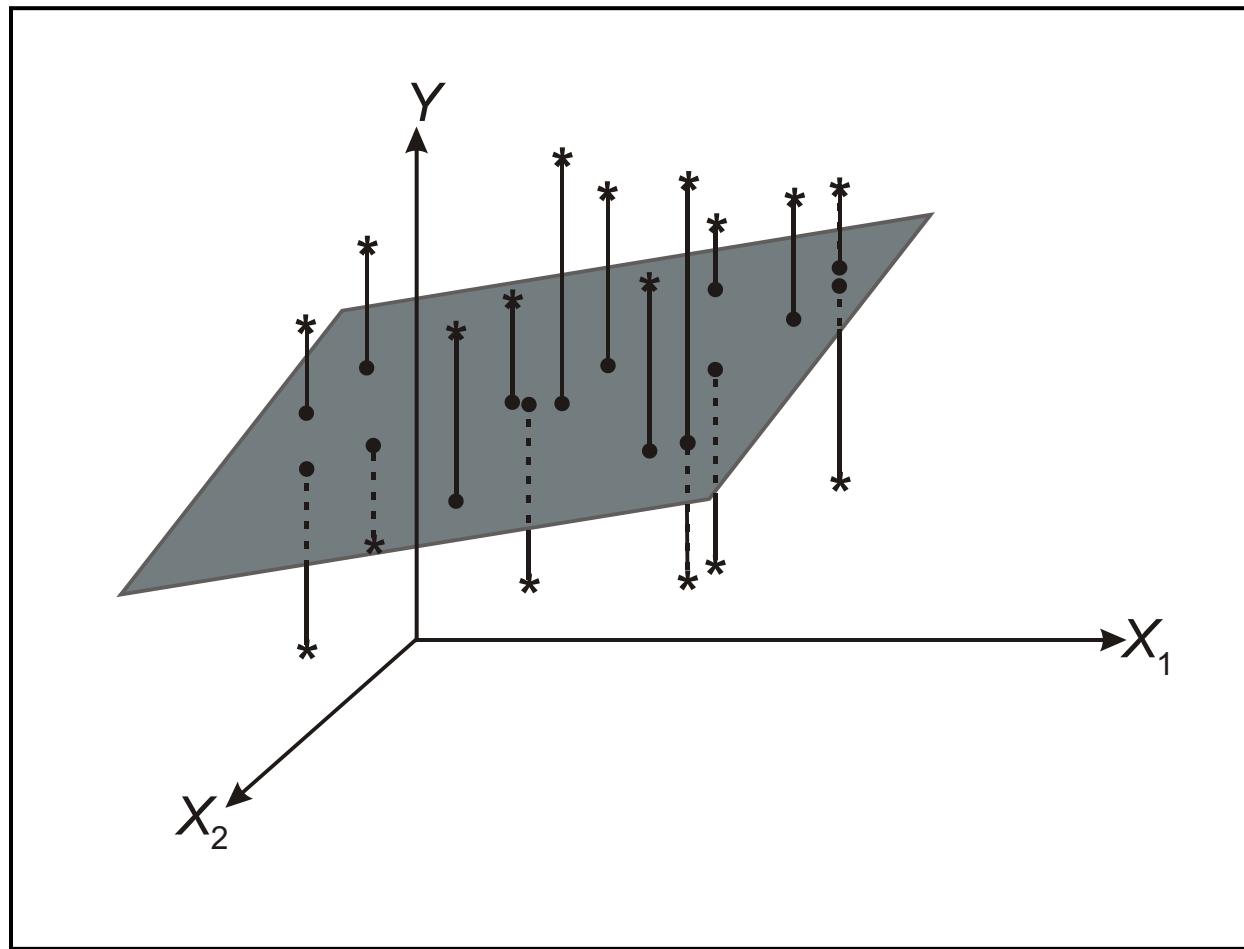
$\varepsilon$  is the error term.

$\beta_0$ ,  $\beta_1$ , and  $\beta_2$  are unknown parameters.

# Picturing the Model: No Relationship



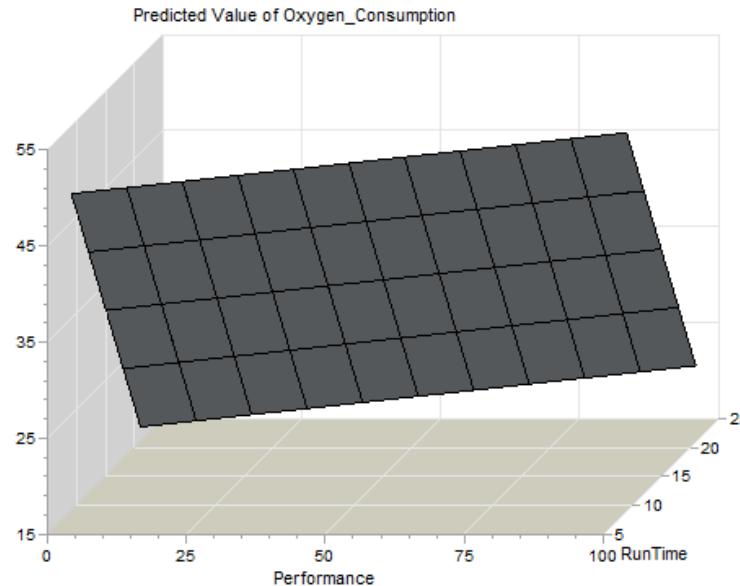
# Picturing the Model: A Relationship



# The Multiple Linear Regression Model

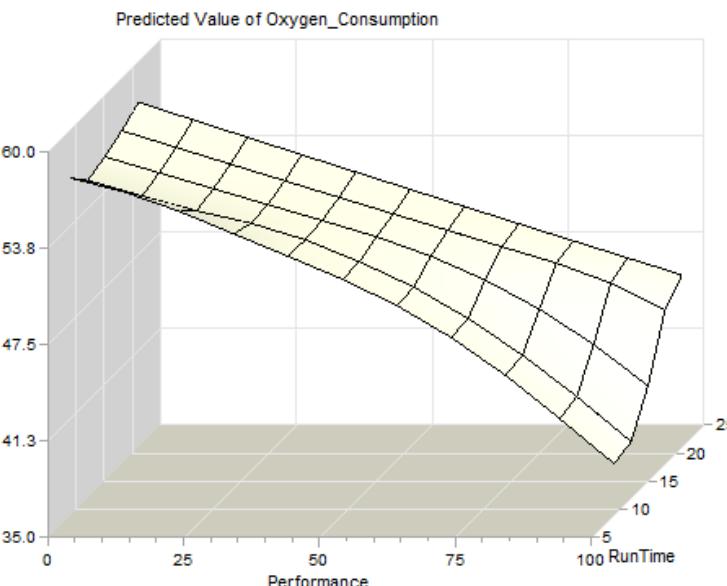
- In general, you model the dependent variable,  $Y$ , as a linear function of  $k$  independent variables,  $X_1$  through  $X_k$ :

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon$$



$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

**Linear Model with  
only Linear Effects**



$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_2 + \beta_4 X_2^2 + \varepsilon$$

**Linear Model with  
Nonlinear Effects**



# The Multiple Linear Regression Model (cont.)

## Matrix Multiplication Approach

inputs		target
Age	Income	Spending
25	25,000	400
35	50,000	500
32	35,000	550

$$Y = \beta_0(1) + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

$$\begin{aligned} [Y] &= [X][\beta] \\ [\beta] &= [X]^{-1}[Y] \end{aligned} \quad \begin{bmatrix} 400 \\ 500 \\ 550 \end{bmatrix} = \begin{bmatrix} 25 & 25000 \\ 35 & 50000 \\ 32 & 35000 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

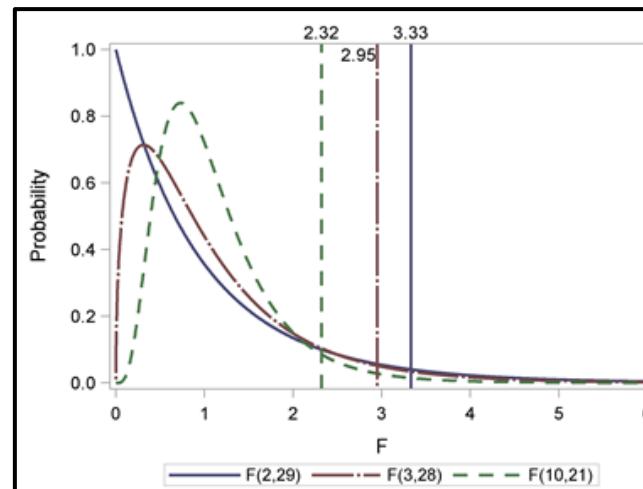
# Model Hypothesis Test (F-Test)

## ■ Null Hypothesis:

- The regression model does not fit the data better than the baseline model.
- $\beta_1 = \beta_2 = \dots = \beta_k = 0$

## ■ Alternative Hypothesis:

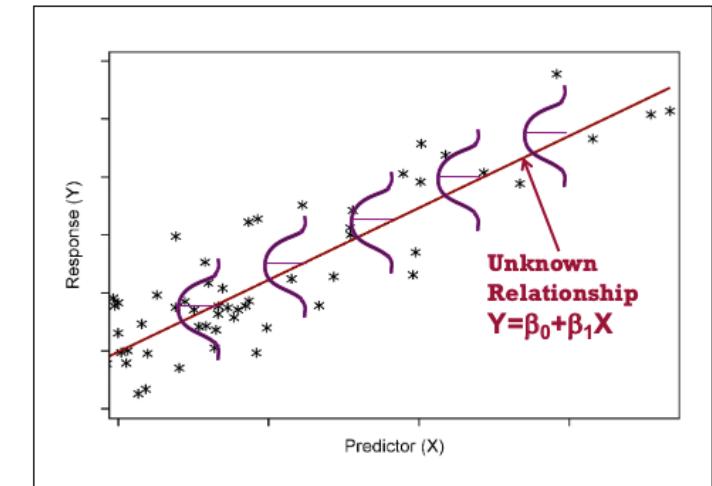
- The regression model does fit the data better than the baseline model.
- Not all  $\beta_i$ s equal zero.



$$F(\text{Model df, Error df}) = MS_M / MS_E$$

# Assumptions for Linear Regression

- The mean of the Ys is accurately modeled by a **linear function** of the X<sub>i</sub>.
  - $(y, x_i) = \text{linear relationship (correlation)}$
- The random error term,  $\varepsilon$ , is assumed to have a **normal distribution** with a mean of zero.
- The random error term,  $\varepsilon$ , is assumed to have a **constant variance**,  $\sigma^2$ .
  - **Not skew**
- The errors are **independent**.



# Multiple Linear Regression versus Simple Linear Regression

## ■ Main Advantage

- Multiple linear regression enables you to **investigate the relationship** among Y and several independent variables simultaneously.

## ■ Main Disadvantages

- Increased complexity makes it **more difficult** to do the following:
  - ascertain which model is “best”
  - interpret the models

# Common Applications of Multiple Regression

36

- Multiple linear regression is a powerful tool for the following tasks:
  - **Prediction** – to develop a model to predict future values of a response variable (Y) based on its relationships with other predictor variables (Xs)
  - **Analytical or Explanatory Analysis** – to develop an understanding of the relationships between the response variable and predictor variables

# Prediction

- The terms in the model, the values of their coefficients, and their statistical significance are of secondary importance.
- The focus is on producing a model that is the best at predicting future values of Y as a function of the Xs. The predicted value of Y is given by this formula:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_k X_k$$

# Analytical or Explanatory Analysis

- The focus is on understanding the relationship between the dependent variable and the independent variables.
- Consequently, the statistical significance of the coefficients is important as well as the magnitudes and signs of the coefficients.

$$\hat{Y} = \underline{\hat{\beta}_0} + \underline{\hat{\beta}_1}X_1 + \dots + \underline{\hat{\beta}_k}X_k$$

$$R^2 = \frac{SSM}{SST} = 1 - \frac{SSE}{SST} \quad \text{Adjusted R Square}$$

$$R_{ADJ}^2 = 1 - \frac{(n-i)(1-R^2)}{n-p}$$

- $i=1$  if there is an intercept and 0 otherwise
- $n$ =the number of observations used to fit the model
- $p$ =the number of parameters in the model

$$R_{ADJ}^2 = 1 - \frac{(100-1)(1-0.7)}{(100-1)} = 0.70; R^2 = 0.7, n = 100, p = 1$$

$$R_{ADJ}^2 = 1 - \frac{(100-1)(1-0.7)}{(100-3)} = 0.69; R^2 = 0.7, n = 100, p = 3$$

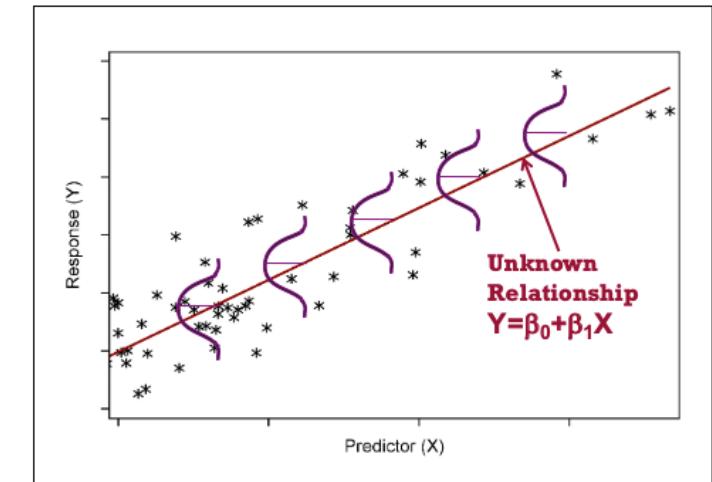
$$R_{ADJ}^2 = 1 - \frac{(100-1)(1-0.7)}{(100-10)} = 0.67; R^2 = 0.7, n = 100, p = 10$$

+

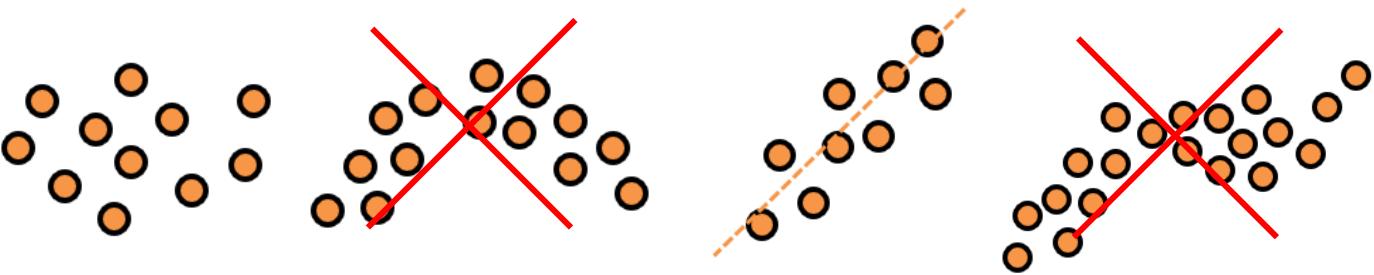
Other topics

# Assumptions for Linear Regression

- The mean of the Ys is accurately modeled by a **linear function** of the X<sub>i</sub>.
  - $(y, x_i) = \text{linear relationship (correlation)}$
  
- The random error term,  $\varepsilon$ , is assumed to have a **normal distribution** with a mean of zero.
- The random error term,  $\varepsilon$ , is assumed to have a **constant variance**,  $\sigma^2$ .
  - **Not skew**
  
- The errors are **independent**.



# Pearson Correlation



- Pearson Correlation, which is the Pearson Product Moment Correlation (PPMC), is used to evaluate **linear relationships** between two **continuous variable**
- Here's the most commonly used formula to find the Pearson correlation coefficient, which can be called Pearson's R:

$$r = \frac{\sum (x_i - \bar{x}_{\text{average}}) (y_i - \bar{y}_{\text{average}})}{\sqrt{\sum (x_i - \bar{x}_{\text{average}})^2 * \sum (y_i - \bar{y}_{\text{average}})^2}}$$

$$\text{VAR}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - E(X))^2 \quad (1)$$

$$\text{COV}(X, Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - E(X))(y_i - E(Y)) \quad (2)$$

$$\text{COR}(X, Y) = \frac{\text{COV}(X, Y)}{\sqrt{\text{VAR}(X)\text{VAR}(Y)}} \quad (3)$$

$$R^2 = 1 - \frac{\text{VAR}(X, Y)_{\text{FittedLine}}}{\text{VAR}(X, Y)_{\text{Mean}}} \quad (4)$$

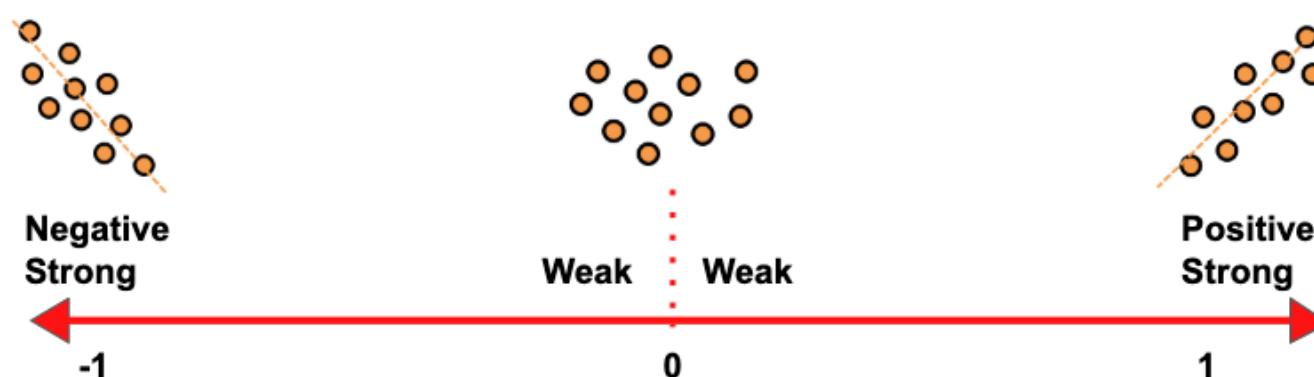
$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

Covariance normalized by Standard Deviation

↓  
Correlation between X and Y  
↓  
Standard deviation of X  
↓  
Standard deviation of Y

# Correlation Coefficient

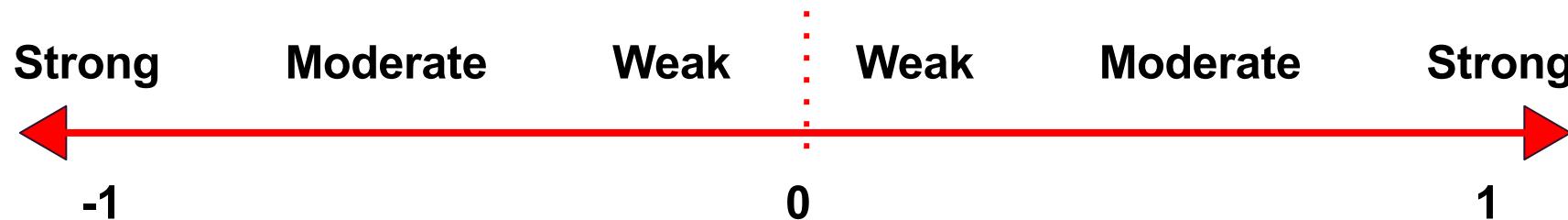
- The numerical measure of the degree of association between two continuous variables is called the **correlation coefficient (r)**.
- The coefficient value is always between **-1 and 1** and it measures both the **strength** and **direction** of the linear relationship between the variables.



# Correlation Coefficient (cont.): Strength

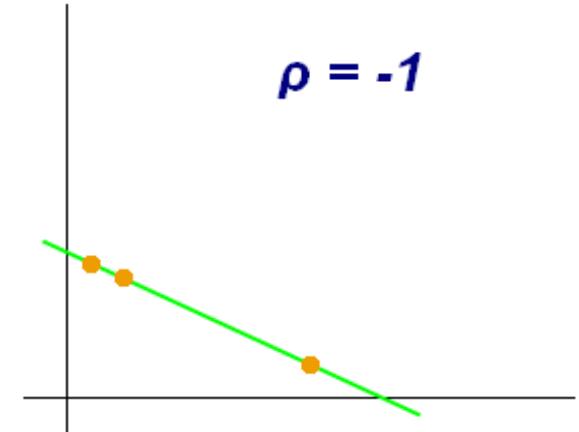
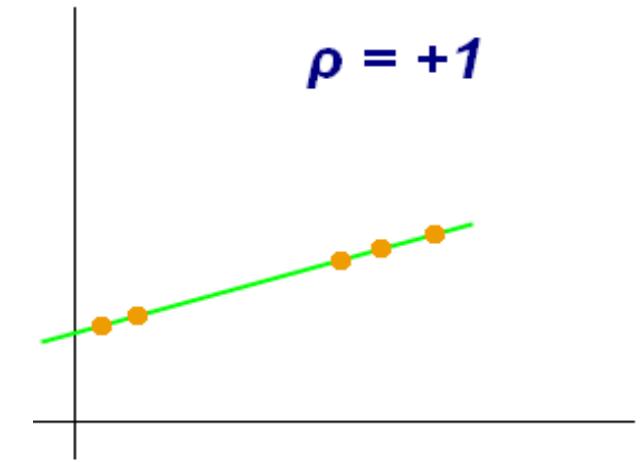
- **Strength**

- The values of **-1 and 1** indicate a perfect **linear relationship** when all the data points fall on a line. Normally, either positive or negative, is **rarely** found.
- A coefficient of **0** indicates no linear relationship between the variables. This is what you are likely to get with two sets of random numbers.
- Values **between 0 and +1/-1** represent a scale of weak, moderate and strong relationships. As the coefficient gets closer to either -1 or 1, the strength of the relationship increases.



# Correlation Coefficient (cont.): Direction

- **Direction**
  - **Positive coefficients** represent **direct** linear association (upward-sloping)
  - **Negative coefficients** represent **inverse** linear association (downward-sloping)

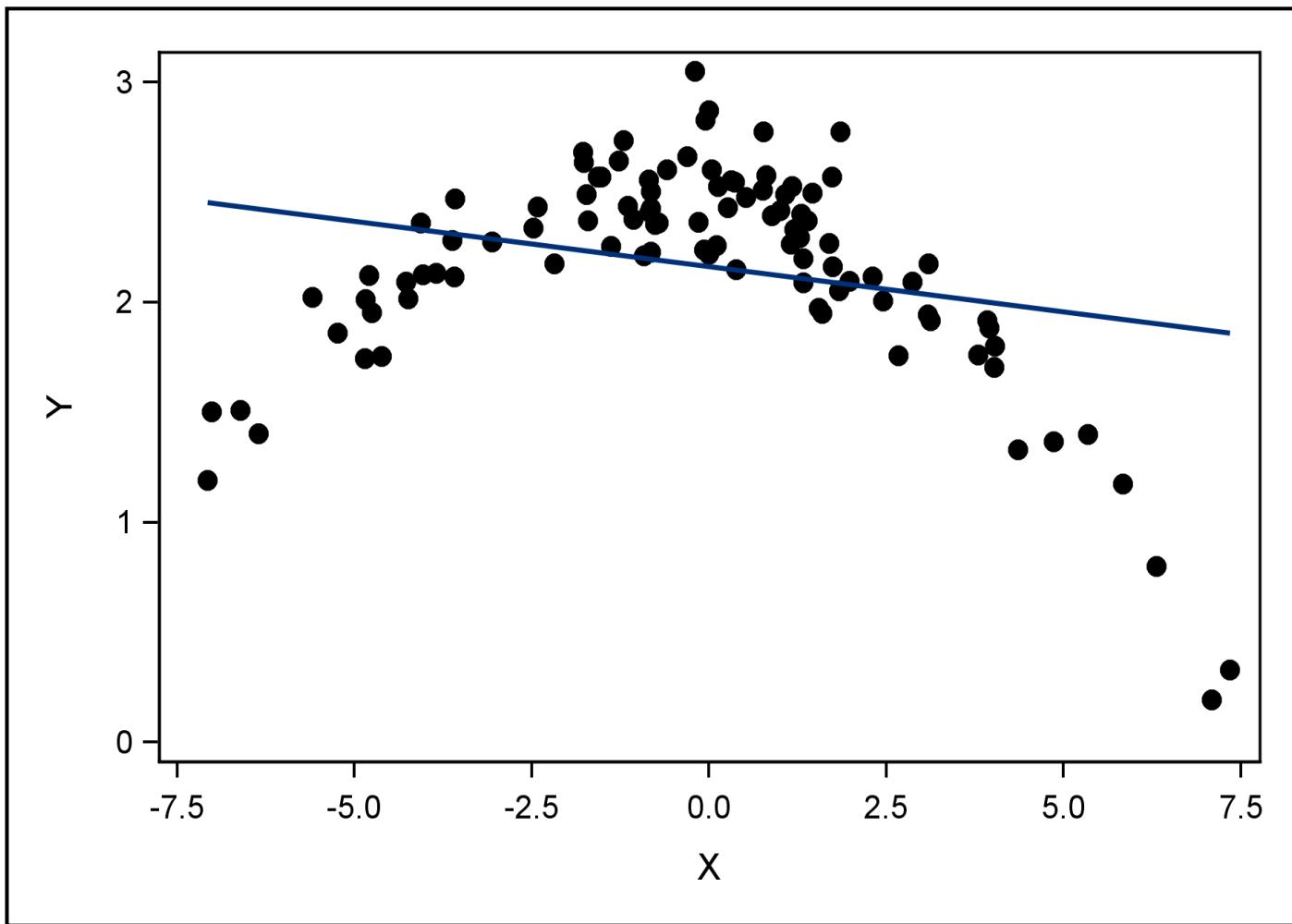


# Hypothesis Test for a Correlation

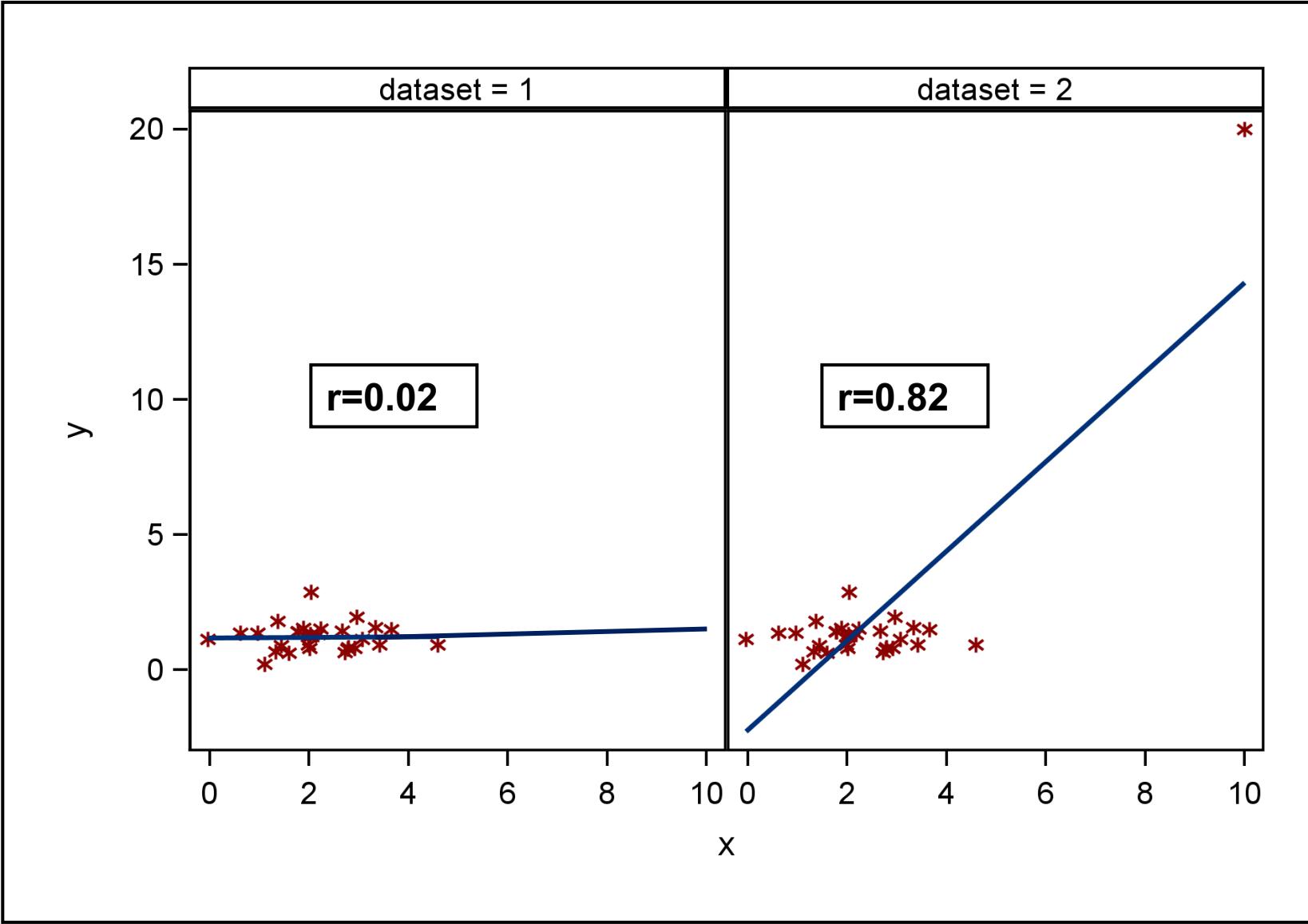
- The parameter representing correlation is  $\rho$ .
- $\rho$  is estimated by the sample statistic  $r$ .
  
- $H_0: \rho=0$
- Rejecting  $H_0$  indicates only great confidence that  $\rho$  is not exactly zero.
- A  $p$ -value does not measure the magnitude of the association.
- Sample size affects the  $p$ -value.
  
- The test focuses on an existence of relationship, **but not the strength**.

# Remark 1: Missing Another Type of Relationship

47



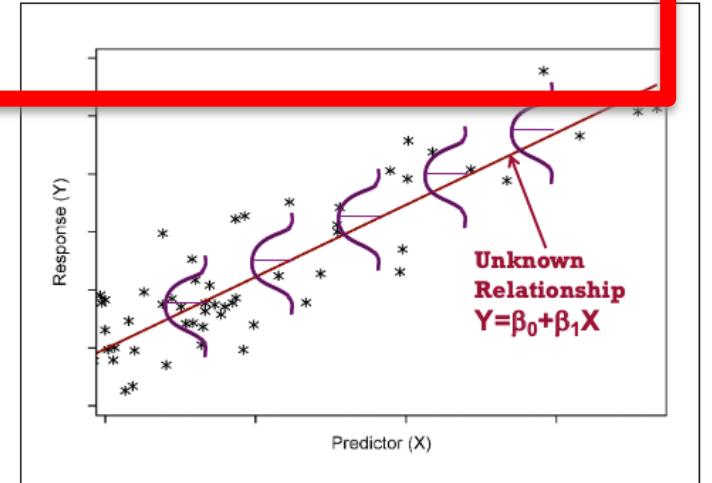
## Remark2: Extreme Data Values



# Assumptions for Linear Regression

- The mean of the Ys is accurately modeled by a **linear function** of the X<sub>i</sub>.
  - (y, x<sub>i</sub>) = **linear relationship (correlation)**

- The random error term,  $\varepsilon$ , is assumed to have a **normal distribution** with a mean of zero.
- The random error term,  $\varepsilon$ , is assumed to have a **constant variance**,  $\sigma^2$ .
  - **Not skew**
- The errors are **independent**.

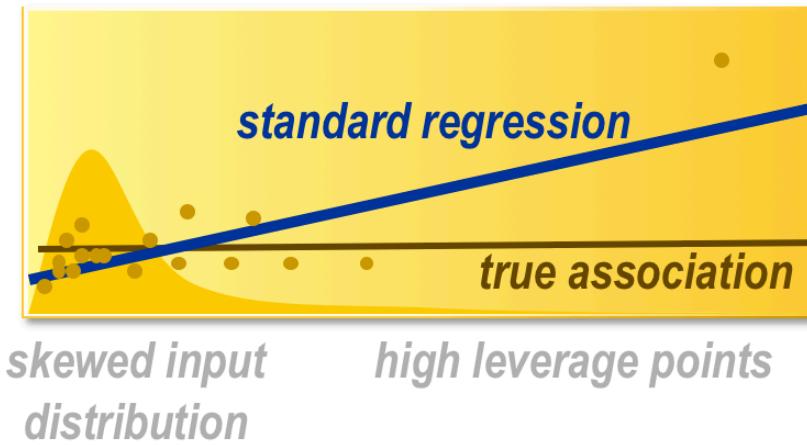
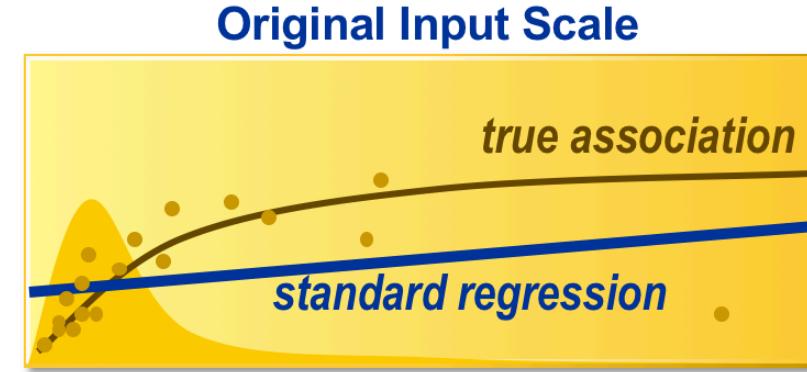




$$\hat{y} = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$



# Outliers & Feature transformation

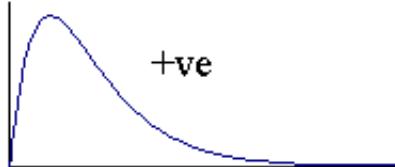


- Skewness

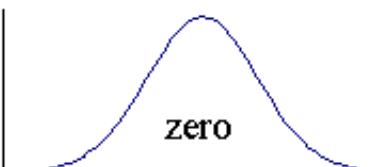
- Example: Salary, Balance in bank account

- Solutions: Log, Binning**

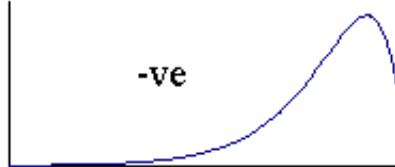
Skewness



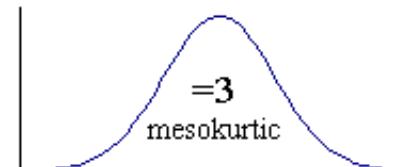
zero



-ve



Kurtosis



>3 leptokurtic





+

Demo





# House Price Prediction: Target=MEDV

## Simple Linear Regression

LinearRegression\_HousePrice.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Comment Share

Table of contents + Code + Text RAM Disk Editing

House Price Prediction (Linear Regression)

Section

House Price Prediction (Linear Regression)

```
[101] 1 # Import libraries necessary for this project
      2 import numpy as np
      3 import pandas as pd
      4 from pandas import set_option
```

# + House Price Prediction: Target=MEDV

## In-Class Activity: Multiple Linear Regression

The screenshot shows a Jupyter Notebook interface with the following details:

- Title:** LinearRegression\_HousePrice.ipynb
- Toolbar:** File, Edit, View, Insert, Runtime, Tools, Help, All changes saved.
- Header:** Comment, Share, User icon.
- Table of Contents:** House Price Prediction (Linear Regression) (selected), Section.
- Section Content:** House Price Prediction (Linear Regression) (with a minus sign icon).
- Image:** A photograph of a residential street with several brick houses under a blue sky with clouds.
- Code Cell (Visible Part):**

```
[101] 1 # Import libraries necessary for this project
2 import numpy as np
3 import pandas as pd
4 from pandas import set_option
```

```

▶ 1 #Correlation with output variable
2 cor_target = abs(cor["MEDV"])
3
4 #Selecting highly correlated features
5 relevant_features = cor_target[cor_target>0.5]
6 relevant_features

▷ RM      0.70
PTRATIO  0.51
LSTAT    0.74
MEDV     1.00
Name: MEDV, dtype: float64

```

Select three inputs  
Write an regression equation  
Report RMSE & R<sup>2</sup> on test

