

# Multiple Linear Regression

Machine Learning

Dr. Adnan Abid

# Regressions

## Simple Linear Regression

$$y = b_0 + b_1 * x_1$$

## Multiple Linear Regression

Dependent variable (DV)

Independent variables (IVs)


$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

# Dummy Variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + ???$$

# Dummy Variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

New York	California

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + ???$$

# Dummy Variables

Profit	R&D Spend	Admin	Marketing	State	New York	California
192,261.83	165,349.20	136,897.80	471,784.10	New York	1	0
191,792.06	162,597.70	151,377.59	443,898.53	California	0	1
191,050.39	153,441.51	101,145.55	407,934.54	California	0	1
182,901.99	144,372.41	118,671.85	383,199.62	New York	1	0
166,187.94	142,107.34	91,391.77	366,168.42	California	0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + ???$$



# Dummy Variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

## Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1$$

1. Add dummy variables for all possible values in the attribute with categorical data, which is STATE in this example.
2. Populate them with One Hot Encoding method
3. Replace the column with Dummy variables

# Dummy Variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

## Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1$$

1. Add dummy variables for all possible values in the attribute with categorical data, which is STATE in this example.
2. Populate them with One Hot Encoding method
3. Replace the column with Dummy variables

1. D1 works like a switch, when it is 1 then it means New York and in case of 0 it is California
2. Thus we may not include the last one, as the information will be complete even when we eliminate it.-

# Dummy Variables

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

## Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1$$



1. D1 works like a switch, when it is 1 then it means New York and in case of 0 it is California
2. Thus we may not include the last one, as the information will be complete even when we eliminate it.

1. Theoretically, we SHOULD NOT include all the dummy variables in the equation.
2. Intuition is that when D1 is 0, then  $b_4 * D_1$  is also 0, and we may think that there is no coefficient for California.
3. Well, the coefficient for California is going to be included in  $b_0$ , i.e. when D1 is 0 the equation becomes an expression for California.
4. When D1 is 1 then the expression becomes the difference between New York and California



# Dummy Variable Trap

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

## Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3$$



$$+ b_4 * D_1 + \underline{b_5 * D_2}$$



Activate Windows  
Go to Settings to activate Windows.

1. The fact that one or more independent variables in a linear regression predict another is called multi-collinearity.
2. If so, the model cannot distinguish the facts of D1 from the facts of D2.
3. This is called Dummy Variable Trap.
4. Thus, you cannot have  $b_0$ ,  $D_1$ , and  $D_2$  in the model at the same time.

# Dummy Variable Trap

					Dummy Variables	
Profit	R&D Spend	Admin	Marketing	State	New York	California
192,261.83	165,349.20	136,897.80	471,784.10	New York	1	0
191,792.06	162,597.70	151,377.59	443,898.53	California	0	1
191,050.39	153,441.51	101,145.55	407,934.54	California	0	1
182,901.99	144,372.41	118,671.85	383,199.62	New York	1	0
166,187.94	142,107.34	91,391.77	366,168.42	California	0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1 + \underline{b_5 * D_2}$$

1. Never include all the dummy variables.
2. Always ignore 1 dummy variable. If you have 10 dummy variables, then include 9 and leave 1; if you have 100 dummy variables include 99 and leave 1.
3. For instance, if you have another variable that tell about the target sector of the organization, then we shall follow the same process all over for this variable too.

# Dummy Variable Trap

Profit	R&D Spend	Admin	Marketing	State
192,261.83	165,349.20	136,897.80	471,784.10	New York
191,792.06	162,597.70	151,377.59	443,898.53	California
191,050.39	153,441.51	101,145.55	407,934.54	California
182,901.99	144,372.41	118,671.85	383,199.62	New York
166,187.94	142,107.34	91,391.77	366,168.42	California

## Dummy Variables

New York	California
1	0
0	1
0	1
1	0
0	1

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3$$

$$+ b_4 * D_1 + \cancel{b_5 * D_2}$$

Always omit one  
dummy variable

# Building A Model

---

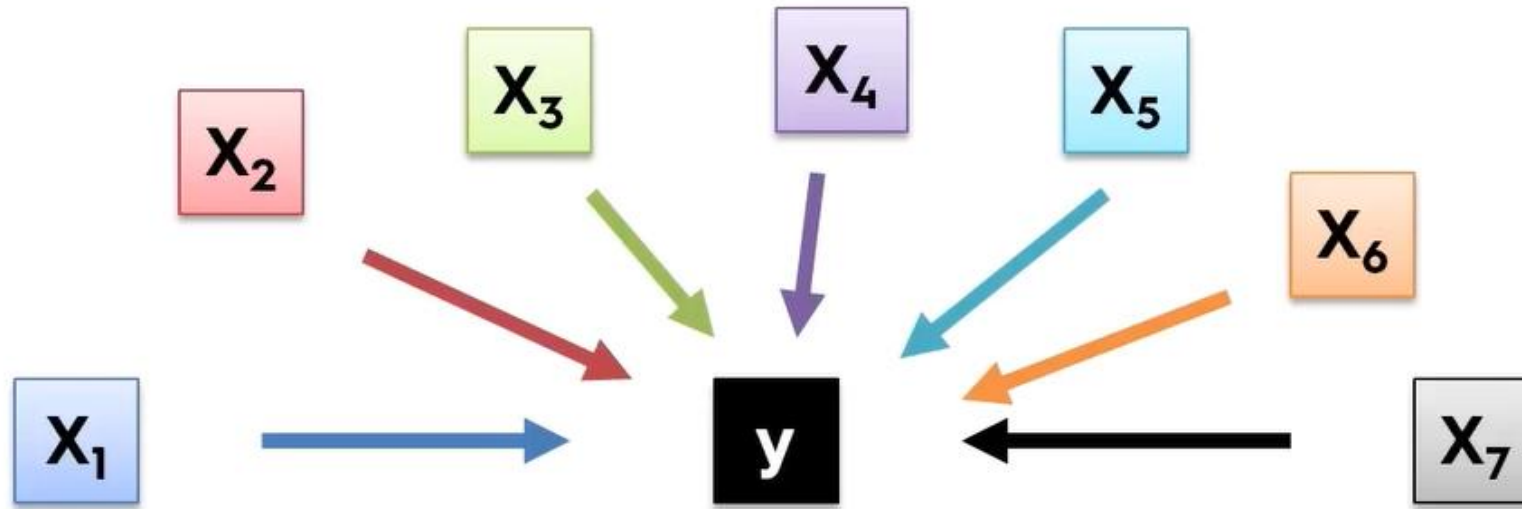
## (Step-By-Step)

Activate Windows  
Go to Settings to activate Windows.



# Building A Model

---



1. There are a lot of variables
2. Should we include all the independent variables?
3. In fact, we need to include the ones which are helpful in prediction, and drop the rest.
4. There are many methods for doing it...

# Building A Model

---

## **5 methods of building models:**

1. All-in
2. Backward Elimination
3. Forward Selection
4. Bidirectional Elimination
5. Score Comparison

# Building A Model

---

## **“All-in” – cases:**

- Prior knowledge; OR
- You have to; OR
- Preparing for Backward Elimination



# Building A Model



## Backward Elimination

**STEP 1:** Select a significance level to stay in the model (e.g.  $SL = 0.05$ )



**STEP 2:** Fit the full model with all possible predictors



**STEP 3:** Consider the predictor with the highest P-value. If  $P > SL$ , go to STEP 4, otherwise go to FIN



**STEP 4:** Remove the predictor



**STEP 5:** Fit model without this variable\*

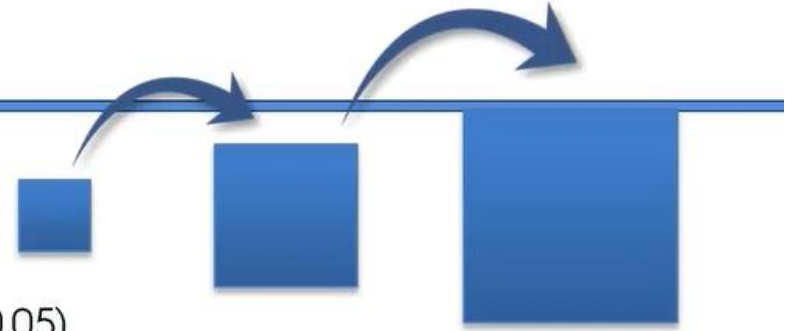


1. Remove the attribute with highest P-value, if the P-value is greater than the threshold Significance Level.
2. Then go to Step 3 and rebuild the model
3. If there is not such variable with P-value higher than SL, the go to FINISH (FIN) state, which means the model is ready and all the remain attributes are statistically significant for prediction.



# Building A Model

## Forward Selection



**STEP 1:** Select a significance level to enter the model (e.g.  $SL = 0.05$ )



**STEP 2:** Fit all simple regression models  $y \sim x_n$ . Select the one with the lowest P-value



**STEP 3:** Keep this variable and fit all possible models with one extra predictor added to the one(s) you already have



**STEP 4:** Consider the predictor with the lowest P-value. If  $P < SL$ , go to STEP 3, otherwise go to FIN

1. Add one variable at a time.
2. Looks reverse of backward elimination, but is much more tedious and complex.
3. Start with a single variable (Simple Regression) like manner.
4. Choose the variable which has lowest SL.
5. Then add each of the remaining variables one by one; and choose the one with lowest SL value.
6. Keep doing it till you find variable with P values  $< SL$ .
7. If you are unable to find any variable which, if added, to the model does not offer  $PL < SL$ .

# Building A Model

## Forward Selection

**STEP 1:** Select a significance level to enter the model (e.g.  $SL = 0.05$ )



**STEP 2:** Fit all simple regression models  $y \sim x_n$ . Select the one with the lowest P-value



**STEP 3:** Keep this variable and fit all possible models with one extra predictor added to the one(s) you already have

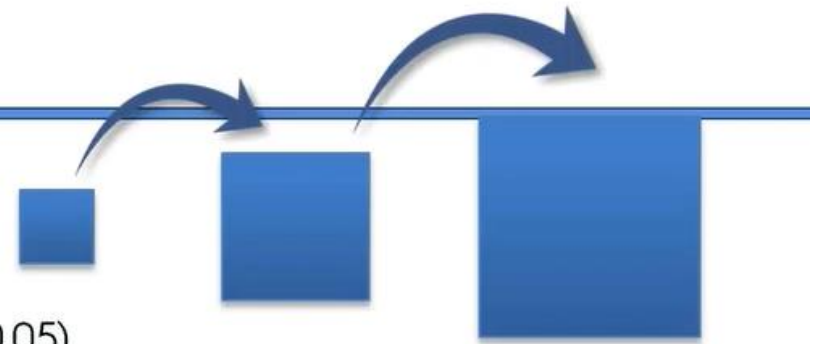


**STEP 4:** Consider the predictor with the lowest P-value. If  $P < SL$ , go to STEP 3, otherwise go to FIN

1. Stop when you are unable to find a variable where  $P < SL$
2. IMPORTANT: Use the previous model, as this attribute does not satisfy the SL criteria. Thus this model should not be included in the model.



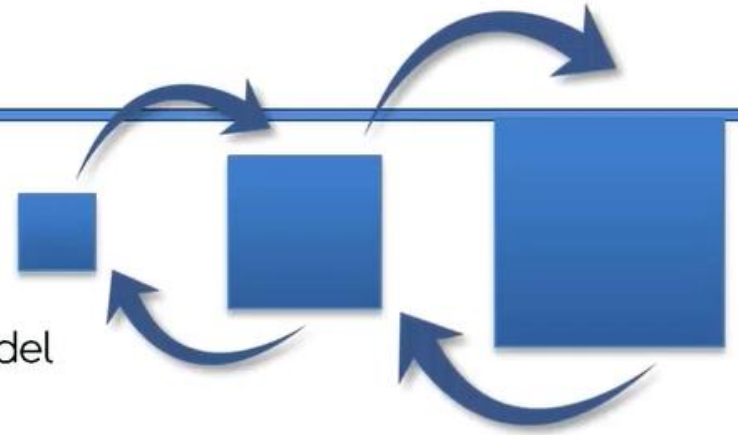
**FIN:** Keep the previous model



# Building A Model

## Bidirectional Elimination

**STEP 1:** Select a significance level to enter and to stay in the model  
e.g.: SLENTER = 0.05, SLSTAY = 0.05



**STEP 2:** Perform the next step of Forward Selection (new variables must have:  $P < \text{SLENTER}$  to enter)

**STEP 3:** Perform ALL steps of Backward Elimination (old variables must have  $P < \text{SLSTAY}$  to stay)

**STEP 4:** No new variables can enter and no old variables can exit

**FIN:** Your Model Is Ready

Activate Windows  
Go to Settings to activate Windows.

# Building A Model

## All Possible Models

**STEP 1:** Select a criterion of goodness of fit (e.g. Akaike criterion)



**STEP 2:** Construct All Possible Regression Models:  $2^N - 1$  total combinations



**STEP 3:** Select the one with the best criterion





# Building A Model

## All Possible Models

**STEP 1:** Select a criterion of goodness of fit (e.g. Akaike criterion)



**STEP 2:** Construct All Possible Regression Models:  $2^N - 1$  total combinations



**STEP 3:** Select the one with the best criterion



**FIN:** Your Model Is Ready



**Example:**  
**10 columns means**  
**1,023 models**

Activate Windows  
Go to Settings to activate Windows.

# Building A Model

---

## **5 methods of building models:**

1. All-in
2. Backward Elimination
3. Forward Selection
4. Bidirectional Elimination
5. Score Comparison

# Polynomial Linear Regression

# Regressions

Simple  
Linear  
Regression

$$y = b_0 + b_1 x_1$$

Multiple  
Linear  
Regression

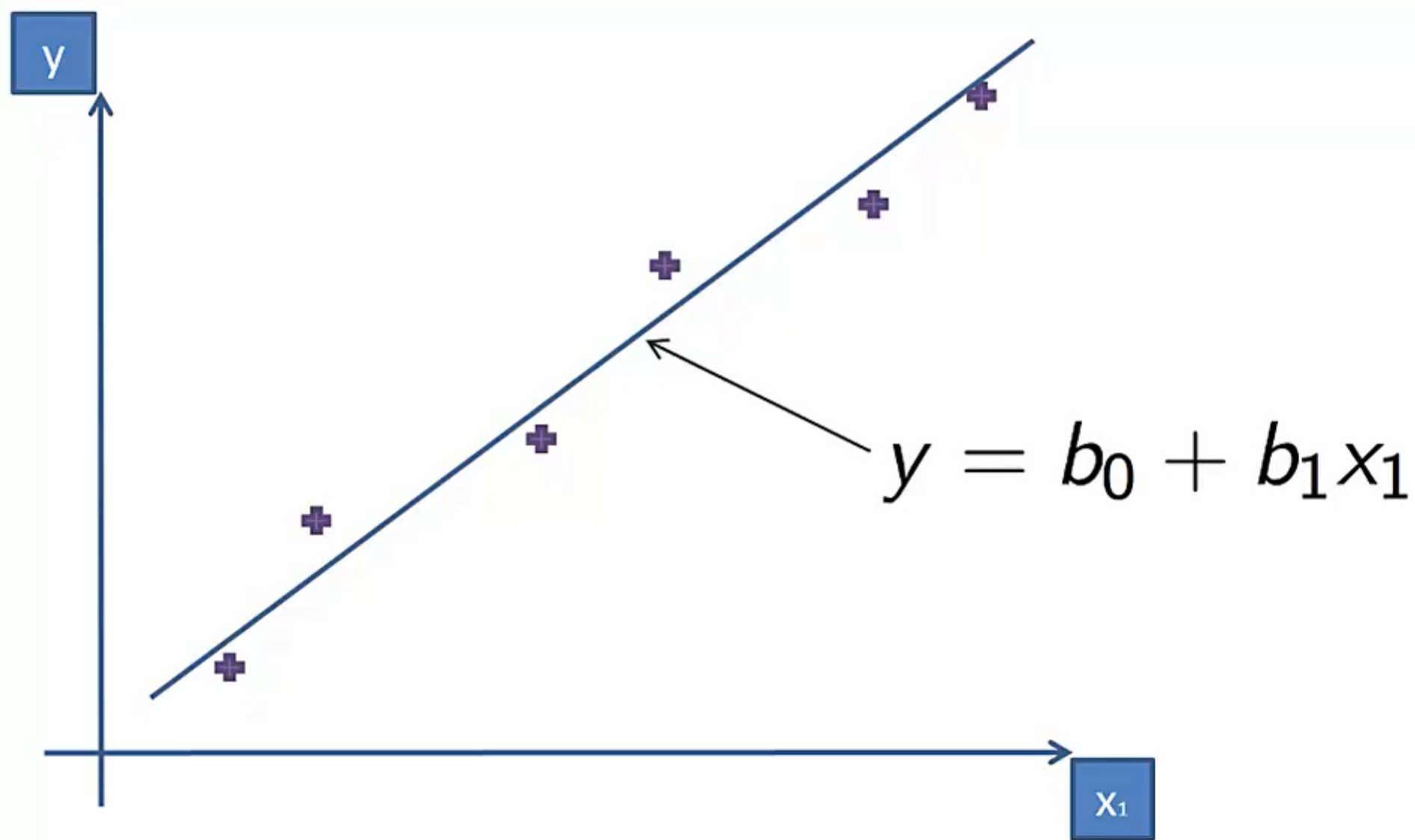
$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

Polynomial  
Linear  
Regression

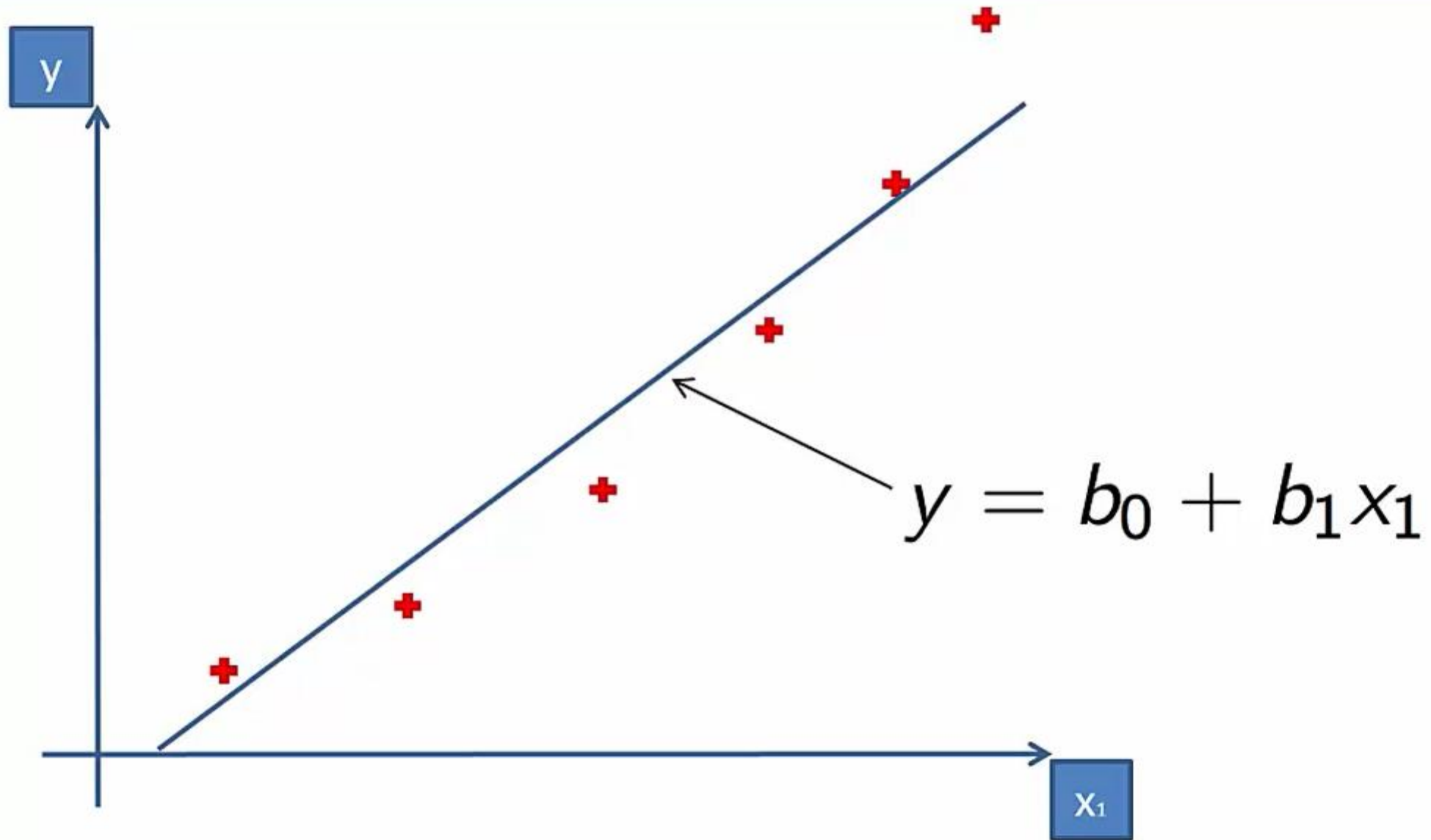
$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \dots + b_n x_1^n$$



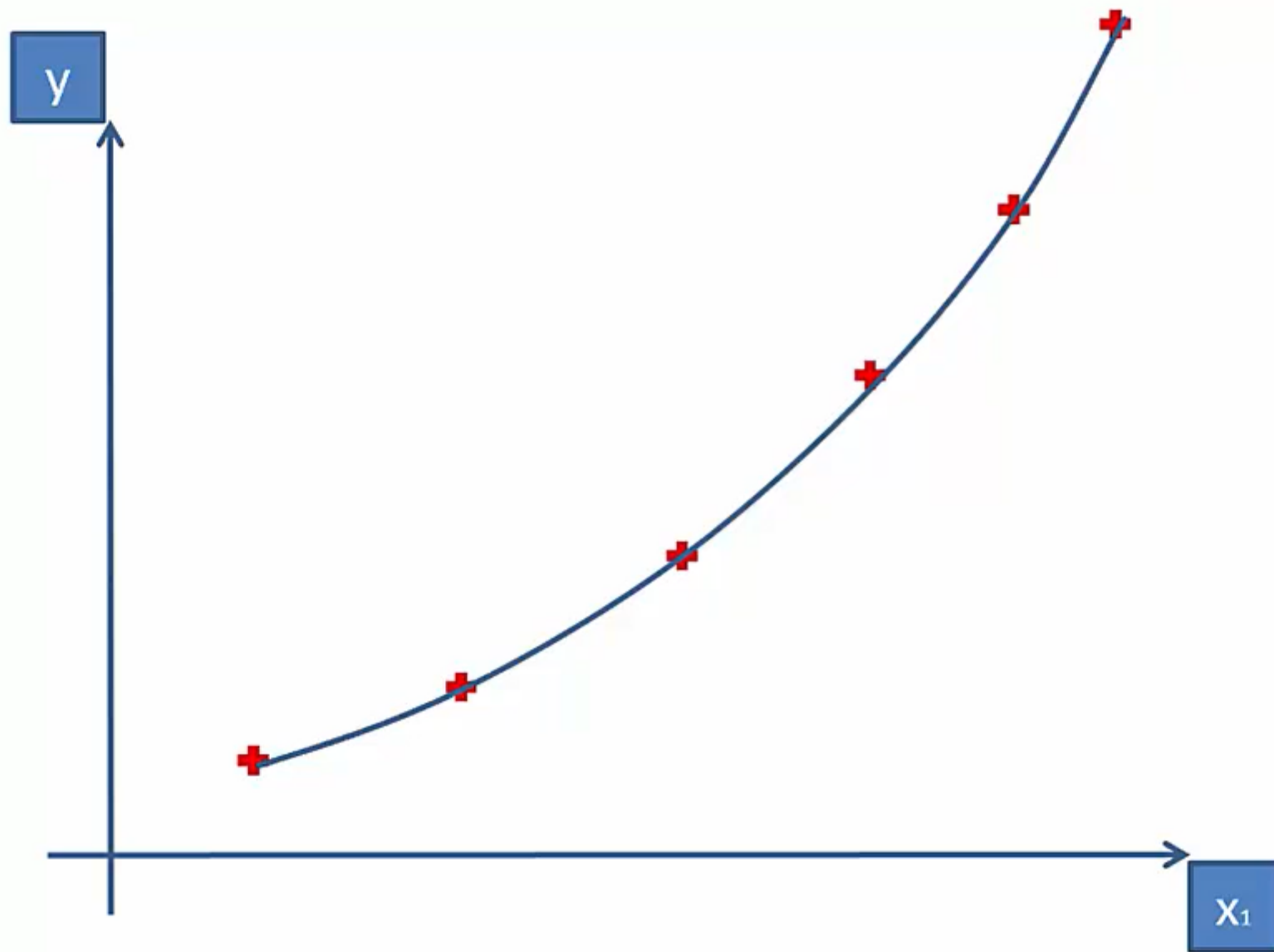
# Simple Linear Regression



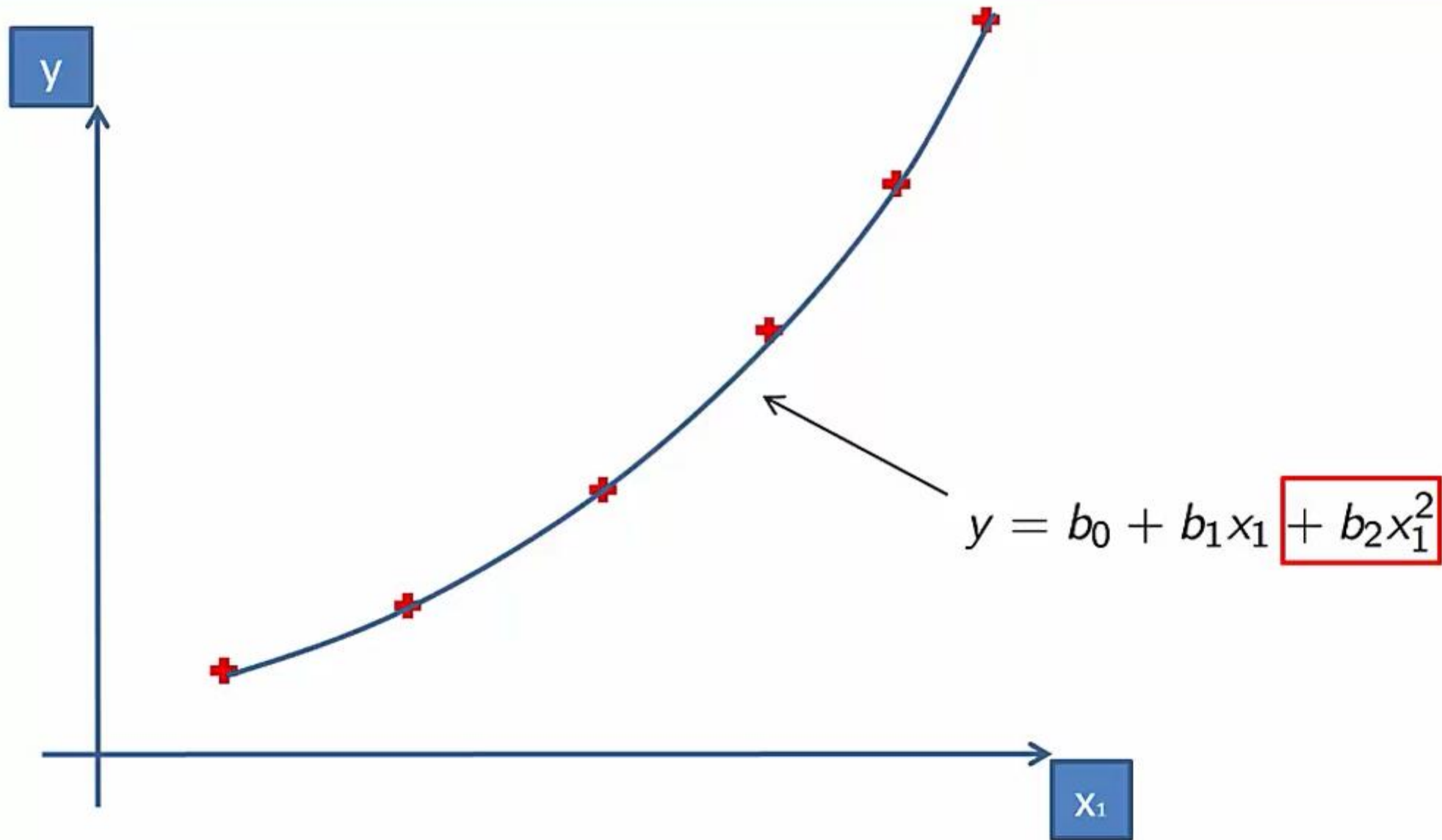
# Simple Linear Regression



# Polynomial Regression



# Polynomial Regression





# Polynomial Regression

Polynomial  
Linear  
Regression

$$y = b_0 + b_1x_1 + b_2x_1^2 + \dots + b_nx_1^n$$