

Microsoft Excel Part 4: Data Analysis

Prakarn Unachak
Department of Computer Science
Faculty of Science
Chiang Mai University

Outline

- Installing Data Analysis ToolPak
- Correlation
- Histogram
- Linear Regression

Department of Computer Science, Faculty of Science, Chiang Mai University

2

Installing Data Analysis ToolPak

1. Check whether if *Data Analysis* command already appears under *Data* tab.
If the command does not appear yet:
2. Click the office button.
3. Click *Excel Options*.
4. Select *Add-ins*.
5. Change Manage: option to *Excel Add-ins*.
6. Click *Go*.

Department of Computer Science, Faculty of Science, Chiang Mai University

3

Installing Data Analysis ToolPak

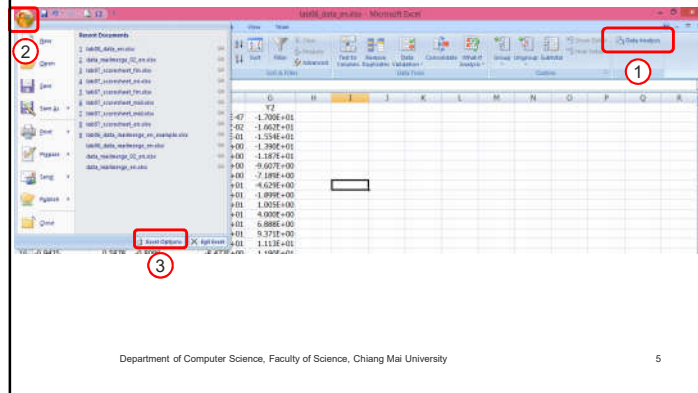
7. Under Add-Ins dialog box, make sure the Analysis ToolPak checkbox is checked.
8. Click OK (a few times).



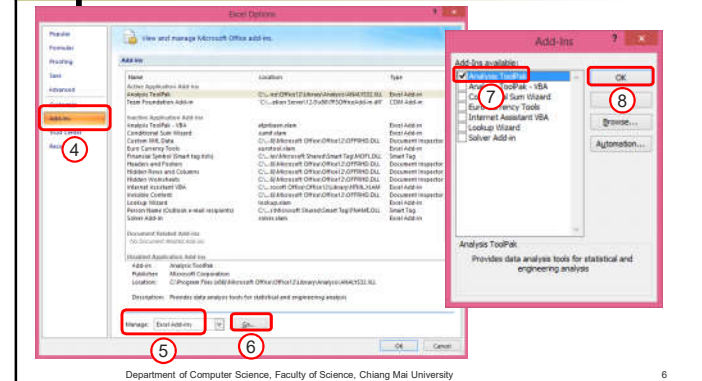
Department of Computer Science, Faculty of Science, Chiang Mai University

4

Installing Data Analysis ToolPak (cont.)



Installing Data Analysis ToolPak (cont.)



CORRELATION

Department of Computer Science, Faculty of Science, Chiang Mai University

7

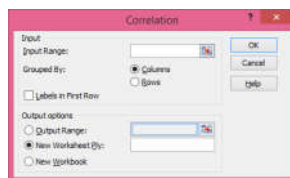
Correlation

- **Correlation** tells you how close two set of data related to each other. The value is between -1 and 1.
- If the two data sets have **positive correlation**, when one increase/decrease, the other will go to the **same** direction.
- If the two data sets have **negative correlation** when one increase/decrease, the other will go to the **other** direction.
- If the correlation is **zero**, there is no relation between two data sets. Whether the one decrease/increase will not predict the other.

Department of Computer Science, Faculty of Science, Chiang Mai University

8

Calculating Correlation



1. Select **Data** → **Data Analysis** → **Correlation**
2. Select **Input Range**
3. Select Options
 - Grouped By: Columns/Rows.
 - Check if first row of data is a label.
4. Select **Output Range**.

Department of Computer Science, Faculty of Science, Chiang Mai University

9

Correlation Results

	A	B	C	D
1	W	X	Y	Z
2	1	1	10	1
3	2	1.5	9	1.5
4	3	2	8	2
5	4	2.5	7	2.5
6	5	3	6	3
7	6	3.5	5	3
8	7	4	4	3.5
9	8	4.5	3	4
10	9	5	2	4.5
11	10	5.5	1	5



	E	F	G	H	I
		W	X	Y	Z
W		1			
X		1	1		
Y		-1	-1	1	
Z		0.994937	0.994937	-0.99494	1

Department of Computer Science, Faculty of Science, Chiang Mai University

10

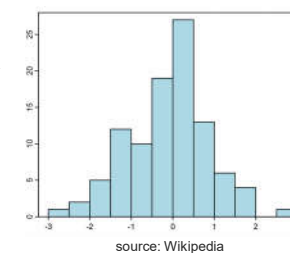
HISTOGRAM

Department of Computer Science, Faculty of Science, Chiang Mai University

11

Histogram

- A **histogram** is a way to display data by grouping each data point into ranges (**bins**) of their values
- Good for showing data distribution



Department of Computer Science, Faculty of Science, Chiang Mai University

12

[Histogram (cont.)]

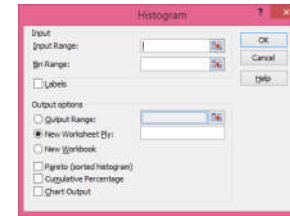
		Bins
	96	20
	35	40
	62	60
	21	80
	87	
	75	
	20	
	83	
	10	
	64	
	20	
	76	
	41	
	85	
	57	
	0	
	67	
	84	
	2	
	19	

- To perform Histogram on Excel, you need **data** and **bins**.
- Bins are cells whose values used to determine which group a data point will go to.

Department of Computer Science, Faculty of Science, Chiang Mai University

13

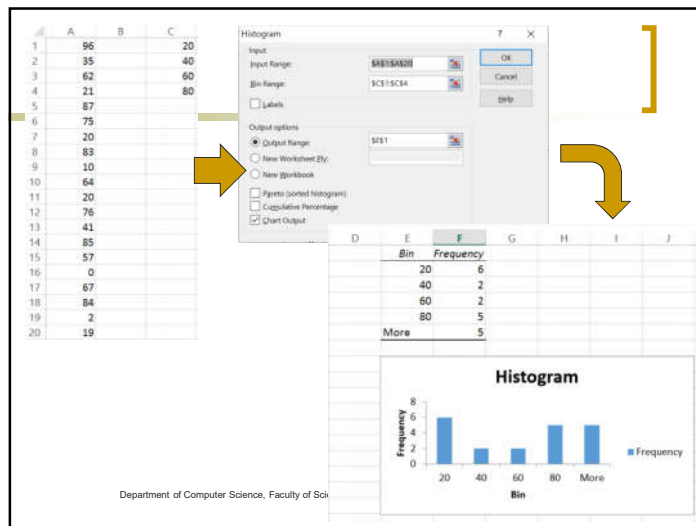
[Calculating Histogram]



1. Select **Data** → **Data Analysis** → **Histogram**
2. Select **Input Range**
3. Select **Bin Range**
4. Select **Output Range**
5. Check other options
 - **Chart Output** if you want charted version of histogram.
6. Click **OK**

Department of Computer Science, Faculty of Science, Chiang Mai University

14



Department of Computer Science, Faculty of Sci

LINEAR REGRESSION

Department of Computer Science, Faculty of Science, Chiang Mai University

16

Linear Regression

- **Linear regression** is an analysis technique used to derive a relationship between dependent variable and one or more independent (explanatory) variables.
- Can be use both to explain data, and predict values of data in the future.

Department of Computer Science, Faculty of Science, Chiang Mai University

17

Linear Regression (cont.)

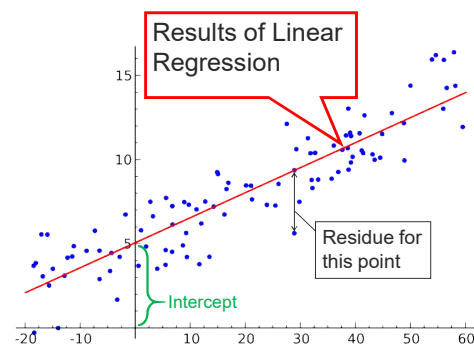
- Linear regression results will be in this form:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon$$
- Y is **dependent variable**.
- X_1, \dots, X_n are **independent** or **explanatory variables**.
- β_0 is an **intercept** or constant values of the equations.
- β_1, \dots, β_n are **coefficients** of each independent variables.
- ε is an **error term** or **residue**. **This is not a constant.**

Department of Computer Science, Faculty of Science, Chiang Mai University

18

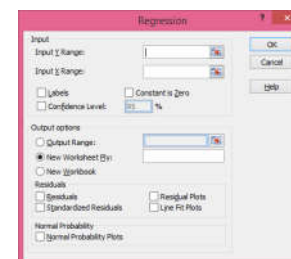
Linear Regression, in Graph



Department of Computer Science, Faculty of Science, Chiang Mai University

19

Linear Regression (cont.)



1. Select **Data** → **Data Analysis** → **Regression**
2. Select **Y** (dependent variable) **Range**
3. Select **X** (independent variable) **Range**
4. Select **Output Range**.
5. Select options.
6. Click **OK**

Department of Computer Science, Faculty of Science, Chiang Mai University

20

Reading the Results

	G	H	I	J	K	L	M	N	O
SUMMARY OUTPUT									
<i>Regression Statistics</i>									
Multiple R	0.99999831								
R Square	0.99999661								
Adjusted R Square	0.99999621								
Standard Error	0.02599426								
Observations	20								

- R-square, adjusted R-square measure accuracy of model as a whole. (The entire equation).
- Closer to 1 is better.

ANOVA					
	df	SS	MS	F	Significance F
Regression	2	3391.344686	1695.67234	2509499	3.1883E-47
Residual	17	0.011486927	0.0006757		
Total	19	3391.356173			
	Coefficients	Standard Error	t Stat	P-value	Lower 95%
Intercept	5.04734708	0.011855951	425.722653	1.1E-35	5.02235524
X1	1.97983667	0.009738943	203.290713	3.2E-30	1.9592893
X2	1.00667981	0.003559745	282.795514	1.2E-32	0.9991694

- t-Stat, P-value measure accuracy for that one variable.
- The higher the t-Stat, the better.
- The lower the P-value, the better

Coefficient for each variable + intercept

The regressed equation is $Y = 1.98 * X1 + 1.01 * X2 + 5.05$

6