

FOSE1025 — Scientific Computing

Week 8 Lecture 1: Summarising and Analysing Data

Diego Mollá

Department of Computer Science
Macquarie University

FOSE1025 2020H1

Programme

- 1 Pivot Tables
- 2 Data Analysis

Reading

- Lecture notes
- <https://www.linkedin.com/learning/excel-pivottables-for-beginners/>

The Scientific Method



Steps of the Scientific Meth...
sciencebuddies.org



Scientific method & variables
slideshare.net



Essays on Scientific Method: exa...
studymoose.com



The scientific met...
khanacademy.org



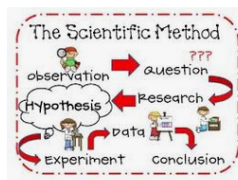
Why does the scientific ...



Formula for Using the Scient...



scientific process for ki...



The Scientific Method

Some results of a Google image search with the words "scientific" and "method" — 1 April 2020.

Excel to Manage Data in Science

We are covering these aspects in FOSE1025:

- Import data from external files (e.g. CSV) — Week 3.
- Explore the data — Week 4.
- Clean the data — Week 6.
- Preprocess, transform the data — Week 7.
- Analyse, summarise, interpret the data — Week 8.

Programme

- 1 Pivot Tables
- 2 Data Analysis

Pivot Tables: A Motivational Example

(data from <https://www.linkedin.com/learning/excel-pivottables-for-beginners>)

- Find the total shopping in each category “Fuel”, etc, of file shopping.csv.
- Find the total shopping of each month.
- What shopping per month and per category??
- Pivot tables can help you generate data for all of above and more.

Date	Buyer	Type	Amt
1-Jan	Mom	Fuel	\$50
2-Jan	Mom	Groceries	\$120
3-Jan	Dad	Cafes	\$10
4-Jan	Dad	Fuel	\$40
4-Jan	Kelly	Groceries	\$129
5-Jan	Mom	Cafes	\$12

A Simple Pivot Table

	F	G	H	I	J	K	L	M	N	O	P	Q	R				
Sum of Amt	Column Labels																
Row Labels	Books	Cafes	Entertainment	Fuel	Groceries	Music	Restaurants	Grand Total									
Jan	169	36		271	209	2147	15	2847									
Feb	476	59		142	202	2820	15	3714									
Mar	160	48		51	329	2348	46	2519	5501								
Apr	418	34		307	100	2985	9	3299	7152								
May	96	63		240	288	2911	14	2136	5748								
Jun	38	145		309	198	2905	86	3352	7033								
Jul	60	33		722	228	2834	6	3419	7302								
Aug	79	38		143	138	3120	17	3651	7186								
Sep	61			163	2377	9	3783	6393									
Oct	39			165	3063	13	3492	6772									
Nov	67			927	117	2373	10	1030	4524								
Dec	328			2627	55	2786	9	5805									
Grand Total	1991	456		5739	2192	32669	249	26681	69977								

PivotTable Fields

FIELD NAME

Search fields

☐ Date

☐ Buyer

☒ Type

Filters

Columns

: Type

Rows

: Months

Values

: Sum of Amt

Drag fields between areas

Anatomy of a Pivot Table

Filters

- What column to use to filter values.
- Only for columns with categorical data.

Rows

- What column to use in the rows of the pivot table.
- Only for columns with categorical data.

Columns

- What column to use in the columns of the pivot table.
- Only for columns with categorical data.

Values

- What value we want to aggregate.
- Only for columns with numerical data.

Pivot Tables to Convert from Long to Wide

Exercise 1 (weather_data.csv)

What is the average precipitation in Antigo?

- Using AVERAGEIFS
- Using a pivot table

Exercise 2 (weather_data.csv)

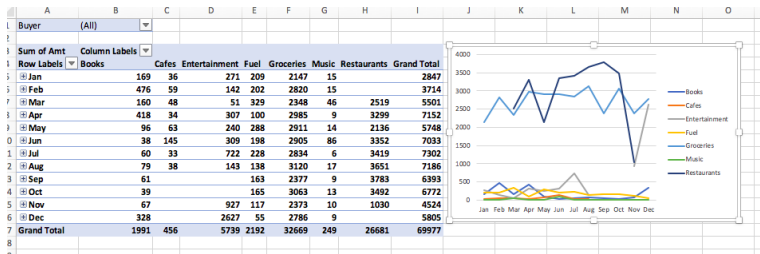
What is the March-2013 average precipitation in Antigo?

- Using AVERAGEIFS
- Using a pivot table

	A	B	C	D	E	F
1		data	date	param	siteid	
2	1	0	1/1/03	Precipitation	ACRE	
3	2	0	2/1/03	Precipitation	Albert Lea	
4	3	11.3199997	3/1/03	Precipitation	Ames	
5	4	0	4/1/03	Precipitation	Antigo	
6	5	3.03999996	5/1/03	Precipitation	Appleton	

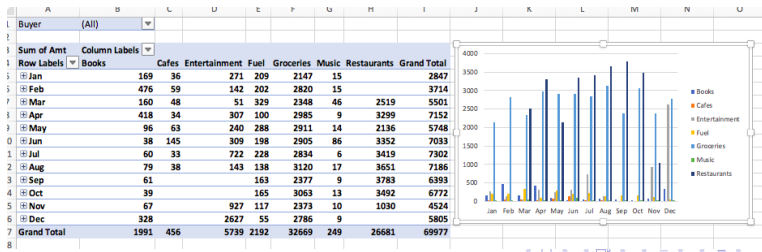
Pivot Tables for Charts

- Pivot tables facilitate the transformation of data for the creation of complex plots.
- In a **multiple chart**, each column of a table is plotted overlaid with the rest. Good for line plots.
- In a clustered chart, each row forms a cluster. Good for bar charts.
- In a stacked chart, columns of a table are plotted one on top of the other.



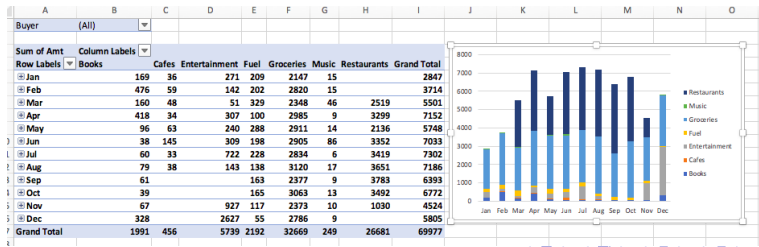
Pivot Tables for Charts

- Pivot tables facilitate the transformation of data for the creation of complex plots.
- In a multiple chart, each column of a table is plotted overlaid with the rest. Good for line plots.
- In a **clustered chart**, each row forms a cluster. Good for bar charts.
- In a stacked chart, columns of a table are plotted one on top of the other.



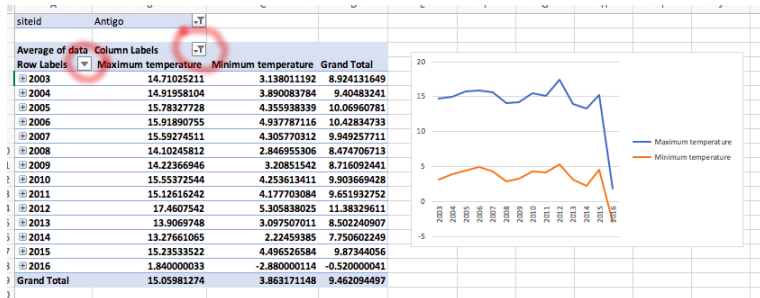
Pivot Tables for Charts

- Pivot tables facilitate the transformation of data for the creation of complex plots.
- In a multiple chart, each column of a table is plotted overlayed with the rest. Good for line plots.
- In a clustered chart, each row forms a cluster. Good for bar charts.
- In a **stacked chart**, columns of a table are plotted one on top of the other.



Pivot Charts: Pivot Tables and Charts!

- Pivot tables are so useful for making charts that there's a tool for that combines both: Pivot charts!
- Exercise: Can you plot (multiple line plot) the maximum and minimum temperature of Antigo as it changes over time? Do not plot precipitation.
 - (hint: you can filter row labels and **column** labels.)



Programme

1 Pivot Tables

2 Data Analysis

- Finding Trends
- Finding Correlations

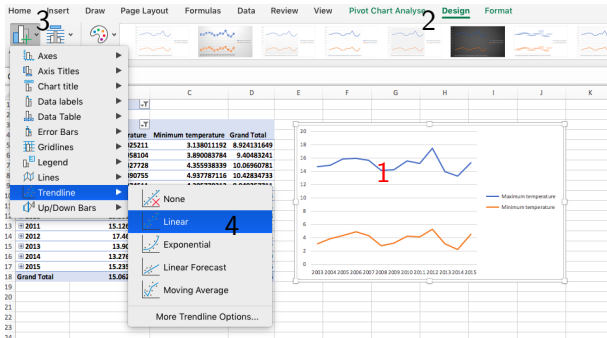
Analysing the Data

- Excel provides various tools for data analysis.
- Understanding most of these tools is beyond the scope of this unit.
- Here we will focus on two goals:
 - Finding trends.
 - Finding correlations.

Adding a Trend Line

- Excel charts support the inclusion of a trend line.
- Select **chart** → Design → Add Chart Element → Trendline.
- Choose the kind of trendline based on what you want to show.

(this figure is based on MS Excel for Mac, Version 16.30, Office 365)

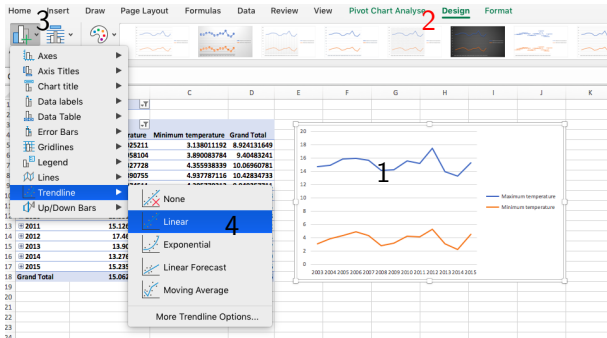


Subscription)

Adding a Trend Line

- Excel charts support the inclusion of a trend line.
- Select chart → **Design** → Add Chart Element → Trendline.
- Choose the kind of trendline based on what you want to show.

(this figure is based on MS Excel for Mac, Version 16.30, Office 365)

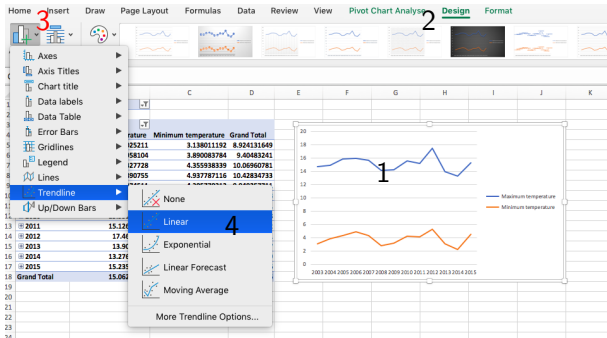


Subscription)

Adding a Trend Line

- Excel charts support the inclusion of a trend line.
- Select chart → Design → **Add Chart Element** → Trendline.
- Choose the kind of trendline based on what you want to show.

(this figure is based on MS Excel for Mac, Version 16.30, Office 365)

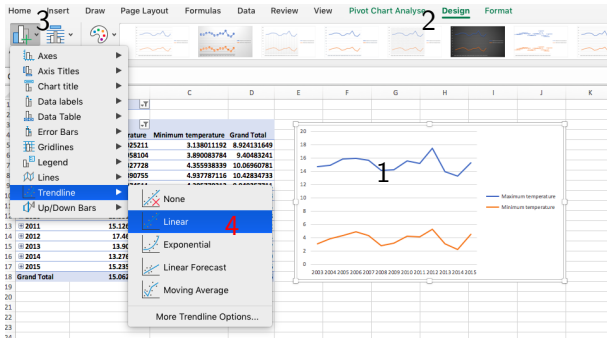


Subscription)

Adding a Trend Line

- Excel charts support the inclusion of a trend line.
- Select chart → Design → Add Chart Element → **Trendline**.
- Choose the kind of trendline based on what you want to show.

(this figure is based on MS Excel for Mac, Version 16.30, Office 365)



Subscription)

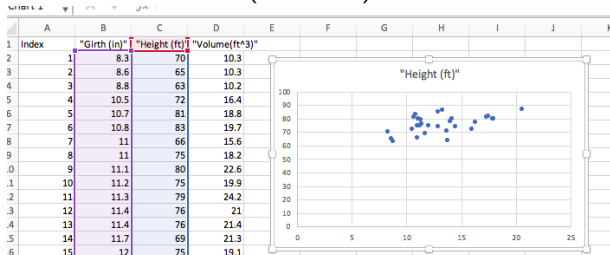
What is Correlation?

- Sometimes two variables are measuring the same property.
 - (each column represents one variable)
 - May happen when multiple agents are providing data.
- You may detect this by observing that the values are the same.
- But sometimes there are minor variations.
- In other cases, two variables are correlated but might not be identical.
 - For example, tree trunk height and girth are correlated.
 - Taller trees will normally have thicker trunks.

Finding Correlations Graphically

- A **scatterplot** can plot one variable against the other.
- If the two variables are not correlated, the scatterplots will look random.
- If the scatterplot has a distinct shape, the two variables are correlated.
- For example, if the shape looks like a line, then the two variables have a **linear correlation**.

(trees.csv)



Finding Correlations on Multiple Columns

- Scatterplots are intuitive but may be cumbersome if you want to check the correlations among many columns.
- E.g. if there are 10 columns you will need to make a plot for each possible pair.
 - This means making $10 \times 9 = 90$ plots.
- There are various formulas that attempt to express the correlation as a number.
- Excel's CORREL function uses one of those formulas.
 - e.g. `=CORREL(B:B,C:C)` computes the correlation between columns B and C.
 - If you want to know what formula Excel uses, look for the “sample correlation coefficient”.
- A number close to 1 (or -1) indicates positive (or negative) correlation.

Correlation Matrix

- Excel's "Data Analysis" tool can compute a correlation matrix.
- Data → Data Analysis → Correlation.

(trees.csv)

The screenshot shows an Excel spreadsheet with the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Index	"Girth (in)"	"Height (ft)"	"Volume(ft³)"									
2	1	8.3	70	10.3									
3	2	8.6	65	10.3									
4	3	8.8	63	10.2									
5	4	10.5	72	16.4									
6	5	10.7	81	18.8									
7	6	10.8	83	19.7									
8	7	11	66	15.6									
9	8	11	75	18.2									
10	9	11.1	80	22.6									
11	10	11.2	75	19.9									
12	11	11.3	79	24.2									
13	12	11.4	76	21									
14	13	11.4	76	21.4									
15	14	11.7	69	21.3									
16	15	12	75	19.1									
17	16	12.9	74	22.2									
18	17	12.9	85	33.8									
19	18	13.3	86	27.4									
20	19	13.3	71	16.7									

The 'Correlation' dialog box is open, showing the following settings:

- Input Range: \$B:\$D
- Grouped By: Columns
- Labels in first row: ☒
- Output Range: \$L\$7
- New Worksheet Ply: ☐
- New Workbook: ☐

(you will observe a strong correlation between girth and volume)

Exercise

- File: shopping.png
- Build the correlation matrix between all types of shopping.
- What are the two most correlated types of shopping?
- Show it clearly by introducing **conditional formatting** that highlights the highest correlations.
 - Home → Conditional Formatting → Colour Scales

	Books	Cafes	Entertainment	Fuel	Groceries	Music	Restaurants
Books	1						
Cafes	-0.289396228	1					
Entertainment	0.160093641	-0.08	1				
Fuel	-0.271487084	0.09	-0.625410842	1			
Groceries	0.09428483	0.19	-0.000504711	-0.2	1		
Music	-0.243270987	0.88	-0.285322756	0.34	0.026135	1	
Restaurants	0.030060483	-0	-0.470731464	-0.1	0.470645	0.071	1

Take-home Messages

- Pivot tables are very powerful to process tables in long format.
- You must be able to use pivot tables for a range of tasks.
- You must be able to create charts based on pivot tables.
- You must be able to show trends by adding trend lines to a plot.
- You must be able to detect whether two variables are correlated.

What's Next

- Week 9 lecture: Ethics related to Scientific Computing.
- Week 9: Submit the project.