

# FOSE1025 — Scientific Computing

## Week 8 Lecture 1: Summarising and Analysing Data

Diego Mollá

FOSE1025 2020H1

### Abstract

This lecture will focus on several approaches for summarising and preparing the data for the final analysis. We will look at pivot tables as a powerful tool to transform and summarising the data. With pivot tables we can convert tables from the long to the wide format. In addition, we can aggregate and filter data and make it ready for insightful analysis and graphic representations. Beside pivot tables, we will look at some specific tools that Excel provides for the analysis of data.

Update April 27, 2020

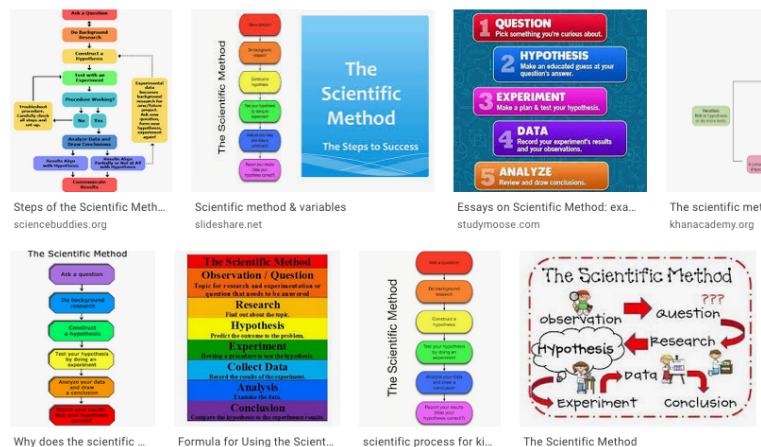
## Contents

<b>1</b>	<b>Pivot Tables</b>	<b>2</b>
<b>2</b>	<b>Data Analysis</b>	<b>5</b>
2.1	Finding Trends . . . . .	5
2.2	Finding Correlations . . . . .	6

## Reading

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

## 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.

# 1 Pivot Tables

## 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
6-Jan	Kelly	Cafes	\$14
7-Jan	Kelly	Books	\$129
7-Jan	Dad	Groceries	\$252
9-Jan	Kelly	Fuel	\$44
10-Jan	Dad	Groceries	\$39
12-Jan	Mom	Books	\$20
13-Jan	Dad	Groceries	\$132
14-Jan	Dad	Groceries	\$79
16-Jan	Kelly	Groceries	\$172
16-Jan	Dad	Music	\$8
18-Jan	Kelly	Fuel	\$30

## A Simple Pivot Table

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	3501											
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	2634	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	2796	9		5805											
Grand Total	1991	456	5739	2192	32669	249	26681	69977											

## 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?

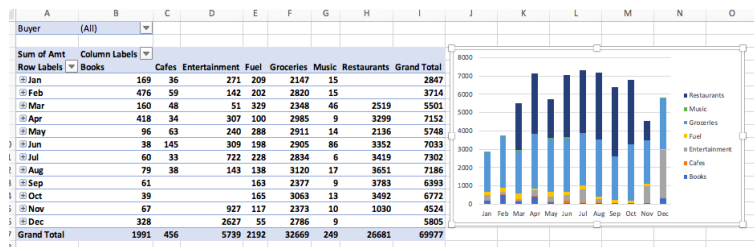
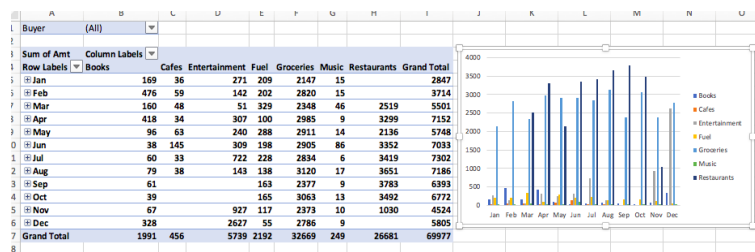
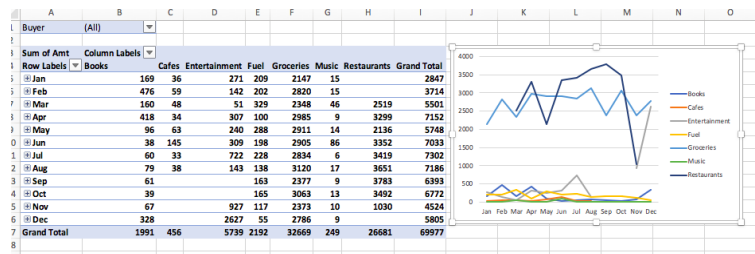
- Using AVERAGEIFS

- 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	AlbertLea	
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	
7	6	0.49000001	6/1/03	Precipitation	Arlington	
8	7	0	7/1/03	Precipitation	Bean&Beet	
9	8	0	8/1/03	Precipitation	Brookings	
10	9	0	9/1/03	Precipitation	Brownstown	
11	10	0	10/1/03	Precipitation	Columbia	
12	11	0	11/1/03	Precipitation	Crookston	
13	12	0	12/1/03	Precipitation	Dekalb	
14	13	0	13/1/03	Precipitation	DixonSprings	

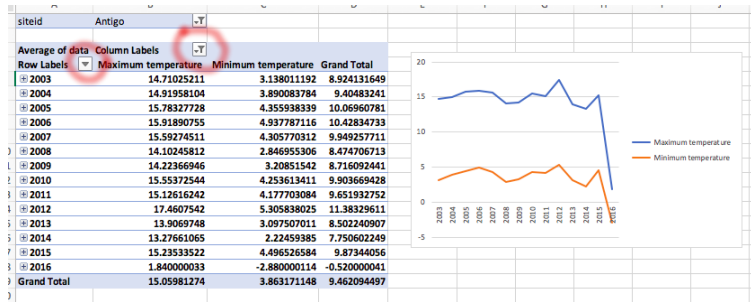
## Pivot Tables for Charts

- 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.)



## 2 Data Analysis

### Analysing the Data

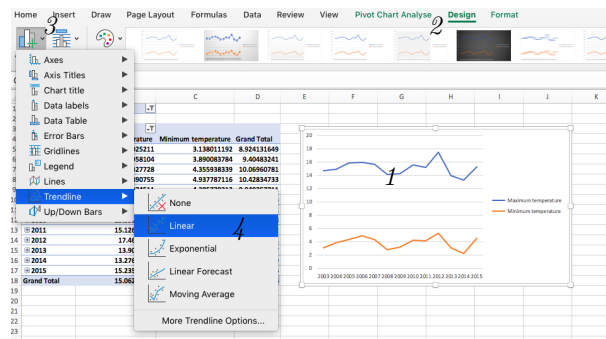
- Excel provides various tools for data analysis.
- Understanding most of these tools is beyond the scope of this unit.
- Here we will focus on two goals:
  - Finding trends.
  - Finding correlations.

### 2.1 Finding Trends

#### 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)



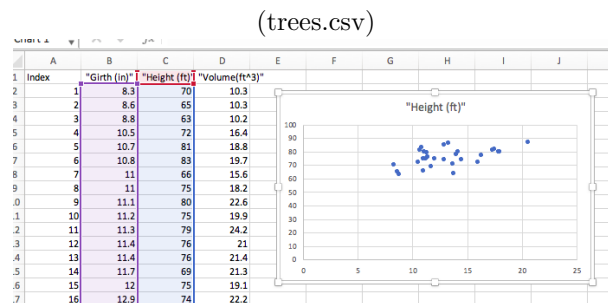
## 2.2 Finding Correlations

### 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*.

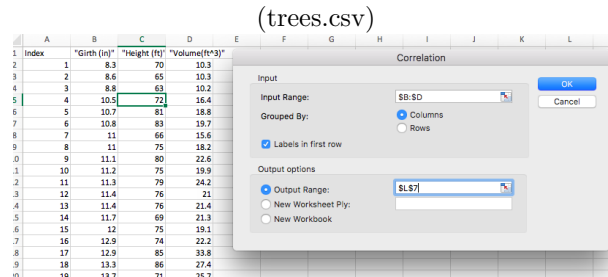


### Finding Correlations on Multiple Columns

- 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.



(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.