

---

# MITS6005

## Big Data

---

***Copyright © 2015 - 2019, Victorian Institute of Technology.***

*The contents contained in this document may not be reproduced in any form or by any means, without the written permission of VIT, other than for the purpose for which it has been supplied. VIT and its logo are trademarks of Victorian Institute of Technology.*

---

# Session 3

## Analyze & Present Data

---

***Copyright © 2015 - 2019, Victorian Institute of Technology.***

*The contents contained in this document may not be reproduced in any form or by any means, without the written permission of VIT, other than for the purpose for which it has been supplied. VIT and its logo are trademarks of Victorian Institute of Technology.*

# *Analyze data*

# Pivot tables

- A pivot table is an aggregation of a source table
- Supports summary statistics:
  - Count
  - Sum
  - Min/Max
  - Average
- This summary is presented in a table format which can be formatted and filtered
- A table with categories down the rows and across the columns is a **cross table**

	A	B	C	D	E	F	G	H	I
1	Rec	Salespers	Code	Item	Region	Month	Year	Units	Sales
2		1 Fred	A	Diary	East	Jan	1991	12	3
3		2 Fred	D	Dictionary	East	Jan	1991	34	40
4		3 Fred	E	Encyclope	East	Jan	1991	52	416
5		4 Fred	N	Novel	East	Jan	1991	34	17
6		5 Fred	A	Diary	East	Feb	1991	6	1
7		6 Fred	D	Dictionary	East	Feb	1991	5	6
8		7 Fred	E	Encyclope	East	Feb	1991	55	440
9		8 Fred	N	Novel	East	Feb	1991	33	16
10		9 Fred	A	Diary	East	Mar	1991	65	19
11		10 Fred	D	Dictionary	East	Mar	1991	34	40
12		11 Fred	E	Encyclope	East	Mar	1991	87	696
13		12 Fred	N	Novel	East	Mar	1991	23	11
14		13 Bert	A	Diary	South	Jan	1991	98	29
15		14 Bert	D	Dictionary	South	Jan	1991	55	66
16		15 Bert	E	Encyclope	South	Jan	1991	21	168
17		16 Bert	N	Novel	South	Jan	1991	11	6

Sum of Sales	Month			
Salesperson	Jan	Feb	Mar	Grand Total
Bert	7508	10360	6281	24149
Bill	5113	8916	7642	21671
Fred	7561	9735	11221	28517
Harry	10513	3583	9452	23548
Grand Total	30695	32594	34596	97885

# Customizing pivot tables

## Filters

Change total sample

## Rows

Can stack multiple fields, e.g. site broken down by team

## Columns

## Options

Details of the way fields are displayed in the table are changed here

	A	B	C	D
1	Status	(All)		
2				
3	Sum of Salary		Cost Centre	
4	Site	Team	CS21	CS22
5	Dundee	ABC Bank		
6		Hi Fi Motors		
7		Local Software		
8		Management		
9	Dundee Total			
10	Newcastle	Broken Cars Insurance		37225
11		Management		
12		The Double Glazing Co		533100
13	Newcastle Total			570325
14	Reading	Broken Pipes	26000	133650
15		Broken Pipes / XYZ Insurance		
16		Local Motors		39000
17		Management	35000	
18		XYZ Insurance	111000	1067250
19	Reading Total		172000	1239900
20	Grand Total		172000	1810225

PivotTable Field List

Choose fields to add to report:

- ☐ Staff Number
- ☒ Cost Centre
- ☒ Site
- ☒ Team
- ☐ Grade
- ☒ Status
- ☐ FTE
- ☐ Hire date
- ☒ Salary

Drag fields between areas below:

Report Filter: Status

Column Labels: Cost Centre

Row Labels: Site, Team

Values: Sum of Salary

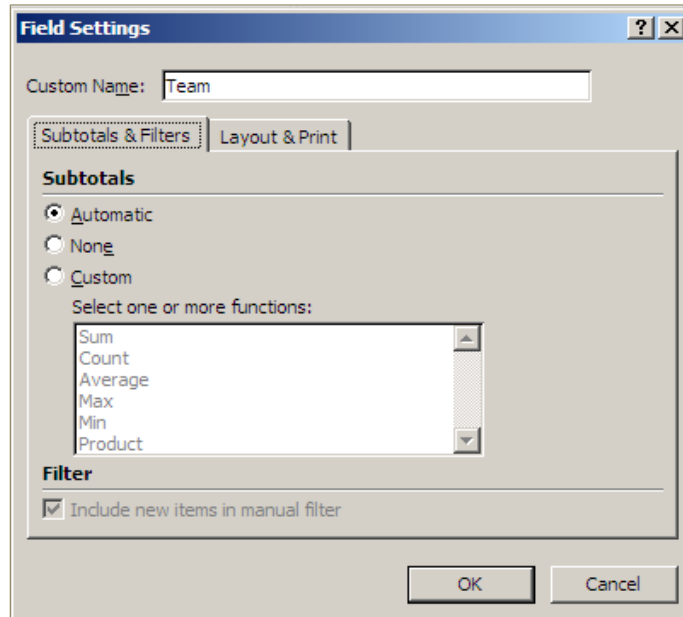
☐ Defer Layout Update

Update

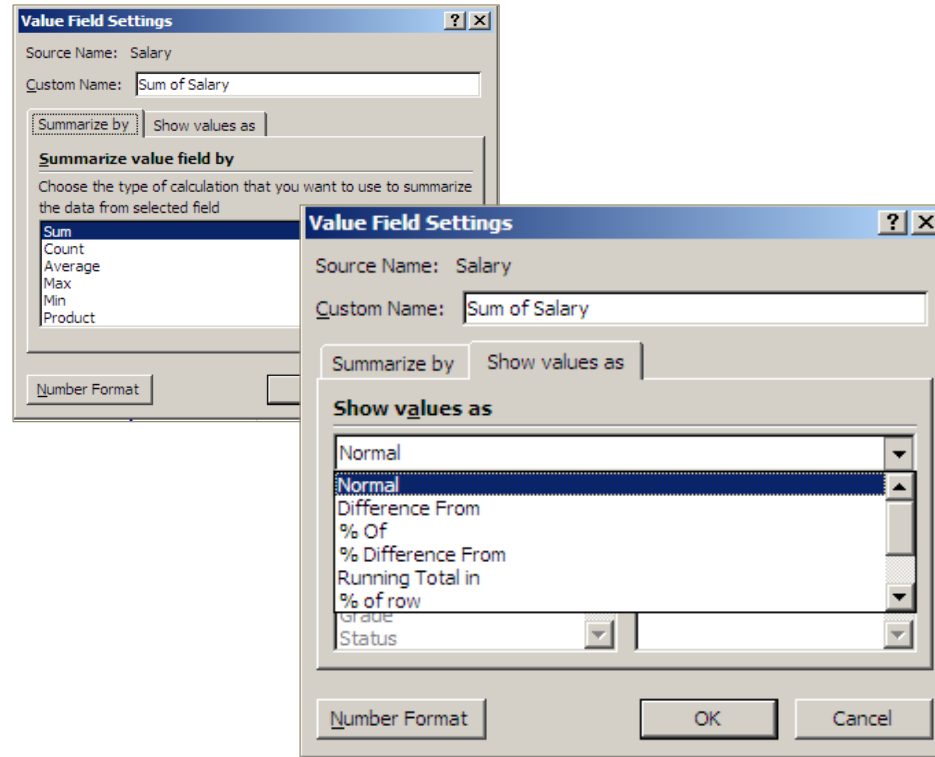
**Cells:** A variable summarised by rows/columns/filters

# Customizing pivot tables

**Rows and Columns** can have subtotals added to them, names changed, and additional options applied to the layout.



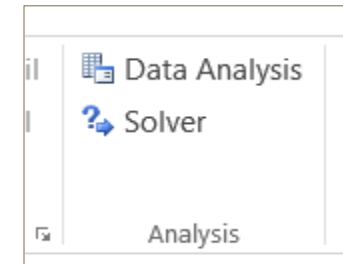
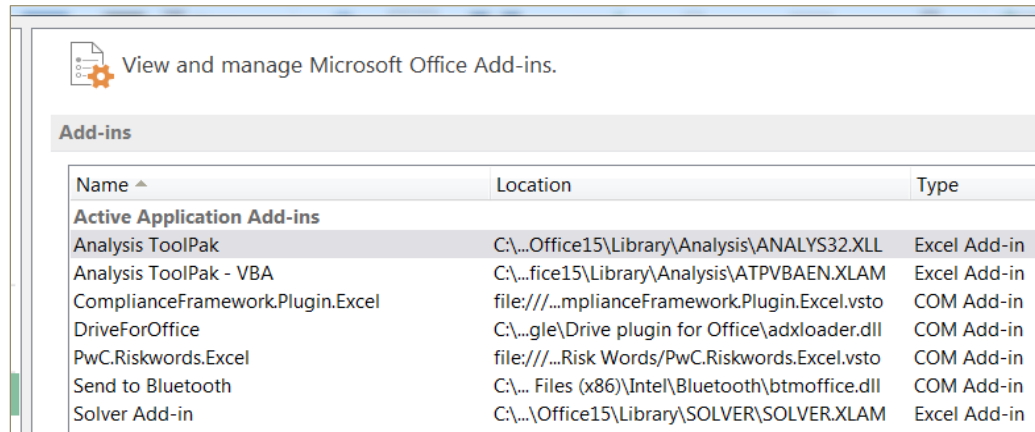
**Value Fields/Cells** can be summarised by different functions and shown as various values



# Analysis ToolPak

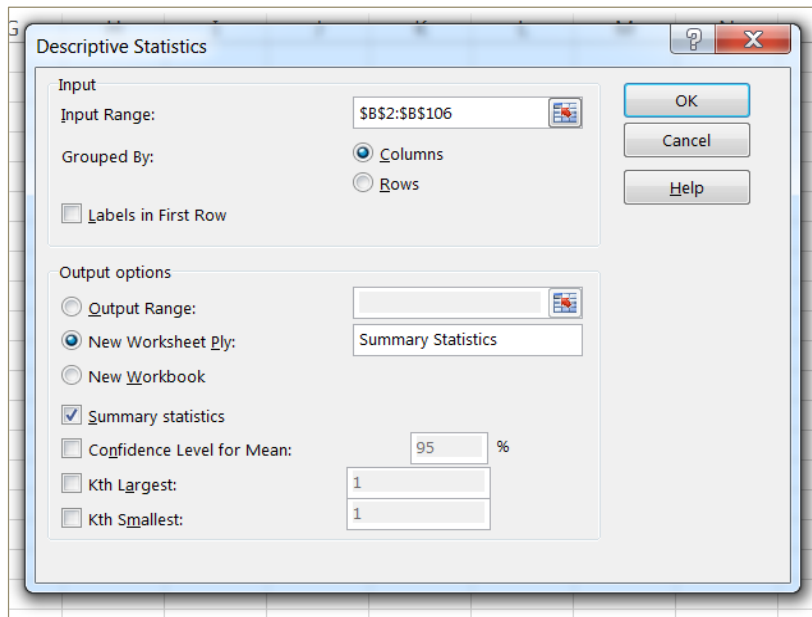
Excel can be used for more advanced analysis than summary statistics and charts

1. Click the File tab, click Options, and then click the Add-Ins category
2. In the Manage box, select Excel Add-ins and then click 'Go'
3. In the Add-Ins box, check the 'Analysis ToolPak' and 'Solver' check boxes, and then click 'OK'
4. Go to the 'Data' tab and look for the Analysis section



# Descriptive statistics

- Instead of using individual functions, the 'Data analysis' tab enables the user to calculate various statistics for a column at the same time



<i>Number of Participants</i>	
Mean	39.69524
Standard Error	5.648011
Median	16
Mode	2
Standard Deviation	57.87489
Sample Variance	3349.502
Kurtosis	6.352246
Skewness	2.515343
Range	271
Minimum	1
Maximum	272
Sum	4168
Count	105
Largest(2)	252
Smallest(2)	1
Confidence Level(95.0%)	11.20022

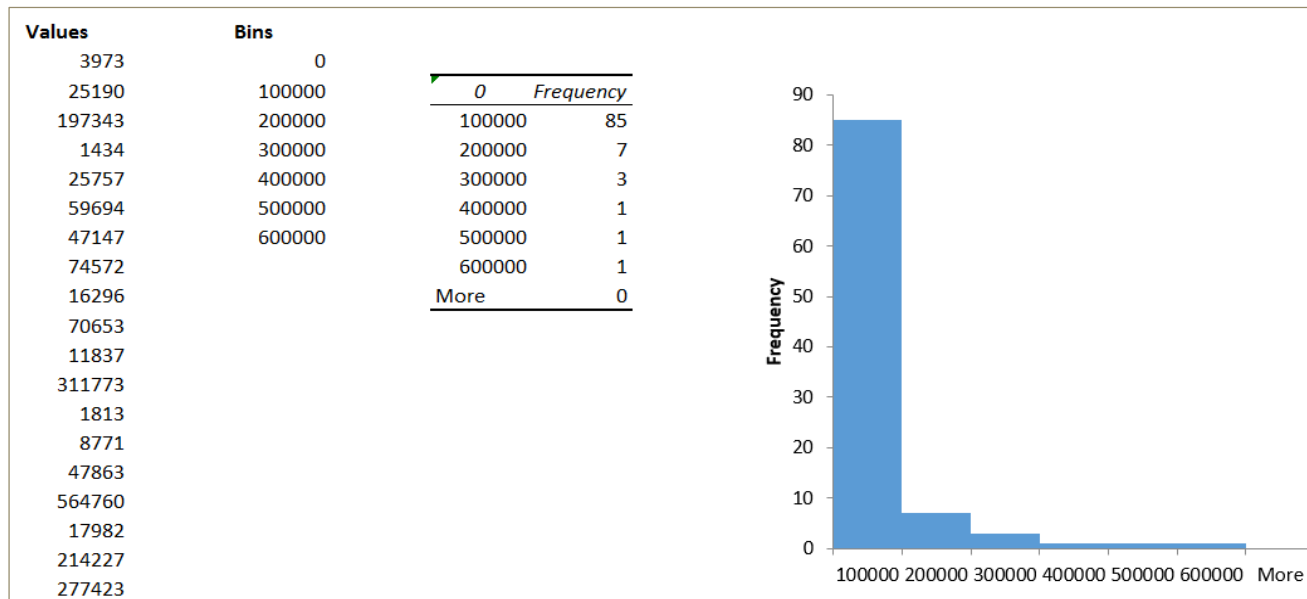


# Histograms

A histogram depicts the frequency or probability distribution of a numeric variable across bins of equal width

Variable distributions may appear statistically “normal” or display non-normal characteristics such as skewness or kurtosis

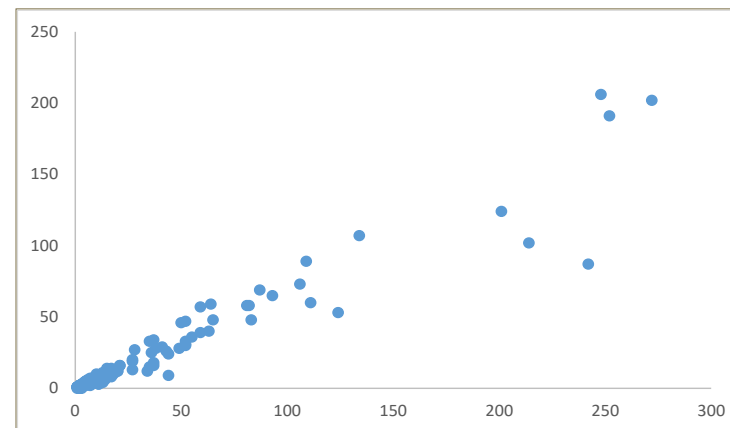
Create a histogram through ‘Data analysis’ by specifying the bins in a separate range



# Correlation

- The CORREL() function (and Data Analysis correlation option) quantify the direction and degree of the relationship between two variables
- Correlation coefficient is between -1 and 1:
  - **Negative values** indicate the variables are inversely related
  - **Positive values** indicate the variables are directly related
  - **Values close to zero** indicate a weak correlation
  - **Values close to 1** indicate a strong correlation
- Correlation is often best visualized with a scatter plot

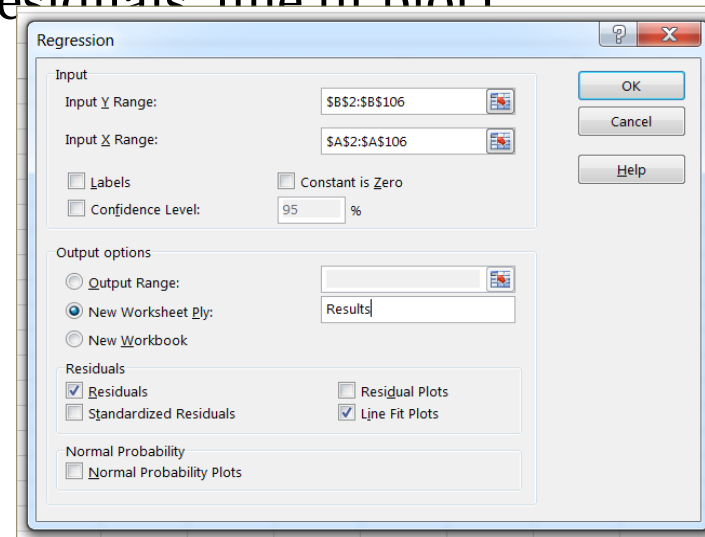
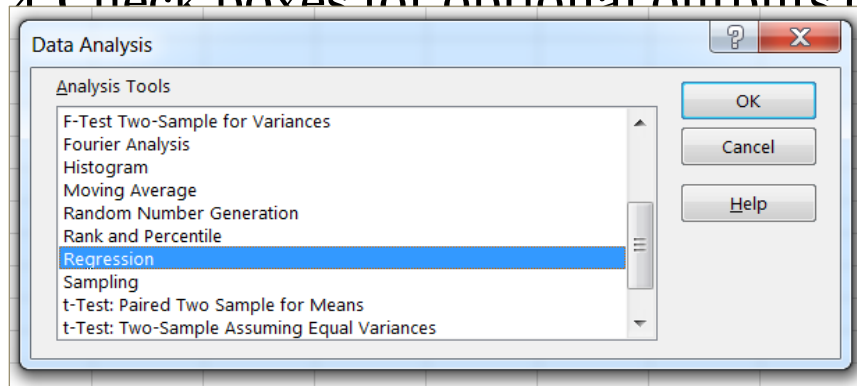
=CORREL(C2:C106,D2:D106)						
B	C	D	E	F	G	
	Column 1	Column 2			Correlation	0.956211
	44	24				
	35	33				
	1	1				
	2	2				
	6	6				
	1	0				
	2	2				
	3	3				
	8	7				
	17	8				



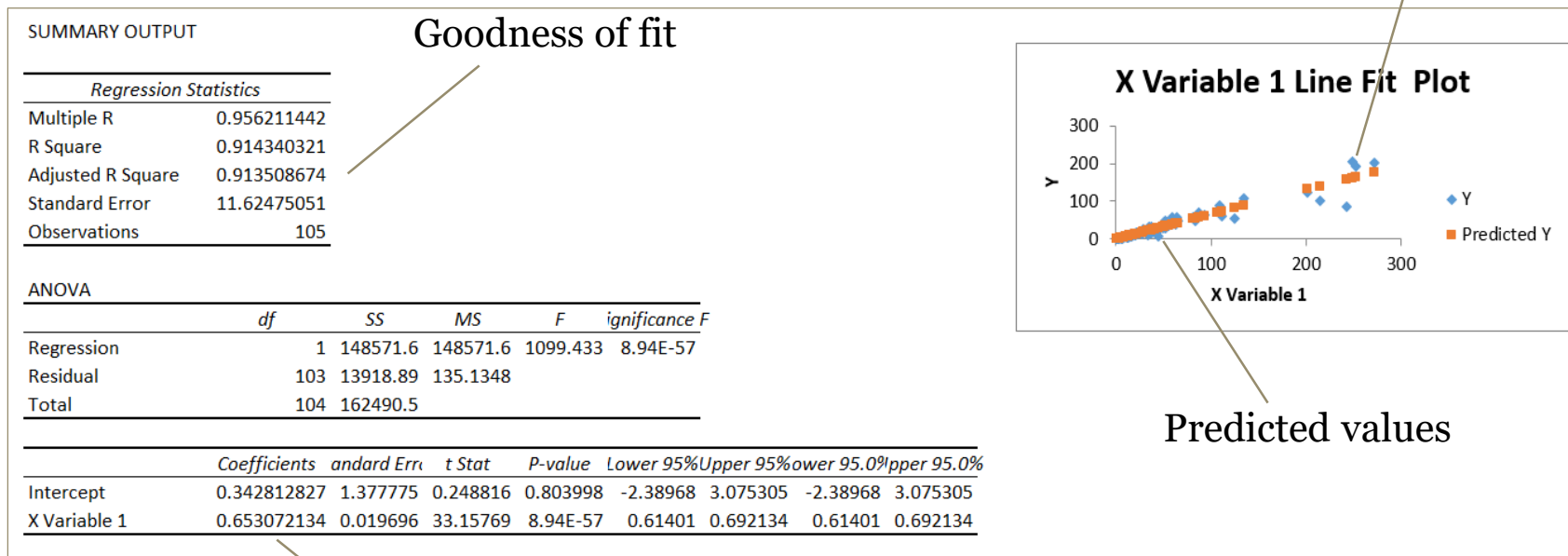
# *Linear regression (univariate)*

In addition to descriptive statistics, Excel can build and evaluate simple predictive models:

1. Select 'Regression' from the Data Analysis section
2. Select your dependent (Y) variable (to be predicted)
3. Select your independent (X) variable (used to predict the Y variable)
4. Check boxes for optional outputs (e.g. residuals, line fit plot)



# Linear regression output



---

## ***Exercise #3***

In 'Payroll Data.xlsx':

1. What is the average percentage of total compensation that is overtime pay, not including staff with no overtime pay?
2. Create a pivot table with the following statistics by Operating Unit:
  - Number of employees
  - Average total compensation
  - Average tenure (in years)
3. Add Job Title to the pivot table across the columns
4. What is the correlation between total compensation and tenure?
5. Create a histogram for total compensation... how would you characterize the distribution? What if you transform the data to reduce any skewness?
6. Create a linear regression model using tenure to predict total compensation... what is the adjusted R-squared value?

# *Present findings*

---

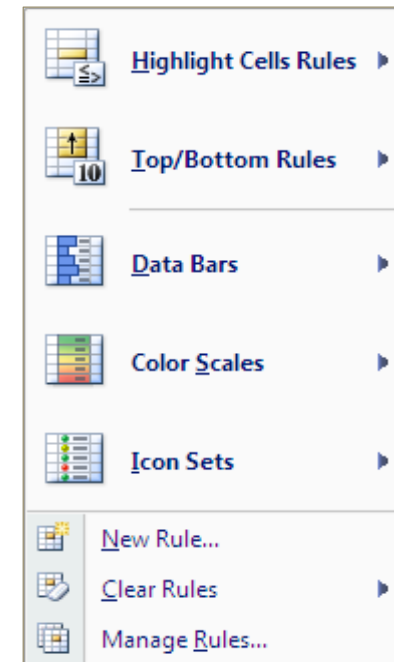
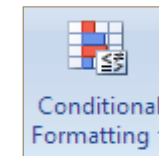
## *Visualisation in excel*

- Visualization is often essential for gaining an understanding of the data and presenting findings to a new audience
- Excel can easily produce a variety of basic charts
- For an effective visualization, always consider:
  - What question do I want to answer?
  - What message do I want my audience to take away?
  - How can I keep it simple?
- We'll focus more on visualization later in the course

# Conditional formatting

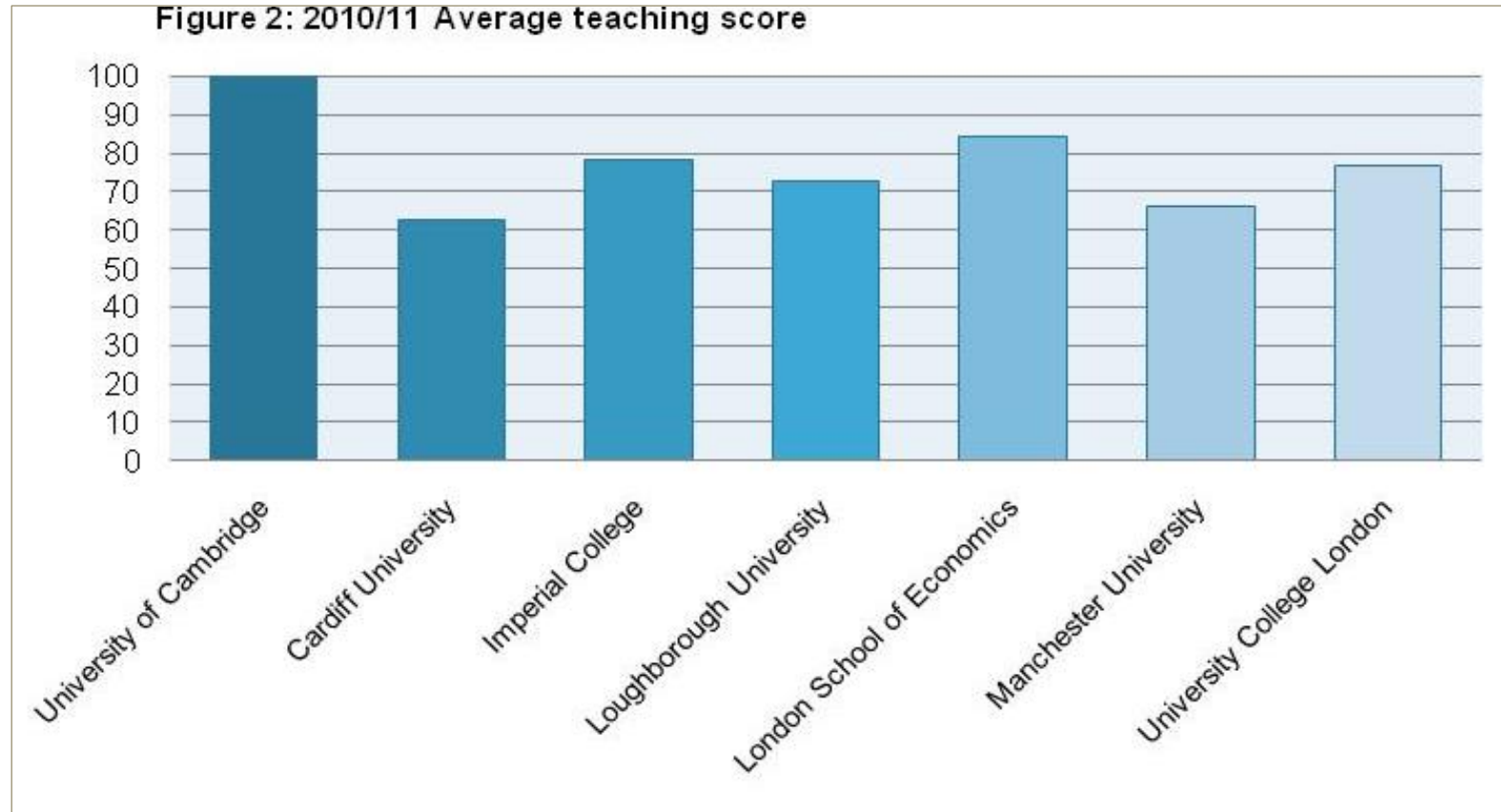
- Conditional formatting adjusts the color of a cell according to the relative magnitude of the values or established rules
- It is commonly used for RAG (Red/Amber/Green) reports, in

Average of Salary	Column Labels				
Row Labels	CT	NJ	NT	NY	WA
Communication & Strategic Sales	\$ 138,128.00	\$ 223,768.50		\$ 27,605.00	
Finance		\$ 94,055.25			
Human Resources		\$ 80,350.00			
IT Services	\$ 168,449.00	\$ 81,555.33	\$ 76,752.00	\$ 75,240.60	
Luxury Treats	\$ 81,241.93	\$ 67,250.20		\$ 152,621.25	\$ 88,688.14
Marketing		\$ 62,922.05		\$ 137,842.00	
Natural Health Foods	\$ 119,963.13	\$ 84,642.50		\$ 165,881.67	
Plant Support Services		\$ 83,708.20		\$ 130,655.32	
Tasty Biscuits	\$ 80,390.89	\$ 58,220.67		\$ 115,482.36	\$ 35,814.25
Tasty Bizkits	\$ 94,501.00			\$ 72,242.00	
Tasty Snacks Food	\$ 105,900.00	\$ 184,037.00		\$ 64,903.33	\$ 93,503.00
Tasty Soups	\$ 112,294.00	\$ 45,017.67		\$ 52,969.63	\$ 120,435.50



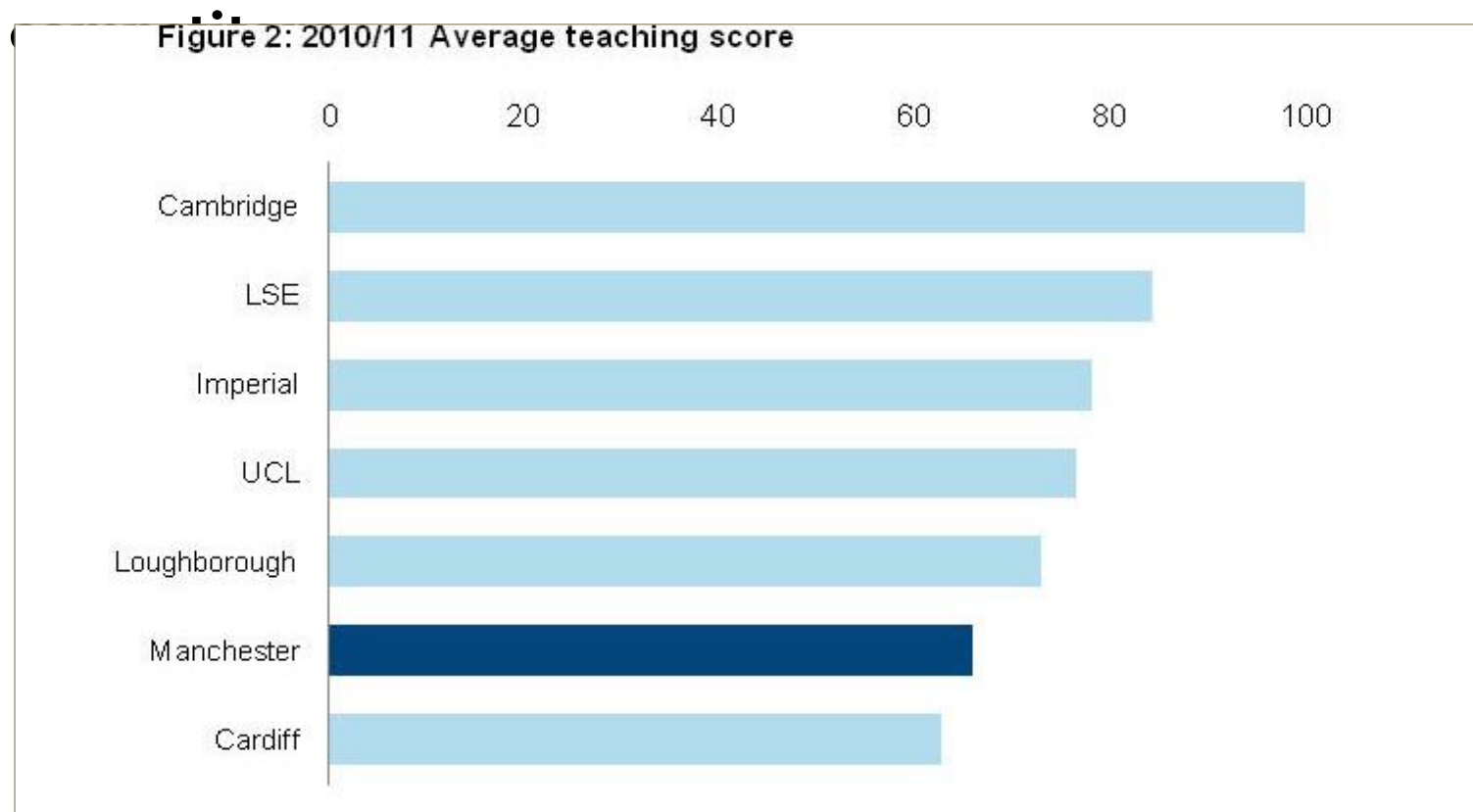


## *What do you think of this visualization?*



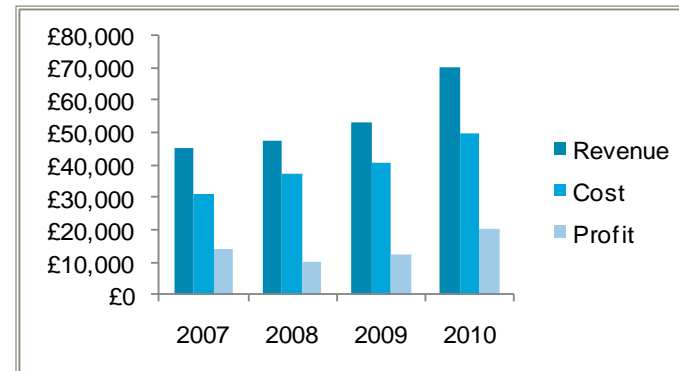
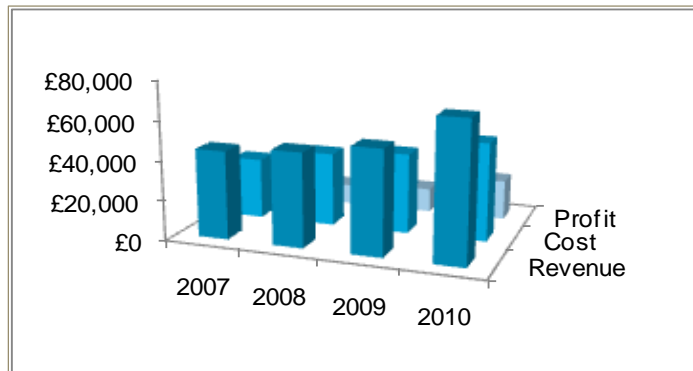
## ***What about this one?***

**Manchester University has low teaching score compared to its**



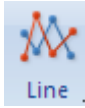
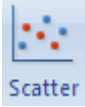


## *Effective visualisation*

- More information, less ink
- Clear out the junk!
- Use color, shape, placement, etc. to draw attention
- Use 2D rather than 3D charts



# Basic charts

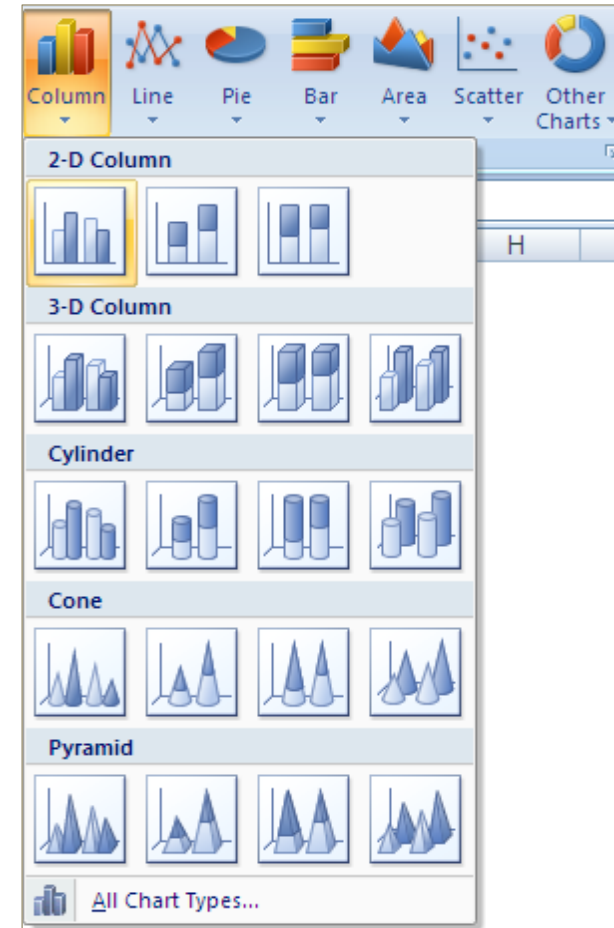
Type of chart		When to use them	Examples
Bar/Column	 	Discrete values	Number of employees at different offices
Line		Continuous values, over time	Annual revenue
Scatter		Continuous values, two variables	Quantity and price for different products

## *Insert chart*

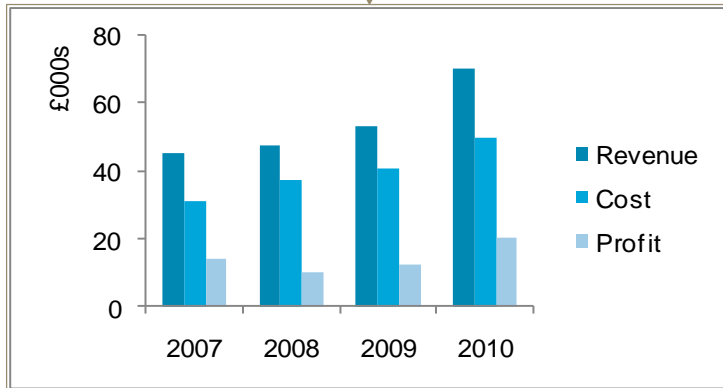
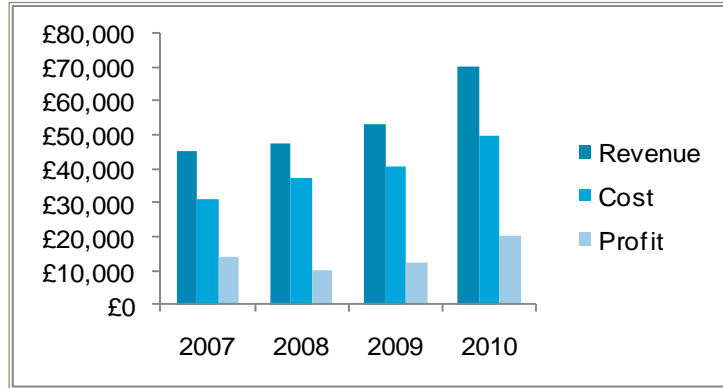
- The first step to setting up your chart is to arrange the data into a table
- Then select the table and select

	A	B	C	D	E
1					
2					
3					
4					
5					
6					
7					
8					
9					

	Revenue	Cost	Profit
2006	£ 44,300	£ 28,500	£ 15,700
2007	£ 46,000	£ 31,100	£ 13,900
2008	£ 47,000	£ 37,600	£ 10,200
2009	£ 53,000	£ 40,800	£ 12,200
2010	£ 70,200	£ 49,800	£ 20,400



# Axis settings



Axis settings menu:

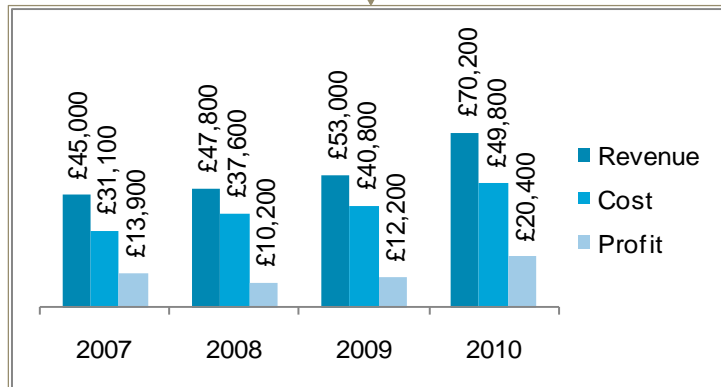
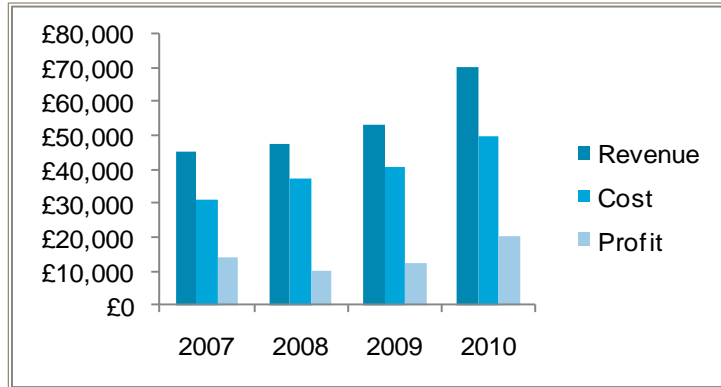
- Primary Horizontal Axis
- Primary Vertical Axis

Options for Primary Vertical Axis:

- None**: Do not display Axis
- Show Default Axis**: Display Axis with default order and labels
- Show Axis in Thousands**: Display Axis with numbers represented in Thousands (highlighted)
- Show Axis in Millions**: Display Axis with numbers represented in Millions
- Show Axis in Billions**: Display Axis with numbers represented in Billions
- Show Axis with Log Scale**: Display Axis using a log 10 based scale

[More Primary Vertical Axis Options...](#)

# Data labels

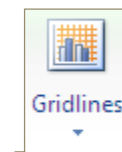
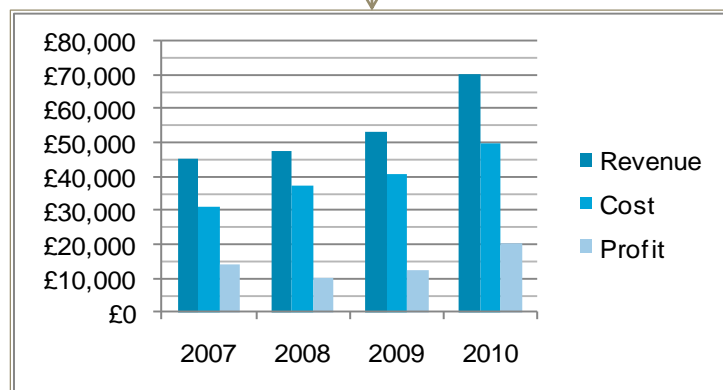
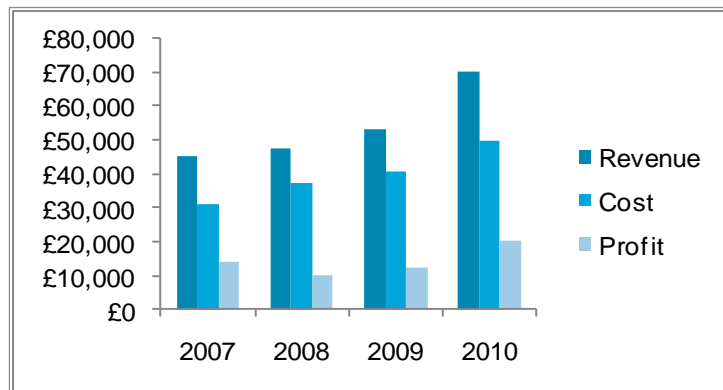


Data Labels ▾

- None**  
Turn off Data Labels for selection
- Center**  
Display Data Labels and position centered on the data point(s)
- Inside End**  
Display Data Labels and position inside the end of data point(s)
- Inside Base**  
Display Data Labels and position inside the base of data point(s)
- Outside End**  
Display Data Labels and position outside the end of data point(s)

[More Data Label Options...](#)

# Gridlines



Gridlines



Primary Horizontal Gridlines ▶



Primary Vertical Gridlines ▶



**None**

Do not display Horizontal Gridlines



**Major Gridlines**

Display Horizontal Gridlines for Major units



**Minor Gridlines**

Display Horizontal Gridlines for Minor units



**Major & Minor Gridlines**

Display Horizontal Gridlines for Major and Minor units

[More Primary Horizontal Gridlines Options...](#)



# Titles

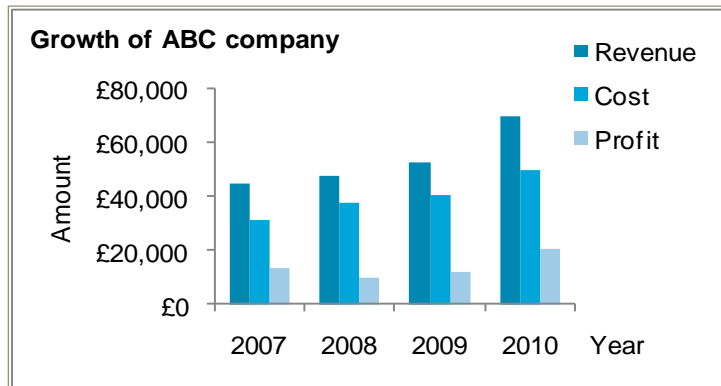
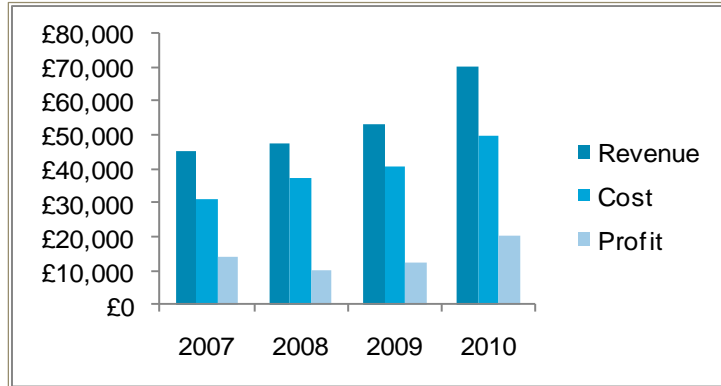


Chart Title ▾

Axis Titles ▾

- Primary Horizontal Axis Title ▶
- Primary Vertical Axis Title ▶

**None**  
Do not display a chart Title

**Centered Overlay Title**  
Overlay centered Title on chart without resizing chart

**Above Chart**  
Display Title at top of chart area and resize chart

[More Title Options...](#)

**None**  
Do not display an Axis Title

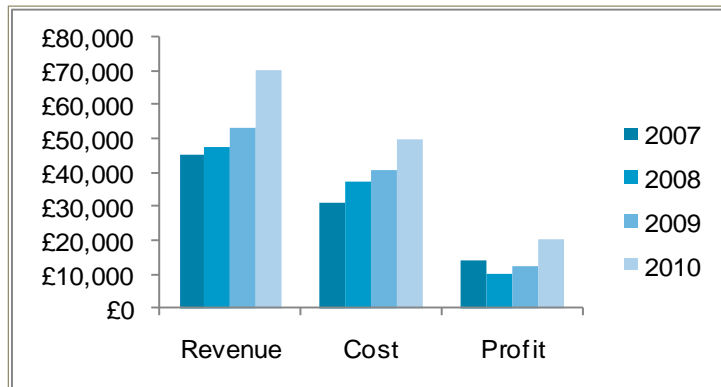
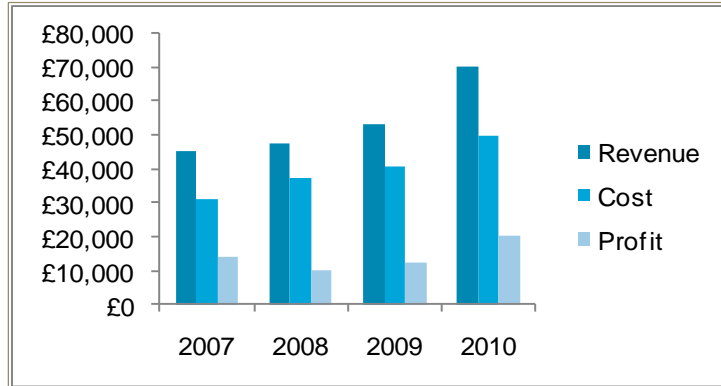
**Rotated Title**  
Display Rotated Axis Title and resize chart

**Vertical Title**  
Display Axis Title with vertical text and resize chart

**Horizontal Title**  
Display Axis Title horizontally and resize chart

[More Primary Vertical Axis Title Options...](#)

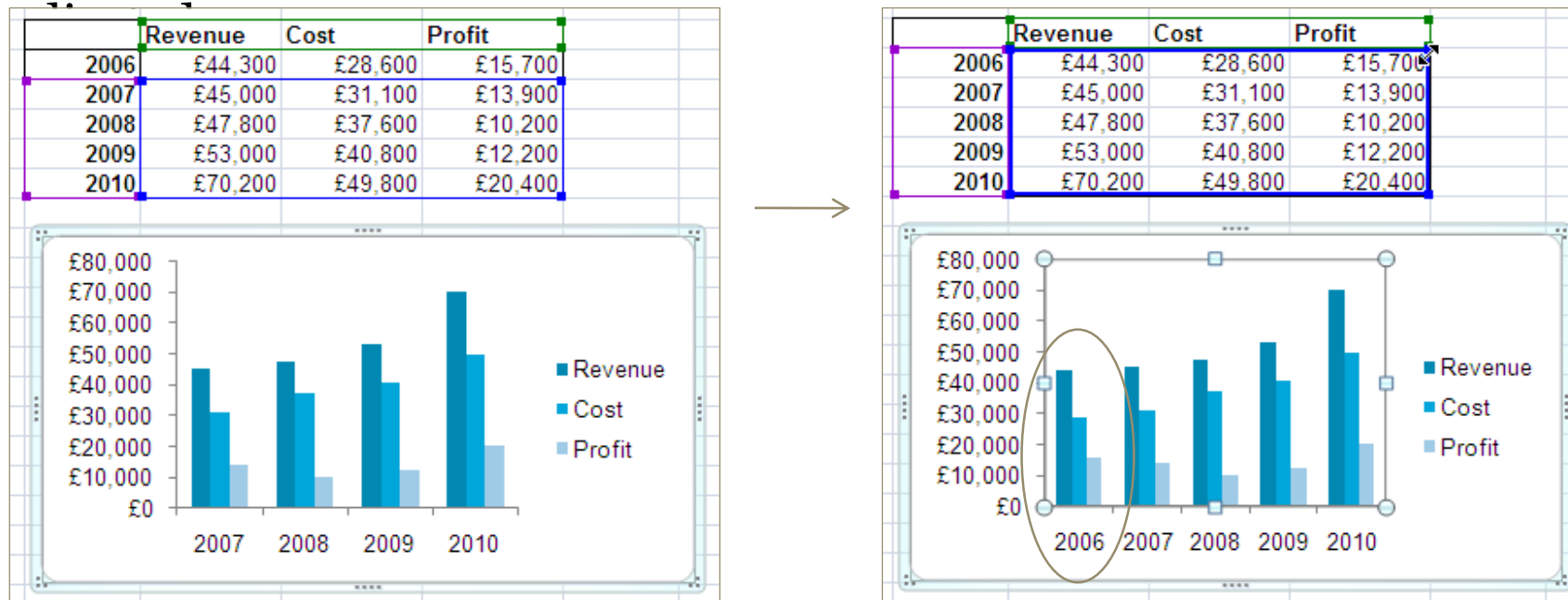
## *Data in rows or columns*



## Changing data range

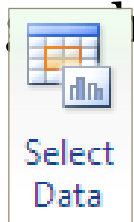
Sometimes we want to change the data range covered by a chart or add an extra series

When you select a data series on a chart, the source data range can be



## Select data source

You can also use the Select Data Source dialog box from the Design tab in the ribbon or by right-clicking the chart to add a new series to a



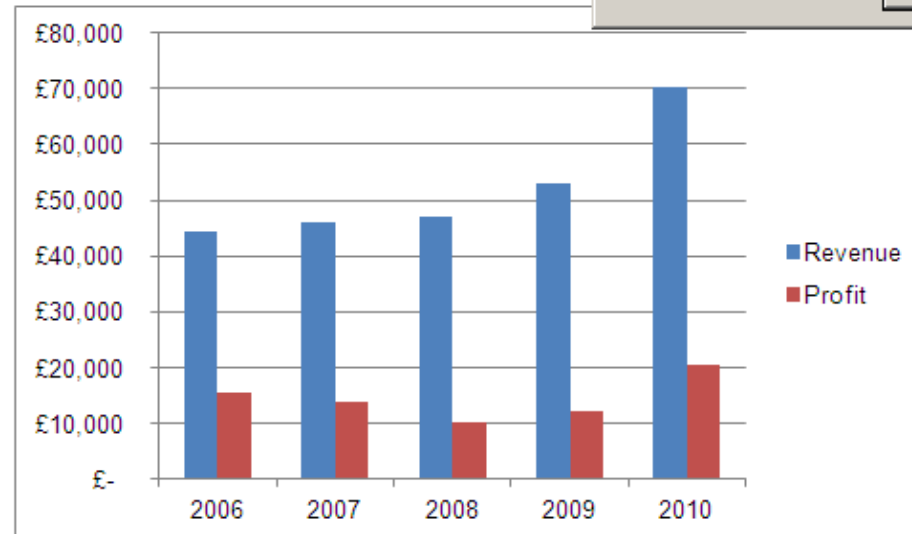
	Revenue	Cost	Profit
2006	£ 44,300	£ 28,500	£ 15,700
2007	£ 46,000	£ 31,100	£ 13,900
2008	£ 47,000	£ 37,600	£ 10,200
2009	£ 53,000	£ 40,800	£ 12,200
2010	£ 70,200	£ 49,800	£ 20,400

**Edit Series** ? X

Series name:  
=Sheet1!\$E\$3 = Profit

Series values:  
=Sheet1!\$E\$4:\$E\$8 = £15,700 , £1...

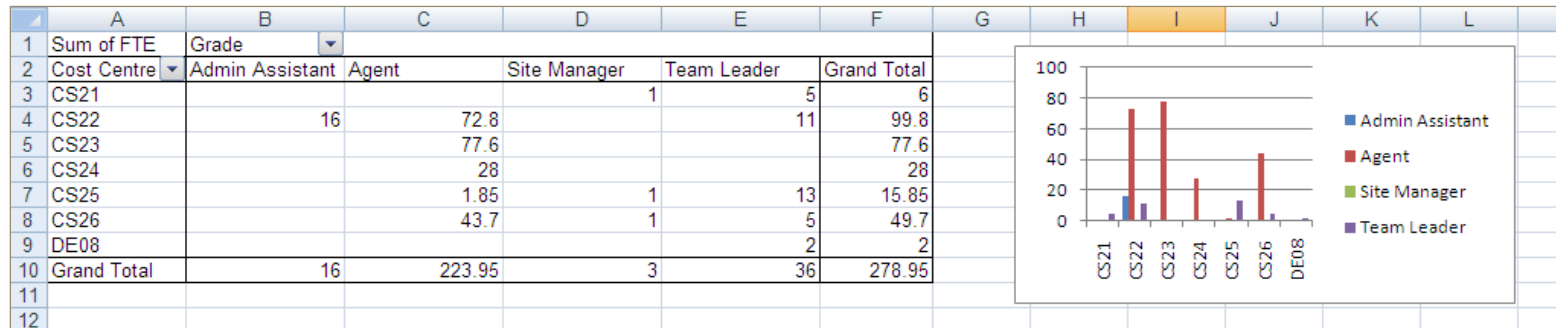
OK Cancel



# Pivot charts

In addition to creating a Pivot table on the data, you can also create a Pivot Chart which is based on the Pivot table itself

Changes made to the chart are replicated on the table and vice versa



---

## ***Exercise #4***

In 'Payroll Data.xlsx':

1. Create a pivot table with average salary by Operating Unit and State and add conditional formatting.
2. Create a bar plot with the number of employees by State
3. Create a scatter plot of salary and tenure
4. Create a line chart to show the cumulative number of employees hired by year.

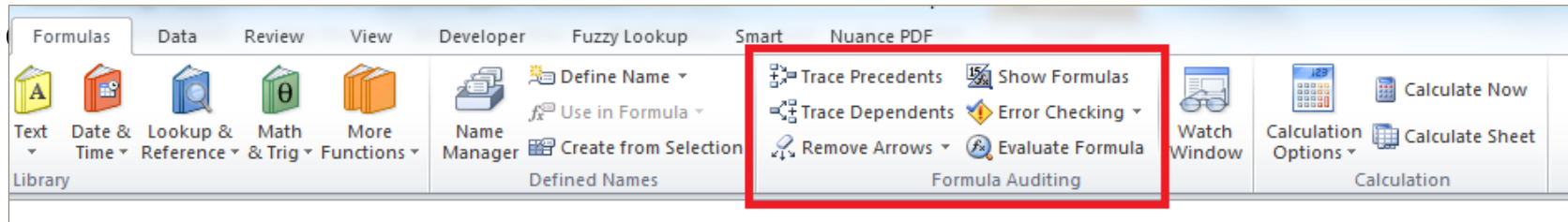
---

## *Preparing a final workbook*

- Getting an Excel file into a presentable state can be a challenge
- Storyboard your workbook and consider how someone would “read” through it
- Keep it simple!
- The following approaches can help with the finishing touches:
  - Formula auditing
  - Page layout
  - Freeze panes
  - Validation
  - Removing gridlines

# Formula auditing

To check that formulas have the right cell references, use Trace Precedents or Trace Dependents to display an arrow between linked



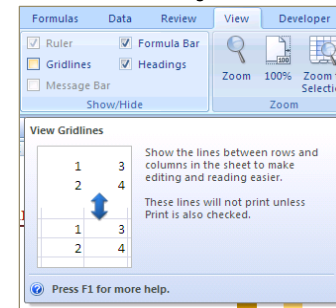
2011 Salary	Bonus	20
38,000	760	
15,300	306	
5,400	108	
12,800	256	
14,325	287	
11,250	225	
10,000	200	
41,250	825	
25,000	500	
14,350	287	
46,500	930	
4,500	90	
13,000	260	
29,500	590	
29,950	599	
12,000	240	
Total Salary		323,125

Salary	DOB	Joined	Service	Regional Salary
39,710	23-Apr-1979	18-Apr-1996	15.4	0
15,989	09-Mar-1981	01-Apr-2000	11.4	15300
5,643	04-Nov-1976	04-Dec-2008	2.7	0
13,376	24-Nov-1983	20-Apr-2001	10.4	0
14,970	30-Aug-1980	30-Sep-2003	7.9	14325
11,756	07-Oct-1983	03-Jul-2008	3.2	11250
10,450	14-Jun-1987	03-Aug-2003	8.1	0
43,106	22-Nov-1977	29-Sep-1994	16.9	41250
26,125	05-Aug-1972	11-Nov-1999	11.8	0
14,996	15-Sep-1985	31-Jul-2002	9.1	0
48,593	12-Oct-1970	23-Oct-1992	18.8	0
4,703	29-Sep-1982	10-Jan-2009	2.6	0
13,585	06-Feb-1988	15-Aug-2004	7.0	0
30,828	26-Aug-1975	23-Sep-1997	13.9	0
31,298	12-Jun-1979	03-Jun-1997	14.2	29950
12,540	29-Jul-1980	15-Aug-2004	7.0	12000
Regional Total Salary				124075

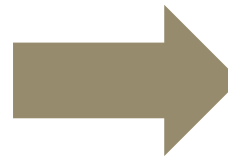


# Removing gridlines

One quick tip for making a workbook look instantly less cluttered is to remove the gridlines



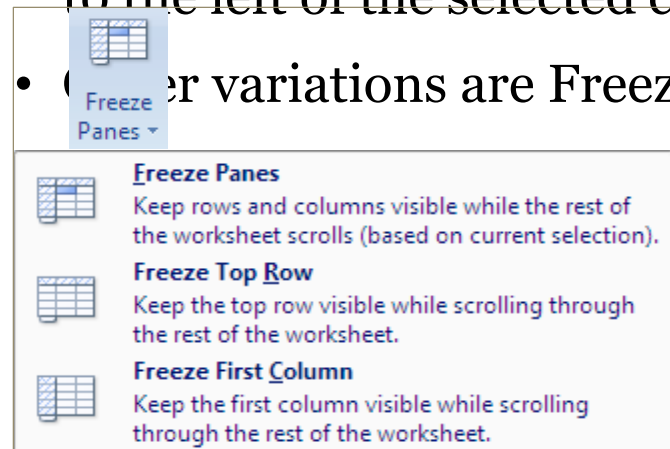
Examples:				
Approx. one third of the world's population is in two countries				
Fig.1 List of countries with largest population				
Rank	Country	Population	% of World Population	
1	China	1,340	19%	
2	India	1,210	17%	
3	US	313	4%	
4	Indonesia	238	3%	
5	Brazil	192	3%	
6	Pakistan	178	3%	
7	Nigeria	162	2%	
8	Russia	143	2%	
9	Bangladesh	142	2%	
10	Japan	128	2%	
11	Mexico	112	2%	
12	Philippines	94	1%	
13	Vietnam	87	1%	
14	Ethiopia	82	1%	
15	Germany	82	1%	
16	Egypt	81	1%	
17	Iran	76	1%	
18	Turkey	74	1%	
19	Thailand	70	1%	
20	The Congo	68	1%	
21	France	65	1%	
22	UK	62	1%	
23	Italy	61	1%	
24	South Africa	51	1%	
25	South Korea	49	1%	



Examples:				
Approx. one third of the world's population is in two countries				
Fig.1 List of countries with largest population				
Rank	Country	Population	% of World Population	
1	China	1,340	19%	
2	India	1,210	17%	
3	US	313	4%	
4	Indonesia	238	3%	
5	Brazil	192	3%	
6	Pakistan	178	3%	
7	Nigeria	162	2%	
8	Russia	143	2%	
9	Bangladesh	142	2%	
10	Japan	128	2%	
11	Mexico	112	2%	
12	Philippines	94	1%	
13	Vietnam	87	1%	
14	Ethiopia	82	1%	
15	Germany	82	1%	
16	Egypt	81	1%	
17	Iran	76	1%	
18	Turkey	74	1%	
19	Thailand	70	1%	
20	The Congo	68	1%	
21	France	65	1%	
22	UK	62	1%	
23	Italy	61	1%	
24	South Africa	51	1%	
25	South Korea	49	1%	

## *Freeze panes*

- As worksheets can get very large, it is important to ensure that the data being viewed on screen at all points has titles and comments representing the appropriate columns
- Freeze panes allows for parts of the Excel document to be frozen, useful to preserve titles and headings
- Freeze Panes freezes all rows above the selected cell and all columns to the left of the selected cell



- Other variations are Freeze Top Row and Freeze First Column

## *Page layout*

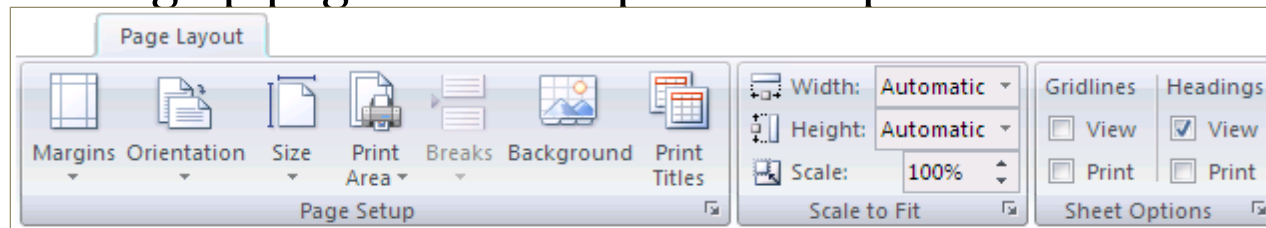
To ensure that the worksheet prints neatly, you may need to change settings on the Page Layout tab:

Changing the page orientation (Portrait vs Landscape)

Adjusting the margins

Specifying a print area (the part of the worksheet that will be printed)

Setting up page breaks at particular points

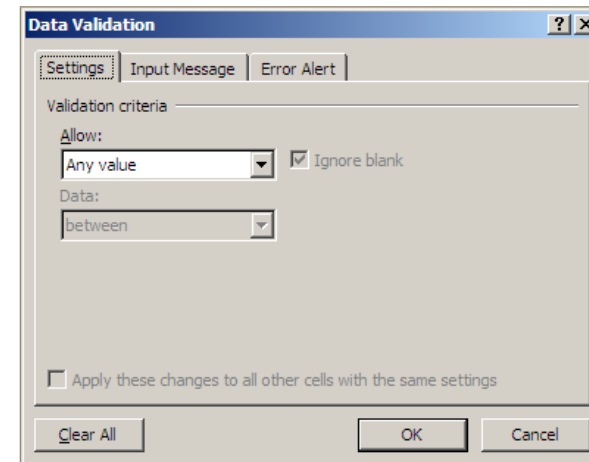
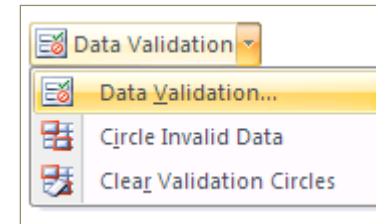


Use Print Preview to check how your worksheet will look when printed

# *Data validation*

Data validation controls what type of data a user can input into a cell, for example:

- Only a number within a certain range
- Only a time/date within a certain period
- Only an item from a predefined list



---

## *More leading practices*

- Save your workbook with A1 as the active cell
- Put a title and description for the workbook in the upper left corner of the first sheet
- Name your spreadsheet tabs so users can easily navigate throughout your workbook
- Don't hide columns or rows; instead, group them
- Consider if hardcoded parameters make sense or should be avoided
- Be careful with merged cells
- Keep source data in the workbook
- Break down complicated formulas

# *Summary*

# Acquire data

Task	Description	Excel
Data access	Connect to a data source	<ul style="list-style-type: none"><li>• File &gt; Open</li><li>• Open in text editor and copy/paste</li></ul>
Importing data	Read the data into an analytical environment	<ul style="list-style-type: none"><li>• Text Import Wizard</li><li>• Data &gt; Text to Columns</li></ul>
Data profiling	Review data dimensions and summary statistics	<ul style="list-style-type: none"><li>• COUNT()</li><li>• MIN(), MAX(), etc.</li></ul>
Data quality assessment	Identify aspects of the data that pose challenges for subsequent analysis	<ul style="list-style-type: none"><li>• Sort</li><li>• Filter</li><li>• COUNTBLANK()</li></ul>
Data simulation	Generate data based on analytical requirements	<ul style="list-style-type: none"><li>• RAND()</li><li>• RANDBETWEEN()</li><li>• CHOOSE()</li></ul>

# Transform data

Task	Description	Excel
Cleaning data	Address data quality issues to facilitate analysis	<ul style="list-style-type: none"><li>• Find/Replace</li></ul>
Changing data types	Convert a value to the appropriate format for analysis	<ul style="list-style-type: none"><li>• Format</li></ul>
Filtering data	Create subsets of records and features based on specified conditions	<ul style="list-style-type: none"><li>• Filter</li><li>• IF()</li></ul>
Deriving data	Create new features from original features	<ul style="list-style-type: none"><li>• MID()</li><li>• FIND()</li><li>• LEN()</li><li>• ROUND()</li><li>• WEEKDAY()</li><li>• ...</li></ul>
Scaling data	Put features with different ranges of values on the same scale while preserving relative values	<ul style="list-style-type: none"><li>• SUM()</li><li>• AVERAGE()</li><li>• EXP()</li></ul>



# Transform data

Task	Description	Excel
Sampling data	Create subsets of records based on a probability distribution	<ul style="list-style-type: none"><li>• RAND()</li><li>• RANDBETWEEN()</li></ul>
Aggregating data	Return a statistic or value for one feature according to different values of another feature	<ul style="list-style-type: none"><li>• Pivot Table</li></ul>
Reshaping data	Change whether values are represented in different records or different features	<ul style="list-style-type: none"><li>• Pivot Table</li></ul>
Concatenating data	Combine data sets through juxtaposition	<ul style="list-style-type: none"><li>• Cut and paste</li></ul>
Merging data	Combine data sets by matching records on a common identifier	<ul style="list-style-type: none"><li>• VLOOKUP()</li><li>• HLOOKUP()</li><li>• INDEX()/MATCH()</li></ul>

# Analyze data

Task	Description	Excel
Summary analysis	Calculate representative statistics for features of interest	<ul style="list-style-type: none"><li>• AVERAGE()</li><li>• MEDIAN()</li><li>• PERCENTILE.INC()</li></ul>
Perform statistical tests	Estimate the probability that the data supports a specific claim	<ul style="list-style-type: none"><li>• Data Analysis Toolpak</li></ul>
Clustering	Identify similar groups of records	
Predictive modeling	Use one set of features to predict the value of another feature	<ul style="list-style-type: none"><li>• Data Analysis Toolpak</li></ul>
Network analysis		

## ***Present findings***

<b>Task</b>	<b>Description</b>	<b>Excel</b>
Data visualization	Display data using lines, shapes, colors, and other abstract representations	• Charts
Dashboarding	Create a collection of dynamic visualizations	• Charts
Exporting data	Produce output from an analytical environment for future use	• File > Save As
Make recommendations	Use results of data analysis to guide decision-making	

