Lecture 6

CORRELATION AND LINEAR REGRESSION

In this lecture we look at relationships between **pairs** of **continuous variables**. These might be

In this lecture we look at relationships between **pairs** of **continuous variables**. These might be

height and weight

In this lecture we look at relationships between **pairs** of **continuous variables**. These might be

- height and weight
- market value and number of transactions

In this lecture we look at relationships between **pairs** of **continuous variables**. These might be

- height and weight
- market value and number of transactions
- temperatures and sales of ice cream

In this lecture we look at relationships between **pairs** of **continuous variables**. These might be

- height and weight
- market value and number of transactions
- temperatures and sales of ice cream

We can examine the relationship between our two variables through a **scatter diagram** (plot one against the other).

In this lecture we look at relationships between **pairs** of **continuous variables**. These might be

- height and weight
- market value and number of transactions
- temperatures and sales of ice cream

We can examine the relationship between our two variables through a scatter diagram (plot one against the other).

Today we will look at how to examine such relationships more formally.

The following data are monthly ice cream sales at Luigi Minchella's ice cream parlour.



Month	Average Temp (°C)	Sales (£000's)
January	4	73
February	4	57
March	7	81
April	8	94
May	12	110
June	15	124
July	16	134
August	17	139
September	14	124
October	11	103
November	7	81
December	5	80

Month	Average Temp (°C)	Sales (£000's)
January	4	73
February	4	57
March	7	81
April	8	94
May	12	110
June	15	124
July	16	134
August	17	139
September	14	124
October	11	103
November	7	81
December	5	80

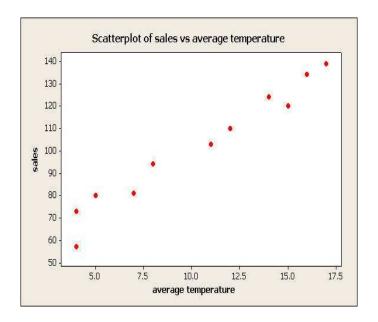
Is there any relationship between temperature and sales?

Month	Average Temp (°C)	Sales (£000's)
January	4	73
February	4	57
March	7	81
April	8	94
May	12	110
June	15	124
July	16	134
August	17	139
September	14	124
October	11	103
November	7	81
December	5	80

Is there any relationship between temperature and sales?

Can you DESCRIBE this relationship?





• The scatter plot goes "uphill", so we say there is a **positive** relationship between temperature and ice cream sales.

- The scatter plot goes "uphill", so we say there is a positive relationship between temperature and ice cream sales.
- As temperatures rise, so do ice cream sales!

- The scatter plot goes "uphill", so we say there is a positive relationship between temperature and ice cream sales.
- As temperatures rise, so do ice cream sales!
- If the scatter plot had shown a "downhill" slope, we would have a **negative** relationship.

- The scatter plot goes "uphill", so we say there is a positive relationship between temperature and ice cream sales.
- As temperatures rise, so do ice cream sales!
- If the scatter plot had shown a "downhill" slope, we would have a **negative** relationship.
- We could draw a straight line through the middle of most of the points. Thus, there is also a linear association present.

- The scatter plot goes "uphill", so we say there is a positive relationship between temperature and ice cream sales.
- As temperatures rise, so do ice cream sales!
- If the scatter plot had shown a "downhill" slope, we would have a negative relationship.
- We could draw a straight line through the middle of most of the points. Thus, there is also a linear association present.
- If we were to draw a line through the points, most points would lie close to this line. The closer the points lie to a line, the stronger the relationship.

We now look at how to **quantify** the relationship between two variables by calculating the **sample correlation coefficient**. This is

We now look at how to **quantify** the relationship between two variables by calculating the **sample correlation coefficient**. This is

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \times S_{YY}}},$$

We now look at how to **quantify** the relationship between two variables by calculating the **sample correlation coefficient**. This is

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \times S_{YY}}},$$

$$S_{XY} = \left(\sum xy\right) - n\bar{x}\bar{y},$$

We now look at how to **quantify** the relationship between two variables by calculating the **sample correlation coefficient**. This is

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \times S_{YY}}},$$

$$S_{XY} = \left(\sum xy\right) - n\bar{x}\bar{y},$$

$$S_{XX} = \left(\sum x^2\right) - n\bar{x}^2,$$

We now look at how to **quantify** the relationship between two variables by calculating the **sample correlation coefficient**. This is

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \times S_{YY}}},$$

$$S_{XY} = \left(\sum xy\right) - n\bar{x}\bar{y},$$

$$S_{XX} = \left(\sum x^2\right) - n\bar{x}^2,$$

$$S_{YY} = \left(\sum y^2\right) - n\bar{y}^2.$$

ullet r always lies between -1 and +1.

- r always lies between -1 and +1.
- If r is close to +1, we have evidence of a strong **positive** (linear) association.

- r always lies between -1 and +1.
- If r is close to +1, we have evidence of a strong **positive** (linear) association.
- If r is close to -1, we have evidence of a strong **negative** (linear) association.

- r always lies between -1 and +1.
- If r is close to +1, we have evidence of a strong **positive** (linear) association.
- If r is close to -1, we have evidence of a strong **negative** (linear) association.
- If r is close to zero, there is probably **no** association!

- r always lies between -1 and +1.
- If r is close to +1, we have evidence of a strong **positive** (linear) association.
- If r is close to -1, we have evidence of a strong **negative** (linear) association.
- If *r* is close to zero, there is probably **no** association!
- A correlation coefficient close to zero does not imply no relationship at all, just no linear relationship.

х	у	x ²	y ²	xy
4	73			

x	у	x ²	y ²	xy
4	73	16		

X	у	x ²	y ²	ху
4	73	16	5329	

x	у	x ²	y ²	ху
4	73	16	5329	292
4	57			

х	у	x ²	y ²	xy
4	73	16	5329	292
4	57	16		

х	у	x ²	y ²	xy
4	73	16	5329	292
4	57	16	3249	·

x	у	x ²	y ²	xy
4	73	16	5329	292
4	57	16	3249	228

х	у	x ²	y ²	xy
4	73	16	5329	292
4	57	16	3249	228
7	81	49	6561	567

Example: ice cream sales

The easiest way to calculate r is to draw up a table!

х	у	x ²	y ²	xy
4	73	16	5329	292
4	57	16	3249	228
7	81	49	6561	567
8	94	64	8836	752

Example: ice cream sales

The easiest way to calculate r is to draw up a table!

х	у	x ²	y ²	xy
4	73	16	5329	292
4	57	16	3249	228
7	81	49	6561	567
8	94	64	8836	752
12	110	144	12100	1320
15	124	225	15376	1860
16	134	256	17956	2144
17	139	289	19321	2363
14	124	196	15376	1736
11	103	121	10609	1133
7	81	49	6561	567
5	80	25	6400	400
120	1200	1450	127674	13362

Notice in the formulae for S_{XY} , S_{XX} and S_{YY} we need the sample means \bar{x} and \bar{y} . Thus,

Notice in the formulae for S_{XY} , S_{XX} and S_{YY} we need the sample means \bar{x} and \bar{y} . Thus,

$$ar{x} = rac{120}{12}$$
 $= 10$ and

Notice in the formulae for S_{XY} , S_{XX} and S_{YY} we need the sample means \bar{x} and \bar{y} . Thus,

$$\bar{x} = \frac{120}{12}$$

$$= 10 \quad \text{and}$$

$$\bar{y} = \frac{1200}{12}$$

$$= 100.$$

$$S_{XY} = \left(\sum xy\right) - n\bar{x}\bar{y}$$

= 13362 - 12 × 10 × 100
= 1362,

$$S_{XY} = \left(\sum xy\right) - n\bar{x}\bar{y}$$

= 13362 - 12 × 10 × 100
= 1362,
 $S_{XX} = \left(\sum x^2\right) - n\bar{x}^2$
= 1450 - 12 × 10 × 10
= 250 and

$$S_{XY} = \left(\sum xy\right) - n\bar{x}\bar{y}$$
= 13362 - 12 × 10 × 100
= 1362,
$$S_{XX} = \left(\sum x^2\right) - n\bar{x}^2$$
= 1450 - 12 × 10 × 10
= 250 and
$$S_{YY} = \left(\sum y^2\right) - n\bar{y}^2$$
= 127674 - 12 × 100 × 100
= 7674.

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \times S_{YY}}}$$

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \times S_{YY}}}$$
$$= \frac{1362}{\sqrt{250 \times 7674}}$$

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \times S_{YY}}}$$

$$= \frac{1362}{\sqrt{250 \times 7674}}$$

$$= 0.983 \text{ (to 3 decimal places)}.$$

- If your calculated correlation coefficient does not lie between
 - -1 and +1, you've done something wrong!

- If your calculated correlation coefficient does not lie between -1 and +1, you've done something wrong!
- Check that your correlation coefficient agrees with your plot

- If your calculated correlation coefficient does not lie between -1 and +1, you've done something wrong!
- Check that your correlation coefficient agrees with your plot
 - an "uphill" slope ties in with a positive correlation coefficient;

- If your calculated correlation coefficient does not lie between -1 and +1, you've done something wrong!
- Check that your correlation coefficient agrees with your plot
 - an "uphill" slope ties in with a positive correlation coefficient;
 - a "downhill" slope ties in with a negative correlation coefficient;

- If your calculated correlation coefficient does not lie between
 -1 and +1, you've done something wrong!
- Check that your correlation coefficient agrees with your plot
 - an "uphill" slope ties in with a positive correlation coefficient;
 - a "downhill" slope ties in with a negative correlation coefficient;
 - a random scattering of points indicates a correlation coefficient close to zero;

- If your calculated correlation coefficient does not lie between -1 and +1, you've done something wrong!
- Check that your correlation coefficient agrees with your plot
 - an "uphill" slope ties in with a positive correlation coefficient;
 - a "downhill" slope ties in with a negative correlation coefficient;
 - a random scattering of points indicates a correlation coefficient close to zero;
 - the closer the points lie to a straight line (either "uphill" or "downhill"), the closer to either +1 or -1 the correlation coefficient will be!

A correlation analysis helps to establish whether or not there is a linear relationship between two variables. However, it doesn't allow us to *use* this linear relationship.

A correlation analysis helps to establish whether or not there is a linear relationship between two variables. However, it doesn't allow us to *use* this linear relationship.

Regression analysis allows us to *use* the linear relationship between two variables. For example, with a regression analysis, we can predict the value of one variable given the value of another.

A correlation analysis helps to establish whether or not there is a linear relationship between two variables. However, it doesn't allow us to *use* this linear relationship.

Regression analysis allows us to *use* the linear relationship between two variables. For example, with a regression analysis, we can predict the value of one variable given the value of another.

To perform a regression analysis, we must assume that

A correlation analysis helps to establish whether or not there is a linear relationship between two variables. However, it doesn't allow us to *use* this linear relationship.

Regression analysis allows us to *use* the linear relationship between two variables. For example, with a regression analysis, we can predict the value of one variable given the value of another.

To perform a regression analysis, we must assume that

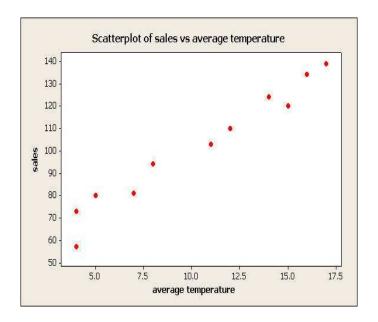
 the scatter plot of the two variables (roughly) shows a straight line, and

A correlation analysis helps to establish whether or not there is a linear relationship between two variables. However, it doesn't allow us to *use* this linear relationship.

Regression analysis allows us to *use* the linear relationship between two variables. For example, with a regression analysis, we can predict the value of one variable given the value of another.

To perform a regression analysis, we must assume that

- the scatter plot of the two variables (roughly) shows a straight line, and
- the **spread** in the Y-direction is roughly constant with X.



• A "line of best fit" can be drawn through the data;

- A "line of best fit" can be drawn through the data;
- This line could then be used to **make predictions** of ice cream sales based on temperature.

- A "line of best fit" can be drawn through the data;
- This line could then be used to make predictions of ice cream sales based on temperature.
- However, everyone's line is bound to be slightly different!

- A "line of best fit" can be drawn through the data;
- This line could then be used to make predictions of ice cream sales based on temperature.
- However, everyone's line is bound to be slightly different!
- And so everyone's predictions will be slightly different!

- A "line of best fit" can be drawn through the data;
- This line could then be used to make predictions of ice cream sales based on temperature.
- However, everyone's line is bound to be slightly different!
- And so everyone's predictions will be slightly different!
- The aim of regression analysis is to find the "best" line which goes through the data.

The simple linear regression model is given by

$$Y = \alpha + \beta X + \epsilon,$$

where

The simple linear regression model is given by

$$Y = \alpha + \beta X + \epsilon,$$

where

• Y is the response variable and

The simple linear regression model is given by

$$Y = \alpha + \beta X + \epsilon,$$

where

- Y is the response variable and
- *X* is the **explanatory variable**.

The simple linear regression model is given by

$$Y = \alpha + \beta X + \epsilon,$$

where

- Y is the response variable and
- *X* is the **explanatory variable**.
- α represents the **intercept** of the regression line (the point where the line "cuts" the Y-axis),

The simple linear regression model is given by

$$Y = \alpha + \beta X + \epsilon,$$

where

- Y is the response variable and
- *X* is the **explanatory variable**.
- α represents the **intercept** of the regression line (the point where the line "cuts" the Y-axis),
- $oldsymbol{\beta}$ represents the **slope** of the regression line (i.e. how steep the line is), and

The simple linear regression model is given by

$$Y = \alpha + \beta X + \epsilon,$$

where

- Y is the response variable and
- *X* is the **explanatory variable**.
- α represents the **intercept** of the regression line (the point where the line "cuts" the Y-axis),
- $oldsymbol{\circ}$ eta represents the **slope** of the regression line (i.e. how steep the line is), and
- \bullet ϵ is known as "random error"

The simple linear regression model is given by

$$Y = \alpha + \beta X + \epsilon,$$

where

- Y is the response variable and
- *X* is the **explanatory variable**.
- α represents the **intercept** of the regression line (the point where the line "cuts" the Y-axis),
- $oldsymbol{\circ}$ eta represents the **slope** of the regression line (i.e. how steep the line is), and
- \bullet ϵ is known as "random error"

In practice, we assume ϵ is zero, and so the only things we need to find are α and β . But how?



 Remember, the aim of regression analysis is to find a line which goes through the middle of the data in the scatter plot, closer to the points than any other line;

- Remember, the aim of regression analysis is to find a line which goes through the middle of the data in the scatter plot, closer to the points than any other line;
- So the "best" line will minimise any gaps between the line and the data!

- Remember, the aim of regression analysis is to find a line which goes through the middle of the data in the scatter plot, closer to the points than any other line;
- So the "best" line will minimise any gaps between the line and the data!
- It turns out that the values of α and β which give this best line are

$$\hat{\beta} = \frac{S_{XY}}{S_{XX}}$$
 and

- Remember, the aim of regression analysis is to find a line which goes through the middle of the data in the scatter plot, closer to the points than any other line;
- So the "best" line will minimise any gaps between the line and the data!
- It turns out that the values of α and β which give this best line are

$$\hat{\beta} = \frac{S_{XY}}{S_{XX}}$$
 and

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}.$$

Example: ice cream sales

For the ice cream sales data, we can find the "best" line using the formulae on the previous slide! Thus

Example: ice cream sales

For the ice cream sales data, we can find the "best" line using the formulae on the previous slide! Thus

$$\hat{\beta} = \frac{S_{XY}}{S_{XX}}$$

$$= \frac{1362}{250}$$

$$= 5.448,$$

and

Example: ice cream sales

For the ice cream sales data, we can find the "best" line using the formulae on the previous slide! Thus

$$\hat{\beta} = \frac{S_{XY}}{S_{XX}}$$

$$= \frac{1362}{250}$$

$$= 5.448,$$

and

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

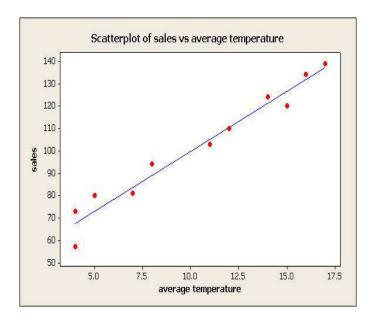
$$= 100 - 5.448 \times 10$$

$$= 45.52.$$

Thus, the regression equation is

Thus, the regression equation is

$$Y = 45.52 + 5.448X + \epsilon$$



Making predictions

We can use our estimated regression equation to make **predictions** of ice cream sales based on average temperature.

For example: Predict monthly sales if the average temperature is 10°C .

For example: Predict monthly sales if the average temperature is 10° C.

We can use take a reading from our graph, or, more accurately, use our regression equation!

$$Y = 45.52 + 5.448X$$
 i.e.
= $45.52 + 5.448 \times 10$
= $45.52 + 54.48$
= 100 ,

For example: Predict monthly sales if the average temperature is 10°C.

We can use take a reading from our graph, or, more accurately, use our regression equation!

$$Y = 45.52 + 5.448X$$
 i.e.
= $45.52 + 5.448 \times 10$
= $45.52 + 54.48$
= 100 ,

i.e. if the average temperature is 10°C , we can expect sales of £100,000.