

Correlation and Regression

Load important statistical concepts directly into your brain

Nikesh Bajaj
nikesh.14730@lpu.co.in
Asst. Prof., ECE Dept.
Lovely Professional University


Two Variable

- Till now, We dealt with one variable
- Let's try to find out,
 - How two things are connected?
 - How they effect each other?

117 By Nikesh Bajaj

Concert and Weather

- *Guys: organizing concerts*
Concert are best in open air
- *Ticket sales in summer look promising*



118 By Nikesh Bajaj

Let's Analyze and predict

- Sunshine and Attendance of audience

Sunshine (hours)	1.9	2.5	3.2	3.8	4.7	5.5	5.9	7.2
Concert attendance (100's)	22	33	30	42	38	49	42	55

- *Scenario:* Temperature is dipping, look like rain, guys want to predict attendance of audience for given hours of sun shine.
- If attendance will be less than 3500, where ticket won't even cover the expenses they will **cancel the event**.
- What you can do with given data?

119 By Nikesh Bajaj

What sort of analysis you suggest?

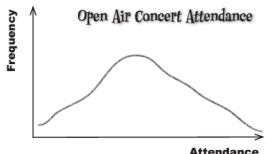
Sunshine (hours)	1.9	2.5	3.2	3.8	4.7	5.5	5.9	7.2
Concert attendance (100's)	22	33	30	42	38	49	42	55

How would you go about modelling the connection between sets of data?

120 By Nikesh Bajaj

Exploring types of Data

- *Univariate Data:* Frequency or probability of one variable, e.g. weight, player's score etc. "One thing"
- It does not tell connection between two
- If



121 By Nikesh Bajaj

Exploring types of Data

- **Bivariate Data:** Values of two variable for each observation.

Sunshine (hours)	1.9	2.5	3.2	3.8	4.7	5.5	5.9	7.2
Concert attendance (100's)	22	33	30	42	38	49	42	55

- **Independent or Explanatory variable**
 - One of variable has been controlled in some way or used to explain other
- **Dependent or Response variable**

■ So Which is which for our example?

122 By Nikesh Bajaj

Visualizing bivariate data

- **Scatter plot or scatter diagram: DOES IT HELPS?**

x (sunshine)	1.9	2.5	3.2	3.8	4.7	5.5	5.9	7.2
y (attendance)	22	33	30	42	38	49	42	55

SO WHAT YOU CAN OBSERVE??

123 By Nikesh Bajaj

Correlation

- Scatter diagram shows the correlation between two variable
- Correlation
 - Linear: If it is straight line, can be others

124 By Nikesh Bajaj

Correlation

Positive linear correlation

Negative linear correlation

No correlation

125 By Nikesh Bajaj

Correlation Coefficient r

- r tells you kind of correlation, positive, negative, perfect or no

r = -1

r = 0

r = 1

126 By Nikesh Bajaj

Computing r

$$r = \frac{b s_x}{s_y}$$

Im the correlation coefficient, r. I say how strong the correlation is between the two variables.

r

$$b = \frac{\sum((x - \bar{x})(y - \bar{y}))}{\sum(x - \bar{x})^2}$$

$$s_x = \sqrt{\frac{\sum(x - \bar{x})^2}{n - 1}} \quad s_y = \sqrt{\frac{\sum(y - \bar{y})^2}{n - 1}}$$

127 By Nikesh Bajaj

Correlation and Causation

- If there is correlation between two variable Does that mean one caused the value of other??

128 By Nikesh Bajaj

Let's See example

- One intern was given many scatter plot of..

129 By Nikesh Bajaj

Correlation and Causation

- If there is correlation between two variable Does that mean one caused the value of other??
- "Not always"
- Let's see example

Coffee shops vs. record shops

130 By Nikesh Bajaj

So for Concert

- Sunshine effect Attendance very much
- Good but
- *What about attendance of 3500 people??*

Concert Attendance and Sunshine

131 By Nikesh Bajaj

Predict the Attendance

- Line of Best Fit

132 By Nikesh Bajaj

Find a Line $y = a + bx$

133 By Nikesh Bajaj

Line with minimum Error $\Sigma(y_i - \hat{y}_i)$

- Error

These are the actual values.

$y = a + bx$

$\Sigma(y_i - \hat{y}_i)$

SSE = $\Sigma(y - \hat{y})^2$

134 By Nikesh Bajaj

Let's find Line $y = a + bx$

- b: Steepness of line, Slope

$$b = \frac{\Sigma((x - \bar{x})(y - \bar{y}))}{\Sigma(x - \bar{x})^2}$$

No need proof right now

x (sunshine)	1.9	2.5	3.2	3.8	4.7	5.5	5.9	7.2
y (attendance)	22	33	30	42	38	49	42	55

- Find b = ?

135 By Nikesh Bajaj

What about 'a' ???

- How to compute?
- $y = a + bx$

$$\bar{y} = a + b\bar{x}$$

$$a = \bar{y} - b\bar{x}$$

136 By Nikesh Bajaj

Solution (*Linear Regression*)

- Line of best fit

attendance (100's) y

sunshine (hours) x

$y = 15.80 + 5.32x$

137 By Nikesh Bajaj

Now Answer Concert guys

- $y = 15.8 + 5.32x$
 - Q1. When predicted sunshine is 6 Hours what would be attendance of audience in concert?
- Q2. what should be sunshine hours for at least audience of 3500

138 By Nikesh Bajaj

Answers

- Ans 1: $y = 47.72$ means 4772 people
- Ans 2: $x = 3.61$ Hours

139 By Nikesh Bajaj

