



`</talentlabs>`

Bivariate Analysis

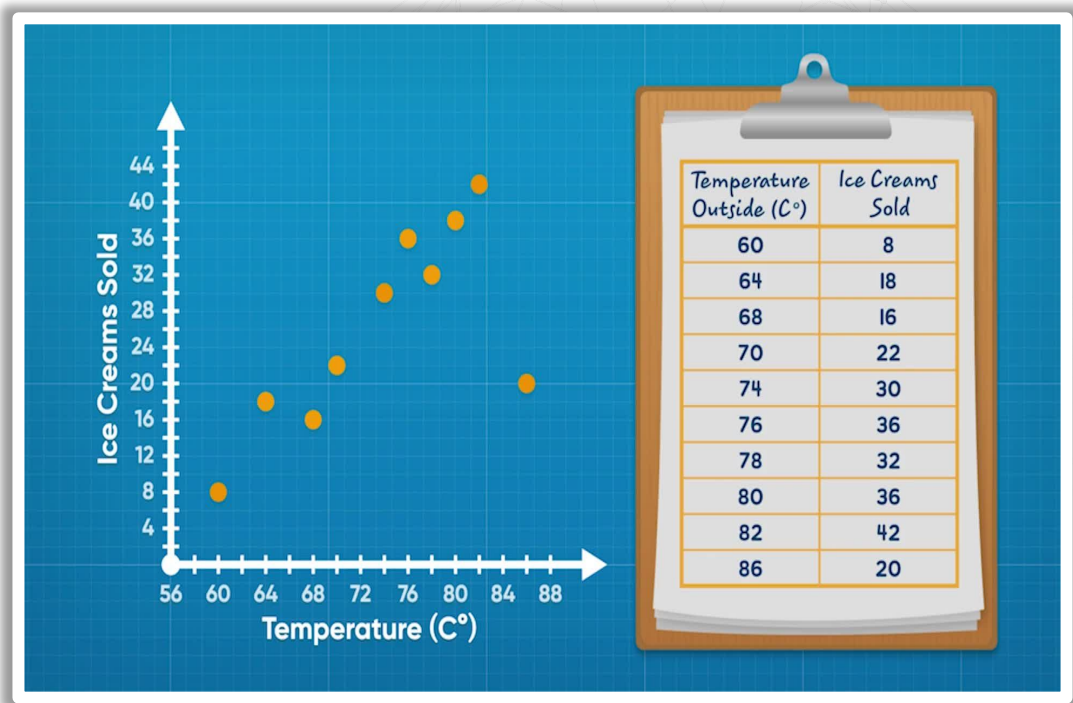
Agenda

1. What is Bivariate analysis and why is it so important?
2. Recap – Covariance and Correlations
3. Quantitative vs Quantitative
4. Quantitative vs Qualitative
5. Qualitative vs Qualitative Analysis
7. Statistical analysis using t-tests and anova



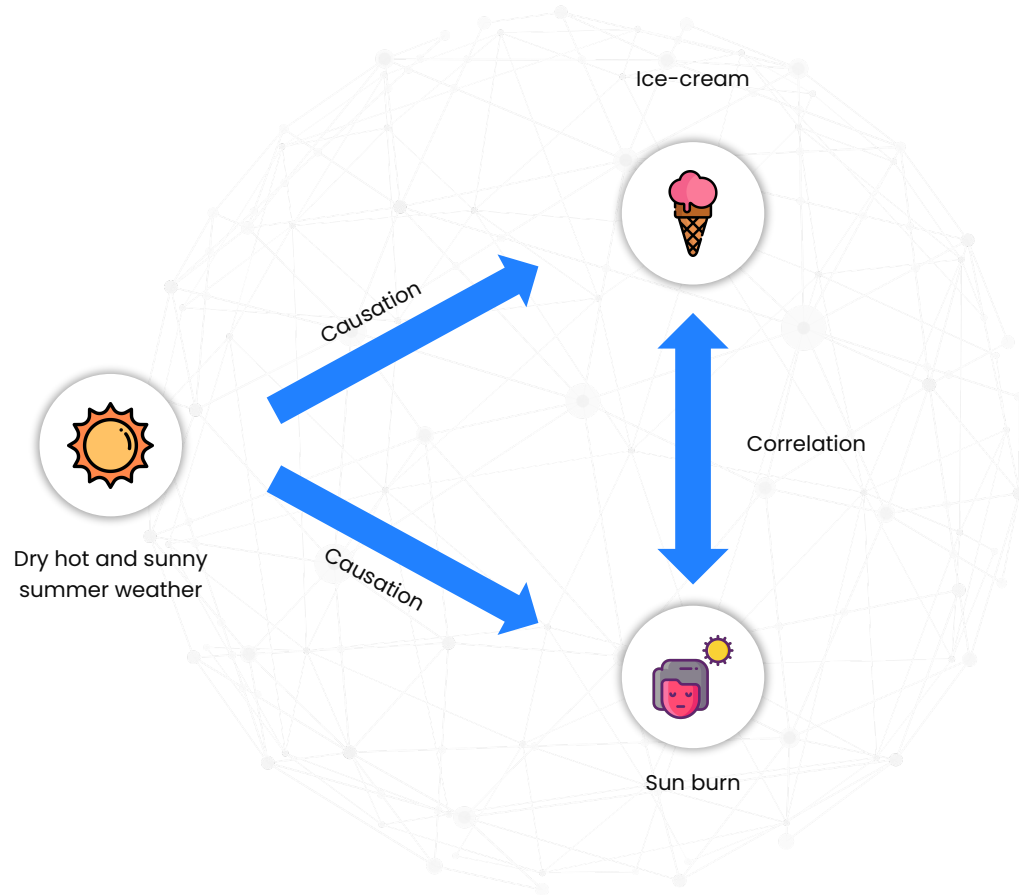
What is Bivariate Analysis?

- It involves the analysis of two variables
- purpose of determining the empirical relationship between the two variables
- Want to test a hypothesis of association or relationship?
- The results from bivariate analysis can be stored in a two-column data table.



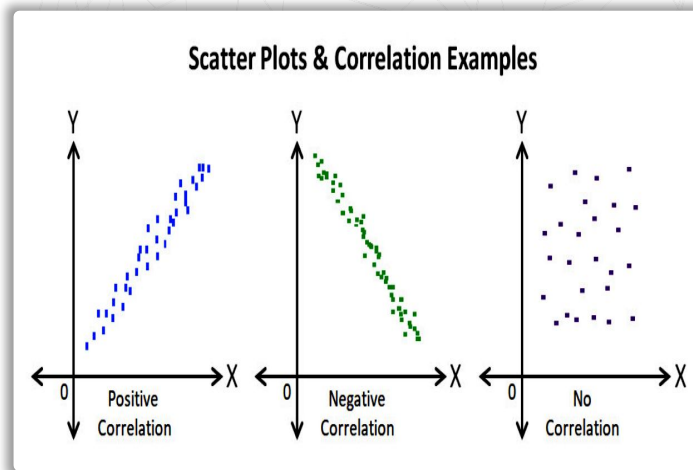
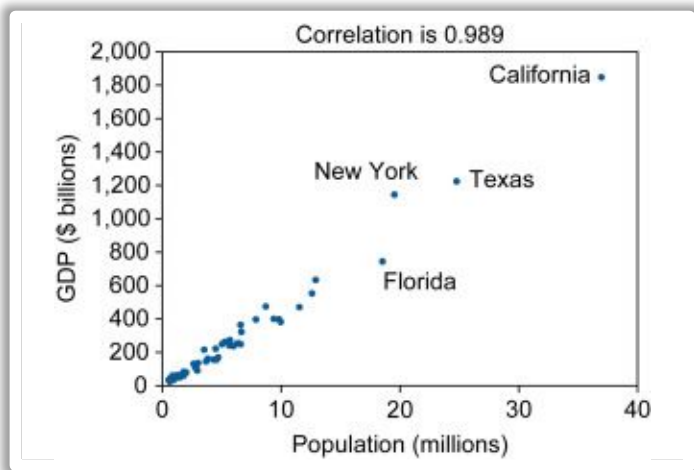
Statistics to represent relationships

- Covariance is a statistic that measures the size of the relationship between two variables.
- Correlation is a measure that determines the strength of the relationship between two variables.
- A myth: correlation = causation.



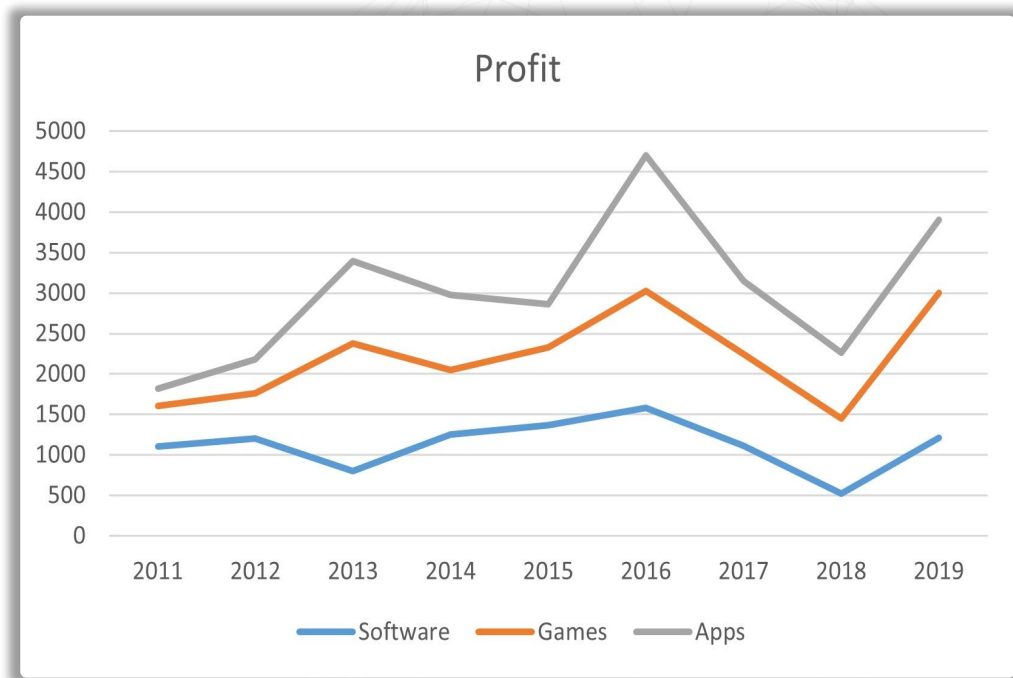
The Scatterplot

- Relationship between numerical or continuous variables.
- **The graphic representation of the relationship between the two variables coming from a bivariate data set.**
- Think of them as the graphic representation of two data sets which have been put into place by dedicating each axis in **the plot to a different variable.**
- **Here, majorly we look into plotting data points without connectivity**



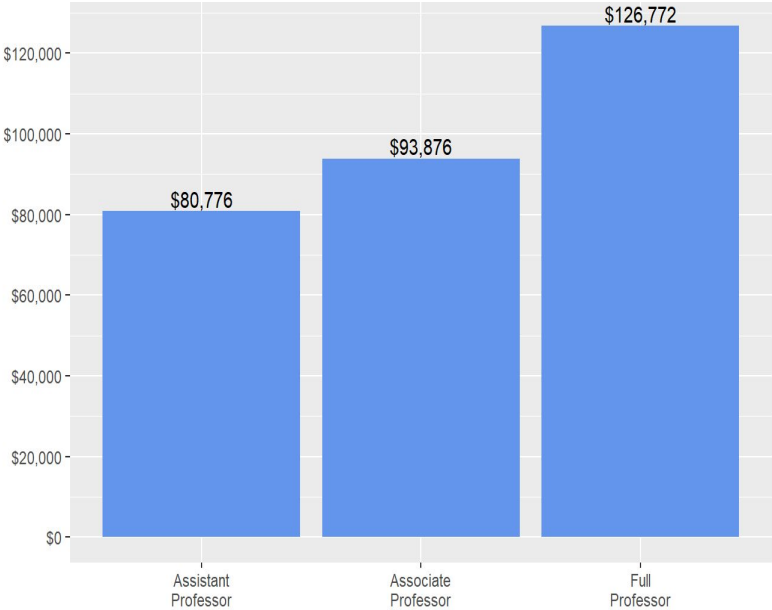
The Line plot

- Relationship between numerical or continuous variables.
- **The graphic representation of the relationship between the two variables coming from a bivariate data set.**
- Here, **majorly we look into plotting data points with connectivity, but we can't see the points themselves.**
- More of trend line plots

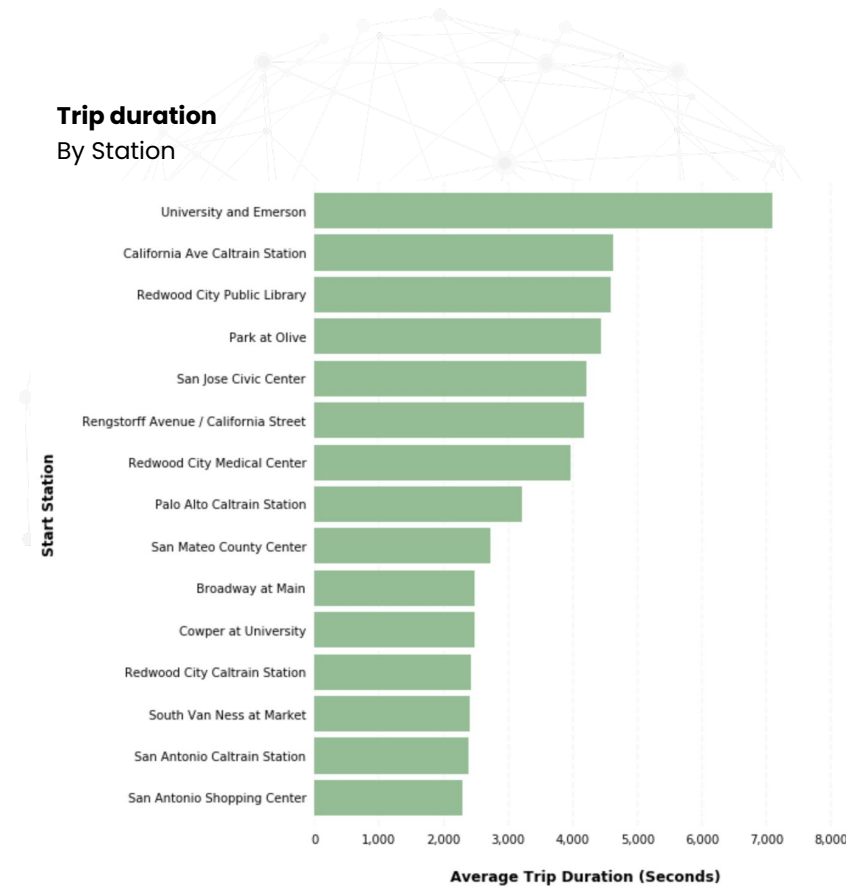


Bar Charts

9 Month average academic salary
By Professor Level

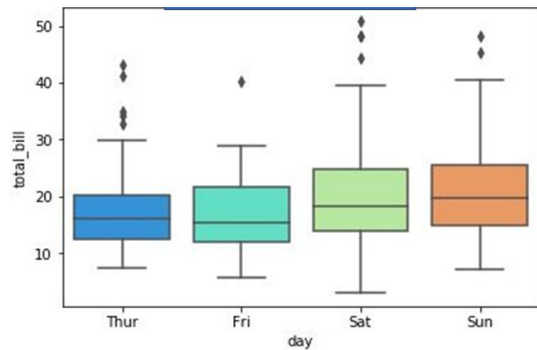


Trip duration
By Station

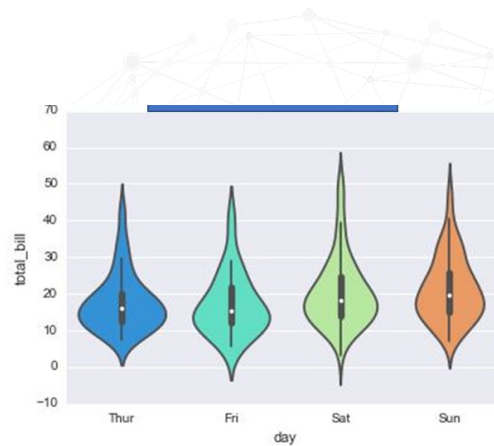


Box plots

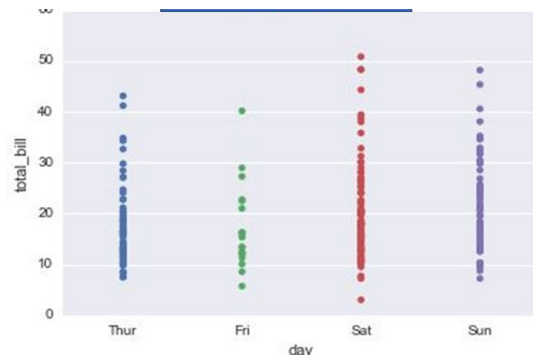
Box plot



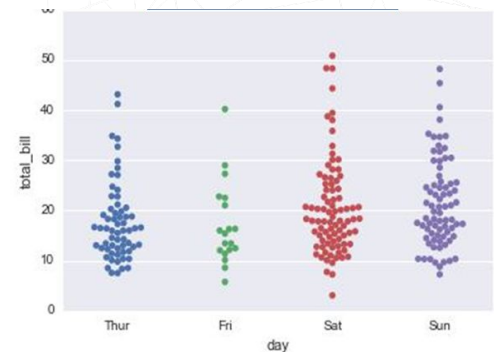
Violin plot



Strip plot



Swarm plot



Boxplot – Two variable example

- What? – I want to know how the distribution of heart rate differs for people resting, walking and running.
- Assumption that with more exercise activity, the median heart rate increases.

Id	Diet	Pulse	Time	Kind
5	Low fat	91	30 min	Rest
21	Low fat	93	1 min	Running
27	No fat	100	1 min	Running
21	Low fat	110	30 min	Running
4	Low fat	80	1 min	Rest



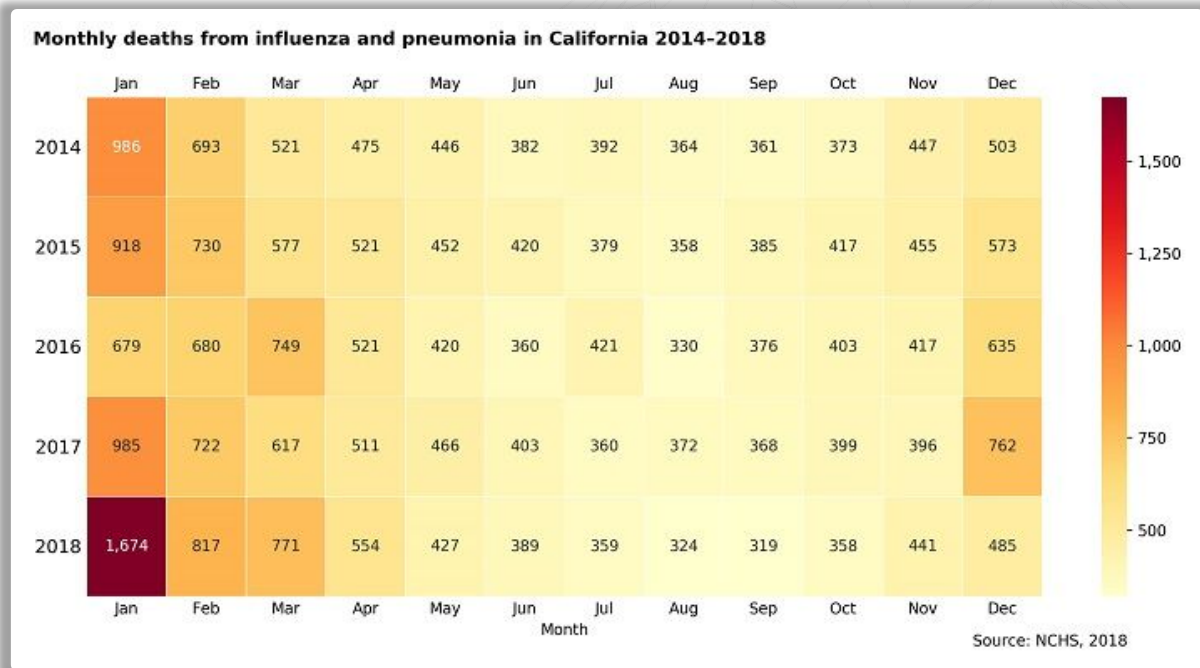
Cross Tabulations

- Both variables are categorical
- It is used to count between categories, or get summaries between two categories.

	Cust	Categ_X	Categ_y	Categ_y	a	b	All
0	1	a	a	Categ_X			
1	1	a	b	a	3	2	5
2	1	b	a	b	2	4	6
3	1	b	b	All	5	6	11
4	2	a	a				
5	3	a	a				
6	3	a	b				
7	3	b	a				
8	3	b	b				
9	4	b	b				
10	5	b	b				

Heat Maps in contingency tables

Contingency table with colors showing strength in counts/frequencies



Bar Charts with two variables

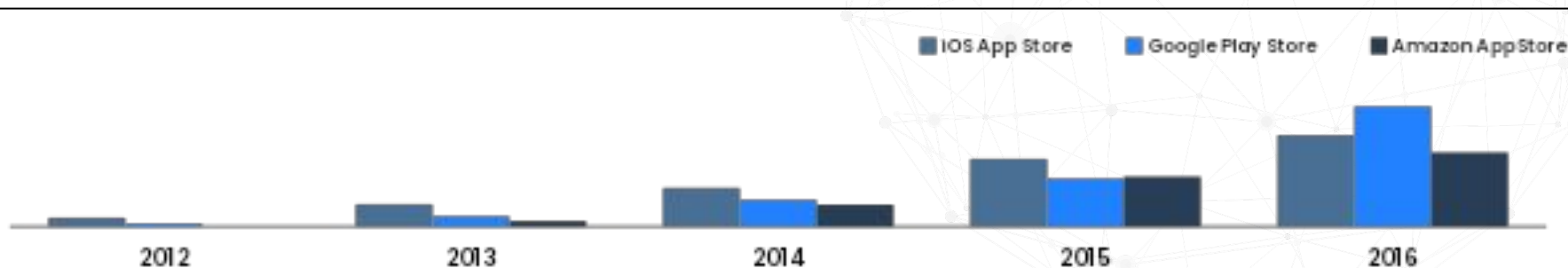
Distribution of population in South Africa

By Provinces



App Publishing trend

2012 – 2016





</talentlabs>

Thank you

Pivot tables

- Wrangle the data around the pivot to analyze better.
- Derived from Excel.
- Shift the data and aggregate it around the pivot.
- Explains better than normal dataframes

