# Exploratory Data Analysis (EDA)

PART TWO

# Exploratory Data Analysis

Utilize data's statistical attributes

- Temporal comparison
- Attributes comparison
- Ranking comparison
- Composition analysis
- Distributions analysis
- Variance analysis
- Correlation analysis
- Geographic analysis

# Distribution analysis

Compare distribution of attributes
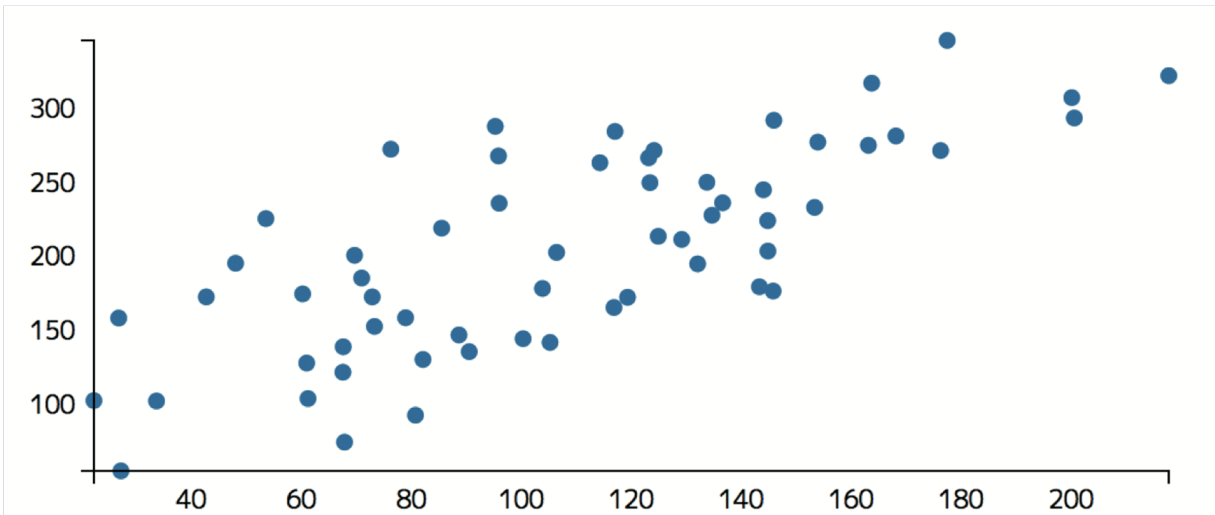- Visualise /compare
- Clusters
- Spread /distribution

# Basic use of various visualisation for Distribution analysis

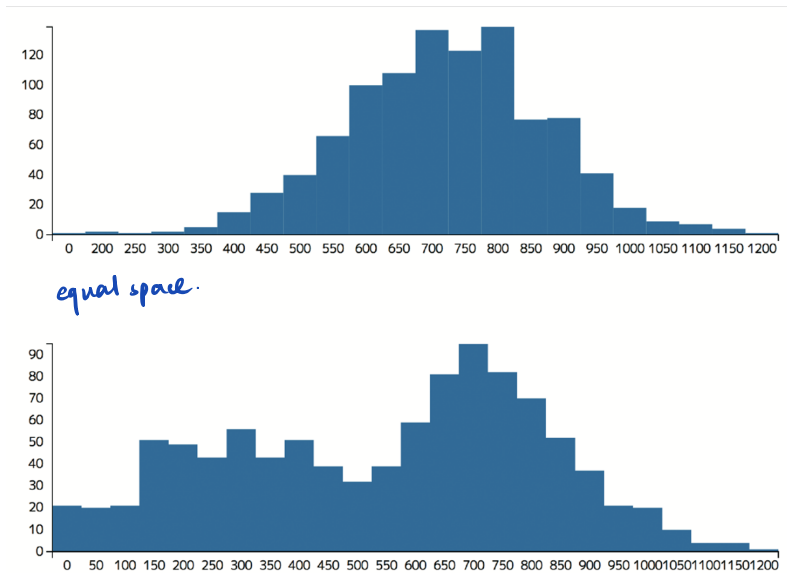| Basic | | Scatter plot |
|---|---|---|
| Overall Distribution | | Histogram |
| Distribution of each attributes | | Box plot |

# Scatter plot



- points

# Histogram



equal space.

- area (width/height)

# Box plot



- line, points (3) : median, max, min

# Variance Analysis

Compare distribution of attributes with respect to the average

- Visualise /compare
- Clusters
- Spread /distribution

# Correlation analysis

Compare distribution of attributes with respect to the average
- Visualise /compare
- Positive /negative correlation amount attributes

# Basic use of various visualisation for Correlation analysis
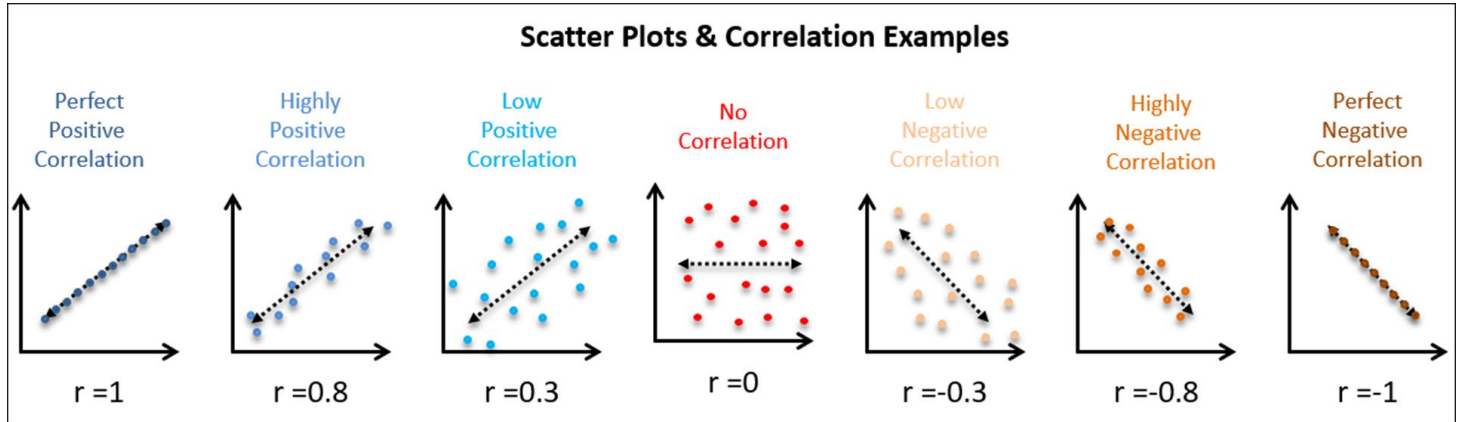
| Two attributes | | Scatter plot |
|---|---|---|
| Three attributes | | Bubble chart |

Icons from Data Visualization, Toshikuni and Watanabe (2019)

# Correlation Scatter plot



Scatter Plots & Correlation Examples

Perfect Positive Correlation — r =1
Highly Positive Correlation — r =0.8
Low Positive Correlation — r =0.3
No Correlation — r =0
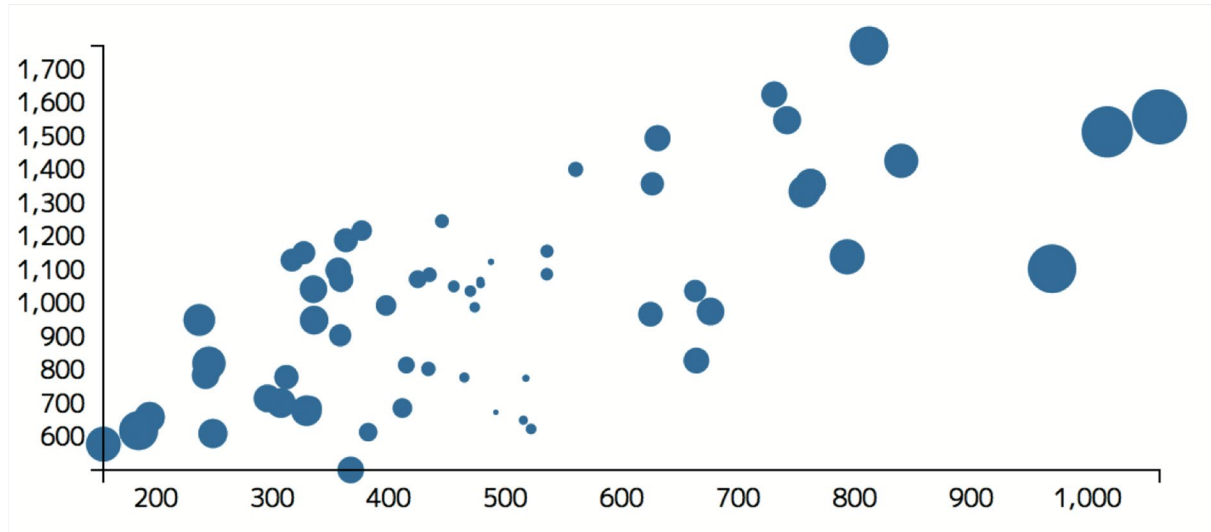Low Negative Correlation — r =-0.3
Highly Negative Correlation — r =-0.8
Perfect Negative Correlation — r =-1

# Bubble chart



- Circle : value ↔ area

# Geographical analysis

- Analyse location/arrangement of attributes based on geo-referenced information

# Basic use of various visualisation for Geographical analysis

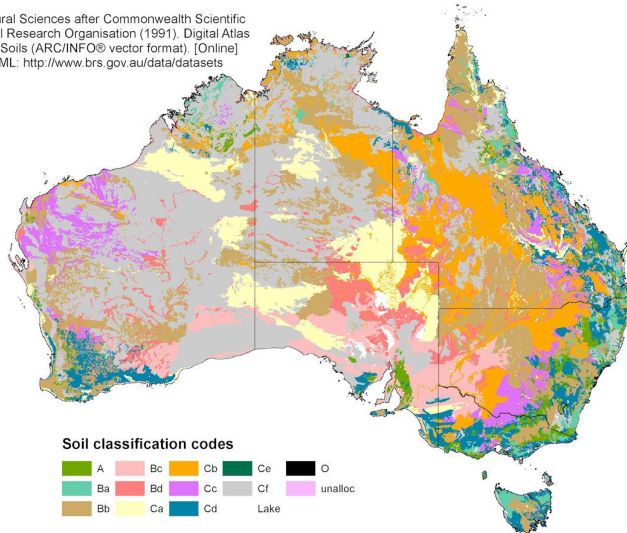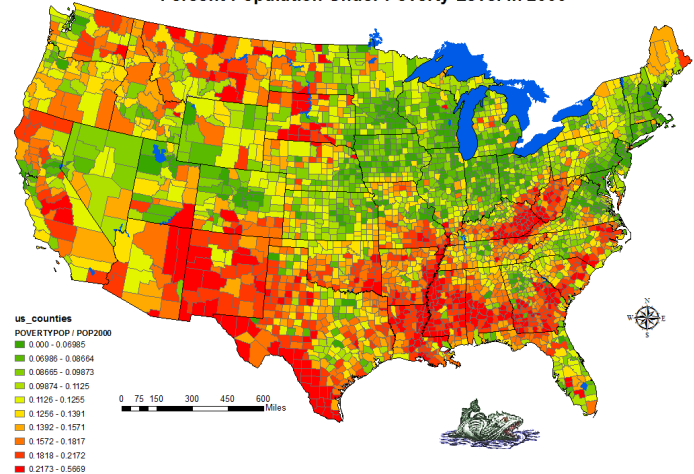| Compare an index among regions |  | Thematic map |
|---|---|---|
| Compare multiple indices among regions |  | Symbol map |

# Thematic Map



Bureau of Rural Sciences after Commonwealth Scientific and Industrial Research Organisation (1991). Digital Atlas of Australian Soils (ARC/INFO® vector format). [Online] Available HTML: http://www.brs.gov.au/data/datasets
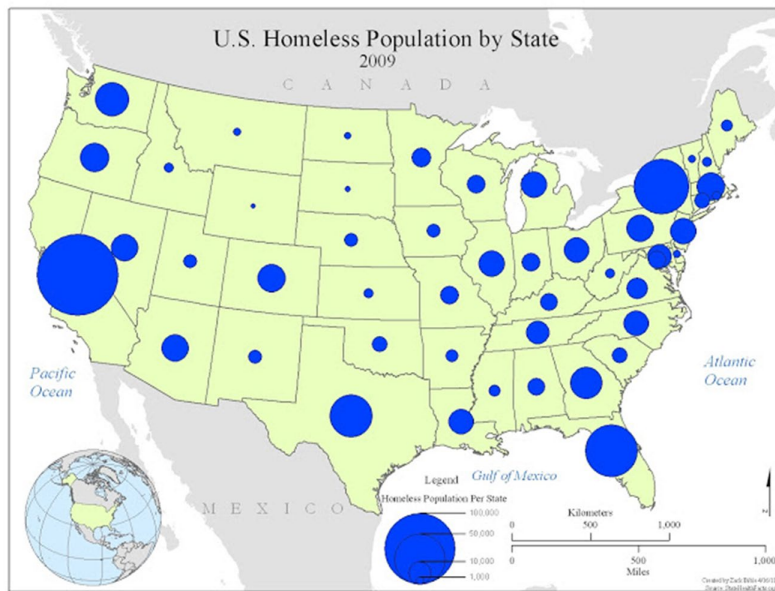
**Soil classification codes**

| | | | | | |
|---|---|---|---|---|---|
| A | Bc | Cb | Ce | O | |
| Ba | Bd | Cc | Cf | unalloc | |
| Bb | Ca | Cd | Lake | | |

**Percent Population Under Poverty Level in 2000**

us_counties
POVERTYPOP / POP2000
- 0.000 - 0.06985
- 0.06986 - 0.08664
- 0.08665 - 0.09873
- 0.09874 - 0.1125
- 0.1126 - 0.1255
- 0.1256 - 0.1391
- 0.1392 - 0.1571
- 0.1572 - 0.1817
- 0.1818 - 0.2172
- 0.2173 - 0.5669

0  75 150    300    450    600
Miles

https://www.globalsecurity.org/jhtml/jframe.html#https://www.globalsecurity.org/military/world/australia/images/australia-soils-1.jpg|||

https://mapgeeks.org/different-types-maps/

# Symbol Map



http://mph720-jennifer-2012.blogspot.com

# Visualising Statistical Feature

# Statistics from Data

What sort of statistical feature can you get from the dataset?
- Average,
- Range (minimum and maximum)
- Median
- Variance
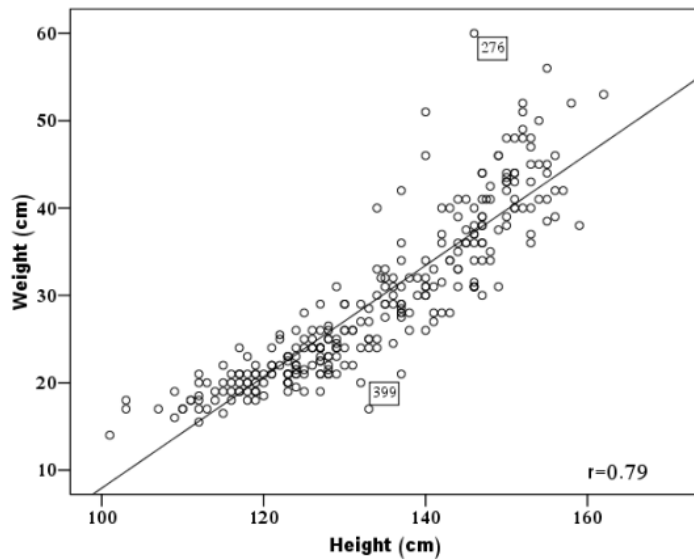- Standard Deviation
- Quartile
- Skewness
- Kurtosis

etc.

# Statistical feature to Visualisation

Statistical values (as indices) can be mapped to graph/chart for visualisation:
- Scatter Plot,
- Histogram
- Probability Plot (Q-Q (quantile-quantile) plot, P-P (Prob-Prov) plot)
- Spaghetti Plot
- Residual Plot
- Box Plot
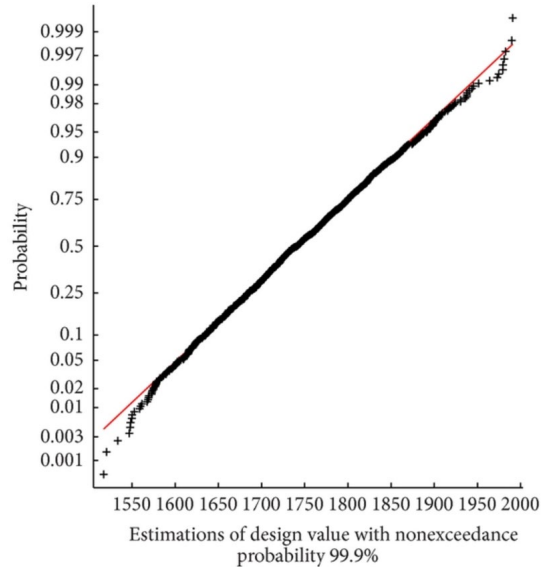- Block Plots
- Biplots

etc.
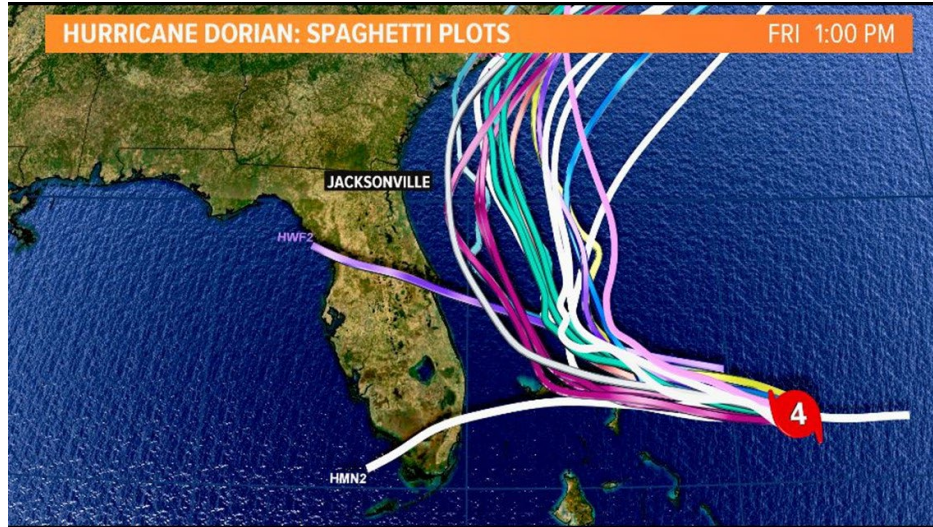
# Scatter Plot

# Histogram



https://pythonclass.in/images/histogram-matplotlib-example.jpg

# Probability plot (QQ plot, PP plot)



Estimations of design value with nonexceedance
probability 99.9%

# Spaghetti Plot

# Residual plot

# Box plot



Guinea Pigs' Tooth Growth

# Biplots (PCA)

how clusters related in original space

new dimension A



f:id:hayataka2049:20180328230804p:plain

sepal length (cm)

petal width (cm)
petal length (cm)

new dimension B

# Multi-dimensional Statistical Features Handling

# How do we handle multi-dimensional statistical features?

- Coordinated Multiview Visualization
- Spatialisation

# Coordinated Multiview Visualization

- A typical visualization used to display statistical features are 2D, or 3D.
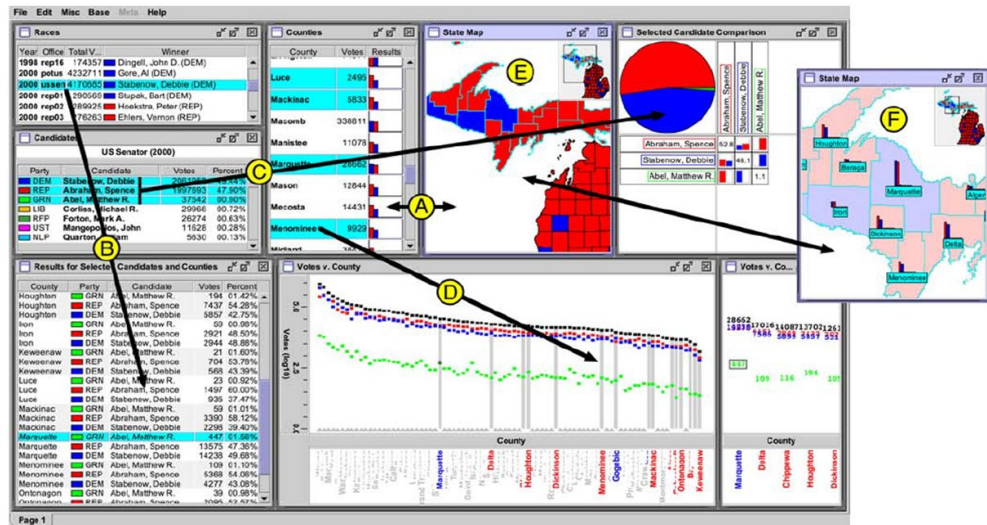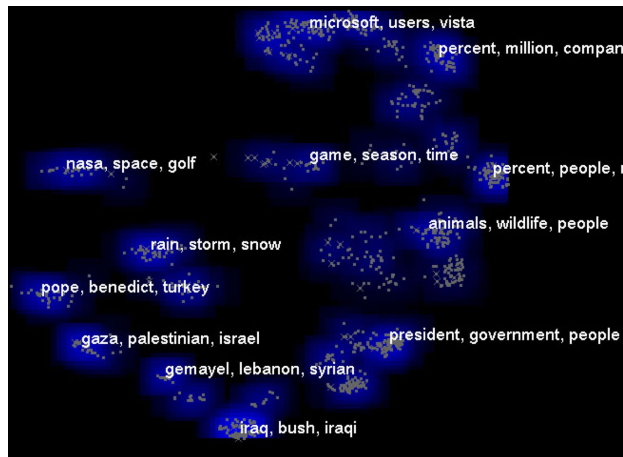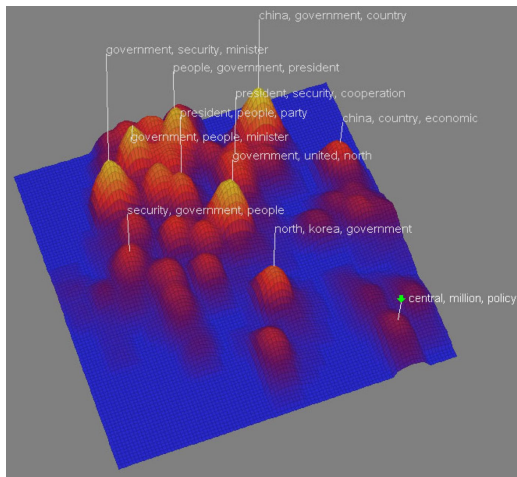


Figure 11: Visualization of election results in Michigan from 1998 to 2002. (A) Shared selection of counties between a table view and a map. (B) Selecting a race causes the election results for that race to be loaded (from a file) and shown throughout the visualization. (C) A pie chart uses a filter to compare results for selected candidates only. (D) A scatterplot highlights selected counties with gray bars. (E) A four-layer scatterplot colors counties by winning candidate party. (F) Semantic zoom labels counties with nested bar plots at sufficient zoom.

Weaver, C. (2004, 10-12 Oct. 2004). Building Highly-Coordinated Visualizations in Improvise. IEEE Symposium on Information Visualization,
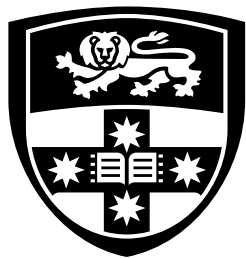
# Spatialisation

- Mapping multidimensional space to a lower dimensional space

# Summary

- Exploratory Visual Data Analysis : try to understand the data through visualization of various statistical data.

- Depending on what sort of analysis you would like to carry out, you should choose appropriate visualization.

- Statistics that can be used for visualization
  - Various statistical features can be mapped to visual attributes to assist EDA processes.

- Other techniques used in EDA
  - Various multi-dimensional scaling methods can be used to place multi-dimensional data point on the 2D/3D space to study the complex data.