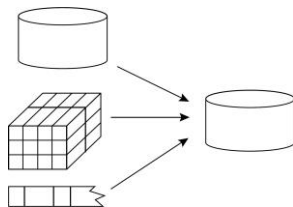**Queen Mary**
University of London

# IOT607U Data Mining

## Week 3: Data exploration and visualization
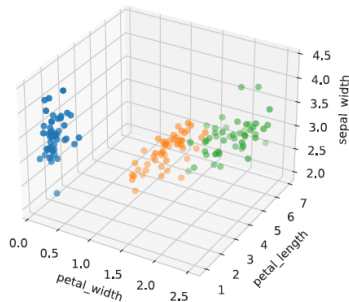
Dr Lin Wang

School of EECS, Queen Mary University of London

- Attributes and Objects
- Characteristics of Data
- **Data Representations**
- Basic Statistical Descriptions of Data
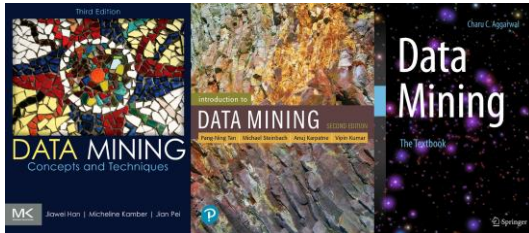- Similarity and Distance

# This week's contents

1. Data Exploration

2. Data Summarization

3. Data Visualization

# Reading

- **Chapter 2.3 of J. Han, M. Kamber, J. Pei, "Data Mining: Concepts and Techniques", 3rd edition, Elsevier/Morgan Kaufmann, 2012**
- Chapter "Data Exploration" of Section 2.3 of P.-N. Tan, M. Steinbach, A. Karpatne, V. Kumar, "Introduction to Data Mining", 2nd edition, Pearson, 2019

# Data Exploration

- Data exploration refers to a preliminary investigation of the data
- This investigation typically has the following goals:
  - Revealing the need for **data pre-processing**
  - Answering simple questions about the data
  - Identifying applicable data analysis techniques
- There are two main methods for data exploration: **data summarisation** and **data visualisation**

## Data Summarization

Summarization techniques

## Data Summarization

- Descriptive statistics
- Central tendency: mean, median, mode
- Dispersion: range, inter-quartile range, variance

# Data Visualization

- **Exploratory data analysis:** what does my data look like?

- **Error detection:** anything weird in the data before you try to use it?

- **Communication:** can you explain the data? Or convince your audience about your point?
    - **e.g. data journalism:** reporting a story with charts.

- Data **collection and cleansing**.

- **Research** on the data, e.g. get descriptive statistics.

- Produce **visualization that conveys those statistics**.

- Assess the **effectiveness of the visualization**.
  - Does it show what we intend to show?
  - Is it easy to interpret for the target reader?
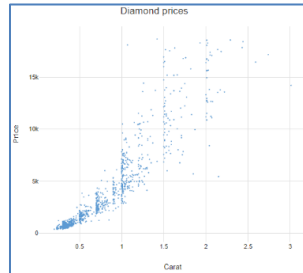  - Could it mislead readers?

- **Scatter plot:** exploring data
- **Time series:** time
- **Bar chart:** comparisons
- **Histogram:** binned frequency
- **Density:** continuous functions
- **Box plot:** quantiles
- **QQ plot:** comparing distributions
- **Other:** displaying multivariate
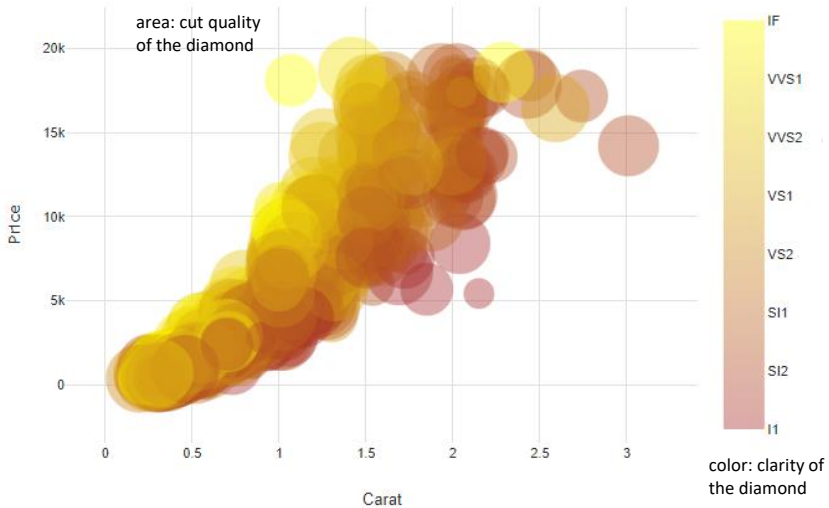
# Data Visualization

## Scatter Plot

- Using Cartesian coordinates to **display values for two variables** for a set of data

- The data is displayed as a collection of points, each having the value of one variable determining the position on the horizontal axis and the value of the other variable determining the position on the vertical axis.
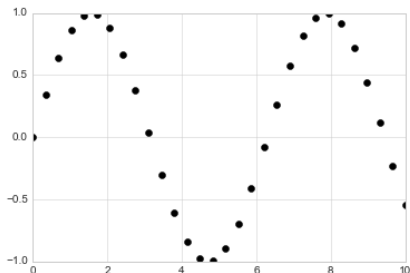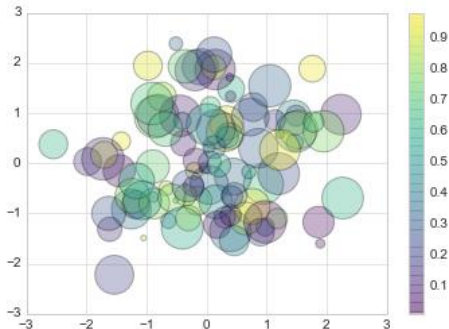
Scatter Plot vs Bubble Plot

BUBBLE PLOT

```python
x = np.linspace(0, 10, 30)
y = np.sin(x)
plt.plot(x, y, 'o', color='black')
```

```
rng = np.random.RandomState(0)
x = rng.randn(100)
y = rng.randn(100)
colors = rng.rand(100)
sizes = 1000 * rng.rand(100)
plt.scatter(x, y, c=colors, s=sizes,
        alpha=0.3, cmap='viridis')
plt.colorbar()
```

**aka bubble plot**

https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.scatter.html
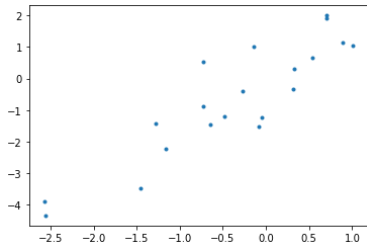
- **Line of best fit:**
  straight line that is the best approximation of the given
  set of data

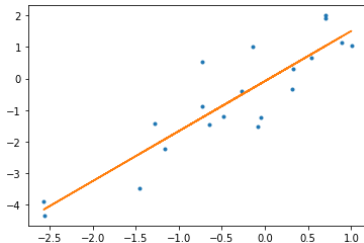- In a scatter plot, what is the linear tendency of my data?

  ```
  numpy.polyfit(x, y, deg)
  ```

```python
x = np.random.randn(20)
y = 2*x + np.random.randn(20)
plt.plot(x,y,'.')
```
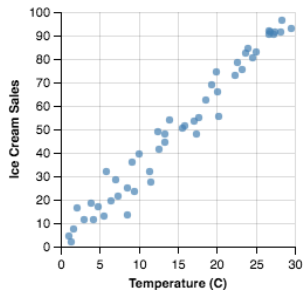
```python
m,b = np.polyfit(x,y,1)
yy = m*x + b
plt.plot(x,yy)
```
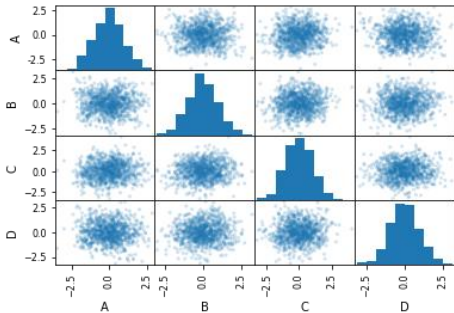
# Scatter Plot: Discussion

- Useful for **exploration of variables** and their **relationships**

- Shows the raw data: **outliers will show up**

- Link to descriptive statistics: provides an immediate feeling of

    - **Range**
    - **Dispersion**
    - **Correlation**

## Scatter Plot Matrix

```python
import pandas as pd
df = pd.DataFrame(np.random.randn(1000, 4),
            columns=['A','B','C','D'])
pd.plotting.scatter_matrix(df, alpha=0.2)
```



|   | A | B | C | D |
|---|---|---|---|---|
| 0 | 1.193096 | -0.946914 | -0.936265 | -1.848978 |
| 1 | -1.813538 | -1.042910 | -2.177161 | -0.230932 |
| 2 | 0.569173 | 0.044870 | 0.271126 | -0.372552 |
| 3 | 0.575154 | -0.738166 | -0.546672 | 0.336331 |
| 4 | 2.454912 | -1.546656 | -0.928218 | -1.181470 |

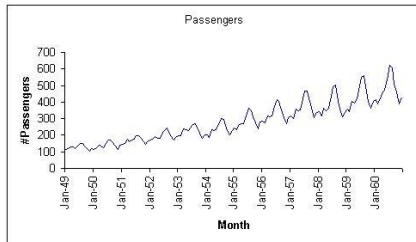# Data Visualization

Time Series

## Time Series: Definition

- **Sequence of data points**, measured typically at **successive points in time** spaced at uniform time intervals

- Time series are used in statistics, signal processing, pattern recognition, econometrics, mathematical finance, weather forecasting, earthquake prediction, electroencephalography, control engineering, astronomy, communications engineering, and largely in any domain of applied science and engineering which involves temporal measurements.

```
plt.plot(x, y)
```



A time series graph of the population of the United States from the years 1900 to 2000



Monthly airline bookings of a flight company

- Trends and patterns over time
- Immediate feeling of
    - Noise vs. signal
    - Trends
    - Anomalies



UK daily cases between 2 April 2020 and 15 July 2021
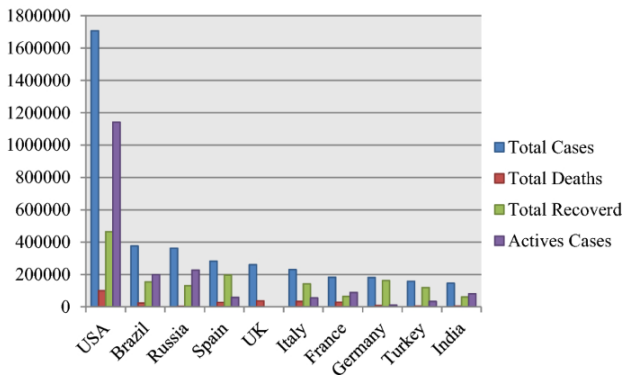
# Data Visualization

Bar Chart

- **Rectangular bars** with lengths proportional to the values they represent.

- Used to show **comparisons among categories**.

- Some bar graphs present bars clustered in groups of more than one (**grouped bar graphs**), and others show the bars divided into subparts to show cumulative effect (**stacked bar graphs**).
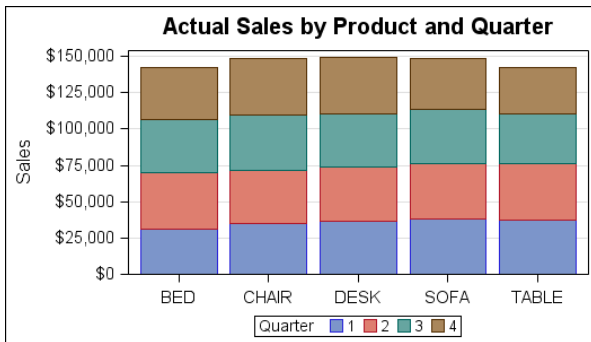

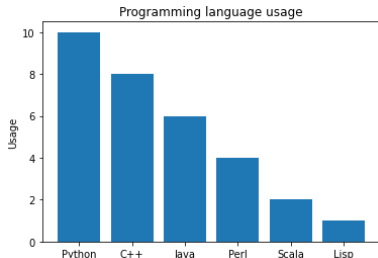Programming language usage

Example of a grouped (clustered) bar chart

Example of a stacked bar chart



Actual Sales by Product and Quarter

# Bar Chart in Python

```python
import matplotlib.pyplot as plt
import numpy as np
objects = ('Python', 'C++', 'Java', 'Perl', 'Scala', 'Lisp')
x_pos = np.arange(len(objects))
usage = [10, 8, 6, 4, 2, 1]
plt.bar(x_pos, usage)
plt.xticks(x_pos, objects)
plt.ylabel('Usage')
plt.title('Programming language usage')
```
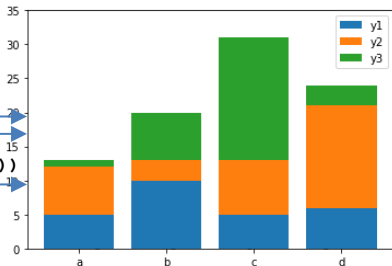


https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.bar.html
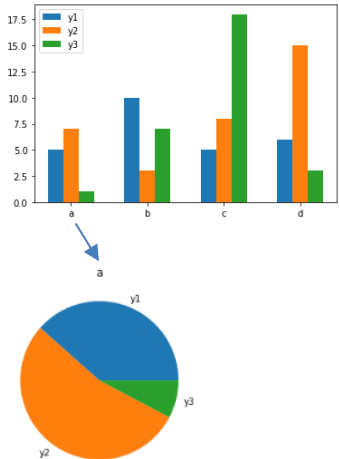
```python
import matplotlib.pyplot as plt
import numpy as np
x = ['a', 'b', 'c', 'd']
y1 = [5, 10, 5, 6]
y2 = [7, 3, 8, 15]
y3 = [1, 7, 18, 3]

w = 0.2
x_pos = np.array(range(len(x)))
plt.bar(x_pos-0.2, y1, width=w)
plt.bar(x_pos, y2, width=w)
plt.bar(x_pos+0.2, y3, width=w)
plt.xticks(x_pos,x)
plt.legend(['y1','y2','y3'])
```

```python
import matplotlib.pyplot as plt
import numpy as np
x = ['a', 'b', 'c', 'd']
y1 = [5, 10, 5, 6]
y2 = [7, 3, 8, 15]
y3 = [1, 7, 18, 3]
plt.bar(x, y1)
plt.bar(x, y2, bottom=y1)
z = list(np.array(y1) + np.array(y2))
plt.bar(x, y2, bottom=z)
```

## PIE CHART IN PYTHON

```python
import matplotlib.pyplot as plt
import numpy as np
x = ['a', 'b', 'c', 'd']
y1 = [5, 10, 5, 6]
y2 = [7, 3, 8, 15]
y3 = [1, 7, 18, 3]

za= [y1[0], y2[0], y3[0]]
plt.pie(za, labels=['y1','y2','y3'])
plt.title(x[0])
```



https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.pie.html

# Data Visualization

Histogram

- Representation of **tabulated frequencies, shown as adjacent rectangles**, erected over discrete intervals (bins), with an area proportional to the frequency of the observations in the interval.

- The **height of a rectangle is also equal to the frequency density of the interval**, i.e., the frequency divided by the width of the interval. The total area of the histogram is equal to the number of data.



**Night price distribution of Airbnb appartements**

# Histogram in Python

```python
import matplotlib.pyplot as plt
x = [21, 22, 23, 4, 5, 6, 77, 8, 9, 10, 31, 32, 3
3, 34, 35, 36, 37, 18, 49, 50, 100]
num_bins = 5

[count, bin_boundary, patch] = plt.hist(x,
 bins=num_bins, facecolor='blue')

#bin_boundary: [4.0, 23.2, 42.4, 61.6, 80.8, 100.0]
#count: [10.0, 7.0, 2.0, 1.0, 1.0]


# can also use np.histogram + bar
[count2, bins2] = np.histogram(x, bins=5)
plt.bar((bins2[0:-
1]+bins2[1:])/2,count2,width=20)
```

- Can be hard to read due to the bin size choice

- No "best" number of bins: different bin sizes can reveal different features of the data

- Square root choice (used by Excel): $k = \sqrt{n}$
- Sturges (normal data): $k = \lceil \log_2 n + 1 \rceil$



$n = 5000$ normally distributed data points

# Data Visualization

Density Plot

## Density: Definition

- Function that describes the **relative likelihood for this variable to take a given value**

- The probability of the random variable falling within a particular range of values is given by the integral of this variable's density over that range—that is, it is given by the **area under the density function** but above the horizontal axis and between the lowest and greatest values of the range. The probability density function is **non-negative** everywhere, and its **integral over the entire space is equal to one**.

$$p(x) \geq 0$$

$$p(a < x < b) = \int_a^b p(x)dx$$

$$\int_{-\infty}^{\infty} p(x)dx = 1$$

$$p(x) \geq 0$$
$$p(a < x < b) = \int_a^b p(x)dx$$
$$\int_{-\infty}^{\infty} p(x)dx = 1$$

## **plt.hist()**

```python
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd

flights = pd.read_csv("http://www.eecs.qmul.ac.uk/~linwang/download/
ecs764/formatted_flights.csv")
JetBlue = flights[flights['name']=='JetBlue Airways']
delay = JetBlue['arr_delay']

hist, bins, path = plt.hist(delay, bins=18, density=True)
bin_centers = (bins[1:]+bins[:-1])*0.5
plt.plot(bin_centers, hist, linewidth=5)
plt.title('JetBlue Airways')
```



JetBlue Airways

## seaborn.distplot()

```python
import seaborn as sns
sns.distplot(delay, hist=True, kde=True, bins=18,
             color = 'darkblue',
             hist_kws={'edgecolor':'black'},
             kde_kws={'linewidth': 5})
plt.title('JetBlue Airways')
```

kernal density estimation (kde)



https://seaborn.pydata.org/generated/seaborn.distplot.html

# Density: Kernal Density Estimation

| Sample | 1 | 2 | 3 | 4 | 5 | 6 |
|--------|------|------|------|-----|-----|-----|
| Value | -2.1 | -1.3 | -0.4 | 1.9 | 5.1 | 6.2 |



Comparison of the histogram (left) and kernel density estimate (right) constructed using the same data. The six individual kernels are the red dashed curves, the kernel density estimate the blue curves. The data points are the rug plot on the horizontal axis.

```
bw_list = [0.1, 0.2, 0.5, 1, 5]
for b in bw_list:
    sns.distplot(delay, bins = 18, hist = False, kde = True,
                 kde_kws={'bw': b})
plt.legend(bw_list)
plt.title('Density plot with varying bandwidth')
```



Density plot with varying bandwidth

# Data Visualization

## Box Plot

## BOX PLOT: DEFINITION

- Graphically depicts **groups of numerical data through their quartiles/percentiles**.

- Box plots display five number summary: "minimum", first quartile (Q1), median, third quartile (Q3), and "maximum".
    - Lower whisker: Minimum - Q1
    - Upper whisker: Q3 - Maximum
    - Small/large values (**outliers**) may be plotted as individual points.

# Box Plot



Symmetric distribution    Skewed distribution

## Box Plot in Python

## `plt.boxplot()`

```python
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
flights = pd.read_csv("http://www.eecs.qmul.ac.uk/~linwang/download/ecs764/formatted_flights.csv")

JetBlue = flights[flights['name']=='JetBlue Airways']['arr_delay']
Delta = flights[flights['name']=='Delta Air Lines Inc.']['arr_delay']
United = flights[flights['name']=='United Air Lines Inc.']['arr_delay']
American = flights[flights['name']=='American Airlines Inc.']['arr_delay']
ExpressJet = flights[flights['name']=='ExpressJet Airlines Inc.']['arr_delay']

delay = [United, JetBlue, ExpressJet, Delta, American]
plt.boxplot(delay)
```

**`plt.boxplot()`**

## **seaborn.violinplot ()**

```
import seaborn as sns
sns.violinplot(data=delay)
```



https://seaborn.pydata.org/generated/seaborn.violinplot.html

## Box Plot: Discussion

- Quick way of examining one or more sets of data visually (exploratory)

- More primitive than histogram or density

- No need for choice of bins and bandwidth

- Compact

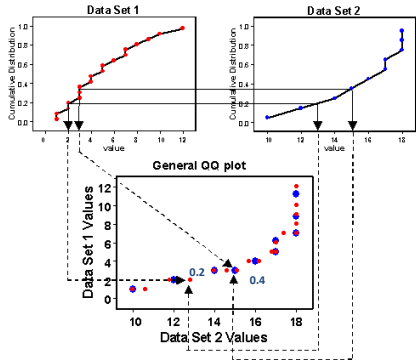- Histogram better for one group while boxplot better for multiple groups

# Data Visualization

QQ Plot

- **Comparing** two probability **distributions** by plotting their **quantiles against each other**

- Method:

  1) The set of intervals for the quantiles is chosen (typically $k/(n + 1)$)

  2) A point (x, y) on the plot corresponds to one of the quantiles of the second distribution (y-coordinate) plotted against the same quantile of the first distribution (x-coordinate).

# QQ Plot in Python

```python
import matplotlib.pyplot as plt
import numpy as np

def standardize(data):
  z = (data - np.mean(data))/np.std(data)
  return z

def get_quantiles(data):
    quantiles = []
    for q in np.arange(0, 1.001, 0.001):
        quantiles.append(np.quantile(data, q))
    return quantiles

def pyqqplot(data1, data2):
    data1 = standardize(data1)
    data2 = standardize(data2)

    q1 = np.array(get_quantiles(data1))
    q2 = np.array(get_quantiles(data2))
    plt.scatter(q1, q2 )

    minim = min(data1.min(), data2.min())
    maxim = max(data1.max(), data2.max())
    plt.plot([minim, maxim], [minim, maxim], 'r-')
```

```python
x1 = np.random.normal(0, 1, 10000 )
x2 = np.random.normal(0, 1, 10000 )
pyqqplot(x1, x2)
```

# QQ Plot: Example

# Data Visualization

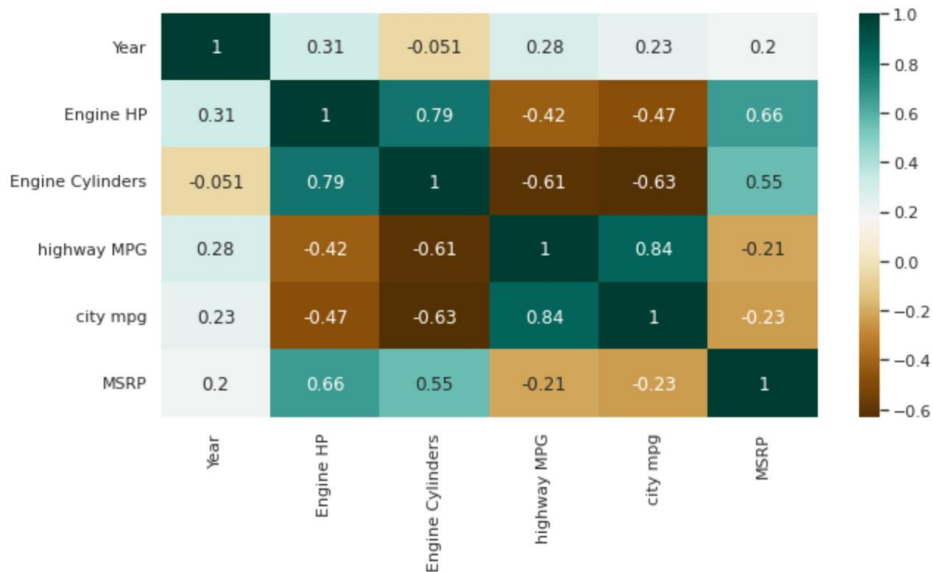## Displaying Multivariate:
Heatmap, Surface, Contour, Quiver

# Heatmap

- A heat map uses a colour to represent the value $f(x, y)$ of each point $(x, y)$ in a scalar field $f : \mathbb{R}^2 \to \mathbb{R}$

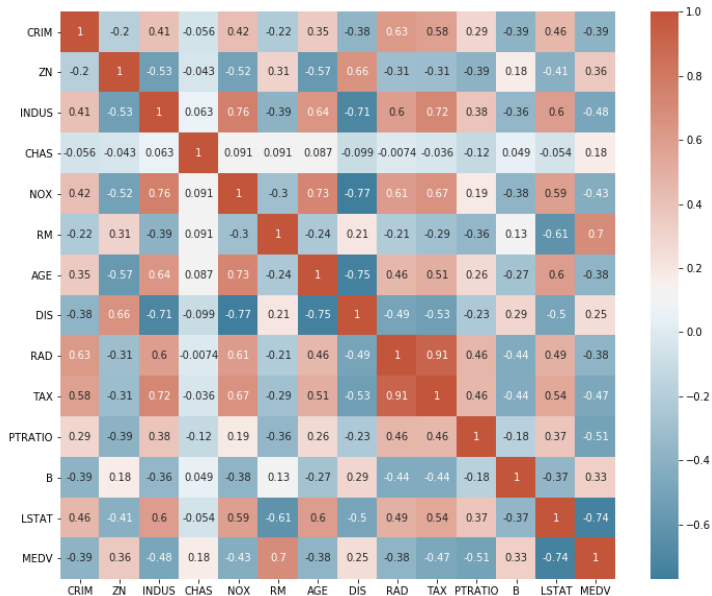# Heatmap Example

# Heatmap Example

## Surface Plot

- A surface plot displace the graph of a bivariate function
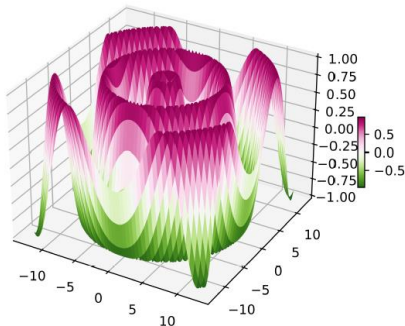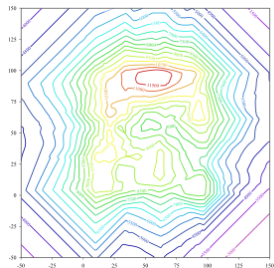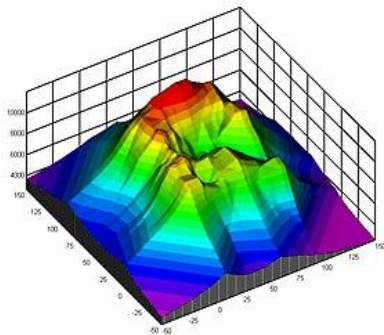- A surface is plotted to fit a set of data (X, Y, Z), where Z if obtained by the function to be plotted Z = f(X, Y).
- Optionally, the plotted values can be color-coded.



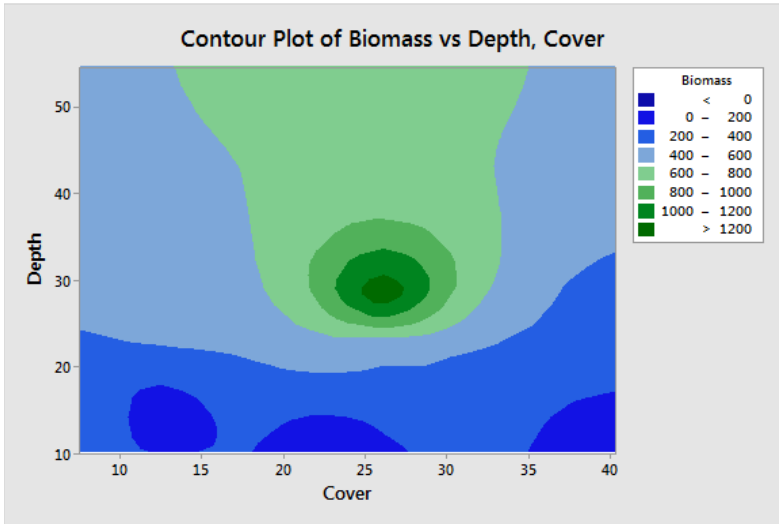https://matplotlib.org/stable/gallery/mplot3d/surface3d.html

# Contour Plot

- A contour plot represents a 3-D surface by plotting lines that connect points with common z-values along a slice.
- For example, you can use a contour plot to visualize the height of a surface in two dimensions.
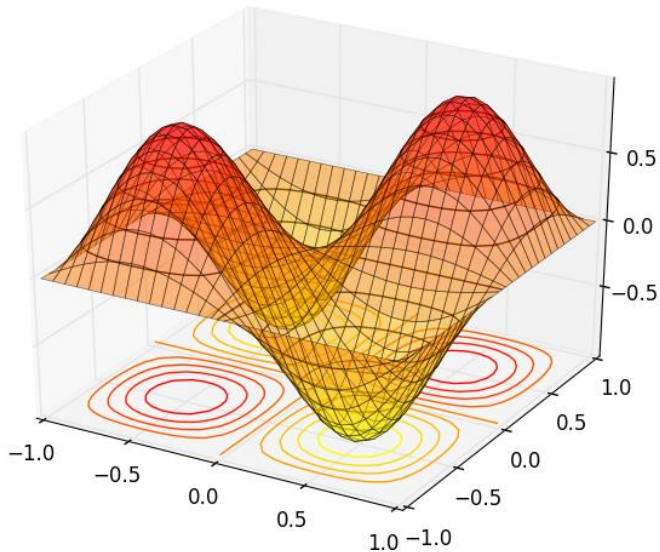


https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.contour.html

## Contour Plot Example

A biologist studies the effect of stream depth and canopy cover on fish biomass.



Contour Plot of Biomass vs Depth, Cover
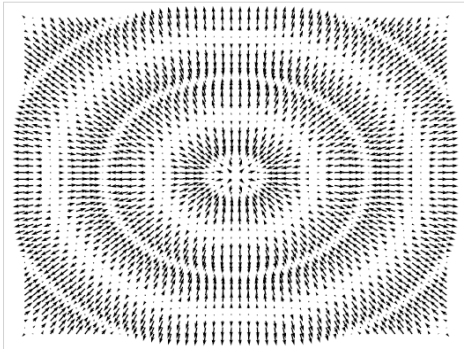
## Vector Field (Quiver)

- In some data, a characteristic may have both magnitude and a direction associated with it (e.g. flow of a substance)
- Vector filed plots (or quiver plots) display both direction and magnitude
- Vector fields can model velocity, magnetic force, fluid motion, and gradients.
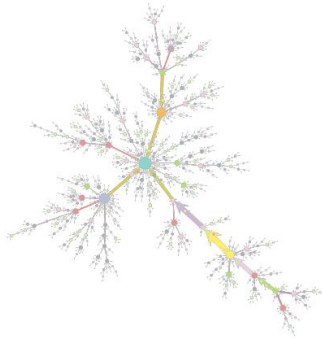
Optical flow in computer vision
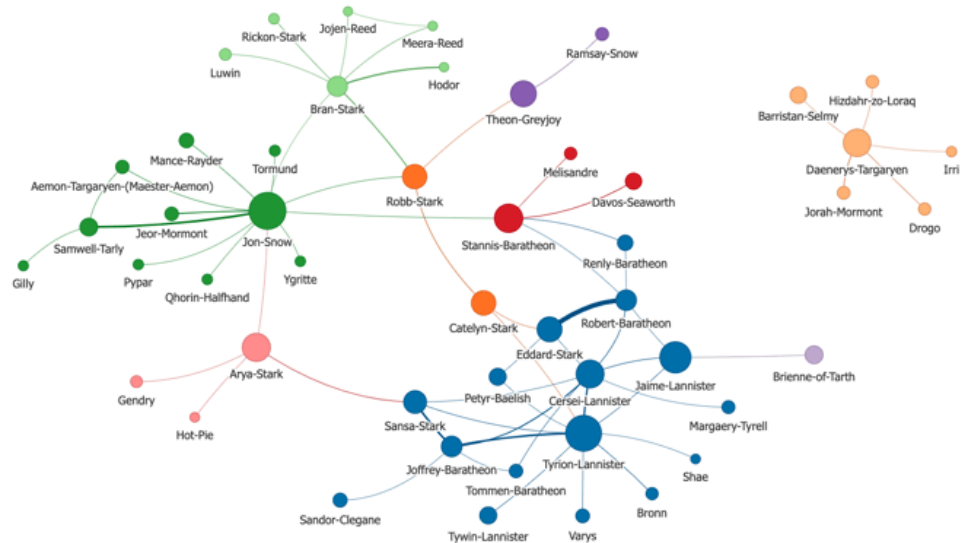
- Graph drawing is concerned with obtaining visual representations of graphs based on nodes and links
- There are many different graph layout algorithms that attempt to maximise different **quality measures**



https://plotly.com/python/network-graphs/
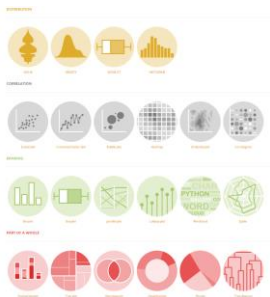
# GRAPH VISUALIZATION



*Game of Thrones characters with more than 60 interactions*

## MORE ON PYTHON GRAPHS

- https://python-graph-gallery.com/

THE PYTHON
GRAPH GALLERY

# Data Visualization

## Visualization principles

- **Scatter plot:** exploring data
- **Time series:** time
- **Bar chart:** comparisons
- **Histogram:** binned frequency
- **Density:** continuous functions
- **Box plot:** quantiles
- **QQ plot:** comparing distributions
- **Others**: displaying multivariate

- Data visualisation is an art

- Many plot types for different purposes: knowing the right tools saves a lot of time

- Important tool for data exploration, especially extracting basic relationships

- **Apprehension**: does the visualisation maximize the understanding of the relationships in the data?
- **Clarity**: are the most important relationships the most visually prominent?
- **Consistency**: are the visual elements consistent with their use in previous visualisations?
- **Efficiency**: are the visual elements economically used and easy to interpret?
- **Necessity**: is the visualisation more useful than alternative ways to represent the data?
- **Truthfulness**: are the visual elements accurate and unambiguous?

**Questions?**
**also please use the forum on QM+**