# Analytics

Chapter 3

# Analytics

Analytics or data analytics is the computer science discipline of analysing and interpreting data.

Most computer science applications amass a lot of data over time.

For example a medical system could store millions of instances of data in seconds.

The goal of Analytics is to process and transform the data, present the data in a meaningful way and finally interpret the data.

# Data Analytics

Terms that are similar:

- Data Analytics
- Data Analysis
- Data Mining
- Data Science

# Data

Data is the term used to typically represent raw data, from sources such as a CSV file, a List or data from a database.

It can be in a form that is of no use to a computer scientist.

Pre-processing is usually the next step.

# Pre processing

One of the first things that computer scientists exam are outliers - for example someone could have accidently typed an extra zero when entering an age resulting in 200 and not 20.

A simple method is to remove that value.

Another pre processing method is to sort the data.

This is useful to get the max and min values or getting ready to count frequency.

# Information

Information usually refers to gathering or presenting the data for the user.

This typically requires presenting the data visually such as a graph or a figure.

# Analysis

Even with information the fine details are not always obvious.

Information is a great aid for presenting your findings, but algorithms can be more accurate when interpreting the data.

Algorithms such as frequency, mean, median, minimum, maximum and mode of data can really add clarity and help interpret the data.

# Data - Python Lists

Lists in Python are fundamental in analytics. Even if our data is stored in a 2D list (like a database), most algorithms such as frequency, mean, median and mode can be calculated on a 1D list.

# Text Files

Text files are one of the most versatile files for computer scientists.

Example of creating a text file and entering data.

```
file = open("myFirstTextFile.txt","w")

file.write("Test text")

file.close()
```

```
file = open("myFirstTextFile.txt","w")
```

The first line allows access to the file.  If the file does not exist, Python automatically creates it for you.

`File` (on the first line)  is just a variable name.  The `open`  method takes two arguments:

The first argument is a string that is the file name.  You also need the file extension.  For a text file, this is .txt.

The second argument is what you are doing with the file.  There are three options:

"w" - write.  This will overwrite any previous contents.

"a" - append.  This adds the text to the end of the current contents.

"r" - read.  This read the contents of the file as a string.

```
file.write("Test text")
```

The second line of code, we pass a string in as an argument to be written to the file.

```
file.close()
```

The third line of code closes the connection to the file.

Python must request to open a file, it must also say when it is complete and release the file.  If you don't then you may get an error when trying to use this file again.

# Reading the file

To read in the same file we just created, we could do the following:

```
file = open("myFirstTextFile.txt", "r")

dataIn = file.read()

file.close()

print(dataIn)
```

# CSV files

While text files have their uses, they are not the most suitable for storing data, especially large amounts of data.

A CSV file is essentially like a .txt file.  It uses the same Unicode/ASCII structure.

CSV stands for "comma separated values"

Computer Scientists like this format for the following reasons:

- Quick to read and write to in Python
- Excel has a row limit of 1,048,576 rows
- The file can be used even if you do not have excel on your computer.

# CSV Structure

The structure of a CSV file is relatively simple but the format can cause some issues:

- Commas separate elements on a row.  These are called delimiters.
- A new line special character "\n" represents the end of a row and the start of a new row.

|   | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | 1 | 2 | 3 | 4 | 5 |
| 2 | 6 | 7 | 8 | 9 | 10 |

1, 2, 3, 4, 5  \n
6, 7, 8, 9, 10 \n

# Reading and Writing to CSV file

Reading and Writing to a csv file can be done the same way as for a text file. You just need to change the file extension.

```python
file = open("myFirstCSVFile.csv", "w")
file.write("1,2,3,4,5")
file.close()
```

```python
file = open("myFirstCSVFile.csv", "r")
dataIn = file.read()
file.close()
print(dataIn)
```

The output is:

```
1,2,3,4,5
```

The code in the previous slide prints a string for us.  This is not really useful if we want to find a total or an average etc.

We know that every element is separated by a comma in a csv file.

We can use this to our advantage and write a loop to go through the string and take out elements with a comma in between them.

This sounds complicated but Python has a method for doing this - it can separate a string into a list based on any specified delimiter (in our case it's a comma)

```
myList = dataIn.split(",")
print(myList)
```

The output is:

```
['1','2','3','4','5']
```

The .split() method takes the delimiter as an argument.

Notice that all the elements in the list are strings.

We can convert each element to an int so that we can do calculations on them. Python has an easy way to do this.

```
myList = dataIn.split(",")

myList = [int(element) for element in myList]

print(myList)
```

The output is (now all integers):

```
[1,2,3,4,5]
```

# Pre-Processing and Transforming Data

Pre-Processing allows to apply techniques to ensure the data is of good quality.

Typical errors in data that require pre-processing are:

- Missing data
- Data out of range
- Non validated data
- Data structure incorrect

# Sorting

Sorting is useful for many reasons.  If you sort, you will know the largest and the smallest are the first and the last value of the list.

```
myList = [1,19,27,8,5,9]

myList.sort()

print("Sorted Ascending:" myList)

myList.sort(reverse = True)

print("Sorted Descending:" myList)
```

The output is:
```
Sorted Ascending   [1, 5, 8, 9, 19, 27]
Sorted Descending  [27, 19, 9, 8, 5, 1]
```

# Removing Missing Data

Sometimes data in a list could be identified as an error and may need to be removed.

For example if a list should only contain positive numbers and a negative number is present.

A pre-processing step could be to remove all negative numbers

```
myList = [1,-19,27,8,-5,9]
for item in myList:
        if item<0:
                myList.remove(    )
print(myList)
```

The output is:
```
[1, 27, 8, 9]
```

# Imputing Missing Data

Imputing data is an option instead of removing missing data.

Sometimes removing data is not an option, especially in a 2D list as there is other data associated with it.

| Name | Age | Claims |
|------|-----|--------|
| Joe Bloggs | 23 | 1 |
| Mary Murphy | 32 | -19 |
| Claire Whelan | 27 | 27 |
| Mike Fahey | 46 | 8 |
| Gillian Marks | 21 | -5 |
| Steph Curry | 31 | 9 |

| Name | Age | Claims |
|------|-----|--------|
| Joe Bloggs | 23 | 1 |
| Mary Murphy | 32 | -19 |
| Claire Whelan | 27 | 27 |
| Mike Fahey | 46 | 8 |
| Gillian Marks | 21 | -5 |
| Steph Curry | 31 | 9 |

[1, -19, 27, 8, -5, 9]

| Name | Age | Claims |
|------|-----|--------|
| Joe Bloggs | 23 | 1 |
| Mary Murphy | 32 | 27 |
| Claire Whelan | 27 | 8 |
| Mike Fahey | 46 | 9 |
| Gillian Marks | 21 | |
| Steph Curry | 31 | |

# Imputing Missing Data

We can go through a list and replace all values with a number.

Imputing is an algorithm that finds the average of all non missing values and then replaces the missing values with the average of the group.

There are two parts:

First find the average of all non missing values

Second replace all error values with the calculated average

**Part A**, find the average of all non-missing values:

```
totalNonMissing = 0
countNonMissing = 0
myList = [1, -19, 27, 8, -5, 9]
for item in myList:
    if item > 0:
        totalNonMissing += item
        countNonMissing += 1
averageNonMissing = totalNonMissing / countNonMissing

print(averageNonMissing)
```

The average value of non-error numbers is:

```
11.25
```

**Part B**, replace all error values with the recently calculated `averageNonMissing`:

```
for counter in range(len(myList)):
    if myList[counter] < 0:
        myList[counter] = averageNonMissing
print(myList)
```

The imputed list:

```
[1, 11.25, 28, 8, 11.25, 9]
```

# Cleaning/Pre-processing Data

In the programs we do, the most common cleaning involves one of the following:

- Remove unwanted characters such as a comma in the middle of a number
- Select only part of the data in a column

Examples of each of these are here:

Remove unwanted characters

Select part of the data in a column