JUST **IT**

# HANDBOOK

**Data Technician bootcamp**

(Data Search, Loading, Manipulation, Visualisation, Analysis)

- Excel
- Tableau
- Power BI
- SQL
- Programming language R

Student:

Alexei Kosyhin

10 Feb 2023

# CONTENT

Page

# Concepts

In this section I will outline everything what we learned about data.

- Data can be defined as an elementary value or the collection of values;
Group Items which have subordinate data items are called Group item;
Record can be defined as the collection of various data items;
File can be defined as various records of one type of entity;
Attribute and Entity represents the class of certain objects. It contains various attributes.

- There are many different sources of data. There are different types of files: MS Word Documents, Spreadsheets, Emails, PDFs, HTML, plaintext files, CSV, JSON,
- Defining Big Data – 4 bigs: volume, velocity, variety, veracity

- Open and Private Data

- Data Analysis Lifecycle: Gathering, Preparing, Choosing a model, Analysing, Presenting results, Making Decisions,

- Data Mapping Techniques: Define datatype and size; Map; Transform; Validate and Test; Deploy; Maintain and update.

- Data Structures: Relational Data  Tables- fields, rows, values;
                                    Python – Strings, lists, tuples, sets, dictionaries, classes, functions …
- Structured and unstructured data

- Evolution to Big Data: Database servers -> Distributed data systems -> Onsite and Cloud computing Solutions

- Basic Data Management Technologies: Flat file database, Relational Database.

- Types of Data Analysis: Scalable Technologies, Business Intelligence.

- Conception of  different types lists

- Purpose of Big Data Analysis: for Descriptive Analysis, Predictive Analysis, Prescriptive Analysis

- Types of Data: categorical data and numerical data. Qualitive categorical data can be nominal data and Ordinal data. Quantitative data can be Interval and Ration

-

# Legislation relating to data security

**Current Legislation:** There are laws designed to protect users and their data from attack and misuse.

- **Computer Misuse Act 1990.** This includes planting a virus which is intended to cause damage, altering data, slowdown computer operation, frequent computer crashes. This particularly includes spreading malicious software like viruses, worms, Trojans, ..

- **Police and Justice Act 2006** (Computer Misuse) this is continuation of the Computer Misuse Act of 1990. This is particularly about hacking – unauthorised access to somebody computer.

- **The Copyright (Computer Programs) Regulation 1992**. From Wikipedia Software copyright is the application of copyright in law to machine-readable software.

- **Data Protection Acts (1998, 2018) and GDPR**. These acts are essentially abiut two points: The ethical use of personal data, and keeping individual's personal data secure. Everyone responsible for using personal data has to follow strict rules called 'data protection principles'. They must make sure the information is: used fairly, lawfully and transparently.

- **Consumer Right Act 2015.** The aim of the 2015 Consumer Rights Act was to aid both consumers and retailers in understanding their rights and responsibilities, and thus to reduce and simplify disputes. It also seeks to encourage business based on fair practices and access to information.

# Introduction to the problem – education via wealth and population

Even though the purpose of this handbook is to demonstrate pure technical skills, I will try to make my data manipulation, visualisation and analysis meaningful and interesting.

I decided to explore a topic about education in the world and how level of education (number of the best universities) depends on population and income. I found a list of 1000 best universities from Wikipedia even though total number of universities in the world is about 25,000. Then I found data of population and revenue by country from Wikipedia. Then I will try to make analysis and make conclusion.

I decided to use data from Wikipedia because Wikepedia is open source and doesn't require special permission for using data, Wikipedia is the most reliable and affordable source of good quality information.

Here are URLs for data from Wikipedia

1) [List of top 1000 universities in the world - Wikipedia](#)
2) [List of countries by GDP (PPP) - Wikipedia](#)
3) [List of countries by population (United Nations) - Wikipedia](#)

In order to provide analysis I will:

a) download the data using Power BI;

b) clean the data (delete empty rows, to replace wrong types e t c) using Power BI and Excel

c) group by the list of universities using Power BI

d) visualize using Excel, Tableau, Power BI

e) merge/ join using Power BI, Microsoft SQL

f) draw some conclusion

Data manipulation – creating tables.





This pic. Shows original data in Wikipedia. This table consists 1000 rows.
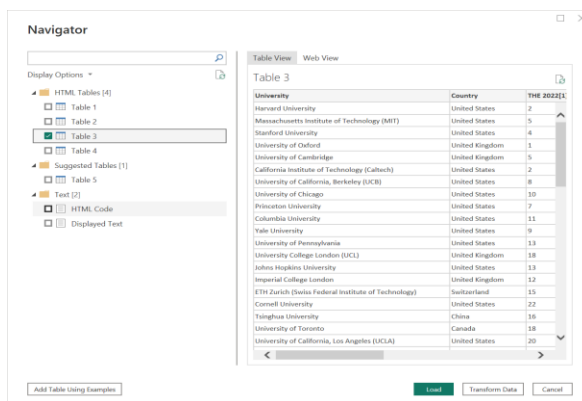
I download this table using Power BI





When I tried to load the table an error appeared. So, instead of load and used transform data. Then I copied the entire table and pasted into excel file.





Once I load the table into excell, I simplified the data. I deleted some columns and left only 2. The table was big, has 1000 rows. I cleaned the table manually. Then I imported the data back to Power BI. The software didn't show an error.

==============================================================

In the similar way I created in Power BI 2 tables – population by country and revenue by country. Now I have 3 tables in Power BI file. The data is cleaned enough for manipulation and analusis.



In order to compare some metrics such as education /number of universities, revenue, population by country we need to transform some tables. For the table of list of universities we need to use "group by" function in Power BI to find number universities in each country.

Using JOIN/ MERGE/ COMBINE in Power BI. For analysis and visualisation it would be nice to merge these 3 tables into one, key index will be "country".



Here I used Power BI Query functions To combine 3 tables into one.

This is the main table from which we will do visualisation and analysis.

However a manager was curious about the data for each continent. For this purpose I will group the table by key "Continents".

**VISUALIZATION using Power BI tools**





Number Universities by Country



Revenue by Country



Population by Country

# Visualisation using Power Bi

(Displaying numbers of universities, population and revenue by continent)

## Number Universities by Continent

44 (4.45%)
259 (26.21%)
380 (38.46%)
294 (29.76%)

**Continent**
- Europe
- Asia
- Americas
- Oceania
- Africa

## Revenue by Continent

2M
26M (1.74%)
(22.32%)
57M (48.71%)
30M (25.97%)

**Continent**
- Asia
- Americas
- Europe
- Africa
- Oceania

## Population by Continent

1bn 0bn
(11.72%) (0.5%)
1bn (15.35%)
4bn (69.61%)

**Continent**
- Asia
- Americas
- Europe
- Africa
- Oceania

# Visualization using TABLEAU.

# Universities



Universities

1.0 ▮▮▮▮▮▮▮ 154.0

# Revenue



Revenue

20,600  23M

# Population



Population

428,963 ▮▮▮▮▮ 1B

# Data Manipulation in R Studio





```
1   library(readr)
2   uni <- read_csv("uni.csv")
3   View(uni)
4   install.packages("tidyverse")
5   sum_uni <- aggregate(University ~ Country, uni, FUN = length)
6   View(sum_uni)
7   write.csv(sum_uni, 'sum_uni.csv')
8
9   library(readr)
10  pop <- read_csv("pop.csv")
11  View(pop)
12
13  join1 <- merge(sum_uni, pop, by = "Country")
14  View(join1)
15
16  library(readr)
17  rev <- read_csv("rev.csv")
18  View(rev)
19
20  join3 <- merge(join1, rev, by = "Country")
21  View(join3)
22
23  write.csv(join3, "join3.csv")
24
```

# Visualization in R Studio

# Analysis/Conclusion

Here I will provide a simple analysis how level education/ number of top universities in each country depends on revenue and population.

From above graphs we can see an obvious correlations 1) between number of universities and revenue, and 2) between number of universities and population.

In general, the more country's avenue or population the more universities in the country.

For example, the USA is the richest country and has the biggest number top universities.

For example, India is the most populated country and has the significant number of universities.

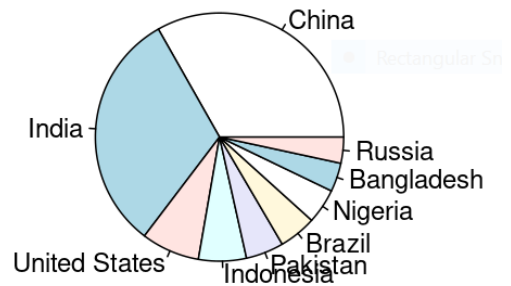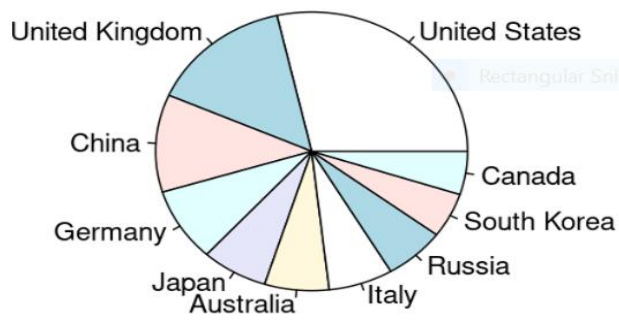There are some exceptions, for example African countries have a very high population, but only few universities.

However other factors should be considered for example an official country language.
I would speculate that the UK has a great number of universities particularly because of English language (everybody wants to learn English in an English speaking country)

# SQL – Big questions

## Task 6

The manager of Northwind Company Limited has marked out certain products for promotional discounts and these would need to be offered as Christmas Discounted Price.

Please create a report containing all products apart from chang, Ikura, Pavlova, and tofu. Ensuring that they relate products with the price ranging from 18 to 32

Please sort the Christmas Discounted Price in descending order



SELECT * FROM [Products]
WHERE (ProductName NOT IN ("Chang", "Ikura", "Pavlova", "Tofu"))
AND (PRICE BETWEEN 18 AND 32);

## Task 7

The manager of Northwind Company Limited wants to embark on promotional activities to focus on their most active customers.

She has therefore asked me to prepare a report that will return to her, the most active customers who have order for goods with quantities that are 50 or above, as well as the quantity ordered



```
Select CU.CustomerID, ODD.Quantity
From Customers as "CU"
Join Orders as "Ord"
On CU.CustomerID = Ord.CustomerID
Join OrderDetails as "ODD"
On Ord.OrderID = ODD.OrderID
Where ODD.Quantity > 49
```

## Task 8

Return to me all orders that have been made at Northwind Company Limited

Please grade them according to the size of quantity ordered e.g. Low value customer to refer to those who order for quantities that 30 or less;  Value customer to refer to those who order for quantities that 50 or less;  Large customer to refer to those who order for quantities that 70 or less;  and finally Premium customers for the others



```
SELECT OrderID, Quantity,
CASE WHEN Quantity  < 30 THEN 'SMALL'
WHEN Quantity < 50 THEN "NORMAL"
WHEN Quantity < 80 THEN "BIG"
ELSE 'EXTRA LARGE'
END AS QuantityText FROM OrderDetails;
```

# Task 9

Return to me the sales of Northwind Company Limited according to their categories.

Please ensure that your report only concentrates on Sales that are in excess of 300

SQL Statement:

```
SELECT ORDERDETAILS.ORDERID,  PRODUCTS.PRODUCTID,
PRODUCTS.PRODUCTNAME, CATEGORIES.CATEGORYNAME,
(ORDERDETAILS.QUANTITY *PRODUCTS.PRICE) AS "SALES"
FROM PRODUCTS
JOIN CATEGORIES, ORDERDETAILS
ON ORDERDETAILS.PRODUCTID = PRODUCTS.PRODUCTID
```

Edit the SQL Statement, and click "Run SQL" to see the result.

**Run SQL »**

Result:

Number of Records: 310

| OrderID | ProductID | ProductName | CategoryName | SALES |
|---------|-----------|-----------------|--------------|-------|
| 10253 | 39 | Chartreuse verte | Beverages | 756 |
| 10255 | 2 | Chang | Beverages | 380 |

SELECT ORDERDETAILS.ORDERID,  PRODUCTS.PRODUCTID,
PRODUCTS.PRODUCTNAME, CATEGORIES.CATEGORYNAME,
(ORDERDETAILS.QUANTITY *PRODUCTS.PRICE) AS "SALES"
FROM PRODUCTS
JOIN CATEGORIES, ORDERDETAILS
ON ORDERDETAILS.PRODUCTID = PRODUCTS.PRODUCTID

AND CATEGORIES.CATEGORYID = PRODUCTS.CATEGORYID
WHERE SALES > 300
ORDER BY CATEGORYNAME

**Question 10**

To specify multiple values (in order words, a shorthand for multiple OR conditions), which of the following Operators is can be used

**a.     IN**
b.     LIKE
c.     BETWEEN
d.     AND


**Question 11**

When an aggregate SQL statement is involved, the SQL Clause which is used instead of a Where Clause is called
a.     WHERE Clause
**b.     Group By Clause**
c.     Case Statement
d.     The LIKE Condition

**Question 12**

The Keyword that eliminates duplicate rows from the results of a SELECT Statement is called
a.     Select *
**b.     DISTINCT**
c.     Avoid Duplicate
d.     The IN Operator


**Question 13**

The Clause which sorts data of a table is called
a.     The Group By Clause
b.     Having Clause
**c.     Order By Clause**
d.     The Where Clause


**Question 14**

We use the ……**LIKE**……… Operator with SELECT to set conditions based on pattern matching
a.     IN
**b.     LIKE**
c.     BETWEEN
d.     NOT