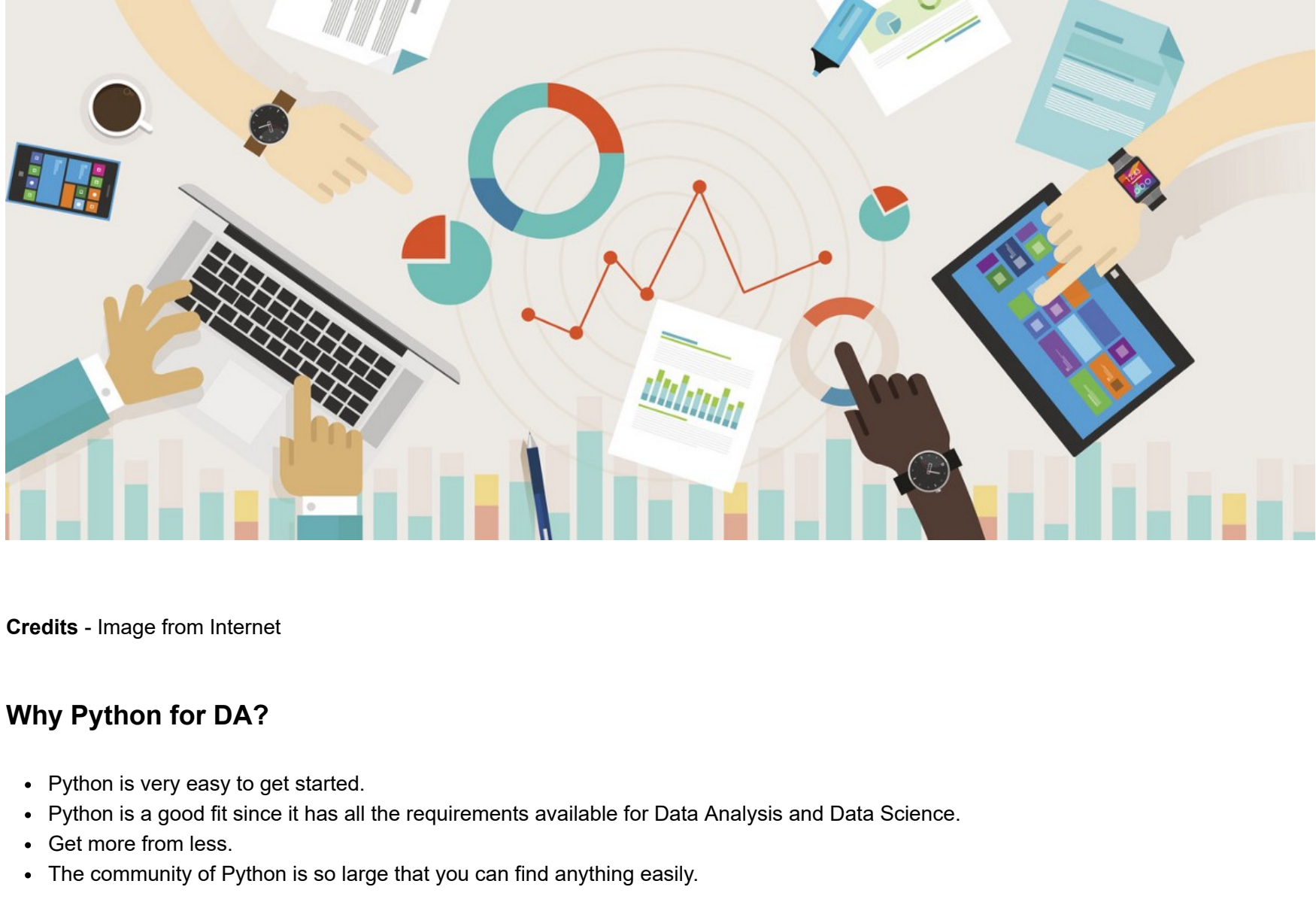


# Intro to Basic understanding of the data using Python



Credits - Image from Internet

## Why Python for DA?

- Python is very easy to get started.
- Python is a good fit since it has all the requirements available for Data Analysis and Data Science.
- Get more from less.
- The community of Python is so large that you can find anything easily.

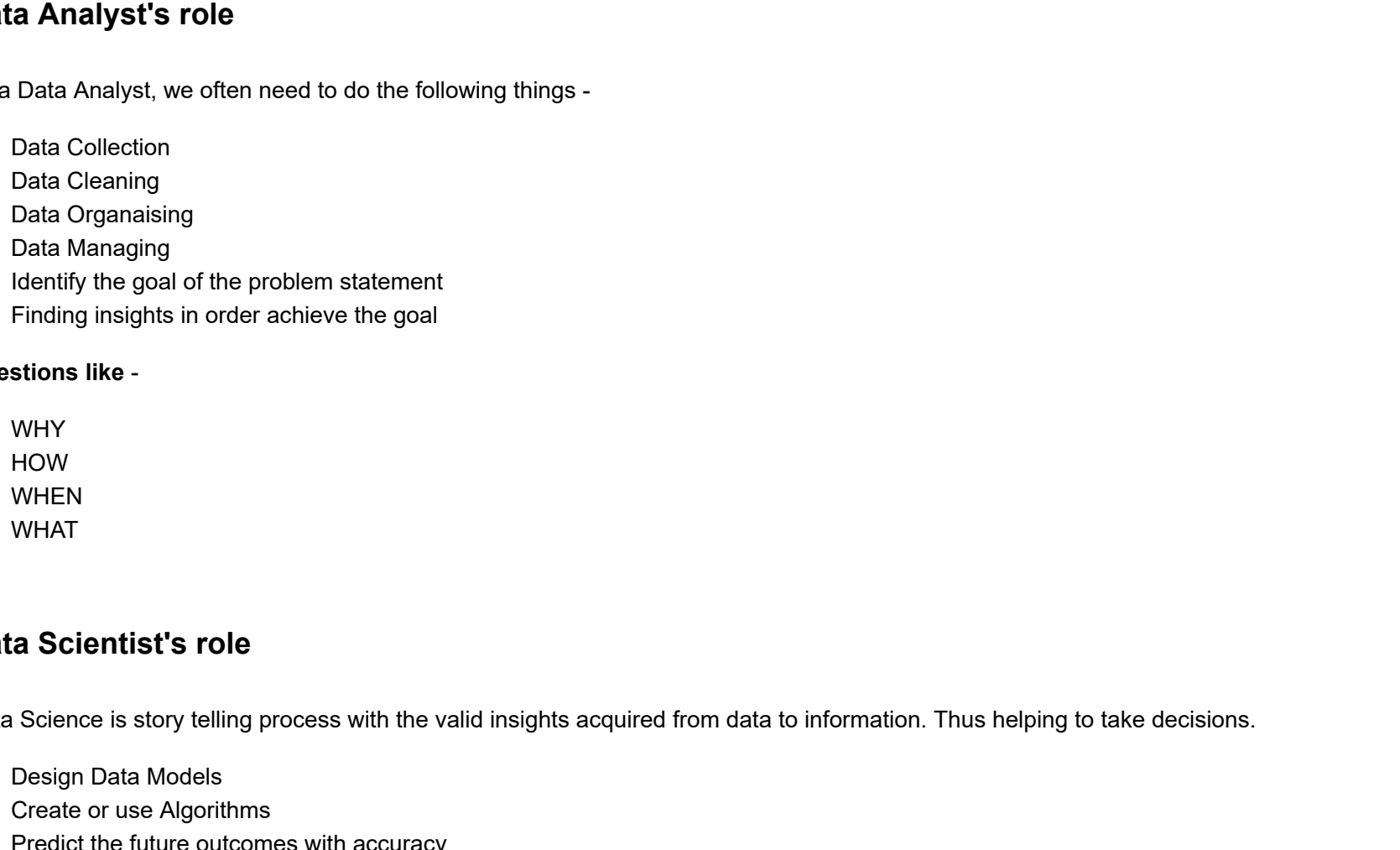
**Note** - Check google trends for comparison.

- Before clicking the below link, please sign-in your google account on web.
- Link → [https://trends.google.com/trends/explore?cat=1299&q=%2Fm%2F05z1\\_,%2Fm%2F0212/m,%2Fm%2F0%2F0%3d/j](https://trends.google.com/trends/explore?cat=1299&q=%2Fm%2F05z1_,%2Fm%2F0212/m,%2Fm%2F0%2F0%3d/j)

“ But when it comes to speed compatibility, Python is slower. ”

- Link → <https://julia.computing.com/blog/2020/06/fast-csv/>

## Data Analysis vs Data Science



## Data Analyst's role

As a Data Analyst, we often need to do the following things -

- Data Collection
- Data Cleaning
- Data Organising
- Data Managing
- Identify the goal of the problem statement
- Finding insights in order achieve the goal

**Questions like** -

- WHY
- HOW
- WHEN
- WHAT

## Data Scientist's role

Data Science is story telling process with the valid insights acquired from data to information. Thus helping to take decisions.

- Design Data Models
- Create or use Algorithms
- Predict the future outcomes with accuracy
- Make decisions from the insights

More information - <https://www.northeastern.edu/graduate/blog/what-does-a-data-scientist-do/>

**Note** - Data Scientist with analytical skills is a Blessing upon the blessed.

```
In [ ]: 
```

```
In [ ]: 
```

## Practise question

1. Collect data from online using Pandas.
2. Check if data cleaning is necessary.

### - yes → Clean the data

- no → Proceed
- 3. Identify the relationship between data variables.
- Apply Correlation
- Plot the relationship

- **Data Source** → [http://wiki.stat.ucla.edu/socr/index.php/SOCR\\_Data\\_Dinov\\_020108\\_HeightsWeights](http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_Dinov_020108_HeightsWeights)

```
In [ ]: 
```

```
In [ ]: 
```

## 1. Collect data from online

1. Pandas is python library mainly used for data analysis.
2. It is similar to doing analysis on Excel.
3. It is one of the best open source libraries available for doing data manipulation and data wrangling.

More information → <https://nandas.pydata.org/>

```
In [1]: import pandas as pd
```

```
In [2]: pip install pandas --user
```

Requirement already satisfied: pandas in c:\users\tapal\appdata\local\programs\python\python38-32\lib\site-packages (1.1.3)Note: you may need to restart the kernel to use updated packages.  
Requirement already satisfied: pytz>=2017.2 in c:\users\tapal\appdata\roaming\python\python38\site-packages (from pandas) (2020.1)  
Requirement already satisfied: python-dateutil>=2.7.3 in c:\users\tapal\appdata\roaming\python\python38\site-packages (from pandas) (2.8.1)  
Requirement already satisfied: numpy>=1.15.4 in c:\users\tapal\appdata\local\programs\python\python38-32\lib\site-packages (from pandas) (1.19.2)  
Requirement already satisfied: six>=1.5 in c:\users\tapal\appdata\roaming\python\python38\site-packages (from python-dateutil>=2.7.3->pandas) (1.15.0)

`read_html()` extracts all the tables from the html page.

```
In [3]: data_source = 'http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_Dinov_020108_HeightsWeights'  
data = pd.read_html(data_source)
```

```
In [4]: type(data)
```

```
Out[4]: list
```

```
In [5]: len(data)
```

```
Out[5]: 3
```

```
In [6]: data[0]
```

```
Out[6]:
```

0 Contents 1 SOCR Data - 25,000 Records of Human...

```
In [7]: data[1]
```

```
Out[7]:
```

	Index	Height(Inches)	Weight(Pounds)
0	1	65.78	112.99
1	2	71.52	136.49
2	3	69.40	153.03
3	4	68.22	142.34
4	5	67.79	144.30
...	...	...	...
195	196	65.80	120.84
196	197	66.11	115.78
197	198	68.24	128.30
198	199	68.02	127.47
199	200	71.39	127.88

200 rows × 3 columns

```
In [8]: data[2]
```

```
Out[8]:
```

	0	1	2	3	4	5	6	7	8	9	10	11
0 (default)	Deutsch	Espanol	Francais	Italiano	Portugues	日本語	България	الامارات العربية المتحدة	Suomi	ਭਾਰਤੀ	Norge	
1	한국어	中文	繁体中文	Русский	Nederlands	Ελληνικά	Hrvatska	Česká republika	Danmark	Polsha	România	Sverige

```
In [9]: df = data[1]
```

```
In [10]: df.head(5)
```

```
Out[10]:
```

	Index	Height(Inches)	Weight(Pounds)
0	1	65.78	112.99
1	2	71.52	136.49
2	3	69.40	153.03
3	4	68.22	142.34
4	5	67.79	144.30

## 2. Check if data cleaning is necessary

Data Cleaning is one of the important aspects in both Data Analysis and Data Science roles.

- It is one of the procedural steps where a data analyst or data scientist spends most of their time.

More information → [https://en.wikipedia.org/wiki/Data\\_cleaning](https://en.wikipedia.org/wiki/Data_cleaning)

### a. Check for any NaN values → Missing values

```
In [11]: df.isnull().any()
```

```
Out[11]:
```

Index False  
Height(Inches) False  
Weight(Pounds) False  
dtype: bool

- Since the dataset is sort of big, we cannot see all the values. Infact we cannot comprehend the actual missing values from the `isna()` dataset.
- In order to get the actual values (indices), the below function can be used.

Above result is clear, every column has `non-nan` values. Hence we can proceed with further steps.

### b. Check for the datatypes from each column

```
In [12]: df.dtypes
```

```
Out[12]:
```

Index int64  
Height(Inches) float64  
Weight(Pounds) float64  
dtype: object

```
In [13]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 200 entries, 0 to 199  
Data columns (total 3 columns):  
# Column Non-Null Count Dtype  
---  
0 Index 200 non-null int64  
1 Height(Inches) 200 non-null float64  
2 Weight(Pounds) 200 non-null float64  
dtypes: float64(2), int64(1)  
memory usage: 4.8 KB
```

Seems like every column has a unique data type. If at all there is then it is required to purify the data - make sure all the values are of same type.

### c. Overall description of the data frame

```
In [14]: df.describe()
```

```
Out[14]:
```

	Index	Height(Inches)	Weight(Pounds)
count	200.000000	200.000000	200.000000
mean	100.500000	67.949800	127.221950
std	57.879185	1.940363	11.960959
min	1.000000	63.430000	97.900000
25%	50.750000	66.522500	119.895000
50%	100.500000	67.935000	127.875000
75%	150.250000	69.202500	136.097500
max	200.000000	73.900000	158.960000

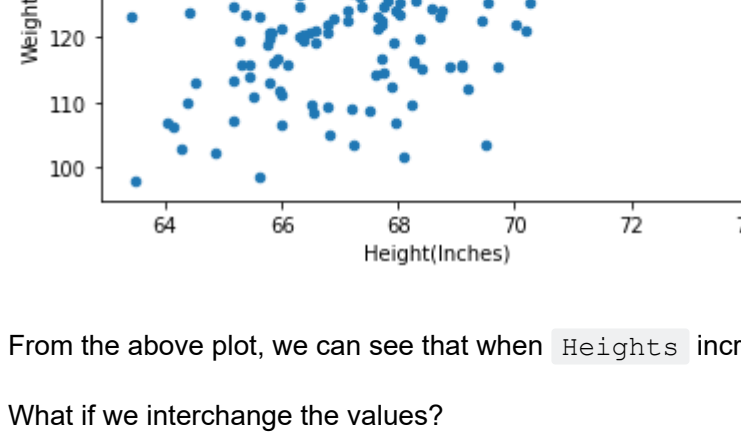
### d. Some visualization to explore more about the data

- We can use pandas plotting functions like `plot()` to explore about the data visually.
- `plot()` can show the following plots -
  - `line` → line plot (default)
  - `bar` → vertical bar plot
  - `barh` → horizontal bar plot
  - `hist` → histogram
  - `box` → boxplot
  - `kde` → Kernel Density Estimation plot
  - `density` → same as 'kde'
  - `area` → area plot
  - `pie` → pie plot
  - `scatter` → scatter plot
  - `hexbin` → hexbin plot

#### Ugly plot example

```
In [15]: df.plot()
```

```
Out[15]: <AxesSubplot: >
```



The above is the plot of all the data variables. This is not something we should do.

#### Plotting without unimportant data variables - excluded Index

```
In [16]: df['Weight(Pounds)']
```

```
Out[16]:
```

0 112.99  
1 136.49  
2 153.03  
3 142.34  
4 144.30  
...  
195 120.84  
196 115.78  
197 128.30  
198 127.47  
199 127.88  
Name: Weight(Pounds), Length: 200, dtype: float64

```
In [17]: df[['Height(Inches)', 'Weight(Pounds)']]
```

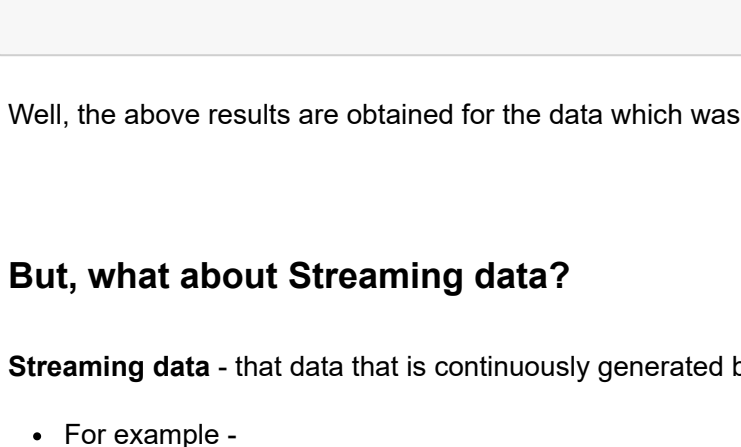
```
Out[17]:
```

	Height(Inches)	Weight(Pounds)
0	65.78	112.99
1	71.52	136.49
2	69.40	153.03
3	68.22	142.34
4	67.79	144.30
...	...	...
195	65.80	120.84
196	66.11	115.78
197	68.24	128.30
198	68.02	127.47
199	71.39	127.88

200 rows × 2 columns

```
In [18]: df[['Height(Inches)', 'Weight(Pounds)']].plot()
```

```
Out[18]: <AxesSubplot: >
```



The above is the plot of both `Heights` and `Weights` from the data frame `df`.

```
In [ ]: 
```

```
In [ ]: 
```

## 3. Relationship between data variables

**Correlation** - one of the statistical measurements applied to find out if any two variables are linrely related.

- If one variables is increasing, then other variable also increases. Vice versa.
- For example
  - If income of an employee increases then the household expenses increase.
  - If income of an employee decreases then the household expenses decrease.
- Scatter plot is really helpful to find the relationship between two variables. With this, it can be easily noticed the linear trend as well.

**Correlation plots** based on the correlation value obtained. → [https://en.wikipedia.org/wiki/Correlation\\_and\\_dependence#/media/File:Correlation\\_examples2.svg](https://en.wikipedia.org/wiki/Correlation_and_dependence#/media/File:Correlation_examples2.svg)

### b. Plot the relationship

```
In [19]: df.plot(x='Height(Inches)', y='Weight(Pounds)', kind='scatter')
```

```
Out[19]: <AxesSubplot: xlabel='Height(Inches)', ylabel='Weight(Pounds)'>
```



From the above plot, we can see that when `Heights` increase, then `Weights` also increased.

What if we interchange the values?

```
In [20]: df.plot(y='Height(Inches)', x='Weight(Pounds)', kind='scatter')
```

```
Out[20]: <AxesSubplot: xlabel='Weight(Pounds)', ylabel='Height(Inches)'>
```



### a. Find the Correlation

Correlation value ranges from `-1` to `1`.

- If the calculated correlation value is -
  - `-1`, then it is perfectly **negative correlation**
  - `1`, then it is perfectly **positive correlation**
  - `< 1`, then it means that **error** in the correlation measurement
  - `> 1`, then it means that **error** in the correlation measurement

More information → <https://www.investopedia.com/terms/c/correlationcoefficient.asp>

```
In [21]: df.corr()
```

```
Out[21]:
```

	Index	Height(Inches)	Weight(Pounds)
Index	1.000000	-0.094260	-0.128882
Height(Inches)	-0.094260	1.000000	0.556865
Weight(Pounds)	-0.128882	0.556865	1.000000

```
In [22]: relation = df.corr()
```

```
relation.style.background_gradient(cmap='Reds')
```

```
Out[22]:
```

	Index	Height(Inches)	Weight(Pounds)
Index	1.000000	-0.094260	-0.128882
Height(Inches)	-0.094260	1.000000	0.556865
Weight(Pounds)	-0.128882	0.556865	1.000000

```
In [ ]: 
```

```
In [ ]: 
```

Well, the above results are obtained for the data which was already stored.

## But, what about Streaming data?

**Streaming data** - that data that is continuously generated by different sources is called streaming data.

- For example -
  - Tesla Self-driving Car generates the data continuously.
  - One tesla car generates 11 TB and 152 TB data per day.
  - Big data problem
- More information → <https://www.luxera.com/blog/autonomous-and-adas-test-cars-produce-over-11-tb-of-data-per-day/>

```
In [ ]: 
```

```
In [ ]: 
```

## Case Study → Activity

1. Select any one of these or you can find your own topic of interest not specifically from below.

- Study to analyse peoples' habits on YouTube platform
- Study to analyse the changes occurred in peoples' life due to Demonization
- Study to analyse the students' overall development due to online education

2. Create a google form where you can have a set of questions and answer options.

- Have atleast 8 to 10 questions
- 3. Collect the data from your friends, families etc (by sharing the link).
- The data will be stored in your drive (in a spreadsheet)

4. Once the data is collected -
  - Create your own data variables from the questions
  - Try to basic analysis like processing and visualization

**Note** - To learn how to create google forms (Questionnaires) and collect the data,

- Please watch this video → <https://www.youtube.com/watch?v=Qy2DlylDUU>