



# FUNDAMENTAL DATA ANALYSIS WITH PYTHON

---

By Muhamad Irfan Fadhullah



# Outline

- Pandas Foundation
- Manipulating Data Using Pandas
- Fandas Functionality
- Data Cleansing
- Data Wrangling
- Data Visualization
- Exploratory Data Analysis

# Pandas Foundation

## Why Pandas?

Pandas is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool built on top of the Python programming language.



### Advantages and Disadvantages of Using Pandas Library



#### Advantages

Less writing and more work done

Excellent data representation

Made for Python

An extensive set of features

#### Disadvantages

Steep learning curve

Difficult syntax

Poor compatibility for 3D matrices

Bad documentation

Pandas





# Manipulating Data Using Pandas

## Examples:

### SELECTION

```
year = datac[['adr', 'reservation_status_date']]
```

### ADDITION

```
df['gender'] = ['male', 'male', 'female']  
df
```

### RENAME

```
cek=high.groupby([high['reservation_status_date'].dt.year.rename('Year'),  
                  high['reservation_status_date'].dt.month.rename('Month')]).agg('count')
```



# Manipulating Data Using Pandas

## DELETION

```
datac=data.drop(columns=['company'],axis=1)
```

## SORTING

```
df.sort_index(ascending =False)
```



# PANDAS FUNCTIONALITY

BASIC FUNCTION

DATA  
ENGINEERING

SIMPLE PLOTTING

## Basic Information

<code>&gt;&gt;&gt; df.shape</code>	(rows,columns)
<code>&gt;&gt;&gt; df.index</code>	Describe index
<code>&gt;&gt;&gt; df.columns</code>	Describe DataFrame columns
<code>&gt;&gt;&gt; df.info()</code>	Info on DataFrame
<code>&gt;&gt;&gt; df.count()</code>	Number of non-NA values

## Summary

<code>&gt;&gt;&gt; df.sum()</code>	Sum of values
<code>&gt;&gt;&gt; df.cumsum()</code>	Cummulative sum of values
<code>&gt;&gt;&gt; df.min()/df.max()</code>	Minimum/maximum values
<code>&gt;&gt;&gt; df.idxmin()/df.idxmax()</code>	Minimum/Maximum index value
<code>&gt;&gt;&gt; df.describe()</code>	Summary statistics
<code>&gt;&gt;&gt; df.mean()</code>	Mean of values
<code>&gt;&gt;&gt; df.median()</code>	Median of values

## CHECK MISSING VALUES

```
df.notnull()
```

## DROP MISSING VALUES

```
df.dropna()
```

## FILLING MISSING VALUES – DIRECT REPLACE

```
df.fillna(ScalarValue)
```

## COMPUTING UNIQUE VALUES

```
df.unique()
```

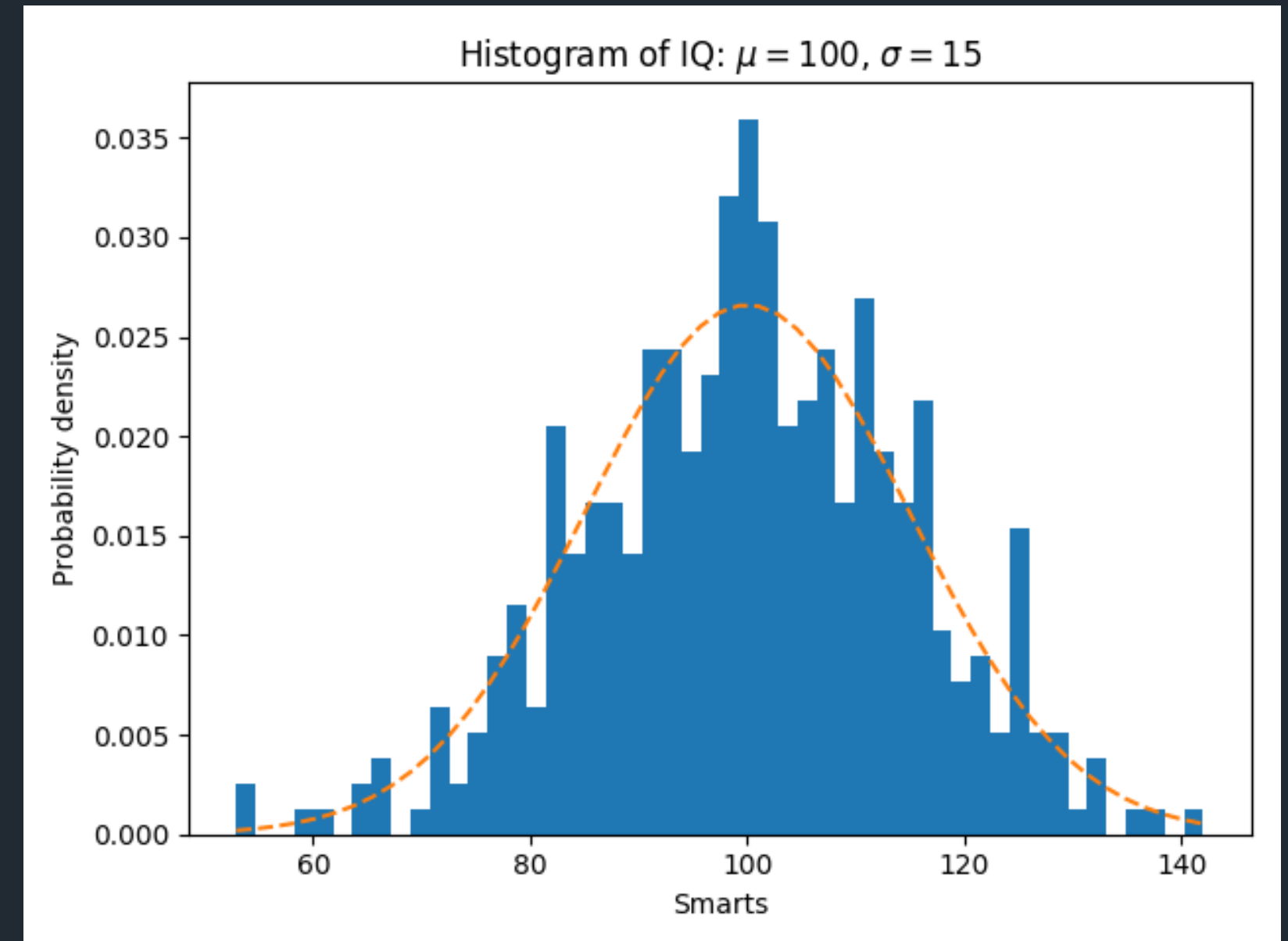
```
df.value_counts()
```



# SIMPLE PLOTTING

Pandas Dataframe offers a range of graphical plotting options. We can plot, box plot, area, scatter plots, stacked charts, barcharts, histograms, etc.

- `df.plot.scatter()` #plots a scatter chart
- `df.plot.line()` # plots a line chart
- `df.boxplot()` # plots a box plot





# CHEAT SHEET PANDAS FROM DATACAMP

## Python For Data Science Cheat Sheet Pandas Basics

Learn Python for Data Science Interactively at [www.DataCamp.com](https://www.datacamp.com)



### Pandas

The Pandas library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language.



Use the following import convention:

```
>>> import pandas as pd
```

### Pandas Data Structures

#### Series

A one-dimensional labeled array capable of holding any data type

a	3
b	-5
c	7
d	4

Index

```
>>> a = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

#### DataFrame

Columns

	Country	Capital	Population
0	Belgium	Brussels	11190846
1	India	New Delhi	1303171035
2	Brazil	Brasilia	207847528

Index

A two-dimensional labeled data structure with columns of potentially different types

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],  
          'Capital': ['Brussels', 'New Delhi', 'Brasilia'],  
          'Population': [11190846, 1303171035, 207847528]}
```

```
>>> df = pd.DataFrame(data,  
                      columns=['Country', 'Capital', 'Population'])
```

### I/O

#### Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)  
>>> df.to_csv('myDataFrame.csv')
```

#### Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')  
>>> df.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
```

Read multiple sheets from the same file

```
>>> xlax = pd.ExcelFile('file.xls')  
>>> df = pd.read_excel(xlax, 'Sheet1')
```

### Asking For Help

```
>>> help(pd.Series.loc)
```

### Selection

Also see NumPy Arrays

#### Getting

```
>>> a['b']  
-5
```

Get one element

```
>>> df[1:]
```

Get subset of a DataFrame

```
Country Capital Population  
1 India New Delhi 1303171035  
2 Brazil Brasilia 207847528
```

### Selecting, Boolean Indexing & Setting

#### By Position

```
>>> df.iloc[[0], [0]]
```

```
'Belgium'
```

Select single value by row & column

```
>>> df.iat[[0], [0]]
```

```
'Belgium'
```

#### By Label

```
>>> df.loc[[0], ['Country']]
```

```
'Belgium'
```

Select single value by row & column labels

```
>>> df.at[[0], ['Country']]
```

```
'Belgium'
```

#### By Label/Position

```
>>> df.ix[2]
```

```
Country Brazil  
Capital Brasilia  
Population 207847528
```

Select single row of subset of rows

```
>>> df.ix[:, 'Capital']
```

```
0 Brussels  
1 New Delhi  
2 Brasilia
```

Select a single column of subset of columns

```
>>> df.ix[1, 'Capital']
```

```
'New Delhi'
```

Select rows and columns

#### Boolean Indexing

```
>>> a[~(a > 1)]
```

```
>>> a[(a <= 1) | (a > 2)]
```

```
>>> df[df['Population'] > 1200000000]
```

Series a where value is not >1  
a where value is <=1 or >2

Use filter to adjust DataFrame

#### Setting

```
>>> a['a'] = 6
```

Set index a of Series a to 6

### Dropping

```
>>> a.drop(['a', 'c'])
```

Drop values from rows (axis=0)

```
>>> df.drop('Country', axis=1)
```

Drop values from columns (axis=1)

### Sort & Rank

```
>>> df.sort_index()
```

```
>>> df.sort_values(by='Country')
```

```
>>> df.rank()
```

Sort by labels along an axis  
Sort by the values along an axis  
Assign ranks to entries

### Retrieving Series/DataFrame Information

#### Basic Information

```
>>> df.shape
```

(rows, columns)

```
>>> df.index
```

Describe index

```
>>> df.columns
```

Describe DataFrame columns

```
>>> df.info()
```

Info on DataFrame

```
>>> df.count()
```

Number of non-NA values

#### Summary

```
>>> df.sum()
```

Sum of values

```
>>> df.cumsum()
```

Cumulative sum of values

```
>>> df.min()/df.max()
```

Minimum/maximum values

```
>>> df.idxmin()/df.idxmax()
```

Minimum/Maximum index value

```
>>> df.describe()
```

Summary statistics

```
>>> df.mean()
```

Mean of values

```
>>> df.median()
```

Median of values

### Applying Functions

```
>>> f = lambda x: x*2
```

```
>>> df.apply(f)
```

```
>>> df.applymap(f)
```

Apply function

Apply function element-wise

### Data Alignment

#### Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> a3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
```

```
>>> a + a3
```

```
a 10.0
```

```
b NaN
```

```
c 5.0
```

```
d 7.0
```

#### Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=0)
```

```
a 10.0
```

```
b -5.0
```

```
c 5.0
```

```
d 7.0
```

```
>>> s.sub(s3, fill_value=2)
```

```
>>> s.div(s3, fill_value=4)
```

```
>>> s.mul(s3, fill_value=3)
```

DataCamp

Learn Python for Data Science Interactively



# DATA CLEANSING



Data cleaning or cleansing is the process of detecting and correcting (or removing) corrupt or inaccurate records from a record set, table, or database and refers to identifying incomplete, incorrect, inaccurate or irrelevant parts of the data and then replacing, modifying, or deleting the dirty or coarse data.

# DATA WRANGLING

Data wrangling is one of the most important components in the data science workflow. It involves the processing of data in various formats like concatenating, grouping, merging, etc. for the purpose of getting them used with another set of data or for analysing.



# **DATA CLEANSING**

## **Common Problems in the Data Cleansing:**

- 1. Duplicate Dataset**
- 2. Missing Data**
- 3. Outliers**
- 4. Data Type**

# METHODS IN COMBINING DATA:

1. JOIN
2. MERGE
3. CONCAT

# DATA WRANGLING

Combining Data

data1		data2	
X1	X2	X1	X3
a	11.432	a	20.784
b	1.303	b	NaN
c	99.906	d	20.784

Merge

```
>>> pd.merge(data1,
              data2,
              how='left',
              on='X1')
```

X1	X2	X3
a	11.432	20.784
b	1.303	NaN
c	99.906	NaN

```
>>> pd.merge(data1,
              data2,
              how='right',
              on='X1')
```

X1	X2	X3
a	11.432	20.784
b	1.303	NaN
d	NaN	20.784

```
>>> pd.merge(data1,
              data2,
              how='inner',
              on='X1')
```

X1	X2	X3
a	11.432	20.784
b	1.303	NaN

```
>>> pd.merge(data1,
              data2,
              how='outer',
              on='X1')
```

X1	X2	X3
a	11.432	20.784
b	1.303	NaN
c	99.906	NaN
d	NaN	20.784

Join

```
>>> data1.join(data2, how='right')
```

Concatenate

Vertical

```
>>> s.append(s2)
```

Horizontal/Vertical

```
>>> pd.concat([s,s2],axis=1, keys=['One','Two'])
>>> pd.concat([data1, data2], axis=1, join='inner')
```

# Data Visualization With Python

## Introduction

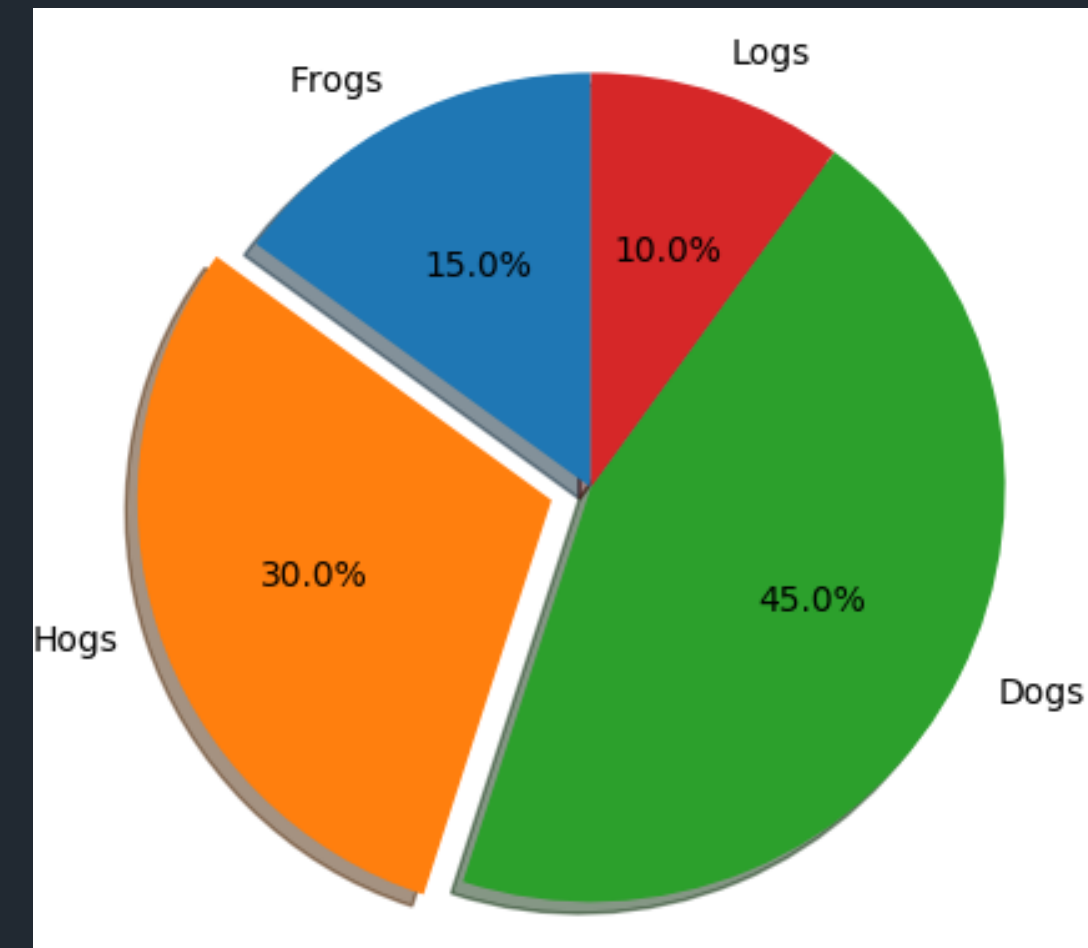


# How to create right Data Visualization

- Understand the context!
- Choose the appropriate plot type. If there are various options, we can try to compare them, and choose the one that fits our model the best.
- When we choose your type of plot, one of the most important things is to label your axis.
- Add a title to make our plot more informative.
- Add labels for different categories when needed.
- Optionally we can add a text or an arrow at interesting data points.
- In some cases we can use some sizes and colors of the data to make the plot more informative.

# Pie Chart

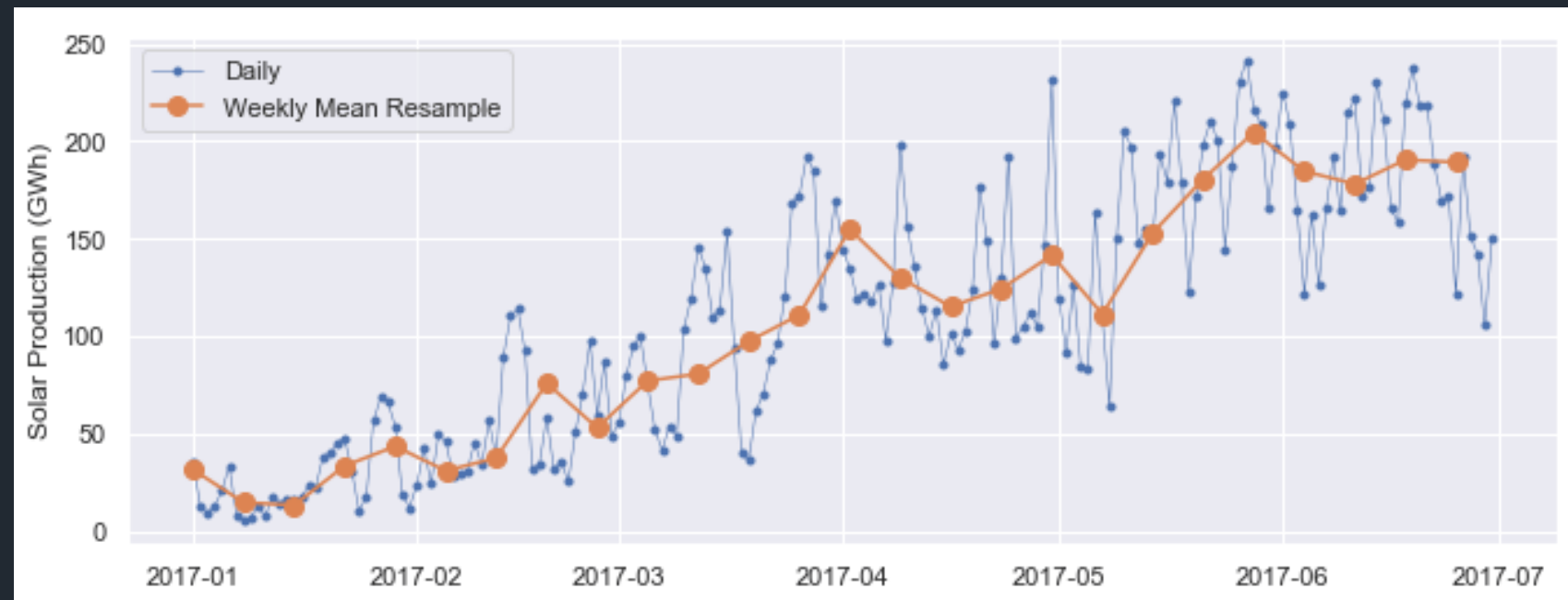
Pie chart should be used seldom as it is difficult to compare sections of the chart. Barplot is used instead as comparing sections is easy. eg: Market share in Films





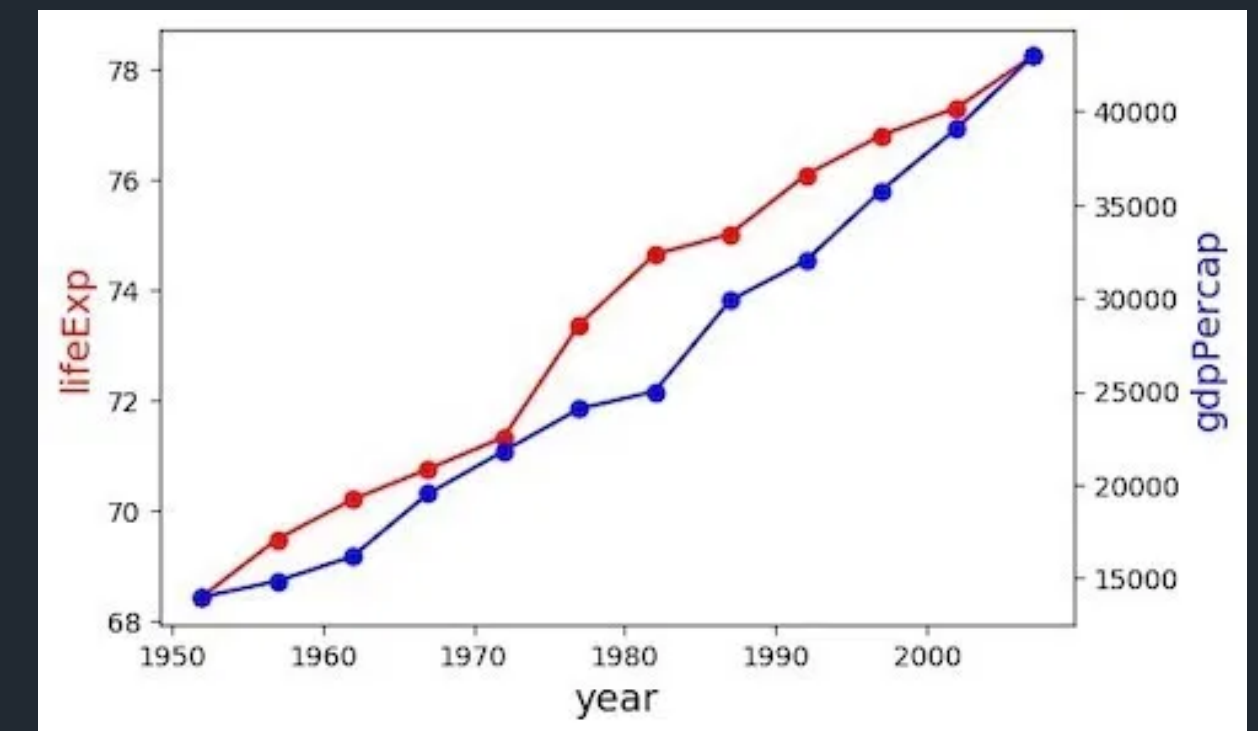
# LINE CHART

## TIME SERIES



It should when we require 2 plots or grouped data in the same direction.  
e.g : Population and GDP data in the same x-axis

## TWIN AXIS

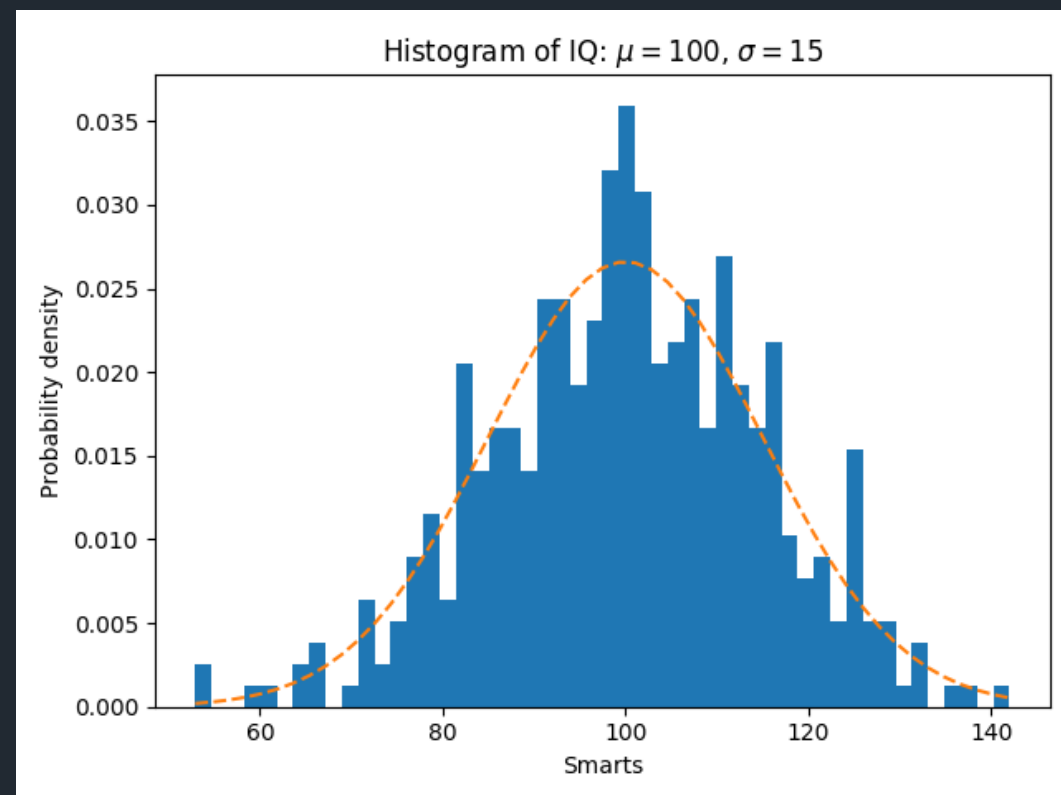


Time Series should be used when single or multiple variables are to be plotted over time. Eg: Stock Market Analysis of Companies, Weather Forecasting.

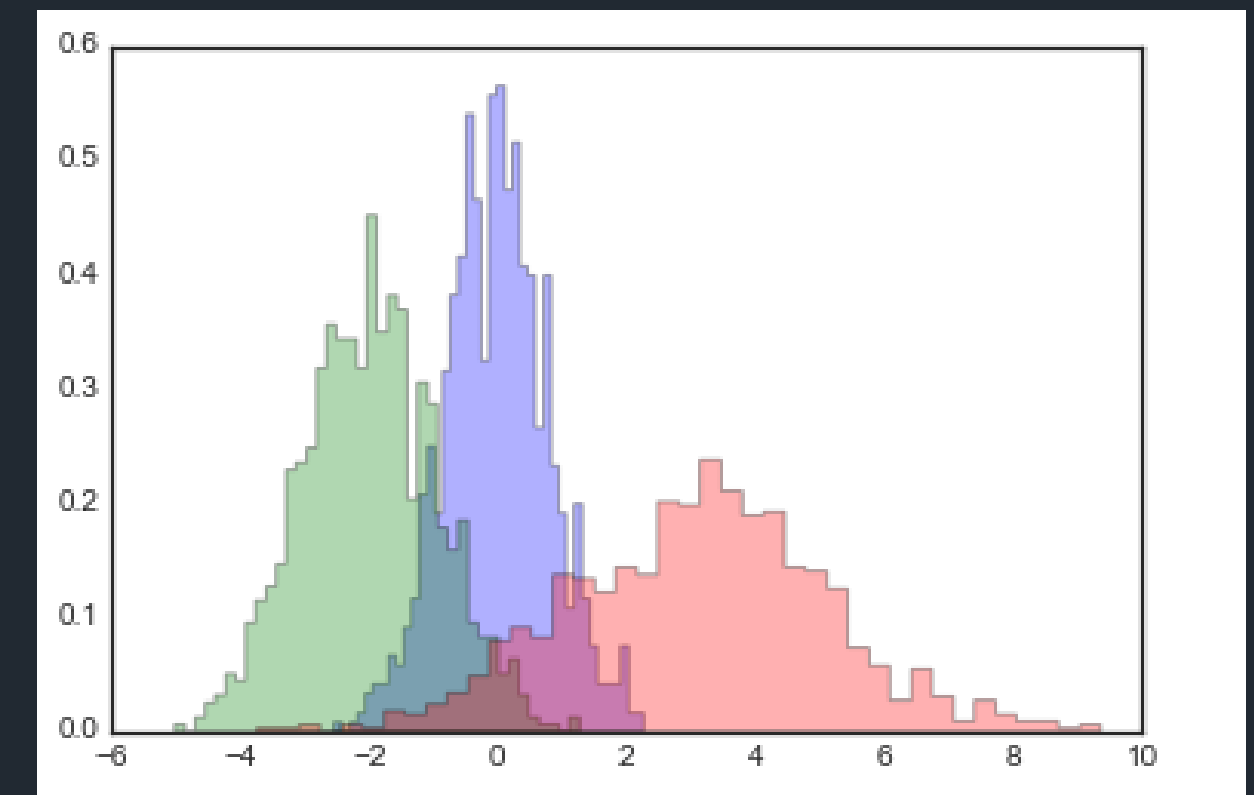
# HISTOGRAM



## SIMPLE



## MULTIPLE HISTOGRAM

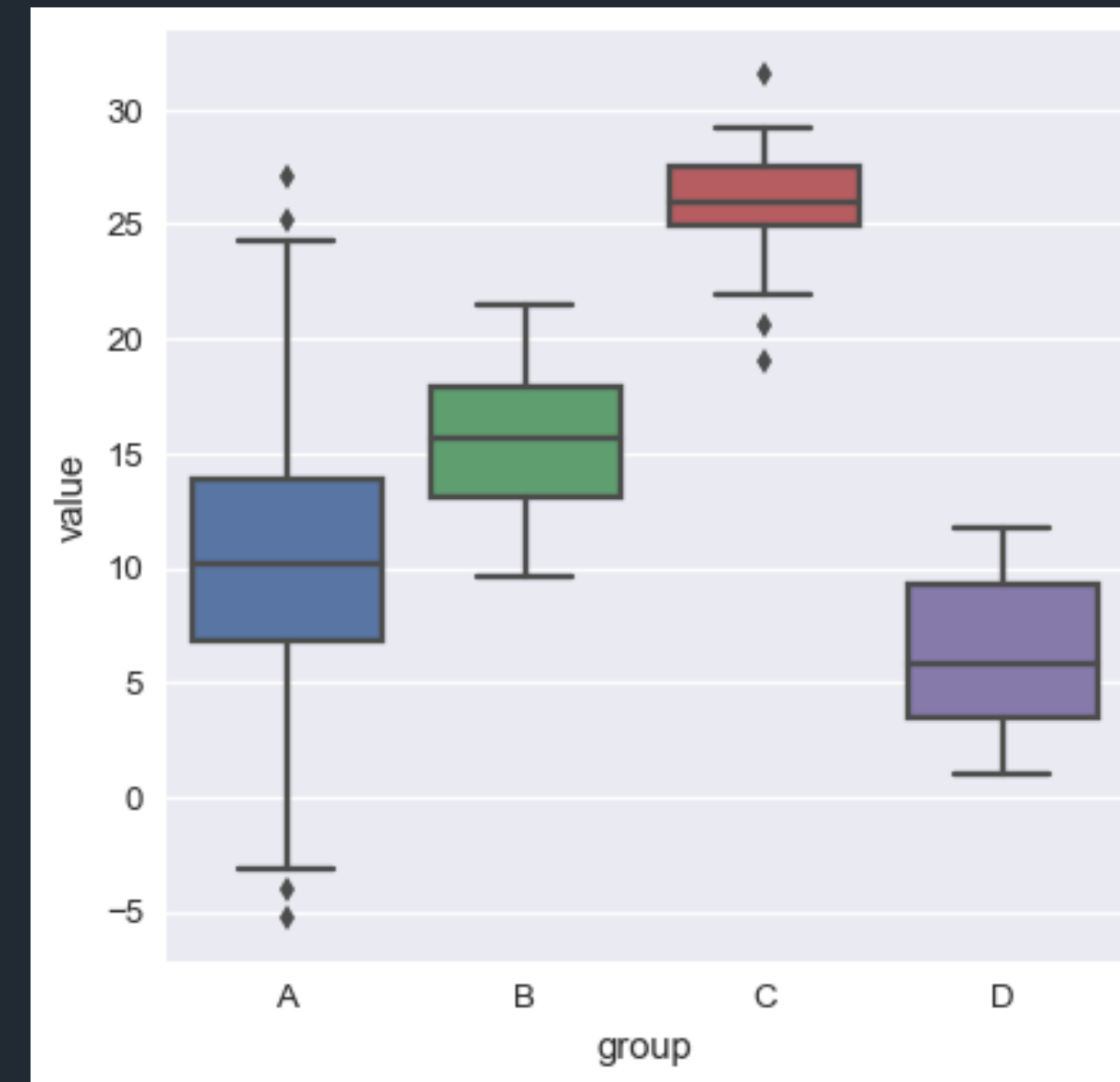


We should use histogram when we need the count of the variable in a plot. eg:  
Number of particular games sold in a store.

When we need to understand the distributions between 2 entity variables. We can see that Grand Canyon has comparably more visitors than Bryce Canyon

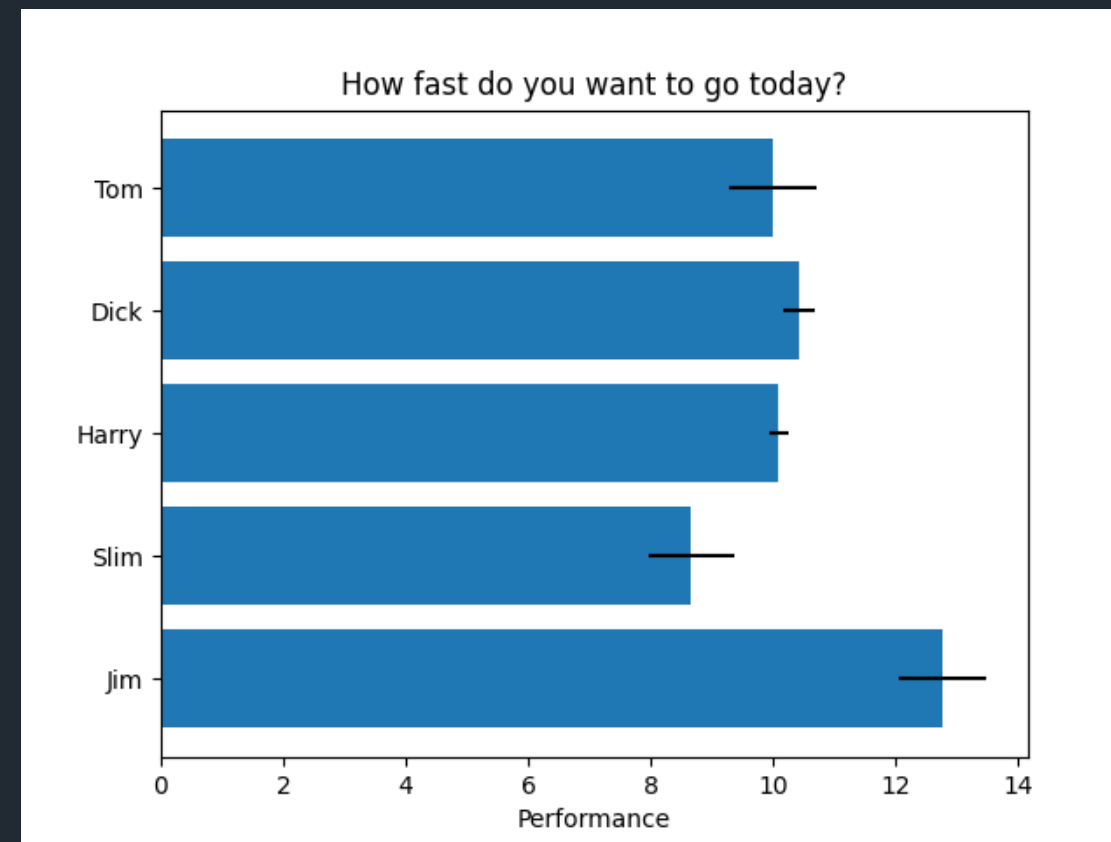
# BOX PLOT

It should be used when we require to use the overall statistical information on the distribution of the data. It can be used to detect outliers in the data. Eg: Credit Score of Customer. We can get the max, min and much more information about the mark



# BAR PLOT

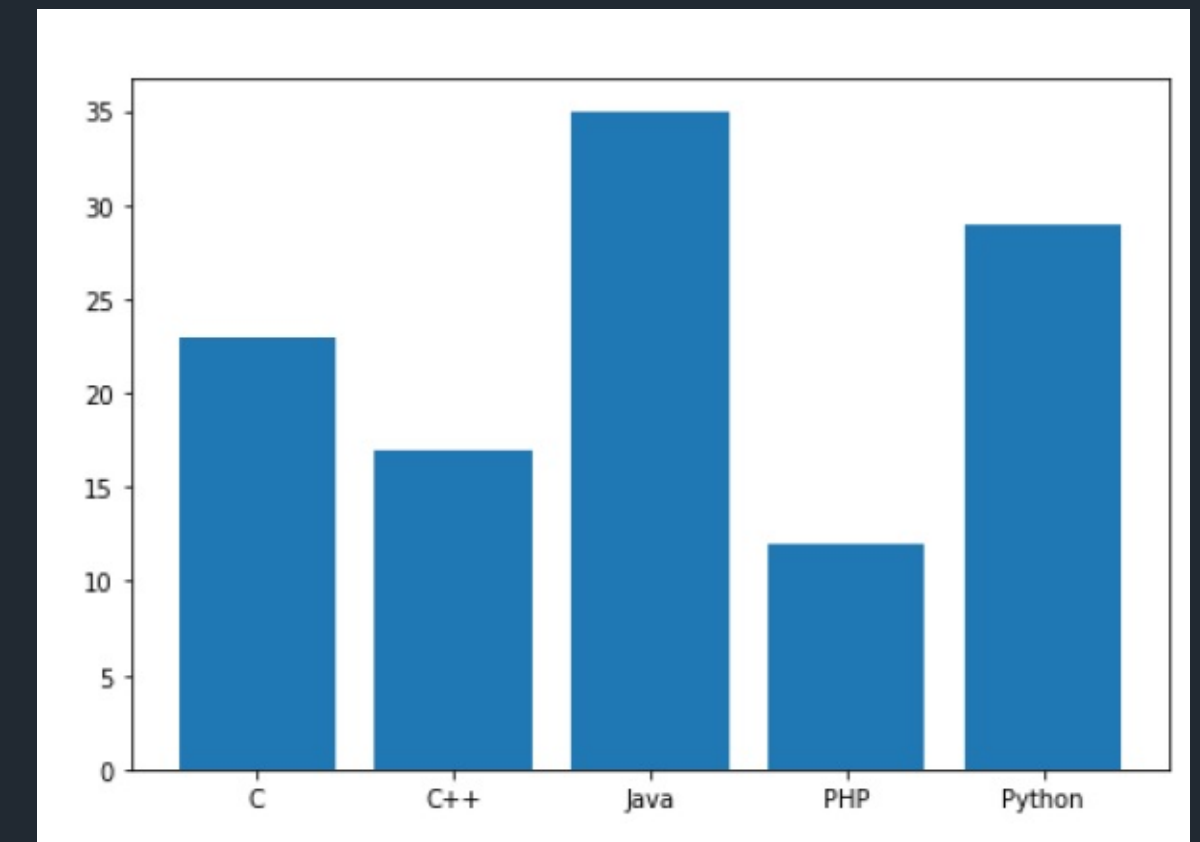
## HORIZONTAL



It is used when to compare between several groups.

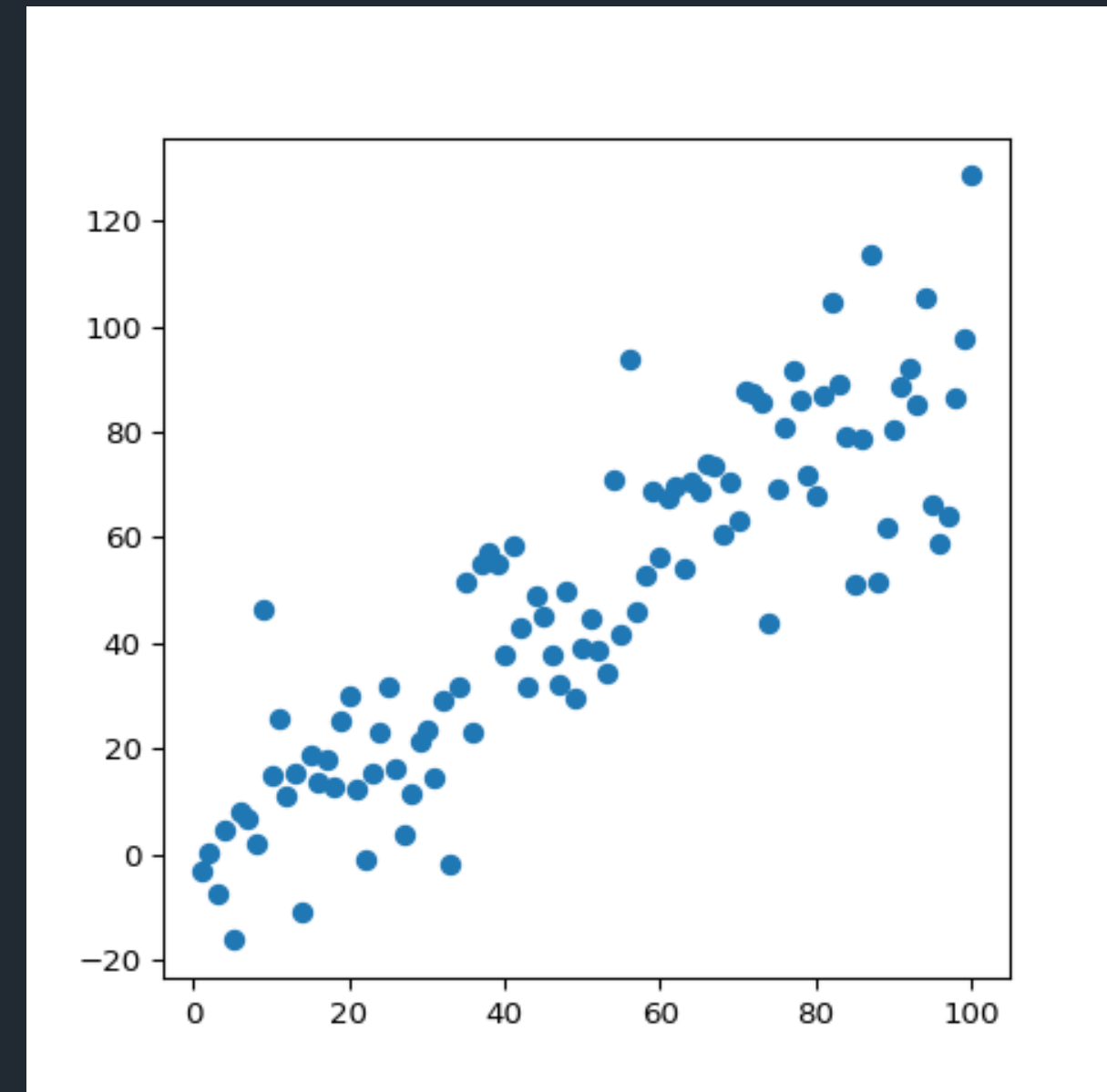
Eg: Student marks in an exam

## VERTIKAL



# SCATTER PLOT

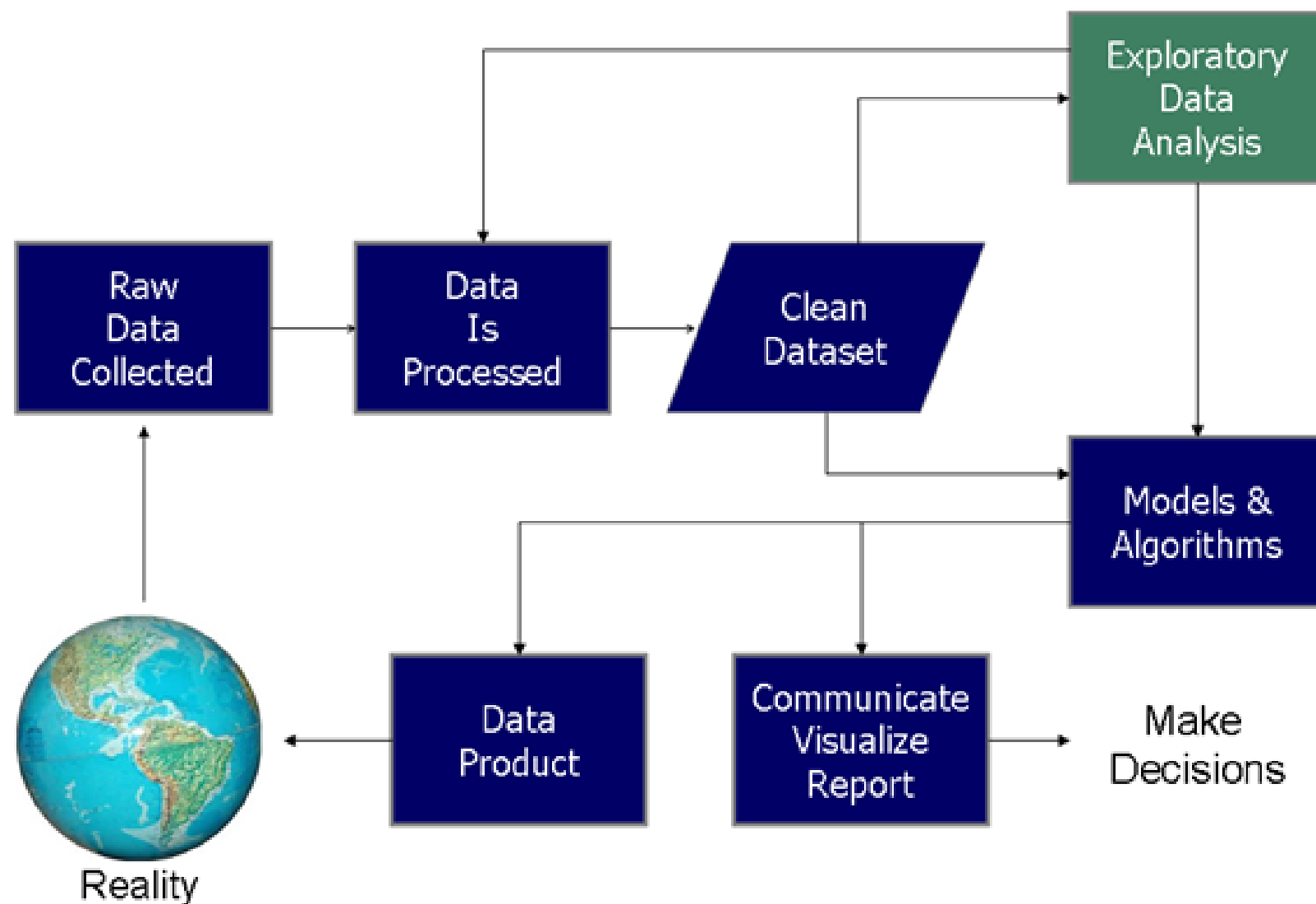
It is used in Machine learning concepts like regression, where  $x$  and  $y$  are continuous variables



# Exploratory Data Analysis



## Data Science Process



Exploratory Data Analysis is an approach to analyzing data sets to summarize their main characteristics, often with visual methods



# OBJECTIVES

## QUICKLY DESCRIBE A DATASET

number of  
rows/columns,  
missing data,  
datatypes, preview  
of the data

## CLEAN CORRUPTED DATA

handle missing  
data, invalid data  
types,  
incorrect values

## VISUALIZE DATA DISTRIBUTIONS

bar chart,  
histogram, box  
plots

## CALCULATE AND VISUALIZE CORRELATION

Calculate and  
visualize correlation  
(relationships)  
between  
variables (heatmap)



# REFERENCES

- <https://iykra.com>
- [datacamp.com](https://datacamp.com)
- [towardsdatascience.com](https://towardsdatascience.com)
- [www.analyticsvidhya.com](https://www.analyticsvidhya.com)
- [pandas.pydata.org](https://pandas.pydata.org)
- [stackoverflow.com](https://stackoverflow.com)
- <https://data-flair.training/blogs>





# Thank You

**Contact Information**

**EMAIL ADDRESS**

irfanfadhullah@gmail.com

**For any questions or concerns**