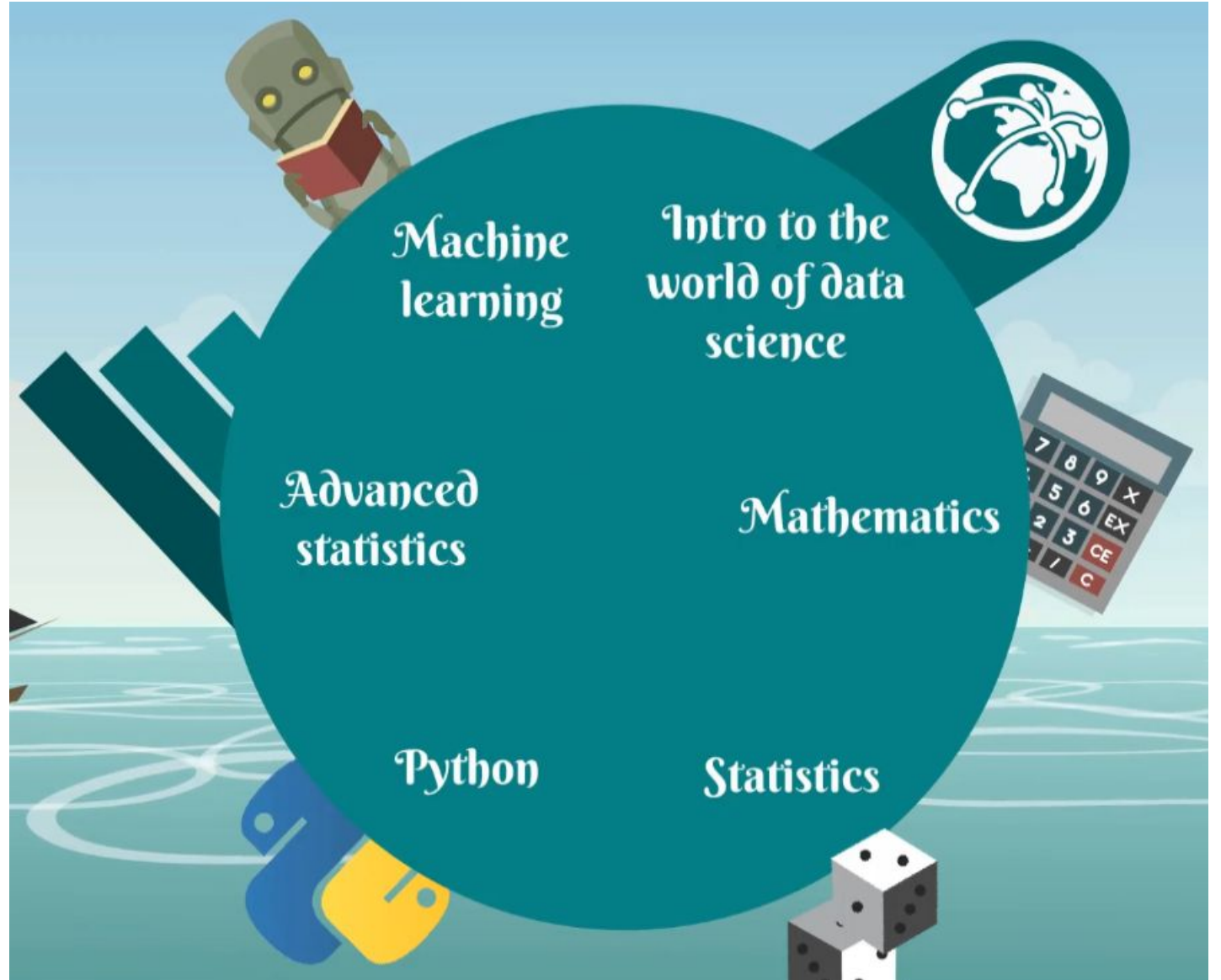


What is Data Science?

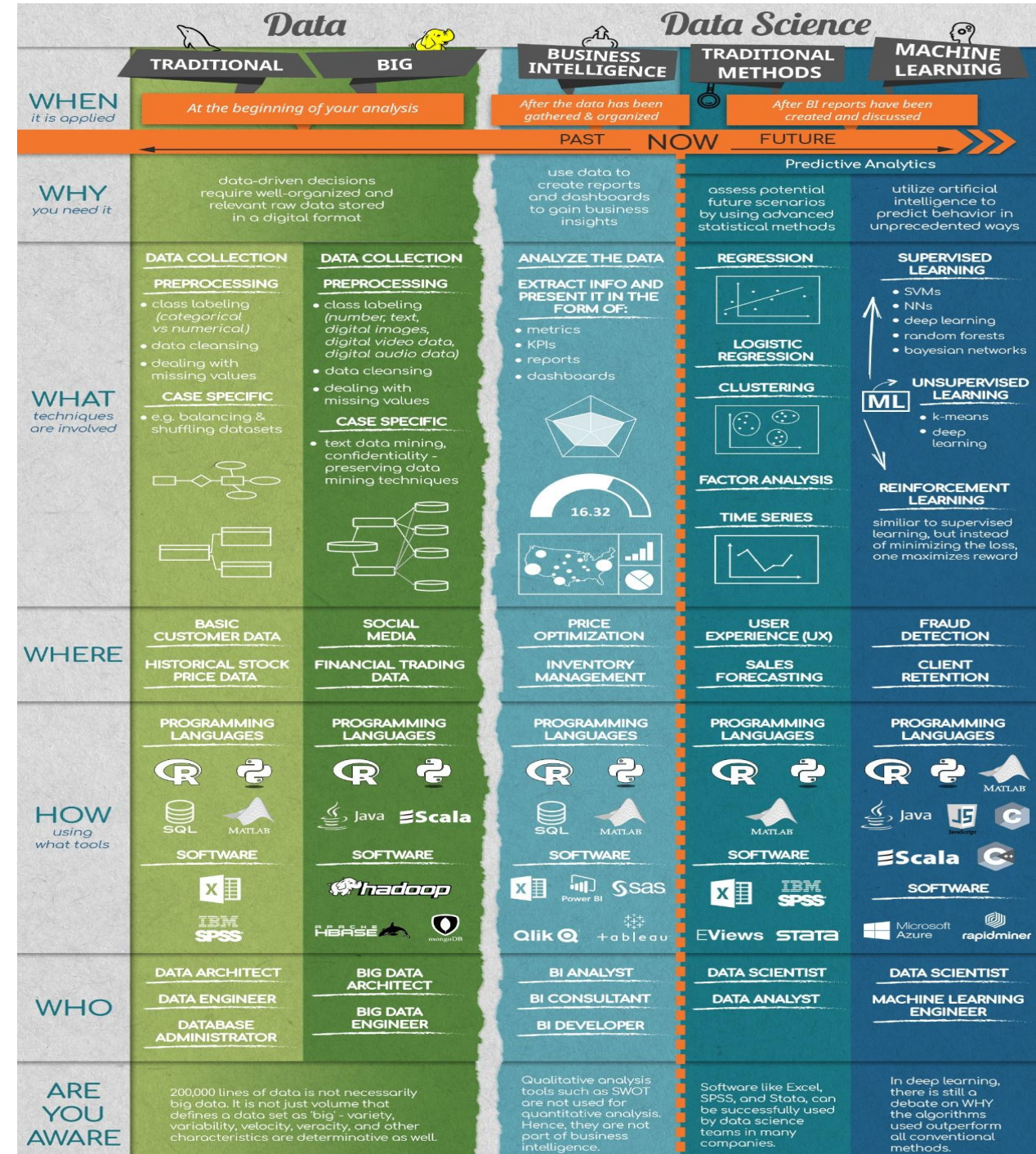
- **Data Science is the science of analysing raw data** using statistics and machine learning techniques with the purpose of drawing insights from the data.
- **Data Science** is used in many industries to allow them to make better business decisions, and in the sciences to test models or theories.
- This requires a process of inspecting, cleaning, transforming, modelling, analyzing, and interpreting raw data.

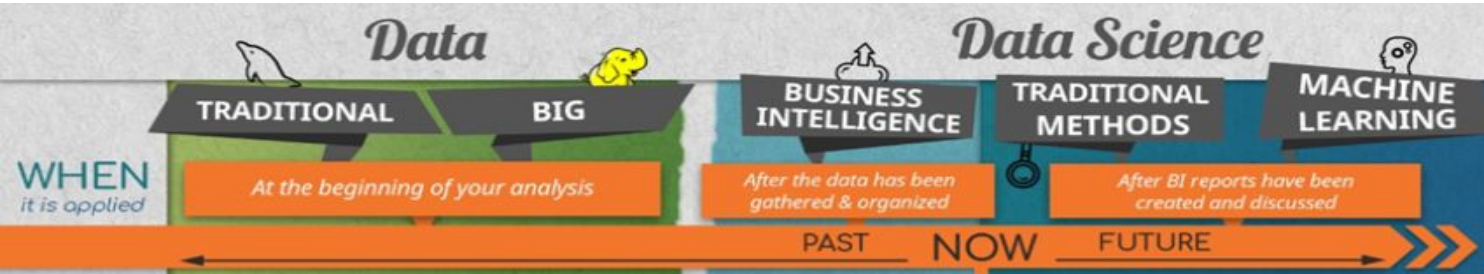


Important Disciplines Under Data Science



A Breakdown of Data Science



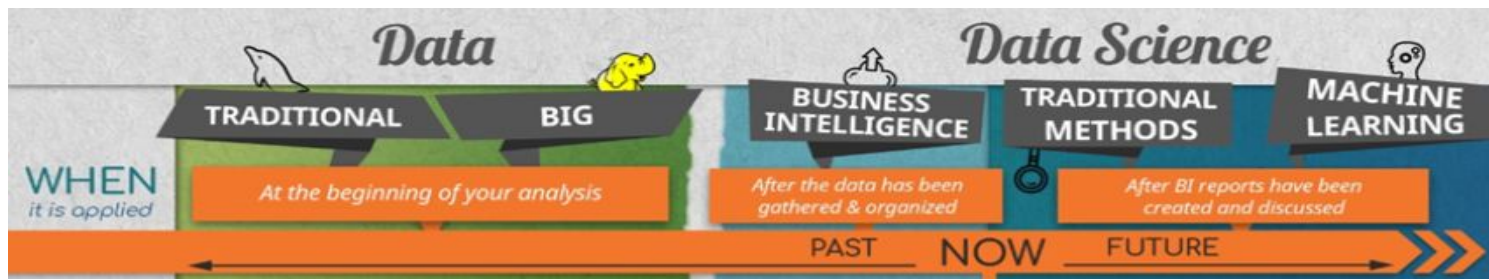


The step by step comparison between the terms and buzzwords related to each discipline

- **DATA**

- Data is defined as information stored in a digital format which can then be used as a base for performing analyses and decision making.
- As you can see there are two types of data.
- **Traditional data and big data:** Dealing with data is the first step when solving business problems or researching.
- Traditional data is the data in the form of tables containing numeric or text values data that is structured and stored in databases which can be managed from one computer.

- **Big data** is a term reserved for extremely large data and it is not just humongous in terms of volume. This data could be in various formats. It can be structured, Semi structure or unstructured. Big data is just that big.
- You will also often see it characterized by the letter V as in big data. They may include the vision you have about big data, the value Big Data carries, the visualization tools you use or the variability and the consistency of big data and so on.
- However, the following are probably the most important criteria. You must remember volume as we already said. Big Data needs a whopping amount of memory space typically distributed between minicomputers. Its size is measured in terabytes, Peta bytes, and even exabytes variety.
- Here we are not talking just about numbers and text. Big data often implies dealing with images, audio files, mobile data and others.
- Velocity when working with big data. One's goal is to make extracting patterns from it as quickly as possible. The progress that has been done in this area is remarkable outputs from huge data sets can be retrieved in real time.



- **Data science** is a broad subject. It's an interdisciplinary field that combines statistical, mathematical, programming, problem solving and data management tools.
- We have divided data science into three segments **business intelligence, traditional methods and machine learning.**

- **Business Intelligence is the discipline includes technology driven tools involved in the process of analysing, understanding, and reporting available past data.** This will result in having reports or dashboards and will help you on your way to making an informed strategic and tactical business decisions.
- You can extract insights and ideas about your business that will help to grow and give you an edge of your competitors giving you added stability.
- Business intelligence means understanding-
 - how your sales grew and why did competitors lose market share,
 - Was there an increase in the price of your products or did you sell a mix of more expensive products?
 - How did your profitability margins behave in the same time frame of a previous year?
 - Were there client accounts that were more profitable?
- This is what BI is all about understanding past business performance in order to improve future performance.
- Once your BI-reports and dashboards are completed and presented it's time to apply one of two types of data science.



- **Traditional methods according to our framework are a set of methods that are derived mainly from statistics and are adapted for business.**
- There is no denying that these conventional data science tools are applicable today. They are perfect for forecasting future performance with great accuracy.
- Regression analysis, cluster analysis and factor analysis all of which are prime examples of traditional methods.

- The last column we will be discussing with **machine learning** in contrast to traditional methods.
- The responsibility is left for the machine through mathematics.
- **A significant amount of computer power in applying AI the machine is given the ability to predict outcomes from data without being explicitly programmed to smell is all about creating algorithms that let machines receive data perform calculations and apply statistical analysis in order to make predictions with unprecedented accuracy.**

The Benefits of Each Discipline



- There are two types of data. **Traditional and big data**. Data driven decisions require well organized and relevant raw data stored in a digital format which can be processed and transformed into meaningful and useful information. It is the material on which you base your analysis. **Without data, a decision maker wouldn't be able to test their decisions and ensure they have taken the right course of action.**
- The data you have describes what happened in the past. It is the job of the **business intelligence** analyst to study the numbers and **explain where and why some things went well and others not so well**. Having the business context in mind the business intelligence analyst will present the data in the form of reports and dashboards.
- What else is needed once the patterns have been interpreted. You can forecast potential future outcomes. The application of any term related to the columns **traditional methods or machine learning** can be said to belong to the field of **predictive analytics**.
- There is a difference between the two. **Traditional methods relate to traditional data**. They were designed prior to the existence of big data where the technology simply wasn't as advanced as it is today. They involve applying statistical approaches to create predictive models.
- If you want to dig deeper however or tackle huge amounts of big data utilizing unconventional methods or AI then you can predict behaviour in unprecedented ways using **machine learning** techniques and tools. Both techniques are useful for different purposes.
- Traditional methods are better suited for traditional data while machine learning will have better results when it comes to tackling big data.

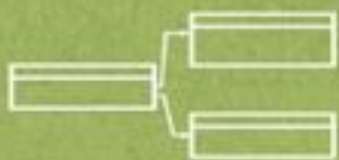
DATA COLLECTION

PREPROCESSING

- class labeling (categorical vs numerical)
- data cleansing
- dealing with missing values

CASE SPECIFIC

- e.g. balancing & shuffling datasets



Techniques for Working with Traditional Data

Data Collection

The gathering of raw data is referred to as **data collection**

Example would be the use of surveys asking people to rate how much they like or dislike a product or experience on a scale of 1 to 10

Preprocessing

Preprocessing is a group of operations that will basically convert your raw data into a format that is more understandable and hence useful for further processing.

For example Customer has entered age as 942 or name as UK then off course this entries are invalid which you need to correct before further processing.

Class Labelling

One technique is **class labelling**. This involves labelling the data point to the correct data type or arranging data by category.

One such category is **numerical**. For example if you are storing the number of goods sold daily t. These are numbers which can be manipulated such as the average number of goods sold per day or month.

The other label is **categorical**. Here you are dealing with information that cannot have mathematical manipulations. For example a person's profession or place of birth

Data Cleansing

The goal of data cleansing is to deal with inconsistent data.

This can come in various forms. Say you are provided with a data set containing the US states and a quarter of the names are misspelled in this situation. Certain techniques must be performed to correct these mistakes

Missing Values

Missing values are another thing you'll have to deal with.

Data cleansing and dealing with missing values are problems that must be solved before you can process the data further

Case Specific

Balancing.

Shuffling Database

Visualisation

E R Diagram

Relational Scema

Real Life Examples of Traditional Data

Customers

customer_id	first_name	last_name	email_address	number_of_complaints
1	John	McKinley	john.mackinley@365careers.com	0
2	Elizabeth	McFarlane	e.mcfarlane@365careers.com	2
3	Kevin	Lawrence	kevin.lawrence@365careers.com	1
4	Catherine	Winnfield	c.winnfield@365careers.com	0

Consider basic Customer data as example the difference between a numerical and categorical variable.

The first column shows the id of the different customers. These numbers however cannot be manipulated. Calculating an average ID is not something that would give you any sort of useful information. This means that even though they are numbers they hold no numerical value and therefore representing categorical data.

Now focus on the last column. This shows how many times that customers filed a complaint. These numbers are easily manipulated. Adding them all together to give a total number of complaints is useful information. Therefore they are numerical data.

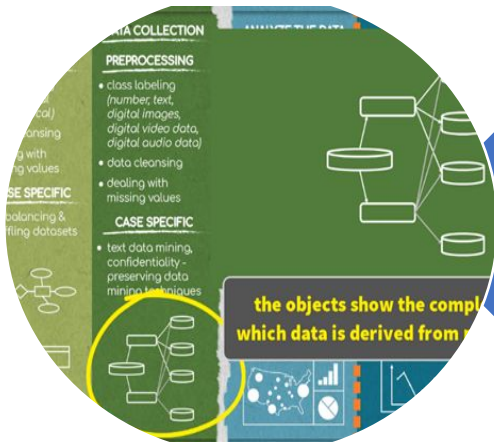
HISTORICAL STOCK PRICE DATA

	PG
Date	
2007-01-03	46.149067
2007-01-04	45.798710
2007-01-05	45.405422
2007-01-08	45.505543
2007-01-09	45.391144
2007-01-10	45.934578
2007-01-11	46.220585
2007-01-12	46.478012
2007-01-16	46.478012
2007-01-17	46.959393
2007-01-18	47.088707

Another example we can look at is daily historical stock price data

There's a column containing the dates of the observations which is considered categorical data and a column containing the stock prices which is numerical data.

Techniques for Working with Big Data



Some as traditional data preprocessing can also be implemented on big data is essential to help organize the data before doing analyses or making predictions as is grouping the data into classes or categories.

While working with Big Data things can get a little more complex. As you have much more variety beyond the simple distinction of numerical and categorical data. Examples of big data can be text data, digital image data, digital video data, digital audio data and more.

Consequently, with a larger amount of data types comes a wider range of data cleansing methods.

There are techniques that verify that a digital image observation is ready for processing and specific approaches exists that can ensure the audio quality of your file is adequate to proceed.

So what about dealing with missing values. This step is a crucial one as big data has big missing values which is a big problem to exemplify.

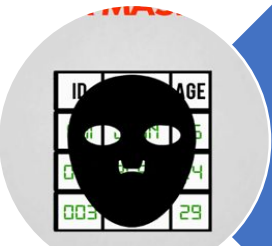


Text data mining represents the process of deriving valuable unstructured data from a text.

Consider you may have a database which has stored information from academic papers about marketing expenditure.

It may contain information from academic papers, blog, articles, online platforms, private Excel files and more.

This means you will need to extract marketing expenditure information from many sources. This technique can find the information you need without much of a problem.



Data masking If you want to maintain a credible business or governmental activity you must preserve confidential information.

However when personal information is shared online it doesn't mean that it can't be touched or used for analysis. Instead you must apply some data masking techniques so you can analyse the information without compromising private details like data shuffling.

Masking can be quite complex. It conceals the original data with random and false data allowing you to conduct analysis and keep all confidential information in a secure place. An example of applying data masking to big data is through what we called **confidentiality preserving data mining techniques**.

Real Life Examples of Big Data



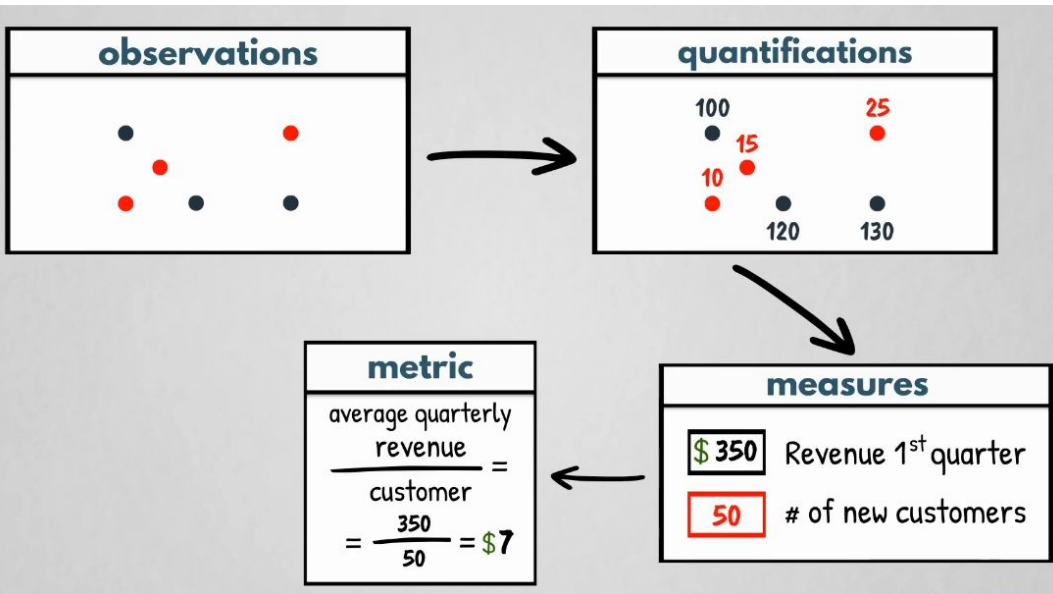
- **Facebook** keeps track of its users names, personal data, photos, videos, recorded messages and so on. This means that their data has a lot of variety in with over 2 billion users worldwide. The volume of data stored on their servers is tremendous.
- Facebook requires real time reporting of the aggregated anonymised voice of its users and it applies many analytical tools for its mobile applications.
- This means the company is investing in boosting its real time data processing powers increasing the velocity of its data set.



- Lets take **financial trading data** for example what happens when we record the stock price every five seconds or every single second.
- We get a data set that is incredibly voluminous requiring significantly more memory disk space in various techniques to extract meaningful information from it. Data like this would also be considered big data.

Business Intelligence (BI) Techniques

- Let's assume your data has been pre-processed and is ready for analysis. It is beautifully organized. This means you are ready to enter the realm of business intelligence.
- The job of a business intelligence analyst requires her to understand the essence of a business and strengthen that business through the power of data.
- So here we have techniques to measure business performance.



Collecting observations

From above diagram, you can observe variables such as sales volume (marked as blue colour dots) or new customers who have enrolled in your web site (marked as red colour dots).

Each monthly revenue or each customer is considered a single observation.

Quantification

However no mathematical manipulations can be applied to these observations. What we must do is quantify that information.

Quantification is the process of representing observations as numbers.

Consider your revenues from new customers for January, February and March were 100, 120 and \$130 respectively while the corresponding number of new customers for the same three months are 10, 15 and 25.

Measure

A measure is the accumulation of observations to show some information.

For example, if you total the revenues of all three months to obtain the value of \$350 that would be a measure of the revenue of the first quarter of that year.

Similarly add together the number of new customers for the same period and you have another measure

Metric

A metric refers to a value that derives from the measures you obtain and aims at gauging business performance or progress to compare.

If a measure is related to something like simple descriptive statistics of past performance a metric has a business meaning attached.

E.g. If you estimate the average quarterly revenue per new customer which equals 350 divided by 50 that is \$7. This is a metric.

KPI

In a real business where the number of observations is significantly larger you can derive thousands of metrics where we can't keep track of all possible metrics we can extract from a data set.

What you need to do is choose the metrics that are tightly aligned with your business objectives. These metrics are called K.P.I.'s, Key Performance Indicators.

Key because they are related to your main business goals. **Performance** because they show how successfully you have performed within a specified time frame and **indicators** because their values or metrics that indicate something related to your business performance.

Real Life Examples of Business Intelligence



- BI can be used for **price optimization**, hotels use price optimization very effectively by raising the price of a room at periods when many people want to visit the hotel. And by reducing it to attract visitors when demand is low they can greatly increase their profits in order to competently apply such a strategy.
- They must extract the relevant information in real time and compare it with historical BI allows you to adjust your strategy to pass data. As soon as it is available.



- Another application of business intelligence is enhancing **inventory management** over and undersupply can cause problems in a business.
- However, implementing effective inventory management means supplying enough stock to meet demand with the minimal amount of waste and cost to do this well.
- You can perform an in-depth analysis of past sales transactions for the purpose of identifying seasonality patterns and the times of the year with the highest cell's.
- Additionally, you could track your inventory to identify the months of which you have over or understocked a detailed analysis can even pinpoint the day or time of day were the need for a good is highest if done right business intelligence will help to efficiently manage your shipment logistics and in turn reduce costs and increase profit.

- So once prepared the BMI reports and dashboards are prepared and the executives have extracted insights about the business what do you do with the information you use it to predict some future values as accurately as possible. That's why at this stage you stop dealing with analysis and start applying analytics more precisely **predictive analytics**.
- We separate predictive analytics into two branches **traditional methods** which comprises classical statistical methods for forecasting and **machine learning**.

Techniques for Working with Traditional Methods

1) Regression

In business statistics, a regression is a model used for quantifying causal relationships among the different variables included in your analysis.

	House Price (\$)	House Size (sq.ft.)
0	1115000	1940
1	860000	1300
2	818400	1420
3	1000000	1680
4	640000	1270
5	1010000	1850
6	600000	1000
7	700000	1100
8	1100000	1600
9	570000	1000
10	860000	2150
11	1085000	1900
12	1250000	2200

In this dataset we have house prices in dollars while the other house sizes measured in square feet .

Every row on the data table is an observation and each can be plotted on this graph as a dot the house size is measured along the horizontal line in its price.

On the vertical line the further to the right an observation is the larger the house size and further up the higher the price so once we've plotted all 20 observations from our dataset our graph will appear like this.

The thing is there is a straight line(red) called a regression line that goes through these dots while being as close as it can be to all of them simultaneously.

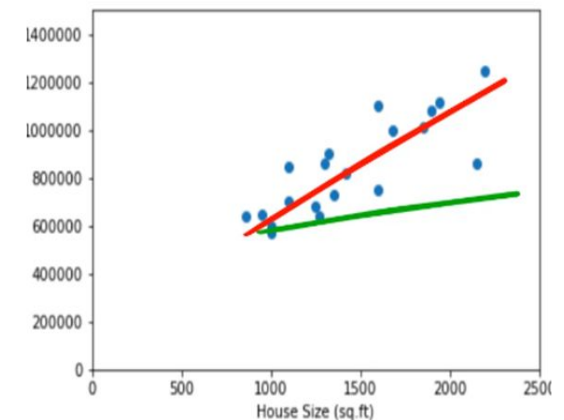
Now imagine we drew another line(green). If you observe altogether dots are closer to the first red line than the second green one. This means that it more accurately represents the distribution of the observations.

So in this case if y signifies the house price then B represents a coefficient which we multiply by x the house size.

$$y = Bx$$

So the equation professionals work with is Y equals B times X and they use the graph as visual support.

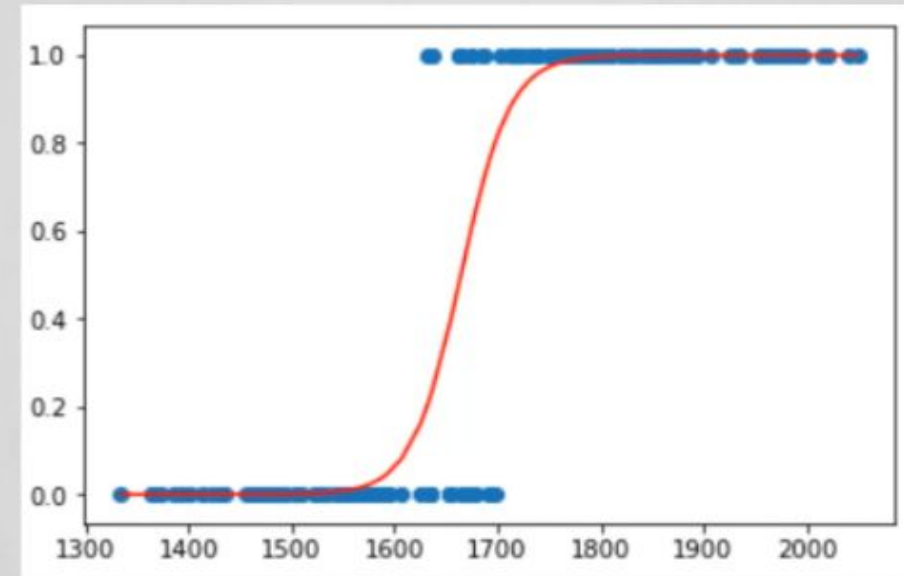
A) Linear Regression



B) Nonlinear Regression – Logistic Regression

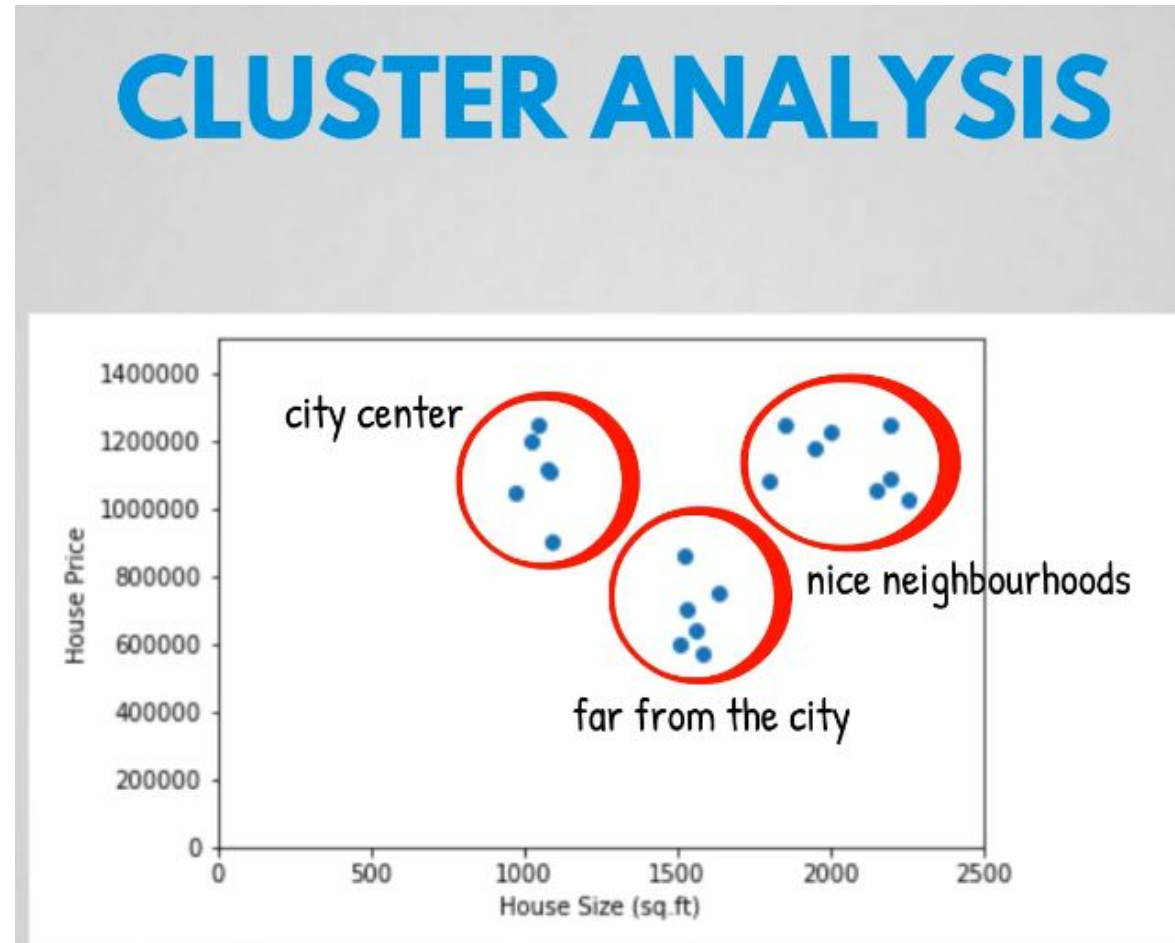
- A **logistic regression** is a common example of a nonlinear model.
- In this case unlike the house prices example the values on the vertical line won't be arbitrary integers.
- They'll be ones or zeros only such a model's useful during a decision-making process.
- Companies apply logistic regression algorithms to filter job candidates during their screening process.
- If the algorithm estimates the probability that a prospective candidate will perform well and the company is above 50 percent it would predict one or a successful application.
- Otherwise it will predict zero therefore the nonlinear nature of the logistic regression is nicely summarized by its Graph very different from the linear regression.

LOGISTIC REGRESSION



2) Cluster Analysis

- Imagine they are derived from research on German house prices. Hence they are dispersed differently.
- When the data is divided into a few groups called clusters you can apply **cluster analysis**.
- This is another technique that will take into account that certain observations exhibit similar house sizes and prices.
- For instance this cluster of observations denotes small houses but with a high price. This could be typical for houses in the city centre.
- The second cluster could represent houses that are far from the city because they are quite big but cost less.
- Finally, the last cluster concerns houses that are probably not in the city centre but are still in nice neighbourhoods. They are big and cost a lot.
- Noticing that **your data can be clustered is important so you can improve your further analysis**. In our example clustering allowed us to conclude that location is a significant factor when pricing a house.



3) Factor Analysis

What about a more complicated study where you consider explanatory variables apart from house size. You might have quantified the location, number of rooms, years of construction and so on which can all affect house price then when thinking about the mathematical expression corresponding to the regression model you won't just have one explanatory variable X you will have x_1, x_2, x_3 and so on. Note that an explanatory variable can also be called a regressor or an independent variable or a predictor variable.

Imagine analysing a survey that consists of 100 questions performing any analysis on 100 different variables is tough. This means you are variables starting from x1 and going all the way up to x100. The good thing is that often different questions are measuring the same issue and this is where **factor analysis** comes in

Assume your survey contained this question on a scale from 1 to 5. How much do you agree with the following statements.

Survey:

1. I like animals
2. I care about animals.
3. I am against animal cruelty.

People are likely to respond consistently to these three questions.

That is however marks five to the first question does the same for the second and third questions as well. In other words, if you strongly agree with one of these three statements you won't disagree with the other two right.

With factor analysis we can combine all the questions into general attitude towards animals. So instead of three variables we now have one in a similar manner.

You can reduce the dimensionality of the problem from 100 variables to 10 which can be used for a regression that will deliver a more accurate prediction.

To sum up we can say that **clustering is about grouping observations together and factor analysis is about grouping explanatory variables together.**

	House Price (\$)	House Size (sq.ft.)	State	Number of Rooms	Year of Construction
0	1116000	1940	IN	8	2002
1	860000	1300	IN	5	1992
2	818400	1420	IN	6	1987
3	1000000	1680	IN	7	2000
4	640000	1270	IN	5	1995
5	1010000	1850	IN	7	1998
6	600000	1000	IN	4	2015
7	700000	1100	LA	4	2014
8	1100000	1600	LA	7	2017
9	570000	1000	NY	5	1997
10	860000	2150	NY	9	1997
11	1085000	1900	NY	9	2000
12	1250000	2200	NY	9	2014

$$y = a + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_n x_n$$

x: explanatory variable

= regressor

= independent variable

= predictor variable

[illegible]

4) Time Series Analysis

You will use this technique especially if you are working in economics or finance in these fields you will have to follow the development of certain values over time such as stock prices or sales volume you can associate **time series** with plotting values against time.

Time will always be on the horizontal line as time is independent of any other variable therefore such a graph can end up depicting a few lines that illustrate the behaviour of your stocks over time.

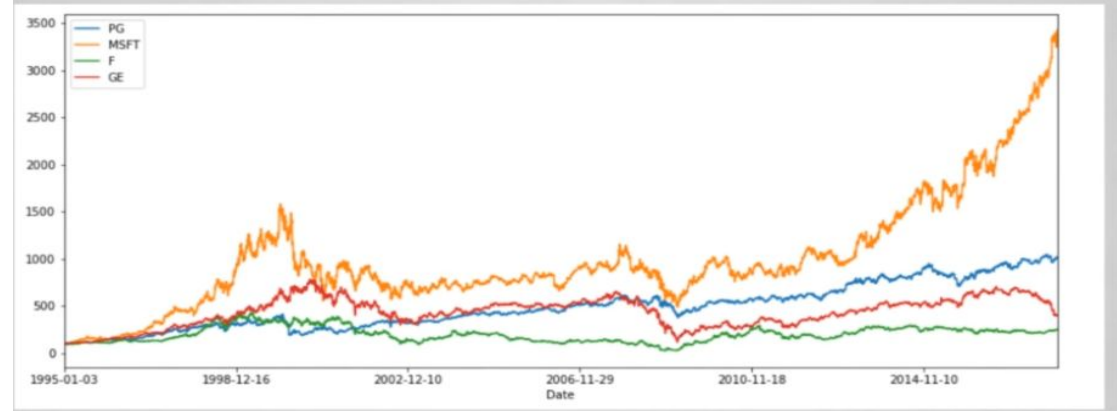
So when you study the visualization you can spot which stock performed well and which did not. We must admit there is a vast variety of methods that professionals can choose from.

TIME SERIES



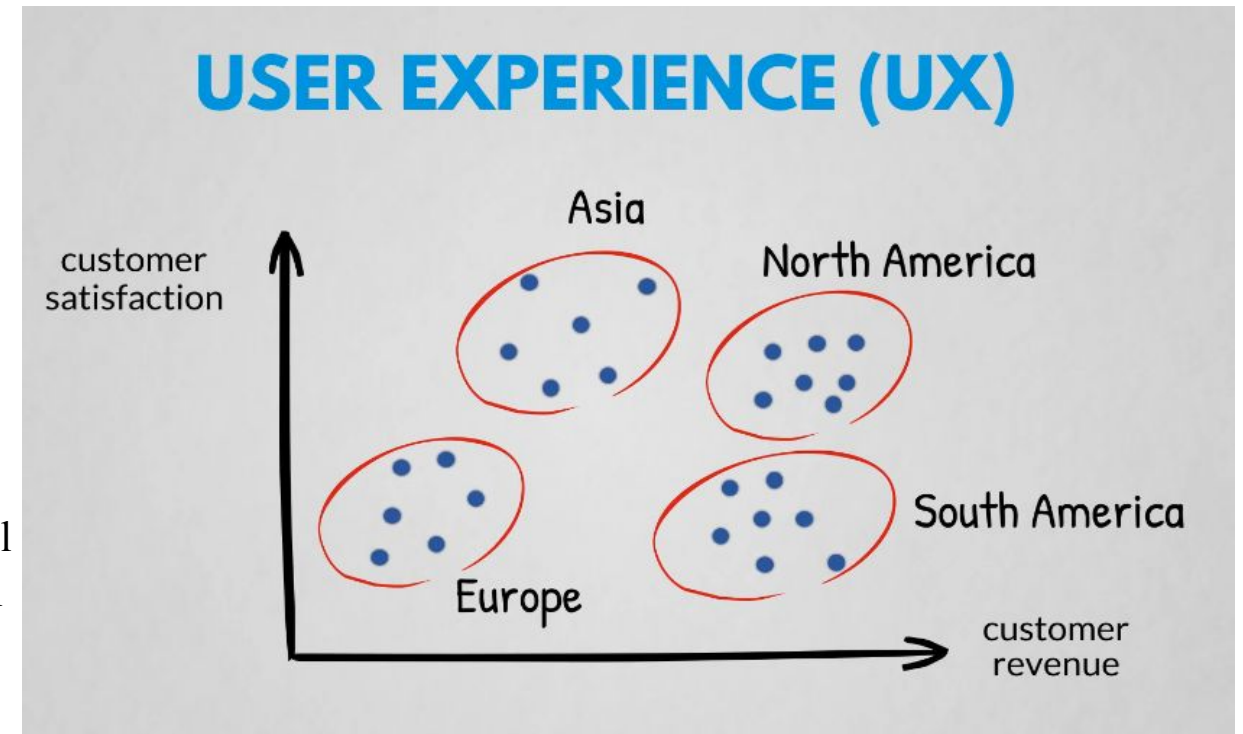
- stock price
- sales volume

TIME SERIES



Real Life Examples of Traditional Methods

- Imagine you are the head of the user experience department of a web site selling goods on a global scale which we often abbreviate as UX.
- So what is your goal as head of UX to maximize user satisfaction right.
- Assume you have already designed and implemented a survey that measures the attitude of your customers towards the latest global product.
- You have launched the graph where you plot your observations will likely appear in the following way when the data is concentrated in such a way.
- You should cluster the observations.
- Remember So when you perform cluster analysis you will find that each cluster represents a different continent.
- This group may refer to the responses gathered from Asia or Europe or South America or from North America.
- Once you realize there are four distinct groups it makes sense to run four separate tests. Obviously, the difference between clusters is too great for us to make a general conclusion.
- Asians may enjoy using your web site one way while Europeans in another. Thus, it would be sensible to adjust your strategy for each of these groups individually.
- Another noteworthy example we can give you is forecasting sales volume every business and financial company does this.
- So, which traditional statistical technique that we discussed would fit the picture here. Time series analysis it is say this was your data until a certain date. What will happen next? How should you expect the cells to be for the year ahead? Will their volume increase or decrease?
- Several types of mathematical and statistical models allow you to run multiple simulations which could provide you with future scenarios based on these scenarios. You can make better predictions and implement adequate strategies awesome job people. You are already acquainted with many of the data science essential terms but not all.

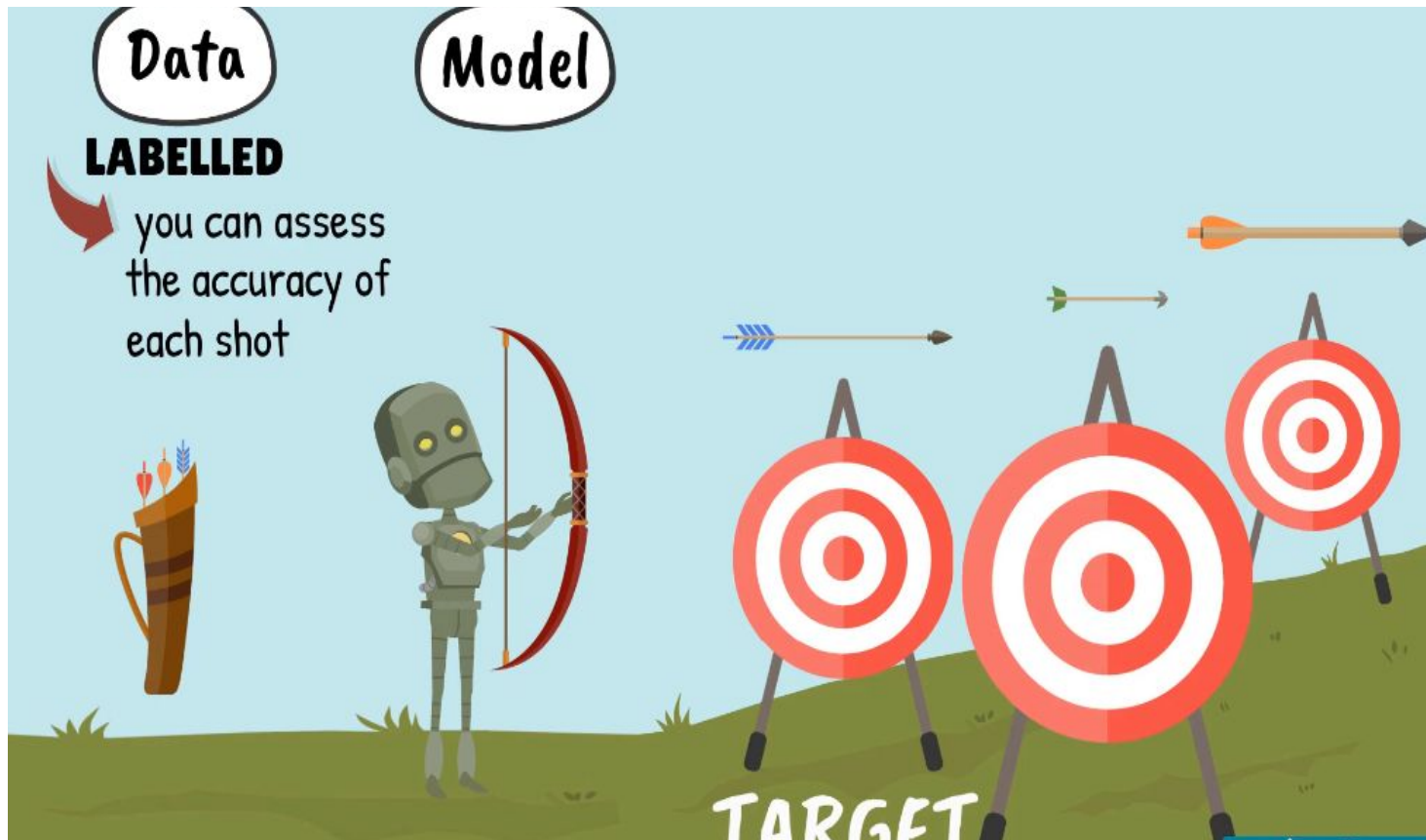


Machine Learning (ML) Techniques

- The core of machine learning is creating an algorithm which a computer then uses to find a model that fits the data as best as possible and makes very accurate predictions based on that and how is that different from conventional methods. We provided with algorithms which give the machine directions on how to learn on its own.
- A machine learning algorithm is like a trial-and-error process. Each consecutive trial is at least as good as the previous one.
- Technically speaking there are four ingredients data, model, objective function and optimization algorithm.
- Example. Imagine a robot holding a bow. We want to find the best way to use that bow to fire accurately. In other words the usage of the bow is our *model*, the best way to learn archery is to train right. We train by taking different arrows and trying to hit the target. So, the quiver of arrows will be or *data* or more precisely the data that the robot will use for training.
- They are all arrows but they have their subtleties. There are straight ones, crooked ones, light ones, heavy ones. So we can safely say the arrows represent different data values.
- We said the robot will be firing at a target. In machine learning or at least in the most common type supervised learning, we know what we are aiming for and we call it a target.
- The *objective function* will calculate how far from the target the robot shots were on average.
- Here comes the fourth ingredient the optimization algorithm.
- It steps on the findings of the objective function and consists of the mechanics that will improve the robot's archery skills somehow. It's posture the way it holds the bow how strong it pulls the bowstring etc. Then the robot will take the exact same data or arrows and fire them once again with its adjusted posture.
- This time the shots will be on average closer to the centre of the target. Normally the improvement will be almost unnoticeable. This entire process could have been hundreds or thousands of times until the robot finds the optimal way to fire this set of arrows and hit the centre every single time.

- Nevertheless, **it is important to remember that while training you won't provide the robot with a set of rules that is you won't have programmed a set of instructions like place the arrow in the middle of the bow pull the bow string and so on.**
- Instead, you will have given the machine a final goal to place the arrow in the centre of the target.
- So you don't care if it places the arrow in the middle or in the bottom of the bow as long as it hits the target.
- **Another important thing is that it won't learn to shoot well right away but after a hundred thousand tries it may have learned how to be the best archer out there.**
- Now there can be infinite possibilities to trial, when will the robots stop training first.
- The robot will learn certain things on the way and will take them into consideration for the next shots at fires for instance if it learns that it must look towards the target it will stop firing in the opposite direction.
- That is the purpose of the [optimization algorithm](#). Second it cannot fire arrows forever.
- However, hitting the centre nine out of 10 times may be good enough. So, we can choose to stop it after it reaches a certain level of accuracy or fires a certain number of arrows.
- So, let us follow the four ingredients at the end of the training. Our robot or model is already trained on this data. With this set of arrows most shots hit the centre so the air or the objective function is quite low or minimized as we like to say the posture the technique and all other factors cannot be improved.
- So, the optimization algorithm has done its best to improve the shooting ability of the machine.
- We own a robot that is an amazing Archer. So, what can you do? Give it a different bag of arrows. If they had seen most types of arrows while training it will do great with the new ones.
- However, if we give it half an arrow or a longer arrow than it has seen it will not know what to do with it.
- In all ordinary cases though we would expect the robot to hit the centre or at least get close.
- The benefit of using machine learning is that the robot can learn to fire more effectively than a human.
- It might even discover that we've been holding bows in a wrong way for centuries to conclude we must say that machine learning is not about robots.

Types of Machine Learning

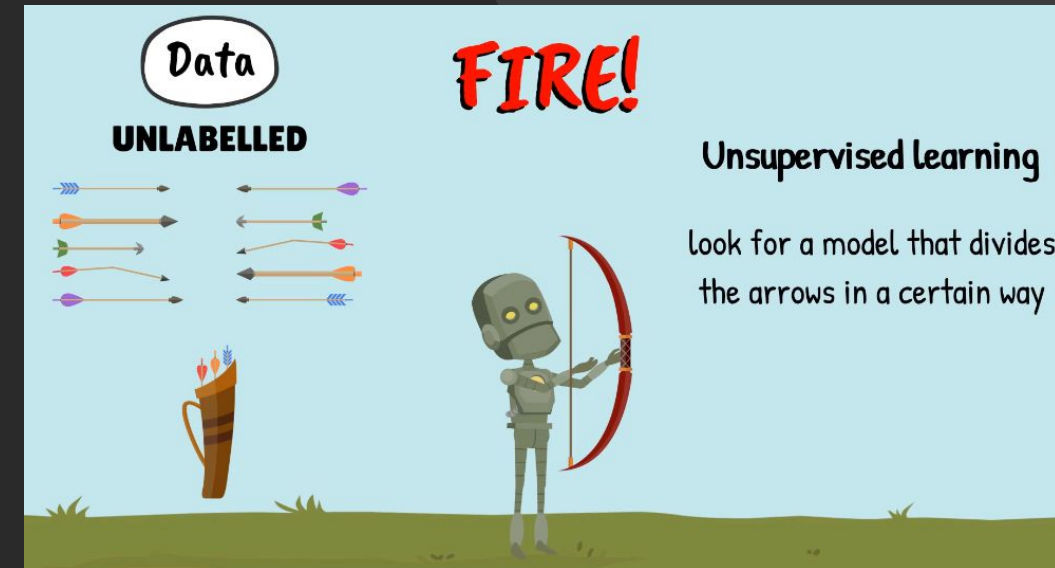


1) Supervised Machine learning

- This name derives from the fact that training an algorithm resembles a teacher supervising her students.
- In Supervised machine learning, it is important to mention you have been dealing with label data. In other words, you can assess the accuracy of each shot.
- Consider previous example, where there isn't a single target different arrows have their own targets.
- Let's check what the robot sees when shooting the ground, a target at a short distance a target at a further distance a target hanging on a tree far behind it a house to the side and the sky.
- *So, having labelled data means the associating or labelling a target to a type of Arrow.*
- You know that with a small arrow the robot is supposed to hit the closest target with a medium arrow it can reach the target located further away while with a larger arrow the target that's hanging on the tree. Finally, a crooked arrow is expected to hit the ground not reaching any target during the training process.
- The robot will be shooting arrows at the respective targets as well as it can. After training is finished.
- Ideally the robot will be able to fire the small arrow at the centre of the closest target the middle arrow at the centre of the one further away and so on.
- To summarize label data means we know the target prior to the shot, and we can associate that shot with the target this way. We're sure where the arrow should hit.
- This allows us to measure the inaccuracy of the shot through the objective function and improve the way the robot shoots through the optimization algorithm. So, what we supervise is the training itself. If a shot is far off from its target, we correct the posture. Otherwise, we don't get.

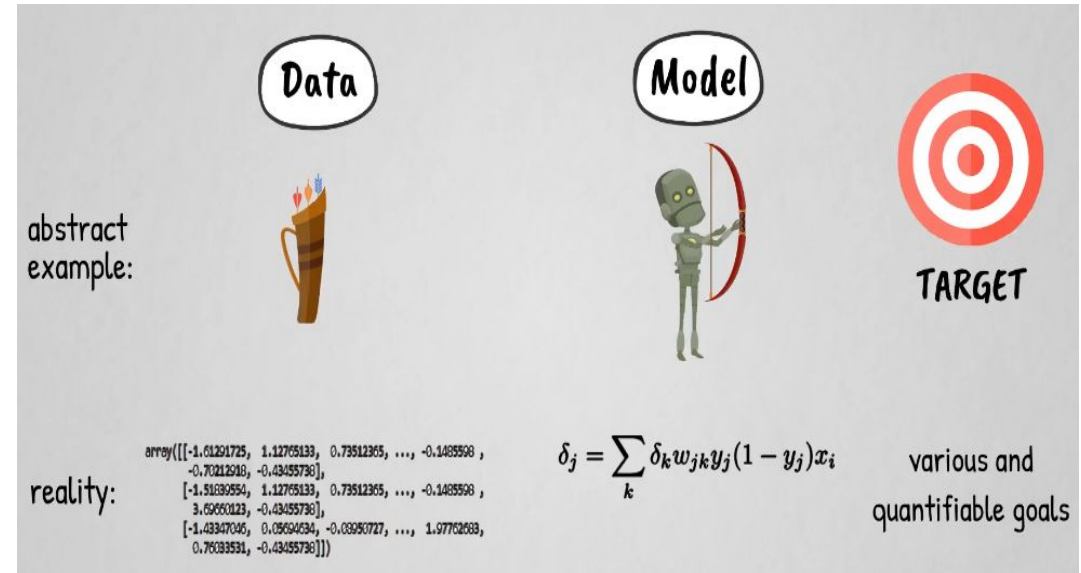
2) Unsupervised Machine learning

- In practice though it might happen that you won't have the time or the resources to associate the arrows with targets before giving them to the robot.
- In that case you could apply the other major type of M-L unsupervised learning here you will just give your robot a bag of arrows with unknown physical properties unlabelled data. This means neither you nor the robot will have separated the arrows into groups.
- Then you'd ask the machine to simply fire in a direction without providing it with targets. Therefore, in this case you won't be looking for a model that helps you shoot better rather you'll be looking for one which divides the arrows in a certain way.
- The robot will see just the ground the tree the House and the sky. Remember there are no targets. So, after firing thousands of shots during the training process we will end up having different types of arrows stuck in different areas.
- For instance, you may identify all the broken arrows by noticing they have fallen on the ground nearby the others you may realise are divided into small medium and large arrows.
- There may be anomalies like crossbow bolts in your bag that after being shot may have accumulated in a pile over here.
- You wouldn't want to use them with a simple bow would you. At the end of the training the robot will have fired so many times that it could discover answers that may surprise you.
- The machine may have managed to split the arrows not into four but into five sized categories due to discovering the crossbow bolt. Or it may have identified that some arrows are going to break soon by placing them in the Broken Arrow pile.
- It is worth mentioning that supervised learning can deal with such problems too and it does very often. However, if you have one million arrows you don't really have the time to assign targets to all of them do you.
- To save time and resources you should apply unsupervised learning.



3) Reinforcement learning

- The third major type of machine learning is called reinforcement learning. This time we introduce a reward system.
- Every time the robot fires an arrow better than before it will receive an award say a chocolate it will receive nothing if it fires worse.
- So instead of minimizing an error we are maximizing a reward or in other words maximizing the objective function.
- If you put yourselves in the shoes of the machine, you'll be reasoning in the following way. I fire an arrow and receive a reward. I'll try to figure out what I did correctly.
- So, I get more chocolate with the next shot or I fire an arrow and don't receive a reward. There must be something I need to improve.
- For me to get some chocolate on my next shot positive reinforcement..



- In addition, don't forget the robot Archer was an abstract depiction of what a machine learning model can do.
- In reality there are robots, but the model will be a highly complex mathematical formula the arrows will be a data set and the goals will be various and quantifiable
- Here are the most notable approaches you will encounter when talking about machine learning support vector machines neural networks deep learning random forest models and Bazy and networks are all types of supervised learning.
- There are neural networks that can be applied to an unsupervised type of machine learning, but K means is the most common unsupervised approach.
- By the way you may have noticed we have placed deep learning in both categories.
- This is a relatively new revolutionary computational approach which is acclaimed as the State-of-the-art email today.
- Describing it briefly we can say it is fundamentally different from the other approaches.
- However, it has a broad practical scope of application in all M-L areas because of the extremely high accuracy of its models.
- Note that deep learning is still divided and supervised, unsupervised and reinforcement, so it solves the same problems but in a conceptually different way.

Real Life Examples of Machine Learning (ML)



The *financial sector and banks* have ginormous data sets of credit card transactions. Unfortunately, banks are facing issues with fraud daily. They are tasked with preventing fraudsters from acquiring customer data and in order to keep customers funds safe they use machine learning algorithms. They take past data and because they can tell the computer which transactions in their history were legitimate and which were found to be fraudulent, they can label the data as such. So through supervised learning they train models that detect fraudulent activity when these models detect even the slightest probability of theft. They flagged the transactions and prevent the fraud in real time. Although no one in the sector has reached a perfect solution.

Another example of using supervise machine learning with label data can be found in *client retention*.

A focus of any business be it a global supermarket chain or an online clothing shop is to retain its customers.

But the larger a business grows the harder it is to keep track of customer trends. A local corner shop owner will recognize and get to know their most loyal customers. They will offer them exclusive discounts to thank them for their custom.

And by doing so keep them returning on a larger scale. Companies can use machine learning and past label data to automate the practice.

And with this they can know which customers may purchase goods from them. This means the store can offer discounts and a personal touch in an efficient way minimizing marketing costs and maximizing profits.





Popular Data Science Tools

DATA VISUALIZATION

Understanding and Visualizing Data

- Data visualization is **the process of translating large data sets and metrics into charts, graphs and other visuals**. The resulting visual representation of data makes it easier to identify and share real-time trends, outliers, and new insights about the information represented in the data.
- Data visualization is the graphical representation of information and data. By using visual elements like charts, graphs, and maps, data visualization tools provide an accessible way to see and understand trends, outliers, and patterns in data.
- In the world of Big Data, data visualization tools and technologies are essential to analyze massive amounts of information and make data-driven decisions.

The advantages and benefits of good data visualization

- Our eyes are drawn to colors and patterns. We can quickly identify red from blue, square from circle. Our culture is visual, including everything from art and advertisements to TV and movies.
- Data visualization is another form of visual art that grabs our interest and keeps our eyes on the message. When we see a chart, we quickly see trends and outliers.
- It's storytelling with a purpose. If you've ever stared at a massive spreadsheet of data and couldn't see a trend, you know how much more effective a visualization can be.
- As the "age of Big Data" kicks into high-gear, visualization is an increasingly key tool to make sense of the trillions of rows of data generated every day.
- Data visualization helps to tell stories by curating data into a form easier to understand, highlighting the trends and outliers.
- A good visualization tells a story, removing the noise from data and highlighting the useful information. However, it's not simply as easy as just dressing up a graph to make it look better or slapping on the "info" part of an infographic.
- Effective data visualization is a delicate balancing act between form and function. The plainest graph could be too boring to catch any notice or it make tell a powerful point; the most stunning visualization could utterly fail at conveying the right message or it could speak volumes.
- The data and the visuals need to work together, and there's an art to combining great analysis with great storytelling.



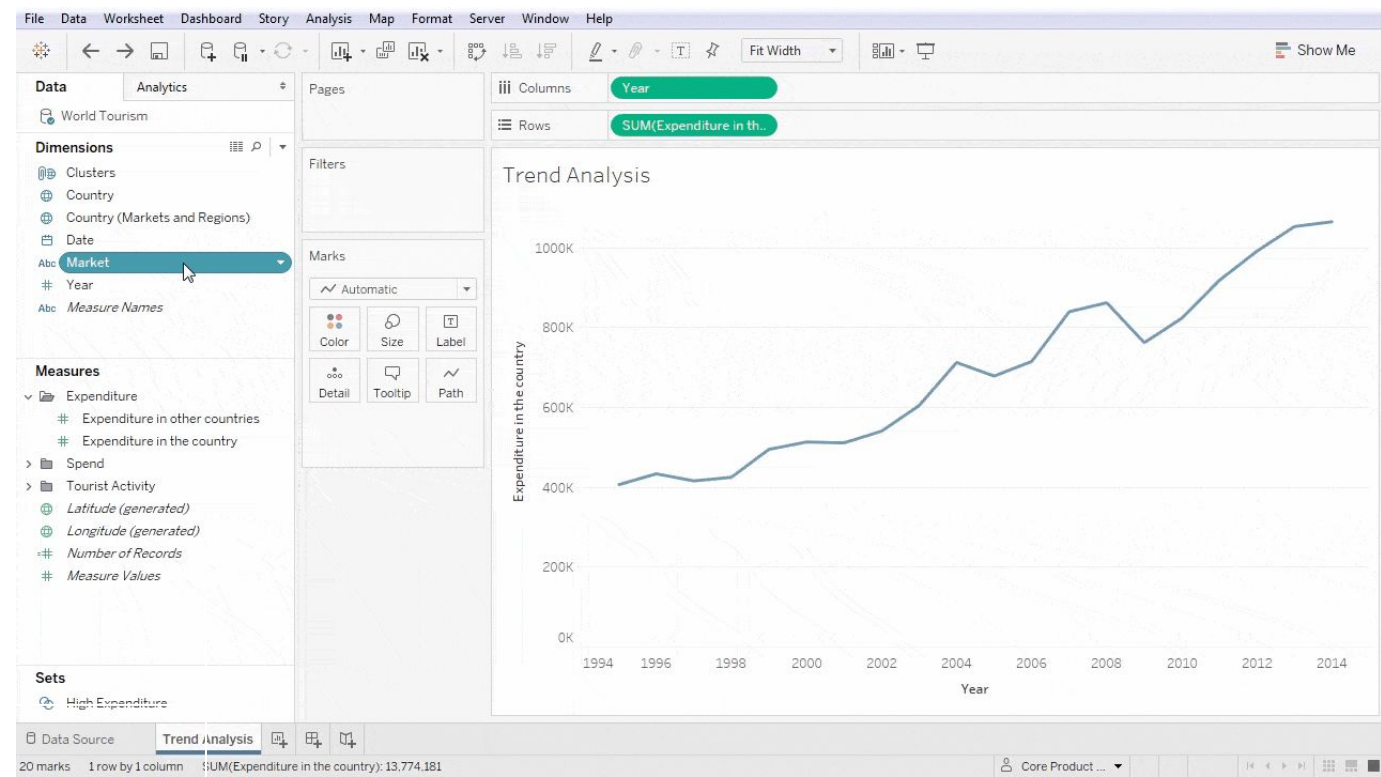
THE DIFFERENT TYPES OF VISUALIZATIONS

When you think of data visualization, your first thought probably immediately goes to simple bar graphs or pie charts. While these may be an integral part of visualizing data and a common baseline for many data graphics, the right visualization must be paired with the right set of information. Simple graphs are only the tip of the iceberg. There's a whole selection of visualization methods to present data in effective and interesting ways. **Common general types of data visualization:**

- Charts
- Tables
- Graphs
- Maps
- Infographics
- Dashboards

More specific examples of methods to visualize data:

- | | | |
|-------------------------|---------------------|---------------------------|
| • Box-and-whisker Plots | • Area Chart | • Bubble Cloud |
| • Histogram | • Matrix | • Bullet Graph |
| • Cartogram | • Circle View | • Dot Distribution Map |
| • Gantt Chart | • Heat Map | • Highlight Table |
| • Network | • Radial Tree | • Scatter Plot (2D or 3D) |
| • Streamgraph | • Text Tables | • Timeline |
| • Treemap | • Wedge Stack Graph | • Word Cloud |



Business Intelligence

What is BI?

- Business intelligence (BI) is a technology-driven process for analyzing data and delivering actionable information that helps executives, managers and workers make informed business decisions.
- As part of the BI process, organizations collect data from internal IT systems and external sources, prepare it for analysis, run queries against the data and create data visualizations, BI dashboards and reports to make the analytics results available to business users for operational decision-making and strategic planning.
- The ultimate goal of BI initiatives is to drive better business decisions that enable organizations to increase revenue, improve operational efficiency and gain competitive advantages over business rivals.
- To achieve that goal, BI incorporates a combination of analytics, data management and reporting tools, plus various methodologies for managing and analyzing data.

Business Intelligence features

Business intelligence can help companies make better decisions by showing present and historical data within their business context. Analysts can leverage BI to provide performance and competitor benchmarks to make the organization run smoother and more efficiently. Analysts can also more easily spot market trends to increase sales or revenue. Used effectively, the right data can help with anything from compliance to hiring efforts. **A few ways that business intelligence can help companies make smarter, data-driven decisions:**

- Identify ways to increase profit
- Analyze customer behavior
- Compare data with competitors
- Track performance
- Optimize operations
- Predict success
- Spot market trends
- Discover issues or problems

Data Visualization

Data visualization is the technique to present the data in a pictorial or graphical format.



Data Visualization



You are a Sales Manager in a leading global organization. The organization plans to study the sales details of each product across all regions and countries. This is to identify the product which has the highest sales in a particular region and up the production. This research will enable the organization to increase the manufacturing of that product in the particular region.

Data Visualization



Data Visualization Considerations

Three major considerations for data visualization are:



Clarity



Accuracy



Efficiency

Ensure the dataset is complete and relevant. This enables the Data Scientist to use the new patterns yield from the data in the relevant places.

Data Visualization Considerations

Three major considerations for data visualization are:



Clarity



Accuracy



Efficiency

Ensure using appropriate graphical representation to convey the right message.

Data Visualization Considerations

Three major considerations for data visualization are:



Clarity



Accuracy



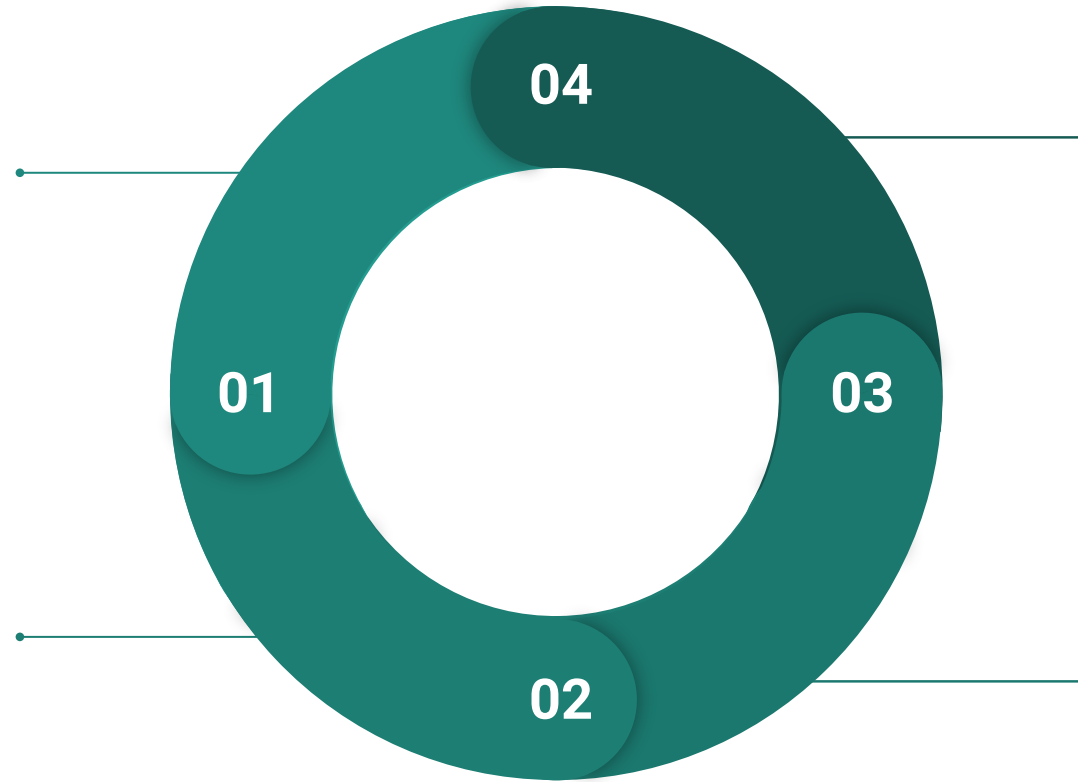
Efficiency

Use efficient visualization technique which highlights all the data points.

DATA VISUALIZATION FACTORS

The visual effect includes the usage of appropriate shapes, colours, and size to represent the analyzed data

The coordinate system helps to organize the data points within the provided coordinates.



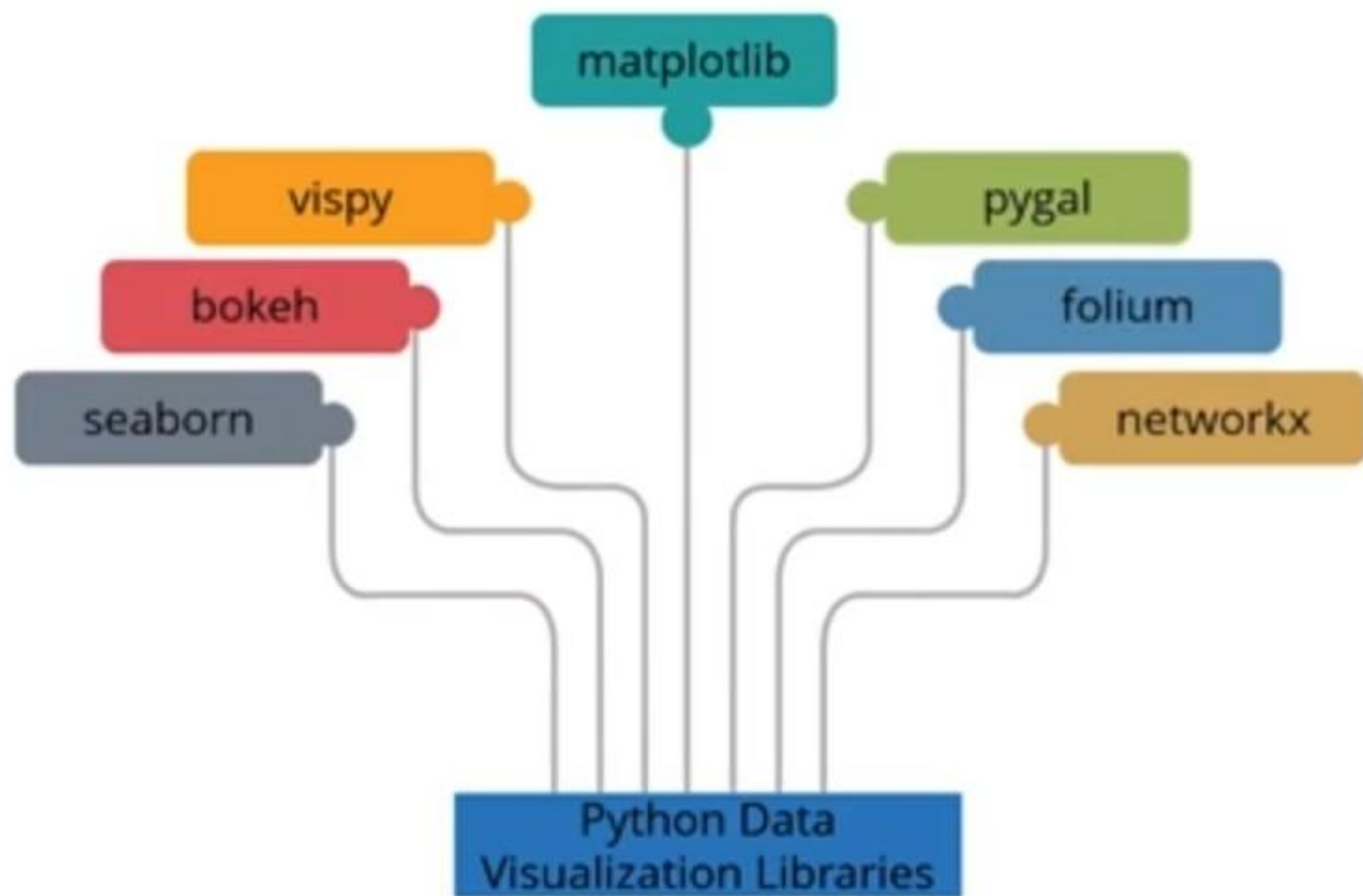
The informative interpretation helps to create visuals in an effective and easily interpretable manner using labels, title, legends and pointers.

The data types and scale choose the type of data such as numeric or categorical.

There are some basic factors that one needs to be aware of before visualizing the data.

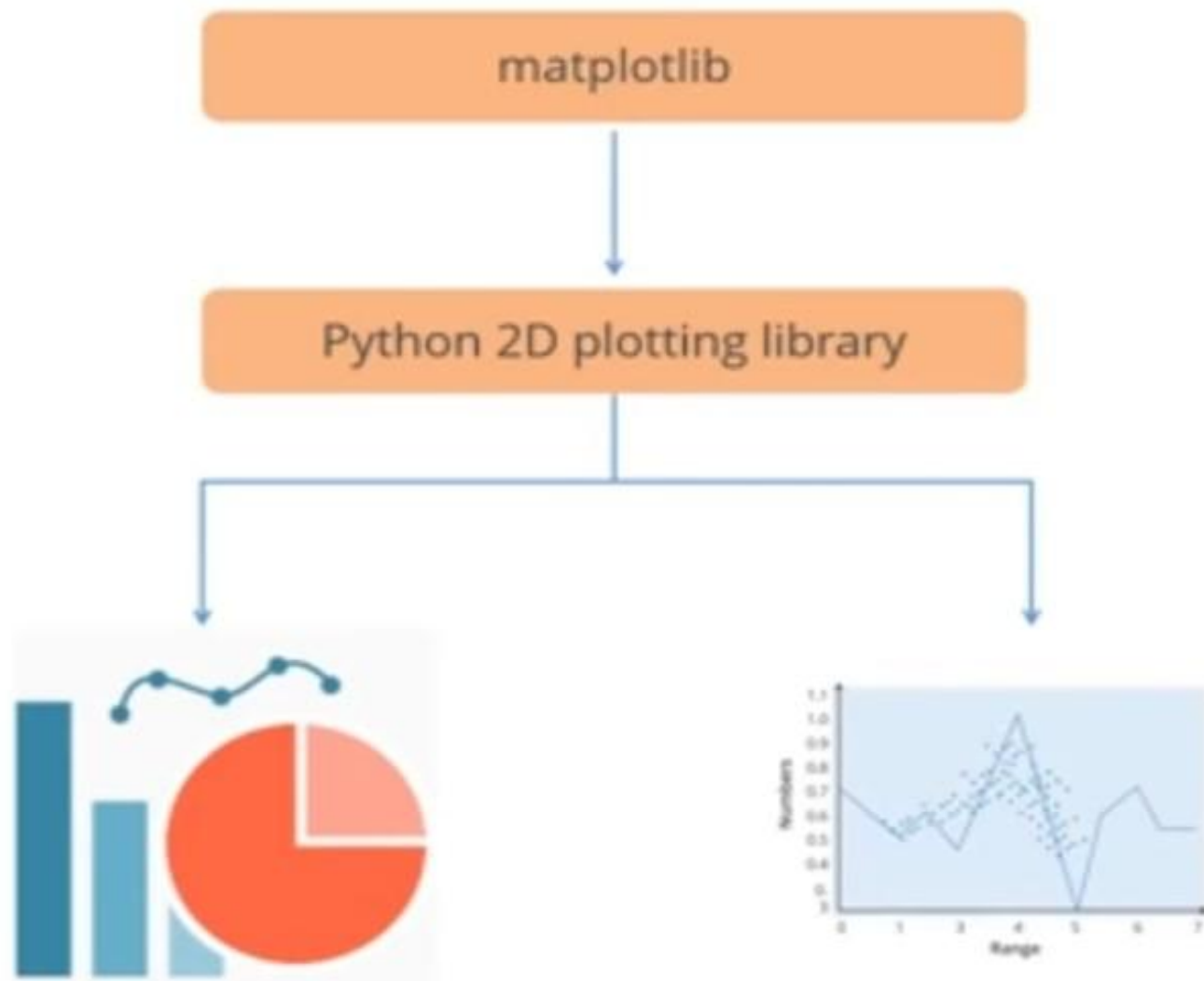
Python Libraries

Many new Python data visualization libraries are introduced recently such as:



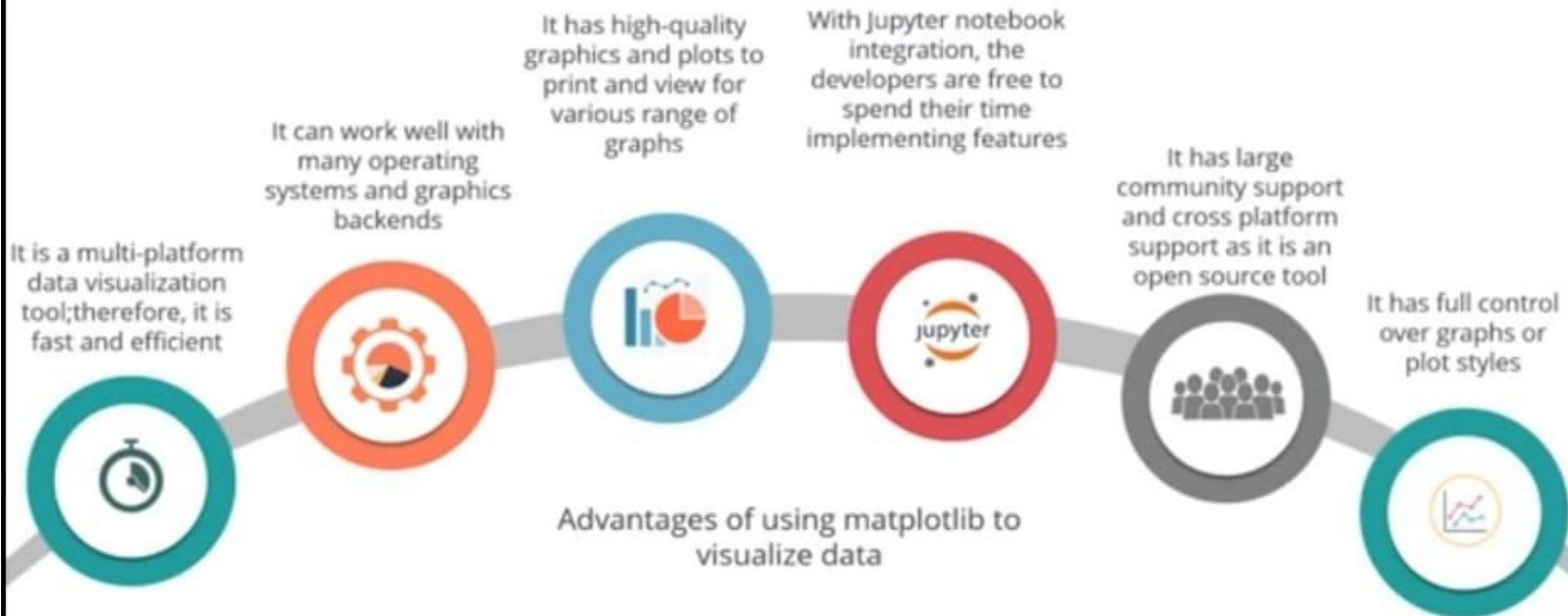
Python Libraries—matplotlib

Using Python's matplotlib, the data visualization of large and complex data becomes easy.



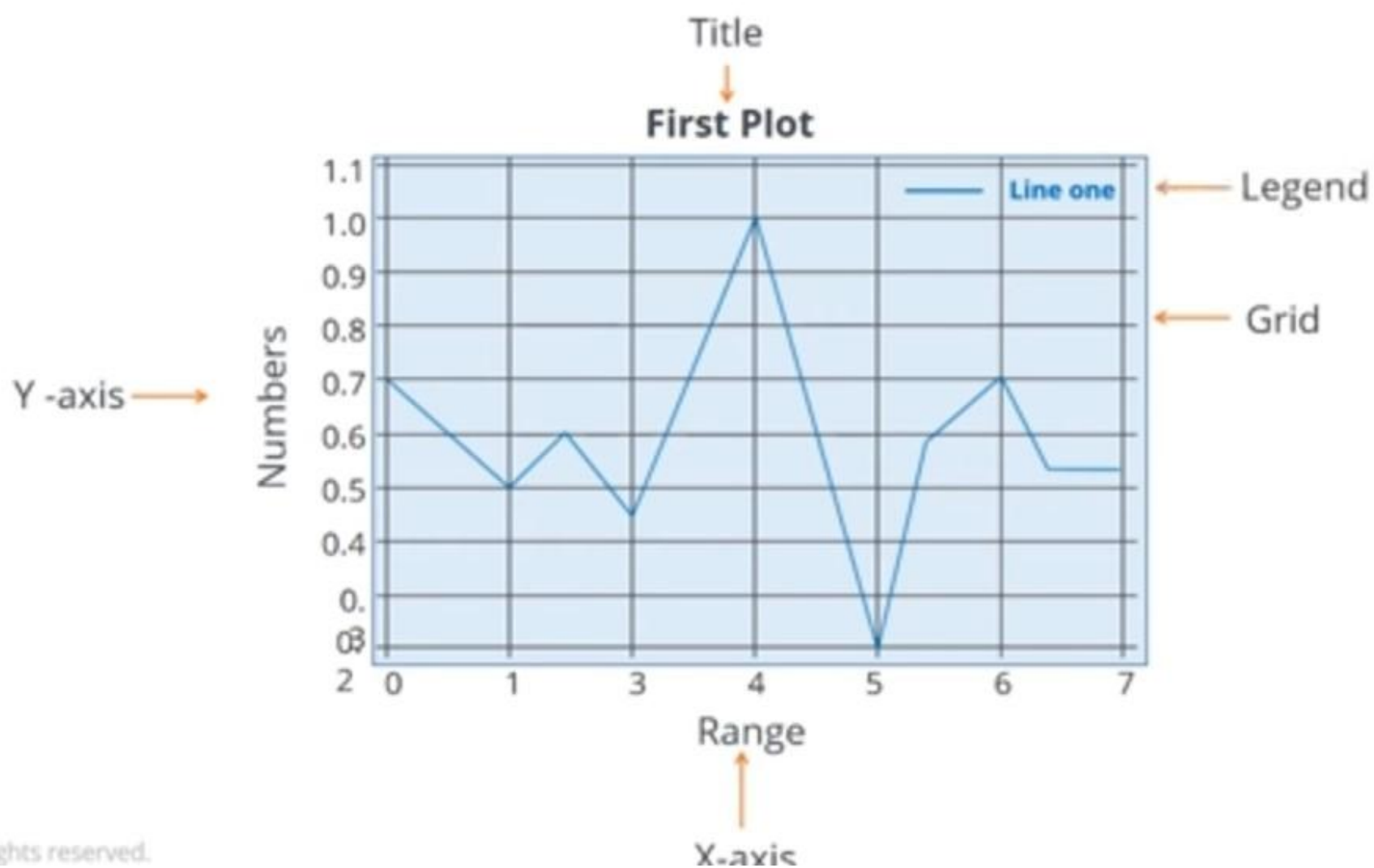
Python Libraries—matplotlib

There are several advantages of using matplotlib to visualize data. They are as follows:



Understanding the Plot

A plot is a graphical representation of data which shows relationship between two variables or the distribution of data.



Steps to Create a Plot

You can create a plot using four simple steps.



