

**JYOTI NIVAS COLLEGE AUTONOMOUS  
POST GRADUATE CENTER**



**DEPARTMENT OF MCA**

**E JOURNAL ON  
DATA ANALYTICS**

**ISSUE: 5**

**Jan 2023**

<b>SL No</b>	<b>TITLE</b>	<b>PAGE NO</b>
1	DATA ANALYTICS IN IOT	1
2	CHI-SQUARE AND ANOVA TESTS	2
3	SPATIAL ANALYTICS FOR GIS DATA	4
4	TABLEAU FOR ANALYTICS PLATFORM	6
5	NATURAL LANGUAGE PROCESSING TOOL KIT	8
6	DATA SCRAPING	9
7	IMPUTATION TECHNIQUES IN DATA ANALYTICS	11
8	REDUCING MANUFACTURING FAILURES USING DATA ANALYTICS	13
9	APACHE SPARK	14
10	SQL FOR DATA ANALYTICS	16
11	FUNNEL ANALYSIS IN MARKETING	18
12	HYPOTHESIS TESTING STATISTICAL ANALYSIS	21
13	MOVIE REVIEW SENTIMENT ANALYSIS	23
14	DATA ANALYTICS USING QUANTUM COMPUTING	24
15	SMART HEALTHCARE SYSTEMS LEVERAGING BIG DATA ANALYTICS	26
16	WEB ANALYTICS FOR DIGITAL MARKETING	28
17	DATA ANALYTICS IN AGILE DATA SCIENCE	30
18	R PROGRAMMING FOR DATA SCIENCE	32
19	SOCIAL MEDIA ANALYTICS	33
20	CORRELATION AND REGRESSION	35
21	AI IN DATA ANALYTICS	38

22	KNIME IN DATA ANALYTICS	40
----	-------------------------	----

## DATA ANALYTICS IN IOT

AMRITHA BALAKRISHNAN (21MCA02)

FEBA BIJU (21MCA16)

Data and IoT continue to be closely related. Data production and consumption are growing exponentially. Numerous IoT-based apps are in use across numerous industries and have been enormously beneficial to their customers. Data analytics come into play since the analysis of the IoT device data is what makes it valuable in the first place. Data analytics (DA) is a procedure used to analyse large and small data sets with a variety of data attributes in order to draw meaningful conclusions and useful information. These findings, which are typically presented as trends, patterns, and statistics, support corporate organizations in actively utilizing data to put into place efficient decision-making procedures.

### **Merging Data Analytics and IoT will Positively Impact Businesses**

Data Analytics has a significant role to play in the growth and success of IoT applications and investments. Analytics tools will allow the business units to make effective use of their datasets. The use of data analytics in IoT investments will allow the business units to gain an insight into customer preferences and choices. This would lead to the development of services and offers as per the customer demands and expectations. This, in turn, will improve the revenues and profits earned by the organizations.

In order to benefit from IoT investments, various forms of data analytics can be leveraged and applied. The list and descriptions of a few of these types are provided below.

**Streaming Analytics:** This form of data analytics is also referred as event stream processing and it analyses huge in-motion data sets. Real-time data streams are evaluated to find urgent problems and take prompt action. This technique can be used for Internet of Things applications based on financial transactions, air fleet tracking, traffic analysis, etc.

**Spatial Analytics:** examines geographic patterns to ascertain the geographical relationships between actual things. This kind of data analytics can be useful for location-based IoT applications, including smart parking apps.

**Time Series Analytics:** This type of data analytics is based upon time-based data that is examined to identify related trends and patterns. This type of data analytics technique can be useful for Internet of Things applications like weather forecasting software and health monitoring devices.

**Predictive Analysis:** This form of data analytics combines descriptive analysis and forecasting. It is used to determine the optimum course of action that can be followed in a specific circumstance. This type of data analytics can be used by commercial IoT applications to draw more accurate conclusions.

Every country's top industry is healthcare, and using data analytics in IoT-based healthcare applications can lead to advances in this field. Using the same, it is possible to lower healthcare expenses, improve remote monitoring and distant health services, and increase diagnosis and treatment.

Therefore, the use of data analytics must be encouraged in the IoT space to boost sales, gain an advantage over competitors, and increase consumer satisfaction. Businesses can combine data analytics with IoT to use data to their advantage by working with the right strategy partner.

## CHI-SQUARE AND ANOVA TESTS

ANDRIA DSOUZA (21MCA04)  
PRIYADARSHINI N (21MCA31)

In this Article, we discuss two different techniques such as **Chi-square** and **ANOVA** Tests. Both are hypothesis testing mainly theoretical.

### **Chi-Square**

Chi-square statistical method commonly used for testing a relationship between categorical variables. In statistics, there are two types of variables: numerical (countable) variables and non-numerical (categorical) variables. The null hypothesis of the Chi-square test is that no relationship exists on the categorical variables in the population and they are the independent variables. The chi-square test can be used to determine whether observed frequencies are significantly different from expected frequencies.

$$\chi_c^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Where,

O = observed score

E = Excepted score

A low value for chi-square means there is a high correlation between your two sets of data.

The hypothesis being tested for chi-square is

**Null:** Variable A and Variable B are independent.

**Alternate:** Variable A and Variable B are not independent.

### **Types of Chi-square**

#### **1. Chi-square goodness of fit test**

It determines if a sample data matches a population.

#### **2. Chi-square test for independence**

Compares two variables in a contingency table to and check they are related or not. In a more general sense, it tests to see whether distributions of categorical variables differ from each other.

### **ANOVA (Analysis of Variance)**

ANOVA is a statistical tool used for comparing the dependent and the independent variables. ANOVA technique that uses a sample of observations to compare the number of means. t is

similar to that of t-test and z-test, which are used to compare mean along with relative variance. However, in ANOVA, it is best suited when two or more populations/samples are compared. An ANOVA test is to find if you need to reject the null hypothesis or accept the alternate hypothesis.

### **Types of ANOVA**

#### **One-way ANOVA**

One-way has one independent variable (with 2 levels). The one-way ANOVA compares the means between the groups you are interested in and determines whether any of those means are statistically significantly different from each other. Specifically, it tests the null hypothesis.

$$H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$$

Where,

$\mu$  = group means

$k$  = number of groups.

#### **Two-way ANOVA**

Two-way ANOVA is used to compare two or more factors (i.e. check the effect of two independent variables on a single dependent variable.) Both types of ANOVA have a single continuous response variable. Use a two way ANOVA when you have one measurement variable (i.e. a quantitative variable) and two nominal variables.

#### **N-way ANOVA**

If more multiple independent variables then use N-way analysis of variance. The N-way ANOVA can show whether there are effects of the independent variable and interactions between them. Interactions are usually seen when one independent variable depends on the second independent variable.

#### **Application**

- Lean-Six Sigma/operational efficiency.
- Comparing the gas mileage of different vehicles, and also the same vehicle under different fuel types.
- Understanding the impact of temperature, pressure, or chemical reaction (power reactors, chemical plants, etc.).
- Understanding the performance, quality, or speed of manufacturing processes based on the number of cells.

# SPATIAL ANALYTICS FOR GIS DATA

**KAVYA SHREE.B (21MCA38)**

**SMRUTHI.R (21MCA40)**

## **Spatial Analysis: What Is It?**

Examining, evaluating, and modelling spatial data characteristics, such as locations, qualities, and relationships that indicate the geometric or geographic properties of the data, is referred to as spatial analysis. In order to assimilate geographic information and determine whether it is appropriate for a target system, it makes use of a number of computer models, analytical methods, and algorithmic approaches.

## **How Does Spatial Analysis Work?**

A key component of the Geographic Information System is spatial analysis (GIS). It is primarily used to produce weather forecasts for a certain region or evaluate the viability of a place for a particular system. It enables users to model problems and find detailed solutions that have geographic features attached to them.

**1. Data collection:** The collection of data is essential for the spatial analysis process. It includes data collection from a broad range of sources, such as remote sensing devices such as LiDAR (light detection and ranging) and airborne systems.

Data collected by such devices is used to create maps that reveal the spatial pattern of entities under consideration, such as a temperature map for various regions. In this case, data refers to high-resolution images or images captured by satellites or aerial systems.

**2. Data analysis:** The collected data is then analyzed using AI and ML solutions to produce results in the second step. Besides that, by analyzing millions of images, one can train ML models to detect objects or structures in a given area.

Colleges, playgrounds, traffic zones, residential areas, and so on are examples of objects. Furthermore, visualization tools can be used to highlight different objects with different shades, shapes, or annotations. Such tools make it simpler to identify objects within huge amounts of data.

**3. Data Presentation:** Post-analysis data presentation can be time-consuming because important elements that reveal the findings must be highlighted. Data visualization tools, which use tables, maps, and bar charts to project relevant data and communicate with stakeholders, make such tasks easier.

Besides that, 3D visualization tools add variables to 2D data and provide a more accurate perspective. Such practices improve planning and implementation strategies, resulting in more effective solutions to the modelled problems.

## **Here are some examples of spatial analysis:**

1. Urban development and planning: Spatial analysis is crucial in urban planning and development activities. Some of the projects that belong to this category.

- Establish resilient cities: Climate change has a negative impact on urban life. City authorities are constantly looking for ways to reduce the impact on urban residents. This is where technologies like GIS, which provides geospatial information, are useful. Policymakers can use GIS while keeping sustainability in mind.

- Monitor the urban heat island (UHI) effect in cities: The urban heat island (UHI) effect refers to the phenomenon in which natural vegetation is removed in order to build apartments and structures that retain heat for longer periods of time. Because it is one of the most serious issues that humanity must resolve, techniques such as spatial analysis can help. Measures such as Satellite data, satellite imagery, thermal remote sensing, and field observational studies can help develop a better understanding of how the UHI effect causes a specific spatial pattern. It is possible to identify the source of the UHI effect and take appropriate action.

2. Administration of public health: Multiple health and govt departments use spatial analysis to manage public health.

- Disease distribution maps: Satellite data is crucial for forecasting disease spread across regions. Such spatial data patterns enable policymakers to prevent disease spread by implementing preventive measures. Moreover, one can combine weather variables such as rainfall or temperature with disease data to better understand how weather affects disease spread or prevalence in various areas. Temperatures and the presence of nearby water bodies, such as lakes and rivers, are typically crucial in determining disease progression in the case of waterborne diseases.
- Vaccination statistics on a map: The COVID-19 pandemic has been one of the world's most significant challenges, especially for the healthcare sector, because vaccinating people was the only way out of this problem. However, how does one keep track of such vaccination campaigns? Governments can handle this scenario effectively by utilising GIS technologies. They can use spatial analysis to track vaccine distribution and ensure uniform distribution across communities. Thus, spatial analysis is critical in such large-scale vaccination campaigns.



# TABLEAU FOR ANALYTICS PLATFORM

CHANDANA C H (21MCA10)

UMA BHUVANESHWARI (21MCA42)

**Tableau** is an end-to-end data analytics platform that allows you to prep, analyze, collaborate, and share your big data insights. Tableau excels in self-service visual analysis. It is one of the most powerful, secure, and flexible end-to-end analytics platform.

Tableau is an excellent data visualization and business intelligence tool used for reporting and analyzing vast volumes of data.

Tableau has a lot of unique, exciting features that make it one of the most popular tools in business intelligence (BI).

As the market-leading choice for modern business intelligence, our analytics platform makes it easier for people to explore and manage data, and faster to discover and share insights that can change businesses and the world.

## Tableau Features:

- Tableau supports powerful data discovery and exploration that enables users to answer important questions in seconds
- No prior programming knowledge is needed; users without relevant experience can start immediately with creating visualizations using Tableau
- It can connect to several data sources that other BI tools do not support. Tableau enables users to create reports by joining and blending different datasets
- Tableau Server supports a centralized location to manage all published data sources within an organization

## Tools of Tableau:



### **Tableau Desktop:**

Tableau Desktop has a rich feature set and allows us to code and customize reports. Right from creating the reports, charts to blending them all to form a dashboard, all the necessary work is created in Tableau Desktop.

### **Tableau Public:**

This Tableau version is specially built for cost-effective users. The word '**Public**' means that the created workbooks cannot be saved locally. They should be kept on the Tableau's public cloud, which can be accessed and viewed by anyone. There is no privacy of the files saved on the cloud, so anyone can access and download the same data

### **Tableau Online:**

It creates a direct link over 40 data sources who are hosted in the cloud such as the **Hive, MySQL, Spark SQL, Amazon Aurora**, and many more.

The software is correctly used to share the workbooks, visualizations, which is created in the admin of the organization has full control over the server. The organization maintains the hardware and the software. Tableau Desktop application over the organization

### **Tableau Reader:**

Tableau Reader is a free tool which allows us to view the visualizations and workbooks, which is created using Tableau Desktop or Tableau Public. The data can be filtered, but modifications and editing are restricted. There is no security in Tableau Reader as anyone can view workbook using Tableau Reader.

### **Tableau Server:**

The software is correctly used to share the workbooks, visualizations, which is created in the Tableau Desktop application over the organization. To share dashboards in the Tableau Server, you should first publish your workbook in the Tableau Desktop. Once the workbook has been uploaded to the server, it will be accessible only to the authorized users.

### **References:**

<https://www.javatpoint.com/tableau-tools>

<https://www.guru99.com/what-is-tableau.html>

<https://www.tableau.com/why-tableau/what-is-tableau>

# NATURAL LANGUAGE PROCESSING TOOL KIT

PRIYANKA N (21MCA32)

VIDYA M (21MCA45)

## INTRODUCTION:-

Natural language processing (NLP) is the ability of a computer program to understand human language as it is spoken and written -- referred to as natural language. It is a component of artificial intelligence (AI).

NLP has existed for more than 50 years and has roots in the field of linguistics. It has a variety of real-world applications in a number of fields, including medical research, search engines and business intelligence.

## Features:-

- \*Text classification
- \*Part-of-speech tagging
- \*Entity extraction
- \*Tokenization
- \*Parsing
- \*Stemming
- \*Semantic reasoning

## Major tools are included:-

**1. NLTK - entry-level open-source NLP Tool:-**Natural Language Toolkit (AKA NLTK) is open-source software powered with Python NLP. From this point, the NLTK library is a standard NLP tool developed for research and education.

NLTK provides users with a basic set of tools for text-related operations. It is a good starting point for beginners in Natural Language processing.

**2. Stanford Core NLP Tool:-** We can say that the Stanford NLP library is a multi-purpose tool for text analysis. Like NLTK, Stanford Core NLP provides many different natural language processing software. But if you need more, you can use custom modules. The main advantage of Stanford NLP tools is scalability. Unlike NLTK, Stanford Core NLP is a perfect choice for processing large amounts of data and performing complex operations

**3. Apache OpenNLP:-**Accessibility is essential when you need a tool for long-term use, which is challenging in the realm of Natural Language Processing open-source tools. Because while being powered with the right features, it could be too complex to use.

Apache OpenNLP is an open-source library for those who prefer practicality and accessibility. Like Stanford CoreNLP, it uses Java NLP libraries with Python decorators.

REFERENCES: - <https://theappsolutions.com/blog/development/nlp-tools>

# DATA SCRAPING

**AASIYA MEHAK M (21MCA01)**

Data scraping refers to the process of acquiring data from websites and online sources for data analytics and data science purposes. It can be done by hand or with automated software, depending on the level of automation desired by the user and the amount of information being scraped. Data scraping can provide valuable insights into both your business's performance and that of your competitors', and it can make it easier to update your products if you scrape competitor information as well as information on your customers and their preferences.



The Five Essential Tools for Data Scraping:

## 1. Web Crawlers:

Web crawlers, or spiders, are computer programs that search through websites. They help you find information that isn't available on a website's homepage. An example would be searching for local restaurants to include on a food delivery service like Seamless or Grub hub. There are lots of web crawlers out there and each has its own particular way of grabbing content from a website.

## 2. Screen Scrapers

Web-Based Screen Scrapers: Web scraping is a popular way to extract data from any web pages. There are many web-based screen scrapers which you can use to scrape through the web pages quickly and easily without needing too much coding knowledge. These screen scrapers are very useful for data extraction purpose. Data Scraper is one such screen scraper for extracting data from different websites like Facebook, YouTube, and Twitter etc.

## 3. Database

The most traditional way to get data is by connecting to a database. Data aggregation tools like SQL, Hive, Pig, etc., make it easy to pull out data sets and combine them into a single table that can be analyzed as one. If you're taking data from a relational database (like MySQL), then there are libraries for just about every language that make connecting with your database and retrieving its information trivial. Just make sure you have permission before accessing someone else's databases!

## 4. APIs

The quickest way to collect data from a variety of sources is to use one or more APIs, which stands for Application Programming Interface. The term may be confusing because it doesn't have anything to do with learning to code. It has everything to do with gathering data. In other words, using an API allows you to connect one application (such as Google Drive) with another (such as your own spreadsheet).

## 5. Ecommerce Sites

If you're looking to scrape data off an ecommerce site, then it's likely that you want to get product information like prices and descriptions. In those cases, there are two popular methods: web scraping and API integration. Both work just fine, but they each have their own strengths and weaknesses. Web scraping is faster but requires more coding; APIs are typically simpler to integrate with but slower. So if speed is your biggest concern, you might lean toward using a web scraper for now; if long-term access is most important, you can use an API instead. Ecommerce websites tend to list their public APIs on their website or with other third-party tools.

### References:

- 5 Popular Ways of Data Scraping in Data Analytics - Machine Learning | AI | Data Science (connectjaya.com)
- <https://www.datamation.com/big-data/data-scraping>

# IMPUTATION TECHNIQUES IN DATA ANALYTICS

ARUNODAYA P (21MCA06)  
LAKSHMI SHRUTHI J (21MCA23)

## **Data Imputation:**

Data Imputation is the process of evaluating the missing values to retain the majority information of the datasets. These methods are used because it would be impractical to remove data from the dataset each time. Missing values can increase the chance of making Type I and Type II errors and reduce statistical power and limit the reliability of confidence intervals.

## **Features:**

1. **Distorts Dataset:** Data distortion is the deviation of data from its true or most accurate representation of the full picture. Bad data may interject incorrect or misguided facts into useful information or worse into models and predictions of customers or business.
2. **Unable to work with the majority of machine learning-related Python libraries:** When utilizing ML libraries, mistakes may occur because there is no automatic handling of these missing data.
3. **Impacts on the Final Model:** Missing data may lead to bias in the dataset, which could affect the final model's analysis.
4. **Desire to restore the entire dataset:** This typically occurs when we don't want to lose any of the data in our dataset because all of it is crucial. Additionally, while the dataset is not very large, eliminating a portion of it could have a substantial effect on the final model.

## **Data Imputation Technique**

1. **Next or Previous Value:** This technique is basically used for time-series data or ordered data. The next or previous value inside the time series is typically substituted for the missing value as part of a common method for imputed incomplete data in the time series. This strategy is effective for both nominal and numerical values.
2. **K Nearest Neighbors:** The objective is to find the k nearest examples in the data where the value in the relevant feature is not absent and then substitute the value of the feature that occurs most frequently in the group.
3. **Maximum or Minimum Value:** The minimum is simply the lowest observation, while the maximum is highest observation. It is the easiest way to determine the minimum and maximum if the data ordered from lowest to highest.
4. **Missing Value Prediction:** Missing data, or missing values, occur when you don't have data stored for certain variables or participants. Data can go missing due to incomplete data entry, equipment malfunctions, lost files, and many other reasons. In any dataset, there are usually some missing data.
5. **Most Frequent Value:** The most frequent value in the column is used to replace the missing values in another popular technique that is effective for both nominal and numerical features.
6. **Average or Linear Interpolation:** Average interpolation means determining a value from the existing values in a given data set. Another way of describing it is the act of inserting or interjecting an intermediate value between two other values.

7. Mean or Moving Average or Median Value: The average or linear interpolation, which calculates between the previous and next accessible value and substitutes the missing value, is similar to the previous/next value imputation but only applicable to numerical data.

8. Fixed Value: Fixed value imputation is a general method that works for all data types and consists of substituting the missing value with a fixed value.

## **REDUCING MANUFACTURING FAILURES USING DATA ANALYTICS**

**HARIPRIYA S (21MCA18)**

**GOOTY SATWIKA (21MCA17)**

The daily challenge of manufacturing products according to deadlines, characteristics and quality standards is one of the most important responsibilities of any company especially those on in the industrial sector. In order to do this, the components of these products and the final products have to be thoroughly tested for the possibility of defects and to ensure optimum performance. Testing is a very complex process which costs a company a fair amount of time, energy and money. To avoid this failure Data Analytics applies various parts of the manufacturing process, that aim to reduce the time and cost required for testing while ensuring that the product quality still remains up to the mark.

### **Techniques and systems to reduce production failures:**

1. Preventive measures: To prevent future problems, companies replace obsolete equipment and renew machines that are failing permanently. However the maintenance team will need to constantly do check-ups and possibly replace equipment, resulting in frequent interruptions of production for this type of technique.
2. Regular inspections: Every company should have a maintenance team that is dedicated, at least on a regular basis, to monitoring the machinery to ensure that it is in good working order at all times.
3. Attend to communication flows: To avoid future errors, it is important to maintain positive and constant communication between production employees and other departments.
4. Maintaining quality controls: Innovation in product design and other aspects are useless if the production line is not regularly controlled.

**Smart Manufacturing** is a set of collaborative manufacturing systems that respond, in real time, to the changing needs of factories, including predictive maintenance. This involves concepts of Big Data, IIoT or Machine Learning are applied to address a more flexible, modular and optimized production. In today's hyper connected and flexible world, companies need systems that help reduce costs, improve efficiency and make the most of production processes in order to succeed. Hence, the best way to increase your business efficiency is by converting it into a smart factory.

### **References:**

- <https://nexusintegra.io/how-to-reduce-production-failures/>
- [https://www.projectpro.io/article/big-data-analytics-projects-for-students-/436#mctoc\\_1f913kev34](https://www.projectpro.io/article/big-data-analytics-projects-for-students-/436#mctoc_1f913kev34)



# APACHE SPARK

**HEMA N (21MCA20)**

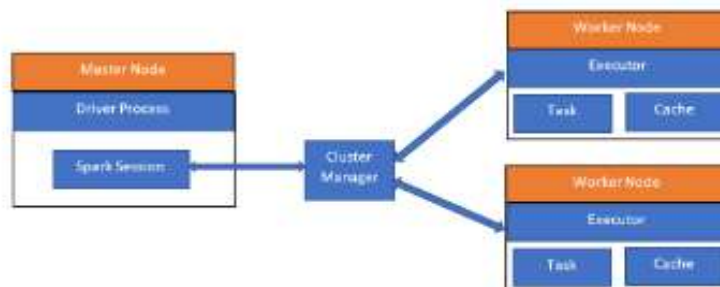
**SHILPA V (21MCA39)**

- Apache Spark is a multi-language engine for executing data engineering, data science, and machine learning on single-node machines or clusters.
- Apache Spark is a unified analytics engine for large-scale data processing. It provides high-level APIs in Java, Scala, Python and R, and an optimized engine that supports general execution graphs. It supports a rich set of higher-level tools including Spark SQL for SQL and structured data processing, pandas API on Spark for pandas workloads, ML lib for machine learning, GraphX for graph processing, and Structured Streaming for incremental computation and stream processing.
- Apache Spark is used in banks, telecommunications companies, games companies, governments, and all of the major tech giants such as Apple, Facebook, IBM, and Microsoft.

## **Features**

- Real-time stream processing.
- Ease of use.
- Lightning-fast processing speed.
- It is flexible.

## **Architecture**



## **Driver process**

Driver process is the heart of a spark application. It sits on master node and maintains the information about the spark application throughout the program execution. Driver program accepts the user's input and then it analyzes, distributes, and schedules the work across the executors running at worker node.

**Executor process**

The executor node is responsible for running the actual task assigned to it by the driver process. It simply executes the given instructions and sends the results back to the driver node.

**References**

<https://sqlrelease.com/big-data-processing-using-apache-spark-introduction>

<https://docs.cloud.sdu.dk/Apps/spark-cluster.html>

<https://www.altexsoft.com/blog/spark/>

# SQL FOR DATA ANALYTICS

**RAKSHITHA J: 21MCA34**

**VAISHNAVI P: 21MCA43**

SQL for Data Analytics is a powerful programming language that helps data analysts interact with data stored in Relational databases. SQL consists of five basic commands to control structure, perform manipulation for transactions, and Data Analytics.

- **Benefits of SQL for Data Analytics:**

- It is easy to understand and learn.
- It is efficient at fast query processing and helps in retrieving big data from multiple databases efficiently.

- **Understanding SQL for Data Analytics:**

## 1. SQL Queries

SQL queries can be classified into five parts as they perform specific roles to execute queries on any RDBMS system, and they are:

- . Data Definition Language (DDL)
- a. Data Manipulation Language (DML)*
- b. Data Query Language (DQL)*
- c. Data Control Language (DCL)*
- d. Transaction Control Language (TCL)*

## 2. SQL Joins

The SQL join clause is used to combine different tables in databases, where JOIN is made using a Primary and Foreign key. There are four major joins which include inner, left, right, and full join used in combination with the 'from' clause.

## 3. SQL Aggregations

It is a process of combining multiple entities can be performed by SQL aggregation query. SQL comes with some standard functions like count, sum, min, max, and avg operation. These functions are often used in conjunction with 'groupby', 'orderby,' and 'having' clauses to

evaluate specific columns.

```
mysql> select sum(Total_amt_spend) as sum_all from ConsumerDetails ;
+-----+
| sum_all |
+-----+
|    12560 |
+-----+
1 row in set (0.00 sec)
```

#### 4. SQL Views and Stored Procedures

SQL views are virtual tables whose content is obtained from an existing table, and it optimizes the database to provide an additional level of security by restricting users from fetching complete information from the database.

Stored procedures are created to process one or more DML operations on a database and are also capable of taking user input to perform a group of SQL commands.

# **FUNNEL ANALYSIS IN MARKETING**

**PRITY KUMARI (21MCA30)**

**RUPA KUMARI (21MCA36)**

## **INTRODUCTION**

Funnel analysis is a powerful analytics method that shows visually the conversion between the most important steps of the user journey. It is a method of understanding the steps required to reach an outcome on a website and how many users get through each of those steps. The set of steps is referred to as a “funnel” because the typical shape visualizing the flow of users is similar to a funnel in your kitchen or garage. Funnel analysis is typically useful if we want to map out a linear user journey.

### **Marketing Funnel Analysis**

- A Marketing Funnel can be seen as a powerful analytics technique that can help businesses understand the customer journey all the way from the first time they visit the website to the last stage until they're ready to make the purchase.
- Marketing Funnel Analysis can be used to understand the steps required to reach a particular outcome on a website and how the number of users varies as they go through these steps.
- This is referred to as a Funnel because a visualization displaying the number of users as they go through different stages until they're ready to make a purchase is in the shape of a Funnel.

With funnel analysis we try to model processes where users can reach their goal via a relatively straightforward path.

### **Steps to Perform Funnel Analysis**

**Step 1: Define Metrics :-** For any useful analysis to be performed, the metrics that will be tracked have to be defined first. The metrics can then be tracked which will result in the collection of relevant data for analysis.

**Step 2: Identify Key Touchpoints: -** To identify the touchpoints, we need to understand what goals users have in mind when they visit your website and how they can achieve those goals.

**Step 3: Perform In-depth Analysis by Segmenting Users:-** Various segments can be made based on factors such as geographic location, gender, age-group, time period, etc., and further Funnel Analysis can be performed on these groups to get deeper insights, allowing you to plan future strategies accordingly.

## **Benefits of Performing Funnel Analysis**

1) Find Pages with High Dropout Rates:- Funnel Analysis can help us visualize the Dropout Rates for each page and the final Conversion Rate. Having an understanding of where our visitors are dropping out can help you plan our future optimization goals accordingly.

2) Determine Source of High-Quality Visitors:-Advanced Analytics tools such as Google Analytics will allow to gain deeper insights by performing advanced Funnel Analysis. This analysis can be used to determine the source of high-quality traffic coming to your website. This information can be used to plan campaigns to increase the volume of traffic from these high-quality sources further.

3) Easy Communication with Team Members:-Funnel Analysis is a simple to understand analytics technique that can be used to communicate observations to the team members and all stakeholders involved allowing them to make data-driven decisions quickly.

## **Conclusion**

Funnel analysis is a great tool to measure the user journey, the conversion and the development of users – step by step.

When we set one up, define 4-8 steps and create a funnel visualization that's easy to understand!

Extract insights! Check at which step most people churn, and which segment is the strongest!

## **References:-**

<https://hevo.com/learn/understanding-marketing-funnel-analysis/>

<https://chartio.com/learn/product-analytics/what-is-a-funnel-analysis/>

## HYPOTHESIS TESTING STATISTICAL ANALYSIS

NIVEDITHA (21MCA26)

POOJA (21MCA28)

Hypothesis Testing is a type of statistical analysis in which you put your assumptions about a population parameter to the test. It is used to estimate the relationship between 2 statistical variables.

The two types of hypothesis testing in statistics.

- Null Hypothesis and Alternate Hypothesis
- Simple and Composite Hypothesis Testing

### **Null Hypothesis and Alternate Hypothesis:**

The Null Hypothesis is the assumption that the event will not occur. A null hypothesis has no bearing on the study's outcome unless it is rejected.

$H_0$  is the symbol for it, and it is pronounced H-naught.

The Alternate Hypothesis is the logical opposite of the null hypothesis. The acceptance of the alternative hypothesis follows the rejection of the null hypothesis.  $H_1$  is the symbol for it.

#### **Example:**

A sanitizer manufacturer claims that its product kills 95 percent of germs on average.

To put this company's claim to the test, create a null and alternate hypothesis.

$H_0$  (Null Hypothesis): Average = 95%.

Alternative Hypothesis ( $H_1$ ): The average is less than 95%.

### **Simple and Composite Hypothesis Testing:**

Simple Hypothesis: A simple hypothesis specifies an exact value for the parameter.

Composite Hypothesis: A composite hypothesis specifies a range of values.

#### **Example:**

A company is claiming that their average sales for this quarter are 1000 units. This is an example of a simple hypothesis.

Suppose the company claims that the sales are in the range of 900 to 1000 units. Then this is a case of a composite hypothesis.

### **One-Tailed and Two-Tailed Hypothesis Testing:**

In a one-tailed test, the critical distribution area is one-sided, meaning the test sample is either greater or lesser than a specific value.

In two tails, the test sample is checked to be greater or less than a range of values in a Two-Tailed test, implying that the critical distribution area is two-sided.

**Type 1 and Type 2 Error:**

A hypothesis test can result in two types of errors.

Type 1 Error: A Type-I error occurs when sample results reject the null hypothesis despite being true.

Type 2 Error: A Type-II error occurs when the null hypothesis is not rejected when it is false, unlike a Type-I error.



## MOVIE REVIEW SENTIMENT ANALYSIS

**RABIYA BASRI (21MCA33)**

**VANDHANA R (21MCA44)**

Sentiment Analysis involves classifying a remark as positive, negative or neutral. Broader classifications can also be done such as “strongly agree, agree, neutral, disagree, strongly disagree”. By analyzing reviews written by individuals who have watched a particular movie, it is possible to get better movie recommendations. With the popularity of social media, it is very easy to express one’s sentiments towards a particular movie. There are many repercussions to these reviews that are expressed.

Sentiment analysis of the reviews can be used to identify interest patterns in the audience, which can accordingly be used to generate recommendations. The use of sarcasm, ambiguity in language use, negations and multipolarity add to the challenges associated with categorizing the reviews based on sentiments. Through this project, it will be possible to understand more about the concept of speech tagging, the difference between stemming and lemmatization and applications of sparse matrices. Application of Naive Bayes model and SVM (Support Vector Machine) can be used for the training model and accordingly making predictions.

### Sentiment Analysis



#### **Features :**

Key features that are essential in a sentiment monitoring tool are :

1. multilingual efficacy
2. precise aspect-based sentiment analysis
3. named entity recognition
4. an effective visualization dashboard.

#### **Advantages of sentiment analysis:**

Sentiment analysis can be used to analyze customer data across geographies, ethnicities, and cultures in order to get the most intrinsic insights about product popularity, competitive brands, success rates of advertising campaigns, and more.

#### **Disadvantages of sentiment analysis:**

Sentiment analysis methods usually do not identify sarcasm, negation, grammar mistakes, misspellings, or irony. Thus, it may not be suitable for analyzing data gathered from social media platforms.

# DATA ANALYTICS USING QUANTUM COMPUTING

**Ashritha G (21MCA07)**

## **What Is Quantum Computing?**

Quantum computing provides high-speed detection, analysis, integration, and diagnosis when dealing with huge, scattered data sets. Quantum computers can quickly find patterns in enormous, unsorted data sets by seeing every item in a massive database at the same time. Quantum computers can execute highly complex computations in seconds, but non-quantum computers can take hundreds of years to do so.

## **Data Analytics Using Quantum Computing:**

- When using large scattered data sets, quantum computing offers high-speed detection, analysis capabilities, integration, and diagnosis.
- Quantum computers can locate patterns quickly in large, unsorted data sets by simultaneously viewing every item in a huge database.

## **Quantum Computing Has The Potential To Revolutionize Data Analysis In Several Ways. Here Are A Few Examples:**

1. **Machine learning:** Quantum computing can be used to speed up the training process for machine learning algorithms. For example, quantum algorithms can be used to perform linear algebra operations, which are the backbone of many machine learning algorithms, much faster than classical computers.
2. **Optimization problems:** Many data analysis problems can be framed as optimization problems, where the goal is to find the optimal solution that maximizes or minimizes a particular objective function. Quantum computing can be used to solve these optimization problems much faster than classical computers.
3. **Database search:** Searching through large databases can be very time-consuming on classical computers. Quantum computers can perform this task much faster by using quantum parallelism to simultaneously search through all possible solutions.
4. **Sampling:** Data analysis often requires sampling, which is the process of selecting a subset of data points from a larger dataset. Quantum computing can be used to perform sampling more efficiently by using quantum algorithms to generate the samples.
5. **Dimensionality reduction:** High-dimensional data can be difficult to analyze because of the large number of variables involved. Quantum computing can be used to perform dimensionality reduction, which is the process of reducing the number of variables in the data, much faster than classical computers.
6. **Predictive Analytics:** Artificial intelligence can be used to extract meaningful historical facts and current data from datasets. When combined with quantum computing, more data is processed, which yields relevant information that can then be used to make predictions.

## **Benefits Of Using Quantum Computing For Data Analysis:**

1. **Speed:** One of the biggest benefits of quantum computing is its ability to perform certain operations much faster than classical computers. For example, quantum algorithms can

perform linear algebra operations, which are the backbone of many data analysis algorithms, much faster than classical algorithms. This means that quantum computing has the potential to dramatically speed up data analysis tasks.

2. **Improved accuracy:** Quantum computing can also improve the accuracy of data analysis by allowing for more complex models to be developed and applied to large datasets. For example, quantum algorithms can be used to perform more advanced forms of dimensionality reduction, which can result in more accurate representations of the data.
3. **Handling large datasets:** Classical computers can struggle to process and analyze large datasets in a reasonable amount of time. Quantum computing has the potential to overcome this challenge by using quantum parallelism to process multiple pieces of data simultaneously. This makes quantum computing well suited for analyzing large datasets, which are becoming increasingly common in many fields.
4. **Optimization problems:** Many data analysis problems can be framed as optimization problems, where the goal is to find the optimal solution that maximizes or minimizes a particular objective function. Quantum computing has the potential to solve these optimization problems much faster than classical computers, allowing for more efficient and accurate solutions to be found.
5. **Quantum simulation:** Quantum computing can also be used to simulate quantum systems, which are commonly used in many fields, including chemistry, materials science, and electronics. This can be useful for analysing data related to these fields and gaining new insights into their underlying processes.

#### **Algorithms Used By Quantum Computing For Data Analysis:**

Quantum computing has the potential to revolutionize data analysis by providing new algorithms that can solve problems faster than classical algorithms. Some of the quantum algorithms that have been developed for data analysis are:

1. **Quantum Principal Component Analysis (PCA):** This is a quantum algorithm for dimensionality reduction, which is used to reduce the number of variables in a dataset. The algorithm can perform PCA much faster than classical algorithms and is particularly useful for large datasets.
2. **Quantum Linear System Algorithm (HHL):** This is a quantum algorithm for solving linear systems of equations, which are commonly used in many data analysis problems. The HHL algorithm can solve linear systems much faster than classical algorithms, making it a valuable tool for data analysis.
3. **Quantum Support Vector Machines (SVMs):** This is a quantum algorithm for classification problems, which are used to categorize data points into different classes based on their features. The quantum SVM algorithm can perform classification much faster than classical algorithms and is particularly useful for large datasets.
4. **Quantum K-Means Clustering:** This is a quantum algorithm for clustering problems, which are used to group data points into clusters based on their similarity. The quantum K-Means algorithm can perform clustering much faster than classical algorithms and is particularly useful for large datasets.
5. **Quantum Grover's Algorithm:** This is a quantum algorithm for searching problems, which are used to find a specific item in a dataset. The Grover's algorithm can search through a dataset much faster than classical algorithms and is particularly useful for large datasets.

# **SMART HEALTHCARE SYSTEMS LEVERAGING BIG DATA ANALYTICS**

**ANANYA G S (21MCA03)  
CHITHRA A V (21MCA12)**

Smart healthcare systems are the new way of thinking about technology in healthcare. These systems use sensors, computers, and other technologies to gather patient data, store them in a database, analyze them, and provide solutions based on those results. A smart health care system is a computerized network that allows information to be shared between various doctors, clinics, and hospitals.

**Smart healthcare systems are based on the following four fundamental pillars:**

- **Connectedness:** Connectedness is the ability to share information between people, devices, and systems.
- **Transparency:** Transparency means accessing the information stored in various systems without barriers or limitations.
- **Security:** Security refers to data protection from hackers or unauthorized users who may try to access your information without permission.
- **Analytics:** Analytics refers specifically to the ability to analyse large amounts of data quickly and effectively to identify patterns that may indicate potential risks.

**Cloud and Big Data analytics in Healthcare Organizations:**

Cloud and big data analytics are bringing a new level of intelligence to healthcare organizations. Analysing massive amounts of data from multiple sources can help doctors, nurses, and other medical professionals make better decisions.

Cloud computing has played a crucial role in transforming healthcare by improving efficiency and reducing costs. Big data analytics lets you look at trends across entire populations, which helps you make better decisions about treatment options and care plans.

**Healthcare applications of big data analytics include:**

- **Patient monitoring systems** – enable continuous monitoring of vital signs such as blood pressure, heart rate, temperature, etc.
- **Alert systems** – It monitors the presence of infectious diseases.
- **Smart medical devices** – Devices like glucose meters that automatically send data on blood sugar levels to doctors.
- **Disease Detection** – The ability to detect disease early and accurately is one of the essential features of smart health care systems.
- **Personalized Care** – One of the biggest challenges in personalized care is figuring out what works for whom and why.

**References:**

<https://intellectdata.com/smart-healthcare-systems-leveraging-big-data-analytics/#:~:text=It%20can%20improve%20treatment%20efficacy,heart%20rate%2C%20temperature%2C%20etc>

# WEB ANALYTICS FOR DIGITAL MARKETING

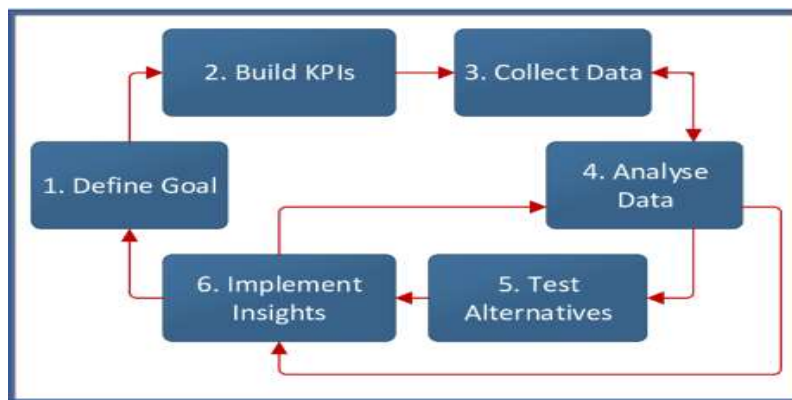
AVULA LAKSHMI (21MCA08)

JAISRI D (21MCA21)

## Introduction to Web Analytics

Web analytics is a subfield of data analytics that specifically deals with the collection, analysis, and interpretation of data from websites and online applications.

The **goal of web analytics** is to provide insights into how people are interacting with a website, including traffic patterns, user behavior, and conversion rates. This information is used to improve website performance, optimize marketing campaigns, and measure the success of digital marketing efforts.



## Process of web analytics

### Types of Web Analytics for Digital Marketing

There are two main types of web analytics for digital marketing:

1. Real-time analytics
2. Historical analytics.

**Real-time analytics** provide insights into user behavior in real-time and can be used to track website performance, monitor user engagement, and identify potential areas of improvement.

**Historical analytics** provide insights into user behavior over a period of time. This data can be used to understand customer behavior, create targeted campaigns, and measure the effectiveness of digital marketing campaigns.

## **Role of Web analytics in digital marketing**

Web analytics plays a critical role in digital marketing, as it provides insights into how website visitors interact with a brand's digital presence. **Some of the key roles of web analytics in digital marketing include:**

- 1. Measurement and Tracking:** Web analytics tools allow marketers to track and measure various metrics, such as website traffic, conversion rates, and user behavior, so that they can understand the impact of their marketing efforts.
- 2. Data-Driven Decision Making:** Web analytics provides the data and insights needed to make informed, data-driven decisions about digital marketing campaigns. This allows marketers to optimize their efforts and improve ROI.
- 3. Customer Insights:** Web analytics provides valuable insights into customer behavior, preferences, and motivations. This information can be used to improve the customer experience, tailor marketing messages, and increase conversion rates.
- 4. Competitive Intelligence:** Web analytics allows marketers to track and compare their website performance against competitors, and identify areas for improvement.
- 5. Optimization and Personalization:** Web analytics can help businesses to personalize the customer experience, by using data and insights to deliver relevant, targeted content and experiences. This can help improve customer engagement and conversion rates.
- 6. Campaign Tracking:** Web analytics allows marketers to track and measure the success of specific marketing campaigns, and identify areas for improvement.

Some common types of data collected in web analytics include:

- Website traffic data, such as page-views, unique visitors, and referrers
- User behavior data, such as clickstream analysis, time on site, and exit pages
- Conversion data, such as conversion rates, goals, and funnels
- Demographic data, such as location, age, and gender

## **References**

1. [https://www.tutorialspoint.com/digital\\_marketing/digital\\_marketing\\_web\\_analytics.html](https://www.tutorialspoint.com/digital_marketing/digital_marketing_web_analytics.html).
2. <https://www.techtarget.com/searchbusinessanalytics/definition/Web-analytics>

# **DATA ANALYTICS IN AGILE DATA SCIENCE**

**JHANCY.S (21MCA22)**

## **INTRODUCTION**

Agile data science is an approach of using data science with agile methodology for web application development. It focusses on the output of the data science process suitable for effecting change for an organization. Data science includes building applications that describe research process with analysis, interactive visualization and now applied machine learning as well. The major goal of agile data science is to document and guide explanatory data analysis to discover and follow the critical path to a compelling product.

## **WHY AGILE FOR DATA SCIENCE**

Data science has to be agile, with agility being defined as the ability to offer actionable insights quickly, iterate on such insights, and validate the results. Agile Data Science is a programming process that deals with the unpredictability of generating analytics platforms from data at scale. Many reputable firms are in need of skilled data scientists who can implement agile approaches in their projects.

## **ADVANTAGES**

- Continuous improvement and delivery.
- Enhanced client satisfaction.
- Enhanced communication.
- Higher quality deliverables with easy shippable.
- Planning a good sprint is prioritized.

## **DISADVANTAGES**

- Being very skilled or organized in Scrum.
- The needs and the scope are prone to rapid change.
- It is harder to quantify data science activities since they are less well-defined.
- Data Science sprints, like engineering sprints, are expected to provide deliverables.
- The developer role may not be well defined.

## **What Can I Use Agile Analytics For?**

One major application for agile analytics is in the financial industry. When looking for evidence of malfeasance or illegal activities, investigators cannot always find a linear paper trail to prove a



crime was committed. Agile analytics allow for a more open analysis of data, which can show different anomalous activities that can be connected, even if not obvious at first. When paired with a powerful data visualization tool, it can be an excellent tool for investigators.

## LATEST TRENDS



## REFERENCES

- <https://eugeneyan.com/writing/data-science-and-agile-what-works-and-what-doesnt/>
- <https://www.sisense.com/glossary/agile-analytics/>

## **R PROGRAMMING FOR DATA SCIENCE**

**ANKITA NAG (21MCA05)**

**DIPIKA (21MCA14)**

R is an open-source programming language that is widely used as a statistical software and data analysis tool. R is an important tool for Data Science. It is highly popular and is the first choice of many statisticians and data scientists.

### **Data Science in R Programming Language**

Data Science has emerged as the most popular field of the 21st century. It is because there is a pressing need to analyze and construct insights from the data. Industries transform raw data into furnished data products. In order to do so, it requires several important tools to churn the raw data. R is one of the programming languages that provide an intensive environment for you to research, process, transform, and visualize information.

### **Features of R – Data Science**

Some of the important features of R for data science application are:

- R provides extensive support for statistical modelling.
- R is a suitable tool for various data science applications because it provides aesthetic visualization tools.
- R is heavily utilized in data science applications for ETL (Extract, Transform, Load). It provides an interface for many databases like SQL and even spreadsheets.
- R also provides various important packages for data wrangling.
- One of the important feature of R is to interface with NoSQL databases and analyze unstructured data.

### **Applications of R for Data Science**

Top Companies that use R for Data Science:

- Google: At Google, R is a popular choice for performing many analytical operations. The Google Flu Trends project makes use of R to analyze trends and patterns in searches associated with flu.
- Facebook Facebook makes heavy use of R for social network analytics. It uses R for gaining insights about the behavior of the users and establishes relationships between them.
- IBM: IBM is one of the major investors in R. It recently joined the R consortium. IBM also utilizes R for developing various analytical solutions. It has used R in IBM Watson – an open computing platform.
- Uber: Uber makes use of the R package shiny for accessing its charting components. Shiny is an interactive web application that's built with R for embedding interactive visual graphics.

### **References:**

<https://www.geeksforgeeks.org/r-programming-for-data-science/>

## **SOCIAL MEDIA ANALYTICS**

**S PAVITHRA (21MCA37)**  
**PADMA PRIYA (21MCA27)**

### **What is Social Media Analytics?**

"The art and science of extracting valuable hidden insights from vast amounts of semi-structured and unstructured social media data to enable informed and insightful decision making" -Gohfar F. Khan

First, social media analytics isn't about brands. It's about people sharing their lives with others they know based on common interests.

"Social media analytics is the collection and analysis of data points that help you measure the performance of your social media accounts"

### **Why Social Media Analytics?**

- Many businesses adopt a brand-centric focus when starting out on their data analytics journey.
- As of 2022, 92% of marketers in companies with more than 100 employees have started using social analytics to better understand the landscape.
- This is crucial, as social media offers brands a huge pool of consumers ripe for brand communication targeted toward relevant interests but consumers resent interruptions.

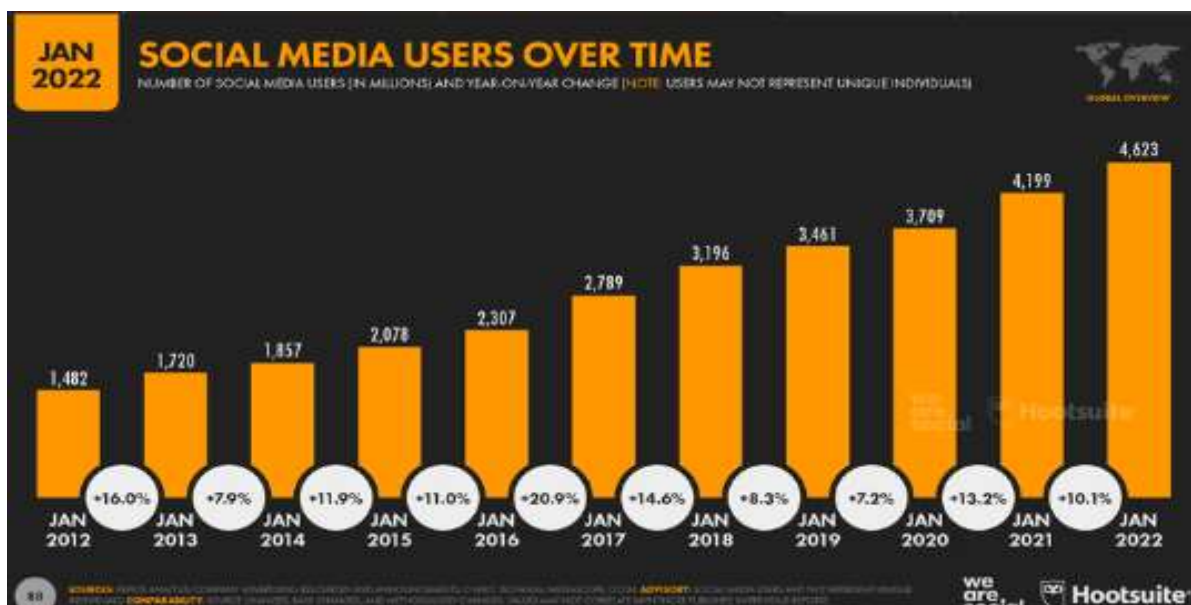
### **Importance:**

- Increase Customer Acquisition
- Protect Brand Health
- Lower Customer Care Costs
- Maximize Product Launches
- Boost Campaign Performance
- Improve Crisis Management

### **Social Media Analytics Involves:**

- Measuring Sentiments
- Measuring Conversation
- Analyses images
- Understand Audiences
- Influencer Analysis
- Social Media Intelligence

### **Users and social media:**



# CORRELATION AND REGRESSION

MEGHANA N (21MCA24)

STEFFI.P (21MCA41)

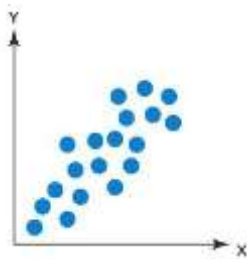
## INTRODUCTION

The most commonly used techniques for investigating the relationship between two quantitative variables are correlation and linear regression. Correlation quantifies the strength of the linear relationship between a pair of variables, whereas regression expresses the relationship in the form of an equation. For example, in patients attending an accident and emergency unit (A&E), we could use correlation and regression to determine whether there is a relationship between age and urea level, and whether the level of urea can be predicted for a given age.

## Correlation in Data Analytics

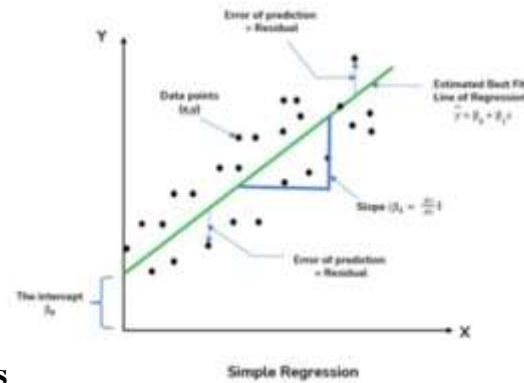
Correlation analysis in research is a statistical method used to measure the strength of the linear relationship between two variables and compute their association. Simply put - correlation analysis calculates the level of change in one variable due to the change in the other. The correlation coefficient 'r' is used to decide the strength of the relationship between two variables, and its value ranges between -1 and 1, where:

- 1 indicates a strong positive relationship.
- -1 indicates a strong negative relationship.
- A result of zero indicates no relationship at all.



**Correlation analysis** is done so as to determine whether there is a relationship between the variables that are being tested. Furthermore, a correlation coefficient such as Pearson's

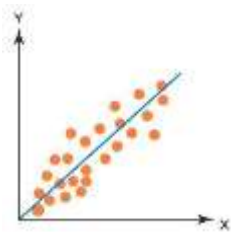
correlation coefficient is used to give a signed numeric value that depicts the strength as well as the direction of the correlation. The scatter plot gives the correlation between two variables x and y for individual data points as shown.



### Regression in Data Analytics

Regression analysis is a statistical technique of measuring the relationship between variables. It provides the values of the dependent variable from the value of an independent variable. The main use of regression analysis is to determine the strength of predictors, forecast an effect, a trend, etc.

This can be shown as a graph with the two variables on the x and y-axis. The independent variable or variables, when changed, they affect the dependent variables, and the regression analysis tries to provide an indication of which particular input variables affect the output most. Furthermore, the governing equation can also quantify this change.



**Regression analysis** is used to determine the relationship between two variables such that the value of the unknown variable can be estimated using the knowledge of the known variables. The goal of linear regression is to find the best-fitted line through the data points. For two variables, x, and y, the regression analysis can be visualized as follows:

### Correlation and Regression

Correlation shows the relationship between the two variables, while regression allows us to see how one affects the other. The data shown with regression establishes a cause and effect, when one changes, so does the other, and not always in the same direction. With correlation, the variables move together.

- Correlation and regression are statistical measurements that are used to quantify the strength of the linear relationship between two variables.
- Correlation determines if two variables have a linear relationship while regression describes the cause and effect between the two.
- Pearson's correlation coefficient and ordinary least squares method are used to perform correlation and regression analysis.

## REFERENCES

<https://www.digitalvidya.com/blog/correlation-and-regression/>

<https://byjus.com/maths/correlation-and-regression/>

# AI IN DATA ANALYTICS

DEVIKA K (21MCA13)

## **Introduction**

Artificial intelligence refers to a set of technologies which enable computers to simulate human intelligence. Examples of AI include speech recognition, such as directing virtual assistants like Alexa to perform tasks, image recognition for identification, and autonomous driving. Big data primarily refers to data sets that are too large or complex to be dealt with by traditional data-processing application software. Some examples are transaction processing systems, customer databases, documents, emails, medical records, internet clickstream logs, mobile apps and social networks.

## **The impact of Big data on AI**

- Big data and AI complement each other. AI becomes better, the more data it is given. It's helping organizations understand their customers a lot better, even in ways that were impossible in the past
- **Detecting anomalies** - AI can analyze artificial intelligence data to detect unusual occurrences in the data. For example, having a network of sensors that have a predefined appropriate range. Anything outside of that range is an anomaly.
- **Probability of future outcome** - Using known conditions that have a certain probability of influencing the future outcome, AI can determine the likelihood of that outcome.
- **AI can recognize patterns** - AI can see patterns that humans don't
- **Data Bars and Graphs** - AI can look for patterns in bars and graphs that might stay undetected by human supervision

## **Relationship between AI and Big Data**

- AI can assist users in all phases of the big data cycle, or the processes involved in the aggregation, storage, and retrieval of diverse types of data from various sources. These include data management, pattern management, context management, decision management, action management, goal management, and risk management.
- Big data and artificial intelligence are also linked in terms of research and technological innovation for each field. Big data technology uses AI theories and methods and AI relies on large volumes of data and the supporting big data technologies to improve and evolve decision-making capabilities.

## **References**

- <https://www.qlik.com/us/augmented-analytics/big-data-ai>



- <https://indatalabs.com/blog/big-data-tech-and-ai>
- The impact of Big Data on AI-Souâd Demigha CRI (University of Paris 1 Sorbonne), Paris, France-2020 International Conference on Computational Science and Computational Intelligence (CSCI)

# KNIME IN DATA ANALYTICS

HARSHITHA PATIL (21MCA19)

MYTHREYI S N (21MCA25)

## **Introduction**

- KNIME (Konstanz Information Miner) is a data analytics platform that provides a visual workflow-based approach to data analysis and modeling. It is a user-friendly and flexible platform that allows data scientists, statisticians, and business analysts to develop and execute complex data analysis tasks without having to write code.
- KNIME supports a wide range of data analysis tasks, including data import and export, data cleaning and preprocessing, data transformation, statistical analysis, machine learning, and data visualization. The platform provides a library of pre-built nodes for these tasks, as well as a rich set of tools for customizing and extending the platform to meet specific needs.
- KNIME provides a visual and flexible platform for data analytics that allows users to develop complex workflows using a drag-and-drop interface. The platform supports a wide range of data analysis tasks and provides a rich set of nodes for data cleaning, transformation, analysis, and reporting.
- Additionally, KNIME integrates with many popular data analytics tools, such as R and Python, allowing users to leverage the power of these tools within the KNIME workflow. This makes it possible to leverage existing analytics scripts and models within the KNIME environment.
- Overall, KNIME is a powerful and flexible platform for data analytics that provides a visual and user-friendly approach to data analysis and modeling. Its wide range of features, rich library of pre-built nodes, and integration with popular data analytics tools make it a popular choice for data scientists and analysts.

The following is a general explanation of how KNIME works in data analytics, based on a typical workflow:

- **Data Import:** The first step in a KNIME workflow is to import data from various sources, including spreadsheets, databases, and cloud services. KNIME supports a wide range of file formats and can automatically detect the format of the data to be imported.
- **Data Cleaning:** After importing the data, the next step is to clean and preprocess the data to ensure that it is ready for analysis. This may involve tasks such as removing missing

values, handling outliers, and converting data into a consistent format. KNIME provides a variety of nodes for data cleaning, including filtering, grouping, and transformation.

- **Data Transformation**: Once the data is cleaned, it may need to be transformed into a different format to support the desired analysis. For example, data may need to be aggregated, pivoted, or reshaped. KNIME provides a variety of nodes for data transformation, including aggregating, pivoting, and reshaping.
- **Data Analysis**: The next step is to perform the actual data analysis. This may involve tasks such as descriptive statistics, data visualization, and predictive modeling. KNIME provides a variety of nodes for data analysis, including statistical analysis, machine learning, and data visualization.
- **Model Deployment**: After building a predictive model, the next step is to deploy it in a production environment. KNIME provides a variety of nodes for model deployment, including scoring, prediction, and model management.
- **Results Reporting**: The final step is to present the results of the analysis, either in the form of a report or as a dashboard. KNIME provides a variety of nodes for results reporting, including reporting, visualization, and dashboarding.
- **KNIME Analytics Platform**: This is the desktop edition of KNIME and is the most widely used edition of the platform. It provides a visual and flexible platform for data analytics that allows users to develop complex workflows using a drag-and-drop interface. The platform provides a rich set of nodes for data import and export, data cleaning and preprocessing, data transformation, statistical analysis, machine learning, and data visualization.
- **KNIME Server**: This is the enterprise edition of KNIME and is designed for organizations that need to scale their data analytics solutions. KNIME Server provides a web-based interface for executing and sharing workflows, as well as a centralized repository for storing and managing workflows and data. It also provides a secure and scalable infrastructure for running data analytics workflows in a production environment.