



Basics of Data Science

Dr.K.Pazhanikumar

Head

Department of Computer Science

S.T.Hindu College

What is Data Science?

Data science is the study of data to extract meaningful insights for business. It is a multidisciplinary approach that combines principles and practices from the fields of mathematics, statistics, artificial intelligence, and computer engineering to analyze large amounts of data.

In [Wikipedia](#), **Data Science** is defined as *a scientific field that uses scientific methods to extract knowledge and insights from structured and unstructured data, and apply knowledge and actionable insights from data across a broad range of application domains.*

Why Data Science?

- Data is the oil for today's world. With the right tools, technologies, algorithms, we can use data and convert it into a distinct business advantage
- Data Science can help us to detect fraud using advanced machine learning algorithms
- It helps us to prevent any significant monetary losses
- Allows to build intelligence ability in machines
- It enables us to take better and faster decisions
- It helps us to recommend the right product to the right customer to enhance your business

Other Related Fields

- Databases
- Big Data
- Machine Learning
- Artificial Intelligence
- Visualization

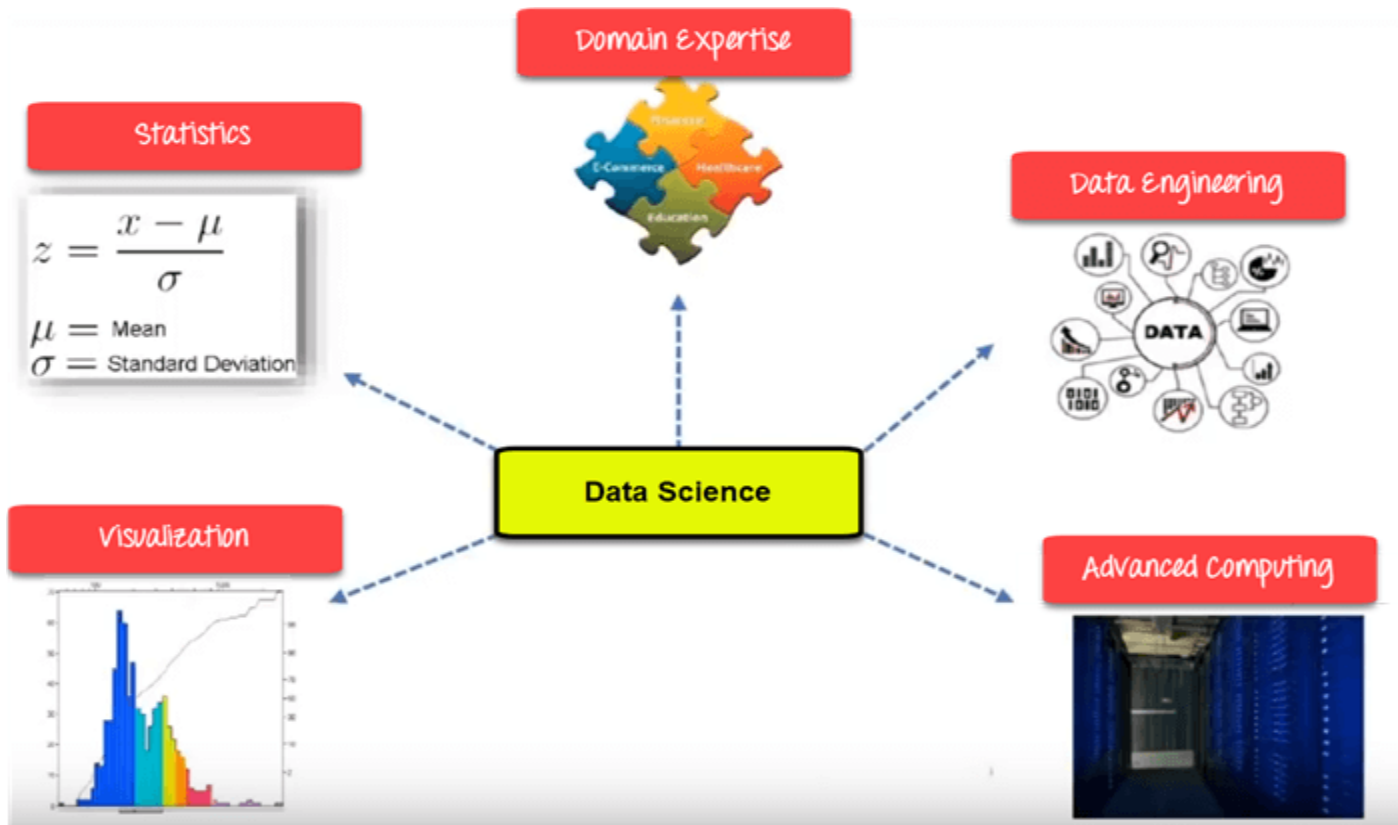
Types of Data

- Structured
- Unstructured
- Semi Structured

Digitalization and Digital Transformation

In the last decade, many businesses started to understand the importance of data when making business decisions. To apply data science principles to running a business, one first needs to collect some data, i.e. translate business processes into digital form. This is known as **digitalization**. Applying data science techniques to this data to guide decisions can lead to significant increases in productivity called **digital transformation**.

Data Science Components



Statistics & Visualization

- Statistics is the most critical unit of Data Science basics, and it is the method or science of collecting and analyzing numerical data in large quantities to get useful insights.
- Visualization technique helps us access huge amounts of data in easy to understand and digestible visuals.

Machine Learning & Deep Learning

- AI enables the machine to think, that is without any human intervention the machine will be able to take its own decision.
- Machine Learning is a subset of Artificial Intelligence that uses statistical learning algorithms to build systems that have the ability to automatically learn and improve from experiences without being explicitly programmed.
- Deep learning is a machine learning technique that is inspired by the way a human brain filters information, it is basically learning from examples. It helps a computer model to filter the input data through layers to predict and classify information.

Data Engineering

- Data engineering refers to the building of systems to enable the collection and usage of data. This data is usually used to enable subsequent analysis and data science; which often involves machine learning.

Advanced Computing

- **It** refers to systems with the ability to process data and perform calculations at high speeds, such as supercomputers.

Data Science Process



Discovery

Discovery step involves acquiring data from all the identified internal & external sources, which helps us answer the business question.

The data can be:

- Logs from web servers
- Data gathered from social media
- Census datasets
- Data streamed from online sources using APIs

Preparation

- Data can have many inconsistencies like missing values, blank columns, an incorrect data format, which needs to be cleaned. We need to process, explore, and condition data before modelling. The cleaner our data, the better are our predictions.

Model Planning

- In this stage, we need to determine the method and technique to draw the relation between input variables. Planning for a model is performed by using different statistical formulas and visualization tools. SQL analysis services, R, and SAS/access are some of the tools used for this purpose.

Model Building

- In this step, the actual model building process starts. Here, Data scientist distributes datasets for training and testing. Techniques like association, classification, and clustering are applied to the training data set. The model, once prepared, is tested against the “testing” dataset.

Operation

- We deliver the final baselined model with reports, code, and technical documents in this stage. Model is deployed into a real-time production environment after thorough testing.

Communicate Results

- In this stage, the key findings are communicated to all stakeholders. This helps us decide if the project results are success or failure based on the inputs from the model.

Tools for Data Science



SQL



MATLAB

Java

- Java can be used for many of the processes:
- Data import and export.
- Cleaning data.
- Statistical analysis.
- Machine learning and Deep learning.
- Deep learning.
- Text analytics (also known as Natural Language Processing or NLP).
- Data visualization.

R

- R is a popular programming language used for statistical computing and graphical presentation.
- Its most common use is to analyze and visualize data.

Why Use R?

- It is a great resource for data analysis, data visualization, data science and machine learning
- It provides many statistical techniques (such as statistical tests, classification, clustering and data reduction)
- It is easy to draw graphs in R, like pie charts, histograms, box plot, scatter plot, etc..
- It works on different platforms (Windows, Mac, Linux)
- It is open-source and free
- It has a large community support
- It has many packages (libraries of functions) that can be used to solve different problems

Python

- Python is open source, interpreted, high level language and provides great approach for object-oriented programming. It is one of the best language used by data scientist for various data science projects/application. Python provide great functionality to deal with mathematics, statistics and scientific function. It provides great libraries to deals with data science application.

SAS

SAS stands for **Statistical Analysis Software**. It was created in the year 1960 by the SAS Institute. From 1st January 1960, SAS was used for data management, and business intelligence, Since then, many new statistical procedures and components were introduced in the software.

Why we use SAS

- Data Management
- Statistical Analysis
- Report formation with perfect graphics
- Business Planning
- Application Development
- Data extraction
- Data transformation
- Data updation and modification

MATLAB

- MATLAB offers a notebook environment, toolboxes, and apps for developing analytic models.
- Using MATLAB we can combine statistics and machine learning with application specific techniques such as signal processing, image processing, text analytics, optimization and controls

Data Science Job Roles

- Data Scientist
- Data Engineer
- Data Analyst
- Statistician
- Data Architect
- Data Admin
- Business Analyst

- **Data Scientist:**

Role: A Data Scientist is a professional who manages enormous amounts of data to come up with compelling business visions by using various tools, techniques, methodologies, algorithms, etc.

Languages: R, SAS, Python, SQL, Hive, Matlab, Pig, Spark

- **Data Engineer:**

Role: The role of a data engineer is of working with large amounts of data. He develops, constructs, tests, and maintains architectures like large scale processing systems and databases.

Languages: SQL, Hive, R, SAS, Matlab, Python, Java, Ruby, C + +, and Perl

- **Data Analyst:**

Role: A data analyst is responsible for mining vast amounts of data. They will look for relationships, patterns, trends in data. Later he or she will deliver compelling reporting and visualization for analyzing the data to take the most viable business decisions.

Languages: R, Python, HTML, JS, C, C++ , SQL

- **Statistician:**

Role: The statistician collects, analyses, and understands qualitative and quantitative data using statistical theories and methods.

Languages: SQL, R, Matlab, Tableau, Python, Perl, Spark, and Hive

- **Data Administrator:**

Role: Data admin should ensure that the database is accessible to all relevant users. He also ensures that it is performing correctly and keeps it safe from hacking.

Languages: Ruby on Rails, SQL, Java, C#, and Python

- **Business Analyst:**

Role: This professional needs to improve business processes. He/she is an intermediary between the business executive team and the IT department.

Languages: SQL, Tableau, Power BI and, Python

Applications of Data Science

- Internet Search
- Recommendation System
- Image and Speech Recognition
- Online Price Comparison

Internet Search

- Every business today uses Data Science to understand what their customers want. In the same way, **Google uses Data Science to understand what its users want to know.** Google search uses Data science technology to search for a specific result within a fraction of a second

Recommendation System

- A recommendation system is an artificial intelligence or AI algorithm, usually associated with machine learning, that uses Big Data **to suggest or recommend additional products to consumers.** These can be based on various criteria, including past purchases, search history, demographic information, and other factors. To create a recommendation system. For example, “suggested friends” on Facebook or suggested videos” on YouTube, everything is done with the help of Data Science.

Image Recognition

- Data science tools with AI has the ability to not just assist users in face recognition but **help detect objects available in the camera.** The tools scan all the objects and attempt to name and identify them. Facebook recognizes our friend when we upload a photo with them, with the help of Data Science.

Speech Recognition

- Speech recognition is an interdisciplinary field that combines computer science and linguistics and provides machines with the ability to recognize human speech and translate it into text. Speech recognizes systems like Siri, Google Assistant, and Alexa run on the Data science technique

Online Price Comparison

- PriceRunner, Jungle, Shopzilla work on the Data science mechanism. Here, data is fetched from the relevant websites using APIs.

Challenges of Data Science Technology

- A high variety of information & data is required for accurate analysis
- Not adequate data science talent pool available
- Management does not provide financial support for a data science team
- Unavailability of/difficult access to data
- Business decision-makers do not effectively use data Science results
- Explaining data science to others is difficult
- Privacy issues
- Lack of significant domain expert
- If an organization is very small, it can't have a Data Science team

Thank u

