# Data Science

## Introduction
## Lecture No. 1

## Instructor: Abdul Rehman

10/23/20
25

# Books & Resources

υ Burkov, A. (2019). *The hundred-page machine learning book* (Vol. 1, p. 32). Quebec City, QC, Canada: Andriy Burkov.

υ Jordan, M. I., & Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, *349*(6245), 255-260.

υ Shalev-Shwartz, S., & Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge university press.
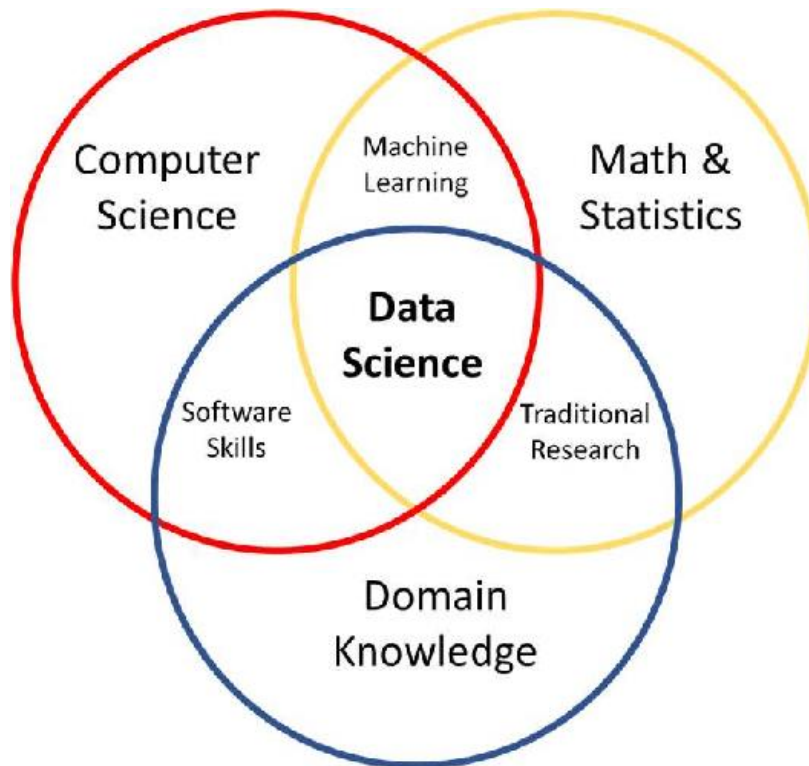
# Lecture Outlines

- ❑ What is Data Science?
- ❑ Key Components of Data Science
- ❑ Key Characteristics of Big Data (5 Vs)
- ❑ What is Machine Learning?
- ❑ Difference between AI and ML
- ❑ Applications of ML
- ❑ Dataset and their types (how to get datasets?)
- ❑ Introduction to Data and their types
- ❑ What are features and labels in dataset
- ❑ Skills Required for a Data Scientist

# What is Data Science?

υ **Data Science** is the process of extracting knowledge and insights from structured and unstructured data using techniques like statistics, mathematics, machine learning, data visualization, and domain knowledge.

υ **Data science** is a multidisciplinary field.

υ It is also used to solve business problem
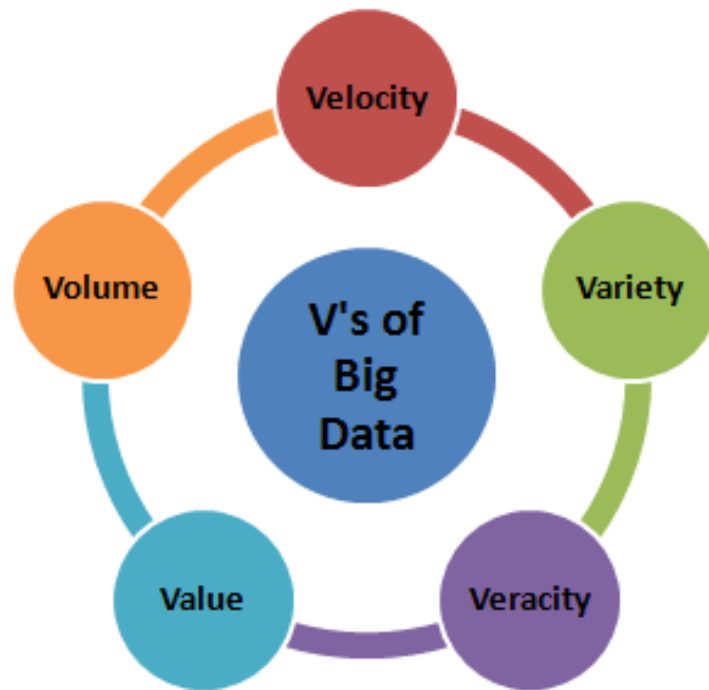
# Data Science is Multidisciplinary Field

# Key Characteristics of Big Data (5 Vs)

ʊ **Volume**: Massive amounts of data (terabytes to petabytes).

ʊ **Velocity**: High-speed data generation (real-time streaming data).

ʊ **Variety**: Different types of data (structured, unstructured, semi-structured).

ʊ **Veracity**: Ensuring data quality and accuracy.

ʊ **Value**: Extracting useful insights for decision-making.

# Key Components of Data Science

υ **Data Collection** – Gathering structured and unstructured data from various sources (databases, APIs, web scraping).

υ **Data Preprocessing** – Cleaning and transforming raw data into a usable format.

υ **Exploratory Data Analysis (EDA)** – Identifying patterns, trends, and outliers using visualization tools.

# Key Components of Data Science

υ **Machine Learning** – Building predictive models using algorithms like regression, classification, and clustering.

υ **Data Visualization** – Presenting findings through graphs, charts, and dashboards (using tools like Matplotlib, Power BI, and Tableau).

υ **Decision Making** – Using insights for business intelligence, automation, and strategy.

# Applications of Data Science

υ **Healthcare**

  υ Disease prediction and diagnosis

  υ Drug discovery and personalized medicine

  υ Medical image analysis (X-rays, MRIs, CTSCAN, etc.)

υ **Finance**

  υ Fraud detection and prevention

  υ Algorithmic trading and investment analysis

υ **Marketing**

  υ Customer segmentation and targeting

  υ Sentiment analysis on social media

  υ Marketing campaign optimization

# Applications of Data Science

υ **E-Commerce**

    υ Product recommendation systems

    υ Customer behavior analysis

υ **Transportation**

    υ Route optimization (like Google Maps)

    υ Autonomous vehicle navigation

υ **Education**

    υ Student performance prediction

υ **Sports Analytics**

    υ Player performance tracking

    υ Game strategy optimization

10/23/2025

# Machine Learning & Artificial Intelligence

ᴜ Machine Learning is the field of study that give computers the capability to learn without being explicitly programed.

ᴜ Machine can learn itself from past data and automatically improve with experience.

ᴜ Artificial Intelligence comprises two words "Artificial" and "Intelligence". Artificial refers to something which is made by humans or a non-natural thing and Intelligence means the ability to understand or think.

# Traditional Programming and Machine Learning



υ **Traditional Programming:** Rules + Data → Output.

υ **Machine Learning:** Data + Output → Rules (model).

# Examples of AI, ML, and DL

**Example Scenario:**

Let's say you want to develop a **self-driving car**:

- **AI:** The car should be able to drive itself intelligently like a human.
- **ML:** The car learns from driving data — when to brake, accelerate, or turn based on previous experiences.
- **DL:** The car uses deep learning models (like Convolutional Neural Networks) to recognize traffic signs, pedestrians, and lane markings from camera images.

# Uses cases of ML & AI

υ **For example:**

1. Fraud Detection (ML)
2. Email Filtering (ML)
3. Medical Diagnosis (ML)
4. Pattern Recognition (ML)
5. Face recognition, (ML)
6. Social media, (ML)
7. Movie recommendation systems, (ML)
8. Virtual assistants, (ML)

υ **For example:**

1. Autonomous vehicles, (AI) [Tesla, Cruise, Waymo, Drones]

2. Chabot, (AI) [Gimini, Tesla, ChatGPT, Siri, Xiaowei, Jeeney AI, etc.]

3. NLP, (AI) [Google Translator]

4. Robotics (AI & ML)

5. Self-driving Car (AI & ML)

# What is a Dataset

❑ A dataset is a collection of data organized in a structured format, often used for analysis, research, or training machine learning models.

❑ Datasets can contain various types of data, such as numerical values, text, images, or even videos, and are typically arranged in rows and columns (similar to a table).

❑ Each row represents

  ❑ an instance/record), and

❑ Each column represents

  ❑ a feature or attribute of the data

# Types of Dataset

- ❑ **Structured datasets**: Organized in a tabular form (like a spreadsheet or database).

- ❑ **Unstructured datasets:** Data without a predefined format, like images, videos, or text documents.

- ❑ **Semi-structured datasets**: Have some organization but aren't as strictly structured, such as JSON or XML files.

# Data Sets Sources

**Datasets:**
1. Kaggle Datasets
2. Amazon Datasets
3. UCI Machine Learning Repository
4. Google Data Search Engine
5. Microsoft Datasets
6. Awesome Public Dataset Collection
7. Scikit Learn Datasets
8. Computer Vision Datasets
9. Government Datasets (such as: https://www.data.gov/us)