# ANNAMACHARYA UNIVERSITY

**(ESTD UNDER AP PRIVATE UNIVERSITIES (ESTABLISHMENT AND REGULATION) ACT, 2016)**

**(UNIVERSITY LISTED IN UGC AS PER THE SECTION 2(f) OF THE UGC ACT, 1956)**

**RAJAMPET, Annamayya District, AP – 516126, INDIA**

# Sub Name: Predictive Analytics
# Sub Code: 23A3261T
# Branch: CSE(DS)
# Year: III B Tech II Sem

### Prepared by

B. Venkatesu Goud,

Assistant Professor, Dept. of AI&DS,

Annamacharya University, Rajampet.

# UNIT I

**Introduction to Predictive Analytics and Business Intelligence, Types of Predictive Models: Classification, Regression, Time Series, Supervised vs Unsupervised Learning, Predictive Modeling Workflow, Applications in Marketing, Finance, Healthcare, Challenges in Predictive Analytics.**

## => Introduction to Predictive Analytics and Business Intelligence

### What is Predictive Analytics?
Predictive analytics is a branch of data science that leverages statistical techniques, machine learning algorithms, and historical data to make data-driven predictions about future outcomes.

### Why Predictive Analytics is important?
Predictive analytics is important for several reasons:
- **Informed Decision-Making:** By anticipating future trends and outcomes, businesses and organizations can make more strategic decisions. Imagine being able to predict customer churn (when a customer stops using your service) or equipment failure before it happens. This allows for proactive measures to retain customers or prevent costly downtime.
- **Risk Management:** Predictive analytics helps identify and mitigate potential risks. For example, financial institutions can use it to detect fraudulent transactions, while healthcare providers can predict the spread of diseases.
- **Optimization and Efficiency:** Predictive models can optimize processes and resource allocation. Businesses can forecast demand and optimize inventory levels, or predict equipment maintenance needs to avoid disruptions.
- **Personalized Experiences:** Predictive analytics allows for personalization and customization. Retailers can use it to recommend products to customers based on their past purchases and browsing behavior.
- **Innovation and Competitive Advantage:** Predictive analytics empowers organizations to identify new opportunities and develop innovative products and services. By understanding customer needs and market trends, businesses can stay ahead of the competition.

### How Predictive Analytics Modeling works?

**1. Define a Problem:**
- Firstly data scientists or data analysts define the problem.
- Defining the problem means clearly expressing the challenge that the organization aims to focus using data analysis.
- A well- defined problem statement helps determine the appropriate predictive analytics approach to employ.

**2. Gather and Organize Data:**
- Once you define a problem statement it is important to acquire and organize data properly.
- Acquiring data for predictive analytics means collecting and preparing relevant information and data from various sources like databases, data warehouses, external

data providers, APIs, logs, surveys, and more that can be used to build and train predictive models.

### 3. Pre-process Data:

- Now after collecting and organizing the data, we need to pre-process data.
- Raw data collected from different sources is rarely in an ideal state for analysis. So, before developing a predictive models, data need to be pre-processed properly.
- Pre-processing involves cleaning the data to remove any kind of anomalies, handling missing data points and addressing outliers that could be caused by errors or input or transforming the data , which can be used for further analysis.
- Pre-processing ensures that data is of high quality and now the data is ready for model development.

### 4. Develop Predictive Models:

- Data scientists or data analysts leverage a range of tools or techniques to develop a predictive models based on the problem statement and the nature of the datasets.
- Now techniques like machine learning algorithms, regression models , decisions trees, neural networks are much among the common techniques for this.
- These models are trained on the prepared data to identify correlations and patterns that can be used for making predictions.

### 5. Validate and Deploy Results:

- After building the predictive model, validation is the critical steps to assess the accuracy and reliability of predictions.
- Data scientists rigorously evaluate the model's performance against known outcomes or test datasets.
- If required, modifications are implemented to improve the accuracy of the model.
- Once the model achieve satisfactory outcomes it can be deployed to deliver predictions to stakeholders.
- This can be done through applications, websites or data dashboards, making the insights easily accessible to decision makers or stakeholders.

## Predictive Analytics Techniques:

Predictive analytical models leverage historical data to anticipate future events or outcomes, employing several distinct types:

- **Classification Models**: These predict categorical outcomes or categorize data into predefined groups. Examples include Logistic Regression, Decision Trees, Random Forest, and Support Vector Machine.
- **Regression Models**: Used to forecast continuous outcome variables based on one or more independent variables. Examples include Linear Regression, Multiple Regression, and Polynomial Regression.
- **Clustering Models**: These group similar data points together based on shared characteristics or patterns. Examples comprise K-Means Clustering and Hierarchical Clustering.
- **Time Series Models**: Designed to predict future values by analyzing patterns in historical time-dependent data. Examples include Autoregressive Integrated Moving Average (ARIMA) and Exponential Smoothing Models.
- **Neural Networks Models**: Advanced predictive models capable of discerning complex data patterns and relationships. Examples encompass Feed Forward Neural Networks, Recurrent Neural Networks, and Convolutional Neural Networks.

## How Businesses Use Analytics?

1. **Operational Efficiency**: Analytics is employed to optimize operational processes and resource allocation, leading to cost savings and improved productivity.

2. **Customer Relationship Management (CRM)**: By analyzing customer data, businesses can gain insights into customer preferences, behavior, and satisfaction levels, allowing for more targeted and effective customer relationship management strategies.
3. **Risk Management**: Analytics is utilized to assess and mitigate various types of risks, such as financial, operational, and cybersecurity risks, enabling businesses to make proactive decisions to safeguard their assets and reputation.
4. **Supply Chain Management**: Analytics helps businesses analyze and forecast demand, optimize inventory levels, and improve logistics and distribution processes, leading to more efficient supply chain management and reduced costs.
5. **Financial Analysis and Planning**: Businesses use analytics for financial forecasting, budgeting, and performance analysis, enabling better financial decision-making and strategic planning.
6. **Human Resources Management**: Analytics assists in workforce planning, talent acquisition, performance management, and employee engagement initiatives, helping businesses attract, retain, and develop top talent.
7. **Product Development and Innovation**: By analyzing market trends, customer feedback, and competitor activities, businesses can identify opportunities for product development and innovation, ensuring the delivery of products and services that meet customer needs and preferences.
8. **Compliance and Regulatory Reporting**: Analytics aids businesses in monitoring compliance with regulations and reporting requirements, facilitating timely and accurate regulatory submissions and reducing the risk of non-compliance penalties.

## Benefits of Using Predictive Analytics

- **Improved Decision Making:** Predictive analytics enables businesses to make informed decisions by analyzing trends and patterns in historical data. This allows organizations to develop market strategies tailored to the insights gained from data analysis, leading to more effective decision-making processes.
- **Enhanced Efficiency and Resource Allocation**: By leveraging predictive analytics, businesses can optimize their operational processes and allocate resources more efficiently. This leads to cost savings, improved productivity, and better utilization of available resources.
- **Enhanced Customer Experience:** Predictive analytics enables businesses to enhance the customer experience by providing personalized product recommendations based on user behavior. By analyzing customer data, businesses can understand individual preferences and tailor their offerings accordingly, leading to increased customer satisfaction and loyalty.

## Applications of Predictive Analytics

Predictive analytics has a vast range of applications across different industries. Here are some key examples:

### Applications of Predictive Analytics in Business

- **Customer Relationship Management (CRM):** Predicting customer churn (customer leaving), recommending products based on past purchases, and personalizing marketing campaigns.
- **Supply Chain Management:** Forecasting demand for products, optimizing inventory levels, and predicting potential disruptions in the supply chain.
- **Fraud Detection:** Identifying fraudulent transactions in real-time for financial institutions and e-commerce platforms.

### Applications of Predictive Analytics in Finance

- **Credit Risk Assessment:** Predicting the likelihood of loan defaults to make informed lending decisions.
- **Stock Market Analysis:** Identifying trends and patterns in stock prices to inform investment strategies.
- **Algorithmic Trading:** Using models to automate trading decisions based on real-time market data.

**Applications of Predictive Analytics in Healthcare**
- **Disease Outbreak Prediction:** Identifying potential outbreaks of infectious diseases to enable early intervention.
- **Personalized Medicine:** Tailoring treatment plans to individual patients based on their genetic makeup and medical history.
- **Readmission Risk Prediction:** Identifying patients at high risk of being readmitted to the hospital to improve patient care and reduce costs.

**Applications of Predictive Analytics in Other Industries**
- **Manufacturing:** Predicting equipment failures for preventive maintenance, optimizing production processes, and improving product quality.
- **Insurance:** Tailoring insurance premiums based on individual risk profiles and predicting potential claims.
- **Government:** Predicting crime rates for better resource allocation and crime prevention strategies.


## The Future of Predictive Analytics

The future of predictive analytics is brimming with exciting possibilities fueled by advancements in technology and a growing focus on responsible use. Here's a glimpse into what we can expect:

- **Enhanced Accuracy and Real-Time Capabilities**
  - **Advanced AI and Machine Learning:** As Artificial Intelligence (AI) and machine learning algorithms become more sophisticated, predictive models will achieve even greater accuracy. This will lead to more reliable and nuanced predictions across various fields.
  - **Real-Time Data Integration:** The increasing availability of real-time data streams will allow models to adapt and update continuously. This ensures predictions stay relevant and reflect the ever-changing dynamics of the world.
- **Prescriptive Analytics Taking Center Stage**
  - **Beyond Predictions to Actionable Insights:** Predictive analytics will evolve beyond just forecasting what will happen. We'll see a rise in prescriptive analytics, which suggests specific actions to optimize outcomes based on predictions.
  - **Decision Support Systems:** Predictive models will be integrated with decision support systems, providing real-time recommendations and guidance to users.
- **Democratization of Predictive Analytics**
  - **Cloud-Based Solutions and User-Friendly Tools:** Cloud-based solutions and user-friendly interfaces will make predictive analytics more accessible to a wider range of organizations, even those without extensive data science expertise.

- o **Rise of Citizen Data Scientists:** With user-friendly tools, more business users will be empowered to leverage the power of predictive analytics for data-driven decision making within their specific roles.
- **Ethical Considerations and Responsible Use**
  - o **Focus on Data Privacy and Security:** As the use of personal data in analytics grows, ensuring data privacy and security will be paramount. Regulations and best practices will continue to evolve to protect individuals.
  - o **Addressing Bias and Fairness:** Mitigating bias in data and algorithms will be crucial to ensure fair and responsible use of predictive analytics across different demographics and social groups.
- **Impact on Society**
  - o **Shaping the Future with Data-Driven Insights:** Predictive analytics will play a significant role in shaping various aspects of society. From personalized healthcare and education to urban planning and environmental sustainability, data-driven insights will guide decision-making for a better future.

## Analytics Vs Machine Learning

- Analytics involves examining data to derive insights and make informed decisions based on historical information.
- Machine learning, a subset of artificial intelligence, focuses on developing algorithms that enable computers to learn from data and make predictions or decisions without explicit programming.
- While analytics often involves descriptive and diagnostic analysis, machine learning emphasizes predictive and prescriptive modeling.
- Analytics typically involves statistical methods and data visualization techniques, while machine learning utilizes algorithms such as decision trees, neural networks, and support vector machines.
- Analytics is broader in scope and encompasses various techniques for data analysis, while machine learning specifically focuses on algorithms that improve with experience and data.
- Both analytics and machine learning play crucial roles in extracting value from data, with analytics providing insights and machine learning enabling automation and prediction.

### Difference between Business Intelligence and Predictive Analytics :

| S.No. | Business Intelligence | Predictive Analytics |
|---|---|---|
| 1. | It is about descriptive analytics or searching at what happened. | It is about discovering hidden patterns the use of complicated algorithms that help to predict future outputs. |
| 2. | With the assist of BI, raw data can be processed to information related to the product, client, region, the quarter for income etc. | Using Predictive analytics, raw data is converted into "cleaned Data" by algorithms. |
| 3. | It helps people to assist in making decisions on what they can do for getting into insights. | In this, the model tells us the best decision in a particular situation based on existing data. |

| S.No. | Business Intelligence | Predictive Analytics |
|-------|----------------------|----------------------|
| 4. | It helps people to get insights to solve any business problem. | It helps to detect the complex problem with the help of algorithms. |
| 5. | It has Ad-hoc reporting technology, alerts technology etc. | It includes technologies like predictive modelling, forecasting etc. |
| 6. | It has data types such as structured data and data sets which can be managed. | It contains both structured and unstructured data and has massive data sets. |
| 7. | Some Common types of questions include How many did we sell? | Common questions in predictive analytics are what will be the optimum case for our business? Why this is happening? |

## => Types of Predictive Models: Classification, Regression, Time Series

**What is Predictive Modeling?**

Predictive modeling is a statistical technique used to predict the outcome of future events based on historical data. It involves building a mathematical model that takes relevant input variables and generates a predicted output variable. Machine learning algorithms are used to train and improve these models to help you make better decisions. Predictive modeling is used in many industries and applications and can solve a wide range of issues, such as fraud detection, customer segmentation, disease diagnosis, and stock price prediction.

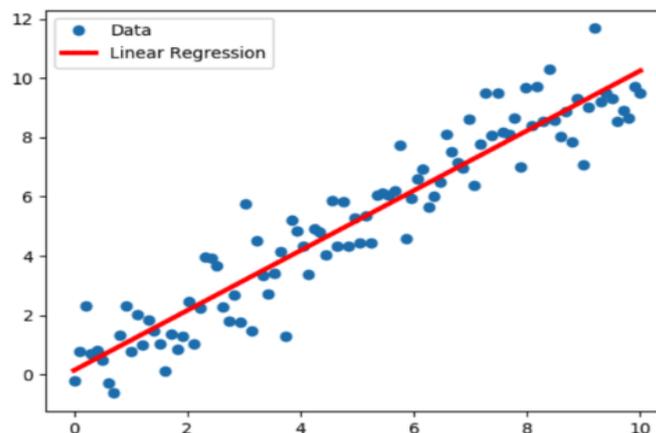**Model Types and Algorithms**

The chart below lists the 7 key types of predictive models and provides examples of predictive modeling techniques or algorithms used for each type. The two most commonly employed predictive modeling methods are regression and neural networks. The accuracy of predictive analytics and every predictive model depends on several factors, including the quality of your data, your choice of variables, and your model's assumptions.

| Predictive Model Types | Predictive Modeling Techniques |
|------------------------|--------------------------------|
| 1. Regression | Linear regression, polynomial regression, and logistic regression. |
| 2. Neural network | Multilayer perceptron (MLP), convolutional neural networks (CNN), recurrent neural networks (RNN), backpropagation, feedforward, autoencoder, and Generative Adversarial Networks (GAN). |
| 3. Classification | Decision trees, random forests, Naive Bayes, support vector machines (SVM), and k-nearest neighbors (KNN). |

| | |
|---|---|
| 4. Clustering | K-means clustering, hierarchical clustering, and density-based clustering. |
| 5. Time series | Autoregressive integrated moving average (ARIMA), exponential smoothing, and seasonal decomposition. |
| 6. Decision tree | Classification and Regression Trees (CART), Chi-squared Automatic Interaction Detection (CHAID), ID3, and C4.5. |
| 7. Ensemble | Bagging, boosting, stacking, and random forest. |

## 1. Regression

Regression models are used to predict a continuous numerical value based on one or more input variables. The goal of a regression model is to identify the relationship between the input variables and the output variable, and use that relationship to make predictions about the output variable. Regression models are commonly used in various fields, including financial analysis, economics, and engineering, to predict outcomes such as sales, stock prices, and temperatures.
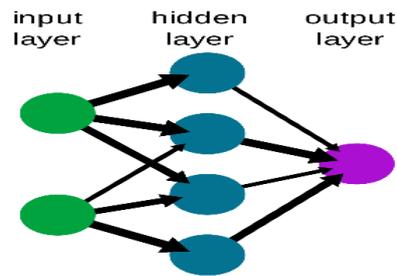


### Regression model algorithms:
**Linear regression** models assume that there is a linear relationship between the input variables and the output variable.
**Polynomial regression** models assume a non-linear relationship between input and output.
**Logistic regression** models are used for binary classification problems, where the output variable is either 0 or 1.

## 2. Neural Network

Neural network models are a type of predictive modeling technique inspired by the structure and function of the human brain. The goal of these models is to learn complex relationships between input variables and output variables, and use that information to make predictions. Neural network models are often used in fields such as image recognition, natural language processing, and speech recognition, to make predictions such as object recognition, sentiment analysis, and speech transcription.
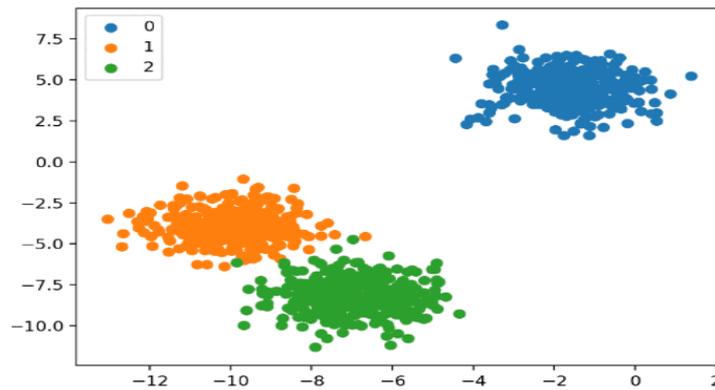
### Neural network model algorithms:

- **Multilayer Perceptron (MLP)** consists of multiple layers of nodes, including an input layer, one or more hidden layers, and an output layer. The nodes in each layer perform a mathematical operation on the input data, with the output of one layer serving as the input for the next layer. The weights between the nodes are adjusted during training using backpropagation to minimize the error between the predicted output and the actual output. MLP is a versatile algorithm that can be used for a wide range of predictive modeling tasks, including classification, regression, and pattern recognition.
- **Convolutional** neural networks (CNN) are commonly used for image recognition tasks, with each layer processing increasingly complex features of the image.
- **Recurrent** neural networks (RNN) are used for sequential data, such as natural language processing, and incorporate feedback loops that allow previous output to be used as input for the next prediction.
- **Long Short-Term Memory** (LSTM) is a type of RNN that addresses the vanishing gradient problem and is particularly useful for learning long-term dependencies in sequential data.
- **Backpropagation** is a common algorithm used to train neural networks by adjusting the weights between nodes in the network based on the error between the predicted output and the actual output.
- **Feedforward** neural networks consist of layers of nodes that process information from previous layers, with each node performing a mathematical operation on the input data.
- **Autoencoder** is used for unsupervised learning, where the network is trained to reconstruct the input data and can be used for tasks such as dimensionality reduction and anomaly detection.
- **Generative Adversarial Networks** (GAN) involves two neural networks, one that generates synthetic data and another that discriminates between real and synthetic data, and is commonly used for tasks such as image generation and data synthesis.

## 3. Classification

Classification models are used to classify data into one or more categories based on one or more input variables. Classification models identify the relationship between the input variables and the output variable, and use that relationship to accurately classify new data into the appropriate category. Classification models are commonly used in fields like marketing, healthcare, and computer vision, to classify data such as spam emails, medical diagnoses, and image recognition.
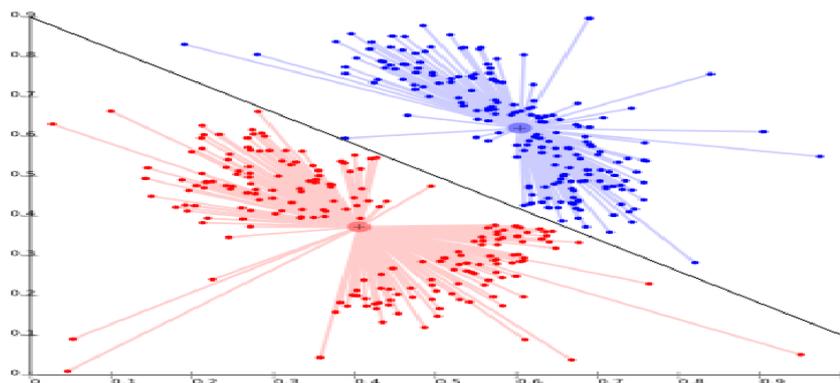
**Classification model algorithms:**

- **Decision trees** are a graphical representation of a set of rules used to make decisions based on a series of if-then statements.
- **Random forests** are an ensemble method that combines multiple decision trees to improve accuracy and reduce errors.
- **Naive Bayes** is a probabilistic model that assumes independence between input variables
- **Support vector machines (SVM)** and k-nearest neighbors (KNN) are distance-based models that use mathematical algorithms to classify data.

# 4. Clustering

Clustering models are used to group data points together based on similarities in their input variables. The goal of a clustering model is to identify patterns and relationships within the data that are not immediately apparent, and group similar data points into clusters. Clustering models are typically used for customer segmentation, market research, and image segmentation, to group data such as customer behavior, market trends, and image pixels.
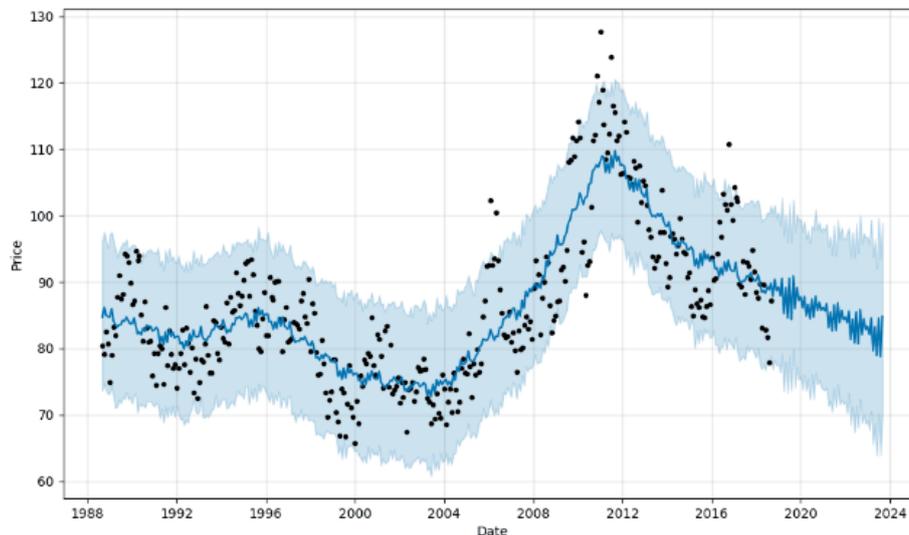


**Clustering model algorithms:**

- **K-means** clustering is a popular method that partitions the data into k clusters based on the distances between data points.
- **Hierarchical** clustering creates a tree-like structure of nested clusters based on the distances between data points.
- **Density-based** clustering groups data points based on their density in a particular area.

# 5. Time series

Time series models are used to analyze and forecast data that varies over time. Time series models help you identify patterns and trends in the data and use that information to make predictions about future values. Time series models are used in a wide variety of fields, including financial analytics, economics, and weather forecasting, to predict outcomes such as stock prices, GDP growth, and temperatures.
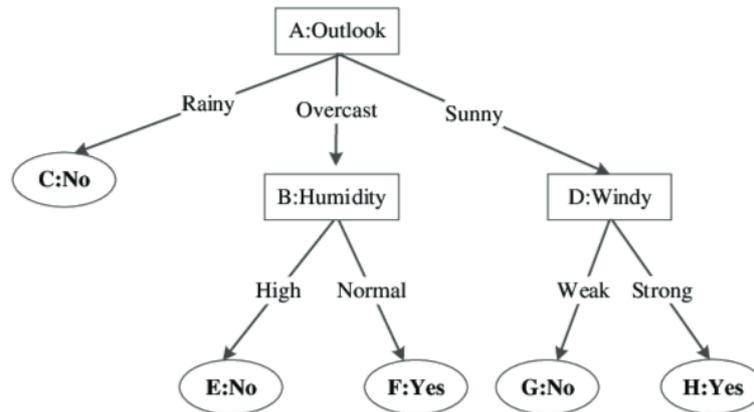


## Time series model algorithms:

- **ARIMA** (autoregressive integrated moving average) algorithms use previous values of a time series to predict future values, taking into account factors such as seasonality, trends, and stationarity.
- **Exponential smoothing** algorithms use a weighted average of past observations to predict future values, and are particularly useful for short-term forecasting.
- **Seasonal decomposition** algorithms decompose the time series into seasonal, trend, and residual components, and then use those components to make predictions.

## 6. Decision Tree

Decision tree models use a tree-like structure to model decisions and their possible consequences. The tree consists of nodes that represent decision points, with branches representing the possible outcomes or consequences of each decision. Each node corresponds to a predictor variable and each branch corresponds to a possible value of that variable. The goal of a decision tree model is to predict the value of a target variable based on the values of the predictor variables. The model uses the tree structure to determine the most likely outcome for a given set of predictor variable values.

Decision tree models can be used for both classification and regression tasks. In a classification tree, the target variable is categorical, while in a regression tree, the target variable is continuous. Decision tree models are easy to interpret and visualize, making them useful for understanding the relationships between predictor variables and the target variable. However, they can be prone to overfitting and may not perform as well as other predictive modeling techniques on complex datasets.
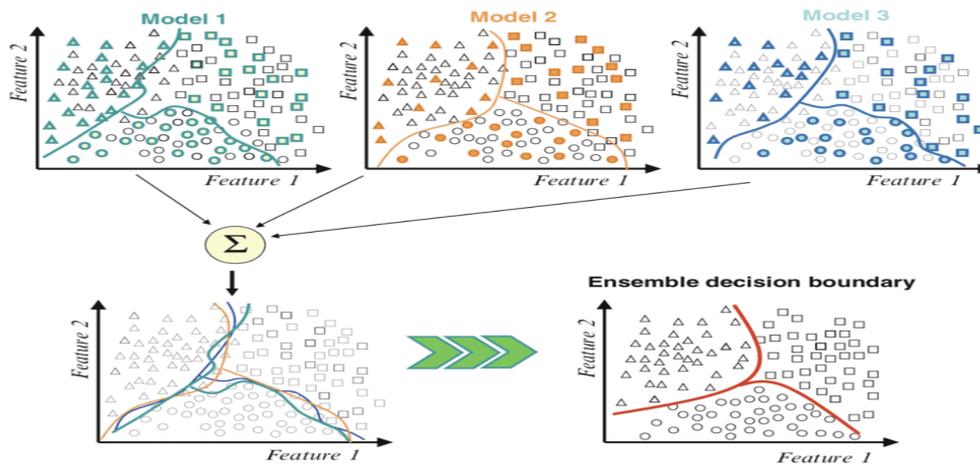
## Decision tree model algorithms:

- **CART** (Classification and Regression Tree) can be used for both classification and regression tasks. It uses Gini impurity as a measure of the quality of a split, aiming to minimize it. CART constructs binary trees, where each non-leaf node has two children.
- **CHAID** (Chi-squared Automatic Interaction Detection) is used for categorical variables and constructs trees based on chi-squared tests to determine the most significant associations between the predictor variables and the target variable. It can handle both nominal and ordinal categorical variables.
- **ID3** (Iterative Dichotomiser 3) is used to build decision trees for classification tasks. It selects the attribute with the highest information gain at each node to split the data into subsets. Information gain is calculated based on the entropy of the subsets.
- **C4.5** is an extension of the ID3 algorithm that can handle both categorical and continuous variables. It uses information gain ratio to select the splitting attribute, which takes into account the number of categories and their distribution in the subsets.

  These algorithms use various criteria to determine the optimal split at each node, such as information gain, Gini index, or chi-squared test.

## 7. Ensemble

Ensemble models combine multiple models to improve their predictive accuracy and stability. By combining multiple models, the errors and biases of individual models are usually reduced, leading to better overall performance. Ensemble models can be used for both classification and regression tasks and are well suited for data mining. They're often used in machine learning or AI competitions and real-world applications where high predictive accuracy is required.
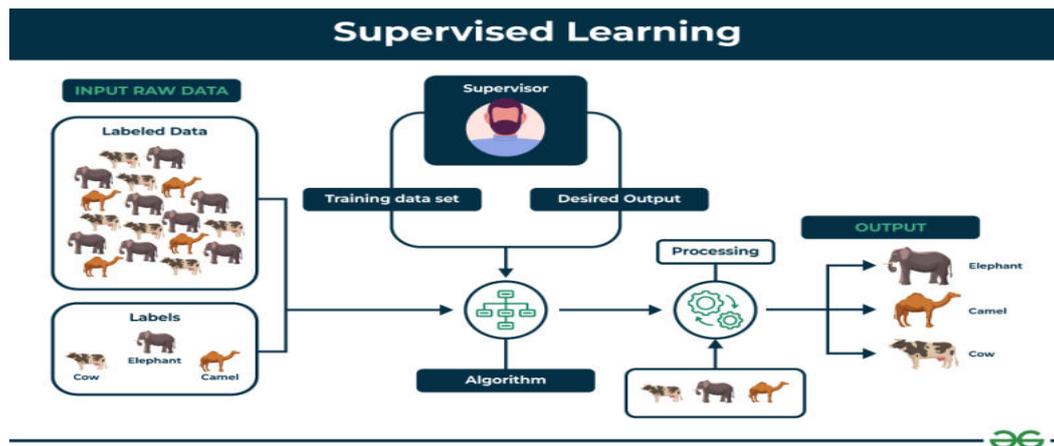
### Ensemble model algorithms:

- **Bagging** (Bootstrap Aggregating) involves creating multiple versions of the same prediction model on different subsets of the training data, and then aggregating their predictions to make the final prediction. Bagging is used to reduce the variance of a single model and improve its stability.
- **Boosting** involves creating multiple weak models sequentially, where each model tries to correct the errors of the previous model. Boosting is used to reduce the bias of a single model and improve its accuracy.
- **Stacking** involves training multiple models and using their predictions as input to a meta-model, which then makes the final prediction. Stacking is used to combine the strengths of multiple models and achieve better performance.
- **Random Forest** is an extension of bagging that uses decision trees as the base models. Random Forest creates multiple decision trees on different subsets of the training data, and then aggregates their predictions to make the final prediction.

## => Supervised vs Unsupervised Learning

### What is Supervised learning?

Supervised learning as the name suggests, works like a **teacher** or **supervisor** guiding the machine. In this approach we teach or train the machine using the labelled data(correct answers or classifications) which means each input has the correct output in the form of answer or category attached to it. After that machine is provided with a new set of examples (data) so that it can analyses the training data and produces a correct outcome from labeled data.

**For example**, a labeled dataset of images of Elephant, Camel and Cow would have each image tagged with either "**Elephant**", "**Camel**" or "**Cow.**"

## Example to Understand

Imagine we have a basket full of different fruits that we want the machine to identify. The machine first looks at the image of a fruit and extracts features like its shape, color and texture. Then it compares these features to the fruits it has already learned during training. If the new fruit's features closely match those of an apple, the machine will predict that the fruit is an apple.

For example, suppose we train the machine by showing it fruits one by one:

- If the fruit is round, has a small depression at the top and is red, it is labeled as an Apple.
- If the fruit is long, curved and greenish-yellow, it is labeled as a Banana.

Now after this training, if we give the machine a new fruit (say a banana) from the basket and ask it to identify it, the machine will use what it has learned during training. It will analyze the shape and color of the new fruit and classify it as a Banana placing it in the correct category. In this way, the machine learns from the training data (the basket with labeled fruits) and applies that knowledge to recognize new, unseen fruits.

## Types of Supervised Learning

Supervised learning is classified into two types of algorithms:

### 1. Regression

A regression is used to predict **continuous** values such as house prices, stock prices or temperature. Regression algorithms learn how to connect input data to a specific number or value.

Some common regression algorithms include:

1. Linear Regression
2. Polynomial Regression
3. Lasso Regression
4. Ridge Regression

### 2. Classification

A classification is used to predict **categorical** values such as whether a customer will buy or not, whether an email is spam or not or whether a medical image shows a tumor or not. Classification algorithms learn how to connect input data to the probability of belonging to different groups or categories.

Some of the most common classification algorithms include:

1. Logistic Regression
2. Support Vector Machines
3. Decision Trees
4. Random Forests
5. Naive Baye

**Applications of Supervised learning**

It can be used to solve variety of problems which includes:

1. **Image classification:** It can automatically classify images into different categories such as animals, objects or scenes helps in the tasks like image search, content moderation and image-based product recommendations.

2. **Medical diagnosis:** It can assist in medical diagnosis by analyzing patient data such as medical images, test results and patient history to identify patterns that suggest specific diseases or conditions.

3. **Fraud detection:** They can analyze financial transactions and identify patterns that shows fraudulent activity which helps financial institutions prevent fraud and protect their customers.

4. **Natural language processing (NLP):** It plays a important role in NLP tasks including sentiment analysis, machine translation and text summarization which enables machines to understand and process human language effectively.

**Advantages of Supervised learning**

1. It learns from labeled examples to make accurate predictions on new, unseen data.

2. With more data and training, these models increases their accuracy which leads to better performance and more reliable predictions.

3. It works well for many tasks from detecting spam emails to predicting house prices as it has the ability to handle various computational challenges.

4. It can handle both classification (sorting data into categories) and regression (predicting numbers) which makes it flexible for different problems.
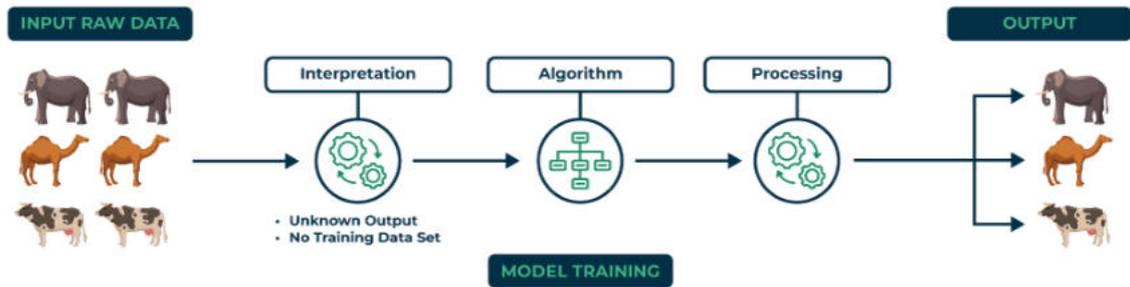
**Disadvantages of Supervised learning**

1. It requires a well-labeled dataset where each input has a corresponding output. Creating such datasets takes a lot of time, money and effort and can sometimes have mistakes, this makes supervised learning hard to use.

2. It works well on many tasks but can struggle with very complex or unstructured problems like understanding patterns or abstract ideas that doesn't relate to what it was trained on.

3. These models can sometimes overfit the training data which means they perform well on training data but poor on new, unseen data.

4. These models often need constant updating with new labeled data to stay accurate as real-world data changes over time.

# What is Unsupervised learning?

Unsupervised learning is a part of machine learning which works differently from supervised because there is no teacher(supervisor) involved to guide the machine. In this approach the machine is given with data that has **no labels or categories**. It analyzes the data on its own to find patterns, groups or relationships without any prior knowledge. The machine learns by discovering hidden structures within the data without being told what the correct output should be.

For example, unsupervised learning can analyze animal data and group the animals by their traits and behavior. These groups might represent different species which allows the machine to organize animals without any prior labels or categories.

**Example to understand**

Imagine we have a machine learning model trained on many unlabeled images of dogs and cats. The model has never seen any labeled example that says "dog" or "cat" before so it doesn't know how these animals look like.

Now, if we give the model a new image that contains both dogs and cats it won't be able to directly label them as "dog" or "cat." It will group parts of the image based on similarities and differences in features like shape or texture. It might separate the image into two groups one with dog-like features and other with cat-like features.

This happens because unsupervised learning doesn't rely on prior knowledge or training with labeled data. It finds patterns and organizes data on its own helps in discovering information that wasn't given before.

## Types of Unsupervised Learning

Unsupervised learning is divided into two categories of algorithms:

### 1. Clustering

A clustering is used to group similar data points together. Clustering algorithms work by repeatedly moving data points closer to to the center of their group (cluster) and farther from points in other groups. This helps the algorithm to create clear and meaningful clusters. Some popular clustering algorithms include:

1. K-means clustering
2. Hierarchical clustering
3. Principal Component Analysis (PCA)
4. Singular Value Decomposition (SVD)
5. Independent Component Analysis
6. Gaussian Mixture Models (GMMs)
7. Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

### 2. Association rule learning

An association rule learning used to find patterns and relationships between different items in a dataset. It looks for rules like "people who buy X often also buy Y".

Some common Association rule learning algorithms include:

1. Apriori Algorithm
2. Eclat Algorithm
3. FP-Growth Algorithm

## Application of Unsupervised learning

Unsupervised learning can be used to solve a variety of problems which includes:

1. **Anomaly detection**: It can identify unusual patterns or behaviors in data helps in the detection of fraud, security breaches or system problems.

2. **Scientific discovery**: It can show hidden relationships and patterns in scientific data which leads to new insights and ideas.
3. **Recommendation systems**: It finds similarities in user behavior and preferences to recommend products, movies or music that align with their interests.
4. **Customer segmentation**: It can identify groups of customers with similar characteristics which allows businesses to target marketing campaigns and improve customer service more effectively.

**Advantages of Unsupervised learning**
1. It doesn't need labeled data so we can start working with large datasets more easily and quickly.
2. This handles large amounts of data and reduce it into simpler forms without losing important patterns which makes it manageable and efficient.
3. It discovers patterns and relationships in the data that were previously unknown which offers valuable insights.
4. By analyzing unlabeled data, it shows meaningful trends and groups that help us to understand our data deeply.

**Disadvantages of Unsupervised learning**
1. Without labeled answers, it's difficult to tell how accurate or effective the model is.
2. Lack of clear guidance can lead to less precise results for complex problems.
3. After grouping the data, we may needs to check and label these groupings which can be time-consuming.
4. Missing data, outliers or noise in the data can easily affect the quality of the results.

## Supervised vs Unsupervised Machine Learning

| Parameters | Supervised machine learning | Unsupervised machine learning |
|---|---|---|
| Input Data | They are trained on labeled data. | They are trained on unlabeled data. |
| Computational Complexity | Simpler method | Computationally complex |
| Accuracy | Highly accurate | Less accurate |
| No. of classes | No. of classes is known | No. of classes is not known |
| Data Analysis | Uses offline analysis | Uses real-time analysis of data |
| Algorithms used | Linear and Logistics regression, KNN Random forest, multi-class classification, decision tree, Support Vector Machine, Neural Network etc. | K-Means clustering, Hierarchical clustering, Apriori algorithm etc. |

| Parameters | Supervised machine learning | Unsupervised machine learning |
|---|---|---|
| Output | Desired output is given. | Desired output is not given. |
| Training data | Use training data to infer model. | No training data is used. |
| Complex model | It is not possible to learn larger and more complex models with supervised learning. | It is possible to learn larger and more complex models with unsupervised learning. |
| Model | We can test our model. | We can not test our model. |
| Supervision | Supervised learning needs supervision to train the model. | Unsupervised learning does not need any supervision to train the model. |
| Classification | Divided into two types: <br> 1. Regression <br> 2. Classification | Divided into two types: <br> 1. Clustering <br> 2. Association |
| Feedback | It has feedback mechanism. | It has no feedback mechanism. |
| Time Consumption | It's more time consuming. | It's less time consuming. |
| Example | Optical character recognition. | Find a face in an image. |

## => **Predictive Modeling Workflow**

Predictive Modeling Workflow is a systematic sequence of steps used to build, evaluate, and deploy predictive models. It ensures that data is properly prepared, appropriate algorithms are selected, and the model performs reliably on real-world data.

**1. Problem Definition**
Clearly define the business or research problem.
- What do we want to predict? (e.g., sales, churn, classification)
- Is the task **classification**, **regression**, or **clustering**?
- What are the success metrics?

## 2. Data Collection

Gather data from various sources:

- Databases
- Excel/CSV files
- APIs / Web scraping
- Sensors / Logs

Key tasks:

- Identify relevant variables
- Understand data availability and limitations

## 3. Data Preprocessing

Prepare raw data for modeling.

- Handle missing values
- Remove duplicates
- Fix inconsistent or incorrect entries
- Encode categorical variables
- Normalize or scale numerical variables

## 4. Exploratory Data Analysis (EDA)

Understand the data through visualization and statistics.

- Discover patterns and relationships
- Identify outliers
- Perform correlation analysis
- Visualize distributions, trends, clusters

Output: insights that influence feature selection and model choice.

## 5. Feature Engineering

Create meaningful input variables.

- Feature creation (ratios, aggregates)
- Feature transformation (log, standardization)
- Feature selection (removing irrelevant features)
- Dimensionality reduction (PCA)

Goal: improve model accuracy and efficiency.

## 6. Model Selection

Choose appropriate algorithms depending on the problem:

- **Classification** → Logistic Regression, Decision Tree, SVM, Random Forest, Neural Networks
- **Regression** → Linear Regression, Ridge/Lasso, Gradient Boosting
- **Time Series** → ARIMA, Prophet, LSTM

## 7. Model Training

Train the model using the training dataset.

- Fit the algorithm on features and target
- Use cross-validation to avoid overfitting

## 8. Model Evaluation

Check how the model performs on unseen data (test set).
**Metrics depend on the task:**

- Classification → Accuracy, Precision, Recall, F1, AUC
- Regression → RMSE, MAE, R²
- Time Series → MAPE, RMSE

## 9. Hyperparameter Tuning
Improve performance by adjusting algorithm parameters.
- Grid Search
- Random Search
- Bayesian Optimization

Goal: find the best model settings.

## 10. Model Deployment
Deploy the model into production.
- Integrate with applications
- Create APIs
- Schedule automated predictions
- Monitor performance over time

## 11. Model Monitoring & Maintenance
Models degrade over time due to data changes (concept drift).
Monitoring includes:
- Tracking accuracy
- Re-training with new data
- Updating features or algorithms

## Summary Diagram
**Define Problem → Collect Data → Preprocess → EDA → Feature Engineering → Model Selection → Training → Evaluation → Tuning → Deployment → Monitoring**

# => Applications in Marketing, Finance, Healthcare

## 1. Predictive Analytics in Marketing

### a) Customer Segmentation
- Clustering algorithms group customers based on behavior, demographics, and preferences.
- Used for personalized marketing.

### b) Customer Churn Prediction
- Identifies customers likely to stop using a product/service.
- Helps companies retain customers by targeting them with offers.

### c) Lead Scoring
- Predicts the probability that a lead will convert into a customer.
- Sales teams prioritize high-scoring leads.

### d) Recommendation Systems
- Suggest products based on customer history.
- Used by Amazon, Flipkart, Netflix, etc.

### e) Marketing Campaign Optimization
- Predicts which customers will respond to promotions.
- Helps allocate budgets more effectively.

### f) Demand Forecasting

- Predicts future product demand using time-series forecasting.
- Prevents overstocking or understocking.

## 2. Predictive Analytics in Finance

### a) Credit Scoring
- Predicts the likelihood of a customer defaulting on a loan.
- Used by banks and credit institutions.

### b) Fraud Detection
- Machine learning models detect unusual patterns indicating fraudulent transactions.
- Real-time detection systems reduce financial losses.

### c) Algorithmic Trading
- Predictive models analyze market trends to make automatic trading decisions.
- Uses time-series forecasting, deep learning, sentiment analysis.

### d) Risk Management
- Estimates financial risks (market risk, credit risk, operational risk).
- Models simulate stress-test conditions.

### e) Customer Lifetime Value (CLV) Prediction
- Predicts the long-term value of a customer.
- Helps in planning cross-selling and up-selling strategies.

### f) Loan Default Prediction
- Banks use predictive analytics to identify risky borrowers.

## 3. Predictive Analytics in Healthcare

### a) Disease Prediction & Diagnosis
- Predicts diseases like diabetes, cancer, heart disease using patient history and lab data.
- Early detection saves lives.

### b) Hospital Readmission Prediction
- Identifies patients at high risk of returning within 30 days.
- Helps improve treatment and reduce costs.

### c) Predictive Patient Monitoring
- Real-time analytics on ICU data predicts deterioration.
- Wearables forecast health events (e.g., seizures, cardiac issues).

### d) Personalized Treatment Plans
- Machine learning recommends treatment based on patient similarity.

### e) Drug Discovery & Development
- Predicts chemical compound effectiveness.
- Reduces time and cost of research.

### f) Resource & Staff Management
- Forecasts patient load, emergency visits, bed occupancy.
- Helps hospitals optimize resource allocation.

## Summary Table

| Industry | Key Applications | Benefits |
|---|---|---|
| Marketing | Segmentation, Churn Prediction, Recommendation Systems, Lead Scoring | Better targeting, higher sales, customer retention |

| Industry | Key Applications | Benefits |
|---|---|---|
| Finance | Fraud Detection, Credit Scoring, Algorithmic Trading, Risk Management | Reduced losses, better decisions, improved profitability |
| Healthcare | Disease Prediction, Patient Monitoring, Readmission Prediction, Personalized Care | Improved patient outcomes, reduced cost, faster diagnosis |

# => **Challenges in Predictive Analytics**

Predictive Analytics offers powerful insights, but real-world implementation comes with technical, data-related, operational, and ethical challenges. Understanding these challenges helps organizations build reliable and responsible predictive models.

**1. Data Quality Issues**
Poor-quality data is the biggest challenge.
- Missing values
- Noise and outliers
- Incorrect or inconsistent data entries
- Duplicate records

**Impact:** Models become inaccurate and unreliable.

**2. Data Availability & Accessibility**
Predictive analytics requires large volumes of historical data.
- Limited or insufficient data
- Difficulty obtaining data from multiple sources
- Data stored in incompatible formats

**Impact:** Restricts model accuracy and scalability.

**3. Data Integration Challenges**
Organizations often have data spread across:
- Databases
- Cloud storage
- Legacy systems
- Real-time streaming sources

Integrating these different systems is difficult and time-consuming.

**4. Feature Engineering Complexity**
Selecting and designing useful features requires:
- Domain knowledge
- Technical skills
- Time-consuming experimentation

**Impact:** Poor feature selection leads to weak model performance.

**5. Model Overfitting & Underfitting**
- **Overfitting:** Model learns noise, performs poorly on new data
- **Underfitting:** Model is too simple, fails to learn patterns

Balancing the complexity of the model is a key challenge.

**6. Model Interpretability**
Advanced models (e.g., deep learning, ensemble methods) are often "black boxes."
- Hard to explain predictions to stakeholders
- Difficult to justify decisions in sensitive areas (finance, healthcare)

Regulators increasingly require explainable AI.

## 7. Scalability & Computational Constraints

Predictive analytics with large datasets requires:

- High computing power
- Efficient storage
- Fast processing algorithms

Small organizations may lack these resources.

## 8. Real-time Prediction Challenges

Implementing real-time analytics needs:

- Stream processing systems
- Low-latency infrastructure
- Continuous updating

Many businesses struggle to handle real-time data efficiently.

## 9. Model Deployment Issues

Even well-performing models face challenges after deployment:

- Integration with existing systems
- API management
- Performance monitoring
- Model drift over time

Ensuring the model remains accurate in production is crucial.

## 10. Data Privacy & Security Concerns

Predictive analytics involves sensitive data:

- Personal information
- Financial records
- Medical data

Challenges include:

- GDPR, HIPAA compliance
- Secure storage and transmission
- Preventing unauthorized access

## 11. Ethical & Bias Issues

Models may unintentionally reproduce biases found in training data.

- Gender bias
- Racial bias
- Socioeconomic bias

**Impact:** Unfair or discriminatory outcomes.

## 12. High Cost & Skill Requirements

Predictive analytics requires:

- Skilled data scientists
- Expensive tools and infrastructure
- Time for model development and tuning

Many organizations struggle with budget and talent shortages.

# UNIT II

**Data Cleaning: Handling Missing, Noisy, and Inconsistent Data, Feature Selection and Dimensionality Reduction (PCA, LDA), Feature Scaling: Normalization, Standardization, Encoding Categorical Variables, Feature Extraction and Construction, Dealing with Imbalanced Datasets.**

## => Data Cleaning: Handling Missing:

### Identify the Missing Data Values

Most analytics projects will encounter three possible types of missing data values, depending on whether there's a relationship between the missing data and the other data in the dataset:

- **Missing completely at random (MCAR):** In this case, there may be no pattern as to why a column's data is missing. For example, survey data is missing because someone could not make it to an appointment, or an administrator misplaces the test results he is supposed to enter into the computer. The reason for the missing values is unrelated to the data in the dataset.
- **Missing at random (MAR):** In this scenario, the reason the data is missing in a column can be explained by the data in other columns. For example, a school student who scores above the cutoff is typically given a grade. So, a missing grade for a student can be explained by the column that has scores below the cutoff. The reason for these missing values can be described by data in another column.
- **Missing not at random (MNAR):** Sometimes, the missing value is related to the value itself. For example, higher income people may not disclose their incomes. Here, there is a correlation between the missing values and the actual income. The missing values are not dependent on other variables in the dataset.

### How to Handle Missing Data Values

Data teams can use a number of strategies to handle missing data. On one hand, algorithms such as random forest and KNN are robust in dealing with missing values.

On the other hand, you may have to deal with missing data on your own. The first common strategy for dealing with missing data is to **delete the rows with missing values.** Typically, any row which has a missing value in any cell gets deleted. However, this often means many rows will get removed, leading to loss of information and data. Therefore, this method is typically not used when there are few data samples.

You can also **impute the missing data.** This can be based solely on information in the column that has missing values, or it can be based on other columns present in the dataset.

Finally, you can use **classification or regression models** to predict missing values.

Let's look at these three strategies in depth:

*1. Missing Values in Numerical Columns*

The first approach is to replace the missing value with one of the following strategies:

- Replace it with a constant value. This can be a good approach when used in discussion with the domain expert for the data we are dealing with.
- Replace it with the mean or median. This is a decent approach when the data size is small—but it does add bias.
- Replace it with values by using information from other columns.

B Venkatesu Goud, Assistant Professor

In the employee dataset subset below, we have salary data missing in three rows. We also have State and Years of Experience columns in the dataset:

| Row no | State | Salary | Yrs of Experience |
|--------|-------|--------|-------------------|
| 1 | NY | 57400 | Mid |
| 2 | TX | | Entry |
| 3 | NJ | 90000 | High |
| 4 | VT | 36900 | Entry |
| 5 | TX | | Mid |
| 6 | CA | 76600 | High |
| 7 | NY | 85000 | High |
| 8 | CA | | Entry |
| 9 | CT | 45000 | Entry |

## Missing values

The first approach is to fill the missing values with the mean of the column. Here, we are solely using the information from the column which has missing values:

| Row no | State | Salary | Yrs of Experience |
|--------|-------|--------|-------------------|
| 1 | NY | 57400 | Mid |
| 2 | TX | 65150 | Entry |
| 3 | NJ | 90000 | High |
| 4 | VT | 36900 | Entry |
| 5 | TX | 65150 | Mid |
| 6 | CA | 76600 | High |
| 7 | NY | 85000 | High |
| 8 | CA | 65150 | Entry |
| 9 | CT | 45000 | Entry |

## Replaced with the mean salary

With the help of a domain expert, we can do little better by using information from other columns in the dataset. The average salary is different for different states, so we can use that to fill in the values. For example, calculate the average salary of people working in Texas and replace the missing data with an average salary of people who typically work in Texas:

| Row no | State | Salary | Yrs of Experience |
|---:|---|---:|---|
| 1 | NY | 57400 | Mid |
| 2 | TX | 45000 | Entry |
| 3 | NJ | 90000 | High |
| 4 | VT | 36900 | Entry |
| 5 | TX | 45000 | Mid |
| 6 | CA | 76600 | High |
| 7 | NY | 85000 | High |
| 8 | CA | 55000 | Entry |
| 9 | CT | 45000 | Entry |

Replaced with mean salary in TX

Replaced with mean salary in CA

What else can we do better? How about making use of the Years of Experience column as well? Calculate the average entry-level salary of people working in Texas and replace the row where the salary is missing for an entry-level person in Texas. Do the same for the mid-level and high-level salaries:

Replaced with mean entry level salary in TX

| Row no | State | Salary | Yrs of Experience |
|---:|---|---:|---|
| 1 | NY | 57400 | Mid |
| 2 | TX | 35000 | Entry |
| 3 | NJ | 90000 | High |
| 4 | VT | 36900 | Entry |
| 5 | TX | 48000 | Mid |
| 6 | CA | 76600 | High |
| 7 | NY | 85000 | High |
| 8 | CA | 43000 | Entry |
| 9 | CT | 45000 | Entry |

Replaced with mean Mid level salary in TX

Replaced with mean Entry level salary in CA

Note that there are some boundary conditions. For example, there might be a row that has missing values in **both** the Salary and Years of Experience columns. There are multiple ways to handle this, but the most straightforward is to replace the missing value with the average salary in Texas.

*2. Predicting Missing Values Using an Algorithm*
Another way to predict missing values is to create a simple regression model. The column to predict here is the Salary, using other columns in the dataset. If there are missing values in the input columns, we must handle those conditions when creating the predictive model. A

simple way to manage this is to choose only the features that do not have missing values, or take the rows that do not have missing values in any of the cells.

## 3. Missing Values in Categorical Columns

Dealing with missing data values in categorical columns is a lot easier than in numerical columns. Simply replace the missing value with a constant value or the most popular category. This is a good approach when the data size small, though it does add bias.

For example, say we have a column for Education with two possible values: High School and College. If there are more people with a college degree in the dataset, we can replace the missing value with College Degree:

| Row no | State | Education |
|---|---|---|
| 1 | NY | High School |
| 2 | TX | |
| 3 | NJ | High School |
| 4 | VT | High School |
| 5 | TX | |
| 6 | CA | College |
| 7 | NY | High School |
| 8 | CA | |
| 9 | CT | College |

## Missing values

We can tweak this more by making use of information in the other columns. For example, if there are more people from Texas with High School in the dataset, replace the missing values in rows for people from Texas with High School.

One can also create a classification model. The column to predict here is Education, using other columns in the dataset. But the most common and popular approach is to model the missing value in a categorical column as a new category called Unknown:

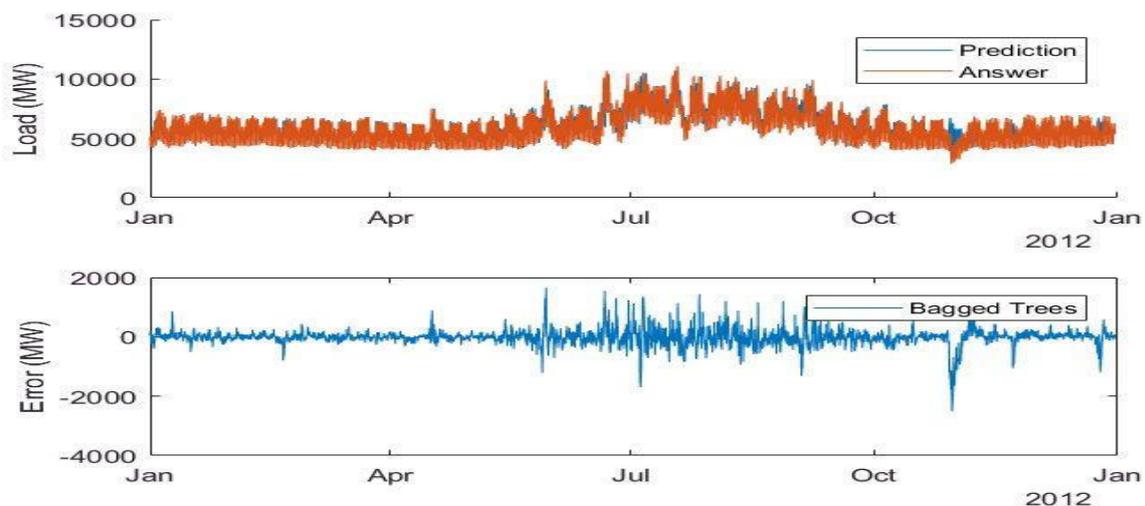| Row no | State | Education |
|---|---|---|
| 1 | NY | High School |
| 2 | TX | Unknown |
| 3 | NJ | High School |
| 4 | VT | High School |
| 5 | TX | Unknown |
| 6 | CA | College |
| 7 | NY | High School |
| 8 | CA | Unknown |
| 9 | CT | College |

## Missing values as a new category

In summary, you'll use different approaches to handle missing data values while data cleaning depending on the type of data and the problem at hand. If you have access to a domain expert, always incorporate their expert advice when filling in the missing values.

Most importantly, no matter the imputation method you choose, always run the predictive analytics model to see which one works best from the standpoint of data accuracy.

# => <u>Noisy, and Inconsistent Data:</u>

Noisy and inconsistent data are common challenges in real-world datasets, referring to errors, outliers, missing values, or format discrepancies that corrupt genuine patterns, leading to inaccurate analysis and poor model performance; they are handled through data pre-processing steps like removing outliers (Z-score, IQR), smoothing (binning, moving averages), imputation (filling missing data), and standardization to clean data for reliable insights.

## Noisy Data



**Definition:**
Noisy data contains **random errors, fluctuations, or outliers** that do not represent the true underlying pattern.

**Sources of Noise**
- Faulty sensors or measurement errors
- Data entry mistakes (typos, wrong units)
- Transmission errors
- Natural variability in real-world data

**Examples**
- A heart-rate sensor briefly recording **300 bpm**
- Customer age entered as **222**
- Sudden spikes/drops in time-series sales data due to logging glitches

**Impact on Predictive Models**
- Misleading patterns learned by the model
- Reduced accuracy and poor generalization
- Overfitting to random fluctuations

**Techniques to Handle Noisy Data**
1. **Outlier Detection & Treatment**
   - Z-score, IQR (Interquartile Range)
   - Isolation Forest, DBSCAN
2. **Smoothing Techniques**
   - Moving average, exponential smoothing

B Venkatesu Goud, Assistant Professor

3. **Filtering**
   - Noise filters for sensor/time-series data
4. **Robust Models**
   - Decision Trees, Random Forests (less sensitive to noise)
5. **Data Validation Rules**
   - Set realistic bounds (e.g., age $\in [0,120]$)

# Inconsistent Data



| Date | StartTime | Duration | Activity | ActivityType | LogType | Steps | Distance | E |
|------|-----------|----------|----------|--------------|---------|-------|----------|---|
| 10/01/17 | 12:21:10 | 5172000 | Walk | 90013 | auto_detect | 5768 | 0 | |
| 10/03/17 | 19:14:12 | 1792000 | Walk | 90013 | auto_detect | 2257 | 0 | |
| 10/03/17 | 23:19:57 | 1787000 | Walk | 90013 | auto_detect | 2334 | 0 | |
| 10/05/17 | 18:37:24 | 1075000 | Walk | 90013 | auto_detect | 1752 | 0 | |
| 10/09/17 | 10:55:51 | 7936000 | Walk | 90013 | auto_detect | 10638 | 0 | |
| 10/10/17 | 11:27:42 | 5446000 | Hike | 90012 | tracker | 8519 | 4.46707 | |
| 10/10/17 | 13:16:01 | 4309000 | Hike | 90012 | tracker | 5477 | 3.71722 | |
| 10/10/17 | 15:18:13 | 7886000 | Hike | 90012 | tracker | 11892 | 7.9415 | |
| 10/12/17 | 7:22:06 | 1178000 | Walk | 90013 | auto_detect | 1247 | 0 | |
| 10/12/17 | 20:43 | 973000 | Walk | 90013 | auto_detect | 1328 | 0 | |
| 10/13/2017 | 1:26:20 | 3533000 | Walk | 90013 | auto_detect | 4812 | 0 | |
| 10/13/2017 | 13:53:22 | 2766000 | Walk | 90013 | auto_detect | 4147 | 0 | |
| 10/14/2017 | 8:32:58 | 1587000 | Walk | 90013 | auto_detect | 2352 | 0 | |
| 10/15/2017 | 18:11:03 | 1740000 | Walk | 90013 | auto_detect | 2965 | 0 | |
| 10/17/2017 | 8:30:12 | 1076000 | Walk | 90013 | auto_detect | 1640 | 0 | |
| 10/21/2017 | 18:01:08 | 2166000 | Run | 90009 | tracker | 5760 | 6.465203 | |
| 10/21/2017 | 19:59:14 | 1178000 | Walk | 90013 | auto_detect | 2015 | 0 | |

1. Inconsistency in Date Format and Distance value          2. Missing data

**Definition:**
Inconsistent data contains **conflicting, contradictory, or non-uniform values** for the same entity or attribute.

**Causes of Inconsistency**
- Multiple data sources with different formats
- Lack of standardization
- Manual data entry variations
- Duplicate records

**Examples**
- Date formats: `01-02-2025` vs `2025/02/01`
- Gender values: `M`, `Male`, `male`, `1`
- Same customer with two different addresses
- Product price differs across databases

**Impact on Predictive Models**
- Confused feature interpretation
- Incorrect labeling and biased predictions
- Increased preprocessing complexity

**Techniques to Handle Inconsistent Data**
1. **Standardization**
   - Uniform formats for dates, units, categories
2. **Data Integration & Schema Matching**
   - Align fields from multiple sources
3. **Deduplication**
   - Remove or merge duplicate records
4. **Constraint Checking**
   - Enforce rules (e.g., DOB < current date)
5. **Master Data Management (MDM)**
   - Maintain a single "source of truth"

# UNIT 13    FEATURE SELECTION AND EXTRACTION

## 13.1  INTRODUCTION

Data sets are made up of numerous data columns, which are also referred to as data attributes. These data columns can be interpreted as dimensions on an n-dimensional feature space, and data rows can be interpreted as points inside that space. One can gain a better understanding of a dataset by applying geometry in this manner. In point of fact, these characteristics are measurements of the same entity. It is possible for their existence in the algorithm's logic to get muddled, which will result in a change to how well the model functions.

Input variables are the columns of data that are fed into a model in order to provide a forecast for a target variable. However, if your data is given in the form of rows and columns, such as in a spreadsheet, then features is another term that can be used interchangeably with input variables. It is possible that the presence of a large number of dimensions in the feature space implies that the volume of that space is enormous. As a result, the points (data rows) in that space reflect a small and non-representative sample of the space's contents. It is possible for the performance of machine learning algorithms to degrade when there are an excessive number of input variables. The existence of an excessive number of input variables has a significant impact on the efficiency with which machine learning algorithms function. when it is used to data that contains a large number of input attributes; this phenomenon is referred to as the "curse of dimensionality." As a consequence of this, one of the most common goals is to cut down on the number of input features. The process of decreasing the number of dimensions that characterise a feature space is referred to as "dimensionality reduction," which is a phrase that was made up specifically to describe this phenomenon.

The usefulness of data mining can be hindered by an excessive amount of information on occasion. There are occasions when only a handful of the columns of data characteristics that have been compiled for the purpose of constructing and testing a model do not offer any information that is significant

to the model. However, there are some that actually reduce the reliability and precision of the model.
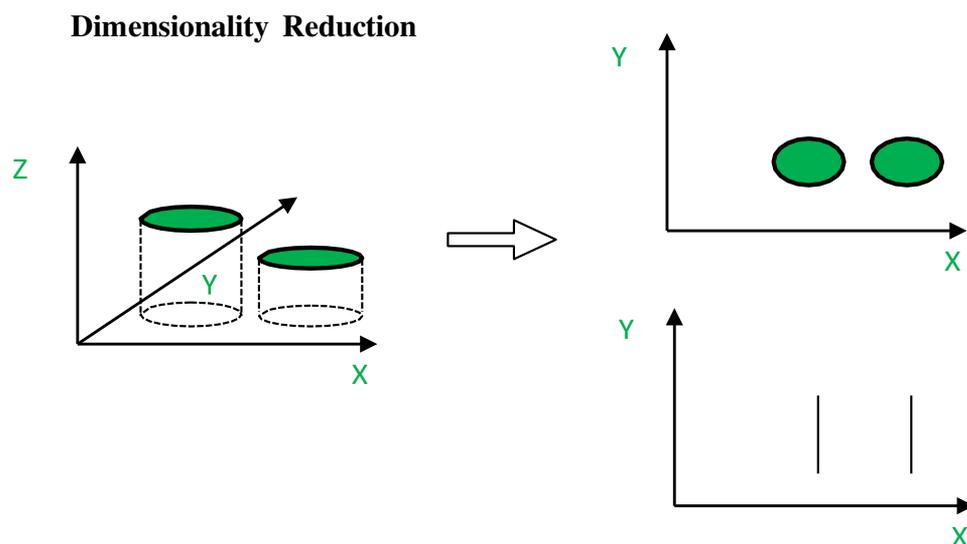
For instance, let's say you want to build a model that can forecast the incomes of people already employed in their respective fields. Therefore, data columns like cellphone number, house number, and so on will not truly contribute any value to the dataset, and they can therefore be omitted. This is because irrelevant qualities introduce noise to the data and affect the accuracy of the model. Additionally, because of the Noise, the size of the model as well as the amount of time and system resources required for model construction and scoring are both increased.

At this point in time, we are required to put the concept of Dimension Reductionality into practise. This can be done in one of two ways: either by selecting features to be extracted or by extracting features to be selected. Both of these approaches are broken down in greater detail below. The step of dimension reduction is one of the preprocessing phases that occurs during the process of data mining. This step is one of the preprocessing steps that may be beneficial in minimising the impacts of noise, correlation, and excessive dimensionality.

Some more examples are presented below to let you understand What does dimensionality reduction have to do with machine learning and predictive modelling?

- A simple issue concerning the classification of e-mails, in which we are tasked with deciding whether or not a certain email constitutes spam. can be brought up as a practical illustration of the concept of dimensionality reduction. This can include elements like whether or not the email has a standard subject line, the content of the email, whether or not it uses a template, and so on. However, some of these features may overlap with one another.

- A classification problem that involves humidity and rainfall can sometimes be simplified down to just one underlying feature as a result of the strong relationship that exists between the two variables. As a direct consequence of this, the number of characteristics could get cut down in some circumstances.

- A classification problem with three dimensions can be difficult to understand, whereas a problem with two dimensions can be translated to a fundamental space with two dimensions, and a problem with one dimension can be mapped to a line with one dimension. This concept is depicted in the diagram that follows, which shows how a three-dimensional feature space can be broken down into two one-dimensional feature spaces, with the number of features being reduced even further if it is discovered that they are related.

In context of dimensionality reduction, various techniques like Principal Component Analysis, Linear Discriminant Analysis, Singular Value Decomposition are frequently used. In this unit we will discuss all the mentioned concepts, related to Dimension reductionality

## 13.2  DIMENSIONALITY REDUCTION

The Data mining and Machine Learning methodologies both have processing challenges when working with big amounts of data (many attributes). In point of fact, the dimensions of the feature space utilised by the approach, often referred to as the model attributes, play the most important function. Processing algorithms grow more difficult and time-consuming to implement as the dimensionality of the processing space increases.

These elements, also known as the model attributes, are the fundamental qualities, and they can either be variables or features. When there are more features, it is more difficult to see them all, and as a result, the work on the training set becomes more complex as well. This complexity was further increased when a significant number of characteristics were linked; hence, the classification became irrelevant as a result. In circumstances like these, the strategies for decreasing the number of dimensions can prove to be highly beneficial. In a nutshell, "the process of making a set of major variables from a huge number of random variables is what is referred to as dimension reduction." When conducting data mining, the step of dimension reduction can be helpful as a preprocessing step to lessen the negative effects of noise, correlation, and excessive dimensionality.

Dimension reduction can be accomplished in two ways :

- Feature selection: During this approach, a subset of the complete set of variables is selected; as a result, the number of conditions that can be utilised to illustrate the issue is narrowed down. It's normally done in one of three ways.:

  - Filter method

  - Wrapper method

  - Embedded method

- Feature extraction: It takes data from a space with many dimensions and transforms it into another environment with fewer dimensions.

## 13.2.1　Feature Selection

It is the process of selecting some attributes from a given collection of prospective features, and then discarding the rest of the attributes that were considered. The use of feature selection can be done for one of two reasons: either to get a limited number of characteristics in order to prevent overfitting or to avoid having features that are redundant or irrelevant. For data scientists, the ability to pick features is a vital asset. It is essential to the success of the machine learning algorithm that you have a solid understanding of how to choose the most relevant features to analyse. Features that are irrelevant, redundant, or noisy can contaminate an algorithm, which can have a detrimental impact on the learning performance, accuracy, and computing cost. The importance of feature selection is only going to increase as the size and complexity of the typical dataset continues to balloon at an exponential rate.
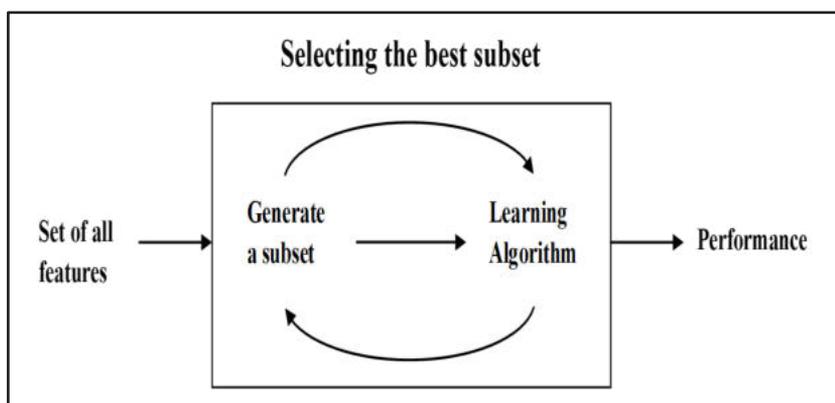
**Feature Selection Methods:** Feature selection methods can be divided into two categories: supervised, which are appropriate for use with labelled data, and unsupervised, which are appropriate for use with unlabeled data. Filter methods, wrapper methods, embedding methods, and hybrid methods are the four categories that unsupervised approaches fall under.:

- **Filter methods:** Filter methods choose features based on statistics instead of how well they perform in feature selection cross-validation. Using a chosen metric, irrelevant attributes are found and recursive feature selection is done. Filter methods can be either univariate, in which an ordered ranking list of features is made to help choose the final subset of features, or multivariate, in which the relevance of all the features as a whole is evaluated to find features that are redundant or not important.

- **Wrapper methods:** Wrapper feature selection methods look at the choice of a set of features as a search problem. Their quality is judged by preparing, evaluating, and comparing a set of features to other sets of features. This method makes it easier to find possible interactions between variables. Wrapper methods focus on subsets of features that will help improve the quality of the results from the clustering algorithm used for the selection. Popular examples are Boruta feature selection and Forward feature selection.

- **Embedded methods:** Embedded feature selection approaches incorporate the feature selection machine learning algorithm as an integral component of the learning process. This allows for simultaneous classification and feature selection to take place within the method. Careful consideration is given to the extraction of the characteristics that will make the greatest contribution to each iteration of the process of training the model. A few examples of common embedded approaches are the LASSO feature selection algorithm, the random forest feature selection algorithm, and the decision tree feature selection algorithm.

Among all approaches the most conventional feature selection is feed forward feature selection.

**Forward feature selection:** The first step in the process of feature selection is to evaluate each individual feature and choose the one that results in the most

effective algorithm model. This is referred to as "forward feature selection." After that step, each possible combination of the feature that was selected and a subsequent feature is analysed, and then a second feature is selected, and so on, until the required specified number of features is chosen. The operation of the forward feature selection algorithm is depicted here in the figure.



The procedure to follow in order to carry out forward feature selection

1. Train the model with each feature being treated as a separate entity, and then evaluate its overall performance.

2. Select the variable that results in the highest level of performance.

3. Carry on with the process while gradually introducing each variable.

4. The variable that produced the greatest amount of improvement is the one that gets kept.

5. Perform the entire process once more until the performance of the model does not show any meaningful signs of improvement.

Here, a fitness level prediction based on the three independent variables is used to show how forward feature selection works.

| ID | Calories_burnt | Gender | Plays_Sport? | Fintess Level |
|---|---|---|---|---|
| 1 | 121 | M | Yes | Fit |
| 2 | 230 | M | No | Fit |
| 3 | 342 | F | No | Unfit |
| 4 | 70 | M | Yes | Fit |
| 5 | 278 | F | Yes | Unfit |
| 6 | 146 | M | Yes | Fit |
| 7 | 168 | F | No | Unfit |
| 8 | 231 | F | Yes | Fit |
| 9 | 150 | M | No | Fit |
| 10 | 190 | F | No | Fit |

So, the first step in Forward Feature Selection is to train n models and judge how well they work by looking at each feature on its own. So, if you have three independent variables, we'll train three models, one for each of these three traits. Let's say we trained the model using the Calories Burned feature and the Fitness Level goal variable and got an accuracy of 87 percent.

| ID | Calories_burnt | Gender | Plays_Sport? | Fintess Level |
|----|----------------|--------|--------------|---------------|
| 1  | 121 | M | Yes | Fit |
| 2  | 230 | M | No  | Fit |
| 3  | 342 | F | No  | Unfit |
| 4  | 70  | M | Yes | Fit |
| 5  | 278 | F | Yes | Unfit |
| 6  | 146 | M | Yes | Fit |
| 7  | 168 | F | No  | Unfit |
| 8  | 231 | F | Yes | Fit |
| 9  | 150 | M | No  | Fit |
| 10 | 190 | F | No  | Fit |
|    |     |   |     |     |

Accuracy = 87%

We'll next use the Gender feature to train the model, and we acquire an accuracy of 80%. –

| ID | Calories_burnt | Gender | Plays_Sport? | Fintess Level |
|----|----------------|--------|--------------|---------------|
| 1  | 121 | M | Yes | Fit |
| 2  | 230 | M | No  | Fit |
| 3  | 342 | F | No  | Unfit |
| 4  | 70  | M | Yes | Fit |
| 5  | 278 | F | Yes | Unfit |
| 6  | 146 | M | Yes | Fit |
| 7  | 168 | F | No  | Unfit |
| 8  | 231 | F | Yes | Fit |
| 9  | 150 | M | No  | Fit |
| 10 | 190 | F | No  | Fit |

Accuracy = 80%

And similarly, the **Plays_sport** variable gives us an accuracy of 85%–

| ID | Calories_burnt | Gender | Plays_Sport? | Fintess Level |
|----|----------------|--------|--------------|---------------|
| 1  | 121 | M | Yes | Fit |
| 2  | 230 | M | No  | Fit |
| 3  | 342 | F | No  | Unfit |
| 4  | 70  | M | Yes | Fit |
| 5  | 278 | F | Yes | Unfit |
| 6  | 146 | M | Yes | Fit |
| 7  | 168 | F | No  | Unfit |
| 8  | 231 | F | Yes | Fit |
| 9  | 150 | M | No  | Fit |
| 10 | 190 | F | No  | Fit |

Accuracy = 85%

At this point, we are going to select the variable that produced the most favourable results. If you take a look at this table, you'll notice that the variable titled "Calories Burned" alone has an accuracy rating of 87 percent, while the variable titled "Gender" has an accuracy rating of 80 percent, and the variable titled "Plays Sport" has an accuracy rating of 85 percent. When these two sets of data were compared, the winner was, unsurprisingly, the number of calories burned. As a direct result of this, we will select this variable.

| Variable used | Accuracy |
|---|---|
| Calories_burnt | 87.00% |
| Gender | 80.00% |
| Plays_Sport? | 85.00% |

The next thing we'll do is repeat the previous steps, but this time we'll just add a single variable at a time. Because of this, it makes perfect sense for us to retain the Calories Burned variable as we proceed to add variables one at a time. Consequently, if we use gender as an illustration, we have an accuracy rate of 88 percent. –

| ID | Calories_burnt | Gender | Plays_Sport? | Fintess Level |
|---|---|---|---|---|
| 1 | 121 | M | Yes | Fit |
| 2 | 230 | M | No | Fit |
| 3 | 342 | F | No | Unfit |
| 4 | 70 | M | Yes | Fit |
| 5 | 278 | F | Yes | Unfit |
| 6 | 146 | M | Yes | Fit |
| 7 | 168 | F | No | Unfit |
| 8 | 231 | F | Yes | Fit |
| 9 | 150 | M | No | Fit |
| 10 | 190 | F | No | Fit |
|  |  |  |  |  |

Accuracy = 88%

We acquire a 91 percent accuracy when we combine Plays Sport with Calories Burnt. The variable that yields the greatest improvement will be kept. That makes natural sense. As you can see, when we combine Plays Sport with Calories Burnt, we get a better result. As a result, we'll keep it and use it in our model. We'll keep repeating the process till all the features are considered in improving the model performance

| ID | Calories_burnt | Gender | Plays_Sport? | Fintess Level |
|---|---|---|---|---|
| 1 | 121 | M | Yes | Fit |
| 2 | 230 | M | No | Fit |
| 3 | 342 | F | No | Unfit |
| 4 | 70 | M | Yes | Fit |
| 5 | 278 | F | Yes | Unfit |
| 6 | 146 | M | Yes | Fit |
| 7 | 168 | F | No | Unfit |

| 8 | 231 | F | Yes | Fit |
|----|-----|---|-----|-----|
| 9 | 150 | M | No | Fit |
| 10 | 190 | F | No | Fit |
| | | | | |
| | | | | |

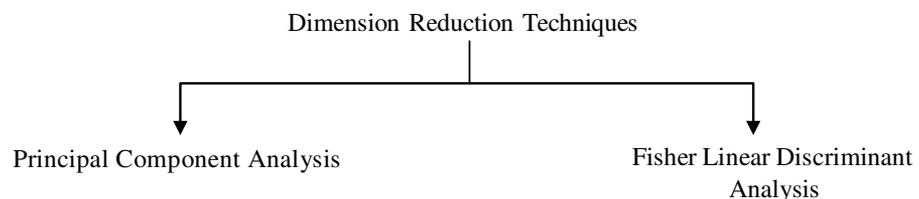Accuracy = 91%

### 13.2.2   Feature extraction:

The process of reducing the amount of resources needed to describe a large amount of data is called "feature extraction." One of the main problems with doing complicated data analysis is that there are a lot of variables to keep track of. A large number of variables requires a lot of memory and processing power, and it can also cause a classification algorithm to overfit to training examples and fail to generalise to new samples. Feature extraction is a broad term for different ways to combine variables to get around these problems while still giving a true picture of the data. Many people who work with machine learning think that extracting features in the best way possible is the key to making good models. The data's information must be shown by the features in a way that fits the needs of the algorithm that will be used to solve the problem. Some "inherent" features can be taken straight from the raw data, but most of the time, we need to use these "inherent" features to find "relevant" features that we can use to solve the problem.

In simple terms *"feature extraction."* can be described as a technique for *Defining a set of features, or visual qualities, that best show the information.* Feature Extraction Techniques such as: PCA, ICA, LDA, LLE, t-SNE and AE. are some of the common examples in machine learning.

**Feature extraction fills the following requirements:**

It takes raw data, called features, and turns them into useful information by reformatting, combining, and changing the primary features into new ones. This process continues until a new set of data is created that the Machine Learning models can use to reach their goals.

**Methods of Dimensionality Reduction :** The following are two well-known and widely used dimension reduction techniques:

Dimension Reduction Techniques

Principal Component Analysis                    Fisher Linear Discriminant Analysis

- (PCA) Principal Component Analysis

- (LDA) Fisher Linear Discriminant Analysis

The reduction of dimensionality can be linear or non-linear, depending on the method used. The most common linear method is called Principal Component Analysis, or PCA.

**Check Your Progress - 1**

Qn1. Define the term feature selection.

Qn2. What is the purpose of feature extraction in machine learning?
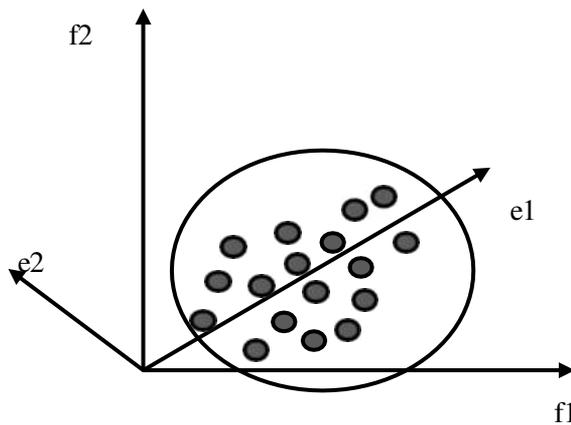
Qn3. Expand the following terms : PCA,LDA,GDA

Qn4. Name components of dimensionality reduction.

## 13.3  PRINCIPAL COMPONENT ANALYSIS

Karl Pearson was the first person to come up with this plan. It is based on the idea that when data from a higher-dimensional space is put into a lower-dimensional space, the lower-dimensional space should have the most variation. In simple terms, principal component analysis (PCA) is a way to get important variables (in the form of components) from a large set of variables in a data set. It tends to find the direction in which the data is most spread out. PCA is more useful when you have data with three or more dimensions.



When applying the PCA method, the following are the primary steps that should be followed:

1.  Obtain the dataset you need.

2.  Calculate the mean of the vectors ().

3.  Deduct the mean of the given data from the total.

4.  Complete the computation for the covariance matrix.

5.  Determine the eigenvectors and eigenvalues of the matrix that represents the covariance matrix.

6.  Creating a feature vector and deciding which components would be the major ones i.e. the principal components.

7.  Create a new data set by projecting the weight vector onto the dataset.As a result, we have a smaller number of eigenvectors, and some data may have been lost in the process. However, the remaining eigenvectors should keep the most significant variances.

**Merits of Dimensionality Reduction**

- It helps to compress data, which reduces the amount of space needed to store it and the amount of time it takes to process it.

- If there are any redundant features, it also helps to get rid of them.

**Limitations of Dimensionality Reduction**

- You might lose some data.

- You might lose some data.

- PCA fails when the mean and covariance are not enough to describe a dataset.

- We don't know how many major parts we need to keep track of, but in practice, we follow some rules.

Below is the practice question for Principal Component Analysis (PCA) :

**Problem-01:** 2, 3, 4, 5, 6, 7; 1, 5, 3, 6, 7, 8 are the given data. Using the PCA Algorithm, calculate the primary component.

**OR**

Consider the two-dimensional patterns (2, 1), (3, 5), (4, 3), (5, 6), (6, 7), (8, 8) and (9, 10). (7, 8).

Using the PCA Algorithm, calculate the primary component.

**OR**

Calculate the principal component of following data-

| Class1 | values | | Class 2 | values |
|--------|--------|---|---------|--------|
| X | 2,3,4 | | X | 5 , 6 , 7 |
| Y | 1,5,3 | | Y | 6 , 7 , 8 |

**Answer :**

**Step-1:** Get data.

The given feature vectors are- x1,x2,x3,x4,x5,x6 with the values given below:

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 3 \\ 5 \end{bmatrix} \begin{bmatrix} 4 \\ 3 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix} \begin{bmatrix} 6 \\ 7 \end{bmatrix} \begin{bmatrix} 7 \\ 8 \end{bmatrix}$$

**Step-2:**

Find the mean vector ($\mu$).

Mean vector ($\mu$) = ((2 + 3 + 4 + 5 + 6 + 7) / 6, (1 + 5 + 3 + 6 + 7 + 8) / 6)= (4.5, 5)

Thus, Mean vector ($\mu$) = $\begin{bmatrix} 4.5 \\ 5 \end{bmatrix}$

## Step-03:

On subtracting mean vector ($\mu$) from the given feature vectors.

* x1 – $\mu$ = (2 – 4.5, 1 – 5) = (-2.5, -4)

same for others

Feature vectors (xi) generated after subtraction are

$$\begin{bmatrix} -2.5 \\ -4 \end{bmatrix} \begin{bmatrix} -1.5 \\ 0 \end{bmatrix} \begin{bmatrix} -0.5 \\ -2 \end{bmatrix} \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \begin{bmatrix} 1.5 \\ 2 \end{bmatrix} \begin{bmatrix} 2.5 \\ 3 \end{bmatrix}$$

## Step-04:

Now to find covariance matrix : Covariance Matrix = $\dfrac{\sum (X_i - \mu)(X_i - \mu)^t}{n}$

$m_1 = (x_1 - \mu)(x_1 - \mu)^t = \begin{bmatrix} -2.5 \\ -4 \end{bmatrix} \begin{bmatrix} -2.5 & -4 \end{bmatrix} = \begin{bmatrix} 6.25 & 10 \\ 10 & 16 \end{bmatrix}$

$m_2 = (x_2 - \mu)(x_2 - \mu)^t = \begin{bmatrix} -1.5 \\ 0 \end{bmatrix} \begin{bmatrix} -1.5 & 0 \end{bmatrix} = \begin{bmatrix} 2.25 & 0 \\ 0 & 0 \end{bmatrix}$

$m_3 = (x_3 - \mu)(x_3 - \mu)^t = \begin{bmatrix} -0.5 \\ -2 \end{bmatrix} \begin{bmatrix} -0.5 & -2 \end{bmatrix} = \begin{bmatrix} 0.25 & 1 \\ 0 & 4 \end{bmatrix}$

$m_4 = (x_4 - \mu)(x_4 - \mu)^t = \begin{bmatrix} 0.5 \\ 2 \end{bmatrix} \begin{bmatrix} 0.5 & 1 \end{bmatrix} = \begin{bmatrix} 0.25 & 1 \\ 1 & 4 \end{bmatrix}$

$m_5 = (x_5 - \mu)(x_5 - \mu)^t = \begin{bmatrix} 1.5 \\ 2 \end{bmatrix} \begin{bmatrix} 1.5 & 2 \end{bmatrix} = \begin{bmatrix} 2.25 & 3 \\ 3 & 4 \end{bmatrix}$

$m_6 = (x_6 - \mu)(x_5 - \mu)^t = \begin{bmatrix} 1.5 \\ 2 \end{bmatrix} \begin{bmatrix} 2.5 & 3 \end{bmatrix} = \begin{bmatrix} 6.25 & 7.5 \\ 7.5 & 9 \end{bmatrix}$

Covariance Marix = $\dfrac{1}{6} \begin{bmatrix} 17.5 & 22 \\ 22 & 34 \end{bmatrix}$

Covariance Matrix = $\begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix}$

## Step-05:

Eigen values and Eigen vectors of the covariance matrix.

$\begin{vmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{vmatrix} - \begin{vmatrix} \lambda & 0 \\ 0 & \lambda \end{vmatrix} = 0$

$\begin{vmatrix} 2.92 - \lambda & 3.67 \\ 3.67 & 5.67 - \lambda \end{vmatrix} = 0$

From here,

$(2.92 - \lambda)(5.67 - \lambda) - (3.67 \times 3.67) = 0$

$16.56 - 2.92\lambda - 5.67\lambda + \lambda^2 - 13.47 = 0$

$\lambda^2 - 8.56\lambda + 3.09 = 0$

Solving this quadratic equation, we get $\lambda = 8.22, 0.38$

Thus, two eigen values are $\lambda\_1=8.22$ and $\lambda\_2=0.38$.

Clearly, the second eigen value is vary small compared to the first eigen value.

So, the second eigen vactor can be left out.

Eigen vector corresponding to the greatest eigen value is the principle component for the given data set.

So, we find the eigen vector corresponding to eigen value $\lambda\_1$.

We use the following equation to find the eigen vector-

$$MX=\lambda X$$

Where-

- M=Covariance Matrix ; X=Eigen vector ,and $\lambda$=Eigen value

Substituting the values in the above equation, we get-

On being substituting the values in the above equation, we get-

$$\begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix} \begin{bmatrix} X1 \\ X2 \end{bmatrix} = 8.22 \begin{bmatrix} X1 \\ X2 \end{bmatrix}$$

Solving these, we get-

$2.92X\_1+3.67X\_2=8.22X\_1$

$3.67X\_1+5.67X\_2=8.22X\_2$

On simplification, we get-

$5.3X_1=3.67X_2$.........(1)

$3.67X_1=2.55X_2$.........(2)

From (1) and (2), X1 =0.69X2

From (2), the eigen vector is-

Eigen Vector: $\begin{bmatrix} X1 \\ X2 \end{bmatrix} = \begin{bmatrix} 2.55 \\ 3.67 \end{bmatrix}$

Thus, PCA for the given problem is

Principle Component: $\begin{bmatrix} X1 \\ X2 \end{bmatrix} = \begin{bmatrix} 2.55 \\ 3.67 \end{bmatrix}$

Lastly, we project the data points onto the new subspace as-

**Problem -02**

Use PCA Algorithm to transform the pattern (2, 1) onto the eigen vector in the previous question.

**Solution-**

The given feature vector is (2, 1) i.e. Given Feature Vector: $\begin{bmatrix} 2 \\ 1 \end{bmatrix}$

The feature vector gets transformed to :

= Transpose of Eigen vector x (Feature Vector – Mean Vector)

$$= \begin{bmatrix} 2.55 \\ 3.67 \end{bmatrix}^T x \left( \begin{bmatrix} 2 \\ 1 \end{bmatrix} - \begin{bmatrix} 4.5 \\ 5 \end{bmatrix} \right) = \begin{bmatrix} 2.55 & 3.67 \end{bmatrix} x \begin{bmatrix} -2.5 \\ -4 \end{bmatrix} = -21.055$$

**Check Your Progress -3**

Qn1. What are the advantages of dimensionality reduction?

Qn2. What are the disadvantages of dimensionality reduction?

# 13.4  LINEAR DISCRIMINANT ANALYSIS

In most cases, the application of logistic regression has been restricted to problems involving two classes of subjects. On the other hand, the Linear Discriminant Analysis is the linear classification method that is recommended to use when there are more than two classes.

The algorithm for linear classification known as logistic regression is known for being both straightforward and robust. On the other hand, there are a few restrictions or faults in the system that highlight the requirement for more complex linear classification algorithms. The following is a list of some of the problems:

- **Binary class Problems.** Concerns regarding the binary class is that the Logistic regression is utilised for issues that involve binary classification or two classes. It is possible to enhance it such that it can manage multiple-class categorization, but in practise, this is not very common.

- **Unstable, but with well-defined classes.** When the classes are extremely distinct from one another, logistic regression may become unstable.

- **It is prone to instability when there are only a few occurrences.** When there are not enough examples from which to draw conclusions about the parameters, the logistic regression model may become unstable.

In view of the limitations of logistic regression that were discussed earlier, the linear discriminant analysis is one of the prospective linear methods that can be used for multi-class classification. This is because it addresses each of the aforementioned concerns in their totality, which is the primary reason for its success (i.e. flaws of logistic regression). When dealing with issues that include binary categorization, two statistical methods that could be effective are logical regression and linear discriminant analysis. Both of these techniques are linear and regression-based.

**Understanding LDA Models :** In order to simplify the analysis of your data and make it more accessible, LDA will make the following assumptions about it:

1. The distribution of your data is Gaussian, and when plotted, each variable appears to be a bell curve.

2. Each feature has the same variance, which indicates that the values of each feature vary by the same amount on average in relation to the mean.

On the basis of these presumptions, the LDA model generates estimates for both the mean and the variance of each class. In the case where there is only one input variable, which is known as the univariate scenario, it is straightforward to think about this.

When the sum of values is divided by the total number of values, we are able to compute the mean value, or mu, of each input, or x, for each class(k), in the following manner.

$$\mu_k = 1/n_k * sum(x)$$

Where,

$\mu_k$ represents the average value of x for class k and

$n_k$ represents the total number of occurrences that belong to class k.

When calculating the variance across all classes, the average squared difference of each individual result's distance from the mean $x - \mu$ is employed.

$$\sigma^2 = 1 / (n\text{-}K) * sum((x - \mu)^2)$$

Where $\sigma^2$ represents the variance of all inputs (x), n represents the number of instances, K represents the number of classes, and μ is the mean for input x.

**Now we will discuss How to use LDA to Make Predictions ?**

LDA generates predictions by calculating the likelihood that each class will be given a fresh batch of data and then extrapolating from there. A forecast is created by selecting the output class that contains the events that are the most likely to occur. The Bayes Theorem is incorporated into the model in order to calculate the probabilities involved. Utilizing the likelihood of each class as well as the probability of data belonging to that class, Bayes's Theorem may be utilised to estimate the probability of the output class (k) given the input class (x). This is accomplished by using the following formula:

$$P(Y=x|X=x) = (\Pi_k k * f_k(x)) / sum(\Pi_l * f_l(x))$$

The base probability of each class (k) that can be found in your training data is denoted by the symbol $\Pi_k$ (e.g. 0.5 for a 50-50 split in a two class problem). This concept is referred to as the prior probability within Bayes' Theorem.

$$\Pi_k = n_k/n$$

The value of f, which represents the estimated likelihood that x is a member of the class, is presented here as f(x). We make use of a Gaussian distribution function for the variable (x). By simplifying the previous equation and then

introducing the Gaussian, we are able to arrive at the equation that is presented below. This type of function is referred to as a discriminate function, and the output classification (y) is determined by selecting the category that contains the greatest value:

$$D_k(x) = x * (\mu_k/\sigma^2) - (\mu_k^2/(2*\sigma^2)) + \ln(\Pi_k)$$

Where, $D_k(x)$ is the discriminating function for class k given input x, and $\mu_k$, $\sigma^2$, and $\Pi_k$ are all estimated from your data.
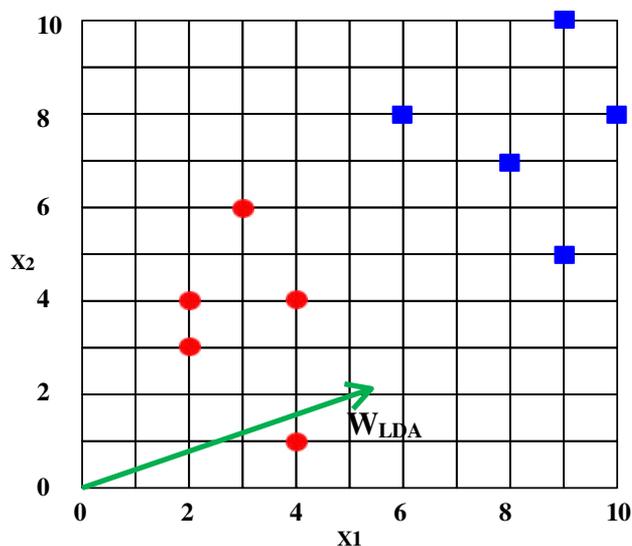
**Now to perform the above task we need to prepare our data first, so the question arises, How to prepare data suitable for LDA?**

In order to prepare data suitable for LDA, one needs to understand following:-

1) **Problems with Classification:** LDA is used to solve classification problems where the output variable is a categorical one. This may seem obvious, as LDA works with both two and more than two classes.

2) **Gaussian Distribution:** The standard way to use the model assumes that the input variables have a Gaussian distribution. Think about looking at the univariate distributions of each attribute and using transformations to make them look more like Gaussian distributions (e.g. log and root for exponential distributions and Box-Cox for skewed distributions).

3) **Remove Outliers:** Think about removing outliers from your data. These things can mess up the basic statistics like the mean and the standard deviation that LDA uses to divide classes.

4) **Same Variance:** LDA assumes that the variance of each input variable is the same. Before using LDA, you should almost always normalise your data so that it has a mean of 0 and a standard deviation of 1.

Below is a practice problems based on Linear Discriminant Analysis (LDA) -

**Problem-2 :** Compute the Linear Discriminant projection for the following two-dimensional datasetX1=(x1,x2)={(4,1),(2,4),(2,3),(3,6),(4,4)} & X2=(x1,x2)={(9,10),(6,8),(9,5),(8,7),(10,8)}

**Solution :**

To understand the working of LDA lets take an example and workout step by step.

Step 1:- Compute the within – Class Scatter matrix ($S_w$), Sw measures how well the data is scattered with in the class , which is done by finding the mean of the classes i.e. $\mu_1$ & $\mu_2$ where, $\mu_{1\&}$ $\mu_2$ are the mean of class $C_1$ & $C_2$ respectively

Now, $S_w$ is given by
$$S_w = S_1 + S_2$$

where,  $S_1$ is the covariance matrix for the class $C_1$ and
$\quad\quad\quad S_2$ is the covariance matrix for the class $C_2$

Let's find the covariance matrices $S_1$ & $S_2$ of each class

$$S1 = \sum_{x\epsilon c_1} (x - \mu_1)\,(x - \mu_1)^T$$

where, $\mu_1$ is the mean of class $C_1$, which is computed by averaging the coordinates of dataset $X_1$

(Coordinates of $X_1$)
$$\mu_1 = \left\{ \frac{4+2+2+3+4}{5},\ \frac{1+4+3+6+4}{5} \right\};$$

$$\mu_1 = [3.00\ ;\ 3.60]$$

Similarly, $\mu_2 = [8.4\ ;\ 7.60]$

$$S_1 = \sum_{x\epsilon w_1} (x - \mu_1)\,(x - \mu_1)^T \quad\quad \mu_1 = [3 \quad\quad 3 \cdot 60]$$

Mean reduced data of class $C_1$ is given by $(x_1 - \mu_1)$

$$(x_1 - \mu_1) = \begin{bmatrix} 1 & -1 & -1 & 0 & 1 \\ -2.6 & 0.4 & -0.6 & 2.4 & 0.4 \end{bmatrix}$$

Now, for each x1 we are going to calculate $(x - \mu_1)\,(x - \mu_1)^T$ . So, we will have "5" such matrices.

Now solving it one by one i.e. for each column of $(x_1 - \mu_1)$ we find $(x_1 - \mu_1)$.
$(x_1 - \mu_1)^T$
$$\begin{bmatrix} 1 \\ -2.6 \end{bmatrix} [1 \quad -2.6] = (x_1 - \mu_1) = \begin{bmatrix} 1 & -2.6 \\ -2.6 & 6.16 \end{bmatrix} \rightarrow ①$$

FIRST MATRIX
Similarly for all columns $\begin{bmatrix} 1 \\ 0.4 \end{bmatrix} [-1 \quad 0.4] = \begin{bmatrix} 1 & -0.4 \\ 0.4 & 0.16 \end{bmatrix} \rightarrow ②$

$$\begin{bmatrix} -1 \\ 0.6 \end{bmatrix} [-1 \quad 0.6] = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 0.36 \end{bmatrix} \rightarrow ③$$

$$\begin{bmatrix} 0 \\ 2.4 \end{bmatrix} [0 \quad 2.4] = \begin{bmatrix} 1 & 0 \\ 0 & 5.76 \end{bmatrix} \rightarrow ④$$

$$\begin{bmatrix} 1 \\ 0.4 \end{bmatrix} [1 \quad 0.4] = \begin{bmatrix} 1 & 0 \\ 0.4 & 0.16 \end{bmatrix} \rightarrow ⑤$$

Adding ①+②+③+④+⑤ and taking average we get covariance matrix $S_1$

$$S_1 = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 2.6 \end{bmatrix}$$

Similarly for class 2, the covariance matrix is given by

$$S_2 = \begin{bmatrix} 1.84 & -0.04 \\ -0.4 & 2.64 \end{bmatrix} \ \& \ \mu_2 = \begin{bmatrix} 8.4 & 7.6 \end{bmatrix}$$

$$S_w = S_1 + S_2$$

$$S_w = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

Step 2. Compute Between class scatter matrix (SB)

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$= \begin{pmatrix} -5.4 \\ -4 \end{pmatrix} (-5.4 - 4 = \begin{pmatrix} 29.16 & 21.6 \\ 21.6 & 16.00 \end{pmatrix}$$

Step 3: Find the best LDA Projection Vector.

The LDA projection vector is a vector on to which all data samples are projected and it carries all necessary features required for data, in fact this projection vector is the Eigen vector, and we know that Eigen vector is the vector that carries good balance between all the features of the dataset. The expression used for it is $S_w^{-1} S_B V = \lambda V$ where V is the projection vector.

Similar to PCA we find this using Eigen vectors having largest Eigen value.

$$S_w^{-1} S_B V = \lambda V \qquad \longrightarrow \text{Projection Vector}$$

i.e. $|S_w^{-1} S_B - \lambda I| = 0$

$$\begin{vmatrix} 11.89 - \lambda & 8.81 \\ 5.08 & 3.76 \end{vmatrix} = 0$$

Substituting in equation a $\lambda$

$$\begin{bmatrix} 11.89 & 8.81 \\ 5.08 & 3.76 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = 15.65 \begin{bmatrix} V_1 \\ V_2 \end{bmatrix}$$

We get $\begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} 0.91 \\ 0.39 \end{bmatrix}$

Or we directly solve

$$\begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = S_w^{-1}(\mu^1 - \mu^2) = \begin{bmatrix} 0.1921 & -0.032 \\ -0.031 & 0.38 \end{bmatrix}\begin{bmatrix} -5.4 \\ -4 \end{bmatrix} = \begin{bmatrix} 0.91 & 0.39 \end{bmatrix}^T$$

Note: $S_w^{-1}$ is found using the formula.

$$\begin{bmatrix} a & b \\ c & c \end{bmatrix}^{-1} = \frac{1}{ad-bc} = \begin{bmatrix} a & -b \\ -c & c \end{bmatrix}$$
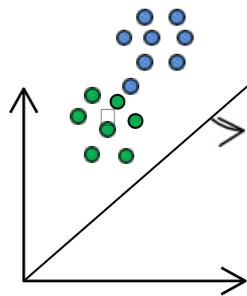
So, $S_w = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$

$$S_w^{-1} = \frac{1}{13.74}\begin{bmatrix} 5.28 & 0.44 \\ 0.44 & 2.68 \end{bmatrix} = \begin{bmatrix} 0.384 & 0.032 \\ 0.032 & 0.192 \end{bmatrix}$$

Step 4 : Dimension reduction

$y = w^T X \longrightarrow$ input data samples i.e. $y = \begin{bmatrix} 0.91 & 0.39 \end{bmatrix}\begin{bmatrix} 4 & 2 & 2 & 3 & 4 \\ 1 & 2 & 3 & 6 & 4 \end{bmatrix}$

Project on vector $\qquad y = \begin{bmatrix} 0.91 & 0.39 \end{bmatrix}\begin{bmatrix} 9 & 0 & 9 & 8 & 10 \\ 10 & 8 & 5 & 7 & 8 \end{bmatrix}$



Projection Vector Corresponding to highest Eigen

**Note :** Linear Discriminant Analysis (LDA) is a technique for dimensionality reduction, and it is used as a pre-processing step for pattern classification and machine learning applications. LDA is similar to PCA but the basic difference is that LDA in addition find the axis that measures the separation between multiple classes. Here the Goal is to project the N-dimensional feature space on to a smaller k-dimensional subspace where k<=n, while maintaining the class discriminator information.

**Check Your Progress -3**

Q1. Define LDA.

.................................................................................................................................

.................................................................................................................................

.................................................................................................................................

Q2. Write any two limitations of LDA.

.................................................................................................................................

.................................................................................................................................

.................................................................................................................................

# 13.5  SINGLE VALUE DECOMPOSITION

The Singular Value Decomposition (SVD) method is a well-known technique for decomposing a matrix into a large number of component matrices. This method is valuable since it reveals many of the interesting and helpful characteristics of the initial matrix. We can use SVD to discover the optimal lower-rank approximation to the matrix, determine the rank of the matrix, or test a linear system's sensitivity to numerical error.

Singular value decomposition is a method of decomposing a matrix into three smaller matrices.

$$A = U\sum V^T$$

Where:

- A : is an m × n matrix

- U : is an m × n orthogonal matrix

- S : is an n × n diagonal matrix

- V : is an n × n orthogonal matrix

Below is a practice problems based on single value decomposition

**Problem-03.**  Find the SVD of the matrix $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ -1 & 1 \end{bmatrix}$

**Solution :** Let $U \sum V^T$ be the Singular Value Decomposition (SVD) of the Matrix A, so need to compute V, $\sum$ and U, the steps are as follows

**Step 1 :** To find the SVD of A , we need to **determine matrix V** then we will find $V^T$

In order to find matrix V firstly we need to **Find $A^T A$**

$$A^T A = \begin{bmatrix} 1 & 0 & -1 \\ 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$$

Then we are required to **find the Eigen values** ($\lambda$) of $A^T A$

The procedure of finding the Eigen values ($\lambda$) of $A^T A$ is as follows:

Let $\lambda^2 + S_1 \lambda + S_2 = 0$ be the characteristic equation of $A^T A$,

Where, $S_1$ = Trace of $A^T A$ = $Tr(A^T A) = (2+3) = 5$

*Note : Trace implies the sum of diagonal elements*

And $S_2$ = Determinant of $A^T A = \begin{vmatrix} 2 & 0 \\ 0 & 3 \end{vmatrix} = 6$

Hence by substituting the values of S1 & S2 in the characteristic equation of $A^T A$ i.e. $\lambda^2 + S_1\lambda + S_2 = 0$ we get $\lambda^2 + 5\lambda + 6 = 0$; solving the characteristic equation

$\lambda^2 + 5\lambda + 6 = 0$ we get the Eigen values $\lambda_1 = 3$ & $\lambda_2 = 2$

Note : for SVD we need to write Eigen values in decreasing order i.e. $\lambda = 3, 2$

Now **Calculation for Eigen vectors** is required for the determined Eigen Values $\lambda = 3, 2$

For $\lambda_1 = 3$ let $X_1 = \begin{bmatrix} a \\ b \end{bmatrix}$ be the Eigen Vector of $A^T A$; i.e. $(A^T A) X_1 = \lambda_1 X_1$ ;

Substituting the values we get

$$\begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = 3 \begin{bmatrix} a \\ b \end{bmatrix}$$

Solving we get

$$-a + 0b = 0 \text{---------------------- (1)}$$

$$0a + 0b = 0 \text{---------------------- (2)}$$

From (1); $-a + 0b = 0$ ; substituting a=0 & b=1 satisfies the equation so the Eigen Vector

$$X_1 = \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Now we need to Normalize this Eigen Vector $X_1$, and for this we need to divide each element of Eigen Vector $X_1$ by the length of Eigen Vector $X_1$.

length of Eigen Vector $X_1 = \sqrt{a^2 + b^2} = \sqrt{0^2 + 1^2} = 1$

Therefore, $N(X_1)$ i.e. Normalized Eigen Vector $X_1 = \begin{bmatrix} 0/1 \\ 1/1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

Similarly, $N(X_2)$ i.e. Normalized Eigen Vector $X_2 = \begin{bmatrix} 1/1 \\ 0/1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Hence matrix $V = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$, and since V is a symmetric matrix, $V^T = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$

**Step-2 : Compute $\Sigma$**

To compute $\Sigma$ We need to consider that the order of $\Sigma$ must be equal to the order of A, and we need to calculate the singular values $\sigma1$ and $\sigma2$ ; where $\sigma1 = \sqrt{\lambda_1}$ & $\sigma2 = \sqrt{\lambda_2}$ i.e. $\sigma1 = \sqrt{3}$ & . $\sigma2 = \sqrt{2}$ . Also, we need to recall the relation between rank of Matrix and Eigen Values that number of non zero Eigen values should be equal to the rank of the matrix. Since number of non zero Eigen values is Two, thus the Rank (R) of the matrix is also two i.e. 2

$$\Sigma = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{bmatrix};$$

Where, $\sum_1$. is a diagonal matrix whose diagonal elements are $\sigma_1$ & $\sigma_2$ ; and Order of $\sum_1$. = R X R where R is the rank of the matrix thus Order of $\sum_1$.= 2 X 2 (since R = 2).

Therefore, $\sum = \begin{bmatrix} \sqrt{3} & 0 \\ 0 & \sqrt{2} \\ 0 & 0 \end{bmatrix}$ ; we are adding zeros in the 3rd row because the order

of $\sum$ must be equal to the order of A

**Step-3 : Compute** U

Calculate $A.A^T = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 2 \end{bmatrix}$

Now, find Eigen Values and Eigen Vectors of $A.A^T = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 2 \end{bmatrix}$, as we did it earlier for the case of matrix V

**find the Eigen values** ($\lambda$) of $AT^A$

The procedure of finding the Eigen values ($\lambda$) of $A^T A$ is as follows

Let. $\lambda^3 - S_1 \lambda^2 + S_2 \lambda - S_3 = 0$ be the characteristic equation of $A^T A$,

where

$S_1$ = Trace of $AA^T$ = $Tr(AA^T)$ = (2+1+2) = 5

Note : Trace implies the sum of diagonal elements

$S_2$ = Sum of Minors of Diagonal Matrix $AA^T$

$S_2 = \begin{vmatrix} 1 & 1 \\ 1 & 2 \end{vmatrix} + \begin{vmatrix} 2 & 0 \\ 0 & 2 \end{vmatrix} + \begin{vmatrix} 2 & 1 \\ 1 & 1 \end{vmatrix} = 1+4+1 = 6$

**And**

$S_3$ = Determinant of $AA^T = \begin{vmatrix} 2 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 2 \end{vmatrix} = 2(1) - 1(2) = 0$

Hence by substituting the values of S1, S2 & S3 in the characteristic equation of $AA^T$

i.e. $\lambda^3 - S_1 \lambda^2 + S_2 \lambda - S_3 = 0$  we get $\lambda^3 - 5 \lambda^2 + 6\lambda - 0 = 0$ as the characteristic equation of $AA^T$

Now, by solving the characteristic equation $\lambda^3 - 5 \lambda^2 + 6\lambda - 0 = 0$ ; we get the

Eigen values $\lambda_1 = 2$ , $\lambda_2 = 0$ & $\lambda_3 = 3$

We know for SVD we need to arrange Eigen Values in decreasing order thus $\lambda = 3, 2, 0$

**Now find Eigen vectors for corresponding Eigen values**

**For** $\lambda = 3$ let $X_1 = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$ be the Eigen vector of $AA^T$; then $(AA^T) X_1 = \lambda X_1$

Thus, $(AA^T) X_1 = \lambda X_1 \Rightarrow \begin{bmatrix} 2 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 2 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix} = 3 \begin{bmatrix} a \\ b \\ c \end{bmatrix}$; solving the expression we get

$$-a + 0b + 0c = 0 \text{-----------------------(1')}$$
$$a - 2b + c = 0 \text{-----------------------(2')}$$
$$0a + b - c = 0 \text{-----------------------(3')}$$

Solving equation (1'),(2') & (3') we get $X_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$

Now we need to Normalize this Eigen Vector $X_1$, and for this we need to divide each element of Eigen Vector $X_1$ by the length of Eigen Vector $X_1$.

length of Eigen Vector $X_1. = \sqrt{a^2 + b^2 + c^2} = \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$

Therefore, N(X1) i.e. Normalized Eigen Vector $X_1 = \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix}$

Similarly $X_2 = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}$ Therefore, N(X2) i.e. Normalized Eigen Vector $X_2 = \begin{bmatrix} 1/\sqrt{2} \\ 0 \\ -1/\sqrt{2} \end{bmatrix}$

Similarly $X_3 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$ Therefore, N(X3) i.e. Normalized Eigen Vector $X_3 = \begin{bmatrix} 1/\sqrt{6} \\ -2/\sqrt{6} \\ 1/\sqrt{6} \end{bmatrix}$

Now, Matrix U can be written by using $N(X_1)$, $N(X_2)$ & $N(X_3)$ i.e.

$$U = \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{2} & 1/\sqrt{6} \\ 1/\sqrt{3} & 0 & -2/\sqrt{6} \\ 1/\sqrt{3} & -1/\sqrt{2} & 1/\sqrt{6} \end{bmatrix}$$

Since, U $\sum$ V$^T$ is the Singular Value Decomposition (SVD) of the Matrix A, so need to substitute values of V$^T$, $\sum$ and U to find U $\sum$ V$^T$

SVD of the matrix A i.e.$= \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ -1 & 1 \end{bmatrix} = U\Sigma V^T = \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{2} & 1/\sqrt{6} \\ 1/\sqrt{3} & 0 & -2/\sqrt{6} \\ 1/\sqrt{3} & -1/\sqrt{2} & 1/\sqrt{6} \end{bmatrix} \begin{bmatrix} \sqrt{3} & 0 \\ 0 & \sqrt{2} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$

**Check Your Progress - 4**

Q1. Define SVD

.........................................................................................................

.........................................................................................................

.........................................................................................................

## 13.6  SUMMARY

In this unit we learned about the concept of Dimension Reductionality, wherein we understood basics of Feature Selection and Feature extraction techniques. Thereafter an explicit discussion of Principal Component Analysis(PCA), Linear Discriminant Analysis(LDA) and Singular Value Decomposition(SVD) was given.

## 13.7  SOLUTIONS TO CHECK YOUR PROGRESS

**Check Your Progress – 1**

Refer Section 13.2 for detailed Solutions

**Ans1.** In machine learning and statistics, feature selection is the process of choosing a subset of relevant features to use when making a model. It is also called variable selection, attribute selection, or variable subset selection.

**Ans2.** The goal of Feature Extraction is to reduce the number of features in a dataset by making new features from the ones that are already there (and then discarding the original features).

Refer Section 13.2 for detailed Solutions

**Ans3.** Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA) and Generalized Discriminant Analysis (GDA)

**Ans4.** Feature selection and feature extraction.

**Check Your Progress  - 2**

Refer Section 13.3 for detailed Solutions

**Ans1.** Advantages of dimensionality reduction

- It reduces the time and space complexity.

- Getting rid of multi-colinearity makes it easier to understand how the machine learning model's parameters work.

- When there are only two or three dimensions, it's easier to see the data.

**Ans2.** Dimensionality Reduction's Drawbacks

- PCA likes to find linear relationships between variables, which isn't always a good thing.

- PCA doesn't work when the mean and covariance aren't enough to describe a dataset.

- In some problems, dimension reduction could also cause data loss.

### Check Your Progress - 3

Refer Section 13.4 for detailed Solutions

**Ans1.** The LDA technique is a multi-class classification method that may be utilised to automatically perform dimensionality reduction. LDA cuts the number of features down from the initial number of features.

LDA projects the data into a new linear feature space, and obviously, the classifier will have a high level of accuracy if the data can be linearly separated

**Ans2.** Some of the limitations of Logistic Regression are as follows:

- Logistic regression is typically applied when attempting to solve problems involving binary or two-class classifications. Problems involving two classes Even while it is possible to extrapolate this information and apply it for multi-class classification, very few people actually do this. On the other hand, Linear Discriminant Analysis is regarded as a superior option whenever multi-class classification is necessary, and in the case of binary classifications, both logistic regression and LDA are utilised in the analysis process.

- Instability with Clearly Delineated Social Groups — Logistic Regression is known to be unreliable in situations when the classes are clearly differentiated from one another.

### Check Your Progress - 4

Refer Section 13.5 for detailed Solutions

**Ans1.** The singular value decomposition is a factorization technique that can be used in linear algebra for real or complex matrices. It extends the application of the eigen decomposition of a square normal matrix to any m x n matrix by using an ortho normal eigen basis. It has something to do with the polar decomposition.
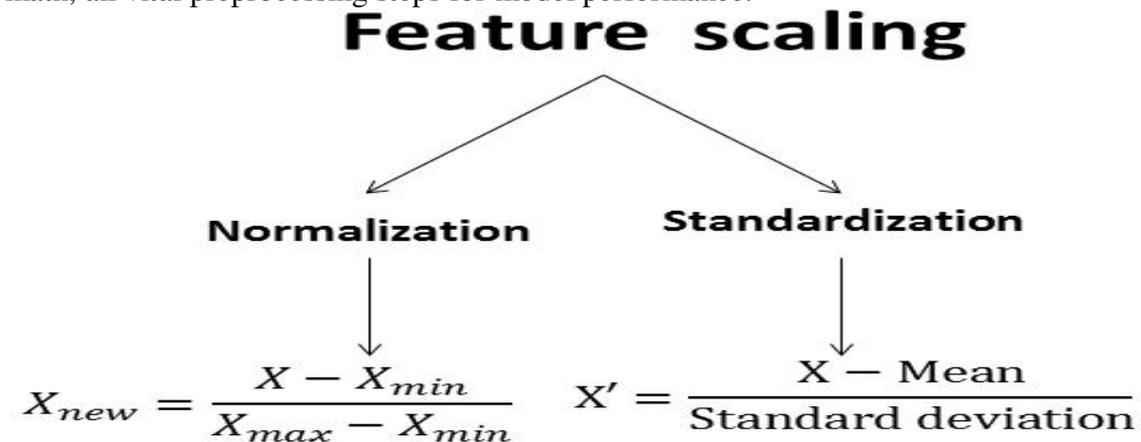
## 13.8 <u>FURTHER READINGS</u>

1. Machine learning an algorithm perspective, Stephen Marshland, 2nd Edition, CRC Press, 2015.

2. Machine Learning, Tom Mitchell, 1st Edition, McGraw- Hill, 1997.

3. Machine Learning: The Art and Science of Algorithms that Make Sense of Data, Peter Flach, 1st Edition, Cambridge University Press, 2012.

# => Feature Scaling: Normalization, Standardization:

Feature scaling (Normalization, Standardization) brings numerical features to a common scale, preventing large-range features from dominating models, while categorical encoding converts text categories into numbers; Normalization scales to [0] (good for non-Gaussian, distance-based), Standardization centers data to mean 0, std dev 1 (good for Gaussian data, robust to outliers), and encoding like One-Hot (nominal) or Label (ordinal) prepares text data for math, all vital preprocessing steps for model performance.

## Feature scaling

Normalization                     Standardization

$$X_{new} = \frac{X - X_{min}}{X_{max} - X_{min}} \qquad X' = \frac{X - Mean}{Standard\ deviation}$$

## Feature Scaling (Numerical Data)

**Why it's needed**: Algorithms like KNN, SVM, Neural Networks, and Linear Regression (using Gradient Descent) are sensitive to feature scales; features with larger ranges can unfairly influence the model, leading to poor performance.

1. **Normalization (Min-Max Scaling)**
- **What it does**: Rescales features to a fixed range, usually [0].
- **Formula**: $(X - X\_min) / (X\_max - X\_min)$ .

- **Best for**: Algorithms that don't assume data distribution (like Neural Networks, KNN) and when outliers are not a major concern, as it's sensitive to them.

2. **Standardization (Z-Score Scaling)**
- **What it does**: Transforms data to have a mean (μ) of 0 and a standard deviation (σ) of 1.
- **Formula**: $(X - \mu) / \sigma$ .

- **Best for**: Algorithms that assume normally distributed data (Gaussian) and when data has outliers, as it's less affected by them than normalization.

3. **Robust Scaling**
- **What it does**: Uses median and Interquartile Range (IQR) to scale, making it robust to outliers.

## Comparison of Various Feature Scaling Techniques
Let's see the key differences across the five main feature scaling techniques commonly used in machine learning pre-processing.

B Venkatesu Goud, Assistant Professor

| Type | Method Description | Sensitivity to Outliers | Typical Use Cases |
|---|---|---|---|
| Absolute Maximum Scaling | Divides values by max absolute value in each feature | High | Sparse data, simple scaling |
| Min-Max Scaling (Normalization) | Scales features to by min-max normalization | High | Neural networks, bounded input features |
| Normalization (Vector Norm) | Scales each sample vector to unit length (norm = 1) | Not applicable (per row) | Direction-based similarity, text classification |
| Standardization (Z-Score) | Centers features to mean 0 and scales to unit variance | Moderate | Most ML algorithms, assumes approx. normal data |
| Robust Scaling | Centers on median and scales using IQR | Low | Data with outliers, skewed distributions |

## => Encoding Categorical Variables:

### Encoding Categorical Variables

**Why it's needed**: Models work with numbers, so text categories must be converted.

1. **One-Hot Encoding (for Nominal Data)**
- **What it does**: Creates new binary columns for each category, marking presence (1) or absence (0).
- **Use when**: Categories have no order (e.g., colors: Red, Blue, Green).
2. **Label Encoding (for Ordinal Data)**
- **What it does**: Assigns a unique integer to each category (e.g., Small=1, Medium=2, Large=3).
- **Use when**: Categories have a meaningful order (e.g., education levels).
3. **Dummy Encoding**
- Similar to One-Hot but drops one column to avoid multicollinearity in regression models.

### Key Takeaway
- **Scale first**: Apply scaling to numerical features.
- **Encode next**: Convert categorical features.

- **Fit on Train, Transform All**: Fit your scaler (e.g., `StandardScaler` ) *only* on the training data, then use it to transform both training and test sets to prevent data leakage.

# => <span style="color:red">Feature Extraction:</span>

Feature extraction is the process of transforming raw data into a simplified and informative set of features or attributes. This reduces data complexity and highlights the most relevant information making it easier for machine learning models to analyze and learn from the data efficiently. It plays an important role in improving model accuracy and reducing computational costs by focusing on essential aspects of the data.

## Importance of Feature Extraction

Feature extraction is important for several reasons:

- **Reduced Computation Cost:** Raw data, especially from images or large datasets can be very complex. Feature extraction makes this data simpler hence reducing the computational resources needed for processing.
- **Improved Model Performance:** By focusing on key features, machine learning models can work with more relevant information leading to better performance and more accurate results.
- **Better Insights:** Reducing the number of features helps algorithms concentrate on the most important data, eliminating noise and irrelevant information which can lead to deeper insights.
- **Prevention of Overfitting:** Models with too many features may become too specific to the training data, making them perform poorly on new data. Feature extraction reduces this risk by simplifying the model.

## Key Techniques for Feature Extraction

There are various techniques for extracting meaningful features from different types of data:

## 1. Statistical Methods

Statistical methods are used in feature extraction to summarize and explain patterns of data. Common data attributes include:
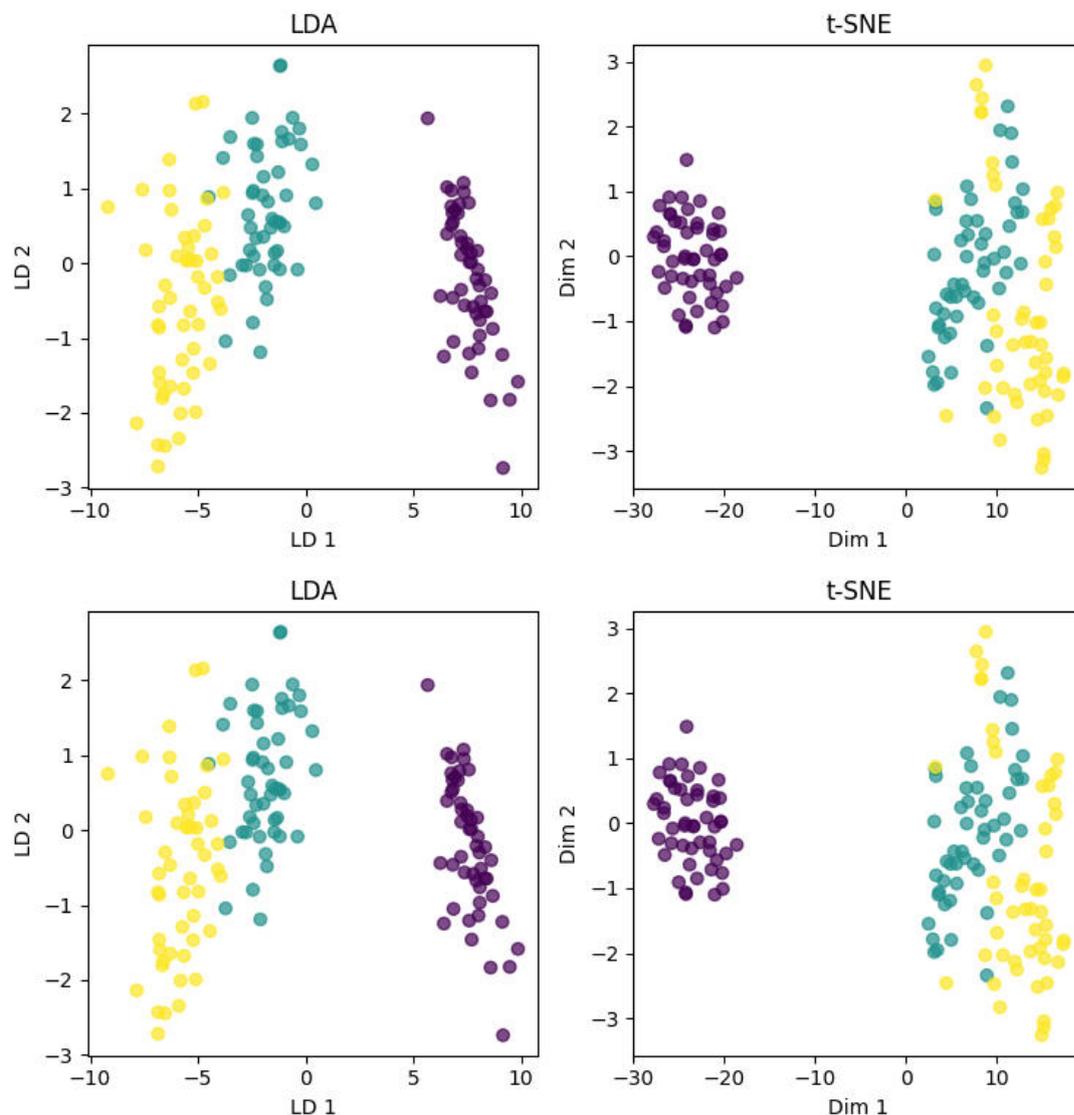


B Venkatesu Goud, Assistant Professor

**Statistical Methods**

- **Mean:** The average value of a dataset.
- **Median:** The middle value when it is sorted in ascending order.
- **Standard Deviation:** A measure of the spread or dispersion of a sample.
- **Correlation and Covariance:** Measures of the linear relationship between two or more factors.
- **Regression Analysis:** A way to model the link between a dependent variable and one or more independent factors.

These statistical methods can be used to represent the center trend, spread and links within a collection.

## 2. Dimensionality Reduction

Dimensionality reduction reduces the number of features without losing important information. Some popular methods are:



- **Principal Component Analysis**: It selects variables that account for most of the data's variation, simplifying the dataset by focusing on the most important components.
- **Linear Discriminant Analysis (LDA):** It finds the best combination of features to separate different classes, maximizing class separability for better classification.

- **t-Distributed Stochastic Neighbor Embedding (t-SNE)**: A technique that reduces high-dimensional data into two or three dimensions ideal for visualizing complex datasets.
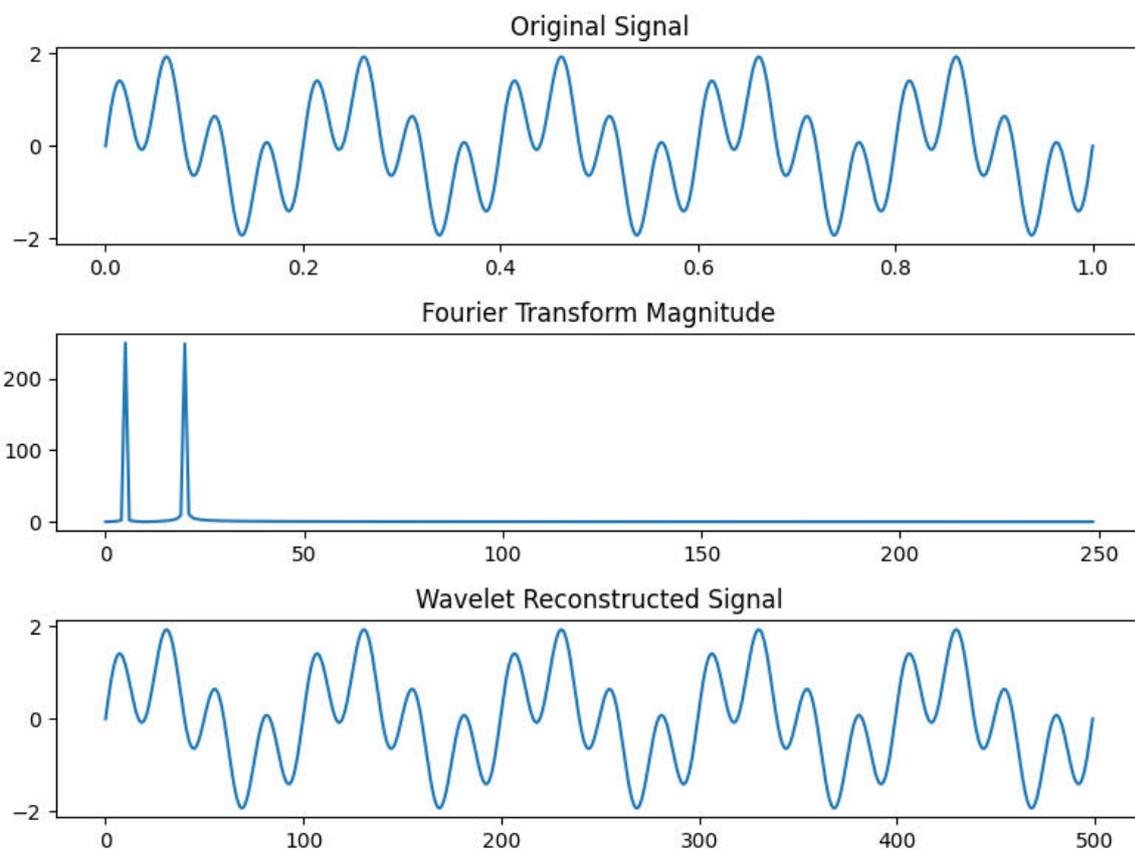
## 3. Feature Extraction for Textual Data

In Natural Language Processing (NLP), we often convert raw text into a format that machine learning models can understand. Some common techniques are:

1. **Bag of Words (BoW):** Represents a document by counting word frequencies, ignoring word order, useful for basic text classification.
2. **Term Frequency-Inverse Document Frequency (TF-IDF)**: Adjusts word importance based on frequency in a specific document compared to all documents, highlighting unique terms.

4. Signal Processing Methods

It is used for analyzing time-series, audio and sensor data:



Signal processing methods

1. **Fourier Transform:**It converts a signal from the time domain to the frequency domain to analyze its frequency components.
2. **Wavelet Transform:**It analyzes signals that vary over time, offering both time and frequency information for non-stationary signals.

## 5. Image Data Extraction

Techniques for extracting features from images:

Original Image     HOG Features     CNN Features (downsampled)

Image Data Extraction

1. **Histogram of Oriented Gradients (HOG):**This technique finds the distribution of intensity gradients or edge directions in an image. It's used in object detection and recognition tasks.
2. **Convolutional Neural Networks (CNN) Features:** They learn hierarchical features from images through layers of convolutions, ideal for classification and detection tasks.

**Choosing the Right Method**

Selecting the appropriate feature extraction method depends on the type of data and the specific problem we're solving. It requires careful consideration and often domain expertise.

- **Information Loss:** Feature extraction might simplify the data too much, potentially losing important information in the process.
- **Computational Complexity:** Some methods, especially for large datasets can be computationally expensive and may require significant resources.

**Feature Selection vs. Feature Extraction**

Since Feature Selection and Feature Extraction are related but not the same, let's quickly see the key differences between them for a better understanding:

| Aspect | Feature Selection | Feature Extraction |
|---|---|---|
| **Definition** | Selecting a subset of relevant features from the original set | Transforming the original features into a new set of features |
| **Purpose** | Reduce dimensionality | Transform data into a more manageable or informative representation |
| **Process** | Filtering, wrapper methods, embedded methods | Signal processing, statistical techniques, transformation algorithms |
| **Output** | Subset of selected features | New set of transformed features |

| Aspect | Feature Selection | Feature Extraction |
|---|---|---|
| Computational Cost | Lower cost | May be higher, especially for complex transformations |
| Interpretability | Retains interpretability of original features | May lose interpretability depending on transformation |

## Applications of Feature Extraction

Feature extraction plays an important role in various fields where data analysis is important. Some common applications include:

- **Computer Vision and Image Processing:** Used in autonomous vehicles to detect road signs and pedestrians by extracting key visual features for safe navigation.
- **Natural Language Processing (NLP):** Powers email spam filtering by extracting textual features to accurately classify messages as spam or legitimate.
- **Biomedical Engineering:** Extracting features from EEG or MRI signals helps diagnose neurological disorders or detect early signs of disease.
- **Industrial and Equipment Monitoring:** Predictive maintenance uses sensor data features to foresee machine failures, reducing downtime and repair costs.
- **Financial and Fraud Detection:** Analyzes transaction patterns to identify fraudulent activities and prevent financial losses.

## Tools and Libraries for Feature Extraction

There are several tools and libraries available for feature extraction across different domains. Let's see some popular ones:

- Scikit-learn: It offers tools for various machine learning tasks including PCA, ICA and preprocessing methods for feature extraction.
- OpenCV: A popular computer vision library with functions for image feature extraction such as SIFT, SURF and ORB.
- TensorFlow / Keras: These deep learning libraries in Python provide APIs for building and training neural networks which can be used for feature extraction from image, text and other types of data.
- PyTorch: A deep learning library enabling custom neural network designs for feature extraction and other tasks.
- NLTK (Natural Language Toolkit): A popular NLP library providing feature extraction methods like bag-of-words, TF-IDF and word embeddings for text data.

## Advantages

Feature extraction has various advantages which are as follows:

- **Simplifies Data:** Reduces complex data into a manageable form for easier analysis and visualization.
- **Boosts Model Performance:** Removes irrelevant data, making algorithms faster and more accurate.
- **Highlights Key Patterns:** Filters out noise to focus on important features for quicker insights.
- **Improves Generalization:** Helps models perform better on new, unseen data by emphasizing informative features.
- **Speeds Up Training and Prediction:** Fewer features mean faster model training and real-time predictions

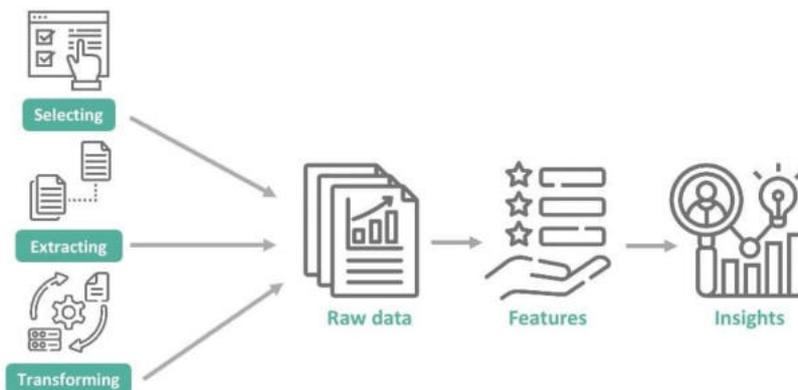B Venkatesu Goud, Assistant Professor

### Challenges
- **Managing High-Dimensional Data:** Extracting relevant features from large, complex datasets can be difficult.
- **Risk of Overfitting or Underfitting:** Too many or too few features can hurt model accuracy and generalization.
- **Computational Costs:** Complex methods may require heavy resources, limiting use with big or real-time data.
- **Redundant or Irrelevant Features:** Overlapping or noisy features can confuse models and reduce efficiency.

## => Feature Construction:

Feature construction in predictive analytics is the creative process of generating new, more informative features from existing data, often using domain knowledge, to help machine learning models better capture complex patterns and relationships for improved accuracy, going beyond simple transformations to extract deeper insights. It involves combining, splitting, or mathematically manipulating raw features (e.g., creating `Total_Spend` from `Quantity` * `Unit_Price`, or extracting `Hour_of_Day` from a timestamp) to reveal hidden predictive power.



What is Feature Engineering?

### Key Aspects & Techniques
- **Domain Expertise:** Leverages human understanding of the problem (e.g., business rules, scientific principles) to design relevant features.
- **Manual vs. Algorithmic:** Can be done manually by data scientists or assisted by algorithms.
- **Types of Construction:**
  - **Interaction Features:** Combining variables (e.g., `Age` * `Income`) to capture synergistic effects.
  - **Mathematical Transformations:** Applying functions like logs or polynomials to numerical data to handle skewness or non-linearity.
  - **Feature Splitting:** Breaking down a feature (e.g., `Name` into `Title` and `Last_Name`).
  - **Aggregation:** Summarizing data (e.g., total purchases per customer).
- **Goal:** To enhance model performance by making underlying patterns explicit, reducing dimensionality, and improving explainability (Explainable AI).

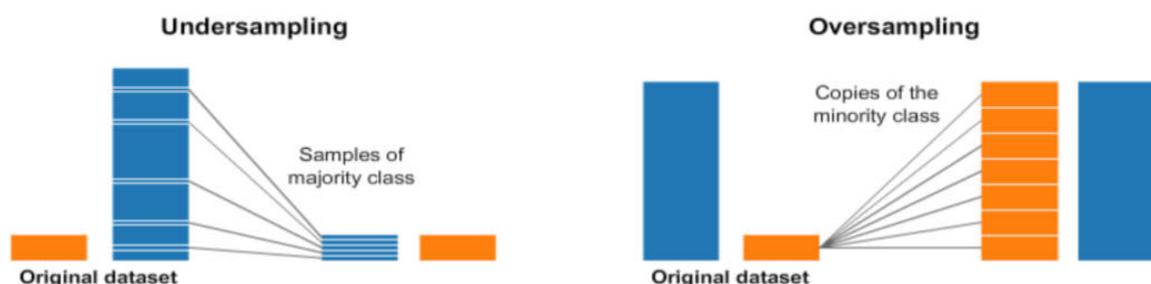B Venkatesu Goud, Assistant Professor

### Examples in Action

- **E-commerce:** Creating `Total_Spend` (Quantity * Unit_Price) or `Days_Since_Last_Purchase`.
- **Time Series:** Converting a full timestamp into `Day_of_Week`, `Month`, or `Is_Weekend`.
- **Health:** Calculating Body Mass Index (BMI) from `Height` and `Weight`.

### Relationship to Feature Engineering

Feature construction is a core part of feature engineering, the broader process that also includes feature selection (choosing features) and feature transformation (scaling, encoding). Construction specifically focuses on *creating* new data representations to unlock predictive insights.

## => <u>Dealing with Imbalanced Datasets:</u>

Dealing with imbalanced datasets in predictive analytics involves using resampling (oversampling minority/undersampling majority), class weighting, synthetic data generation (like SMOTE), or anomaly detection to prevent models from favoring the majority class, along with using metrics beyond accuracy (like Precision, Recall, F1-Score) for better evaluation. These strategies balance class representation, improve minority class prediction, and enhance overall model robustness for real-world applications like fraud detection.



### Data-Level Techniques (Resampling)

- **Oversampling:** Increases minority class samples by duplicating existing ones or creating synthetic data (e.g., **SMOTE**) to reduce overfitting risks compared to simple duplication.
- **Undersampling:** Reduces majority class samples, useful when data is abundant but risks information loss.
- **Hybrid Methods:** Combines oversampling and undersampling for balanced results, such as SMOTE with Tomek links.

### Algorithm-Level Techniques (Cost-Sensitive Learning)

- **Class Weighting:** Assigns higher penalties (costs) to misclassifying the minority class, forcing the model to pay more attention to it (e.g., `class_weight='balanced'` in scikit-learn models).

- **Ensemble Methods:** Uses specialized algorithms like BalancedBaggingClassifier or Balanced Random Forest, which integrate resampling within the bagging process to build stronger, balanced models.

## Evaluation & Approach Adjustments

- **Use Appropriate Metrics:** Rely on Precision, Recall, F1-Score, ROC-AUC, or PR-AUC instead of just accuracy, as accuracy can be misleading with imbalanced data.
- **Anomaly Detection:** Treat the minority class as an anomaly if it's extremely rare, focusing on outlier detection rather than standard classification.

## Key Considerations

- **Choose based on Use Case:** The best method depends on dataset size and specific problem; no single technique is universally superior.
- **Experiment:** Try different rebalancing ratios and combine techniques for optimal performance.

B Venkatesu Goud, Assistant Professor

# UNIT III

**Linear Regression and Polynomial Regression, Logistic Regression for Binary Classification, Decision Trees and Random Forest, k-Nearest Neighbors (k-NN) and Naïve Bayes, Support Vector Machines (SVM), Model Selection and Comparison.**

## => <span style="color:red">**Linear Regression and Polynomial Regression:**</span>

Regression analysis is a cornerstone technique in data science and machine learning, used to model the relationship between a dependent variable and one or more independent variables. Among the various types of regression, Linear Regression and Polynomial Regression are two fundamental approaches.



Linear Regression and Polynomial Regression

**What is Linear Regression?**
Linear regression is a statistical method used to model the relationship between a dependent variable and one or more independent variables by fitting a linear equation to observed data. The equation for simple linear regression is:
$y = a + bx$
- where **y** is the dependent variable,
- **x** is the independent variable and
- **a,b** are coefficients.

Linear regression is ideal when the relationship between variables is linear.

**What is Polynomial Regression?**
Polynomial regression is an extension of linear regression that models the relationship between the dependent variable and the independent variable(s) as an n-th degree polynomial. The equation for polynomial regression is:
$y = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \ldots + a_n x^n$
- where **y** is dependent variable
- **x** is the independent variable and

B Venkatesu Goud, Assistant Professor

- **a0,a1,a2,a3,an** are the coefficients.

Polynomial Regression is useful for modeling non-linear relationships where the data points form a curve.

**Key Differences Between Linear and Polynomial Regression**

| Model | Nature of Relationship | Model Complexity | Flexibility |
|---|---|---|---|
| **Linear Regression** | Assumes a straight-line relationship between dependent and independent variables. Suitable for linear trends. | Simpler, easier to interpret, fewer parameters, less prone to overfitting. | Limited to linear relationships, may underfit non-linear data. |
| **Polynomial Regression** | Can model non-linear relationships by fitting a polynomial equation to the data. Ideal for complex patterns. | More complex, higher-degree polynomial, more prone to overfitting. | Offers greater flexibility, can model curves and intricate patterns, suitable for non-linear trends. |

**Understanding Practical Examples for Linear and Polynomial Regression**

**Real-Life Linear Regression Examples**
1. **Real Estate Pricing Prediction**:
   - **Problem**: Predict the selling price of houses based on features like size, location, and number of bedrooms.
   - **Why Linear Regression**: The relationship between house features and price is often linear, making linear regression suitable for a first approximation.
2. **Sales Forecasting for a Retail Store**:
   - **Problem**: Estimate next month's sales based on historical sales data, taking into account factors like advertising budget, seasonality, and store location.
   - **Why Linear Regression**: It provides a straightforward model to understand how different factors linearly impact sales, aiding in budget planning and marketing strategies.

**Real-Life Polynomial Regression Examples**
1. **Agricultural Yield Prediction Based on Environmental Conditions**:
   - **Problem**: Predict the crop yield based on variables such as temperature, rainfall, and soil quality, where the relationship between these factors and yield is not linear.
   - **Why Polynomial Regression**: Environmental factors often have a non-linear impact on crop yields. Polynomial regression can model these complex relationships more effectively than linear regression.
2. **Modeling Electricity Consumption in Relation to Temperature**:
   - **Problem**: Forecast the electricity consumption of a city based on the temperature, where consumption increases during extreme cold and hot temperatures but drops at moderate temperatures.

- **Why Polynomial Regression**: The relationship between temperature and electricity consumption is likely to be non-linear (U-shaped curve), making polynomial regression a better fit for capturing these dynamics.

**When to Use Linear Regression vs. Polynomial Regression**

Choosing between linear and polynomial regression depends on the nature of your data and the relationship between the variables you are analyzing. Here are some scenarios to help you decide when to use each method:

**Linear Regression**
- When the relationship between variables is linear.
- When simplicity and interpretability are crucial.
- With smaller datasets to avoid overfitting.
- For initial analysis to understand basic trends.

**Scenario: Predicting house prices based on square footage and location:**

Why Use Linear Regression: The relationship between house prices and their size/location is often linear. As the size increases, the price generally increases proportionally. Linear regression provides a straightforward model that is easy to interpret and works well with this type of data.

**Polynomial Regression**
- When the relationship between variables is non-linear.
- To capture more complex relationships in large datasets.
- When flexibility is needed to fit a wider range of data shapes.
- With careful consideration of the polynomial degree to avoid overfitting.

**Scenario: Modeling the growth rate of bacteria over time:**

Why Use Polynomial Regression: The growth rate of bacteria often follows a non-linear pattern, such as an S-curve or exponential growth followed by a plateau. Polynomial regression can capture this complex relationship by fitting a curve to the data, which linear regression cannot do.

**Advantages and Disadvantages of Regression Models**

**Advantages and Disadvantages of Linear Regression**
**Advantages:**
- **Simplicity**: Easy to implement and interpret.
- **Efficiency**: Requires fewer computational resources.
- **Robustness**: Less prone to overfitting with large datasets.

**Disadvantages:**
- **Limited Flexibility**: Cannot model non-linear relationships.
- **Underfitting**: May not capture the complexity of the data if the true relationship is non-linear.

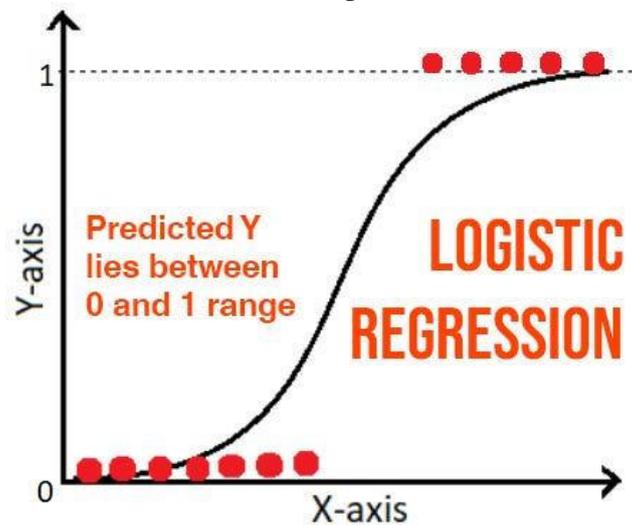**Advantages and Disadvantages of Polynomial Regression**
**Advantages:**
- **Flexibility**: Can model a wide range of relationships.
- **Better Fit**: Can capture non-linear trends in the data.

**Disadvantages:**
- **Complexity**: More complex and harder to interpret.
- **Overfitting**: Prone to overfitting, especially with higher-degree polynomials.
- **Sensitivity to Outliers**: More sensitive to outliers compared to linear regression.

## => **Logistic Regression for Binary Classification:**

Logistic regression is a statistical method employed to model the relationship between a binary dependent variable and one or more independent variables.



- Unlike linear regression, which predicts continuous outcomes, logistic regression estimates the probability that a given input point belongs to a particular category.

- This is achieved using the logistic function, also known as the sigmoid function, which maps any real-valued number to a value between 0 and 1.

**Logistic Regression Equation and Assumptions**

**Logistic Regression Equation**

- The logistic regression model is based on the following equation:

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n)}}$$

**where :**

- $P(Y=1|X)$ is the probability of the outcome being 1

- $\beta_0$ is the intercept

- $\beta_1, \beta_2, \ldots, \beta_n$ are the coefficients

- $X_1, X_2, \ldots, X_n$ are the input features

The result is passed through the sigmoid function:

B Venkatesu Goud, Assistant Professor

$$Sigma(z) = \frac{1}{1+e^{-z}}$$

- This converts the linear combination of features into a **probability between 0 and 1**.

**Assumptions of Logistic Regression**

- **Binary Output:** Logistic regression assumes the outcome is binary (0 or 1).

- **Independence of Observations:** Each observation should be independent of others.

- **No Multicollinearity:** Independent variables should not be highly correlated.

- **Linearity of Logit:** There should be a linear relationship between the log odds and the input variables.

- **Large Sample Size:** A bigger dataset improves the performance and reliability.

### Types of Logistic Regression with Examples

While binary classification is the most common, logistic regression comes in different forms depending on the number of categories in the output.

**1. Binary Logistic Regression**

- Used when the outcome has two classes.

- Example: Predicting if a customer will churn or not.

**2. Multinomial Logistic Regression**

- Used when the outcome has three or more unordered categories.

- Example: Predicting which product category a customer will choose (Electronics, Clothing, Groceries).

**3. Ordinal Logistic Regression**

- Used when the output variable has ordered categories.

- Example: Predicting customer satisfaction (Low, Medium, High).

| Type | Target Variable Type | Example Use Case |
|---|---|---|
| Binary | Two categories | Spam/Not Spam |
| Multinomial | Three or more (unordered) | Product preference prediction |
| Ordinal | Three or more (ordered) | Rating satisfaction level (Poor, Fair, Good) |

### Application of Logistic Regression in Binary Classification

Here are some real-world examples where logistic regression is widely used:

**1. Email Spam Detection**

Classifies whether an email is spam or not spam based on keywords, sender info, etc.

**2. Credit Scoring**

Predicts whether a person will default on a loan based on income, age, credit history, etc.

**3. Customer Churn Prediction**

Identifies whether a customer is likely to leave a service or not.

**4. Disease Diagnosis**

Predicts whether a patient has a disease (1) or not (0) based on symptoms, lab results, etc.

**5. Marketing Campaign Effectiveness**

Checks if a user will respond to a marketing campaign based on demographics and past behavior.

### Logistic Regression vs. Linear Regression

| Feature | Logistic Regression | Linear Regression |
|---|---|---|
| Output | Probability (0 to 1) | Continuous values |
| Use Case | Classification | Regression |
| Function Used | Sigmoid Function | Identity function |
| Best For | Binary/Multiclass problems | Predicting numerical outcomes |

### Tools and Libraries for Logistic Regression

You can implement logistic regression using various tools:

- Python: scikit-learn, statsmodels

- R: glm function

- Excel: Add-ins like XLSTAT

- SAS/SPSS: Built-in statistical tools

**Advantages of Logistic Regression**

- Simple and easy to implement

- Works well with linearly separable classes

- Fast training and prediction

- Interpretable model – easy to explain to stakeholders

**Limitations of Logistic Regression**

- Doesn't work well with non-linear data

- Sensitive to outliers

- Assumes linearity in the logit

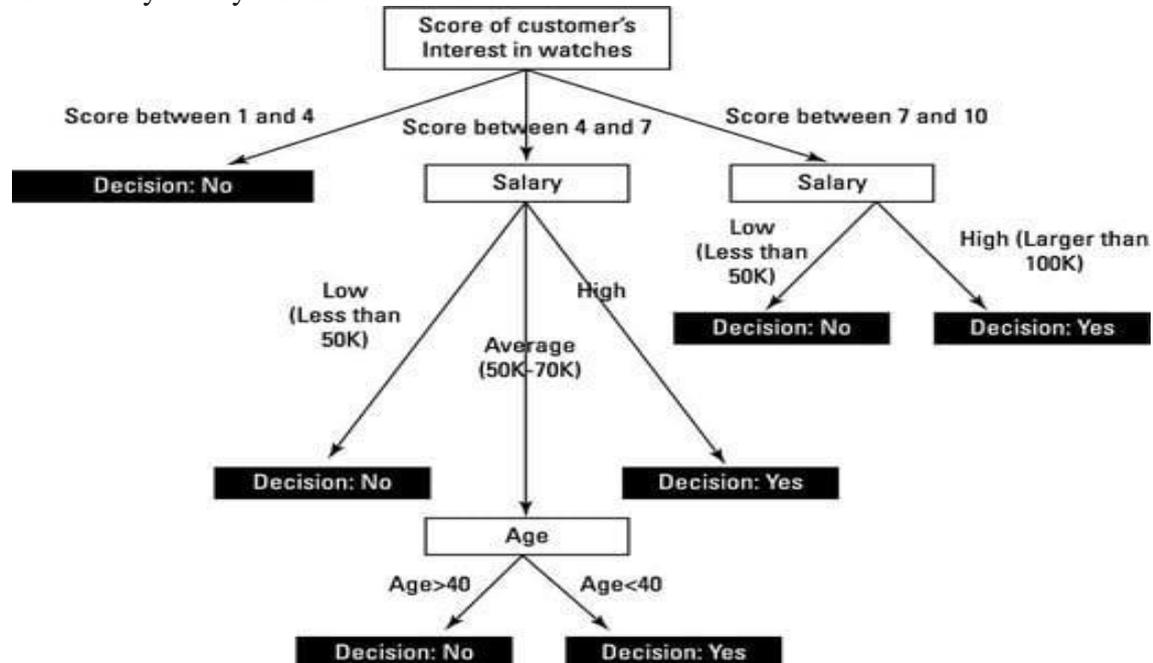- May underperform with high-dimensional or complex data

# => Decision Trees and Random Forest:

## What Is a Decision Tree?

A Decision Tree is a predictive model that maps observations about data points to conclusions about their target value. Visually, it resembled a tree structure where each node represents a feature (attribute), each branch indicates a decision rule, and each leaf node

represents a conclusion. It's like a flowchart that guides you to the right decision based on cumulative information.

You can think of a Decision Tree as a series of questions that lead to specific outcomes. For example, if you're trying to decide whether to play outside based on the weather, your first question might be, "Is it raining?" Depending on the answer, you follow a different branch that leads you to your final decision.



## How Do Decision Trees Work?

At its core, a Decision Tree splits the data into subsets based on the value of input features. Here's a simplified step-by-step breakdown of the process:

1. **Selecting the Best Feature to Split**: Algorithms like Gini impurity or information gain are used to decide which feature will give the best separation of classes in the data.
2. **Creating Branches**: Once the best feature is found, branches are created based on the values of that feature, leading to more questions.
3. **Recursive Splitting**: This process continues recursively, creating new nodes and branches until a stopping criterion is reached—such as a maximum depth or a minimum number of samples in a node.
4. **Making Predictions**: Finally, predictions are made by following the branches of the tree based on new input data and reaching a leaf node that contains the predicted outcome.

## Advantages of Decision Trees

1. **Easy to Interpret**: One of the greatest strengths of Decision Trees lies in their interpretability. You can visualize how decisions are made, making it easier for stakeholders to understand.
2. **Requires Little Data Preparation**: Decision Trees can handle both numerical and categorical data and do not require scaling.
3. **Exploits Non-linear Relationships**: These trees can model complex relationships in data, which linear models may miss.

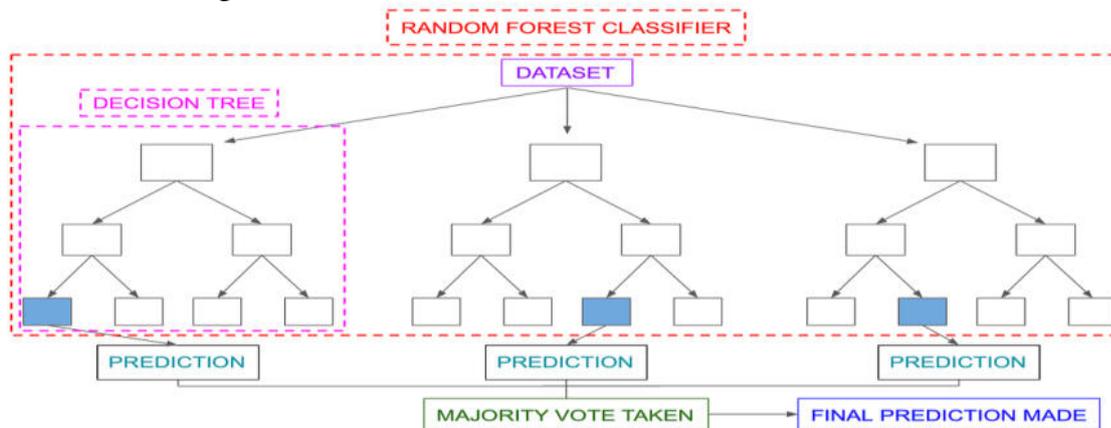## Disadvantages of Decision Trees

1. **Overfitting**: One of the most significant challenges with Decision Trees is their tendency to overfit. This happens when a model learns not only the underlying pattern but also the noise in the training data.

B Venkatesu Goud, Assistant Professor

2. **Instability**: Small changes in the data can lead to very different tree structures, which makes Decision Trees somewhat unreliable.
3. **Bias towards Dominant Classes**: In cases of imbalanced datasets, Decision Trees may be biased towards the more frequent class.
Introduction to Random Forests

## What Is a Random Forest?

A Random Forest is an ensemble learning method that combines the predictions of multiple Decision Trees to improve the overall accuracy and robustness of predictions. Think of it as a committee of Decision Trees that votes to reach a final decision.

The idea is simple but powerful—by averaging the results across many trees, you reduce the risk of overfitting and create a more stable model.



## How Do Random Forests Work?

1. **Bootstrapping**: Random Forests employ a technique called bootstrapping, which involves sampling subsets of the original dataset with replacement to train multiple trees.
2. **Feature Randomness**: At each split in the trees, only a random subset of the features is considered for splitting, which enhances the diversity among trees.
3. **Building Multiple Decision Trees**: Each tree is constructed independently. After all trees are built, they collaborate to make predictions through voting (in classification) or averaging (in regression).
4. **Final Prediction**: The collective output of all trees is what you end up with, resulting in a more balanced and accurate prediction.

## Advantages of Random Forests

1. **High Accuracy**: By aggregating multiple trees, Random Forests often achieve higher accuracy than single Decision Trees.
2. **Reduced Overfitting**: The combination of multiple models minimizes the chance of overfitting, making it a more reliable choice.
3. **Handles Missing Values**: Random Forests can handle missing data well, making them adaptable to various scenarios.
4. **Feature Importance**: You can evaluate the importance of different features in predicting outcomes, which can guide further analysis.

## Disadvantages of Random Forests

1. **Complexity**: While Random Forests are powerful, they can lack the interpretability of Decision Trees. It's often challenging to understand how the final decision was made.
2. **Resource-Intensive**: Training a large number of trees requires more computational power and memory, which might not be feasible for all users.

B Venkatesu Goud, Assistant Professor

3. **Longer Training Time**: Compared to a single Decision Tree, Random Forests have longer training times due to the multiple models they build.

## => <span style="color:red">k-Nearest Neighbors (k-NN) and Naïve Bayes, Support Vector Machines (SVM):</span>

**Naïve Bayes**
Naïve Bayes classifiers are probabilistic models based on Bayes' theorem, assuming conditional independence between features. Despite the "naïve" assumption, they perform well in many practical applications, especially text classification.

$$P(C \mid X) = \frac{P(X \mid C)P(C)}{P(X)}$$

Where $C$ is the class, $X$ is the feature vector, $P(C)$ is the prior probability of the class, and $P(X \mid C)$ is the likelihood of features given the class.

**Advantages:**
- Computationally efficient
- Performs well with high-dimensional data
- Requires small training datasets

**Limitations:**
- Assumption of feature independence rarely holds in practice
- Less effective for correlated features
- Probabilistic outputs may be less intuitive for some applications

*Example:* In email spam detection, Naïve Bayes evaluates the probability of a message being spam based on word occurrences. Despite correlations between words, the model achieves high accuracy due to its probabilistic framework.

**Real-World Application:**
Consider a telecommunications company predicting customer churn. Features include contract type, monthly charges, tenure, and customer complaints. Logistic regression provides interpretable coefficients, indicating which factors increase churn probability. Decision trees segment customers into high-risk groups. Random forests enhance prediction accuracy by aggregating multiple trees. Naïve Bayes is employed for automated email classification, detecting customers at risk of churn through interaction patterns and textual feedback.

**Advantages of Classification Models in Predictive Analytics:**
- Ability to handle categorical and numerical data
- Capture linear and non-linear relationships
- Provide probabilistic and interpretable outputs (logistic regression, Naïve Bayes)
- Robust predictive power with ensemble methods (random forests)

**Challenges:**
- Overfitting in complex trees
- Selecting appropriate hyperparameters
- Handling imbalanced datasets (requires preprocessing strategies like SMOTE)

- Ensuring interpretability in ensemble methods

k-NEAREST NEIGHBORS (k-NN) AND SUPPORT VECTOR MACHINES (SVM)
Classification and regression tasks in predictive analytics often require models that can handle non-linear relationships, high-dimensional data, and complex decision boundaries. Two powerful algorithms that address these challenges are k-Nearest Neighbors (k-NN) and Support Vector Machines (SVM). While both are supervised learning algorithms, they operate on different principles: k-NN is an instance-based, non-parametric method relying on feature similarity, whereas SVM is a margin-based method seeking an optimal separating hyperplane. This part explores these models in depth, covering their theoretical foundations, mathematical formulations, practical applications, evaluation methods, advantages, limitations, and real-world case studies.

### k-Nearest Neighbors (k-NN)
The k-Nearest Neighbors algorithm predicts the class of a new observation by examining the 'k' closest data points in the feature space. The proximity between points is typically measured using distance metrics such as Euclidean distance, Manhattan distance, or Minkowski distance. For regression tasks, k-NN predicts the target as the average of the k nearest neighbors' values. For classification, it assigns the majority class among the neighbors.

### Mathematical Formulation:
For a given query point $x_q$, the Euclidean distance to a training point $x_i$ in n-dimensional space is:

$$d(x_q, x_i) = \sqrt{\sum_{j=1}^{n} (x_{qj} - x_{ij})^2}$$

The algorithm then selects the k points with the smallest distances and determines the predicted outcome:
- **Classification:** Majority voting among k neighbors
- **Regression:** Mean (or weighted mean) of neighbor values

### Key Considerations in k-NN:
1. Choice of k: A small k may lead to overfitting, while a large k can oversmooth decision boundaries.
2. Distance metric: Euclidean is common, but Manhattan or Minkowski may be better for certain data distributions.
3. Feature scaling: Essential, as distance calculations are sensitive to feature magnitudes.

*Example:* Predicting whether a customer will buy a product based on age, income, and previous purchase behavior. For a new customer, k-NN identifies the closest k customers in the training set and assigns the majority purchase behavior.
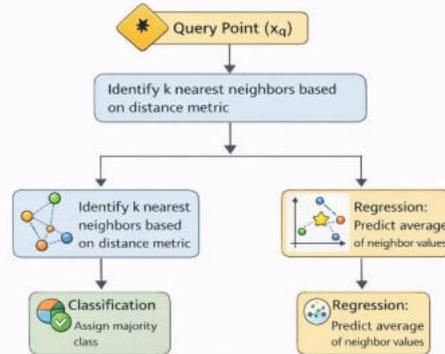
### Advantages:
- Simple to understand and implement
- Non-parametric: No assumptions about underlying data distribution
- Can model non-linear decision boundaries naturally

### Limitations:
- Computationally expensive for large datasets

B Venkatesu Goud, Assistant Professor

- Sensitive to irrelevant or noisy features
- Requires effective feature scaling
- Performance degrades in high-dimensional spaces (curse of dimensionality)

**Text-Based Diagram Representation:**



This figure illustrates the k-Nearest Neighbors (k-NN) algorithm, a widely used method in both classification and regression tasks. The process begins with a query point ($x_q$), representing the data instance for which a prediction is required. The algorithm first identifies the k nearest neighbors of this query point using a chosen distance metric, such as Euclidean or Manhattan distance, which measures similarity between data points. Once the nearest neighbors are determined, the workflow splits into two distinct predictive paths:

1. Classification: In this path, the query point is assigned the majority class among its k neighbors. This approach allows the model to categorize the instance based on the most common label in its local neighborhood.
2. Regression: Here, the prediction is made by computing the average of the neighbor values, providing a continuous output. This approach leverages the proximity of neighboring points to estimate the target value for the query point.

Overall, this diagram captures the key steps of k-NN, emphasizing its simplicity, reliance on distance-based similarity, and dual applicability for both categorical and continuous prediction tasks. The visualization clarifies how local information around a query point is utilized to make predictions, making it suitable for academic explanations of instance-based learning.

**Real-World Case Study:**

A retail company wants to predict customer loyalty. Features include purchase frequency, average spending, and product categories. k-NN identifies the k most similar customers and predicts whether a new customer will become loyal based on neighbor behavior. Feature scaling ensures that high-spending customers do not disproportionately influence the distance metric.

**Support Vector Machines (SVM)**
SVM is a supervised learning algorithm that constructs a hyperplane or set of hyperplanes in high-dimensional space to separate classes with maximum margin. The hyperplane represents

a decision boundary that best separates data points of different classes. SVM can handle linear and non-linear separations using kernel functions, which implicitly map input features to higher-dimensional spaces.

**Mathematical Formulation:**

For linearly separable data, SVM seeks to maximize the margin $M$ between classes:

$$\text{Maximize } M = \frac{2}{\parallel \mathbf{w} \parallel}$$

Subject to constraints:

$$y_i(\mathbf{w} \cdot x_i + b) \geq 1, i = 1,2,\ldots,n$$

Where:
- $x_i$ are feature vectors
- $y_i \in \{-1,1\}$ are class labels
- $\mathbf{w}$ is the normal vector to the hyperplane
- $b$ is the bias term

For non-linear data, kernel functions such as polynomial, radial basis function (RBF), or sigmoid map features to a higher-dimensional space where a linear separation is possible.
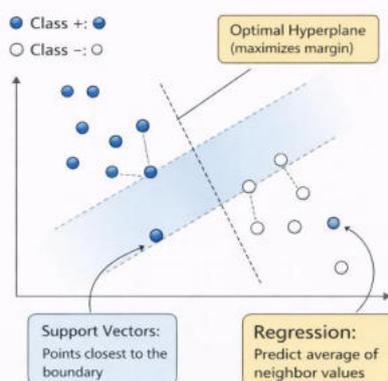
**Advantages of SVM:**
- Effective in high-dimensional spaces
- Robust to overfitting, especially with regularization
- Can model complex non-linear boundaries using kernel tricks
- Provides clear geometric interpretation of classification boundaries

**Limitations:**
- Computationally intensive for very large datasets
- Choice of kernel and hyperparameters significantly impacts performance
- Less interpretable than simpler models like logistic regression
- Sensitive to noise and overlapping classes

**Diagram Representation:**



This figure illustrates the fundamental concepts of a Support Vector Machine (SVM) for binary classification. The data points belong to two classes: Class +, represented by filled circles (●), and Class -, represented by empty circles (○). The SVM algorithm identifies an optimal hyperplane, which is the decision boundary that maximizes the margin between the two classes. The support vectors are the data points that lie closest to this hyperplane, and

they are crucial because they define the position and orientation of the boundary. When a new data point is introduced, its classification depends on which side of the hyperplane it falls. Points on one side are assigned to Class +, while points on the other side are assigned to Class -.

This diagram effectively conveys the geometric intuition behind SVM, highlighting how the margin and support vectors determine the classification, and making it suitable for academic explanations in courses on machine learning or predictive analytics.

**Real-World Example:**

In healthcare, SVM predicts the presence of a disease based on patient biomarkers. Features include lab results, demographic data, and lifestyle indicators. By using an RBF kernel, SVM captures complex non-linear interactions between biomarkers to distinguish between healthy and diseased patients. Support vectors identify critical boundary cases influencing classification.

**Evaluation Metrics for k-NN and SVM:**
- Confusion matrix (True Positives, False Positives, True Negatives, False Negatives)
- Accuracy, Precision, Recall, F1 Score
- ROC-AUC curve for threshold-independent evaluation

**Challenges:**
1. k-NN: Computational cost for large datasets, sensitivity to high dimensionality, and noise
2. SVM: Kernel selection, hyperparameter tuning, and interpretability
3. Both require careful feature scaling and pre-processing

**Case Study:**

A bank wants to detect fraudulent transactions. k-NN identifies similarity between new and historical transactions, predicting fraud based on nearest neighbours. SVM separates fraudulent and legitimate transactions using a non-linear kernel, capturing complex interactions among transaction amount, location, merchant type, and time. Pre-processing steps including scaling and handling class imbalance are critical for accurate detection.

## => **Model Selection and Comparison:**

**The Model Selection Framework**
Choosing the right model is not just about accuracy; it involves balancing complexity, data size, and the nature of the problem.[1]
- **Define the Task:**
  - **Regression:** Predicting continuous values (e.g., house prices).[2]
  - **Classification:** Grouping data into categories (e.g., spam vs. not spam).[3]
  - **Clustering:** Finding hidden patterns without labels (e.g., customer segments).
- **Assess Data Characteristics:**

- **Size:** Small datasets often favor simpler models (Linear Regression, Naïve Bayes) to avoid overfitting.[4] Large datasets can leverage Deep Learning.[5]
- **Linearity:** If the relationship is straight-line, use Linear models.[6] If it's complex, use Tree-based models (Random Forest) or Neural Networks.
- **Interpretability vs. Performance:** In regulated industries like healthcare or finance, an interpretable model (Decision Tree) is often preferred over a "black box" model (Neural Network), even if the latter is slightly more accurate.

## Model Comparison Techniques

To compare models fairly, you must evaluate them on data they haven't seen before.[7]

### Resampling Methods

- **Train/Test Split:** Dividing your data (typically **80/20**) to train on one set and validate on the other.[8]
- **K-Fold Cross-Validation:** The data is split into [9]$k$ subsets.[10] The model is trained [11]$k$ times, each time using a different subset as the test set and the rest as training data.[12] The results are averaged for a robust score.

### Comparison Metrics

| Metric Type | Key Metrics | Best Used For... |
|---|---|---|
| **Classification** | Accuracy, Precision, Recall, F1-Score | Determining how well the model categorizes data. |
| **Regression** | RMSE, MAE, $R^2$ | Measuring the distance between predicted and actual values. |
| **Probabilistic** | Log Loss, AUC-ROC | Evaluating the model's confidence in its predictions. |

## The Trade-off: Bias vs. Variance

A critical part of comparison is understanding why a model fails.

- **High Bias (Underfitting):** The model is too simple and misses the trend (e.g., using a straight line for a curved trend).[13]
- **High Variance (Overfitting):** The model is too complex and "memorizes" the noise in the training data, failing to generalize to new data.[14]

### Summary Checklist for Comparison

1. **Check for Overfitting:** Is the training accuracy much higher than the test accuracy?
2. **Computational Cost:** How long does the model take to train and predict?
3. **Stability:** Does the model performance change drastically with small changes in the data?

# UNIT IV

**Training, Testing, and Validation Sets, Cross-Validation Techniques (k-Fold, Stratified, LOOCV), Evaluation Metrics: Accuracy, Precision, Recall, F1 Score, ROC-AUC, Confusion Matrix and Classification Report, Bias-Variance Trade-off and Overfitting, Hyperparameter Tuning: Grid Search, Random Search.**

## TRAINING, TESTING, AND VALIDATION SETS

A fundamental principle in predictive modeling is to separate data into distinct subsets to train and evaluate models effectively. This prevents overfitting, where a model performs exceptionally well on training data but poorly on unseen data. The standard approach involves dividing the dataset into a **training set**, a **validation set**, and a **test set**.

**Training Set:** The training set is used to fit model parameters. It contains the majority of the data and provides the examples from which the algorithm learns relationships between features and the target variable. Proper training requires representative sampling to capture the underlying patterns of the population.

**Validation Set:** The validation set is used for model selection and hyperparameter tuning. For example, in decision trees, the maximum depth or minimum samples per leaf can be optimized by evaluating performance on the validation set. Using a validation set ensures that hyperparameters are chosen without bias from the test set.

**Test Set:** The test set is reserved for final evaluation after model training and tuning are complete. It provides an unbiased assessment of model generalizability. The model is evaluated only once on the test set to estimate its performance in real-world scenarios.

**Data Splitting Strategies:**
1. **Holdout Method:** Randomly splits the data into training, validation, and test sets. Common ratios include 70:15:15 or 60:20:20. While simple, this method may introduce variance depending on the random split.
2. **Stratified Sampling:** Ensures that class proportions are maintained across splits, particularly important for imbalanced datasets.
3. **Temporal Splits:** For time series data, the split respects chronological order, using past data for training and future data for testing.

*Example:* In a retail churn prediction model, 70% of customer data is used for training, 15% for validation to tune model hyperparameters, and 15% for testing to estimate real-world predictive performance. Stratified sampling ensures that churners and non-churners are proportionally represented in all subsets.
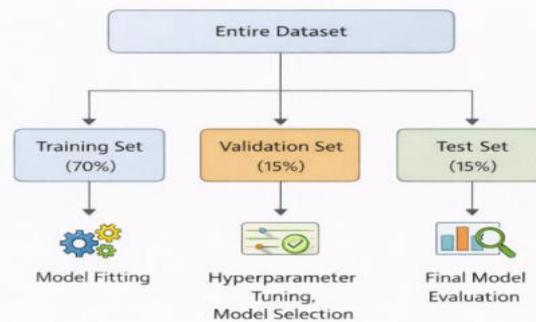
**Advantages:**
- Reduces overfitting
- Provides unbiased model evaluation
- Enables systematic hyperparameter tuning

**Challenges:**
- Small datasets may yield unstable estimates due to limited validation data
- Random splits may introduce sampling bias
- Ensuring consistent pre-processing across subsets is critical to prevent data leakage

**Diagram Representation:**

This figure illustrates the standard procedure for splitting a dataset in machine learning. The entire dataset is divided into three subsets: training set, validation set, and test set. The training set, typically comprising 70% of the data, is used for model fitting, allowing the algorithm to learn patterns from the data. The validation set, usually 15% of the dataset, is employed for hyperparameter tuning and model selection, helping to optimize the model's performance and avoid overfitting. Finally, the test set, also generally 15% of the data, is used for final model evaluation, providing an unbiased assessment of the model's predictive performance on unseen data.

This workflow ensures that the model is trained, tuned, and evaluated properly, maintaining generalizability and robustness. The diagram visually distinguishes each subset and clearly depicts their roles in the machine learning pipeline, making it suitable for academic explanations of model development best practices.

## CROSS-VALIDATION TECHNIQUES

Cross-validation provides a more robust evaluation than a single train-test split, particularly for small or medium-sized datasets. It involves partitioning the data into k subsets (folds), training the model on k-1 folds, and testing on the remaining fold. This process is repeated k times, and performance metrics are averaged to obtain a more reliable estimate of model generalization.

**k-Fold Cross-Validation:**
- Divide data into k equally sized folds.
- Train on k-1 folds and validate on the remaining fold.
- Repeat for all folds and compute average metrics.

*Example:* With k=5, the dataset is split into 5 folds. Fold 1 is used for validation while folds 2–5 are used for training. The process repeats until each fold serves as validation once.

**Stratified k-Fold Cross-Validation:**
- Ensures class distributions are preserved in each fold.
- Essential for imbalanced classification tasks, preventing folds from containing disproportionate minority or majority class instances.

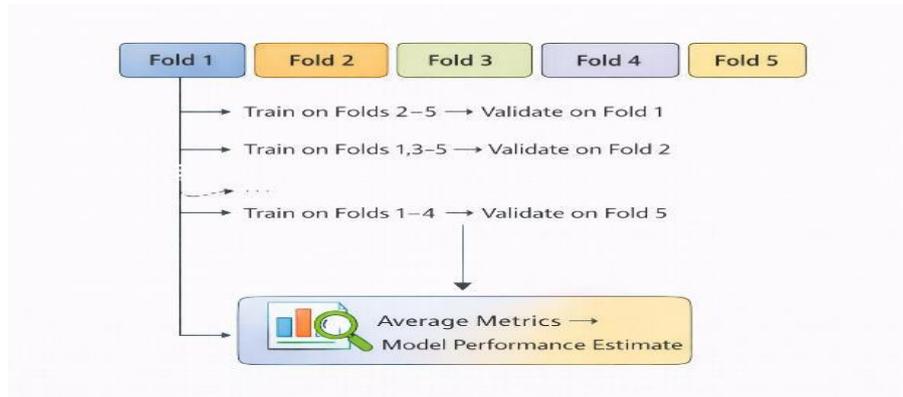**Leave-One-Out Cross-Validation (LOOCV):**
- Each observation is used once as validation while all others form the training set.
- Provides an almost unbiased estimate but is computationally expensive for large datasets.

**Advantages of Cross-Validation:**
- Reduces variance in performance estimates
- Provides robust evaluation for small datasets
- Enables hyperparameter tuning without overfitting

**Limitations:**
- Computationally intensive for large datasets
- Stratification required for imbalanced classes
- LOOCV may have high variance for noisy data



This figure illustrates the k-Fold Cross Validation (CV) process, a widely used technique to evaluate the generalization performance of machine learning models. In this method, the dataset is divided into k equally sized folds. For each iteration, one fold is used as the validation set, while the remaining k-1 folds serve as the training set. For example, in a 5-fold CV, Fold 1 is held out for validation while Folds 2–5 are used for training. The process repeats for each fold: Fold 2 is validated while Folds 1, 3–5 are used for training, and so on, until each fold has served once as the validation set. After all iterations, the performance metrics from each fold are averaged to provide a robust estimate of the model's predictive performance.

This diagram effectively demonstrates how k-Fold Cross Validation maximizes the use of available data for both training and validation, reduces variance in performance estimates, and helps in model selection and hyperparameter tuning. It is particularly useful for small to medium-sized datasets where a separate validation set may reduce the training data excessively.

## EVALUATION METRICS FOR REGRESSION AND CLASSIFICATION

Selecting appropriate evaluation metrics is crucial for assessing predictive model quality. Different metrics capture different aspects of performance and must align with business or domain objectives.

**Regression Metrics:**
1. **Mean Absolute Error (MAE):** Measures average absolute deviations between actual and predicted values.
2. **Mean Squared Error (MSE):** Penalizes larger errors more heavily by squaring deviations.
3. **Root Mean Squared Error (RMSE):** Square root of MSE, interpretable in original units.
4. **R-squared ($R^2$):** Fraction of variance explained by the model, indicating goodness-of-fit.

**Classification Metrics:**
1. **Accuracy:** Fraction of correctly classified instances.

2. **Precision:** Proportion of true positives among predicted positives.
3. **Recall (Sensitivity):** Proportion of true positives detected among actual positives.
4. **F1 Score:** Harmonic mean of precision and recall, balancing false positives and negatives.
5. **Confusion Matrix:** Detailed view of true/false positives and negatives.
6. **ROC-AUC:** Threshold-independent measure of classifier discriminative ability.

*Example:* In a medical diagnosis model, recall may be prioritized to ensure that as many patients with a disease are correctly identified, even at the expense of some false positives.

**Advantages:**
- Tailored metrics provide insights aligned with domain objectives
- Facilitate comparison of multiple models
- Highlight model strengths and weaknesses

**Challenges:**
- Choosing the right metric for the task
- Balancing trade-offs between metrics (e.g., precision vs. recall)
- Interpretation complexity for multi-class classification

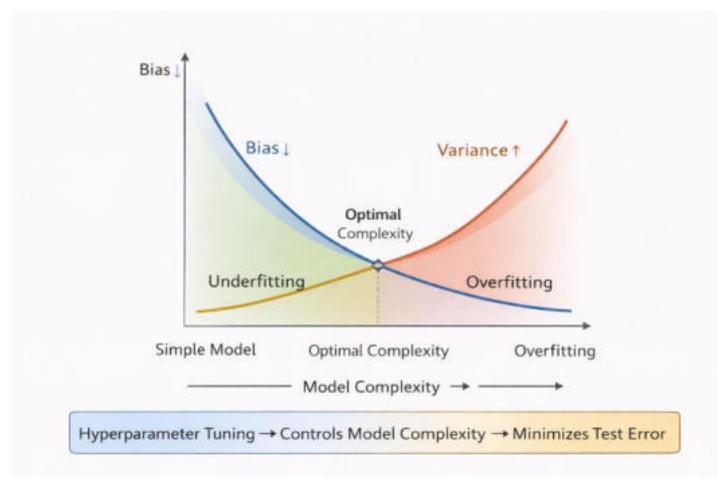BIAS-VARIANCE TRADE-OFF AND HYPERPARAMETER TUNING

Predictive models must balance bias (error due to oversimplification) and variance (error due to overfitting). High-bias models underfit data, failing to capture patterns. High-variance models overfit training data, performing poorly on unseen data.

**Hyperparameter Tuning:**
- **Grid Search:** Exhaustively searches all combinations of hyperparameters.
- **Random Search:** Randomly samples hyperparameter space, often more efficient than grid search.
- **Bayesian Optimization:** Uses probabilistic models to intelligently explore hyperparameter space.

*Example:* Tuning the maximum depth of decision trees or the regularization parameter in SVM ensures optimal generalization without overfitting.

**Diagram Representation:**



This figure illustrates the bias-variance tradeoff, a fundamental concept in machine learning that describes the relationship between model complexity and predictive performance. On the x-axis, model complexity ranges from simple models (low capacity) to highly complex

models (high capacity). On the y-axis, bias decreases as model complexity increases, while variance increases.

- **Underfitting (Simple Model):** When the model is too simple, it cannot capture the underlying patterns in the data, leading to high bias and poor training performance.
- **Optimal Complexity:** At this point, the model balances bias and variance, achieving minimal test error and generalizing well to unseen data.
- **Overfitting (Complex Model):** Excessively complex models fit the training data very closely, capturing noise rather than true patterns. This results in low bias but high variance, causing poor performance on new data.

The figure also emphasizes the role of hyperparameter tuning, which adjusts model parameters to control complexity and minimize test error. This diagram effectively conveys the need to strike a balance between underfitting and overfitting for robust model performance, making it suitable for academic explanations of model evaluation and optimization.

**Challenges:**
- Computationally expensive for large datasets
- Requires robust validation strategy
- Risk of overfitting validation set if tuned excessively

# UNIT V

**Ensemble Learning: Bagging, Boosting (AdaBoost, XGBoost), Predictive Analytics with Time Series (ARIMA, Prophet), Deep Learning for Predictive Modeling (ANNs, LSTM), Use of Predictive Analytics in IoT, Retail, and Healthcare, Ethics and Privacy in Predictive Analytics, Building and Deploying End-to-End Predictive Systems.**

## ENSEMBLE LEARNING: BAGGING AND BOOSTING

Ensemble learning refers to combining multiple individual models (base learners) to create a single, robust predictive model. By leveraging the strengths of multiple models, ensemble methods reduce variance, improve accuracy, and mitigate overfitting. Two widely used ensemble techniques are **bagging** and **boosting**.

### Bagging (Bootstrap Aggregating)

Bagging involves training multiple models on different bootstrap samples of the training data. Each model independently predicts outcomes, and final predictions are aggregated via majority voting (classification) or averaging (regression). Random Forest is a prominent example of bagging applied to decision trees.

*Mathematical*                                                        *Concept:*

Given a training set of size $N$, bagging creates $B$ bootstrap samples by sampling with replacement. Each sample trains a base model $h_b(x)$, and predictions are combined as:

- Classification:

$$\hat{y} = \text{mode}(h_1(x), h_2(x), \ldots, h_B(x))$$

- Regression:

$$\hat{y} = \frac{1}{B} \sum_{b=1}^{B} h_b(x)$$

**Advantages:**
- Reduces variance without increasing bias
- Robust to overfitting
- Handles high-dimensional datasets effectively

**Limitations:**
- Requires more computational resources
- Less interpretable than single models
- Dependent on diversity among base models

### Boosting

Boosting sequentially trains models, each focusing on correcting errors of the previous ones. AdaBoost and XGBoost are widely used boosting algorithms. Boosting assigns higher weights to misclassified instances, encouraging subsequent models to focus on difficult examples.

*Mathematical Concept (AdaBoost):*

1. Initialize weights $w_i = 1/N$ for all training examples
2. Train a weak learner $h_t(x)$
3. Compute weighted error $\epsilon_t = \sum_i w_i\, I(y_i \neq h_t(x_i))$
4. Compute model weight $\alpha_t = 0.5\ln\dfrac{1-\epsilon_t}{\epsilon_t}$
5. Update example weights:

$$w_i \leftarrow w_i \cdot e^{-\alpha_t y_i h_t(x_i)}$$

**Advantages:**
- Reduces both bias and variance
- Focuses on difficult examples
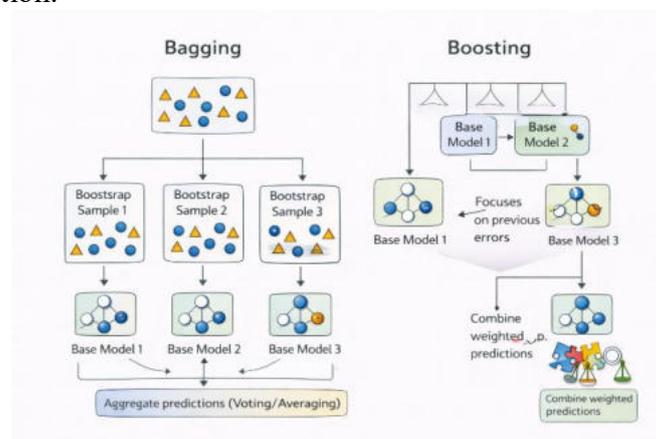- Often achieves high accuracy on structured datasets

**Limitations:**
- Sensitive to noisy data and outliers
- Computationally intensive
- Risk of overfitting with too many iterations

**Case Study:**
In a credit scoring application, Random Forest (bagging) predicts loan defaults by aggregating multiple decision trees. XGBoost (boosting) further improves predictions by focusing on customers misclassified in previous iterations, capturing subtle patterns in income, repayment history, and transaction behavior. Feature importance analysis identifies key variables influencing default risk.

Diagram Representation:



This figure illustrates the two main types of ensembles learning techniques: Bagging and Boosting.

Bagging (Bootstrap Aggregating): The process begins by generating multiple bootstrap samples from the original dataset. Each bootstrap sample trains a separate base model, such as a decision tree. After all models are trained independently, their predictions are aggregated using methods like majority voting for classification or averaging for regression. Bagging reduces variance and improves model stability without significantly increasing bias, making it effective for high-variance models.

Boosting: In contrast, boosting builds models sequentially, where each new model focuses on the errors made by previous models. Each base model is assigned a weight based on its

accuracy, and the final prediction is a weighted combination of all models. Boosting reduces bias and improves predictive accuracy but can increase variance if overfitting occurs.

The diagram effectively contrasts the two approaches: Bagging emphasizes parallel model training and aggregation, while Boosting emphasizes sequential error correction and weighted combination, making it suitable for academic explanations of ensemble techniques in predictive analytics.

## PREDICTIVE ANALYTICS WITH TIME SERIES

Time series data, characterized by sequential dependencies, requires specialized predictive models. Unlike independent observations, time series exhibit trends, seasonality, and autocorrelation. Predictive modeling techniques include **ARIMA**, **Seasonal ARIMA (SARIMA)**, and **Prophet**, each designed to handle different aspects of temporal data.

**ARIMA (Autoregressive Integrated Moving Average)**

ARIMA models forecast based on three components:

- **Autoregression (AR):** Regression of the target on its own previous values
- **Integration (I):** Differencing to achieve stationarity
- **Moving Average (MA):** Modeling residual errors of past observations

Model parameters $(p, d, q)$ represent AR order, differencing order, and MA order, respectively. ARIMA is effective for univariate stationary series and short-term forecasting.

**Prophet**

Prophet is a decomposable time series model combining trend, seasonality, and holiday effects. It is robust to missing data, outliers, and non-linear trends, making it suitable for business applications such as sales forecasting and web traffic prediction.

*Example:* A retail chain forecasts weekly sales for multiple stores. ARIMA captures local trends in individual stores, while Prophet models overall trends with seasonality and promotional events.

**Challenges in Time Series Predictive Analytics:**

- Handling missing or irregularly spaced data
- Capturing complex seasonality patterns
- Incorporating exogenous variables (promotions, holidays)
- Avoiding overfitting to historical fluctuations

## DEEP LEARNING FOR PREDICTIVE MODELING: ANNS AND LSTM

**Artificial Neural Networks (ANNs)**

ANNs are highly flexible models inspired by biological neurons. They consist of input layers, hidden layers, and output layers, with weighted connections and non-linear activation functions. ANNs can model complex, non-linear relationships in high-dimensional datasets.

**Key Concepts:**

- Forward propagation: Computes predictions from inputs through layers
- Backpropagation: Optimizes weights by minimizing loss function
- Activation functions: Sigmoid, ReLU, Tanh, enabling non-linear transformations

*Example:* Predicting customer lifetime value based on demographics, purchase behavior, and engagement metrics. ANNs capture intricate interactions between multiple features to forecast long-term revenue.

B Venkatesu Goud, Assistant Professor

**Long Short-Term Memory (LSTM) Networks**
LSTM networks are specialized recurrent neural networks (RNNs) capable of modeling sequential data with long-term dependencies. They include memory cells, input/output gates, and forget gates, allowing retention and selective updating of relevant information.
*Applications:*
- Stock price prediction
- Sales forecasting
- Sensor data analysis in IoT

**Advantages:**
- Handles non-linear, complex relationships
- Captures temporal dependencies in sequential data
- Scalable to high-dimensional inputs

**Limitations:**
- Requires large datasets for effective training
- Computationally intensive
- Difficult to interpret compared to traditional models

## APPLICATIONS OF PREDICTIVE ANALYTICS IN INDUSTRY

**Healthcare:** Predictive models assist in disease diagnosis, treatment planning, and patient risk stratification. Time series models forecast patient vitals in ICU, while ensemble learning predicts disease progression based on multi-modal patient data.

**Retail:** Predictive analytics enhances demand forecasting, inventory management, and personalized recommendations. Deep learning and time series models capture seasonal trends, promotional impacts, and customer purchase behavior.
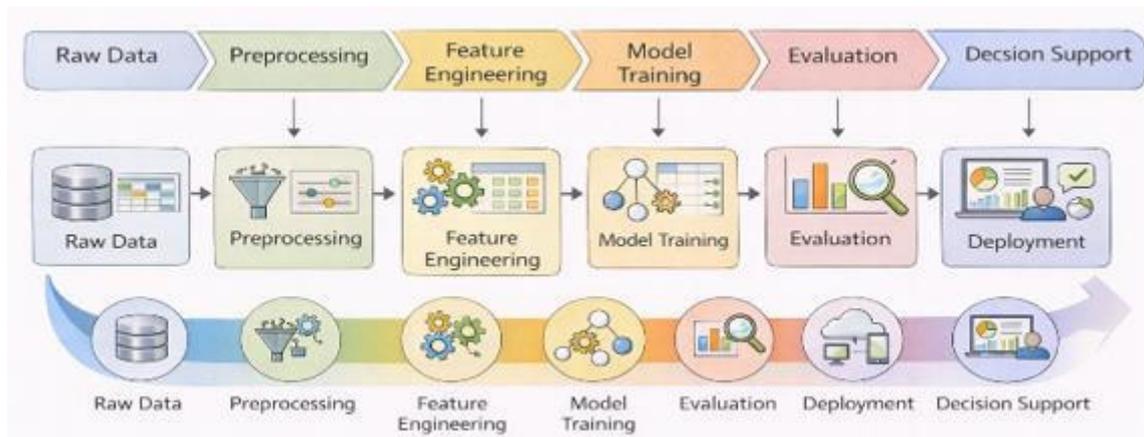
**IoT and Smart Devices:** Sensor data from IoT devices is analyzed to predict equipment failures, optimize energy consumption, and detect anomalies. LSTM networks and ensemble models handle continuous data streams in real-time applications.

**Ethics and Privacy:** Predictive analytics must consider ethical implications, data privacy, and regulatory compliance. Misuse of predictive models can lead to biased decisions, discrimination, or privacy breaches.

**Building End-to-End Predictive Systems:**
1. **Data Ingestion:** Collect raw data from multiple sources
2. **Data Pre-processing:** Cleaning, feature engineering, scaling, and encoding
3. **Modeling:** Select and train appropriate models (regression, classification, deep learning, or ensembles)
4. **Evaluation:** Apply cross-validation, metrics, and hyperparameter tuning
5. **Deployment:** Integrate models into applications or dashboards for real-time predictions

**Diagram Representation:**

This figure illustrates the complete workflow of a predictive analytics project, starting from raw data and ending with decision support. The process begins with Raw Data, which is collected from various sources and forms the basis for analysis. The data undergoes pre-processing, where it is cleaned, normalized, and transformed to ensure quality and consistency. Next, feature engineering is performed to create or select meaningful variables that improve model performance.

The prepared dataset is then used for model training, where machine learning algorithms learn patterns and relationships within the data. Following training, the model undergoes evaluation using appropriate metrics to measure its predictive performance and generalizability. Once validated, the model is deployed into operational systems or applications to generate predictions on new data. Finally, these predictions inform decision support, enabling stakeholders to make data-driven decisions effectively.

This diagram effectively summarizes the predictive analytics lifecycle, highlighting the sequential and interdependent nature of each stage, and serves as a foundational visual for academic and professional learning in data science and machine learning.