

## Programs Offered

### Post Graduate Programmes (PG)

- Master of Business Administration
- Master of Computer Applications
- Master of Commerce (Financial Management / Financial Technology)
- Master of Arts (Journalism and Mass Communication)
- Master of Arts (Economics)
- Master of Arts (Public Policy and Governance)
- Master of Social Work
- Master of Arts (English)
- Master of Science (Information Technology) (ODL)
- Master of Science (Environmental Science) (ODL)

### Diploma Programmes

- Post Graduate Diploma (Management)
- Post Graduate Diploma (Logistics)
- Post Graduate Diploma (Machine Learning and Artificial Intelligence)
- Post Graduate Diploma (Data Science)

### Undergraduate Programmes (UG)

- Bachelor of Business Administration
- Bachelor of Computer Applications
- Bachelor of Commerce
- Bachelor of Arts (Journalism and Mass Communication)
- Bachelor of Arts (General / Political Science / Economics / English / Sociology)
- Bachelor of Social Work
- Bachelor of Science (Information Technology) (ODL)



**AMITY** UNIVERSITY  
DIRECTORATE OF  
DISTANCE & ONLINE EDUCATION

Amity Helpline: 1800-102-3434 (toll-free), 0120-4614200

For Distance Learning Programmes: [dladmissions@amity.edu](mailto:dladmissions@amity.edu) | [www.amity.edu/addoe](http://www.amity.edu/addoe)

For Online Learning programmes: [elearning@amity.edu](mailto:elearning@amity.edu) | [www.amityonline.com](http://www.amityonline.com)

Foundations of Machine Learning

# Foundations of Machine Learning

AMITY



**AMITY** UNIVERSITY  
DIRECTORATE OF  
DISTANCE & ONLINE EDUCATION

# **Foundations of Machine Learning**



**AMITY** | DIRECTORATE OF DISTANCE &  
UNIVERSITY ONLINE EDUCATION

© Amity University Press

**All Rights Reserved**

No parts of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise without the prior permission of the publisher.

**SLM & Learning Resources Committee**

*Chairman* : Prof. Abhinash Kumar

*Members* : Dr. Divya Bansal  
Dr. Coral J Barboza  
Dr. Monica Rose  
Dr. Winnie Sharma

*Member Secretary* : Ms. Rita Naskar

# Contents

	Page No.
<b>Module - I: Introduction to Machine Learning</b>	<b>01</b>
1.1 Machine Learning	
1.1.1 Brief History and Importance of Machine Learning	
1.1.2 Role of Data and Connection to the Knowledge /Experience in Learning	
1.1.3 Show Successful Examples of Machine Learning in Industry/Working	
1.1.4 Motivation and Role of Machine Learning in Computer Science and Problem Solving	
1.1.5 Connect Machine Learning to the Broader Theme of Computer Science	
1.2 Types of Machine Learning	
1.2.1 Supervised Machine Learning	
1.2.2 Unsupervised Machine Learning	
1.2.3 Reinforcement Learning	
1.3 Machine Learning Algorithms	
1.3.1 Statistical View of Machine Learning	
1.3.2 Multiple Linear Regression	
1.4 ML Applications in Retail	
1.4.1 How Machine Learning Improves Pricing and Promotions	
1.4.2 Using ML to Justify Promotions	
1.4.3 ML Helps Retailers Avoid Promotional Waste	
1.4.4 ML Determines Optimal Product Prices	
1.4.5 ML Saves Time and Effort	
1.4.6 ML Reduces Promo Risks and Improve Business Bottom Line	
<b>Module - II: Exploratory Data Analysis</b>	<b>61</b>
2.1 Introduction to Data and Data Distribution	
2.1.1 What is Exploratory Data Analysis (EDA)	
2.1.2 Processes in EDA	
2.2 Data Types	
2.2.1 Handling of Data Types	
2.2.2 Univariate and Bivariate Analysis	
2.2.3 Visualising the Data	
2.3 Feature Extraction	
2.3.1 Feature Extraction: Part I	
2.3.2 Feature Extraction: Part II	
2.4 Dimensionality Reduction	
2.4.1 Need for Dimensionality Reduction	
2.4.2 Introduction of Dimensionality Reduction	
2.4.3 Introduction to Principle Component Analysis (PCA)	
2.4.4 Disadvantage of PCA	
2.5 Data Visualisation	
2.5.1 Understanding Data Visulisation Concepts	
<b>Module - III: Supervised Learning</b>	<b>118</b>
3.1 Introduction to Supervised Learning	

- 3.1.1 Introductory Concepts of Supervised Learning
- 3.2 Introduction to Linear Regression
  - 3.2.1 Linear Regression: Important Terms
  - 3.2.2 Linear Regression: Assumptions
  - 3.2.3 Gradient Descent
  - 3.2.4 Linear Regression Using Gradient Descent
  - 3.2.5 Linear Regression: Real Life Events
- 3.3 Logistic Regression
  - 3.3.1 Logistic Regression Using Gradient Descent

#### **Module - IV: Unsupervised Learning**

**159**

- 4.1 Introduction to Clustering
  - 4.1.1 Introduction to Clustering
  - 4.1.2 Types of Clustering and Clustering Algorithms
  - 4.1.3 K-means Clustering
- 4.2 Evaluation of Clustering
  - 4.2.1 Understanding Clustering Types
  - 4.2.2 Anomaly Detection
- 4.3 K-means Clustering
  - 4.3.1 K-means Clustering
  - 4.3.2 Steps to Implement K-means
- 4.4 Principal Component Analysis
  - 4.4.1 What is Principle Component Analysis (PCA)
  - 4.4.2 Why Do We Need PCA
  - 4.4.3 Basic Terminologies of PCA
  - 4.4.4 How Does PCA Work
  - 4.4.5 Advantages and Disadvantage of PCA
  - 4.4.6 Application for Principal Component Analysis

#### **Module - V: Deep Learning**

**214**

- 5.1 Concepts of Deep Learning
  - 5.1.1 Introduction to Deep Learning: Part I
  - 5.1.2 Introduction to Deep Learning: Part II
- 5.2 Introduction to Neural Networks
  - 5.2.1 The Neuron Model
  - 5.2.2 The Neural Network and Its Utility in Modelling and Solving Problems
  - 5.2.3 Connect to the Biological Motivations and Parallelism
  - 5.2.4 Popular CNN (Convolutional Neural Network) Architectures
  - 5.2.5 RNNs (Recurrent Neural Networks)
- 5.3 Neural Networks: Its Types and Applications
  - 5.3.1 Perceptron
  - 5.3.2 Feed Forward Neural Network
  - 5.3.3 Multilayer Perceptron
  - 5.3.4 Convolutional Neural Network

# Module - I: Introduction to Machine Learning

Notes

## Learning Objectives

At the end of this module, you will be able to:

- Recognise the history and importance of machine learning
- Explain role of data and connection to knowledge/experience in learning
- Infer successful examples of machine learning in industry/working
- Summarise the role of machine learning in computer science and problem solving
- Discuss supervised, unsupervised and reinforcement machine learning
- Describe multiple linear regression and statistical view of machine learning
- Know how machine learning improves pricing and promotions
- Infer how machine learning helps retailers avoid promotional waste
- Explain how machine learning determines optimal product prices
- Discuss how machine learning saves time and effort
- Analyse how machine learning reduces promotional risks and improves business bottom line

## Introduction

Machine Learning (ML) is a specific domain within the field of artificial intelligence (AI) that is dedicated to the development of algorithms and models that enable computers to acquire knowledge and generate predictions or judgements by analysing and interpreting data. The process entails the creation of systems that possess the ability to enhance their performance autonomously over time, without the need for explicit programming tailored to each individual task. Machine learning (ML) has found widespread applications in diverse domains, including but not limited to finance, healthcare, marketing and robotics.

## Types of Machine Learning:

- **Supervised Learning:** Supervised learning involves the acquisition of knowledge by a model through the use of annotated training data, enabling the model to generate predictions based on the provided labels. Typical tasks encompass classification and regression.
- **Unsupervised Learning:** In the context of unsupervised learning, the model acquires knowledge of patterns through the analysis of unlabeled data. Clustering and dimensionality reduction are often encountered activities within this particular domain.
- **Reinforcement Learning:** The process of reinforcement learning entails the training of a computational model to make optimal judgements by use of iterative interactions with an external environment. The model is provided with feedback in the form of either rewards or punishments.

### 1.1 Machine Learning

In recent years, Machine Learning has emerged as a prominent component of information technology, assuming a significant but often inconspicuous role in our daily lives. Given the continuous growth of available data, it is reasonable to anticipate that

## Notes

the widespread adoption of smart data analysis will become increasingly essential for technological advancement. Machine learning can manifest itself in various forms. In the field of machine learning, a significant aspect is the process of transforming a diverse array of problems into a more focused collection of prototypes. A significant portion of the field of machine learning is dedicated to addressing these challenges and offering robust assurances for the derived solutions.

### 1.1.1 Brief History and Importance of Machine Learning

Machine learning (ML) plays a significant role in the pursuit of harnessing artificial intelligence (AI) technology. Machine learning is commonly referred to as AI due to its capacity for learning and decision-making, however it is more accurately classified as a subset of AI. It constituted a significant aspect of the evolutionary trajectory of artificial intelligence. Subsequently, it diverged and underwent independent evolutionary development. Machine learning has emerged as a crucial tool in the realm of cloud computing and e-commerce, finding applications in various state-of-the-art technologies.

Presented here is a concise overview of the historical development of machine learning and its significance in the realm of data management. Machine learning has become an indispensable component of contemporary commercial and research endeavours for numerous organisations in the present era. Algorithms and neural network models are employed to facilitate the gradual enhancement of computer systems' performance. Machine learning algorithms possess the ability to autonomously construct a mathematical model by utilising a set of sample data, commonly referred to as training data. This enables them to make judgements without the need for explicit programming instructions.

The significance of machine learning lies in its ability to provide organisations with insights into customer behaviour trends and company operational patterns, while also facilitating the creation of novel products. Machine learning has become an integral component of the operational strategies employed by prominent contemporary corporations, including Facebook, Google and Uber. The utilisation of machine learning has emerged as a notable factor that distinguishes firms in terms of their competitiveness.

The initial development of an electric circuit-based neural network can be attributed to Warren McCulloch and Walter Pitts. The objective of the network was to address a challenge initially presented by John von Neumann and other researchers, which involved devising a means for computers to establish communication among themselves. The initial prototype demonstrated the feasibility of computer-to-computer communication in the absence of human intervention. This particular event holds significance due to its role in facilitating the advancement of machine learning.

Despite the ongoing evolution of machine learning and the emergence of numerous new technologies, its utilisation remains prevalent across diverse industries.

Machine learning encompasses numerous practical applications that yield tangible business outcomes, including significant time and cost reductions. These outcomes possess the capacity to profoundly influence the future trajectory of your organisation. Significant advancements are observed in the customer service business, namely with the implementation of machine learning techniques. This has resulted in enhanced productivity and efficiency, enabling individuals to do tasks at a faster pace.

Virtual Assistant solutions utilise machine learning to automate operations that would often require the intervention of a human agent, such as password changes or account balance inquiries. This allows for the allocation of agent time, which can be dedicated

to the provision of customer care that is most effectively performed by humans, such as intricate decision-making and personalised interactions that are not as readily executed by automated systems. At Interactions, the process is enhanced by removing the need to determine whether a request should be directed to a human or a machine. This is achieved through the utilisation of a distinctive technology called Adaptive Understanding, wherein the machine acquires knowledge of its own limitations and defers to humans when it lacks confidence in delivering the accurate solution.

Machine learning is a field that encompasses both cybernetics and computer science, often known as Control Science and Computer Science. It has garnered significant attention from professionals and the general public in recent times. In recent years, the field of computer science has experienced notable advancements, such as the introduction of Graphics Processing Units (GPUs) and the subsequent enhancements in computational capabilities, as well as the development of specialised software enabling the processing of large datasets. As a result, machine learning has increasingly become associated with the domain of computer science.

Historically, it can be observed that learning algorithms that exhibit convergence and a satisfactory rate of convergence in the learning process have emerged from the field of cybernetics/control. The following is an examination of the historical perspective on machine learning from the standpoint of a control theorist. The author possesses a background in mathematics and has acquired expertise in the field of pattern recognition and control through the utilisation of adaptation and learning methodologies during the 1960s. The modern conception of machine learning is commonly attributed to Frank Rosenblatt, a psychologist affiliated with Cornell University. Drawing inspiration from the functioning of the human nervous system, Rosenblatt led a team that developed a machine capable of letter recognition. The machine, referred to as "perception" by its developer, employed a combination of analogue and discrete signals, incorporating a threshold component to convert analogue impulses into discrete counterparts.

The aforementioned system served as the prototype for contemporary artificial neural networks (ANN) and its learning model closely resembled the learning models created in psychology for animals and humans.

Rosenblatt conducted the initial mathematical investigations on the perceptron. The Novikoff theorem, which provides the necessary conditions for the convergence of a perceptron learning algorithm within a finite number of iterations, has gained increased recognition. Additional research was conducted on the architecture and learning capabilities of multilayer neural networks. In the year 1980, Kunihiko Fukushima introduced a hierarchical multilayered convolutional neural network, which is commonly referred to as the neocognitron developed by Fukushima.

The creation of the back propagation learning algorithm in the mid-eighties had a notable impact, as demonstrated by various authors. Rumelhart's theories, which were initially offered in the early 1960s, have been discussed by Dreyfus, Widrow, Lehr and Werbos. In contrast to the conventional gradient descent method, which updates all parameters simultaneously based on the error, the back propagation algorithm follows a different approach. It first propagates the error term from the output layer back to the layer where parameter updates are required. Subsequently, it employs the chain rule to update the parameters in accordance with the propagated error. Several limitations of the back propagation algorithm were identified in the studies conducted by Brady et al., as well as Gori and Tesi. This approach may encounter limitations when attempting to classify instances that are not linearly separable within the given classes.

## Notes

### The Origins and History of Machine Learning

**1950** – In the year under consideration, Alan Turing, a highly esteemed and famous mathematician and computer scientist from Britain, devised the Turing test. The purpose of the test was to ascertain whether a computer possesses intelligence that is comparable to that of a person. To successfully achieve a passing grade on the examination, the computer must possess the capability to persuade a human interlocutor into perceiving it as a fellow human being. With the exception of a computer program that purportedly emulated a 13-year-old Ukrainian boy and reportedly passed the Turing test, no other endeavours have achieved similar success thus far.

**1967** – The nearest neighbour method was initially developed in the current year. This technology enables computers to initiate the utilisation of rudimentary pattern recognition. The technique presented herein offers a solution for the optimisation problem of determining an efficient route for a travelling salesman, commencing from an arbitrarily selected city, while guaranteeing the inclusion of all specified cities within the minimum possible time.

**1981** – The notion of explanation-based learning (EBL) was established by Gerald DeJong. In this particular mode of learning, the computer undertakes an analysis of the training data and thereafter formulates a general rule that may be applied by excluding data that appears to lack significance.

**1985** – The NetTalk program, developed by Terry Sejnowski, possesses the ability to acquire word pronunciation in a manner akin to that observed in infants throughout the process of language learning. The primary objective of the artificial neural network was to create a simpler model that could effectively demonstrate the intricate nature of acquiring cognitive abilities comparable to those exhibited by humans.

**1990s** – In the 1990s, there was a notable transition in the field of machine learning from a knowledge-driven strategy to a data-driven approach. Programs were developed by scientists and researchers to enable computers to analyse vast quantities of data and derive conclusions based on the obtained results. As a result of these advancements, the IBM Deep Blue computer emerged and achieved victory over Garry Kasparov, the reigning world chess champion, in 1997.

**2006** – The phrase “deep learning” was introduced by Geoffrey Hinton in this year. The individual employed the terminology to elucidate a novel category of algorithms that enable computers to visually see and differentiate items or text within photos or videos.

**2010** – In this year, Microsoft introduced the Kinect technology, which possesses the capability to track up to 20 distinct human traits at a frequency of 30 instances per second. The Microsoft Kinect system facilitated user-machine interaction through the utilisation of gestures and movements.

**2011** – The year in review has presented notable developments in the field of machine learning. To begin with, IBM’s Watson successfully outperformed human contestants on the game show Jeopardy. Additionally, Google has created Google Brain, a system that utilises a deep neural network capable of acquiring the ability to identify and classify various things, with a specific focus on cats.

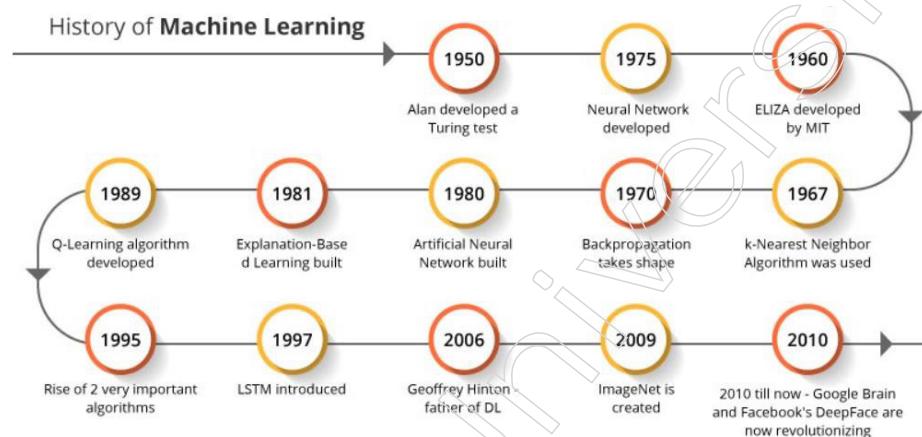
**2012** – The machine learning algorithm, built by Google X lab, possesses the capability to independently browse through YouTube movies and accurately detect the presence of feline subjects.

**2014** – Facebook has developed a software algorithm known as DeepFace, which possesses the capability to identify and authenticate individuals in photographs with a level of accuracy comparable to that of human beings.

**2015** – Amazon introduced its exclusive machine learning platform, so enhancing the accessibility of machine learning and elevating its prominence within the realm of software development. In addition, Microsoft has developed the Distributed Machine Learning Toolkit, a software tool that facilitates the effective distribution of machine learning tasks over several computing devices.

**2016** – At the aforementioned year, Google's artificial intelligence algorithms achieved the remarkable feat of surpassing a skilled human player at the strategic Chinese board game, Go. The game of Go is widely regarded as the most intricate board game in the world. The AlphaGo system, developed by Google, achieved a perfect record of five victories in the competition, thereby garnering significant attention and placing artificial intelligence in the spotlight.

**2020** – The groundbreaking natural language processing algorithm GPT-3 was announced by Open AI. It possesses an exceptional capacity to generate text that closely resembles human language when provided with a prompt. Currently, GPT-3 is widely regarded as the largest and most advanced language model globally, as it utilises 175 billion parameters and Microsoft Azure's AI supercomputer for its training process.



Some important applications in which machine learning is widely used are given below:



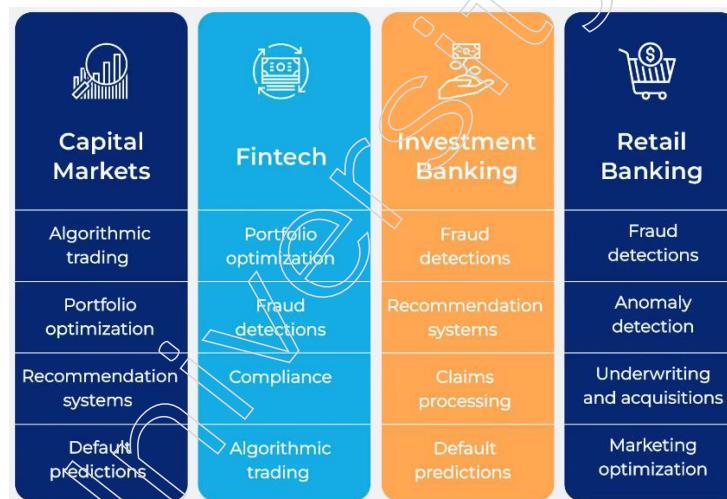
Figure: Machine Learning in Healthcare

**Healthcare:** The healthcare industry extensively utilises machine learning techniques. The analysis of data points and subsequent consequence suggestions are beneficial for healthcare researchers. Natural language processing (NLP) has been essential in providing precise insights to enhance patient outcomes. Moreover, the field of machine learning has made significant advancements in enhancing treatment

## Notes

methodologies through the analysis of external data pertaining to patients' illnesses, such as X-ray images, ultrasound scans and CT scans. Natural Language Processing (NLP), medical imaging and genetic information are prominent domains within the field of machine learning that have demonstrated significant advancements in enhancing diagnostic, detection and predictive systems within the healthcare industry.

**Banking and Finance:** The field of Banking and Finance include the application of Machine Learning, which is a branch of Artificial Intelligence (AI) that utilises statistical models to generate precise predictions. Machine learning has proven to be highly beneficial in the banking and finance industry, offering several advantages across multiple areas. These include fraud identification, portfolio management, risk management, chatbots, document analysis, high-frequency trading, mortgage underwriting, anti-money laundering (AML) detection, anomaly detection, risk credit score detection and know your customer (KYC) processing, among others. Therefore, machine learning is extensively utilised in the banking and financial industry with the aim of minimising both errors and time consumption.



**Image Recognition:** Image recognition is a prevalent use of machine learning that is utilised to identify images on the internet. Moreover, a multitude of social media platforms, like Facebook, employ picture recognition technology to facilitate the process of tagging individuals in uploaded images. This functionality is commonly referred to as "auto friend tagging suggestion."

Moreover, in contemporary times, nearly all mobile gadgets are equipped with captivating face detection capabilities. This feature allows users to enhance the security of their mobile data by implementing face unlocking. Consequently, unauthorised individuals attempting to access the mobile device will be unable to do so without successful face recognition.

**Speech Recognition:** Speech recognition is considered to be a significant accomplishment within the realm of machine learning applications. The feature allows users to do searches by utilising voice commands, hence eliminating the need for manual text input. By utilising voice recognition technology, this system has the capability to do searches for information and items across several platforms like as YouTube, Google, Amazon and others. The aforementioned technology is commonly known as voice recognition.

The technique described involves the conversion of spoken commands into written text, sometimes referred to as 'Speech to text' or 'Computer speech recognition'. Several

notable instances of speech recognition systems include Google Assistant, Siri, Cortana and Alexa, among others.

**Product Recommendation:** Machine learning has emerged as a significant accomplishment, facilitating digital advertising for e-commerce and entertainment enterprises such as Flipkart, Amazon, Netflix, among others. When individuals conduct a search for a particular product, they then encounter targeted advertisements for that same product while engaging in internet browsing activities on the same web browser.

This capability is facilitated through the utilisation of machine learning algorithms, which analyse users' interests and past experiences to provide tailored product recommendations. For example, when conducting a search for a laptop on the Amazon platform, numerous other computers with similar categories and criteria are also displayed. Likewise, when utilising the Netflix platform, users are presented with personalised suggestions for a variety of entertainment options, including series, movies and more. Therefore, the achievement of this outcome can also be facilitated through the utilisation of machine learning techniques.

**Self-driving cars:** Self-driving automobiles provide a very captivating implementation of machine learning techniques. Machine learning is an integral component in the production of autonomous vehicles. The system employs an unsupervised learning approach to train vehicle models in order to effectively identify and classify pedestrians and various objects during the course of driving. Tata and Tesla are widely recognised automotive manufacturers that have garnered significant attention for their respective endeavours in the development of autonomous vehicles. Therefore, this significant transformation occurs within the context of the technology era, facilitated by the utilisation of machine learning techniques.



**Credit card fraud detection:** The identification of credit card fraud has emerged as a significant concern due to its vulnerability to exploitation by online hackers. The prevalence of online and digital payment systems has led to a concomitant rise in the associated risks pertaining to credit and debit card usage. Machine Learning is a valuable tool for developers in the identification and examination of fraudulent activities within online transactions. This study proposes a unique approach for detecting fraud in Streaming Transaction Data. The primary aim is to analyse historical transaction information of consumers and identify their behavioural patterns. Moreover, cardholders are categorised into different groups based on their transaction amounts, allowing for the extraction of distinct behavioural patterns within each group. Therefore, the utilisation of the Aggregation Strategy and Feedback Mechanism in machine learning presents a pioneering method for detecting credit card fraud.

## Notes

**Stock Marketing and Trading:** Machine learning plays a significant role in the domain of stock marketing and trading by leveraging historical trends and past experiences to forecast market risks. Share marketing, also known as marketing risk, can be mitigated to a certain degree by the application of machine learning techniques. These techniques enable the prediction of data pertaining to marketing risk. The long short-term neural memory network (LSTM) within the field of machine learning is employed for the purpose of forecasting stock market patterns.

The process of converting text from one language to another is known as language translation. This practice is commonly used to facilitate communication between individuals.

**Data analysis:** In contemporary times, organisations are producing an unprecedented amount of data, necessitating the utilisation of machine learning algorithms to effectively analyse and interpret this data. Machine learning has the potential to expedite the process of extracting the most useful insights from data sets by assuming the labor-intensive task of comprehensively examining the entire dataset. Machine learning algorithms have been developed to support managerial decision-making processes and enhance the efficiency of teams across several departments, including sales, marketing and production, by facilitating faster data analysis.

**Identification and Prevention of Fraudulent Activities:** The increasing reliance of consumers on internet platforms for buying has provided cybercriminals with several opportunities to engage in fraudulent activities. Numerous organisations implement a variety of online security measures, with machine learning emerging as the most promising approach. Machine learning algorithms are employed to discern fraudulent transactions, such as money laundering, from lawful ones. Machine learning methods facilitate the analysis of distinct attributes within a given dataset, enabling the construction of a robust model that serves as a foundation for scrutinising individual transactions for potential indications of fraudulence. In this manner, organisations have the ability to intervene in the process prior to its completion, thereby mitigating potential complications.

**Natural language processing (NLP):** Natural Language Processing (NLP) has enabled the resolution of various tasks such as tech assistance, help desks, customer care and numerous others through the use of machine learning algorithms and their proficiency in processing natural language. The potential for computers to supplant human agents arises from the capabilities of Natural Language Processing (NLP), which facilitates automated translation between computer and human languages. Machine learning-driven technologies such as chatbots and virtual assistants prioritise various aspects of human language, including context, jargon, meaning and other intricate details, in order to enhance their human-like conversational abilities.

### 1.1.2 Role of Data and Connection to the Knowledge/Experience in Learning

Background The field of Machine Learning pertains to the inquiry of constructing a computer system that autonomously enhances its performance and acquires knowledge through experiential learning. Currently, machine learning is experiencing significant growth within the realm of technology, since it resides at the convergence of computer science and statistics, serving as a fundamental component of Artificial Intelligence. The misinterpretation of Machine Learning as Artificial Intelligence (AI) is a prevalent occurrence. Artificial Intelligence (AI) is a discipline that falls within the expansive realm of computer science, encompassing both contemporary advancements and longstanding

principles.

The overarching domain of artificial intelligence endeavours to comprehend intelligent creatures and replicate human intelligence within a computational system. The pursuit of intelligence is evident not just in philosophy and psychology but also in the field of computer science, particularly in the subject of artificial intelligence (AI). In this domain, the primary objective is to construct intelligent entities while simultaneously seeking to comprehend their underlying mechanisms.

The acquisition of knowledge is the fundamental principle underlying human cognitive abilities. Machine learning is often regarded as a means to attain intelligent machines through their capacity to acquire knowledge and adjust their behaviour in response to their surroundings. The field of machine learning has had significant progress in recent years, which may be attributed in part to the emergence of vast quantities of data being collected.

Data analytics is widely recognised as a crucial function inside organisations, as it enables the extraction of valuable insights pertaining to business value and customer behaviour from the vast amount of data that organisations gather or generate. Despite the promising prospects offered by data, there are inherent repercussions associated with the proliferation of big and diverse datasets originating from many sources. One of the primary challenges in data analytics lies in the capacity to effectively analyse the vast scale and intricate nature of complicated data sets. Machine learning is a method that can provide assistance in carrying out these duties, since it serves as an enhancement to the cognitive capabilities of individuals.

Large data sets bring both challenges and opportunities for machine learning systems. This is because such data helps algorithms identify intricate patterns, enabling more precise and fast predictions than before. AI and Machine Learning (ML) are emerging technologies that will dramatically impact the digital business landscape in the future decade.

AI and ML are disruptive because to sophisticated algorithms, large datasets for algorithmic training, and parallel computing advances. AI and ML include artificial neural networks, deep learning, and natural language processing. These traits enable the creation of advanced systems that can understand, learn, predict, adapt, and maybe operate autonomously.

Companies can use data more efficiently to solve problems with learning systems. Learning systems may solve problems and help companies create new products and services. Consumer products giant Wise Inc. operates worldwide. Wise Inc. prioritises digital technology integration to implement the 2020 strategic business plan.

The organization's meaning of "digital first" is unclear. Since there is no framework, each department must determine the meaning and ramifications of a digital-first approach in their setting. They must also determine how they can help implement the company's 2020 strategic business strategy.

Big data's volume, diversity, and velocity characteristics explain data management's complexity. Social media, ecommerce, sensors, and computing gadgets generate unprecedented amounts of data. This rapid data growth complicates data management.

Different technologies store or generate different amounts of data, which is called "volume". Data volume varies according to temporal fluctuations and data type. Variety means different data types. Technology phenomena like social media platforms have generated new types of data, increasing data diversity. Structured, semi-structured, and unstructured data should be considered.

## Notes

Standardised data in spreadsheets or relational databases with labels is called structured data. Semi-structured data meets application structural requirements but deviates from relational database standardisation. XML is a popular standard for hierarchical data organisation. It lets users exchange data includes language, which is used to exchange data on the Web.

Unstructured data, which includes text, photos, audio, and video, is difficult to analyse. Speed is how fast data is generated, processed, and evaluated. Due to the time and cost of data processing, batch-processing has been the main approach for analysis. Smartphones and other new computing gadgets have increased data creation. Real-time analytics are needed due to the increasing growth of data.

The utilisation of data is an essential element within the domain of Machine Learning. The term “training data” pertains to the collection of observations or measurements that are utilised in the process of training a machine-learning model. The performance of a machine-learning model is heavily influenced by the quality and quantity of data that is accessible for both training and testing purposes. Data can manifest in several formats, including numerical, categorical, or time-series data and can originate from a multitude of sources, such as databases, spreadsheets, or application programming interfaces (APIs). Machine learning algorithms leverage data to acquire knowledge of patterns and correlations among input variables and target outputs, enabling them to make predictions or perform classification tasks.

Data is typically divided into two types:

- Labeled data
- Unlabeled data

Labelled data comprises a label or target variable that the model endeavours to predict, whereas unlabeled data lacks a label or target variable. In the field of machine learning, the data utilised is commonly classified into two main types: numerical and categorical. Numerical data include quantifiable values that may be systematically arranged and assessed, such as age or salary. Categorical data encompasses discrete values that correspond to distinct categories, such as gender or the specific variety of fruit.

Data can be partitioned into two distinct sets, namely the training set and the testing set. The training dataset is utilised for the purpose of training the model, whereas the testing dataset is employed to assess the model’s performance. Ensuring that the data is divided in a manner that is both random and representative has significant importance.

Data preparation plays a crucial role in the machine learning pipeline. This stage encompasses many tasks such as data cleansing and normalisation, addressing missing values and doing feature selection or engineering.

### **Data:**

An unprocessed fact, value, text, sound, or picture that lacks interpretation and analysis might be considered as raw data. The significance of data cannot be overstated in the domains of Data Analytics, Machine Learning and Artificial Intelligence. The absence of data renders the training of any model impossible, hence rendering futile all contemporary research and automation endeavours. Large corporations are investing significant financial resources in order to get a substantial amount of specific data.

The data has undergone interpretation and manipulation, resulting in relevant inferences for the consumers.

**Knowledge:**

The amalgamation of deduced knowledge, personal encounters, educational pursuits and discerning perspectives. The outcomes involve the development of consciousness or the construction of conceptual frameworks for individuals or organisations.



Figure: Data connection Process

**How do we split data in Machine Learning:**

**Training Data:** The portion of data utilised for training our model. The provided information constitutes the dataset that the model actively perceives and acquires knowledge from.

**Validation Data:** The portion of data utilised for regular model evaluation involves fitting the training dataset and optimising the associated hyper parameters, which are predetermined parameters prior to the model's learning process. This data is important during the training phase of the model.

**Testing Data:** After the completion of training, the utilisation of testing data allows for an impartial assessment of our model. When the inputs of the Testing data are provided, our model will provide predictions for certain values without having access to the actual outcome. Following the process of prediction, the evaluation of our model is conducted by a comparison between the model's output and the actual output contained inside the testing data. This is the methodology employed to assess the extent to which our model has acquired knowledge from the input data provided during the training phase.

Training data refers to the dataset employed for the purpose of training an algorithm or machine learning model, with the objective of enabling the model to make predictions pertaining to the desired outcome it has been designed to forecast. When employing supervised learning or a hybrid methodology that incorporates supervised learning, the dataset will undergo a process of data labelling or annotation to enhance its quality and usefulness.

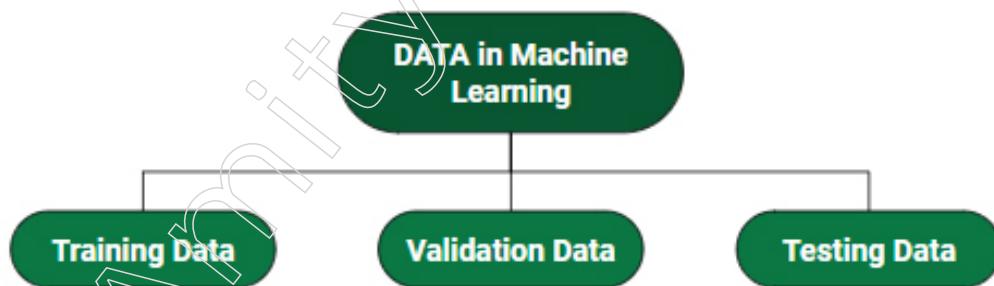


Figure: Data in Machine learning

The utilisation of test data is essential for evaluating the performance metrics, such as accuracy and efficiency, of the algorithm employed in machine training. The utilisation of test data facilitates the evaluation of a model's predictive capabilities for novel responses, drawing upon its training. Both the training and test datasets play a crucial role in enhancing and validating machine learning models.

## Notes

Training data refers to the dataset utilised for training a machine learning algorithm or model with the objective of achieving accurate predictions for a specific outcome or desired answer. In the context of supervised learning, the process of training data necessitates the involvement of a human agent who actively selects and assigns labels to the pertinent features within the dataset. These labelled features are subsequently utilised to facilitate the training of the machine learning model.

Unsupervised learning is a machine learning approach that leverages unlabeled data to identify patterns, such as making conclusions or clustering data points. Semi-supervised learning is a hybrid approach that integrates both supervised and unsupervised learning techniques.

### 1.1.3 Show Successful Examples of Machine Learning in Industry/ Working

The five primary uses of machine learning encompass fraud detection, virtual personal assistants, product suggestions, speech recognition and consumer segmentation. Machine learning has proven to be a valuable tool in the realm of manufacturing, as it can be effectively employed for many purposes such as quality control, automation and customisation. One potential use of machine learning involves the utilisation of this technology to identify and diagnose flaws or imperfections in goods prior to their distribution to end-users. Moreover, it has the capability to streamline monotonous activities, such as those encountered in assembly line operations.

Machine learning, a contemporary technological advancement, has significantly improved several industrial and professional procedures, as well as our everyday existence. Machine learning is a specialised field within the realm of artificial intelligence (AI) that centres on the utilisation of statistical methods to develop intelligent computer systems capable of learning from existing databases.

Machine learning enables computer systems to effectively leverage client data for various purposes. The system functions based on its pre-established programming and is capable of adapting to novel circumstances or modifications. Algorithms possess the capability to dynamically adjust their operations based on the input data, hence acquiring behaviours that were not explicitly pre-programmed.

The acquisition of reading skills and the ability to discern context enables a digital assistant to effectively analyse emails and extract pertinent information. An essential aspect of this learning process involves the capacity to generate forecasts on prospective client behaviours. This approach facilitates a deeper understanding of clients, enabling businesses to move beyond mere responsiveness and adopt a proactive stance. Deep learning is a subfield within the broader domain of machine learning. Essentially, it is a multi-layered artificial neural network. Single-layer neural networks have the capability to generate approximate predictions. The use of additional layers has the potential to enhance optimisation and improve accuracy. Machine learning holds significant relevance across several sectors and industries, exhibiting the potential for continuous expansion and development. Presented here are six empirical instances that exemplify the practical applications of machine learning.

**Image recognition:** Image recognition is a widely recognised and prevalent instance of machine learning in practical applications. The process involves the identification of an object within a digital image, utilising the pixel intensity present in either monochromatic or colour images.

**Examples:** Classify an x-ray image as either indicative of cancer or non-cancerous.

In the realm of social media, the practice of assigning a name to a photographed face, commonly referred to as “tagging,” is a prevalent phenomenon. The proposed approach involves the segmentation of a single letter into smaller pictures in order to facilitate the recognition of handwriting.

Facial recognition inside an image is a common use of machine learning. By utilising a database containing information about individuals, the system is capable of identifying shared characteristics and subsequently associating them with corresponding facial features. This is frequently employed within the realm of law enforcement.

**Speech recognition:** Machine learning algorithms have the capability to convert spoken language into written text. Several software systems have the capability to transform real-time spoken language and pre-recorded speech into a written document. The speech can also be divided into segments based on intensities observed in time-frequency bands.

**Medical diagnosis:** The utilisation of machine learning techniques has the potential to contribute significantly to the process of disease diagnosis. Numerous medical practitioners employ chatbots equipped with speech recognition skills to identify and analyse patterns in symptoms.

### Examples

- Aiding in the development of a diagnostic plan or proposing a course of therapy
- The fields of oncology and pathology employ machine learning techniques to discern malignant tissue.
- This study aims to conduct an analysis of various physiological fluids. The utilisation of facial recognition software and machine learning in the context of rare diseases facilitates the examination of patient photographs and the identification of traits that exhibit a correlation with uncommon genetic disorders.

**Statistical arbitrage:** Arbitrage is a computerised trading approach employed in the field of finance with the purpose of effectively managing a substantial quantity of financial instruments. The employed approach utilises a trading algorithm to assess a collection of assets by means of economic indicators and correlations.

### Example

- Trading algorithms that examine the microstructure of a market
- Analyse a lot of data
- Find prospects for real-time arbitrage
- The arbitrage approach is optimised via machine learning to improve outcomes.

**Predictive Analytics:** The data that is accessible can be classified using machine learning techniques into distinct categories, which can afterwards be refined and delineated by rules developed by analysts. Once the classification process has been completed, the analysts possess the ability to ascertain the probability of a problem occurring.

### Examples

- Determining whether a transaction is legitimate or fraudulent.
- Enhance prediction methods to determine the likelihood of a defect.
- One of the most promising applications of machine learning is predictive analytics. It can be used for anything, including pricing for real estate and product development.

## Notes

**Extraction:** From unstructured data, machine learning can extract structured information. Organisations gather enormous amounts of client data. The process of automatically annotating datasets for predictive analytics tools is done using a machine learning algorithm.

### 1.1.4 Motivation and Role of Machine Learning in Computer Science and Problem Solving

#### Manual data entry

The presence of inaccurate and duplicated data poses significant challenges for organisations seeking to automate their processes. However, the utilisation of machine learning (ML) algorithms and predictive modelling algorithms can greatly alleviate this issue. ML programs leverage the acquired data to enhance the process with each subsequent calculation, enabling machines to acquire the capability to execute laborious tasks such as documentation and data entry. Consequently, knowledge workers can allocate more of their time towards engaging in higher-value problem-solving activities. Arria, an artificial intelligence (AI) firm, has developed a natural language processing technology that examines textual content and discerns the connections between concepts in order to generate reports.

#### Detecting Spam

The problem of spam identification was one of the initial challenges addressed by machine learning. Four years prior, email service providers employed pre-established rule-based methodologies for the purpose of spam elimination. However, the current state of spam filters involves the utilisation of machine learning algorithms to autonomously generate new rules. Due to the implementation of neural networks within its spam filters, Google currently maintains a spam rate of 0.1 percent, a notable achievement. The spam filters of the system incorporate brain-like “neural networks” that possess the ability to acquire knowledge in order to identify unsolicited and fraudulent emails. This is achieved by the analysis of patterns and regulations across a vast network of interconnected computers. In conjunction with spam detection, social media platforms are employing machine learning techniques to discern and mitigate instances of abusive content.

#### Financial Analysis

Machine learning has the potential to be employed in financial research due to its ability to handle vast amounts of data, its quantitative character and its capacity to provide precise historical information. The utilisation of machine learning (ML) in the field of finance encompasses various applications, such as algorithmic trading, portfolio management, fraud detection and loan underwriting. Potential future uses of machine learning (ML) in the field of finance encompass the utilisation of chatbots and conversational interfaces to enhance customer service, as well as to bolster security measures and conduct sentiment analysis.

#### Predictive Maintenance

The utilisation of artificial intelligence (AI) and machine learning (ML) holds potential for the manufacturing industry to identify significant trends within industrial data. The implementation of corrective and preventive maintenance practises incurs significant expenses and exhibits inefficiencies. Predictive maintenance serves to mitigate the likelihood of unforeseen problems and diminishes the extent of superfluous preventive maintenance endeavours.

## Corrective, Preventive and Predictive Maintenance

In the context of predictive maintenance, it is possible to construct a machine learning architecture that incorporates many components. These components include historical device data, a flexible analytical environment, a workflow visualisation tool and an operations feedback loop. The Azure Machine Learning platform offers a demonstration of predictive maintenance modelling by utilising a set of simulated aircraft engine run-to-failure events.

It is postulated that the asset exhibits a pattern of degradation that is characterised by progression through time. The aforementioned pattern is evident in the measurement of the sensor of the asset. To forecast forthcoming failures, machine learning algorithms acquire knowledge about the correlation between sensor readings and the variations in sensor values with respect to past occurrences of failures.

## Image Recognition

Computer vision is a field of study that involves the extraction of numerical or symbolic data from photographs and high-dimensional datasets. The subject matter encompasses machine learning, data mining, knowledge discovery in databases and pattern recognition. The application of picture recognition technology in several industries such as healthcare, autonomous cars and marketing campaigns presents promising opportunities for businesses. Baidu has successfully created a prototype known as DuLight, designed to assist individuals with visual impairments. This innovative device utilises computer vision technology to see and analyse the surrounding environment, subsequently providing auditory descriptions and explanations through an earpiece. Marketing efforts that utilise image recognition technology, such as L'Oreal's Makeup Genius, have been found to effectively enhance social sharing and user engagement.

Machine Learning (ML) plays a significant role in computer science and problem-solving, transforming the way we approach complex tasks and challenges. Here's how ML contributes to these areas:

- **Automation and Efficiency:** ML algorithms can automate repetitive tasks and processes that are time-consuming for humans. This leads to increased efficiency and frees up human resources for more creative and strategic work.
- **Pattern Recognition:** ML excels at identifying patterns in large datasets that might be difficult or impossible for humans to discern. This ability to recognize patterns is used in image and speech recognition, natural language processing, and fraud detection, among others.
- **Predictive Analytics:** ML models can make predictions based on historical data. This has applications in financial forecasting, weather prediction, stock market analysis, and more.
- **Classification and Categorization:** ML algorithms can classify and categorize data into distinct classes or groups. This is used in applications like email spam detection, sentiment analysis, and medical diagnosis.
- **Recommendation Systems:** ML powers recommendation engines that suggest products, services, or content based on user preferences. This is commonly seen in streaming services, online shopping platforms, and social media.
- **Natural Language Processing (NLP):** ML techniques enable computers to understand, interpret, and generate human language. This is used in chatbots, language translation, sentiment analysis, and more.

## Notes

- **Computer Vision:** ML algorithms enable computers to interpret and understand visual information from the world. This has applications in facial recognition, autonomous vehicles, object detection, and medical imaging.
- **Anomaly Detection:** ML can identify unusual or anomalous behavior in data, making it useful for fraud detection, network security, and fault detection in industrial systems.
- **Optimization:** ML can find optimal solutions to complex optimization problems. This is used in supply chain management, resource allocation, and scheduling.
- **Personalization:** ML helps tailor user experiences by analyzing user behavior and preferences. This leads to personalized content, recommendations, and marketing strategies.

### 1.1.5 Connect Machine Learning to the Broader Theme of Computer Science

Machine learning (ML) is a subfield of artificial intelligence (AI) that facilitates the autonomous acquisition of knowledge by computers through the analysis of training data, leading to iterative enhancements in performance, all without the need for explicit programming. Machine learning algorithms have the capability to identify patterns within datasets and acquire knowledge from these patterns, hence enabling them to generate predictions autonomously.

Machine Learning is a prominent and disruptive topic within the discipline of Computer Science, which encompasses a diverse set of techniques and algorithms that enable computers to acquire knowledge from data and enhance their performance on a given task without the need for explicit programming. It connects with other domains within the field of Computer Science, exerting a significant impact on their development and direction. Machine Learning is intricately linked to the overarching domain of Computer Science due to its fundamental principles and applications within this field.

**Artificial Intelligence (AI):** Machine Learning is a fundamental element of Artificial Intelligence (AI), facilitating the replication of human-like learning and decision-making mechanisms by computers. Artificial intelligence (AI) comprises a diverse range of methodologies, such as Machine Learning, Natural Language Processing (NLP), Computer Vision, Robotics and other related approaches. Machine Learning algorithms serve as the fundamental framework for numerous artificial intelligence (AI) applications and systems.

**Data Science:** Machine learning plays a key role in the field of data science. Data Science encompasses the process of collecting valuable insights and knowledge from extensive datasets through the use of several methodologies, including data gathering, data cleansing, data analysis, data visualisation and predictive modelling approaches. Machine Learning models are employed for the purpose of constructing predictive and descriptive models based on data.

**Pattern Recognition:** Machine Learning algorithms are specifically engineered to identify and comprehend patterns and correlations present in datasets. This statement is consistent with the domain of Pattern Recognition in the discipline of Computer Science, which is concerned with the detection of consistent patterns in data and the subsequent use of these patterns for decision-making purposes.

**Computational Statistics:** Machine Learning draws extensively from principles and methodologies in the field of Statistics. The utilisation of techniques including regression, classification, clustering and hypothesis testing is considered vital in both

areas. Machine Learning algorithms frequently utilise statistical concepts to create predictions and draw inferences from data.

**Algorithm Design and Analysis:** Machine Learning encompasses the process of developing and refining algorithms with the objective of acquiring knowledge from data in a manner that is both efficient and effective. This observation is relevant to the overarching topic of algorithm design and analysis within the field of Computer Science, whereby scholars investigate the effectiveness, accuracy and capacity for growth of algorithms.

**Computer Vision:** Computer Vision is a specialised domain within the discipline of Computer Science that is dedicated to facilitating the capacity of computers to analyse and comprehend visual data derived from photos or videos. The field of Computer Vision has been significantly transformed by the advent of Machine Learning methods, with deep learning in particular playing a pivotal role. This has led to groundbreaking advancements in tasks such as image recognition, object detection and image segmentation.

**Natural Language Processing (NLP):** Natural Language Processing (NLP) pertains to the study and analysis of the dynamic interplay between computer systems and human language. Machine Learning is of utmost importance in Natural Language Processing (NLP) applications such as sentiment analysis, language translation and speech recognition. In these jobs, computers acquire knowledge from extensive text and speech datasets by identifying and understanding patterns.

**Reinforcement Learning:** The field of Natural Language Processing (NLP) encompasses the examination and evaluation of the dynamic interaction between computer systems and human language. Machine Learning plays a crucial role in the realm of Natural Language Processing (NLP) applications, encompassing tasks such as sentiment analysis, language translation and audio recognition. In these occupations, computers assimilate information through the analysis of vast text and speech databases, wherein they discern and comprehend patterns.

**Database Systems:** Machine Learning models frequently depend on extensive datasets for the purpose of training. This topic pertains to the intersection of database design and management within the field of Computer Science, where the optimisation of data storage, retrieval and querying processes holds significant importance.

**Ethical and Social Implications:** Ethical considerations are becoming more crucial as AI and machine learning develop. The ethical and social ramifications of AI and machine learning, including questions of bias, justice, transparency and privacy, are studied in computer science. Overall, Machine Learning is a key factor in the development of numerous ground-breaking computer science applications and improvements. Its linkages to numerous subfields serve as an example of how it enhances and has an impact on the larger field of computing and technology.

## 1.2 Types of Machine Learning

The ability for a machine to automatically learn from data, enhance performance based on prior experiences and make predictions is known as machine learning. A collection of algorithms used in machine learning operate on vast amounts of data. These algorithms are fed data to train them and after training, they develop a model and carry out a certain task.

### 1.2.1 Supervised Machine Learning

Typically, supervised learning starts with a pre-existing set of data and a

## Notes

predetermined comprehension of how that data is categorised. The goal of supervised learning is to identify data patterns that can be used in an analytics process. The labelled features in this data define the data's meaning. For instance, you could have millions of photographs of animals with descriptions of each one and you could then develop a machine learning program that can tell one animal from another.

You may create hundreds of categories of various species by categorising this information about animal types. The users who are training the modelled data to meet the specifics of the labels have a good understanding of the data because the attributes and meaning have been determined. Regression occurs when the label is continuous; classification occurs when the data comes from a finite collection of values. Regression applied for supervised learning essentially aids in your understanding of the relationship between variables. Forecasting the weather is an illustration of supervised learning. Regression analysis is a technique used in weather forecasting to produce a prediction of the weather by taking into consideration both the current conditions and known historical weather patterns.

The learning process of a basic machine learning model consists of two distinct steps, namely training and testing. The learning model is developed through the training process, wherein samples from the training data are utilised as input. During this process, the learning algorithm or learner acquires knowledge of the features. During the process of testing, the learning model utilises the execution engine to generate predictions for the test or production data. The output of the learning model, known as tagged data, represents the ultimate prediction or classified information.

Given the objective of inducing the computer to acquire a constructed classification system, supervised learning (as depicted in the upper figure) is the prevailing technique employed in addressing classification challenges.

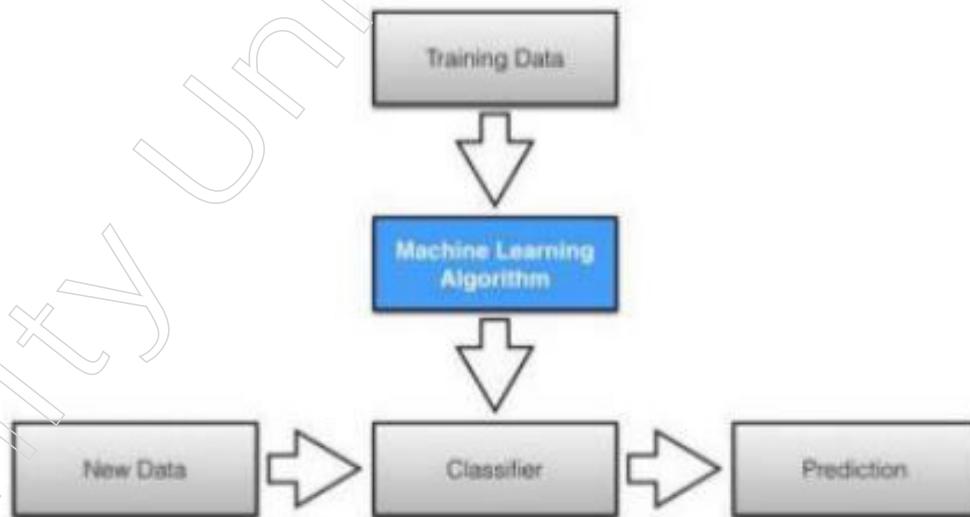


Figure: Supervised learning process

The probability associated with input in supervised learning is often not explicitly given, even in cases where the anticipated outcome is already known. The aforementioned technique facilitates the creation of a dataset with distinct features and corresponding labels. The primary objective is to develop an estimator capable of predicting the label of an object using a given collection of features.

The learning algorithm is provided with a set of features and their corresponding desired outputs as input. It proceeds to learn by comparing the desired outputs with the

actual outputs, thereby identifying any discrepancies or defects. Subsequently, the model is modified in accordance with the necessary adjustments. In the presence of all input variables, the utilisation of a model is unnecessary. However, in instances where certain input values are missing, it is not possible to make any inferences on the corresponding outputs.

Supervised learning is widely recognised as the predominant approach for training neural networks and decision trees. The relationship between these two variables is contingent upon the data that is supplied through the pre-established classification process.

Moreover, this acquired knowledge is utilised in algorithms that predict potential future occurrences based on past data. There are a plethora of practical applications for this form of learning, including an application capable of discerning the species of an iris based on measurements of its floral structure.

Supervised learning problems can be categorised into two distinct groups: classification and regression. In the context of regression, the label is characterised by a continuous variable, but in classification, the label is represented by a discrete variable.

As illustrated previously, the algorithm differentiates between the observed data  $X$  and the training data, which often refers to structured data that is supplied to the model during the training process. The construction of the prediction model is achieved by employing the supervised learning technique during the entirety of this process. The trained model will then proceed to predict the most likely labels for a new batch of samples  $X$  in the testing set. Different kinds of supervised learning can be established by classifying the type of goal variable  $y$ .

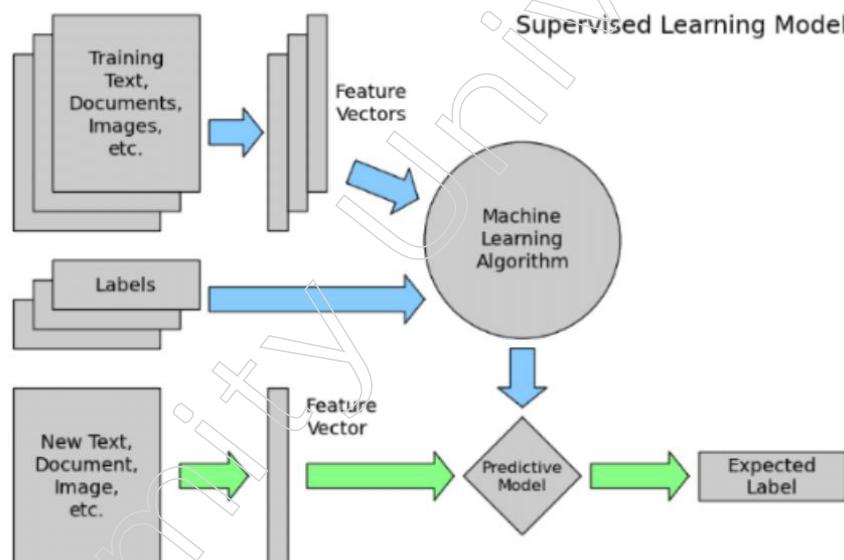


Figure: Supervised Learning Model

- The task of predicting  $y$  is known as classification if it contains values that fall into a predetermined set of category outcomes (integers).
- Regression is the process of predicting  $y$  when  $y$  has floating point values.
- Categories of Supervised Machine Learning

The challenges associated with supervised machine learning can be categorised into two distinct groups, as outlined below:

## Notes

- ❖ Classification
- ❖ Regression

### Classification

Classification methods are utilised to address classification problems that involve a categorical output variable, such as "Yes" or "No," Male or Female, Red or Blue and so on. The classification algorithms are responsible for predicting the categories that are present in the dataset. The contemporary utilisation of categorization systems encompasses several applications such as spam detection, email filtering and other similar instances.

### Regression

Regression problems characterised by a linear association between the input and output variables are addressed through the application of regression methodologies. These models are utilised for the purpose of predicting variables that have continuous outputs, such as market movements, weather forecasts and so on.

### Advantages Supervised Learning

- The utilisation of a labelled dataset in supervised learning enables accurate identification of object types.
- These algorithms become valuable in creating predictions regarding outcomes by leveraging historical performance.

### Disadvantages

- These algorithms are incapable of solving complex issues.
- If there is a discrepancy between the test data and the training data, it could potentially lead to the prediction of an inaccurate outcome.
- The user's text lacks sufficient information to be rewritten in an academic manner. The computational aspect of the algorithm's training process is characterised by a significant duration.

### Applications of Supervised Learning

Following are a few typical applications of supervised learning:

- **Image Segmentation:** Algorithms based on Supervised Learning are utilised for image segmentation. With the use of pre-established labels, picture classification is carried out in this method on various image data.
- **Medical Diagnosis:** The employment of supervised algorithms in the medical field is very common. It is done utilising historical data with labels for disease conditions and medical photos. The machine can diagnose a disease for new patients using such a procedure.
- **Fraud Detection:** Classification algorithms developed through supervised learning are used to find fraudulent transactions, fraudulent clients, etc. In order to find the patterns that could point to potential fraud, historical data is used.
- **Spam detection:** Classification algorithms are employed in spam detection and filtering. These algorithms determine whether an email is spam or not. The spam folder receives the spam emails.
- **Speech Recognition:** Algorithms for supervised learning are also applied in this field. Voice data was used to train the algorithm and it may be used to identify a

variety of things, including voice-activated passwords, speech commands, etc.

### Supervised learning algorithms

#### Decision Tree

The partition of the instance space by a decision tree represents a classifier. The decision tree is made up of nodes that make up a “root tree,” which is a distributed tree with a base node that has no outgoing edges.

Every other node has just one incoming edge. The term “internal node” or “test node” refers to the node with outgoing edges. The remaining nodes are referred to as leaves. Each test node divides the instance space into two or more sub-spaces in a decision tree in accordance with a certain discrete function of the input values. In the simplest scenario, each test just takes into account one attribute, dividing the instance space in accordance with the attribute's value. When a property has a numeric value, the condition relates to a range.

Each leaf is given a class that corresponds to the ideal target value. The leaf could store a probability vector that represents the likelihood that the target attribute will have a particular value. According to the results of the tests conducted along the way, the instances are classified by moving them from the base of the tree down to the leaf.

Figure describes a straightforward application of the decision tree. Each node is given a name for the attribute it is testing and each branch is given a label for the attribute's corresponding value.

With the help of this classifier, the analyst is able to anticipate some potential customers' reactions while also comprehending the general behavioural traits of the population of potential customers.

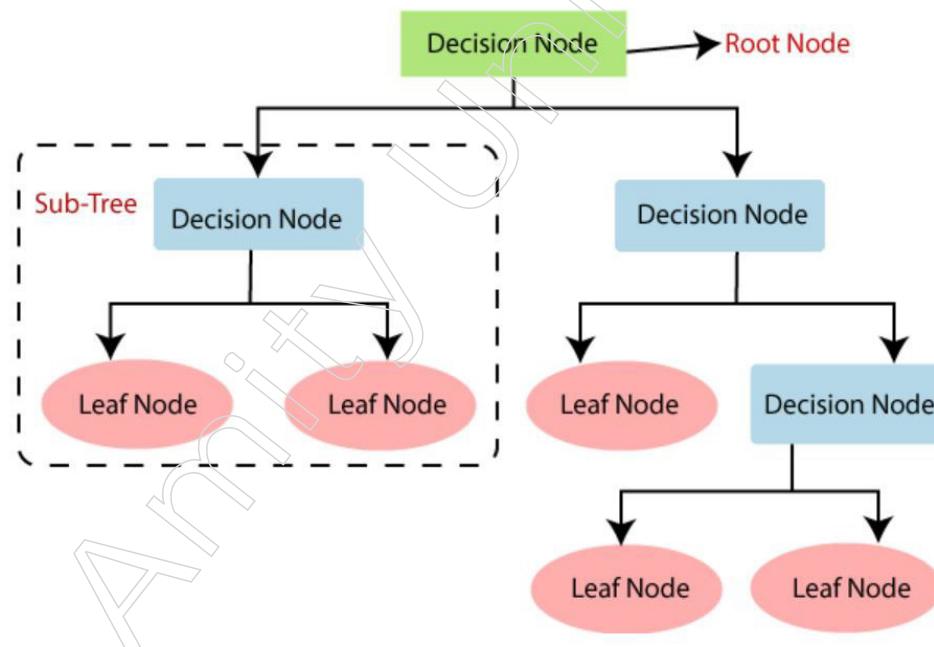


Figure: Decision tree example

Decision trees can be geometrically seen as collections of hyper planes, each orthogonal to one of the axes, in the case of numerical characteristics. Less complicated decision trees are preferred by decision-makers because they may be viewed as being more thorough.

Notes

## Notes

### Linear regression

Finding correlations and dependencies between variables is the aim of the linear regression algorithm, which belongs to the family of regression algorithms. It represents the modelling relationship between one or more (a D-dimensional vector) explanatory variables (also independent variables, input variables, features, observed data, observations, attributes, dimensions, data point, etc.) denoted X and a continuous scalar dependent variable y (also label or target in machine learning terminology). While classification is another field where the objective is to predict a label from a finite set, regression analysis aims to predict a continuous target variable. The multiple regression model that uses a linear combination of input variables has the following structure:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + e$$

Supervised learning algorithms include linear regression among their subcategories. In other words, we use the model to predict labels on sets of unlabeled data (testing data) after training it on a set of labelled data (training data).

The model (red line) is constructed using training data (blue points), where each point has a known label (y axis), to match the points as exactly as possible by minimising the value of a selected loss function, as shown in the above Figure. When we just know the value of x and wish to anticipate the value of y, we can use the model to forecast unknown labels.

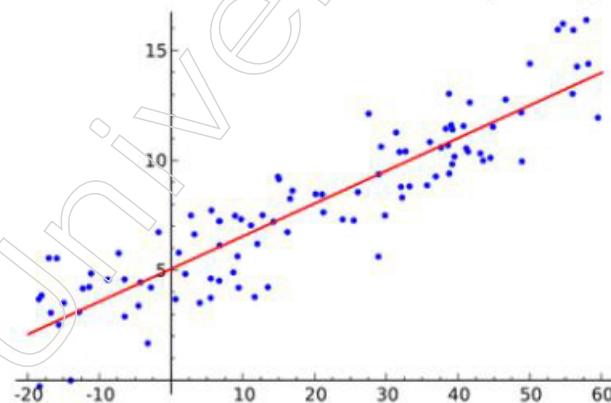


Figure: Visual representation of the linear regression

### Naive Bayes

Another supervised learning technique, as well as the statistical method for classification, is the Bayesian approach. Assumes an underlying probabilistic model, which allows for the principled capture of model uncertainty through the determination of outcome probabilities. The Bayesian classification's main benefit is its ability to address prediction issues.

This categorization can incorporate observable data and offers useful learning techniques. Bayesian classification offers an insightful viewpoint for comprehending and assessing learning systems. It determines the explicit probability for each hypothesis and accounts for input data noise.

Consider the general probability distribution  $P(x_1, x_2)$  for two values. Without sacrificing generality, the Bayes rule yields the following equation:

$$P(x_1, x_2) = P(x_1|x_2)P(x_2)$$

Similar, the following equation is obtained if there is a second class variable:

$$P(x_1, x_2 | c) = P(x_1 | x_2, c)P(x_2 | c)$$

The following results are obtained if the scenario is generalised from the case of two variables to a conditional independence assumption for a set of variables  $x_1, \dots, x_N$  conditional on another variable  $c$ :

$$P(x | c) = \prod_{i=1}^N P(x_i | c)$$

### Logistic Regression

Similar to the naive Bayes model, logistic regression extracts a series of weighted features from the input, converts them to logs and then combines them linearly. This process involves multiplying each feature by a weight and adding the results.

The main distinction between naive Bayes and logistic regression is that the former is a generative classifier, whilst the latter is a discriminative one.

By fitting data to a logistic function, the type of regression known as “logistic regression” can forecast the likelihood that an event will occur. Logistic regression uses a number of predictor variables that may be numerical or categorical, just like many other types of regression analysis.

The definition of the logistic regression theory is:

$$h_\theta(x) = g(\theta^T x)$$

When the definition of the function  $g$  is a sigmoid function

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

As seen in Figure below, the sigmoid function has unique characteristics that lead to values in the range  $[0, 1]$ .

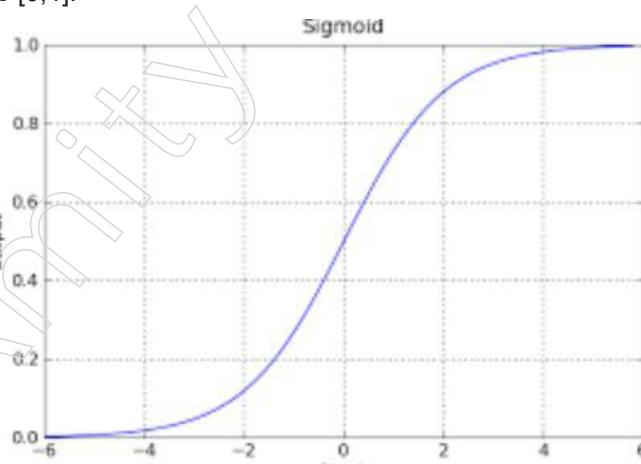


Figure: Visual representation of the Logistic Function

The cost function for logistic regression is given as:

**Notes**

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [-y^{(i)} \log(h_\theta(x^{(i)})) - (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))]$$

We will utilise a built-in function called `fmin_bfgs` in machine learning to determine the minimum of this cost function. This method, given a fixed dataset (of  $x$  and  $y$  values), provides the optimum parameters for the logistic regression cost function. The parameters are the starting values of the parameters that need to be optimised and a function computes the logistic regression cost and gradient with respect to for the dataset with  $x$  and  $y$  values given the training set and a specific. The training data's decision boundary will be visualised using the final value of.

### 1.2.2 Unsupervised Machine Learning

Unsupervised learning is a type of machine learning in which models are not supervised using training datasets, as the name implies. Instead, models themselves decipher the provided data to reveal hidden patterns and insights. It is comparable to the learning process that occurs in the human brain while learning something new. Unsupervised learning can be summed up as a sort of machine learning in which models are taught with unlabeled datasets and then allowed to act on that data unsupervised.

Because unlike supervised learning, we have the input data but no corresponding output data, unsupervised learning cannot be used to solve a regression or classification problem directly. Finding the underlying structure of a dataset, classifying the data into groups based on similarities and representing the dataset in a compressed format are the objectives of unsupervised learning.

**Example:** Consider providing an input dataset including photographs of several breeds of cats and dogs to the unsupervised learning algorithm. The algorithm is never trained on the provided dataset, thus it has no knowledge of its characteristics. The unsupervised learning algorithm's job is to let the image features speak for themselves. This work will be carried out by an unsupervised learning algorithm, which will cluster the image collection into groups based on visual similarities.



Figure: Example of Unsupervised machine Learning

Unsupervised learning is a machine learning approach that is employed to derive insights from datasets comprising input data without accompanying labelled replies. Unsupervised learning algorithms do not incorporate any form of classification or

categorization within the observed data. The absence of output values precludes the estimate of functions.

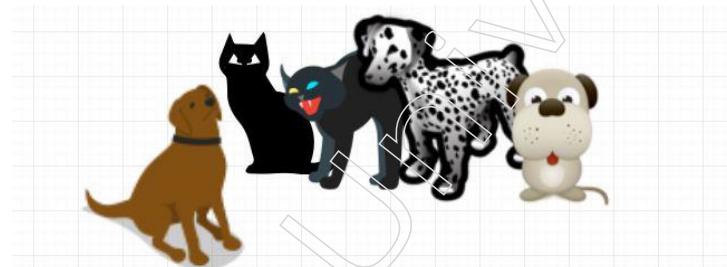
The absence of labels on the examples provided to the learner hinders the ability to assess the accuracy of the algorithm's generated structure. Cluster analysis is widely recognised as the prevailing technique in the realm of unsupervised learning. This method is mostly employed for the purpose of doing exploratory data analysis, with the objective of uncovering latent patterns or establishing groupings within a given dataset.

### Example

Please take into account the following data pertaining to patients who are seeking treatment at a medical facility. The dataset comprises information regarding the gender and age of the patients.

Unsupervised learning refers to the process of training a machine by utilising unclassified and unlabeled information, enabling the algorithm to autonomously process and analyse the data without explicit supervision. The objective of the machine is to categorise unorganised data based on similarities, patterns and distinctions, without any pre-existing data training.

In contrast to supervised learning, unsupervised learning does not involve the provision of a teacher, hence resulting in the absence of training for the machine. Consequently, the machine is limited to autonomously discovering the concealed patterns within unannotated data. For example, consider a scenario in which an image is presented, containing both dogs and cats, which the model has not previously encountered.



Therefore, due to the machine's lack of knowledge on the characteristics of dogs and cats, it is not possible to classify it as 'dogs and cats'. However, it is possible to classify them based on their similarities, patterns and distinctions. For instance, the aforementioned image may be readily divided into two distinct halves. The initial section may encompass a compilation of photographs featuring dogs, whereas the subsequent segment may encompass a compilation of photographs featuring cats. In this particular instance, there was a lack of prior knowledge or exposure to any form of instructional material or illustrative instances.

This feature enables the model to autonomously identify previously unnoticed patterns and information. The primary focus of this study pertains to data that lacks explicit categorization.

Unsupervised learning can be categorised under two distinct algorithmic approaches:

- **Clustering:** When you wish to find the underlying groupings in the data, such as categorising clients based on their shopping habits, you have a clustering problem.
- **Association:** You want to find rules that broadly characterise your data, such as "people who buy X also tend to buy Y," in an association rule learning task.

Unsupervised learning can be further divided into the following two categories:

## Notes

### 1. Clustering

When looking for the innate groups in the data, we employ the clustering technique. It is a method of clustering the items so that those that share the most similarities stay in one group and share little to none in common with those in other groups. Grouping clients based on their purchase habits is an illustration of the clustering method.

The following list of well-liked clustering algorithms includes

- ❖ K-Means Clustering algorithm
- ❖ Mean-shift algorithm
- ❖ DBSCAN Algorithm
- ❖ Principal Component Analysis
- ❖ Independent Component Analysis

### 2. Association

An unsupervised learning method called association rule learning identifies intriguing relationships between variables in a sizable dataset. This learning algorithm's primary goal is to identify the dependencies between data items and then map the variables in a way that maximises profit. This algorithm is mostly used in continuous production, web usage mining, market basket analysis, etc.

Apriori, Eclat and FP-growth algorithms are a few of the well-known algorithms for learning association rules.

#### Advantages

- These algorithms, as opposed to supervised ones, can be utilised for more challenging problems because they operate on unlabeled datasets.
- For a variety of jobs, unsupervised techniques are preferred since it is simpler to obtain the unlabeled dataset than the labelled dataset.

#### Disadvantages:

- Since the dataset is not labelled and the algorithms are not trained using the exact output in advance, the output of an unsupervised method may be less accurate.
- Working with unsupervised learning is more challenging since it uses a dataset that is not mapped to the output and is unlabeled.

#### Applications of Unsupervised Learning

- Network Analysis: In document network analysis of text data for scholarly papers, unsupervised learning is utilised to detect plagiarism and copyright.
- Suggestion Systems: Recommendation systems frequently construct suggestion applications for various online applications and e-commerce websites using unsupervised learning techniques.
- Unsupervised learning is frequently used for anomaly detection, which can find out-of-the-ordinary data items in a collection. It is employed to find erroneous transactions.
- Singular Value Decomposition: To extract specific data from the database, SVD, or singular value decomposition, is utilised. Taking information on each user who is present in a specific location.

### 1.2.3 Reinforcement Learning

By executing actions and observing the outcomes of those actions, an agent learns how to behave in a given environment via reinforcement learning, a feedback-based machine learning technique. The agent receives compliments for each positive activity and is penalised or given negative feedback for each negative action.

In contrast to supervised learning, reinforcement learning uses feedback to autonomously train the agent without the use of labelled data. The agent can only learn from its experience because there is no labelled data. A particular class of problems, such as those in robotics, gaming and other long-term endeavours, are solved using RL. The agent engages with the environment and independently explores it. In reinforcement learning, an agent's main objective is to maximise positive reinforcement while doing well.

The agent learns through hit-and-miss and depending on its experience, it develops the skills necessary to carry out the mission more effectively. In light of this, we may state that "Reinforcement learning is a type of machine learning method where an intelligent agent interacts with the environment and learns to act within that." One example of reinforcement learning is how a robotic dog learns to move his arms.

It is a fundamental component of artificial intelligence and the idea of reinforcement learning is the basis for all AI agents. In this case, there is no need to pre-program the agent because it learns on its own without assistance from humans.

**Example:** Let's say an AI agent is present in a maze setting and his objective is to locate the diamond. The agent interacts with the environment by taking certain actions and depending on those activities, the agent's state is altered and it also receives feedback in the form of rewards or penalties.

The agent keeps performing these three tasks, which help him learn about and investigate his surroundings.

The agent gains knowledge of which behaviours result in positive feedback or rewards and which behaviours result in negative feedback penalties. The agent receives good points for rewards and negative points for penalties.

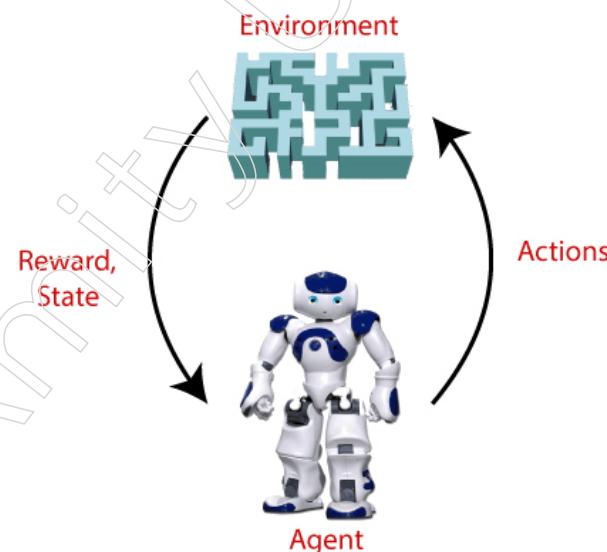


Figure: Reinforcement learning

Reinforcement learning operates through a feedback-driven procedure, wherein an artificial intelligence agent (a software component) autonomously explores its

## Notes

environment through trial and error, executing actions, acquiring knowledge from past encounters and enhancing its overall performance. The reinforcement learning agent is motivated by a system where it receives rewards for doing desirable actions and incurs punishments for undesirable actions.

Consequently, the primary objective of the agent is to optimise its behaviour in order to maximise the cumulative rewards it obtains. In the context of reinforcement learning, it is important to note that unlike supervised learning, the availability of labelled data is absent. Instead, agents rely solely on their experiences to acquire knowledge and improve their performance.

The concept of reinforcement learning pertains to the challenge of enabling an agent to make decisions and take actions within its environment in a manner that optimises the accumulation of rewards. In the context of machine learning, the learner (referring to the program) is not provided with explicit instructions regarding the actions to be taken, which is commonly observed in other machine learning methodologies.

Instead, the learner is required to autonomously explore and experiment with different actions in order to determine which activities result in the highest level of reward. In instances that are particularly intriguing and demanding, actions have the potential to impact not only the current outcome but also subsequent scenarios, so influencing all subsequent rewards.

The process of reinforcement learning exhibits similarities to human learning, as exemplified by a youngster acquiring knowledge through experiential encounters in their daily routines. One instance of reinforcement learning involves engaging in a game, where the game itself serves as the environment. The actions taken by an agent during each iteration establish distinct states, while the objective of the agent is to achieve a favourable score. The agent is provided with feedback in the form of both punitive measures and incentives.

Reinforcement learning is utilised in various domains, including but not limited to game theory, operations research, information theory and multi-agent systems, due to its distinctive operational characteristics. The formalisation of a reinforcement learning problem often involves the utilisation of a Markov Decision Process (MDP). In the context of Markov Decision Processes (MDPs), the agent engages in ongoing interaction with the environment by executing actions. With each action, the environment reacts and produces a new state.

Reinforcement learning constitutes a distinct domain within the broader field of Machine Learning. The concept pertains to the process of implementing appropriate measures in order to optimise the attainment of benefits within a specific context. It is utilised by diverse software and machines to determine the optimal behaviour or trajectory to be pursued in a given circumstance. Reinforcement learning and supervised learning exhibit distinct characteristics.

In supervised learning, the training data is accompanied by an answer key, enabling the model to be trained using the correct answers. Conversely, reinforcement learning lacks an explicit answer key, as the reinforcement agent autonomously determines the appropriate actions to accomplish the assigned task. When there is a lack of a training dataset, the system is compelled to acquire knowledge through its own experiences.

Reinforcement Learning (RL) is a field of study that focuses on the scientific principles behind the process of decision making. The objective is to acquire the most favourable actions within a given context in order to maximise the overall reward. In the field of Reinforcement Learning (RL), data is gathered through the utilisation of machine

learning systems that employ a trial-and-error approach. The inclusion of data is not a component typically observed in the input for both supervised and unsupervised machine learning methodologies.

Reinforcement learning employs algorithms that acquire knowledge from observed outcomes and subsequently determine the optimal course of action to be taken. Following each step, the algorithm is provided with feedback that aids in its determination of the correctness, neutrality, or incorrectness of the choice it executed. Utilising this approach seems to be advantageous for automated systems tasked with making numerous incremental decisions in the absence of human intervention.

Reinforcement learning is an autonomous and self-directed system that acquires knowledge via a process of trial and error. The agent engages in actions with the objective of optimising rewards, or more precisely, it acquires knowledge through practical experience to attain the most favourable results.

### Categories of Reinforcement Learning

Reinforcement learning is categorized mainly into two types of methods/algorithms:

- **Positive Reinforcement Learning:** Positive reinforcement learning involves the use of a stimulus or reward to enhance the likelihood of a desired behaviour recurring in the future. It amplifies the efficacy of the agent's behaviour and has a beneficial influence on it.
- **Negative Reinforcement Learning:** Negative reinforcement learning operates in a manner that is diametrically opposed to positive reinforcement learning. By avoiding the negative condition, there is an increased likelihood that the exact behaviour will be repeated.

### Elements of Reinforcement Learning

- **Policy:** Policy refers to the prescribed set of rules or guidelines that dictate the behaviour of a learning agent within a specific timeframe. The concept being described is a cognitive process that involves the translation of perceived environmental states into corresponding actions to be executed in response to those states.
- **Reward function:** The reward function is employed to establish an objective in a problem of reinforcement learning. A reward function is a mathematical function that assigns a numerical value to represent the quality or desirability of a particular situation within an environment.
- **Value function:** Value functions define the criteria for determining what is considered beneficial or advantageous in the context of long-term outcomes. The state's value refers to the cumulative benefit an agent can anticipate to accrue in the future, commencing from such state.

### Model of the environment: Models are used for planning.

- **The issue of credit allocation:** Reinforcement learning algorithms develop the ability to produce an internal value for the effectiveness of each intermediate state in achieving the goal. The agent is the name given to the learning decision-maker. The environment, which is made up of everything around the agent, is interacted with by the agent.
- The reinforcement learning problem model involves an agent that interacts with its surroundings continuously. A series of time steps determine how the agent and the environment interact. The agent chooses an action after receiving information about

## Notes

- the environment and a scalar numerical reward for the prior action at each time step t.
- The relationship between a learning agent and its environment is defined by a formal framework in terms of states, actions and rewards in reinforcement learning. This framework aims to offer a straightforward method of expressing key aspects of the artificial intelligence challenge.

### Real-world Use cases of Reinforcement Learning

- Video games:** Real-time learning algorithms are widely used in game applications. It is utilised to perform at a superhuman level. The video games AlphaGO and AlphaGO Zero are examples of well-known RL algorithms.
- Resource Management:** The work “Resource Management with Deep Reinforcement Learning” demonstrated how to employ RL in computers to automatically train and plan resources to wait for various jobs in order to minimise average job lag.
- Robotics:** Many applications of RL are found in robotics. In the industrial and manufacturing sectors, robots are deployed and reinforcement learning is used to increase their power. The development of intelligent robots utilising AI and machine learning technologies is a goal shared by many sectors.
- Text Mining:** The Salesforce organisation is currently implementing text mining, one of the great applications of NLP, with the aid of reinforcement learning.

### Advantages

- It aids in the resolution of complex real-world issues that are challenging to resolve using conventional methods.
- Since the RL learning paradigm is comparable to how people learn, the results are generally very accurate.
- Aids in obtaining long-term effects.

### Disadvantage

- For straightforward tasks, RL algorithms should not be used.
- RL algorithms call for a lot of processing and data.
- An overload of states brought on by excessive reinforcement learning may degrade the outcomes.

## 1.3 Machine Learning Algorithms

An area of artificial intelligence (AI) called machine learning (ML) enables computers to “self-learn” from training data and get better over time without having to be explicitly programmed. Detecting patterns in data and learning from them allows machine learning algorithms to develop their own predictions. In short, algorithms and models for machine learning gain knowledge through experience.

In conventional programming, an engineer for computers creates a set of instructions that tell a machine how to change input data into a desired output. The majority of instructions follow an IF-THEN structure: when particular criteria are met, the program performs a particular action.

In contrast, machine learning is a process that is automated and gives computers the ability to solve issues with little to no human involvement and make decisions based on prior experiences.

Although the terms artificial intelligence and machine learning are frequently used synonymously, they are actually two distinct ideas. Machine learning, a subset of AI that enables intelligent systems to autonomously learn new things from data, is the notion that encompasses robots making decisions, learning new skills and solving problems in a comparable fashion to people. AI is the more general term.

You can give machine learning algorithms examples of labelled data (referred to as training data) instead of programming them to carry out specific tasks, which enables them to calculate, analyse data and recognise patterns automatically.

Simply said, machine learning is a sophisticated labelling machine, according to Google's Chief Decision Scientist. Machines may be trained to label items like apples and pears by showing them examples of fruit; eventually, if they have learnt from suitable and accurate training examples, they will begin labelling apples and pears on their own.

Massive amounts of data can be used to put machine learning to work and it is considerably more accurate than people. It can assist you in saving time and money on tasks and analyses, such as reducing customer frustration to increase customer happiness, automating support ticket processing and data mining from internal sources and across the internet.

### 1.3.1 Statistical View of Machine Learning:

An important tool for studying this data and identifying trends is statistics for machine learning. It gives you the right direction for using, analysing and presenting the raw data that is successfully used in disciplines like computer vision and voice analysis, which enables you to discover previously hidden patterns.

#### Statistics for machine learning:

In the initial stages of the machine learning algorithm, statistics are introduced. Since it provides the basis for putting statistical principles into practice and, eventually, understanding conclusions from it, it aids us in dealing with the data. To put it simply, it is a branch of mathematics that is utilised for data collection, organisation and analysis. However, we must be familiar with its subcategories in order to comprehend how statistics actually works. These are listed below:

#### Descriptive statistics:

With the use of graphs and numbers, the data is arranged and summarised in this process. Examples include graphs, pie charts, histograms and others. Both population data and sample data can be used to perform it.

#### Inferential statistics:

This draws conclusions from data. On the sample data, a number of tests are run, including data manipulation, data visualisation and more, to reach various conclusions. Now that we have a fundamental understanding of statistics principles, the picture would look as follows if we laid down the points of distinction between statistics and machine learning.

#### A brief on machine learning v/s statistics

The idea that machine learning and statistics are the same thing was spread through a number of ambiguous remarks. They are not, however, the same. In the following section, we will discuss the stark contrast between statistics and machine learning.

## Notes

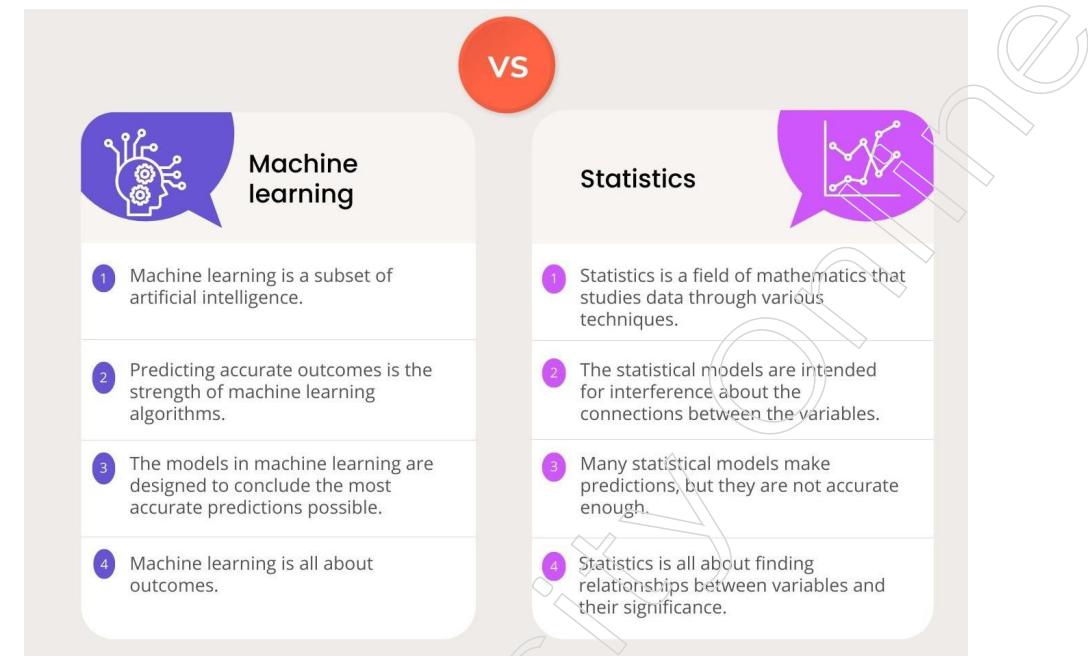


Figure: Difference between machine learning and Statistics

The statistical view is one of the key viewpoints used to comprehend and evaluate the learning process in the context of machine learning. It emphasises the connection between data and models and sees machine learning as an issue of statistical inference. By identifying patterns and relationships in the data, this point of view seeks to create predictions or conclusions based on the data that has been observed. Here are some essential ideas for understanding machine learning from a statistical perspective:

### Data:

Data is the main component of the statistical perspective. It alludes to the set of instances or examples that are employed in the development and assessment of machine learning models. These examples frequently include input features (also known as attributes or variables) and the accompanying labels or results (in supervised learning or unsupervised learning, respectively).

### Model:

The link between input features and output labels is represented mathematically by a model. It displays the underlying dependencies and patterns in the data. Finding a model that performs well on new data, makes precise predictions, or recognises significant patterns in the data is the aim of the learning process.

### Parameters:

Model parameters frequently need to be calculated from the data throughout the learning process in statistical learning. Model fitting or parameter estimation refers to the process of determining the ideal values for certain parameters.

### Probability:

When looking at machine learning from a statistical perspective, probability theory is crucial. We can use it to measure uncertainty, evaluate the likelihood of various outcomes and make inferences about the model's effectiveness.

**Objective function:**

An objective function, often referred to as a loss function or cost function, is used in many machine learning methods to gauge how well the model's predictions match the actual labels. The goal function directs learning by demonstrating how closely the model matches the intended result.

**Optimization:**

The learning process includes optimisation, which seeks to identify the model's parameters that minimise the objective function after an objective function has been created. Various optimisation strategies are commonly used to achieve this.

**Overfitting and Underfitting:**

From a statistical perspective, a significant obstacle lies in achieving an optimal equilibrium between a model that effectively captures the inherent patterns within the data, while avoiding the pitfalls of overfitting (excessive adherence to the training data) or underfitting (inadequate capture of the underlying patterns).

In general, the statistical perspective of machine learning offers a robust framework for comprehending the process of learning, evaluating the effectiveness of models and making informed judgements guided by probabilities and the presence of uncertainty. It serves as the foundation for numerous classical and contemporary machine learning algorithms and methodologies.

### 1.3.2 Multiple Linear Regression

Multiple linear regression (MLR), commonly referred to as multiple regression, is a statistical methodology that uses numerous explanatory variables to forecast the result of a response variable. The objective of multiple linear regression is to establish a mathematical model that represents the linear association between the explanatory variables, also known as independent variables and the response variables, also referred to as dependent variables. Several regression can be seen as an expansion of ordinary least-squares (OLS) regression, as it encompasses the inclusion of several explanatory variables.

Multiple Linear Regression is a significant regression approach that aims to represent the linear association between a solitary dependent continuous variable and multiple independent variables. The multiple linear regression model is a flexible statistical model used to assess the associations between a continuous dependent variable and multiple independent variables. Predictors encompass a range of variables, including continuous, categorical and derived fields, hence facilitating the analysis of non-linear interactions. The model exhibits linearity since it comprises of additive components, whereby each component represents a predictor that is multiplied by an estimated coefficient. In the model, it is customary to include a constant (intercept) term.

Linear regression is a statistical technique employed to derive meaningful interpretations from charts that exhibit a minimum of two continuous variables, wherein one variable is designated as the target variable and the other as the predictor variable. Furthermore, it is possible to include a categorical predictor and two auxiliary continuous variables in a chart, which may then be utilised to construct a suitable regression model.

- Multiple linear regression (MLR), commonly referred to as multiple regression, is a statistical methodology that uses numerous explanatory variables to forecast the result of a response variable.

## Notes

- Many regression analysis is an advanced statistical technique that expands upon the linear regression model, often known as ordinary least squares (OLS) regression, by incorporating many explanatory variables.
- Multiple Linear Regression (MLR) is widely employed in the fields of econometrics and financial reasoning.

### Formula and Calculation of Multiple Linear Regressions

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon$$

Where, for i=n observations:

$y_{(i)}$  = Dependent variable

$x_{(i)}$  = explanatory variables

$\beta_{(i)}$  = y-intercept (constant term)

$\beta_{(p)}$  = slope coefficients for each explanatory variable

$\epsilon$  = the model's error term (also known as the residuals)

The concept of simple linear regression pertains to a mathematical model employed by analysts or statisticians to generate predictions for a specific variable, relying on the available information pertaining to another variable. Linear regression is applicable exclusively in situations where there are two variables that are continuous in nature, namely an independent variable and a dependent variable. The independent variable refers to the parameter utilised in the computation of the dependent variable or outcome. The multiple regression model encompasses many explanatory variables.

The multiple regression model is founded upon a set of underlying assumptions:

- A direct correlation exists between the dependent variables and the independent variables.
- The independent variables have a low degree of correlation among themselves.
- The observations  $y_i$  are selected in a manner that is both independent and random from the population.
- The residuals are expected to exhibit a normal distribution, characterised by a mean of zero and a variance denoted by  $\sigma$ .

The coefficient of determination, commonly referred to as R-squared, is a statistical measure employed to assess the extent to which the variability in the dependent variable can be accounted for by the variability in the independent variables. The coefficient of determination ( $R^2$ ) consistently increases when additional predictors are incorporated into the multiple linear regression (MLR) model, irrespective of the potential lack of association between these predictors and the outcome variable.

$R^2$ , in alone, does not possess the capability to determine the inclusion or exclusion of predictors in a model. The coefficient of determination,  $R^2$ , is bounded between 0 and 1. A value of 0 signifies that none of the independent variables can predict the outcome, while a value of 1 shows that the outcome can be accurately predicted by the independent variables without any error.

In the context of understanding the outcomes of multiple regression analysis, it is important to note that beta coefficients retain their validity under the assumption of maintaining all other variables constant, sometimes referred to as "all else equal." The results obtained from a multiple regression analysis can be presented in two different formats: a horizontal equation or a vertical table.

### Example of How to Use Multiple Linear Regressions

For instance, an analyst may want to ascertain the impact of market fluctuations on the pricing dynamics of ExxonMobil (XOM). In this scenario, the linear equation will have the S&P 500 index as the independent variable, also known as the predictor and the price of XOM as the dependent variable.

In actuality, there are several factors that can be used to forecast the result of an event. The price fluctuations of ExxonMobil, as an illustrative example, are influenced by factors beyond the overall market performance. Additional factors, such as the fluctuations in oil prices, interest rates and the price dynamics of oil futures, have the potential to exert an influence on the valuation of XOM and the stock prices of other oil corporations. Multiple linear regression is employed to comprehend relationships involving more than two variables.

Multiple linear regression (MLR) is a statistical technique employed to establish a mathematical association between multiple independent variables and a dependent variable. In alternative terminology, multiple linear regression (MLR) investigates the relationship between numerous independent variables and a single dependent variable. After establishing the predictive capacity of each independent factor on the dependent variable, the data pertaining to various variables can be utilised to generate a precise estimation of the extent to which they influence the result variable. The model establishes a linear relationship that serves as the most accurate approximation for all the discrete data points. With regard to the multiple linear regression (MLR) equation provided above, in the specific example under consideration:

- $y_i$  = dependent variable—the price of XOM
- $x_{i1}$  = interest rates
- $x_{i2}$  = oil price
- $x_{i3}$  = value of S&P 500 index
- $x_{i4}$  = price of oil futures
- $B_0$  = y-intercept at time zero
- $B_1$  = regression coefficient that measures a unit change in the dependent variable when  $x_{i1}$  changes - the change in XOM price when interest rates change
- $B_2$  = coefficient value that measures a unit change in the dependent variable when  $x_{i2}$  changes—the change in XOM price when oil prices change

The least-squares estimates— $B_0$ ,  $B_1$ ,  $B_2$ ... $B_p$ —are usually computed by statistical software. The regression model can include as many variables as necessary and each independent variable is distinguished by a number—1, 2, 3, 4, etc. Using data on numerous explanatory variables, an analyst can predict a result using the multiple regression model!

Even so, the model is not always 100% accurate because each data point may vary somewhat from the model's outcome prediction. To account for such minute fluctuations, the residual value, E, which is the difference between the actual result and the anticipated result, is incorporated in the model.

If we use a statistics computation program to run our XOM price regression model, the results are as follows:

XOM Price = 75 - 1.5 interest rates + 7.8 oil price + 3.2 S&P 500 + 5.7 oil futures
R - Sq = 86.5%

## Notes

According to an analyst's interpretation of this production, if all other factors remain constant, the price of XOM will rise by 7.8% if the market price of oil rises by 1%. The model also predicts that after a 1% increase in interest rates, the price of XOM will fall by 1.5%. According to R<sup>2</sup>, fluctuations in the interest rate, oil price, oil futures and S&P 500 index may account for 86.5% of changes in Exxon Mobil's stock price.

### The Difference Between Linear and Multiple Regression

The response of a dependent variable in response to a change in several explanatory factors is compared using ordinary linear squares (OLS) regression. A dependent variable, however, is rarely explained by just one factor. In order to explain a dependent variable using several independent variables, an analyst employs multiple regression in this situation. There are both linear and nonlinear multiple regressions.

The foundation of multiple regressions is the presumption that the connection between the dependent and independent variables is linear. Additionally, it presupposes that the independent variables have little to no association.

Rarely can one variable fully explain a dependent variable. In these circumstances, a researcher does multiple regression, which tries to explain a dependent variable using multiple independent variables. However, the model makes the assumption that there are no significant connections among the independent variables.

## 1.4 ML Applications in Retail

Machine learning is used in the retail industry to process large datasets, identify relevant metrics, recurring patterns, anomalies, or cause-and-effect relationships among variables and thereby gain a deeper understanding of the dynamics guiding the sector and the environments where retailers operate. As more retail data is analysed, machine learning algorithms become more effective at discovering new correlations and better organising the business environment they are studying.

The Fortune Business Insights study emphasised the potential application of machine learning to develop forecasting models and deliver valuable information on constantly changing consumer preferences. Machine learning is now the largest segment of the global AI industry for retail, with a predicted growth from \$5.84 billion in 2021 to \$18.33 billion in 2028. A machine learning model for retailers might swiftly analyse a sizable amount of complex data and transform it into insights that could be used to:

- The precise prediction of future need.
- Enhancement of inventory management.
- The process of discerning consumer demands by employing suitable segmentation techniques.
- Enhancing the distinctiveness of product offerings.
- Determine the optimal pricing strategy to enhance sales.

Machine learning is being employed by a multitude of enterprises to augment the customer experience and enhance sales performance. The following are the prevailing machine learning applications in the retail industry:

- Auto-pricing
- Price optimization
- Demand prediction
- Customer segmentation

- Logistics support (Supply chain management)
- Personalized offers
- Predictive analytics
- Market basket analysis
- Churn rate prediction
- Location optimization
- Classification of customer reviews
- Fraud detection
- Determining customer lifetime value (CLTV)
- Document work automation
- Cross-sell prediction
- Merchandising

#### 1.4.1 How Machine Learning Improves Pricing and Promotions

The majority of pricing decisions were made by traditional marketers intuitively, with little regard for customer behaviour, market trends, the effects of promotions during holidays, or how these factors affected how sensitively the goods responded to price. Big data technologies are being used by most firms to optimise pricing decisions due to developments in powerful computing that enable analysis of massive volumes of data over time.

This is done to ensure that maximum clearance revenue margin targets are fulfilled while attempting to offer more competitive prices. Businesses must choose the optimal price for their goods in order to meet objectives like increasing revenue and profitability. Product price is the single most crucial factor in determining sales and revenue. Therefore, businesses must choose the optimal pricing for their products in order to maximise revenue. Our model reads and evaluates data of a product that was acquired from retail sources using an OLS linear regression model in order to train and build a demand curve. The optimal price at which the product will sell for the retailer to make the most money can be predicted using the concept of price elasticity.

Pricing and promotions across a range of businesses have been significantly impacted by machine learning. Machine learning may assist organisations in optimising pricing tactics, enhancing promotional efforts and boosting total revenue and customer happiness by utilising sophisticated algorithms and data analytics. Following are a few ways machine learning enhances pricing and promotions:

##### Adaptive Pricing:

Dynamic pricing, which adjusts prices in real-time based on variables including demand, competition, inventory levels, time of day and customer behaviour, is made possible by machine learning. As a result, pricing are constantly optimised to increase sales and profit margins. The emergence of internet sales channels in recent years has created a wealth of comprehensive client information.

The emergence of online sales channels in recent years has given sellers access to a wealth of specific client data. Customer demographics (postal code, date of birth, education level and income status), previous purchasing patterns and social media activity are a few examples of this data. Many vendors are now able to dynamically improve their pricing decisions by using sales data on these client qualities. For internet sellers, the availability of such information presents both distinct obstacles and

## Notes

opportunities. In specifically, how can a seller use dynamic learning to simultaneously apply this information in pricing decisions to maximise income over time? Customer traits have an impact on product demand.

In order to answer this query, we look into personalised dynamic pricing with learning, or dynamic pricing with limited knowledge of the impact of client characteristics on demand. To be more specific, we take into account a seller who provides customised prices to various clients whose distinctive traits are encoded as various vectors, often known as features or the feature vector.

Although the seller is unsure of the connection between client attributes and product demand, over time, sales observations can teach them about it. We stress that by personalization we do not mean client segmentation, which, to the best of our knowledge, appears to be the prevalent method for customised pricing in practice. Instead, we mean customised pricing at the individual level. We thus reserve the terms “personalised pricing” for price discrimination at the most granular level conceivable with the information supplied and “customised pricing” for price discrimination on at least one dimension of client attributes.

Pricing identical products differently for various clients may appear to be a dubious practice at first glance. However, excluding distinction that contravenes antitrust or price-fixing rules, price discrimination based on client attributes is a long-standing practice that is allowed.

**Demand Forecasting:** To estimate future demand for goods and services, machine learning models can analyse historical data and current market patterns. With this knowledge, companies may set price and promotional plans that correspond to expected demand, lowering the possibility of overstocking or understocking.

Pricing and promotions that are specifically tailored to each consumer can be developed by firms thanks to machine learning. Companies can offer customised discounts and promotions to boost customer loyalty and engagement by studying client data, preferences, purchase histories and behaviour.

**Competitive Pricing Analysis:** Machine learning algorithms can automatically alter prices to remain competitive by keeping track on competitors' pricing practices. This enables companies to keep their market share while responding swiftly to changes in the marketplace.

**Optimised Promotional Timing:** Based on past data and customer behaviour, machine learning can determine the best times to conduct promotions and discounts. This guarantees that marketing initiatives have the most possible impact, resulting in higher conversion rates and greater income.

**A/B Testing:** Machine learning can make it easier to test various pricing and promotion methods effectively. Businesses can determine the best strategies for distinct market segments by analysing the performance of various pricing models and advertising activities.

Analysis of the price sensitivity of different client segments and product categories is possible with machine learning models. Businesses can use this information to determine which items can sustain price increases and which might need more cautious price modifications.

**Preventing Price Discrimination and Collusion:** Machine learning algorithms are able to identify potential instances of price discrimination and collusion, ensuring that pricing practises comply with legal requirements and moral principles.

**Price Optimisation for Product Bundles:** Machine learning can identify the most appealing product bundles and their pricing to increase sales and satisfy customers. This encourages cross-selling opportunities and enables businesses to provide enticing packages.

**Repricing in Real-Time:** In markets that are undergoing rapid change, machine learning can help companies adjust their prices in real-time to reflect changes in the supply, demand, or competitive environment.

Businesses may increase their profitability and growth by utilising machine learning to optimise revenue, acquire useful insights into their pricing and promotional tactics and provide customers with a more individualised and engaging experience.

### 1.4.2 Using ML to Justify Promotions

The secret to success for retailers is to be close to their customers and to provide them with exactly what they need, when they need it. Customers increasingly anticipate more tailored offers and promotions depending on their demands thanks to technological improvements. When there is a steady stream of both new and returning clients, retail enterprises may survive in the market. Because of this, merchants and companies generally go to considerable efforts to satisfy their clients. Additionally, since satisfying a customer's needs is one of the best business methods, retailers are constantly coming up with innovative promos and exclusive deals to interact with their clientele.

Emerging technologies like artificial intelligence (AI) and machine learning (ML), which are opening new realities to eCommerce and merchants as a whole, complement this desire to keep customers satisfied. Retailers are hurriedly implementing these contemporary technologies in an effort to enhance consumer experience, reduce operational friction and boost financial performance.

AI and ML have been used to enhance product suggestions and ad targeting. These technologies are now being utilised to send promotions and special offers, which is a significant component of product pricing.

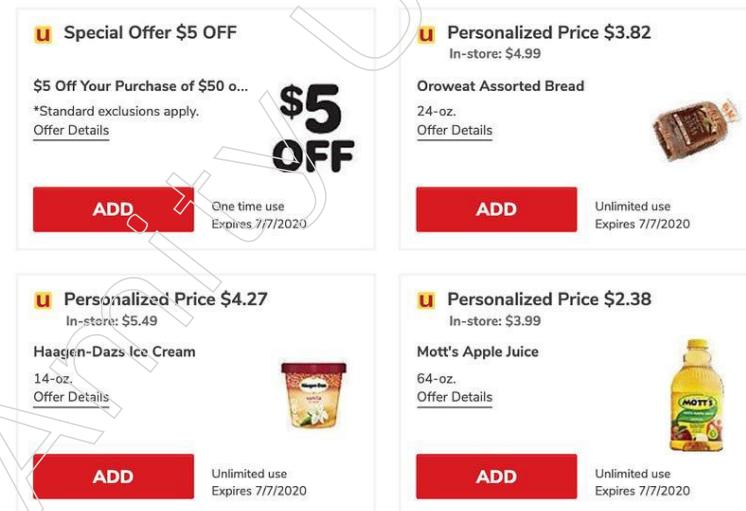


Figure: Promotions poster

### How Machine Learning Improves Pricing and Promotions

Both sides mention how ML is used to send promotions and exclusive deals. With the goal of assuring customer happiness, ML forecasts the most effective promotions

## Notes

for each client category. These promotions also have a favourable overall effect on the company's bottom line.

Machine learning, for the benefit of beginners, is "teaching a computer system, how to make accurate predictions when fed with data." So, after being "trained," an ML model can identify patterns in data without being explicitly programmed. When a shop uses an ML model, it can continuously integrate with fresh product and market data, identify patterns and make precise forecasts about demand and prices. This is its main advantage. Typically, the model created by a skilled data scientist gives the store the freedom to experiment with various approaches and factors in order to find the best pricing and promotional offers.

Here are a few particular examples of how machine learning is assisting merchants with promotions and product pricing:

### ML Is Used To Justify Promotions

Only when discounts can be recovered by the company in terms of economic value, such as improved turnover, increased client acquisition, brand loyalty, etc., are promotions justified. Managers can forecast things more precisely with ML algorithms by:

- The optimal discount value for a product
- The best time to start or end a promotion
- The best parameters for the promotion
- The overall impact of the promotion on the business balance sheet

The crucial phrase here is precision. ML models understand both past and present market trends as well as other underlying aspects to produce more reliable suggestions, whereas managers without ML seek to answer these business problems by looking at historical data.

### ML Helps Retailers Avoid Promotional Waste

Giving a promotion for a product to a customer who doesn't need them sounds like a horrible idea, unless it's a loyalty program, right? A Forrester Study found that 52% of consumers said they received weekly to monthly promotional offers on goods they would have happily paid full price for. By examining buyer data and the anticipated effects of promotions across numerous channels, AI and ML-based models can assist businesses in plugging these loss opportunities.

### ML Determines Optimal Product Prices

ML assists the retailer in determining the optimal price for a product instead of the traditional markdowns by taking into account a variety of variables, including but not limited to total expenditure, transaction volume, profits, anticipated demand, promotions data, demographics and rival prices. If applied further, machine learning (ML) can power dynamic pricing, an automated system in which prices adjust in real-time in response to some of these parameters.

### ML Gives the Managers the Tools To Meet Several KPIs

An ML-enabled pricing system is the optimum option whether the objective is to maximise profit per unit, boost overall revenue through increased turnover, or a mix of objectives. As ML becomes more advanced, marketing possibilities with clean data increase, unlocking growth prospects with important KPIs. The ML model can forecast the appropriate pricing for sales, marketing, or revenue rise depending on the factors.

### ML Enables Personalization of Promotions

A machine learning model that analyses the purchasing behaviour of individual consumers has the capability to suggest promotional offers to purchasers who exhibit similar tendencies. As an example, the model has the potential to provide a cart-based promotion that would encourage customers to surpass specific criteria in order to receive incentives.

### ML Saves Time and Effort

The machine learning model conducts an analysis of extensive data volumes and complex sets of variables, which would pose insurmountable challenges for conventional systems to handle. Machine learning models are capable of integrating new information and making predictions without the need for explicit programming, as they undergo a process of training.

### ML Reduces Promo Risks and Improve Business Bottom Line

Precise machine learning models significantly mitigate the risks encountered by retailers when implementing promotional campaigns. The algorithm facilitates the prediction of potential promotion outcomes, so aiding the product manager in assessing the optimal promotion and product pricing strategies to enhance the financial performance of the organisation. By implementing the model, the store is able to do hypothetical analysis on revenue and profit for several situations prior to initiating a sales promotion. The utilisation of Machine Learning (ML) for the purpose of justifying promotions is a multifaceted and delicate subject matter. Machine learning (ML) has the potential to provide valuable insights and impartial analysis. However, it is crucial to prioritise fairness, openness and ethical issues throughout the ML process. The following points should be taken into consideration:

**Data Collection and Bias:** It is imperative to ensure the accuracy, relevance and representativeness of the data utilised for machine learning analysis in order to assess the employees' performance. The presence of bias within historical data has the potential to sustain and perpetuate inequitable practises, thereby underscoring the importance of recognising and addressing any biases.

**Feature Selection:** Select appropriate and equitable criteria for evaluating employee performance. It is advisable to refrain from utilising attributes that have the potential to be discriminatory or unrelated to the determination of promotion.

**Algorithm Selection:** Choose machine learning algorithms that possess transparency and interpretability. The utilisation of black-box algorithms may lead to favourable outcomes; yet, it can pose difficulties in comprehending the rationale behind the promotions.

**Model Training and Validation:** It is essential to appropriately train and validate the machine learning model in order to ascertain its correctness and ability to generalise when presented with unseen data.

**Human-in-the-loop:** ML should not be employed as a fully autonomous system but rather as a tool for decision-making. The final say should always rest with human managers, who can also leverage ML insights to make better judgements.

**Regular Audits:** Audit the ML model on a regular basis to find and fix any biases or problems that might develop.

**Communicate Transparently:** Be honest and forthright with your staff when it comes to the usage of ML for promotions. Employees are entitled to information about how their performance is assessed.

## Notes

**Feedback Mechanism:** Create a feedback process so that employees may inquire about and comprehend the elements that went into their promotion choice.

**Accountability:** Design, implementation and upkeep of the ML-based promotion system should be delegated to certain people or groups. Continuous Improvement: Evaluate and enhance the ML model and promotion process continuously to keep pace with evolving needs and prevent stagnation.

**Legal and Ethical Compliance:** Make sure that all applicable laws and regulations are followed during the promotion process. Anti-discrimination laws and data privacy regulations are included in this.

Keep in mind that machine learning (ML) is a tool to aid decision-making; it should not, however, take the role of human judgement and empathy in assessing an employee's potential and appropriateness for advancement. To develop a fair and just promotion process, strive for a balance between data-driven insights and human understanding.

### 1.4.3 ML Helps Retailers Avoid Promotional Waste

The practice of teaching computer systems to learn from data and make predictions or judgements without explicit programming is known as machine learning, which is a subset of artificial intelligence (AI).

The phrase "ML helps retailers avoid promotional waste" alludes to the use of machine learning techniques to maximise the efficacy of promotions and marketing plans in the retail sector. This is how it can be done:

**Demand Prediction:** Machine Learning models can reliably predict future demand by analysing previous sales data, consumer behaviour and other pertinent aspects. Retailers may plan their promotions more efficiently and avoid costly over-promotion at times of already high demand by knowing when and how much demand is likely to increase.

**Customer Segmentation:** Machine learning algorithms are able to divide up a customer base based on preferences, purchasing patterns and demographics. This gives businesses the ability to customise promotions to particular client groups, ensuring that the proper offers reach the right audience and lowering the likelihood that some segments would be promoted with irrelevant products.

**Personalization:** Based on a customer's previous contacts with the brand, ML can be used to tailor promotional offers for specific customers. Promotional waste can be reduced and conversion rates can rise when customers respond favourably to personalised offers.

**Price Optimisation:** Using data from rival pricing, demand variations and consumer behaviour, machine learning can assist merchants in fine-tuning their price strategies. Prices that are appropriate prevent promotions from being too generous and causing unneeded losses.

**Inventory management:** ML algorithms can help merchants more effectively manage their inventory levels. Retailers can avoid overstocking or understocking situations by effectively forecasting customer demand and changing inventory levels accordingly. This reduces the need for clearance sales and waste from out-of-date or expired goods.

**Dynamic promos:** Using ML, merchants can design dynamic, real-time promos that adapt to shifting market conditions and consumer behaviour rather than static ones. This adaptability guarantees that marketing initiatives stay effective and pertinent.

**Promotion Analysis:** Retailers can use machine learning to analyse the effectiveness of previous promotions. Retailers should prevent wasting promotional funds by identifying successful promotional methods and eliminating those that did not yield meaningful returns by evaluating the impact of each promotion.

Overall, applying machine learning to the retail industry enables companies to make data-driven decisions, maximise marketing initiatives and lessen waste brought on by ineffective or poorly focused advertising. Retailers profit from increased ROI and consumers gain from receiving more pertinent and interesting offers.

Giving a promotion for a product to a customer who doesn't need them sounds like a horrible idea, unless it's a loyalty program, right? A Forrester Study found that the majority of customers acknowledged receiving weekly to monthly promotional offers for goods they would have happily paid full price for. By examining buyer data and the anticipated effects of promotions across numerous channels, AI and ML-based models can assist businesses in plugging these loss opportunities.

There is an old saying that goes, "Retailers are famous for being first at being second"; in other words, they have a bad reputation for being rather slow to adopt new technology and explore new avenues. However, there are a few notable exceptions.

First, some of the most cutting-edge technology developments in fields like recommendation engines, dynamic pricing, gamification and loyalty program models have been driven by pure-play online retailers. Second, when it comes to applying AI and machine learning in the real world and generating quantifiable business effect and ROI, the retail industry has been a surprise early adopter.

Some of the applications are very obvious and catch people's attention, such as in-store robots that scan the shelves and communicate with customers or online chatbots that improve customer service by suggesting relevant products to a customer's choices or accelerating the checkout process. Self-learning technology is used by many robots that are used in manufacturing for retail stores or in distribution centres. Other autonomous delivery technologies, such as drones and self-driving delivery robots, are powered by AI.

Price, promotion and markdown optimisation is arguably one of the most established and high-impact areas of AI in retail. More than ten years ago, second-generation price and promotion technologies hit the market. These new approaches overcame the non-productized science, whose model ran out of steam and produced a "black box" that had made retailers sceptical about adoption or accepting recommendations. In contrast to the earlier first-generation solutions, these new approaches gave retailers transparency into how the science was weighing various parameters.

Perhaps more importantly, the younger generation adopted AI and ML, allowing the algorithms to continuously learn (and unlearn when necessary) and evolve at the pace of changes in the market, the competitive environment and consumer behaviour. These algorithms became ever more capable and accurate as they chowed down on real-world data and matured natively outside of a controlled R&D environment. The retail price and promotion landscape is currently undergoing further change across all channels thanks to AI and ML. Retailers are so keen to embrace AI in pricing and promotions because it can give a win-win situation—a positive, demonstrable influence on the retailer's bottom line while also giving customers prices they perceive as fair and non-arbitrary on the things they care about most.

According to a recent global study by Forrester Consulting for Revionics, "Retail Success Requires Personalised, AI-Driven Pricing Strategies," 76% of retailers think AI-driven pricing would benefit customers. The good news is that consumers concur, with

## Notes

a different survey indicating that a resounding 78% of consumers believe it is fair to use data science to increase and decrease pricing as long as they are provided with prices they are ready to pay.

### How it works

In what specific ways do AI and ML effect price and promotion optimisation, then? With regard to pricing, these skills can assist retailers in offering targeted, more customised prices and offers that account for consumer sensitivity and competitive elasticity all the way down to the level of individual store items. In this perspective, machine learning can be thought of as a collection of various methods and strategies for problem-solving. Similar to how a specialised item can only be used effectively by a competent technician or artisan, highly skilled data scientists use their understanding of the tools in their ML toolkit to successfully use them in the area of real-world retail pricing and promotions.

Retail users can “play with the knobs” to adjust the dials to achieve their pricing strategies, for instance by applying science to different pricing strategies on goods that are successful at generating margins against those that are successful at generating traffic or transactions. These weightings and the best pricing recommendations are given to users as a result of the crucial transparency. The science analyses margins, competitive prices, pricing elasticity, brand sensitivity, good-better-best linkages, price tiers, price families and more while taking basket affinities into account.

Dynamic pricing is the practice of offering fresh price recommendations at the speed that a retailer desires in each of their channels while simultaneously responding in real-time to price adjustments made by other shops online. Online retailers benefit most from dynamic pricing, but retailers have also adopted it through in-store channels employing Electronic Shelf Labels (ESLs) or by simply giving the most important price adjustments priority.

### Promotions: Discovering what to avoid

On the promotion front, science can quickly examine past promotions and the strategies that went along with them, assisting retailers in rapidly halting negative margin leakage simply by understanding what not to do. Retailers have saved more than \$60M just by taking this one action. Moving forward, science can suggest more effective promotions that accomplish category strategy, take halo and cannibalization effects into account and satisfy financial goals.

According to a report from October 2018, 52% of consumers receive weekly or monthly promotional offers from merchants for goods they would have happily paid full price for. Fortunately, scientific study of promotion performance helps identify ineffective promotions. Stopping poor promotions can instantly save retailers millions of dollars. To be more proactive, retailers can use prescriptive AI-based analytics to recommend promotions with the best channel, vehicle and offer to deliver carefully crafted promotions that meet the retailer’s strategic goals. These recommendations take into account consumer buying influences, forecast demand, promotional vehicle impact, cross-item effects and vendor fund influence.

Retailers who use AI and ML-based price and promotion skills have a distinct advantage when it comes to pleasing their customers with thoughtfully chosen prices and promotions. Retailers who do not make use of these capabilities run the danger of alienating customers, wasting limited resources and damaging their brand. In today’s competitive retail environment, failing to take use of proven scientific capabilities can spell the difference between success and failure.

#### 1.4.4 ML Determines Optimal Product Prices

Machine learning has the potential to extend its utility beyond the development of precise pricing models for businesses. Machine learning (ML) algorithms can be effectively utilised in the context of price optimisation to precisely anticipate customer responses to specific pricing and forecast product demand. Price optimisation using machine learning incorporates a comprehensive analysis of various factors, enabling the generation of appropriate pricing recommendations for a wide range of products. This approach takes into account the primary objective of the retailer, such as enhancing sales or maximising profit margins, in order to facilitate more lucrative pricing decisions by pricing managers.

The utilisation of machine learning in pricing optimisation effectively mitigates the inherent risks associated with price adjustments, owing to its predictive capabilities. Retail teams have the ability to employ machine learning techniques in order to experiment with different promotional and pricing strategies. This allows them to gain insights into the potential impact of these initiatives, transforming their informed hypotheses into empirically supported scientific findings.

In essence, the utilisation of machine learning for price optimisation offers pricing teams the ability to determine the most suitable price for a product, taking into account various factors such as sales, revenue increase and promotional strategies.

As depicted in the aforementioned image, the use of machine learning in pricing exhibits unparalleled potential when compared to a singular reliance on human intervention. Competera's predictive algorithms include the ability to simultaneously analyse 60 pricing and non-pricing elements, hence enabling price managers to conserve 4 hours during each repricing cycle. Furthermore, our price optimisation software facilitates the transition of retail teams from SKU-based pricing to portfolio-level pricing, without any restrictions on the number of categories or goods being monitored. The utilisation of machine learning algorithms for the purpose of optimising the pricing process is imperative for pricing teams operating within established retail organisations that possess a substantial product inventory, necessitating frequent repricing activities. With the increasing prevalence of technology in various industries, the proficiency in handling machine learning (ML)-enabled software is poised to become an essential component of the job responsibilities for price or category managers. There is no feasible alternative to this approach, as it provides pricing managers with an unparalleled level of accuracy and efficiency in making decisions for a wide range of products.

The task of optimising prices in a way that simultaneously maximises revenues and avoids deterring customers from making purchases has consistently posed a difficulty for merchants. Retailers have always employed conventional, rule-based approaches to effectively handle pricing optimisation. The aforementioned approaches encompass the human examination of customer and market data, wherein pricing managers employ basic mathematical models (e.g., linear regression) to determine the impact of price adjustments on both profit margins and customer price sensitivity.

The study conducted provides the basis for the establishment of price rules, which serve to describe the rationale employed in determining pricing strategies. The pricing rules are subsequently kept within price optimisation systems, serving as the foundation for automated price adjustments. It is imperative to conduct regular checks and monitoring, such as AB testing, to ensure that the strategies are in line with prevailing market conditions and are functioning optimally. Effectively managing this process incurs significant time and effort for firms. Moreover, a considerable number of organisations lack an effective monitoring system.

## Notes

The conventional approaches to price optimisation have had additional effects due to the rapid growth of e-commerce and the total digitization of the market in recent times. Consequently, there has been a substantial surge in the volume of consumer and sales-related data that merchants are required to manage. The ever-expanding quantity of data now accessible presents a growing difficulty for merchants in accurately and consistently assessing it. In general, the prevailing market conditions in which corporations are presently engaged in competition are progressively growing more intricate. Consequently, conventional approaches to pricing optimisation are no longer sufficient in assisting merchants with price setting. Nevertheless, recent advancements in pricing optimisation technology have emerged, enabling retailers to fully leverage their data and efficiently establish prices that optimise their profitability. Machine learning is the fundamental component that underlies the functionality of this technology.

### Machine learning-based pricing

The utilisation of machine learning technology in pricing optimisation has emerged as a significant transformative force, offering merchants a means to effectively tackle the numerous obstacles they currently encounter. Firstly, it should be noted that machine learning algorithms possess the capability to analyse considerably larger datasets and consider a significantly greater number of variables compared to traditional pricing methods. In the past, pricing managers were required to engage in manual processes to ascertain pricing regulations. In contrast, machine learning models employ algorithms that iteratively learn from their outcomes in a somewhat automated manner. Retailers can utilise machine learning algorithms to establish pricing strategies in alignment with sales aims. The task can be accomplished in a totally automated manner, with significantly higher precision and a fraction of the required labour.

Pricing tools that utilise machine learning algorithms are not solely focused on acquiring knowledge, but also exhibit an enhanced ability to identify the most advantageous price thresholds for businesses. These tools gradually refine their pricing recommendations by effectively discerning the delicate balance between excessively low and excessively high price points.

Moreover, these pricing solutions that utilise machine learning techniques have the capability to incorporate essential internal data and influential external data within their algorithms. The ability of modern technologies to handle greater and more diversified datasets, in conjunction with their capacity to analyse such information, enables them to exhibit a high level of precision when determining prices in connection to these influential data points. The algorithms assess various factors, which may include:

- ❖ Historical sales and transaction data
- ❖ Seasonal changes
- ❖ Weather conditions
- ❖ Inventory levels
- ❖ Product features
- ❖ Marketing campaigns

Based on the available data, pricing solutions that utilise machine learning algorithms can effectively determine price elasticity, which quantifies the extent to which demand will vary in response to alterations in prevailing conditions. Subsequently, the software makes appropriate modifications to the costs. These techniques can also ascertain the products that exhibit a relatively consistent demand, rendering them appropriate for margin optimisation. Conversely, they can identify products that

significantly contribute to overall sales, necessitating cautious adjustments. The utilisation of machine learning algorithms for pricing strategies is increasingly prevalent within the retail industry. Based on a study conducted by IBM, it is said that a significant majority of the studied organisations, specifically 73 percent, have expressed their intention to implement intelligent automation in order to enhance their pricing and promotional strategies by the conclusion of the year 2021. In order to maintain competitiveness, merchants must prioritise the adoption of machine learning-based technology for their pricing strategies.

### Machine Learning in Retail: Dealing with Data

The volume of data in the retail industry is expanding at an unprecedented rate, rendering manual management increasingly impractical. Therefore, enterprises that aim to enhance their operational efficiency are increasingly adopting algorithms to generate pricing recommendations and forecast sales, so enabling managers to allocate their time and efforts towards strategic endeavours.

Prior to implementation, algorithms must acquire knowledge through the analysis of historical and competitive data. During the learning phase, the model or algorithm examines each individual variable that influences sales, including factors such as pricing and traffic. After the completion of the training phase, during which the algorithm attains a high level of accuracy in its predictions, then validated by empirical evidence, the model becomes prepared for a pilot implementation. Should the store express contentment with the outcomes, the model may then be employed for subsequent applications.

Frequently, merchants encounter challenges with data that is either insufficient, arduous to retrieve, or improperly formatted in an incorrect manner. In the subsequent discussion, we will explore the methodologies employed by machine learning algorithms to address the issue of inadequate data within the retail industry.



Figure: Machine Learning in Retail: Dealing with Data

### How does machine learning based pricing work?

What is the precise mechanism behind machine learning-based pricing? Although it may appear unfamiliar to certain individuals, the process of establishing an effective price optimisation program based on machine learning is actually rather basic. The procedure is as follows:

#### Data collection and data cleansing

In order to construct a machine learning model, it is important to acquire several forms of data. Specifically, in the realm of price optimisation, the database may assume the following structure:

- Transactional data:
  - ❖ A list of goods offered at various costs.
  - ❖ Product descriptions Information about each product in the catalogue (category, brand, size, colour, etc.)
- Cost data:

## Notes

- ❖ Sourcing cost
- ❖ Shipping cost
- ❖ Cost of returns
- ❖ Marketing cost
- Competitor data:
  - ❖ Competitor prices for comparable products.
- Inventory and delivery data:
  - ❖ Data on inventory levels
  - ❖ Product availabilities

### Price history

Not every sector or corporation will require or be able to use all of this information. For instance, many retailers lack a ‘clean’ price history. However, pricing based on machine learning is able to get the most information out of the data that is available. Most often, this results in a material betterment of the current situation (for example, greater profit). Additionally, businesses are understandably quite cautious when using personal data. The good news is that no personal data must be handled in order to optimise prices at the product level.

After that, the information must be error-free and ready for processing. The need to combine data from several sources and different forms makes this step difficult. Therefore, the process should be carried out by skilled data scientists who make sure the data is totally and accurately translated into an algorithm.

### Training of the Machine Learning Model

The machine learning model’s training comes next. The model first examines each variable to identify any potential impacts of pricing adjustments on sales. By doing this, the machine learning model discovers connections and patterns on its own that human analysts would typically overlook. These serve as the foundation for estimates of sales and profitability and are put into the algorithm for determining optimal prices.

The basic model is tested after it is constructed and it can also be manually optimised on a regular basis. The algorithm independently learns and enhances its outcomes with each adjustment. The algorithm’s accuracy can be improved further by adding additional data sets. While the software’s efficiency keeps rising over time, the training effort gradually declines.

### Optimization based on price elasticity predictions

Once created, a machine learning model can define the best prices to achieve particular business goals and calculate the price elasticity for thousands of products in a matter of minutes.

Because they can more accurately predict the possible impact on demand and sales, internal marketing and product teams can use these estimates to experiment more daringly with entry-level prices and discounts. They can now reason based on findings from the machine learning algorithm rather than depending solely on intuition and past experience. They have more room to operate as a result, which typically results in higher sales and profits. With good reason, retailers are embracing machine learning and AI-enabled technology to provide them a competitive edge. Particularly in the area of pricing, machine learning-based approaches are proving to be advantageous and

significant. It's challenging to accurately estimate the adoption rate of machine learning-based pricing, but e-commerce companies in particular have done so and reaped the rewards. However, a Business Insider analysis found that 72% of retail organisations intend to invest in AI and machine intelligence in 2021, likely including price optimisation.

Many well-known retailers are already using machine intelligence to their advantage. These include well-known companies including British supermarket chain Morrisons, American electronics giant Monoprice and fashion shop Bonprix.

Other well-known multinational brands that use machine learning-based pricing are:

- **Zara:** The apparel brand uses AI to set its entry-level rates and lets them automatically adjust to trends. According to Ghemawat and Nueno, as opposed to 30 to 40% at other European stores, Zara simply needs to sell 15 to 20% of its products at a discount.
- **Ralph Lauren:** With the help of markdown discounts, Ralph Lauren and Michael Kors are able to sell fewer items while also improving inventory management and sales.
- **Boohoo, Shein:** Despite having cheap entry costs, fast fashion retailers are known for using machine learning to achieve their business objectives.
- **Amazon:** Amazon sets the prices for its products using ML. The business's ML models consider a wide range of variables, including previous sales data, client demographics and rival pricing. As a result, Amazon is able to establish pricing that maximise revenues while still achieving its corporate goals.
- **Walmart:** Walmart is also utilising ML to optimise prices. Walmart uses its machine learning (ML) algorithms to establish prices that are both competitive with those of its competitors and profitable.
- **Netflix:** Netflix using ML to determine the cost of its subscription packages. The business's ML models consider a range of elements, including client demographics, viewing preferences and competitiveness. This enables Netflix to establish prices that are reasonable for its users and bring in enough money to pay its expenses.

### Benefits of machine learning based price optimization

Traditional price optimisation has limits because of its rule-based methodology, as was previously mentioned. Due to the rigidity of these pricing regulations and the fact that they only consider a small portion of price-relevant elements, businesses systematically lose money. These fundamental constraints and errors can be removed using machine learning, enabling retailers to fully utilise their data and increase their revenues. Some applications where machine learning-based technologies are advantageous are:

#### Analysis of huge and complex data sets

Traditional pricing optimisation relies on formulas from mathematics that are too simple for today's complicated market settings and the vast amounts of data that customers are now producing. The effectiveness of the price adjustment rules and the pricing manager in charge of them also influence the outcomes of such procedures. With conventional price optimisation, simple human mistake might lead to the misinterpretation of important variables or the neglect of crucial changes. All of this results in price optimisations that are not as effective as they may be.

However, machine learning models are trained to recognise even subtle correlations in pricing that is based on machine learning. More than any human could reasonably manage, they are also uniquely equipped to manage the enormous amounts of internal and external data that can affect price decisions.

## Notes

### Improved pricing in diversified assortments

Retailers often oversee huge product assortments that span numerous different categories. Due to the limitations of conventional price optimisation tools, automated changes based on pricing rules can have a detrimental effect on the sales of specific products. That or expensive manual tweaks are needed to monitor and maintain prices at the product or category level when using price automation. In contrast, technologies for price optimisation based on machine learning can adjust pricing down to the level of a specific product, leading to changes that go beyond assortment or category-level adjustments.

### No more price reductions at the expense of profit

Retailers employ price reductions, such discounts or coupons, to boost sales and get rid of outdated merchandise in their warehouses. However, if a blanket discount, such as 30% off the entire collection, is given, a store may be selling many items that might have been offered at a much lower price, losing profit. This “bulldozing” approach to pricing is typical of traditional methods, in contrast to machine-learning based pricing, which can be considerably more precise and is therefore more advantageous to retailers.

### Avoidance of loss of revenue due to traditional pricing

Traditional pricing is typically too insensitive to effectively assess how different variables, such as the time of year, the weather and other elements that affect consumer behaviour, may affect sales. Inaccuracies in the traditional approach to defining pricing guidelines, such as evaluating some elements either excessively favourably or unfavourably, might result in incorrect price setting and, eventually, a loss of sales.

Because of this, pricing regulations frequently do not operate for extremely price-sensitive products aggressively enough. Price guidelines can't change to safeguard profits on particular products because they are static. In contrast, pricing strategies based on machine learning function according to company objectives. They are able to “know” how pricing changes will affect important business KPIs by training and learning using both historical and current data and they can make adjustments accordingly.

### Improved calculation of willingness to buy

An effective pricing strategy depends on being able to accurately predict people's purchasing intentions. Traditional pricing techniques are unable to provide accurate forecasts in this case. By precisely predicting how price elasticity will change from large data variables, machine learning-based price optimisation seeks to maximise revenues.

### Consideration of a wider range of influential factors

Numerous dynamic elements have an impact on a product's price elasticity. The following is a drawback of conventional optimisation tools: To account for recent changes in the market and the state of the competition, the database and pricing rules must be manually updated on a regular basis. Price optimisation also needs to manually take into account changes in corporate strategy. Pricing based on machine learning operates considerably more autonomously and continuously gains new knowledge from fresh data.

In fact, businesses can benefit greatly from the use of machine learning to determine the best rates for their products. Analysis of a variety of parameters, including market demand, competition prices, production costs, historical sales data, customer behaviour and external economic conditions are all part of the complicated work of pricing optimisation. Businesses can set prices to maximise revenue and profit by using machine

learning techniques to model and anticipate how changes in these variables will affect sales and earnings.

Here's how machine learning can be applied to determine optimal product prices:

- Data gathering: Compile pertinent information on your product, the market, your competitors, customer behaviour and other elements that affect pricing choices.
- Finding the important factors that influence pricing decisions is a feature engineering task. These might include previous sales information, clientele statistics, economic indicators and more.
- Model selection: Pick the best machine learning technique to create a pricing model. Neural networks, decision trees, random forests and linear regression are examples of common algorithms.
- Instruction of the Model: Utilise relevant attributes and previous pricing information to train the selected model. The model will understand the connections between different variables and how they affect sales and profitability.
- Cross-validation is a technique that can be used to test and validate a model's performance. By doing this, the model is guaranteed to generalise effectively to new data.
- The model can be used to anticipate how changes in price will effect sales and profits once it has been trained and validated.
- Use optimisation algorithms to determine the price that maximises sales, profits, or any other pertinent company goal. Iterative procedures to polish the pricing approach may be necessary.
- Implement dynamic pricing techniques that let you change prices in real-time in response to shifts in demand, supply, or rival prices. This necessitates ongoing monitoring of pertinent data and pricing adjustments.
- Monitoring and Refinement: Continue to keep an eye on the effects of your pricing choices while gathering fresh information. You can hone your pricing model over time to increase accuracy and accommodate shifting market circumstances.

It's vital to remember that, even though machine learning can automate and provide useful insights for pricing decisions, human skills and subject-matter knowledge remain essential. There is a chance that machine learning models won't fully account for the complexities of human behaviour, market dynamics and outside influences. Therefore, the best method for deciding on the best product prices is frequently a blend of data-driven insights and human judgement.

#### 1.4.5 ML Saves Time and Effort

In order to process large datasets, identify pertinent metrics, recurrent patterns, anomalies, or cause-and-effect relationships among variables and thus gain a deeper understanding of the dynamics guiding this industry and the contexts where retailers operate, machine learning is used in the retail industry. Machine learning algorithms improve their performance when they find new correlations and better clarify the business context they are analysing as more retail data is processed.

Two methods are commonly used to make use of such features. First, machine learning can be used to power augmented analytics systems that, in comparison to conventional statistical analysis techniques, will go much deeper into data, identify even the most minute correlations between data points and be better able to handle new trends and constantly changing phenomena.

## Notes

Second, pattern recognition in machine learning opens the door for important advancements in the field of artificial intelligence (AI) known as “cognitive technologies,” which enable machines to mimic some of the intrinsic human skills. This includes computer vision solutions that use algorithms to find visual patterns and link them to particular objects, as well as natural language processing systems that use machine learning to recognise the linguistic patterns of human communication in order to understand and duplicate them. The capacity for retailers to use machine learning and benefit from it in a range of business operations and scenarios is effectively translated by the aforementioned. Examples include:

- Market and consumer analytics are used to anticipate retail trends, such as shifts in product demand and to create effective marketing, pricing and restocking plans.
- A fully customised shopping experience that includes recommendation engines, niche marketing, dynamic pricing and special offers based on consumer requirements.
- Through chatbots, virtual assistants and contextual purchasing, interactive solutions for digital retailers can simulate the conventional in-store experience in a virtual setting.
- Using anticipatory shipping, intelligent route planning, self-driving cars, or drones, machine learning-enhanced logistics can speed up the delivery of goods.
- Use video surveillance to monitor your assets, employees and customers. Anomaly detection based on machine learning is used to look for symptoms of fraud.

### The state of the market

Nowadays, machine learning algorithms are used in most AI-powered retail software solutions in one way or another. In fact, machine learning is one of the key factors driving the growth of the global AI market in this industry, as highlighted by Mordor Intelligence in its 2021 Artificial Intelligence in Retail Market report, as it enables market players to give end-users a more individualised and interactive purchasing experience.

On the other hand, the same problem was covered in a 2020 study by Fortune Business Insight that highlighted the value of machine learning in developing forecasting models and providing useful insights into constantly shifting consumer preferences. With a predicted increase from \$3.75 billion in 2020 to \$31.18 billion in 2028, machine learning is now the largest subset of the worldwide AI market in retail.

In this context, it's important to point out that the increasing use of machine learning and associated technologies in retail hasn't only been a reactionary move to more effectively address the problems described in our introduction. In fact, it has also been a proactive effort to take advantage of the opportunities presented by artificial intelligence, which ultimately served as a catalyst for general dynamics the retail industry had already undergone in previous years, such as the transition from a purely brick-and-mortar model to the coexistence of in-store shopping and e-commerce.

These potent machines' function in digital marketplaces is to identify the type of customer we are dealing with and match them with the appropriate goods, which is analogous to the job of a human sales assistant in a real store. Advisory systems do this by:

- Segmenting individuals based on their personal information, such as browsing behaviours, past purchases made on the platform, comments made on social media, past purchases and more.
- Adding details about the product and contextual factors, including larger market trends or geography, to this study.

- Instead of letting customers get lost in a vast maze of products, offer them personalised choices that fit their consumer archetype and then fine-tune those recommendations based on user input.

### Targeted marketing

Another technology that makes use of machine learning's promise in e-commerce but is also easily capable of boosting in-store sales and tailored advertising. Its basic mechanism resembles that of recommendation systems in certain ways. In order to investigate pertinent indicators and reveal their relationships, a machine learning-based predictive analytics system can collect and evaluate user data from social media or e-commerce platforms. These may consist of:

- Behavioural aspects, such as the duration of time spent on a website, the rate at which users navigate away from the site without interacting further and the rate at which users successfully complete desired actions or conversions.
- Demographic variables, such as age and gender, are factors that are commonly used in academic research to categorise and analyse individuals within a population.
- Psychographic measurements pertaining to individuals' interests and personality traits.
- Geospatial information such as the user's city and region.
- Based on the given characteristics and their interdependencies.

### Contextual shopping

The potential for marketers to leverage machine learning in the retail industry extends beyond the realm of targeted advertising. Contextual shopping is an additional technique that can be employed to reduce the length of the virtual pathway that links clients with the specific products they want. This software solution, which incorporates machine learning algorithms and computer vision technology, possesses a high level of interactivity. It is designed to identify and emphasise the products featured in online content on prominent social media platforms. Consequently, users are able to conveniently access your digital store and make purchases by simply clicking on the desired item.

### Chatbots and virtual shopping assistants

Chatbots are an example of machine learning in retail that focuses on interaction and contextual shopping. Natural language processing, a separate form of cognitive technology, is what gives them the ability to help clients around-the-clock with a range of tasks. This can entail providing users with assistance in locating a product they need, sending updates about new collections, proposing related products based on predictions made by recommendation engines and so forth. Such adaptable and untiring virtual assistants are widely available online and have already been used on the websites of numerous companies, including Victoria's Secret and Burberry. Additionally, they've been integrated into social media sites like Facebook, which introduced its Messenger bots in 2016 and enabled companies to partially automate their customer support processes.

### Dynamic pricing

The predictive capacities of machine learning extend beyond discerning clients' product preferences to include their price elasticity. One technique that has greatly profited from the implementation of machine learning algorithms is dynamic pricing. This approach involves making dynamic adjustments to product prices in response to changing circumstances.

## Notes

Machine learning systems can effectively assess various parameters to accurately predict the potential impact of price fluctuations on demand, also known as price elasticity. By doing so, these systems can suggest appropriate price changes or promotions that can optimise profit while minimising the risk of customer churn. This automated approach is advantageous as it allows for the consideration of a wide range of parameters that would be challenging to evaluate manually. As an illustration, one such approach is to employ web crawling techniques to retrieve data on the costs of comparable products sold by rival companies, as well as information on their ongoing promotional campaigns and special offers. Additionally, examining the pricing trends of certain products over time can also be facilitated by this method.

Machine learning algorithms have the capability to consider broader market trends, the relationship between product demand and supply and the presence of products with the lowest prices within their respective categories. This often leads to an uneven distribution of demand. These findings possess significant utility in formulating a markdown price plan and facilitating the disposal of outdated inventory at the conclusion of a particular season.

### Demand prediction for inventory management

As previously indicated, it is evident that prices have a significant influence on product demand. Consequently, the integration of machine learning in the retail sector offers the potential to elucidate the dynamics of this relationship, so enabling the optimisation of pricing strategies. However, it is important to note that the dynamics of demand are far more intricate due to the influence of numerous variables operating on a broad scale.

- Business endeavours, including the aforementioned pricing strategies, promotional activities and the management of product assortment and display within a retail establishment.
- The study focuses on consumer purchasing behaviours with regards to certain temporal factors such as weekdays, seasonality and celebratory occasions, shown by the demand for chocolate hearts during Valentine's Day.
- Various external factors can serve as motivators for customers to venture outside their homes, such as favourable weather conditions or the presence of local activities.

The monitoring of various parameters to anticipate fluctuations in product demand is a crucial element of effective stock and supply chain management. This practice enables retailers to maintain adequate stock levels. Deloitte has emphasised that retailers may incur losses of up to \$1 trillion annually due to stockouts. Additionally, Deloitte reported on the effective implementation of machine learning algorithms by a prominent German sportswear retailer to optimise restocking processes across both online and physical stores. In order to effectively monitor these parameters, retailers should continuously supply machine learning-based predictive analytics systems with a steady flow of information derived from various data channels. These channels encompass both virtual sources such as social media and ecommerce platforms, as well as data obtained from physical stores. These may encompass "digital sources" such as social media, electronic commerce platforms and the incorporation of voice commerce in conjunction with data gathered from brick-and-mortar establishments.

### Delivery optimization

The turnover of products is a dynamic process and it is as crucial to focus on enhancing their exit strategy from the store as much as their restocking process. The

utilisation of machine learning in route planning within the transportation sector is a significant application of artificial intelligence. This approach capitalises on the integration of many interconnected devices, including sensors and cameras, through the Internet of Things. These devices enable the collection of real-time weather and traffic data, which is then utilised to determine the most efficient route for deliveries.

Anheuser-Busch, as an illustration, implemented machine learning techniques for the purpose of optimising routes in a trial initiative including two cities in the United States, with the objective of enhancing the efficiency of daily deliveries. The algorithms took into account variables such as weather conditions and the driver's level of expertise in order to provide the optimal delivery time for each individual consumer. Following a period of several months during which the pilot program was implemented, Anheuser-Busch conducted an evaluation of the program's efficacy, taking into account several success indicators such as driver satisfaction and working hours. Based on this assessment, the company made the decision to expand the utilisation of this navigation strategy to encompass all of its wholesalers operating within the United States.

#### 1.4.6 ML Reduces Promo Risks and Improve Business Bottom Line

Machine learning (ML) has the potential to mitigate promotional risks and enhance a business's financial performance through various mechanisms. There are several methods by which machine learning (ML) can get these advantages:

- Personalised Targeting: Machine learning algorithms possess the capability to scrutinise extensive quantities of client data, encompassing purchase history, browsing behaviour and demographics, with the objective of generating customer profiles that exhibit enhanced accuracy. By utilising this data, businesses have the ability to strategically focus their promotional efforts and tailor offers to particular client segments, so enhancing the likelihood of achieving desired outcomes and mitigating the potential for irrelevant or unsuccessful promotions.
- The utilisation of machine learning algorithms enables the analysis of past sales data, facilitating the identification of patterns that might inform optimal timing for promotional campaigns. By strategically selecting the appropriate timing, organisations have the ability to circumvent the initiation of promotional activities during periods of low demand. This approach allows for the optimisation of promotional impact during moments of heightened demand, thereby mitigating the potential for promotions to yield subpar results.
- Dynamic pricing refers to the practice of using machine learning algorithms to analyse current market data, competition pricing strategies and customer demand in order to make real-time adjustments to prices. This facilitates enterprises to provide competitive pricing while upholding profit margins, hence mitigating the likelihood of overpricing or underpricing products during promotional activities.
- Fraud Detection and Prevention: Promotional events are associated with an elevated susceptibility to fraudulent actions, including the creation of counterfeit accounts, misuse of coupons and perpetration of transactional fraud. Machine learning (ML) has the potential to aid in the identification of fraudulent behaviour by examining trends and abnormalities within client data. This analytical approach can effectively reduce risks and safeguard the financial performance of the firm.
- Demand forecasting is a process that utilises machine learning techniques to generate precise predictions on the demand for particular items or services. This is achieved by analysing historical data and market patterns. This enables organisations to strategically plan promotions, so optimising supply availability and mitigating the potential risks associated with excessive inventory or stock shortages.

## Notes

- Recommendation engines, which are powered by machine learning algorithms, have the capability to provide customers with suggestions for related items or services. These recommendations are generated by analysing the customers' interests and historical purchase history. By providing pertinent suggestions during promotional activities, enterprises can enhance prospects for cross-selling and enhance their total sales success.
- Churn prediction and customer retention involve the utilisation of machine learning algorithms to forecast customer churn. This enables businesses to take proactive measures by targeting customers who are at danger of churning with customised promotions or offers, so incentivizing them to remain loyal. Consequently, this strategy effectively reduces customer attrition and enhances the overall financial performance of the firm.
- Machine learning (ML) has the capability to analyse advertising data and discern the most efficient channels and advertising campaigns. By using strategies to optimise advertising expenditure, firms can mitigate the potential for resource wastage on inefficient advertising campaigns and enhance the return on investment (ROI) derived from promotional endeavours.
- Supply chain optimisation is the utilisation of machine learning (ML) techniques to enhance the efficiency of supply chain operations, resulting in cost reduction and mitigation of inventory holding risks. This practice guarantees the availability of products during critical periods, such as promotional events, hence mitigating the risk of financial losses resulting from stockouts.

In general, the utilisation of machine learning in promotional tactics enables firms to make decisions based on data, mitigate risks and get more favourable results, ultimately resulting in enhanced financial performance.

### Summary

- In 1950's the term "artificial intelligence" (AI) is coined and researchers begin exploring the idea of creating machines that can mimic human intelligence.
- Machine learning has become a transformative technology that influences nearly every aspect of modern life. Its applications are vast and growing, making it a crucial tool for businesses, researchers and individuals seeking to leverage data and automation for improved decision-making and problem-solving.
- Machine learning (ML) plays a crucial role in computer science by enabling computers to learn from data and perform tasks that were previously considered the domain of human intelligence. It enhances various aspects of computer science and drives innovations across multiple subfields.
- Supervised machine learning is widely used in real-world applications, including image recognition, speech recognition, medical diagnosis, credit scoring and more. It forms the foundation of many machine learning techniques and serves as a basis for further exploration into more complex algorithms and methodologies.
- Unsupervised machine learning is a category of machine learning where the algorithm learns patterns from unlabeled data without explicit target labels. The goal of unsupervised learning is to find hidden structures, groupings, or relationships within the data. Unlike supervised learning, there are no predetermined outcomes to predict. Unsupervised learning is particularly useful for exploring and understanding the underlying characteristics of data
- Reinforcement Learning (RL) is a category of machine learning that focuses on training agents to make sequences of decisions in an environment to maximize a

cumulative reward. Unlike supervised learning where the model learns from labeled data and unsupervised learning where the goal is to find patterns in data, RL involves learning through interaction and feedback from the environment.

- Multiple Linear Regression is a statistical technique used in machine learning and statistics to model the relationship between multiple independent variables (also called features or predictors) and a single dependent variable (also called the target or outcome). It's an extension of simple linear regression, which deals with only one independent variable. Multiple linear regression aims to find the best-fitting linear equation that explains how the independent variables collectively influence the dependent variable.
- Machine learning has significantly improved the pricing and promotions strategies for businesses across various industries. By leveraging the power of data analysis, predictive modelling and optimization algorithms, machine learning enables businesses to make more informed decisions about pricing, discounts and promotions.

### Glossary

- ❖ ML: Machine Learning
- ❖ GPU: Graphics Processing Units
- ❖ ANN: Artificial Neural Networks
- ❖ KNN: K-Nearest Neighbour
- ❖ AI: Artificial Intelligence
- **Labelled Data:** Labelled data comprises a label or target variable that the model endeavours to predict, whereas unlabeled data lacks a label or target variable.
- **Supervised Learning:** Supervised learning involves the acquisition of knowledge by a model through the use of annotated training data, enabling the model to generate predictions based on the provided labels. Typical tasks encompass classification and regression.
- **Unsupervised Learning:** In the context of unsupervised learning, the model acquires knowledge of patterns through the analysis of unlabeled data. Clustering and dimensionality reduction are often encountered activities within this particular domain.
- **Reinforcement Learning:** The process of reinforcement learning entails the training of a computational model to make optimal judgements by use of iterative interactions with an external environment. The model is provided with feedback in the form of either rewards or punishments.
- **A/B Testing and Experimentation:** Machine learning assists in designing and analyzing A/B tests for pricing and promotion changes. It helps measure the impact of different strategies on customer behavior and revenue.

### Check Your Understanding

1. What event is considered the birth of the field of artificial intelligence (AI) and marks the beginning of machine learning research?
  - a) The creation of the first computer
  - b) The introduction of deep learning algorithms
  - c) The Dartmouth Workshop in 1956
  - d) The development of the first neural network

**Notes**

2. In which decade did machine learning gain momentum with advancements in computational power and algorithmic techniques?
  - a) 1960s
  - b) 1970s
  - c) 1980s
  - d) 1990s
3. Which historic milestone showcased the potential of machine learning by defeating human champions in the board game Go?
  - a) Deep Blue defeating Garry Kasparov in chess
  - b) Watson winning Jeopardy!
  - c) AlphaGo defeating Lee Sedol in Go
  - d) DeepMind's victory in the Turing Test
4. What is the primary goal of supervised machine learning?
  - a) Clustering data points into groups
  - b) Finding hidden patterns in unlabeled data
  - c) Predicting outcomes based on labeled data
  - d) Identifying anomalies in a dataset
5. Which of the following is an example of a supervised learning problem?
  - a) Image segmentation
  - b) Customer segmentation
  - c) Email spam detection
  - d) Principal Component Analysis (PCA)
6. In supervised learning, what are the labeled data points used for?
  - a) Feature extraction
  - b) Training the model
  - c) Evaluating model performance
  - d) Identifying outliers
7. What are the two main types of supervised learning tasks?
  - a) Regression and classification
  - b) Clustering and dimensionality reduction
  - c) Feature selection and feature scaling
  - d) Reinforcement learning and unsupervised learning
8. What is the primary objective of unsupervised machine learning?
  - a) Predicting outcomes based on labeled data
  - b) Grouping similar data points into clusters
  - c) Maximizing rewards through interaction with an environment
  - d) Determining the optimal policy for decision-making
9. Which of the following is an example of an unsupervised learning task?
  - a) Handwriting recognition
  - b) Sentiment analysis
  - c) Customer churn prediction
  - d) Clustering customer segments
10. In unsupervised learning, what are the data points typically lacking?
  - a) Labels
  - b) Features
  - c) Clusters
  - d) Relationships
11. Which key concept in unsupervised learning involves reducing the dimensionality of the data while preserving important information?
  - a) Clustering
  - b) Regression
  - c) Principal Component Analysis (PCA)
  - d) Reinforcement Learning

12. What is the primary goal of reinforcement learning?
- a) Grouping similar data points into clusters
  - b) Predicting outcomes based on labeled data
  - c) Maximizing cumulative rewards through interaction with an environment
  - d) Classifying data into predefined categories
13. Which of the following is an essential element in reinforcement learning?
- a) Labeled training data
  - b) Unlabeled training data
  - c) Interaction with an environment
  - d) Static dataset
14. In reinforcement learning, what do rewards signify?
- a) Losses incurred by the agent
  - b) Anomalies in the environment
  - c) Positive or negative feedback for actions taken
  - d) The probability of taking a particular action
15. Which type of learning involves the trade-off between exploration and exploitation?
- a) Supervised learning
  - b) Unsupervised learning
  - c) Reinforcement learning
  - d) Semi-supervised learning
16. What is the primary objective of multiple linear regression?
- a) Classifying data into predefined categories
  - b) Finding hidden patterns in unlabeled data
  - c) Predicting a single dependent variable using multiple independent variables
  - d) Clustering data points into groups
17. In multiple linear regression, how many independent variables are considered?
- a) One
  - b) Two
  - c) Three or more
  - d) It varies based on the problem
18. What is the purpose of the coefficients in multiple linear regression?
- a) They represent the intercept term only.
  - b) They determine the relationship between the independent and dependent variables.
  - c) They are used for clustering data points.
  - d) They have no significance in regression.
19. How is multiple linear regression different from simple linear regression?
- a) Multiple linear regression has only one independent variable.
  - b) Multiple linear regression can handle only categorical variables.
  - c) Multiple linear regression deals with multiple dependent variables.
  - d) Multiple linear regression involves multiple independent variables.
20. How does machine learning benefit retailers?
- a) By predicting stock market trends
  - b) By automatically generating code for online stores
  - c) By optimizing pricing, inventory and demand forecasting
  - d) By providing social media marketing services

**Notes**

## Notes

### Exercise

1. Write a short on history of machine learning.
2. What is the role of data and connection to the knowledge?
3. Define machine learning.
4. Explain different types of machine learning.
5. What is the role of machine learning in computer science and problem solving?
6. Define multiple linear regressions.
7. How machine learning improves pricing and promotions?
8. Define ML applications in retail.

### Learning Activities

1. How does machine learning contribute to problem-solving within the field of computer science? Provide examples of specific problem domains where machine learning techniques have been successfully applied to improve solutions and outcomes.
2. Explain how unsupervised learning techniques can be used to identify anomalies or outliers in datasets.

### Check Your Understanding- Answers

- |       |       |       |       |
|-------|-------|-------|-------|
| 1. c  | 2. d  | 3. c  | 4. c  |
| 5. c  | 6. b  | 7. a  | 8. b  |
| 9. d  | 10. a | 11. c | 12. c |
| 13. c | 14. c | 15. c | 16. c |
| 17. c | 18. b | 19. d | 20. c |

### Further Readings and Bibliography

1. Pattern Recognition and Machine Learning by Christopher M. Bishop
2. Machine Learning: A Probabilistic Perspective by Kevin P. Murphy
3. Elements of Statistical Learning by Trevor Hastie, Robert Tibshirani and Jerome Friedman
4. Machine Learning for Retail by Vineet Vashishta and Shubhendu Trivedi

# Module - II: Exploratory Data Analysis

Notes

## Learning Objectives

At the end of this module, you will be able to:

- Define exploratory data analysis (EDA)
- Analyse the different types of data types
- Explain feature extraction
- Analyse dimensionality reduction
- Analyse the process and tools for data visualisation

## Introduction

Exploratory data analysis (EDA) refers to the preliminary examination and investigation of data sets in order to uncover patterns, relationships and potential insights. It involves the utilisation of various statistical techniques, visualisations and computational tools to get a deeper understanding. What were the origins of this phenomenon? What are the origins and locations of its inception? How does this methodology differ from conventional and Bayesian data analysis methods? Is exploratory data analysis (EDA) synonymous with statistical graphics? Statistical graphics play a significant part in exploratory data analysis (EDA). Is statistical graphics synonymous with exploratory data analysis (EDA)?

The present topic addresses these inquiries and their associated inquiries. This topic addresses the aforementioned inquiries and offers the essential context for understanding the assumptions, ideas and procedures of exploratory data analysis (EDA). The EDA method can be characterised as an overarching framework rather than a rigid set of techniques. It embodies a particular mindset or philosophy regarding the appropriate manner in which a data analysis should be conducted.

Although EDA and statistical graphics are often used interchangeably, it is important to note that they are not identical. Statistical graphics encompasses a variety of visually-oriented strategies that are centred around the characterization of data from a singular perspective. Exploratory Data Analysis (EDA) includes a broader scope, as it pertains to a data analysis methodology that defers the conventional assumptions regarding the data's adherence to a specific model. Instead, EDA adopts a more direct approach by permitting the data to manifest its inherent structure and model. Exploratory Data Analysis (EDA) encompasses more than a mere compilation of techniques.

### 2.1 Introduction to Data and Data Distribution

#### Data

Data can be defined as factual information, such as measurements or statistics, that serves as a foundation for reasoning, discussion, or calculation. Data can be defined as information that is specifically organised for the purpose of analysis.

- Factual information, such as measurements or statistics, serves as the foundation for reasoning, discussion and calculation.
- The output of a sensing device or organ consists of a combination of helpful information as well as irrelevant or redundant information, which necessitates processing in order to derive meaning.

## Notes

- Digital data refers to information that is represented in numerical form and may be transported or processed using digital technology. Drawing from the aforementioned definitions, a pragmatic perspective on the definition of data posits that data encompasses numerical values, textual characters, visual representations, or any other means of recording information, presented in a format that enables evaluation for the purpose of reaching a conclusion or making a specified course of action. It is widely held that data lacks inherent meaning and only acquires significance and transforms into information through interpretation.

### Data Distribution

The concept of data distribution encompasses the complete set of potential values for a given variable, while also providing a measure of the relative frequency or probability associated with each value. Distributions are defined as populations that exhibit a dispersion of data points. The determination of the distribution of a population is of utmost significance as it enables the appropriate application of statistical methodologies during its analysis.

Data distributions play a key role in the field of statistics. Assume an engineer acquires a sample size of 500 data points within a manufacturing facility. The lack of value for management arises from the absence of data categorization or organisation that would render it useful.

Data distribution methods are utilised to arrange raw data into various graphical representations, such as histograms, box plots and run charts, among others. These methods serve the purpose of offering valuable insights and information.

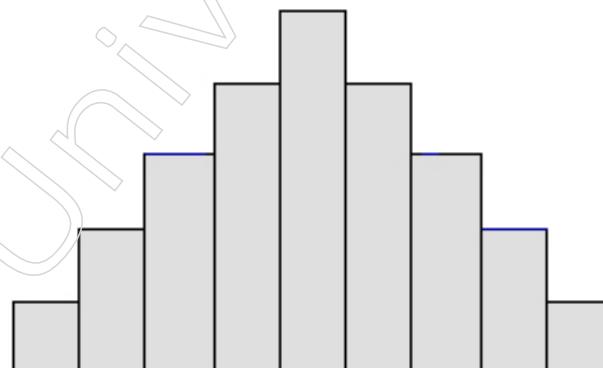


Figure: Histogram

The primary benefit of data distribution is in its capacity to calculate the likelihood of a particular observation occurring within a given sample space. A probability distribution is a mathematical construct that quantifies the likelihood of many potential outcomes in a given test or experiment. Random variables are employed in order to establish distinct categories for various forms of data, which are normally classified as either discrete or continuous. These categorizations are crucial for making informed decisions, since they are contingent upon the specific models being utilised. Various statistical methods, such as mean, mode, range and probability, can be employed within the context of random variable categorization.

**Types of Distribution:** Distributions are typically categorised according to the nature of the data.

- Typically discreet
- Continuous

**Discrete Distributions:**

A discrete distribution arises when dealing with countable data that has a limited number of potential values. In addition, it is possible to provide discrete distributions in tabular form, where the values of the random variable are enumerable. For instance, engaging in the act of rolling dice or selecting several heads.

**Probability mass function (pmf):**

A probability mass function (PMF) is a mathematical function that assigns probabilities to discrete random variables. It is sometimes referred to as a discrete density function.

Simply Discrete = counted

Different types of discrete distributions are:

**Binomial distribution:**

The binomial distribution quantifies the likelihood of obtaining a specific number of successful or unsuccessful outcomes in a given experiment for each trial.

Characteristics can be categorised into two distinct and comprehensive classes, namely the count of successes/failures and the count of accepted/rejected instances, which adhere to a binomial distribution.

**Ex: Tossing a coin:** The probability of a coin landing on heads is 1/2, whereas the probability of it landing on tails is also 1/2.

**Poisson distribution:**

The Poisson distribution is a discrete probability distribution that quantifies the probability of a certain number of events occurring during a specified time period, assuming that the events occur sequentially and in a clearly defined manner.

The Poisson distribution is associated with characteristics that possess the potential to assume high numerical values in theory, but in practise tend to exhibit smaller values.

**Ex:** Number of defects, errors, accidents, absentees, etc.

**Hypergeometric distribution:**

The hypergeometric distribution is a discrete probability distribution that quantifies the likelihood of observing a particular number of successes in a series of trials ( $n$ ) without replacement, given a sufficiently large population size ( $N$ ). In alternative terms, the process of sampling without replacement.

The hypergeometric distribution exhibits similarities to the binomial distribution. However, a fundamental distinction lies in the fact that the probability of success remains constant for all trials in the binomial distribution, whereas this is not the case for the hypergeometric distribution.

**Geometric distribution:**

The geometric distribution is a discrete probability distribution that quantifies the probability of the occurrence of the first success.

The concept can be further expanded to encompass the negative binomial distribution.

**Ex:** A marketing representative employed by an advertising agency employs a

## Notes

random selection process to identify hockey players from multiple institutions until he successfully identifies a hockey player who has participated in the Olympic Games.

### **Continuous Distributions:**

A continuous distribution is characterised by an unlimited number of data points that are displayed on a continuous measuring scale. A continuous random variable is a type of random variable characterised by an unlimited and uncountable collection of possible values. Instead of just tallying, it quantifies a phenomenon; conventionally, it is characterised using probability density functions (PDFs).

### **Probability density function (pdf):**

The probability density function characterises the statistical distribution of a random variable. The typical approach involves organising data into a grouped frequency distribution. Therefore, the probability density function perceives it as the inherent structure of the distribution. The concept of "simply continuous" refers to a variable that has the ability to assume a wide range of values. Different types of continuous distributions are:

### **Normal Distribution:**

Normal distributions are commonly referred to be Gaussian distributions by certain individuals. The distribution has a symmetrical bell-shaped curve, characterised by a greater frequency (probability density) in the vicinity of the central value. The frequency exhibits a significant decline as values deviate from the central value in either direction.

In essence, it is anticipated that a normal distribution will exhibit symmetrical properties, wherein the occurrences of values deviating from the desired mean are equally likely to occur on both sides. In a normal distribution, the mean, median and mode exhibit equivalence.

### **Lognormal distribution:**

A lognormal distribution is observed in a continuous random variable  $x$  when its natural logarithm,  $\ln(x)$ , conforms to a normal distribution. As the sample size increases, the sum distribution of random variables tends to approximate a normal distribution, irrespective of the underlying distribution of the individual variables. The aforementioned situation is also applicable to the concept of multiplication. The location parameter represents the central tendency of the data set, which is obtained by applying a logarithmic transformation. Similarly, the scale parameter corresponds to the dispersion of the data set, calculated as the standard deviation following the same transformation.

### **F distribution:**

The F distribution is commonly employed in statistical analysis to assess the equality of variances between two populations that follow a normal distribution. The F distribution is a non-symmetric probability distribution characterised by a lower bound of 0, while lacking an upper bound. Significantly, the curve asymptotically approaches 0 without intersecting the horizontal axis.

### **Chi Square distribution:**

The chi-square distribution is obtained by squaring and summing independent variables that follow a conventional normal distribution.

**Ex:** If  $Z$  represents the standard normal random variable, then:

$y = Z_{12} + Z_{22} + Z_{32} + Z_{42} + \dots + Z_n$

The chi-square distribution exhibits symmetry and is constrained to values greater than or equal to zero. As the degrees of freedom rise, the form of the distribution tends to approximate that of a normal distribution.

### **Exponential distribution:**

The exponential distribution is a commonly employed continuous probability distribution that is frequently utilised in several fields. Frequently employed for the purpose of simulating objects characterised by a consistent rate of failure. The exponential distribution exhibits a close relationship with the Poisson distribution. The failure rate is constant as the form parameters remain same.

**Ex:** The lifetime of a bulb, the time between fires in a city.

### **T Student distribution:**

The t-distribution, also known as the Student's t-distribution, is a probability distribution that exhibits a bell-shaped curve and is symmetric with respect to its mean. Hypothesis testing and the construction of confidence intervals for means are frequently employed in statistical analysis.

The Student's t-distribution is employed as a substitute for the regular normal distribution when the value of the standard deviation is not known. Similar to the normal distribution, the distribution of averages of random variables tends to approximate a normal distribution, irrespective of the underlying distribution of the individual variables.

### **Weibull Distribution:**

- The primary objective of utilising the Weibull distribution is to effectively represent and analyse time-to-failure data. This distribution has extensive application in several fields such as reliability engineering, medical research and statistical analysis.
- The form of an object can vary based on its shape, size and position factors. The influence of the shape parameter  $\beta$  on the Weibull distribution.
- For example, when the shape parameter  $\beta$  is equal to 1, the distribution becomes indistinguishable from the exponential distribution.

### **Non-normal distributions:**

It is commonly assumed that during the process of conducting a hypothesis test, the data being analysed is a representative sample derived from a normal distribution. However, it is important to note that this assumption does not hold universally.

Put otherwise, it is possible for data to deviate from a normal distribution. Therefore, nonparametric tests are employed in situations where there is no underlying assumption about the distribution of the population.

### **Odd Distributions**

#### **Bivariate Distribution:**

The probability distribution of a random variable can be categorised into two main types: continuous distributions, such as the normal, chi-square and exponential distributions and discrete distributions, such as the binomial and geometric distributions.

The bivariate distribution refers to the likelihood of an event occurring when there are two independent random variables present. This distribution can take on either a continuous or discrete form. The bivariate distribution is distinguished by its characteristic as the joint distribution of two variables.

## Notes

### Bi-modal:

- A bi-modal distribution is characterised by the presence of two distinct modes. In essence, it is quite probable that two outcomes will be compared in order to assess the relative merits of their respective regions.
- Two sources of data are inputted into a unified process panel.

### 2.1.1 What is Exploratory Data Analysis (EDA)

Exploratory data analysis (EDA) is a widely employed technique in the field of data science, utilised for the purpose of analysing and investigating data sets in order to derive meaningful insights and summarise their primary characteristics. This approach frequently incorporates the use of data visualisation methods to facilitate the exploration and comprehension of the data.

Exploratory data analysis, commonly referred to as EDA, holds paramount significance as the initial phase in the examination of experimental data. There are several primary justifications for employing exploratory data analysis (EDA).

- Identification of errors.
- Verification of assumptions.
- Initial identification and evaluation of suitable models.
- Evaluating the relationship and estimated magnitude of connections between independent and dependent variables.
- Identifying the associations between the explanatory variables. Exploratory data analysis, in its broadest sense, encompasses several methods of analysing data that do not rely on formal statistical models and inference techniques.

Exploratory Data Analysis, or EDA, is a technique for analysing and comprehending data by employing a variety of methods, such as data transformation, summary statistics and visualisation, to abstract its essential features. EDA is typically carried out before formal modelling or hypothesis testing in order to obtain a sense of the data and identify any potential difficulties or concerns that need to be addressed.

In order to help with future analysis or decision-making, it seeks to spot patterns, connections and trends in the data. EDA can be used to analyse data of many forms, including text, category and numerical data. To find and fix data mistakes, as well as to visualise the essential data qualities, this is often done before data analysis.

Exploratory Data Analysis (EDA) is a method or philosophy for data analysis that makes use of a range of methods, most of which are graphical, to:

The objectives of this analysis are to enhance understanding of a given data set, reveal the fundamental structure within it, identify significant variables, identify outliers and anomalies, validate underlying assumptions, construct concise models and establish the most favourable factor settings.

The EDA methodology is exactly that—an approach—rather than a collection of methods or a set of principles for doing data analysis.

The aim to understand the engineering/scientific process underlying the data is EDA's overarching purpose. EDA is proactive and futuristic in contrast to summary statistics, which are passive and historical. EDA uses the data as a “window” to look into the process that produced the data in an effort to “understand” the process and make it better going forward. Summary statistics have an archive function in the area of research and production, but the EDA approach has an exponentially bigger use.

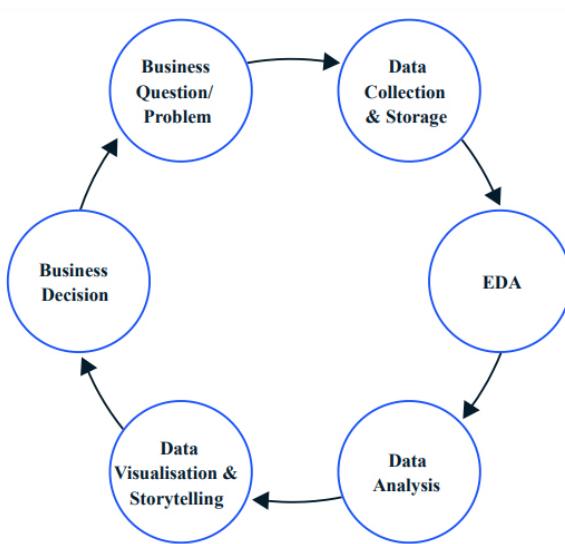


Figure: Process of EDA

EDA is a methodical way to comprehending data archiving. By skillfully manipulating data sources, data scientists can utilise it to find patterns, identify anomalies, test theories, or confirm presumptions.

#### Typical data format and the types of EDA:

In the context of experimental data storage, it is common practise to employ a rectangular array. This array consists of rows that correspond to each experimental subject and columns that represent subject identification, outcome variables and explanatory variables. Each column inside the dataset contains either the levels for a categorical variable or the numerical values associated with a specific quantitative variable.

Individuals often lack proficiency in assessing a column of numerical values or an entire spreadsheet in order to discern crucial elements of the data. Individuals may see the act of examining numerical data as arduous, lacking in appeal, or overwhelming. Various techniques have been developed to facilitate exploratory data analysis.

The majority of these solutions operate by partially concealing certain aspects of the data while simultaneously enhancing the visibility of others.

Exploratory data analysis is commonly cross-classified into two groups. One initial differentiation can be made between graphical and non-graphical methodologies. Furthermore, it is worth noting that each technique can be classified as either multivariate or univariate. Graphical procedures are effective in presenting data in a visual or pictorial format, whereas non-graphical methods commonly include the aggregation of summary statistics.

Multivariate methods are employed to analyse the interrelationships among two or more variables concurrently, whereas univariate approaches concentrate on the examination of individual variables (data columns) in isolation. The multivariate exploratory data analysis (EDA) we will do mostly focuses on bivariate analysis, examining two elements simultaneously. However, there may be instances where our analysis includes three or more variables.

#### Univariate non-graphical EDA:

Our observations for a single attribute, like as age, gender, speed at a task, or

## Notes

response to a stimulus, are represented by the data that result from performing a specific measurement on all of the subjects in a sample. These observations should be viewed as a “sample distribution” of the variable, which in turn roughly corresponds to the variable’s “population distribution.” Univariate non-graphical EDA often aims to increase understanding of the “sample distribution” and draw some broad conclusions about which population distribution(s) are compatible with it. This approach also includes outlier detection.

### Categorical Data:

The range of values and the frequency (or relative frequency) of occurrence for each value are the only attributes of a categorical variable that are of interest. Using the methods in the second half of this section, it is occasionally suitable to treat ordinal variables as quantitative variables. Therefore, tabulation of the frequencies in some form, typically accompanied with determination of the proportion (or percent) of data that falls in each category, is the only usable univariate non-graphical technique for categorical variables.

A real population of all students enrolled in the Fall semester exists, for instance, if we categorise courses by College at Carnegie Mellon University as H&SS, MCS, SCS and “other”. For the purpose of conducting a memory experiment, a random sample of 20 students might be chosen and the sample’s “measurements” may be listed as H&SS, H&SS, MCS, other, other, SCS, MCS, other, H&SS, MCS, SCS, other, MCS, MCS, H&SS, MCS, other and H&SS, SCS. Our EDA would appear like follows:

Statistic/College	H&SS	MCS	SCS	other	Total
Count	5	6	4	5	20
Proportion	0.25	0.30	0.20	0.25	1.00
Percent	25%	30%	20%	25%	100%

Note that the total count (frequency) is helpful in confirming that we have one observation for each participant we recruited. (Losing data is a frequent error and EDA is particularly useful for identifying errors. Furthermore, if we accurately calculate the proportions, we should anticipate that they sum up to 1.00. You won’t need both proportion (relative frequency) and percent after you become used to it because you’ll think of them as interchangeable terms.

### Importance of EDA in data science :

The phase of exploratory data analysis holds significant importance within the data science process as it facilitates a comprehensive understanding of the data at a more profound level for data scientists. This inquiry aims to elucidate the significance of exploratory data analysis (EDA) within the field of data science through the delineation of its objectives.

- Performing an exploratory data analysis (EDA) might serve to validate the feasibility of the gathered data within the specific business challenge under consideration. If the data or the technique employed by the data analysts is found to be inadequate, it may be necessary to make modifications.
- Data quality concerns, such as duplicates, missing data, inaccurate numbers, data kinds and anomalies, can be identified and resolved by the use of this technique.
- Exploratory data analysis (EDA) is an essential component in the process of obtaining valuable insights from data, as it uncovers significant statistical metrics such as the mean, median and standard deviation.

- The user's text does not contain any information to rewrite. Frequently, certain values exhibit substantial deviations from the established normative set of values; these outliers necessitate validation prior to the commencement of data analysis. If left unaddressed, these factors have the potential to significantly disrupt the analysis process, resulting in inaccuracies and misinterpretations. Therefore, a primary goal of exploratory data analysis (EDA) is to identify and detect outliers and abnormalities within the dataset.
- Exploratory Data Analysis (EDA) is a valuable technique that enables data scientists to get insights into the behaviour of variables when they are combined. By visualising and analysing the data, EDA assists in identifying patterns, correlations and interactions among these variables. The provided information proves to be valuable in the development of artificial intelligence models.
- Exploratory Data Analysis (EDA) facilitates the identification and removal of unnecessary columns, as well as the derivation of novel variables. Therefore, it can aid in the identification of the most significant features for predicting the target variable, facilitating the selection of features to be used in the modelling process.
- EDA can be utilised to determine suitable modelling strategies by considering the attributes of the data.

### 2.1.2 Processes in EDA

Proficiency in various tools and computer languages is essential for the execution of exploratory data analysis (EDA). In the given instance, an exploratory data analysis (EDA) would be conducted utilising the Python programming language within the open-source web-based tool known as Jupyter Notebook. The exploratory data analysis (EDA) process can be succinctly summarised into three sequential parts, namely:

- Understanding the data
- Cleaning the data
- Analysis of the relationship between variables

Let us understand the process of Exploratory Data Analysis (EDA) step-by-step:

#### Understanding the data:

- Import necessary libraries: The initial step is importing the necessary libraries. The code block employs the Pandas library for data reading and manipulation, while the Pandas-profiling module is utilised for exploratory data analysis (EDA). The Iris dataset is loaded using the datasets module from the scikit-learn toolkit.
- import pandas as pd
- import pandas\_profiling
- from sklearn import datasets

Subsequently, it is necessary to import the dataset. In this study, we will utilise the multivariate dataset known as the Iris dataset. `iris = datasets.load_iris()`.

#### The EDA cycle typically involves the following steps:

1. **Data Collection:** Obtain the raw data from various sources, such as databases, spreadsheets, APIs, or other data files.
2. **Data Cleaning:** Clean the data to ensure its accuracy and reliability. This step involves handling missing values, dealing with outliers, correcting inconsistencies, and converting data types if needed.

## Notes

3. **Data Exploration:** Explore the dataset's basic statistics, such as mean, median, standard deviation, and range. Create visualizations like histograms, scatter plots, and box plots to understand the distribution of data and relationships between variables.
4. **Feature Engineering:** If necessary, create new features (variables) from the existing data that might provide more meaningful insights or improve the performance of machine learning models. This could involve transformations, aggregation, or extracting relevant information from raw data.
5. **Hypothesis Generation:** Formulate hypotheses or questions about potential patterns or relationships in the data. These hypotheses guide further analysis and investigation.
6. **Data Visualization:** Create more advanced visualizations like heatmaps, correlation matrices, and time series plots to uncover deeper insights and relationships between variables.
7. **Statistical Analysis:** Perform statistical tests and calculations to validate or reject hypotheses. Common tests include t-tests, ANOVA, correlation analysis, and regression analysis.
8. **Interpretation and Insights:** Based on the results of the analysis, draw conclusions, and generate actionable insights. Determine if the initial hypotheses are supported by the data and identify any unexpected findings.
9. **Iteration and Refinement:** EDA is an iterative process. If new questions arise or if the initial analysis suggests further exploration is needed, you can refine your approach and perform additional analyses.
10. **Communication:** Present your findings and insights to stakeholders, team members, or clients. Visualization tools, reports, and presentations can help effectively communicate complex results.

Throughout the EDA cycle, it's important to maintain a balance between automated analysis and human intuition. While automated tools can assist with some steps, human intuition and domain knowledge play a significant role in formulating hypotheses and interpreting results.

### Converting to Pandas Data Frame:

The scikit-learn dataset is loaded into memory as a Bunch object, which bears resemblance to a dictionary data structure. In order to utilise the dataset in the context of Pandas, it is necessary to transform it into a Pandas DataFrame.

```
iris_data = pd.DataFrame(iris.data, columns=iris.feature_names) iris_data['target'] = iris['target']
```

### Checking data attributes:

Examining the characteristics of the data, such as its dimensions or the quantity of rows and columns inside the dataset, is consistently regarded as a prudent methodology. To assess the structure of the data, execute the provided code snippet.

```
iris_data.shape
```

The column names of the DataFrame can be inspected by accessing the columns attribute.

```
iris_data.columns
```

To inspect the initial records in a DataFrame, one may execute the following code if the dataset is of substantial size: The function `iris_data.head()` is used to display the first few rows of the `iris_data`

**Cleaning the data:**

After completing the process of scanning the attributes of the data, it is possible to perform any required alterations to the dataset, such as altering the names of the columns or rows. It is important to refrain from altering the variables of the dataset, as doing so can have a substantial impact on the ultimate outcome.

**Check for null values:**

In order to cleanse the data, it is imperative to initially examine the variables for any instances of null values. The presence of null values in any of the variables within a dataset has the potential to impact the outcomes of the analysis. When encountering missing data in a dataset, it is necessary to address this issue utilising various methods such as imputation, deletion of observations or variables, or employing models that are capable of handling missing data.

**Dropping the redundant data and removing outliers:**

Subsequently, in the event that any superfluous data is identified within the dataset, which does not contribute to the resultant output, it is also possible to eliminate such data from the table. All columns and rows within the iris dataset under consideration are deemed significant. The data would not be discarded. During this stage, it is imperative to identify any potential outliers within the dataset.

**Analysis of the relationship between variables:**

The ultimate stage in the process of exploratory data analysis (EDA) involves the examination and assessment of the interrelationships among variables. This encompasses the subsequent components:

**Correlation analysis:** The analyst does a computation of the correlation matrix among variables in order to ascertain the presence of strong correlations between them.

**Visualization:** The data analyst generates visual representations in order to investigate the interplay between factors. This encompasses many visual representations such as scatter plots, heatmaps and other similar graphical tools.

**Hypothesis testing:** The data analyst use visualisations to examine the interaction among various components. This comprises a variety of visual representations, including scatter plots, heatmaps and other comparable graphical tools.

**Types of EDA techniques:**

Various methodologies for exploratory data analysis can be employed to acquire a deeper understanding of the data. There are several prevalent forms of exploratory data analysis (EDA) that are frequently employed in various domains.

**Univariate non-graphical:**

Univariate non-graphical exploratory data analysis is a fundamental technique for reviewing data, wherein a single variable is used to analyse the information. Univariate non-graphical exploratory data analysis (EDA) entails the determination of the underlying distribution or pattern within the dataset, while providing objective information regarding the population.

The aforementioned technique encompasses the analysis of several characteristics of the distribution of a population, such as measures of dispersion, measures of central tendency, skewness and kurtosis.

## Notes

- The central tendency refers to the average or middle value of a distribution. One frequently used metric of central tendency is the mean, which is typically followed by the median and mode. The median is often favoured as a measure of central tendency in cases when the distribution is skewed or when there are concerns regarding the presence of outliers.
- The user's text does not contain any information to rewrite. The concept of spread refers to the extent to which the values of information deviate from the central tendency. The standard deviation and variance are two important measures of dispersion. The variance is calculated as the average of the squared differences between each data point and the mean, while the standard deviation serves as the fundamental measure from which the variance is derived.
- Skewness and kurtosis are two more useful univariate characteristics of the distribution. Skewness is a statistical measure that quantifies the degree of asymmetry in a distribution, whereas kurtosis is a measure that compares the peakedness of a distribution to a standard dispersion.

### Multivariate non-graphical:

Multivariate non-graphical exploratory data analysis (EDA) is a method employed to examine the association between many variables by means of cross-tabulation or statistical analysis. The utilisation of this technique seems to be advantageous in the identification of patterns and correlations among variables. This analytical approach proves to be particularly valuable in cases where a dataset contains several variables and there is a desire to examine their interrelationships.

Cross-tabulation is an advantageous expansion of tabulation specifically designed for the analysis of categorical data. Cross-tabulation is the preferred method of analysis when there are two variables being considered. In order to do this task, it is necessary to construct a bivariate table including two columns, each representing a distinct variable and two rows, each representing a different variable. Subsequently, proceed to populate the respective counts with all participants exhibiting identical pairs of levels.

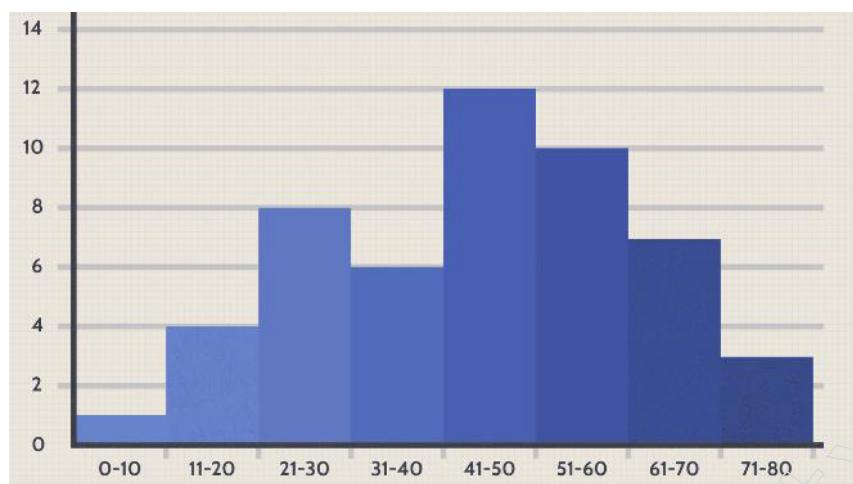
Statistics are generated for each level of every categorical variable, as well as for a single quantitative variable. These statistics are then compared across all categorical variables. The primary objective of multivariate non-graphical exploratory data analysis (EDA) is to discern and comprehend the associations between variables. The analysis of the interrelationships among variables enables the identification of patterns and trends that may not be readily apparent when considering individual variables in isolation.

### Univariate graphical:

A univariate graphical exploratory data analysis (EDA) technique utilises a range of graphs to acquire understanding of the distribution of a single variable. These graphical tools facilitate the rapid comprehension of the shapes, core trends, spreads, modalities, skewness and outliers present in the data under investigation. The subsequent methods are frequently employed in exploratory data analysis for univariate data visualisation.

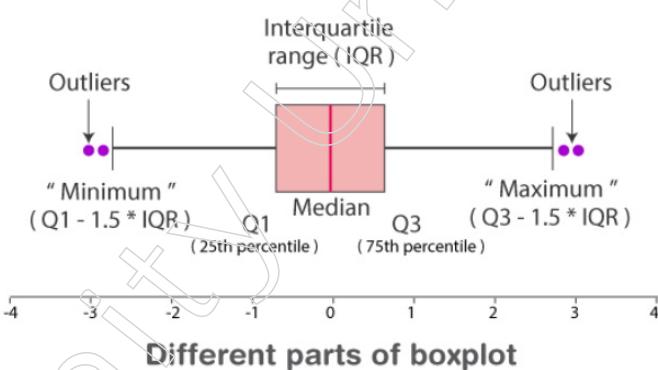
- Histogram:** The graph being referred to is a fundamental visualisation commonly employed in exploratory data analysis (EDA). A histogram is a graphical representation in the form of a bar plot that illustrates the frequency or proportion of occurrences within certain intervals, sometimes referred to as bins, for a given variable. The vertical dimension of each bar corresponds to the frequency or relative frequency of observations that lie within each interval. Histograms offer a perceptible understanding of the configuration and dispersion of the distribution, together with the identification of

any exceptional values.



## Notes

- Stem-and-leaf plots:** The stem-and-leaf plot is a graphical representation that serves as an alternative to a histogram, as it presents individual data values alongside their corresponding magnitudes. The stem-and-leaf plot is a graphical representation where the data values are divided into two parts: the stem, which represents the leading numbers and the leaf, which represents the trailing digits. The aforementioned plot facilitates the visual depiction of the distribution of the data, hence enabling the identification of characteristics such as symmetry and skewness.
- Boxplots:** Boxplots, alternatively referred to as box-and-whisker plots, offer a graphical representation that succinctly summarises the central tendency, dispersion and outliers of a distribution. The boxplot visually depicts the interquartile range (IQR) of the data, with the median line positioned within the box.



The whiskers of the boxplot are defined as the lines that extend from the lower and upper quartiles to the smallest and biggest observations, respectively, within a distance of 1.5 times the interquartile range (IQR) from the box. Outliers are data points that fall outside the range defined by the whiskers.

- Quantile-Normal Plots:** A quantile-normal plot, commonly referred to as a Q-Q plot, is a graphical tool used to evaluate the distribution of data by comparing the actual values against the expected values derived from a normal distribution. The Q-Q plot is a graphical tool used to compare the observed data with the quantiles of a normal distribution. If the data is regularly distributed, it is expected that the points will be aligned in a linear fashion. In the event that the data exhibits non-normality, the plot will effectively display any presence of skewness, kurtosis, or outliers.

## Notes

### Multivariate graphical:

A multivariate graphical exploratory data analysis (EDA) use visual representations to illustrate the links among two or more datasets. When exploring relationships among variables that extend beyond two, this methodology is employed to acquire a more comprehensive comprehension of the data. The utilisation of a grouped barplot is a prevalent approach in multivariate data visualisation, wherein each group corresponds to a distinct level of a variable and each bar signifies the respective quantity associated with it. Multivariate graphics can be visually depicted using various types of plots, including scatterplots, run charts, heat maps, multivariate charts and bubble charts.

- Scatterplots serve as visual depictions that illustrate the correlation between two variables of a quantitative nature. The process involves the representation of one variable on the horizontal x-axis and another variable on the vertical y-axis. In the context of the plot, it can be observed that each individual point serves as a representation of a specific observation. Scatterplots facilitate the identification of outliers or patterns within a dataset, as well as the determination of the direction and magnitude of the link between two variables.
- The user's text does not provide any information to be rewritten in an academic manner. A run chart is a graphical representation in the form of a line graph that illustrates the temporal variation of data. The tool in question is a straightforward yet potent instrument for the purpose of monitoring and analysing data alterations and observing patterns over time. Run charts are a valuable tool for identifying and analysing patterns, fluctuations and changes in a given process as it evolves over a period of time.
- A multivariate chart is a graphical representation that depicts the correlation between various factors and their corresponding responses. A multi-variable scatterplot is a graphical representation that illustrates the interrelationships among multiple variables concurrently. A multivariate chart is a graphical representation that illustrates the interrelationships between multiple variables, enabling the identification of patterns or clusters within the dataset.
- A bubble chart is a form of data visualisation that presents many circles, often known as bubbles, within a two-dimensional plot. The magnitude of each circular entity corresponds to the numerical representation of a tertiary variable. Bubble charts are frequently employed in the analysis of data sets containing three variables, as they offer a straightforward means of visually representing the interrelationships among these variables.

### Visualization techniques in EDA:

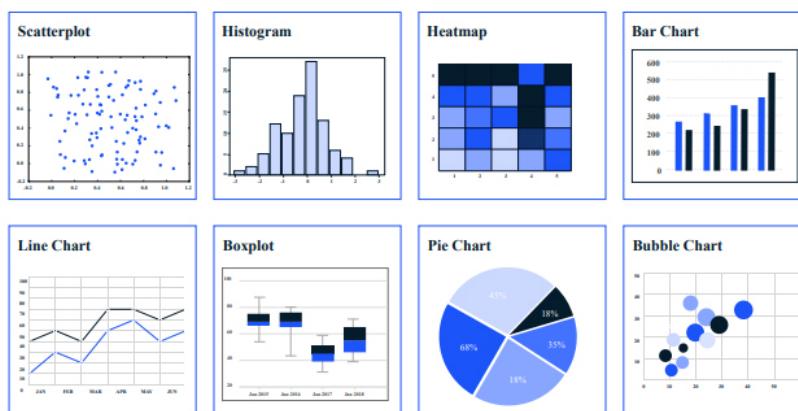


Figure: Visualization techniques in EDA

Visualisation techniques are crucial in exploratory data analysis (EDA) as they facilitate the visual exploration and comprehension of intricate data structures and linkages. Several visualisation techniques commonly employed in exploratory data analysis (EDA) include:

1. Histograms are visual depictions that illustrate the distribution of numerical values. Visualising the frequency distribution aids in comprehending the central tendency and dispersion of the data.
2. Boxplots, also known as box-and-whisker plots, are graphical representations that display the distribution of a quantitative variable. This visualisation technique facilitates the identification of outliers and provides insights into the distribution of the data by visually representing its quartiles.
3. Heatmaps are visual representations of data that utilise colours to depict different values. Visualisations are frequently employed to present intricate data sets, offering a convenient and expeditious means of visually representing patterns and trends within extensive quantities of data.
4. Bar charts are graphical representations that display the distribution of a categorical variable. The purpose of this tool is to visually represent the distribution of frequencies within a dataset, facilitating comprehension of the relative frequencies associated with each category.
5. Line charts are graphical representations that display the trajectory of a quantitative variable as it evolves over a specific period. The purpose of utilising this tool is to visually represent the temporal evolution of the data and to discern any discernible patterns or trends.
6. Pie charts are a type of graphical representation that effectively displays the relative distribution of a category data. The purpose of this tool is to visually represent the relative proportions of each category and facilitate comprehension of the distribution of the data.

#### **Exploratory Data Analysis Tools:**

##### **Spreadsheet software:**

Spreadsheet software, such as Microsoft Excel, Google Sheets, or Libre Office Calc, is frequently employed for exploratory data analysis (EDA) due to its user-friendly interface, straightforward functionality and rudimentary statistical analysis capabilities. These tools enable users to effectively organise, categorise, modify and conduct fundamental statistical analysis on data, such as computing measures of central tendency (mean and median) and variability (standard deviation).

##### **Statistical software:**

Specialised statistical software, such as R or Python, along with their diverse libraries and packages, provide a range of advanced statistical analysis capabilities, encompassing regression analysis, hypothesis testing and time series analysis. The software provides users with the capability to create personalised functions and do intricate statistical analysis on extensive datasets.

##### **Data visualization software:**

Visualisation software such as Tableau and Power BI empowers users to generate interactive and dynamic data visualisations. These technologies facilitate the identification of patterns and relationships within the data, so enabling users to make more informed decisions. In addition, a diverse range of charts and graphs are provided,

## Notes

including the capability to generate dashboards and reports. The programme facilitates the seamless sharing and dissemination of data, rendering it highly advantageous for collaborative endeavours and presentations.

### Programming languages:

Programming languages such as R, Python, Julia and MATLAB possess robust numerical computation capabilities and afford users access to a diverse range of statistical analytic tools. These programming languages possess the capability to create tailored functions to address specific analytical requirements, making them particularly advantageous in the context of handling extensive datasets. In addition to facilitating the automation of repetitive processes, they also provide increased flexibility in the management and manipulation of data.

### Business Intelligence (BI) tools:

SAP Business Objects, IBM Cognos, and Oracle BI offer data exploration, dashboards, and reports. Users can visually depict and analyse data from databases and spreadsheets. The company provides data preparation and quality management tools for data-driven decision-making in numerous corporate settings.

**Data mining tools:** KNIME, RapidMiner, and Weka do data preprocessing, clustering, classification, and association rule mining. These methods are useful for finding patterns and correlations in large datasets and building prediction models. Data mining is used in finance, healthcare, and retail.

**Cloud-based tools:** Cloud platforms like Google Cloud, Amazon Web Services (AWS), and Microsoft Azure offer several data analysis capabilities. These systems provide a flexible foundation for data storage and manipulation, as well as a variety of tools for analysis and visualisation. Cloud-based solutions are ideal for large and complex datasets due to their high-performance computing resources and flexibility to modify operations to project needs.

**Text analytics tools:** RapidMiner and SAS Text Analytics are used to analyse unstructured data like text documents and social media posts. Natural language processing (NLP) methods including sentiment analysis, entity recognition, and topic modelling extract useful information from textual input. Text analytics is used in marketing, customer service, and political analysis.

**Geographic Information System (GIS) tools:** GIS tools like ArcGIS and QGIS are used to analyse and visualise geospatial data. Users can map data and do spatial analysis, including finding patterns and trends in geographical data and running spatial queries. GIS tools are used in urban planning, environmental management, and transportation.

- Exploratory data analysis (EDA) is essential before any other data analysis jobs. This tool helps data scientists and analysts understand and gain insights from data. The eventual analysis must be checked for biases or errors by identifying missing or faulty data. Using data cleaning and preprocessing procedures during exploratory data analysis (EDA), analysts can ensure data accuracy and reliability.
- Exploratory Data Analysis (EDA) can help identify key features for machine learning models by easing feature selection. This procedure improves model performance. Exploratory data analysis (EDA) finds anomalies, patterns, and links in datasets. This competence can help organisations make educated judgements and acquire a competitive edge in the ever-changing technology landscape.

## 2.2 Data Types

John Tukey advocated for the promotion of exploratory data analysis (EDA) as a means to encourage statisticians to thoroughly investigate data, perhaps leading to the formulation of hypotheses that could prompt more data gathering and experimental endeavours. Exploratory Data Analysis (EDA) is primarily concerned with the meticulous examination of assumptions necessary for the process of model fitting and hypothesis testing. Additionally, it performs tests to handle missing values and applies necessary modifications to variables.

Exploratory Data Analysis (EDA) facilitates the development of a comprehensive comprehension of the data, as well as the identification of potential challenges related to either the information itself or the underlying processes. The methodology employed in this study adopts a scientific approach to elucidate the narrative embedded within the dataset.

### 2.2.1 Handling of Data Types

The initial stage in managing data types involves accurately identifying them. Common data types encompass various categories, such as numerical (including integer and float), categorical (such as string and enum), datetime and boolean.

**Converting Data Types:** On occasion, it is possible for the data to be stored in an incorrect format. An instance of a potential error in data storage involves the inadvertent conversion of numerical data into string format. In instances of this nature, it becomes imperative to do data type conversion in order to ascertain that the data is appropriately formatted for analytical purposes.

**Handling Missing Values:** The presence of missing data is a prevalent concern in databases. The approach to managing missing values may vary depending on the specific data type. When dealing with numerical data, it is common practise to substitute missing values with either the mean or median. Conversely, when working with categorical data, it is often appropriate to assign a distinct category to reflect missing values.

**Dealing with Categorical Data:** Categorical variables necessitate specific consideration. The categories can be classified as either nominal, which do not possess an intrinsic order, or ordinal, which are ordered. The numerical representation of categorical data can vary depending on its nature, with potential methods including one-hot encoding, label encoding, or ordinal encoding.

**Encoding Text Data:** In order to facilitate analysis, textual data necessitates preprocessing and transformation. Commonly employed methods in the field of natural language processing (NLP) involve the utilisation of techniques such as tokenization, stemming and vectorization (e.g., TF-IDF) to facilitate the conversion of textual data into a numerical representation.

**Handling DateTime Data:** The extraction of useful information, such as year, month, day, hour, etc., from DateTime data frequently necessitates the process of parsing. This facilitates the examination and categorization of data based on temporal factors.

**Normalization and Scaling:** In the context of numerical data, it is frequently imperative to employ normalisation or scaling techniques in order to guarantee that the values are uniformly distributed over comparable scales. Two commonly used strategies in data normalisation are Min-Max scaling and Z-score normalisation.

**Checking Data Integrity:** It is imperative to ensure that the data included within

## Notes

each column strictly conforms to its designated data type. An instance of this is when a column is designated as “numerical,” it is expected that it will exclusively contain values that are numeric in nature.

**Visualization:** Diverse data kinds require the utilisation of distinct visualisation strategies. For example, histograms, scatter plots and box plots are commonly employed to visually depict numerical data, whereas bar charts and pie charts are often preferred for the representation of categorical data.

**Handling Outliers:** The presence of outliers can have a substantial impact on exploratory data analysis (EDA). A comprehensive grasp of data types is crucial in the process of identifying the appropriate approach for dealing with outliers. For instance, the presence of outliers in categorical data may suggest the occurrence of errors, whereas in numerical data, they may represent genuine extreme values.

### 2.2.2 Univariate and Bivariate Analysis

#### Univariate Analysis:

The term “uni” refers to “one,” while “variate” pertains to “variable.” Consequently, in the context of univariate analysis, there exists a sole dependent variable. The primary aim of univariate analysis is to extract, delineate and summarise the data, as well as examine the inherent patterns within it. The dataset is analysed by examining each variable individually.

There are two types of variables that can be distinguished: categorical variables and numerical variables. Univariate analysis allows for the identification of several patterns, including measures of central tendency such as the mean, mode and median. Additionally, dispersion measures such as the range and variance, quartiles such as the interquartile range and the standard deviation can also be readily detected by this analysis.

#### Univariate data can be described through:

**Frequency Distribution Tables:** The frequency distribution table provides a representation of the frequency or number of occurrences of a particular event or observation within a given dataset. The provision of a concise overview of the data facilitates the identification of patterns.

#### Example:

The list of IQ scores is: 118, 139, 124, 125, 127, 128, 129, 130, 130, 133, 136, 138, 141, 142, 149, 130, 154.

IQ Range	Number
118-125	3
126-133	7
134-141	4
142-149	2
150-157	1

- Bar charts are often used for the purpose of comparing categories or groupings of data. They offer a straightforward visual representation that facilitates the comparison process. Tracking changes over time can be beneficial. Visualising discrete data is more effective.

- Histograms are a type of graphical representation that shares similarities with bar charts. They are utilised to visually depict categorical variables in relation to the corresponding data categories. Histograms visually represent categories as bins, which provide information about the frequency or count of data points within specific ranges. Visualising continuous data is more effective.
- Pie charts are mostly utilised for the purpose of comprehending the manner in which a collective entity is subdivided into smaller constituent parts. The entirety of the pie symbolises a complete set of 100 percent, while the individual slices of the pie serve to indicate the proportional magnitude of each respective category.
- The term “frequency” refers to the number of times an event or phenomenon occurs within a Polygons, such as frequency polygons, are employed in the realm of data analysis to facilitate the comparison of datasets or the visual representation of cumulative frequency distributions.

### Bivariate Analysis:

The term “bi” denotes the quantity of two, whereas “variate” refers to a variable. Consequently, in this context, there are two variables under consideration. The analysis pertains to causality and the interrelationship between the two variables. The field of bivariate analysis encompasses three distinct forms of analysis.

#### Bivariate Analysis of two Numerical Variables : (Numerical-Numerical)

A scatter plot is a graphical representation that displays individual data points as dots. These graphical representations facilitate the identification of potential relationships between two variables. The observed pattern provides information regarding the nature (linear or non-linear) and magnitude of the association between two variables.

Linear correlation refers to the degree of association between two numerical variables, indicating the strength of their linear link. In the absence of a correlation between the two variables, there is no discernible inclination for one variable to vary in conjunction with changes in the values of the second variable.

#### Bivariate Analysis of two categorical Variables: (Categorical-Categorical)

**Chi-square Test:** The chi-square test is employed to assess the relationship between categorical variables. The calculation is derived from the discrepancy between the anticipated frequencies and the observed frequencies within one or more categories of the frequency table. A probability value of zero signifies a state of complete dependence between two category variables, while a probability value of one shows that the two categorical variables are entirely independent.

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares (MS)	F
Within	$SSW = \sum_{j=1}^k \sum_{i=1}^l (X_i - \bar{X}_j)^2$	$df_w = k - 1$	$MSW = \frac{SSW}{df_w}$	$F = \frac{MSB}{MSW}$
Between	$SSB = \sum_{j=1}^k (\bar{X}_j - \bar{X})^2$	$df_b = n - k$	$MSB = \frac{SSB}{df_b}$	
Total	$SST = \sum_{j=1}^n (\bar{X}_j - \bar{X})^2$	$df_t = n - 1$		

Figure: Analysis of Variance

## Notes

**ANALYSIS OF VARIANCE (ANOVA):** The analysis of variance (ANOVA) test is employed to ascertain the presence of a statistically significant disparity among the means of multiple groups that are distinct from one another. This analytical approach is suitable for doing a comparative analysis of the mean values of a quantitative variable across multiple categories of a qualitative variable.

### 2.2.3 Visualising the Data

The ability to visually represent and analyse data is increasingly recognised as a crucial competency in contemporary business environments. The utilisation of data visualisations and dashboards enables executives to efficiently make informed decisions by promptly accessing the necessary information. This constitutes a vital aspect in the development of business intelligence, which has subsequently evolved into the field commonly referred to as business analytics.

Data visualisation refers to the visual depiction of information and data. Data visualisation tools utilise visual components like as charts, graphs and maps to facilitate the comprehension of trends, outliers and patterns within data, hence offering a user-friendly means of analysis. In the realm of big data, the utilisation of data visualisation tools and technologies is vital for the examination of vast quantities of information and the formulation of decisions based on data.

Data visualisation refers to the visual depiction of data and information through the utilisation of various visual tools such as charts, graphs, maps and similar visual aids. The utilisation of visualisations facilitates the comprehension of patterns, trends and outliers within a given data set. Data visualisation is a method of presenting data in a form that is easily understandable to both the general public and specific audiences who may lack technical expertise. As an illustration, the health agency inside a governmental entity may furnish a cartographic representation of locations that have undergone vaccination. The primary objective of data visualisation is to facilitate the process of informed decision-making by presenting data in a visually appealing manner, hence enhancing its interpretability and significance. Benefits of data visualization: Data visualisation is a multifaceted instrument that is utilised in several sectors, including but not limited to public policy, finance, marketing, retail, education, sports, history and numerous other disciplines. There are several advantages associated with data visualisation.

- **Storytelling:** The aesthetic elements of colours and patterns found in various forms of expression, such as clothing, artwork, culture and architecture, possess a captivating allure for individuals. The narrative encapsulated within the data can also be represented through the utilisation of hues and patterns.
- **Accessibility:** Data is disseminated in a manner that is comprehensible and easily accessible to a diverse array of individuals.
- **VisualiseRelationships:** The utilisation of graphs or charts facilitates the identification of connections and patterns within a dataset.
- **Exploration:** The increased availability of data has expanded the possibilities for research, collaboration and the formulation of practical judgements.

Types of data visualization: The process of data visualisation can range from utilising basic graphical representations such as bar graphs or scatter plots to more sophisticated techniques that yield enhanced analytical capabilities. The following are few prevalent categories of data visualisations:

- **Table:** Data is displayed in rows and columns in a table, which is a simple task to accomplish in an Excel or Word document.

- **Chart or graph:** Data are commonly depicted on a two-dimensional coordinate system, consisting of an x-axis and a y-axis. This graphical representation employs various visual elements such as bars, points, or lines to effectively convey and compare data. The data is presented in a tabular manner. An infographic is a specialised form of graphical representation that combines textual information and visual elements to effectively convey factual data.
  - ❖ **Gantt chart:** A bar chart known as a Gantt chart, which is used primarily in project management, shows a timeline and tasks.
  - ❖ **Pie chart:** The percentages that make up the “slices” of a pie in a pie chart all add up to 100%.
- **Geospatial visualization:** Data is visually represented in the form of maps, utilising various shapes and colours to effectively convey the correlation between distinct geographical areas. This can be achieved through the utilisation of techniques such as choropleth or heat maps.
- **Dashboard:** Data and visualisations are commonly exhibited, typically for the goal of facilitating comprehension and presentation of data by analysts, particularly in the context of business operations.
- **Data visualization examples:** Data visualisation technologies can be utilised to generate various types of charts and graphs in order to visually represent significant data. The following examples illustrate instances of data visualisation in practical contexts:
- **Data science:** Data scientists and researchers are equipped with programming languages or tools, such as Python or R, that provide them with access to libraries. These libraries enable them to analyse data sets, discern trends and gain insights. Tools assist data professionals in enhancing their efficiency through the utilisation of various visual elements such as colour coding, graphs, lines and shapes to facilitate the organisation and analysis of research data.
- **Marketing:** The utilisation of tracking data, such as web traffic and social media analytics, enables marketers to conduct an analysis of customer behaviour in terms of product discovery and adoption patterns, distinguishing between early adopters and laggard buyers. Charts and graphs serve as effective tools for synthesising data, enabling marketers and stakeholders to gain a comprehensive understanding of prevailing trends.
- **Finance:** Investors and advisors that specialise in the trading of stocks, bonds, dividends and other commodities engage in the analysis of price fluctuations over time in order to assess the investment potential of these assets for both short-term and long-term durations. Line graphs are utilised by financial analysts to visually represent data, allowing them to switch between different time periods such as months, years and even decades.
- **Health policy:** Choropleth maps, characterised by color-coded divisions based on geographical areas such as nations, states, or continents, can be employed by policymakers. One potential application of these maps is to illustrate the variations in cancer or Ebola fatality rates throughout different regions of the globe.

### 2.3 Feature Extraction

The process of feature extraction is conducted, resulting in three distinct outcomes for each categorization, which are subsequently compared with one another. The optimal feature extraction approach is selected for each of the two categories and subsequently

## Notes

integrated into the entirety of the system. The utilised methods encompass the histogram of oriented gradients (HOG), features derived from the discrete cosine transform and features obtained from a pre-trained convolutional neural network. Various feature extraction approaches possess distinct advantages, hence justifying their selection.

The Histogram of Oriented Gradients (HOG) method is frequently employed in the field of object detection. This approach relies on the utilisation of gradient information to effectively characterise images. The utilisation of gradients in image analysis allows for the extraction of valuable information pertaining to edges and corners. Consequently, the Histogram of Oriented Gradients (HOG) method is highly advantageous in accurately defining the content present inside a picture.

The selection of the method for extracting features from the Discrete Cosine Transform (DCT) domain is based on the fact that the resulting features are specifically designed to characterise various quality parameters in a picture. The last approach involves utilising characteristics derived from a Convolutional Neural Network (CNN). This CNN is trained on a substantial collection of photos for the purpose of object recognition.

The objective is to enable the network to generalise its learning to other tasks and datasets that it has not specifically been trained on. The selection of this particular strategy is based on its shown efficacy in effectively addressing a wide range of general tasks.

### 2.3.1 Feature Extraction: Part I

The process of extracting features plays a pivotal role in numerous machine learning and data analysis endeavours. The process entails converting unprocessed data into a condensed and significant format, hence enhancing the efficiency of machine learning algorithms. The objective of feature extraction is to effectively capture pertinent information from the input data, while simultaneously lowering its dimensionality and minimising the impact of noise.

**Principal Component Analysis (PCA):** main Component Analysis (PCA) is a widely used method in the field of dimensionality reduction. Its primary objective is to identify the main components, which are orthogonal directions, that account for the highest amount of variance present in the dataset. The major components can be expressed as linear combinations of the original attributes.

By utilising the process of picking the highest N principal components, it is possible to significantly decrease the dimensionality of the data while still preserving a substantial portion of the crucial information.

**Independent Component Analysis (ICA):** Independent Component Analysis (ICA) is a dimensionality reduction technique that aims to identify statistically independent components from a given set of mixed data. Blind source separation challenges, such as the separation of mixed audio sources, frequently employ this technique.

The utilisation of Independent Component Analysis (ICA) is a statistical and computational method commonly employed in the field of machine learning. Its purpose is to disentangle a multivariate signal by isolating its independent non-Gaussian components. The Independent Component Analysis (ICA) framework operates under the assumption that the observed data can be represented as a linear combination of independent signals that are not distributed according to a Gaussian distribution.

The objective of Independent Component Analysis (ICA) is to identify a linear mapping of the data that yields a collection of components that are statistically independent from each other.

- Independent Component Analysis (ICA) is a robust and versatile technique that finds extensive use in various domains, including signal processing, image analysis and data reduction. The utilisation of Independent Component Analysis (ICA) has been observed across various disciplines, encompassing finance, biology and neuroscience.
- The fundamental concept underlying Independent Component Analysis (ICA) is to ascertain a collection of fundamental functions that can be employed to depict the data that has been observed. The selection of these basis functions is made such that they exhibit statistical independence and non-Gaussian characteristics. Once the identification of these basis functions has been accomplished, they can be employed to partition the observed data into its distinct and autonomous components.
- Independent Component Analysis (ICA) is frequently employed in tandem with many machine learning methodologies, including clustering and classification. For instance, Independent Component Analysis (ICA) can be employed as a preprocessing technique for data prior to conducting clustering or classification. Alternatively, ICA can also be utilised to extract features that are then employed in these aforementioned tasks.

Independent Component Analysis (ICA) possesses several constraints, which encompass the assumption that the underlying sources exhibit non-Gaussian distribution and are blended in a linear manner. Moreover, Independent Component Analysis (ICA) can impose a significant computing burden and may encounter convergence challenges if the data is not adequately pre-processed. Notwithstanding these constraints, Independent Component Analysis (ICA) continues to be a robust and extensively employed methodology in the fields of machine learning and signal processing.

#### Advantages of Independent Component Analysis (ICA):

- Ability to separate mixed signals:** Independent Component Analysis (ICA) is a robust and effective technique utilised for the purpose of decomposing mixed signals into their constituent independent components. This phenomenon has proven to be beneficial in a diverse range of applications, including but not limited to signal processing, image analysis and data compression.
- Non-parametric approach:** Independent Component Analysis (ICA) is classified as a non-parametric methodology due to its characteristic of not relying on assumptions on the underlying probability distribution of the data.
- Unsupervised learning:** Independent Component Analysis (ICA) is a form of unsupervised learning methodology, wherein it is capable of being employed on data sets without the requirement of annotated instances. This attribute renders it advantageous in scenarios where annotated data is unavailable.
- Feature extraction:** Independent Component Analysis (ICA) possesses the capability to do feature extraction, hence enabling the identification of significant characteristics within the dataset. These extracted features can subsequently be employed for many purposes, including classification tasks.

#### Disadvantages of Independent Component Analysis (ICA):

- The assumption made by ICA is that the underlying sources are non-Gaussian, although this assumption may not always hold true. If the underlying sources exhibit a Gaussian distribution, independent component analysis (ICA) may not yield satisfactory results.

## Notes

- The linear mixing assumption posits that Independent Component Analysis (ICA) considers the sources to be mixed in a linear manner, although this assumption may not always hold true. In cases where the sources are combined in a nonlinear manner, the effectiveness of Independent Component Analysis (ICA) may be compromised.
- ICA can incur significant computing costs, particularly when applied to datasets of considerable size. The application of ICA to real-world problems might be challenging.
- One potential challenge in utilising Independent Component Analysis (ICA) is the occurrence of convergence issues, wherein the algorithm may encounter difficulties in reaching a solution. Consequently, there may be instances where ICA fails to identify a suitable solution. Complex datasets with numerous sources can be a significant challenge in various contexts.

### Restrictions on ICA

- The assumption used in ICA is that the independent components produced are statistically independent from one another.
- The independent components produced by the Independent Component Analysis (ICA) must exhibit a non-Gaussian distribution.
- The user's text does not provide any information to rewrite in an academic manner. The quantity of autonomous elements produced by the Independent Component Analysis (ICA) is equivalent to the quantity of observed mixtures.

### Difference between PCA and ICA are as follows:

Principal Component Analysis	Independent Component Analysis
It reduces the dimensions to avoid the problem of overfitting.	It decomposes the mixed signal into its independent sources' signals.
It deals with the Principal Components.	It deals with the Independent Components.
It focuses on maximizing the variance.	It doesn't focus on the issue of variance among the data points.
It focuses on the mutual orthogonality property of the principal components.	It doesn't focus on the mutual orthogonality of the components.
It doesn't focus on the mutual independence of the components.	It focuses on the mutual independence of the components.

**Autoencoders:** Autoencoders are a class of neural networks that possess the capability to do unsupervised learning and reduce the dimensionality of data. Autoencoders are composed of two main components: an encoder and a decoder. The encoder is responsible for mapping the input data to a representation with lower dimensions. On the other hand, the decoder is responsible for reconstructing the original data from the lower-dimensional representation. The primary objective of autoencoders is to acquire a condensed and semantically significant depiction of the input data.

**Feature Selection:** The process of feature selection entails the identification and selection of a subset of the original features, taking into consideration their significance or pertinence to the given job. This procedure has the potential to mitigate overfitting, enhance the interpretability of the model and expedite the training process. Several commonly used approaches for feature selection in academic research include univariate feature selection, recursive feature elimination and L1 regularisation.

**Bag of Words (BoW):** The Bag-of-Words (BoW) methodology is a widely employed method for extracting features in the field of natural language processing (NLP). The representation of a text document is achieved by considering it as an unordered set, sometimes referred to as a bag, of its constituent words, without taking into account grammatical structure or the sequence in which the words appear.

The frequency or occurrence of each word in the document is considered as a feature and the resultant feature vector represents the content of the document.

**Word Embeddings:** Word embeddings refer to compact vector representations of words, which are typically acquired by the utilisation of methodologies such as Word2Vec, GloVe, or fastText, by training on extensive collections of texts. These word embeddings effectively capture semantic links, enabling natural language processing (NLP) models to comprehend word context and meaning more efficiently compared to the Bag-of-Words (BoW) approach.

**Histogram of Oriented Gradients (HOG):** The Histogram of Oriented Gradients (HOG) method is a widely employed methodology for extracting features in computer vision applications, particularly in the context of object detection. The method assesses the directional changes in pixel intensities inside a picture and subsequently depicts these changes as histograms. The feature vector obtained represents both the shape and edge characteristics of the item.

### 2.3.2 Feature Extraction: Part II

In the subsequent section of Feature Extraction, further investigation will be conducted into additional feature extraction methodologies that are frequently employed in diverse machine learning and data analysis endeavours.

- Word2Vec, which embeds words as vectors in a continuous space, is popular. A shallow neural network uses contextual usage in a large text corpus to represent words. The word embeddings capture semantic relationships between words, allowing models to understand parallels and analogies.
- Global Vectors for Word Representation (GloVe) is a word embedding method that factors the word co-occurrence matrix from a large textual corpus. Similar to Word2Vec, GloVe generates word embeddings that capture semantic relationships and are useful in NLP applications.
- TF-IDF is a popular text analysis approach. The metric measures the importance of a term in a text compared to a corpus of many texts. High-frequency terms in a text but low-frequency terms in the corpus have high TF-IDF scores. This approach vectorizes textual data, focusing on document-specific terms.
- Mel-Frequency Cepstral coefficients (MFCC) are commonly used in speech and audio processing for feature extraction. The discrete cosine transform (DCT) is applied to the log-mel-spectrogram after extracting the short-term power spectrum from an audio source. The Mel-frequency cepstral coefficients (MFCCs) represent the fundamental features of the audio signal, making them useful for speech recognition and audio classification.
- A popular computer vision method for extracting unique features from images is the Scale-Invariant Feature Transform (SIFT). Its main uses are object detection and image matching. The system uses gradient orientation histograms to identify visible components in an image and describe their local appearance. SIFT characteristics are resistant to scale, rotation, and illumination conditions, making them ideal for image analysis applications.

## Notes

- A typical object detection method in computer vision is the Histogram of Oriented Gradients (HOG). The distribution of local gradient directions in an image is computed to provide a feature vector that accurately characterises the image's structural and form features.
- The topic is colorhistograms. Colour histograms are used to show a picture's colour frequency. The number of pixels in each bin is counted in the colour space. Colour histograms are useful in image retrieval and object recognition because they use colour information.
- The Bag of Visual Words (BoVW) computer vision method extends the Bag of Words approach. Image representation is achieved by grouping it into visual words. These words are obtained by analysing the image's keypoints and local appearance descriptors like SIFT or HOG. Picture categorization and object recognition are common computer vision problems using the Bag-of-Visual-Words (BoVW) technique.
- Academic research focuses on VGG and ResNet features. VGG and ResNet are popular deep convolutional neural network designs. Instead of using the network for classification, one can use intermediary layer activations (features) to represent the input image. The high-level attributes above capture the image's visual content and can be fed to other machine learning models. It is vital to remember that the data and problem determine the feature extraction technique. The tasks and data kinds affect the efficiency of different techniques. Thus, experimentation and data analysis are essential for making educated decisions.
- Therefore, it is crucial to engage in experimentation and develop a comprehensive understanding of the data in order to make well-informed decisions.

### Feature Extraction Process:

Data can manifest in several formats, including textual, numerical, visual (pictures) and audiovisual (videos). An instance of this scenario involves a client details form in which certain fields have not been completed and are left blank. The term "missing data" is used to refer to such data. In the majority of instances, data can consist of missing data, unstructured data, or data that exhibits irregular organisation. In the field of data visualisation, it is necessary to do data cleaning prior to data processing in order to ensure its suitability for future analysis.

Data cleansing has been a longstanding practise in the field of databases and is an essential component of the extract, transform, load (ETL) process, frequently employed in data warehouses. This process involves extracting data from one or multiple sources, transforming it into the appropriate format and structure, which includes the cleansing of the data and ultimately loading it into a designated destination, such as a singular database or file. This refined dataset can then be utilised for :

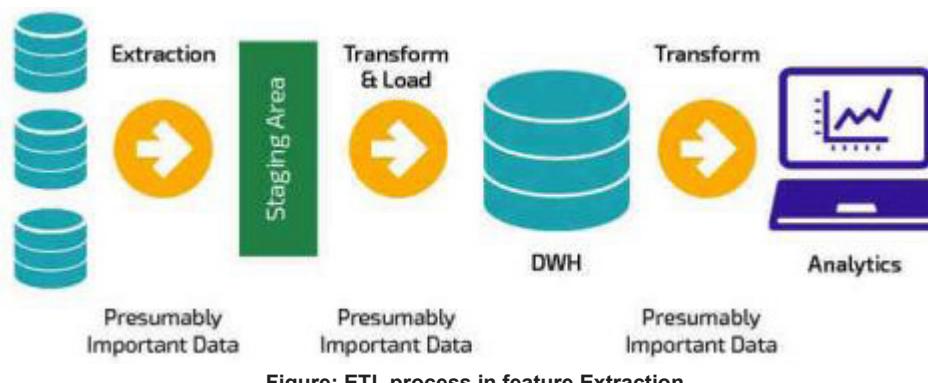


Figure: ETL process in feature Extraction

### Extraction, Transformation and Load (ETL):

#### Extraction:

The initial stage of the Extract, Transform, Load (ETL) process involves the extraction of data. During this stage, data is extracted from several source systems, which may exist in different formats such as relational databases, NoSQL databases, XML files and flat files. The extracted data is then loaded into the staging area. The extraction of data from diverse source systems and its subsequent storage in the staging area prior to being loaded into the data warehouse is a crucial step. This is primarily due to the fact that the extracted data exists in multiple formats and is susceptible to corruption. Therefore, the direct loading of data into the data warehouse may result in potential harm, making the process of rollback significantly more challenging. Hence, this stage holds significant importance inside the ETL process.

#### Transformation:

Transformation is the second stage of the ETL process. In this step, the extracted data is subjected to a set of rules or functions to transform it into a single standard format. It might entail the following procedures or tasks:

- **Filtering** – The process of selectively incorporating specific attributes into the data warehouse.
- **Cleaning** – The process involves replacing the missing variables with predetermined default values, as well as standardising the representation of geographical references such as U.S.A, United States and America to the standardised abbreviation USA.
- **Joining** – consolidating several qualities into one.
- **Splitting** – The process of dividing a singular attribute into numerous attributes.
- **Sorting** – The process of arranging tuples based on a specific attribute, typically referred to as the key attribute.

#### Loading:

The loading phase is the third and ultimate step inside the Extract, Transform, Load (ETL) process. During this stage, the data that has been transformed is ultimately loaded into the data warehouse. The frequency at which data is updated in the data warehouse varies, with some instances involving frequent updates and others occurring at regular intervals over longer periods of time. The rate and duration of loading are contingent upon the specific requirements of each system, hence exhibiting variability.

## 2.4 Dimensionality Reduction

Dimensionality refers to the quantity of input features, variables, or columns that are present inside a certain dataset. The act of reducing these features is commonly referred to as dimensionality reduction. The dataset encompasses a substantial quantity of input information across diverse scenarios, hence rendering the predictive modelling assignment more intricate.

In situations when the training dataset contains a large number of features, it becomes challenging to visualise or make predictions. Consequently, the utilisation of dimensionality reduction techniques becomes necessary in order to address this issue.

## Notes

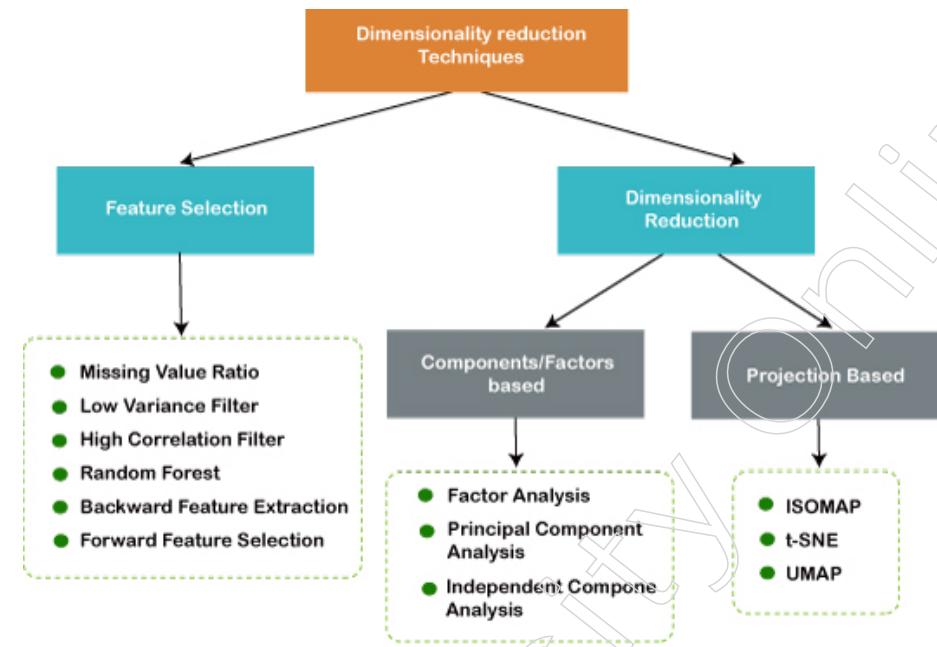


Figure: Dimensionality Reduction

Dimensionality reduction techniques can be defined as methods used to transform datasets with higher dimensions into datasets with less dimensions, while ensuring that the resulting dataset retains similar information. These strategies are commonly employed to improve the accuracy of predictive models while addressing classification and regression challenges. High-dimensional data is frequently utilised in various domains, including but not limited to speech recognition, signal processing and bioinformatics. Additionally, it has the potential to be utilised for purposes such as data visualisation, noise reduction and cluster analysis.

### 2.4.1 Need for Dimensionality Reduction

Dimensionality reduction is a technique employed in the field of Exploratory Data Analysis (EDA) that facilitates the rapid visualisation of data with a high number of dimensions. Additionally, it offers the potential to uncover concealed systematic patterns within a given dataset. Linear dimensionality reduction approaches, such as Principal Component Analysis (PCA), are often regarded as the benchmark in various domains of data science. However, their effectiveness appears to be limited when it comes to analysing non-linear data with large dimensions. In this particular scenario, the utilisation of non-linear dimensionality reduction techniques such as t-distributed Neighbour Embedding and Uniform Manifold Approximation and Projection (UMAP) has gained significant popularity. These methods have demonstrated exceptional efficacy in the analysis of high-dimensional data, thereby establishing themselves as cutting-edge approaches in this field.

Some needs of applying dimensionality reduction technique to the given dataset are given below:

- The reduction of feature dimensions leads to a corresponding decrease in the amount of space needed to store the dataset.
- Reduced feature dimensions result in decreased computation training time.
- The reduction of feature dimensions in the dataset facilitates efficient data visualisation.

- The process of removing redundant features, if they exist, is accomplished by addressing the issue of multicollinearity.
- Dimensionality reduction refers to the procedure of decreasing the number of features present in a dataset, with the objective of preserving a substantial amount of information.
- One potential approach to address the aforementioned objectives is simplifying the model, enhancing the efficacy of the learning algorithm, or facilitating the visualisation of the data.
- Various methods for reducing dimensionality are available, including principal component analysis (PCA), singular value decomposition (SVD) and linear discriminant analysis (LDA).
- Every method employed in this study projects the dataset onto a space with fewer dimensions, while ensuring the retention of crucial information.
- Dimensionality reduction is commonly conducted as part of the pre-processing phase prior to constructing a model, with the aim of enhancing performance.
- It is imperative to acknowledge that dimensionality reduction approaches have the potential to eliminate valuable information, thus necessitating caution during their application.

#### Methods of Dimensionality Reduction:

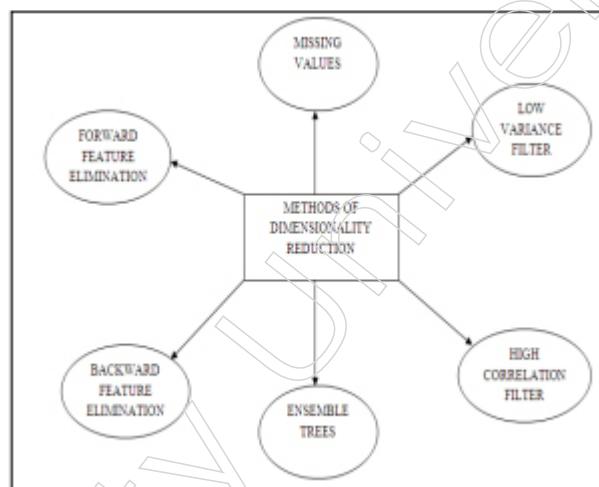


Figure: Methods of Dimensionality Reduction

#### Missing values:

The data column exhibiting a substantial number of missing values is computed and subsequently eliminated. The calculation of missing values is performed using either a statistical node or a group node. Statistical nodes are employed for data analysis, while the group node is utilised in the methodologies of dimensionality reduction. The gathered missing values that exceed the specified threshold are subsequently eliminated. If the determined threshold value is greater than the initial values, the reduction will be characterised by a high level of intensity.

The calculation of missing data is performed using the formula depicted in Figure. The KNIME tool's missing value node computes the ratio of missing values by dividing the number of missing values by the total number of rows.

## Notes

RATIO OF MISSING VALUES=NUMBER OF  
MISSING VALUES/TOTAL NUMBER OF  
ROWS

Figure: Formulae for Missing Values

### Low Variance Filter:

The low variance filter is a method used to assess the variability of data in order to determine the amount of information contained within a given data column. In the scenario where the column cells converge to a constant value, the variance would approach zero, rendering the column ineffective in distinguishing between distinct data groups. The data column with values below the specified threshold are filtered and thereafter eliminated. The data undergoes a filtering process in order to identify instances of low variance. The initial step in the data normalisation process involves applying normalisation techniques to all columns in order to identify and reduce the presence of lower variance.

### High Correlation Filter:

The input features frequently exhibit correlation, indicating that they are interdependent and contain redundant information. A data column that has a strong correlation with another data column is unlikely to contribute significantly novel information to the existing set of input attributes. In order to decrease the presence of correlated columns, the linear correlation node is utilised to quantify the correlation between pairs of columns. In the event that a strongly correlated column exists within a pair, the corresponding data is eliminated. The correlation filter utilises a correlation matrix as its input. The process of filtering highly connected data columns necessitates the presence of homogeneous data ranges, which can be achieved by utilising a Normalise node.

### Ensemble Trees

Ensemble trees, also referred to as random forests, are a popular machine learning technique. In order to achieve optimal classification performance, the process of feature selection is conducted. One method for reducing dimensionality involves creating a comprehensive set of trees based on a target characteristic, then subsequently utilising the usage statistics of each variable to identify the most informative subset of features. The score is determined based on the applicants' level and serves to indicate the relative traits that are most indicative. An ensemble of shallow trees has been built, where each tree is trained on a subset of all characteristics. If the specific attribute is chosen as the optimal split, the informative features will be preserved.

### Backward Feature Elimination:

The Backward Feature Elimination procedure involves reducing the dimensionality of a dataset by iteratively removing features in the context of a certain machine learning method. The proposed approach utilises an iterative process, wherein the chosen classification algorithm is executed for a certain number of input features. Next, one input feature is excluded and the model is trained using  $(n-1)$  input features repeatedly. The input feature exhibiting a high error rate is identified and subsequently removed from the dataset. Subsequently, the aforementioned steps are iteratively performed for a range of values, specifically  $n-2$ ,  $n-3$  and so forth, in order to ascertain the presence of a higher mistake rate. The Backward Feature Elimination Filter provides a visual representation of

the number of features retained at each iteration, along with the accompanying error rate. This methodology is only applicable to a limited dataset size.

## Notes

### Forward feature elimination:

Forward feature elimination is a method that is similar to backward feature elimination. In this approach, a forward feature construction loop is used to create a set of classifiers. These classifiers are built by incrementally adding input features. The forward feature loop commences with the inclusion of a single feature and thereafter incorporates an additional feature on each iteration. Both the forward and backward processes are characterised by high computing costs and significant expenses. By executing the optimisation loop, the optimal cutoff values were identified for each of the six dimensionality reduction techniques, as well as for the most effective model, based on criteria such as the lowest number of columns and highest accuracy.

**Techniques in Dimensionality Reduction:** Dimensionality reduction encompasses two distinct primary approaches. Two common techniques in machine learning are feature selection and feature reduction.

**Feature Selection:** In the fields of machine learning and statistics, feature selection (FS) is a crucial process that involves the selection of a subset of pertinent features (also referred to as variables, predictors) for the purpose of constructing a model. This process is also known as variable selection, attribute selection, or variable subset selection. This technique encompasses three distinct rationales: The interpretation of models is to enhance their efficiency and simplicity for users, save training time and minimise variation to improve generalisation capabilities. Feature selection is a technique that is commonly employed in domains characterised by a high number of features.

It involves the selection of specific features based on their relevance to the target function. A feature selection algorithm can be conceptualised as the amalgamation of a search methodology for generating novel feature subsets, coupled with an assessment metric that assigns scores to the various feature subsets. Each subset of the feature is individually examined in order to minimise the error or noise rate. Feature selection can be categorised into three distinct kinds, namely embedding, filter and wrapper methods. The categorization of feature selection methods is depicted in figure below.

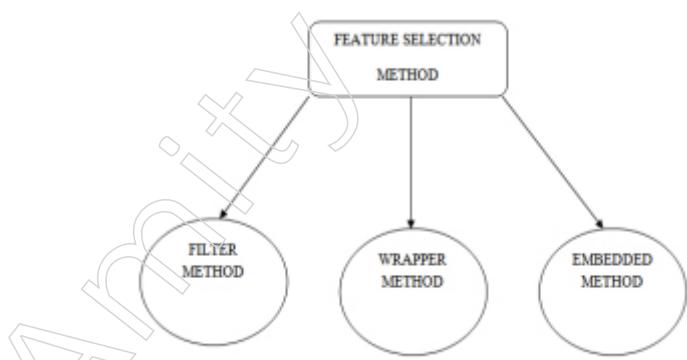


Figure : Feature Selection Methods

### Filter Method

The selection of variables in the filter approach is determined by the model employed. These predictions are solely derived from overarching characteristics, such as the level of correlation with the variable being predicted. The filter approach encompasses a limited number of variables that are deemed to be of particular interest.

## Notes

Additional variables will be employed for the purposes of either classification or regression models. These solutions demonstrate high efficiency while minimising time consumption. The aforementioned method is exclusively employed as a preprocessing technique, wherein the selection of redundant variables is made without considering their link with other variables. In this study, all the attributes are considered and subsequently, the optimal subsets are selected. These selected subsets are then subjected to various algorithms in order to evaluate their performance.

### Wrapper Method

In contrast to the filter approach, the wrapper method determines the subset of variables that account for the interaction between the variables. There exist two drawbacks associated with the wrapper approach, namely:

- The number of observations are inadequate when the over fittings are increased.
- Larger the variable, the computation time is increased.

The subset feature is formed by the use of the search strategy and subsequently, the performance is assessed. The wrapper approach employs a search algorithm to determine the optimal subset. Alternative search-based strategies utilise targeted projection pursuit, a method that identifies low-dimensional projections of the data with high scores. These projections highlight the features that exhibit the greatest projections in the lower-dimensional space, thereby enabling their selection.

**Embedded Method:** The utilisation of embedded approaches has been suggested as a means to mitigate the complexity associated with learning classification tasks. The formation of the embedded class involves the utilisation of filter and wrapper classes, which offer various benefits. Several embedded feature selection methods have been introduced in recent years; however, a unified theoretical framework has not yet been constructed. The learning algorithm has its own variable selection mechanism.

**Feature Reduction:** Feature reduction is a technique used to transform data from a high-dimensional space to a lower-dimensional space. The process of data transformation can be carried out using either a linear or non-linear approach. The primary method employed for reducing dimensionality is a linear strategy that involves mapping the data onto a lower-dimensional space.

This mapping is designed to maximise the variance of the data in the resulting low-dimensional representation. To diminish the number of features, the correlation coefficient is computed by the determination of the Eigen vectors and Eigen values. In this model, the selection of relevant features for the class is achieved through the elimination of redundant features.

Initially, the input consists of features, which are subsequently subjected to a process of feature extraction. This extraction involves the identification and removal of superfluous features, as well as the extraction of comparable features through the utilisation of algorithms such as PCA, SVD, LDA and others.

The filters employ criteria that do not rely on machine learning, such as a relevance index derived from correlation coefficients or test statistics. Subsequently, the filters are employed for the purpose of approximating the value of the feature.

### 2.4.2 Introduction of Dimensionality Reduction

Dimensionality reduction refers to the procedure of decreasing the quantity of features, also known as dimensions, inside a dataset, while striving to preserve the

maximum amount of pertinent information. There are several motivations for undertaking this task, including the simplification of a model, the enhancement of a learning algorithm's performance and the facilitation of data visualisation.

Various methods exist for reducing the dimensionality of data, such as principal component analysis (PCA), singular value decomposition (SVD) and linear discriminant analysis (LDA). Each strategy employs a distinct methodology to project the data onto an area with fewer dimensions, while simultaneously retaining crucial information.

Dimensionality reduction is a widely employed method in data analysis that aims to decrease the number of features present in a dataset, while simultaneously endeavouring to preserve the most significant information. In essence, the procedure entails the conversion of data with a high number of dimensions into a space with less dimensions while retaining the fundamental characteristics of the initial data.

Within the field of machine learning, the term "high-dimensional data" pertains to datasets that possess a substantial quantity of features or variables. The phenomenon known as the curse of dimensionality is a prevalent issue encountered in the field of machine learning, wherein the efficacy of a model diminishes with the escalation of the number of characteristics.

The reason for this phenomenon is because as the number of characteristics increases, the model's complexity also increases, hence making it more challenging to identify an optimal solution. Moreover, the presence of high-dimensional data can also result in overfitting, a phenomenon characterised by the model excessively conforming to the training data and subsequently failing to effectively generalise to novel data instances. The utilisation of dimensionality reduction techniques can effectively address these issues by decreasing the intricacy of the model and enhancing its ability to generalise. Dimensionality reduction encompasses two primary methodologies: feature selection and feature extraction.

**Feature Selection:** The process of feature selection entails the identification and selection of a subset of the original features that exhibit the highest degree of relevance to the specific problem being addressed. The objective is to decrease the number of dimensions in the dataset while preserving the essential characteristics. Feature selection encompasses several techniques, such as filter methods, wrapper methods and embedding methods. Filter techniques include ranking the features according to their relevance to the target variable. Wrapper methods, on the other hand, utilise the model's performance as the criterion for selecting features. Embedded methods, meanwhile, integrate feature selection into the model training process.

**Feature Extraction:** The process of feature extraction entails the generation of novel characteristics through the amalgamation or alteration of the initial features. The objective is to generate a collection of characteristics that effectively encapsulates the fundamental nature of the initial data within a reduced-dimensional framework.

Various techniques exist for feature extraction, such as principal component analysis (PCA), linear discriminant analysis (LDA) and t-distributed stochastic neighbour embedding (t-SNE). Principal Component Analysis (PCA) is a widely utilised method that facilitates the projection of the original features into a space with reduced dimensions, all the while striving to retain the maximum amount of variance.

#### Advantages of Dimensionality Reduction:

- One of the benefits of this technology is its ability to facilitate data compression, resulting in a reduction of required storage space.

## Notes

- The utilisation of this technique results in a reduction in computing time. Additionally, it aids in the elimination of any redundant characteristics that may be present.
- **Improved Visualization:** Visualising high dimensional data can be challenging due to the complexity of representing several dimensions. However, employing dimensionality reduction techniques can facilitate the visualisation of such data in two or three dimensions, hence enhancing comprehension and analysis.
- **Overfitting Prevention:** The utilisation of high-dimensional data in machine learning models has the potential to result in overfitting, hence causing a decline in the overall generalisation performance. The application of dimensionality reduction techniques can effectively mitigate data complexity, hence mitigating the risk of overfitting.
- **Feature Extraction:** Dimensionality reduction techniques are employed to extract significant features from datasets with a high number of dimensions. This process becomes valuable in the context of feature selection for machine learning models.
- **Data Preprocessing:** Before using machine learning methods, dimension reduction can be used as a preprocessing step to lower the dimensionality of the data and thus increase the performance of the model.
- **Improved Performance:** Dimensionality reduction techniques can be employed to enhance the efficacy of machine learning models by mitigating the intricacy of the data. Consequently, this leads to a reduction in noise and extraneous information within the data set.

### Disadvantages of Dimensionality Reduction

- The potential consequence of this situation is the possibility of experiencing a certain degree of data loss.
- Principal Component Analysis (PCA) has a tendency to identify linear relationships among variables, which can be considered undesirable in certain cases.
- Principal Component Analysis (PCA) may not be effective in situations where the mean and covariance of datasets do not provide sufficient information for their characterization.
- In practical applications, it is common to apply certain heuristics to determine the number of major components to retain, as the optimal number may not be known.
- **Interpretability:** The interpretability of the reduced dimensions may provide challenges, as it may be arduous to discern the connection between the original features and the reduced dimensions.
- Overfitting can occur in some scenarios when dimensionality reduction is employed, particularly when the selection of the number of components is based on the training data.
- One potential concern with certain dimensionality reduction approaches is their susceptibility to outliers, as this can lead to a distorted portrayal of the underlying data.
- The computational complexity of certain dimensionality reduction methods, such as manifold learning, can be significant, particularly when used to datasets of considerable size.

### 2.4.3 Introduction to Principle Component Analysis (PCA)

The aforementioned technique was initially proposed by Karl Pearson. The prerequisite for successful mapping of data from a higher dimensions space to a lower

dimensional space is that the variance of the data in the lower dimensional space must be maximised.

Principal component analysis (PCA) is a robust technique rooted in the principles of linear algebra, which may be effectively employed in various domains including data analysis, computer graphics and data compression. Principal Component Analysis (PCA) is a widely utilised technique in data analysis for identifying and extracting fundamental features or structures that may not be immediately apparent in a given dataset.

Sensors are employed in scientific studies for the purpose of data acquisition. In order to effectively capture the crucial information required for comprehending the underlying physics, it is imperative to strategically position sensors at precise places, the specifics of which are typically undisclosed. Principal Component Analysis (PCA) can thereafter be employed to represent the data in a novel basis that captures the most prominent components of the underlying data, which would have been measured had the ideal sensor positions been known in advance.

Principal component analysis (PCA) is a robust technique rooted in the principles of linear algebra, which may be effectively employed in various domains including data analysis, computer graphics and data compression. Principal Component Analysis (PCA) is a widely employed technique in data analysis that enables the identification and extraction of crucial elements or underlying structures that may not be immediately apparent in a given dataset.

In experimental settings, sensors are employed to capture and document data. In order to effectively comprehend the fundamental principles of physics, it is imperative to strategically position sensors at precise places, the specifics of which are typically undisclosed. Principal Component Analysis (PCA) can thereafter be employed to represent the data in a novel basis that captures the most prominent components of the underlying data, assuming that the ideal sensor placements were determined in advance.

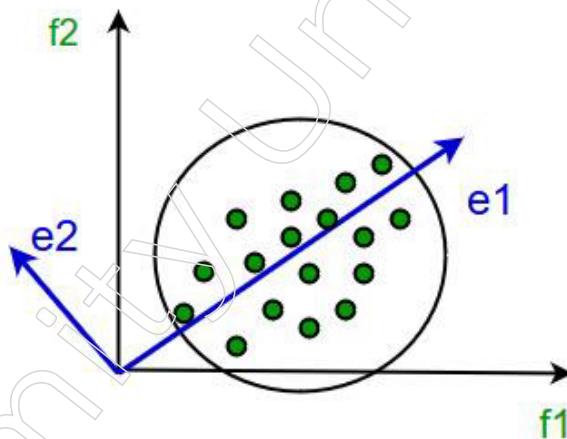


Figure:Principal Component Analysis

Principal Component Analysis (PCA) is a statistical method employed for the purpose of reducing the dimensionality of data, specifically by reducing the number of features present in a dataset. This reduction is achieved by identifying and selecting the most significant features that effectively capture the maximum amount of information pertaining to the dataset.

The selection of characteristics is based on their ability to produce variety in the output. The original characteristics of the dataset are transformed into principal components, which are obtained through linear combinations of the pre-existing features.

## Notes

The primary factor that contributes to the most variability is the initial Principal Component. The characteristic that accounts for the second biggest amount of variability is referred to as the second Principal Component and so forth.

Principal Component Analysis (PCA) is a technique utilised to extract significant elements, referred to as components, from a vast array of features present within a given dataset. Principal Component Analysis (PCA) is a statistical technique used to identify the principal directions of maximum variance within a dataset that has a high number of dimensions.

By projecting the data onto a lower-dimensional subspace, PCA aims to preserve the majority of the information contained in the original dataset. By employing a technique of projecting our data into a lower-dimensional space, we are effectively decreasing the number of dimensions in our feature space.

The intrinsic dimensionality of a data set  $X$ , denoted as  $X \subseteq R^l$ , is defined as  $m < l$ , indicating that  $X$  can be accurately represented using  $m$  independent parameters. The user's text: "I think the government should provide free healthcare for all citizens." Let us consider the generation of vectors in  $X$  as functions dependent on  $m$  random variables. Let  $x$  be a function of  $u_1, u_2, \dots, u_m$ , where each  $u_i$  is a real number for  $i = 1, 2, \dots, m$ . The observation vectors will be situated on a manifold, the specific shape of which is determined by the vector-valued function  $g : R^m \rightarrow R^l$ . Let us consider the given form of the function  $g$ , which can be expressed as  $x = [r \cos \theta, r \sin \theta]$ .

The variables  $T$  and  $r$  are constant, but the variable  $\theta$  belongs to the interval  $[0, 2\pi]$ . In this particular scenario, a single parameter is adequate for describing the data. When a minor amount of noise is introduced, the data tends to form clusters around the circumference. From a statistical perspective, this suggests that there is a correlation between the data.

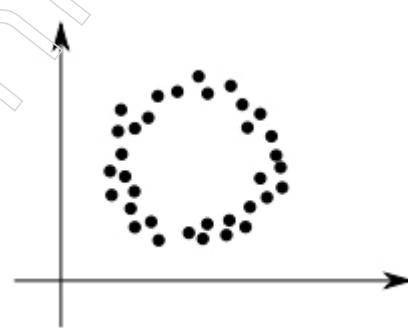


Figure: A sample of points that can be well described using a single parameter  $\theta$ .

It is postulated that the data under consideration are produced by a system or process that is influenced by a limited number of latent variables, which are not directly observable. The objective is to get knowledge about the underlying structure. Consider a vector of  $l$  elements, denoted as  $x_n \in R^l$ , where  $n = 1, 2, \dots, N$ . This vector is assumed to have a mean of zero. If the vector does not have a zero mean, the mean is subtracted. Here,  $n$  represents the  $n$ -th observation out of  $N$  total samples.

Principal Component Analysis (PCA) is a technique used to identify a subspace of dimension  $m$  that allows for optimal retention of statistical variation in the data following projection into this subspace, where  $m$  is less than or equal to  $l$ . The subspace under consideration possesses  $m$  axes that are mutually orthogonal to each other. The computations are performed in such a way that maximises the variance of the data once it has been projected onto the subspace.

Principal Component Analysis (PCA) does not directly augment the variance of the data. Rather, it performs a rotation of the data points, aligning them with the principal components, which are the directions that exhibit the most spread or variance.

main Component Analysis (PCA) use the Eigen decomposition of the covariance matrix to ascertain the main components. Pairs of eigenvalues and eigenvectors are a fundamental concept in linear algebra, whereby each eigenvector is associated with a specific eigenvalue. In the case of a  $n \times n$  covariance matrix, the number of eigenvectors will be equal to  $n$ . Eigenvectors are employed in order to comprehend the variance, or spread, inside a given dataset.

Specifically, they allow us to determine the variables that exhibit greater variance, with the magnitude of the eigenvector corresponding to the extent of that particular direction. When the Eigenvalues are arranged in a descending sequence, the Eigenvector corresponding to the first Eigenvalue represents the first principal component. Similarly, the second Eigenvector linked with the second Eigenvalue represents the second main component and so forth.

#### **PCA involves the following steps:**

- The task at hand involves the construction of the covariance matrix for the given data set.
- Please calculate the eigenvectors of the given matrix.
- The utilisation of eigenvectors that correspond to the greatest eigenvalues enables the reconstruction of a significant portion of the variance included in the initial data set.

Consequently, a reduced set of eigenvectors remains, potentially resulting in data loss during the procedure. However, it is crucial to preserve the significant differences through the retention of the remaining eigenvectors.

Principal components analysis (PCA) is a member of a group of methodologies designed to handle high-dimensional data by leveraging the interrelationships among variables to express it in a more manageable, lower-dimensional format, while minimising the loss of information. Principal Component Analysis (PCA) is considered to be a straightforward and resilient method for conducting dimensionality reduction. Additionally, this transformation holds the distinction of being one of the most ancient techniques, having undergone numerous rediscoveries across various disciplines.

Consequently, it has acquired alternative names such as the Karhunen-Loëve transformation, the Hotelling transformation, the method of empirical orthogonal functions and singular value decomposition. The technique will be referred to as Principal Component Analysis (PCA).

#### **Advantages of Principal Component Analysis:**

One of the main advantages of using PCA for data visualization is that it can handle high-dimensional data sets that are difficult to plot or interpret otherwise. PCA can reduce the complexity and noise of the data and highlight the most important features and relationships. For example, you can use PCA to visualize the similarities and differences between different groups of customers, products, or genes, based on multiple attributes. Another advantage of PCA is that it is a linear and unsupervised method, which means that it does not require any prior assumptions or labels about the data. It can also be easily implemented and interpreted using standard tools and libraries.

- **Removes Correlated Features:** In practical situations, it is frequently observed

## Notes

that datasets often contain a large number of characteristics, typically numbering in the thousands. Running the algorithm on all features may lead to a decrease in algorithm efficiency and hinder the ability to effectively visualise the numerous features in a graphical representation. It is important to undertake feature reduction in the dataset. It is necessary to ascertain the correlation among the features, namely the variables that exhibit correlation.

- The manual identification of correlations among thousands of features is an arduous, time-consuming task that is practically difficult to do. This process can be highly frustrating for researchers. Principal Component Analysis (PCA) efficiently performs this task. Upon the use of Principal Component Analysis (PCA) on the dataset, it is observed that all the Principal Components exhibit independence from each other. There is no discernible association between them.

- **Improves Algorithm Performance:** The presence of several features will lead to a significant deterioration in the performance of your algorithm. Principal Component Analysis (PCA) is a widely employed technique in the field of Machine Learning to enhance algorithmic efficiency by eliminating associated variables that do not significantly contribute to decision-making processes. The training duration of the algorithms experiences a notable decrease when the number of features is reduced.

If the input dimensions are excessively large, employing Principal Component Analysis (PCA) as a means to expedite the procedure is a justifiable decision.

- **Reduces Overfitting:** Overfitting primarily arises when the dataset has an excessive number of variables. Principal Component Analysis (PCA) serves the purpose of mitigating the problem of overfitting by diminishing the dimensionality of the feature space.
- **Improves Visualization:** The comprehension and visualisation of data in high dimensions pose significant challenges. Principal Component Analysis (PCA) is a mathematical technique used to reduce the dimensionality of high-dimensional data to a lower dimension, often two dimensions. This reduction allows for easier visualisation of the data.

The utilisation of a 2D Scree Plot enables the identification of Principal Components that exhibit a greater degree of variance and exert a more substantial influence in comparison to other Principal Components. The IRIS dataset, even in its most basic form, consists of four dimensions, making it challenging to visualise. Principal Component Analysis (PCA) can be employed to reduce the data to a two-dimensional representation, hence enhancing the visualisation capabilities.

Let us contemplate a scenario in which we are presented with a set of 50 features, denoted as  $p = 50$ . A total of 1225 scatter plots can be generated to analyse the associations between variables, which can be calculated using the formula  $p(p-1)/2$ . Performing exploratory analysis on this dataset would be a laborious task. Hence, it is imperative to employ Principal Component Analysis (PCA) as a means to mitigate this issue.

### 2.4.4 Disadvantage of PCA

Principal Component Analysis (PCA) also exhibits several limits and downsides when employed for the purpose of data visualisation. One of the primary drawbacks associated with Principal Component Analysis (PCA) is its potential to result in information and detail loss during the process of dimensionality reduction. The process described has the potential to result in the simplifying or distortion of the data, hence impeding the identification of outliers or abnormalities.

For instance, it is possible to overlook certain intricate patterns or trends that can only be discerned in dimensions beyond the conventional ones. One additional drawback of Principal Component Analysis (PCA) is its susceptibility to scaling and outliers, which might have an impact on the integrity and robustness of the outcomes. Hence, it is imperative to meticulously preprocess and normalise the data prior to implementing Principal Component Analysis (PCA).

- Independent variables become less interpretable: Upon applying Principal Component Analysis (PCA) to the dataset, the original features will be transformed into Principal Components. Principal components refer to the linear combinations of the initial features. Principal components have lower readability and interpretability compared to the original features.
- When applying primary Component Analysis (PCA) to our dataset, the original features undergo a transformation, resulting in the creation of primary components. These components are derived from linear combinations of the original data's characteristics. However, it is imperative to determine the most significant aspects within the given data collection. Answering this question may pose challenges following the computation of Principal Component Analysis (PCA). Biplots are commonly employed as a valuable tool for facilitating the interpreting process.
- Data standardization is must before PCA: It is imperative to perform data standardisation prior to the use of Principal Component Analysis (PCA) in order to ensure that the ideal Principal Components may be accurately identified.

For example, when a feature set contains data that is measured in units such as Kilogrammes, Light years, or Millions, the variance scale within the training set is significantly large. When Principal Component Analysis (PCA) is done to a given feature set, the resulting loadings for features that exhibit considerable variance will likewise be substantial. Therefore, the principal components will exhibit a bias towards traits that possess a large degree of volatility, consequently yielding inaccurate outcomes.

In order to ensure standardisation, it is necessary to convert all category data into numerical features prior to the application of PCA.

Principal Component Analysis (PCA) is influenced by the scale of the features within the dataset. Therefore, it is necessary to standardise or normalise the features prior to using PCA. The StandardScaler module from the Scikit-learn library is employed for standardising features in a dataset. It is essential to standardise the features of a dataset onto a unit scale, with a mean of 0 and a standard deviation of 1. This standardisation is necessary to ensure the optimal performance of various Machine Learning algorithms.

- The principal component analysis (PCA) algorithm is utilised to ascertain the principal directions that correspond to the greater changes within a dataset. Prior to computing the principal components, it is necessary to standardise all variables such that they possess a mean of zero and a normal deviation of one, as the variance of a variable is assessed on its own squared scale. Alternatively, the principal component analysis (PCA) may be biassed towards variables with bigger scales.
- Information Loss: While Principal Components aim to capture the highest amount of variance within a dataset's features, it is important to exercise caution when selecting the number of Principal Components. In doing so, there is a risk of overlooking certain information in comparison to the original set of features. The utilisation of principle Component Analysis may result in a certain degree of information loss if an inadequate number of principle components is chosen, hence failing to adequately account for the variance present in the dataset.

## Notes

### 2.5 Data Visualisation

**Visualize:** The process of creating a cognitive representation, visualisation, or depiction of an object or concept that is not physically observable or currently in existence, with the intention of rendering it perceptible to the mind or imagination.

- Visualisation refers to the application of computer graphics to provide visual representations that facilitate the comprehension of intricate and frequently extensive data representations.
- Visual data mining refers to the systematic exploration and extraction of valuable knowledge from extensive datasets through the utilisation of visualisation techniques.

#### Data visualization Definition:

Data visualisation refers to the visual depiction of data and information. Data visualisation tools utilise visual components such as charts, graphs and maps to facilitate the comprehension of data trends, outliers and patterns in a user-friendly manner. Moreover, it offers a commendable means for employees or business proprietors to effectively communicate data to individuals lacking technical expertise, hence minimising any potential uncertainty.

In the realm of Big Data, the utilisation of data visualisation tools and technologies is crucial for the examination of vast quantities of information and the facilitation of data-informed decision-making. Data visualisation is a technique employed to analyse and derive meaningful understanding from vast quantities of data. This is accomplished by utilising visual representations, frequently interactive in nature, of unprocessed data.

#### Why Data Visualization Is Important:

The significance of data visualisation is straightforward: It facilitates people in observing, engaging with and comprehending data in a more efficient manner. The accurate representation of one's level of knowledge using visual means can facilitate a shared understanding among individuals, regardless of the simplicity or intricacy of the subject matter. The enhancement of data comprehension is expected to yield benefits across several professional sectors. The comprehension of data holds advantages for several fields within STEM, as well as for those in the public sector, commerce, marketing, history, consumer products, services, sports and other relevant domains.

Data visualisation has undeniable practical and tangible implications in real-world contexts, despite the tendency to express admiration for it (especially on the Tableau website). The ubiquity of visualisation renders it a very advantageous skill to acquire within professional contexts. The utilisation of information, whether through a dashboard or a slide presentation, correlates positively with enhanced effectiveness.

The concept of a citizen data scientist is gaining increasing prominence in contemporary discourse. The required skill sets in a data-driven environment are undergoing continuous evolution. The increasing value of professionals who possess the ability to analyse data for decision-making purposes and effectively communicate data-driven narratives through visual representations is evident.

The field of data visualisation occupies a central position between analysis and visual storytelling. Traditional education tends to emphasise a clear separation between creative narrative and technical analysis. However, contemporary professional environments increasingly value individuals who possess the ability to navigate both domains.

### Different Types of Visualizations:

When contemplating data visualisation, one's initial inclination likely gravitates towards rudimentary bar graphs or pie charts. Although visualisations are often considered essential for representing data and serve as a standard foundation for various data visualisations, it is crucial to ensure that the appropriate visualisation is combined with the relevant information. Basic graphs represent just a fraction of the larger whole. There exists a wide array of visualisation techniques that can be employed to effectively and engagingly portray data.

### General Types of Visualizations:

- A chart is a visual representation of data that is given in a tabular or graphical manner, where the information is organised along two axes. The information can be visually represented through several means such as a graph, diagram, or map.
- A table is a collection of numerical values organised in a grid-like structure, with rows and columns.
- A graph is a visual representation of data points, lines, segments, curves, or areas that depict the relationship between specific variables. Typically, graphs are constructed with two axes intersecting at a right angle to facilitate the comparison of these variables.
- Geospatial refers to a method of visualisation that utilises maps to represent data, employing various shapes and colours to depict the correlation between different data points and their respective geographical locations.
- An infographic is a visual representation of data that combines images and text. Typically employs graphical representations such as charts or diagrams.
- Dashboards refer to a compilation of visual representations and data that are consolidated in a single location, serving the purpose of facilitating the analysis and presentation of data.

### More specific examples:

- **Area Map:** Area maps are a type of geospatial visualisation that is employed to depict specific values assigned to various regions on a map, such as a country, state, county, or any other geographical area. Choropleths and isopleths are two prevalent categories of area maps.
- **Bar Chart:** Bar charts are graphical representations that depict the relative magnitudes of numerical values in comparison to one another. The magnitude of each variable is visually represented by the length of the corresponding bar.
- **Box-and-whisker Plots:** The provided visualisations depict a variety of intervals (represented by the box) spanning a given metric (shown by the bar).
- **Bullet Graph:** A graphical representation consisting of a bar placed on a background to indicate the level of progress or performance achieved in relation to a specific objective, typically represented by a line on the graph.
- **Gantt Chart:** Gantt charts, commonly employed in the field of project management, serve as graphical representations of timeframes and tasks through the use of bar charts.
- **Heat Map:** Geospatial visualisation in the form of maps is a prevalent technique that employs distinct colours to represent specific data values, however it is not limited to temperature data.

## Notes

- **Highlight Table:** A type of tabular representation that employs the use of colour to classify comparable data, hence facilitating enhanced readability and intuitive comprehension for the viewer.
- Histogram: A form of graphical representation known as a histogram, which partitions a continuous variable into distinct intervals or bins in order to facilitate the examination of its distribution.
- Pie Chart: A circular diagram composed of triangle segments that visually represents data in terms of proportions relative to a whole unit.
- Tree map: A graphical representation that depicts various interconnected values through the arrangement of nested rectangles.

### The advantages and benefits of good data visualization:

- The human visual system is naturally inclined to be attracted to various hues and patterns. The ability to distinguish between red and blue, as well as between square and circular, can be rapidly accomplished. Contemporary society exhibits a visual-oriented culture, wherein various forms of visual media such as art, commercials, television and cinema play a significant role.
- The user's text does not contain any information to rewrite. Data visualisation is a distinct manifestation of visual art that captivates our attention and sustains our focus on the conveyed content.
- The user's text does not contain any information to rewrite. Upon observing a chart, individuals are able to promptly identify patterns and anomalies. The process of internalising visual stimuli occurs rapidly when they are perceptible to our senses.

The act of storytelling serves a certain intention or objective. If one has encountered the situation of observing a substantial dataset in the form of a spreadsheet and being unable to discern any discernible pattern, it becomes evident how significantly more impactful a visualisation can be in such circumstances.

**Enhanced Consensus:** In the realm of business, it is sometimes necessary to assess the performances of two entities or scenarios throughout several time periods. One commonly employed approach is the systematic collection of extensive data pertaining to the given conditions, followed by a thorough analysis of the data. Undoubtedly, this task will require a significant amount of time.

**An Enhanced Approach:** This method effectively addresses the challenge of incorporating the information from both perspectives within the visual framework. Undoubtedly, this will provide a more comprehensive understanding of the circumstances. For example, Google trends aid in the comprehension of data pertaining to popular industries or queries through visual or graphical representations.

The act of sharing data in a straightforward manner involves the transmission of information, enabling organisations to establish a new mode of communication. Instead of disseminating complex information, the utilisation of visual data facilitates the transmission of more easily comprehensible and assimilable information.

**Sales Investigation:** Through the utilisation of data visualisation, a salesperson can easily appreciate the sales trajectory of various products. By utilising information perception tools such as heat maps, one can gain insights into the factors that contribute to the increase in company figures, as well as the factors that lead to a decline in business figures. Information representation plays a crucial role in enhancing comprehension of patterns and various variables, such as customer preferences, repeat customers, geographical influences and more.

**Exploring the Connections Between Events:** A business is subject to several factors that exert effect upon it. Establishing a correlation between these factors or situations fosters leaders' comprehension of the difficulties associated with their firm. The online business industry is not a novel concept in contemporary times. On specific joyous occasions, such as Christmas or Thanksgiving, the graphical representations of internet-based entities experience an increase. In this context, it can be stated that if an online organisation is generating a revenue of \$1 million in a certain quarter and experiences subsequent growth, it can promptly identify the factors contributing to this increase.

**Exploring Openings and Patterns:** Business leaders have the ability to access a vast amount of information, allowing them to have a comprehensive understanding of the various trends and opportunities that exist in their environment. By employing information representation techniques, experts are able to identify patterns in the behaviour of their customers, so enabling them to analyse trends and identify commercial opportunities.

#### **Disadvantages of Data Visualization:**

- The process of creating visualisations can be characterised as time-consuming, particularly when confronted with extensive and intricate datasets. The presence of this factor has the potential to impede the efficiency of the machine learning process and diminish overall output.
- Potential for Misleading: Although data visualisation has the potential to facilitate the identification of patterns and relationships within data, it is important to acknowledge that incorrect execution can lead to misleading interpretations. Visualisations have the potential to generate perceptions of patterns or trends that may lack empirical basis, hence resulting in erroneous deductions and suboptimal decision-making.
- The task of interpretation can pose challenges. Certain forms of visualisations, particularly those incorporating three-dimensional or interactive components, may provide challenges in terms of interpretation and comprehension. This phenomenon has the potential to induce ambiguity and misconstrual of the collected data.
- This approach may not be applicable to all forms of data. Certain categories of data, such as textual or auditory data, may not be inherently suitable for visualisation. In such instances, it may be more suitable to employ different approaches of analysis.
- The accessibility of the content may not be available to all users. Certain individuals may experience visual impairments or other limitations that pose challenges or render them incapable of comprehending visual representations. In certain instances, it may be imperative to employ alternate approaches for data presentation in order to guarantee inclusivity.

#### **Effective visualization:**

The greatest way to convey information from the increasingly massive and complicated datasets in the natural and social sciences is through effective visualisation. But as power grows, business analysts frequently must sift through a bewildering assortment of visualisation alternatives.

Data visualisation utilises visual data to effectively and expeditiously communicate information to diverse audiences. Moreover, this practise can aid organisations in discerning the variables that exert influence on consumer behaviour, pinpointing areas that necessitate enhancement or further scrutiny, enhancing stakeholders' retention of data, ascertaining optimal timings and locations for the sale of specific products and projecting sales volumes.

## Notes

The significance of data visualisation stems from the cognitive processes involved in human information processing. When visualising extensive and intricate data sets, individuals often find studying graphs and charts to be more enjoyable compared to spreadsheets and reports. Data visualisation is a highly effective and expeditious method for conveying ideas to a wide audience. By implementing a minor adjustment, one can experiment with a novel structure. A way of encoding numerical, relational, or spatial information into visuals is known as data visualisation. The graphical display of information and data using visual components including charts, graphs, maps and dashboards is known as data visualisation. Data visualisation focuses on how to present data to the appropriate audience at the appropriate moment so that they can most effectively get insights.

There are many different commercial and free data visualisation tools on the market. Tableau, Qlikview, Sisense, GoogleData Studio, Zoho Analytics, Fusioncharts, Highcharts, Datawrapper, Plotty, Microsoft PowerBI and IBM Watson Analytics are a few of the widely used data visualisation technologies.

### Benefits of data visualization:

Data visualisation technologies provide innovative methods to significantly enhance the capacity to comprehend concealed information inside extensive sets of company data. There are several key benefits associated with data visualisation for decision makers and their organisations. These advantages can be summarised as follows:

- Enhanced Assimilation of Business Information: The human visual system processes visual and picture information better than textual and numerical information. However, most business intelligence reports for senior executives have static tables and charts, which do not visually engage the audience.

Data visualisation allows the collection of massive amounts of operational and corporate data.

Data visualisation uses heat maps, fever charts, and other visually immersive graphics to help decision makers understand complex data sets.

- Quick Access to Relevant Business Insights: Visual data discovery helps companies find and access relevant information faster than competitors, boosting productivity. A recent study found that corporate managers who use visual data discovery tools are 28% more likely to receive timely information than those who use controlled reporting and dashboards.

Business intelligence users, specifically 48%, in organisations using visual data discovery techniques can independently access the needed information without relying on IT personnel.

- **Determine patterns in business operations:** Data visualisation lets users see complex patterns. It helps demonstrate the relationship between business and operations and evaluate performance measures. Data visualisation helps identify operational changes that may have affected business performance by examining the relationship between everyday job activities and firm success.

Data visualisation technology let decision-makers quickly identify and understand changes in customer behaviour or market conditions. The business pattern revealed helps decision makers identify the causes of positive or negative performance swings and take corrective action.

- **Rapid Identification of Latest Trends:** Modern organisations can collect massive amounts of customer and market data. Depending on their capacity to find new

ways to make money and grow their businesses, this plethora of data can offer company executives significant insights. Decision-makers may easily understand client behaviour and market conditions across data sets by visualising data.

Accurate consumer sentiment analysis Data visualisation can help companies comprehend customer sentiment and other data, revealing new service opportunities. These useful insights enable firms to seize on new business opportunities to stay ahead of the competition.

- **Direct Interaction with Data:** Data visualisation helps organisations manage and interact with their data directly and effectively. An advantage of data visualisation is its capacity to convey actionable insights. Data visualisation tools let people interact with and change data, unlike static one-dimensional tables and charts.
- **Predictive Sales Analysis:** Sales leaders can utilise real-time data visualisation to conduct advanced predictive analytics on their sales figures. This enables them to access current sales data and analyse the factors contributing to underperforming products and lagging sales. One potential factor contributing to this phenomenon could be the availability of discounts provided by other companies.
- **Drill-Down Sales Analysis:** Heat map data visualisation helps corporate executives identify successful and unsuccessful product categories. This method also lets executives go deeper into the data to find the factors affecting sales. For instance, the statistics may show that pet-care products perform poorly while higher-income customers dominate sales.

These data can be used to strategically target promotions to this consumer demographic to boost conversion rates and revenue.

- **Easy Comprehension of Data:** Data visualisation helps firms assess and present enormous datasets clearly. This applies to entertainment, current events, finance, and politics. It also fosters deep comprehension, encouraging people to make good decisions and act quickly when needed.
- **Customized Data-Visualization:** Data visualisation can not only display data in a graphical style, but also allow users to change its structure, remove extraneous pieces, and apply filters for more thorough information. This feature attracts business executives and improves communication. Compared to traditional data presenting methods, it is also beneficial.

#### Categories of Data Visualization:

The utilisation of data visualisation plays a crucial role in market research, as it allows for the representation of both numerical and categorical data.

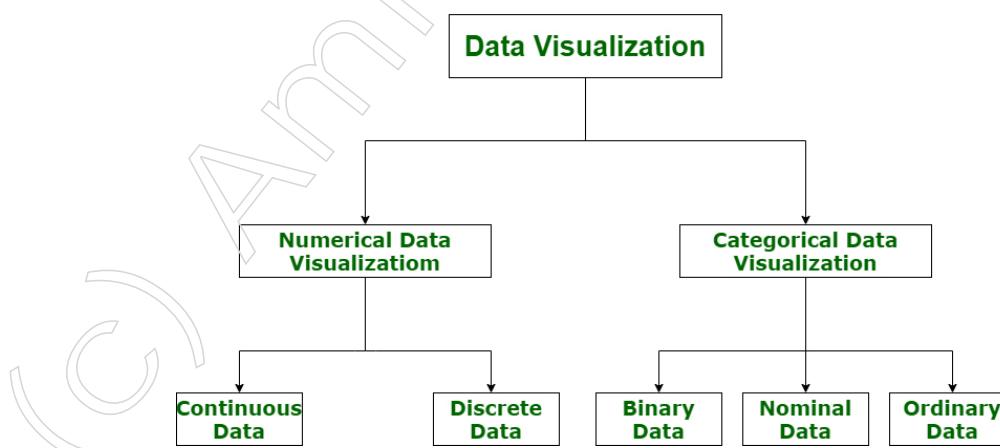


Figure: Data visualization category

## Notes

This practise enhances the effectiveness of insights and mitigates the potential for analysis paralysis. Data visualisation can be classified into the subsequent categories:

**Numerical Data:** Numerical data is commonly referred to as quantitative data in academic discourse. Numerical data refers to data that quantitatively depicts various attributes, such as the height, weight and age of individuals. The utilisation of numerical data visualisation represents the most straightforward method for visually representing data. The primary purpose of this tool is to facilitate the comprehension and use of extensive datasets and unprocessed numerical information by individuals, hence enhancing their ability to get actionable insights from such data. Numerical data is classified into two distinct categories:

- **Continuous Data:** It can be narrowed or categorized (Example: Height measurements).
- **Discrete Data:** This type of data is not “continuous” (Example: Number of cars or children’s a household has).

Charts and numerical values are commonly employed as visualisation approaches to depict numerical data visualisation. Some examples of visual representations commonly used in data analysis are pie charts, bar charts, averages and scorecards.

**Categorical Data:** Categorical data is alternatively referred to as qualitative data. Categorical data refers to data that is typically representative of distinct groupings. The dataset has categorical variables that serve as indicators for various attributes, such as an individual's rating or gender. The visualisation of categorical data primarily involves the representation of significant themes, the establishment of relationships and the provision of contextual information. Categorical data is divided into three distinct categories:

- **Binary Data:** Positioning is used to categorise in this (for example, “Agrees” or “Disagrees”).
- **Nominal Data:** Here, categorization is done according to characteristics (for instance, male or female).
- **Ordinal Data:** Here, classification is done based on traits (such male or female, for example).

Graphics, Diagrams and Flowcharts are the types of visualisation techniques that are used to display categorical data. Word clouds, sentiment mapping and Venn diagrams are a few examples.

### 2.5.1 Understanding Data Visulisation Concepts

The graphic display of information and data is known as data visualisation. In order to assist people comprehend and evaluate the underlying patterns, trends and insights within the data, it entails the use of visual components such as charts, graphs and maps. Complex datasets can be made more approachable, comprehensible and useful through effective data visualisation. To help you better comprehend data visualisation, consider the following essential ideas:

**Data:** Data can be defined as the unprocessed information or factual observations that are gathered, quantified, or documented. Data can manifest in diverse formats, encompassing numerical values, textual content, visual representations and even audiovisual recordings. Visualization:Visualisation refers to the process of presenting data in a graphical format, wherein raw data is converted into visual elements such as charts, graphs, or maps. The purpose of visualisation is to facilitate comprehension and enhance understanding of the data.

**Visual Encoding:** Visual encoding refers to the procedure of associating data attributes, such as numerical values or categories, with visual qualities like position, size, colour and shape inside a data visualisation. The selection of suitable visual encodings plays a critical role in the effective communication of information.

**Charts and Graphs:** These are the most common forms of data visualizations. Some popular types include:

- **Bar Charts:** Used to compare different categories or discrete data points.
- **Line Charts:** Useful for showing trends and patterns over time or continuous data.
- **Pie Charts:** Suitable for displaying proportions or percentages of a whole.
- **Scatter Plots:** Shows the relationship between two variables and their correlation.
- **Heatmaps:** Visualize data in a tabular format with color-coded cells to represent values.

**Data Abstraction:** Data abstraction is a process that aims to simplify intricate data sets in order to emphasise fundamental patterns and trends, while ensuring that crucial information is not compromised or omitted. Aggregation, filtering and summarization are often employed procedures for abstracting data.

**Data Storytelling:** Data visualisation is an influential instrument for narrative construction. This tool facilitates analysts in effectively communicating data insights by employing persuasive techniques, so guiding the audience through the data and aiding their comprehension of the main findings.

**Color Theory:** The selection of colours plays a crucial role in the field of data visualisation. Colours can serve as a means of representing distinct categories or emphasising particular data points. Nevertheless, it is vital to guarantee that the utilisation of colours is both intuitive and significant.

**Data Labels and Annotations:** Labels and annotations play a crucial role in enhancing the interpretability of data visualisations by providing contextual information, hence facilitating proper comprehension for the audience.

**Interactivity:** Interactive data visualisations enable users to actively engage with and manipulate data in real-time, facilitating a dynamic and immersive exploration of the information presented. Users have the ability to delve deeper into specific details, apply filters to manipulate data and modify visual representations, so augmenting their comprehension and gaining valuable insights.

**Data Visualization Tools:** A wide array of data visualisation tools are already accessible, encompassing both basic spreadsheet-based options and advanced software capable of managing large-scale datasets. Several widely used technologies in data visualisation include Tableau, Power BI, matplotlib, D3.js, ggplot2, among others.

**Data Visualization Best Practices:** In the process of constructing data visualisations, adherence to established guidelines is crucial in order to guarantee the attainment of clarity, precision and efficacy. Several recommended practices encompass selecting suitable chart types, accurately labelling axes and data points, minimising visual clutter and tailoring the design to the intended audience.

#### **Data Visualization Tools:**

There are many different commercial and free data visualisation tools on the market. Tableau, Qlikview, Sisense, Looker, Google Data Studio, Zoho Analytics, Fusioncharts, Highcharts, Datawrapper, Klipfolio, Kibana, Chartio, Plotly, Infogram, Visme, Geckoboard,

## Notes

AnyChart, D3.js, Microsoft PowerBI, IBM Watson Analytics and SAP Analytics Cloud are a few of the widely used data visualisation tools.

The features, advantages and disadvantages few important tools in the market are elaborated below:

### **Google Data Studio:**

The Google Account is required to utilise the free data visualisation tool Google Data Studio. Utilising data from Google Analytics, Google Sheets, Google Ads, Google Search Console, YouTube and MySQL, Google Data Studio is a tool for data visualisation. The templates in Google Data Studio make it simple and quick to set up reports and dash boards. Google Data Studio is the newest product and a component of Google's analytics solutions. Being relatively new to the market, it aims to distinguish itself from numerous rivals by usability, an understated yet elegant design, creative problem-solving and simple, routine ways to share dashboards (much like sharing papers). There is no desktop version of this solution; it is entirely web-based.

With the Google Analytics Solutions data toolkit, a software suite for data analysis and facilitating data-driven solutions, Google aims to hit the sweet spot on the market by not only promoting one BI tool but also all of their other tools for working with data. Raw data can be transformed using Google Data Studio to be shown in interactive visualisations that are assembled into dashboards. Additionally, the tool is perfectly suited for use with data sources that are exclusive to Google. It offers straightforward data connectors that make accessing the data simple. Last but not least, one of the best features relates to Google Data Studio's methods for team collaboration, which enable a group of developers to work together on a same issue. Similar to Google Docs, Data Studio enables users to view and change the dashboard.

The sharing of reports is simple and works similarly to Google Drive. The control of access levels functions similarly, allowing you to invite others to read a report or a folder of reports via email or a shareable link and to specify whether they can alter the report after viewing it. The total scope of Google Data Studio is still constrained. When compared to other tools, it falls short when it comes to calculating functions, customising images and making reports interactive.

#### **Advantages:**

- The development of this particular software occurred within the Google Analytics suite and it is designed to seamlessly interface with other pertinent Google products.
- The system is characterised by its simplicity in all significant elements, rendering it user-friendly.
- The cooperation skills are exceptional.

#### **Disadvantages:**

- Less adaptable than its rivals.
- Less capacity to add unique visuals; instead, one can only slightly alter those that are already there.
- There is no support for interaction.
- Lack of ability to combine and blend data.
- Only small changes can be made to the data; it must be suitable for visualisation.

### **Tableau:**

The aforementioned data visualisation tool is widely recognised and holds a prominent position in the market. It is utilised to visually represent and examine data in a manner that is simply comprehensible. This technology is widely employed by numerous firms globally, as it is highly effective in facilitating business insight and analysis. This feature enables the utilisation of real-time data sets and prioritises the allocation of time towards data analysis.

The software's widespread adoption across various industries can be attributed to its user-friendly interface and its capacity to generate dynamic visual representations. The technology is highly adept at managing large and rapidly evolving datasets commonly employed in big data operations, encompassing artificial intelligence and machine learning applications.

Tableau offers many licencing plans that are comparable to those found in other business intelligence (BI) tools. Tableau provides a range of three unique products that exhibit significant variations in terms of pricing. There are three distinct products offered by Tableau, namely Tableau Desktop, Tableau Online and Tableau Server, each of which is priced differently.

The platform facilitates integration with a wide range of diverse and sophisticated data sources, including various structured file formats (such as CSV, JSON, XML, MS Excel), relational and non-relational database systems (such as Postgre SQL, MySQL, SQL Server, Mongo DB) and cloud-based systems (such as AWS, Oracle Cloud, Google Big Query, Microsoft Azure).

Tableau differentiates itself from its competitors with the inclusion of a unique functionality known as Data Blending. One notable characteristic of this technology is its capacity for real-time collaboration, rendering it a significant asset for both business and non-commercial entities. There exist multiple methods for disseminating Tableau reports, including publishing them to a Tableau Server, utilising the email functionality of Tableau Reader and openly releasing the Tableau workbook while granting access to individuals through a shared link. The wide range of choices available allows for significant adaptability and alleviates numerous limitations.

Tableau provides a wide range of visualisation capabilities that possess unique qualities, facilitating intelligent methods of data exploration and profound understanding. Tableau offers a diverse range of visualisation forms, such as word clouds and bubble charts, which include distinct qualities that facilitate enhanced comprehension. The utilisation of tree diagrams and tree maps offers contextual information to the visual representations. The latter is commonly employed to represent categorical data, with a specific emphasis on the most pertinent aspects of the information.

Tableau dashboards exhibit a remarkable degree of flexibility. The fundamental characteristics of this system enable the user to arrange the dashboard in a customised manner without any overlapping elements. This feature is particularly useful for optimising the ergonomics of the screen space. Tableau is characterised by its user-friendly interface and accessible learning curve, making it an approachable tool for individuals with varying levels of technical expertise in visualisation workflows. It aims to empower users, including those without prior exposure to technical intricacies, by providing comprehensive functionality. This objective is achieved by the implementation of an intuitive user interface, where all necessary functions can be accessed within a maximum of two clicks. Additionally, the system has robust filters and drill-down options that are conveniently located and user-friendly. Furthermore, actions inside the system are thoroughly documented and clearly labelled, enhancing usability and ease of navigation.

## Notes

### **Qlikview:**

This programme serves as a prominent contender in the market for data visualisation and stands as the primary challenger to Tableau. One notable benefit of the tool is its extensive range of capabilities, which may be highly customised to suit individual needs. QlikView not only possesses robust data visualisation capabilities, but also offers advanced business intelligence analytics and enterprise reporting functionalities, complemented by a streamlined and uncluttered user interface. QlikView is widely recognised as a platform with high cost implications within the business intelligence (BI) domain. QlikView is a data-centric solution that places emphasis on the end-user as the primary recipient of information. The platform enables users to investigate and uncover information within their data using a workflow that mirrors the one employed by developers during data processing. In order to preserve the linkage between data, this software endeavours to uphold flexibility in its methodology for data exploration and visualisation. This feature enhances the accessibility of your data to end-users who are searching for specific information, regardless of any obstacles or discrepancies in the provenance of the relevant items.

QlikView exhibits a high degree of flexibility. This feature facilitates the configuration and fine-tuning of various attributes pertaining to individual objects, enabling users to personalise the visual appearance and user experience of visualisations and dashboards.

The software offers a significant degree of flexibility, which is complemented by an integrated ETL (Extract, Transform, Load) Engine. This engine facilitates the execution of standard data cleansing procedures. Nevertheless, it is possible that the associated expenses could be substantial.

### **Unique Features:**

- Individuality and adaptability.
- Comprehensive collection of features for designing complex dashboards.
- Automated manipulation of data associations.
- Facilitates faster data exploration and queries by storing data in memory.

### **Advantages:**

- An appealing user interface.
- Filtering for any type of visuals is simple to set up and both graphs and tables render quickly.
- The ability to mail reports in the practical PDF format.

### **Disadvantages:**

- During filtering, some data elements may be inadvertently combined.
- There is no way to combine bookmark results.
- Issues in using it as a business tool.

### **Common features of Data Visualization Tool:**

The utilisation of Data Visualisation tools enables enterprises, organisations and firms to present data in a structured and organised manner. This approach facilitates ease of interpretation, while also ensuring that the data is meaningful and conducive to informed decision-making. The process involves the identification of patterns, the reduction of noise and the elimination of insignificant values from the data in order to generate insights that may be effectively acted upon.

In order to optimise the utilisation of Data Visualisation, organisations must carefully choose a tool that offers a diverse range of features and capabilities. This section will address some key qualities that a data visualisation tool should include. Experts suggest that organisations should take into account a range of characteristics and capabilities when they decide to employ visualisation tools for big data.

#### **Clear and Customizable Dashboard:**

The dashboard is undoubtedly an essential component of any data visualisation tool. Similar to how one glance at the dashboard of a car provides all the necessary information, such as the speed, an indicator, a light, a seat belt and fuel, etc. Similar to this, a visualisation dashboard should be able to quickly convey all the important data.

A decent data visualisation dashboard should do several tasks at once. To start, it need to be fantastic-looking. It must be concise, with bright colour pops and enough white space. Too little colour or too much white can be intimidating. A balance should be achieved on the dashboard. All relevant data should be appropriately summarised by the dashboard. The most important Key Process Indicators (KPIs) you are trying to measure, the crucial trends you are keeping an eye on, or any other datum that is crucial to your organisation should be clearly visualised on the dashboard so that you can get a broad overview shortly after loading the dashboard.

The data displayed on the dashboard should be easy to understand and comprehend at a glance. Customizability is a crucial feature that every dashboard must have. Your business can be monitoring dozens of distinct datasets at any given time. The dashboard should give you the option to choose which datasets are shown prominently. Because each team has a different set of priorities, the data visualisation tool must be completely customizable.



Figure: Sample Dashboard prepared using Microsoft Power BI visualization tool

### **Embeddability:**

The smooth integration of visual reports into various applications is crucial for fully harnessing the potential of data visualisation. In order to optimise team productivity, enhance collaboration and provide seamless information sharing across diverse platforms, it is imperative that the data visualisation software possesses the capability to seamlessly integrate numerous media formats, such as graphs and charts, into a multitude of mobile and online applications.

The preservation of visual report quality and richness should not be compromised

## Notes

during the process of transferring it to an alternative application. The reports ought to retain their interactive nature, facilitating more exploration and analysis of the material.

It is not necessary for every department to engage in the analysis of all the high-level data that your tool gathers. The majority of users express a preference for integrating only a specific portion of the data into their respective applications in a seamless manner. The individual or organisation requires prompt and practical insights that may be implemented to enhance the effectiveness of their duties and campaigns. An effective data visualisation tool should possess the capability to be easily embedded.

### Performance:

If the utilisation of visualisation tools for big data hinders workers' workflow, it is probable that their adoption and usage would decrease. While a brief delay of a few seconds may not hold much significance in many scenarios, it has the potential to discourage users who are responsible for assessing numerous selections within a single day. Various features can contribute to enhancing performance, such as the utilisation of prompts, the implementation of data optimisation settings and the incorporation of dynamic loading options.

An additional aspect to contemplate in terms of performance is the tool's capacity to execute calculations on graphics processing units (GPUs). The increasing size of data sets has posed challenges in effectively processing substantial volumes of data using conventional systems. The utilisation of GPUs in conjunction with direct memory access (DMA) can significantly enhance the speed and efficiency of processing extensive datasets. This facilitates the construction of high-definition visualisations on the server side, which are subsequently delivered by the application through a web-based interface.

### Interactive Reporting:

The data visualisation tool should produce visual reports that exhibit a high degree of interactivity, facilitating effortless exploration of patterns and revelations. Interactive data visualisation plays a crucial role in discerning patterns and constructing narratives using data. This encompasses functionalities such as filtering, segmenting and exploring data at high rates, enabling users to efficiently analyse large datasets and promptly obtain answers to their inquiries.

Professionals in the field of data analysis and decision-making must possess the capability to gather data from many sources and merge datasets in order to provide reports that offer valuable insights. The tool should include the capability to enable the viewing of reports in many formats, while also allowing for the selective highlighting of different sections at distinct intervals.

Customization of industry-specific key performance indicators (KPIs) is essential in order to offer specialised insights. In order to provide these functionalities, it is imperative that the business intelligence and data visualisation tool possesses a high degree of interactivity.

### Data Collection and Sharing:

The management of importing raw data into the visualisation tool and afterwards exporting visual reports in diverse formats is a task that necessitates firm oversight and customization according to its preferences. Certain datasets can be directly inputted into the tool in their unprocessed state, while others necessitate prior aggregation due to their excessive size. In certain instances, data may be derived from a singular source, however in other cases, it necessitates aggregation from multiple sources and

subsequent visualisation through a suitable technology. Additionally, the tool should include the capability to distribute the reports across team members and other relevant individuals involved in the project. It is imperative that the reports possess the capability to be exported to various other programmes.

#### **Geo-tagging and Location Intelligence:**

The increasing interconnectedness of global business necessitates the incorporation of location intelligence within data visualisation tools. From where is the data sourced? Which states or regions demonstrate higher levels of engagement with the services and which places exhibit a greater need for improvement? The capacity to arrange multiple types of data in a chronological and spatial manner holds significance for enterprises that require monitoring key performance indicators (KPIs) depending on geography.

#### **Collaboration:**

Big data visualisation technologies with real-time collaboration features enable staff to discuss their findings in more depth. Instead than needing them to email each other static files and screenshots, this allows employees to collaborate in real time on current data.

#### **Streaming data support:**

Enterprises are currently confronted with the challenge of managing substantial quantities of intricate, real-time data originating from diverse sources. Numerous visualisation tools employ older back ends that rely on structured batch data analysis. Analysing data of an extreme kind in real-time is a considerable challenge. The inclusion of support for streaming data might facilitate the exploration of various visualisation applications that involve data derived from social media, internet of things devices and mobile applications.

#### **Artificial Intelligence Integration:**

Visualisation tools for big data are now exploring the integration of machine learning, deep learning and natural language processing techniques in order to enhance the ease of analysis, exploration, prediction and action recommendation.

The selection of an appropriate data visualisation tool is a significant undertaking, since it entails considerable financial investment and exerts a substantial influence on the formulation of corporate strategies. A tool capable of delivering visually engaging and precise reports has the potential to enhance decision-making, strategic planning and Key Performance Indicator (KPI) monitoring for company professionals. The selection of a suitable tool should be based on the specific features that hold the most importance for the business, ensuring that it provides the necessary representations.

#### **Summary**

- Exploratory Data Analysis (EDA) is a critical initial step in the data analysis process that involves visually and statistically summarizing, understanding and visualizing the main characteristics of a dataset. EDA helps analysts and data scientists to gain insights into the data, identify patterns, detect anomalies and inform subsequent analysis or modelling.
- Feature extraction is a fundamental process in data analysis and machine learning, where raw data is transformed into a more compact and informative representation. Feature extraction involves selecting or creating a subset of relevant features (variables) from the original dataset.

## Notes

- Univariate analysis focuses on a single variable in the dataset. The goal is to understand the distribution, central tendency, dispersion and other properties of that variable.
- Bivariate analysis involves the study of two variables to understand their relationship or association. It helps uncover patterns, dependencies, correlations and potential causal relationships between the two variables.
- Dimensionality reduction is a technique used in data analysis and machine learning to reduce the number of features (dimensions) in a dataset while preserving as much relevant information as possible. High-dimensional datasets can pose challenges such as increased computational complexity, noise sensitivity and overfitting.
- Principal Component Analysis (PCA) is a dimensionality reduction technique used in statistics and machine learning to transform high-dimensional data into a lower-dimensional representation while retaining as much of the original data's variance as possible. It is widely used for data visualization, noise reduction and feature extraction.

### Glossary

- ❖ EDA: Exploratory Data Analysis
- ❖ PCA: Principal Component Analysis
- ❖ KPI: Key Performance Indicator
- ❖ PDF: Probability Density Functions
- ❖ AWS: Amazon Web Services
- ❖ NLP: Natural Language Processing
- ❖ GIS: Geographic Information System
- ❖ HOG: Histogram of Oriented Gradients
- ❖ CNN: Convolutional Neural Network
- ❖ ICA: Independent Component Analysis
- ❖ SVD: Singular Value Decomposition
- ❖ LDA: Linear Discriminant Analysis
- **Histogram:** The graph being referred to is a fundamental visualisation commonly employed in exploratory data analysis (EDA). A histogram is a graphical representation in the form of a bar plot that illustrates the frequency or proportion of occurrences within certain intervals, sometimes referred to as bins, for a given variable.
- **Boxplots:** Boxplots, alternatively referred to as box-and-whisker plots, offer a graphical representation that succinctly summarises the central tendency, dispersion and outliers of a distribution.

### Check Your Understanding

1. What is the primary goal of Exploratory Data Analysis (EDA)?
  - a. To build predictive models directly.
  - b. To summarize data using only descriptive statistics.
  - c. To visually represent data without performing any calculations.
  - d. To understand the main characteristics of the data and generate hypotheses.
2. Which visualization is commonly used in EDA to display the distribution of a continuous variable?

3. In EDA, which statistical measure is used to quantify the strength and direction of a linear relationship between two continuous variables?

  - Bar chart
  - Pie chart
  - Box plot
  - Contingency table

4. Which step in EDA involves handling missing values, detecting outliers and ensuring data quality?

  - Median
  - Variance
  - Pearson's correlation coefficient
  - Mode

5. Univariate analysis focuses on:

  - Exploring the relationship between two variables.
  - Analyzing a single variable in isolation.
  - Identifying outliers and anomalies.
  - Comparing multiple variables simultaneously.

6. Which of the following is an appropriate visualization for univariate analysis of a continuous variable?

  - Scatter plot
  - Bar chart
  - Histogram
  - Box plot

7. Bivariate analysis involves:

  - Analyzing a single variable.
  - Exploring the distribution of a single variable.
  - Studying the relationship between two variables.
  - Creating summary statistics for a dataset.

8. Correlation coefficients are used in bivariate analysis to measure:

  - The strength and direction of a linear relationship between two continuous variables.
  - The difference between two categorical variables.
  - The impact of an outlier on a scatter plot.
  - The variance of a single variable.

9. What is the primary goal of feature extraction in data analysis and machine learning?

  - To create new data points from scratch.
  - To visualize data in higher dimensions.
  - To reduce the number of features while retaining relevant information.
  - To add noise to the dataset for better generalization.

10. Which of the following is a dimensionality reduction technique used for feature extraction?

  - One-hot encoding
  - Normalization
  - Principal Component Analysis (PCA)
  - Cross-validation

**Notes**

11. In the context of Natural Language Processing (NLP), which technique is commonly used for feature extraction from text data?
  - a. Clustering
  - b. Singular Value Decomposition (SVD)
  - c. K-Means
  - d. Decision Trees
12. How does feature extraction differ from feature selection?
  - a. Feature extraction creates new features, while feature selection eliminates existing ones.
  - b. Feature extraction increases dimensionality, while feature selection reduces dimensionality.
  - c. Feature extraction retains relevant information, while feature selection removes irrelevant features.
  - d. Feature extraction is used for categorical data, while feature selection is used for numerical data.
13. Dimensionality reduction is primarily used to:
  - a. Increase the number of features in a dataset.
  - b. Improve visualization of high-dimensional data.
  - c. Add noise to the data for better generalization.
  - d. Make the data more complex.
14. Which of the following is a feature selection technique?
  - a. Principal Component Analysis (PCA)
  - b. Singular Value Decomposition (SVD)
  - c. Recursive Feature Elimination (RFE)
  - d. Independent Component Analysis (ICA)
15. In Principal Component Analysis (PCA), the principal components are:
  - a. Linear combinations of the original features.
  - b. Transformed feature space with a higher dimension.
  - c. A subset of the original features.
  - d. Created using deep learning techniques.
16. What is a key consideration when applying dimensionality reduction techniques?
  - a. Maximizing the number of features retained.
  - b. Avoiding any loss of information.
  - c. Evaluating the impact on model performance.
  - d. Using all available features without any preprocessing.
17. Data visualization is primarily used for:
  - a. Adding complexity to data analysis.
  - b. Reducing the need for statistical analysis.
  - c. Displaying patterns, trends and insights in data.
  - d. Replacing the need for data preprocessing.
18. Which type of data visualization is most suitable for showing the distribution of a single numerical variable?

## Notes

## Exercise

1. What is exploratory data analysis (EDA)? Also define the process for EDA.
  2. Define univariate and bivariate analysis.
  3. Explain feature extraction.
  4. What is the need for dimensionality reduction?
  5. Define principle component analysis.
  6. Define various data visualisation concepts.

## Learning Activities

1. Provide a step-by-step outline of how you would perform EDA on a real-world dataset.
  2. Provide examples of real-world scenarios where dimensionality reduction can be beneficial.

### **Check Your Understanding- Answers**

- |       |       |       |       |
|-------|-------|-------|-------|
| 1. d  | 2. c  | 3. c  | 4. c  |
| 5. b  | 6. c  | 7. c  | 8. a  |
| 9. c  | 10. c | 11. b | 12. c |
| 13. b | 14. c | 15. a | 16. c |
| 17. c | 18. c | 19. c | 20. b |

#### **Further Readings and Bibliography**

1. Exploratory Data Analysis by John W. Tukey
  2. Applied Data Science by Kelleher, Mac Namee and D'Arcy
  3. Dimensionality Reduction by C. Burges, A. Smola and P. Bartlett
  4. Applied Multivariate Statistical Analysis by Richard A. Johnson and Dean W. Wichern

## Module - III: Supervised Learning

### Learning Objectives

At the end of this module, you will be able to:

- Exploratory Data Analysis by John W. Tukey
- Applied Data Science by Kelleher, Mac Namee and D'Arcy
- Dimensionality Reduction by C. Burges, A. Smola and P. Bartlett
- Applied Multivariate Statistical Analysis by Richard A. Johnson and Dean W. Wichern

### Introduction

Supervised learning is a well studied domain within the realm of machine learning, attracting significant scholarly attention and research efforts. A variety of methodologies in supervised learning have been effectively utilised in the examination and manipulation of multimedia data. The key characteristic that sets supervised learning apart is the existence of annotated training data. The term “supervisor” denotes an entity that offers direction to the learning system in determining the suitable labels to apply to training samples. In the context of classification problems, it is frequently seen that these labels are utilised to represent class labels.

Supervised learning techniques involve the construction of models based on a provided dataset for training purposes. These models are then utilised to classify unlabeled input. The current research offers an examination of supervised learning, situated within the theoretical framework of risk minimization theory. This work provides a thorough analysis of support vector machines and nearest neighbour classifiers, which are widely acknowledged as the primary supervised learning techniques employed in the domain of multimedia research.

Supervised Learning is a machine learning paradigm that entails acquiring knowledge on the input-output correlation of a system through the utilisation of a pre-defined collection of paired input-output training instances. The label of the input data or the supervision is often used to denote the outcome of a given input. Therefore, it is common to refer to a training sample that includes both the input and its corresponding output as labelled training data or supervised data. Sometimes, it is alternatively referred to as Learning with a Teacher, Learning from Labelled Data, or Inductive Machine Learning. The primary goal of supervised learning is to develop an artificial system that can effectively learn and understand the correlation between input and output variables. This acquired knowledge enables the system to make precise predictions about the output when provided with new and unseen inputs. If the outcome comprises a finite set of unique values that symbolise the categorization labels of the input.

### 3.1 Introduction to Supervised Learning:

In supervised learning, a mapping between a set of input variables  $X$  and an output variable  $Y$  is learned and this mapping is then used to forecast the results for hypothetical data. The most crucial machine learning methodology is supervised learning, which is equally crucial for handling multimedia data. Here, we concentrate on supervised learning methods based on kernels. We examine support vector machines, which are now the most popular supervised learning method, especially for

handling multimedia data. We also discuss closest neighbor classifiers, which are roughly categorized as a kernel-based approach.

Because of the emphasis on similarity being acceptable for multimedia data and the wide variety of similarity assessment techniques available, nearest neighbor techniques are popular in the field of multimedia. In order to wrap up this overview of supervised learning, we'll also talk about the ensemble concept, a key tactic for improving a classifier's stability and accuracy that entails replacing a single classifier with a group of classifiers. An overview of the fundamentals of statistical learning theory serves as the unit's introduction since it offers a broad framework for evaluating learning algorithms and practical tools for handling applications in the real world. Before introducing various methods, we introduce fundamental statistical learning concepts and theorems.

## Notes

### 3.1.1 Introductory Concepts of Supervised Learning:

Supervised learning is a machine learning paradigm wherein machines utilise appropriately labelled training data to anticipate the output. This training data is employed to train the machines. The concept of "labelled data" pertains to input data that has been previously allocated the corresponding output. Supervised learning involves the utilisation of training data as a form of guidance for computers, enabling them to accurately predict the desired output. It utilises the same concept that a student would acquire knowledge under the guidance of an instructor.

The process of supervised learning entails providing the machine learning model with appropriate input data as well as corresponding output data. The objective of a supervised learning algorithm is to identify a mapping function that establishes a relationship between the input variable ( $x$ ) and the output variable ( $y$ ). Supervised learning exhibits practical utility in various domains, including but not limited to risk assessment, picture categorization, fraud detection and spam filtering.

#### How Supervised Learning Works:

Supervised learning involves the training of models using a dataset that is annotated with labels, enabling the model to acquire knowledge about different types of input. After the completion of the training process, the model undergoes testing using a designated test dataset, which is a subset of the original training set. Subsequently, the model generates predictions based on this evaluation.

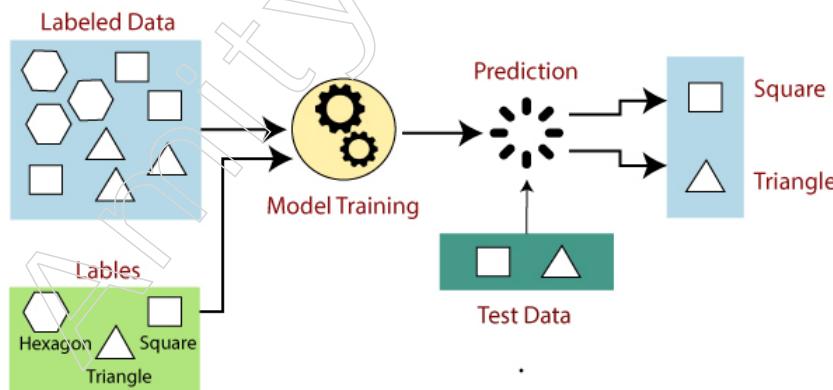


Figure: Working of Supervised Learning

Consider a dataset comprising several forms, encompassing square, rectangle, triangle and polygon. The initial phase entails training the model for every individual shape.

## Notes

- A shape with four sides, all of which are of equal length, is classified as a square.
- A form with three sides will be classified as a triangle.
- A form with six equal sides is classified as a hexagon.
- Following the completion of the training phase, the model is subjected to testing using the designated test set. The primary objective of the model is to accurately discern and classify the shape.

The computer has undergone comprehensive training including various shapes. Upon encountering a novel shape, it employs a classification process based on the number of sides to determine the appropriate categorization and then generate a prediction.

### Steps Involved in Supervised Learning:

- The initial step is to select the type of training dataset.
- Acquire the training data through the utilisation of labels.
- The process of partitioning the training dataset into training, test and validation datasets is undertaken.
- The training dataset's input features should possess comprehensive information to facilitate accurate prediction of the output.
- Selecting the optimal approach for the model, such as a decision tree or a support vector machine, is crucial.
- Utilise the algorithmic procedure using the provided practise dataset. Validation sets, which are subsets of training datasets, are sometimes necessary for the purpose of controlling parameters.
- Utilise the test set in order to ascertain the accuracy of the model. If the model accurately predicts the outcome, it might be deemed as being accurate.

**Types of supervised Machine learning Algorithms:** Supervised learning can be further categorized into two distinct types:

- Regression
- Classification

**Regression:** Regression techniques are employed when there exists a discernible association between the input variable and the output variable. This technique is employed to make predictions about continuous variables, such as weather patterns, market trends and other similar phenomena. In the context of regression analysis, it is common to assign a continuous numerical value to each observed data point. One possible instance that may be used to illustrate this concept is the process of estimating the weight of a bovine creature, utilizing various measured attributes such as the animal's height and length. The following are several widely used regression methods that fall within the domain of supervised learning.

Regression analysis is a statistical method used to examine the association between a dependent variable (also known as the target variable) and one or more independent variables (also known as predictors). This methodology is employed for the purpose of making predictions, modelling time series data and determining the causal relationship between variables. The examination of the correlation between reckless driving behavior and the frequency of road accidents caused by a driver is most effectively conducted using regression analysis.

### Why Do We Use Regression Analysis?

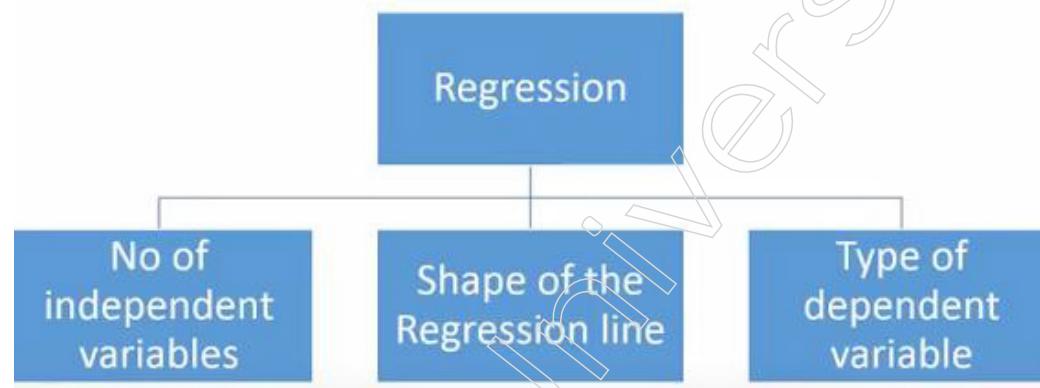
There exist numerous advantages associated with the utilization of regression analysis. The following points can be identified:

- It serves to demonstrate the substantial associations between a dependent variable and an independent variable.
- It serves to quantify the magnitude of the influence of several independent variables on a dependent variable.

Regression analysis enables the comparison of variables assessed on disparate scales, such as the impact of price fluctuations and the quantity of promotional endeavors. These advantages assist market researchers, data analysts and data scientists in the process of eliminating and evaluating the optimal collection of variables for constructing predictive models.

### **Types of Regression Techniques:**

There exists a diverse range of regression procedures that can be utilized for the purpose of making predictions. The strategies primarily rely on three metrics, namely the number of independent variables, the type of dependent variables and the shape of the regression line. The subsequent sections will provide a detailed discussion of these topics.



If you feel the need to apply a combination of the aforementioned parameters that hasn't been used before, you can even concoct brand-new regressions for the inventive ones. Let's first learn the most widely utilized regressions before you begin.

### **Linear Regression:**

This modelling technique is widely recognized and acknowledged. Linear regression is commonly one of the initial subjects that individuals choose to study when acquiring knowledge in the field of predictive modelling. This technique involves a continuous dependent variable, independent variable(s) that can be either continuous or discrete and a linear regression line.

Linear regression is a statistical method that aims to construct a mathematical link between a dependent variable (Y) and one or more independent variables (X). This relationship is represented by a straight line, commonly referred to as the regression line, which is determined by a process of finding the best fit.

The equation  $Y=a+b*X + e$  represents a mathematical model, where Y is the dependent variable, X is the independent variable, a represents the intercept of the line, b represents the slope of the line and e represents the error term. The aforementioned equation possesses the capability to forecast the value of the target variable by utilizing the provided predictor variable(s).

## Notes

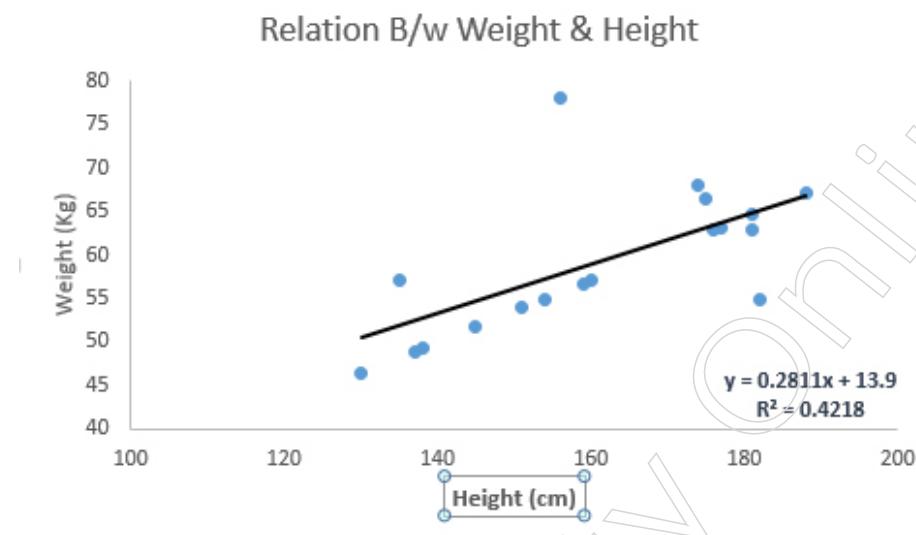


Figure: Linear Regression

The distinction between simple linear regression and multiple linear regression is that the former has more than one independent variable, whilst the latter has only one.

How to obtain best fit line (Value of a and b): The Least Square Method can be employed to readily achieve this goal. The strategy employed for fitting a regression line that is most frequently utilized is known as the most common approach. The method employed in this calculation involves reducing the sum of the squares of the vertical deviations between each data point and the line, resulting in the determination of the best-fit line for the observed data. The squaring of deviations ensures that positive and negative numbers do not cancel out when they are summed.

### Important Points:

- It is important for there to exist a linear correlation between the independent and dependent variables.
- Multiple regression is susceptible to issues such as multicollinearity, auto correlation and heteroskedasticity.
- Linear regression is highly susceptible to the influence of outliers. Multicollinearity has the potential to significantly impact the regression line and subsequently the projected values. Specifically, it can lead to an increase in the variance of the coefficient estimates and render them highly susceptible to even modest alterations in the model. The outcome of the analysis indicates that the coefficient estimates exhibit instability.
- When faced with several independent variables, researchers have several approaches to select the most significant ones. These include forward selection, backward elimination and stepwise selection.

### Logistic Regression:

Logistic regression is a statistical method employed to estimate the likelihood of an event occurring, specifically the events of success and failure. Logistic regression is the appropriate statistical technique to employ when the dependent variable exhibits a binary nature, characterized by two possible outcomes such as 0/1, True/False, or Yes/No. The value of Y varies between 0 and 1 and can be expressed using the following equation. The odds, denoted as  $p/(1-p)$ , are the ratio of the chance of an event occurring to the probability of the event not occurring. The natural logarithm of the chances can be

expressed as the natural logarithm of the probability divided by one minus the probability. The equation  $\text{logit}(p) = \ln(p/(1-p)) = b_0 + b_1X_1 + b_2X_2 + b_3X_3 + \dots + b_kX_k$  is commonly used in academic research. The variable  $p$  represents the probability associated with the occurrence of the attribute under consideration. One pertinent inquiry to be posed in this context is, "What is the rationale behind the utilization of logarithms in the equation?"

Given that we are dealing with a binomial distribution as the dependent variable, it is necessary to select an appropriate link function that is most suitable for this particular distribution. The function in question is the logit function. The parameters in the aforementioned equation are selected with the objective of maximizing the likelihood of observing the sample values, as opposed to minimizing the sum of squared errors, as typically done in ordinary regression.

### **Polynomial Regression:**

A regression equation can be classified as a polynomial regression equation when the exponent of the independent variable exceeds 1. The provided equation can be identified as a polynomial equation, namely in the form of  $y = a + b*x^2$ . In the context of regression analysis, the optimal fitting line does not adhere to a linear trajectory. The curve can be described as a suitable mathematical model that aligns with the given data points.

#### **Important Points:**

- While the desire to fit a higher degree polynomial may exist in order to obtain a smaller error, doing so may lead to over-fitting. Plot the relationships consistently to check the fit and pay close attention to whether the curve matches the nature of the issue.
- Pay close attention to shapes and trends that bend toward the ends and determine whether they make sense. Higher polynomials may result in strange extrapolation outcomes.

### **Ridge Regression:**

When the data exhibits multicollinearity—highly connected independent variables—Ridge Regression is the method utilized. Even though the least squares estimates (OLS) in multicollinearity are unbiased, their enormous variances cause the observed value to differ much from the true value. Ridge regression lowers the standard errors by biasing the regression estimates to some extent.

#### **Important Points:**

- Ridge regression lowers the value of coefficients but doesn't approach zero, suggesting no feature selection feature.
- The assumptions of this regression are similar to least squared regression with the exception that normality is not to be assumed.
- L2 regularization is used in this regularization technique.

**How to Select the Right Regression Model:** Typically, life exhibits a semblance of simplicity when one possesses a limited repertoire of one or two skills. One training institute advises its students to utilize linear regression when dealing with continuous outcomes. If the data follows a binary distribution, logistic regression is an appropriate statistical technique to employ. Nevertheless, as the number of accessible possibilities increases, the task of selecting the most appropriate one gets increasingly challenging.

A comparable scenario arises in the context of regression models. When considering various regression models, it is crucial to select the most appropriate technique depending on the nature of the independent and dependent variables, the dimensionality of the data

## Notes

and other pertinent properties of the dataset. Outlined below are the fundamental factors that one should adhere to while selecting an appropriate regression model:

- Creating a prediction model inevitably involves data investigation. Before choosing the appropriate model, you should first establish the link and influence of the various factors.
- The statistical significance of parameters, R-square, Adjusted r-square, AIC, BIC and error term can be utilised to assess and compare the goodness of fit across different models. Another criterion is the Mallow's Cp. By doing a comparison between your model and all possible submodels, or a strategically selected subset of them, you are effectively assessing the presence of any potential bias in your model.
- The user's text does not provide any information to be rewritten in an academic manner. Cross-validation is often regarded as the most effective approach for evaluating prediction models. The data collection process is divided into two distinct groups, namely the training and validation sets. The prediction accuracy can be assessed by calculating the simple mean squared difference between the observed and anticipated numbers.
- The utilisation of the automatic model selection method is not recommended in cases where the data set contains numerous confounding variables, as it is undesirable to incorporate all of these variables simultaneously into a model. Furthermore, the outcome will be contingent upon the objective you are aiming to achieve. In certain instances, it is possible for a less powerful model to exhibit greater ease of use compared to a more statistically significant model.
- The user's text does not contain any information. Regression regularisation techniques, such as Lasso, Ridge and ElasticNet, have been found to be useful in situations when the variables in the dataset exhibit high dimensionality and multicollinearity.

**Classification:** Classification methods are employed when the output variable is categorical, consisting of two classes, such as Yes-No, Male-Female, True-False and so on. The categorization or classification of data into distinct groups or classes, or the assignment of a discrete target value, represents another prominent instance of supervised activities.

Various other challenges, such as time series forecasting, medical risk assessment and pixel-wise image segmentation, can be formulated and understood as regression or classification problems. The consideration of target values enables the establishment and evaluation of precise quality criteria, such as the projected mean square error (MSE) in regression or the number of misclassifications in a particular test dataset.

When contrasting unsupervised tasks with supervised learning, it becomes evident that the latter exhibits a greater level of well-definedness. The presence of well-defined quality standards in the training data allows for the utilisation of objective functions that are naturally applicable, enabling effective guidance of the learning process.

A number of difficulties, including the choice of an appropriate model, must be carefully addressed in supervised learning as well. Systems that are incompatible, oversimplified, or excessively complex might make learning more difficult. Similar to how training process specifics can have a significant impact on performance. The success of supervised training also depends on the actual depiction of the observations and the choice of the proper features.

In the following, we will mostly consider a prototypical work flow of supervised learning where:

- In a training phase, a model or hypothesis about the target rule is developed by analyzing a set of labeled samples. This could be accomplished, for example, by adjusting a feed-forward neural network's weights.
- After training, the working phase allows the learnt model, such as the network, to be used with new data.

### Applications:

Supervised Learning facilitates the acquisition of knowledge by a machine pertaining to human or object behavior in certain activities. The acquired information can afterwards be utilized by the machine to execute analogous actions on these assignments.

Given the potential for computer gear to execute input-output mappings with greater speed and durability than humans, machines that possess a proficient supervised learning capability have the ability to outperform humans in specific jobs with increased efficiency and accuracy.

However, the current Supervised Learning algorithms are still unable to replicate the complex learning capabilities of humans due to constraints in hardware, software and algorithm designs. Risk assessment is a crucial process in the financial services and insurance sectors, aimed at mitigating the risk exposure of firms. To do this, supervised learning techniques are employed.

Image categorization is a prominent application that showcases the utilization of supervised machine learning. As an illustration, Facebook has the capability to identify an individual's buddy inside a photograph sourced from a collection of images that have been labeled with tags.

Fraud detection refers to the process of determining the authenticity of user transactions. Visual recognition refers to the capacity of a machine learning model to accurately discern and classify various elements such as objects, locations, individuals, activities and visual representations.

### General issues of supervised learning algorithms:

The method of inductive machine learning involves acquiring a collection of rules based on cases, or examples, within a training set. In a broader sense, it entails constructing a classifier that has the ability to generalize from novel occurrences. The application of supervised machine learning to a real-world situation is depicted in the Figure below.

The initial stage involves the collection of the dataset. If an expert with the necessary expertise is present, they may provide recommendations for the domains that offer the highest level of information. If not, then the most straightforward approach is that of "brute-force," wherein all accessible measurements are taken in the expectation that the appropriate (informative, relevant) traits can be identified. Nevertheless, a dataset obtained by the "brute-force" approach is not inherently appropriate for the process of induction. In the majority of instances, the data set is characterized by the presence of noise and missing feature values, necessitating extensive pre-processing procedures.

The subsequent stage involves the preparation and preprocessing of the data. Researchers have a range of approaches at their disposal to address missing data, which can be selected based on the specific circumstances. Hodge and Austin (year) have recently conducted a comprehensive assessment on modern methodologies

## Notes

employed for the detection of outliers, also referred to as noise. The researchers have recognized the advantages and cons of these strategies.

The utilization of instance selection is not solely limited to addressing noise, but also serves as a means to manage the impracticality of learning from excessively big datasets. The process of instance selection in these datasets can be viewed as an optimization problem, wherein the objective is to preserve the quality of data mining while simultaneously lowering the size of the sample.

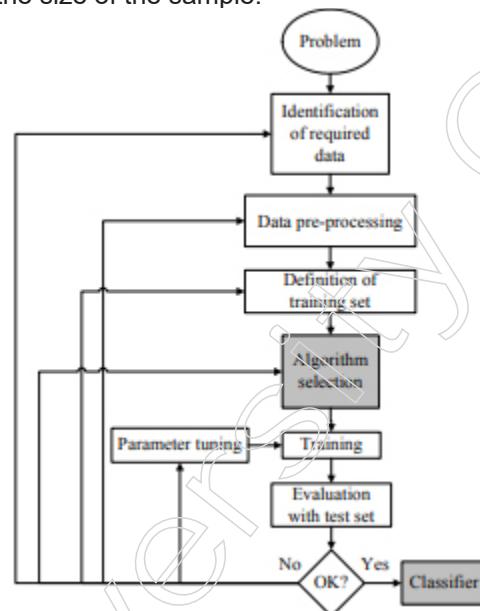


Figure: The process of supervised ML

The act of discovering and eliminating as many unnecessary and redundant features as you can is known as feature subset selection. As a result, the data's dimensions are reduced, making it possible for data mining algorithms to work more quickly and efficiently. The accuracy of supervised ML classification models is frequently unreasonably influenced by the fact that numerous features depend on one another. By building new features from the existing feature set, this issue can be solved.

The term "feature construction/transformation" refers to this method. The development of clearer and more precise classifiers could be facilitated by these recently developed features. Additionally, the identification of significant traits helps people comprehend the acquired topic and the created classifier more clearly.

### Algorithm selection:

A crucial decision is which precise learning algorithm to employ. When the results of preliminary testing are deemed good, the classifier (which converts unlabeled instances into classes) is ready for regular use. Prediction accuracy (the proportion of right predictions divided by the total number of predictions) is frequently used as the basis for the classifier's evaluation.

The accuracy of a classifier is calculated using at least three different methods. One method is to divide the training set in half, with one third used for training and the other third for performance estimation. Another method is called cross-validation and it divides the training set into equal-sized, mutually exclusive subsets. For each subset, the classifier is trained using the union of all the other subsets. The classifier's error rate can therefore be estimated using the average error rate of each subset. Cross validation is a specific example of leave-one-out validation. One instance is present in each test subset.

The computational cost of this method of validation is higher, but it is beneficial when the most precise estimate of a classifier's error rate is needed.

#### Algorithm process:

Statistical comparisons of trained classifier accuracies on certain datasets are a frequent way to compare supervised machine learning methods. If there is an adequate amount of data available, we can sample a number of training sets of size N, apply the two learning algorithms to each one and then use a sizable test set to measure the accuracy of each pair of classifiers. The variance of these differences is an estimate of the variance of the classifier in the entire set and their average represents an estimate of the predicted difference in generalization error across all feasible training sets of size N.

The paired t-test will be employed to examine the null hypothesis that there exists no significant difference in means between the classifiers. This test may provide two distinct types of faults. The phenomenon in which a test erroneously rejects the null hypothesis, so indicating a significant difference despite the absence of one, is referred to as a type I error. The probability of accepting the null hypothesis in the presence of a true difference is sometimes referred to as a type II error. The probability of committing a Type I mistake in the test will be approximately equal to the chosen level of significance.

Given the desire for ease in replicating published experimental findings, it is preferable for the outcomes of the test to be unaffected by the particular partitioning resulting from the randomization technique. Nevertheless, the implementation of partitioning in practical scenarios sometimes involves a certain level of sensitivity. In order to assess the replicability of our findings, it is necessary to conduct multiple iterations of the same test on identical data, with diverse random partitionings. Typically, this involves repeating the test ten times and observing the frequency at which consistent results are obtained.

Supervised categorization is a prevalent activity carried out by Intelligent Systems. Consequently, a multitude of techniques have been developed in the field of Artificial Intelligence, including Logical/Symbolic methods, Perceptron-based methods and Statistical approaches such as Bayesian Networks and Instance-based methods.

**The Simple Linear Regression Model:** The most elementary deterministic mathematical association between two variables, denoted as x and y, is a linear relationship expressed as  $y = \beta_0 + \beta_1 x$ . The aim of this section is to construct a linear probabilistic model that is equal to the original model. If two variables exhibit probabilistic dependence, then, given a constant value of x, there exists inherent uncertainty regarding the value of the second variable. The equation  $Y = \beta_0 + \beta_1 x + \varepsilon$  is postulated, with  $\varepsilon$  representing a stochastic variable. Two variables have a linear relationship "on average" when, for a constant value of x, the observed value of Y deviates from its predicted value due to a random component, or in other words, random error is present.

**A Linear Probabilistic Model:** The model equation  $Y = \beta_0 + \beta_1 x + \varepsilon$  incorporates parameters  $\beta_0$ ,  $\beta_1$  and  $\sigma^2$ , where Y represents the dependent variable and x represents the independent variable. It is important to note that for any given fixed value of x, the dependent variable Y is a random variable that is associated with x through this model equation. The variable  $\varepsilon$  in the model equation represents the "error," which is a random variable. It is assumed to have a symmetric distribution with an expected value of  $E(\varepsilon) = 0$  and a variance of  $V(\varepsilon) = \sigma^2$ . At this point, no assumptions have been made about the specific distribution of  $\varepsilon$ .

- **X:** the independent, predictor, or explanatory variable (usually known). NOT RANDOM.

## Notes

- **Y:** The dependent or response variable. For fixed  $x$ ,  $Y$  will be random variable.
  - **$\epsilon$ :** The random deviation or random error term. For fixed  $x$ ,  $\epsilon$  will be random variable.
- The coordinates  $(x_1, y_1), \dots, (x_n, y_n)$  obtained from  $n$  distinct observations will exhibit dispersion around the actual regression line.

### 3.2 Introduction to Linear Regression:

**Regression:** Regression analysis is a statistical technique that use one or more independent variables to elucidate the association between a dependent variable, also known as the target variable and independent variables, also referred to as predictor variables. Regression analysis allows us to gain a deeper understanding of how the value of the dependent variable varies in relation to an independent variable, while keeping other independent variables constant. It forecasts real, continuous values like temperature, age, salary and cost, among others. Using the example below, we may comprehend the notion of regression analysis:

#### Linear regression:

A straightforward method for supervised learning, predictive analysis and quantitative response is linear regression. It is among the simplest approaches to take into account.

Linear regression is a highly robust statistical methodology. A considerable number of individuals possess a certain level of acquaintance with regression analysis from their exposure to news articles, wherein scatterplots are often accompanied by superimposed straight lines. Linear models are commonly employed for the purposes of prediction or to assess the presence of a linear association between two numerical variables. There exist two variables that can be accurately represented by a linear connection. The mathematical expression representing a straight line is given by:

$$y = 5 + 57.49x$$

If there were a perfect linear relationship, you could determine the precise value of  $y$  by only knowing the value of  $x$ . In practically any natural process, this is implausible. For instance, if we took family income as  $x$ , this value would offer some beneficial insight into the potential financial aid  $y$  that a college might provide for a potential student. Even when comparing pupils whose families had similar financial circumstances, there would yet be variation in financial support.

The assumption underlying linear regression is that a straight line may adequately represent the connection between two variables,  $x$  and  $y$ :

$$y = \beta_0 + \beta_1 x$$

where  $0$  and  $1$  are two model parameters for a linear model and  $0$  and  $1$  are represented by the Greek letter beta. (This usage of the  $0$  and  $1$  Linear model parameters of is unrelated to the we used to represent the likelihood of a Type II error.)

These parameters are calculated using data and their point estimates are denoted by the letters  $b_0$  and  $b_1$ . Typically, when  $x$  is used to predict  $y$ ,  $x$  is referred to as the explanatory or predictor variable, while  $y$  is referred to as the outcome.

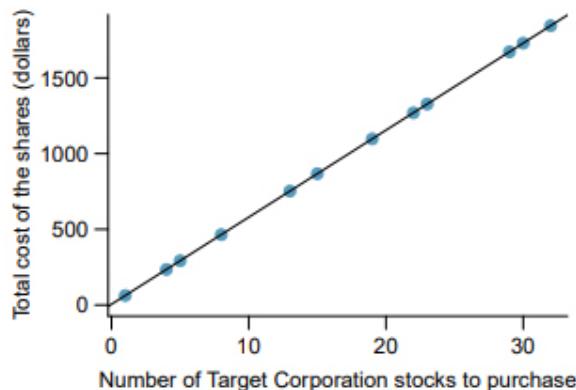


Figure: two variables perfectly with straight line

**Figure shows:** Twelve distinct buyers concurrently submitted requests to a trading company for the acquisition of Target Corporation stock, with the aggregate cost of the shares being duly recorded. Due to the utilization of a linear method for cost computation, the linear fit is deemed to be flawless.

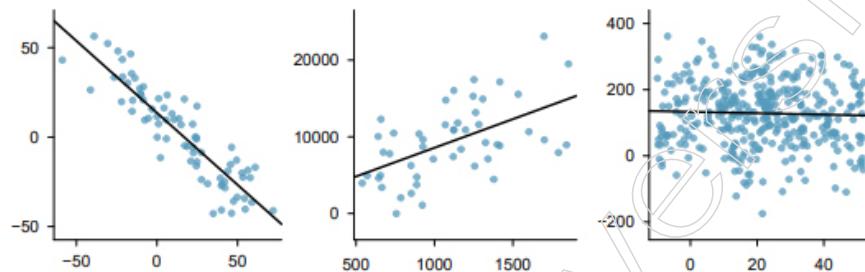
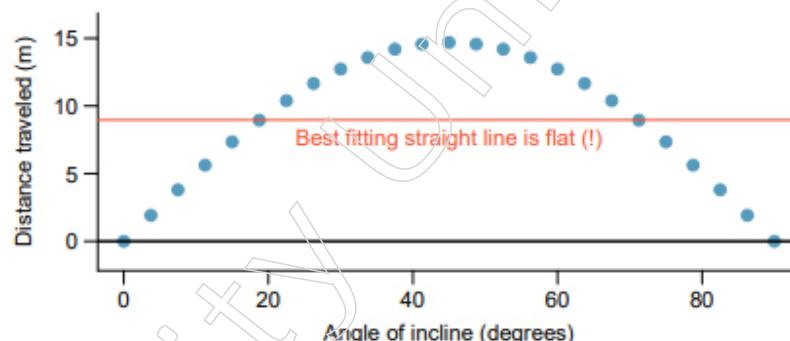


Figure: Three data sets where a linear model may be useful even though the data do not all fall exactly on the line.

Figure: A linear model is not useful in this nonlinear case.  
These data are from an introductory physics experiment.

### 3.2.1 Linear Regression: Important Terms:

#### Linear Regression:

Linear regression is a statistical technique that aims to establish a mathematical model representing the association between two variables by employing a linear equation to best fit the observed data. In this context, it is customary to designate one variable as the independent or explanatory variable, while the other is designated as the dependent variable. An instance where a modeler may seek to establish a connection between the weights and heights of individuals is through the utilization of a linear regression model. Prior to fitting a linear model to the observed data, it is imperative for the modeler to ascertain the presence or absence of a link between the variables under consideration.

## Notes

This observation does not necessarily indicate a causal relationship between the variables (e.g., better SAT scores do not causally lead to higher college grades), but rather suggests a notable correlation between the two variables.

The utilization of a scatterplot can prove to be a valuable instrument in assessing the degree of correlation between two variables. If there is no apparent correlation between the suggested independent and dependent variables (i.e., the scatter plot does not exhibit any discernible upward or downward patterns), it is unlikely that applying a linear regression model to the dataset will yield a meaningful model. The correlation coefficient is a significant quantitative measure of the relationship between two variables. It is a numerical number ranging from -1 to 1, representing the degree of association detected in the data for the two variables. The equation of a linear regression line can be expressed as  $Y = a + bX$ , where  $X$  represents the explanatory variable and  $Y$  represents the dependent variable. The line's slope is denoted as  $b$ , whereas the intercept, representing the value of  $y$  when  $x$  equals zero, is referred to as  $a$ .

- Linear regression is a statistical regression approach commonly employed in predictive analysis.
- Regression is a commonly used algorithm that effectively models the relationship between continuous variables.
- This technique is commonly employed in the context of regression problems in machine learning.
- Linear regression, as its name suggests, illustrates a linear relationship between the independent variable ( $X$ -axis) and the dependent variable ( $Y$ -axis).
- Simple linear regression is the term used to describe a form of linear regression where there is a sole input variable ( $x$ ). Moreover, the aforementioned form of regression analysis is sometimes referred to as multiple linear regression when there exists a multitude of input variables.
- The provided visual representation can be utilised to elucidate the interconnection among the variables within the linear regression model. In this study, we are doing an estimation of an employee's compensation by considering their level of experience over time.
- Below is the mathematical equation for Linear regression
  - ❖  $Y = aX + b$  Here,  $Y$  = dependent variables (target variables),  
 $X$  = Independent variables (predictor variables),  
 $a$  and  $b$  are the linear coefficients.

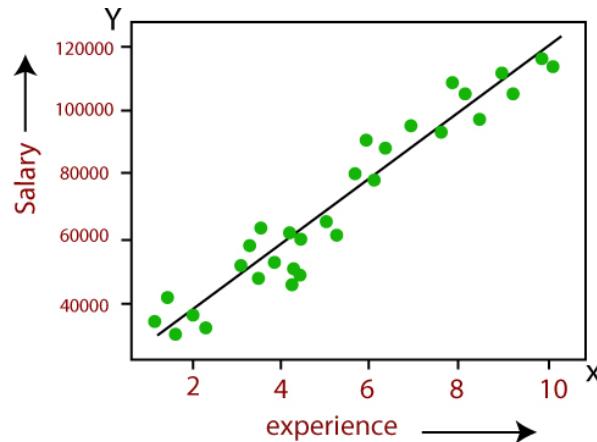


Figure: linear regression model

**Least-Squares Regression:** The method of least-squares is often utilized as the primary approach for fitting a regression line. The proposed approach aims to determine the optimal line that best fits the given data by minimizing the total of the squared vertical deviations between each data point and the line. In cases when a data point sits perfectly on the fitted line, its vertical deviation is considered to be zero. Due to the process of squaring the deviations before summing them, the occurrence of cancellations between positive and negative values is eliminated.

**Example:** The dataset titled “Televisions, Physicians and Life Expectancy” encompasses many variables, including the population-to-television ratio and the population-to-physician ratio for a total of 40 countries. Given that both variables likely represent the measure of economic prosperity inside each respective country, it is justifiable to suggest that there exists a positive correlation between them.

After excluding 8 countries with missing data from the dataset, the remaining 32 nations exhibit a correlation coefficient of 0.852 between the variables representing the number of individuals per television set and the number of individuals per physician. The coefficient of determination, denoted as  $r^2$ , is calculated as the square of the correlation coefficient.

In this case, the  $r^2$  value is determined to be 0.726. This result signifies that about 72.6% of the variability observed in one variable can be accounted for by the other variable. Let us assume that we opt to examine the number of individuals per television set as the independent variable, while considering the number of individuals per physician as the dependent variable.

The results obtained via the utilization of the MINITAB “REGRESS” command are as follows: The equation representing the regression model is as follows:

$$\text{People.Phys.} = 1019 + 56.2 \text{ People.Tel.}$$

One can plot the generated regression line over the actual data points to assess how well the model fits the observed data. In this illustration, the plot is shown on the right, with the dependent variable (number of people per doctor) on the y-axis and the explanatory variable (number of people per television set) on the x-axis. There are a few data points that are far from the main cluster of the data, despite the fact that the majority of the data points are concentrated in the lower left quadrant of the plot (showing that there are relatively few people per television and each doctor). These are referred to as outliers and depending on where they are, they might have a significant effect on the regression line.

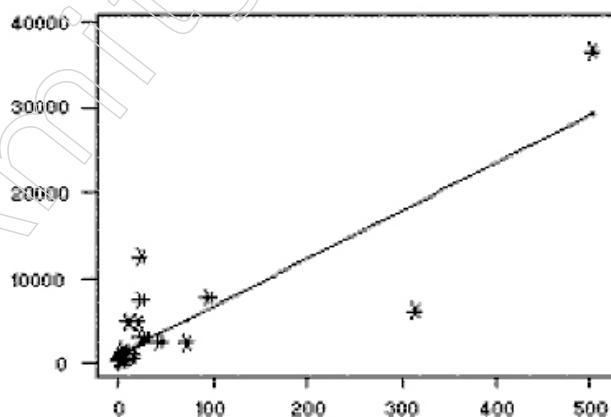


Figure: Regression line

**Outliers and Influential Observations:** An outlier refers to a data point that

## Notes

deviates significantly from the regression line and exhibits a substantial residual value following the construction of a regression line for a given dataset. Possible sources of inaccuracies in the data points could include the presence of erroneous or unreliable information, as well as the existence of a regression line that exhibits a poor fit with the observed data. An observation is deemed important if it is situated at a considerable distance from the other data points along the horizontal axis. This discrepancy arises due to the potential significant impact of these points on the slope of the regression line.

The regression equation is now People after this significant observation has been eliminated. Physical is 1650 plus 21.3 People.Tel. The r<sup>2</sup> value drops to 0.182 because the correlation between the two variables has fallen to 0.427. Less than 20% of the fluctuation in the number of individuals per physician may be explained by the number of people watching television when this significant observation is taken into account. The new model also exhibits influential observations and it is important to examine their effects.

**Residuals:** Once a regression model has been applied to a given dataset, the modeller has the ability to assess the validity of the assumption that a linear relationship exists. This may be achieved by examining the residuals, which represent the discrepancies between the observed values and the values predicted by the fitted line.

The presence of a probable non-linear relationship between the variables can be demonstrated by creating a scatter plot where the residuals are plotted on the y-axis and the explanatory variable is plotted on the x-axis. This visual representation may lead the researcher to investigate the presence of hidden variables.

The discrepancy between an observed response variable value and the value suggested by the regression line is known as a residual.

$$\text{Residual} = \text{Observed Value} - \text{Predicted Value}$$

- The response variable is predicted using regression coefficients based on the explanatory variable.
- Sometimes the anticipated score and the observed score for a certain person are very different

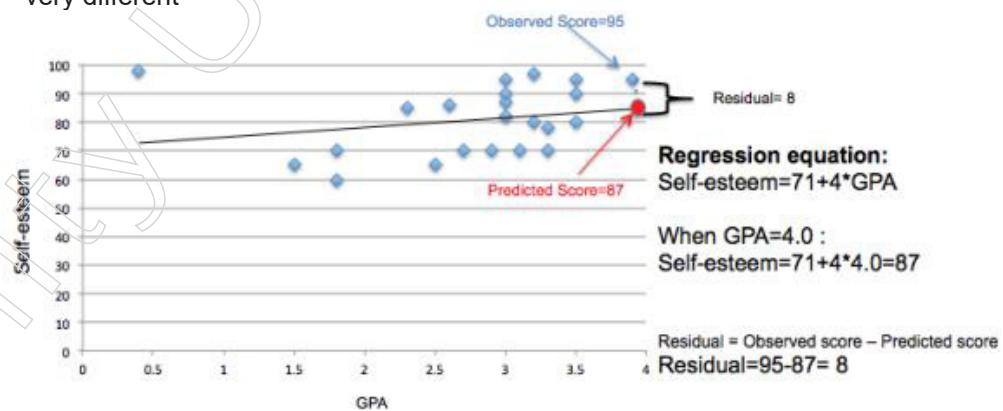


Figure: Residual

- Using the regression equation, a person with a GPA of 4.0 would be anticipated to have a self-esteem score of 87 in the example above. However, this person's observed score is 95. The observed score and the projected score are separated by 8 points. The residual refers to this variation.
- The residuals are minimal when there is a strong correlation between two variables because the data points are closer to the regression line.

**Notes**

- For data points above the line, residuals are positive, while for points below the line, residuals are negative.
- All residuals have a mean that is always equal to zero.
- Outliers and Influential Points:
- An outlier refers to an observation that deviates significantly from the rest of the observations.
- Large residuals are observed in outliers along the y-axis.
- In the context of least-squares regression, it is commonly observed that outliers in the x direction exert a significant influence on the resulting regression line. Consequently, the exclusion of these data points would result in a substantial alteration of the equation representing the line.
- Regression coefficients are derived through the utilization of the means, standard deviations and correlation between the two variables. Consequently, these coefficients exhibit a high degree of sensitivity towards outliers.
- Hence, the positioning of the regression line may undergo substantial alterations in the presence of outliers, or instances of extreme values, within the dataset, particularly when such extreme values are seen for the explanatory variable.
- Upon removal of the outlier, as depicted in the graph on the right, the slope of the regression line exhibits an increase and aligns more closely with the remaining data points.
- Given the proximity of the right line to the data points, it may be inferred that the residuals have a lesser magnitude.

**Lurking Variables:** In the event that a link between an explanatory variable and a dependent variable exhibits non-linear tendencies, it is prudent to take into account additional influential variables. The concept of a hidden variable arises when the association between two variables is substantially influenced by the existence of an unaccounted third variable in the modelling process. Given the potential influence of time-related factors, such as political or economic cycles, it is common practice to employ a time series plot to effectively detect the existence of lurking variables.

**Extrapolation:** It is imperative to carefully observe the range of the data while fitting a linear regression model to a given dataset. Using a regression equation to predict values above the given range is frequently deemed inappropriate and can result in implausible outcomes. This phenomenon is commonly referred to as extrapolation. One illustrative instance is a linear model that establishes a connection between weight gain and age in the context of young children. The application of such a model to individuals in the adult or adolescent age range would be deemed illogical, as the association between age and weight increase lacks uniformity across all age cohorts.

**Important Terms of Linear Regression:**

Linear regression is a statistical technique employed to establish a mathematical model that describes the association between a dependent variable and one or more independent variables. The objective is to determine the optimal linear regression model that accurately describes the data points in a scatter plot. Below are many key words that are commonly connected with linear regression:

- The dependent variable, sometimes referred to as the response variable or target variable, represents the variable that is being sought to be predicted or explained through the use of independent variables.

## Notes

- The independent variables, often referred to as predictor variables or characteristics, are the variables utilized to elucidate or forecast the value of the dependent variable.
- Simple linear regression refers to the linear regression model in which there exists just one independent variable.
- Multiple linear regression refers to the statistical model used when there are two or more independent variables. In this model, the relationship between the dependent variable and the independent variables is represented by a linear equation.
- A linear equation is a mathematical expression that represents a straight line and is commonly employed in the linear regression model to depict the association between the dependent and independent variables. The equation is structured as follows:  $Y$  equals the sum of  $b_0, b_1X_1, b_2X_2$  and so on, up to  $b_nX_n$ . In this equation,  $b_0$  represents the intercept, while  $b_1, b_2, \dots, b_n$  denote the coefficients associated with the independent variables.
- The intercept ( $b_0$ ) refers to the value of the dependent variable ( $Y$ ) when all independent variables are assigned a value of zero. The concept of the intercept refers to the estimated value of the dependent variable in a regression model when all independent variables have a value of zero.
- The coefficients ( $b_1, b_2, \dots, b_n$ ) in the linear regression model denote the parameters that signify the alteration in the dependent variable when there is a one-unit modification in the corresponding independent variable, while keeping all other variables constant.
- Residuals refer to the disparities that exist between the actual values of the dependent variable and the anticipated values generated by the linear regression model. Residuals are employed as a means of evaluating the adequacy of the model's fit.
- The Least Squares Method is a statistical methodology employed to estimate the coefficients in linear regression models. This method aims to minimize the sum of the squared residuals, which are the differences between the observed and predicted values.
- The R-squared ( $R^2$ ) is a statistical metric used to assess the degree of appropriateness of the linear regression model in explaining the variation in the dependent variable. The term "it" refers to the coefficient of determination, which quantifies the fraction of the variation in the dependent variable that can be accurately predicted by the independent variables.
- The adjusted R-squared is a variant of the R-squared metric that accounts for the presence of irrelevant independent variables, hence enhancing its reliability as an indicator of model fit in the context of multiple linear regression.
- Linear regression is predicated on a set of assumptions that must be satisfied in order for the model to be considered valid. These assumptions encompass linearity, independence of errors, constant variance of errors and normality of errors.
- Outliers are data points that exhibit substantial deviation from the overall pattern of the data and has the potential to exert a substantial impact on the regression line.
- Multicollinearity refers to a scenario in which two or more independent variables exhibit a high degree of correlation, resulting in coefficient estimates that are unstable and untrustworthy.

Understanding these terms is crucial for effectively using and interpreting linear regression models in various applications.

### 3.2.2 Linear Regression: Assumptions:

**Assumptions of Linear Regression:** Predictive modelling requires regression analysis, but simply running a line of code or checking R<sup>2</sup> and MSE values is insufficient. The plot() function in R creates four charts that provide insightful information about the data. Unfortunately, a lot of beginners have trouble understanding these plots.

To clarify key regression assumptions, corrections for errors and the importance of these plots. Your regression models will benefit substantially from understanding these principles. For accurate and trustworthy regression findings, it is crucial to identify and rectify assumptions of linear regression that have been violated because doing so can result in inaccurate or inefficient estimates. The following premises underlie linear regression:

- **Linearity:** The dependent and independent variables have a straight line connection.
- **Independence:** There is no interdependence between the observations.
- **Homoscedasticity:** Across all levels of the independent variables, the error variance is constant.
- **Normality:** A normal distribution characterizes the errors.
- **No multicollinearity:** There is little correlation between the independent variables.
- **No endogeneity:** There is no association between the independent variables and the mistakes.

Let's look at the important assumptions in regression analysis: The association between the dependent variable, which is the response variable and the independent variable(s), also known as the predictor variable(s), should exhibit linearity and additivity. A linear connection is characterised by a consistent change in the response variable Y for every unit change in the predictor variable X<sub>1</sub>, independent of the specific value of X<sub>1</sub>. An additive relationship suggests that the influence of X<sub>1</sub> on Y remains constant regardless of the presence or influence of other variables.

- The residual (error) terms should exhibit no correlation with each other. Autocorrelation refers to the absence of this particular phenomena.
- The presence of correlation among the independent variables is considered undesirable. Multicollinearity is a phenomena that is not observed in empirical data.
- The constant variance assumption must hold for the error terms. The aforementioned phenomenon is commonly known as homoskedasticity. Heterogeneous skewness refers to the existence of non-uniform variance.
- The error terms need to be evenly spaced apart.

#### What Happens When You Violate the Assumptions of Linear Regression:

**Linear and Additive:** If a non-linear, non-additive data set is fitted with a linear model during the regression process, the resulting model would be inadequate in mathematically representing the underlying trend. Moreover, this will result in imprecise predictions for a collection of unrecorded information.

**How to check:** Please search for graphical representations comparing fitted values with residual values. In order to incorporate the non-linear impact, it is possible to introduce polynomial terms (X, X<sup>2</sup>, X<sup>3</sup>) into the model.

**Autocorrelation:** The presence of correlation in error terms has a substantial negative impact on the accuracy of the model. Time series models are commonly used when the value of a future instant is dependent on the value of the preceding instant. If

## Notes

the error terms are correlated, the estimated standard errors tend to underestimate the true standard error.

In the event of such an occurrence, it can be expected that the proximity between confidence intervals and prediction intervals will be reduced. A confidence interval with a confidence level smaller than 0.95 would possess a reduced likelihood of encompassing the true values of the coefficients. In order to enhance comprehension of tight prediction intervals, let us employ an illustrative approach.

For instance,  $X^1$  has a least squares coefficient of 15.02 and a standard error of 2.08 (without accounting for autocorrelation). However, when autocorrelation is present, the standard error drops to 1.20. As a result, from (12.94, 17.10), the prediction interval is reduced to (13.82, 16.22). Lower standard errors would also result in p-values that were less accurate than they actually were. This will lead us to draw the wrong conclusion that a parameter is statistically significant.

**How to check:** Please search for the Durbin-Watson (DW) statistic. The value must be within the range of 0 and 4. A Durbin-Watson (DW) statistic value of 2 shows the absence of autocorrelation. A DW value between 0 and 2 suggests the presence of positive autocorrelation, while a DW value between 2 and 4 suggests the presence of negative autocorrelation. Additionally, one can observe a residual vs time plot to identify any seasonal or linked patterns among the residual values.

**Multicollinearity:** This phenomenon occurs when there is a significant correlation between the independent variables, ranging from moderate to strong. The identification of the true association between a predictor and response variable in a model containing correlated variables poses a challenging task. In alternative terms, the task of determining the specific variable that is exerting influence on the forecast of the response variable gets increasingly difficult. An additional factor to consider is that the presence of associated predictors often leads to an increase in standard errors.

Moreover, when there are substantial standard errors, the confidence interval widens, leading to less precise predictions of the slope parameters. Furthermore, in the presence of correlated predictors, the incorporation of additional predictors into the model has an impact on the estimated regression coefficient of a correlated variable. In the event that such a circumstance arises, it is imperative to acknowledge that the conclusion drawn regarding the extent of impact exerted by a variable on the target variable may be rendered erroneous. The estimated regression coefficients of the model would undergo changes if a single connected variable were to be removed from it.

**How to check:** Scatter plots are commonly employed to visualise the correlation between variables. The utilisation of the VIF factor can also be employed. A Variance Inflation Factor (VIF) score equal to or greater than 10 signifies the presence of substantial multicollinearity, whereas a score of 4 or less indicates the absence of multicollinearity. Furthermore, it is important to utilise a correlation table in order to achieve the desired objective.

**Normal Distribution of error terms:** The width or narrowness of confidence intervals can be affected by the presence of non-regularly distributed error components. The presence of confidence interval instability poses a significant challenge in accurately estimating coefficients using the minimization of least squares. The existence of a non-normal distribution suggests the presence of outliers or anomalous data points that necessitate thorough investigation to develop a more precise model.

The parametric technique is regression. Parametric implies it makes presumptions about the data in order to analyze it. Regression has a restrictive aspect because of its

parametric side. When using data sets that don't support its assumptions, it doesn't produce satisfactory results. Therefore, it's crucial to verify these assumptions in order to conduct a successful regression analysis. Continue reading to find out more about how linear regression and polynomial regression work. The utilisation of linear regression analyses aims to assess the extent to which one or more predictor variables effectively explain the dependent (or criterion) variable. The regression is underpinned by five key ideas.

In order to do linear regression, it is important to have a linear association between the independent and dependent variables. Given the susceptibility of linear regression to the influence of outliers, it is imperative to actively identify and examine potential outliers. Scatter plots provide a robust means of assessing the linearity assumption. The subsequent two illustrations demonstrate two scenarios whereby linearity is either lacking or very minimally evident. Furthermore, it is a prerequisite for conducting a linear regression analysis that all variables exhibit a multivariate normal distribution. The most effective approach for validating this assumption is through the utilisation of a histogram or a Q-Q plot. The verification of normalcy can be accomplished through the utilisation of a goodness of fit test, such as the Kolmogorov-Smirnov test. In cases when the data does not exhibit normal distribution, the utilisation of a non-linear transformation, such as a log-transformation, may offer a potential solution to address this issue.

Moreover, it is important to note that linear regression relies on the assumption that the dataset exhibits low or no multicollinearity. Multicollinearity arises when the independent variables exhibit a high degree of correlation with each other.

#### Three main criteria can be used to test for multicollinearity:

- **Correlation matrix** – When calculating the matrix of Pearson's Bivariate Correlation for all independent variables, it is necessary for the correlation coefficients to be less than 1.
- **Tolerance** – Tolerance is a statistical measure that quantifies the impact of one independent variable on all other independent variables. It is derived by an initial linear regression analysis. The concept of tolerance in the context of first step regression analysis is mathematically stated as  $T = 1 - R^2$ . When the value of  $T$  is less than 0.1, it is possible that multicollinearity exists in the data. However, when the value of  $T$  is less than 0.01, it may be concluded that multicollinearity is present in the data.
- **Variance Inflation Factor (VIF)** – The variance inflation factor (VIF) in linear regression is mathematically defined as the reciprocal of the tolerance ( $T$ ), where  $VIF = 1/T$ . When the Variance Inflation Factor (VIF) exceeds a value of 5, it suggests the possible presence of multicollinearity. Similarly, when the VIF exceeds a value of 10, it indicates the presence of multicollinearity among the variables.

In the event of the presence of multicollinearity in the dataset, a potential resolution could involve the process of centering the data. This entails the subtraction of the mean of the variable from each individual score. Nevertheless, the most direct approach to address the problem is to remove independent variables that have high VIF values.

In order to do linear regression analysis, it is imperative that the data exhibits minimal or negligible autocorrelation. Autocorrelation occurs when the residuals exhibit dependence among themselves. For example, the volatility of stock prices can often be attributed to the reliance of present prices on historical prices.

#### 3.2.3 Gradient Descent:

In the process of training deep neural networks using gradient descent, it is necessary to choose a step size denoted as  $\alpha$  for the optimizer. When the value of  $\alpha$

## Notes

is too little, the optimization process experiences a significant decrease in speed. Conversely, if  $\alpha$  is too large, the optimizer fails to reach convergence. The selection of an adequate  $\alpha$  is a job of optimization that machine learning practitioners encounter on a daily basis. Why is the application of gradient descent not considered in this context as well? In order to accomplish this, it is necessary to calculate the derivative of the loss function with respect to both the weights of the neural network and the parameter  $\alpha$ .

In this study, Baydin et al. utilize a concept introduced by Almeida et al. to elucidate the effective computation of “hypergradients” through the manual differentiation of conventional optimizer update methods with regard to the step size hyperparameter. The implementation of online learning rate adaptation facilitates the improvement of convergence, particularly in cases where the starting  $\alpha$  value is sub-optimal.

Nevertheless, the aforementioned approach exhibits three constraints.

- Firstly, the process of manually distinguishing optimizer update rules is laborious and susceptible to errors, necessitating repetition for each optimizer variant.
- Secondly, the method solely optimizes the step size hyperparameter, neglecting other hyperparameters like the momentum coefficient.
- Lastly, the method introduces an additional hyperparameter, namely the hyper-step-size, which also requires tuning.

By swapping out manual differentiation for automatic differentiation (AD), which automatically computes proper derivatives without any additional human work and naturally generalizes to other hyperparameters (such as momentum coefficient) for no additional cost, we are able to overcome all three restrictions. As for the hyperparameters, we see that AD may be used to optimize them as well as the hyper-hyperparameters, hyper-hyper-hyperparameters and so forth. In fact, we can design towers of recursive optimizers of any height and they get more and more resistant to initial hyperparameter selection. As a result, these “hyperoptimizers” lighten the workload of the people in charge of fine-tuning the hyperparameters.

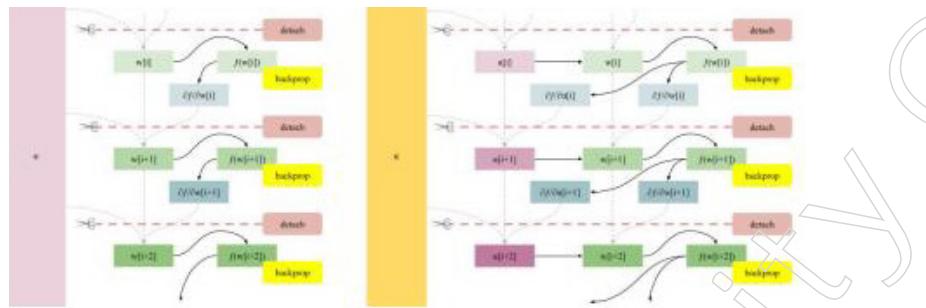
### Computing the step-size update rule automatically:

To facilitate the automatic computation of hypergradients, it is necessary to begin by providing a concise overview of the mechanics underlying reverse-mode automatic differentiation (AD). In the context of differentiable programming systems, it is common for reverse-mode automatic differentiation (AD) to involve the construction of a computation graph during the forward calculation of a function. As an illustration, when a user performs the computation of the function  $f(w_i)$ , the system internally maintains a Directed Acyclic Graph (DAG). In this DAG, the weights  $w_i$  are represented as the leaves, intermediate computations are represented as internal nodes and the ultimate loss is represented as the root.

The process of backpropagation involves the propagation of gradients across a computation graph, commencing with the root node. As the gradients descend, they are deposited in each internal node until they reach the leaf nodes, where the weights  $w_i$  accrue their respective gradients  $\partial f(w_i)/\partial w_i$ . After calculating the gradient of the function  $f$  with respect to the weights  $w_i$  during the backward pass, we proceed to update the weights for the next iteration of gradient descent. This update is performed by setting the new weights  $w_{i+1}$  equal to the current weights  $w_i$  multiplied by a learning rate  $\alpha$  and the computed gradient  $\partial f(w_i)/\partial w_i$ . This process is then repeated for the subsequent step of gradient descent.

A crucial factor to take into account at this juncture is the necessity of “detaching” the

weights from the computation graph prior to the subsequent iteration of this procedure. This entails forcefully converting the weights into leaves of the graph by eliminating any inbound edges. The impact of the “detach” operation is seen in Figure X. If this operation were omitted, the backpropagation process in the subsequent iteration would extend beyond the current weights and into the computation graph of the previous iteration. As the number of steps done increases, the computation graph exhibits linear growth. This is due to the linear relationship between backpropagation and the size of the graph. Consequently, the total training process becomes quadratic-time and intractable.



A: Computation graph of SGD with a single fixed hyperparameter  $\alpha$ .  
B :Computation graph of SGD with a continuouslyupdated hyperparameter  $\alpha_i$ .

**Figure: Visualizing the computation graphs of SGD and HyperSGD.**

The computational overhead of the hyperoptimizer is minimal because we just slightly expand the computation graph in order to evaluate the optimizer.

**Extending to other optimizers:** In contrast to other research, our approach enables the concurrent optimization of all hyperparameters of these optimizers without incurring additional costs. The approach employed in the implementation involves treating them in a manner analogous to alpha. The evaluation conducted in Section illustrates the clear advantages of pursuing this approach. Nevertheless, it is crucial to consider two significant nuances: Initially, it is necessary to ensure that hyperparameters such as 1 and 2 adhere rigidly to the domain  $(0, 1)$ .

To achieve this, we use a scaled sigmoid function to restrict the “raw” values within this domain. In the absence of this particular phase, there is a possibility of inadvertently modifying these values beyond their respective domains. Furthermore, it should be noted that the Adam optimizer incorporates the term  $p v^i$ , which exhibits continuity but lacks differentiability when  $v^i = 0$ . The failure of backpropagation on the first step can be attributed to the division by zero, which occurs as a result of Adam’s usual initialization of  $v^0$  as 0. To address this issue, the problem can be resolved by initializing the variable  $v^0$  to the value of  $\epsilon$  instead of 0.

**Stacking Hyperoptimizers Recursively:** At this point, it is only logical to wonder if the hyperoptimizer itself can be tuned, i.e., if a hyper-hyperoptimizer can modify the hyper-hyperparameters. Baydin et al. proposed that this process might be repeated indefinitely to produce an optimization algorithm that is extremely resistant to the human-selected hyperparameter. Without knowing the precise order of the stacked optimizers beforehand, it is impossible to compute the gradients of these higher-order hyperparameters manually and doing so would also be very time-consuming and error-prone. However, this vision can be realized because AD can calculate these gradients automatically. Let’s review our earlier HyperSGD implementation to achieve this. There is a chance for recursion here, as you may factor out the change to alpha by calling SGD.step, where the hyperparameter for SGD is kappa.

## Notes

This approach is identical to the one above in terms of functionality because SGD is already cautious to properly detach its argument (usually  $w$ , but in this instance  $). Any optimizer that follows this protocol is sufficient, in fact.$

A level-2 hyperoptimizer, HyperSGD(0.01, HyperSGD(0.01, SGD(0.01))) can then be created by feeding HyperSGD itself recursively as the optimizer after this refactoring. Similar to this, we can envision larger buildings or towers that combine various optimizers, like Adamoptimized-by-SGD-optimized-by-Adam. The question of whether this procedure actually makes the problem of hyperparameter optimization worse by adding more hyperparameters is understandable. Baydin et al. expected that the resulting algorithms would become less sensitive to the human-selected hyperparameters as the towers of hyperoptimizers became taller.

### What Is Gradient Descent In Machine Learning:

An optimization method called gradient descent can be used to locate the local minimum of a differentiable function. In machine learning, gradient descent is simply used to identify the parameters (coefficients) of a function that minimize a cost function as much as is practical.

**Gradient:** A gradient calculates how much a function's output will vary if its inputs are slightly altered. — MIT's Lex Fridman. The gradient is a measure that quantifies the rate of change of all weights with respect to the change in error. The concept of a gradient bears resemblance to the notion of a function's slope. A larger gradient is observed when the slope is steeper and the model exhibits faster learning capabilities. Nevertheless, the learning process of the model ceases when the slope reaches zero. In the field of mathematics, a gradient refers to a partial derivative with respect to its inputs.

Consider a blindfolded individual who wishes to ascend a hill with the fewest number of steps possible. He might begin climbing the hill by making really large steps in the direction that is the steepest, which he is allowed to do as long as he is not very close to the summit. But in order to prevent overshooting it, he will take fewer and smaller steps as he approaches the top. The gradient can be used to mathematically characterize this process.

Imagine that the top-down illustration of our slope shows the climber's steps as red arrows. Consider a gradient in this context as a vector that specifies the length and direction of the steepest step the blindfolded man can take.

**How Gradient Descent Works:** Gradient drop can be compared to hiking down to the bottom of a valley rather than up a hill. Because this minimization procedure minimizes a specified function, it is a better analogy. The gradient descent technique operates as shown in the equation below, where  $b$  is our climber's future location and  $a$  is his current position. The gradient descent algorithm's minimization phase is denoted by a minus sign. The gradient term ( $f(a)$ ) is just the direction of the steepest drop and the gamma in the middle is a waiting factor.

$$b=a-\gamma \nabla f(a)$$

The direction of the steepest slope is what this formula essentially instructs us to move to next. To fully put the idea into perspective, let's look at another example.

Consider a machine learning scenario where you want to train your algorithm using gradient descent to reduce your cost-function  $J(w, b)$ , attain its local minimum and do so. The cost function  $J(w, b)$  is depicted on the vertical axes in the picture below, with the horizontal axes showing the parameters ( $w$  and  $b$ ). As a function, gradient descent

is convex. We are aware that we need to identify the values of  $w$  and  $b$  that fall inside the range of the cost function's minimum (shown by the red arrow). We initialize  $w$  and  $b$  with some random numbers to get the search for the ideal values started. Gradient descent then begins there (about at the top of our example) and proceeds step by step in the steepest downward direction (i.e., from the top to the bottom of the illustration) until it reaches the location where the cost function is as m as feasible.

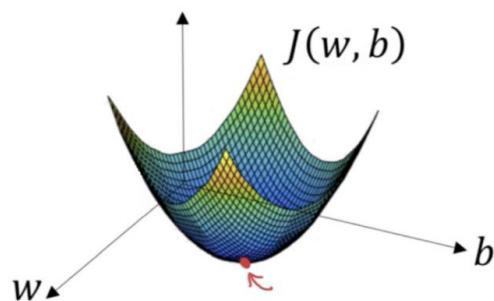


Figure : Gradient descent is a convex function.

**Gradient Descent Learning Rate:** The learning rate, a critical parameter in the optimisation algorithm known as gradient descent, governs the rate at which the algorithm converges towards the optimal weights. It controls the magnitude of the steps taken by the algorithm in the direction of the local minimum.

The learning rate must be carefully chosen to ensure that the gradient descent algorithm converges to the local minimum, without being too little or too large. This observation has significance since the size of the steps taken during the optimisation process may hinder the attainment of the local minimum. This is due to the possibility of the algorithm oscillating between the convex function of the gradient descent, as depicted in the left figure. When the learning rate is adjusted to a sufficiently small value, gradient descent will eventually converge to the local minimum, but potentially requiring a considerable amount of time.

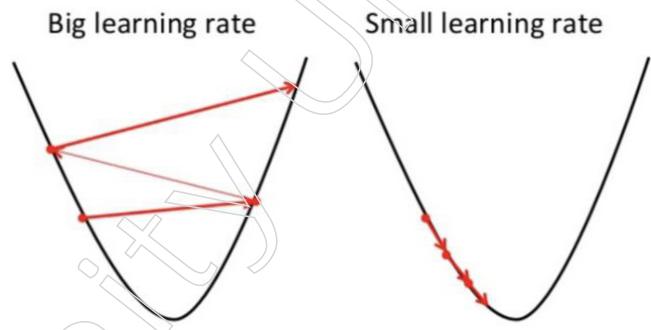


Figure: Gradient Descent Learning Rate

**How to Solve Gradient Descent Challenges:** Plotting the cost function as the optimization process proceeds is an excellent approach to ensure that the gradient descent algorithm functions correctly. Place the cost function's value on the y-axis and the number of iterations on the x-axis. This gives you a tool to quickly determine how appropriate your learning rate is after each gradient descent iteration and enables you to examine the value of your cost function. Simply experiment with various values and plot them collectively. Such a plot is shown in the image on the left below and the difference between good and terrible learning rates is shown in the image on the right.

## Notes

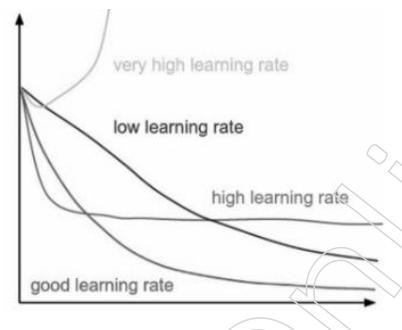
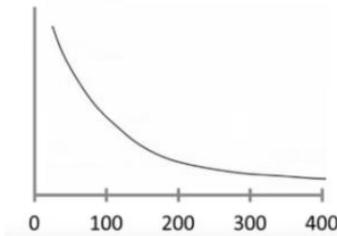


Figure: difference between good and bad learning rates

Every iteration of the gradient descent process should result in a smaller cost function. Gradient descent has converged when it is no longer able to reduce the cost-function and stays essentially at the same level. It can occasionally vary greatly how many iterations of gradient descent are required for convergence. It is difficult to predict in advance the number of iterations necessary to get convergence because it could take 50, 60,000, or even 3 million.

However, you must first specify a threshold for the convergence, which is equally difficult to estimate. Some algorithms can automatically notify you if gradient descent has converged. This makes the chosen convergence test simple graphs.

Monitoring gradient descent using plots has the additional benefit of making it simple to identify when it isn't functioning properly, such as when the cost function is rising. When applying gradient descent, an excessive learning rate is typically the cause of an expanding cost-function.

Reduce the learning rate if the learning curve appears to be bouncing up and down without actually reaching a lower position. Additionally, to see which learning rate performs the best when using gradient descent to begin with a problem, try 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, etc. as the learning rates.

**Types of Gradient Descent:** There are three popular types of gradient descent that mainly differ in the amount of data they use:

**Batch Gradient Descent:** Batch gradient descent, often referred to as vanilla gradient descent, involves evaluating the error for each individual case in the training dataset. However, the model remains unchanged until all training instances have been examined. The complete process is commonly known as a training phase and has cyclical characteristics.

The computing efficiency of batch gradient descent, which yields a steady error gradient and a stable convergence, are some of its benefits. One drawback is that the stable error gradient might occasionally lead to a state of convergence below the maximum that the model is capable of. Additionally, the algorithm must have access to and be able to store the complete training dataset in memory.

**Stochastic Gradient Descent:** In contrast, the stochastic gradient descent (SGD) algorithm updates the parameters individually for each training example inside the dataset, one at a time. The potential for increased speed in stochastic gradient descent (SGD) compared to batch gradient descent is contingent upon the specific problem at hand. One advantage is that the periodic updates provide us with a rather precise understanding of the rate of progress.

Nevertheless, the batch gradient descent method is more computationally efficient compared to frequent updates. The frequency at which these updates occur may also

result in gradients that contain noise, potentially leading to fluctuations in the error rate rather than a progressive decrease.

**Mini-Batch Gradient Descent:** The favoured strategy is mini-batch gradient descent due to its incorporation of both batch gradient descent and stochastic gradient descent. The process involves partitioning the training dataset into smaller, more manageable subsets and subsequently updating each subset individually. This approach achieves a trade-off between the efficacy of batch gradient descent and the durability of stochastic gradient descent.

Mini-batch sizes commonly vary between 50 and 256 in the context of machine learning. However, it is important to note that, similar to other machine learning methodologies, there is no universally prescribed standard for mini-batch sizes as their determination is contingent upon the specific application at hand. The most widely utilised kind of gradient descent in the field of deep learning, this approach is considered the optimal choice for training neural networks.

#### Advantages of Gradient Descent:

- Flexibility is a notable characteristic of Gradient Descent, as it can be applied to a wide range of cost functions and is capable of effectively addressing non-linear regression problems.
- Scalability is a notable characteristic of Gradient Descent, as it possesses the ability to efficiently handle enormous datasets by updating the parameters for each training example individually.
- **Convergence:** The process of Gradient Descent has the potential to converge towards the global minimum of the cost function, given that the learning rate is suitably chosen.

#### Disadvantages Of Gradient Descent:

- The sensitivity to the learning rate is an important consideration in the context of Gradient Descent. It is crucial to select an appropriate learning rate, as a high value can result in the algorithm overshooting the minimum, while a low value can lead to slow convergence of the algorithm.
- One potential drawback of Gradient Descent is its slow convergence, which can be attributed to the fact that it updates the parameters for each training sample individually, thereby necessitating a larger number of iterations to reach the minimum.
- One potential limitation of the Gradient Descent algorithm is its susceptibility to becoming trapped in local minima when the cost function exhibits several local minima.
- The updates in Gradient Descent exhibit noise and possess a significant degree of variance, hence potentially compromising the stability of the optimization process and resulting in oscillatory behavior around the minimum.

### 3.2.4 Linear Regression Using Gradient Descent:

The supervised learning algorithm of linear regression in machine learning uses both input and output variables to teach the mapping function,  $y=f(x)$ , from the input to the output. In order to accurately forecast the output for fresh input data ( $x$ ), it is important to accurately identify the mapping function.

#### Linear Regression with Gradient Descent.

The first method uses the gradient descent algorithm and Python's numpy package to identify the model parameters' optimal values. Use the well-known Scikit-learn

## Notes

Machine Learning library for the second strategy. As a result, we can easily train the model. In the third method, the parameters' optimal values are determined analytically using the normal equation, linear algebra and numpy ndarray objects. After training, there is almost no difference. All of these algorithms produce models that are remarkably similar and make predictions in the same manner.

### Linear Regression with Gradient Descent:

To discover the optimal values for the model parameters, use the gradient descent approach with Python's numpy package. The well-known Scikit-learn Machine Learning library is not utilized here. Therefore, this is a fantastic approach to gain a deeper understanding of the underlying workings of linear regression.

The gradient descent algorithm: The above cost function can be minimized and the optimal values for the linear regression model parameters can be discovered using the optimization process known as the gradient descent algorithm. Although we apply it here with just one predictor, this can also be used for more generic issues like multiple linear regression issues with 100 or 1000 predictors. The gradient descent algorithm is as follows for simple linear regression with two parameters:

```
repeat until convergence {
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ 
    (for  $j = 1$  and  $j = 0$ )
}
```

$:=$  notation: Assignment operator. In programming, it is just  $=$  notation.  $\alpha$ : Learning rate. It is a fixed value.

**Derivative term:** The partial derivative of the cost function  $J(\theta_0, \theta_1)$  for  $j=0$  and 1. After we partially differentiate the cost function  $J(\theta_0, \theta_1)$  for  $j=0$  and 1, the algorithm is:

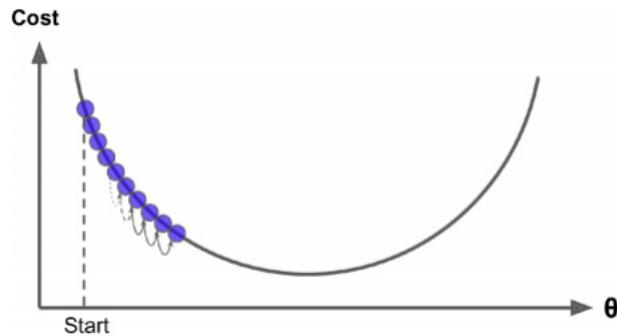
```
repeat until convergence {
     $\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$ 
     $\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$ 
}
```

Now we apply gradient descent to our basic linear regression model by following the steps below.

- **Step 1:** Begin by selecting initial estimates for the parameters  $\theta_0$  and  $\theta_1$ , typically denoted as  $\theta_0 = 0$  and  $\theta_1 = 0$ .
- **Step 2:** Select an appropriate value for the parameter  $\alpha$ .
- **Step 3:** Concurrently adjust the parameters  $\theta_0$  and  $\theta_1$  in order to minimize the cost function  $J(\theta_0, \theta_1)$  until reaching the desired minimum.

**Choosing (Learning rate):** In order to achieve convergence within a suitable timeframe, it is advisable to make appropriate adjustments to the value of  $\alpha$  in the gradient descent process. If the optimization process fails to converge or requires excessive time to reach the minimal value, it suggests that there may be an issue with the chosen  $\alpha$ .

If the value of  $\alpha$  is sufficiently small, the process of gradient descent may exhibit a slower rate of convergence. The attainment of the minimum will necessitate a greater number of iterations, thereby resulting in a longer duration of time.



The learning rate is too small

Figure: Small Learning Rate

In the event where  $\alpha$  is excessively big, the gradient descent algorithm may exhibit the phenomenon of overshooting the minimum. Under such circumstances, it is possible for the process to exhibit a lack of convergence or even demonstrate divergence. In the event that this situation occurs during the execution of the Python code, the returned values for  $\theta_0$  and  $\theta_1$  are NaN. In this scenario, it is advisable to reduce the value of  $\alpha$  and re-execute the algorithm.

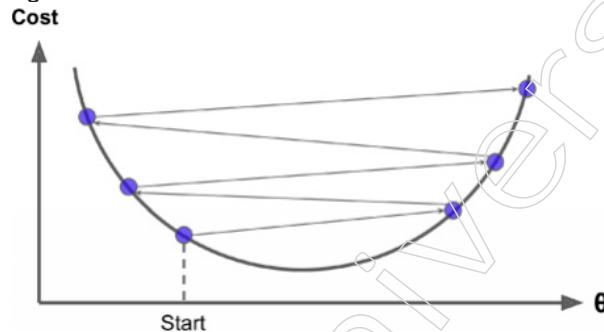


Figure : Large Learning Rate

By selecting an appropriate value for the parameter  $\alpha$  and executing a sufficient number of iterations, the algorithm will ultimately converge to the global minimum. To ascertain this, it is most effective to generate a graphical representation of the cost function values over iterations of the gradient descent algorithm. One should observe a monotonically decreasing function, which is characterized by a consistent decrease or constancy without any instances of increase.

Values of Cost Function over iterations of Gradient Descent

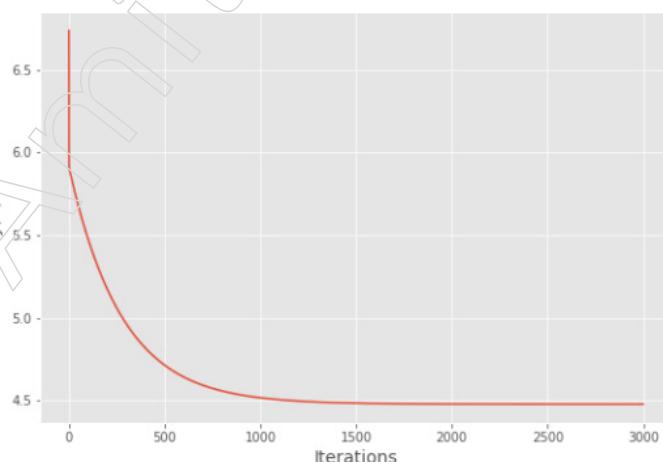


Figure: cost function over iterations of gradient descent

## Notes

In the context of calculus, it can be observed that the derivative of a function, in this case the cost function, will consistently yield a value of zero at the minimum point. Consequently, the cost function's value remains unaltered at this critical point. The minimum cost at the last iteration corresponds to the optimal  $\theta$  values.

**Global minimum and local minimum:** The above graphic illustrates the presence of both local and global minima in various cost functions.

Fortunately, the mean squared error (MSE) cost function for a linear regression model is a convex function that possesses a single global minimum and lacks any local minima. Gradient Descent exhibits the property of convergence towards the global minimum, provided that sufficient time is allowed for the process and the learning rate is appropriately chosen.

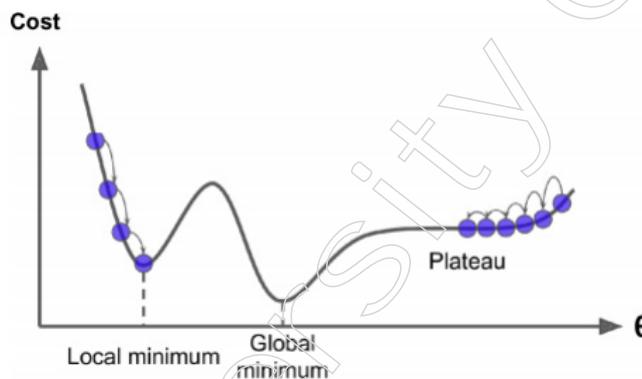


Figure: Global minimum and local minimum

**Feature Scaling:** In the event that the predictors exhibit significant variations in their scales, it is imperative to perform feature scaling before executing the gradient descent algorithm. If the algorithm is executed without feature scaling, the convergence to the global minimum will be significantly delayed, resulting in a prolonged computational time.

- There is no requirement to do feature scaling when there is only a single predictor variable.
- It is advisable to perform feature scaling when the predictors exhibit significant differences in scales, such as a range of 1:100 versus 1:1000.
- Feature scaling should be applied exclusively to the predictors. There is no requirement to implement feature scaling for the response variable, even in cases where it exhibits a significant disparity in scale.
- A frequently employed method for feature scaling is mean normalization, which entails subtracting the mean value of an input variable and dividing the resulting value by the standard deviation of that variable.

$$x_i := \frac{x_i - \mu_i}{s_i} \text{ That's it.}$$

### 3.2.5 Linear Regression: Real Life Events:

#### Linear Regression in Real Life:

Linear regression is frequently employed by businesses to gain insights into the correlation between expenditures on advertising and generated income. As an illustration, one may employ a basic linear regression model to examine the relationship between advertising expenditure, serving as the independent variable and revenue, serving as the dependent variable. The regression model can be represented by the following equation:

$\text{revenue} = \beta_0 + \beta_1(\text{ad spending})$

- The coefficient  $\beta_0$  can be interpreted as the predicted value of total income in the absence of any advertising expenditure.
- The coefficient  $\beta_1$  denotes the mean alteration in total income when advertising expenditure is augmented by a single unit (e.g., one dollar).
- If the coefficient  $\beta_1$  is negative, it would indicate an inverse relationship between ad expenditure and revenue, suggesting that an increase in ad spending is connected with a decrease in income.
- If the value of  $\beta_1$  approaches 0, it would indicate that there is a minimal impact of advertising expenditure on revenue.
- If the coefficient  $\beta_1$  is positive, it indicates a positive relationship between advertising expenditure and revenue.
- The user's text does not contain any information to rewrite. The decision to increase or decrease advertising expenditures by a firm is contingent upon the magnitude of  $\beta_1$ .

#### Linear Regression Real Life Example #2:

Linear regression is frequently employed by medical researchers to get insights into the correlation between the administration of medicine dosages and the corresponding blood pressure levels exhibited by patients.

As an illustration, scientists may deliver different doses of a specific medication to individuals and subsequently monitor the corresponding effects on their blood pressure. A potential approach could involve employing a basic linear regression model, where the dosage is considered as the independent variable and the blood pressure is regarded as the dependent variable. The regression model can be represented by the following equation:

$\text{blood pressure} = \beta_0 + \beta_1(\text{dosage})$

- The coefficient  $\beta_0$  can be interpreted as the anticipated blood pressure level in the absence of any medication.
- The coefficient  $\beta_1$  can be interpreted as the mean alteration in blood pressure that occurs on average when the dosage is incremented by a single unit.
- If the coefficient  $\beta_1$  is negative, it would indicate that there is a negative relationship between dosage and blood pressure, implying that an increase in dosage is correlated with a reduction in blood pressure.
- If the value of  $\beta_1$  approaches zero, it would indicate that there is no significant association between an increase in dosage and changes in blood pressure.
- If the coefficient  $\beta_1$  is positive, it would indicate a positive relationship between dosage and blood pressure, suggesting that an increase in dosage is related with an increase in blood pressure.
- The user's text does not contain any information to rewrite. Researchers may opt to adjust the amount administered to a patient based on the magnitude of  $\beta_1$ .

#### Linear Regression Real Life Example #3:

Linear regression is frequently employed by agricultural experts to assess the impact of fertilizer and water on crop yields.

As an illustration, researchers may employ varying quantities of fertilizer and water

## Notes

across distinct agricultural plots to examine their impact on crop productivity. One potential approach would involve employing a multiple linear regression model, wherein the predictor variables of fertilizer and water are utilized, while the response variable is represented by crop yield. The regression model can be represented by the following equation:

$$\text{crop yield} = \beta_0 + \beta_1(\text{amount of fertilizer}) + \beta_2(\text{amount of water})$$

- The coefficient  $\beta_0$  can be interpreted as the anticipated crop yield in the absence of any fertilizer or water inputs.
- The coefficient  $\beta_1$  can be interpreted as the average effect on crop yield when the amount of fertilizer is increased by one unit, while holding the amount of water constant.
- The coefficient  $\beta_2$  signifies the mean alteration in crop output resulting from a one-unit increase in water, while holding the amount of fertilizer constant.
- The adjustment of fertilizer and water quantities to optimize crop yield is contingent upon the values assigned to  $\beta_1$  and  $\beta_2$  by the scientists.

### Linear Regression Real Life Example #4:

Linear regression is frequently employed by data scientists working with professional sports teams to assess the impact of various training regimens on player performance.

As an illustration, data scientists inside the National Basketball Association (NBA) may undertake an analysis to examine the impact of varying frequencies of weekly yoga sessions and weightlifting sessions on players' scoring performance. One possible approach would involve employing a multiple linear regression model, wherein the predictor variables of yoga sessions and weightlifting sessions are utilized to explain the response variable of total points scored. The regression model can be represented by the following equation:

$$\text{points scored} = \beta_0 + \beta_1(\text{yoga sessions}) + \beta_2(\text{weightlifting sessions})$$

- The coefficient 0 would indicate the predicted number of points for a player who does neither yoga nor weightlifting.
- If the number of weekly weightlifting sessions stays the same while the number of yoga sessions increases by one, the coefficient 1 would represent the average change in points scored.
- The coefficient 2 would represent the average change in points scored when the number of weekly weightlifting sessions increases by one.
- The data scientists may advise a player to take part in more or less weekly yoga and weightlifting sessions in order to maximize the number of points they score, depending on the values of 1 and 2.

The following represents some real-world examples / use cases where linear regression models can be used:

- Linear regression models are commonly utilised by organisations for the purpose of sales forecasting. This might potentially be advantageous for the purposes of strategic planning and financial allocation. Algorithms such as Amazon's item-to-item collaborative filtering employ customers' previous purchases to predict their future buying behaviour.
- Cash forecasting is a common practise in several organisations, wherein linear regression is utilised to anticipate the future cash position. The effective

management of spending and the maintenance of sufficient funds to cover unanticipated expenses are of utmost importance.

- Survey data analysis often involves the utilisation of linear regression as a valuable tool for examining the data. The capacity to assess factors such as consumer satisfaction and product preferences might yield advantages for companies. For example, a commercial entity may employ linear regression analysis to ascertain the probability of consumers recommending their product to others.
- Stock market predictions: Numerous corporations utilise linear regression models as a means to anticipate forthcoming stock performance. This objective is achieved by the identification and analysis of patterns present in historical data pertaining to stock prices and trends.
- Consumer behaviour forecasting is a valuable use of linear regression in which businesses may utilise this statistical technique to predict many aspects, such as the average customer spending. Regression models can also be utilised to forecast consumer behaviour. This approach has the potential to yield benefits in areas such as targeted marketing and product development. Walmart, as an example, employs linear regression to predict the potential high demand for various products nationwide.
- The investigation of relationships between variables can be conducted using linear regression, which enables the identification of links among different variables. One such approach is to utilise linear regression analysis in order to determine the influence of temperature on the sales of ice cream.

### 3.3 Logistic Regression:

Logistic regression, commonly referred to as the logistic model or logit model, is a statistical technique that involves fitting data to a logistic curve. This method is utilised to assess the association between several independent factors and a categorical dependent variable, while also estimating the probability of an event taking place. There are two distinct forms of logistic regression models, namely binary logistic regression and multinomial logistic regression.

Binary logistic regression is commonly used when the dependent variable is dichotomous and the independent variables might be either continuous or categorical. The multinomial logistic regression model is applicable in cases where the dependent variable consists of multiple categories and is not limited to binary outcomes. Consider the utilisation of blood cholesterol levels as a predictive indicator for the progression of coronary heart disease (CHD), serving as an illustrative example. Elevated blood cholesterol levels are associated with an increased risk of coronary heart disease (CHD).

The association between blood cholesterol levels and the risk of coronary heart disease (CHD) exhibits a little decrease at both the lower and higher ends of the spectrum. Moreover, it is important to note that the relationship between CHD and serum cholesterol levels is not characterised by a linear pattern. This pattern is considered typical due to the inherent constraint that probabilities must fall within the range of 0 to 1. The relationship can be described using a 'S'-shaped curve. The logistic function, upon which the logistic regression model is founded, provides estimates throughout the interval of 0 to 1. It also presents an aesthetically pleasing S-shaped representation of the cumulative impact of many risk factors on the likelihood of an event occurring. The aforementioned characteristics contribute to the widespread popularity of the logistic model.

## Notes

### Advantages and Disadvantages of Logistic Regression

#### Advantages:

- The efficacy and ease of implementation and analysis of logistic regression training are noteworthy.
- The statement refrains from making any assumptions on the distribution of classes in feature space.
- The inclusion of multiple classes, namely through the use of multinomial regression, along with the adoption of a natural probabilistic framework for class predictions, can be seamlessly incorporated.
- In addition to providing a measure of the appropriateness of a predictor (as indicated by the size of its coefficient), it also indicates the magnitude and direction (positive or negative) of its relationship.
- The system rapidly categorises unfamiliar data.
- The model exhibits strong performance in cases when the dataset can be linearly split and demonstrates high accuracy across a range of straightforward datasets.
- The utilisation of model coefficients enables the determination of the statistical significance of a given attribute.
- While the probability is lower, logistic regression might still experience overfitting when dealing with high-dimensional datasets. In order to mitigate the issue of over-fitting in these instances, it may be advisable to consider the utilisation of regularisation techniques such as L1 and L2 approaches.

#### Disadvantages:

- The utilisation of logistic regression is not recommended when the number of data points is smaller than the number of features, since this may lead to the problem of overfitting.
- It establishes linear boundaries.
- The primary limitation of logistic regression is in the assumption of a linear relationship between the dependent variable and the independent variables. This method is solely applicable to the prediction of discrete functions. Consequently, the discrete set of numbers is limited to the dependent variable in logistic regression.
- Logistic regression has a linear decision boundary, hence limiting its ability to effectively handle non-linear problems. Linearly separable data is a rarity in real-world scenarios.
- Logistic regression necessitates the presence of average or negligible multicollinearity among independent variables.
- Establishing complex correlations might be challenging when utilising logistic regression. The efficacy of this approach can be easily outperformed by more powerful and compact algorithms such as neural networks.
- The relationship between independent and dependent variables in linear regression is characterised by linear correlation. In the context of Logistic Regression, it is a requirement that the independent variables exhibit a linear connection with the logarithm of the odds ratio ( $\log(p/(1-p))$ ).

**Application of logistic regression:** Logistic regression is a widely employed statistical technique utilized for the purpose of binary classification problems. In such tasks, the objective is to make predictions regarding one of two potential outcomes,

such as Yes/No, True/False, or 0/1, by leveraging one or more predictor variables. It is extensively utilized across diverse disciplines, including machine learning, statistics and the social sciences. Logistic regression is a statistical modelling technique that is commonly used in the field of machine learning. It is a type of regression analysis that is specifically designed for predicting binary outcomes, where the dependent variable is categorical and has only two possible values. The underlying

**Sigmoid Function (Logistic Function):** The logistic regression technique utilizes the logistic or sigmoid function to establish a mathematical representation of the association between the predictor variables and the likelihood of a specific outcome taking place. The sigmoid function is a mathematical function that assigns a value within the range of [0, 1] to any real number. This property makes it well-suited for the representation of probabilities. The sigmoid function is formally defined as:

$$\sigma(z) = 1 / (1 + e^{-z})$$

where  $z$  is the linear combination of predictor variables and their associated weights.

**Hypothesis Representation:** The hypothesis function  $h(z)$  in logistic regression is formulated as the application of the sigmoid function to the linear combination of the input features ( $x$ ) and their respective weights ( $\theta$ ).

$$h(x; \theta) = \sigma(\theta^T * x)$$

where: The function  $h(x; \theta)$  denotes the estimated probability of the positive class (1) based on the input features  $x$  and the model parameters  $\theta$ . The symbol  $\theta^T$  represents the transpose of the parameter vector  $\theta$ . It signifies the dot product between  $\theta^T$  and  $x$ .

**Cost Function:** The cost function, which is alternatively referred to as the log loss or cross-entropy loss, is employed to quantify the disparity between the predicted probabilities and the true class labels inside the training dataset. The objective is to minimize the cost function throughout the training process in order to acquire the most optimal model parameters.

**Training:** In the training phase, the logistic regression model is applied to the training data by employing an optimization procedure, such as gradient descent, in order to determine the most favourable values for the model parameters  $\theta$ .

**Decision Boundary:** After the completion of the training process, the model becomes capable of making predictions on the likelihood of the positive class for novel input data. In order to generate binary predictions, such as 0 or 1, a decision threshold is typically employed on the anticipated probabilities, with a common value being 0.5. If the calculated probability is equal to or greater than the specified threshold, the sample is categorized as belonging to the positive class. Conversely, if the probability is below the threshold, the sample is categorized as belonging to the negative class.

The logistic regression model can be expanded to accommodate multi-class classification tasks by employing methodologies such as one-vs-rest (OvR) or softmax regression. It is imperative to acknowledge that logistic regression is predicated on the assumption of a linear association between the predictor variables and the log-odds of the outcome. Consequently, its performance may be suboptimal when confronted with datasets that exhibit pronounced nonlinearity. In instances of this nature, it may be more appropriate to employ more intricate models such as decision trees, support vector machines (SVM), or neural networks.

**When to use logistic regression:** The application of logistic regression is utilized to make predictions for a categorical dependent variable. In alternative terms, logistic regression is employed when the anticipated outcome is of a categorical nature, such

## Notes

as binary choices between yes or no, true or false, or 0 or 1. The projected probability or output of logistic regression is dichotomous, with no intermediate values. Regarding predictor variables, they have the potential to belong to any of the subsequent categories:

- **Continuous data:** Information that can be scaled infinitely. Any value between two numbers may be used for it. Examples are the number of pounds or the Fahrenheit scale.
- **Discrete, nominal data:** data that falls under specific categories. Here's a simple illustration: blonde, black, or brown hair.
- **Discrete, ordinal data:** Information that conforms to some sort of scale. Giving your level of satisfaction with a good or service on a scale of one to five is an example.

### 3.3.1 Logistic Regression Using Gradient Descent:

To identify the ideal parameters that minimize the cost function, logistic regression is frequently trained using optimization methods like gradient descent. Here is a detailed explanation of how gradient descent is used to train logistic regression.

**Sigmoid functions:** A continuous function with a "S" shape and a domain over all R is known as a sigmoid function or logistic function. The range, though, only extends from 0 to 1. The graph for it is shown in Figure

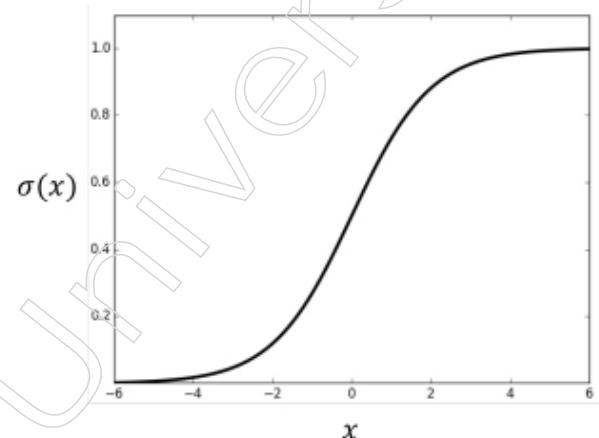


Figure: Sigmoid functions

The logistic regression model, as was already said, employs the sigmoid function to translate the output to a probability value between 0 and 1.

$$\sigma(z) = 1 / (1 + e^{-z})$$

The objective of binary logistic regression is to develop a classifier capable of making a binary determination on the class of a novel input observation. In this section, we present the sigmoid classifier, which will aid in our decision-making process. In the context of our study, we will examine a solitary input observation denoted as  $x$ . This observation will be represented as a vector of features  $[x_1, x_2, \dots, x_n]$ . The specific features will be discussed in detail in the subsequent subsection. The output of the classifier, denoted as  $y$ , can take on two possible values: 1, indicating that the observation belongs to the class, or 0, indicating that the observation does not belong to the class.

We seek to determine the probability  $P(y = 1|x)$  that this particular observation belongs to the class. The decision at hand pertains to distinguishing between "positive sentiment" and "negative sentiment". The characteristics in question are the word counts

within a given document.  $P(y = 1|x)$  denotes the likelihood that the document exhibits positive sentiment, whereas  $P(y = 0|x)$  represents the probability of the document displaying negative emotion. The task at hand is addressed by logistic regression through the process of acquiring knowledge from a training set, wherein a vector of weights and a bias term are learned. Each weight, denoted as  $w_i$ , is a real-valued number that corresponds to one of the input features,  $x_i$ . The weight, denoted as  $w_i$ , signifies the significance of the input feature in determining the classification outcome. It can assume positive values, indicating support for the instance being categorized as belonging to the positive class, or negative values, indicating support for the instance being classified as belonging to the negative class.

Therefore, it is reasonable to anticipate that in a sentiment analysis task, the term "awesome" would possess a significantly positive weight, whereas the term "abysmal" would exhibit a very negative weight. The bias term, commonly referred to as the intercept, is an additional real number that is included in the weighted inputs. In order to arrive at a judgment regarding a test instance, subsequent to the acquisition of weight values during the training phase, the classifier initially performs a multiplication operation between each feature  $x_i$  and its corresponding weight  $w_i$ . The resultant products are then aggregated, followed by the addition of the bias term  $b$ . The resultant numerical value, denoted as  $z$ , represents the aggregated weight of the evidence supporting the classification.

**Designing features:** Features are often generated by carefully analyzing the training dataset, taking into consideration linguistic intuitions and the existing body of linguistic literature pertaining to the specific topic. Conducting a meticulous error analysis on the training set or devset of a preliminary iteration of a system frequently yields valuable insights into its features. In certain activities, the construction of intricate features that encompass a combination of basic features can be particularly advantageous. In the previous section, we observed a characteristic related to period disambiguation. Specifically, it was noted that when the word "St." was preceded by a capitalized word, the likelihood of the period being the conclusion of the sentence decreased.

In the context of logistic regression and naive Bayes, the creation of combination features or feature interactions necessitates manual design. The phenomenon of feature interactions refers to the occurrence of unexpected and often undesirable behaviors that arise when different features or components of a system interact. In numerous activities, particularly those involving feature values that can make reference to specific words, a substantial quantity of features is required. Frequently, these are generated automatically using feature templates, which are abstract definitions of features. As an illustration, a collection of bigram template features designed for period disambiguation could entail the creation of a feature for each consecutive pair of words that precedes a period within the training dataset.

Consequently, the feature space exhibits sparsity due to the need of creating a feature only if the corresponding n-gram is present in the given place within the training set. The feature is often generated by applying a hash function to the string descriptions. The user's description of a feature, specifically "bigram(American breakfast)", is transformed into a unique integer denoted as  $i$ , which subsequently serves as the feature number  $f_i$ .

**Scaling input features:** In scenarios where input characteristics exhibit significant disparities in their value ranges, it is customary to do rescaling in order to establish comparable ranges. The process of standardizing input values involves centering them to achieve a mean of zero and a standard deviation of one.

## Notes

This transformation is commonly referred to as the z-score. In the context provided, the mean ( $\mu_i$ ) of the values of feature  $x_i$  across the  $m$  observations in the input dataset and the standard deviation ( $\sigma_i$ ) of the values of features  $x_i$  across the input dataset, allow for the replacement of each feature  $x_i$  with a new feature  $x'_i$ , which is computed in the following manner.

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)} \quad \sigma_i = \sqrt{\frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2}$$

$$x'_i = \frac{x_i - \mu_i}{\sigma_i}$$

**Processing Many Examples at Once:** The equations for logistic regression for a singular example have been presented. However, in practical applications, it is necessary to process a whole test set consisting of numerous samples. Let us consider a test set including  $m$  test cases, each of which we aim to classify. The nomenclature used on page 2 will be maintained, where a superscript value enclosed in parentheses denotes the index of an example within a given set of data, be it for training or testing purposes. In this scenario, each test example  $x(i)$  is associated with a feature vector  $x(i)$ , where  $i$  ranges from 1 to  $m$ .

**Hypothesis Representation:** The hypothesis function in logistic regression is defined as:

$$h(x; \theta) = \sigma(\theta^T * x) \text{ where,}$$

- $h(x; \theta)$  is the predicted probability of the positive class given input features  $x$  and model parameters  $\theta$ .
- $\theta^T$  is the transpose of the parameter vector  $\theta$ .
- denotes the dot product between  $\theta^T$  and  $x$ .

**Cost Function:** The cost function, alternatively referred to as the log loss or cross-entropy loss, is employed to quantify the disparity between the predicted probability and the true class labels inside the training dataset. The cost function is written as follows for a single training example  $(x, y)$ , where  $y$  represents the actual class label (0 or 1).

$$J(\theta) = -[y * \log(h(x; \theta)) + (1 - y) * \log(1 - h(x; \theta))]$$

The comprehensive cost function for the complete training dataset is computed as the mean of the individual costs over all training cases.

**Gradient Descent:** Finding the ideal parameters for logistic regression will minimize the cost function  $J()$ . The cost function is minimized by using the iterative optimization technique known as gradient descent to update the parameters in the direction of the sharpest descent (negative gradient).

**Training Process:** Gradient descent is used to update the parameters frequently during the training process until convergence. The gradient of the cost function with respect to each parameter is determined in each iteration using the current values of. Then, using the gradient descent update rule, the parameters are changed. The procedure continues until either a predetermined number of iterations is reached or the cost function converges to a minimum (i.e., changes get very small).

**Decision Boundary:** The learned parameters determine the decision boundary after the model has been trained. In the feature space, it symbolizes the line dividing the positive class from the negative class.

**Prediction:** The trained logistic regression model applies the hypothesis function

$h(x; \cdot)$  to the new input features  $x$  to make predictions on new data and the output probability is compared to a threshold (often 0.5).

The fundamental optimization method of gradient descent is employed in many machine learning methods, including logistic regression. It's important to keep in mind that there are gradient descent variants, such stochastic gradient descent (SGD) and mini-batch gradient descent, which can be quicker and more effective for large datasets.

### Summary

- Supervised learning is a type of machine learning where an algorithm learns from labeled training data to make predictions or decisions. In supervised learning, the algorithm is provided with input-output pairs, where the input (also called features) is the data the algorithm uses to make predictions and the output (also called labels or targets) is the expected outcome.
- Linear regression is a fundamental statistical and machine learning technique used for modelling the relationship between a dependent variable (also called the target or response) and one or more independent variables (also called predictors or features). It assumes a linear relationship between the variables and aims to find the best-fitting linear equation that describes this relationship.
- Logistic regression is a statistical and machine learning technique used for binary classification problems, where the goal is to predict one of two possible outcomes (usually coded as 0 and 1) based on input features. Despite its name, logistic regression is used for classification, not regression. It models the probability that a given input belongs to a particular class and provides a decision boundary to separate the classes.

### Glossary

- ML: Machine Learning
- MSE: mean square error
- KNN: k-Nearest Neighbors
- VIF: Variance Inflation Factor
- AD: Automatic Differentiation
- SGD: Stochastic Gradient Descent
- Polynomial Regression: A regression equation is a polynomial regression equation if the power of independent variable is more than 1.
- Clustering: Clustering is a method of grouping the objects into clusters such that objects with most similarities remains into a group and has less or no similarities with the objects of another group.
- Underfitting and Overfitting: If our algorithm works well with the training dataset but not well with test dataset, then such problem is called Overfitting. And if our algorithm does not perform well even with training dataset, then such problem is called underfitting.

### Check Your Understanding

1. What is the main objective of supervised learning?
  - a. To discover hidden patterns in unstructured data.
  - b. To classify data into different categories without labeled examples.
  - c. To make predictions or decisions based on labeled training data.

**Notes**

- d. To perform unsupervised clustering of data points.
2. In supervised learning, what are the two main types of tasks?
  - a. Classification and regression.
  - b. Clustering and dimensionality reduction.
  - c. Feature extraction and model validation.
  - d. Anomaly detection and time series analysis.
3. Which of the following is an example of a classification problem?
  - a. Predicting the price of a house.
  - b. Detecting fraudulent credit card transactions.
  - c. Estimating the temperature on a given day.
  - d. Analyzing the sentiment of customer reviews.
4. In a linear regression model, what is being predicted?
  - a. Categorical labels.
  - b. Probabilities.
  - c. Clusters.
  - d. A continuous numeric value.
5. What is the primary purpose of evaluating a supervised learning model?
  - a. To determine the number of clusters in the data.
  - b. To visualize high-dimensional data.
  - c. To preprocess the data for better accuracy.
  - d. To assess the model's performance and generalization on new data.
6. Linear regression is used primarily for:
  - a. Classifying data into multiple categories.
  - b. Predicting a continuous numeric value.
  - c. Clustering data points into groups.
  - d. Reducing the dimensionality of data.
7. In simple linear regression, how many independent variables are used to predict the dependent variable?
  - a. None, as it only predicts the target variable directly.
  - b. One.
  - c. Two or more.
  - d. It depends on the problem complexity.
8. The goal of linear regression is to:
  - a. Minimize the number of features.
  - b. Maximize the variance of the target variable.
  - c. Minimize the difference between actual and predicted values.
  - d. Maximize the complexity of the model.
9. What is the purpose of the coefficient in linear regression?
  - a. It represents the dependent variable.
  - b. It scales the independent variable.
  - c. It defines the equation of the regression line.
  - d. It determines the number of iterations in training.

10. When evaluating a linear regression model, which metric is commonly used to measure the goodness of fit?
  - a. Mean Squared Error (MSE).
  - b. Accuracy.
  - c. Precision.
  - d. F1-score.
11. Gradient descent is an optimization algorithm used primarily for:
  - a. Clustering data points into groups.
  - b. Solving linear regression problems.
  - c. Finding eigenvalues of a matrix.
  - d. Feature extraction.
12. In gradient descent, what does the “gradient” refer to?
  - a. The direction of steepest ascent.
  - b. The second derivative of the cost function.
  - c. The rate of learning in the algorithm.
  - d. The number of iterations.
13. Which term best describes the learning rate in gradient descent?
  - a. The slope of the cost function.
  - b. The size of the training dataset.
  - c. The rate at which the algorithm converges.
  - d. The step size taken towards the minimum.
14. Which variant of gradient descent takes a fixed step size for each iteration?
  - a. Batch Gradient Descent.
  - b. Mini-Batch Gradient Descent.
  - c. Stochastic Gradient Descent.
  - d. Adaptive Gradient Descent.
15. In gradient descent, what happens if the learning rate is too large?
  - a. The algorithm converges faster.
  - b. The algorithm gets stuck in a local minimum.
  - c. The cost function increases.
  - d. The algorithm stops iterating.
16. Logistic regression is primarily used for:
  - a. Clustering data points into groups.
  - b. Predicting a continuous numeric value.
  - c. Solving classification problems.
  - d. Reducing dimensionality.
17. In logistic regression, the sigmoid function is used to:
  - a. Calculate the mean of the input features.
  - b. Normalize the input features.
  - c. Model the probability of the positive class.
  - d. Reduce the number of features.
18. The output of logistic regression is:
  - a. A continuous numeric value.
  - b. A probability between 0 and 1.
  - c. A categorical label.
  - d. A cluster identifier.
19. What is the purpose of the logistic regression decision boundary?
  - a. It separates the data into clusters.
  - b. It defines the linear equation of the model.
  - c. It indicates the midpoint of the sigmoid function.
  - d. It helps classify data into different classes.

**Notes**

**Notes**

20. In binary classification using logistic regression, the decision boundary is typically set at:
- 0.5 probability.
  - 0 probability.
  - 1 probability.
  - 1 probability.

**Exercise**

- Explain the various concepts of supervised learning.
- What do you mean by linear regression? Also define the important terms used in linear regression.
- Define gradient descent.
- Define logistic regression.

**Learning Activities**

- Provide a step-by-step explanation of how linear regression works, including the formulation of the linear model, estimation of coefficients and prediction process.
- Discuss the importance of evaluating a logistic regression model and describe common evaluation metrics used for binary classification.

**Check Your Understanding- Answers**

- |       |       |       |       |
|-------|-------|-------|-------|
| 1. c  | 2. a  | 3. b  | 4. d  |
| 5. d  | 6. b  | 7. b  | 8. c  |
| 9. c  | 10. a | 11. b | 12. a |
| 13. d | 14. a | 15. c | 16. c |
| 17. c | 18. b | 19. d | 20. a |

**Further Readings and Bibliography**

- Introduction to Linear Regression Analysis by Douglas C. Montgomery, Elizabeth A. Peck and G. Geoffrey Vining
- Linear Regression Analysis by George A. F. Seber and Alan J. Lee
- Supervised Learning Algorithms by S. Sathiya Keerthi and C. Chandrasekaran
- Pattern Recognition and Machine Learning by Christopher M. Bishop

# Module IV: Unsupervised Learning

Notes

## Learning Objectives

At the end of this module, you will be able to:

- Define clustering and clustering algorithms
- Analyse the clustering algorithms in clustering
- Explain K-means clustering and analyse the steps to implement K-means Clustering
- Explain principal component analysis
- Analyse the application for principal component analysis

## Introduction

Perhaps one of the most extensively researched issues in the data mining and machine learning communities is the clustering issue. Over the course of five decades, scholars from several fields have explored this issue. Numerous different problem domains, including text, multimedia, social networks and biological data, are among the problem domains where clustering is applied. Additionally, the issue could arise in a variety of situations, including flowing or unclear data. The basic techniques for clustering vary considerably depending on the problem scenario and data area.

### 4.1 Introduction to Clustering

Because it has so many uses in summarization, learning, segmentation and target marketing, the subject of data clustering has received a lot of attention in the literature on data mining and machine learning. Clustering is a condensed model of the data that can be understood as either a summary or a generating mode in the absence of specific labeled information.

#### 4.1.1 Introduction to Clustering

Clustering is a widely employed methodology in the fields of machine learning and data analysis, aimed at the grouping of data points or objects that exhibit similarities. Its primary objective is to uncover patterns, correlations and structures inherent within a given dataset. The objective of clustering is to divide a given dataset into distinct subsets, commonly known as clusters, in such a way that the data points within each cluster exhibit greater similarity among themselves compared to those belonging to different groups. Clustering is classified as an unsupervised learning methodology, indicating that it does not necessitate the utilisation of labelled data or pre-established categories. Instead, it discerns patterns within the data by leveraging natural similarities.

The basic problem of clustering may be stated as follows: Divide a set of data points into as many groups that are comparable as you can. Please take note that this is a very basic formulation and the problem definition can vary significantly based on the particular model utilized. For instance, a generative model might employ a probabilistic generative mechanism to determine similarity, whereas a distance-based approach would quantify similarity using a conventional distance function. The specific data type also significantly affects how the problem is defined. The following list includes some typical application domains where the clustering issue appears:

Clustering often serves as a critical intermediate step for several core data mining challenges, such as classification or outlier analysis, due to its ability to function as a

## Notes

form of data summarization. A concise summary of the data often proves beneficial in obtaining application-specific insights across different domains.

Collaborative filtering algorithms employ clustering to summarise people who exhibit similar interests. Collaborative filtering is implemented by utilising the ratings that are exchanged among different users. In several contexts, this can be employed to provide suggestions. Customer segmentation refers to the process of grouping customers based on similar characteristics or behaviours. In the context of this application, it involves aggregating customers with equivalent attributes in the data. This approach bears some resemblance to collaborative filtering, a technique commonly used in recommendation systems. One key differentiation from collaborative filtering lies in the utilisation of arbitrary qualities pertaining to the objects, rather than relying solely on rating data, for the goal of grouping.

The field of data summarization is closely intertwined with various clustering approaches, particularly in the context of dimensionality reduction. These strategies can be classified as a form of data summarization. Data summarization enables the development of succinct data representations that facilitate efficient processing and interpretation across diverse applications.

The detection of trends in social networking platforms can be accomplished through the utilisation of dynamic and streaming techniques. In the context of these applications, the data is dynamically streamed and organised into groups, enabling the identification of major patterns of change. Examples of streaming data include multidimensional data, text streams, streaming time-series data and trajectory data. Clustering methodologies can be employed to discern significant patterns and occurrences within the dataset.

Multimedia Data Analysis encompasses a diverse array of document kinds, such as photographs, audio recordings and videos. There are multiple applications for the identification of comparable segments, including the discovery of similar musical samples or visual imagery. The data often exhibits multimodality and encompasses diverse formats. These variables exacerbate the complexity of the issue.

The proliferation of biological data has become increasingly prevalent in recent years due to the success of the human genome project and advancements in the collection of diverse gene expression data. Biological data is commonly structured in the form of networks or sequences. Clustering algorithms provide valuable insights into the prevalent patterns and significant trends present within the data.

The identification of major communities in a social network is a fundamental aspect of social network analysis in various applications. This process involves analysing the structure of the network to determine the significant communities within it. Community detection plays a crucial role in social network analysis due to its ability to enhance comprehension of the network's community structure. The process of summarising social networks, which has proven to be beneficial in various applications, can also be effectively utilised in the context of clustering. The aforementioned list of applications is not comprehensive, nevertheless it offers a comprehensive insight of the diverse array of problems that clustering approaches may address. The research conducted within the domain of data clustering frequently encompasses several prominent areas of investigation. The utilisation of various strategies, such as probabilistic techniques, distance-based techniques, spectral techniques, density-based techniques and dimensionality-reduction based techniques, in the clustering process is not surprising due to the prevalence of clustering as a commonly studied subject. Each of these approaches possesses distinct advantages and disadvantages and can prove efficacious in diverse

scenarios and problem domains. Big data, streaming data and high dimensional data each provide distinct challenges and often require specialised methodologies.

The focus of this approach is on the types of data being used. Different apps offer a range of data kinds accompanied by diverse functionalities. A social network platform will offer a combination of document-based and structural data, whereas an ECG machine will generate a sequence of time data points that exhibit strong interdependence. Categorical data, time series data, discrete sequences, network data and probabilistic data are among the most commonly encountered instances in various academic disciplines. The selection of clustering methods is evidently influenced by the inherent characteristics of the data. Moreover, the differentiation between different types of attributes, such as behavioural and contextual attributes, renders certain data types more complex than others.

Further insights can be gained from exploring variations in clustering. Numerous insights have been generated for different variations in clustering. In order to get further understanding, individuals may utilise methodologies such as visual analysis, guided analysis, ensemble analysis, or multiview analysis. Moreover, the matter of cluster validation holds significant importance in terms of acquiring detailed insights into the efficacy of the clustering process.

#### 4.1.2 Types of Clustering and Clustering Algorithms

**Common Techniques Used in Cluster Analysis:** Following is a discussion of various clustering techniques and algorithms:

##### Feature Selection Methods:

The feature selection phase is a crucial preprocessing step that is required to improve the underlying clustering's quality. Given that some features could be noisier than others, not all features are equally important for locating clusters. Therefore, using a preprocessing phase to remove the noisy and pointless information from contention is frequently beneficial. Dimensionality reduction and feature selection go hand in hand. Original subsets of the features are chosen during feature selection. To further improve the feature selection effect in dimensionality reduction, linear combinations of features may be used in methods like principal component analysis. The former has the advantage of better interpretability, but the later has the advantage of requiring fewer modified directions throughout the representation process.

For more accurate locality-specific insights, feature selection can also be added directly to the clustering method. When several attributes are pertinent to various data locales, this is especially helpful. The failure of global feature selection methods is the driving force behind high dimensional subspace clustering algorithms. As seen in several examples of real data, certain points are linked in relation to a particular set of dimensions, while others are correlated in relation to various dimensions. Therefore, it might not always be possible to reduce the number of dimensions without also losing a significant amount of information. Therefore, the ideal method for attaining this goal is to use local feature selection and incorporate the feature selection process into the algorithm. These local feature selections, also known as local dimensionality reduction, can be applied to the dimensionality reduction issue.

##### Probabilistic and Generative Models:

The main goal of probabilistic models is to simulate data generated by generative processes. The Expectation Maximization (EM) algorithm is used to estimate the

## Notes

parameters of this model after presuming a particular type of generative model (for example, a mixture of Gaussians). The parameters are estimated using the existing data set so that they have the best possible chance of fitting the generative model.

We next calculate the generative probability (or fit probabilities) of the underlying data points using this model as a guide. Anomalies will have very low fit probability, whereas data points that fit the distribution well will have high fit probabilities. A mixture-based generative model's general premise is to suppose that the data were produced using the following procedure, which involved mixing  $k$  distributions with probability distributions  $G_1$  through  $G_k$ :

- In order to select one of the  $k$  distributions, it is necessary to choose a data distribution with a prior probability  $\alpha_i$ , where  $i$  belongs to the set  $\{1\dots k\}$ . Let us consider the scenario where the  $r$ th option is selected.
- Produce a data point derived from the " $G_r$ ".

**The probability distribution**  $G_r$  is chosen from a wide range of potential choices. Be aware that this process is generative and that it necessitates the identification of a number of parameters, including the prior probabilities and the model parameters for each distribution  $G_r$ . Depending on whether the prior probabilities are stated as part of the issue setting or whether interattribute correlations are expected within a component of the mixture, models with varying degrees of flexibility may be created. Note the circular relationship between the model parameters and the likelihood that data points will be assigned to clusters. Therefore, iterative techniques are preferred to get rid of this circularity. The EM technique is usually used to solve the generative models. It begins with a random or heuristic initialization and then iteratively employs two phases to fix the circular computation:

- **(E-Step)** Using the parameters of the current model, ascertain the anticipated likelihood of clustering data points.
- **(M-Step)** Utilizing the assignment probabilities as weights, determine the ideal model parameters for each mixture.

As long as the generative model for each component is properly chosen for the specific mixture component  $G_r$ , one good feature of EM-models is that they can be adapted to new types of data rather easily. Here are a few instances:

- A Gaussian mixture model may be used to model each component  $G_r$  in numerical data.
- A Bernoulli model may be used for  $G_r$  to represent the creation of discrete values for categorical data.
- A Hidden Markov Model (HMM) can be applied to sequence data to model the development of a sequence. It's interesting to note that an HMM is a specific type of mixture model in which the various parts of the mixture depend on one another through transitions. So, it is possible to think of the clustering of sequence data using a combination of HMMs as a two-level mixture model.

Because they attempt to comprehend the underlying mechanism that causes a cluster to be produced, generative models are among the most fundamental clustering techniques. By taking into account specific circumstances in terms of prior probabilities or mixture parameters, a number of intriguing links between other clustering techniques and generative models can be discovered. For instance, the particular scenario in which all mixture components are assumed to have the same radius along all dimensions and where each prior probability is fixed to the same value simplifies to a soft version of the k-means algorithm.

**Distance-Based Algorithms:**

It is possible to demonstrate that many specialized generative algorithms can be reduced to distance-based methods. This is due to the fact that generative models' mixture components frequently employ a distance function within the probability distribution. In terms of the Euclidian distance from the mixture's mean, the Gaussian distribution, for instance, indicates the probability of data creation. As a result, it can be demonstrated that a generative model with the Gaussian distribution and the k-means method have a very tight relationship. In fact, it can be demonstrated that many distance-based algorithms are condensations or simplifications of various generative models.

Distance-based approaches are frequently preferred due to their simplicity and adaptability in a wide range of situations.

**Distance-based algorithms can be generally divided into two types:**

**Flat:** In this scenario, the data is often partitioned into multiple clusters at once using partitioning representatives. The behavior of the underlying algorithm is controlled by the selection of the partitioning representative and distance function. The data points are assigned to the partitioning representatives that are closest to them in each iteration and the representation is then modified in accordance with the data points allotted to the cluster. Comparing this to the iterative structure of the EM algorithm, in which soft assignments are carried out in the E-step and model parameters (corresponding to cluster representatives) are modified in the M-step, is instructive. The following are a few typical techniques for generating the partitions:

- **k-Means:** The partitioning representatives in these methods are the average of each cluster. It should be noted that the partitioning representation is formed as a function of the underlying data and is not pulled from the original data set. To calculate distances, one uses the Euclidian distance. The k-Means approach is thought to be one of the most straightforward and traditional methods for data clustering [46] and is also maybe one of the most widely utilized methods in actual implementations.
- **k-Medians:** Instead than using the mean to produce the dividing representative, these methods use the median along each dimension. The partitioning representatives are not chosen from the original data set, much like with the k-Means technique. The median of a group of values is typically less sensitive to extreme values in the data, making the k-Medians technique more resilient to noise and outliers. Also to be mentioned is the fact that the term "k-Medians" is occasionally misused in the literature because it can also refer to a k-Medoid technique in which the partitioning representations are taken directly from the original data. It should be highlighted that the k-Medians and k-Medoid approaches should be regarded as separate methodologies despite the overload and misunderstanding in the study literature.
- **k-Medoids:** These techniques sample the partitioning representative from the initial data. These methods are especially helpful when the data points to be grouped are arbitrary objects because it is frequently meaningless to discuss the functions of such items. For instance, discussing the mean or median of a group of network or discrete sequence objects may not be pertinent. In these situations, partitioning representatives are selected from the data and their quality is enhanced using iterative procedures. To see if the quality of the clustering improves, one of the representatives is swapped out for a representative from the most recent data in each iteration. As a result, this strategy might be thought of as a form of hill climbing.

## Notes

Compared to the k-Means and k-Medoids approaches, these methods typically require a lot more iterations.

**Hierarchical:** These techniques use a dendrogram to depict the clusters hierarchically and at various levels of granularity. These representations can be categorized as either agglomerative or divisive depending on whether the hierarchical structure was formed top-down or bottom-up.

- **Agglomerative:** These techniques employ a bottom-up methodology whereby we begin with the individual data points and subsequently merge clusters to produce a tree-like structure. There are several options available for combining these clusters, each of which offers a distinct trade-off between quality and efficiency. Single-linkage, all-pairs, centroid and sampled-linkage clustering are a few examples of these options. The smallest distance between any two points in two clusters is utilized in single-linkage clustering. In all-pairs linkage, the average between all pairs is used, but in sampled linkage, the average distance is calculated using a selection of the data points in the two clusters. The distance between the centroids is employed in centroid-linkage. Some of these versions suffer from the drawback of chaining, where larger clusters are predisposed to have closer distances to other points and will thus draw an increasing number of points. This behavior is particularly prone to affect single linkage clustering. Many data domains, including network clustering, are particularly prone to this tendency.
- **Divisive:** These techniques employ a top-down strategy to successively divide the data points into a tree-like structure. The partitioning at each phase can be carried out using any flat clustering algorithm. Greater flexibility in the tree's hierarchical structure and the degree of balance between the various clusters is made possible by dividing partitions. It is not necessary to have a tree with all of the nodes' depths precisely balanced or one with exactly two degrees on each branch. This makes it possible to build a tree structure with a variety of node depth and node weight (number of data points in a node) tradeoffs. For instance, in a top-down approach, the largest cluster can be preferred for division at each level if the node weights of the various branches of the tree are out of balance. In order to construct evenly distributed clusters in huge social networks, where the problem of cluster imbalance is particularly acute, such a method is utilized in METIS. Although METIS is not a distance-based algorithm, same broad guidelines also apply to them.

Since they may be used to practically any data type as long as a suitable distance function is developed for that data type, distance-based approaches are widely employed in the literature. By discovering a distance function for that data type, the clustering issue can then be simplified to a problem of distance function discovery. As a result, the design of distance functions has itself grown to be a significant area of data mining research. Dedicated techniques have frequently been developed for certain data domains, including time series or categorical data. Of course, the quality of the distance functions decreases in many domains, such as high-dimensional data, due to the presence of numerous irrelevant dimensions and may exhibit mistakes as well as concentration effects, which lessen the statistical significance of the findings from data mining. In these situations, one can either use projections to directly locate the clusters in pertinent subsets of attributes or the redundancy in larger portions of the pairwise distance matrix to abstract away the noise in the distance computations.

### Density- and Grid-Based Methods:

Two groups that attempt to granularly explore the data space at high levels are density- and grid-based approaches. The density at any given location in the data space

is expressed either as the number of data points in the locality's predetermined volume or as a smoother kernel density estimate. Typically, a post processing step is used to "put together" the dense parts of the data space into an arbitrary shape after the data space has been investigated at a very high level of granularity. In density-based approaches known as "grid-based methods," the many data space regions that are investigated are organized into a grid-like structure. Grid-like structures are frequently more practical due to the simplicity with which the various dense pieces can be assembled at the post-processing stage. Since lower dimensional grids define clusters on subsets of dimensions, such grid-like methods can likewise be employed in the context of high-dimensional methods.

These methods have the significant benefit of being able to reconstruct the whole form of the data distribution since they explore the data space at a high level of granularity. DBSCAN and STING are two traditional approaches for density-based methods and grid-based methods, respectively. Because density-based approaches are inherently predicated on data points in a continuous domain, this presents a significant issue. As a result, unless an embedding strategy is applied, they frequently cannot be used in a discrete or noneuclidian space meaningfully. Consequently, without specialized modifications, many arbitrary data types, such as time-series data, are not quite as simple to employ with density-based approaches. Due to the larger number of cells in the underlying grid structure and the sparsity of the data in the underlying grid, density computations become more challenging to describe as dimensionality increases.

#### Leveraging Dimensionality Reduction Methods:

In that they try to exploit the similarity and correlations between dimensions to minimize the dimensionality of representation, dimensionality reduction approaches are closely related to feature selection and clustering. As a result, dimensionality reduction techniques are frequently thought of as a vertical kind of clustering, clustering the data's columns rather than its rows using either correlation or proximity analysis. It follows that the question of whether it is possible to complete these procedures simultaneously by grouping the data's rows and columns together naturally emerges. According to the theory, conducting row and column clustering simultaneously is probably more advantageous than doing each of these tasks separately. Numerous techniques, including matrix factorization, spectral clustering, probabilistic latent semantic indexing and co-clustering, have been developed as a result of this broader notion. Although some of these techniques, such spectral clustering, are slightly different, they nonetheless share the same fundamental idea. These techniques are also closely related to projected clustering techniques, which are frequently employed in database literature for high dimensional data.

#### Some common models will be discussed below:

- **Generative Models for Dimensionality Reduction:** These models incorporate the relationships between the data points and dimensions using a generative probability distribution. A generalized Gaussian distribution, for instance, can be thought of as a combination of arbitrarily correlated (directed) clusters, the parameters of which can be discovered using the EM technique. Of fact, given to the greater number of parameters involved in the learning process, it is frequently difficult to perform this robustly as dimensionality increases. It is generally known that techniques like EM are very susceptible to overfitting, a condition in which the number of parameters drastically rises. This is due to the fact that EM approaches attempt to maintain all information in terms of soft probabilities rather than relying on nonparametric

## Notes

methods to make difficult decisions on the selection of points and dimensions. Nevertheless, generative models have been used to successfully learn numerous special situations for various data kinds.

- **Matrix Factorization and Co-Clustering:** Another class of dimensionality reduction techniques that is frequently utilized is matrix factorization and co-clustering. Although it is theoretically feasible to apply these approaches to other types of matrices, they are typically applied to data that is represented as sparse nonnegative matrices. The added interpretability provided by nonnegative matrix factorization techniques, which allow a data point to be described as a nonnegative linear combination of the concepts in the underlying data, is what really draws people to this approach. Methods for nonnegative matrix factorization are closely connected to co-clustering, which simultaneously groups a matrix's rows and columns.
- **Spectral Methods:** Using the similarity (or distance) matrix of the underlying data instead of the original points and dimensions, spectral approaches are an intriguing tool for dimensionality reduction. Of course, this has advantages and disadvantages of its own. The main benefit is that arbitrary objects can now be used for dimensionality reduction instead of only data points that are represented in a multi-dimensional environment. In reality, in addition to achieving the dimensionality reduction, spectral methods also accomplish the twin goal of embedding these objects into a Euclidean space.

In order to conduct clustering on arbitrary objects, such as node sets in a graph, spectral approaches are very well-liked. The drawback of spectral approaches is that the time complexity of even building the similarity matrix rises with the square of the amount of data points because they operate on a  $n \times n$  similarity matrix. Furthermore, unless just a very small number of these eigenvectors are required, the process of figuring out this matrix's eigenvectors can be exceedingly costly. Another drawback of spectral approaches is that, unless the data points were included in the initial sample from which the similarity matrix was derived, it is considerably more challenging to represent data points in lower dimensions. Unless the data is exceedingly noisy and high dimensional, using such a huge similarity matrix for multidimensional data is really pointless.

### The High Dimensional Scenario:

Due to the significant differences in the behavior of the data variables across different regions of the data, the high dimensional scenario presents specific difficulties for cluster analysis. This creates a variety of difficulties in clustering, closest neighbor search and outlier analysis, among other data mining issues. It should be noted that a significant portion of the algorithms for these problems rely on distances as a crucial subroutine. The distances, however, appear to lose statistical significance and effectiveness as dimensionality rises due to irrelevant features.

According to the theory, as dimensions increase, a lower and smaller percentage of the qualities are frequently still relevant, which causes distances to get blurrier and concentration effects to become stronger due to the averaging tendency of the irrelevant attributes. Concentration effects describe how all pairwise distances between data points become comparable when multiple features are noisy and uncorrelated due to their additive effects. For distance-based clustering algorithms (and many other algorithms), the noise and concentration effects present two challenges:

- The distance representation may become inaccurate as a result of rising noise from irrelevant attributes, failing to accurately reflect the actual distance between data objects.

- If employed directly with distances that have not been appropriately denoised, the concentration effects from the irrelevant dimensions cause a drop in the statistical significance of the results from distance-based methods.

### Scalable Techniques for Cluster Analysis:

Data collection has been easier and easier in a wide range of situations thanks to advancements in hardware and software technology. For instance, individuals may carry mobile or wearable sensors in social sensing applications, which could lead to the continuous collection of data over time. This presents a variety of difficulties when real-time analysis and insights are needed. Because the data are frequently too vast to be collected under constrained resource restrictions, this is known as the "streaming scenario," in which it is assumed that just one pass is permitted over the data stream. Even when the data are collected offline, there are still several scalability problems when trying to use the vast amounts of data in a distributed scenario with a big data framework or integrate them with conventional database systems. Thus, depending on the nature of the underlying data, varied degrees of problems are feasible.

Each of these issues is discussed below:

- **I/O Issues in Database Management:** When a data mining technique is used in conjunction with a conventional database system, the most fundamental scalability problems appear. In these situations, it can be demonstrated that the main bottleneck results from the I/O times needed to access the database items. As a result, algorithms that sequentially scan the data are frequently more useful than those that randomly read the data records. In the database literature, a variety of traditional approaches were put up to deal with these scaling problems.

The flat partitioning techniques, which use sequential scans over the database to assign data points to representatives, are the simplest to adapt to this situation. CLARANS, which uses a k-medoids methodology to determine these representatives, was one of the first approaches in this direction. Because each iteration needs testing new partitioning representatives through an exchange procedure (from the current set), the k-medoids method can still be highly computationally demanding. The number of runs across the data set will rise if many iterations are necessary.

- **Streaming Algorithms:** Due to real-time analysis requirements, underlying data evolution and concept drift, the streaming scenario presents unique challenges for clustering algorithms. While database-centric algorithms only need a few passes over the data, streaming algorithms just need one pass because no storage of the data is possible. In addition to this difficulty, real-time analysis is frequently required and it must properly take into account changing patterns in the data.

Almost all streaming approaches require a summarization method to produce intermediate representations that can be used for clustering in order to accomplish these objectives. One of the earliest approaches in this direction creates and maintains the clusters from the underlying data stream using a microclustering methodology. For the microclusters, summary statistics are kept in order to facilitate efficient grouping. A pyramidal time frame is used in conjunction with this to capture the changing characteristics of the underlying data stream. Other data types, including discrete data, massive-domain data, text data and graph data, can also be extended to via stream clustering. The distributed setting also presents a number of particular difficulties.

- **The Big Data Framework:** While big data framework makes use of advancements in storage technology to actually store the data and process it, streaming algorithms operate under the presumption that the data are too massive to be stored explicitly.

## Notes

Even while the data can be explicitly saved, it is frequently difficult to analyse them and get insights from them, as the discussion that follows will demonstrate. This is due to the fact that when data grows in size, a distributed file system must be employed to store it and distributed processing techniques are necessary to maintain adequate scalability.

The problem here is that it is frequently too expensive to shuffle the data between several workstations in order to derive integrated insights from them if big parts of the data are available on different machines. As a result, in order to reduce transmission costs, it is preferable to communicate intermediate insights in distributed infrastructures. This can occasionally present difficulties for application programmers in terms of tracking the locations of various data components and the exact sequencing of connections to reduce expenses.

### 4.1.3 K-means Clustering

In this section, we will discuss the K-Means clustering algorithm, which is considered as the initial partitional clustering technique. One of the most simplistic and efficient clustering methodologies ever introduced in the field of data clustering research. After providing a comprehensive explanation of the approach, we will highlight many critical factors that exert a substantial influence on the ultimate clustering outcome. This section will additionally encompass an exploration of other common iterations of the K-Means algorithm.

#### K-Means Clustering:

K-means clustering is widely recognised as the predominant method for partitional clustering. K representative locations are selected as the initial centroids. Each point is then assigned to the closest centroid depending on the chosen proximity metric. The centroids of each cluster undergo updates subsequent to the formation of the clusters. Subsequently, the algorithm proceeds to repetitively execute these two acts until the centroids exhibit no more changes or a distinct relaxed convergence criterion is met.

The K-means clustering algorithm is a heuristic approach that is guaranteed to converge to a local minimum solution. However, it has been proven to be computationally challenging to minimise its objective function, as it falls under the class of NP-Hard problems. Typically, a less stringent criterion can be employed, thereby relaxing the criteria for convergence. The iterative process must be iterated until the point at which 1% of the points undergo a change in their cluster memberships, in accordance with the governing rule. A comprehensive rationale for the mathematical convergence of the K-means algorithm may be obtained.

#### Algorithm K-Means Clustering:

- 1: Select K points as initial centroids.
- 2: repeat
- 3: Form K clusters by assigning each point to its closest centroid.
- 4: Recompute the centroid of each cluster.
- 5: until convergence criterion is met.

The foundational K-Means approach is delineated in an algorithm. The 3-means algorithm was used to the Fisher Iris dataset at different stages, as depicted in the picture provided. In the initial iteration, three centroids are randomly initialised.

The centroids undergo displacement in subsequent iterations until reaching a state of convergence. The K-means algorithm has the capability to calculate the nearest

centroid by employing a range of proximity measurements. The decision can have a significant impact on the quality of the final solution and the assignment of centroids. In this context, there are numerous sorts of metrics that can be utilised, namely Manhattan distance (L1 norm), Euclidean distance (L2 norm) and cosine similarity. The Euclidean distance metric is commonly employed as the primary measure for K-means clustering.

It is possible to obtain different clusterings by varying the values of K and the proximity measure. The objective function employed by K-means is commonly referred to as the Residual Sum of Squares (RSS) or Sum of Squared Errors (SSE). The subsequent mathematical expression represents the Sum of Squared Errors (SSE) or Residual Sum of Squares (RSS).

Let's define the clustering produced after applying K-means clustering as  $C = \{C_1, C_2, \dots, C_k, \dots, C_K\}$  given a dataset  $D = \{x_1, x_2, \dots, x_N\}$  consists of N points. Equation, where  $c_k$  is the centroid of cluster  $C_k$ , defines the SSE for this clustering. Finding a clustering that reduces the SSE score is the goal. The K-means algorithm's iterative assignment and update stages are designed to reduce the SSE score for the specified set of centroids.

$$\begin{aligned} SSE(C) &= \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - c_k\|^2 \\ c_k &= \frac{\sum_{x_i \in C_k} x_i}{|C_k|} \end{aligned}$$

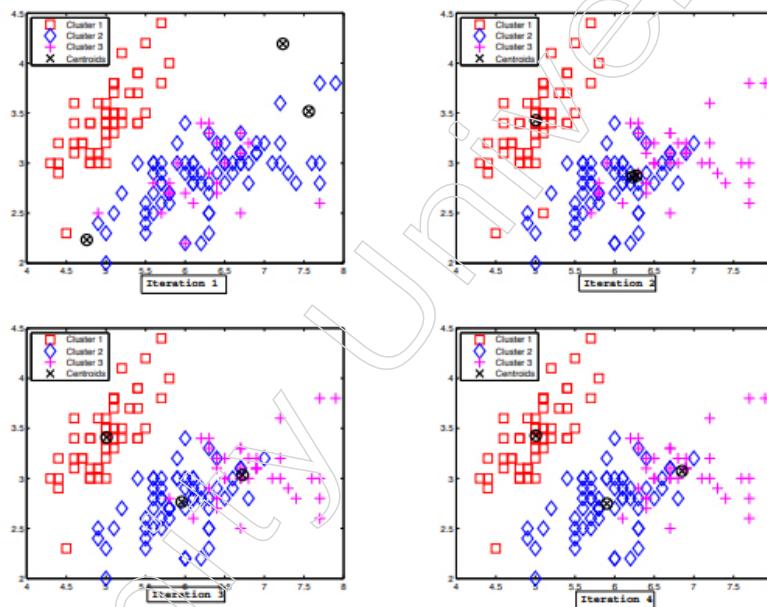


Figure: An illustration of 4 iterations of K-means over the Fisher Iris dataset.

## 4.2 Evaluation of Clustering:

The evaluation of clustering models is challenging since there are no specific labels or well-established ground truths that can be used to examine the clustering findings. Therefore, a variety of measures have been created to evaluate the efficacy of clustering approaches based on their characteristics and objectives. Metrics that are frequently used include:

**Silhouette Score:** Each data point's silhouette score assesses how well it fits into the cluster to which it has been assigned, taking into account its proximity to other data

## Notes

points in that cluster as well as to data points in other clusters. A score of 1 indicates that the data point is correctly classified, while a number of -1 indicates that it is not. The score for a silhouette ranges from -1 to 1.

**Calinski-Harabasz Index:** Greater clustering performance is implied by a higher index value. The ratio of between-cluster variation to within-cluster variance is assessed using the Calinski-Harabasz index.

**Davies-Bouldin index:** Since the Davies-Bouldin index measures the average similarity between each cluster and its closest comparable cluster, a lower value indicates better clustering performance.

**Rand Index:** Better clustering performance is indicated by a higher Rand index. It measures how similar the expected grouping and the actual clustering are.

**Adjusted Mutual Information (AMI):** A higher index indicates better grouping capabilities. The AMI measures the mutual information between the ground truth clustering and the predicted clustering after chance correction.

**Choosing the Right Evaluation Metric:** The best assessment method to use will depend on the clustering problem's goals and nature. The Calinski-Harabasz index or the silhouette score can be useful if the purpose of clustering is to put similar data points together. However, the Rand index or AMI would be a better choice if the clustering results needed to be compared to ground truth clustering. Therefore, when choosing the evaluation metric, it is crucial to take the clustering problem's goals and limitations into account.

**Evaluating the Stability of Clustering Results:** Since the algorithm's parameters and the starting conditions could have an impact on the outcomes, clustering presents certain difficulties. To assess the sustainability of the clustering findings, it is crucial to repeatedly run the clustering approach using various random initializations or settings. Metrics like the Jaccard index or information variance can be used to assess the stability of the clustering results.

**Visualizing the Clustering Results:** Visualization of the clustering results might help to reveal the patterns and structure of the data. One method to view the clustering results is to use scatter plots or heat maps, where each data point is shown as a point or a cell with a color-coded depending on its cluster assignment. Dimensionality reduction methods like principle component analysis (PCA) or t-SNE can be used to map the high-dimensional data into a lower-dimensional space and reveal the clusters. Additionally, cluster analysis software packages usually provide visualization tools like dendograms or silhouette plots to let users explore the results of clustering.

### 4.2.1 Understanding Clustering Types:

The two most extensively researched clustering algorithms are partitional and hierarchical clustering. These algorithms have been extensively utilised in various applications mostly due to their comparative simplicity and ease of implementation when compared to other clustering approaches. Partitional clustering algorithms aim to discover the groupings present in the data by iteratively improving the quality of the partitions through the optimisation of a certain objective function. In order to determine the prototype points that act as the representative elements for each cluster, these methods commonly necessitate the utilisation of certain user-defined configurations. These methods are alternatively referred to as prototype-based clustering algorithms due to this rationale. In contrast, hierarchical clustering methods address the clustering problem by constructing a dendrogram, which is a binary tree-based data structure. By partitioning the dendrogram at different levels, several clustering solutions can be

obtained for the same dataset without the need to repeat the clustering process. Bottom-up clustering and top-down clustering are two distinct approaches used to accomplish hierarchical grouping. Both of these methods employ the dendrogram concept to cluster the data, but the outcomes can vary significantly based on the chosen clustering criterion.

### **Partitional Clustering:**

In this section, the K-Means clustering algorithm will be discussed as the initial partitional clustering technique. One of the most notable and efficient clustering approaches ever proposed in the field of data clustering research. After providing a comprehensive explanation of the approach, we will highlight many crucial factors that exert a substantial influence on the ultimate clustering outcome. K-Means The most used partitional clustering algorithm is clustering. As the initial centroids, K representative points are selected. The closest centroid is then assigned to each point using the selected proximity metric.

Data is divided into non-overlapping clusters by partitioning techniques, with each data point belonging to exactly one cluster. The K-Means algorithm is a well-known illustration of a partitioning algorithm. K-Means seeks to reduce the sum of squared distances between each data point and the centroid of the cluster to which it is assigned. Although K-Means is effective and simple to use, it may have trouble with complex or non-linear cluster shapes.

### **Hierarchical Clustering:**

The drawbacks of flat or partitional-based clustering techniques were addressed by the development of hierarchical clustering algorithms. To get a clustering solution, partitional methods typically need a user-predefined parameter K and they are frequently nondeterministic in nature. To create a more predictable and adaptable technique for clustering the data objects, hierarchical algorithms were created. Agglomerative and divisive clustering methods are two types of hierarchical methods. Agglomerative techniques begin by combining two clusters at a time to create singleton clusters, which have just one data object per cluster. This creates a bottom-up hierarchy of the clusters. On the other hand, dividing methods begin with all the data objects in a large macro-cluster and constantly divide them into two groups, creating a top-down hierarchy of clusters.

### **Agglomerative Clustering:**

An agglomerative hierarchical clustering algorithm involves the following fundamental steps. The dissimilarity matrix is first built using a specific closeness metric and all the data points are then visually depicted at the dendrogram's base. At each level, the clusters that are the closest to one another are combined and the dissimilarity matrix is then updated appropriately. Up until the final maximal cluster—which comprises all of the data objects in a single cluster—is formed, this agglomerative merging process is continued. This would be the peak of our dendrogram and the point at which the merging procedure would be finished.

### **Divisive Clustering:**

Divisive hierarchical clustering is a top-down method where all the data points are started at the root and are split recursively to create the dendrogram. This approach has the advantage of being more effective than agglomerative clustering, particularly when a whole hierarchy down to the level of individual leaves is not required. Since it incorporates all information before dividing the data, it can be regarded as a global approach.

## Notes

### Density-Based Clustering:

By counting the number of points in a fixed-radius neighborhood, DBSCAN calculates the density. Two points are regarded as connected if they are in the same neighborhood as one another. A point is referred to as a core point if it is surrounded by at least MinPts points in the neighborhood of radius Eps, meaning that the density in the area must be higher than a certain level. If a point q is located in the Eps-neighborhood of a core point p and is directly density-reachable from p, then density-reachability is determined by the transitive closure of direct density-reachability. If there is a third location o from which both p and q are density-reachable, then two points p and q are said to be "density-connected." The collection of density-connected points that are maximal in terms of density-reachability are then referred to as a cluster. The set of points in a database that do not belong to any of its clusters is referred to as noise. Finding every cluster in a given database that respects the parameters Eps and Min Pts is the goal of density-based clustering.



Figure:Density-reachability and connectivity

### Grid-Based Clustering:

Grid-based grouping methods are an effective approach for efficiently mining large multidimensional data sets. The algorithms partition the data space into a discrete set of grid-like cells, so establishing a grid structure. Subsequently, they proceed to form clusters by aggregating the grid-like cells. Clusters are defined as regions that exhibit a higher density of data points per unit area compared to the surrounding environment. The utilisation of grids as a method for organising the feature space was initially proposed by Warnekar and Krishna, as exemplified in the GRIDCLUS algorithm. Subsequently, the adoption of this approach became more widespread following the release of STING, CLIQUE and WaveCluster algorithms.

The primary advantage of grid-based clustering is a significant decrease in time complexity, particularly when dealing with large datasets. Grid-based algorithms employ clustering methods that focus on the vicinity surrounding the cells representing the data points, rather than the precise locations of the data points themselves. In the majority of applications, the performance of grid-based approaches is dramatically improved due to the reduced number of cells compared to the amount of data points. Grid-based clustering methods typically consist of five regularly featured phases:

- dividing the data area into a certain number of cells in order to create the grid layout.
- figuring out how many cells there are in each unit.
- arranging the cells in descending order of density.
- locating cluster nodes.
- crossing of nearby cells

The majority of grid-based clustering methods may also be referred to as density-based because it is sometimes necessary to compute cell density in order to sort cells and choose cluster centers. To arrange cells according to their density, several grid-

based clustering methods additionally incorporate hierarchical clustering or subspace clustering.

## Notes

### Data problems include the following for grid-based clustering:

- **Non-Uniformity:** For extremely erratic data distributions, using a single rigid, uniform grid may not be adequate to achieve required clustering quality or efficiency.
- **Locality:** The efficiency of grid-based clustering is constrained by predetermined cell sizes, cell boundaries and the density threshold for significant cells if there are local differences in the shape and density of the distribution of data points.
- **Dimensionality:** Grid-based techniques may not be scalable for clustering very high-dimensional data because performance depends on the size of the grid structures and the size of grid structures may increase dramatically with more dimensions. Additionally, with a grid-based clustering approach, there are components of the “curse of dimensionality” such as noise filtering and choosing the most pertinent qualities that become more challenging with more dimensions.

### Spectral Clustering:

We give a brief overview of the family of spectral clustering techniques, which has grown in popularity recently. Numerous publications in this field of study have been published, beginning with the foundational studies in. Because spectral clustering does not make assumptions about the shapes of clusters, it is able to solve problems in much more complex situations than “traditional clustering algorithms” like k-means and generative mixture models, which always produce clusters with convex geometric shape. The inherent difficulties in the Expectation Maximization (EM) framework, which is frequently used to develop a mixed model for clustering, are another drawback of the earlier approaches. Since this approach essentially involves finding local minima iteratively, getting a suitable solution necessitates several restarts.

With conventional techniques, it can be infamously difficult to cluster several data sets. The toy data sets in the figure below show how challenging they are for conventional clustering techniques. These data sets were developed in order to produce clusters of various shapes. Algorithms that implicitly assume particular cluster geometries cannot produce satisfactory results on such datasets. For instance, the underlying clusters assume a convex form to the Euclidian distance metrics. It goes without saying that such suppositions can affect how well data sets are clustered in general. It will become clear that the spectral clustering method can successfully handle such data sets.

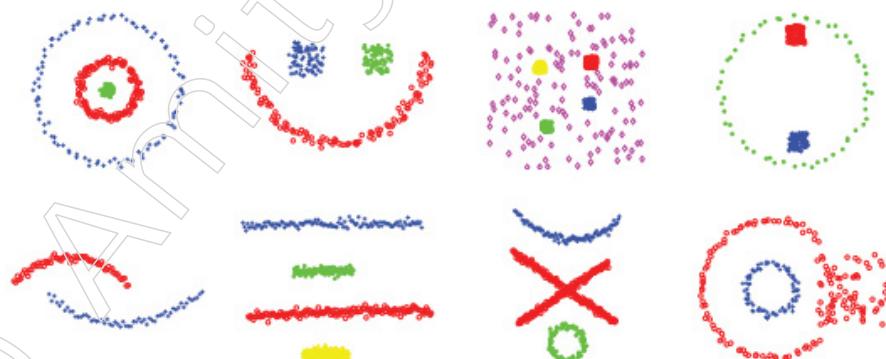


Figure : Self-tuning spectral clustering on toy dataset.

### The spectral clustering family can be viewed as a three-step algorithm:

- ❖ Building a similarity graph for each data point is the initial stage.

## Notes

- ❖ The usage of the eigenvectors of the graph Laplacian embeds the data points in a space where the clusters are more “obvious.”
- ❖ Finally, the embedding is divided using a traditional clustering approach like k-means.

### Fuzzy Clustering:

Data points can be included in several clusters with varying degrees of membership thanks to fuzzy clustering. Fuzzy clustering assigns membership scores, which indicate the likelihood that a data point belongs to each cluster, as opposed to partitioned clustering, which assigns each data point to a single cluster.

### 4.2.2 Anomaly Detection

The identification of atypical events or observations that exhibit statistical divergence from the remaining observations is sometimes referred to as anomaly detection. The observed behaviour, characterised as “anomalous,” typically suggests the presence of an underlying problem, such as credit card theft, server malfunction, or a cyberattack. Anomalies can be classified into three general groups.

- **Point Anomaly:** A tuple in a dataset is referred to as a point anomaly if it deviates significantly from the rest of the data.
- **Contextual Anomaly:** Contextual anomalies are observations that are anomalous because of the context in which they were made.
- **Collective Anomaly:** A collection of data instances aid in the discovery of an anomaly.

The two categories of anomaly detection are global and local detection. To assess whether a point is an outlier, the model compares it to every other point in the data set, which is known as global anomaly detection. Local detection does not compare every point to the full data set, in contrast to global detection. Instead, just the immediate neighborhood is taken into account when determining whether or not a point is anomalous. Since it would also find data points that have a normal value when compared to the entire data set, but an odd value when compared to its neighbors, the latter type typically detects more anomalies than the former. Anomalies are classified using a variety of metrics. Some models employ isolation, some use data density and still others use closeness as a measurement.

Finding abnormalities in logs could aid in finding process problems before they become more serious. Trimma wants to know why a process is taking longer than normal in this situation. A lengthy procedure could mean that anything in their system is acting strangely. It's possible that someone altered the code in a way that had an impact on the procedure, or that some of their servers are operating at a lesser level of performance than others. Another possibility is that their customer did something incorrectly, such sending a data set with a different size or in a different format. Their analysts and engineers may look at prospective anomalies to determine what caused them before they become recurring problems.

**Supervised Anomaly Detection:** In order to build a prediction model to categorize upcoming data points, this method needs a labeled dataset with both normal and anomalous examples. The most popular methods for this purpose include K-Nearest Neighbors Classifier, Support Vector Machine learning and Supervised Neural Networks.

**Unsupervised Anomaly Detection:** This method makes two assumptions about the data: first, that only a small portion of the data is anomalous and second, that any

anomaly is statistically distinct from the normal samples. The data is then clustered using a similarity metric based on the aforementioned presumptions and the data points that are far from the cluster are regarded as anomalies.

#### Anomaly detection techniques of Unsupervised Machine Learning:

Unsupervised machine learning methods don't need any practice data. When the data is unlabeled, unsupervised analysis is advantageous since it involves less effort to prepare the data set. Unsupervised anomaly detection methods make the assumption that normal data points occur far more frequently than anomalous data points. Less common data points are categorized as anomalies using this supposition. Instead of labeling each data point, unsupervised techniques assign each one a score. According to this rating, the data point is likely an abnormality. The model that is employed determines how the score is determined.

The contamination rate  $c$ , or estimated fraction of anomalies in the data set, is one metric necessary for unsupervised anomaly identification. The ground truth is unknown in unsupervised machine learning because the data are not labeled. As a result, different methods must be used to assess the contamination rate. An assessment was made by looking at the data points labeled as anomalies because rarity in the data set was one of an anomaly's features. Domain specialists at Trimma found that the data points labeled as anomalies had much too many false-positives with  $c = 10\%$ , the default value of the models. The lower threshold, which was used to classify anomalies, had a divergence of 0.08% from the category median. The threshold was reduced by about 10% as a result of halving  $c$ , which is more accurate but, in Trimma's opinion, still not a significant enough variation. Once more, halving  $c$  produced  $c = 2.5\%$ , giving a lower threshold of about 25%. Since it was decided that this contamination rate best defined an anomaly in this case, it was adopted as the final contamination rate in all of the models. The contamination rate correlates with the lower threshold, as shown in the figure below. The size of the smallest anomaly discovered increased with reduced contamination rates. Due to Trimma's preference for discovering anomalies that have increased in duration, just the increase threshold was taken into account during this process.

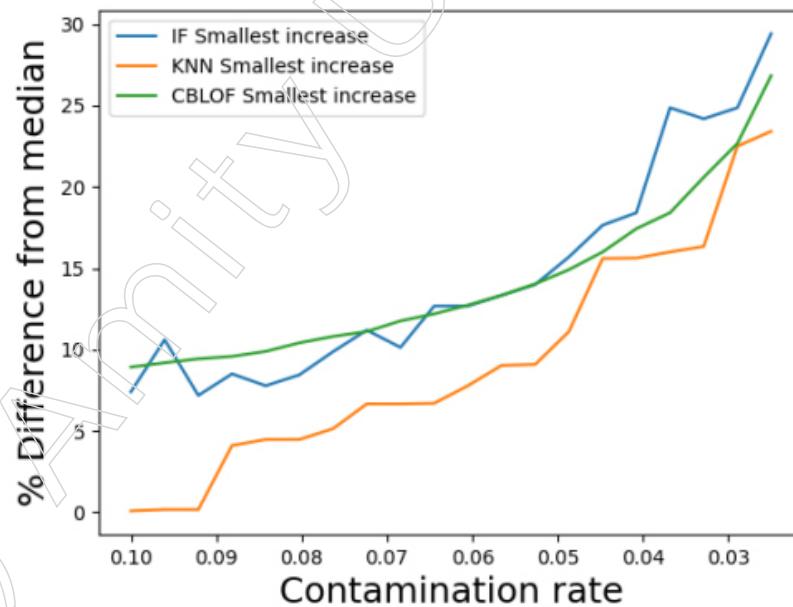


Figure: Lower threshold for each model with various contamination rates.

## Notes

### Isolation Forest (IF) used in Anomaly detection:

**Isolation:** A tree-ensemble technique called Forest is used to find anomalies. By creating a profile of typical data points, other techniques typically identify data points that do not meet this profile as anomalies. Contrary to this, IF actively seeks out abnormal data points. This is accomplished by building an ensemble of decision trees, which are then split recursively on a random feature between its minimum and maximum value. Values that became separated along a shorter journey have a higher likelihood of becoming unusual. The rationale behind why this method works is provided by Liu et al.

Since the pathways for anomalies are noticeably shorter thanks to this random partitioning

- Shorter pathways in a tree structure and
- instances with distinguishing attribute values are more likely to be separated in early partitioning result from fewer instances of anomalies.

The number of trees utilized is one of the variables because it is an ensemble of trees. Each tree employs the aforementioned recursive splitting and random data partitioning. The time it takes to isolate each data point is then determined on average. Data points with lower average split requirements are more likely to be anomalous. The size of the subsample is the other factor. Only a portion of the total data collection is used by IF. This makes IF scale nicely and helps with problems like masking or swamping. Swamping is the problem that arises when a model incorrectly classifies typical data points as abnormalities. When there are groups of anomalies, this is known as masking and it makes them more difficult to find since it would take more splits to isolate any one of them. These problems are less severe because of the reduced sample size. The sample size is set to  $(256, n)$ , where  $n$  is the size of the data set, by default. The number of isolation trees used in the ensemble and the number of features employed both have a significant impact on the execution time.

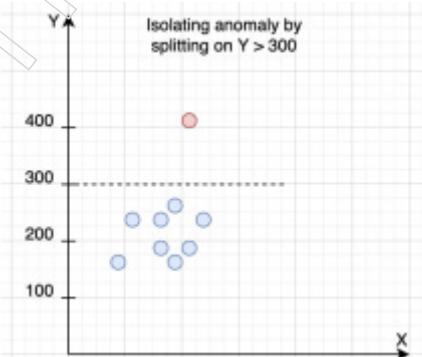


Figure: Example of splitting a data set to isolate a normal data point.

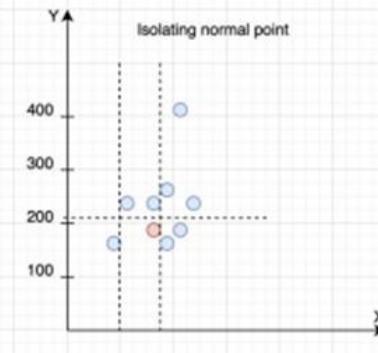


Figure: Example of splitting a data set to isolate a normal data point.

A data collection is divided to isolate several points in the both of figure above. As demonstrated, isolating an anomalous point requires fewer splits than isolating a typical point. Generally speaking, this is the case, however it's possible that a typical point requires fewer splits to be isolated than an abnormal point does. However, when more than one tree is utilized, this problem is lessened because more than one tree is used.

#### **Cluster-based Local Outlier Factor (CBLOF) used in Anomaly Detection:**

**Cluster-based:** A distance-based anomaly detection method called the local outlier factor is closely related to clustering methods. This will allow for the use of a single functionality to simultaneously cluster and detect outliers. According to the original algorithm, data are first clustered using a clustering algorithm. These clusters are subsequently classified as large or tiny clusters. He et al. uses the following concept to distinguish between large and small clusters:

**Definition:** (Cluster, both large and little) Assume that  $C = \{C_1, C_2, \dots, C_k\}$  is the collection of clusters in the progression of  $|C_1| \geq |C_2| \geq \dots \geq |C_k|$ . If one of the following formulas applies, we define  $b$  as the boundary of the large and tiny cluster given two numerical parameters  $\alpha$  and  $\beta$ .

$$(|C_1| + |C_2| + \dots + |C_b|) \geq |D| * \alpha / |C_b| \geq \beta$$

The set of small cluster is therefore defined as  $SC = \{C_j | j > b\}$  and the set of large cluster as  $LC = \{C_i | i \leq b\}$ . In this case, the parameter is utilized to determine the overall size of the huge clusters. The large clusters should collectively encompass 90% of the data points if is set to 90%. Determines the size difference between the Large and Small clusters. A large cluster, for instance, must be at least five times bigger than a small cluster. CBLOF uses a mix of distance and size to determine outlierness.

Assuming  $p$  is a data point, determine its anomaly score as follows: First, determine the distance from  $p$  to the center of the closest major cluster. Second, determine the size of the cluster that  $p$  belongs to so that it can be weighted. The cluster sizes are not used for weighting in the CBLOF implementation in PyOD. The claim is that this has an impact on performance since it prevents the detection of anomalies near small clusters. Instead of the Squeezing algorithm used in the original technique, this implementation employs Kmeans++ clustering.

#### **K-nearest neighbors (kNN) used in Anomaly Detection:**

kNN is a straightforward technique that bases classification on a data point's  $k$  nearest neighbors. The only variables needed are  $k$ , the number of neighbors to take into account and contamination rate. For instance, when classifying a point  $p$  with  $k = 5$ , the point's 5 closest neighbors are taken into account.

The majority class is picked as  $p$ 's class by examining the classes of the  $k$  neighbors. By locating the k-nearest neighbors, kNN can be utilized in anomaly detection to find widespread anomalies. There are several approaches to determine how anomalous a data point is. The distance to its  $k$ th nearest neighbor can be used as one method. As a result, data points that are scarce and distant from other data points are farther away from their  $k$ th nearest neighbor and are therefore more likely to be anomalous. Utilizing the average distances to its  $k$ -nearest neighbors is an additional choice. The advantage of this method is that it also considers the point's local density.

### **4.3 K-means Clustering:**

The clustering difficulties are addressed using a direct unsupervised learning strategy. The topic of interest is K-means clustering. The utilised methodology employs

## Notes

a systematic procedure to partition a provided dataset into many clusters, each of which is represented by the preassigned symbol “k.” The clusters are subsequently represented as individual points and all observations or data points are linked to the nearest cluster, calculated and refined. The aforementioned process is subsequently replicated utilising the updated modifications and so forth, until the intended result is achieved. K-means clustering is employed in various domains such as search engine optimisation, market segmentation, statistics and astronomy.

### Factors Affecting K-Means:

The following are the main variables that can affect how well the K-means algorithm performs:

- Choosing the initial centroids.
- Estimating the number of clusters K.

Now go over a few strategies that have been put out in the literature to address each of these concerns.

**Popular Initialization Methods:** A straightforward initialization technique that MacQueen developed selects K seeds at random. The literature frequently employs this approach because it is the easiest. The following list includes some well-liked K-means initialization techniques that have been effective in enhancing clustering performance.

**Hartigan and Wong:** This strategy argues that the locations that are well spaced and have a high number of points within their surrounding multi-dimensional sphere can be ideal candidates for beginning points by using the notion of nearest neighbor density. Equation is used to get the average pair-wise Euclidean distance between locations. In order to keep  $d_1$  apart from all earlier seeds, subsequent points are picked in the sequence of decreasing density. Please take note that we continue to use the earlier introduced notation for the formulae supplied below.

$$d_1 = \frac{1}{N(N-1)} \sum_{i=1}^{N-1} \sum_{j=i+1}^N \|x_i - x_j\|$$

**Milligan:** This method utilizes the results of the Ward's method using the agglomerative hierarchical clustering (with the aid of the dendrogram). By calculating the sum of squared errors to determine how far apart two clusters are from one another, Ward's approach selects the initial centroids. The greedy approach of Ward's method limits the amount of agglomerative growth that can occur.

**Bradley and Fayyad:** Pick random subsamples from the data and use random seeds to apply K-means clustering to all of these subsamples. After gathering the centroids from each of these subsamples, a fresh dataset made up of only these centroids is produced. These centroids are used as the beginning seeds in the clustering of this new dataset. The best seed set is guaranteed by the minimum SSE attained.

**K-Means++:** The initial centroids for K-means clustering are carefully chosen by the K-means++ algorithm. The first centroid is originally chosen at random as part of the algorithm's straightforward probability-based approach. The centroid that is farthest away from the currently chosen centroid is chosen as the next centroid. Based on a weighted likelihood score, this choice is made. Once we obtain K centroids, we do K-means clustering on these centroids to continue the selection process.

**Estimating the Number of Clusters:** One of the main difficulties with K-means clustering is estimating the right number of clusters (K). In the literature, several

academics have suggested fresh approaches to this problem. We'll briefly go through a few of the more popular techniques.

**Calinski–Harabasz Index:** Equation, where N is the number of data points, defines the Calinski–Harabasz index. By maximizing the function specified in Equation, the cluster count is determined. With K clusters, B(K) and W(K) represent the between and within cluster sums of squares, respectively.

**Gap Statistic:** This technique generates B distinct datasets, each with the same range of values as the original data. For each of them, a different number of clusters is used to calculate the within cluster sum of squares. The both-uniform dataset's within-cluster sum of squares is denoted as W \* b (K).

$$\text{Gap}(K) = \frac{1}{B} \times \sum_b \log(W_b^*(K)) - \log(W(K))$$

The amount of clusters selected corresponds to the minimal value of K that fulfills Equation.

$$\text{Gap}(K) \geq \text{Gap}(K+1) - s_{k+1}$$

where Sk+1 is an estimate of the standard deviation of the formula  $\log(W * b (K + 1))$ .

**Duda and Hart:** By selecting the appropriate cut-off level in the dendrogram, this estimation technique halts the hierarchical clustering procedure. The following techniques are frequently used to cut a dendrogram:

- cutting the dendrogram at a user-specified threshold of similarity, where the level of similarity is predetermined.
- cutting the dendrogram where the biggest gap exists between two successive merges.
- stopping the procedure when the density of the merged cluster falls below a predetermined level.

**Newman and Girvan:** The dendrogram is treated as a graph in this approach and a betweenness score—which will be used as a measure of dissimilarity between the edges—is suggested. The process begins by determining the betweenness score of each graph edge. The edge with the highest betweenness rating is then eliminated. The remaining edges' betweenness scores are then recalculated until the complete set of connected components is acquired. This technique yields a set whose cardinality can be used as a good approximation for K.

**Silhouette Coefficient:** This formulation takes both the intra- and inter-cluster distances into account. The average distance to every point in the same cluster is first determined for a specific point xi. It has been set to ai. The average distance of xi to all the data points in each cluster is then calculated for each cluster that does not contain xi. It has been set to bi. The silhouette coefficient of a point is calculated using these two numbers. The average silhouettes width for all the points in the dataset is the average of all the silhouettes in the dataset. One can calculate the average silhouette coefficient across all points to assess the quality of a clustering.

$$S = \frac{\sum_{l=1}^N \frac{b_l - a_l}{\max(a_l, b_l)}}{N}$$

**Newman and Girvan:** The dendrogram is treated as a graph in this approach and a betweenness score—which will be used as a measure of dissimilarity between the

## Notes

edges—is suggested. The process begins by determining the betweenness score of each graph edge. The edge with the highest betweenness rating is then eliminated. The remaining edges' betweenness scores are then recalculated until the complete set of connected components is acquired. This technique yields a set whose cardinality can be used as a good approximation for K.

**ISODATA:** For clustering the data using the nearest-centroid method, ISODATA was suggested. In this approach, the dataset is initially processed through K-means to identify clusters. Then, if a cluster has less points than a certain threshold or a distance below a specific threshold, it is combined. Similar to this, if a cluster's internal standard deviation is higher than a user-defined threshold, the cluster is divided.

### 4.3.1 K-means Clustering:

The K-means algorithm's straightforward structure makes it easy to alter and add additional effective algorithms on top of it. Some of the K-means algorithm modifications that have been suggested are based on

- Applying some sort of feature transformation technique (Weighted K-means, Kernel K-means),
- picking improved initial centroid estimates (Intelligent K-means, Genetic K-means) and
- selecting various representative prototypes for the clusters (K-medoids, K-medians, K-modes).

The most well-known K-means clustering variations that have been put out in the partitional clustering literature will be covered in this section.

#### Variations of K-Means:

**K-Medoids Clustering:** Compared to Kmeans, the clustering algorithm K-medoids is more robust to outliers. Finding a clustering solution that minimizes a predefined objective function is the aim of K-medoids, just as K-means. The K-medoids approach, which is more resistant to data noise and outliers, selects the actual data points as the prototypes. Instead of focusing on SSE, the K-medoids algorithm seeks to reduce the absolute error criteria. The K-medoids technique works in an iterative fashion, much like the K-means clustering algorithm, until each representative object is truly a medoid of the cluster. Algorithm provides the fundamental K-medoids clustering algorithm.

1. Choose K points as the starting representative objects.
2. repeat
3. The allocation of each point to the cluster is determined by its proximity to the representative item.
4. Select an object  $x_i$  in a random manner, without considering its representativeness.
5. Calculate the overall cost S of replacing the model item  $m$  with  $x_i$ .
6. Swap  $m$  with  $x_i$  to create the new set of K representative objects if  $S < 0$ .
7. until Convergence criterion is met.

In some circumstances, the K-medoids clustering method considers replacing a representative point  $m$  with an arbitrary random point  $x_i$ . The membership of the points that were initially a part of  $m$  is checked after this phase. One of the two possibilities exists for the change in membership of these points. These locations can now be nearer to any of the previous sets of representative points or to  $x_i$  (the new representative point).

The absolute error criteria for K-medoids is used to determine the cost of swapping. This cost of switching is determined for each assignment operation and is a component of the overall cost function.

A variation of the K-medoids clustering technique called Partitioning Around Medoids (PAM) algorithm is suggested [26] to address the issue of carrying out many swap operations while getting the final representative points for each cluster. On the dissimilarity matrix of a given dataset, this algorithm functions. PAM swaps all the nonmedoid points and medoids repeatedly until convergence in order to minimize the objective function. K-medoids are more reliable than K-means, but because of their higher processing cost, they are not appropriate for huge datasets. The Clustering LARge Application (CLARA) algorithm was also proposed using PAM in conjunction with a sampling technique. To arrive at the collection of ideal medoids, CLARA takes into account numerous samples and applies PAM to each one.

**K-Medians Clustering:** In K-means clustering, the mean of the cluster is found. In K-medians clustering, the median of each cluster is found. The K-medians clustering method chooses K cluster centres by making sure that the sum of the distances between each point and the closest cluster centre is as small as possible. In the K-means algorithm, the square of the L2-norm is used as a distance measure. In the K-medians algorithm, the L1-norm is used as a distance measure. The criterion function of the K-medians method is as follows:

$$S = \sum_{k=1}^K \sum_{x_i \in C_k} |x_{ij} - med_{kj}|$$

where  $med_{kj}$  is the median for the  $j$ th attribute in the  $k$ th cluster  $C_k$  and  $x_{ij}$  is the  $j$ th attribute of the instance  $x_i$ . K-medians are more resilient to outliers than K-means. Finding the median point subsets that reduce the cost of assigning data points to the closest medians is the aim of the K-Medians clustering. The algorithm's general structure is comparable to that of K-means. The following two steps are repeated until convergence:

- The medians are recalculated using the median of each unique characteristic after each datapoint is allocated to its nearest median.

**K-Modes Clustering:** The inability of K-means to deal with nonnumerical properties is one of its biggest drawbacks. Categorical data can be changed into new feature spaces using certain data transformation techniques and the K-means algorithm can then be used to extract the final clusters from this transformed space. However, this approach has shown to be incredibly unsuccessful and does not result in high-quality clusters. It has been noted that when working with categorical data, the SSE function and mean usage are inappropriate. As a result, the K-modes clustering algorithm has been suggested to address this issue.

When processing categorical data, the nonparametric clustering algorithm K-modes optimizes a matching metric (L0 loss function) without utilizing any explicit distance metrics. The loss function in this instance is a particular instance of the common  $L_p$  norm where  $p$  goes to zero. The loss function in K-modes clustering operates as a metric and uses the number of mismatches to estimate the similarity between the data points, in contrast to the  $L_p$  norm, which determines the distance between the data point and centroid vectors. Algorithm includes a detailed description of the K-modes algorithm. This is an optimization problem, just like K-means and this approach cannot ensure an overall optimal answer.

## Notes

### K-Modes Clustering algorithm:

1. Choose K initial modes.
2. Repeat
3. To form K clusters, the data points are assigned to the cluster with the closest mode using the matching metric.
4. Recalculate the modal values of the clusters.
5. Until the convergence requirement is satisfied.

**Fuzzy K-Means Clustering:** Fuzzy C-Means clustering is another name for this. In complex datasets with overlapping clusters, it is not possible to perform hard allocation of points to clusters. It is possible to utilize a fuzzy clustering technique to extract these overlapping structures. The membership of points to various clusters in the fuzzy C-means clustering method (FCM) can range from 0 to 1. In Equation, the SSE function for FCM is given.

$$\begin{aligned} SSE(C) &= \sum_{k=1}^K \sum_{x_i \in C_k} w_{xik}^\beta \|x_i - c_k\|^2 \\ w_{xik} &= \frac{1}{\sum_{j=1}^K \left( \frac{|x_i - c_j|}{|x_i - c_k|} \right)^{\frac{2}{\beta-1}}} \\ c_k &= \frac{\sum_{x_i \in C_k} w_{xik}^\beta x_i}{\sum_{x_i \in C_k} w_{xik}} \end{aligned}$$

The membership weight of point  $x_i$  that belongs to  $C_k$  in this case is  $w_{xik}$ . The fuzzy C-means update phase use this weight. For  $C_k$  (represented by  $c_k$ ), the weighted centroid based on the fuzzy weights is determined. The fundamental method operates similarly to K-means in that it iteratively minimizes the SSE before updating  $w_{xik}$  and  $c_k$ . Up to the centroids' convergence, this process is continued. The FCM method, like K-means, is sensitive to outliers and the results it produces will match the local minimum of the objective function. Rough C-means and Possibilistic C-means are two more extensions of this method that have been documented in the literature.

**X-Means Clustering:** X-means is a clustering technique that can effectively estimate K's value. To find the sets of centroids among the existing ones that can be separated to better suit the data, it employs a technique called blacklisting. The Akaike or Bayesian Information Criterion is used to make decisions in this situation. By first narrowing the search space with a heuristic in this algorithm, the centroids are selected. The experiment's K values are chosen between a predefined lower bound and higher bound.

The effectiveness of the model is then evaluated for various  $K$  in the confined space using a certain model selection criterion. The maximum likelihood estimates and the Gaussian probabilistic model are used to build this model selection criterion. The model with the highest model score is the one for which the best K value applies. This algorithm's main objective is to provide a scalable K means clustering algorithm when the number of data points increases while effectively estimating K.

**Intelligent K-Means Clustering:** The idea behind the Intelligent K-means (IK-means) clustering technique is that a point is more interesting the further it is from the centroid. IK-means chooses the points that are farthest from the centroid and correlate to the greatest amount of data dispersion using the fundamental principles of principal

component analysis (PCA). Anomaly pattern clusters are the clusters created from such locations. Algorithm provides the IK-means clustering algorithm.

1. Determine the centroid for the provided dataset, denoted as  $cg$ , by the process of calculating the centre of gravity.
2. Repeat
3. Construct a centroid, denoted as  $c$ , that is located at the maximum distance from the centre of gravity ( $cg$ ).
4. Construct a cluster, denoted as  $Siter$ , consisting of data points that exhibit a closer proximity to the centroid  $c$  compared to the centroid  $cg$ . This is achieved by allocating all the remaining data points, denoted as  $xi$ , to  $Siter$  if the distance between  $xi$  and  $c$ , represented as  $d(xi,c)$ , is less than the distance between  $xi$  and  $cg$ , represented as  $d(xi,cg)$ .
5. The centroid of the region known as  $Siter$  has been updated to  $sg$ .
6. Let  $cg$  be equal to  $sg$ .
7. Eliminate tiny clusters, if present, by using a predetermined threshold.
8. Until the stopping requirement is satisfied.

We can choose K in IK-means in a variety of ways, some of which are comparable to how we choose K in the previously discussed K-means. It is possible to use a structural-based method that contrasts within-cluster coherence with between-cluster dissociation. K can also be calculated using common hierarchical clustering techniques that create a dendrogram. IK-means can be thought of as a deterministic algorithm, in contrast to K-means, which is thought to be nondeterministic. When clusters are dispersed across the dataset rather than compactly formed in a single area, IK-means can be particularly effective in extracting them. Before using K-means, the initial centroid seed selection can be done using IK-means clustering. Only the top centroids will remain after the IK-means for further selection. Small aberrant pattern clusters have already been trimmed, thus they won't offer any candidate centroids.

**Bisecting K-Means Clustering:** A divisive hierarchical clustering technique called bisecting K-means clustering utilizes Kmeans repeatedly on the parent cluster C to find the optimal split to produce the two child clusters C1 and C2. When K-means are bisected to get the best split, uniform-sized clusters are produced. Algorithm contains the algorithm for bisecting K-means clustering.

1. Repeat.
2. Select the parent cluster to be divided; then, repeat.
3. Pick at random two centroids from C.
4. Apply a predetermined distance measure to the remaining points and place them in the closest subcluster.
5. Recalculate the centroids, then keep assigning clusters until convergence.
6. Use the centroids to calculate inter-cluster dissimilarity for the two subclusters until I iterations are finished.
7. Select the subcluster centroids with the highest inter-cluster dissimilarity.
8. The centroids for these centroids are C1 and C2.
9. Up until K clusters have been obtained, choose the larger of C1 and C2 and establish it as the parent cluster.

## Notes

**Kernel K-Means Clustering:** The final clusters in Kernel K-means clustering are obtained when the data is projected onto the high-dimensional kernel space. The kernel function is used to first map the data points in the input space onto a high-dimensional feature space. Polynomial, Gaussian and sigmoid kernels are a few significant kernel functions. Equation contains the formula for the kernel K-means' SSE criterion as well as the cluster centroid. Additionally, the following is the formula for the kernel matrix K for any two points  $x_i, x_j \in C_k$ :

$$SSE(C) = \sum_{k=1}^K \sum_{x_i \in C_k} \|\phi(x_i) - c_k\|^2$$

$$c_k = \frac{\sum_{x_i \in C_k} \phi(x_i)}{|C_k|}$$

$$K_{x_i x_j} = \phi(x_i) \cdot \phi(x_j)$$

The sole differentiating factor between the conventional K-means criterion and the novel kernel K-means criterion is in the utilisation of the projection function  $\phi$ . In order to compute the Euclidean distance between a point and the cluster centroid in the high-dimensional feature space using kernel K-means, it is sufficient to have knowledge of the kernel matrix K. Hence, the clustering process can be performed for the data points  $x_i$  and  $x_j$  belonging to cluster  $C_k$ , without relying on the explicit individual projections  $\phi(x_i)$  and  $\phi(x_j)$ . The production of the kernel matrix from the kernel function for the given data necessitates a higher computing cost compared to the utilisation of K-means. A variant of the technique called Weighted Kernel K-means has also been developed. Spectral clustering is a well-researched variant of kernel K-means clustering.

**Mean Shift Clustering:** Mean shift clustering is a well-liked nonparametric clustering method that has been applied in numerous pattern recognition and computer vision applications. Through a convergence routine, it seeks to identify the modes that are present in the data. Finding the local maxima or modes existing in the data distribution is the main objective of the mean shift process. The mean shift clustering approach is built on the Parzen window kernel density estimation technique. Each point is the starting point and from there, a gradient ascent process is used until convergence. The mean shift vector might describe a path leading to a stationary point of the estimated density because it always points in the direction of the highest increase in density.

Such stationary locations are the local density maxima (or modes). One of the popular clustering techniques that belongs to the group of mode-finding techniques is the mean shift algorithm. Some of the fundamental mathematical formulae used in the mean shift clustering algorithm are provided here. With kernel  $K(x)$  and window radius  $h$ , the multivariate Parzen window kernel density estimate  $f(x)$  is produced for a set of N data points  $x_i$ , where  $i = 1, \dots, N$  on a d-dimensional space  $R^d$ . It is provided by

$$f(x) = \frac{1}{Nh^d} \sum_{i=1}^N K\left(\frac{x-x_i}{h}\right)$$

$$m_h(x) = \frac{\sum_{i=1}^N x_i \cdot g\left(\left\|\frac{x-x_i}{h}\right\|^2\right)}{\sum_{i=1}^N g\left(\left\|\frac{x-x_i}{h}\right\|^2\right)}$$

### Algorithm: Mean Shift Clustering

- 1: Select K random points as the modes of the distribution.

- 2: repeat
- 3: For each mode  $x$  calculate the mean shift vector  $mh(x)$ .
- 4: Update the point  $x = mh(x)$ .
- 5: until Modes become stationary and converge.

**Weighted K-Means Clustering:** The Weighted K-means (WK-means) Algorithm augments the basic K-means algorithm with a feature weighting method. The weights for various features are automatically learned in this iterative optimization process. Standard K-means ignores the relative relevance of different features and treats them all identically. Equation contains the modified SSE function that the WK-means clustering method optimized. Here, the clusters are numbered starting with  $k = 1$  and the characteristics start with  $v = 1, \dots, M$ . There is also a user-defined parameter  $\beta$  that uses the effect of the feature weights on the clustering. The clusters are designated by the letters  $C = C_1, \dots, C_k, \dots, C_K$ , where  $c_k$  is the cluster  $C_k$ 's  $M$ -dimensional centroid and  $c_{kv}$  is the centroid's  $v$ th feature value. In WKmeans, feature weights are modified in accordance with  $w_v$ .  $D_v$  is the total of the feature  $v$  within-cluster variances, weighted by the cluster cardinalities.

Similar to K-means clustering, the WK-means clustering technique also uses feature weights to adjust the distance metric. The weights and centroids for  $M$  features are initialized in line 1. The nearest centroids for each location are allocated in lines three and four and the weighted centroid is computed. After that, there is a weight update step that constrains the total of weights as illustrated in the equation below:

$$SSE(C, w) = \sum_{k=1}^K \sum_{x_i \in C_k} \sum_{v=1}^M s_{xik} w_v^\beta (x_{iv} - c_{kv})^2$$

$$w_v = \frac{1}{\sum_{u \in V} \left[ \frac{D_u}{D_v} \right]^{\beta-1}}$$

$$d(x_i, c_k) = \sum_{v=1}^M w_v^\beta (x_{iv} - c_{kv})^2$$

$$\begin{cases} s_{xik} \in (0, 1) \\ \sum_{k=1}^K s_{xik} = 1 \\ \sum_{v=1}^M w_v = 1 \end{cases}$$

The centroids are brought together by repeating these processes. Comparatively speaking, this algorithm is more computationally expensive than K-means. This algorithm has convergence problems, just like K-means. The Intelligent Weighted K-means algorithm is produced by combining Intelligent K-means (IK-means) with WK-means.

**Genetic K-Means Clustering:** The issue of K-means convergent to a local minimum is problematic. Stochastic optimization techniques that excel at preventing the convergence to a local optimal solution can be used to address this issue. It has been demonstrated that genetic algorithms (GA) converge to a global optimum. These algorithms develop across generations, creating a new population from the existing one throughout each generation by using a variety of genetic operators such natural selection, crossover and mutation. They create a fitness function, select the fittest person from each generation based on the probability score and use them to construct the following population using the mutation operator. With the help of GA, the issue of local optima can be efficiently overcome, giving rise to the Genetic K-means algorithm (GKA).

## Notes

The initial dataset for GKA is composed of a population of strings of group numbers that were used to code the data when it was first transformed. Major steps in the GKA algorithm include the following:

- **Initialization:** To start the procedure, choose a random population. This is comparable to the K-means random centroid initialization phase.
- **Selection:** Using the probability calculation provided in Equation below, select the population's fittest members.

where  $F(s_i)$  denotes the string  $s_i$ 's fitness score in the population. To determine the quality of the solution, a fitness function is further created. Similar to the SSE of K-means is this fitness function.

- **Mutation:** Similar to the K-means assignment stage, where points are given to the centroids that are closest to them and then the centroids are updated at the conclusion of iteration, is this. Iteratively applying the selection and mutation phases results in convergence. A proof of the GA's convergence is provided and a detailed discussion of the exact GKA algorithm's pseudocode is included.

### 4.3.2 Steps to Implement K-means

The unsupervised learning algorithm is a computational method used to extract patterns and structures from unlabeled data without the need for explicit supervision or guidance. The K-Means algorithm is a popular clustering technique used in machine learning and data mining. The process of clustering involves the partitioning of an unlabeled dataset into many clusters. The variable  $K$  denotes the number of predetermined clusters that are to be generated as a component of the process. For instance, if  $K$  is equal to 2, then two clusters will be formed. Similarly, if  $K$  is equal to 3, then three clusters will be formed and so forth.

- The algorithm employed is an iterative approach that partitions the dataset lacking labels into  $k$  different clusters. Each cluster consists of a single dataset and exhibits a common set of properties.
- The technique enables the partitioning of data into many clusters and offers an efficient approach for autonomously discerning these clusters within an unlabeled dataset, obviating the necessity for any form of training.
- A centroid is assigned to each cluster in the method due to its centroid-based nature. The main objective of this approach is to minimise the overall distances between each individual data point and the clusters to which they belong.
- The technique commences by taking an unlabeled dataset as its input, partitions it into  $k$  clusters and subsequently iterates this process until all available clusters have been utilised. The determination of the value of  $k$  is a prerequisite in this algorithm.

**The k-means clustering technique is mostly utilised for two key functions:**

- The algorithm employs an iterative approach to select the optimal values for  $K$  centre points or centroids.
- Each individual data point is paired with the k-center that is closest to it. A cluster is comprised of data points that exhibit proximity to a designated k-center. Consequently, every cluster exhibits a unique set of characteristics that differentiate it from other clusters, while yet sharing certain similarities among its constituent data points.

**The below diagram explains the working of the K-means Clustering:**

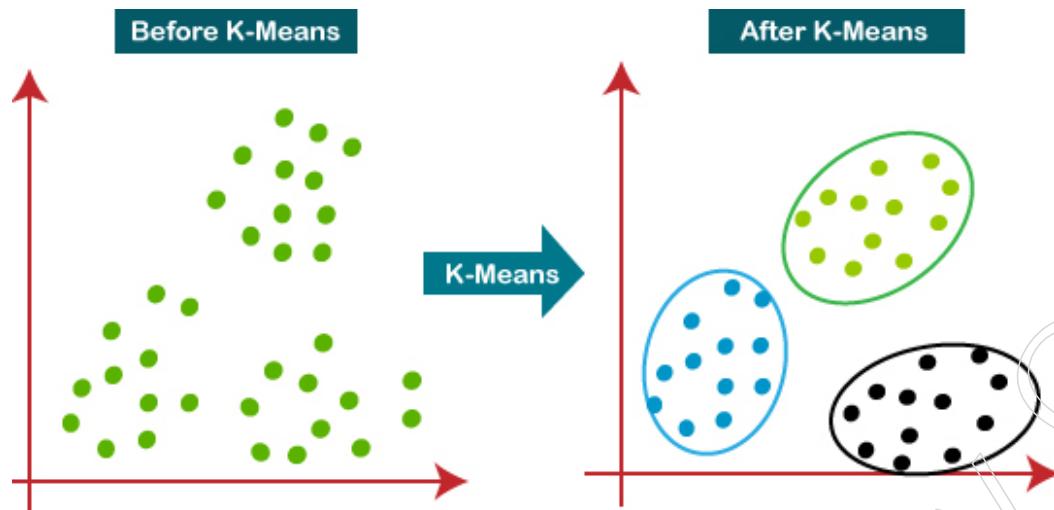


Figure : K-means Clustering

How does the K-Means Algorithm Work: Steps to Implement K-means: The following steps illustrate how the K-Means algorithm functions:

The first step in the clustering process involves selecting a value, denoted as  $K$ , which determines the number of clusters to be formed.

**Step 2:** Randomly select  $K$  locations or centroids. (It may differ from the supplied dataset).

In the third step of the process, each data point is assigned to the centroid that is closest to it, resulting in the formation of  $K$  clusters as predetermined.

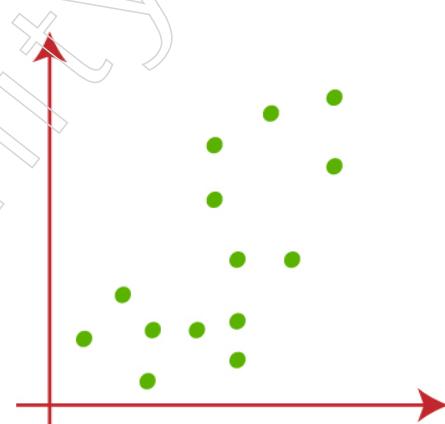
In the fourth step, the variance is computed and a new centroid is assigned to each cluster.

In accordance with phase 5 of the algorithm, it is necessary to iterate through the third phase, wherein each datapoint is reassigned to the closest centroid of its respective cluster.

**Step 6:** In the event that a reassignment takes place, proceed to step 4; otherwise, proceed to the FINISH step.

In step 7, the model has been prepared.

Let's understand the above steps by considering the visual plots: Consider that there are two variables,  $M_1$  and  $M_2$ . The following shows the  $x$ - $y$  axis scatter plot of these two variables:

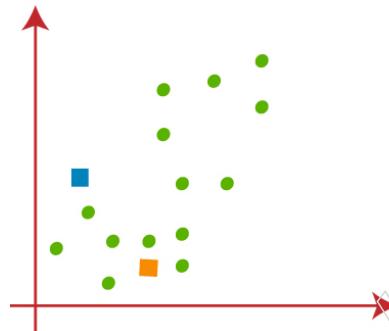


In order to categorise the dataset and partition it into distinct clusters, we shall employ the utilisation of a cluster number denoted as  $k$ , where  $k$  is equal to 2. In this particular case, our objective is to partition these datasets into two separate clusters.

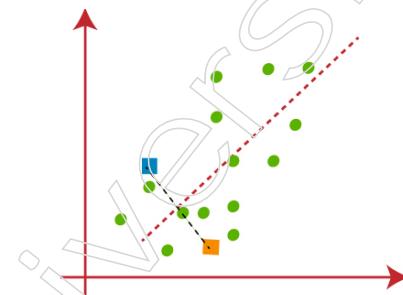
Notes

## Notes

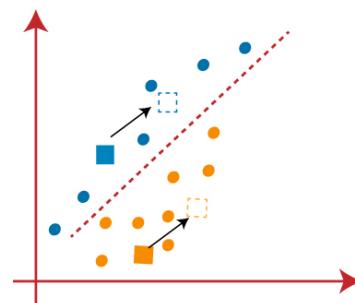
In order to establish the cluster, it is necessary to randomly select  $k$  points or centroids. These points may consist of any other points or points derived from the dataset. Consequently, we have elected to designate the two aforementioned points as  $k$  points, despite their absence from our dataset. Please consider the photograph provided below:



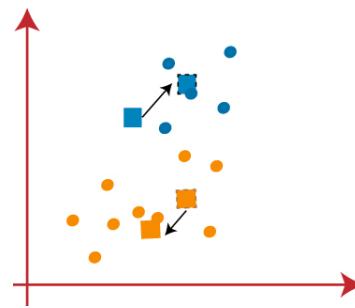
The data points of the scatter plot will be reassigned to their nearest K-points or centroid. The calculation will be performed utilising mathematical principles previously acquired in order to ascertain the distance between two given places. The median will be drawn between the two centroids. Consider the visual representation depicted below.



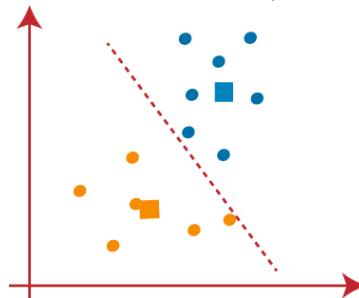
The process will be repeated by selecting an alternative centroid in order to identify the cluster that is nearest in proximity. The determination of the centres of gravity for these centroids will be conducted in order to facilitate the selection of the new centroids. This process will be carried out in the following manner:



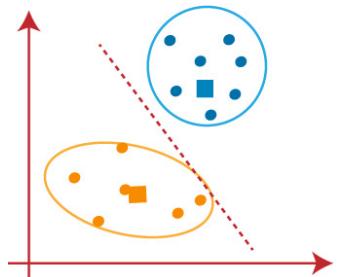
The centroids will be depicted in the subsequent image upon repetition of the procedure to ascertain their centres of gravity.



The process of redrawing the median line and reassigning the data points will be repeated in light of the updated centroids. Therefore, the depiction will be as follows:



As we can see in the image above, our model is constructed since there are no data points that are dissimilar on either side of the line. Take a look at the image below:



Now that our model is complete, we may delete the posited centroids, leaving the two final clusters as displayed in the image below:

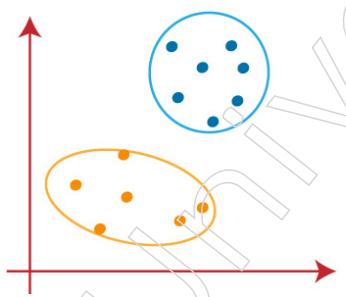


Figure: Two Final Clusters

The efficiency of the K-means clustering technique is predicated on its ability to generate highly efficient clusters. However, the process of calculating the optimal number of clusters is a significant challenge. This paper examines the optimal approach for identifying the appropriate number of clusters, commonly referred to as the value of K. There exist multiple methodologies for determining the optimal number of clusters. As an example, consider the following method:

**Elbow Method:** The Elbow approach is a commonly employed methodology for identifying the optimal number of clusters. The concept of WCSS value is utilised in this methodology. The abbreviation "WCSS" refers to the concept of total variations inside a cluster, where "WCSS" stands for within Cluster Sum of Squares. The subsequent equation may be employed to ascertain the numerical value of within-cluster sum of squares (WCSS) for a partition consisting of three clusters:

$$\text{WCSS} = \sum_{\text{Pi in Cluster1}} \text{distance}(\text{P}_i, \text{C}_1)^2 + \sum_{\text{Pi in Cluster2}} \text{distance}(\text{P}_i, \text{C}_2)^2 + \sum_{\text{Pi in Cluster3}} \text{distance}(\text{P}_i, \text{C}_3)^2$$

In the above formula of WCSS,

$\sum_{\text{Pi in Cluster1}} \text{distance}(\text{P}_i, \text{C}_1)^2$ : For the other two variables, it is the sum of the squares of the distances between each data point and its centroid within a cluster1.

## Notes

- Various methods, including the Manhattan distance and the Euclidean distance, can be employed to compute the degree of separation between the data points and the centroid.
- The elbow technique employs the following steps to ascertain the optimal number of clusters:
- The K-means clustering algorithm is applied to a specific dataset, with the K parameter varying from 1 to 10.
- Compute the within-cluster sum of squares (WCSS) value for each given value of K.
- The curve is generated by connecting the estimated within-cluster sum of squares (WCSS) values with the corresponding K-cluster count.
- When a bend has a sharp tip or a plot point that resembles an arm, it is said to possess the highest K value.

The elbow method is so named because the graph depicts a sharp bend that resembles an elbow. The elbow method's graph resembles the image below:

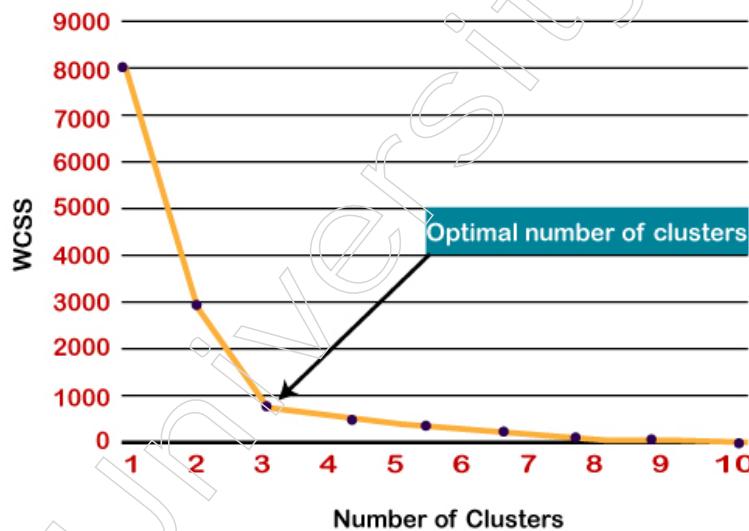


Figure: elbow method

## 4.4 Principal Component Analysis:

The primary objective of principal component analysis (PCA) is to reduce the dimensionality of a dataset consisting of numerous interrelated variables, while retaining a significant portion of the dataset's variation. To do this, a series of newly introduced variables known as principal components (PCs) is employed. These PCs are uncorrelated and arranged in a manner that ensures the initial few components capture the majority of the variance present in all of the original variables.

The contemporary definition of the introductory Principal Component Analysis (PCA) encompasses the commonly accepted process of obtaining Principal Components (PCs) by utilising eigenvectors derived from a covariance matrix.

### 4.4.1 What is Principle Component Analysis (PCA):

Principal Component Analysis (PCA) is an unsupervised learning approach commonly employed in the field of machine learning for the purpose of dimensionality reduction. The application of orthogonal transformation is a statistical procedure that converts observations of correlated features into a collection of linearly uncorrelated data.

The Principal Components refer to the modified characteristics. This particular technology is extensively utilised for the purposes of exploratory data analysis and predictive modelling. This approach involves the identification of meaningful patterns within a given dataset through the reduction of variances.

In general, Principal Component Analysis (PCA) aims to identify the subspace with the minimum dimensionality onto which the data of higher dimension can be projected.

Principal Component Analysis (PCA) operates by considering the variance of each feature, as attributes with significant variance indicate a distinct separation across classes, resulting in a reduction in dimensionality. Practical applications of Principal Component Analysis (PCA) encompass image processing, movie recommendation systems and power allocation optimisation in numerous communication channels. By employing a feature extraction technique, the model selectively retains the significant variables while disregarding the less relevant ones.

**The PCA algorithm is based on some mathematical concepts such as:**

- Variance and Covariance
- Eigenvalues and Eigen factors

**Some common terms used in PCA algorithm:**

- **Dimensionality:** The term “quantity” refers to the number of features or variables present in the dataset under consideration. The ease of determining this is facilitated by the number of columns in the dataset.
- **Correlation:** The term “correlation” refers to the extent to which two variables exhibit a relationship with each other. For example, when there is a modification in one variable, there is a corresponding alteration in the other variable as well. The correlation coefficient is bounded within the range of -1 to +1. In this context, a value of -1 signifies an inverse correlation between the variables, while a value of +1 signifies a positive correlation between the variables.
- **Orthogonal:** It establishes that there is no correlation between the variables and as a result, there is none.
- **Eigenvectors:** If a non-zero vector  $v$  is given along with a square matrix  $M$ , those are called eigenvectors. In the event where  $Av$  is  $v$ 's scalar multiple,  $v$  will then be an eigenvector.
- **Covariance Matrix:** The term “covariance matrix” refers to a matrix that contains the covariance between two variables.

**Principal Components in PCA:** The Principal Components refer to the modified attributes or outcomes obtained by Principal Component Analysis (PCA), as mentioned earlier. The number of PCs in the dataset is either equal to or less than the number of initial characteristics included. The subsequent paragraphs outline certain attributes of these fundamental constituents.

- The primary determinant must be the linear combination of distinct characteristics.
- Due to the orthogonality of these components, there exists no discernible correlation between any pair of variables.
- As the index progresses from 1 to  $n$ , the significance of each constituent diminishes, resulting in PC-1 assuming the highest level of relevance while PC-N assumes the lowest level of importance.

**Derivation of Principal Components:**

## Notes

Let us consider a vector, denoted as  $x$ , consisting of  $p$  random variables. In this context, we are particularly interested in examining the structure of the covariances or correlations between the  $p$  variables, as well as the variances of these random variables.

Examining merely the  $p$  variances and the  $1/2 p(p - 1)$  correlations or covariances may often yield limited utility, particularly when  $p$  is small or the structure is rudimentary. An alternative approach is to identify a limited number of derived variables ( $\leq p$ ) that capture the majority of the information contained in the variances, correlations, or covariances.

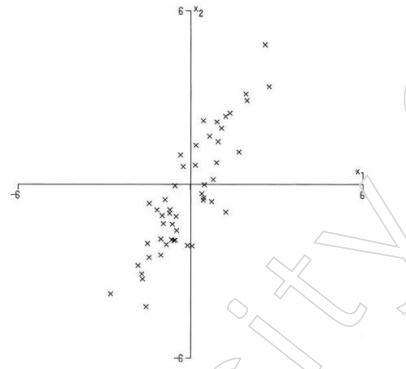


Figure: Plot of 50 observations on two variables  $x_1, x_2$ .

While covariances and correlations are taken into consideration, PCA focuses mostly on variances. The first stage is to find a linear function  $\alpha'_1 x$  of the  $x$  elements with the greatest variance, where  $\alpha'_1$  is a vector of  $p$  constants, such as  $\alpha'_{11}, \alpha'_{12}, \dots, \alpha'_{1p}$  and represents transposition.

$$\alpha'_1 x = \alpha'_{11} x_1 + \alpha'_{12} x_2 + \dots + \alpha'_{1p} x_p = \sum_{j=1}^p \alpha'_{1j} x_j.$$

Find a linear function with maximum variance that is uncorrelated with  $\alpha'_1 x$ ,  $\alpha'_2 x, \dots, \alpha'_{k-1} x$ , then look for a linear function with maximum variance that is uncorrelated with  $\alpha'_k x$ ,  $\alpha'_1 x$  and so on. The  $k$ th PC is the  $k$ th derived variable,  $kx$ . It is possible to find up to  $p$  PCs, but it is envisaged that, generally speaking,  $m$  PCs, where  $m < p$ , will account for the majority of the variation in  $x$ . In several cases later in the book, the complexity reduction gained by changing the original variables to PCs will be shown, but it will be helpful here to first explore the implausible but straightforward scenario when  $p = 2$ . The fact that the data can be plotted precisely in two dimensions is, of course, an advantage of  $p = 2$ .

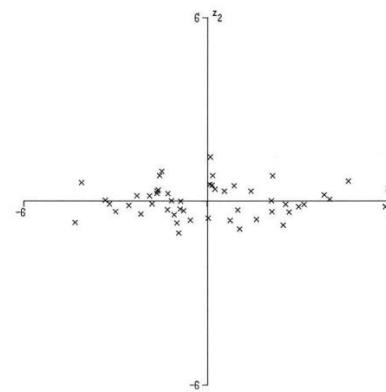
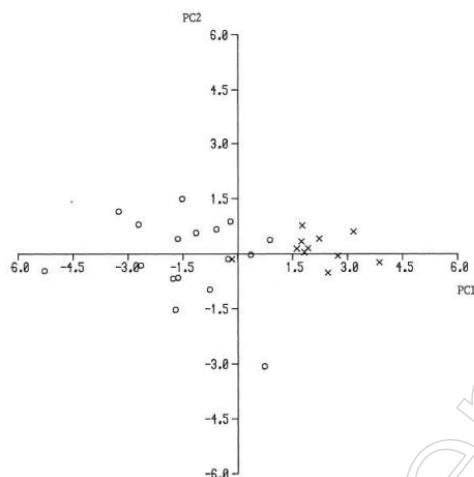


Figure: Plot of the 50 observations from above Figure with respect to their PCs  $z_1, z_2$ .

The direction of  $z_1$  demonstrates a higher degree of variance compared to the original variables, while  $z_2$  has a comparatively low level of fluctuation. In a broader sense, if a set of  $p$  ( $> 2$ ) variables exhibits substantial correlations, the majority of the

variation in the original variables can be explained by the top few principal components. On the other hand, the most recent principal components exhibit minimal variation in terms of orientations, specifically, they demonstrate nearly consistent linear correlations among the initial variables. The provided figure depicts a graphical representation of the values associated with the initial two principal components in a scenario involving seven variables. This serves as an introductory illustration to foreshadow the multitude of cases that will be presented in subsequent sections of the book. This dataset comprises seven anatomical measurements collected from a sample of 28 students, consisting of 11 females and 17 males.



**Figure: Plots of 28 pupils' anatomical measures in relation to their first two personal computers.**  
○ indicates women; ✕ indicates men.

The 2-dimensional representation of the data in the aforementioned Figure is a relatively accurate depiction of the 28 observations' placements in 7-dimensional space. The first PC, which, as we shall see later, can be taken as a measure of the overall size of each student, performs a decent job of differentiating the males and women in the sample, as is also evident from the image.

#### 4.4.2 Why Do We Need PCA:

Machine learning frequently performs effectively when a computer is educated on a sizable, well-organized dataset. Principal component analysis (PCA) is one method for addressing the problem of dimensionality in machine learning. Typically, having enough data allows us to build a prediction model that is more accurate since we have more data to utilize to train the computer. However, working with a large data set has its own disadvantages. The biggest trap is the dimensionality curse.

The title of an unpublished Harry Potter book refers to the curse of dimensionality rather than what happens when your data has too many features and perhaps not enough data points. The dimensionality curse can be beaten through dimensionality reduction. 50 variables could be decreased to 40, 20, or even 10. Here, dimensionality reduction has the greatest effects. When dealing with high-dimensional data, overfitting problems will appear and dimensionality reduction will be utilized to deal with them. avoiding information loss and improving interpretability. locates key qualities to help. helps find a linear combination of different sequences.

A common method in data analysis and machine learning is principal component analysis. Dimensionality reduction and data preparation are its main applications.

## Notes

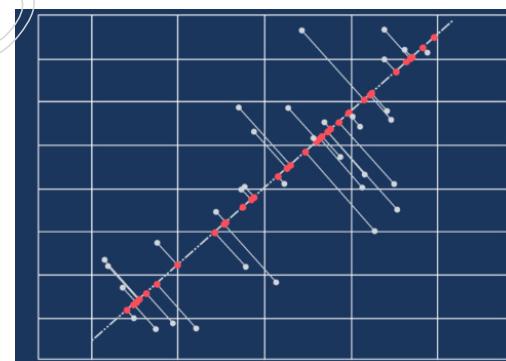
### The Benefits of PCA (Principal Component Analysis):

- Unsupervised learning techniques like PCA have several benefits. For instance, by reducing the dimensionality of the data, PCA enables us to increase the generalizability of machine learning models. The “curse of dimensionality” is thus more manageable as a result.
- Performance of algorithms is frequently influenced by the dimension of the data. High-dimensional data models may operate very slowly or even crash, consuming a lot of server resources. At a negligible loss of model accuracy, PCA can assist us in enhancing performance.
- Other benefits of PCA include its ability to extract distinct, uncorrelated features from the data and its ability to reduce data noise. PCA also makes it possible to examine clustering/classification methods and visualise data.

**A Closer Look at PCA:** Let's look at the model to better understand PCA. It presumes that the data actually exists as a low-dimensional space embedded in a high-dimensional representation. When  $L$  is a low-rank matrix,  $X$  is the original data and  $V$  is a projection operator, we presume that  $L \approx XV$ .

- We want to condense the  $n$  features in our original data to  $k$  features, where  $k \ll n$ . In order to do this, we must project the data onto a different vector space.
- There are numerous intriguing representations of PCA, but we'll focus on one that feels the most natural to us:
- We want to minimize the squared reconstruction error between our original data points and their projection onto some  $k$ -dimensional vector space in order to determine the Principal Components (PCs).
- We must mathematically identify a projection matrix  $V$  that resolves the ensuing optimization issue: Our initial data is  $\text{argmin } VTV = \|k\|X - XVV^T\|^2$ . In  $k$ -dimensional space,  $XV$  is the projection, while  $XVV^T$  is the reconstruction in the original space.

The following image serves as an example of projection from a two-dimensional space to a one-dimensional one:



**In the image above:**

- Our original data vectors are shown as blue dots.
- their projections are as the red spots.
- Reconstruction errors are represented by the blue lines.
- The ideal projection  $V$  that we discovered after solving our optimization problem is represented by the dashed black line.

Now that the data has been successfully projected from a two-dimensional ( $x, y$ ) space to a one-dimensional (line) space.

### There are several reasons why we need PCA:

**Dimensionality Reduction:** The number of features (dimensions) in real-world datasets, particularly those employed in machine learning applications, can be rather high. High-dimensional data may increase computing difficulty, necessitate more storage space and suffer from the dimensionality curse. By lowering the number of features while maintaining as much data as feasible, PCA helps to streamline the data and boost productivity.

Principal component analysis (PCA) is one of the most often used techniques for minimising the number of linear dimensions. A projection-based method is used to convert the data, putting it on a set of orthogonal (perpendicular) axes. When mapping data from a higher-dimensional space to data in a lower-dimensional space, PCA operates under the assumption that the variance or spread of the data in the lower-dimensional space should be at its maximum.

**Data Visualization:** Data visualization in high-dimensional settings (three dimensions or more) is difficult. We can display and study the data more readily, finding patterns, clusters and correlations by using PCA to reduce data to two or three dimensions.

- The well-known iris dataset has four dimensions. We are unable to show data with more than three features, unfortunately.

### We projected the data to a 2-dimensional space using PCA:

- This is particularly useful for displaying statistics to different employees inside your company. Additionally, it enables the visualization and examination of the operation of clustering and classification methods.

### Improve Classification Accuracy:

Shows that employing PCA, the accuracy of the tree classifier C4.5 has improved. The authors were able to increase model precision from 33% to 100% while also increasing model accuracy from 86% to 91% by conducting an experiment on a medical dataset from the UCI Machine Learning repository. six different machine learning classifiers' performance. While reducing the amount of characteristics from over 1000 to only a few principal components, they were still able to enhance accuracy.

**Feature Extraction:** The original characteristics are changed by PCA into new, principal component-level orthogonal features. The directions with the largest variance in the data are represented by these components, which are linear combinations of the original attributes. We can concentrate on the most informative principal components because they are rated in terms of relevance.

Consider that the first principal component vector in the previous example is (0.905, 0.423) when choosing features. With a ratio of roughly 2:1, this indicates that the projection is a linear mixture of the two features. This information could be used to do feature selection. The essential features are those that we observe receiving the majority of the projection's weight in the basic principal components. That knowledge could be used to run models on our original data (without PCA treatment) and keep the model's interpretability.

**Noise Reduction:** Many datasets may have features that are noisy or contain unrelated data. By concentrating on the principle components that carry the most signal and disregarding the principal components that primarily contribute to noise, PCA can assist eliminate or lessen the influence of noise. Let's look at the following example to show how noise can be reduced:

## Notes

- The distribution of the data is dominated more by noise on the y-axis than it is along the x-axis.
- This noise is removed by transforming the data into a 1-dimensional space.
- A “cleaner” signal can be obtained by using PCA to remove noise from data that represents a signal (image, audio).

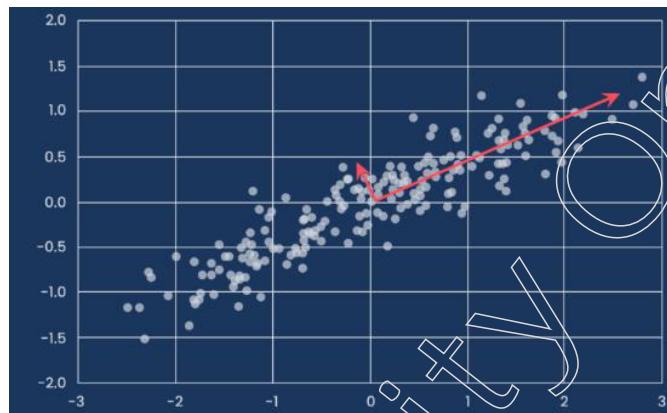


Figure: Reduce Noise in Data

**Multicollinearity:** Multicollinearity refers to a dataset's features having a high degree of correlation with one another. Regression and other statistical studies may suffer as a result, producing unstable model estimates. By converting the correlated characteristics into uncorrelated principle components, PCA helps reduce multicollinearity.

**Preprocessing for Machine Learning:** Before using machine learning techniques, PCA is frequently employed as a data preprocessing step. PCA helps speed up the training and testing of machine learning models while keeping the key information by reducing the dimensionality of the input.

**Computational Efficiency:** Processing highly dimensional data can be computationally expensive. PCA may considerably speed up the calculations required for a variety of data processing jobs by lowering the number of features.

### 4.4.3 Basic Terminologies of PCA

It will be necessary to have some basic mathematical background knowledge to comprehend the principal components analysis procedure. It is more crucial to comprehend the rationale for the usage of a particular mathematical technique and what the operation's outcome reveals about our data than it is to remember the precise workings of the technique. Even though not all of these strategies are employed in PCA, those that do serve as the foundation upon which the most crucial procedures are built. Included a section on statistics that looks at measurements of the distribution of the data, or how it is spread out. The second portion, which deals with matrix algebra, examines eigenvectors and eigenvalues, two crucial matrices' characteristics that are essential to PCA.

**Standard Deviation:** We require a data set in order to comprehend standard deviation. The majority of the time, statisticians are focused on sampling a population. A sample is a subset of the population that statisticians measure and we may use election polls as an example to illustrate the difference. The wonderful thing about statistics is that it allows you to determine what would most likely be the measurement if you used the full population by simply measuring a sample of it (in this case, by doing a phone poll or something similar).

I'm going to make the assumption in this section of the statistics that our data sets are samples of a larger population. A reference pointing to more details concerning samples and populations can be found later in this section. Here is a set of examples:

$$X = \{1 2 4 6 1 2 1 5 2 5 4 5 6 8 6 7 6 5 9 8\}$$

One could just refer to the complete group of numbers by the symbol. Subscripts can be used on the symbol to designate a specific number if one wishes to refer to a specific number in the data set. Example relates to the number 4, which is the third number in. Not like you might see in some textbooks, that is the first number in the sequence. The number of elements in the set will also be denoted by the symbol. We can calculate a variety of information about a data set. We could figure out the sample mean, for instance. The following formula is providedas it is assumed that the reader is aware of what the mean of a sample is:

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

**Variance:** Another indicator of the distribution of data in a data set is variance. In actuality, it resembles the standard deviation very closely. It goes like this:

$$s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{(n - 1)}$$

There is no square root in the formula for variance, therefore you will see that this is just the standard deviation squared in both the symbol (/ 4) and the expression. The typical notation for sample variance is / 4. These two measurements both reflect the data's distribution. The most popular metric is standard deviation, but variance is sometimes employed. Variance has been included in addition to standard deviation in order to establish a strong foundation for the covariance section that follows.

**Covariance:** The last two measurements we examined are only one dimensional. These data sets could include things like the room's population's heights, the results of the most recent COMP101 exam, etc. However, a lot of data sets have more than one dimension and the goal of statistically analyzing these data sets is typically to determine whether the dimensions are related in any way. For instance, we might include both the average height of each student in a class and the grade they earned for that particular paper in our data set. Then, using statistical analysis, we could determine whether a student's height has any bearing on their grade. You can only compute the standard deviation for each dimension of the data set separately from the other dimensions since standard deviation and variance only apply to one dimension. To determine how much the dimensions deviate from the mean relative to one another, it is helpful to have a comparable measure.

**Such a metric is covariance:** The covariance between two dimensions is always measured. The variance can be calculated by calculating the covariance between one dimension and itself. As a result, you could calculate the covariance between the dimensions <and =, <and > and =and > if you had a 3-dimensional data set (<, =, >). You might determine the variance of the<, = and> dimensions, respectively, by measuring the covariance between <and<, or = and =, or > and >. The covariance formula is quite similar to the variance formula. The variance formula might alternatively be expressed as follows:

$$var(X) = \frac{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}{(n - 1)}$$

where I've just made the square word larger to include both portions. Given that information, the covariance formula is as follows:

## Notes

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(n - 1)}$$

Remember that covariance is always calculated between two dimensions when looking at the covariance matrix. There are multiple covariance measurements that can be determined if the data set has more than two dimensions. Could you, for instance, calculate B&CD from a set of three-dimensional data (dimensions <, = and >)? 0 < E = 2, 0 B&CD? B&CD and 0 < E > 2? 0 = E > 2. In fact, you can determine! PO Q!SR4DT O U4 various covariance values for a set of -dimensional data.

**Matrix Algebra:** The background information for the matrix algebra needed for PCA is provided in this section. We'll focus on a matrix's eigenvalues and eigenvectors in particular. It is once more presumed that you are familiar with matrices.

$$2 \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 6 \\ 4 \end{pmatrix}$$

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 6 \\ 4 \end{pmatrix} = \begin{pmatrix} 24 \\ 16 \end{pmatrix} = 4 \times \begin{pmatrix} 6 \\ 4 \end{pmatrix}$$

Figure: Example of one non-eigenvector and one eigenvector

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 11 \\ 5 \end{pmatrix}$$

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 12 \\ 8 \end{pmatrix} = 4 \times \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

Figure: Example of how a scaled eigenvector is still an eigenvector

As you are aware, if two matrices have the same size, you can multiply them together. One particular example of this is eigenvectors. In Figure, there are two times a matrix and a vector are multiplied. Unlike in the second example, where the example is exactly 4 times the original vector, the resultant vector in the first example is not an integer multiple of the original vector. Why is that so? In two dimensions, the vector is, well, a vector. An arrow pointing from the origin to the point is represented by the vector d (from the second example multiplication). One might imagine the other matrix, the square one, as a transformation matrix. The result of multiplying the matrix to the left of a vector is another vector that has undergone a transformation from its initial state. The eigenvectors result from the nature of the transformation. Consider a transformation matrix that multiplied on the left to reflect the line's vectors.

**Eigenvalues:** Eigenvalues and eigenvectors have a close relationship. Take note of how, in both instances, the original vector's scaling after being multiplied by the square matrix was the same. The eigenvalue for the eigenvector in the preceding case was 4, which was the value. Before multiplying the eigenvector by the square matrix, we could take any multiple of it and the result would always be a scaled vector that was four times as large. As a result, it is clear that eigenvectors and eigen values are always paired. The eigen values are frequently provided together with the eigen vectors when using a sophisticated programming package to construct your eigen vectors.

**Exercises For the following square matrix:**

For the following square matrix:

$$\begin{pmatrix} 3 & 0 & 1 \\ -4 & 1 & 2 \\ -6 & 0 & -2 \end{pmatrix}$$

Decide which, if any, of the following vectors are eigenvectors of that matrix and give the corresponding eigenvalue.

$$\begin{pmatrix} 2 \\ 2 \\ -1 \end{pmatrix} \begin{pmatrix} -1 \\ 0 \\ 2 \end{pmatrix} \begin{pmatrix} -1 \\ 1 \\ 3 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}$$

#### 4.4.4 How Does PCA Work:

##### Working of PCA:

A statistical technique called principal component analysis, or PCA, aims to clarify the covariance structure of a collection of data. It enables us to pinpoint the primary directions in which the data vary. Consider the triangles in figure, for instance, as a two-variable data set that we measured using the X-Y coordinate system. The U axis displays the primary direction of variation in the data and the V axis, which is perpendicular to it, displays the secondary direction. The U V axis system provides us with a concise representation when we set it at the data's mean.

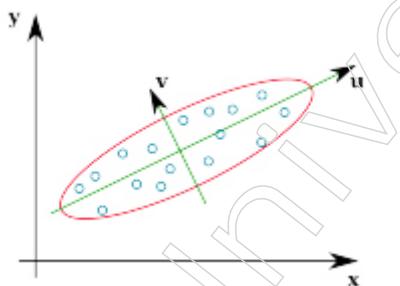


Figure: PCA for Data Representation

The data is de-correlated if each (X, Y) coordinate is converted into its corresponding (U, V) value, which results in zero covariance between the U and V variables. principal component analysis identifies the axis system bounded by the primary directions of variance for a given collection of data. The major components are the directions U and V.

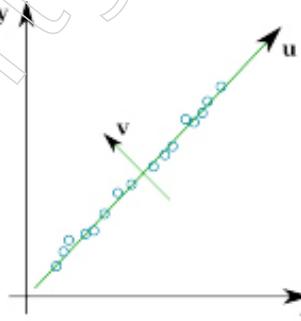


Figure: PCA for Dimension Reduction

We can anticipate a data set's variance to be regularly distributed if it results from a natural property or a random experimental error. In this instance, a hyper-ellipse (in this case, a two-dimensional ellipse) is used to illustrate the notional extent of the normal distribution. The hyper ellipse contains data points deemed to be members of a class. It

## Notes

can be considered a class boundary because it is drawn at a distance beyond which a point's likelihood of belonging to the class is low.

We can use PCA to reduce the dimensionality of a data collection if the variance in the data is due to another relationship. Think about two variables that are almost linearly connected, as in figure 2. Similar to figure 1, the U axis displays the primary direction of variation in the data while the V axis displays the secondary direction. The V coordinates are all, however, very close to 0 in this instance. For instance, we could infer that the only reason they are not zero is because of experimental noise. As a result, we can represent the data set using only variable U and ignore variable V in the U-V axis system.

**Computing the Principal Components:** Computing the data covariance matrix's eigenvectors and eigenvalues yields the primary components in terms of computation. Finding the axis system in which the covariance matrix is diagonal is the equivalent of this approach. The direction with the greatest variation is represented by the eigenvector with the largest eigenvalue; the next highest variation is represented by the eigenvector with the second largest eigenvalue; and so on. We will provide a brief overview of eigenvectors and eigenvalues so that you can see how the computation is carried out.

**Let A be an  $n \times n$  matrix. The eigenvalues of A are defined as the roots of:**

$$\text{determinant}(A - \lambda I) = |(A - \lambda I)| = 0$$

where I is the identity matrix of size  $n \times n$ . The characteristic equation, often known as the characteristic polynomial, has  $n$  roots. Make  $\lambda$  eigenvalue for A. Then a vector x exists in such a way that  $Ax = \lambda x$

Associated with the eigenvalue  $\lambda$ , the vector x is referred to as an eigenvector of A. In the equation above, x does not have a single solution. It can be scaled to any magnitude because it is merely a direction vector. To locate a copy of Lecture 15 of DOC493: Intelligent Data Analysis and Probabilistic Inference Inorder to get a single numerical solution for x, we must assign one of its components to an arbitrary value, let's say 1. This results in a collection of concurrent equations that we must then solve for the remaining components. If a solution is not found, the procedure is repeated with a different element. In most cases, we normalize the final values to have length 1, or  $x \cdot x^T = 1$ . Assuming we have a  $3 \times 3$  matrix A with eigenvalues  $\lambda_1, \lambda_2$  and  $\lambda_3$  and eigenvectors  $x_1, x_2$  and  $x_3$ , respectively, then:

$$Ax_1 = \lambda_1 x_1 \quad Ax_2 = \lambda_2 x_2 \quad Ax_3 = \lambda_3 x_3$$

**Using the eigenvectors as a matrix's columns results in:**

$$A \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}$$

writing:

$$\Phi = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \quad \Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}$$

The matrix equation is presented to us. The equation  $A\Phi = \Phi\Lambda$  is expressed. The eigenvectors have been appropriately scaled to have a magnitude of one and are mutually perpendicular. Therefore, the following relationships hold:  $\Phi\Phi^T = \Phi^T\Phi = I$ , which signifies that  $\Phi^T A \Phi = \Lambda$  and A can be expressed as  $\Phi\Lambda\Phi^T$ . Now, let us consider

the connection between this concept and the covariance matrix employed in Principal Component Analysis (PCA).

Consider a square matrix of size  $n \times n$  denoted as  $\Sigma$ , which represents the covariance matrix. For  $\Sigma$  to be equal, it is necessary to have an orthogonal matrix  $\Phi$  of size  $n \times n$ , where the columns of  $\Phi$  represent the eigenvectors of  $\Sigma$ . Additionally, a diagonal matrix  $\Lambda$  with the eigenvalues of  $\Sigma$  as its diagonal elements must exist, such that the equation  $\Phi^T \Sigma \Phi = \Lambda$  holds true.

In the example of figure, the matrix of eigenvectors can be viewed  $\Phi$  as a linear transformation that converts data points in the  $[X, Y]$  axis system into the  $[U, V]$  axis system. The data points are transformed into a data set with uncorrelated variables in the general situation by the linear transformation  $\Phi$  provided by. The data's correlation matrix in the new coordinate system is  $\Lambda$  and all of the off-diagonal entries have zeros.

#### PCA Transformation:

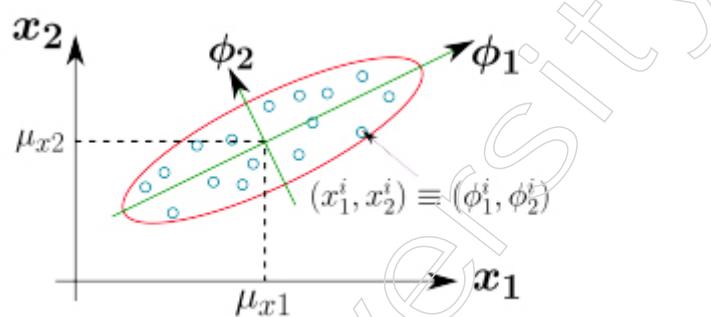


Figure : The PCA Transformation

Figure provides a two-dimensional geometric representation of the process. We determine the mean values of the variables  $(\mu_{x1}, \mu_{x2})$  and the covariance matrix  $\Sigma$ , which in this case is a  $2 \times 2$  matrix, using all of the data points. The direction vectors  $\phi_1$  and  $\phi_2$  are obtained by computing the eigenvectors of the covariance matrix. We generate a transformation matrix by adding the two eigenvectors as columns to the matrix  $\Phi = [\phi_1, \phi_2]$ , which converts our data points from the  $[x_1, x_2]$  axis system to the axis  $[\phi_1, \phi_2]$  system with the formula:

$p\phi = (px - \mu_x) \cdot \Phi$ ,  $px$  represent an arbitrary point in the axis system  $[x_1, x_2]$ , where  $\mu_x = (\mu_{x1}, \mu_{x2})$  is the mean of the data. Additionally, let  $p\phi$  be the coordinate of the point in the axis system  $[\phi_1, \phi_2]$ .

**Face Recognition:** Principal component analysis has been widely applied in face recognition, mostly to lower the number of variables. Let's take a look at the 2D scenario where we have an input image and want to compare it to a collection of database photos to determine which one is the best match. We presume that all of the photographs are equally framed and have the same resolution. (i.e., the faces are displayed in the photographs at the same size and location). Since each pixel may be thought of as a different variable, we have a very high dimensional problem that PCA can help to simplify. In a formal context, it may be stated that an input image containing a total of  $n$  pixels can be mathematically represented as a single point inside the image space.

The image space, in turn, refers to an  $n$ -dimensional space that is specifically utilised in the field of image recognition. The row vector is formed by concatenating the rows of pixels in the picture. Therefore, for an image of moderate size, such as one with a resolution of 128 by 128 pixels, the resulting vector will have a dimension of 16384. The

## Notes

individual ordinates of this point correspond to the intensity values of each pixel in the image, forming a row vector denoted as  $\mathbf{p}_x = (i_1, i_2, i_3, \dots, i_n)$ . As an illustration,

$$\begin{bmatrix} 150 & 152 & \cdots & 151 \\ 131 & 133 & \cdots & 72 \\ \vdots & \vdots & \ddots & \vdots \\ 144 & 171 & \cdots & 67 \end{bmatrix} \quad 128 \times 128$$

becomes the vector: [150, 152, ..., 151, 131, 133, ..., 72, ..., 144, 171, ..., 67]16K

This many variables is obviously much more than is necessary to solve the problem. The majority of the image pixels will have strong correlations. For instance, adjacent background pixels are perfectly connected if the background pixels are all equal. As a result, we must think about how to reduce the number of variables.

**Dimension Reduction:** Let's take a look at a scenario where there are  $N$  photos, each with  $n$  pixels. Our complete set of data can be expressed as a  $N \times n$  data matrix  $D$ . One image from our data collection is represented by each row in  $D$ . For instance, we might have:

$$D = \begin{bmatrix} 150 & 152 & \cdots & 254 & 255 & \cdots & 252 \\ 131 & 133 & \cdots & 221 & 223 & \cdots & 241 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 144 & 171 & \cdots & 244 & 245 & \cdots & 223 \end{bmatrix} \quad N \times n$$

The initial step of Principal Component Analysis (PCA) involves shifting the origin of the coordinate system to the mean of the dataset. In this particular application, the process of obtaining the mean image  $\mu$  involves the computation of the average of the columns of matrix  $D$ . Next, the mean image is subtracted from each image in the data set (i.e., each row of  $D$ ) in order to get the mean-centered data vector, denoted as  $U$ . Let us consider the hypothetical scenario in which the image has been transformed to have a mean-centered distribution.

$$[120 \ 140 \ \cdots \ 230 \ 230 \ \cdots \ 240]$$

$$U = \begin{bmatrix} 30 & 12 & \cdots & 24 & 25 & \cdots & 12 \\ 11 & -7 & \cdots & -9 & -7 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 24 & 31 & \cdots & 14 & 15 & \cdots & -17 \end{bmatrix} \quad N \times n$$

Computing the covariance matrix from the mean-centered data matrix is a straightforward process. The equation  $\Sigma = U^T U / (N - 1)$  is a mathematical expression with dimensions  $n \times n$ . The eigenvectors and eigenvalues of  $\Sigma$  are computed using the conventional method described earlier. In order to fulfill the given condition  $\Sigma = \Phi \Lambda \Phi^T$ , it is necessary to determine the values of  $\Phi$  and  $\Lambda$ . If the eigenvectors are normalized, the set of vectors  $\Phi$  forms an orthonormal basis for the system.

$$\forall \phi_i \phi_j \in \Phi, \phi_i \cdot \phi_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

The user's text is incomplete and does not provide enough information to be rewritten in an academic The axis system serves as a means by which data can be represented in a concise manner. Size reduction can be accomplished by opting to represent our data with a reduced number of dimensions. Typically, the selection of eigenvectors from the set of  $m$  (where  $m$  is less than or equal to  $n$ ) is based on the  $m$

biggest eigenvalues of  $\Sigma$ . In the context of facial recognition systems, it is common for the value of  $m$  to be relatively small, typically ranging from 20 to 50. The PCA projection can be performed by composing the  $\varphi$  vectors into a  $n \times m$  matrix  $\Phi_{pca} = [\varphi_1, \varphi_2, \varphi_3 \dots \varphi_m]$ . Given an image represented by a set of pixel values  $p_x = (i_1, i_2, i_3, \dots, i_n)$ , it is possible to determine a corresponding point in the Principal Component Analysis (PCA) space by performing the following computation:  $p_\varphi = (p_x - \mu_x) \cdot \Phi_{pca}$

The vector  $p_\varphi$ , which has  $m$  dimensions, is sufficient for representing the image. A significant decrease in data size has been accomplished, since the value of  $n$  is often no less than 16K and  $m$  can be as tiny as 20. The storage of database images in the Principal Component Analysis (PCA) space enables efficient searching for the most similar match to a given test image within the database. The inverse transform can be utilized to rebuild any given image.

$$p_x = p_\varphi \cdot \Phi_{pca}^T + \mu_x$$

It can be demonstrated that the selection of the  $m$  eigenvectors of  $\Sigma$ , which possess the highest eigenvalues, results in the minimization of the mean square reconstruction error across all possible selections of  $m$  orthonormal bases. It is evident that there is a need to minimize the value of  $m$  while ensuring accurate recognition and reconstruction. However, it is important to note that the optimal value of  $m$  is contingent upon the specific characteristics of the data. The decision-making process can be influenced by the proportion of total variance explained by the selected  $m$  main components. The assessment of this can be conducted by examining the eigenvalues. Consider the notation  $\sum_{j=1}^n \lambda_j$  to represent the sum of all  $n$  eigenvalues. The  $\Sigma$  symbol used in this context represents summation, rather than covariance. The fraction of the variance explained by the  $i$ -th eigenvector can be mathematically represented as:

$$r_i = 100 \times \frac{\lambda_i}{\sum_{j=1}^n \lambda_j}$$

Subsequently, the selection of the parameter  $m$  might be made in accordance with a heuristic criterion that is contingent upon the specific application. One possible approach to ensure that a minimum percentage of the entire variance, such as 95%, is accounted for, is by setting the condition  $\sum_{j=1}^m r_j \geq 95$ . An alternative approach would involve the exclusion of eigenvectors whose corresponding eigenvalues contribute less than 1% to the overall variance. Regrettably, a precise criterion for determining the appropriate amount of eigenvectors to maintain remains elusive.

**Few samples and many variables:** In the domain of face recognition, it is common to encounter a substantial number of variables, denoted as  $n$ , which can reach a magnitude of 16,000. However, the corresponding number of images, denoted as  $N$ , is considerably smaller. This implies that despite the covariance matrix having a dimension of  $n \times n$ , its rank can only reach a maximum of  $N - 1$ . As a result, the number of non-zero eigenvalues will be restricted to at most  $N - 1$ . The rank of the data may decrease if the data points are not linearly independent. In the event that all data points are perfectly aligned, a single eigenvector will possess a non-zero eigenvalue. This eigenvector will align with the line connecting the data points, resulting in a covariance matrix rank of 1. Due to the considerable magnitude of the variable count, the computation of the eigenvectors of  $\Sigma$  poses challenges. Considerable amounts of calculation time can be saved in the process of obtaining the  $N - 1$  non-zero eigenvalues by utilizing a technique initially introduced by Kohonen and Lowe.

The technique known as Principal Component Analysis (PCA) is alternatively referred to as the Karhunen-Lowe transform or the Householder transform in certain

## Notes

contexts. The precise application being referred to is a mathematical technique known as singular value decomposition (SVD). In this application, a  $n \times m$  matrix is transformed into a diagonal form. The concept under consideration has a strong correlation with the mathematical technique known as single value decomposition (SVD), which serves as a broader method for diagonalization applicable to matrices of non-square dimensions.

**Correspondence in PCA:** The pixels in 2D face images are frequently out of alignment. In other words, a certain pixel  $[x_i, y_i]$  may represent a portion of the cheek in one image, a portion of the hair in another, etc. This presents significant challenges for both face reconstruction and face recognition. This implies that each linear combination of eigenfaces represents a composition of face components rather than a real face. Only when the eigenfaces are combined in precisely the appropriate ratios to recreate one of the original face images from the training set is a real face produced. It is conceivable to construct a correspondence between each point on the surface map, though, if we depict a face in three dimensions. We can arrange the data so that every anatomical feature, such as the tip of the nose, has the same index across all of the distinct face data sets. Figure depicts two distinct faces with matching surface points.

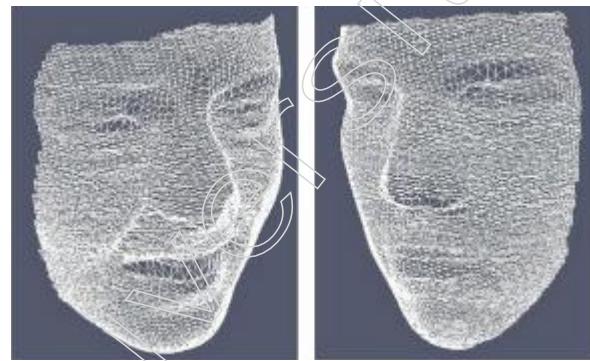


Figure : 3D surfaces in correspondence.

On 3D surface maps of faces (and other anatomical structures), PCA can be applied. The surface points' unique coordinates make up the variables. A data set has the following form:

$$[x_1, y_1, z_1, x_2, y_2, z_2, x_3, \dots, x_N, y_N, z_N].$$

The face maps in the illustration have around 5,000 surface points, therefore there are 15,000 variables. Any logical concatenation of the eigenvectors will serve as a valid face if the points of each subject are in correspondence and the various subjects are aligned as closely as possible in 3D prior to computing the PCA. Accordingly, PCA on 3D surface maps has the capacity to create and modify faces that are dissimilar to any examples in the training set. This has many of uses in the media and in movies. An active shape model is the collection of eigenvectors that in this fashion represents surface data. A 3D surface map can have a texture applied to it to give it a genuine facial appearance. The PCA can also contain the texture values as variables (much like the pixel values in 2D). Such models are referred to as active appearance models.

### 4.4.5 Advantages and Disadvantage of PCA

#### Advantages of Principal Component Analysis:

Overfitting, which occurs when there are too many variables in the data set, is one of the key problems when studying a high-dimensional data set. Such an overfit can be avoided by reducing the dimensionality of the data set using PCA. Finding the

contribution of each independent variable to the projected outcome can be challenging when there is multicollinearity, or significant correlation between independent variables. By lowering the multicollinearity in the dataset, PCA can make it easier to interpret the model. Machine learning algorithms will converge more quickly when we use the primary components of the data set rather than all the variables. The training time of the algorithms will shorten with less features. It can be challenging to comprehend and display a high-dimensional data set. Our high-dimensional data may be transformed into a low-dimensional data set using the PCA, allowing for considerably better visualization. If we want to extract the key characteristics from our enormous data set, doing a PCA can be a really smart idea.

One of the key advantages of PCA for time series data is that it can assist you in denoising and streamlining your data. You can concentrate on the vital aspects of the data and disregard the unnecessary ones by keeping only the most crucial elements. This can facilitate analysis and speed up the process while also enhancing model performance and accuracy. The ability to find underlying patterns and structures in your data, such as seasonality, cycles, or outliers, is another benefit of PCA. You can acquire insights into the dynamics and temporal behavior of your data by visualizing the principal components.

#### Take a look at some of the advantages of PCA:

**Removes Correlated Features:** In a real-world scenario, it is common to encounter datasets that contain a substantial number of features, often numbering in the thousands. Running the algorithm on every feature is not advisable due to the potential negative impact on algorithmic efficiency and the challenges associated with visualising a large number of characteristics in graphical representations. In order to reduce the number of variables in your dataset, it is advisable to decrease the number of characteristics. It is imperative to determine the relationships among the features, namely the correlated variables. The presence of several features poses significant challenges in manually identifying correlations, resulting in difficulties, frustration and time constraints. Principal Component Analysis (PCA) is a method that can be utilised to effectively address this issue. After the application of Principal Component Analysis (PCA) on your dataset, it can be observed that all of the Principal Components (PCs) are mutually independent. There exists no discernible correlation between the two entities.

**Improves Algorithm Performance:** The performance of the algorithm will be significantly compromised by an excessive number of characteristics. Principal Component Analysis (PCA) is a frequently utilised methodology in machine learning for enhancing the performance of algorithms by removing irrelevant factors that do not contribute to the decision-making process. The reduction in the number of characteristics results in a significant decrease in the training timeframes of the algorithms. Hence, employing Principal Component Analysis (PCA) as a means to enhance the efficiency of the method is a logical choice when the number of input dimensions is excessively large.

**Reduces Overfitting:** The primary factor contributing to overfitting is the presence of an excessive number of variables within the dataset. Hence, the utilisation of Principal Component Analysis (PCA) serves to mitigate the issue of overfitting by constraining the number of features.

**Improves Visualization:** The comprehension and visualisation of data in high-dimensional spaces pose significant difficulties. Principal Component Analysis (PCA) is a technique used to transform high-dimensional data into a lower-dimensional representation, often two dimensions, facilitating easier visualisation.

To determine which principal components have a higher variation and more influence

## Notes

than other principal components, we can utilize a 2D scatter plot. Even the most basic IRIS dataset has four dimensions, making it challenging to visualize. For easier visualization, we can reduce it to a 2-dimensional space using PCA.

### **Disadvantages of Principal Component Analysis:**

The PCA algorithm determines the greater variations' directions. Prior to determining the principle components, all the variables should have a mean of 0 and a standard deviation of 1, as the variance of a variable is assessed on its own squared scale. Otherwise, the PCA would be dominated by the variables with bigger scales. The original features of our data set will be converted into principal components, which are the linear combinations of the original characteristics, when we apply principal component analysis to our data set. But which aspects of the data set are the most important? Following the PCA, it may be challenging to provide an answer to this question. Usually, biplots are useful for such kind of interpretation.

PCA might lose some information and interpretation when applied to time series data, which is one of its key downsides. You alter the variables' initial significance and scale by turning the data into principle components. This can make it more difficult to comprehend and interpret the results, especially if you need to contextualize them in real-world situations. The assumption that the data are linearly connected and stationary, or that they don't change over time, is another drawback of PCA. For some time series data, though, which can have nonlinear and nonstationary features, this might not be the case. PCA might not accurately reflect the data's underlying structure and variability in this situation.

- PCA has its own shortcomings, just like any other framework. First, PCA presupposes a linear relationship between the variables. PCA will yield inaccurate findings if the data is embedded on a nonlinear manifold.
- PCA is also outlier-sensitive. Such data inputs may result in results that greatly deviate from the data's intended projection.
- When it comes to interpretability, PCA has drawbacks. Features lose their original meaning because we are changing the data. In situations when the data's interpretability is crucial, this can be an issue. However, there are circumstances in the feature selection example that we previously stated when we can still only partially interpret the model.
- Finally, only continuous, non-discrete data are appropriate for PCA. PCA is a poor choice if some of our characteristics are categorical.

### **Using the Principal Component Analysis method can also have some disadvantages:**

The interpretability of independent variables is reduced as Principal Components (PCA) is utilised to transform the information, so transforming the original features into principal components (PCs). The major components are formed by combining your distinct characteristics in a linear manner. The readability and comprehensibility of original features surpass that of Principal Components.

Data standardisation is a necessary step to be undertaken prior to performing Principal Component Analysis (PCA). Prior to employing Principal Component Analysis (PCA), it is imperative to normalise the data. Failure to do so would hinder PCA's ability to discern the optimal Principal Components. For example, when a feature set includes data expressed in units such as kilogrammes, light years, or millions, the variance scale of the training set becomes significantly enormous.

If Principal Component Analysis (PCA) is used to a feature set with high variation, it will result in significant loadings for those features. Consequently, the principal components will exhibit a preference for traits characterised by significant volatility, thereby leading to the generation of erroneous conclusions. In order to utilise PCA, it is necessary to convert all categorical attributes into numerical features for the purpose of normalisation.

You must scale the features in your data before performing PCA because scale has an impact on PCA. Use Scikit's Standard Scaler. Learn to standardize dataset features onto unit scale (mean = 0 and standard deviation = 1), which is necessary for many machine learning algorithms to operate at their best.

**Information Loss:** Principal Components attempt to account for as much variance among the characteristics in a dataset as possible, but if the number of principal components is not carefully chosen, it may leave out some information from the original list of features.

#### 4.4.6 Application for Principal Component Analysis:

**PCA Application example:** 2020 China's COVID-19 in Ottawa Province 2020 China's COVID-19 in Ottawa Province:

**Data collection:** We chose 10 typical indicators, including "Time to Return to Work," to evaluate the effects of various measures implemented by various local governments on epidemic prevention and control. The "Mode of Travel" The phrase "Measures for Returning Persons" The phrase "Requirements for Body Temperature Measurement for Returning Persons" The "Isolation Period for Returning Persons" Number of phone calls made looking for returning people Requirements for daily necessities How often residents go shopping for necessities Relatives and friends' visit Urban control Residents of Ottawa Province, where the epidemic condition is severe in Canada, received 885 surveys in total. There were 834 valid questionnaires and the response percentage was 94.20%. 58% of the valid questionnaires were filled out by men and 42% by women.

**Calculation by principal component analysis:** The scores of each indicator in nine regions are averaged and the correlation matrix is computed (as shown in the table). The ten indices are represented by the numbers X 1, X 2, X 3, X 4,..... X 10 accordingly.

Table : Correlation matrix:

	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>	X <sub>7</sub>	X <sub>8</sub>	X <sub>9</sub>	X <sub>10</sub>
X <sub>1</sub>	1.000	.341	-.026	.594	-.229	.129	-.062	-.146	.243	-.308
X <sub>2</sub>	.341	1.000	-.770	-.297	-.392	-.294	-.452	.370	.061	.084
X <sub>3</sub>	-.026	-.770	1.000	.459	.612	.627	.094	-.286	.200	-.073
X <sub>4</sub>	.594	-.297	.459	1.000	.097	.014	.183	-.655	.130	-.378
X <sub>5</sub>	-.229	-.392	.612	.097	1.000	.641	-.601	-.279	-.264	-.252
X <sub>6</sub>	.129	-.294	.627	.014	.641	1.000	-.446	.267	.117	.112
X <sub>7</sub>	-.062	-.452	.094	.183	-.061	-.446	1.000	-.209	.323	.150
X <sub>8</sub>	-.146	.370	-.286	-.655	-.279	.267	-.209	1.000	-.005	.436
X <sub>9</sub>	.243	.061	.200	.130	-.264	.117	.323	-.005	1.000	.678
X <sub>10</sub>	-.308	.084	-.073	-.378	-.252	.112	.150	.436	.678	1.000

It is clear from the output findings' correlation matrix that there is a strong relationship between the ten metrics. Therefore, additional research can be conducted using the principal component analysis method.

**Determination of the number of principal components:** The sum of numerous common factor variances is known as the common factor variance. The better the overall

## Notes

effect, the higher the rate of cumulative contribution, representativeness, or interpretation of the extracted common factors for the original variables.

	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$X_6$	$X_7$	$X_8$	$X_9$	$X_{10}$
Initial	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
Extract	0.931	0.939	0.940	0.897	0.917	0.917	0.915	0.727	0.886	0.881

The data in the last column of the above Table are between 70% and 95%, which indicates that the extracted principal components have a high degree of explanation for each variable, according to the common factor variance of the output findings.

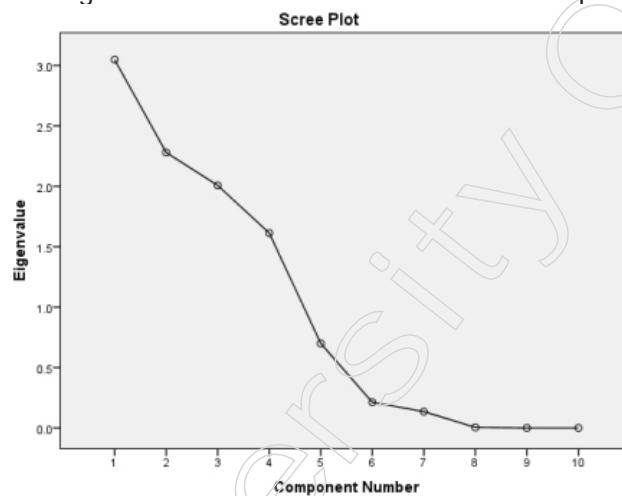


Figure: Principal component line chart.

A broken-line map known as a gravel map places feature roots in size order from large to tiny. This paper extracts four main components because the first four eigen values change visibly and the fifth eigen value becomes softer. The sum of the variances of the first four major components accounted for 89.488% of the overall variance in the output result graph. As a result, the first four main components may essentially preserve the data from the original indicators, resulting in a decrease in the number of indicators from the original 10 to 4.

The association's strength is indicated by the correlation coefficient's absolute value. Each primary component's original variables, as represented by this data, can be identified. The primary component is identified on this basis. In this example, the first principal component—also known as “the degree of intervention for the returned people from Wuhan”—mainly refers to the actions taken for the returned Wuhan residents, the requirements for taking their temperatures and the isolation period. The second important factor, which may be referred to as “the degree of intervention for the returned people from Wuhan,” mostly deals with the purchasing requirements for daily essentials and the frequency of going out to acquire necessities. The third primary factor, which can also be called “policy implementation strength,” generally relates to the Chinese New Year’s visits from family and friends and the city’s level of control. The term “policy strictness” can be used to describe the fourth primary component, which primarily focuses on the return to work time, transport mode and phone call volume.

**Principal component expression.** The component score coefficient matrix can be used to determine the principal component expression. The four main components produced in this example from the original variables after normalization are written as. The following is the primary component expression:

**Notes**

$$\begin{aligned}
 f_1 &= 0.029zx_1 - 0.240zx_2 + 0.289zx_3 + 0.201zx_4 + 0.246zx_5 \\
 &\quad + 0.174zx_6 - 0.023zx_7 - 0.187zx_8 - 0.032zx_9 - 0.134zx_{10}. \\
 f_2 &= -0.195zx_1 + 0.038zx_2 + 0.032zx_3 - 0.277zx_4 + 0.253zx_5 \\
 &\quad + 0.270zx_6 - 0.315zx_7 + 0.251zx_8 - 0.119zx_9 + 0.089zx_{10}. \\
 f_3 &= -0.149zx_1 - 0.206zx_2 + 0.197zx_3 - 0.069zx_4 - 0.073zx_5 \\
 &\quad + 0.120zx_6 + 0.253zx_7 + 0.106zx_8 + 0.365zx_9 + 0.401zx_{10}. \\
 f_4 &= 0.494zx_1 + 0.293zx_2 + 0.033zx_3 + 0.197zx_4 - 0.017zx_5 \\
 &\quad + 0.275zx_6 - 0.230zx_7 + 0.103zx_8 + 0.319zx_9 + 0.094zx_{10}.
 \end{aligned}$$

To obtain the score for each principal component and the overall score for each city, the standardized data for each city is inserted into the principal component expression. Table displays the scoring results.

**Table: Scores by region**

City	Component 1	Component 2	Component 3	Component 4	Total score
Ezhou	0.771	1.441	-0.185	2.059	4.086
Huanggang	-1.060	0.300	0.360	-0.345	-0.745
Huangshi	-0.095	-1.086	-1.902	0.425	-2.657
Jingmen	0.050	1.662	-1.010	-1.695	-0.993
Jingzhou	0.682	-1.058	-0.354	-0.154	-0.884
Suizhou	-1.548	0.059	0.940	0.387	-0.162
Xiangyang	-0.652	-0.370	0.453	0.133	-0.437
Xiaogan	0.134	-0.775	0.334	-0.222	-0.528
Yichang	1.718	-0.173	1.364	-0.588	2.321

The table shows that Ezhou and Yichang perform the best in terms of policy intervention.

**Summary**

- Clustering is a fundamental technique in unsupervised machine learning that involves grouping similar data points together based on certain criteria. The primary goal of clustering is to discover patterns or structure within a dataset without using predefined labels or classes. It is particularly useful for exploratory data analysis, pattern recognition and data compression. Clustering is commonly used in various fields, such as marketing, biology, image analysis and customer segmentation. It helps in gaining insights into the inherent structure of data, identifying outliers and organizing large datasets into meaningful groups.
- K-means clustering is one of the most widely used clustering algorithms. It is a partitioning algorithm that aims to divide a dataset into K clusters, where each data point belongs to the cluster with the nearest mean (centroid). K-means aims to minimize the within-cluster sum of squares, which is the sum of squared distances between each data point and its assigned centroid. However, K-means is sensitive to the initial placement of centroids and may converge to local optima. Multiple runs with different initializations can help mitigate this issue. Remember that the choice of K (number of clusters) is crucial and may require domain knowledge or techniques such as the elbow method or silhouette analysis to determine an appropriate value. K-means is efficient and works well on large datasets, but it assumes that clusters are spherical and equally sized, which might not always be the case.
- Principal Component Analysis (PCA) is a dimensionality reduction technique widely

## Notes

used in various fields, including machine learning, statistics and signal processing. Its primary objective is to transform high-dimensional data into a lower-dimensional representation while preserving as much of the original variability as possible. PCA achieves this by identifying a set of new orthogonal axes, called principal components, that capture the most significant patterns or directions of variance in the data. It's important to note that while PCA is a powerful technique, it does have some limitations. For instance, it assumes that the principal components represent meaningful patterns in the data, which may not always be the case. Additionally, the interpretability of the transformed features can be challenging, as they are combinations of the original features. Overall, PCA is a versatile tool for reducing the dimensionality of high-dimensional data while retaining as much valuable information as possible.

- PCA works by transforming the original high-dimensional data into a new coordinate system defined by a set of orthogonal axes called principal components. These principal components are linear combinations of the original features and are chosen in such a way that they capture the maximum variance in the data. By retaining only a subset of the principal components, PCA effectively reduces the dimensionality of the data while preserving as much information as possible. The new representation is particularly useful for visualization, analysis and subsequent machine learning tasks.
- In summary, PCA is a powerful tool for dimensionality reduction, noise reduction and visualization. However, its application requires careful consideration of the trade-offs between reducing dimensionality and preserving meaningful information in the data.

### Glossary

- Expectation Maximization (EM): The Expectation Maximization (EM) algorithm is used to estimate the parameters of probabilistic models after presuming a particular type of generative model (for example, a mixture of Gaussians).
- Streaming Scenario: “streaming scenario,” in which it is assumed that just one pass is permitted over the data stream. Even when the data are collected offline, there are still several scalability problems when trying to use the vast amounts of data in a distributed scenario with a big data framework or integrate them with conventional database systems.
- Residual Sum of Squares (RSS): The Residual Sum of Squares (RSS) or Sum of Squared Errors (SSE) is the name of the objective function used by K-means.
- Calinski-Harabasz Index: Greater clustering performance is implied by a higher index value. The ratio of between-cluster variation to within-cluster variance is assessed using the Calinski-Harabasz index.
- Davies-Bouldin index: Since the Davies-Bouldin index measures the average similarity between each cluster and its closest comparable cluster, a lower value indicates better clustering performance.
- Rand Index: Better clustering performance is indicated by a higher Rand index. It measures how similar the expected grouping and the actual clustering are.
- Adjusted Mutual Information (AMI): A higher index indicates better grouping capabilities. The AMI measures the mutual information between the ground truth clustering and the predicted clustering after chance correction.

### Check Your Understanding

1. What is the primary goal of clustering in machine learning?

- a) Supervised classification      b) Dimensionality reduction  
c) Discovering patterns and structure within data  
d) Data transformation
2. Which clustering algorithm is known for dividing data points into dense regions separated by sparse areas?
- a) K-means clustering      b) DBSCAN  
c) Agglomerative clustering      d) Spectral clustering
3. In K-means clustering, what does the "K" represent?
- a) Number of iterations      b) Number of clusters  
c) Number of features      d) Number of data points
4. Which evaluation metric measures the compactness of clusters in K-means?
- a) Silhouette score      b) F1-score  
c) Precision      d) Recall
5. Hierarchical clustering is a type of clustering that:
- a) Divides data into dense regions separated by sparse areas.  
b) Uses probabilistic methods to assign data points to clusters.  
c) Focuses on discovering patterns in time-series data.  
d) Builds a tree-like structure of clusters by merging or splitting.
6. Anomaly detection is primarily concerned with identifying:
- a) Outliers or abnormal data points      b) Centroids of clusters  
c) Principal components      d) Mean-shift vectors
7. Which step in K-means involves updating the cluster centroids?
- a) Initialization      b) Assignment  
c) Evaluation      d) Update
8. What does PCA stand for in machine learning?
- a) Principal Component Arrangement      b) Primary Classification Algorithm  
c) Principal Component Analysis      d) Pattern Coordination Approach
9. What is the primary goal of PCA?
- a) Maximizing the variance in the data      b) Creating new data points  
c) Transforming categorical data  
d) Reducing the number of features to zero
10. What problem does PCA help address in machine learning?
- a) Overfitting      b) Underfitting  
c) Dimensionality reduction      d) Feature scaling
11. In PCA, the directions that capture the most variance in the data are called:
- a) Principal axes      b) Eigenvalues  
c) Covariance matrices      d) Centroids
12. Which step in PCA involves finding the eigenvalues and eigenvectors of the covariance matrix?

**Notes**

## Notes



## Exercise

1. Define clustering and its types.
  2. Write short note on K-means clustering.

3. What do you meant by anomaly detection?
4. Define the various steps to implement k-means clustering
5. What is Principle Component Analysis (PCA)?
6. Why Do We Need PCA?
7. How Does PCA Work? Also explain the advantages and disadvantage of PCA.

**Notes**

### Learning Activities

1. Describe strategies to mitigate the sensitivity of K-means to initial centroid selection."
2. How does Principal Component Analysis (PCA) help in dimensionality reduction and data compression?

### Check Your Understanding- Answers

1 c	2 b	3 b	4 a
5 d	6 a	7 d	8 c
9 a	10 c	11 a	12 d
13 c	14 d	15 d	16 a
17 d	18 b	19 a	20 b

### Further Readings and Bibliography

1. "Pattern Recognition and Machine Learning" by Christopher M. Bishop.
2. "Unsupervised Machine Learning: A Practical Introduction to Clustering and Dimensionality Reduction" by AleksandarMilanović.
3. "Introduction to Machine Learning with Python: A Guide for Data Scientists" by Andreas C..
4. "Hands-On Unsupervised Learning Using Python: How to Build Applied Machine Learning Solutions from Unlabeled Data" by Ankur A. Patel.
5. "Applied Unsupervised Learning with Python" by Kelleher, Mac Namee and D'Arcy.
6. "Unsupervised Learning Algorithms" by Keshav Chhabra and SiddharthShekar.

## Module - V: Deep Learning

### Learning Objectives

At the end of this module, you will be able to:

- Define the concepts of deep learning
- Analyse the neural network and its utility in modelling and solving problems
- Explain biological motivations and parallelism
- Explain recurrent neural networks
- Analyse the multilayer perceptron

### Introduction

Deep learning is a specific field within the realm of machine learning that focuses on the application of neural networks with multiple layers, generally known as deep neural networks, to effectively model and solve complex problems. Neuromorphic computing is a specialised subfield within the realm of artificial intelligence that aims to emulate the structural and functional characteristics of neural networks observed in the human brain.

The discipline of deep learning has attracted significant attention and recognition due to its ability to independently gather and represent features from raw data. The aforementioned feature has greatly contributed to its outstanding performance across many domains, encompassing, but not restricted to, image and audio recognition, as well as the processing of natural language.

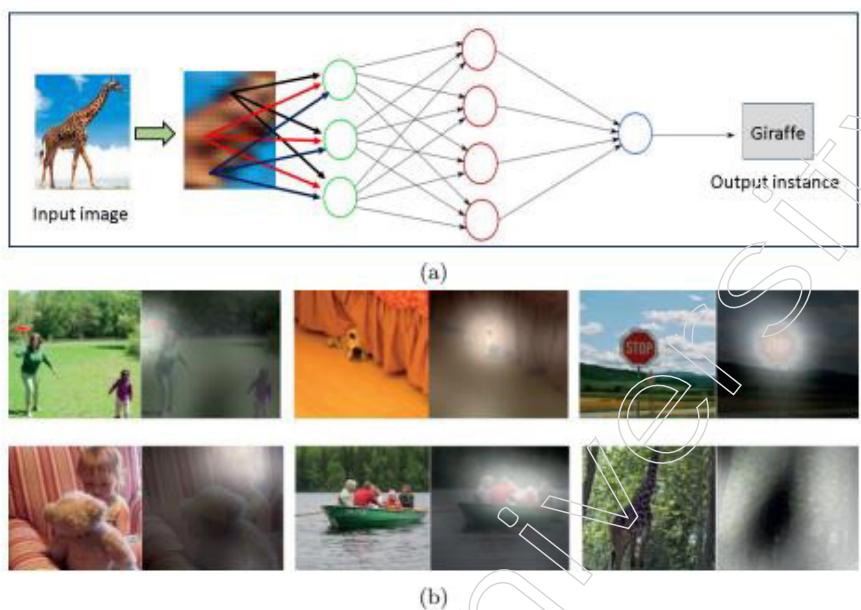
### 5.1 Concepts of Deep Learning

Deep learning is a specific field within the realm of machine learning that focuses on the application of neural networks with multiple layers, generally known as deep neural networks, to effectively model and solve complex problems. Neuromorphic computing is a specialised subfield within the realm of artificial intelligence that aims to emulate the structural and functional characteristics of neural networks observed in the human brain. The discipline of deep learning has attracted significant attention and recognition due to its ability to independently gather and represent features from raw data. The aforementioned feature has greatly contributed to its outstanding performance across many domains, encompassing, but not restricted to, image and audio recognition, as well as the processing of natural language.

#### 5.1.1 Introduction to Deep Learning: Part I

In a startling fashion, Artificial Intelligence (AI), in the numerous areas that integrate it, has begun to faithfully imitate human behavior and cognition at the beginning of the twenty-first century. As a result, notable advancements have been made in a variety of fields, including computer-assisted medical diagnosis, DNA sequence classification, data mining, artificial vision, voice recognition, written language analysis, virtual games, robotics and any others where reasoning is a key component. Deep learning is a relatively new field of AI that was developed a few decades ago with the primary goal of enabling intelligent agents to make their own decisions, a feat that is currently only imagined in science fiction. In this way, multiple Artificial Neural Networks (ANN) methodologies are utilized in deep learning with the aim of providing an agent with a personality similar to a human. Deep Neural Networks, Convolutional Neural Networks and Deep Belief Networks are a few of these methods.

A discipline known as “Deep Learning” uses numerous hidden layers of an ANN to achieve deep abstractions in order to find patterns. The border of the region of interest in a picture is selected below figure as an example of how to achieve abstraction in a hidden layer. Depth is achieved by repeating abstraction in as many hidden layers as needed. Since there is no metric that determines how many to use, both the number of hidden layers and that of neurons range from 1 to n. Instead, this attempts to resolve the agreement to the problem, the dimensions and properties of the data set and based on the experience of who implements the ANN. Our ANN will produce the desired result, which corresponds to a classification that informs us that a “giraffe” would be represented in an input image like the one in the figure below. The representation of the input and output images for deep learning is also shown in the same graphic.



**Figure: a) Abstractions using the Deep learning approach,  
b) Input and output images before and after processing with Deep learning.**

A deep model is simply a complicated tensorial computation that may be ultimately broken down into common mathematical operations from analysis and linear algebra. Over time, the field has produced a sizable number of high-level modules with a distinct semantic and intricate models that combine these modules and have been successful in a variety of application domains. Deeper structures, or lengthy compositions of mappings, result in better performance, according to empirical data and theoretical findings.

Machine thought has long been a goal of inventors. At least since the days of ancient Greece, this yearning has existed. The entities Galatea, Talos and Pandora can be interpreted as examples of artificial life, while the mythological figures Pygmalion, Daedalus and Hephaestus can be regarded as legendary pioneers in the realm of innovation.

More than one hundred years prior to the construction of the first programmable computer, individuals pondered the possibility of machines attaining intelligence. The field of artificial intelligence (AI) is now experiencing significant growth, characterised by a wide range of practical applications and continuous research endeavours. Intelligent software is employed for the purpose of automating monotonous operations, comprehending voice or visuals, facilitating medical diagnosis and providing foundational support for scientific study. A compilation of rigorous mathematical concepts can be employed to elucidate scenarios that pose intellectual challenges for humans but are

## Notes

comparatively straightforward for computers. The early stages of artificial intelligence witnessed prompt and effective resolution of such difficulties. The primary obstacle faced by artificial intelligence pertains to addressing problems that are readily and instinctively resolved by humans, such as recognising spoken words or identifying faces in images. However, these tasks are challenging for people to articulate using formalised language or explanations.

One potential solution to this challenge involves enabling computers to acquire knowledge through experiential learning and develop a comprehensive understanding of the universe by organising information in a hierarchical structure, wherein each concept is defined in relation to more fundamental notions. This approach mitigates the necessity for human operators to explicitly specify each fragment of knowledge that the computer requires, as it acquires knowledge through experiential learning. The hierarchical organisation of concepts facilitates the acquisition of sophisticated ideas by computers through the construction of complex ideas from simpler ones.

Paradoxically, tasks that are characterised by formality and abstraction, which pose significant cognitive challenges for individuals, are very straightforward for computers to accomplish. Computers have consistently outperformed even the most skilled human chess players, but only in recent times have they started to approach the level of proficiency exhibited by average individuals in tasks such as speech recognition and object detection.

In order to navigate their daily existence, individuals must possess a substantial breadth of knowledge about the world. Much of the acquired knowledge in this context exhibits characteristics of irrationality and intuition, rendering its formalisation a formidable task. For computers to exhibit intelligent behaviour, it is imperative that they possess the capability to acquire equivalent knowledge. One of the primary challenges in the field of artificial intelligence pertains to the incorporation of unstructured knowledge into machine systems.

The efficacy of these elementary machine learning algorithms is heavily influenced by the manner in which the data is represented. For example, the AI system does not conduct a physical examination of the patient when making a recommendation for a caesarean birth through the utilisation of logistic regression. In contrast, the medical practitioner furnishes the system with a range of relevant particulars, including the presence or absence of a scar on the uterus. A feature refers to any data element that is observed within the patient representation. Logistic regression can be employed to identify the associations between various patient features and distinct outcomes. Nevertheless, it is important to note that it does not have any influence on the definition of the features. In contrast to the standardised reports generated by medical professionals, it is unlikely that logistic regression can generate meaningful predictions based just on an MRI scan of the patient. The connection between specific MRI scan pixels and potential delivery difficulties is quite low.

This reliance on representations is a widespread phenomenon that may be seen in both everyday life and computer science. If the data collection is appropriately structured and indexed, computer science procedures like querying a collection of data can move forward considerably faster. Arithmetic with Roman numerals takes far longer for people to complete than with Arabic numerals. It is not surprising that the representation used has a significant impact on how well machine learning algorithms function. Take a look at Figure below for a clear visual example. By creating the ideal set of features to extract for a task and then feeding those features to a straightforward machine learning algorithm, many artificial intelligence tasks can be addressed. A measurement of the speaker's

vocal tract size, for instance, can be a useful attribute for speaker identification from sound. As a result, it strongly suggests that the speaker is either a man, woman, or little child. It might be challenging to determine which features should be retrieved for certain activities. Let's say we wanted to create a program that could identify cars in pictures. We would prefer to use the presence of a wheel as a feature since we are aware that autos have wheels. Unfortunately, it is challenging to precisely characterize a wheel's appearance in terms of pixel values. A wheel has a straightforward geometric shape, but its appearance can be complex by shadows falling on it, sunlight reflecting off its metal components, the car's fender or something in the foreground blocking part of the wheel and so on.

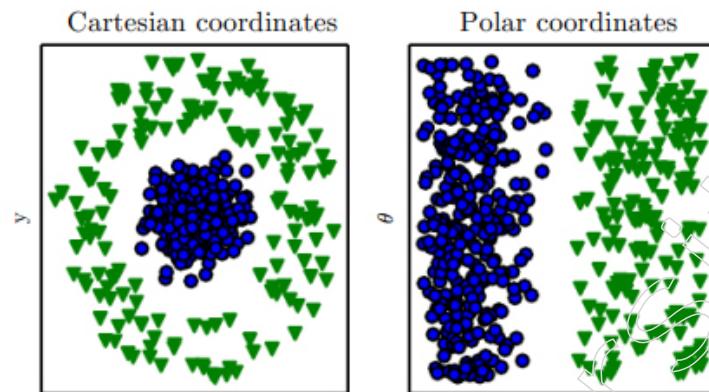


Figure: Example of different representations

Let's say we wish to use a scatterplot to divide two categories of data by drawing a line between them. It is impossible to depict part of the data in the plot on the left using Cartesian coordinates. In the figure to the right, we use polar coordinates to describe the data, which makes it easier to solve the problem using a vertical line.

Using machine learning to identify both the representation itself and the mapping from representation to output is one approach to solving this issue. Representation learning is the name given to this strategy. When compared to representations created by hand, learned representations frequently produce significantly higher performance. Additionally, they enable AI systems to quickly adapt to new jobs with a minimum of human involvement. For a simple task, a representation learning algorithm can find a decent set of features in a matter of minutes; for a difficult task, it may take hours or even months. It takes a lot of human time and effort to manually build features for a complex task; for a whole community of scholars, it can take decades.

The auto encoder is the classic illustration of a representation learning algorithm. An auto encoder combines a function called an encoder, which changes the input data's representation, with a function called a decoder, which changes the new representation back to the original format. Auto encoders are trained to make the new representation have a variety of desirable features, but they are also trained to preserve as much information as possible when an input is passed through the encoder and subsequently the decoder. Different autoencoder types strive to attain various features.

Our objective while creating features or algorithms for learning features is typically to identify the sources of variation that account for the observed data. In this context, we simply refer to individual sources of impact as "factors"; the factors are typically not multiplied together. These factors are frequently not clearly observable quantities. Instead, they might exist in the physical world as invisible forces or invisible things that have an impact on observable numbers.

## Notes

They might also be mental constructs that help to explain or infer the causes of observed facts by providing helpful simplifications. They can be viewed as ideas or abstractions that aid in our understanding of the data's extensive range of variability. The age, sex, accent and words that the speaker uses are all sources of diversity when listening to a speech tape. The position of the car, its color and the sun's angle and brightness are all variables to consider while examining a photograph of a car.

The fact that many of the sources of variation have an impact on every single piece of data we are able to examine is a significant source of difficulties in many real-world applications of artificial intelligence. At night, it's possible that the individual pixels in an image of a red automobile are almost completely black. Depending on the angle of view, the silhouette of the car changes. The majority of applications call for us to separate the sources of variation and toss out the ones we don't care about. Of course, separating such high-level, abstract traits from raw data can be exceedingly challenging. Many of these sources of variance, such as a speaker's accent, cannot be found without a comprehensive understanding of the data that is almost human-level. At first look, representation learning does not appear to be helpful when getting a representation is nearly as challenging as solving the original problem.

By incorporating representations that are expressed in terms of other, simpler representations, deep learning addresses this fundamental issue in representation learning. The machine can create complicated notions from simpler ones thanks to deep learning. Figure below demonstrates how a deep learning system can combine simpler notions like corners and contours, which are in turn described in terms of edges, to represent the concept of an image of a human.

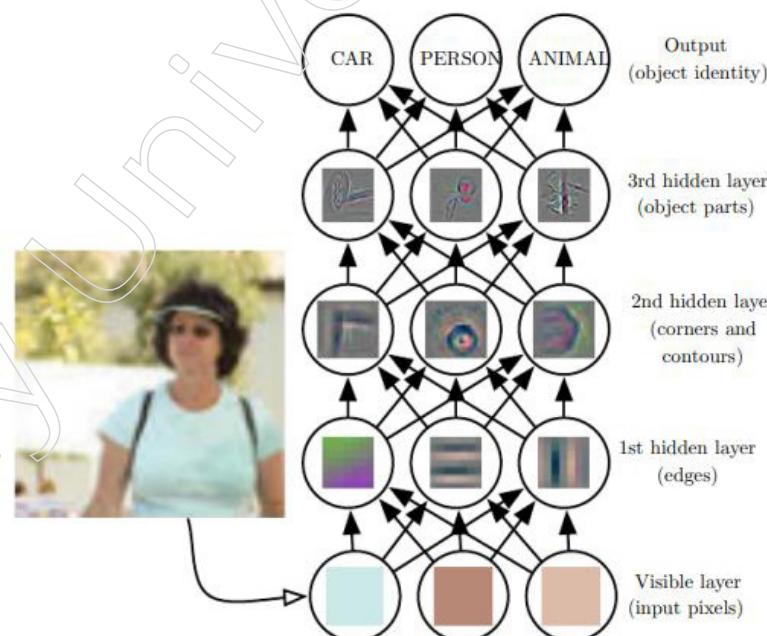


Figure: Illustration of a deep learning model

Computers face difficulties in comprehending raw sensory input data, such as a picture represented as a set of pixel values. Mapping a set of pixels to an object identity requires a substantial amount of effort. If approached directly, the task of learning or evaluating this mapping appears to be insurmountable. The solution to this difficulty is achieved through the utilisation of deep learning techniques. Deep learning effectively breaks down the intricate mapping necessary into multiple layers, each responsible for a distinct and simpler mapping inside the model.

The visible layer, referred to as such due to its inclusion of observable variables, serves as the interface for input presentation. The features of the image are subsequently retrieved through a series of hidden layers. The levels in question are commonly denoted as "hidden" due to the absence of explicit values in the observed data. Instead, the model is tasked with selecting the most suitable concepts to elucidate the relationships present in the observed data.

The visual representations in these images depict the characteristics represented by each hidden unit. The initial layer has the ability to rapidly detect edges by means of comparing the luminosity of adjacent pixels within a particular pixel set. The second latent layer may rapidly identify corners and extended contours, which are perceived as aggregations of edges, based on the edge representations provided by the first latent layer. The third latent layer possesses the capability to identify entire components of certain objects by the identification of specific clusters of corners and contours, based on the visual description provided by the second latent layer in terms of corners and contours. Finally, by utilising this depiction of the image in relation to its constituent elements, it becomes feasible to discern the entities depicted inside the visual representation.

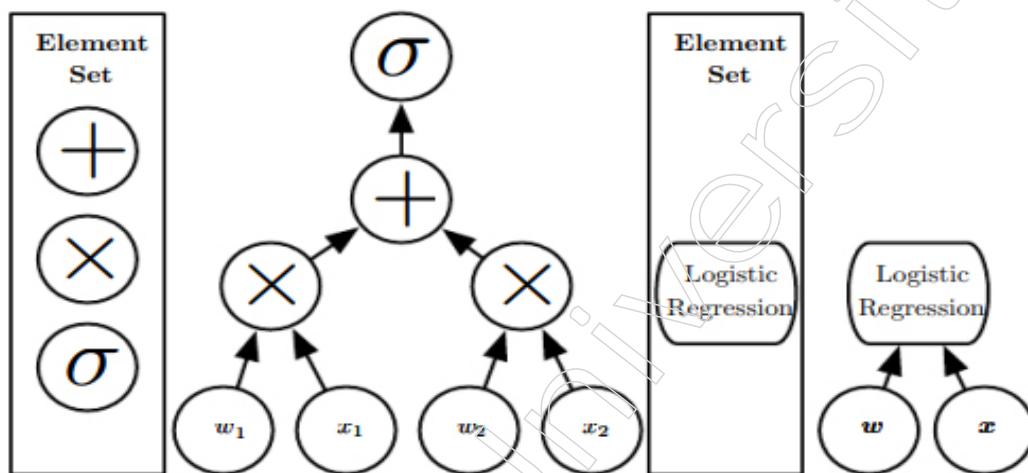


Figure: The depiction of computational networks, which map an input to an output, involves nodes that execute various operations.

Depth, which also depends on how a computational step is defined, is the distance that takes a calculation from input to output. The outcomes of a logistic regression model, where  $(w^T x)$  is the input and  $\sigma(w^T x)$  is the logistic sigmoid function, are displayed in these graphs. If we use logistic sigmoids, addition and multiplication as the building elements of our computer language, this model has depth three. If we consider logistic regression as a component in and of itself, then this model has depth 1.

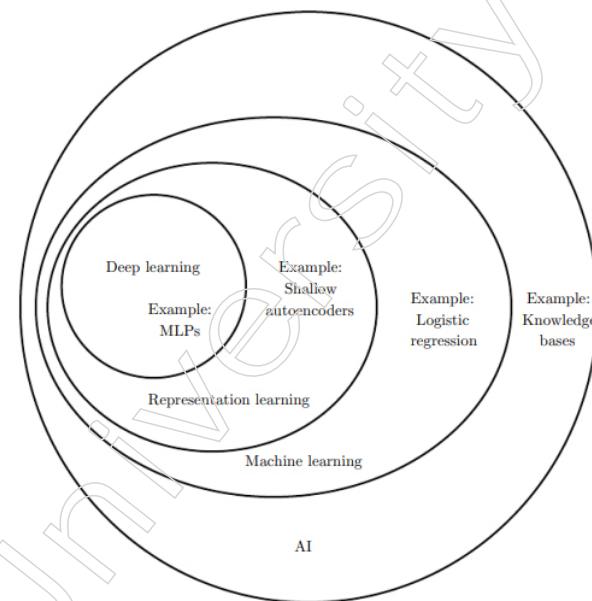
Not all of the data in an activation layer necessarily indicates reasons of variation that explain the input from a deep learning perspective. The representation also keeps track of state information so that a programme can run that can comprehend the input. This state information may be compared, like a counter or pointer in a normal computer programme. It doesn't directly affect the substance of the input, but it helps organise how the model processes data.

There are two techniques to gauge a model's depth. The first point of view centres on how many sequential operations are necessary to evaluate the architecture. A flow chart that illustrates how to calculate each output of the model given its inputs can be used to compare this to the length of the longest path. Just as two equivalent computer programmes will have various lengths based on the language the programme is written

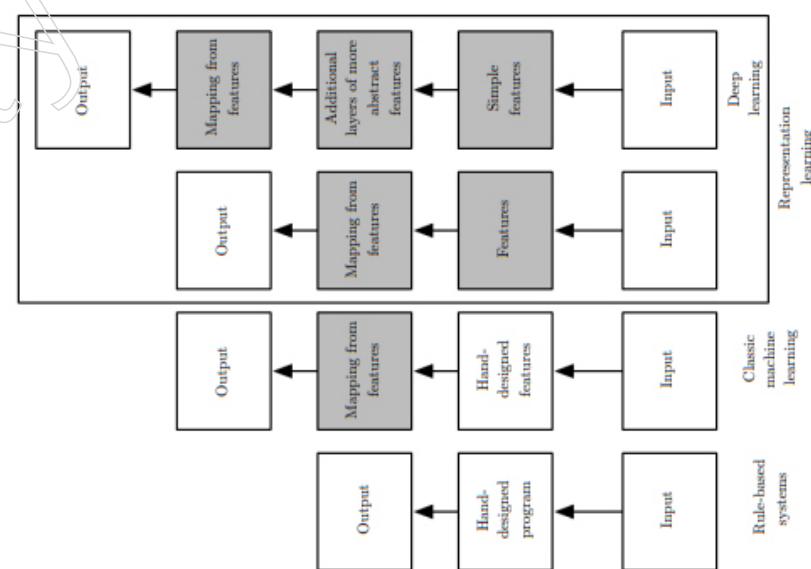
## Notes

in, the same function may be represented as a flowchart with varying depths depending on the functions we enable to be utilised as individual steps in the flow chart above. This language option can lead to two different measurements for the same architecture, as seen in Figure.

Deep learning is an AI methodology, to sum up. It specifically falls under the umbrella of machine learning, a technique that enables computer systems to improve over time and with experience. Machine learning is the only feasible approach for developing AI systems that can work in difficult real-world situations. Deep learning is a particular kind of machine learning that reaches a significant level of power and flexibility by learning to represent the world as a layered hierarchy of concepts, with each concept defined in relation to simpler concepts and more abstract representations computed in terms of less abstract ones. The image below shows the relationship between these several AI disciplines. The following graphic displays a high-level schematic of each function.



**Figure: A Venn diagram illustrating how deep learning is a subset of machine learning, which is employed in many but not all methods to artificial intelligence. An illustration of an AI technology may be found in each section of the Venn diagram.**



**Figure: Flowcharts showing how the different parts of an AI system.**

### 5.1.2 Introduction to Deep Learning: Part II

Historical Trends in Deep Learning: Deep learning is easy to comprehend in the perspective of history. Instead of giving a thorough history of deep learning, we highlight a few significant trends:

- Deep learning, despite its extensive and significant historical background, has been referred to by multiple titles that encompass diverse philosophical perspectives and its level of acceptance and recognition has experienced fluctuations throughout time.
- The increased availability of training data has led to the growing benefits of deep learning.
- The size of deep learning models has increased over time as computer hardware and software infrastructure has become more advanced.
- Over time, deep learning has more accurately and efficiently handled increasingly complex applications.

Deep learning, despite its rich and substantial historical foundation, has been characterised by several designations that span a range of philosophical viewpoints. Furthermore, its level of acceptance and acknowledgment has undergone variations throughout time. The proliferation of training data has resulted in the expanding advantages of deep learning.

The difficulties in feature engineering, which are essential for symbolic-based machine learning, can be overcome via deep learning. The amazing thing about deep learning is that it eliminates the need for computer programming entirely because the models can pick up the features automatically. Programmers simply need to provide the computer with a learning algorithm, expose it to terabytes of input data to train it and then let the machine figure out how to recognize the required items on its own. In other words, these computers can now learn on their own.

As a result, deep learning is a very potent technique for contemporary machine learning. On almost every metric, deep learning techniques outperform conventional symbolic-based machine learning techniques. The amount of money being invested, the number of people choosing deep learning as their field of study and the number of top technological businesses making AI the centerpiece of their long-term strategies are all indications that deep learning is on the rise. It has the ability to alter how individuals live their daily lives while revolutionizing many facets of machine perception. Some people even think that AI might one day be programmed to resemble human common sense.

#### Increasing Dataset Sizes:

Although the initial studies with artificial neural networks were carried out in the 1950s, one might question why deep learning has just lately come to be acknowledged as an important technology. Deep learning has been successfully applied in business since the 1990s, but until recently it was frequently thought of as more of an art than a technology and something that only a specialist could utilize. It's true that using a deep learning algorithm effectively requires some level of competence. Fortunately, as the amount of training data increases, the level of skill needed decreases.

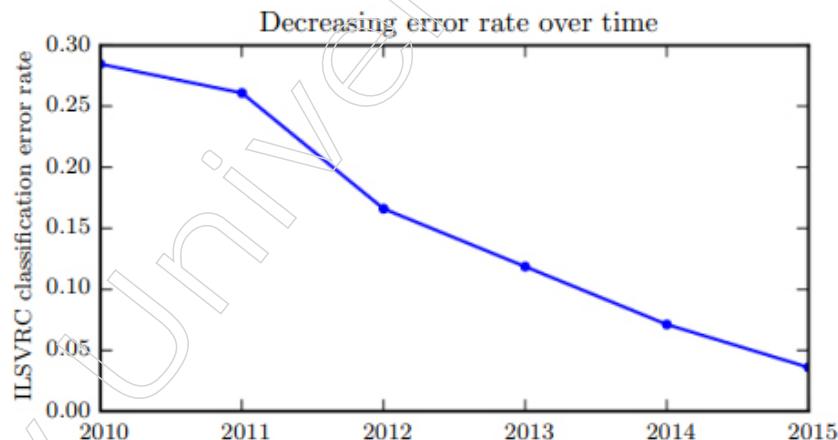
Though the models we train with these algorithms have undergone improvements that make it easier to train very deep architectures, the learning algorithms used to achieve human performance on hard tasks today are essentially identical to the learning algorithms that struggled to handle toy problems in the 1980s. The most significant innovation is that we can now provide these algorithms the tools they require to be successful.

## Notes

### Increasing Accuracy, Complexity and Real-World Impact

The initial deep models were designed to identify specific items in extremely small, closely cropped images. Since then, the size of the images that neural networks could analyze has gradually increased. Modern object recognition networks interpret detailed, high-resolution images without needing to crop the image in close proximity to the object to be identified. Similarly, while these contemporary networks routinely detect at least 1,000 different categories of items, the early networks could only distinguish two types of objects (or, in some circumstances, the absence or presence of a single kind of object).

The Image Net Large-Scale Visual Recognition Challenge (ILSVRC), which is held annually, is the biggest object recognition competition. The victory of a convolutional network in this challenge, surpassing its competitors by a substantial margin, signified a pivotal moment in the exponential growth of deep learning. The outcome of this triumph resulted in a notable decrease in the top-5 error rate, which was lowered from 26.1% to 15.3%. This signifies that the convolutional network generates a list of potential categories for each image and the accurate category was included in this list for all test samples, except for a mere 15.3%. Subsequently, deep convolutional networks have consistently emerged as the victors in these competitions. At the time of composing this text, the progress made in deep learning has resulted in a reduction of the most recent top-5 error rate in this competition to 3.6%, as depicted in the accompanying Figure.



**Figure:** Deep networks have consistently won the competition every year and produced ever-lower mistake rates since they achieved the scale required to participate in the ImageNet Large Scale Visual Recognition Challenge.

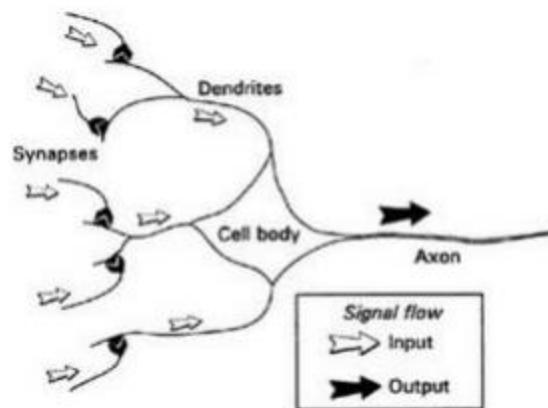
Deep learning is a method of machine learning that has been developed over the past several decades and mainly relies on our understanding of the human brain, statistics and applied arithmetic. Its popularity and utility have greatly increased in recent years, largely as a result of more potent computers, greater datasets and methods for training deeper networks. The coming years will present many difficulties and chances for deep learning to advance and reach new heights.

## 5.2 Introduction to Neural Networks

Let's start with a working definition of what a "neural network" is, then move on to straightforward, straightforward explanations of some of the major elements in the definition:

A neural network refers to a network composed of interconnected processing units, also known as nodes, which exhibit a functional resemblance to the action of biological

neurons in animals. The network's processing capacity resides in the interunit connection strengths, also known as weights, which are obtained through a process of adaptation or learning from a given set of training patterns. In order to facilitate further elaboration, it is imperative to commence by briefly revisiting certain fundamental principles in the field of neuroscience. The human brain, consisting of approximately 100 billion neurons or nerve cells, is represented in a highly stylized fashion in this illustration.



**Figure: Essential components of a neuron shown in stylized form**

Neurons utilise electrical signals, characterised as transient voltage fluctuations or "spikes," to facilitate intercellular communication. Synaptic electrochemical junctions, located on dendritic cell branches, facilitate interneuronal connectivity.

Typically, individual neurons possess numerous connections with other neurons, leading to a continuous stream of messages that ultimately converge at the cell body. In this context, the elements are merged or unified, resulting in the neuron generating a voltage impulse, commonly referred to as "firing," when the resultant signal surpasses a pre-established threshold. The axon, a dendritic extension, serves as a conduit for transmitting this information to adjacent neurons.

Certain incoming signals exhibit an inhibitory influence, impeding the firing process, whereas others possess an excitatory impact, promoting the generation of impulses. The synaptic connections of individual neurons, characterised by their type (excitatory or inhibitory) and intensity, are believed to play a crucial role in shaping the distinct information processing capabilities of each neuron.

The designation of connectionist systems arises from the recognition of the importance placed on interneuron connections. Consequently, the broader investigation of this approach is commonly known as connectionism. The objective is to integrate this particular design and processing technique into neural networks. The utilisation of the aforementioned terminology is commonly observed while examining neural networks within the framework of psychologically inspired models of human cognitive processes.

The nodes or units in our initial description serve as synthetic counterparts to biological neurons, with a representative illustration shown in Figure below. In the context of neuronal communication, it is customary to assign a weight to each input signal prior to its transmission to the cell body. This is due to the fact that synapses, which serve as the connections between neurons, are typically represented by a singular numerical value denoting their strength or efficacy. In this context, the activation of a node is generated through the summation of weighted signals, employing fundamental arithmetic operations. The comparison of activation to a threshold occurs in a node known as a threshold logic unit (TLU), as illustrated in the accompanying figure. If the activation

## Notes

surpasses the threshold, the TLU generates a high-valued output, typically represented as "1". Otherwise, it produces a zero output. Exclusively positive weights have been employed. The diagram utilises arrow widths to symbolise signals and multiplication symbols enclosed in circles to represent weights. The values assigned to the circles are intended to be proportional to the size of the symbol. The TLU, or Threshold Logic Unit, serves as a fundamental model of a synthetic neuron.

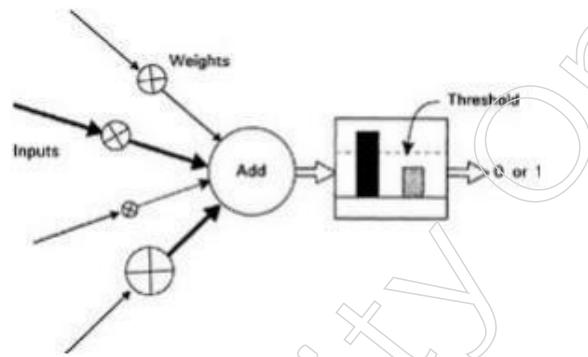


Figure : Simple artificial neuron.

Any artificial neuronal network shall be referred to as a "network" in this paper. This could be as basic as a single node or as complex as a big group of nodes, each of which is connected to every other node in the network. The figure below illustrates one kind of network.

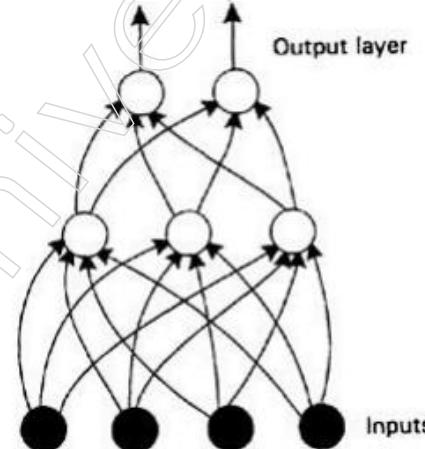


Figure: Simple example of neural network.

in the current configuration, the weights are incorporated implicitly in all connections, although the representation of each node remains limited to a circular shape. Every signal is generated from an input and undergoes two successive nodes prior to reaching an output, at which point it undergoes no further modifications. The placement of the nodes follows a tiered framework.

The feed forward structure, among several alternatives, is frequently employed to categorise an input pattern into one of multiple categories, relying on the resulting output pattern. The output layer, situated at the highest position in the picture, may consist of 26 nodes, each representing a letter of the alphabet. These nodes are utilised to determine the specific letter category to which the input character belongs. For example, if the input is a representation of the light and dark patterns in an image of handwritten letters.

To do this, one output node would be assigned to each class and it would be necessary for only one of these nodes to fire whenever a pattern from the relevant class

was provided at the input. So much for the fundamental structural components and how they work. Returning to our working definition, take note of the significance placed on experience-based learning. Real neurons' synaptic strengths may be altered in specific situations to allow each neuron to respond differently or adapt to its unique sensory input. The equivalent in artificial neurons is the alteration of the weight values. The "knowledge" the network is meant to have is stored in its weights, which develop through a process of adaptation to stimulus from a set of pattern instances. There are no computer programs involved in this information processing.

An input pattern is supplied to the net in one training paradigm known as supervised learning, which is employed in conjunction with nets of the kind shown in Figure above and its response is then compared with a goal output. For instance, if we were using our prior letter recognition example, we might enter "A" and compare the network output to the classification code for A. The weights are then adjusted based on the discrepancy between the two output patterns.

Every specific recipe for change functions as a learning principle. A new pattern is provided, the output is compared to the target and fresh alterations are applied after the necessary weight updates have been made. Iteratively repeating this series of events increases the likelihood that the network's behavior will eventually converge and each pattern's response will be near to the matching target. The entire procedure, including any pattern presentation ordering, process termination requirements, etc., makes up the training algorithm.

### Why study neural networks?

Neural networks are commonly utilised in statistical analysis and data modelling as an alternative to conventional nonlinear regression or cluster analysis techniques. Consequently, these methods are commonly utilised in situations that pertain to categorization or forecasting. Illustrative instances encompass textual character recognition, image and speech recognition, as well as human sectors of expertise such as medical diagnosis, oil exploration geology and forecasting financial market indicators.

Neural networks are perceived by engineers and computer scientists as a form of parallel distributed computing, presenting a viable alternative to the conventional algorithmic approaches that have historically dominated the field of machine intelligence. This particular problem also pertains to the field of classical artificial intelligence (AI). Practitioners in this domain often prioritise the ease of implementing solutions in digital hardware or the efficacy and accuracy of particular approaches, rather than emphasising biological realism.

Nets, referred to as computer models of the animal brain, have garnered attention from neuroscientists and psychologists due to their abstraction of key properties of nerve tissue believed to be vital for information processing. There exists a degree of scepticism among neuroscientists over the final effectiveness of these simplified models, since they argue that a more comprehensive level of information is necessary to provide a thorough understanding of the mechanisms behind brain function. The artificial neurons employed in connectionist models often exhibit significant simplifications compared to their biological counterparts. The ultimate outcome remains uncertain, however, notable advancements have been made in the realm of replicating brain functions by the utilisation of knowledge pertaining to the interconnectedness of actual neurons, referred to as local "circuits."

- Neural networks were first developed in the early 1940s. In the late 1980s, it saw a rise in popularity. This came about as a result of new methods and innovations

## Notes

being discovered as well as general advancements in computer hardware technology.

- Some NNs are models of biological neural networks, but historically, the field was inspired by the desire to create artificial systems capable of complex, possibly “intelligent,” computations similar to those the human brain performs, thereby improving our understanding of the brain.
- Most NNs have a “training” rule of some kind. In other words, NNs “learn” from examples in the same way that kids pick out dogs from pictures of dogs and show some ability to generalize outside of the training data.

### NNs vs Computers:

Digital Computers	Neural Networks
Reasoning deductively. To produce output, we use input data and established rules.	Reasoning inductively. We build the rules using training examples as input and output data.
Centralized, synchronous and serial computation all exist.	Collective, asynchronous and parallelism all describe computation.
Memory is physically stored, addressable by location and packetized.	Memory is internalized, dispersed and addressable by content.
not tolerant of faults. It stops working when one of the transistors fails.	Redundancy, fault tolerance and responsibility sharing.
Exact.	Inexact.
Static connectivity.	Dynamic connectivity.
Provided applicable and provided the rules and input data are clearly established..	Relevant when data is noisy or incomplete or when rules are ambiguous or complex.

Neural networks are designed to mimic the information processing mechanisms of the human brain, enabling them to replicate fundamental cognitive capabilities. Due to its efficient processing and rapid reaction capabilities, it is employed for a diverse range of real-time tasks. The construction of an artificial neural network encompasses various components that are derived from the biological nervous system. An artificial neural network is composed of numerous interconnected processing components, sometimes referred to as Nodes. These nodes establish a connection link with other nodes in order to facilitate their interconnection.

The connection link is characterised by the presence of weights, which serve to encapsulate information pertaining to the incoming signal. The weights undergo updates with each iteration and input. The weights of the neural network and its architectural configuration are commonly referred to as the “trained neural network” once all training data instances have been processed. This method is commonly referred to as the training of neural networks. In order to address the specific concerns delineated in the problem statement, the employed neural network model is utilised. Artificial neural networks have the capability to address a wide range of tasks, including but not limited to classification problems, pattern recognition and data clustering.

Artificial neural networks are used because they can learn quickly and adapt. They are able to use the training data they receive to learn “how” to tackle a particular problem. After learning, you can utilize it to address that particular problem fast, effectively and

accurately. Air Traffic Control, Optical Character Recognition, which is utilized by various scanning apps like Google Lens, Voice Recognition, etc. are some examples of real-world neural network applications.

Notes

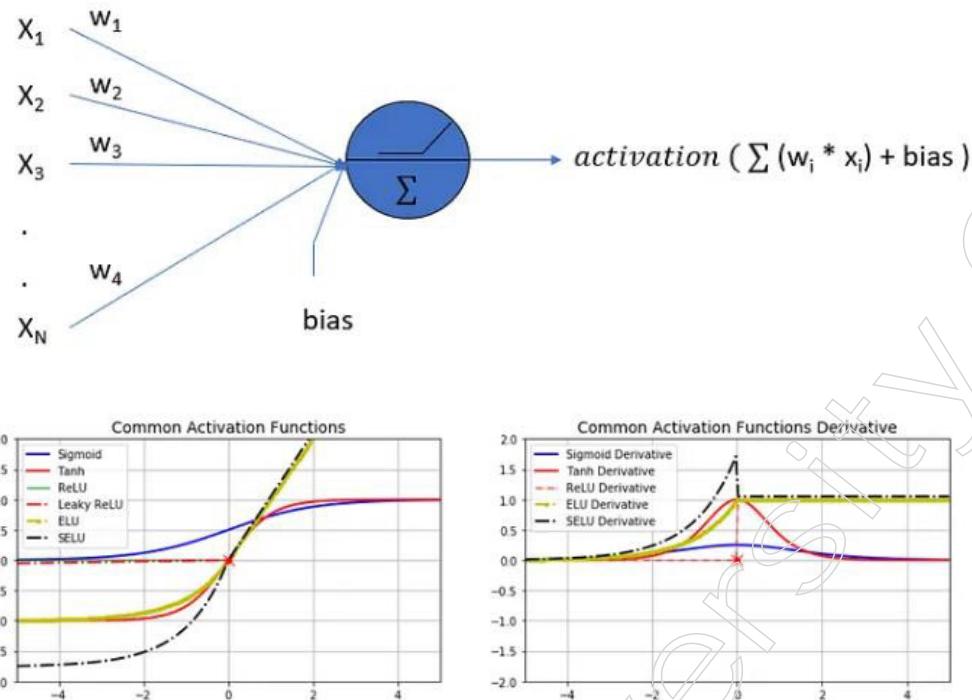


Figure: Common activation functions in neural networks

The activation function is a crucial hyper-parameter during the training of neural networks and we must decide which activation function to employ in both the hidden layers and the output layer. The activation function determines whether or not to fire a neuron by calculating the bias term and the weighted total of the input. The input signal is subjected to an activation function, which may be a linear or non-linear transformation and the output is given as input to the next layer of neurons.

#### Basic building blocks of neural network training:

It is vital to explain the fundamental building blocks of neural network training before going on to other activation functions and their variants. These blocks can be further broken down into the forward pass, the network output error measurement and the backward pass. The training cases are first given to the neural network in forward pass. A forward cascade of calculations is then performed utilizing the current set of weights across the layers to get the network's output prediction. Second, the discrepancy between the desired output and the actual projected output—the output error of the network—is measured. Thirdly, we back-propagate through each layer to determine the contribution of each connection to the error before adjusting the connection weights to lower the error.

#### Humans Versus Computers: Stretching the Limits of Artificial Intelligence:

One additional advantage is that neural networks provide a rapid means of adjusting the complexity of a model by incorporating or removing neurons from the architecture based on the quantity of training data or computational resources at hand. The increasing availability of data and the enhanced computational powers of modern computers have surpassed the limitations of traditional machine learning methods, rendering them

## Notes

incapable of fully harnessing the current possibilities. This phenomenon elucidates a substantial proportion of the recent achievements observed in neural networks. The depicted figure illustrates the given condition.

	Artificial Intelligence	Human Intelligence
<b>Processing</b>	Based on algorithms and mathematical models	Based on cognitive processes and biological structures
<b>Learning</b>	Based on data and feedback loops	Based on experience, intuition, and creativity
<b>Speed</b>	Can process data and perform tasks much faster than humans	Slower than AI in processing large amounts of data, but can make complex decisions quickly
<b>Adaptability</b>	Can quickly adapt to new data and situations	Can adapt to new situations, learn from experience, and make decisions based on context
<b>Emotions</b>	Lacks emotions and empathy	Capable of feeling emotions and empathy
<b>Creativity</b>	Limited ability to be creative or think outside of the box	Capable of creativity, imagination, and innovation
<b>Ethics</b>	Does not have a moral code or conscience	Has a moral code and conscience that guides decision-making
<b>Physical Limitations</b>	Does not have physical limitations, can operate 24/7	Limited by physical capabilities and requires rest and maintenance

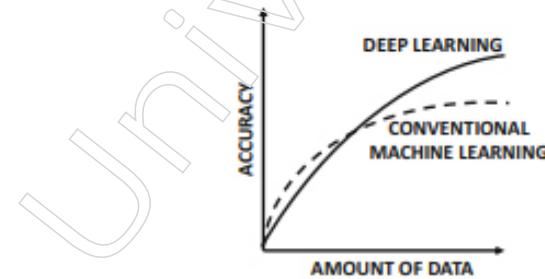


Figure: An example of comparing the precision of a big neural network and a standard machine learning algorithm. When there is enough data and computer power, deep learning becomes more appealing than traditional methods. There has been a “Cambrian explosion” in the use of deep learning technology in recent years as a result of an increase in computer power and data accessibility.

### 5.2.1 The Neuron Model:

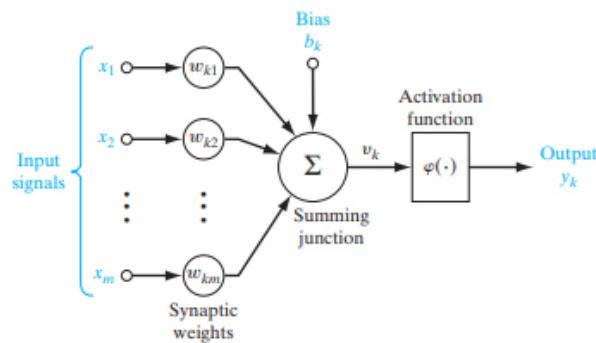


Figure : Nonlinear model of a neuron, labeled k.

A neuron is an information-processing cell that is essential to a neural network's functionality. The model of a neuron, which serves as the inspiration for creating a large family of neural, is shown in the image below. Here, we list the neural model's fundamental three components:

A compilation of synapses, also referred to as connecting points, wherein each synapse possesses a unique weight or strength. At the input of synapse  $j$ , which is linked to neuron  $k$ , the signal  $x_j$  is multiplied by the synaptic weight  $w_{kj}$ . It is advisable to pay attention to the notation used for the subscripts of the synaptic weight  $w_{kj}$ .

The terminal region of the synapse, which is associated with the weight, is commonly known as the. The synaptic weight of an artificial neuron exhibits a range of values that can fluctuate between negative and positive, in contrast to the synaptic weight of a biological synapse in the brain.

- A linear combiner is a device that adds input signals and weights them according to the strength of the synapses in each neuron.
- a neuron's output amplitude is limited by an activation function. The output signal's permitted amplitude range is squashed (limited) to a specific value by the activation function, which is also known as a squashing function.
- Typically, the closed unit interval  $[0,1]$  or, alternatively,  $[-1,1]$  is used to represent the normalized amplitude range of a neuron's output.
- The externally applied bias, represented by  $b_k$ , is likewise included in the neural model of the previous Figure. Depending on whether it is positive or negative, the bias  $b_k$  has the effect of either raising or reducing the net input of the activation function.
- We may write the following two equations to represent the neuron  $k$  shown in the previous Figure in mathematical terms:

$$u_k = \sum_{j=1}^m w_{kj} x_j \quad y_k = \varphi(u_k + b_k)$$

The input signals are denoted as  $x_1, x_2, \dots$  and  $x_m$ , while the corresponding synaptic weights of neuron  $k$  are represented as  $w_{k1}, w_{k2}, \dots$  and  $w_{km}$ . The output of the linear combiner resulting from the input signals is denoted as  $u_k$  and  $b_k$  represents the bias. The activation function is denoted as  $(\cdot)$  and  $y_k$  represents the output signal of the neuron. In the depicted model illustrated in the aforementioned Figure, the utilisation of bias  $b_k$  results in the application of an affine transformation on the output  $u_k$  of the linear combiner, as exemplified by:

$$v_k = u_k + b_k$$

In particular, the relationship between the induced local field, or activation potential,  $v_k$  of neuron  $k$  and the linear combiner output  $u_k$  is altered in the way depicted in Figure below and these two terms will now be used interchangeably. This depends on whether the bias  $b_k$  is positive or negative. Note that the graph of  $v_k$  versus  $u_k$  no longer crosses through the origin as a result of this affine change.

## Notes

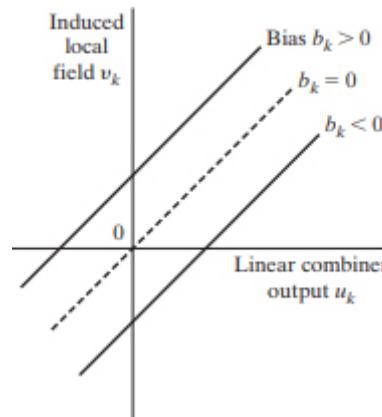


Figure : Affine transformation produced by the presence of a bias; note that  $v_k = b_k$  at  $u_k = 0$ .

An external parameter of neuron k is the bias ( $b_k$ ). We can explain its existence using the equation above. The combination of the two equations can be expressed equivalently as follows:

$$v_k = \sum_{j=0}^m w_{kj} x_j \quad y_k = \varphi(v_k)$$

As a result, we can reformulate the neuron k model, as seen in Figure below. The influence of the bias is taken into consideration in this figure in two ways:

- introducing a fresh input signal fixed at +1 and
- a fresh synaptic weight corresponding to the bias  $b_k$ . Despite the visual differences between the two Figure models, they are technically equal.

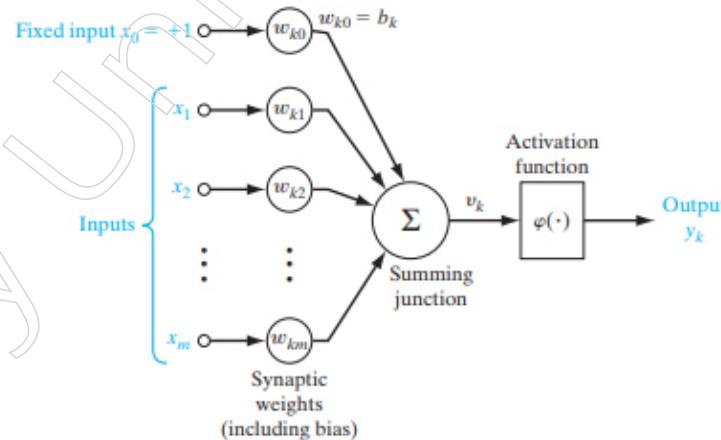


Figure: Another nonlinear model of a neuron; accounts for the bias  $b_k$  .

**Types of Activation Function:** The output of a neuron is defined by the activation function, indicated as ( $v$ ), in terms of the induced local field  $v$ . Following, we define two fundamental categories of activation functions:

**Threshold Function:** Figure following describes this type of activation function.,

$$\varphi(v) = \begin{cases} 1 & \text{if } v \geq 0 \\ 0 & \text{if } v < 0 \end{cases}$$

we have:

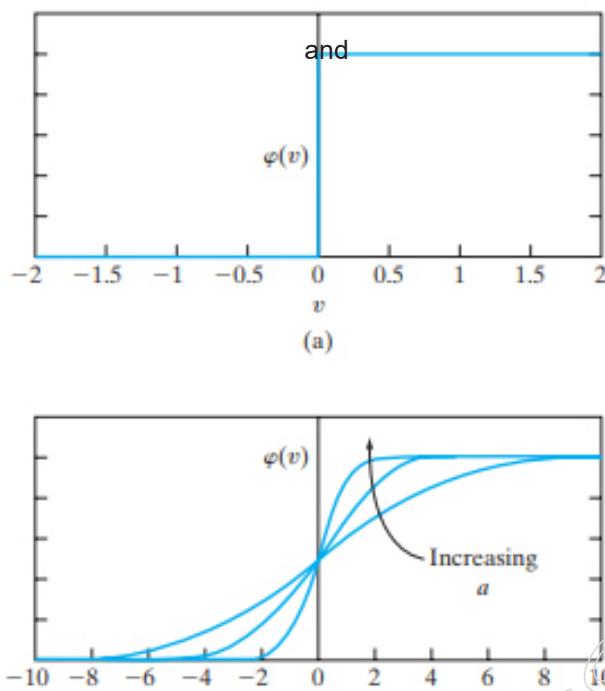


Figure : (a) Threshold function. (b) Sigmoid function for varying slope parameter a.

Such a neuron is known as the McCulloch-Pitts model in neural computation in honor of the groundbreaking work carried out by McCulloch and Pitts. In this paradigm, a neuron's output has a value of 1 if its induced local field is nonnegative and a value of 0 otherwise. The McCulloch-Pitts model's all-or-none attribute is described by this claim.

**Sigmoid Function:** The most typical type of activation function utilized in the creation of neural networks is the sigmoid function, whose graph is "S"-shaped. It is described as a strictly rising function with a delicately balanced linear and nonlinear behavior. The logistic function is an illustration of a sigmoid function and is defined as follows:

$$\varphi(v) = \frac{1}{1 + \exp(-av)}$$

where  $a$  is the sigmoid function's slope parameter. We can get sigmoid functions with various slopes by changing the parameter  $a$ , as shown in Fig. below. In actuality,  $a/4$  is the slope at the origin. The sigmoid function reduces to a threshold function in the limit as the slope parameter approaches infinity. In contrast to threshold functions, which only consider values between 0 and 1, sigmoid functions consider a continuous range of values between 0 and 1.

In the sense that its input-output behavior is exactly defined for all inputs, the neural model illustrated in the aforementioned Fig. is deterministic. It is preferable to use a stochastic neural model as the foundation for the analysis in some neural network applications.

The activation function of the McCulloch-Pitts model is given a probabilistic interpretation using an analytical trace table approach. A neuron is only allowed to exist in one of two states, let's say +1 or -1. A neuron's decision to fire—that is, to change from a "off" state to a "on" state—is probabilistic. Let  $x$  represent the neuron's current state and  $P(v)$  represent the likelihood that it will fire, with  $v$  standing for the neuron's induced local field. We may then type:

**Notes**

$$x = \begin{cases} +1 & \text{with probability } P(v) \\ -1 & \text{with probability } 1 - P(v) \end{cases}$$

The sigmoid function is a common option for  $P(v)$ :

where  $T$  is a pseudotemperature that is utilized to regulate the noise level and, consequently, the firing uncertainty. However, it is crucial to understand that  $T$  is not the actual temperature of a neural network, whether it be a biological or synthetic neural network. As was already mentioned, we should instead consider  $T$  to be merely a variable that regulates the thermal fluctuations that represent the effects of synaptic noise. The stochastic neuron given by the two equations above simplifies to the McCulloch-Pitts model, which is a noiseless (i.e., deterministic) form, when.

### 5.2.2 The Neural Network and Its Utility in Modelling and Solving Problems

**The Neural Network:** The understanding that the human brain computes entirely differently from the traditional digital computer has driven research on artificial neural networks, sometimes known as “neural networks,” from their beginnings. The brain is a very sophisticated, parallel, nonlinear computer (information processing system). It has the ability to organize the neurons that make up its structural components such that it can carry out some computations (including pattern recognition, perception and motor control) much more quickly than the current fastest digital computer. Take human vision as an example, which is an example of an information-processing activity.

The visual system’s job is to give us a picture of our surroundings and, more importantly, to give us the knowledge we need to interact with that environment. To be more precise, activities of considerably lower complexity take a very long time on a powerful computer, whereas the brain routinely completes perceptual identification tests (e.g., recognizing a familiar face embedded in an unfamiliar scene) in about 100-200 ms.

Another illustration would be a bat’s sonar. An active echolocation system is sonar. Bat sonar not only communicates information about the distance to a target (such as a flying bug), but also about the relative velocity, size and size of numerous elements on the target, as well as the azimuth and elevation of the target. A plum-sized brain performs the intricate neural computations required to retrieve all of this information from the target echo. An engineer working on radar or sonar would be jealous of how easily and successfully an echolocating bat can pursue and catch its prey.

So how does a bat’s brain or a human’s brain accomplish this? A brain has significant structure from birth and the capacity to develop its own set of behavioral guidelines through what we typically refer to as “experience.” In fact, experience is acquired over time. The human brain is hardwired in large part during the first two years after birth, but development continues even after that point.

A massively parallel distributed processor made up of simple processing units, known as a neural network, has a built-in inclination for storing and making use of experiencing information.

**It resembles the brain in two respects:**

- Through a process of learning, the network picks up knowledge from its surroundings.
- The learned information is stored in synaptic weights, which are the strengths of interneurons.

A learning algorithm refers to the computational procedure employed to facilitate the execution of the learning process. The intention of this process is to systematically modify the synaptic weights of the network in order to attain a predetermined design objective. The prevailing methodology employed in the creation of neural networks involves the manipulation of synaptic weights. This approach closely aligns with the principles of linear adaptive filter theory, a well recognised and effectively employed framework across various academic domains. The potential for neural networks to modify their own design arises from the capacity of neurons in the human brain to undergo cell death and the formation of new synaptic connections.

#### Benefits of Neural Networks:

It is clear that a neural network gets its processing power from two sources: first, its highly parallel distributed topology and second, its capacity for learning and generalization. Generalization is the neural network's ability to generate plausible results for inputs that weren't present during training (learning). With the help of these two information-processing abilities, neural networks are able to identify reliable approximations to difficult complex (large-scale) problems. However, in real-world applications, neural networks cannot deliver the answer on their own. In place of that, a consistent system engineering methodology needs to incorporate them. In particular, an interesting difficult problem is broken down into a number of relatively easy tasks and neural networks are given a subset of those tasks that are compatible with their natural skills. But it's vital to understand that before we can create a computer architecture that resembles the human brain, we still have a long way to go (if ever).

#### Neural networks offer the following useful properties and capabilities:

**Nonlinearity:** Linear and nonlinear artificial neurons can also be implemented. A neural network possesses inherent nonlinearity due to its composition of interconnected nonlinear neurons. Moreover, the characteristic of nonlinearity in this context is distinct in that it is distributed across the entire network. The feature of nonlinearity holds significant importance, particularly when the input signal (e.g., a speech signal) is inherently nonlinear due to the physical mechanism creating it.

**Input–Output Mapping:** A common learning paradigm known as supervised learning or learning with a teacher includes changing the synaptic weights of a neural network using a set of labeled training examples, or task examples. A distinct input signal and a corresponding desired (target) response make up each sample. In order to minimize the difference between the desired response and the actual response of the network produced by the input signal in accordance with an appropriate statistical criterion, the network's synaptic weights (free parameters) are modified. This is done by presenting the network with an example chosen at random from the set. Until the network reaches a stable state where there are no more noticeable changes in the synaptic weights, the network training is repeated for numerous examples in the set.

Training examples that have already been used may be used again, but in a new sequence. As a result, the network builds an input-output mapping specific to the given challenge in order to learn from examples. A similar strategy is similar to the study of nonparametric statistical inference, a field of statistics that focuses on model-free estimates, or, from a biological perspective, Tabula rasa education. In this context, the term "nonparametric" refers to the absence of any prior assumptions regarding a statistical model or the input data. Take a look at a pattern classification task as an illustration. In this challenge, you must assign a physical object or event represented by an input signal to one of several predetermined categories (classes). To "estimate"

## Notes

arbitrary decision boundaries without using a probabilistic distribution model in the input signal space for the pattern-classification task is the aim of a nonparametric approach to this problem. A similar perspective is implicitly adopted by the supervised learning paradigm, which draws parallels between the input-output mapping performed by a neural network and nonparametric statistical inference.

**Adaptivity:** In neural networks, synaptic weights can automatically change in response to environmental changes. For example, a neural network can be swiftly retrained to adapt to minor changes in the environmental conditions that it operates in after being trained to function in a specific context. When working in a nonstationary environment (one whose statistics change over time), a neural network may be taught to modify its synaptic weights in real time.

A neural network is a useful tool for adaptive pattern classification, adaptive signal processing and adaptive control due to its adjustable capabilities and natural architecture for pattern classification, signal processing and control applications.

As a general rule, it can be asserted that a system's performance is expected to be more dependable when the system is required to operate in a nonstationary environment the more adaptable the system is built to be while still retaining stability. But it's important to remember that adaptivity doesn't always lead to robustness; in fact, it might even have the opposite effect.

For example, a quick-changing adaptive system may be more prone to respond to fictional disruptions, which could cause a sudden fall in system performance. To fully use adaptivity, the major temporal constants of the system must be long enough for it to ignore spurious disturbances and short enough for it to respond to significant environmental changes.

**Evidential Response:** A neural network can be created to provide information regarding the confidence in the judgment made as well as which specific pattern to choose when it comes to pattern classification. If ambiguous patterns do appear, this later knowledge may be used to eliminate them and boost the network's classification performance.

**Uniformity of Analysis and Design:** In essence, neural networks are versatile data processors. We say this in the sense that neural networks are applied in all domains using the same nomenclature. This characteristic can appear in various ways:

- All neural networks contain neurons in one way or another as a basic component.
- This similarity enables the sharing of theories and learning algorithms across many neural network applications.
- It is possible to create modular networks by seamlessly integrating components.

**VLSI Implementability:** A neural network has the potential to be quick for the calculation of some tasks due to its massively parallel nature. A neural network is highly suited for use with very-large-scale-integrated (VLSI) technology because of the same characteristic. The ability to capture extremely complicated activity in a highly hierarchical manner is one of the advantages of VLSI.

**Contextual Information:** A neural network's physical composition and level of activation serve as a representation of knowledge. The overall activity of every neuron in the network has the capacity to influence every individual neuron. As a result, a neural network handles contextual information organically.

**Fault Tolerance:** In terms of resilient computation, a hardware-implemented neural network has the potential to be inherently fault tolerant, meaning that performance

declines smoothly under challenging operating conditions. For instance, recall of a stored pattern is of worse quality if a neuron or one of its connecting pathways is destroyed. However, because the network's data is spread, significant harm must be done before the network's general responsiveness is substantially compromised. Therefore, a neural network in theory displays a gentle performance decline as opposed to catastrophic collapse. Robust computation has some empirical support, but this support is typically unregulated. It could be required to make corrections when creating the algorithm to train the network in order to ensure that the neural network is, in fact, fault resistant.

**Neurobiological Analogy:** The brain, which serves as living proof that fault-tolerant parallel processing is not only physically possible but also quick and effective, serves as the inspiration for the creation of a neural network. (Artificial) neural networks are used as a research tool by neurobiologists to interpret neurobiological events. Engineers, on the other hand, look to neurobiology for novel solutions to challenges that are more complex than those based on traditional hardwired design methods. The following two examples, which are each representative of one of these two points of view:

- In comparison to neural network models based on recurrent networks, linear system models of the vestibulo-ocular reflex (VOR) are more straightforward. The oculomotor system includes the vestibulo-ocular reflex. By rotating the eyes counterclockwise to the head, VOR works to keep the retinal image (i.e., visual) stable. Premotor neurons in the vestibular nucleus, which receive and interpret head rotation information from vestibular sensory neurons, mediate the VOR by sending the results to the motor neurons in the eye muscles. Because both the VOR's input (head rotation) and output (eye rotation) can be precisely set, it is highly suited for modelling. It is also a very straightforward reflex and the constituent neurons' neurophysiological characteristics have been extensively discussed. The vestibular nuclei's premotor neurons are the most intricate and fascinating of the three different neuronal kinds. The VOR has previously been modeled using control theory and lumped, linear system characteristics. Although these models helped to explain some of the VOR's general characteristics, they provided little information about the characteristics of the individual neurons that make up the VOR. Neural network modelling has significantly improved this scenario. Many of the static, dynamic, nonlinear and dispersed elements of signal processing by the neurons that mediate the VOR, especially the vestibular nuclei neurons, can be reproduced and explained using recurrent network models of VOR.
- More than any other area of the brain, the retina is where we start to make connections between the initial cerebral images and the outside world, which is represented by a visual sense, as well as between that sense's actual physical image projected onto a variety of receptors. The retina is a delicate layer of neural tissue that covers the back of the eye. It is the job of the retina to transform an optical image into a neural image that may then be sent down the optic nerve to a variety of locations for additional analysis. As seen by the synaptic structure of the retina, this is a difficult undertaking. The conversion of an optical image into a neural image occurs in the retinas of all vertebrates in three steps: phototransduction by a layer of receptor neurons; transmission of the resulting signals (produced in response to light) by chemical synapses to a layer of bipolar cells; and transmission of these signals, again by chemical synapses, to output neurons known as ganglion cells.

In synaptic stages, there exist specialised neurons known as horizontal cells and amacrine cells, which are laterally coupled. The primary role of these neurons is to modulate the transmission of synaptic signals between different layers. Inter-plexiform

## Notes

cells serve as centrifugal elements responsible for conveying impulses from the inner synaptic layer to the outer synaptic layer. Scientists have successfully developed electrical chips that have a striking resemblance to the structural characteristics of the human retina. The term “neuromorphic integrated circuits” refers to electronic chips.

The neuromorphic image sensor has an array of photoreceptors interconnected with analogue circuitry, whereby each pixel is represented. The artificial system exhibits characteristics similar to the human retina, since it possesses the ability to perceive edges, respond to small variations in luminosity and identify patterns of motion. One notable advantage of the neurobiological analogy, as exemplified by neuromorphic integrated circuits, is its capacity to instill optimism and conviction, while also offering some degree of substantiation, regarding the potential beneficial effects of applying knowledge about neurobiological structures to the advancement of electronics and VLSI technology for neural network development.

### 5.2.3 Connect to the Biological Motivations and Parallelism

**Evolution and Parallel Processing:** Evolution has evolved creatures in the biological world to maximize their chances of survival and procreation. In biology, parallelism is demonstrated by the simultaneous occurrence of several processes within an organism. For instance, the body's cells do a variety of tasks simultaneously, including digestion, breathing and circulation. This simultaneous processing improves flexibility and efficiency.

**Neural Networks and Brain Parallelism:** One of the most intricate biological structures, the human brain, demonstrates astounding parallelism. Neurons process information simultaneously, enabling quick perception, thought and action. In the field of machine learning, artificial neural networks were created as a result of this parallelism notion. These networks enable computers to carry out tasks like image recognition and natural language processing by simulating the connectivity of neurons.

**Genetic Parallelism and Diversity:** The parallel evolution of various species or populations in response to various environmental conditions gives rise to genetic diversity. This parallelism is essential for the survival of all life on Earth because it enables species to adapt to shifting environmental conditions. Similar to this, parallelism is used in computers to handle difficult operations by breaking them down into smaller, more manageable components that may be performed concurrently, improving overall efficiency.

**Ecosystems and Distributed Computing:** The effectiveness of parallelism in nature is shown by ecosystems. In parallel, several species interact and contribute to the efficiency of the ecosystem. Breaking down workloads into smaller subtasks and executing them on different machines is known as distributed computing, a parallel processing strategy in technology. This is similar to how several species work together to maintain the balance of an ecosystem.

**Biological Hierarchies and Parallel Architectures:** From cells to tissues to organs to entire animals, biological systems frequently display hierarchical order. Parallelism is made easier by this hierarchical structure because different levels can carry out separate tasks at the same time. Similar to this, parallel architectures in computing are created with numerous layers, each of which increases the system's overall processing capacity and effectiveness.

**Survival Strategies and Task Parallelism:** Different parallelism-based survival strategies have been created by biological creatures. Animals, for example, may find food, fend off predators and reproduce all at once. Task parallelism in computing refers

to the division of a task into smaller, concurrently executable tasks. This method boosts efficiency and speed in a manner similar to how nature multitasks.

**Notes**

### 5.2.4 Popular CNN (Convolutional Neural Network) Architectures:

#### Convolutional Neural Network:

Artificial intelligence is rapidly narrowing the disparity in skills between robots and humans. A multitude of scholars and non-professionals are engaged in diverse facets of the field of artificial intelligence with the aim of developing remarkable innovations. The discipline of computer vision is a remarkable area of study. The primary goal of computer vision is to emulate the perceptual abilities of humans in artificial systems. Examples of well-known computer vision applications include image detection, image tagging, image recognition, image classification, image analysis, video analysis and natural language processing.

The emergence and advancement of deep learning techniques in the field of computer vision has garnered significant interest and engagement from numerous researchers over the course of several years. CNN is the predominant choice for constructing the majority of computer vision algorithms. Convolutional neural networks employ deep learning techniques to assign learnable biases and weights to various objects within an input image, enabling their differentiation.

CNN requires less preprocessing than other techniques. The best learning algorithm for understanding visual content is therefore a CNN. Additionally, it has proven to be exceptionally good at segmenting, retrieving and classifying images. People outside of academia are now interested in CNN because of its success. Microsoft, Google, AT&T, NEC and Facebook are just a few of the businesses working to expand CNN architecture.

Additionally, they possess research teams that are actively engaged in investigating state-of-the-art convolutional neural network (CNN) architectures. Currently, the majority of leading participants in image processing and computer vision competitions predominantly utilise models based on deep Convolutional Neural Networks (CNNs). There exist numerous modifications of the fundamental Convolutional Neural Network (CNN) design. An examination of the introduction to CNN, the evolution of CNN over time, various design aspects employed by CNN and an architectural analysis of each form of CNN, including an assessment of their respective strengths and weaknesses.

This is the first review that almost covers every aspect of CNN for computer vision, including its background, CNN architecture designs, advantages and disadvantages, applications and the work that needs to be done in the future.

- This review helps readers decide wisely on their future research in the field of CNN for computer vision;
- The limits of the current CNN architectures are included in this manuscript, which inspires the construction of a new architecture. Additionally, it outlines the benefits and drawbacks of practically all common CNN variations;
- This survey paper's fascinating section classifies the CNN architecture into eight types based on its implementation criteria;
- In addition, other CNN applications are described so that readers might use CNN in areas other than computer vision;
- It offers a detailed overview of upcoming research directions in the field of CNN for computer vision.

## Notes

This research paper is organized with the first portion providing a thorough understanding of CNN. The parallels between CNN and the visual brain of the ape are then explained.

**CNN components:** The following figure displays different CNN elements. It's crucial to comprehend the many CNN components and their uses in order to learn about the developments in CNN architecture:

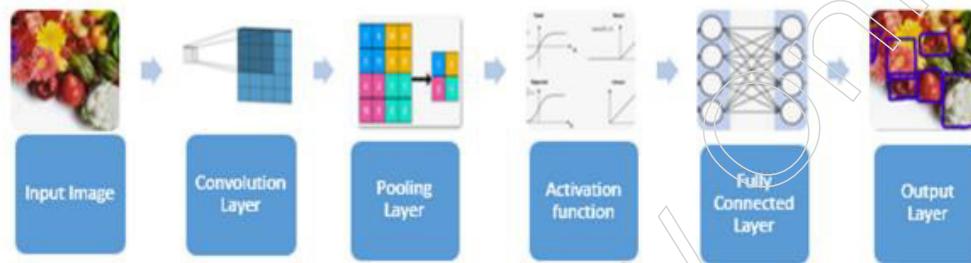


Figure: various CNN components.

**Input Image:** The building blocks of a computer image are called pixels. They serve as the binary representation of the visual data. The digital image is composed of a series of pixels with a range of 0-255 that are arranged in a matrix-like pattern. The brightness and hue of each pixel are specified by its pixel value. In the first second after seeing an image, human brains digest a huge amount of data. In order to span the whole visual field, each neuron in the human brain has its own receptive field and is connected to other neurons. Each neuron in the biological vision system responds to stimuli in the receptive field, which makes up a very small percentage of the visual field. Similar to this, each neuron in CNN only processes data inside its specific receptive area. Before moving on to more complicated patterns, like faces and objects, the CNN layers are designed to recognize simpler patterns first, such lines and curves. Therefore, it is conceivable to assert that utilizing a CNN might provide machines eyesight.

**Convolution Layer:** In the CNN design, the convolution layer is a crucial layer. As illustrated in Figure below, it utilizes a  $3 \times 3$  or  $5 \times 5$  filter and accepts images as input.

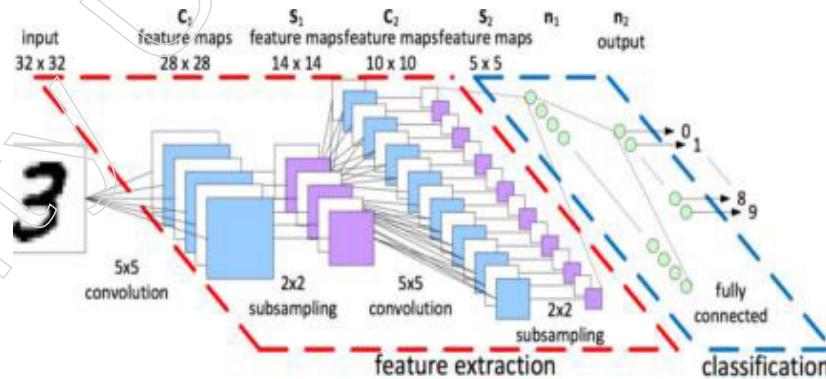


Figure: Convolution layer.

The input image, which is shown in blue in the figure below, is covered by the green filter one pixel at a time, beginning at the top left. The filter multiplies its values by the values that overlap with the image as it passes over it, then adds all of the values to produce a single value output for each overlap until the entire image has been visited.

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

1	0	1
0	1	0
1	0	1

Figure: Input and filter image

In instances where an image possesses many channels, such as the RGB (red, green and blue) model, it is observed that the kernel's depth aligns with the depth of the input image. The process of matrix multiplication is performed on the stacks  $K_n$  and  $I_n$  ( $[K_1, I_1]$ ,  $[K_2, I_2]$ ,  $[K_3, I_3]$ ) as depicted in the accompanying diagram. The outcomes of this operation are subsequently merged with the bias term to produce a dense channel with a single depth.

Each every neuron inside the output matrix possesses an underlying receptive field. The initial ConvLayer is responsible for capturing low-level information, like gradient direction, edges, colour and similar characteristics. Through the incorporation of multiple layers, the architectural structure adapts to the overarching characteristics, resulting in a neural network that possesses a comprehensive comprehension of the images inside the dataset. The figures depict the sequential stages of the convolution process.

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

Figure: Calculation of filter slides over input image.

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

$$1 \times 1 + 1 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 1 + 1 \times 0 + 0 \times 1 + 0 \times 0 + 1 \times 1$$

Figure: First step of convolution

**Feature Extraction:** CNN is renowned for its capacity to automatically extract attributes. The RGB picture matrix calculation is shown in the figure below. CNN typically uses padding to prevent the size of the feature maps from decreasing at each layer, which is not what is desired. The process results in two different kinds of results:

- a kind where the input's dimensions are reduced relative to the twisted feature;
- a type where the dimensionality is either kept or improved, rather than diminished. To accomplish this purpose, padding is used.

## Notes

## Notes

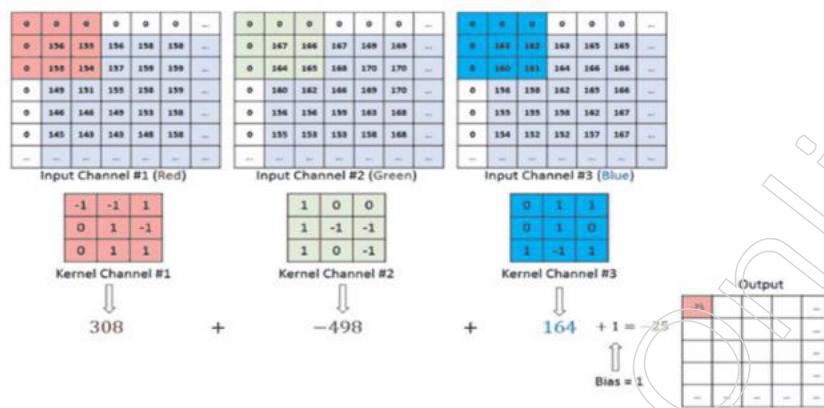


Figure :Matrix calculation.

For instance, the complex matrix is found to have dimensions of  $5 \times 5 \times 1$  when the  $5 \times 5 \times 1$  picture is reinforced into a  $7 \times 7 \times 1$  image and then applied to the  $3 \times 3 \times 1$  kernel over it, as shown in Figure below. It signifies that the input and output images share the same dimensions and padding. If the same process is carried out without padding, the output could include an image with smaller dimensions. Consequently, a  $5 \times 5 \times 1$  image will change into a  $3 \times 3 \times 1$  image.

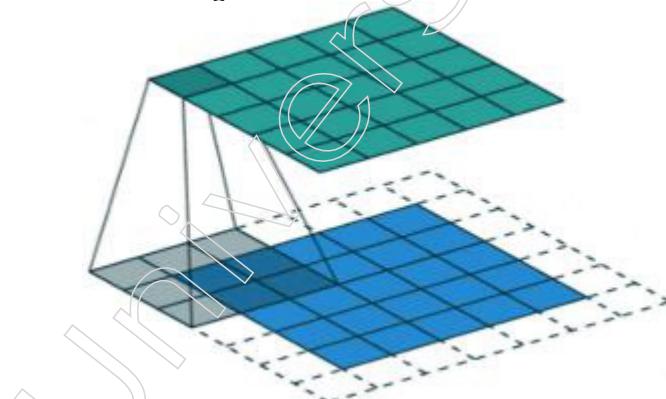


Figure: Padding.

During the forwarding pass, the kernel traverses the picture's width and height. It produces a graphic illustration of the relevant receptive region. A two-dimensional representation of the image that displays the kernel's response at each spatial position of the image is created as a result: an activation map. The kernel's size when it slips is measured in strides. Assume the input picture has the dimensions  $W \times W \times D$ . The output volume can be determined using the following formula if the quantity of kernels with spatial dimensions of  $F$ , stride  $S$  and padding  $P$  is unknown:

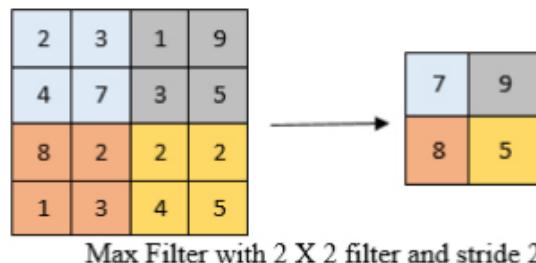
$$W_{out} = \frac{W - F + 2P}{S} + 1$$

This will result in a result of size  $W_{out} \times W_{out} \times D_{out}$ .

**Pooling Layer:** Following the acquisition of the feature maps, a pooling (subsampling) layer must be added to CNN with a convolution layer. The pooling layer's function is to reduce the spatial size of the convolved feature. The amount of computing resources needed to process the data is decreased as a result of the dimensionality reduction. This helps to maintain the model's practical training and the extraction of leading characteristics that are rotational and positional invariant. Pooling shortens the

training period and avoids over-fitting. Maximum pooling and average pooling are the two types of pooling.

- **Maximum Pooling:** The input of the pooling layer is a tensor. In the context of maximum pooling, as depicted in the accompanying Figure, a kernel with dimensions  $n \times n$  (specifically  $2 \times 2$  in the aforementioned instance) is traversed across the matrix. Subsequently, the maximum value inside the kernel is identified and placed in the corresponding position within the output matrix.



- **Average Pooling:** In the process of average pooling, a kernel with dimensions  $n \times n$  is systematically moved across the input matrix. At each position, the average of all the values within the kernel is computed and this average is then placed in the corresponding position of the output matrix. The aforementioned process is iterated for each individual channel within the input tensor. The output tensor has been obtained as a consequence. It is imperative to bear in mind that throughout the process of pooling, the dimensions of the image are reduced in terms of height and width, while the number of channels, representing the depth, remains unchanged. The pooling layer is responsible for computing a summary statistic of the neighbouring outputs, which is then used to occasionally replace the network output.
- As a result, it helps to decrease the spatial dimension of the representation, which lowers the amount of computation and weights needed. Each slice of the representation undergoes the pooling process separately. According to Figure below, the rectangle neighborhood average, the rectangle neighborhood L2 norm and a weighted average based on the distance from the central pixel are all examples of pooling functions. Maximum pooling, which presents the neighborhood's most noteworthy output, is the most popular technique.

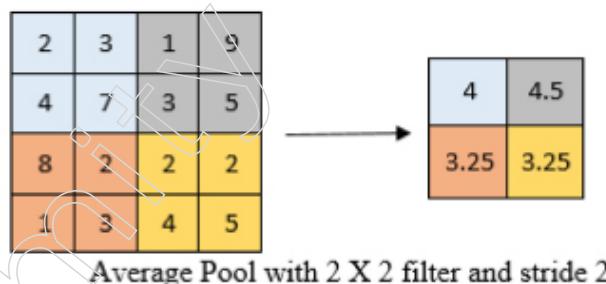


Figure: Average pooling

**Nonlinearity Layer (Activation Function):** In CNN layers, the activation function is crucial. An additional mathematical function known as an activation function receives the filter's output. Rectified linear unit, or ReLu, is the most often used activation function in CNN feature extraction. Utilizing the activation function is mostly done to interpret the output of neural networks, such as yes or no. Depending on the activation function, the output values are mapped between -1 and 1 or 0 and 1, etc. There are two different categories for the activation functions:

## Notes

## Notes

- **Linear Activation Function:** This uses the function  $F(x) = CY$  as the linear activation function. The output signal is proportional to the input after being multiplied by constant  $c$  (the weight of each neuron). Given that it simply provides a yes or no response and not a range of options, the linear function may be preferable to the step function.
- **Non-linear Activation:** Functions of Non-linear Activation Today's neural networks employ non-linear activation functions. In order to learn and model complex data, such as photos, videos, sounds and non-linear or high-dimensional data sets, the model must be able to construct intricate mappings between the network's inputs and outputs.

**Fully Connected Layer:** As seen in the figure below, a completely linked layer is nothing more than a feed-forward neural network. The very bottom layers of the network are where you'll find fully connected layers. The output layer of the last pooling or convolutional layer is flattened before being supplied as input to a fully connected layer. When the output is flattened, all of the values that were obtained after the last pooling or convolutional layer are unrolled into a vector (3D matrix).

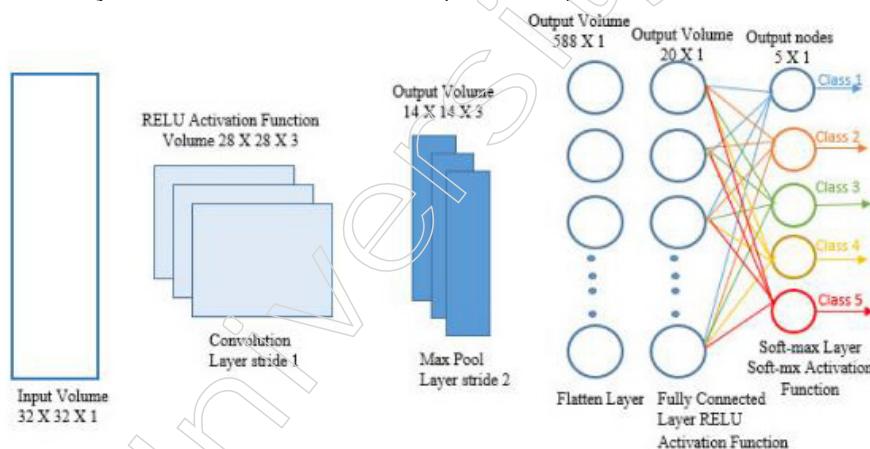


Figure: Fully connected Layer

A quick method for learning nonlinear combinations of high-level features represented by the output of the convolutional layer is to add an FC layer. The FC layer is learning a potentially nonlinear function in that area.

**Architectural Evolution of CNNs:** The many architectural categories of CNN variants are described in the following figure. These categories are all explained in detail in this section.

**Spatial Exploitation-Based CNNs:** Biases, weights, the number of layers, neurons, activation function, stride, filter size, learning rate and other factors are among the many variables in a CNN model. Convolutional procedures take into account the vicinity (locality) of input pixels, allowing for the investigation of various correlation levels using various filter sizes. Different filter sizes correspond to various granularity levels; typically, small filters collect fine-grained information, while large filters retrieve coarse-grained information. In order to improve performance, researchers started using spatial filters in the early 2000s. It was discovered that a spatial filter and network learning are related. During this time, numerous research showed that CNN might perform better on coarse and fine-grained details by modifying filters.

**CNN Based on Depth:** The primary concept underlying deep CNN design is that the network may successfully approximate the goal function with the aid of more

mappings (nonlinear) and more sophisticated feature hierarchies. The depth of the network has been a crucial factor in supervised training. Deep networks are more adept at representing particular function classes than shallow systems. A theorem known as “universal approximation” was first forth in 2001. It described how any function might be approximated by a single hidden layer. However, this results in an unrealistically large number of neurons and an exponentially high cost of computation. In 2011, Bengio and Delalleau proposed that deeper networks could preserve the network’s theatrical impact for less money. Bengio showed this empirically in 2013 and concluded that deep networks are computationally more efficient for complex activities. In the ILSVRC-2014 competition, VGG and Inception fared the best, supporting the idea that depth is a crucial factor in controlling a network’s capacity for learning.

**CNNs with Multiple Paths:** Deep CNNs are frequently effective at difficult tasks. Performance concerns, explosion problems, or gradient fading might occasionally affect them and are caused by increasing the depth rather than overfitting. An rise in test error and training error is a result of the vanishing gradient problem. For deep learning networks, the theory of cross-layer connectivity or multi-path was put out. By avoiding some intermediate levels, shortcut connections or multiple paths can connect one layer to another analytically, enabling a tailored information flow between the layers. Cross-layer connection is used to divide the network into the various components. The vanishing gradient issue is resolved by these paths by extending the gradient to lower layers.

**Feature-Map Exploitation Based CNNs:** Due to its ability to perform automatic feature extraction and hierarchical learning, CNN has been a popular choice for MV tasks. The choice of features has a significant impact on how well classification, segmentation and detection modules work. CNN uses a kernel, often referred to as a mask and associated weights to dynamically choose features. Additionally, numerous stages of feature extraction are carried out, allowing for a variety of features (referred to as feature maps or channels in CNN). Some of the feature maps, however, don’t perform well or at all in object discrimination. Over-fitting of the network may result from excessive feature sets producing a noise effect. This suggests that selection of feature maps can be quite important in enhancing network generalization in addition to network engineering.

**Multi-Connection Depending on the Width:** The primary focus of CNN developments from 2012 to 2015 was on maximizing the depth and effectiveness of connections for network regularization. In 2019, Kawaguchi revealed that the network’s width is also quite important. This suggests that width, in addition to depth, is crucial when creating learning philosophies. It is demonstrated that in order to maintain a universal approximation property while also gaining in depth, neural networks with ReLU activation functions must be wide enough. The failure of several layers to learn useful features is a serious problem with deep neural network topologies. Even though adding more layers (increasing depth) may enable the learning of different feature representations, this does not always improve the NN’s capacity for learning. Furthermore, if a deep network’s maximum width does not exceed its input dimension, it cannot arbitrarily approximate a class of continuous functions on

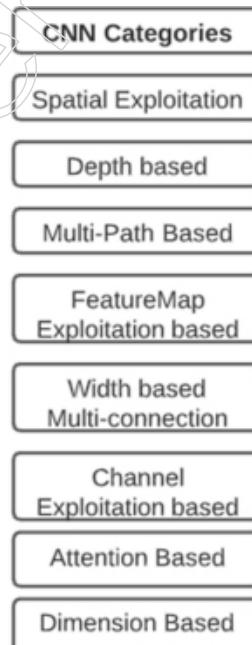


Figure: CNN variants categories.

## Notes

a small set. To overcome this problem, the research's emphasis shifted from deep and narrow designs to wide and thin architectures.

**Exploitation-Based Feature-Map (ChannelFMap) CNNs:** As a result of its ability to perform automatic feature extraction and hierarchical learning, CNN has drawn a lot of interest from people working on computer vision issues. The choice of features has a significant impact on how well classification, segmentation and detection modules work. CNN chooses features in a dynamic manner by varying the weights associated with a kernel, sometimes referred to as a mask. Additionally, CNN uses numerous feature extraction processes to mine different kinds of characteristics. Some feature maps, however, are either insignificant or have no impact on object discrimination. Massive feature sets could produce a noise effect that overfits the network. This suggests that, in addition to network engineering, the choice of feature mappings can be crucial in enhancing network generalization. In the literature, feature maps and channel terminology are sometimes used interchangeably.

**CNNs That Are Based on Attention:** Different levels of abstraction are crucial in defining the discrimination power of the NN. Learning involves using several abstraction hierarchies that are centered on properties important for picture localisation and recognition. The attention effect is what the visual system of the human is called. Any scene can be seen by a human by integrating quick glances at it and concentrating on context-relevant elements. This method enhances the visual structure capture process by concentrating on specific areas and taking into account several interpretations of objects at a given location. RNN and LSTM use an interpretation that is more or less similar. Attention modules are used as a progressive element in RNN and LSTM networks and new tasters are weighted according to how frequently they appeared in preceding rounds. Many researchers employ the convolutional neural network's attention principle to enhance representation and get around computational constraints. This idea of attention also helps CNN develop the ability to detect objects even against cluttered backgrounds and challenging circumstances.

**Dimension-Based CNN:** The traditional convolutions layer simultaneously encodes channel-wise and spatial information, but it is computationally costly. The development of separable (or depth-wise separable) convolutions, which independently encode spatial and channel-wise information using point-wise and depth-wise convolutions, increased the efficiency of conventional convolutions. The point-wise convolutions become a computational bottleneck since this factorization, despite being much more effective, exerts a heavy computational weight on them.

### 5.2.5 RNNs (Recurrent Neural Networks):

Recurrent Neural Networks (RNNs) are a particular sort of neural network design that are mostly employed for pattern recognition in a stream of input. Such information can include handwriting, genomes, language, or numerical time series that are frequently created in professional settings (such as stock markets or sensors). However, if those are respectively divided into a number of patches and handled as a sequence, they are likewise relevant to photographs. RNNs are used at a higher level in speech recognition, language modelling and text generation, image description and video tagging. The way information moves through the network distinguishes Recurrent Neural Networks from Feedforward Neural Networks, sometimes referred to as Multi-Layer Perceptrons (MLPs). While RNNs have cycles and broadcast data back into itself, feed forward networks just transmit data via the network. As a result, they are able to expand the capability of feed forward networks such that they take into account both the current input  $X_t$  and past inputs

X0:t1. The figure below provides a high-level visualization of this disparity. Keep in mind that one Hidden Layer block H is used to aggregate the option of having several hidden levels here. There are certainly other hidden layers that may be added to this block.

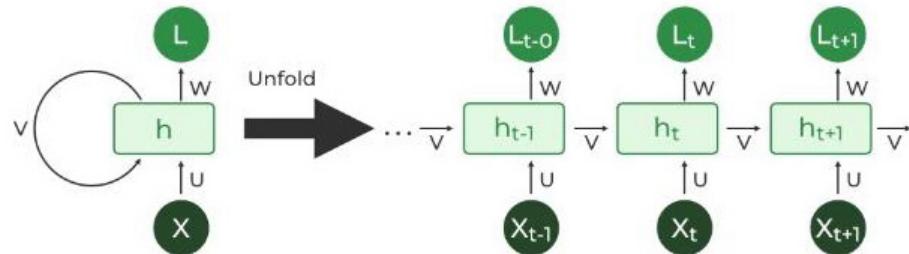


Figure: Recurrent neural network

The mathematical notation proposed in the aforementioned study can be employed to visually represent the process of transferring data from the preceding iteration to the concealed layer. In this context, the notation  $H_t \in \mathbb{R}^{n \times h}$  and  $X_t \in \mathbb{R}^{n \times d}$  is employed, where n represents the quantity of samples, d denotes the number of inputs per sample and h signifies the number of hidden units. In addition, we utilise a bias parameter  $b_h \in \mathbb{R}^{1 \times h}$ , a hidden-state-to-hidden-state matrix  $W_{hh} \in \mathbb{R}^{h \times h}$  and a weight matrix  $W_{xh} \in \mathbb{R}^{d \times h}$ . To facilitate the utilisation of gradients in the backpropagation process, the aforementioned data is subsequently fed into an activation function, commonly characterised by a logistic sigmoid or hyperbolic tangent (tanh) function. The combination of these notations yields the concealed variable, as described in the first equation and the resultant variable, as indicated in the second equation.

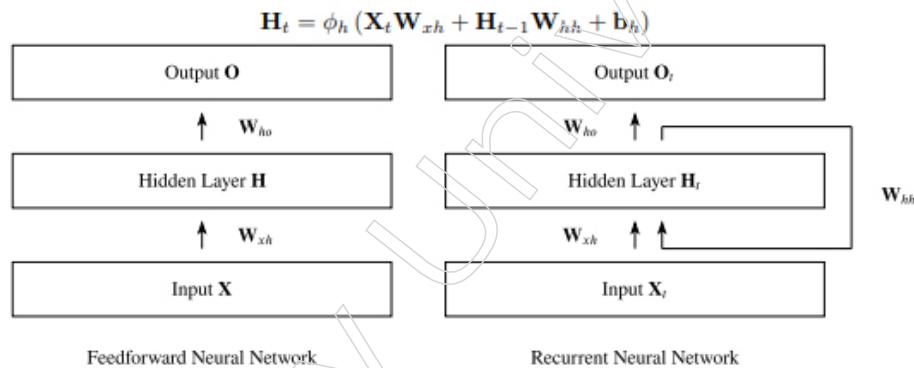


Figure: Visualisation of differences between Feedforward NNs und Recurrent NNs

$$O_t = \phi_o (H_t W_{ho} + b_o)$$

Since  $H_t$  includes  $H_{t-1}$  recursively and this process happens for each time step, the RNN contains traces of all hidden states that came before  $H_{t-1}$  in addition to  $H_{t-1}$ . We can readily observe the distinction we previously outlined if we contrast that notation for RNNs with a syntax comparable to that for feedforward neural networks. The computation for the hidden variable is shown in the equation below and the output variable is shown in the equation that follows.

$$H = \phi_h (X W_{xh} + b_h)$$

$$O = \phi_o (H W_{ho} + b_o)$$

**Architecture Of Recurrent Neural Network:** The input and output architecture of RNNs is identical to that of other deep neural network architectures. But there are variations in how information moves from input to output. In RNN, the weight across

Notes

## Notes

the network is constant, in contrast to Deep Neural Networks where we have separate weight matrices for each Dense network. For each input  $X_i$ , the state hidden state  $H_i$  is calculated. By use the formulas below:

$$\begin{aligned} h &= \sigma(UX + Wh_{-1} + B) \\ Y &= O(Vh + C) \text{ Hence} \\ Y &= f(X, h, W, U, V, B, C) \end{aligned}$$

In this case,  $S$  is the state matrix and element  $s_i$  represents the network's state at timestep  $i$ .  $W$ ,  $U$ ,  $V$ ,  $c$  and  $b$  are network parameters that are shared throughout timesteps. How RNN Function: The recurrent neural network consists of a single fixed activation function unit at each time step. The hidden state of each unit is denoted as its internal state. The concealed state at a specific time step represents the network's existing knowledge of the past. At each time step, the hidden state is modified to incorporate the network's evolving comprehension of the preceding events. The provided recurrence relation is utilised for updating the hidden state.

**The formula for calculating the current state:**

$$h_t = f(h_{t-1}, x_t)$$

**where:**  $h_t$  -> current state , $h_{t-1}$  -> previous state , $x_t$  -> input state

**Formula for applying Activation function(tanh):**

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

**where:**  $W_{hh}$  -> weight at recurrent neuron , $W_{xh}$  -> weight at input neuron.

**The formula for calculating output:**

$$y_t = W_{hy}h_t$$

**where :** $y_t$  -> output , $W_{hy}$  -> weight at output layer

**Backpropagation Through Time (BPTT):** Because the neural network in RNN is ordered, each variable is computed one at a time in a specific order, such as first  $h_1$ , then  $h_2$ , then  $h_3$  and so on. Therefore, we will sequentially perform backpropagation to each of these concealed temporal stages. Figure below:  $L(\theta)$ (loss function) is dependent on  $h_3$ , which is dependent on  $h_2$  and  $W$ , which is dependent on  $h_1$  and  $W$ , which is dependent on  $h_0$  and  $W$ , where  $h_0$  represents a constant beginning condition.

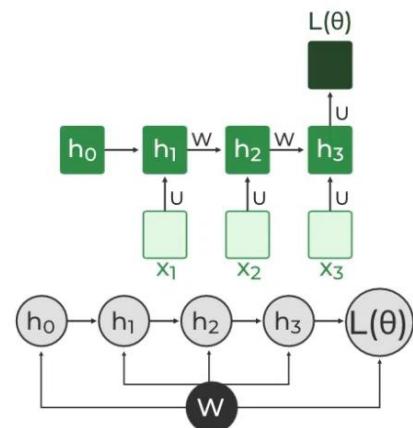


Figure: Backpropagation Through Time (BPTT) In RNN.

**Types Of RNN:** Based on the quantity of inputs and outputs in the network, there are four different types of RNNs.

**One to One:** This kind of RNN, commonly referred to as a “vanilla neural network,” functions similarly like any other straightforward neural network. There is only one input and one output in this neural network.

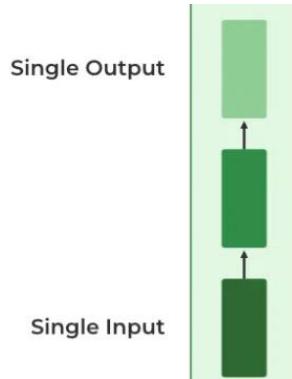
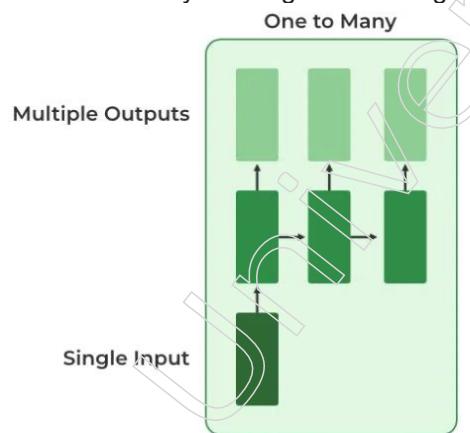


Figure:One to One RNN.

**One To Many:** This kind of RNN has a single input and numerous associated outputs. One of the most popular applications of this network is picture captioning, in which we anticipate a sentence with many words given an image.



**Many to One:** This kind of network only produces one output after receiving several inputs at various network states. Use of this kind of network is made for issues like sentimental analysis. When many words are provided as input, we merely forecast the sentiment of the sentence to be the output.

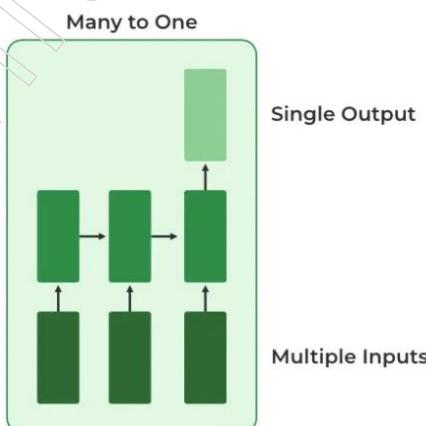


Figure:Many to One RNN

## Notes

**Variation Of Recurrent Neural Network (RNN):** Several new, more sophisticated RNNs have been developed to address issues like disappearing gradient and explosive gradient descent. Some of these include;

- **Bidirectional Neural Network (BiNN):** An adaptation of a recurrent neural network known as a bidirectional neural network (BiNN) combines the outputs of both directions to produce the input. In cases like NLP tasks and Time-series analysis issues, where the context of the input is more crucial, BiNN is helpful.
- **Long Short-Term Memory (LSTM):** Short-Term Long-Term The read-write-forget (RWF) concept describes how memory operates. Given an input of information, the network reads and writes the data that will be most valuable in forecasting the output before forgetting the data that won't be. Three new gates are added to the RNN to accomplish this. Only the chosen information is transmitted via the network in this way.

### 5.3 Neural Networks: Its Types and Applications

Artificial neural networks (ANNs) and simulation neural networks (SNNs) are alternative terms used to refer to neural networks, which are a subset of machine learning and serve as the foundational component of deep learning methodologies. The structure and nomenclature of its organisation are derived from the human brain, emulating the mechanisms by which actual neurons engage in communication. Computers can utilise this approach to develop an adaptive system that facilitates continuous progress by instructing them based on errors. Artificial neural networks are employed in order to tackle challenging tasks such as document summarization or facial recognition.

Neural networks provide the capability to effectively categorise and cluster data, so augmenting the existing data management and storage infrastructure with an additional layer of grouping and categorization. These algorithms aid in the process of data categorization by organising unclassified data into groups or clusters, using similarities observed among sample inputs, with the use of a labelled dataset for training purposes.

**Major Categories of Neural Networks:** Some of the main types of neural networks include the following:

**Classification:** Neural networks frequently perform well in classification tasks, which necessitate labeled datasets for supervised learning. For instance, while recognizing visual patterns in hundreds of photos, neural networks can assign labels fast and consistently. Through practice, they develop the ability to solve challenging, confusing challenges. The brain network develops an innate ability to recognize the most important factors. As a result, the data scientist is not obliged to offer attributes that would allow cats and dogs to be distinguished from one another.

**Sequence learning:** Data sequences are used as input or output in the machine learning category known as "sequence learning." Sequential learning can be done with text streams, audio files, video clips, measurements and more.

**Function approximation:** Function approximation is a strategy that uses previous or current observations from the domain to approximate an unknown underlying function. Artificial neural networks can learn to approximate a function.

#### 5.3.1 Perceptron:

The term "perceptron" is extensively employed in the fields of machine learning and artificial intelligence, enjoying widespread usage among practitioners. Acquiring knowledge in the fields of machine learning and deep learning, which encompass a

collection of weights, input values or scores and a threshold, constitutes the initial stage of this undertaking. The perceptron is a fundamental element within an artificial neural network. The Perceptron was originally developed by Frank Rosenblatt during the mid-20th century with the purpose of doing certain computations to discern the capabilities of input data or commercial intelligence.

A linear machine learning algorithm called the perceptron is utilized for supervised learning for different binary classifiers. With the help of this algorithm, neurons can learn new elements and process them one at a time while preparing. In-depth knowledge of the perceptron will be covered in this tutorial's section on "Perceptron in Machine Learning," along with a brief overview of its fundamental operations. Let's begin with a brief overview of the perceptron.

**What is the Perceptron model in Machine Learning:** For the supervised learning of various binary classification tasks, the perceptron machine learning algorithm is used. Moreover, a perceptron can be classified as an artificial neuron or unit within a neural network, which contributes to the field of business intelligence by facilitating the identification and analysis of particular input data calculations. The perceptron model is widely regarded as one of the most effective and straightforward types of artificial neural networks. Nevertheless, the system employs binary classifiers within a supervised learning framework. The structure of this system can be conceptualised as a single-layer neural network, consisting of four primary components: input values, weights and bias, net sum and an activation function.

**Basic Components of Perceptron:** The perceptron model, which is a binary classifier consisting of three key components, was developed by Mr. Frank Rosenblatt. The following items are enumerated below:

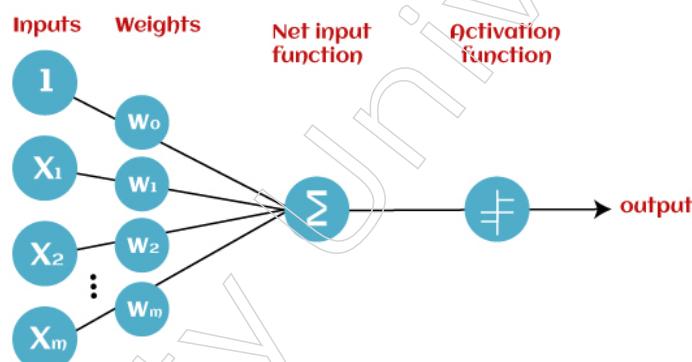


Figure: Basic Components of Perceptron

- **Input Nodes or Input Layer:** The primary component of the Perceptron system responsible for receiving the initial data for further processing is the aforementioned component. Each input node has a numerical value that is authentic.
- **Weight and Bias:** The weight parameter quantifies the degree of connectivity between the units. This particular attribute holds significant importance for components of the Perceptron. The influence of the input neuron on the output is directly proportional to its weight. Moreover, the point of intersection in a linear equation might be conceptualised as a form of inherent prejudice.
- **Activation Function:** The aforementioned aspects are of paramount importance and serve as determining factors in the firing or non-firing of a neuron. The activation function can be conceptualised primarily as a step function. Various activation mechanisms:
  - Sign function, Step function and Sigmoid function

## Notes

The data scientist creates the desired outputs by using the activation function to make a judgment call based on multiple issue statements. By examining whether the learning process is slow or has disappearing or exploding gradients, it is possible to distinguish between different activation functions (such as Sign, Step and Sigmoid) in perceptron models.

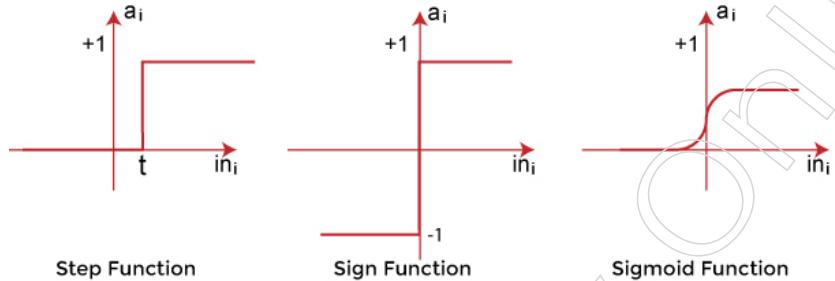


Figure : Activation Function structure

### How does Perceptron work:

Perceptrons are single-layer neural networks that are used in machine learning. The composition of these entities consists of four primary components: input values (Input nodes), weights and bias, net sum and an activation function. To compute the weighted sum, the perceptron model performs a multiplication operation between each input value and its corresponding weight. The weighted sum is further processed using the activation function 'f' to obtain the desired output. The activation function, commonly known as the step function, is denoted by the sign 'f'.

The activation function, sometimes referred to as the step function, plays a critical role in facilitating the transfer of output values within the specified range of (0,1) or (-1,1). It is important to note that the strength of a node can be assessed based on the magnitude of its input weight. In a similar vein, the inclusion of a bias value in an input facilitates the manipulation of the curve of the activation function.

### Perceptron model works in two important steps as follows:

**Step-1:** To obtain the weighted total, multiply all input values by the relevant weight values in the first step, then add the results. The weighted sum can be determined mathematically in the manner shown below:

$\sum w_i * x_i = x_1 * w_1 + x_2 * w_2 + \dots + x_n * w_n$  This weighted sum should be supplemented with a unique term called bias 'b' to enhance the model's performance.

$$\sum w_i * x_i + b$$

**Step-2:** The weighted sum discussed above is combined with an activation function in the second phase to provide output that can either be binary or continuous.

$$Y = f(\sum w_i * x_i + b)$$

**Types of Perceptron Models:** Perceptron models are classified into two categories based on the layers. These are listed below:

- Single-layer Perceptron Model
- Multi-layer Perceptron model

**Single Layer Perceptron Model:** One of the simplest types of artificial neural networks (ANN) is this one. A threshold transfer function and a feed-forward network are both included in a single-layered perceptron model. The single-layer perceptron model's primary goal is to examine things that can be linearly separated into binary categories.

- The weight parameters of a single layer perceptron model are initially allocated random values as the model's processes do not rely on recorded data. Moreover, it calculates the sum of all the inputs in terms of weight. The model exhibits activity by generating an output of +1 when the cumulative sum of all inputs over a specified threshold.
- The model's performance is considered satisfactory if the output matches the predetermined or threshold value, while the weight requirement remains constant. Nevertheless, this model exhibits certain discrepancies that manifest when varying weight input values are employed. To ensure accurate output and minimise errors, it is necessary to make several modifications to the input weights.
- "Single-layerperceptrons can only learn linearly separable patterns."

**Multi-Layered Perceptron Model:** The multi-layer perceptron model possesses a greater number of hidden layers in comparison to the single-layer perceptron model, while maintaining a same model structure. The backpropagation algorithm, also referred to as the multi-layer perceptron model, operates in a two-step process as outlined below:

- Forward Stage:** During the forward propagation stage, the activation functions are initiated at the input layer and propagate through the subsequent layers until reaching the output layer.
- Backward Stage:** During the backward step, the weight and bias parameters are modified in order to align with the requirements of the model. At this juncture, the disparity between the desired output and the observed output originated at the output layer and culminated at the input layer.

Consequently, a comparative analysis has been conducted between a multi-layered perceptron model and other artificial neural networks featuring different layers, whereby the activation function is non-linear. This comparison is made in relation to a single layer perceptron model.

In the context of deployment, alternative activation functions, such as sigmoid, TanH andReLU, can be employed instead of linear activation functions. The multi-layer perceptron model possesses the ability to effectively process both linear and non-linear patterns, hence exhibiting enhanced processing capability. Additionally, it may use logic gates like AND, OR, XOR, NAND, NOT, XNOR and NOR.

### Applications:

- Data Compression:** Data compression is the process of encoding, rearranging, or otherwise changing data to make it smaller. It entails re-encoding data with fewer bits than the original representation in its most basic form.
- Streaming Encoding:** Faster training is achieved by using an encoding technique that whitens the real-valued input data sent to the first hidden units of a fully connected neural network.

**Perceptron Function:** The output of a perceptron function,  $f(x)$ , is obtained by multiplying the input,  $x$ , by the learned weight coefficient,  $w$ . We can formulate it mathematically as follows: If  $w \cdot x + b > 0$ ,  $f(x) = 1$ ; otherwise,  $f(x) = 0$ ; where ' $w$ ' denotes a vector of real-valued weights, ' $b$ ' denotes the bias and ' $x$ ' is a vector of input  $x$  values.

**Characteristics of Perceptron:** The following traits apply to the perceptron model:

- A machine learning algorithm called perceptron is used to learn binary classifiers under supervision. The weight coefficient is automatically learned in a perceptron.
- Weights are first multiplied by input features to determine whether to fire the neuron or not.

## Notes

- To determine if the weight function is larger than zero, the activation function applies a step rule.
- In order to distinguish between the two linearly separable classes, +1 and -1, the linear decision boundary is drawn.
- It must have an output signal if the total of all input values is greater than the threshold value; otherwise, no output will be displayed.

**Limitations of Perceptron Model:** A perceptron model has the following drawbacks:

- The hard limit transfer function of a perceptron limits the output to a binary number (0 or 1) alone.
- Only sets of input vectors that can be linearly separated can be classified using perceptrons. Non-linear input vectors are difficult to correctly categorize.

**Future of Perceptron:** The Perceptron model has a very promising and important future since it aids in the interpretation of data by creating intuitive patterns and using them afterwards. The future of perceptron technology will continue to support and facilitate analytical behavior in machines that will, in turn, increase the efficiency of computers. Machine learning is a rapidly expanding field of Artificial Intelligence that is constantly evolving and in the developing stage. With the aid of artificial neurons, the perceptron model is constantly improving and solving complicated issues effectively.

### 5.3.2 Feed Forward Neural Network:

This topic pertains to feedforward neural networks, which are sometimes known as deep feedforward networks or multi-layer perceptrons. For example, these networks form the basis for convolutional and recurrent neural networks, which are extensively employed in computer vision applications. We will endeavour to render the significant principles easily comprehensible and remember, while avoiding excessive mathematical intricacies. Deep learning technologies are widely employed in several domains such as mobile applications, machine translation and search engines. The mechanism of action involves stimulating the brain to discern and establish patterns from diverse sets of input. This remarkable technology significantly relies on the utilisation of feedforward neural networks due to their efficacy in assisting programmers with tasks such as non-linear regression, function approximation, as well as pattern recognition and classification.

A feedforward neural network is a class of artificial neural networks characterised by the absence of cyclic connections among the nodes. The reverse of a feed forward neural network is a recurrent neural network, wherein specific routes are cyclically traversed. The feed forward model is considered the most elementary form of neural network due to its unidirectional processing of input. While the transmission of data may occur through multiple subterranean nodes, it consistently progresses in a unidirectional manner, without any retrograde movement.

#### How does a Feed Forward Neural Network work:

The single layer perceptron is a prevalent illustration of a feed-forward neural network in its fundamental form. Multiple inputs are introduced into the layer inside this model and subsequently multiplied by the corresponding weights. The summation of the weighted input values yields a cumulative total. The resulting number is commonly 1 and in the event that the cumulative aggregate of the values falls below the specified threshold, the output value is assigned as -1. The threshold is commonly established at a value of zero. The single layer perceptron plays a vital role as a feed forward neural network model in classification problems. Single-layer perceptrons may incorporate certain elements of artificial intelligence.

The neural network has the capability to evaluate the outputs of its nodes in relation to the desired values through a mechanism referred to as the delta rule. This property facilitates the network in adjusting its weights in order to generate output values that exhibit enhanced accuracy. The process of learning and training leads to the phenomenon known as gradient decline. The process of updating weights in multi-layered perceptrons, commonly referred to as back-propagation, exhibits a high degree of similarity. Under these conditions, the hidden layers of the network are modified based on the output values produced by the top layer.

**Feed forward Neural Network's Layers:** The following are the components of a feedforward neural network:

- The layer responsible for receiving information consists of neurons that are situated within it. The subsequent tier obtains the information subsequent to that. The quantity of variables inside the dataset is equivalent to the overall count of neurons present in the input layer.
- The neural network architecture includes a hidden layer, which functions as an intermediary layer positioned between the input and output layers. A multitude of neurons undergo modifications to the inputs inside this particular layer. Subsequently, communication is established with the output layer.
- The ultimate layer, referred to as the output layer, is determined by the architectural design of the model. Furthermore, given your understanding of the intended outcome, the output layer corresponds to the anticipated characteristic.
- The application of weights on neurons serves to quantify the strength of the connections between them. A weight possesses a numerical value that is within the range of 0 to 1.

#### Applications of Feed Forward Neural Networks:

Despite their simplicity, Feed Forward Neural Networks offer advantages in certain machine learning applications because of their streamlined architecture. Utilising a conservative mediator for moderation, an approach could involve deploying multiple feed-forward neural networks in isolation from each other. Analogous to the human brain, this method employs a multitude of individual neurons to effectively process and understand tasks of greater complexity. The integration of the individual findings from each network can be performed in order to get a consolidated and cohesive output.

Neural networks find extensive use across diverse domains. The area units for several of them are indicated as follows:

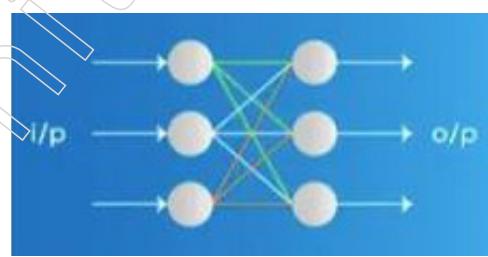


Figure: Feed Forward Neural Network

- **Pattern Recognition:** Pattern recognition refers to the utilisation of a machine learning algorithm for the purpose of identifying and discerning patterns. Pattern recognition refers to the process of categorising data by utilising either preexisting knowledge or statistical information obtained from patterns and their representations.

## Notes

- Computer Vision: Artificial intelligence (AI) encompasses a specialised area known as computer vision, which empowers computers and systems to extract relevant information from digital photos, videos and other visual inputs. This extracted data serves as a basis for subsequent actions or recommendations made by the AI system.
- Physiological feedforward system: The concept of feedforward management is shown by the typical preemptive regulation of heart rate by the central autonomic system before engaging in physical activity.
- Gene regulation and feedforward: A recurrent motif may be observed in these widely recognised networks and empirical evidence has demonstrated that this motif serves as a feedforward mechanism for detecting enduring alterations in the atmosphere.
- Automating and managing machines
- Parallel feedforward compensation with derivative: This is a more recent method of transforming the minimum part of an open-loop transfer system from the non-minimum part.

### **Advantages:**

- With a small middleman to ensure moderation, a number of Feedforward networks can function independently.
- The streamlined architecture of feed forward neural networks can improve machine learning.
- With a moderated intermediary, multiple networks in feed-forward networks operate independently.
- Several neurons are required in the network for complex tasks.
- In contrast to perceptrons and sigmoid neurons, which are otherwise complex, neural networks can handle and process nonlinear data with ease.
- Decision boundaries are a challenging problem that a neural network handles.
- The neural network architecture can change based on the data. For instance, recurrent neural networks (RNNs) excel at processing text and voice, while convolutional neural networks (CNNs) excel at processing images.
- GPUs are necessary for neural networks to handle large datasets for high computational and hardware performance.

### **Disadvantages:**

- Insufficient for deep learning.
- to optimize more variables.
- losing knowledge of the neighborhood.
- Translation invariance isn't the issue.

### **5.3.3 Multilayer Perceptron:**

A multilayer perceptron consists of an input layer, an output layer and one or more hidden layers, with each layer comprising several interconnected neurons. In contrast to the Perceptron, which requires neurons to have an activation function that enforces a threshold, such as ReLU or sigmoid, the neurons in a Multilayer Perceptron have the flexibility to apply any arbitrary activation function.

The multilayer perceptron is considered the most prevalent and extensively utilised type of neural network. In the vast majority of instances, the transmission of signals occurs from the input to the output within the network. The absence of a loop and the lack of influence of a neuron's output on itself are notable characteristics. The architectural design in question is commonly referred to as "feedforward." The phrase "hidden" pertains to strata that are not inherently connected to their immediate environment. There is ongoing discussion in the literature over the classification of the input layer as an independent layer inside the network, as its primary function is to transmit input signals to the higher layers without engaging in any input processing. The inputs are organised into clusters inside the input layer. However, for the sake of our analysis, we will focus just on the layers composed of individual neurons. Moreover, there are feedback networks that possess the ability to send impulses bidirectionally due to the presence of reaction links inside the network.

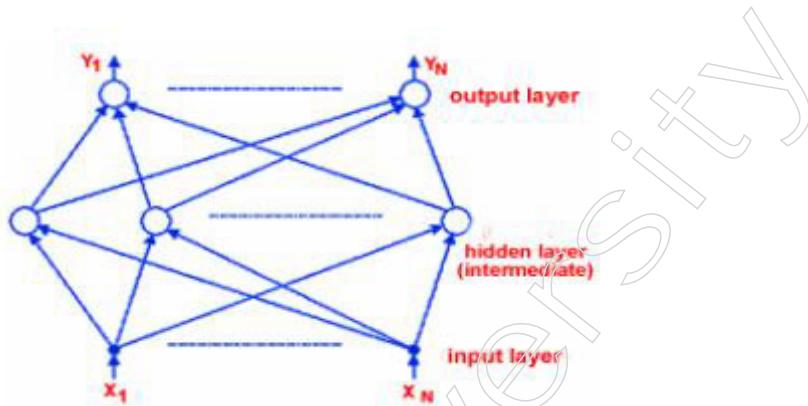


Figure: The multilayer perceptron.

These networks possess significant influence and exhibit a high degree of complexity. The network exhibits dynamic behaviour, continuously adapting until it finds a state of equilibrium. Furthermore, with every alteration in input, the network actively seeks a new state of equilibrium. The necessity to enhance the intricacy of choice regions prompted the incorporation of many layers. A perceptron consisting of only one layer and one input generates decision regions that are semi-planar.

The network's output has the ability to estimate convex decision regions, which are formed by the intersection of semi-planes generated by the neurons. This is achieved by introducing an additional layer where each neuron acts as a standard perceptron for the outputs of the neurons in the preceding layer. An alternate option that can be considered is a three-layer perceptron, as depicted in the accompanying illustration.

It was shown that, if the activation functions of neurons are linear, multilayer networks do not offer an improvement in processing power when compared to networks with a single layer since a linear function of a linear function is likewise a linear function. Non-linear activation functions are precisely where the multilayer perceptron's power lies. Except for polynomial functions, almost every non-linear function can be utilized for this. Currently, the single-pole (or logistic) sigmoid, as shown in Figure below, is the function that is most frequently utilized.

$$f(s) = \frac{1}{1 + e^{-s}}$$

## Notes

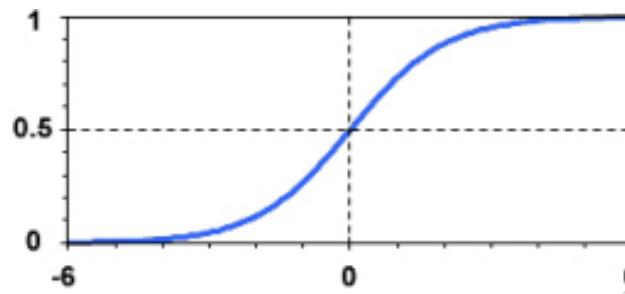


Figure : Sigmoid single-pole activation function.

The bipolar sigmoid function (the hyperbolic tangent), for  $a=2$ , is depicted in Figure below:

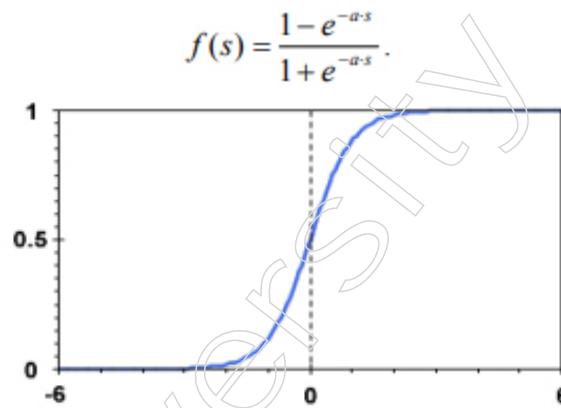


Figure : Sigmoid single-pole activation function.

This supports the multilayer perceptron's ability to serve as a universal approximator. Additionally, it was shown that neural networks can calculate specific polynomial expressions by applying the Stone-Weierstrass theorem to them: if there are two networks that calculate precisely two functions,  $f_1$  and  $f_2$ , then there is a larger network that calculates precisely a polynomial expression of  $f_1$  and  $f_2$ .

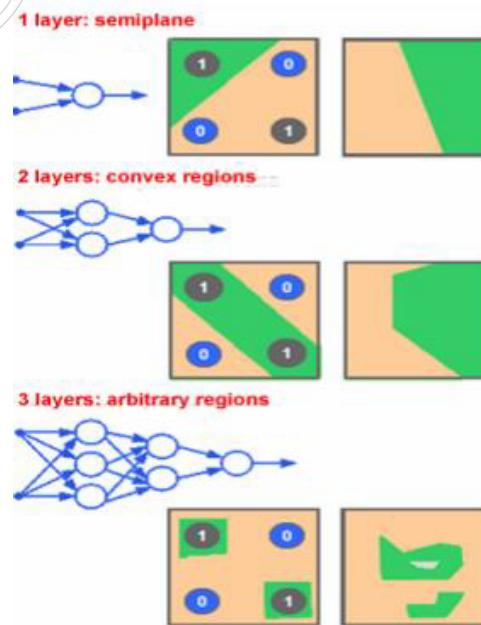


Figure: Decision regions of multilayer perceptrons.

The most popular and well-known type of neural network is a multi-perceptron, which is made up of trained units like those in the figure below. Each of these units creates a weighted input total to which a constant is added. Following that, this amount is processed via a nonlinear function that is frequently referred to as the activation function. The majority of units are interconnected in a “feed forward” way, i.e. interconnections that form an aloop, as seen in the following figure.

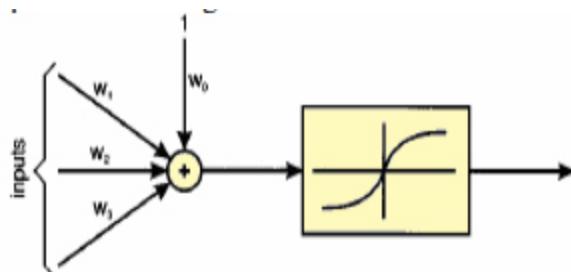


Figure: Example network “feed forward”. Each circle represents a unit of the type shown in Figure . Each connection between units is a share. Each unit also has an entry in the diagonal are not shown.

### Working of MultiLayer Perceptron Neural Network:

- The characteristic of the dataset is represented by the input node.
- Each input node sends the hidden layer’s hidden layer the vector input value.
- Each edge in the hidden layer has a weight that is multiplied by the input variable. The sum of all the production values from the concealed nodes is calculated. to produce the result.
- In the buried layer, the active nodes are recognized using the activation function.
- The output layer receives the output.
- Determine the discrepancy between output as planned and output as achieved at the output layer.
- Following the computation of the anticipated output, the model employs backpropagation.

### Advantages of MultiLayer Perceptron Neural Network:

- Non-linear issues can be solved with ease by MultiLayer Perceptron Neural Networks.
- It is capable of handling big datasets and difficult issues.
- This paradigm is used by developers to address the fitness issue with neural networks.
- It uses backpropagation to increase accuracy and lower prediction error.
- The Multilayer Perceptron Neural Network accurately predicts the outcome after model training.

### Disadvantages of MultiLayer Perceptron Neural Network:

- This neural network uses a lot of computation, which occasionally raises the model’s overall cost.
- Only when the model has received perfect training will it function well.
- Due of the close connections in this model, there are more parameters and nodes are redundant.

## Notes

### 5.3.4 Convolutional Neural Network:

Convnets, also known as convolutional neural networks (CNNs), are a particular type of feedforward neural network. In that they are composed of neurons with trainable weights and biases, they are quite similar to the neural networks discussed above. The fundamental distinction is that we can encode certain features in the CNN design since it implicitly assumes that the input is image-like. Convolutions in particular capture translation invariance (i.e., filters are position independent). As a result, the forward function is more effective, there are much less parameters, the network is easier to improve and the dependence on the amount of the data is reduced.

Unlike conventional neural networks, CNNs' layers feature neurons arranged in a few different dimensions, including channels, width, height and number of filters in the most basic 2D example. Similar to an MLP, a convolution neural network is made up of a series of layers, each of which modifies the activations or outputs of the layer before it using a different differentiable function. The convolution layer, pooling layer and fully connected layers are the most typical building blocks you will find in most CNN architectures. There are other layers used in CNNs as well and they will be covered in following sections. These layers essentially function as dimensionality reduction, feature extractors and classification layers, respectively. The whole convolutional layer of a CNN is created by stacking these CNN layers.

#### General Model of Convolution Neural Network:

**General Model:** The conventional artificial neural network (ANN) model typically consists of multiple hidden layers, in addition to a single input layer and a single output layer. A designated neuron receives an input vector X and produces an output vector Y by the application of a function F, as denoted by the general equation presented herein.

$$F(X, W) = Y$$

where W stands for the weight vector, which symbolizes how strongly neurons in two adjacent layers are connected. The weight vector that was created can now be utilized to classify images. The classification of images based on pixels has been extensively studied in the literature. Contextual information, such as the image's shape, gives better results or outperforms, nonetheless. CNN is a model that is attracting interest due to its ability to classify data based on context. The graphic below describes the CNN model in its entirety. Convolution layer (a), pooling layer (b), activation function (c) and fully connected layer (d) make up a standard CNN model. Below is an illustration of each component's functionality:

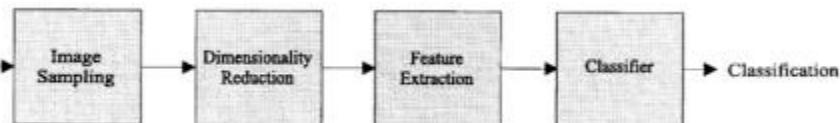


Figure: Elementary constituents of CNN

**Convolution Layer:** The input layer of the system receives an image that requires classification and the predicted class label is generated by utilising features extracted from the image. The receptive field refers to the specific connection established between an individual neuron in the subsequent layer and a subset of neurons in the preceding layer. The retrieval of local details from the input image is accomplished by utilising the concept of receptive field.

The formation of a weight vector is a result of the connection between a neuron's receptive field and a specific place in the preceding layer, which in turn relates to the

neurons in the subsequent layer. Given that the neurons on a two-dimensional plane possess identical weights, it becomes feasible to discern similar characteristics that manifest across multiple locations within the input data. The representation of this concept is depicted in the figure provided below.

The generation of the feature map is achieved by systematically moving the weight vector, also known as the filter or kernel, along the input vector. The convolution operation refers to the procedure of horizontally and vertically moving the filter. The aforementioned methodology involves the generation of N filters and N feature maps through the extraction of N distinct features from the input image.

These features are subsequently positioned on a singular layer, thereby representing each unique feature. The local receptive field phenomenon results in a significant reduction in the number of trainable parameters. The output value  $a_{ij}$  in the subsequent layer for a given point  $(i,j)$  is computed using the convolution process, as described by the formula provided below:

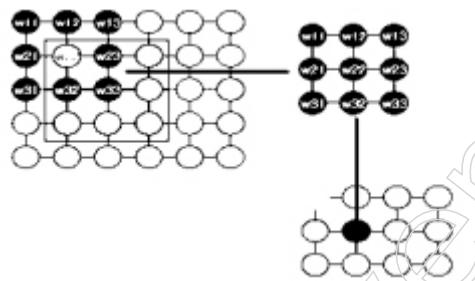


Figure: Receptive field of particular neuron in the next layer.

$$a_{ij} = \sigma((W * X)_{ij} + b)$$

In this context,  $X$  is the input received by the layer,  $W$  denotes a filter or kernel that is applied to the input,  $b$  represents the bias term,  $*$  signifies the convolution operation and symbolises the nonlinearity incorporated into the network.

**Pooling Layer:** Once a feature has been found, the exact positioning of this feature becomes less significant. As a result, the pooling or sub-sampling layer is typically positioned following the convolution layer. The utilisation of the pooling technique offers notable benefits, such as the introduction of translation invariance and a substantial reduction in the number of trainable parameters. As seen in the visual representation provided, a certain window is selected for the purpose of conducting the pooling procedure. Subsequently, the input elements encompassed within said window are subsequently sent through the pooling function.

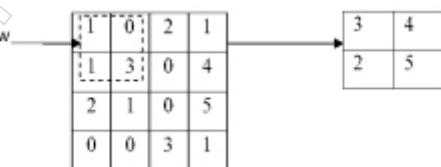


Figure: Pooling operation performed by choosing a  $2 \times 2$  window.

The pooling function generates an additional output vector. There exist a limited number of pooling techniques, including average pooling and max-pooling. Among them, max-pooling is the most commonly employed approach, known for its significant reduction in map size. The error does not propagate to the winning unit during computation mistakes, as it does not partake in the forward flow.

## Notes

**Fully Connected Layer:** The fully connected network in traditional models is analogous to the fully connected layer. The fully connected layer receives the output of the first phase, which comprises repetitive convolution and pooling and computes the dot product of the weight vector and the input vector to produce the final output. Gradient descent, sometimes referred to as batch mode learning or the offline approach, lowers the cost function by estimating the cost throughout the full training dataset. It changes the parameters only once every epoch, or complete traversal of the training dataset. Although it produces global minima, the size of the training dataset has a significant impact on how long it takes to train the network. Stochastic gradient descent was used to replace this method of decreasing the cost function.

**Activation Function:** The utilisation of the sigmoid activation function in conventional machine learning techniques is extensively documented in the academic literature. The utilisation of Rectified Linear Unit (ReLU) has demonstrated its superiority over its predecessor in terms of creating non-linearity, primarily due to two key factors. To begin with, the computation of the partial derivative of the Rectified Linear Unit (ReLU) function is straightforward. Furthermore, it has been observed that saturating non-linearities, such as the sigmoid function represented by  $\frac{1}{1+e^{-x}}$ , exhibit poorer computational speed compared to non-saturating non-linearities, such as the Rectified Linear Unit (ReLU), even when considering the training time.

Third, ReLU prevents gradients from going away. However, a big gradient that is flowing through the network reduces ReLU efficiency and an update in weight prevents activation of the neuron, which results in the Dying ReLU problem, a significant difficulty that is frequently experienced. Leaky ReLU, where  $\alpha$  is a tiny constant, can be used to solve this problem. If  $x>0$ , the function activates as  $f(x)=x$  and if  $x<0$ , the function activates as  $x+\alpha$ .

**Architectures Of Convolution Neural Network:** In CNN, numerous architectures have been created and put into use. Below are succinct descriptions of those architectures:

**LeNet Architecture:** The ability of multi-layer networks to learn from extremely complicated and high-dimensional input makes them suited for picture recognition tasks. The following paragraph provides a summary of the LeNet architecture, which was proposed in 1998 and leverages datasets. The figure below shows the LeNet5 architecture. It comprises eight layers, five of which are convolutional and three of which are fully linked. A plane has 25 inputs per unit. Units in the first hidden layer receive input from the  $5 \times 5$  area, a small portion of the input image that is transferred to the first hidden layer. Receptive field of the unit refers to this particular section of the input image. In a plane, every unit has the same weight vector. The unit's output is kept in the same spot on the feature map.

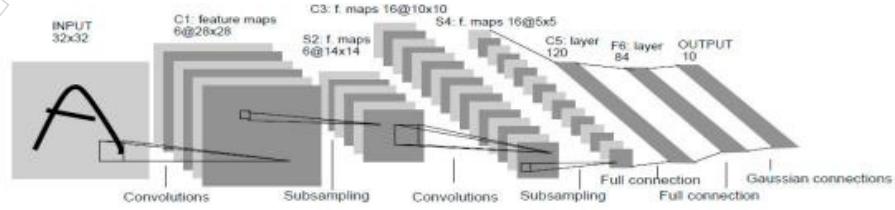


Figure: Architecture of LeNet5, a CNN where each box represents a different feature map.

The adjacent units in the feature map are the outcome of the adjacent units in the preceding layer. Contiguous receptive fields overlap as a result of this. Convolution layer is the first layer and it is made up of neurons that produce sigmoid activation when applied to the weighted sum. As seen in figure, if a  $5 \times 5$  area is selected as an input and a horizontal shift on the area is carried out, it will result in the overlap of four rows and

five columns while computing contiguous units in the feature map. Different feature maps are produced as a result of applying various weight vectors to the same input image.

The produced feature maps can be used to extract various features. The fundamental characteristic of CNN is that a small change in the input has no impact on the feature map. Since it is not important to have a precise position for a feature in an image, subsampling is done to lower the precision value. Sub-sampling has been shown in the second layer of the previous graphic, as seen. The amount of feature maps obtained through subsampling is equal to that of feature maps obtained through convolution. Here, the average of the four inputs has been computed for the sub-sampling layer  $2 \times 2$  area, multiplied by the trainable coefficient, added the trainable bias and then sent to the sigmoid function. As the spatial resolution is reduced layer by layer, an increase in the number of feature maps may be seen. The back propagation approach is used to carry out the learning.

**AlexNet Architecture:** This piece provides a succinct description of the AlexNet architecture, a modified version of LeNet that has been proposed. For the purpose of classifying 1.2 million high resolution photos into 1000 different classes, it is composed of three completely connected layers and five convolution layers. To train the network more quickly, non-saturating neurons and effective GPU implementation are used. As a result, a huge network is needed to enable the network to classify objects from millions of photos, which may ultimately result in a high demand for training a very large number of weights and the overfitting issue.

This issue was solved by using the dropout method. In this method, the neurons that have a probability of 0.5 are damped and do not participate in forward or backward propagation. Overfitting is significantly reduced since the neurons that depend on these damped neurons are forced to learn the most robust features entirely on their own. The dropout method doubles the amount of iterations needed to converge. It takes five to six days to train the network using two GTX 580 3GB GPUs. The main characteristics of this design included the addition of ReLU non-linearity in CNN, which caused the convergence rate to rise quickly.

**GoogleNet Architecture:** GoogleNet is a popular concept that was proposed. After winning the ILSVRC14 competition, it became enlightened. The main objective was to create a model with a smaller budget that could consume less memory, power and trainable parameters. The number of trainable parameters employed in the network was dramatically decreased by the model. The following is a general description of the architecture. In essence, it employs 12 million fewer parameters than the model put forward by. The architecture attempted to create a network that could more accurately identify the items in an image.

This might be accomplished by expanding the network's size, which would increase the number of layers, but a key negative of this idea was that doing so would increase the number of parameters that would need to be trained, which would create the issue of overfitting. Another significant drawback is that as the number of filters increases, so does the computation, which raises overhead. Implementing a sparse matrix was the suggested approach. In order to create the best network topology, the highly correlated units join to form a cluster in the previous layer and send input to the next layer. Even though the computations are sped up by 100 times, the overhead of cache misses still exists when using the non-uniform sparse matrix. Using highly optimized numerical libraries to achieve faster computations is also ineffective. As a result, the state of the art relies on uniform sparse matrices.

## Notes

### Advantages of Convolutional Neural Network (CNN):

- One of the primary characteristics of Convolutional Neural Networks (CNNs) is their ability to efficiently process images. This phenomenon can be attributed to the utilisation of a technique known as convolution, wherein a filter is applied to an image to extract relevant features for the specific task at hand. Convolutional neural networks (CNNs) are capable of processing information in a more efficient and expedient manner compared to other methods. This is achieved by the reduction of data volume that necessitates processing.
- One advantage of Convolutional Neural Networks (CNNs) is their ability to achieve high accuracy rates. This phenomenon can be attributed to the practise of analysing extensive datasets, which enables individuals to develop the ability to identify complex patterns inside photographs. This characteristic renders them very suitable for tasks such as facial recognition or object identification, as they may undergo training to achieve precise recognition of certain items or properties.
- Moreover, Convolutional Neural Networks (CNNs) exhibit robustness to noise, hence facilitating their ability to identify patterns in images that have been compromised or distorted. This phenomenon can be attributed to the process of extracting features from images through the utilisation of several layers of filters, hence endowing them with enhanced resistance to noise compared to alternative techniques.
- Transfer learning is a notable advantage of convolutional neural networks (CNNs). This capability allows CNNs to be initially trained for a specific task and subsequently utilised for another task without requiring much additional guidance or instruction. This phenomenon can be attributed to the fact that convolutional neural networks (CNNs) possess the ability to extract attributes that are often sufficiently generalizable to be applied across a wide range of tasks, hence rendering them a versatile tool for several applications.
- Automated feature extraction – The feature extraction process is automated by CNNs, allowing them to learn to recognize patterns in images without the need for manual feature engineering. Because the CNN can be trained to recognize the relevant features, they are perfect for jobs where the features that are crucial to the task are not known in advance.

### Disadvantages of Convolutional Neural Network (CNN):

- One of the primary limitations of Convolutional Neural Networks (CNNs) is their significant computational requirements. The reason for this phenomenon is that Convolutional Neural Networks (CNNs) may have numerous layers and parameters, resulting in a substantial requirement for computational resources and memory during both training and operation. Due to this factor, its suitability for use in certain applications with limited resources may be questionable.
- One challenge that arises when working with tiny datasets is that Convolutional Neural Networks (CNNs) require large datasets in order to achieve high levels of accuracy. This phenomenon can be attributed to the extensive examination of various instances of patterns in images prior to acquiring the ability to identify them. The phenomenon of overfitting can occur in the CNN model when the size of the dataset is insufficient, leading to an excessive specialisation of the model to the training data and subsequently poor performance when presented with new, unseen data.
- In order to achieve high accuracy rates, Convolutional Neural Networks (CNNs)

require the use of large datasets. This phenomenon can be attributed to the extensive analysis of various patterns in images that individuals undertake prior to acquiring the ability to identify them. The phenomenon of overfitting can occur in the CNN model when the size of the dataset is insufficient, leading to the model being excessively specialised to the training data and therefore exhibiting poor performance when presented with new, unseen data. One further limitation of Convolutional Neural Networks (CNNs) is their lack of interpretability. Therefore, understanding the decision-making process employed by CNN is a significant challenge. This issue might provide a challenge in applications where it is crucial to comprehend the underlying reasoning behind a specific decision.

- **Vulnerability to adversarial attacks** – Adversarial assaults, which entail purposefully modifying the input data to trick the CNN into making false judgments, are another danger to CNNs. In applications like driverless vehicles, where safety is a top priority, this can be a major issue.
- **Limited ability to generalize** – Finally, CNNs can only a limited extent generalize to novel circumstances. As a result, they might not perform well when presented with images that are substantially dissimilar from those in the training dataset. This may present a challenge in applications where the CNN must handle a wide range of image types.

**Applications of Convolutional Neural Networks:** The evident contenders encompass conventional convolutional neural network (CNN) implementations observed in everyday scenarios, such as software for speech recognition, image categorization and facial identification. These beliefs are very prevalent among individuals without specialised knowledge, such as ourselves and significantly influence our everyday experiences, particularly inside image-centric social media sites such as Instagram. Presented here is a compilation of several prominent applications associated with CNN (Convolutional Neural Networks) as developed and utilised by CNN (Cable News Network).

**Understanding Gray Areas:** The goal of adding gray areas to CNNs is to present a far more accurate representation of reality. CNNs currently observe a true and false value for each inquiry, operating largely just like a machine. However, we recognize as humans that there are countless shades of gray in the real world. It will be easier for the computer to comprehend and absorb fuzzier logic if it can see the gray region that exists in human thought and that we try to avoid. This will enable CNN to have a more complete picture of what people see.

**Advertisin:** With the introduction of programmatic buying and data-driven tailored advertising, CNNs have already made a significant effect in the advertising industry.

**Other Interesting Fields:** With the development of autonomous automobiles, mimicking human behavior robots, aids for human genome mapping projects, earthquake and natural catastrophe forecasting and perhaps even self-diagnosis of medical issues, CNNs are positioned to be the technology of the future. Therefore, you wouldn't even need to make a trip to the clinic or make an appointment with a doctor to be sure your severe cold symptoms are just the flu and not the signs of a rare disease. The diagnosis of brain cancer is one issue on which researchers are focusing using CNNs. More lives that are affected by brain cancer may be saved if the disease is detected sooner.

## Summary

- Deep learning is a subfield of machine learning that focuses on using artificial neural networks to model and solve complex problems. It has gained significant attention

## Notes

due to its ability to automatically learn and extract features from data, enabling the creation of highly accurate predictive models. Deep learning has revolutionized various fields, including image and speech recognition, natural language processing, game playing and more. Deep learning's ability to automatically learn from data has led to significant breakthroughs and has reshaped the landscape of artificial intelligence. It continues to advance with research into more effective architectures, training algorithms and applications.

- Neural networks are computational models inspired by the human brain's structure and functionality. They are used in various machine learning tasks, including pattern recognition, classification, regression and more. Neural networks consist of interconnected layers of nodes (neurons) that process and transform data to make predictions or decisions. At the core of a neural network is the neuron, a computational unit that processes input data and produces an output. The neuron model, also known as the perceptron model, consists of several components. The neuron model is a basic building block that forms the foundation of artificial neural networks, a neural network is composed of multiple layers of interconnected neurons.
- Neural networks have the ability to model highly complex relationships in data, making them suitable for tasks like image recognition, language processing and more. Their utility lies in their capacity to automatically learn and extract features from the data, eliminating the need for manual feature engineering. Through a process called backpropagation, neural networks adjust their weights and biases during training to minimize the difference between predicted and actual outputs. Neural networks have achieved remarkable successes in various fields, including image classification, speech recognition, natural language processing and playing games. Their versatility and ability to capture intricate patterns make them a powerful tool for solving a wide range of problems.
- The design of artificial neural networks is inspired by the biological neurons and neural networks in the human brain. While artificial neural networks are not direct replicas of biological systems, there are some connections and motivations. While the analogy between artificial and biological neural networks is not perfect, these connections provide a conceptual framework for understanding the design and behavior of deep learning models. Recurrent Neural Networks (RNNs) are designed for sequence data, where the order of elements matters. RNNs are used in tasks like natural language processing, speech recognition and time series analysis, where the order of data is crucial for understanding and making predictions. These architectures showcase the diversity and specialization of neural networks for various types of data and tasks. Understanding their strengths and design principles is essential for effective application in solving real-world problems.
- The perceptron is the fundamental building block of neural networks and is the simplest form of an artificial neuron. It was introduced by Frank Rosenblatt in the 1950s. The perceptron takes a set of inputs, applies weights to them, sums the weighted inputs, adds a bias term and then passes the result through an activation function. A feed-forward neural network, also known as a single-layer neural network, consists of an input layer, one or more hidden layers and an output layer. Each neuron in a layer is connected to all neurons in the subsequent layer, but there are no connections within the same layer. Convolutional Neural Networks (CNNs) are specialized for processing grid-like data, particularly images. They excel at detecting patterns and features within images due to their unique architecture. CNNs are widely used in image classification, object detection, image generation and various computer vision tasks.

## Glossary

- **ANNs:** Artificial neural networks
- **DNNs:** Deep neural networks
- **CNNs:** Convolutional neural networks
- **RNNs:** Recurrent neural networks
- **DBNs:** Deep belief networks
- **GPUs:** Graphics Processing Units

## Check Your Understanding

1. What is the primary advantage of deep learning models over traditional machine learning models?
  - a) Simplicity and ease of implementation
  - b) Faster training times
  - c) Ability to automatically learn features from data
  - d) Reduced need for labeled training data
2. What is the role of the “hidden layers” in a deep neural network?
  - a) They automatically learn features from the data.
  - b) They process input data and produce final predictions.
  - c) They connect the input layer to the output layer.
  - d) They prevent overfitting by adding noise to the data.
3. What does the term “backpropagation” refer to in the context of deep learning?
  - a) A training technique for reinforcement learning models.
  - b) A method for adjusting neural network parameters using gradients.
  - c) A type of activation function used in convolutional layers.
  - d) An approach to ensemble learning using multiple models.
4. What is a key challenge addressed by regularization techniques in deep learning?
  - a) Increasing the number of layers in neural networks.
  - b) Reducing the depth of neural networks.
  - c) Minimizing the number of training epochs.
  - d) Avoiding overfitting and improving generalization.
5. What is the primary purpose of neural networks in machine learning?
  - a) To capture and learn patterns from data
  - b) To perform mathematical operations on data
  - c) To automate repetitive tasks in data preprocessing
  - d) To generate random numbers for statistical analysis
6. In the context of a neuron model, what is the purpose of the activation function?
  - a) To determine the initial weight of the neuron
  - b) To add a constant bias to the neuron’s output
  - c) To introduce non-linearity into the neuron’s output
  - d) To connect the neuron to other neurons in the network

**Notes**

7. What is a key advantage of using neural networks for modelling and solving complex problems?
  - a) They require minimal computational resources for training.
  - b) They eliminate the need for labeled training data.
  - c) They can only be applied to problems with a small number of variables.
  - d) They can automatically learn and extract features from data.
8. What is a key similarity between artificial neural networks and biological neural networks?
  - a) Artificial neural networks consist of neurons that communicate through physical connections.
  - b) Both types of networks have hierarchical structures with hidden layers.
  - c) Biological neural networks process information sequentially.
  - d) Artificial neural networks can perform autonomic functions.
9. Which popular CNN architecture introduced the concept of residual blocks?
  - a) LeNet
  - b) VGGNet
  - c) AlexNet
  - d) ResNet
10. What is the primary advantage of Recurrent Neural Networks (RNNs) in handling sequence data?
  - a) RNNs have the ability to capture temporal dependencies and context.
  - b) RNNs require less training data compared to other types of networks.
  - c) RNNs can process sequences in parallel.
  - d) RNNs eliminate the need for activation functions.
11. Which type of neural network is best suited for processing grid-like data, such as images?
  - a) Multilayer Perceptron
  - b) Recurrent Neural Network
  - c) Convolutional Neural Network
  - d) Radial Basis Function Network
12. In the context of a perceptron, what role does the activation function play?
  - a) It determines the number of hidden layers in the network.
  - b) It adjusts the weights of the inputs during training.
  - c) It introduces non-linearity to the perceptron's output.
  - d) It controls the flow of data in feedback loops.
13. What is the primary characteristic of a feed-forward neural network?
  - a) It has connections that loop back from later layers to earlier layers.
  - b) It processes data in a sequential manner.
  - c) It only has an input layer and an output layer.
  - d) It lacks feedback connections between layers.

14. What distinguishes a multilayer perceptron from a single-layer perceptron?
- Multilayer perceptrons have multiple input layers.
  - Multilayer perceptrons can have one or more hidden layers.
  - Multilayer perceptrons are used exclusively for regression tasks.
  - Multilayer perceptrons don't use activation functions.
15. Which neural network architecture is specifically designed for processing grid-like data, such as images?
- Recurrent Neural Network (RNN)
  - Multilayer Perceptron (MLP)
  - Convolutional Neural Network (CNN)
  - Radial Basis Function Network (RBFN)
16. What are the disadvantages of CNN?
- High computational requirements
  - Difficulty with small datasets
  - Vulnerability to adversarial attacks
  - All of the above
17. What is the working of MultiLayer Perceptron Neural Network?
- The characteristic of the dataset is represented by the input node.
  - Each input node sends the hidden layer's hidden layer the vector input value.
  - Each edge in the hidden layer has a weight that is multiplied by the input variable. The sum of all the production values from the concealed nodes is calculated. to produce the result.
  - All of the above
18. Which neural network is a type of artificial neural network in which there is no cycle in the connections between the nodes?
- Feed forward
  - Recurrent
  - Convolutional
  - Artificial
19. Which among the following areas Deep Learning has made substantial progress?
- Convolutional neural networks (CNNs)
  - Recurrent neural networks (RNNs)
  - Deep belief networks (DBNs)
  - All of the above
20. What is the full form of ILSVRC?
- Imagine Net Large-Scale Visual Recognition Challenge
  - Image Net Large-Scale Visual Recognition Challenge
  - Image Net Large-Scale Video Record Challenge
  - Image Net Larger-Scale Visual Recorder Challenge

**Notes**

**Notes****Exercise**

1. Define deep learning
2. The Neuron Model
3. What do you meant by the neural network and its utility in modelling and solving problems?
4. Define RNNs.
5. Define perceptron and multilayer perceptron.
6. Explain feed forward neural network and convolutional neural network.

**Learning Activities**

1. What are some challenges and considerations when training deep neural networks on small datasets?
2. Discuss the role of gradient descent in optimizing the neural network's weights and its impact on the training process

**Check Your Understanding-Answers**

- |       |       |       |       |
|-------|-------|-------|-------|
| 1. c  | 2. a  | 3. b  | 4. d  |
| 5. a  | 6. c  | 7. d  | 8. b  |
| 9. d  | 10. a | 11. c | 12. c |
| 13. d | 14. b | 15. c | 16. d |
| 17. d | 18. a | 19. d | 20. b |

**Further Readings and Bibliography**

1. "Deep Learning" by Ian Goodfellow, Yoshua Bengio and Aaron Courville.
2. "Neural Networks and Deep Learning: A Textbook" by Charu C. Aggarwal.
3. "Deep Learning for Computer Vision" by Rajalingappa Shanmugamani.
4. "Deep Learning with PyTorch" by Eli Stevens, Luca Antiga and Thomas Viehmann.