

HOUSE PRICE PREDITION PROJECT

A PROJECT REPORT

In partial fulfillment of the requirements for the award of the degree
MCKV Institute of Engineering
(An Autonomous Institute under UGC Act, 1956 Approved by AICTE Affiliated to
Maulana Abul Kalam Azad University of Technology, West Bengal)

BACHELORE OF TECHNOLOGY

Under the guidance of
SOURAV GOSWAMI SIR



**MCKV Institute of Engineering 243, G.T. Road(N)s
Liluah, Howrah – 711204**

(Note: All entries of the Performa of approval should be filled up with appropriate and complete information. Incomplete Performa of approval in any respect will be summarily rejected.)

1. Title of the Project: HOUSE PRICE PREDICTION PROJECT
2. Project Member: GOURAB BASAK , ANUPAM GHOSH
3. Name of the guide: Mr. SOURAV GOSWAMI SIR
4. Address: Ardent Computech Pvt. Ltd (An ISO 9001:2008 Certified)
CF-137, Sector - 1, Salt Lake City, Kolkata – 700064

Project version control history

Version	Primary Author	Description of Version	Starting Date	Date completed
Final	Gourab Basak	Project Report	1 st July 2025	10 th July 2025
Final	Anupam Ghosh	Project Report	1 st July 2025	10 th July 2025

Signature of Students

Date:

Signature of Approver

Date:

MR. Sourav Goswami Sir

DECLARATION

We hereby declare that the project work being presented in the project proposal entitled “CREDIT CARD FRAUD DETECTION” in partial fulfillment of the requirements for the award of the degree of BACHELOR OF TECHNOLOGY at ARDENT COMPUTECH PVT. LTD, SALLAKE, KOLKATA, WEST BENGAL, is an authentic work carried out under the guidance of MR. Sourav Goswami. The matter embodied in this project work has not been submitted elsewhere for the award of any degree of our knowledge and belief.

Date:

Name of the Student: Gourab Basak ; Anupam Ghosh

Signature of the student

Gourab Basak ; Anupam Ghosh

CERTIFICATE

This is to certify that this proposal of minor project entitled “CREDIT CARD FRAUD DETECTION” is a record of Boniface work, carried out by ADITYA RAJ, ROHIT KUMAR THAKUR under my guidance at ARDENT COMPUTECH PVT LTD. In my opinion, the report in its present form is in partial fulfillment of the requirements for the award of the degree of BACHELOR OF TECHNOLOGY and as per regulations of the ARDENT®. To the best of my knowledge, the results embodied in this report, are original in nature and worthy of incorporation in the present version of the report.

Guide / Supervisor

MR. SOURAV GOSWAMI SIR

ACKNOWLEDGEMENT

Success of any project depends largely on the encouragement and guidelines of many others. I take this sincere opportunity to express my gratitude to the people who have been instrumental in the successful completion of this project work.

I would like to show our greatest appreciation to Mr. Sourav Goswami sir, Professor at ARDENT®. I always feel motivated and encouraged every time by his valuable advice and constant inspiration; without his encouragement and guidance this project would not have materialized.

Words are inadequate in offering our thanks to the other trainees, project assistants and other members at ARDENT®, for their encouragement and cooperation in carrying out this project work. The guidance and support received from all the members and who are contributing to this project, was vital for the success of this project.

ABSTRACT

Artificial intelligence is finding its way into ever more areas of life. The latest craze is AI chips and related applications on the smartphone. However, technology began as early as the 1950s with the Dartmouth Summer Research Project on Artificial Intelligence at Dartmouth College, USA. The beginnings go even further back to the work of Alan Turing - which goes back to the well-known Turing test -, Allen Newell and Herbert A. Simon. With the chess computer Deep Blue from IBM, which succeeded in 1996 as the first machine to beat the then-reigning chess world champion Garry Kasparov in a match, the artificial intelligence managed to get into the focus of the world public. In data centres and on mainframes, AI algorithms have been used for many years.

Table of Contents

Introduction	9
History of Artificial Intelligence	10
Milestones for AI History	10 - 11
What is an Artificial Intelligence	11 - 12
Goals of AI	12 - 14
Application of Artificial Intelligence (AI)	15
Machine Learning	Abstract
.....	16
Introduction to Machine Learning.....	16 - 17
WHAT IS MACHINE LEARNING?.....	18
Training & Test Data	19
ML Categories	20
Supervised Learning	20
Unsupervised Learning	21
Reinforcement Learning	22
The Machine Learning Tool Box	23 - 26
DATA SCRUBBING	27
Feature Selection.....	27 - 28
Missing Data.....	29
SETTING UP YOUR DATA	29 - 30
Cross Validation	31
Machine Learning Algorithms	
Supervised Learning	32
Unsupervised Learning.....	33
Semi-Supervised Learning	33
Reinforcement Learning	33
Linear Regression	34
Logis c Regression.....	34 - 35
Decision Tree	36
Support Vector Machine (SVM).....	37

Naïve Bayes.....	38
KNN Algorithm (K- Nearest Neighbour).....	39
K means Algorithm.....	40
Random Forest.....	40
Algorithms	41
Linear Regression Algorithm.....	41
Logis c Regression Algorithm	41
Decision Tree Algorithm	41
Naïve Bayes.....	42
KNN Algorithm	42
K means Clustering.....	42
Random Forest.....	42
Implementation of the Project	
Data Visualization	
Data Pre - processing:	
Feature Engineering	
Model Creation	
CONCLUSION.....	58
Acknowledgement.....	58

Introduction

In recent years, incredible progress has been made in computer science and AI. Watson, Siri or Deep Learning show that AI systems are now delivering services that must be considered intelligent and creative. And there are fewer and fewer companies today that can do without artificial intelligence if they want to optimize their business or save money. AI systems are undoubtedly very useful. As the world becomes more complex, we need to leverage our human resources and high-quality computer systems help. This also applies to applications that require intelligence. The other side of the AI medal is: The possibility that a machine might possess intelligence scares many. Most people believe that intelligence is something unique, which is what distinguishes Homo sapiens. But if intelligence can be mechanized, what is unique about humans and what sets it apart from the machine? The quest for an artificial copy of man and the complex of questions involved are not new. The reproduction and imitation of thought already occupied our ancestors. From the sixteenth century, it was teeming with legends and the reality of artificial creatures. Homunculi, mechanical automata, the golem, the Mälzel chess automaton, or Frankenstein were all imaginative or real attempts in the past centuries to artificially produce intelligences and to imitate what is essential to us. The idea of making inanimate objects into intelligent beings by giving life a long time is fascinating the mind of mankind. Ancient Greeks had myths about robotics, and Chinese and Egyptian engineers made automatons. We can see the traces of the beginning of modern artificial intelligence as an attempt to define the classical philosophers' system of human thought as a symbolic system. However, the field of artificial intelligence was not formally established until 1956. In 1956, a conference "Artificial Intelligence" was held for the first time in Hanover, New Hampshire, at Dartmouth College. Cognitive scientist Marvin Minsky at MIT and other scientists participating in the conference were quite optimistic about the future of artificial intelligence. As Minsky stated in his book "AI: The Tumultuous Search for Artificial Intelligence": "In a generation, the problem of artificial intelligence creation will be solved at a significant level." One of the most important visionaries and theoreticians was Alan Turing (1912-1954): in 1936, the British mathematician proved that a universal calculator - now known as the Turing machine - is possible. Turing's central insight is that such a machine is capable of solving any problem as long as it can be represented and solved by an algorithm. Transferred to human intelligence, this means that if cognitive processes can be algorithm can be broken down into finite well-defined individual steps, they can be executed on one machine. A few decades later, the first practical digital computers were actually built. Thus, the "physical vehicle" for artificial intelligence was available. The electromechanical machine of Turing, considered a precursor of modern computers, managed to unlock the code used by the German submarines in the Atlantic. His work at Bletchley Park is considered key to the end of World War II. His work at Bletchley Park, an isolated country house north of London, was made public in the 1970s, when the role of the brilliant mathematician in the war was revealed. The cryptographers who worked helped shorten World War II by about two years, by deciphering around 3,000 German military messages a day. Turing's team deciphered the 'Enigma' code, which the Germans considered unbreakable, and designed and developed Colossus, one of the first programmable computers.

History of Artificial Intelligence

To be informed about the history of artificial intelligence, it is necessary to go back to previous dates in Milat. In the Ancient Greek era, it is proven that various ideas about humanoid robots have been carried out. An example of this is Daedalus, who is said to have ruled the mythology of the wind, to try to create artificial humans. Modern artificial intelligence has begun to be seen in history with the aim of defining philosophers' system of human thought. 1884 is very important for artificial intelligence. Charles Babbage, on this date, has worked on a mechanical machine that will exhibit intelligent behavior. However, as a result of these studies, he decided that he would not be able to produce a machine that would exhibit as intelligent behaviors as a human being, and he took his work suspended. In 1950, Claude Shannon introduced the idea that computers could play chess. Work on artificial intelligence continued slowly until the early 1960s. The emergence of artificial intelligence officially in history dates back to 1956. In 1956, a conference artificial intelligence session at Dartmouth College was introduced for the first time. Marvin Minsky stated in his book "Stormed Search for Artificial Intelligence " that "the problem of artificial intelligence modelling within a generation will be solved ". The first artificial intelligence applications were introduced during this period. These applications are based on logic theorems and chess game. The programs developed during this period were distinguished from the geometric forms used in the intelligence tests; which has led to the idea that intelligent computers can be created.

Milestones for AI History

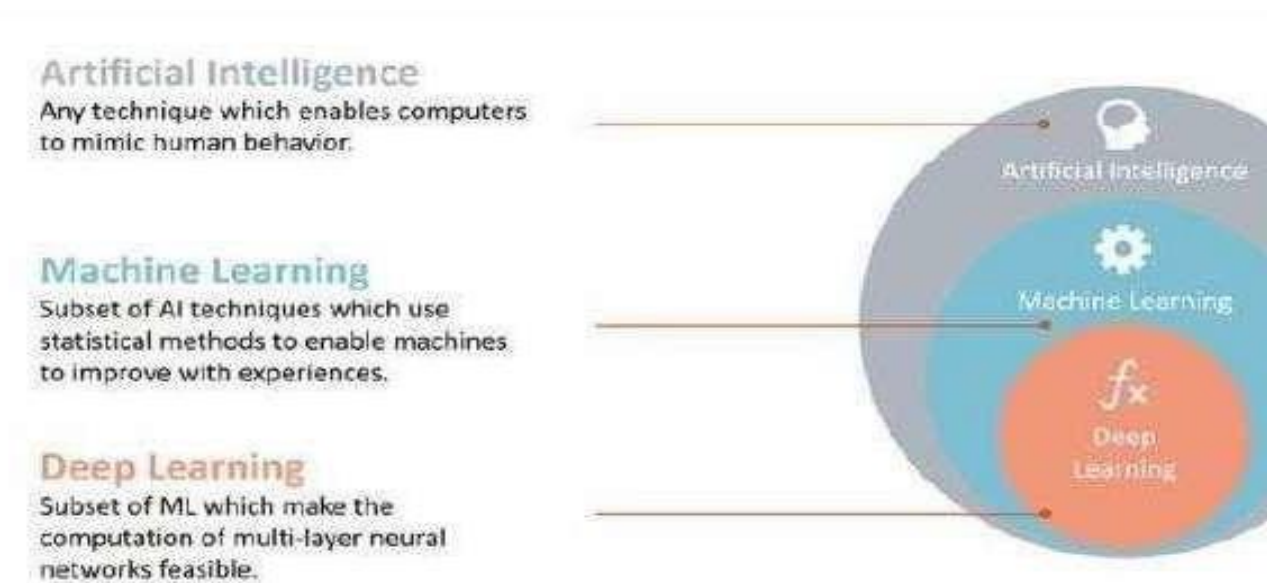
In 1950, Alan Turing created a test to determine whether a machine was intelligent. This test shows the intelligence given to computers. The intelligence level of the machines that passed the test at that time was considered adequate. LISP (List Processing Language), developed by John McCarthy in 1957, is a functional programming language. developed for artificial intelligence. One of the rather old and powerful programming languages, LISP is a language that allows you to create flexible programs that represent basic operations with list structure. Between 1965 and 1970, it could be called a dark period for artificial intelligence. The developments on artificial intelligence in this period are too few to be tested. The hasty and optimistic attitude due to the unrealistic expectations that have emerged has led to the idea that it will be easy to uncover the machines with intelligence. But this period was named as a dark period on behalf of artificial intelligence because it did not succeed with the idea of creating intelligent machines by simply uploading data. Between 1970 and 1975, artificial intelligence gained momentum. Thanks to the success achieved in artificial intelligence systems that have been developed and developed on subjects such as disease diagnosis, the basis of today's artificial intelligence has been established.

During the period 1975-1980 they developed the idea that they could benefit artificial intelligence through other branches of science such as psychology. Artificial Intelligence began to be used in large projects with practical applications in the 1980s. The next time the daylight is passed, the artificial intelligence has been adapted to solve real life problems. Even when the needs of users are already met with traditional

methods, the use of artificial intelligence has reached to a much wider range thanks to more economical software and tools.

The concept of Machine Intelligence emerging with various code algorithms and data studies reveals that all the technological devices produced from the first computers to today's smart phones are developed on the basis of people. The artificial intelligence, which was developed very slowly in the old periods but so important steps as the day- to-day, reveals how much progress has been made with the emergence of gifted robots today.

What is an Artificial Intelligence?



According to the father of Artificial Intelligence, John McCarthy, it is “The science and engineering of making intelligent machines, especially intelligent computer programs”.

Artificial intelligence is the general name of the technology for the development of machines, which are start at the beginning and could not be maintained as expected on both sides. In the symbolic artificial intelligence studies, robots cannot give exactly the expected responses and answers to the questions of the people, whereas on the cybernetic artificial intelligence side, the artificial neural networks do not give the expectation and the works on the two sides cannot be successful with literally. Artificial intelligence has led to the emergence of specialized artificial intelligence exercises that will continue with a single purpose, rather than different branches and minds, after failures in Symbolic and Cybernetic artificial intelligence studies developed on different sides. While the concept of artificial intelligence has stimulated artificial intelligence studies, the fact that artificial intelligence

products do not have enough knowledge about what is being worked on has brought about various problems. However, the artificial intelligence developers who brought rational solutions to the problems that have arisen, have reached to commercial level of artificial intelligence, and the artificial intelligence industry that emerged in the coming periods has shown that the achievement of successful works is achieved with billion-dollar billets.

Recent developments in artificial intelligence studies have revealed the importance of language. As anthropology, Human Science studies show, people have begun to hold the language in front of the artificial intelligence studies in recent years because people think with language and put out various functions. Later, a number of artificial intelligences marking languages appeared with the language studies that were behind those who supported Symbolic Artificial Intelligence studies. Today, artificial intelligence studies carried out by Symbolic artificial intelligence developers have benefited from artificial intelligence languages and have made it possible to show even robots that can speak.

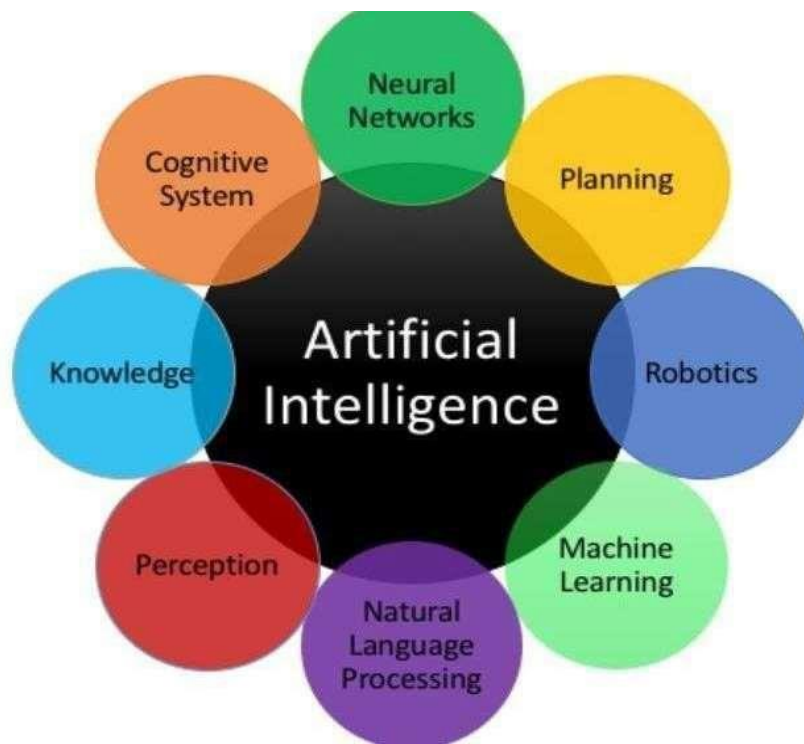
Goals of AI

To Create Expert Systems – The systems which exhibit intelligent behaviour, learn, demonstrate, explain, and advice its users.

To Implement Human Intelligence in Machines – Creating systems that understand, think, learn, and behave like humans.

The overall research goal of artificial intelligence is to create technology that allows computers and machines to function in an intelligent manner. The general problem of simulating (or creating) intelligence has been broken down into sub-problems.

These consist of particular traits or capabilities that researchers expect an intelligent system to display. The traits described below have received the most attention. Erik Sandwell emphasizes planning and learning that is relevant and applicable to the given situation.



- Reasoning, problem solving: Early researchers developed algorithms that imitated step-by-step reasoning that humans use when they solve puzzles or make logical deductions. By the late 1980s and 1990s, AI research had developed methods for dealing with uncertain or incomplete information, employing concepts from probability and economics. For difficult problems, algorithms can require enormous computational resources—most experience a “combinatorial explosion”: the amount of memory or computer time required becomes astronomical for problems of a certain size. The search for more efficient problem-solving algorithms is a high priority.

- Knowledge representation: Knowledge representation and knowledge engineering are central to AI research. Many of the problems machines are expected to solve will require extensive knowledge about the world. Among the things that AI needs to represent are: objects, properties, categories and relations between objects; situations, events, states and time; causes and effects; knowledge about knowledge (what we know about what other people know); and many other, less well researched domains. A representation of “what exists” is an ontology: the set of objects, relations, concepts and so on that the machine knows about. The most general are called upper ontologies, which attempt to provide a foundation for all other knowledge.

- **Planning:** Intelligent agents must be able to set goals and achieve them. They need a way to visualize the future—a representation of the state of the world and be able to make predictions about how their actions will change it—and be able to make choices that maximize the utility (or “value”) of available choices. In classical planning problems, the agent can assume that it is the only system acting in the world, allowing the agent to be certain of the consequences of its actions. However, if the agent is not the only actor, then it requires that the agent can reason under uncertainty. This calls for an agent that can not only assess its environment and make predictions, but also evaluate its predictions and adapt based on its assessment.

- **Learning:** Machine learning, a fundamental concept of AI research since the field’s inception, is the study of computer algorithms that improve automatically through experience. Unsupervised learning is the ability to find patterns in a stream of input. Supervised learning includes both classification and numerical regression. Classification is used to determine what category something belongs in, after seeing a number of examples of things from several categories. Regression is the attempt to produce a function that describes the relationship between inputs and outputs and predicts how the outputs should change as the inputs change.

- **General intelligence:** Many researchers think that their work will eventually be incorporated into a machine with artificial general intelligence, combining all the skills mentioned above and even exceeding human ability in most or all these areas. A few believe that anthropomorphic features like artificial consciousness or an artificial brain may be required for such a project.

Application of Artificial Intelligence (AI)

AI has been dominant in various fields such as:-

Gaming – AI plays crucial role in strategic games such as chess, poker, tic-tac-toe, etc., where machine can think of large number of possible positions based on heuristic knowledge.

Natural Language Processing – It is possible to interact with the computer that understands natural language spoken by humans.

Expert Systems – There are some applications which integrate machine, software, and special information to impart reasoning and advising. They provide explanation and advice to the users.

Vision Systems – These systems understand, interpret, and comprehend visual input on the computer.

For example: A spying aeroplane takes photographs, which are used to figure out spatial information or map of the areas. Doctors use clinical expert system to diagnose the patient. Police use computer software that can recognize the face of criminal with the stored portrait made by forensic artist.

Speech Recognition – Some intelligent systems are capable of hearing and comprehending the language in terms of sentences and their meanings while a human talks to it. It can handle different accents, slang words, noise in the background, change in human's voice due to cold, etc.

Handwriting Recognition – The handwriting recognition software reads the text written on paper by a pen or on screen by a stylus. It can recognize the shapes of the letters and convert it into editable text.

Intelligent Robots – Robots are able to perform the tasks given by a human. They have sensors to detect physical data from the real world such as light, heat, temperature, movement, sound, bump, and pressure. They have efficient processors, multiple sensors and huge memory, to exhibit intelligence. In addition, they are capable of learning from their mistakes and they can adapt to the new environment.

Machine Learning

Abstract Machine

Learning has always been an integral part of artificial intelligence and its methodology has evolved in the concert with the major concerns in the field. In response of increasing of ever increasing of the volume of knowledge in modern AI systems, many researchers have turned their attention into machine learning as to overcome the knowledge acquisition bottleneck.

Together with many other disciplines, machine learning methods have been widely employed in bioinformatics. The difficulties and cost of biological analyses have led to the development of sophisticated machine learning approaches for this application area. In this chapter, we first review the fundamental concepts of machine learning such as feature assessment, unsupervised versus supervised learning and types of classification. Then, we point out the main issues of designing machine learning experiments and their performance evaluation. Finally, we introduce some supervised learning methods.

Introduction to Machine Learning

Machines have come a long way since the Industrial Revolution. They continue to fill factory floors and manufacturing plants, but now their capabilities extend beyond manual activities to cognitive tasks that, until recently, only humans were capable of performing. Judging song competitions, driving automobiles, and mopping the floor with professional chess players are three examples of the specific complex tasks machines are now capable of simulating. But their remarkable feats trigger fear among some observers. Part of this fear nestles on the neck of survivalist insecurities, where it provokes the deep-seated question of what if? What if intelligent machines turn on us in a struggle of the fittest? What if intelligent machines produce offspring with capabilities that humans never intended to impart to machines? What if the legend of the singularity is true? The other notable fear is the threat to job security, and if you're a truck driver or an accountant, there is a valid reason to be worried. According to the British Broadcasting Company's (BBC) interactive online resource Will a robot take my job?, professions such as bar worker (77%), waiter (90%), chartered accountant (95%), receptionist (96%), and taxi driver (57%) each have a high chance of becoming automated by the year 2035. [1] But research on planned job automation and crystal ball gazing with respect to the future evolution of machines and artificial intelligence

(AI) should be read with a pinch of skepticism. AI technology is moving fast, but broad adoption is still an uncharted path fraught with known and unforeseen challenges. Delays and other obstacles are inevitable. Nor is machine learning a simple case of flicking a switch and asking the machine to predict the outcome of the Super Bowl and serve you a delicious martini. Machine learning is far from what you would call an out-of-the-box solution.

Machines operate based on statistical algorithms managed and overseen by skilled individuals—known as data scientists and machine learning engineers. This is one labor market where job opportunities are destined for growth but where, currently, supply is struggling to meet demand. Industry experts lament that one of the biggest obstacles delaying the progress of AI is the inadequate supply of professionals with the necessary expertise and training.

According to Charles Green, the Director of Thought Leadership at Belatrix Software: “It’s a huge challenge to find data scientists, people with machine learning experience, or people with the skills to analyze and use the data, as well as those who can create the algorithms required for machine learning. Secondly, while the technology is still emerging, there are many ongoing developments. It’s clear that AI is a long way from how we might imagine it.”

To build and program intelligent machines, you must first understand classical statistics. Algorithms derived from classical statistics contribute the metaphorical blood cells and oxygen that power machine learning. Layer upon layer of linear regression, k-nearest neighbors, and random forests surge through the machine and drive their cognitive abilities. Classical statistics is at the heart of machine learning and many of these algorithms are based on the same statistical equations you studied in high school. Indeed, statistical algorithms were conducted on paper well before machines ever took on the title of artificial intelligence.

WHAT IS MACHINE LEARNING?

In 1959, IBM published a paper in the IBM Journal of Research and Development with an, at the time, obscure and curious title. Authored by IBM's Arthur Samuel, the paper investigated the use of machine learning in the game of checkers "to verify the fact that a computer can be programmed so that it will learn to play a better game of checkers than can be played by the person who wrote the program." [3] Although it was not the first publication to use the term "machine learning" per se, Arthur Samuel is widely considered as the first person to coin and define machine learning in the form we now know today. Samuel's landmark journal submission, *Some Studies in Machine Learning Using the Game of Checkers*, is also an early indication of homo sapiens' determination to impart our own system of learning to man-made machines.

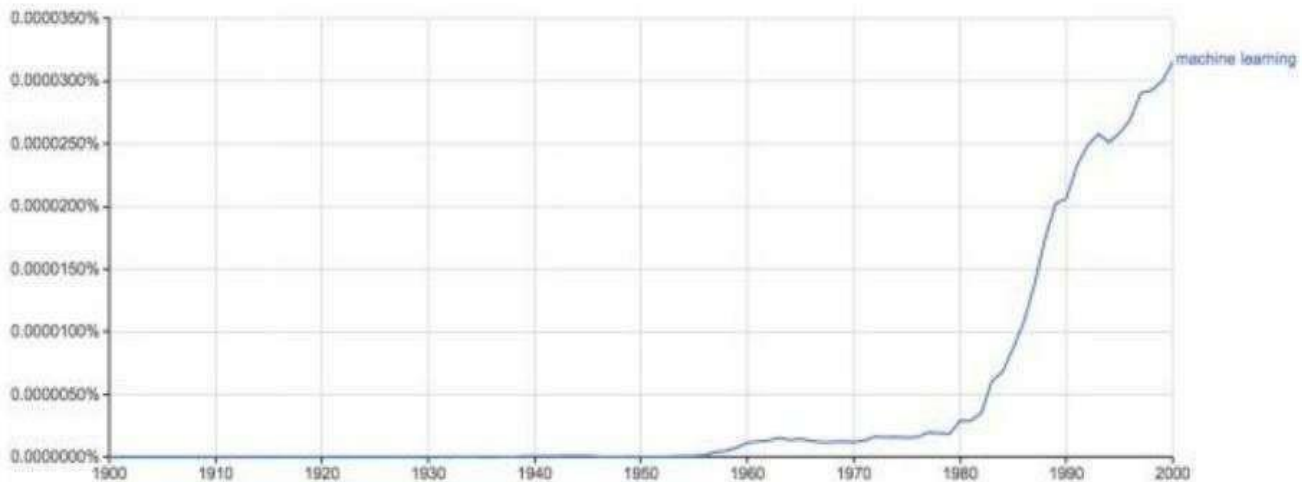


Figure 1: Historical mentions of "machine learning" in published books. Source: Google Ngram Viewer, 2017

Arthur Samuel introduces machine learning in his paper as a subfield of computer science that gives computers the ability to learn without being explicitly programmed. [4] Almost six decades later, this definition remains widely accepted. Although not directly mentioned in Arthur Samuel's definition, a key feature of machine learning is the concept of self-learning. This refers to the application of statistical modelling to detect patterns and improve performance based on data and empirical information; all without direct programming commands. This is what Arthur Samuel described as the ability to learn without being explicitly programmed. But he doesn't infer that machines formulate decisions with no upfront programming. On the contrary, machine learning is heavily dependent on computer programming. Instead, Samuel observed that machines don't require a direct input command to perform a set task but rather input data.

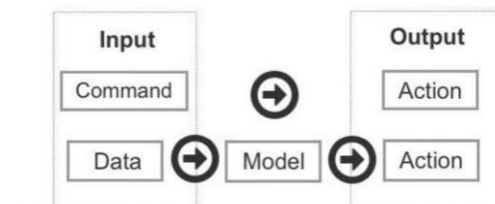


Figure 2: Comparison of Input Command vs Input Data

Training & Test Data

In machine learning, data is split into training data and test data. The first split of data, i.e. the initial reserve of data you use to develop your model, provides the training data. In the spam email detection example, false positives similar to the PayPal auto-response might be detected from the training data. New rules or modifications must then be added, e.g., email notifications issued from the sending address “payments@paypal.com” should be excluded from spam filtering. After you have successfully developed a model based on the training data and are satisfied with its accuracy, you can then test the model on the remaining data, known as the test data. Once you are satisfied with the results of both the training data and test data, the machine learning model is ready to filter incoming emails and generate decisions on how to categorize those incoming messages. The difference between machine learning and traditional programming may seem trivial at first, but it will become clear as you run through further examples and witness the special power of self-learning in more nuanced situations. The second important point to take away from this chapter is how machine learning fits into the broader landscape of data science and computer science. This means understanding how machine learning interrelates with parent fields and sister disciplines. This is important, as you will encounter these related terms when searching for relevant study materials—and you will hear them mentioned ad nauseam in introductory machine learning courses. Relevant disciplines can also be difficult to tell apart at first glance, such as “machine learning” and “data mining.” Let’s begin with a high-level introduction. Machine learning, data mining, computer programming, and most relevant fields (excluding classical statistics) derive first from computer science, which encompasses everything related to the design and use of computers. Within the all-encompassing space of computer science is the next broad field: data science. Narrower than computer science, data science comprises methods and systems to extract knowledge and insights from data through the use of computers.

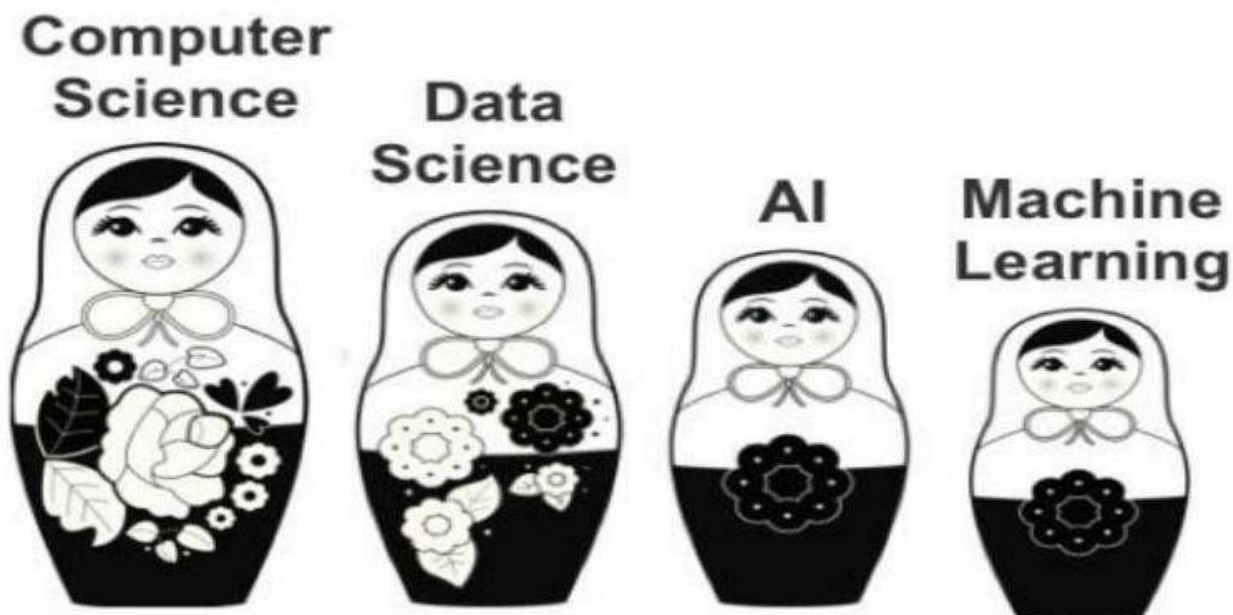


Figure 3: The lineage of machine learning represented by a row of Russian matryoshka dolls

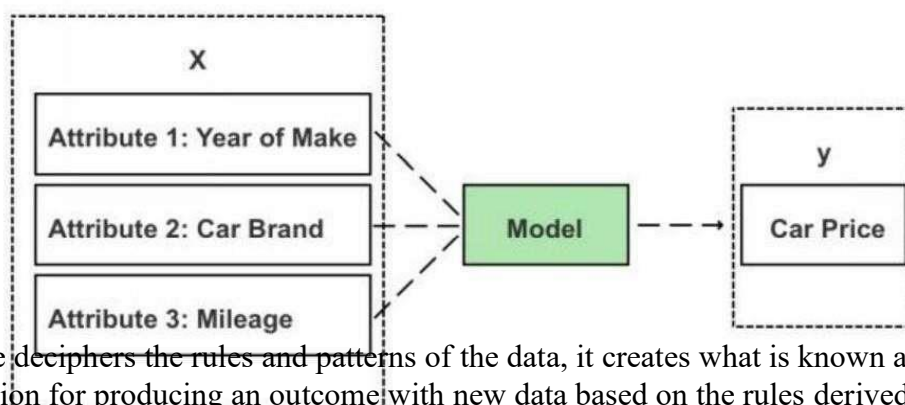
ML Categories

Machine learning incorporates several hundred statistical-based algorithms and choosing the right algorithm or combination of algorithms for the job is a constant challenge for anyone working in this field. But before we examine specific algorithms, it is important to understand the three overarching categories of machine learning. These three categories are supervised, unsupervised, and reinforcement.

Supervised Learning

As the first branch of machine learning, supervised learning concentrates on learning patterns through connecting the relationship between variables and known outcomes and working with labeled datasets. Supervised learning works by feeding the machine sample data with various features (represented as “X”) and the correct value output of the data (represented as “y”). The fact that the output and feature values are known qualifies the dataset as “labeled.” The algorithm then deciphers patterns that exist in the data and creates a model that can reproduce the same underlying rules with new data. For instance, to predict the market rate for the purchase of a used car, a supervised algorithm can formulate predictions by analyzing the relationship between car attributes (including the year of make, car brand, mileage, etc.) and the selling price of other cars sold based on historical data.

Given that the supervised algorithm knows the final price of other cars sold, it can then work backward to determine the relationship between the characteristics of the car and its value



After the machine deciphers the rules and patterns of the data, it creates what is known as a model: an algorithmic equation for producing an outcome with new data based on the rules derived from the training data. Once the model is prepared, it can be applied to new data and tested for accuracy. After the model has passed both the training and test data stages, it is ready to be applied and used in the real world. In Chapter 13, we will create a model for predicting house values where y is the actual house price and X are the variables that impact y, such as land size, location, and the number of rooms. Through supervised learning, we will create a rule to predict y (house value) based on the given values of various variables (X).

Examples of supervised learning algorithms include regression analysis, decision trees, k- nearest neighbours, neural networks, and support vector machines.

Unsupervised Learning

In the case of unsupervised learning, not all variables and data patterns are classified. Instead, the machine must uncover hidden patterns and create labels through the use of unsupervised learning algorithms. The k-means clustering algorithm is a popular example of unsupervised learning. This algorithm groups data points that are found to possess similar features

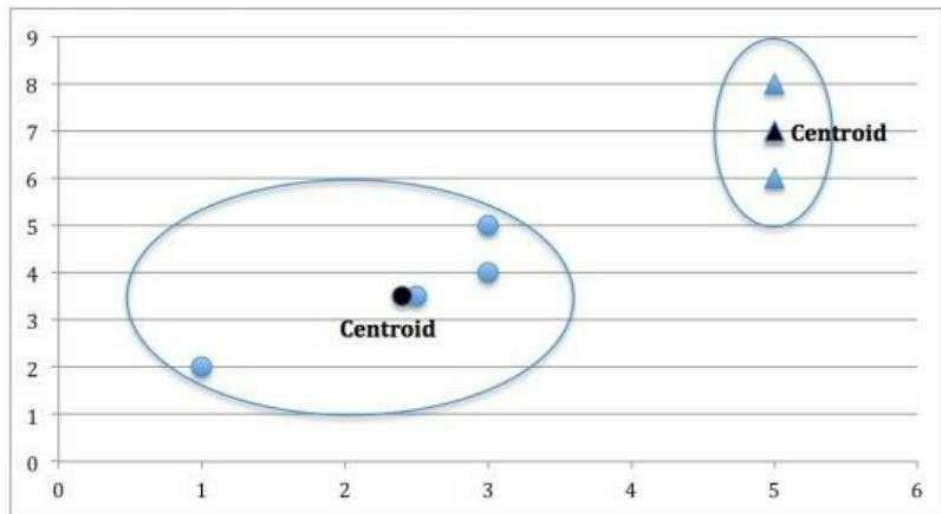


Figure 1: Example of k-means clustering, a popular unsupervised learning technique

If you group data points based on the purchasing behavior of SME (Small and Medium-sized Enterprises) and large enterprise customers, for example, you are likely to see two clusters emerge. This is because SMEs and large enterprises tend to have disparate buying habits. When it comes to purchasing cloud infrastructure, for instance, basic cloud hosting resources and a Content Delivery Network (CDN) may prove sufficient for most SME customers.

The advantage of unsupervised learning is it enables you to discover patterns in the data that you were unaware existed—such as the presence of two major customer types. Clustering techniques such as k-means clustering can also provide the springboard for conducting further analysis after discrete groups have been discovered.

In industry, unsupervised learning is particularly powerful in fraud detection—where the most dangerous attacks are often those yet to be classified. One real-world example is DataVisor, who essentially built their business model based on unsupervised learning.

Reinforcement Learning

Reinforcement learning is the third and most advanced algorithm category in machine learning. Unlike supervised and unsupervised learning, reinforcement learning continuously improves its model by leveraging feedback from previous iterations. This is different to supervised and unsupervised learning, which both reach an indefinite endpoint after a model is formulated from the training and test data segments. Reinforcement learning can be complicated and is probably best explained through an analogy to a video game. As a player progresses through the virtual space of a game, they learn the value of various actions under different conditions and become more familiar with the field of play. Those learned values then inform and influence a player's subsequent behavior and their performance immediately improves based on their learning and past experience. Reinforcement learning is very similar, where algorithms are set to train the model through continuous learning. A standard reinforcement learning model has measurable performance criteria where outputs are not tagged— instead, they are graded. In the case of self-driving vehicles, avoiding a crash will allocate a positive score and in the case of chess, avoiding defeat will likewise receive a positive score. A specific algorithmic example of reinforcement learning is Q-learning. In Qlearning, you start with a set environment of states, represented by the symbol 'S'. In the game Pac-Man, states could be the challenges, obstacles or pathways that exist in the game. There may exist a wall to the left, a ghost to the right, and a power pill above—each representing different states. The set of possible actions to respond to these states is referred to as "A." In the case of Pac-Man, actions are limited to left, right, up, and down movements, as well as multiple combinations thereof. The third important symbol is "Q." Q is the starting value and has an initial value of "0." As Pac-Man explores the space inside the game, two main things will happen: 1) Q drops as negative things occur after a given state/action 2) Q increases as positive things occur after a given state/action In Q-learning, the machine will learn to match the action for a given state that generates or maintains the highest level of Q. It will learn initially through the process of random movements (actions) under different conditions (states). The machine will record its results (rewards and penalties) and how they impact its Q level and store those values to inform and optimize its future actions. While this sounds simple enough, implementation is a much more difficult task and beyond the scope of an absolute beginner's introduction to machine learning. Reinforcement learning algorithms aren't covered in this book, however, I will leave you with a link to a more comprehensive explanation of reinforcement learning and Q-learning following the Pac-Man scenario.

The Machine Learning Tool Box

A handy way to learn a new subject area is to map and visualize the essential materials and tools inside a toolbox. If you were packing a toolbox to build websites, for example, you would first pack a selection of programming languages. This would include frontend languages such as HTML, CSS, and JavaScript, one or two backend programming languages based on personal preferences, and of course, a text editor. You might throw in a website builder such as WordPress and then have another compartment filled with web hosting, DNS, and maybe a few domain names that you've recently purchased. This is not an extensive inventory, but from this general list, you can start to gain a better appreciation of what tools you need to master in order to become a successful website developer. Let's now unpack the toolbox for machine learning. Compartment 1: Data In the first compartment is your data. Data consists of the input variables needed to form a prediction. Data comes in many forms, including structured and non-structured data. As a beginner, it is recommended that you start with structured data. This means that the data is defined and labeled (with schema) in a table, as shown here:

A tabular (table-based) dataset contains data organized in rows and columns. In each column is a feature. A feature is also known as a variable, a dimension or an attribute—but they all mean the

Date	Bitcoin Price	No. of Days Transpired
19-05-2015	234.31	1
14-01-2016	431.76	240
09-07-2016	652.14	417
15-01-2017	817.26	607
24-05-2017	2358.96	736

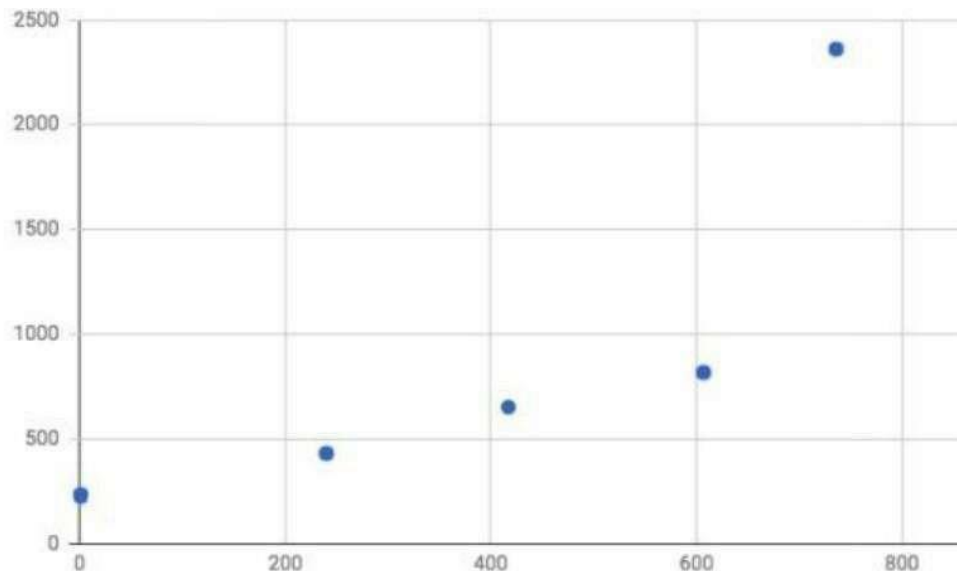
same thing. Each individual row represents a single observation of a given feature/variable. Rows are sometimes referred to as a case or value, but in this book, we will use the term “row.”

Matrices			
Vector			
	Feature 1	Feature 2	Feature 3
Row 1			
Row 2			
Row 3			
Row 4			

Each column is known as a vector. Vectors store your X and y values and multiple vectors (columns) are commonly referred to as matrices. In the case of supervised learning, y will already exist in your dataset and be used to identify patterns in relation to independent variables (X). The y values are commonly expressed in the final column, as shown in Figure 2

	Vector	Matrices		
	Maker (X)	Year (X)	Model (X)	Price (y)
Row 1				
Row 2				
Row 3				
Row 4				

Next, within the first compartment of the toolbox is a range of scatterplots, including 2-D, 3-D, and 4-D plots. A 2-D scatterplot consists of a vertical axis (known as the y-axis) and a horizontal axis (known as the x-axis) and provides the graphical canvas to plot a series of dots, known as data points. Each data point on the scatterplot represents one observation from the dataset, with X values plotted on the x-axis and y values plotted on the y-axis.



	Independent Variable (X)	Dependent Variable (y)
Row 1	1	243.31
Row 2	240	431.76
Row 3	417	653.14
Row 4	607	817.26
Row 5	736	2358.96

Figure 3: Example of a 2-D scatterplot. X represents days passed since the recording of Bitcoin prices and y represents recorded Bitcoin price.

Compartment 2:

Infrastructure The second compartment of the toolbox contains your infrastructure, which consists of platforms and tools to process data. As a beginner to machine learning, you are likely to be using a web application (such as Jupyter Notebook) and a programming language like Python. There are then a series of machine learning libraries, including NumPy, Pandas, and Scikit-learn that are compatible with Python. Machine learning libraries are a collection of pre-compiled programming routines frequently used in machine learning. You will also need a machine from which to work, in the form of a computer or a virtual server. In addition, you may need specialized libraries for data visualization such as Seaborn and Matplotlib, or a standalone software program like Tableau, which supports a range of visualization techniques including charts, graphs, maps, and other visualizations. With your infrastructure sprayed out across the table (hypothetically of course), you are now ready to get to work building your first machine learning model. The first step is to crank up your computer. Laptops and desktop computers are both suitable for working with smaller datasets. You will then need to install a programming environment, such as Jupyter Notebook, and a programming language, which for most beginners is Python. Python is the most widely used programming language for machine learning because: a) It is easy to learn and operate, b) It is compatible with a range of machine learning libraries, and c) It can be used for related tasks, including data collection (web scraping) and data piping (Hadoop and Spark). Other go-to languages for machine learning include C and C++. If you're proficient with C and C++ then it makes sense to stick with what you already know. C and C++ are the default programming languages for advanced machine learning because they can run directly on a GPU (Graphical Processing Unit). Python needs to be converted first before it can run on a GPU, but we will get to this and what a GPU is later in the chapter. Next, Python users will typically install the following libraries: NumPy, Pandas, and Scikit-learn. NumPy is a free and open-source library that allows you to efficiently load and work with large datasets, including managing matrices. Scikit-learn provides access to a range of popular algorithms, including linear regression, Bayes' classifier, and support vector machines. Finally, Pandas enables your data to be represented on a virtual spreadsheet that you can control through code. It shares many of the same features as Microsoft Excel in that it allows you to edit data and perform calculations. In fact, the name Pandas derives from the term "panel data," which refers to its ability to create a series of panels, similar to "sheets" in Excel. Pandas is also ideal for importing and extracting data from CSV files.

```
# Preview dataframe
df.head(n=5)
```

Out[31]:

	Suburb	Address	Rooms	Type	Price	Method	SellerG	Date	Distance	Postcode	Bedroom2	Bathroom	Car	Landsize
0	Abbotsford	68 Studley St	2.0	h	NaN	SS	Jellis	3/09/2016	2.5	3067.0	2.0	1.0	1.0	126.0
1	Abbotsford	85 Turner St	2.0	h	1480000.0	S	Biggin	3/12/2016	2.5	3067.0	2.0	1.0	1.0	202.0
2	Abbotsford	25 Bloomburg St	2.0	h	1035000.0	S	Biggin	4/02/2016	2.5	3067.0	2.0	1.0	0.0	196.0
3	Abbotsford	18/859 Victoria St	3.0	u	NaN	VB	Rounds	4/02/2016	2.5	3067.0	3.0	2.0	1.0	0.0
4	Abbotsford	5 Charles St	3.0	h	1466000.0	SP	Biggin	4/03/2017	2.5	3067.0	3.0	2.0	0.0	134.0

Figure 4: Previewing a table in Jupyter Notebook using Pandas

In summary, users can draw on these three libraries to:

1) Load and work with a dataset via NumPy.

2) Clean up and perform calculations on data, and extract data from CSV files with Pandas.

3) Implement algorithms with Scikit-learn. For students seeking alternative programming options (beyond Python, C, and C++), other relevant programming languages for machine learning include R, MATLAB, and Octave. R is a free and open-source programming language optimized for mathematical operations, and conducive to building matrices and statistical functions, which are built directly into the language libraries of R. Although R is commonly used for data analytics and data mining, R supports machine learning operations as well. MATLAB and Octave are direct competitors to R. MATLAB is a commercial and propriety programming language. It is strong in regards to solving algebraic equations and is also a quick programming language to learn. MATLAB is widely used in electrical engineering, chemical engineering, civil engineering, and aeronautical engineering. However, computer scientists and computer engineers tend not to rely on MATLAB as heavily and especially in recent times. In machine learning, MATLAB is more often used in academia than in industry. Thus, while you may see MATLAB featured in online courses, and especially on Coursera, this is not to say that it's commonly used in the wild. If, however, you're coming from an engineering background, MATLAB is certainly a logical choice. Lastly, Octave is essentially a free version of MATLAB developed in response to MATLAB by the open-source community.

Compartment 3:

Algorithms Now that the machine learning environment is set up and you've chosen your programming language and libraries, you can next import your data directly from a CSV file. You can find hundreds of interesting datasets in CSV format from [kaggle.com](https://www.kaggle.com/). After registering as a member of their platform, you can download a dataset of your choice. Best of all, Kaggle datasets are free and there is no cost to register as a user. The dataset will download directly to your computer as a CSV file, which means you can use Microsoft Excel to open and even perform basic algorithms such as linear regression on your dataset. Next is the third and final compartment that stores the algorithms. Beginners will typically start off by using simple supervised learning algorithms such as linear regression, logistic regression, decision trees, and k-nearest neighbors. Beginners are also likely to apply unsupervised learning in the form of kmeans clustering and dimensionality algorithms.

DATA SCRUBBING

Much like many categories of fruit, datasets nearly always require some form of upfront cleaning and human manipulation before they are ready to digest. For machine learning and data science more broadly, there are a vast number of techniques to scrub data. Scrubbing is the technical process of refining your dataset to make it more workable. This can involve modifying and sometimes removing incomplete, incorrectly formatted, irrelevant or duplicated data. It can also entail converting text-based data to numerical values and the redesigning of features. For data practitioners, data scrubbing usually demands the greatest application of time and effort.

Feature Selection

To generate the best results from your data, it is important to first identify the variables most relevant to your hypothesis. In practice, this means being selective about the variables you select to design your model. Rather than creating a four-dimensional scatterplot with four features in the model, an opportunity may present to select two highly relevant features and build a two-dimensional plot that is easier to interpret. Moreover, preserving features that do not correlate strongly with the outcome value can, in fact, manipulate and derail the model's accuracy. Consider the following table excerpt downloaded from kaggle.com documenting dying languages.

Name in English	Name in Spanish	Countries	Country Code
South Italian	Napolitano-calabres	Italy	ITA
Sicilian	Siciliano	Italy	ITA
Low Saxon	Bajo Sajón	Germany, Denmark, Netherlands, Poland, Russian Federation	DEU, DNK, NLD, POL, RUS
Belarusian	Bielorruso	Belarus, Latvia, Lithuania, Poland, Russian Federation, Ukraine	BRB, LVA, LTU, POL, RUS, UKR
Lombard	Lombardo	Italy, Switzerland	ITA, CHE
Romani	Romani	Albania, Germany, Austria, Belarus, Bosnia and Herzegovina, Bulgaria, Croatia, Estonia, Finland, France, Greece, Hungary, Italy, Latvia, Lithuania, The former Yugoslav Republic of Macedonia, Netherlands, Poland, Romania, United Kingdom of Great Britain and Northern Ireland, Russian Federation, Slovakia, Slovenia, Switzerland, Czech Republic, Turkey, Ukraine, Serbia, Montenegro	ALB, DEU, AUT, BRB, BIH, BGR, HRV, EST, FIN, FRA, GRC, HUN, ITA, LVA, LTU, MKD, NLD, POL, ROU, GBR, RUS, SVK, SVN, CHE, CZE, TUR, UKR, SRB, MNE
Yiddish	Yiddish	Israel	ISR
Gondi	Gondi	India	IND

Let's say our goal is to identify variables that lead to a language becoming endangered. Based on this goal, it's unlikely that a language's "Name in Spanish" will lead to any relevant insight. We can therefore go ahead and delete this vector (column) from the dataset. This will help to prevent overcomplication and potential inaccuracies, and will also improve the overall processing speed of the model. Secondly, the dataset holds duplicate information in the form of separate vectors for "Countries" and "Country Code." Including both of these vectors doesn't provide any additional insight; hence, we can choose to delete one and retain the other.

Another method to reduce the number of features is to roll multiple features into one. In the next table, we have a list of products sold on an e-commerce platform. The dataset comprises four buyers and eight products. This is not a large sample size of buyers and products—due in part to the spatial limitations of the book format. A real-life e-commerce platform would have many more columns to work with, but let's go ahead with this example.

	Protein Shake	Nike Sneakers	Adidas Boots	Fitbit	Powerade	Protein Bar	Fitness Watch	Vitamins
Buyer 1	1	1	0	1	0	5	1	0
Buyer 2	0	0	0	0	0	0	0	1
Buyer 3	3	0	1	0	5	0	0	0
Buyer 4	1	1	0	0	10	1	0	0

In order to analyze the data in a more efficient way, we can reduce the number of columns by merging similar features into fewer columns. For instance, we can remove individual product names and replace the eight product items with a lower number of categories or subtypes. As all product items fall under the single category of "fitness," we will sort by product subtype and compress the columns from eight to three. The three newly created product subtype columns are "Health Food," "Apparel," and "Digital."

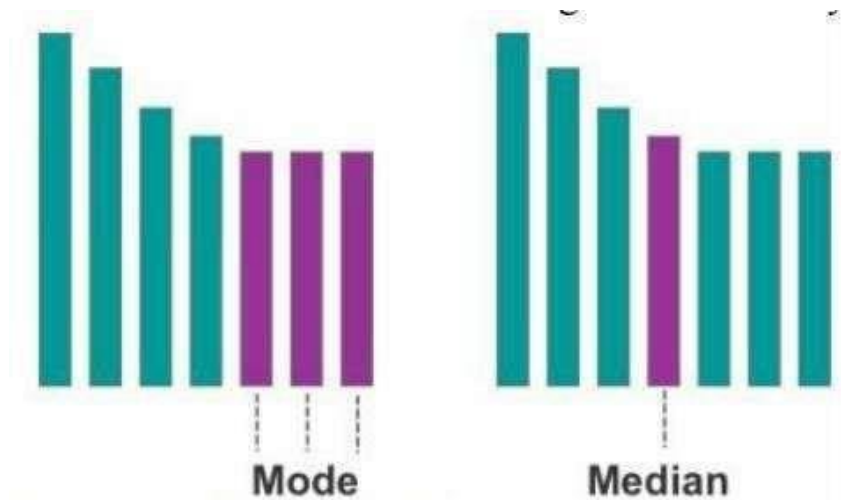
	Health Food	Apparel	Digital
Buyer 1	6	1	2
Buyer 2	1	0	0
Buyer 3	8	1	0
Buyer 4	12	1	0

This enables us to transform the dataset in a way that preserves and captures information using fewer variables. The downside to this transformation is that we have less information about relationships between specific products. Rather than recommending products to users according to other individual products, recommendations will instead be based on relationships between product subtypes.

Missing Data

Dealing with missing data is never a desired situation. Imagine unpacking a jigsaw puzzle that you discover has five percent of its pieces missing. Missing values in a dataset can be equally frustrating and will ultimately interfere with your analysis and final predictions.

There are, however, strategies to minimize the negative impact of missing data. One approach is to approximate missing values using the mode value. The mode represents the single most common variable value available in the dataset. This works best with categorical and binary variable types.



The second approach to manage missing data is to approximate missing values using the median value, which adopts the value(s) located in the middle of the dataset. This works best with integers (whole numbers) and continuous variables (numbers with decimals). As a last resort, rows with missing values can be removed altogether. The obvious downside to this approach is having less data to analyse and potentially fewer comprehensive results.

SETTING UP YOUR DATA

Once you have cleaned your dataset, the next job is to split the data into two segments for testing and training. It is very important not to test your model with the same data that you used for training. The ratio of the two splits should be approximately 70/30 or 80/20. This means that your training data should account for 70 percent to 80 percent of the rows in your dataset, and the other 20 percent to 30 percent of rows is your test data. It is vital to split your data by rows and not columns.

		Variable 1	Variable 2	Variable 3
Training Data	Row 1			
	Row 2			
	Row 3			
	Row 4			
	Row 5			
	Row 6			
	Row 7			
Test Data	Row 8			
	Row 9			
	Row 10			

Figure 1: Training and test partitioning of the dataset 70/30

Before you split your data, it is important that you randomize all rows in the dataset. This helps to avoid bias in your model, as your original dataset might be arranged sequentially depending on the time it was collected or some other factor. Unless you randomize your data, you may accidentally omit important variance from the training data that will cause unwanted surprises when you apply the trained model to your test data. Fortunately, Scikit-learn provides a built-in function to shuffle and randomize your data with just one line of code. After randomizing your data, you can begin to design your model and apply that to the training data. The remaining 30 percent or so of data is put to the side and reserved for testing the accuracy of the model. In the case of supervised learning, the model is developed by feeding the machine the training data and the expected output (y). The machine is able to analyze and discern relationships between the features (X) found in the training data to calculate the final output (y). The next step is to measure how well the model actually performs. A common approach to analyzing prediction accuracy is a measure called mean absolute error, which examines each prediction in the model and provides an average error score for each prediction. In Scikit-learn, mean absolute error is found using the `model.predict()` function on X (features). This works by first plugging in the y values from the training dataset and generating a prediction for each row in the dataset. Scikit-learn will compare the predictions of the model to the correct outcome and measure its accuracy. You will know if your model is accurate when the error rate between the training and test dataset is low. This means that the model has learned the dataset's underlying patterns and trends.

Once the model can adequately predict the values of the test data, it is ready for use in the wild. If the model fails to accurately predict values from the test data, you will need to check whether the training and test data were properly randomized. Alternatively, you may need to change the model's hyperparameters. Each algorithm has hyperparameters; these are your algorithm settings. In simple terms, these settings control and impact how fast the model learns patterns and which patterns to identify and analyze.

Cross Validation

Although the training/test data split can be effective in developing models from existing data, a question mark remains as to whether the model will work on new data. If your existing dataset is too small to construct an accurate model, or if the training/test partition of data is not appropriate, this can lead to poor estimations of performance in the wild. Fortunately, there is an effective workaround for this issue. Rather than splitting the data into two segments (one for training and one for testing), we can implement what is known as cross validation. Cross validation maximizes the availability of training data by splitting data into various combinations and testing each specific combination. Cross validation can be performed through two primary methods. The first method is exhaustive cross validation, which involves finding and testing all possible combinations to divide the original sample into a training set and a test set. The alternative and more common method is non-exhaustive cross validation, known as k-fold validation. The k-fold validation technique involves splitting data into k assigned buckets and reserving one of those buckets to test the training model at each round. To perform k-fold validation, data are first randomly assigned to k number of equal sized buckets. One bucket is then reserved as the test bucket and is used to measure and evaluate the performance of the remaining (k-1) buckets.

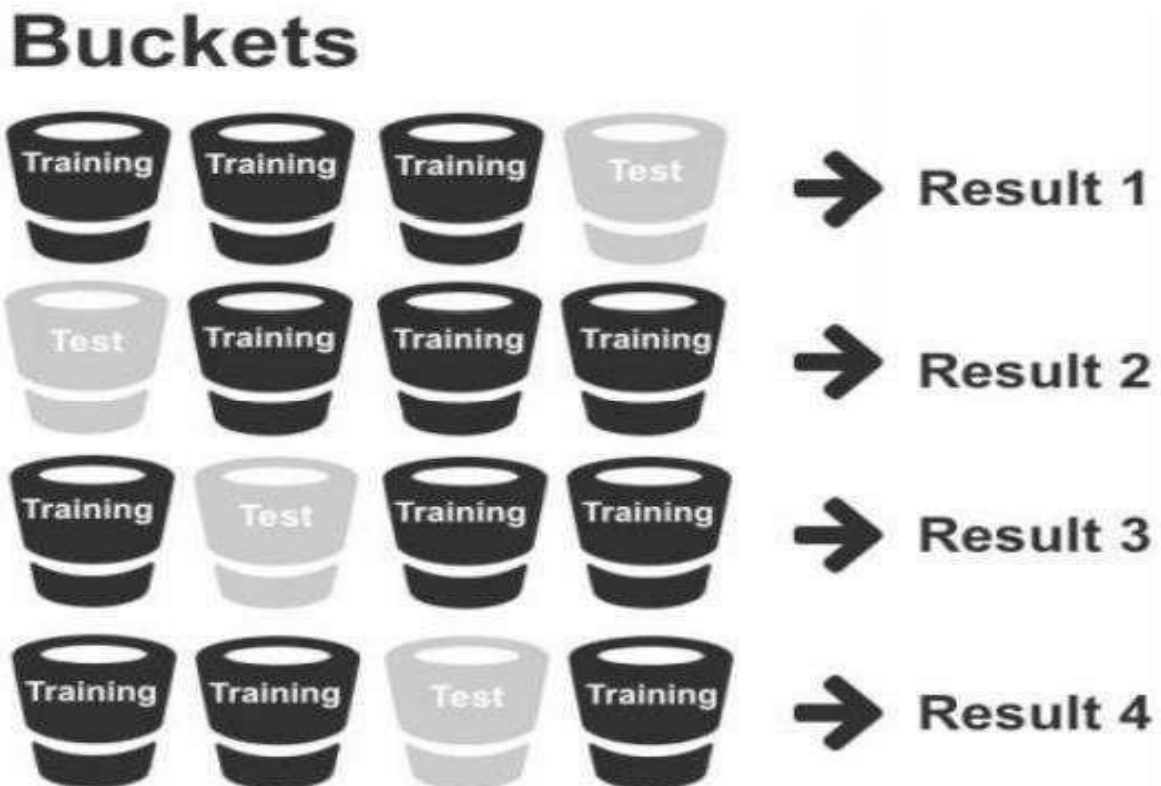


Figure 2: k-fold validation

Machine Learning Algorithms

Supervised Learning

In supervised learning, the target is to infer a function or mapping from training data that is labelled. The training data consist of input vector X and output vector Y of labels or tags. A label or tag from vector Y is the explanation of its respective input example from input vector X . Together they form a training example. In other words, training data comprises training examples. If the labelling does not exist for input vector X , then X is unlabelled data. Why such learning is called supervised learning? The output vector Y consists of labels for each training example present in the training data. These labels for output vector are provided by the supervisor. Often, these supervisors are humans, but machines can also be used for such labelling. Human judgments are more expensive than machines, but the higher error rates in data labelled by machines suggest superior priority of human judgment. The manually labelled data is a precious and reliable resource for supervised learning. However, in some cases, machines can be used for reliable labelling

Example

Table 1.1 demonstrates five unlabeled data examples that can be labeled based on different criteria. The second column of the table titled, “Example judgement for labeling” expresses possible criterion for each data example. The third column describes possible labels after the application of judgment. The fourth column informs which actor can take the role of the supervisor. In all first four cases described in Table 1.1, machines can be used, but their low accuracy rates make their usage questionable. Sentiment analysis, image recognition, and speech detection technologies have made progress in past three decades but there is still a lot of room for improvement before we can equate them with humans’ performance. In the fifth case of tumor detection, even normal humans cannot label the X-ray data, and expensive experts’ services are required for such labeling.

Two groups or categories of algorithms come under the umbrella of supervised learning. They are

1. Regression

2. Classification

Table 1.1 Unlabeled Data Examples along with Labeling Issues

<i>Unlabeled Data Example</i>	<i>Example Judgment for Labeling</i>	<i>Possible Labels</i>	<i>Possible Supervisor</i>
Tweet	Sentiment of the tweet	<i>Positive/negative</i>	Human/machine
Photo	Contains <i>house</i> and <i>car</i>	<i>Yes/No</i>	Human/machine
Audio recording	The word <i>football</i> is uttered	<i>Yes/No</i>	Human/machine
Video	Are weapons used in the video?	<i>Violent/nonviolent</i>	Human/machine
X-ray	Tumor presence in X-ray	<i>Present/absent</i>	Experts/machine

Unsupervised Learning

In unsupervised learning, we lack supervisors or training data. In other words, all what we have is unlabeled data. The idea is to find a hidden structure in this data. There can be a number of reasons for the data not having a label. It can be due to unavailability of funds to pay for manual labeling or the inherent nature of the data itself. With numerous data collection devices, now data is collected at an unprecedented rate. The variety, velocity, and the volume are the dimensions in which Big Data is seen and judged. To get something from this data without the supervisor is important. This is the challenge for today's machine learning practitioner

Semi-Supervised Learning

In this type of learning, the given data are a mixture of classified and unclassified data. This combination of labeled and unlabeled data is used to generate an appropriate model for the classification of data. In most of the situations, labeled data is scarce and unlabeled data is in abundance (as discussed previously in unsupervised learning description). The target of semi-supervised classification is to learn a model that will predict classes of future test data better than that from the model generated by using the labeled data alone. The way we learn is similar to the process of semi-supervised learning. A child is supplied with

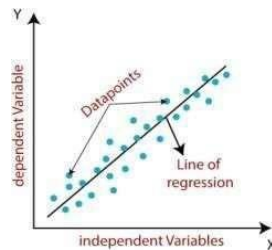
1. Unlabeled data provided by the environment. The surroundings of a child are full of unlabeled data in the beginning.
2. Labeled data from the supervisor. For example, a father teaches his children about the names (labels) of objects by pointing toward them and uttering their names

Reinforcement Learning

The reinforcement learning method aims at using observations gathered from the interaction with the environment to take actions that would maximize the reward or minimize the risk. In order to produce intelligent programs (also called agents), reinforcement learning goes through the following steps: 1. Input state is observed by the agent. 2. Decision making function is used to make the agent perform an action. 3. After the action is performed, the agent receives reward or reinforcement from the environment. 4. The state-action pair information about the reward is stored. Using the stored information, policy for particular state in terms of action can be fine-tuned, thus helping in optimal decision making for our agent.

Linear Regression

Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis.....Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (x) variables, hence called as linear regression.



Here,

Y= Dependent Variable (Target Variable)

X= Independent Variable (predictor Variable)

a_0 = intercept of the line (Gives an additional degree of freedom)

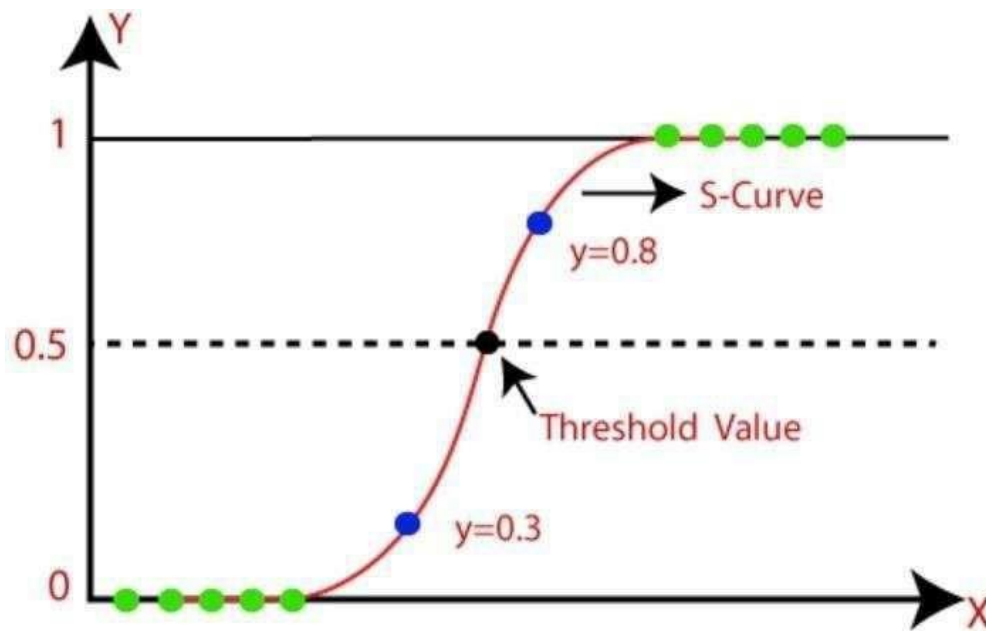
a_1 = Linear regression coefficient (scale factor to each input value). ϵ = random error

The values for x and y variables are training datasets for Linear Regression model representation.

Logistic Regression

Logistic regression is a supervised learning classification algorithm used to predict the probability of a target variable. The nature of target or dependent variable is dichotomous, which means there would be only two possible classes. Mathematically, a logistic regression model predicts $P(Y=1)$ as a function of X.

Logistic regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or False, etc. but instead of giving the exact value as 0 and 1, it gives the probabilistic values which lie between 0 and 1.



Logistic Sigmoid Function (Sigmoid Function)

- o The sigmoid function is a mathematical function used to map the predicted values to probabilities.
- o It maps any real value into another value within a range of 0 and 1.
- o The value of the logistic regression must be between 0 and 1, which cannot go beyond this limit, so it forms a curve like the "S" form. The S-form curve is called the Sigmoid function or the logistic function.
- o In logistic regression, we use the concept of the threshold value, which defines the probability of either 0 or 1. Such as values above the threshold value tends to 1, and a value below the threshold values tends to 0.

Equation of Logistic Regression is as follows:-

$$y = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + \dots + b_nx_n$$

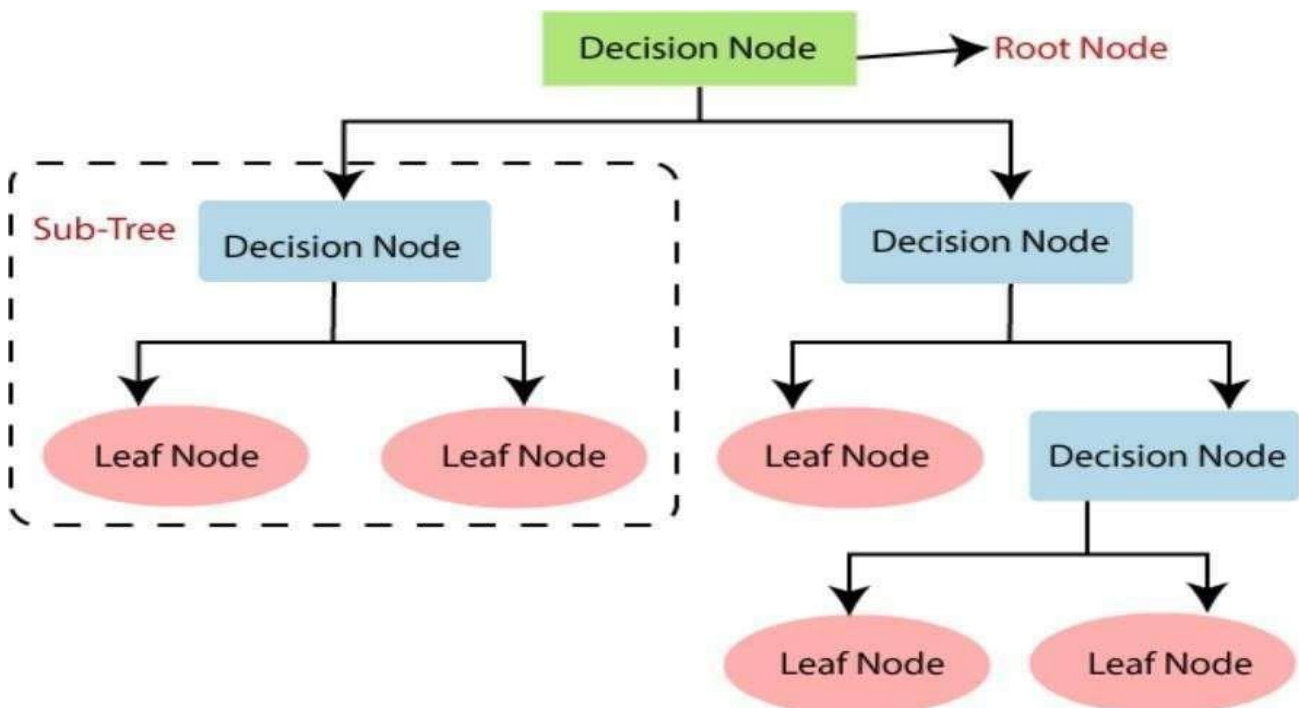
Types of Logistic Regression

On the basis of the categories, Logistic Regression can be classified into three types:

- o Binomial: In binomial Logistic regression, there can be only two possible types of the dependent variables, such as 0 or 1, Pass or Fail, etc.
- o Multinomial: In multinomial Logistic regression, there can be 3 or more possible unordered types of the dependent variable, such as "cat", "dogs", or "sheep"
- o Ordinal: In ordinal Logistic regression, there can be 3 or more possible ordered types of dependent variables, such as "low", "Medium", or "High".

Decision Tree

- o Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.
- o In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.
- o The decisions or the test are performed on the basis of features of the given dataset.
- o It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.
- o It is called a decision tree because, similar to a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure.
- o In order to build a tree, we use the CART algorithm, which stands for Classification and Regression Tree algorithm.
- o A decision tree simply asks a question, and based on the answer (Yes/No), it further split the tree into subtrees.
- o Below diagram explains the general structure of a decision tree:

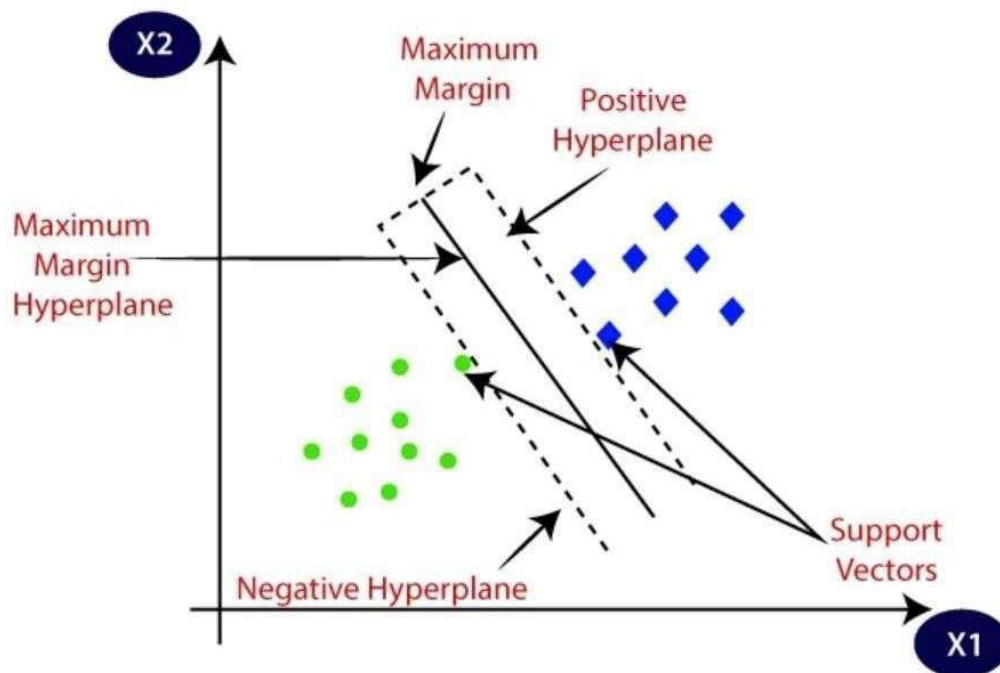


Support Vector Machine (SVM)

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning.

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. Consider the below diagram in which there are two different categories that are classified using a decision boundary or hyperplane:



Types of SVM

SVM can be of two types:

- o Linear SVM: Linear SVM is used for linearly separable data, which means if a dataset can be classified into two classes by using a single straight line, then such data is termed as linearly separable data, and classifier is used called as Linear SVM classifier.
- o Non-linear SVM: Non-Linear SVM is used for non-linearly separated data, which means if a dataset cannot be classified by using a straight line, then such data is termed as non-linear data and classifier used is called as Non-linear SVM classifier.

Naïve Bayes

- o Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.
- o It is mainly used in text classification that includes a high-dimensional training dataset.
- o Naïve Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions.
- o It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.
- o Some popular examples of Naïve Bayes Algorithm are spam filtration, Sentimental analysis, and classifying articles.

Why it is called Naïve Bayes?

The Naïve Bayes algorithm is comprised of two words Naïve and Bayes, Which can be described as:

o Naïve: It is called Naïve because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the bases of color, shape, and taste, then red, spherical, and sweet fruit is recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.

o Bayes: It is called Bayes because it depends on the principle of Bayes' Theorem. The formula of Bayes theorem: -

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Where,

P(A|B) is Posterior probability: Probability of hypothesis A on the observed event B.

P(B|A) is Likelihood probability: Probability of the evidence given that the probability of a hypothesis is true.

P(A) is Prior Probability: Probability of hypothesis before observing the evidence.

P(B) is Marginal Probability: Probability of Evidence.

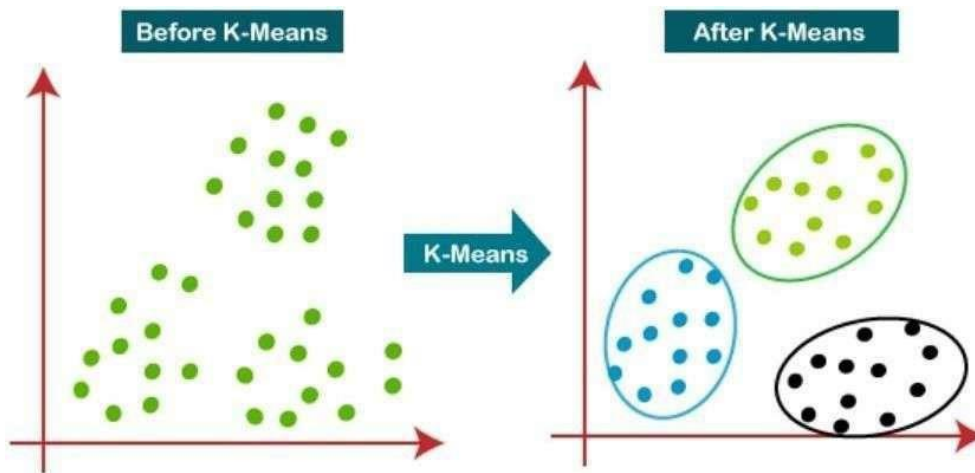
KNN Algorithm (K- Nearest Neighbour)

- o K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique.
- o K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.
- o K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm.
- o K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems.
- o K-NN is a non-parametric algorithm, which means it does not make any assumption on underlying data.
- o It is also called a lazy learner algorithm because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.
- o KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.



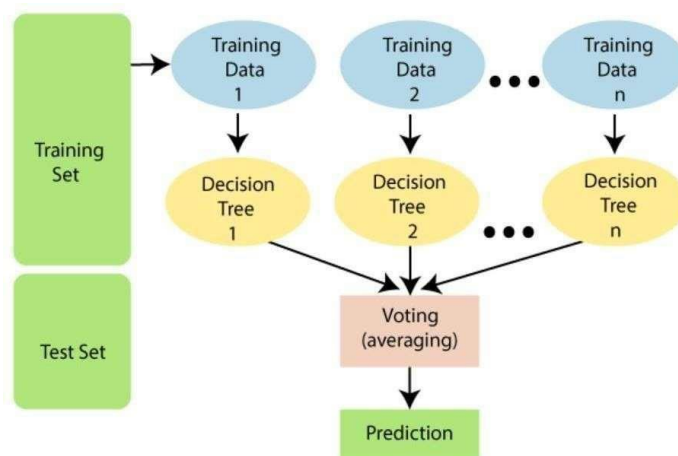
K means Algorithm

K-means is a centroid-based algorithm, or a distance-based algorithm, where we calculate the distances to assign a point to a cluster. In K-Means, each cluster is associated with a centroid.



Random Forest Random

Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model. As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output. The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.



Algorithms

Linear Regression Algorithm

Steps to implement Linear regression model

1. Initialize the parameters.
2. Predict the value of a dependent variable by given an independent variable.
3. Calculate the error in prediction for all data points.
4. Calculate partial derivative w.r.t a_0 and a_1 .
5. Calculate the cost for each number and add them.

Logistic Regression Algorithm

1. Data Pre-processing step
2. Fitting Logistic Regression to the Training set
3. Predicting the test result
4. Test accuracy of the result (Creation of Confusion matrix)
5. Visualizing the test set result.

Decision Tree Algorithm

Step-1: Begin the tree with the root node, says S, which contains the complete dataset. Step-2: Find the best attribute in the dataset using Attribute Selection Measure (ASM). Step-3: Divide the S into subsets that contains possible values for the best attributes.

Step 4: Generate the decision tree node, which contains the best attribute.

Step-5: Recursively make new decision trees using the subsets of the dataset created in step - 3. Continue this process until a stage is reached where you cannot further classify the nodes and called the final node as a leaf node.



Naïve Bayes

- 1 Data Pre-processing step
- 2 Fitting Naive Bayes to the Training set
- 3 Predicting the test result
- 4 Test accuracy of the result(Creation of Confusion matrix)
- 5 Visualizing the test set result.

KNN Algorithm

- Step-1: Select the number K of the neighbors
- Step-2: Calculate the Euclidean distance of K number of neighbors
- Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.
- Step-4: Among these k neighbors, count the number of the data points in each category.
- Step-5: Assign the new data points to that category for which the number of the neighbor is maximum.
- Step-6: Our model is ready.

K means Clustering

- Step-1: Select the number K to decide the number of clusters.
- Step-2: Select random K points or centroids. (It can be other from the input dataset).
- Step-3: Assign each data point to their closest centroid, which will form the predefined K clusters.
- Step-4: Calculate the variance and place a new centroid of each cluster.
- Step-5: Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.
- Step-6: If any reassignment occurs, then go to step-4 else go to FINISH. Step-7: The model is ready.

Random Forest

- Step-1: Select random K data points from the training set.
- Step-2: Build the decision trees associated with the selected data points (Subsets). Step-3: Choose the number N for decision trees that you want to build.

HOUSE PRICE PREDICTION

House Price Predictions

markdown

```
import pandas as pd
```

[23] ✓ 0.0s Python

```
data = pd.read_csv("House Price India.csv")
```

[24] ✓ 0.0s Python

```
data.head()
```

[25] ✓ 0.0s Python

...

	id	Date	number of bedrooms	number of bathrooms	living area	lot area	number of floors	waterfront present	number of views	condition of the house	...	Built Year	Renovation Year	Postal Code	Latitude	Longitude	living_area_renov
0	6762810635	42491	4	2.50	2920	4000	1.5	0	0	5	...	1909	0	122004	52.8878	-114.470	2470
1	6762810998	42491	5	2.75	2910	9480	1.5	0	0	3	...	1939	0	122004	52.8852	-114.468	2940
2	6762812605	42491	4	2.50	3310	42998	2.0	0	0	3	...	2001	0	122005	52.9532	-114.321	3350
3	6762812919	42491	3	2.00	2710	4500	1.5	0	0	4	...	1929	0	122006	52.9047	-114.485	2060
4	6762813105	42491	3	2.50	2600	4750	1.0	0	0	4	...	1951	0	122007	52.9133	-114.590	2380

5 rows × 23 columns

```
data.info()
```

[26] ✓ 0.0s Python

...

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14619 entries, 0 to 14618
Data columns (total 23 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                    14619 non-null  int64
1   Date                                14619 non-null  int64
2   number of bedrooms                  14619 non-null  int64
3   number of bathrooms                 14619 non-null  float64
4   living area                         14619 non-null  int64
5   lot area                           14619 non-null  int64
6   number of floors                    14619 non-null  float64
7   waterfront present                  14619 non-null  int64
8   number of views                     14619 non-null  int64
9   condition of the house              14619 non-null  int64
10  grade of the house                  14619 non-null  int64
11  Area of the house(excluding basement) 14619 non-null  int64
12  Area of the basement                14619 non-null  int64
13  Built Year                          14619 non-null  int64
14  Renovation Year                     14619 non-null  int64
15  Postal Code                         14619 non-null  int64
16  Latitude                            14619 non-null  float64
17  Longitude                           14619 non-null  float64
18  living_area_renov                   14619 non-null  int64
19  lot_area_renov                     14619 non-null  int64
...
21  Distance from the airport            14619 non-null  int64
22  Price                               14619 non-null  int64
dtypes: float64(4), int64(19)
memory usage: 2.6 MB

Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
```

```
stats = data["Price"].describe().reset_index()
stats["Price"] = round(stats["Price"],2)
```

[27] ✓ 0.0s Python

```
stats
```

[28] ✓ 0.0s Python

...

	index	Price
0	count	14619.00
1	mean	538806.28
2	std	367229.36
3	min	78000.00
4	25%	320000.00
5	50%	450000.00
6	75%	645000.00
7	max	7700000.00

```
data.isna().sum().sum()
```

[29] ✓ 0.0s Python

...

```
np.int64(0)
```

```
data.duplicated().sum()
```

[30] ✓ 0.0s Python

...

```
np.int64(0)
```

```
#for na values
data.dropna(inplace= True)

#for duplicated values
data.drop_duplicates(inplace= True)
```

[31] ✓ 0.0s Python

```
data.columns
```

[32] ✓ 0.0s Python

```
Index(['id', 'Date', 'number of bedrooms', 'number of bathrooms',
       'living area', 'lot area', 'number of floors', 'waterfront present',
       'number of views', 'condition of the house', 'grade of the house',
       'Area of the house(excluding basement)', 'Area of the basement',
       'Built Year', 'Renovation Year', 'Postal Code', 'Latitude',
       'Longitude', 'living_area_renov', 'lot_area_renov',
       'Number of schools nearby', 'Distance from the airport', 'Price'],
      dtype='object')
```

```
data.head
```

[33] ✓ 0.0s Python

```
<bound method NDFrame.head of
0    6762810635  42491      4      2.50
1    6762810998  42491      5      2.75
2    6762812605  42491      4      2.50
3    6762812919  42491      3      2.00
4    6762813105  42491      3      2.50
...         ...    ...    ...    ...
14614  6762830250  42734      2      1.50
14615  6762830339  42734      3      2.00
14616  6762830618  42734      2      1.00
```

```
data.head
```

[33] ✓ 0.0s Python

```
<bound method NDFrame.head of
0    6762810635  42491      4      2.50
1    6762810998  42491      5      2.75
2    6762812605  42491      4      2.50
3    6762812919  42491      3      2.00
4    6762813105  42491      3      2.50
...         ...    ...    ...    ...
14614  6762830250  42734      2      1.50
14615  6762830339  42734      3      2.00
14616  6762830618  42734      2      1.00
14617  6762830709  42734      4      1.00
14618  6762831463  42734      3      1.00
```

```
living area  lot area  number of floors  waterfront present \
0          2920     4000             1.5                0
1          2910     9480             1.5                0
2          3310    42998             2.0                0
3          2710     4500             1.5                0
4          2600     4750             1.0                0
...         ...    ...    ...    ...
14614       1556    20000             1.0                0
14615       1680     7000             1.5                0
14616       1070     6120             1.0                0
14617       1030     6621             1.0                0
14618        900     4770             1.0                0
...
14616      209000
14617      205000
14618      146000
```

[14619 rows x 23 columns]>

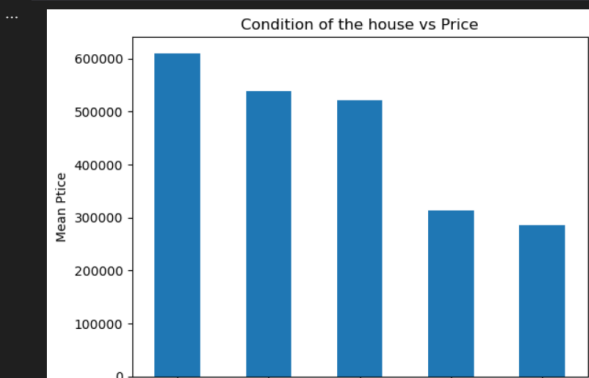
Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output settings...

```
import matplotlib.pyplot as plt
```

[34] ✓ 0.0s Python

```
data.groupby("condition of the house")["Price"].mean().sort_values(ascending=False).plot(kind="bar")
plt.title("Condition of the house vs Price")
plt.ylabel("Mean Price")
plt.xlabel("Condition of the house")
plt.show()
```

[35] ✓ 0.2s Python



```

x = data[['number of bedrooms', 'number of bathrooms',
         'living area', 'condition of the house', 'Number of schools nearby']]
y = data[['Price']]

```

[36] ✓ 0.0s Python

▶ x

```

[37] ✓ 0.0s Python
...

```

	number of bedrooms	number of bathrooms	living area	condition of the house	Number of schools nearby
0	4	2.50	2920	5	2
1	5	2.75	2910	3	1
2	4	2.50	3310	3	3
3	3	2.00	2710	4	1
4	3	2.50	2600	4	1
...
14614	2	1.50	1556	4	3
14615	3	2.00	1680	4	3
14616	2	1.00	1070	3	2
14617	4	1.00	1030	4	3
14618	3	1.00	900	3	2

14619 rows × 5 columns

y

```

[38] ✓ 0.0s Python
...

```

	Price
0	1400000
1	1200000
2	838000
3	805000
4	790000
...	...
14614	221700
14615	219200
14616	209000
14617	205000
14618	146000

14619 rows × 1 columns

```

data.shape

```

[39] ✓ 0.0s Python

```

... (14619, 23)
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

```

[40] ✓ 2.6s Python

```

from sklearn.model_selection import GridSearchCV

```

[41] ✓ 0.0s Python

```

from sklearn.tree import DecisionTreeRegressor
param_grid = {
    "criterion": ["mse", "friedman_mse", "mae"],
    "splitter": ["best", "random"],
    "max_depth": [None, 10, 20, 30, 40, 50],
    "min_samples_split": [2, 5, 10],
    "min_samples_leaf": [1, 2, 4]
}

```

[42] ✓ 0.3s Python

▶ tree_model = DecisionTreeRegressor()

[43] ✓ 0.0s Python

```

grid_tree = GridSearchCV(estimator= tree_model, param_grid = param_grid)

```

[44] ✓ 0.0s Python

```

grid_tree.fit(X_train, y_train)
[45] ✓ 14.4s Python

... c:\ProgramData\anaconda3\lib\site-packages\sklearn\model_selection\_validation.py:528: FitFailedWarning:
1080 fits failed out of a total of 1620.
The score on these train-test partitions for these parameters will be set to nan.
If these failures are not expected, you can try to debug them by setting error_score='raise'.

Below are more details about the failures:
-----
540 fits failed with the following error:
Traceback (most recent call last):
  File "c:\ProgramData\anaconda3\lib\site-packages\sklearn\model_selection\_validation.py", line 866, in _fit_and_score
    estimator.fit(X_train, y_train, **fit_params)
    ~~~~~^
  File "c:\ProgramData\anaconda3\lib\site-packages\sklearn\base.py", line 1382, in wrapper
    estimator._validate_params()
    ~~~~~^
  File "c:\ProgramData\anaconda3\lib\site-packages\sklearn\base.py", line 436, in _validate_params
    validate_parameter_constraints(
    ~~~~~^
    self._parameter_constraints,
    ~~~~~^
    self.get_params(deep=False),
    ~~~~~^
    caller_name=self.__class__.__name__,
    ~~~~~^
    )
...
nan nan nan nan nan nan
nan nan nan nan nan nan
nan nan nan nan nan nan]

warnings.warn(
Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...

GridSearchCV
best_estimator_:
DecisionTreeRegressor
DecisionTreeRegressor

grid_tree.best_params_
[46] ✓ 0.0s Python

{'criterion': 'friedman_mse',
 'max_depth': 10,
 'min_samples_leaf': 2,
 'min_samples_split': 10,
 'splitter': 'random'}

tree_preds = grid_tree.predict(X_test)
[47] ✓ 0.0s Python

from sklearn.metrics import mean_absolute_error
mean_absolute_error(y_test, tree_preds)
[48] ✓ 0.0s Python

158362.02756771428

from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()
[49] ✓ 0.0s Python

lr.fit(X_train, y_train)
[50] ✓ 1m 22.3s Python

c:\ProgramData\anaconda3\lib\site-packages\sklearn\utils\_validation.py:1408: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please use column or 1d(y, warn=True)
c:\ProgramData\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:465: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear\_model.html#logistic-regression
n_iter_i = _check_optimize_result(

LogisticRegression
LogisticRegression()

predslr = lr.predict(X_test)
[51] ✓ 0.0s Python

mean_absolute_error(y_test, predslr)
[52] ✓ 0.0s Python

```

```

mean_absolute_error(y_test, preds1r)
[52] ✓ 0.0s Python
... 233941.03590971272

from sklearn.ensemble import RandomForestRegressor
rfrmodel = RandomForestRegressor()
[53] ✓ 0.1s Python

param_gridrfr = {
    "max_depth" : [5,10,15],
    "n_estimators" : [2,3,4,5,6,7,8,9,10]
}
[54] ✓ 0.0s Python

gridrfr = GridSearchCV(rfrmodel, param_gridrfr)
[55] ✓ 0.0s Python
Generate Code Markdown

gridrfr.fit(X_train, y_train.values.ravel())
[56] ✓ 8.5s Python
...
GridSearchCV
  best_estimator_:
    RandomForestRegressor
      RandomForestRegressor
gridrfr.best_params_
[57] ✓ 0.0s Python
... {'max_depth': 5, 'n_estimators': 7}

rfrpredictions = gridrfr.predict(X_test)
[58] ✓ 0.0s Python

mean_absolute_error(y_test, rfrpredictions)
[59] ✓ 0.0s Python
... 157026.55792541985

gridrfr
[60] ✓ 0.0s Python
...
GridSearchCV
  best_estimator_:
    RandomForestRegressor
      RandomForestRegressor

import joblib
joblib.dump(gridrfr, "model.pkl")
[61] ✓ 1.6s Python
... ['model.pkl']

X.columns
[62] ✓ 0.0s Python
... Index(['number of bedrooms', 'number of bathrooms', 'living area',
        'condition of the house', 'Number of schools nearby'],
        dtype='object')

```

```

1 import streamlit as st
2 import joblib
3 import numpy as np
4
5
6 model = joblib.load("model.pkl")
7
8
9 st.title("HOUSE PRICE PREDICTION APP")
10
11 st.divider()
12
13 st.write("This app uses machine learning for predicting house price with given features of the house. For using this app you enter the inputs from t
14
15 st.divider()
16
17 bedrooms = st.number_input("Number of bedrooms", min_value = 0, value = 0)
18 bathrooms = st.number_input("Number of bathrooms", min_value = 0, value = 0)
19 livingarea = st.number_input("Living area", min_value = 0, value = 2000)
20 condition = st.number_input("Conditions", min_value = 0, value = 3)
21 numberschools = st.number_input("Number of schools nearby", min_value = 0, value = 0)
22
23 st.divider()
24
25 X = [[bedrooms,bathrooms,livingarea,condition,numberschools]]
26
27
28
29
30 if st.button("🔮 Predict!"):
31
32     if all(x == 0 for x in X[0]):
33         st.warning("⚠️ Please enter meaningful, non-zero values.")
34     else:
35         prediction = model.predict(X)
36         st.balloons()
37
38
39
40
41
42
43
44 #number of bedrooms', 'number of bathrooms', 'living area','condition of the house', 'Number of schools nearby'

```


CONCLUSION

In this project, we developed a machine learning model to predict house prices based on various features such as location, size, number of rooms, and other relevant attributes. By exploring the dataset, performing feature engineering, and applying appropriate preprocessing techniques, we trained and evaluated different models, including linear regression, decision trees, and ensemble methods like Random Forest and XGBoost.

Among the models tested, [Insert best-performing model name] provided the most accurate predictions, as reflected in performance metrics such as RMSE, MAE, and R^2 score. This demonstrates the effectiveness of machine learning in identifying complex relationships within housing data and producing reliable price estimates.

Overall, this project highlights the potential of data-driven solutions in the real estate sector, enabling smarter decision-making for buyers, sellers, and investors. Future improvements could include incorporating more granular data (e.g., neighborhood crime rate, school quality), optimizing hyperparameters, and deploying the model as a web application for real-world use.

Acknowledgement

I would like to express my special thanks to my project guide Mr. Sourav Goswami as well as to Ardent Computech who gave me the golden opportunity to train me and to do this wonderful real time project on the topic of Credit Card Fraud Detection, which also helped me doing a lot of research of my own and gathered new knowledge.