

MACHINE LEARNING

(Mobile Price Classification)

*Summer Internship Report Submitted in partial
fulfillment of the requirement for undergraduate degree*

of

Bachelor of Technology

In

COMPUTER SCIENCE AND ENGINEERING

By

**MALLUPEDDI SAI LOHITH
221710304034**

Under the Guidance of

Assistant Professor



Department Of Electronics and Communication
Engineering GITAM School of Technology

GITAM (Deemed to be
University) Hyderabad-502329

July 2020

DECLARATION

I submit this industrial training work entitled “**Mobile Price Classification**” to GITAM (Deemed To Be University), Hyderabad in partial fulfillment of the requirements for the award of the degree of “**Bachelor of Technology**” in “**Computer Science and Engineering**”. I declare that it was carried out independently by me under the guidance of **Mr** Asst. Professor, GITAM (Deemed To Be University), Hyderabad, India.

The results embodied in this report have not been submitted to any other University or Institute for the award of any degree or diploma.

Place: Hyderabad

Name: Mallupeddi Sai Lohith

Date:

Student RollNo:221710304034



GITAM (DEEMED TO BE UNIVERSITY)

Hyderabad-502329, India

Dated:

CERTIFICATE

This is to certify that the Industrial Training Report entitled “**Mobile Price Classification**” is being submitted by Mallupeddi Sai Lohith (221710304034) in partial fulfillment of the requirement for the award of **Bachelor of Technology in Computer Science and Engineering** at GITAM (Deemed to Be University), Hyderabad during the academic year 2018-19

It is faithful record work carried out by her at the **Computer Science and Engineering Department**, GITAM University Hyderabad Campus under my guidance and supervision.

Mr.

Assistant Professor
Department of CSE

Dr. S.Phani Kumar

Professor and HOD
Department of CSE

ACKNOWLEDGEMENT

Apart from my effort, the success of this internship largely depends on the encouragement and guidance of many others. I take this opportunity to express my gratitude to the people who have helped me in the successful completion of this internship.

I would like to thank respected **Dr. N. Siva Prasad**, Pro Vice Chancellor, GITAM Hyderabad and **Dr. N. Seetharamaiah**, Principal, GITAM Hyderabad

I would like to thank respected **Dr.S.Phani Kumar** , Head of the Department of Computer Science and Engineering for giving me such a wonderful opportunity to expand my knowledge for my own branch and giving me guidelines to present a internship report. It helped me a lot to realize of what we study for.

I would like to thank the respected faculties **Mr.** who helped me to make this internship a successful accomplishment.

I would also like to thank my friends who helped me to make my work more organized and well-stacked till the end.

MALLUPEDDI SAI LOHITH
221710304034

ABSTRACT

Machine learning algorithms are used to predict the values from the data set by splitting the data set in to train and test and building Machine learning algorithms models of higher accuracy to predict the values is the primary task to be performed on data set.

To classify “If the mobile with given features will be Economical or Expensive” is the main motive of this project. Real Dataset is collected from website. I have adapted the view point of looking at features of the dataset, for deep understanding of the problem. I have taken the stance of a seller and reasoned out the various factors of choice of the mobile’s customer buys differs from other companies. Different Scaling algorithms are used to identify and normalize the range of independent variables or features of dataset. It is used in data preprocessing and used to normalize minimum computational complexity. Different classifiers are used to achieve as higher accuracy as possible. Results are compared in terms of highest accuracy achieved and minimum features selected. Conclusion is made on the base of best selection algorithm and best classifier for the given dataset. This work can be used in any type of marketing and business to find optimal product (with minimum cost and maximum features). To classify the mobile price range.

Table of Contents

1.MACHINE LEARNING	1
1.1. Introduction:	1
1.2. Importance of Machine Learning:.....	1
1.3. Uses of Machine Learning:	2
1.4. Types of Machine Learning:	5
1.4.1. Supervised Learning:	6
1.4.2. Unsupervised Learning:.....	7
1.4.3. Reinforcement Learning:	7
2.DEEP LEARNING	8
2.1 Deep Learning and It's Importance:.....	8
2.2 Uses of Deep Learning:.....	9
2.3 Relation between Data Mining, Machine Learning and Deep Learning:	11
3.PYTHON	13
3.1 Introduction:	13
3.2 Setup of Python:	13
3.2.1 Installation(using python IDLE):.....	14
3.2.2 Python Installation using Anaconda:	15
3.3 Features:	16
3.4 Variable Types	17
3.4.1 Python Numbers:	17
3.4.2 Python Strings:	18
3.4.3 Python lists:	18
3.4.4 python tuples:	18

3.4.5 python Dictionary:	19
3.5 Functions:	19
3.5.1 Defining a Function:	19
3.5.2 Calling a Function:	19
3.6 OOPs Concepts:	20
3.6.1 Class:	20
4.Mobile Price Classification.....	22
4.1 Project Requirements:	22
4.1.1. Packages used:	22
4.1.2. Versions of the packages:	23
4.1.3. Algorithms used:.....	23
4.2. Problem Statement:	23
4.3. Dataset Description:	24
4.4. Objective of the Case Study:.....	25
5.DATA PREPROCESSING/FEATURE ENGINEERING AND EDA	26
5.1 Statistical Analysis:	26
5.2 Data Type Conversions:	27
5.3 Detection of Outliers:	28
5.4 Handling Missing Values:	29
5.5 Encoding Categorical Data:	30
5.6 Generating Plots:	30
5.6.1. Visualize the data between Target and the Features:	30
5.6.1. Visualize the data between all the Features:	32
6. FEATURE SELECTION	35
6.1 Select relevant features for the analysis:	35

6.2 Drop irrelevant features:.....	36
6.3 Train-Test-Split:	36
6.4 Feature Scaling:.....	37
7. MODEL BUILDING AND EVALUATION	40
7.1 Brief about the algorithms used:	40
7.2 Train the Models:	42
7.3 Make Predictions:.....	43
7.4 Validate the Models:	45
7.5 Parameter Tuning:	47
i. Grid search:	47
iii. Bayesian optimization:	48
7.6 Confusion Matrix in Algorithm:	49
7.7 Classification Report:	52
7.8 Roc Curve:.....	53
7.9 Predictions from raw data:	54
8. Conclusion	58
9. REFERENCES	59

LIST OF FIGURES:

FIG 1.2 USAGE OF MACHINE LEARNING IN DIFFERENT FIELDS	2
FIG 1.3 USES OF MACHINE LEARNING.....	5
FIG 1.4 TYPES OF ML	6
FIG 2.1 DEEP NEURAL NETWORK	9
FIG 2.2 THE DEEP LEARNING PROCESS.....	10
FIG 2.3 RELATION BETWEEN DM,ML,DL	11
FIG 2.3 PROCESS IN MACHINE LEARNING AND DEEP LEARNING.....	11
FIG 3.2.1 : PYTHON DOWNLOAD.....	14
FIG 3.2.1.1 PYTHON INSTALLATION.....	14
FIG 3.2.1.2 IDLE	15
FIG 3.2.2.1 AFTER INSTALLATION	16
FIG 3.2.2.2 JUPYTER NOTEBOOK.....	16
FIG 3.6.1 CLASS DEFINING	20
FIG 3.6.1.1 EXAMPLE OF CLASS	21
FIG 4.1.1 PACKAGES.....	23
FIG 4.1.2 VERSIONS	23
FIG 5.1 LOADING DATA SET	26
FIG 5.1.1 STATISTICAL DATA	27
FIG 5.2 DATATYPES	28
FIG 5.3 DETECTION OF OUTLIER USING BOXPLOT.....	29
FIG 5.5 NUMERICAL IN THE DATASET	30
FIG 5.6.1 CORRELATION	30
FIG 5.6.1.1 CO-RELATION GRAPH	31
FIG 5.6.1.2 RAM VS PRICE.....	32
FIG 5.6.1.3 BLUETOOTH,WIFI.....	32
FIG 5.6.2.1 3G SUPPORT PHONES	32
FIG 5.6.2.2 DUAL SIMS	33
FIG 5.6.2.3 BATTERY POWER	33
FIG 5.6.2.4 4G SUPPORT	34
FIG 6.1 TOP FEATURES IN DATASET	36

FIG 6.2 CORRELATION BETWEEN ATTRIBUTES	36
FIG 6.4.2 SCALED X TRAIN	38
FIG 6.4.3 SCALED X TEST	38
FIG 7.1.1 LOGISTIC REGRESSION FOR TRAIN DATA.....	40
FIG 7.1.2 IMPORTING DECISION PACKAGES	41
FIG 7.1.3 RANDOM FORESTS	42
FIG 7.2.1 TRAINING DATA	43
FIG 7.3.2 PREDICTING TEST IN LOGISTIC REGRESSION	44
FIG 7.3.3 PREDICTING TEST IN DECISION TREE	44
FIG 7.3.4 PREDICTING TEST IN RANDOM FOREST	44
FIG 7.4.1 TEST AND TRAIN ACCURACY IN LOGISTIC REGRESSION	45
FIG 7.4.2 TEST VS TRAIN IN LOGISTIC REGRESSION	45
FIG 7.4.3 TEST AND TRAIN ACCURACY IN DECISION TREE	46
FIG 7.4.4 TEST VS TRAIN IN DECISION TREE	46
FIG 7.4.5 TEST ACCURACY IN RANDOM FOREST	46
FIG 7.4.6 TEST VS TRAIN IN RANDOM FOREST	47
FIG 7.5.1 GRID SEARCH CV	49
FIG 7.5.2 ACCURACY OF TEST USING GRID SEARCH CV I.E ACCURACY IS INCREASED IN DECISION TREE	49
FIG 7.6.1 CONFUSION MATRIX OF TEST DATASET IN RANDOM FOREST.....	50
FIG 7.6.2 HEAT MAP OF CONFUSION MATRIX	50
FIG 7.6.3 ACCURACY	51
FIG 7.6.4 RECALL	51
FIG 7.6.5 PRECISION.....	51
FIG 7.7.1 TEST CLASSIFICATION REPORT	52
FIG 7.9.1 BEST ACCURACY MODEL AMONG ALGORITHMS	55
FIG 7.9.2 LOADING TEST DATASET	55
FIG 7.9.3 PREDICTING TEST DATASET USING RANDOM FOREST.....	56
FIG 7.9.4 ADDED NEW COLUMN AS PRICE RANGE IN TEST DATA SET	56
FIG 7.9.5 RANGE OF TARGET IN TRAIN DATASET.....	56
FIG 7.9.6 OUTPUT OF PREDICTED TEST DATASET	57

1.MACHINE LEARNING

1.1. Introduction:

Over the past two decades Machine Learning has become one of the mainstays of information technology and with that, a rather central, albeit usually hidden, part of our life. With the ever increasing amounts of data becoming available there is good reason to believe that smart data analysis will become even more pervasive as a necessary ingredient for technological progress.

Human designers often produce machines that do not work as well as desired in the environments in which they are used. In fact, certain characteristics of the working environment might not be completely known at design time. Machine learning methods can be used for on-the-job improvement of existing machine designs. The amount of knowledge available about certain tasks might be too large for explicit encoding by humans. Machines that learn this knowledge gradually might be able to capture more of it than humans would want to write down. Environments change over time. Machines that can adapt to a changing environment would reduce the need for constant redesign. New knowledge about tasks is constantly being discovered by humans. Vocabulary changes. There is a constant stream of new events in the world. Continuing redesign of AI systems to conform to new knowledge is impractical, but machine learning methods might be able to track much of it.

1.2. Importance of Machine Learning:

Machine learning is a branch of artificial intelligence that aims at enabling machines to perform their jobs skillfully by using intelligent software. The statistical learning methods constitute the backbone of intelligent software that is used to develop machine intelligence. Because machine learning algorithms require data to learn, the discipline must have connection with the discipline of database. Similarly, there are familiar terms such as Knowledge Discovery from Data (KDD), data mining, and pattern recognition. One wonders how to view the big picture in which such connection is illustrated.

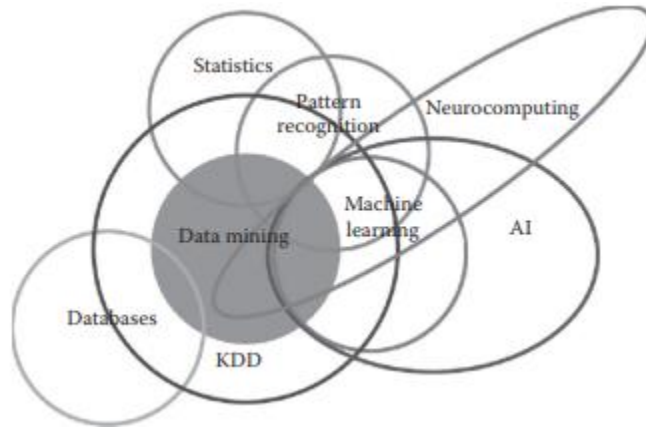


Fig 1.2 usage of Machine learning in different fields

There are some tasks that humans perform effortlessly or with some efforts, but we are unable to explain how we perform them. For example, we can recognize the speech of our friends without much difficulty. If we are asked how we recognize the voices, the answer is very difficult for us to explain. Because of the lack of understanding of such phenomenon (speech recognition in this case), we cannot craft algorithms for such scenarios. Machine learning algorithms are helpful in bridging this gap of understanding.

The idea is very simple. We are not targeting to understand the underlying processes that help us learn. We write computer programs that will make machines learn and enable them to perform tasks, such as prediction. The goal of learning is to construct a model that takes the input and produces the desired result. Sometimes, we can understand the model, whereas, at other times, it can also be like a black box for us, the working of which cannot be intuitively explained. The model can be considered as an approximation of the process we want machines to mimic. In such a situation, it is possible that we obtain errors for some input, but most of the time, the model provides correct answers. Hence, another measure of performance (besides performance of metrics of speed and memory usage) of a machine learning algorithm will be the accuracy of results.

1.3. Uses of Machine Learning:

Artificial Intelligence (AI) is everywhere. Possibility is that you are using it in one way or the other and you don't even know about it. One of the popular applications of AI is Machine Learning (ML), in which computers, software, and devices perform via cognition (very similar to

human brain). Herein, we share few examples of machine learning that we use everyday and perhaps have no idea that they are driven by ML. These are some the uses and applications of ML

i.Virtual Personal Assistants:

Siri, Alexa, Google Now are some of the popular examples of virtual personal assistants. As the name suggests, they assist in finding information, when asked over voice. All you need to do is activate them and ask “What is my schedule for today?”, “What are the flights from Germany to London”, or similar questions. For answering, your personal assistant looks out for the information, recalls your related queries, or send a command to other resources (like phone apps) to collect info. You can even instruct assistants for certain tasks like “Set an alarm for 6 AM next morning”, “Remind me to visit Visa Office day after tomorrow”.

Machine learning is an important part of these personal assistants as they collect and refine the information on the basis of your previous involvement with them. Later, this set of data utilized to render results that are tailored to your preferences.

Virtual Assistants are integrated to a variety of platforms. For example:

- Smart Speakers: Amazon Echo and Google Home
- Smartphones: Samsung Bixby on Samsung S8
- Mobile Apps: Google Allo

ii.Predictions while Commuting:

Traffic Predictions: We all have been using GPS navigation services. While we do that, our current locations and velocities are being saved at a central server for managing traffic. This data is then used to build a map of current traffic. While this helps in preventing the traffic and does congestion analysis, the underlying problem is that there are less number of cars that are equipped with GPS. Machine learning in such scenarios helps to estimate the regions where congestion can be found on the basis of daily experiences.

Online Transportation Networks: When booking a cab, the app estimates the price of the ride. When sharing these services, how do they minimize the detours? The answer is machine learning. Jeff Schneider, the engineering lead at Uber ATC reveals in a an interview that they use ML to define price surge hours by predicting the rider demand. In the entire cycle of the services, ML is playing a major role.

iii. Social Media Services:

From personalizing your news feed to better ads targeting, social media platforms are utilizing machine learning for their own and user benefits. Here are a few examples that you must be noticing, using, and loving in your social media accounts, without realizing that these wonderful features are nothing but the applications of ML.

- **People You May Know:** Machine learning works on a simple concept: understanding with experiences. Facebook continuously notices the friends that you connect with, the profiles that you visit very often, your interests, workplace, or a group that you share with someone etc. On the basis of continuous learning, a list of Facebook users are suggested that you can become friends with.

- **Face Recognition:** You upload a picture of you with a friend and Facebook instantly recognizes that friend. Facebook checks the poses and projections in the picture, notice the unique features, and then match them with the people in your friend list. The entire process at the backend is complicated and takes care of the precision factor but seems to be a simple application of ML at the front end.

- **Similar Pins:** Machine learning is the core element of Computer Vision, which is a technique to extract useful information from images and videos. Pinterest uses computer vision to identify the objects (or pins) in the images and recommend similar pins accordingly.

iv. Search Engine Result Refining:

Google and other search engines use machine learning to improve the search results for you. Every time you execute a search, the algorithms at the backend keep a watch at how you respond to the results. If you open the top results and stay on the web page for long, the search engine assumes that the the results it displayed were in accordance to the query. Similarly, if you reach the second or third page of the search results but do not open any of the results, the search engine estimates that the results served did not match requirement. This way, the algorithms working at the backend improve the search results.

v. Product Recommendations:

You shopped for a product online few days back and then you keep receiving emails for shopping suggestions. If not this, then you might have noticed that the shopping website or the app recommends you some items that somehow matches with your taste. On the basis of your behaviour with the website/app, past purchases, items liked or added to cart, brand preferences etc., the product recommendations are made.

vi. Online Fraud Detection:

Machine learning is proving its potential to make cyberspace a secure place and tracking monetary frauds online is one of its examples. For example: Paypal is using ML for protection against money laundering. The company uses a set of tools that helps them to compare millions of transactions taking place and distinguish between legitimate or illegitimate transactions taking place between the buyers and sellers.



Fig 1.3 Uses of Machine learning

1.4. Types of Machine Learning:

There are 3 types of Machine learning which are widely used in today's world these are:

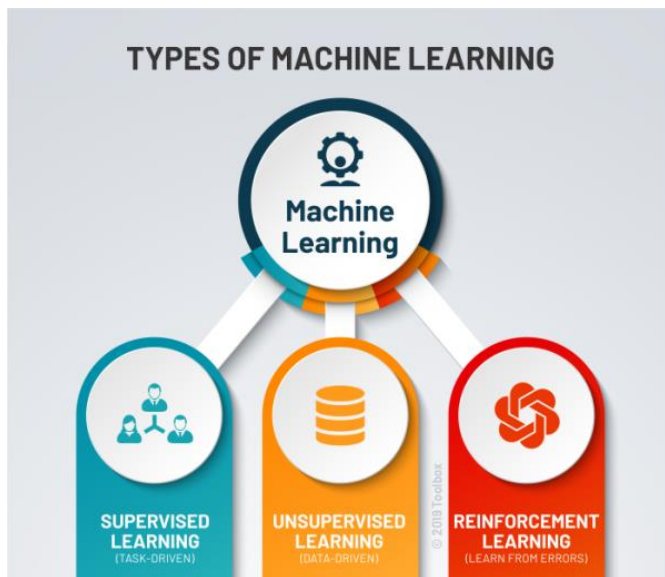


Fig 1.4 types of ML

1.4.1. Supervised Learning:

Supervised learning is one of the most basic types of machine learning. In this type, the machine learning algorithm is trained on labeled data. Even though the data needs to be labeled accurately for this method to work, supervised learning is extremely powerful when used in the right circumstances. In supervised learning, the ML algorithm is given a small training dataset to work with. This training dataset is a smaller part of the bigger dataset and serves to give the algorithm a basic idea of the problem, solution, and data points to be dealt with. The training dataset is also very similar to the final dataset in its characteristics and provides the algorithm with the labeled parameters required for the problem. The algorithm then finds relationships between the parameters given, essentially establishing a cause and effect relationship between the variables in the dataset. At the end of the training, the algorithm has an idea of how the data works and the relationship between the input and the output.

1.4.2. Unsupervised Learning:

Unsupervised machine learning holds the advantage of being able to work with unlabeled data. This means that human labor is not required to make the dataset machine-readable, allowing much larger datasets to be worked on by the program. In supervised learning, the labels allow the algorithm to find the exact nature of the relationship between any two data points. However, unsupervised learning does not have labels to work off of, resulting in the creation of hidden structures. Relationships between data points are perceived by the algorithm in an abstract manner, with no input required from human beings. The creation of these hidden structures is what makes unsupervised learning algorithms versatile. Instead of a defined and set problem statement, unsupervised learning algorithms can adapt to the data by dynamically changing hidden structures. This offers more post-deployment development than supervised learning algorithms.

1.4.3. Reinforcement Learning:

It directly takes inspiration from how human beings learn from data in their lives. It features an algorithm that improves upon itself and learns from new situations using a trial-and-error method. Favorable outputs are encouraged or 'reinforced', and non-favorable outputs are discouraged or 'punished'. Based on the psychological concept of conditioning, reinforcement learning works by putting the algorithm in a work environment with an interpreter and a reward system. In every iteration of the algorithm, the output result is given to the interpreter, which decides whether the outcome is favorable or not. In case of the program finding the correct solution, the interpreter reinforces the solution by providing a reward to the algorithm. If the outcome is not favorable, the algorithm is forced to reiterate until it finds a better result. In most cases, the reward system is directly tied to the effectiveness of the result. In typical reinforcement learning use-cases, such as finding the shortest route between two points on a map, the solution is not an absolute value. Instead, it takes on a score of effectiveness, expressed in a percentage value. The higher this percentage value is, the more reward is given to the algorithm. Thus, the program is trained to give the best possible solution for the best possible reward.

2.DEEP LEARNING

2.1 Deep Learning and It's Importance:

Deep learning algorithms run data through several “layers” of neural network algorithms, each of which passes a simplified representation of the data to the next layer. Most machine learning algorithms work well on datasets that have up to a few hundred features, or columns.

Basically deep learning is itself a subset of machine learning but in this case the machine learns in a way in which humans are supposed to learn. The structure of deep learning model is highly similar to a human brain with large number of neurons and nodes like neurons in human brain thus resulting in artificial neural network. In applying traditional machine learning algorithms we have to manually select input features from complex data set and then train them which becomes a very tedious job for ML scientist but in neural networks we don't have to manually select useful input features, there are various layers of neural networks for handling complexity of the data set and algorithm as well. In my recent project on human activity recognition , when we applied traditional machine learning algorithm like K-NN then we have to separately detect human and its activity also had to select impactful input parameters manually which became a very tedious task as data set was way too complex but the complexity dramatically reduced on applying artificial neural network, such is the power of deep learning. Yes it's correct that deep learning algorithms take lots of time for training sometimes even weeks as well but its execution on new data is so fast that its not even comparable with traditional ML algorithms. Deep learning has enabled Industrial Experts to overcome challenges which were impossible, a decades ago like Speech and Image recognition and Natural Language Processing. Majority of the Industries are currently depending on it , be it Journalism, Entertainment, Online Retail Store, Automobile, Banking and Finance, Healthcare, Manufacturing or even Digital Sector. Video recommendations, Mail Services, Self Driving cars, Intelligent Chat bots, Voice Assistants are just trending achievements of Deep Learning.

Furthermore, Deep learning can most profoundly be considered as future of Artificial Intelligence due to constant rapid increase in amount of data as well as the gradual development in hardware field as well, resulting in better computational power.

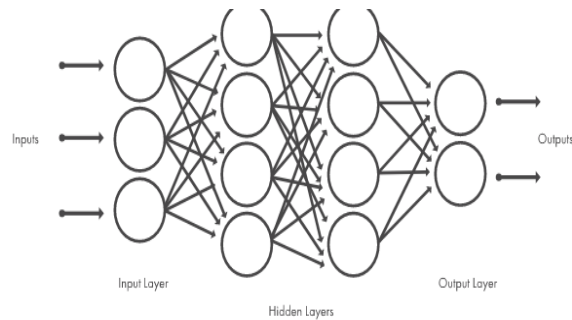


Fig 2.1 Deep Neural network

2.2 Uses of Deep Learning:

i. Translations:

Although automatic machine translation isn't new, deep learning is helping enhance automatic translation of text by using stacked networks of neural networks and allowing translations from images.

ii. Adding color to black-and-white images and videos:

It is used to be a very time-consuming process where humans had to add color to black-and-white images and videos by hand can now be automatically done with deep-learning models.

iii. Language recognition:

Deep learning machines are beginning to differentiate dialects of a language. A machine decides that someone is speaking English and then engages an AI that is learning to tell the differences between dialects. Once the dialect is determined, another AI will step in that specializes in that particular dialect. All of this happens without involvement from a human.

iv. Autonomous vehicles:

There's not just one AI model at work as an autonomous vehicle drives down the street. Some deep-learning models specialize in streets signs while others are trained to recognize pedestrians. As a car navigates down the road, it can be informed by up to millions of individual AI models that allow the car to act.

v. Computer vision:

Deep learning has delivered super-human accuracy for image classification, object detection, image restoration and image segmentation—even handwritten digits can be recognized.

Deep learning using enormous neural networks is teaching machines to automate the tasks performed by human visual systems.

vi. Text generation:

The machines learn the punctuation, grammar and style of a piece of text and can use the model it developed to automatically create entirely new text with the proper spelling, grammar and style of the example text. Everything from Shakespeare to Wikipedia entries have been created.

vii. Deep-learning robots:

Deep-learning applications for robots are plentiful and powerful from an impressive deep-learning system that can teach a robot just by observing the actions of a human completing a task to a housekeeping robot that's provided with input from several other AIs in order to take action. Just like how a human brain processes input from past experiences, current input from senses and any additional data that is provided, deep-learning models will help robots execute tasks based on the input of many different AI opinions.

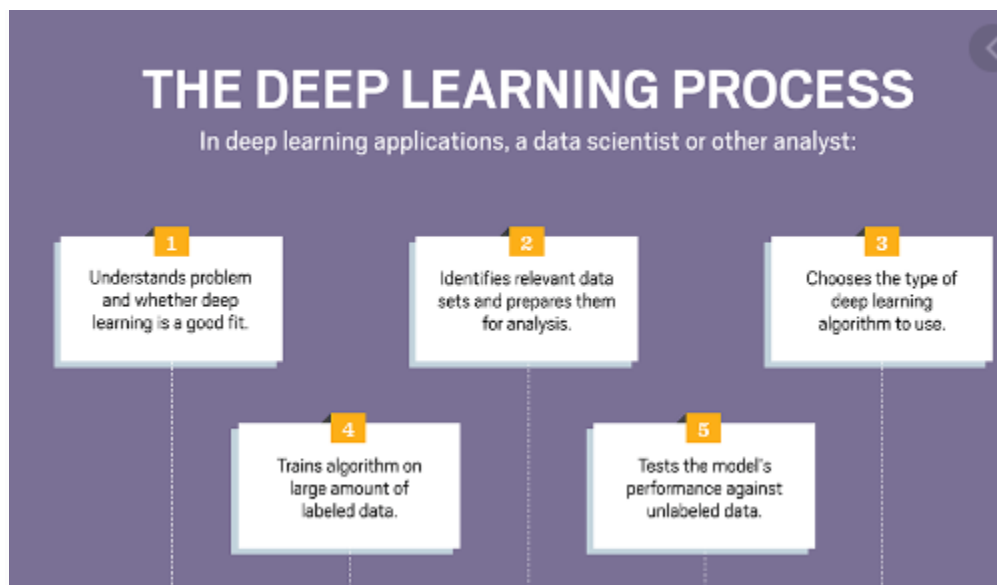


Fig 2.2 The deep learning process

2.3 Relation between Data Mining, Machine Learning and Deep Learning:

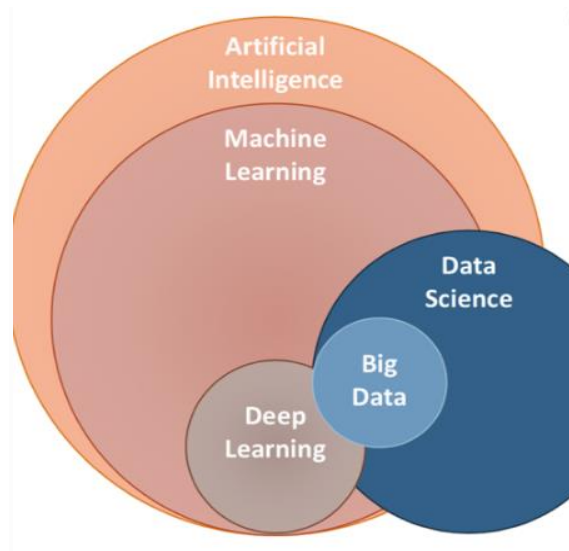


Fig 2.3 Relation between DM,ML,DL

The deep learning, data mining and machine learning share a foundation in data science, and there certainly is overlap between the two. Data mining can use machine learning algorithms to improve the accuracy and depth of analysis, and vice-versa; machine learning can use mined data as its foundation, refining the dataset to achieve better results.

You could also argue that data mining and machine learning are similar in that they both seek to address the question of how we can learn from data. However, the way in which they achieve this end, and their applications, form the basis of some significant differences.

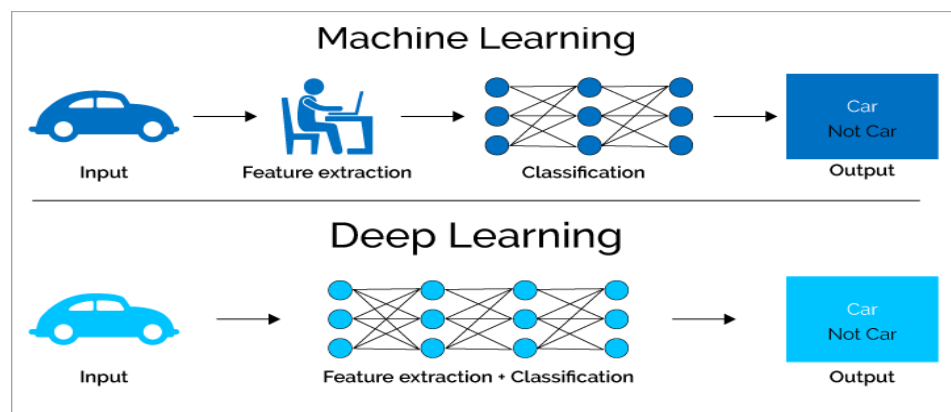


Fig 2.3 process in machine learning and deep learning

Mobile Price Classification

Machine Learning comprises of the ability of the machine to learn from trained data set and predict the outcome automatically. It is a subset of artificial intelligence.

Deep Learning is a subset of machine learning. It works in the same way on the machine just like how the human brain processes information. Like a brain can identify the patterns by comparing it with previously memorized patterns, deep learning also uses this concept.

Deep learning can automatically find out the attributes from raw data while machine learning selects these features manually which further needs processing. It also employs artificial neural networks with many hidden layers, big data, and high computer resources.

Data Mining is a process of discovering hidden patterns and rules from the existing data. It uses relatively simple rules such as association, correlation rules for the decision-making process, etc. Deep Learning is used for complex problem processing such as voice recognition etc. It uses Artificial Neural Networks with many hidden layers for processing. At times data mining also uses deep learning algorithms for processing the data.

3.PYTHON

3.1 Introduction:

Python is a widely used general-purpose, high level programming language. It was created by Guido van Rossum in 1991 and further developed by the Python Software Foundation. It was designed with an emphasis on code readability, and its syntax allows programmers to express their concepts in fewer lines of code. Python is a programming language that lets you work quickly and integrate systems more efficiently.

Python is dynamically typed and garbage-collected. It supports multiple programming paradigms, including structured, object-oriented, and functional programming. Python is often described as a "batteries included" language due to its comprehensive standard library.

Python is an interpreted, object-oriented, high-level programming language with dynamic semantics. Its high-level built-in data structures, combined with dynamic typing and dynamic binding, make it very attractive for Rapid Application Development, as well as for use as a scripting or glue language to connect existing components together. Python's simple, easy-to-learn syntax emphasizes readability and therefore reduces the cost of program maintenance. Python supports modules and packages, which encourages program modularity and code reuse. The Python interpreter and the extensive standard library are available in source or binary form without charge for all major platforms, and can be freely distributed.

3.2 Setup of Python:

- Python distribution is available for a wide variety of platforms. You need to download only the binary code applicable for your platform and install Python.
- The most up-to-date and current source code, binaries, documentation, news, etc., is available on the official website of Python <https://www.python.org/>

3.2.1 Installation(using python IDLE):

- To start, go to python.org/downloads and then click on the button to download the latest version of Python
- We can download python IDLE in windows,mac and linux operating systems also.

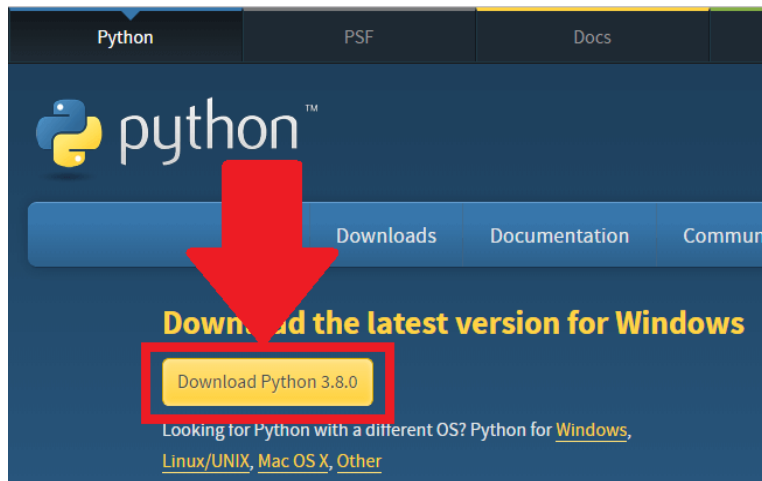


Figure 3.2.1 : Python download

- Run the .exe file that you just downloaded and start the installation of Python by clicking on Install Now
- We can give environmental variable i.e path after completion of downloading

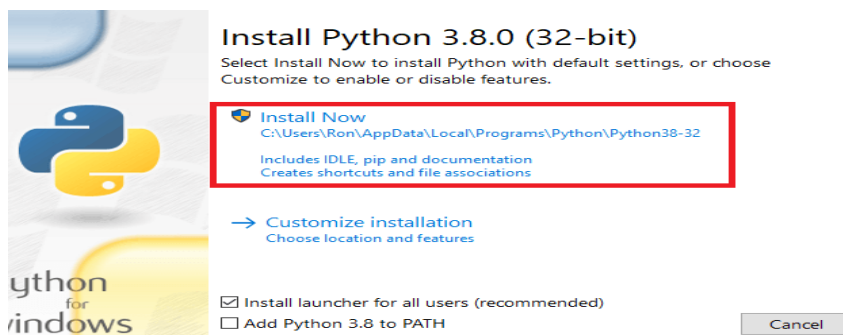


Fig 3.2.1.1 python installation

- When python is installed, a program called IDLE is also installed along with it. It provides a graphical user interface to work with python.

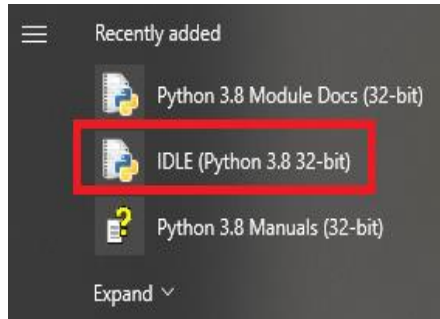


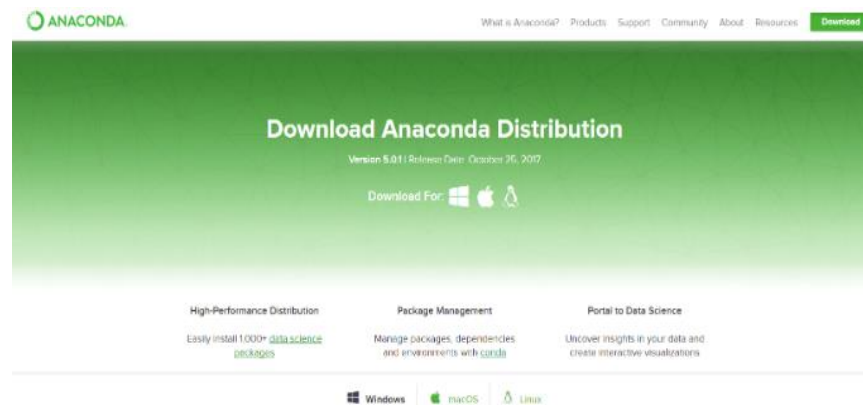
Fig 3.2.1.2 IDLE

3.2.2 Python Installation using Anaconda:

- Anaconda is a free open source distribution of python for large scale data processing, predictive analytics and scientific computing.
- Conda is a package manager quickly installs and manages packages.

Anaconda for Windows installation:

- i. Go to the following link: [Anaconda.com/downloads](https://anaconda.com/downloads)



- ii. Download python 3.4 version for (32-bitgraphic installer/64 -bit graphic installer)
- iii. Select path(i.e. add anaconda to path & register anaconda as default python 3.4)
- iv. Click finish
- v. Open jupyter notebook

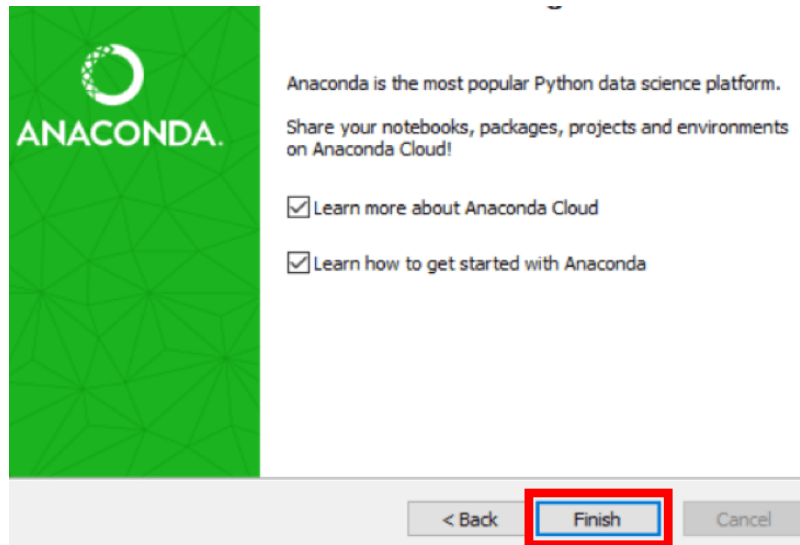


Fig 3.2.2.1 After installation

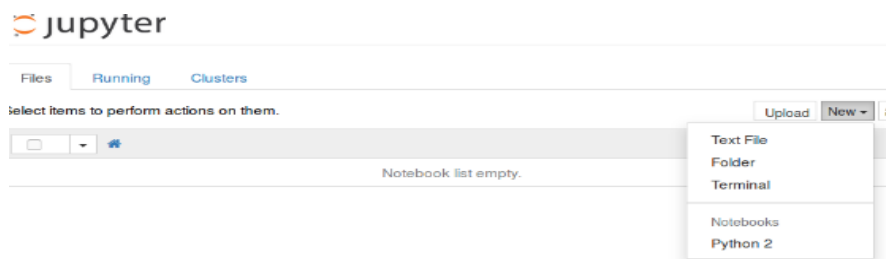


Fig 3.2.2.2 jupyter notebook

3.3 Features:

- i. **Readable:** Python is a very readable language.
- ii. **Easy to Learn:** Learning python is easy as this is a expressive and high level programming language, which means it is easy to understand the language and thus easy to learn
- iii. **Cross platform:** Python is available and can run on various operating systems such as Mac, Windows, Linux, Unix etc. This makes it a cross platform and portable language.
- iv. **Open Source:** Python is a open source programming language.
- v. **Large standard library:** Python comes with a large standard library that has some handy codes and functions which we can use while writing code in Python.
- vi. **Free:** Python is free to download and use. This means you can download it for free and use it in your application. Python is an example of a FLOSS (Free/Libre Open Source

Software), which means you can freely distribute copies of this software, read its source code and modify it.

vii. **Supports exception handling:** If you are new, you may wonder what is an exception? An exception is an event that can occur during program execution and can disrupt the normal flow of program. Python supports exception handling which means we can write less error prone code and can test various scenarios that can cause an exception later on.

viii. **Advanced features:** Supports generators and list comprehensions. We will cover these features later.

ix. **Automatic memory management:** Python supports automatic memory management which means the memory is cleared and freed automatically. You do not have to bother clearing the memory.

3.4 Variable Types:

Variables are nothing but reserved memory locations to store values. This means that when you create a variable you reserve some space in memory. Based on the data type of a variable, the interpreter allocates memory and decides what can be stored in the reserved memory. Therefore, by assigning different data types to variables, you can store integers, decimals or characters in these variables.

Python has five standard data types –

- Numbers
- Strings
- Lists
- Tuples
- Dictionary

3.4.1 Python Numbers:

Number data types store numeric values. They are immutable data types, means that changing the value of a number data type results in a newly allocated object.

Python supports four different numerical types –

- **int (signed integers)** – They are often called just integers or ints, are positive or negative whole numbers with no decimal point.

- **long (long integers)** – Also called longs, they are integers of unlimited size, written like integers and followed by an uppercase or lowercase L.
- **float (floating point real values)** – Also called floats, they represent real numbers and are written with a decimal point dividing the integer and fractional parts. Floats may also be in scientific notation, with E or e indicating the power of 10 ($2.5e2 = 2.5 \times 10^2 = 250$).

3.4.2 Python Strings:

In Python, Strings can be created by simply enclosing characters in quotes. Python does not support character types. These are treated as length-one strings, and are also considered as substrings. Substrings are immutable and can't be changed once created. Strings are the ordered blocks of text that are enclosed in single or double quotations. Thus, whatever is written in quotes, is considered as string. Though it can be written in single or double quotations, double quotation marks allow the user to extend strings over multiple lines without backslashes, which is usually the signal of continuation of an expression, e.g., 'abc', "ABC".

3.4.3 Python lists:

- List is a collection data type in python. It is ordered and allows duplicate entries as well. Lists in python need not be homogeneous, which means it can contain different data types like integers, strings and other collection data types. It is mutable in nature and allows indexing to access the members in a list.
- To declare a list, we use the square brackets.
- List is like any other array that we declare in other programming languages. Lists in python are often used to implement stacks and queues. The lists are mutable in nature. Therefore, the values can be changed even after a list is declared.

3.4.4 python tuples:

- A tuple is a collection of objects which ordered and immutable. Tuples are sequences, just like lists. The differences between tuples and lists are, the tuples cannot be changed unlike lists and tuples use parentheses, whereas lists use square brackets. Creating a tuple is as simple as putting different comma-separated values. Optionally you can put these comma-separated values between parentheses also

3.4.5 python Dictionary:

- It is a collection data type just like a list or a set, but there are certain features that make python dictionary unique. A dictionary in python is not ordered and is changeable as well. We can make changes in a dictionary unlike sets or strings which are immutable in nature. Dictionary contains key-value pairs like a map that we have in other programming languages. A dictionary has indexes. Since the value of the keys we declare in a dictionary are always unique, we can use them as indexes to access the elements in a dictionary.

3.5 Functions:

3.5.1 Defining a Function:

- Function blocks begin with the keyword `def` followed by the function name and parentheses `()`.
- Any input parameters or arguments should be placed within these parentheses. You can also define parameters inside these parentheses.
- The first statement of a function can be an optional statement - the documentation string of the function or docstring.
- The code block within every function starts with a colon `(:)` and is indented.
- The statement `return [expression]` exits a function, optionally passing back an expression to the caller. A return statement with no arguments is the same as `return None`.

3.5.2 Calling a Function:

- Defining a function only gives it a name, specifies the parameters that are to be included in the function and structures the blocks of code.
- Once the basic structure of a function is finalized, you can execute it by calling it from another function or directly from the Python prompt

3.6 OOPs Concepts:

3.6.1 Class:

- Python is an object oriented programming language. Unlike procedure oriented programming, where the main emphasis is on functions, object oriented programming stresses on objects..
- An object is simply a collection of data (variables) and methods (functions) that act on those data. Similarly, a class is a blueprint for that object.
- We can think of class as a sketch (prototype) of a house. It contains all the details about the floors, doors, windows etc. Based on these descriptions we build the house. House is the object.
- As many houses can be made from a house's blueprint, we can create many objects from a class. An object is also called an instance of a class and the process of creating this object is called **instantiation**.
- Like function definitions begin with the def keyword in Python, class definitions begin with a class keyword.
- The first string inside the class is called docstring and has a brief description about the class

```
class MyNewClass:  
    '''This is a docstring. I have created a new class'''  
    pass
```

Fig 3.6.1 Class defining

- As soon as we define a class, a new class object is created with the same name. This class object allows us to access the different attributes as well as to instantiate new objects of that class.

```
class Person:
    "This is a person class"
    age = 10

    def greet(self):
        print('Hello')

# Output: 10
print(Person.age)

# Output: <function Person.greet>
print(Person.greet)

# Output: 'This is my second class'
print(Person.__doc__)
```

Fig 3.6.1.1 Example of class

4.Mobile Price Classification

4.1Project Requirements:

4.1.1. Packages used:

- **Numpy:** In Python we have lists that serve the purpose of arrays, but they are slow to process.NumPy aims to provide an array object that is up to 50x faster than traditional Python lists.The array object in NumPy is called `ndarray`, it provides a lot of supporting functions that make working with `ndarray` very easy.Arrays are very frequently used in data science, where speed and resources are very important.
 - **Pandas:** Pandas is an open-source Python Library providing high-performance data manipulation and analysis tool using its powerful data structures. The name Pandas is derived from the word Panel Data – an Econometrics from Multidimensional data. Python with Pandas is used in a wide range of fields including academic and commercial domains including finance, economics, Statistics, analytics, etc.
 - **Seaborn:** Seaborn is a Python data visualization library based on `matplotlib`. It provides a high-level interface for drawing attractive and informative statistical graphics.
- Mathplotlib :** Matplotlib is one of the most popular Python packages used for data visualization. It is a cross-platform library for making 2D plots from data in arrays. Matplotlib is written in Python and makes use of NumPy, the numerical mathematics extension of Python. It provides an object-oriented API that helps in embedding plots in applications using Python GUI toolkits such as PyQt, WxPython or Tkinter. It can be used in Python and IPython shells, Jupyter notebook and web application servers also.Matplotlib has a procedural interface named the Pylab, which is designed to resemble MATLAB, a proprietary programming language developed by MathWorks. Matplotlib along with NumPy can be considered as the open source equivalent of MATLAB.

1. LOAD required packages

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
```

Fig 4.1.1 packages

4.1.2. Versions of the packages:

The versions of the packages are found by following command

```
#versions of packages
import numpy
import matplotlib
print('numpy:',numpy.__version__)
print('pandas:',pd.__version__)
print('seaborn:',sns.__version__)
print('matplotlib:',matplotlib.__version__)

numpy: 1.18.5
pandas: 1.0.5
seaborn: 0.10.1
matplotlib: 3.2.2
```

Fig 4.1.2 versions

4.1.3. Algorithms used:

Here , 3 algorithms are used they are:

- Logistic regression
- Decision Tree
- Random Forest

4.2. Problem Statement:

Bob has started his own mobile company. He wants to give tough fight to big companies like Apple,Samsung etc.

Mobile Price Classification

He does not know how to estimate price of mobiles his company creates. In this competitive mobile phone market you cannot simply assume things. To solve this problem he collects sales data of mobile phones of various companies.

Bob wants to find out some relation between features of a mobile phone(eg:- RAM,Internal Memory etc) and its selling price. But he is not so good at Machine Learning. So he needs your help to solve this problem.

In this problem you do not have to predict actual price but a price range indicating how high the price is.

4.3. Dataset Description:

In this data the columns contains:

- **id:** ID
- **battery_power:** Total energy a battery can store in one time measured in mAh
- **blue:** Has bluetooth or not
- **clock_speed:** speed at which microprocessor executes instructions
- **dual_sim:** Has dual sim support or not
- **fc:** Front Camera mega pixels
- **four_g:** Has 4G or not
- **int_memory:** Internal Memory in Gigabytes
- **m_dep:** Mobile Depth in cm
- **mobile_wt:** Weight of mobile phone
- **n_cores:** Number of cores of processor
- **pc:** Primary Camera mega pixels
- **px_height:** Pixel Resolution Height
- **px_width:** Pixel Resolution Width
- **ram:** Random Access Memory in Megabytes
- **sc_h:** Screen Height of mobile in cm
- **sc_w:** Screen Width of mobile in cm
- **talk_time:** longest time that a single battery charge will - last when you are

- **three_g:** Has 3G or not
- **touch_screen:** Has touch screen or not
- **wifi:** Has wifi or not

4.4. Objective of the Case Study:

Objective of problem is you do not have to predict actual price but a price range indicating how high the price is.

5.DATA PREPROCESSING/FEATURE ENGINEERING AND EDA

5.1 Statistical Analysis:

Pandas in python provide an interesting method `read_csv()`. The `read_csv` function reads the entire dataset from a comma separated values file and we can assign it to a DataFrame to which all the operations can be performed. It helps us to access each and every row as well as columns and each and every value can be access using the dataframe. Any missing value or NaN value have to be cleaned.

```
df_train=pd.read_csv('https://raw.githubusercontent.com/lohith4034/AI/ML/master/aiml%20project/train.csv')
```

```
df_train.head()
```

	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_cores	pc	px_height	px_width	ram	sc_h	sc_w	talk_ti
0	842	0	2.2	0	1	0	7	0.6	188	2	2	20	756	2549	9	7	
1	1021	1	0.5	1	0	1	53	0.7	136	3	6	905	1988	2631	17	3	
2	563	1	0.5	1	2	1	41	0.9	145	5	6	1263	1716	2603	11	2	
3	615	1	2.5	0	0	0	10	0.8	131	6	9	1216	1786	2769	16	8	
4	1821	1	1.2	0	13	1	44	0.6	141	2	14	1208	1212	1411	8	2	

Fig 5.1 loading data set

Descriptive statistics include those that summarize the central tendency, dispersion and shape of a dataset's distribution, excluding nan values. Analyzes both numeric and object series, as well as DataFrame column sets of mixed data types. The output will vary depending on what is provided.

For numeric data, the result's index will include count, mean, std, min, max as well as lower, 50 and upper percentiles. By default the lower percentile is 25 and the upper percentile is 75. The 50 percentile is the same as the median.

For object data (e.g. strings or timestamps), the result's index will include count, unique, top, and freq. The top is the most common value. The freq is the most common value's frequency. Timestamps also include the first and last items.

If multiple object values have the highest count, then the count and top results will be arbitrarily chosen from among those with the highest count.

Mobile Price Classification

For mixed data types provided via a `DataFrame`, the default is to return only an analysis of numeric columns. If the dataframe consists only of object and categorical data without any numeric columns, the default is to return an analysis of both the object and categorical columns. If `include='all'` is provided as an option, the result will include a union of attributes of each type.

	count	mean	std	min	25%	50%	75%	max
battery_power	2000.0	1238.51850	439.418206	501.0	851.75	1226.0	1615.25	1998.0
blue	2000.0	0.49500	0.500100	0.0	0.00	0.0	1.00	1.0
clock_speed	2000.0	1.52225	0.816004	0.5	0.70	1.5	2.20	3.0
dual_sim	2000.0	0.50950	0.500035	0.0	0.00	1.0	1.00	1.0
fc	2000.0	4.30950	4.341444	0.0	1.00	3.0	7.00	19.0
four_g	2000.0	0.52150	0.499662	0.0	0.00	1.0	1.00	1.0
int_memory	2000.0	32.04650	18.145715	2.0	16.00	32.0	48.00	64.0
m_dep	2000.0	0.50175	0.288416	0.1	0.20	0.5	0.80	1.0
mobile_wt	2000.0	140.24900	35.399655	80.0	109.00	141.0	170.00	200.0
n_cores	2000.0	4.52050	2.287837	1.0	3.00	4.0	7.00	8.0
pc	2000.0	9.91650	6.064315	0.0	5.00	10.0	15.00	20.0
px_height	2000.0	645.10800	443.780811	0.0	282.75	564.0	947.25	1960.0
px_width	2000.0	1251.51550	432.199447	500.0	874.75	1247.0	1633.00	1998.0
ram	2000.0	2124.21300	1084.732044	256.0	1207.50	2146.5	3064.50	3998.0
sc_h	2000.0	12.30650	4.213245	5.0	9.00	12.0	16.00	19.0

Fig 5.1.1 Statistical data

Observations:-

- 1.Max Battery power is 2000mah
- 2.75% of the phones have Dual sim availability
- 3.50% of the mobile phones has 32gb of memory
- 4.most of the mobile phones are screen touch enabled and supports 3g , 4g and are wifi enabled

5.2 Data Type Conversions:

When doing data analysis, it is important to make sure you are using the correct data types; otherwise you may get unexpected results or errors. In the case of pandas, it will correctly infer data types in many cases and you can move on with your analysis without any further thought on the topic. Despite how well pandas works, at some point in your data analysis processes, you will likely need to explicitly convert data from one type to another. This article will discuss the basic

Mobile Price Classification

pandas data types (aka dtypes), how they map to python and numpy data types and the options for converting from one pandas type to another.

```
df_train.dtypes
battery_power    int64
blue             int64
clock_speed      float64
dual_sim         int64
fc              int64
four_g          int64
int_memory       int64
m_dep           float64
mobile_wt        int64
n_cores          int64
pc              int64
px_height        int64
px_width         int64
ram             int64
sc_h            int64
sc_w            int64
talk_time       int64
three_g         int64
touch_screen     int64
wifi            int64
price_range      int64
dtype: object
```

Fig 5.2 datatypes

5.3 Detection of Outliers:

Perhaps the most common or familiar type of outlier is the observations that are far from the rest of the observations or the center of mass of observations. This is easy to understand when we have one or two variables and we can visualize the data as a histogram or scatter plot, although it becomes very challenging when we have many input variables defining a high-dimensional input feature space. In this case, simple statistical methods for identifying outliers can break down, such as methods that use standard deviations or the interquartile range. It can be important to identify and remove outliers from data when training machine learning algorithms for predictive modeling. Outliers can skew statistical measures and data distributions, providing a misleading representation of the underlying data and relationships. Removing outliers from training data prior to modeling can result in a better fit of the data and, in turn, more skillful predictions. Thankfully, there are a variety of automatic model-based methods for identifying outliers in input data. Importantly, each method approaches the definition of an outlier in slightly different ways, providing alternate approaches to preparing a training dataset that can be evaluated and compared, just like any other data preparation step in a modeling pipeline.

There are no outliers in a data set.

Mobile Price Classification

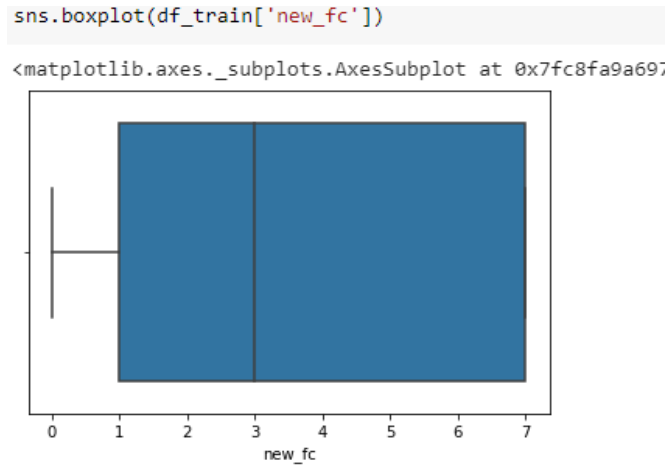


Fig 5.3 detection of outlier using boxplot

5.4 Handling Missing Values:

There are a number of schemes that have been developed to indicate the presence of missing data in a table or DataFrame. Generally, they revolve around one of two strategies: using a mask that globally indicates missing values, or choosing a sentinel value that indicates a missing entry. In the masking approach, the mask might be an entirely separate Boolean array, or it may involve appropriation of one bit in the data representation to locally indicate the null status of a value. In the sentinel approach, the sentinel value could be some data-specific convention, such as indicating a missing integer value with -9999 or some rare bit pattern, or it could be a more global convention, such as indicating a missing floating-point value with NaN (Not a Number), a special value which is part of the IEEE floating-point specification.

Fig 5.4 There are no missing values in a dataset

```
df_train.isnull().sum()
battery_power    0
blue             0
clock_speed      0
dual_sim         0
fc               0
four_g           0
int_memory       0
m_dep            0
mobile_wt        0
n_cores          0
pc               0
px_height        0
px_width         0
ram              0
sc_h             0
sc_w             0
talk_time        0
three_g          0
touch_screen     0
wifi             0
price_range      0
dtype: int64
```


5.5 Encoding Categorical Data:

Categorical Variables are of two types: Nominal and Ordinal

- Nominal: The categories do not have any numeric ordering in between them. They don't have any ordered relationship between each of them. Examples: Male or Female, any colour
- Ordinal: The categories have a numerical ordering in between them. Example: Graduate is less than Post Graduate, Post Graduate is less than Ph.D. customer satisfaction survey, high low medium
- Categorical data can be handled by using dummy variables, which are also called as indicator variables.

In the given dataset I have not used any encoding because dataset is numerical

```
df_train.select_dtypes(include=['int', 'float']).columns #numerical
Index(['battery_power', 'blue', 'clock_speed', 'dual_sim', 'fc', 'four_g',
      'int_memory', 'm_dep', 'mobile_wt', 'n_cores', 'pc', 'px_height',
      'px_width', 'ram', 'sc_h', 'sc_w', 'talk_time', 'three_g',
      'touch_screen', 'wifi', 'price_range'],
      dtype='object')
```

Fig 5.5 numerical in the dataset

5.6 Generating Plots:

5.6.1. Visualize the data between Target and the Features:

i. Correlation:

df_train.corr() #correlation												
	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_cores	pc	px_height
battery_power	1.000000	0.011252	0.011482	-0.041847	0.033334	0.015665	-0.004004	0.034085	0.001844	-0.029727	0.031441	0.014907
blue	0.011252	1.000000	0.021419	0.035198	0.003593	0.013443	0.041177	0.004049	-0.008605	0.036161	-0.009952	-0.006872
clock_speed	0.011482	0.021419	1.000000	-0.001315	-0.000434	-0.043073	0.006545	-0.014364	0.012350	-0.005724	-0.005245	-0.014522
dual_sim	-0.041847	0.035198	-0.001315	1.000000	-0.029123	0.003187	-0.015679	-0.022142	-0.008979	-0.024658	-0.017143	-0.020875
fc	0.033334	0.003593	-0.000434	-0.029123	1.000000	-0.016560	-0.029133	-0.001791	0.023618	-0.013356	0.644595	-0.009990
four_g	0.015665	0.013443	-0.043073	0.003187	-0.016560	1.000000	0.008690	-0.001823	-0.016537	-0.029706	-0.005598	-0.019236
int_memory	-0.004004	0.041177	0.006545	-0.015679	-0.029133	0.008690	1.000000	0.006886	-0.034214	-0.028310	-0.033273	0.010447
m_dep	0.034085	0.004049	-0.014364	-0.022142	-0.001791	-0.001823	0.006886	1.000000	0.021756	-0.003504	0.026282	0.025263
mobile_wt	0.001844	-0.008605	0.012350	-0.008979	0.023618	-0.016537	-0.034214	0.021756	1.000000	-0.018989	0.018844	0.000935
n_cores	-0.029727	0.036161	-0.005724	-0.024658	-0.013356	-0.029706	-0.028310	-0.003504	-0.018989	1.000000	-0.001193	-0.006872
pc	0.031441	-0.009952	-0.005245	-0.017143	0.644595	-0.005598	-0.033273	0.026282	0.018844	-0.001193	1.000000	-0.018465

Fig 5.6.1 correlation

Mobile Price Classification

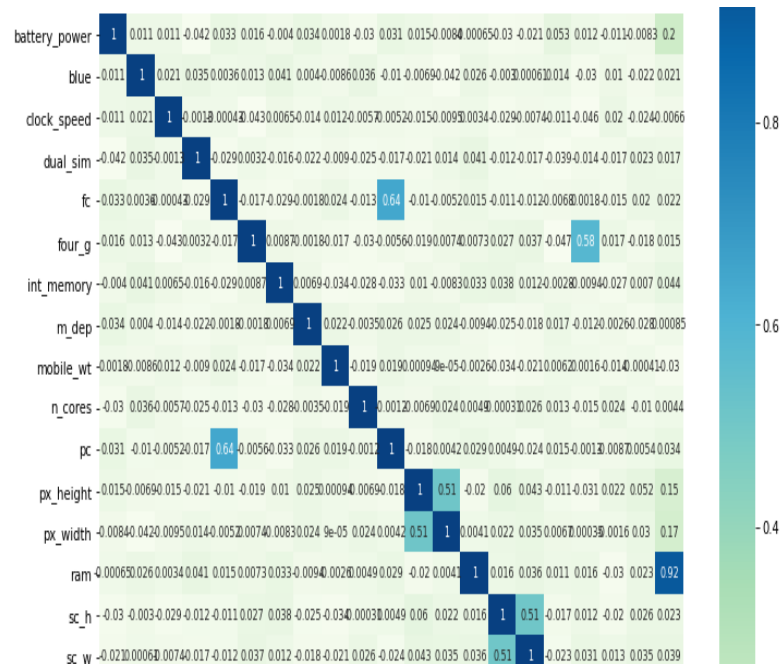


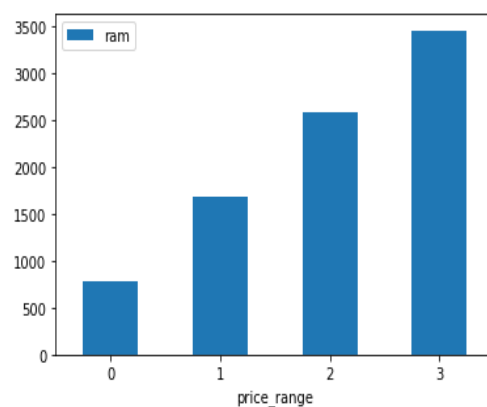
Fig 5.6.1.1 co-relation graph

We can see these attributes having relationship with each other:

- Price range vs. ram: high positive correlation
- fc vs. pc: positive correlation
- four_g vs. three_g: positive correlation
- pc_height vs. pc_width: positive correlation

ii. How does ram is affected by price :

```
df_train.groupby('price_range').mean()['ram'].plot(kind = 'bar', legend = True)
```



Mobile Price Classification

Fig 5.6.1.2 ram vs price

Here, we can see :

- The price_range 3 having more RAM
- The price_range 0 having low RAM

iii. Bluetooth, Wifi vs price:

```
sns.catplot('price_range', col='blue', hue='wifi', data=df_train, kind='count', col_wrap=2)  
<seaborn.axisgrid.FacetGrid at 0x7f9def49c208>
```

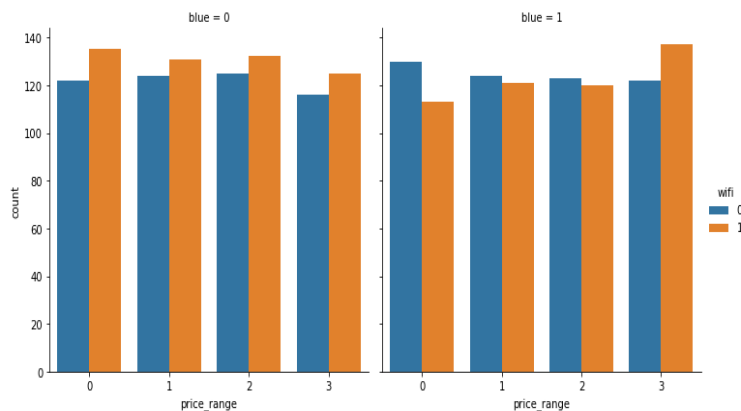


Fig 5.6.1.3 bluetooth, wifi

In the catplot we observe that:

- Bluetooth and Wifi seem to not have a significant affect to phone price since they have similar distribution in every price range.

5.6.1. Visualize the data between all the Features:

i. % of Phones which support 3G:

```
labels = ["3G-supported", 'Not supported']  
values = df_train['three_g'].value_counts().values  
fig1, ax1 = plt.subplots()  
ax1.pie(values, labels=labels, autopct='%1.1f%%', shadow=True, startangle=90)  
plt.show()
```

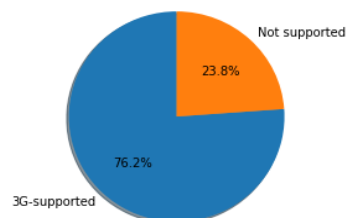


Fig 5.6.2.1 3g support phones

Mobile Price Classification

Here, in the piechart we observe:

- 76.2% phones support 3-G Network
- 23.8% phones does not support 3-G Network

ii.Supports Dualsim or not:

```
x=df_train['dual_sim'].value_counts()
labels='Supports Dualsim: '+str(x[1]),'Does not support Dualsim:- '+str(x[0])
sizes=[x[1],x[0]]
fig1, ax1 = plt.subplots()
ax1.pie(sizes,labels=labels)
ax1.axis('equal')
plt.show()
```

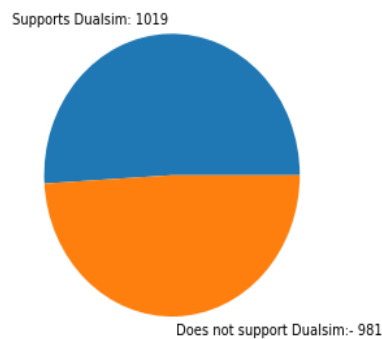


Fig 5.6.2.2 dual sims

Observations:-

- 1019 mobile phones are Dualsim Enabled
- 981 phones does not support Dualsim feature

iii.Battery power vs Price Range:

```
sns.boxplot(x="price_range", y="battery_power", data=df_train)
<matplotlib.axes._subplots.AxesSubplot at 0x7f9def1655c0>
```

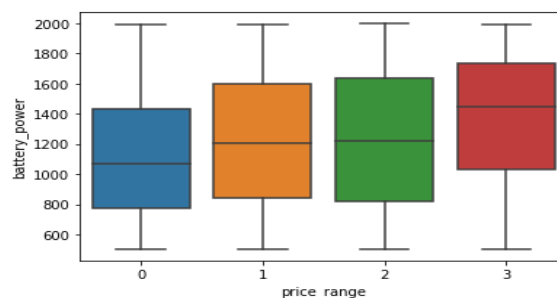


Fig 5.6.2.3 battery power

In box plot we observe that:

Mobile Price Classification

- when price_range is 3 then battery power greater than 1600 MAH
- when price_range is 2 then battery power greater is 1600 MAH
- when price_range is 1 then battery power is 1500 MAH
- when price_range is 0 then battery power is 1400 MAH

iv. % of Phones which support 4G:

```
labels4g = ["4G-supported", 'Not supported']  
values4g = df_train['four_g'].value_counts().values  
fig1, ax1 = plt.subplots()  
ax1.pie(values4g, labels=labels4g, autopct='%1.1f%%', shadow=True, startangle=90)  
plt.show()
```

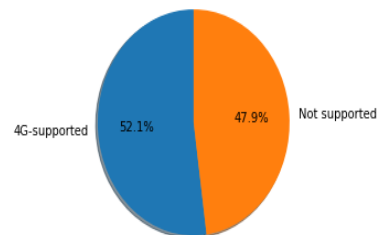


Fig 5.6.2.4 4G support

Here, in the piechart we observe:

- 52.1% phones support 4-G Network
- 47.9% phones does not support 4-G Network

v.Observations:

- the most influential variable is ram
- most of the variables have very little correlation to price range
- having 3G and 4G is somewhat correlated
- there is no highly correlated inputs in our dataset, so there is no multicollinearity problem.
- ram has direct impact on the price range of the phones
- Features like
 1. Dual sim
 2. wifi
 3. 4g
 4. 3g
 5. Touch screen
- these have more impact on the phones prices

6. FEATURE SELECTION

6.1 Select relevant features for the analysis:

Feature Selection is the process where you automatically or manually select those features which contribute most to your prediction variable or output in which you are interested in. Having irrelevant features in your data can decrease the accuracy of the models and make your model learn based on irrelevant features.

- **Reduces Overfitting:** Less redundant data means less opportunity to make decisions based on noise.
- **Improves Accuracy:** Less misleading data means modeling accuracy improves.
- **Reduces Training Time:** fewer data points reduce algorithm complexity and algorithms train faster.

Feature Selection Methods:

I will share 3 Feature selection techniques that are easy to use and also gives good results.

1. Univariate Selection
2. Feature Importance
3. Correlation Matrix with Heatmap

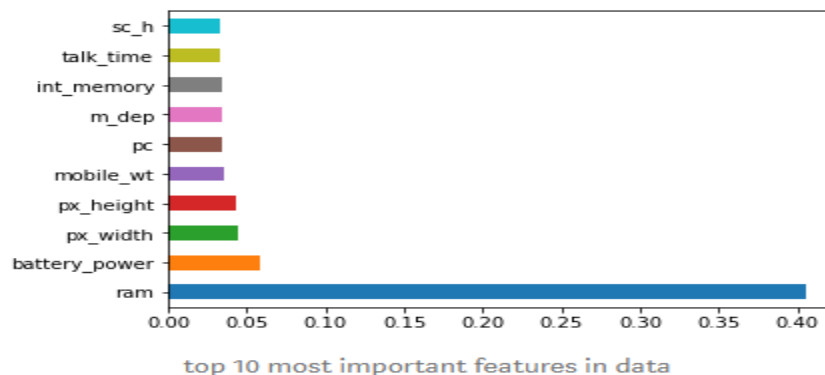


Fig 6.1 top features in dataset

6.2 Drop irrelevant features:

There are no irrelevant features in my dataset

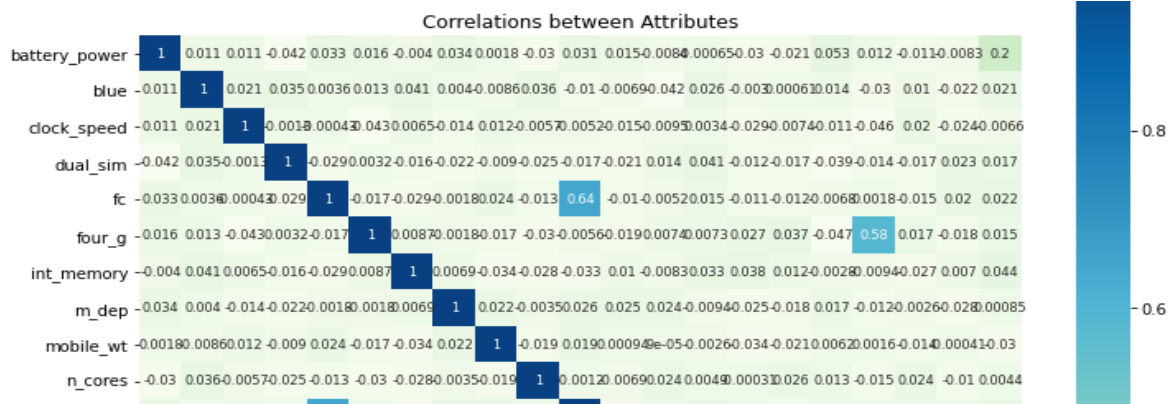


Fig 6.2 correlation between attributes

6.3 Train-Test-Split:

One of the first decisions to make when starting a modeling project is how to utilize the existing data. One common technique is to split the data into two groups typically referred to as the training and testing sets. The training set is used to develop models and feature sets; they are the substrate for estimating parameters, comparing models, and all of the other activities required to reach a final model. The test set is used only at the conclusion of these activities for estimating a final, unbiased assessment of the model's performance. It is critical that the test set not be used prior to this point. Looking at the test sets results would bias the outcomes since the testing data will have become part of the model development process.

There are a number of ways to split the data into training and testing sets. The most common approach is to use some version of random sampling. Completely random sampling is a straightforward strategy to implement and usually protects the process from being biased towards any characteristic of the data. However this approach can be problematic when the response is not evenly distributed across the outcome. A less risky splitting strategy would be to use a stratified random sample based on the outcome. For

Mobile Price Classification

classification models, this is accomplished by selecting samples at random within each class. This approach ensures that the frequency distribution of the outcome is approximately equal within the training and test sets. When the outcome is numeric, artificial strata can be constructed based on the quartiles of the data. For example, in the Ames housing price data, the quartiles of the outcome distribution would break the data into four artificial groups containing roughly 230 houses. The training/test split would then be conducted within these four groups and the four different training set portions are pooled together (and the same for the test set).

```
y = df_train['price_range']
X = df_train.drop('price_range', axis = 1)

y.unique() # We have four price ranges as target values
array([1, 2, 3, 0])
```

Fig 6.3.1 dividing input and output and target values for output

Importing packages:

```
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 101)

print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)

(1600, 20)
(400, 20)
(1600,)
(400,)
```

Fig 6.3.2 importing packages and splitting with random state=101 and test_size=0.2

6.4 Feature Scaling:

It is a step of Data Pre Processing which is applied to independent variables or features of data. It basically helps to normalise the data within a particular range. Sometimes, it also helps in speeding up the calculations in an algorithm. Real world dataset contains features that highly vary

Mobile Price Classification

in magnitudes, units, and range. Normalisation should be performed when the scale of a feature is irrelevant or misleading and not should Normalise when the scale is meaningful.

The algorithms which use Euclidean Distance measure are sensitive to Magnitudes. Here feature scaling helps to weigh all the features equally.

Formally, If a feature in the dataset is big in scale compared to others then in algorithms where Euclidean distance is measured this big scaled feature becomes dominating and needs to be normalized.

$$z = \frac{x - \mu}{\sigma}$$

Fig 6.4.1 formula for scaling

Importing package:

```
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
scaler = StandardScaler().fit(X_train)
```

```
scaled_X_train = pd.DataFrame(scaler.fit_transform(X_train))
scaled_X_train
```

	0	1	2	3	4	5	6	7	8	
0	0.984823	-0.986343	0.345546	0.990050	-0.066906	-1.053953	-1.156036	0.653456	-0.451426	0.6469
1	-0.673224	1.013846	-1.250227	-1.010051	-0.766578	0.948808	1.044235	-1.416900	0.114337	-1.102
2	-0.400668	1.013846	0.222795	0.990050	-0.999802	-1.053953	1.429282	-0.036663	-0.140256	1.084
3	1.545834	-0.986343	-1.250227	-1.010051	-0.300130	0.948808	1.704316	-0.726781	-0.819173	-1.102
4	-1.359156	-0.986343	-1.250227	-1.010051	-0.766578	-1.053953	0.494167	1.688634	1.613611	1.084

Fig 6.4.2 scaled X train

```
scaled_X_test = pd.DataFrame(scaler.fit_transform(X_test))
scaled_X_test
```

	0	1	2	3	4	
0	0.342699	-1.005013	-1.264622	-1.056599	1.232482	0.999999
1	-1.402675	-1.005013	-1.264622	-1.056599	-0.308120	-1.000000
2	-1.455565	-1.005013	-0.167336	0.946433	1.232482	0.999999

Fig 6.4.3 scaled X test

Mobile Price Classification

1. **K-Means** uses the Euclidean distance measure here feature scaling matters.
2. **K-Nearest-Neighbours** also require feature scaling.
3. **Principal Component Analysis (PCA)**: Tries to get the feature with maximum variance, here too feature scaling is required.
4. **Gradient Descent**: Calculation speed increase as Theta calculation becomes faster after feature scaling

7. MODEL BUILDING AND EVALUATION

7.1 Brief about the algorithms used:

i. 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.

In simple words, the dependent variable is binary in nature having data coded as either 1 (stands for success/yes) or 0 (stands for failure/no).

Mathematically, a logistic regression model predicts $P(Y=1)$ as a function of X . It is one of the simplest ML algorithms that can be used for various classification problems such as spam detection, Diabetes prediction, cancer detection etc.

Before diving into the implementation of logistic regression, we must be aware of the following assumptions about the same –

- In case of binary logistic regression, the target variables must be binary always and the desired outcome is represented by the factor level 1.
- There should not be any multi-collinearity in the model, which means the independent variables must be independent of each other.
- We must include meaningful variables in our model.
- We should choose a large sample size for logistic regression.

```
from sklearn.linear_model import LogisticRegression

lr1 = LogisticRegression(multi_class = 'ovr', solver = 'sag', max_iter = 10000)

lr1.fit(scaled_X_train, y_train)

LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, l1_ratio=None, max_iter=10000,
                    multi_class='ovr', n_jobs=None, penalty='l2',
                    random_state=None, solver='sag', tol=0.0001, verbose=0,
                    warm_start=False)
```

Fig 7.1.1 logistic regression for train data

ii. Decision Tree:

A decision tree is a tree-like graph with nodes representing the place where we pick an attribute and ask a question; edges represent the answers to the question; and the leaves represent the actual output or class label. They are used in non-linear decision making with simple linear decision surface. Decision trees classify the examples by sorting them down the tree from the root to some leaf node, with the leaf node providing the classification to the example. Each node in the tree acts as a test case for some attribute, and each edge descending from that node corresponds to one of the possible answers to the test case. This process is recursive in nature and is repeated for every subtree rooted at the new nodes.

Decision Tree Analysis is a general, predictive modelling tool that has applications spanning a number of different areas. In general, decision trees are constructed via an algorithmic approach that identifies ways to split a data set based on different conditions. It is one of the most widely used and practical methods for supervised learning. Decision Trees are a non-parametric supervised learning method used for both classification and regression tasks. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features.

The decision rules are generally in form of if-then-else statements. The deeper the tree, the more complex the rules and fitter the model.

```
from sklearn.tree import DecisionTreeClassifier

dt = DecisionTreeClassifier(random_state=101, criterion='entropy')
dt_model = dt.fit(X_train, y_train)

y_train_pred_dt = dt.predict(X_train)

metrics.confusion_matrix(y_train, y_train_pred_dt)

array([[406,  0,  0,  0],
       [ 0, 398,  0,  0],
       [ 0,  0, 380,  0],
       [ 0,  0,  0, 416]])
```

Fig 7.1.2 importing decision packages

iii. Random forest:

Random forest is a flexible, easy to use machine learning algorithm that produces, even without hyper-parameter tuning, a great result most of the time. It is also one of the most used

Mobile Price Classification

algorithms, because of its simplicity and diversity (it can be used for both classification and regression tasks). In this post we'll learn how the random forest algorithm works, how it differs from other algorithms and how to use it.

Random forest is a supervised learning algorithm. The "forest" it builds, is an ensemble of decision trees, usually trained with the “bagging” method. The general idea of the bagging method is that a combination of learning models increases the overall result.

Random forest is a supervised learning algorithm. The "forest" it builds, is an ensemble of decision trees, usually trained with the “bagging” method. The general idea of the bagging method is that a combination of learning models increases the overall result.

Random forest builds multiple decision trees and merges them together to get a more accurate and stable prediction.

One big advantage of random forest is that it can be used for both classification and regression problems, which form the majority of current machine learning systems.

```
#Import the RFC from sklearn
from sklearn.ensemble import RandomForestClassifier

# initialize the object for RFC
rf = RandomForestClassifier(n_estimators = 100, random_state=101, criterion = 'entropy', oob_score = True)
model_rf = rf.fit(X_train, y_train)
```

```
y_pred_train_1113 = rf.predict(X_train)

from sklearn.metrics import confusion_matrix, classification_report
print(classification_report(y_train, y_pred_train_1113))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	406
1	1.00	1.00	1.00	398
2	1.00	1.00	1.00	380
3	1.00	1.00	1.00	416

Fig 7.1.3 Random Forests

7.2 Train the Models:

Splitting the data : after the preprocessing is done then the data is split into train and test sets

- In Machine Learning in order to access the performance of the classifier. You train the classifier using 'training set' and then test the performance of your classifier on unseen 'test set'.

Mobile Price Classification

An important point to note is that during training the classifier only uses the training set . The test set must not be used during training the classifier. The test set will only be available during testing the classifier.

- training set - a subset to train a model.(Model learns patterns between Input and Output)
- test set - a subset to test the trained model.(To test whether the model has correctly learnt)
- The amount or percentage of Splitting can be taken as specified
- First we need to identify the input and output variables and we need to separate the input set and output set
- In scikit learn library we have a package called model_selection in which train_test_split method is available .we need to import this method
- This method splits the input and output data to train and test based on the percentage specified by the user and assigns them to four different variables.

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 101)
```

Fig 7.2.1 Training data

7.3 Make Predictions:

- Then we have to test the model for the test set ,that is done as follows
- We have a method called predict , using this method we need to predict the output for input test set and we need to compare the out but with the output test data
- If the predicted values and the original values are close then we can say that model is trained with good accuracy

```
y_test_pred_lr = lr1.predict(scaled_X_test)

confusion_matrix = metrics.confusion_matrix(y_test, y_test_pred_lr)
confusion_matrix
array([[93,  1,  0,  0],
       [10, 71, 21,  0],
       [ 0, 22, 82, 16],
       [ 0,  0,  0, 84]])

print(classification_report(y_test, y_test_pred_lr))
```

	precision	recall	f1-score	support
0	0.90	0.99	0.94	94
1	0.76	0.70	0.72	102
2	0.80	0.68	0.74	120
3	0.84	1.00	0.91	84

Fig 7.3.1 Predicting test in logistic regression scaled

```
y_pred_lr_11 = lr.predict(X_test)

y_pred_lr_12 = lr.predict(X_train)

#test
confusion_matrix = metrics.confusion_matrix(y_test, y_pred_lr_11)
confusion_matrix
array([[82, 11, 1, 0],
       [25, 49, 23, 5],
       [ 0, 21, 54, 45],
       [ 0, 4, 8, 72]])
```

Fig 7.3.2 Predicting test in logistic regression

```
y_pred_dt = dt.predict(X_test)

print(metrics.confusion_matrix(y_test, y_pred_dt))

[[ 86  8  0  0]
 [ 7 84 11  0]
 [ 0 10 100 10]
 [ 0  0  6 78]]
```

Fig 7.3.3 Predicting test in Decision Tree

```
# Prediction on test data(unseen data)
y_pred_rf = rf.predict(X_test)
print(metrics.confusion_matrix(y_test, y_pred_rf))

[[ 87  7  0  0]
 [ 7 85 10  0]
 [ 0 14 100  6]
 [ 0  0  3 81]]
```

```
print(classification_report(y_test, y_pred_rf))
```

	precision	recall	f1-score	support
0	0.93	0.93	0.93	94
1	0.80	0.83	0.82	102

Fig 7.3.4 Predicting test in Random Forest

7.4 Validate the Models:

Model validation is the process of evaluating a trained model on test data set. This provides the generalization ability of a trained model. Here I provide a step by step approach to complete first iteration of model validation in minutes.

- The model are validate after completion of training and testing the model.
- Checking the accuracy scores as metrics to validate the models
- We have to check the accuracy among the models and validate best model among those.

```
acc_lr = metrics.accuracy_score(y_train, y_train_pred_lr) #train  
acc_lr  
0.86875
```

```
acc_lr = metrics.accuracy_score(y_test, y_test_pred_lr) #test  
acc_lr  
0.825
```

Fig 7.4.1 Test and train accuracy in logistic regression

```
models = ['training', 'testing']  
acc_scores = [0.86, 0.82]  
plt.bar(models, acc_scores, color=['lightblue', 'pink'])  
plt.ylabel("accuracy scores")  
plt.title("train vs test")  
plt.show()
```

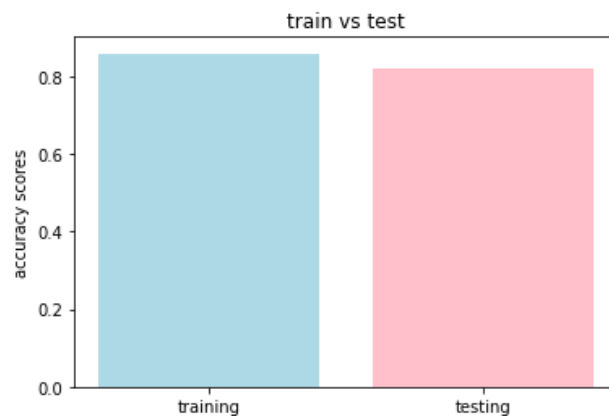


Fig 7.4.2 test vs train in logistic regression

Mobile Price Classification

Here, training accuracy is 86%

testing accuracy is 82%

```
acc_dt = metrics.accuracy_score(y_test, y_pred_dt)
acc_dt #test
```

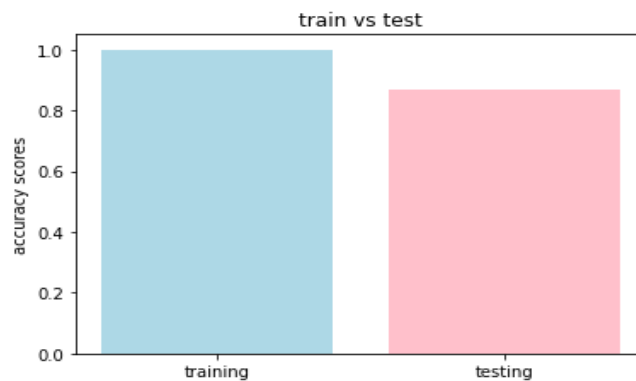
0.87

```
accuracy_score(y_train,y_train_pred_dt) #train
```

1.0

Fig 7.4.3 test and train accuracy in decision tree

```
models = ['training','testing']
acc_scores = [1,0.87]
plt.bar(models, acc_scores, color=['lightblue', 'pink' ])
plt.ylabel("accuracy scores")
plt.title("train vs test")
plt.show()
```



- Here, Training accuracy in decision tree model is 100%
- Here, Testing accuracy in decision tree model is 87%

Fig 7.4.4 test vs train in decision tree

```
acc_rf = metrics.accuracy_score(y_test, y_pred_rf)
acc_rf
```

0.8825

Fig 7.4.5 test accuracy in random forest

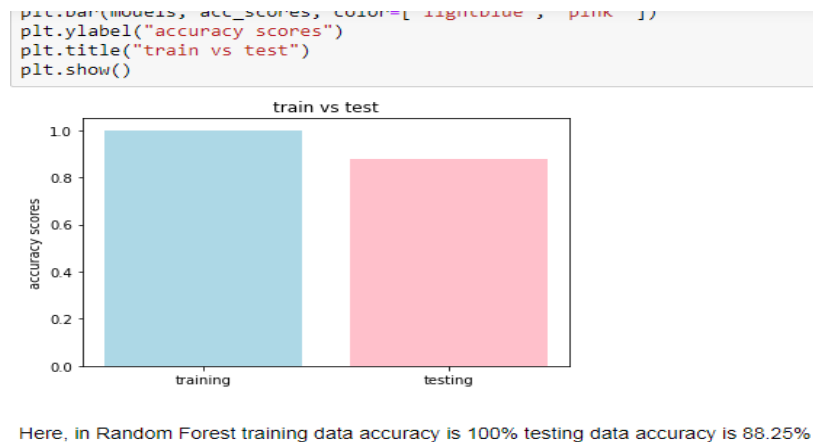


Fig 7.4.6 Test vs train in random forest

7.5 Parameter Tuning:

In machine learning, hyperparameter optimization or tuning is the problem of choosing a set of optimal hyperparameters for a learning algorithm. A hyperparameter is a parameter whose value is used to control the learning process. By contrast, the values of other parameters (typically node weights) are learned.

The same kind of machine learning model can require different constraints, weights or learning rates to generalize different data patterns. These measures are called hyperparameters, and have to be tuned so that the model can optimally solve the machine learning problem. Hyperparameter optimization finds a tuple of hyperparameters that yields an optimal model which minimizes a predefined loss function on given independent data. The objective function takes a tuple of hyperparameters and returns the associated loss. Cross-validation is often used to estimate this generalization performance.

i. Grid search:

The traditional way of performing hyperparameter optimization has been grid search, or a parameter sweep, which is simply an exhaustive searching through a manually specified subset of the hyperparameter space of a learning algorithm. A grid search algorithm must be guided by some performance metric, typically measured by cross-validation on the training set or evaluation on a held-out validation set.

Since the parameter space of a machine learner may include real-valued or unbounded value spaces for certain parameters, manually set bounds and discretization may be necessary before applying grid search.

For example, a typical soft-margin SVM classifier equipped with an RBF kernel has at least two hyperparameters that need to be tuned for good performance on unseen data: a regularization constant C and a kernel hyperparameter γ . Both parameters are continuous, so to perform grid search, one selects a finite set of "reasonable" values for each, say

ii. Random search:

Random Search replaces the exhaustive enumeration of all combinations by selecting them randomly. This can be simply applied to the discrete setting described above, but also generalizes to continuous and mixed spaces. It can outperform Grid search, especially when only a small number of hyperparameters affects the final performance of the machine learning algorithm. In this case, the optimization problem is said to have a low intrinsic dimensionality. Random Search is also embarrassingly parallel, and additionally allows the inclusion of prior knowledge by specifying the distribution from which to sample.

iii. Bayesian optimization:

Bayesian optimization is a global optimization method for noisy black-box functions. Applied to hyperparameter optimization, Bayesian optimization builds a probabilistic model of the function mapping from hyperparameter values to the objective evaluated on a validation set. By iteratively evaluating a promising hyperparameter configuration based on the current model, and then updating it, Bayesian optimization, aims to gather observations revealing as much information as possible about this function and, in particular, the location of the optimum. It tries to balance exploration (hyperparameters for which the outcome is most uncertain) and exploitation (hyperparameters expected close to the optimum). In practice, Bayesian optimization has been shown to obtain better results in fewer evaluations compared to grid search and random search, due to the ability to reason about the quality of experiments before they are run.

Gridsearch CV to increase accuracy of testing in decision tree

```
from sklearn.tree import DecisionTreeClassifier
dt=DecisionTreeClassifier()
grid_params={
    'criterion':['gini', 'entropy'],
    'splitter':['best','random'],
    'max_depth': range(1,11,2),
    'min_samples_leaf':range(1,6,3)
}
from sklearn.model_selection import GridSearchCV
grid=GridSearchCV(estimator=dt,param_grid=grid_params)
grid.fit(X_train,y_train)
grid.best_params_

{'criterion': 'entropy',
 'max_depth': 9,
 'min_samples_leaf': 1,
 'splitter': 'best'}

#Building the model with these params
final_model=DecisionTreeClassifier(criterion='entropy',
                                   max_depth=9,
                                   min_samples_leaf=6,
                                   splitter='best')
final_model.fit(X_train,y_train)
```

Fig 7.5.1 grid search cv

```
acc_dt = metrics.accuracy_score(y_test, pred_test)
acc_dt #test

0.8725
```

In grid based Cv Accuracy of decision tree is

- Here, Testing accuracy in decision tree model is 87.7%

Fig 7.5.2 accuracy of test using grid search cv i.e accuracy is increased in decision tree

7.6 Confusion Matrix in Algorithm:

A confusion matrix is a summary of prediction results on a classification problem. The number of correct and incorrect predictions are summarized with count values and broken down by each class. This is the key to the confusion matrix. The confusion matrix shows the ways in which your classification model is confused when it makes predictions. It gives us insight not only into the errors being made by a classifier but more importantly the types of errors that are being made.

```
# Prediction on test data(unseen data)
y_pred_rf = rf.predict(X_test)
clas_test_rf=metrics.confusion_matrix(y_test, y_pred_rf)
clas_test_rf
```

```
array([[ 87,   7,   0,   0],
       [  7,  85,  10,   0],
       [  0,  14, 100,   6],
       [  0,   0,   3,  81]])
```

Fig 7.6.1 Confusion Matrix of test dataset in Random Forest

```
sns.heatmap(clas_test_rf,annot=True,annot_kws={"size": 16}, fmt='g') #test heatmap
```

<matplotlib.axes._subplots.AxesSubplot at 0x7fbdfaa23ac8>

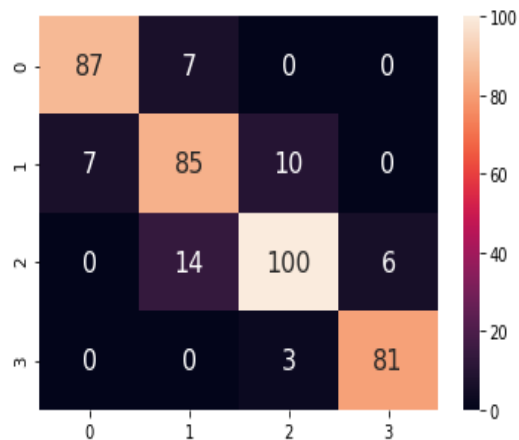


Fig 7.6.2 Heat Map of confusion matrix

- **Positive (P)** : Observation is positive.
- **Negative (N)** : Observation is not positive.
- **True Positive (TP)** : Observation is positive, and is predicted to be positive.
- **False Negative (FN)** : Observation is positive, but is predicted negative.
- **True Negative (TN)** : Observation is negative, and is predicted to be negative.
- **False Positive (FP)** : Observation is negative, but is predicted positive.

Classification Rate/Accuracy:

Classification Rate or Accuracy is given by the relation:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Fig 7.6.3 Accuracy

Recall:

$$\text{Recall} = \frac{TP}{TP + FN}$$

Fig 7.6.4 Recall

Precision:

$$\text{Precision} = \frac{TP}{TP + FP}$$

Fig 7.6.5 Precision

To get the value of precision we divide the total number of correctly classified positive examples by the total number of predicted positive examples. High Precision indicates an example labelled as positive is indeed positive (a small number of FP).

High recall, low precision: This means that most of the positive examples are correctly recognized (low FN) but there are a lot of false positives.

Low recall, high precision: This shows that we miss a lot of positive examples (high FN) but those we predict as positive are indeed positive (low FP).

7.7 Classification Report:

The classification report shows a representation of the main classification metrics on a per-class basis. This gives a deeper intuition of the classifier behavior over global accuracy which can mask functional weaknesses in one class of a multiclass problem. Visual classification reports are used to compare classification models to select models. The metrics are defined in terms of true and false positives, and true and false negatives. Positive and negative in this case are generic names for the classes of a binary classification problem. In the example above, we would consider true and false occupied and true and false unoccupied. Therefore a true positive is when the actual class is positive as is the estimated class. A false positive is when the actual class is negative but the estimated class is positive.

```
print(classification_report(y_test, y_pred_rf)) #test classification report
```

	precision	recall	f1-score	support
0	0.93	0.93	0.93	94
1	0.80	0.83	0.82	102
2	0.88	0.83	0.86	120
3	0.93	0.96	0.95	84
accuracy			0.88	400
macro avg	0.89	0.89	0.89	400
weighted avg	0.88	0.88	0.88	400

Fig 7.7.1 Test Classification Report

- **Precision:** Precision is the ability of a classifier not to label an instance positive that is actually negative. For each class it is defined as the ratio of true positives to the sum of true and false positives.
- **Recall:** Recall is the ability of a classifier to find all positive instances. For each class it is defined as the ratio of true positives to the sum of true positives and false negatives.
- **F1-Score :** The F_1 score is a weighted harmonic mean of precision and recall such that the best score is 1.0 and the worst is 0.0. Generally speaking, F_1 scores are lower than accuracy measures as they embed precision and recall into their computation. As a rule of thumb,

the weighted average of F_1 should be used to compare classifier models, not global accuracy.

- **Support:** Support is the number of actual occurrences of the class in the specified dataset. Imbalanced support in the training data may indicate structural weaknesses in the reported scores of the classifier and could indicate the need for stratified sampling or rebalancing. Support doesn't change between models but instead diagnoses the evaluation process.

7.8 Roc Curve:

The Receiver Operator Characteristic (ROC) curve is an evaluation metric for binary classification problems. It is a probability curve that plots the TPR against FPR at various threshold values and essentially separates the 'signal' from the 'noise'. The Area Under the Curve (AUC) is the measure of the ability of a classifier to distinguish between classes and is used as a summary of the ROC curve.

- The higher the AUC, the better the performance of the model at distinguishing between the positive and negative classes.
- When $AUC = 1$, then the classifier is able to perfectly distinguish between all the Positive and the Negative class points correctly. If, however, the AUC had been 0, then the classifier would be predicting all Negatives as Positives, and all Positives as Negatives.
- When $0.5 < AUC < 1$, there is a high chance that the classifier will be able to distinguish the positive class values from the negative class values. This is so because the classifier is able to detect more numbers of True positives and True negatives than False negatives and False positives.
- In a ROC curve, a higher X-axis value indicates a higher number of False positives than True negatives. While a higher Y-axis value indicates a higher number of True positives than False negatives. So, the choice of the threshold depends on the ability to balance between False positives and False negatives.

- There are useful statistics that can be calculated from this curve, like the Area Under the Curve (AUC) and the Youden Index. These tell you how well the model predicts and the optimal cut point for any given model
- ROC curves typically feature true positive rate on the Y axis, and false positive rate on the X axis. This means that the top left corner of the plot is the “ideal” point - a false positive rate of zero, and a true positive rate of one. This is not very realistic, but it does mean that a larger area under the curve (AUC) is usually better.

```
from sklearn.metrics import roc_curve, auc
# get false and true positive rates
fpr, tpr, thresholds = roc_curve(y_test, y_pred_proba_2[:,0], pos_label=0)
# get area under the curve
roc_auc = auc(fpr, tpr)
# PLOT ROC curve
plt.figure(dpi=90)
plt.plot(fpr, tpr, lw=1, color='green')
plt.title('ROC Curve for Logistic regression with scaling classifier')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate (Recall)')
plt.xlim([-0.05, 1.05])
plt.ylim([-0.05, 1.05])
plt.show()
```

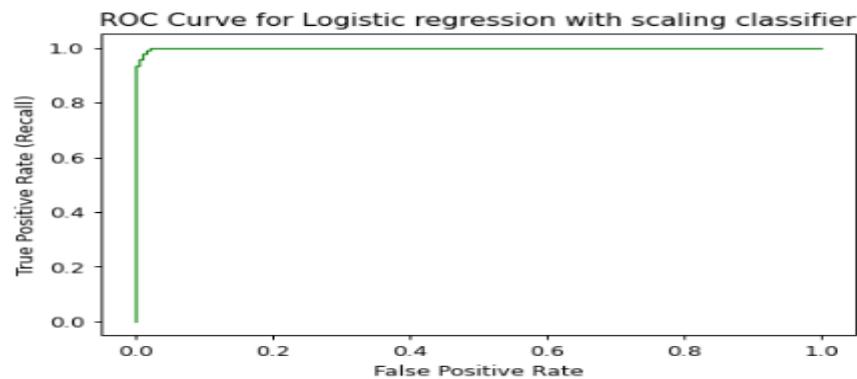


Fig 7.8.1 Roc curve for logistic regression

7.9 Predictions from raw data:

- After completion of accuracy we have to predict best model among them onto raw data
- Load the test dataset to predict
- Check the rows and columns and also shape

Mobile Price Classification

- After completion predict using best algorithm onto train dataset
- Add new column of predicted data i.e output
- Check the relation between test and train dataset to verify output

Best among scaled Logistic regression,logistic regression,decision tree,Random Forest

```
models = ['scaled Logistic regression','logistic regression', 'decision tree','Random Forest']
acc_scores = [0.82, 0.64, 0.87,0.882]
plt.bar(models, acc_scores, color=['lightblue', 'pink', 'lightgrey'])
plt.ylabel("accuracy scores")
plt.title("Which model is the most accurate?")
plt.show()
```

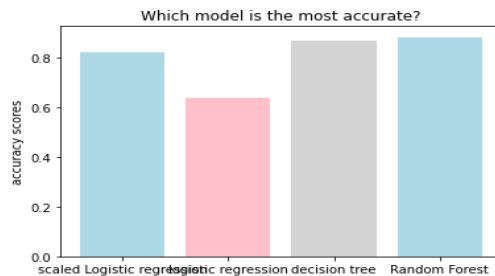


Fig 7.9.1 best accuracy model among algorithms

- Random forest shows best Accuracy among models so we choose Random forest for testing.

```
df_test=pd.read_csv('https://raw.githubusercontent.com/lohith4034/AIML/master/aiml%20project/test.csv')
```

```
df_test.head()
```

	id	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_cores	pc	px_height	px_width	
0	1	1043	1	1.8	1	14	0	5	0.1	193	3	16	226	1412	3
1	2	841	1	0.5	1	4	1	61	0.8	191	5	12	746	857	3
2	3	1807	1	2.8	0	1	0	27	0.9	186	3	4	1270	1366	2
3	4	1546	0	0.5	1	18	1	25	0.5	96	8	20	295	1752	3
4	5	1434	0	1.4	0	11	1	49	0.5	108	6	18	749	810	1

Fig 7.9.2 Loading Test dataset

- After loading train dataset using best algorithm predict the train dataset

Mobile Price Classification

```
predicted_price_range = rf.predict(df_test)
```

Fig 7.9.3 Predicting test dataset using random forest

```
df_test['price_range'] = predicted_price_range #new column price range is added
```

```
df_test.head() #new column price range is added
```

	battery_power	blue	clock_speed	dual_sim	fc	four_g	int_memory	m_dep	mobile_wt	n_cores	pr
0	1043	1	1.8	1	14	0	5	0.1	193	3	10
1	841	1	0.5	1	4	1	61	0.8	191	5	10
2	1807	1	2.8	0	1	0	27	0.9	186	3	10
3	1546	0	0.5	1	18	1	25	0.5	96	8	20
4	1434	0	1.4	0	11	1	49	0.5	108	6	10

Fig 7.9.4 Added new column as price range in test data set

- After adding new column verify the test and train dataset

```
df_train.groupby('price_range').size() #target variable size
```

```
price_range
0      500
1      500
2      500
3      500
dtype: int64
```

Fig 7.9.5 range of target in train dataset

Mobile Price Classification

```
df_train['price_range']
```

0	1
1	2
2	2
3	2
4	1
...	...
1995	0
1996	2
1997	3
1998	0
1999	3

Name: price_range, Length: 2000, dtype: int64

Fig 7.9.6 output of predicted test dataset

```
pd.DataFrame({'batterypower' : df_test['battery_power'],'price_range' : predicted_price_range})
```

	batterypower	price_range
0	1043	3
1	841	3
2	1807	2
3	1546	3
4	1434	1
...
995	1700	2
996	609	1
997	1185	0
998	1533	2
999	1270	2

1000 rows x 2 columns

Fig 7.9.7 example of predicted test dataset

Observations:

- If battery power is 1043 then price range is classified has 3
- It is similar to train dataset

8.Conclusion

It is concluded after performing thorough Exploratory Data analysis which include Statistics models which are computed to get accuracy and also Heat maps which are computed to get a clear understanding of the data set and its come to point of getting the solution for the problem statement being , that the bob should open a mobile shop by classified features of a mobile phone(eg:- RAM,Internal Memory etc) and its selling price and gives tough fight to big companies like Apple,Samsung etc, by classified selling price range .

9.REFERENCES

- https://en.wikipedia.org/wiki/Machine_learning
- <https://www.kaggle.com/iabhishekoofficial/mobile-price-classification/kernels>
- <https://mode.com/blog/python-data-visualization-libraries/>
- <https://www.hackerearth.com/practice/machine-learning/machine-learning-algorithms/ml-decision-tree/tutorial/>
- <https://builtin.com/data-science/random-forest-algorithm>
- <https://towardsdatascience.com/supervised-machine-learning-model-validation-a-step-by-step-approach-771109ae0253>
- **GIT HUB Link:** <https://github.com/lohith4034/AIML/tree/master/aiml%20project>