MACHINE LEARNING

(Predicting Chance of Admission)

Summer Internship Report Submitted in partial

fulfillment of the requirement for undergraduate degree

of

Bachelor of Technology

In

COMPUTER SCIENCE AND ENGINEERING

By

NEELA MADHU KUMAR

221710304038

Under the Guidance of

Assistant Professor



Department Of Computer Science and Engineering GITAM School of Technology GITAM (Deemed to be University) Hyderabad-502329 July 2020

DECLARATION

I submit this industrial training work entitled "Predicting Chance of Admission" 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: N. Madhu Kumar

Date: 12-07-2020 Student Roll No: 221710304038



ii

GITAM (DEEMED TO BE UNIVERSITY)

Hyderabad-502329, India

Dated:12-7-2020

CERTIFICATE

This is to certify that the Industrial Training Report entitled "**Predicting Chance of Admission**" is being submitted by Neela Madhu Kumar (221710304038) in partial fulfilment 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 2020-21

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

Dr. S.Phani Kumar

Professor and HOD

Assistant Professor

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 competition 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 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.

Neela Madhu Kumar 221710304038

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 predict "The Chance of Admission in to a University" 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. 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 Regressors 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 regressor for the given dataset. This work can be used to student to know the eligibility of admission in to a university

Table of Contents

1.MACHINE LEARNING	9
1.1. Introduction	9
1.2. Importance of Machine Learning	9
1.3. Uses of Machine Learning	10
1.4. Types of Machine Learning	13
1.4.1. Supervised Learning	14
1.4.2. Unsupervised Learning	14
1.4.3. Reinforcement Learning	15
2.DEEP LEARNING	16
2.1 Deep Learning and It's Importance	16
2.2 Uses of Deep Learning	17
2.3 Relation between Data Mining, Machine Learning and Deep	
Learning	19
3.PYTHON	21
3.1 Introduction	21
3.2 Setup of Python	21
3.2.1 Installation(using python IDLE)	22
3.2.2 Python Installation using Anaconda	23
3.3 Features	24
3.4 Variable Types	25
3.4.1 Python Numbers	25

3.4.2 Python Strings	26
3.4.3 Python lists	26
3.4.4 python tuples	26
3.4.5 python Dictionary	26
3.5 Functions	27
3.5.1 Defining a Function	27
3.5.2 Calling a Function	27
3.6 OOPs Concepts	27
3.6.1 Class	27
4. Predicting Chance of Admission	29
4.1 Project Requirements	29
4.1.1. Packages used	29
4.1.2. Versions of the packages	30
4.1.3. Algorithms used	30
4.2. Problem Statement	31
4.3. Dataset Description	31
4.4. Objective of the Case Study	31
5.DATA PREPROCESSING/FEATURE ENGINEERING AND EDA	32
5.1 Statistical Analysis	32
5.2 Data Type Conversions	34
5.3 Detection of Outliers	34

5.4 Handling Missing Values	35
5.5 Encoding Categorical Data	36
5.6 Generating Plots	37
5.6.1. Visualize the data between Target and the Features	37
5.6.2. Visualize the data between all the Features	42
6. FEATURE SELECTION	45
6.1 Select relevant features for the analysis	45
6.2 Drop irrelevant features	45
6.3 Train-Test-Split	45
6.4 Feature Scaling	47
7. MODEL BUILDING AND EVALUATION	48
7.1 Brief about the algorithms used	48
7.2 Train the Models	50
7.3 Make Predictions	51
7.4 Validate the Models	52
7.5 Parameter Tuning	55
i. Grid search:	55
ii.Bayesian optimization:	56
7.6 Predictions from raw data	58
8.Conclusion	62
9.REFERENCES	63

List of Figures:	
FIG 1.2.1USAGEOFMACHINELEARNINGINDIFFERENTFIELD	10
FIG 1.3.1 USESOFMACHINELEARNING	13
FIG 1.4.1 TYPESOFML	14
FIG 2.1.1 DEEPNEURALNETWORK	17
FIG 2.2.1 THEDEEP LEARNING PROCESS	18
FIG 2.3.1 RELATION BETWEEN DM,ML,DL	19
FIG2.3 PROCESS INMACHINELEARNINGANDDEEPLEARNING	19
FIG 3.2.1: PYTHON DOWNLOAD	22
FIG 3.2.1.1 PYTHONINSTALLATION	22
FIG 3.2.1.2 IDLE	23
FIG 3.2.2.1 AFTERINSTALLATION	24
FIG 3.2.2.2 JUPYTER NOTEBOOK	24
FIG 3.6.1 CLASS DEFINING	28
FIG 3.6.1.1 EXAMPLE OF CLASS	28
FIG 4.1.1 PACKAGES	30
FIG 4.1.2 VERSIONS	30
FIG 5.1.1 LOADINGDATASET	32
FIG 5.1.2 STATISTICALDATA	33
FIG 5.2.1 DATATYPES	34
FIG 5.4.1 There are no missing values in a Data set	35
FIG 5.5.1 NUMERICALINTHEDATASET	36
FIG 5.6.1 CORRELATION	37
FIG 5.6.1.1 CO-RELATIONGRAPH	38
FIG 5.6.1.2 CGPA vs chance of Admit	39
FIG 5.6.1.3 GRE vs chance of Admit	40
FIG 5.6.1.4 TOEFEL vs chance of Admin	42
FIG 5.6.2.1 Research	43

FIG 5.6.2.2 University Rating	44
FIG 6.3.1 Dividing Input and Output and Target values of Output	46
FIG 6.3.2 Importing Packages and Splitting	46
FIG 7.1.1 Linear Regression for Train Data	48
FIG 7.1.2 IMPORTINGDECISIONPACKAGES	49
FIG 7.1.3 RANDOMFORESTS	50
FIG 7.2.1 TRAININGDATA	51
FIG 7.3.1 PREDICTINGTESTINLINEARREGRESSION	51
FIG 7.3.2 PREDICTINGTESTINRANDOMFOREST	51
FIG 7.3.3 PREDICTINGTESTINDECISIONTREE	51
FIG 7.4.1 TESTANDTRAINACCURACYINLINEARREGRESSION	52
FIG 7.4.2 TESTVSTRAININLINEARREGRESSION	53
FIG 7.4.3 TESTANDTRAINSCOREINDECISIONTREE	53
FIG 7.4.4 TESTVSTRAININDECISIONTREE	54
FIG 7.4.5 TRAIN AND TEST SCOREINRANDOMFOREST	54
FIG 7.4.6 TESTVSTRAININRANDOMFOREST	55
FIG 7.5.1 Grid Search CV	58
FIG 7.5.2 Accuracy of test using Grid Search CV	58
FIG 7.6.1 BESTACCURACYMODELAMONGALGORITHMS	59
FIG 7.6.2 PREDICTMODELFORUNKNOWNLINEARREGRESSION	60

CHAPTER 1

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.

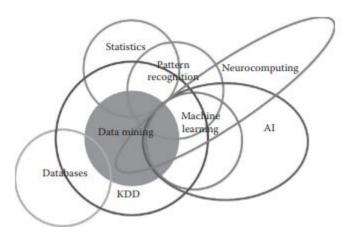


Fig 1.2.1 usage of Machine learning in different fields

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.

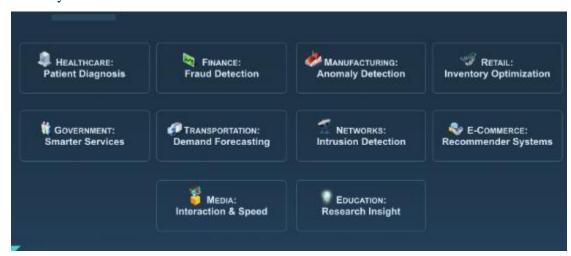


Fig 1.3.1 Uses of Machine learning

1.4. Types of Machine Learning:

There are 3 types of Machine learning which are widely used in todays world these Are:

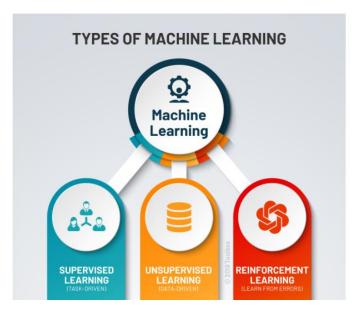


Fig 1.4.1 types of ML

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.

CHAPTER 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.

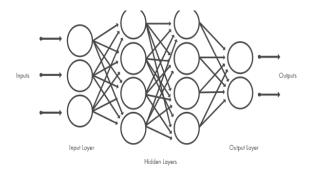


Fig 2.1.1 Deep Neural network

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.

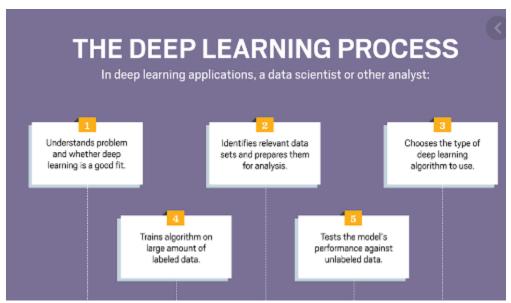


Fig 2.2.1 The deep learning process

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

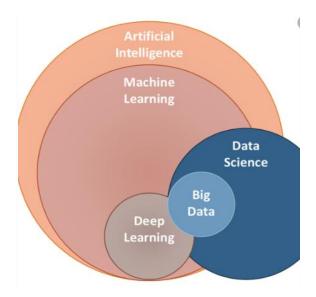


Fig 2.3.1 Relation between DM,ML,DL

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.

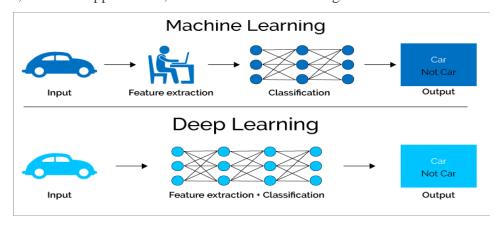


Fig 2.3.1 process in machine learning and deep learning

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.

CHAPTER 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 <u>python.org/downloads</u> 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



- 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

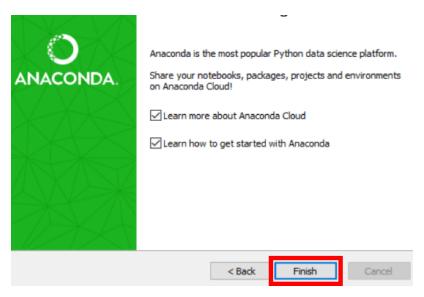


Fig 3.2.2.1 After installation



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 exception 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

CHAPTER 4

Predicting Chance of Admission

4.1 Project 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 that traditional Python lists. The array object in NumPy is called ndata, 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, WxPythonotTkinter. 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.

LOAD required packages

```
#importing the Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
matplotlib inline
import seaborn as sns
```

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: 1.18.1

pandas: 1.0.1

seaborn: 0.10.0

matplotlib: 3.1.3
```

Fig 4.1.2 versions

4.1.3 Algorithms used:

Here, 3 algorithms are used they are:

- Linear regression
- Decision Tree
- Random Forest

4.2 Problem Statement:

400 applicants have been surveyed as potential students for UCLA. The university weighs certain aspects of a student's education to determine their acceptance.

The objective is to explore what kind of data is provided, determine the most important factors that contribute to a student's chance of admission, and select the most accurate model to predict the probability of admission.

4.3 Dataset Description:

Given dataset containing various attributes of student information. The dataset that I am working is downloded from: https://www.kaggle.com/adepvenugopal/predicting-graduate-admissions-using-ml/data

This dataset is created for prediction of graduate admissions and the dataset link is below:

- Features in the dataset:
- GRE Scores (290 to 340)
- TOEFL Scores (92 to 120)
- University Rating (1 to 5)
- Statement of Purpose (1 to 5)
- Letter of Recommendation Strength (1 to 5)
- Undergraduate CGPA (6.8 to 9.92)
- Research Experience (0 or 1)
- Chance of Admit (0.34 to 0.97)

4.4 Objective of the Case Study:

The objective is to explore what kind of data is provided, determine the most important factors that contribute to a student's chance of admission, and select the most accurate model to predict the probability of admission.

Chapter 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.

```
#Read the Data
df = pd.read_csv("Admission_Predict.csv")
df.head()
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	1	337	118	4	4.5	4.5	9.65	1	0.92
1	2	324	107	4	4.0	4.5	8.87	1	0.76
2	3	316	104	3	3.0	3.5	8.00	1	0.72
3	4	322	110	3	3.5	2.5	8.67	1	0.80
4	5	314	103	2	2.0	3.0	8.21	0	0.65

Fig 5.1.1 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.

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
Serial No.	400.0	200.500000	115.614301	1.00	100.75	200.50	300.2500	400.00
GRE Score	400.0	316.807500	11.473646	290.00	308.00	317.00	325.0000	340.00
TOEFL Score	400.0	107.410000	6.069514	92.00	103.00	107.00	112.0000	120.00
University Rating	400.0	3.087500	1.143728	1.00	2.00	3.00	4.0000	5.00
SOP	400.0	3.400000	1.006869	1.00	2.50	3.50	4.0000	5.00
LOR	400.0	3.452500	0.898478	1.00	3.00	3.50	4.0000	5.00
CGPA	400.0	8.598925	0.596317	6.80	8.17	8.61	9.0625	9.92
Research	400.0	0.547500	0.498362	0.00	0.00	1.00	1.0000	1.00
Chance of Admit	400.0	0.724350	0.142609	0.34	0.64	0.73	0.8300	0.97

Fig 5.1.2 Statistical data

Observations:-

- The Average Score for GRE is 316.62
- The Average Score for TOEFL is 107.29
- The Average CGPA is 8.59
- The Average Number for Statement of Purpose is 3.0
- The Average Number for Recommendation letters among the students is 3.0
- The Average Number of Research done by students is 1.0
- The Average University Ratings of Different Students is 3.0

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 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.

:	Serial No.	int64
	GRE Score	int64
	TOEFL Score	int64
	University Rating	int64
	SOP	float64
	LOR	float64
	CGPA	float64
	Research	int64
	Chance of Admit	float64
	dtype: object	

Fig 5.2.1 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 is 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.

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.

df.isnull().sum()							
Serial No.	0						
GRE Score	0						
TOEFL Score	0						
University Rating	0						
SOP	0						
LOR	0						
CGPA	0						
Research	0						
Chance of Admit	0						
dtype: int64							

Fig 5.4.1 There are no missing values in a dataset

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.select_dtypes(include=['int','float']).columns #numerical
Index(['SOP', 'LOR ', 'CGPA', 'Chance of Admit '], dtype='object')
```

Fig 5.5.1 numerical in the dataset

5.6 Generating Plots:

5.6.1 Visualize the data between Target and the Features:

i. Correlation:

#finding correlation
df.corr()

	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
GRE Score	1.000000	0.835977	0.668976	0.612831	0.557555	0.833060	0.580391	0.802610
TOEFL Score	0.835977	1.000000	0.695590	0.657981	0.567721	0.828417	0.489858	0.791594
University Rating	0.668976	0.695590	1.000000	0.734523	0.660123	0.746479	0.447783	0.711250
SOP	0.612831	0.657981	0.734523	1.000000	0.729593	0.718144	0.444029	0.675732
LOR	0.557555	0.567721	0.660123	0.729593	1.000000	0.670211	0.396859	0.669889
CGPA	0.833060	0.828417	0.746479	0.718144	0.670211	1.000000	0.521654	0.873289
Research	0.580391	0.489858	0.447783	0.444029	0.396859	0.521654	1.000000	0.553202
Chance of Admit	0.802610	0.791594	0.711250	0.675732	0.669889	0.873289	0.553202	1.000000

Fig 5.6.1 correlation

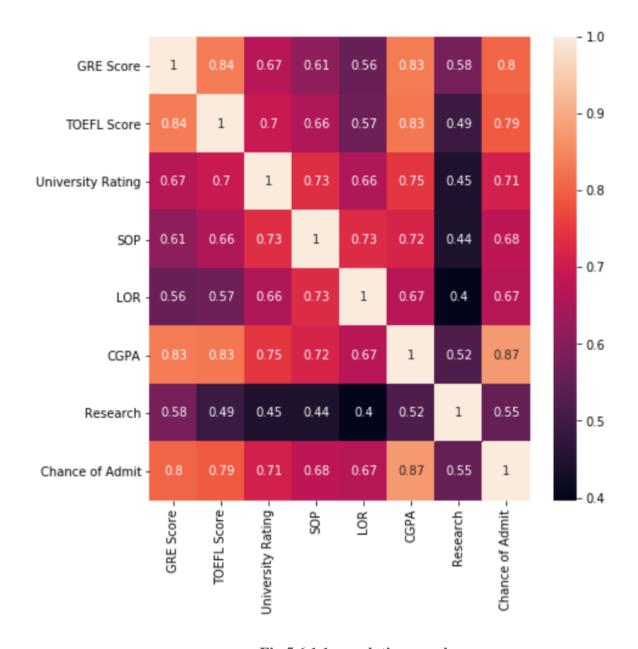


Fig 5.6.1.1 co-relation graph

The top three features that affect the Chance to Admit are:

1.CGPA

2.GRE Score

3.TOEFL Score

1. CGPA:

- The Cumulative Grade Point Average is a 10 point grading system.
- From the data shown below, it appears the submissions are normally distributed. With a mean of 8.6 and standard deviation of 0.6.

CGPA vs Chance of Admit:

• It appears as applicant's CGPA has a strong correlation with their chance of admission.

```
plt.figure(figsize=(20,6))
plt.subplot(1,2,1)
sns.distplot(df['CGPA'])
plt.title('CGPA Distribution')

plt.subplot(1,2,2)
sns.regplot("CGPA","Chance of Admit",data =df)
plt.title('CGPA vs Chance of Admit')
```

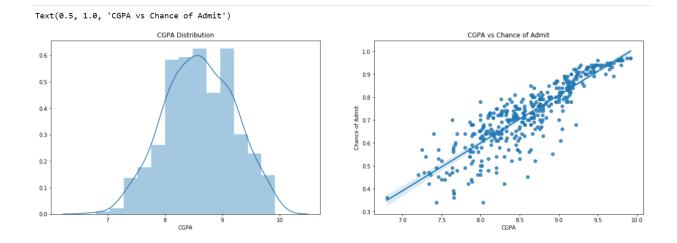


Fig 5.6.1.2 CGPA vs Chance of Admit

2. GRE Score:

The Graduate Record Examination is a standarized exam, often required for admission to graduate and MBA programs globally. It's made up of three components:

- 1. Analytical Writing (Scored on a 0-6 scale in half-point increments)
- 2. Verbal Reasoning (Scored on a 130-170 scale)
- 3. Quantitative Reasoning (Scored on a 130-170 scale)

In this dataset, the GRE Score is based on a maximum of 340 points. The mean is 317 with a standard deviation of 11.5.

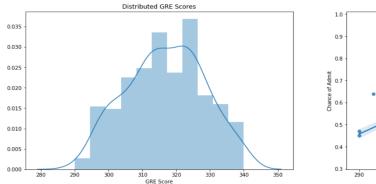
GRE vs Chance of Admit:

GRE scores have a strong correlation with the chance of admission however not as strong as one's CGPA.

```
plt.figure(figsize=(20,6))
plt.subplot(1,2,1)
sns.distplot(df['GRE Score'])
plt.title('Distributed GRE Scores')

plt.subplot(1,2,2)
sns.regplot(df['GRE Score'], df['Chance of Admit'])
plt.title('GRE Scores vs Chance of Admit')
```





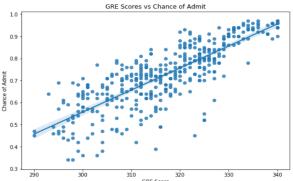


Fig 5.6.1.3 GRE vs Chance of Admit

3. TOEFL Score:

The Test of English as a Foreign Language is a standarized test for non-native English speakers that are choosing to enroll in English-speaking universities.

The test is split up into 4 sections:

- 1.Reading
- 2.Listening
- 3.Speaking
- 4.Writing

All sections are scored out of 30, giving the exam a total score of 120 marks. In this dataset, the TOEFL scores have a mean of 107 and a standard deviation of 6.

TOEFL Score vs Chance of Admit:

Like GRE scores, the scores received for the TOEFL strongly correlate to an applicants chance of admission.

```
plt.figure(figsize=(20,6))
plt.subplot(1,2,1)
sns.distplot(df['TOEFL Score'])
plt.title('Distributed TOEFL Scores')

plt.subplot(1,2,2)
sns.regplot(df['TOEFL Score'], df['Chance of Admit'])
plt.title('TOEFL Scores vs Chance of Admit')
```

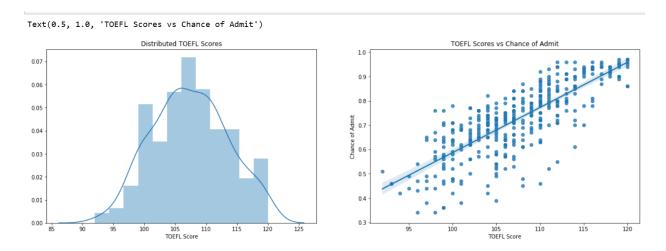


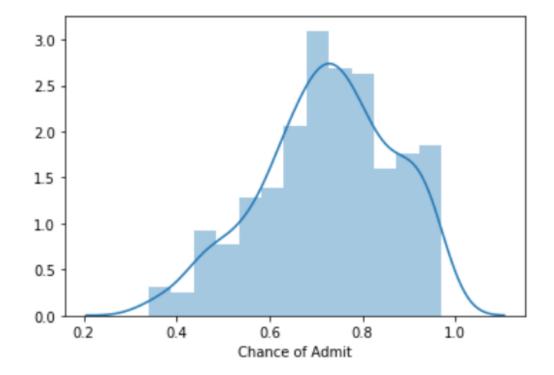
Fig 5.6.1.4 TOEFL vs Chance of Admit

5.6.2 - Visualize the data between all the Features:

i. Visualizing the Chance of Admit:

#Visualization of Distribution of Chance of Admit
sns.distplot(df['Chance of Admit'])

<matplotlib.axes._subplots.AxesSubplot at 0x17edee36188>



ii. Research:

• It seems the majority of applicants have research experience. However, this is the least important feature, so it doesn't matter all too much if an applicant has the experience or not.

```
fig, ax = plt.subplots(figsize=(8,6))
sns.countplot(df['Research'])
plt.title('Research Experience')
plt.ylabel('Number of Applicants')
ax.set_xticklabels(['No Research Experience', 'Has Research Experience'])
```

[Text(0, 0, 'No Research Experience'), Text(0, 0, 'Has Research Experience')]

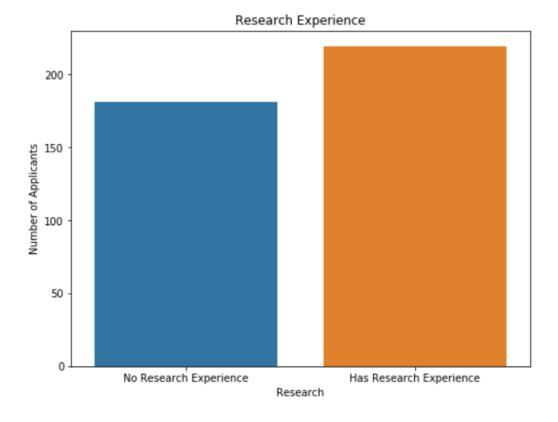


Fig 5.6.2.1 Research

iii. University Rating:

- Let's see the distribution of applicants coming from each kind of university.
- Most applicants come from a tier 3 and tier 2 university.

```
plt.subplots(figsize=(8,6))
sns.countplot(df['University Rating'])
plt.title('University Rating')
plt.ylabel('Number of Applicants')
```

Text(0, 0.5, 'Number of Applicants')

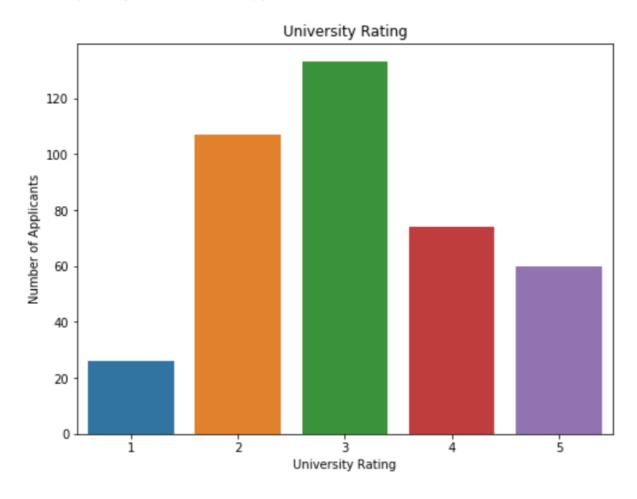


Fig 5.6.2.2 University Rating

Chapter 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

Note: No feature scaling is required for this model

6.2 Drop irrelevant features:

There are no irrelevant features in my dataset

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 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).

```
#splitting of input and output
X=df.drop(columns = {'Chance of Admit'})
y=df['Chance of Admit']
```

Fig 6.3.1 dividing input and output and target values for output

Importing packages:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.2, random_state=42)
```

Fig 6.3.2 importing packages and spliting with random state=42 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 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

Scaling is not required in this model because, it is Regression problem. If we perform scaling the values may not differ in high range that is slight changes may occur.

Chapter 7

MODEL BUILDING AND EVALUATION

7.1 Brief about the algorithms used:

i.Linear Regression:

Linear regression is probably one of the most important and widely used regression techniques. It's among the simplest regression methods. One of its main advantages is the ease of interpreting results

Problem Formulation:

When implementing linear regression of some dependent variable y on the set of independent variables $\mathbf{x} = (x_1, ..., x_r)$, where r is the number of predictors, you assume a linear relationship between y and \mathbf{x} : $y = \beta_0 + \beta_1 x_1 + \cdots + \beta_r x_r + \varepsilon$. This equation is the **regression equation**. $\beta_0, \beta_1, ..., \beta_r$ are the **regression coefficients**, and ε is the **random error**.

Linear regression calculates the **estimators** of the regression coefficients or simply the **predicted weights**, denoted with b_0 , b_1 , ..., b_r . They define the **estimated regression function** $f(\mathbf{x}) = b_0 + b_1x_1 + \cdots + b_rx_r$. This function should capture the dependencies between the inputs and output sufficiently well.

The **estimated** or **predicted response**, $f(\mathbf{x}_i)$, for each observation i = 1, ..., n, should be as close as possible to the corresponding **actual response** y_i . The differences $y_i - f(\mathbf{x}_i)$ for all observations i = 1, ..., n, are called the **residuals**. Regression is about determining the **best predicted weights**, that is the weights corresponding to the smallest residuals.

```
#Build the model on Training data
from sklearn.linear_model import LinearRegression
lm=LinearRegression()
lm.fit(X_train,y_train)
```

LinearRegression(copy X=True, fit intercept=True, n jobs=None, normalize=False)

Fig 7.1.1 linear 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 the 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.

Fig 7.1.2 importing decision packages

random_state=0, splitter='best')

min_weight_fraction_leaf=0.0, presort='deprecated',

iii. Random forest:

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

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'. 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=42)
```

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

```
: #predicting on test data
y_test_pred=lm.predict(X_test)
```

Fig 7.3.1 Predicting test in linear regression

```
#Predicting on the test data
y_predict = rfc.predict(X_test)
```

Fig 7.3.2 Predicting test in Random Forest

```
#predicting on test data
y_testpredictd = dec_tree.predict(X_test)
```

Fig 7.3.3 Predicting test in Decision Tree

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 R2 scores as metrics to validate the models
- Checking for root mean square error
- We have to check the accuracy among the models and validate best model among those.

```
#checking the R2 score, RMSE
#performance score of train data
from sklearn.metrics import r2_score, mean_squared_error
lr_r2 = r2_score(y_train, y_train_pred)
lr_mse = mean_squared_error(y_train, y_train_pred)
lr_rmse = np.sqrt(lr_mse)
print('Linear Regression R2 Score: {0} \nLinear Regression MSE: {1}, \nLinear Regression RMSE: {2}'.format(lr_r2, lr_mse, lr_rms
Linear Regression R2 Score: 0.7951946003191085
Linear Regression MSE: 0.0038704218434173236.
Linear Regression RMSE: 0.062212714483595104
#performance score of test data
from sklearn.metrics import r2_score, mean_squared_error
lrt_r2 = r2_score(y_test, y_test_pred)
lrt_mse = mean_squared_error(y_test, y_test_pred)
lrt_rmse = np.sqrt(lrt_mse)
print('Linear Regression RZ Score: {0} \nLinear Regression MSE: {1}, \nLinear Regression RMSE: {2}'.format(1rt_r2, 1rt_mse, 1rt_r
Linear Regression R2 Score: 0.8212082591486991
Linear Regression MSE: 0.004617003377285011,
Linear Regression RMSE:0.06794853476922819
```

Fig 7.4.1 Test and train score in linear regression

```
models=['training','testing']
r2_scores = [0.7951946003191085,0.8212082591486991]
plt.bar(models,r2_scores,color=['purple','pink'])
plt.ylabel("R2 scores")
plt.title("train vs test of Linear Regression")
plt.show()
```

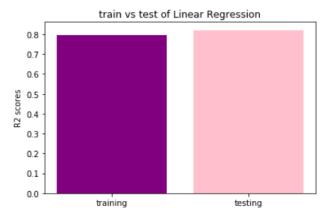


Fig 7.4.2 test vs train in linear regression

Here, Training score in model is 79% And Testing score in model is 82%

Here the unseen data (test data) works better than the train data

```
#performance score of test data
from sklearn.metrics import r2_score, mean_squared_error
dctr2 = r2_score(y_test, y_testpredictd)
dctmse = mean_squared_error(y_test, y_testpredictd)
dctrmse = np.sqrt(dctmse)
print('R2 Score: {0} ,\nMSE: {1}, \nRMSE:{2}'.format(dctr2, dctmse,dctrmse))

R2 Score: 0.7302776364703218 ,
MSE: 0.0069651375249012715,
RMSE:0.08345739946164912
```

```
#performance score of train data
from sklearn.metrics import r2_score, mean_squared_error
dcr2 = r2_score(y_train, y_tpredictd)
dcmse = mean_squared_error(y_train, y_tpredictd)
dcrmse = np.sqrt(dcmse)
print('R2 Score: {0} ,\nMSE: {1}, \nRMSE:{2}'.format(dcr2, dcmse,dcrmse))
```

R2 Score: 0.8812650216479114 , MSE: 0.0022438590706477537, RMSE: 0.04736938959547351

Fig 7.4.3 test and train score in decision tree

```
models=['training','testing']
r2_scores = [0.8812650216479114,0.7302776364703218]
plt.bar(models,r2_scores,color=['purple','pink'])
plt.ylabel("R2 scores")
plt.title("train vs test of Decision Tree")
plt.show()
```

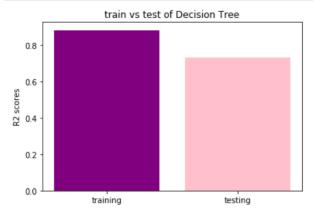


Fig 7.4.4 test vs train in decision tree

Here, the training score is 88% and testing score is 73%.

rfctmse = mean_squared_error(y_test, y_predict)

```
#performance score of train data
from sklearn.metrics import r2_score, mean_squared_error
rfcr2 = r2_score(y_train, y_train_predrfc)
rfcmse = mean_squared_error(y_train, y_train_predrfc)
rfcrmse = np.sqrt(rfcmse)
print('R2 Score: {0} \nMSE: {1}, \nRMSE:{2}'.format(rfcr2, rfcmse, rfcrmse))

R2 Score: 0.906390038230045
MSE: 0.0017690453540795885,
RMSE:0.042060020852105966

#performance score of test data
from sklearn.metrics import r2_score, mean_squared_error
rfctr2 = r2_score(y_test, y_predict)
```

R2 Score: 0.8183993845316527 MSE: 0.004689537955960314, RMSE:0.06848020119684459

rfctrmse = np.sqrt(rfctmse)

Fig 7.4.5 train and test score in random forest

print('R2 Score: {0} \nMSE: {1}, \nRMSE:{2}'.format(rfctr2, rfctmse, rfctrmse))

```
models=['training','testing']
r2_scores = [0.906390038230045,0.8183993845316527]
plt.bar(models,r2_scores,color=['purple','pink'])
plt.ylabel("R2 scores")
plt.title("train vs test of Random Forest")
plt.show()
```



Fig 7.4.6 Test vs train in random forest

Here, the training score is 90% and the testing score is 81%.

7.5 Parameter Tuning:

In machine learning, hyperparameter optimization or tuning is the problem of choosing a set of optimal <u>hyperparameters</u> 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:

mam 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

```
grid_param = {
    'criterion': ['mse'],
    'max_depth' : range(2,20,1),
    'min_samples_leaf' : range(1,10,1)
}
```

```
#Import the GridSearchCV
from sklearn.model_selection import GridSearchCV

# initialization of GridSearch with the parameters- ModelName and the dictionary of parameters

clf = DecisionTreeRegressor()
grid_search = GridSearchCV(estimator=clf, param_grid=grid_param)

# applying gridsearch onto dataset
grid_search.fit(X_train, y_train)
```

```
GridSearchCV(cv=None, error_score=nan,
             estimator=DecisionTreeRegressor(ccp_alpha=0.0, criterion='mse',
                                             max_depth=None, max_features=None,
                                             max_leaf_nodes=None,
                                             min_impurity_decrease=0.0,
                                             min_impurity_split=None,
                                             min_samples_leaf=1,
                                             min samples split=2,
                                             min_weight_fraction_leaf=0.0,
                                             presort='deprecated',
                                             random_state=None,
                                             splitter='best'),
             iid='deprecated', n_jobs=None,
             param_grid={'criterion': ['mse'], 'max_depth': range(2, 20),
                          'min_samples_leaf': range(1, 10)},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
             scoring=None, verbose=0)
```

Fig 7.5.1 grid search cv

```
#performance score of test data after using GridsearchCV
from sklearn.metrics import r2_score, mean_squared_error
gridr2 = r2_score(y_test, gridtest)
gridmse = mean_squared_error(y_test, gridtest)
gridrmse = np.sqrt(gridmse)
print('R2 Score: {0} \nMSE: {1}, \nRMSE:{2}'.format(gridr2, gridmse, gridrmse))

R2 Score: 0.7939106007401951
MSE: 0.005321920620463799,
RMSE: 0.07295149498443332
```

In Grid search CV R2 score of decision tree is

- Here, Testing R2 score in decision tree is 79%
- score as increased from 73% to 79% after using parameter tunning

Fig 7.5.2 Accuracy of test using grid search cv

7.6 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
- After completion predict using best algorithm onto train dataset
- Check the relation between test and train dataset to verify output

```
model = ['linear regression','Random Forest','Decision Tree']
R2_scores = [0.8212082591486991,0.8183993845316527,0.793910600740195]
plt.bar(model, R2_scores, color=['blue', 'pink', 'lightgrey'])
plt.ylabel("R2 scores")
plt.title("Which model is the most accurate?")
plt.show()
```

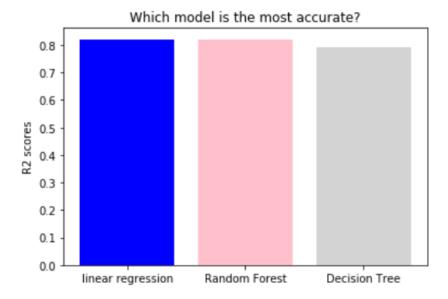


Fig 7.6.1 best accuracy model among algorithms

- Linear Regression shows the best score among model so we choose Linear Regression for predicting the Chance of Admit in this model
- After loading train dataset using best algorithm predict the train dataset

```
#predicting the model for unknown features using Linear Regression
un=lm.predict([[315,112,2,3.0,2.5,8.42,1]])
un
```

array([0.7025534])

```
1 un1=lm.predict([[300,128,4,3.5,3.0,7.89,0]])
2 un1
```

array([0.65768669])

Fig 7.6.2 Predicting model for unknown using Linear regression

Metrics:

i. R2-Score:

- The r2_score function computes the coefficient of determination, usually denoted as R2.
- It represents the proportion of variance (of y) that has been explained by the independent variables in the model. It provides an indication of goodness of fit and therefore a measure of how well unseen samples are likely to be predicted by the model, through the proportion of explained variance.
- If y^i is the predicted value of the i-th sample and yi is the corresponding true value for total n samples, the estimated R² is defined as:

$$R^2(y,\hat{y}) = 1 - rac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - ar{y})^2}$$

ii. Mean Squared Error:

• The mean_squared_error function computes mean square error, a risk metric corresponding to the expected value of the squared (quadratic) error or loss.

• If y^i is the predicted value of the i-th sample, and yi is the corresponding true value, then the mean squared error (MSE) estimated over nsamples is defined as

$$ext{MSE}(y, \hat{y}) = rac{1}{n_{ ext{samples}}} \sum_{i=0}^{n_{ ext{samples}}-1} (y_i - \hat{y}_i)^2.$$

Chapter 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 is based on important factors that contribute to a student's chance of admission, and selected the most accurate model and predicted the probability of admission.

Linear Regression shows the best score among model so we choose Linear Regression for predicting the Chance of Admit in this model

Chapter 9

REFERENCES

- https://en.wikipedia.org/wiki/Machine_learning
- https://www.kaggle.com/adepvenugopal/predicting-graduate-admissions-using-ml/data
- 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

Github Link: https://github.com/neelamadhu/AIML/tree/master/AIMLproject