# PROPOSED SOLUTION

## INTRODUCTION

Student admission for the Master's degree program consists of different criteria/scores which is taken into consideration before admitting the student to the degree program. This process is elaborative and requires lot of thought processing and analysis by the selection committee before choosing the right applicants to the Master's degree program. The purpose of this analysis is to demonstrate the top contributing scores which helps the student to get the admission into the Master's degree program. What factors contributes to successful admission to a Master's degree program? The analysis might seem straight forward but caution has to be exercised to consider the scores like GRE, TOEFL, university rating, SOP, LOR and CGPA and any outliers should not impact the decision making process.

## ABSTRACT

The primary purpose is to discuss the prediction of student admission to university based on numerous factors and using logistic regression. Many prospective students apply for Master's programs. The admission decision depends on criteria within the particular college or degree program. The independent variables in this study will be measured statistically to predict graduate school admission. Exploration and data analysis, if successful, would allow predictive models to allow better prioritization of the applicants screening process to Master's degree programme which in turn provides the admission to the right candidates

## Methodology

### Business Understanding

Initially good amount of time was spent on understanding the problem statement by understanding the concerns of students regarding the current application process, the objectives of the research were defined in this process.

### Data Understanding

Data required for the research was collected from multiple data sources. Different features of the data were analyzed based on their importance and relevance. Data-set would be explained in more detail further.

### Data Preparation

In this phase, the data from multiple data sources were integrated into a final data-set. Further the data was cleaned by removing unwanted columns, performing transformation and cleaning activities on the data.

### Modelling

Multiple machine learning models were developed to predict the likelihood of success of the student's application in a particular university. The user interface was developed to allow the users to access these models.

### Evaluation

Models developed were evaluated based on their performance and accuracy. More information will be presented in the evaluation section of the paper.

# Algorithms

Multiple machine learning algorithms were used for this research, K- Nearest Neighbour and Multivariate Logistic Regression algorithms were used to predict the likelihood of the students getting admission into university based on their profile. Decision Tree algorithm was used to predict the rank of the college that would be suitable for the students based on their profile and suggest the list of universities accordingly.

### K-Nearest Neighbours

It is an algorithm which is used widely for classification and regression problems. Due to its simplicity and effectiveness, it is easy to implement and understand. It is a supervised machine learning algorithm that uses available data to create the model and further that model can be applied to classify the new data. The class of new data is determined by the class of its neighbours. Distance is calculated between the unseen data sample and the all other data samples already present in the data-set. Depending on the value of K, that many nearest neighbours are selected and their class is identified. The class of neighbours which has majority is assigned to the class of the new data sample. Generally, Euclidean distance is used to calculate the distance between the records. Multiple values of K should be tried and tested, and the value of K at which best performance is observed must be selected for the model.

### Logistic Regression

Logistic regression algorithm is used to identify the probability of occurrence of an event based on single predictor variable. Multivariate Logistic regression can be used to determine the probability of the occurrence of an event based on multiple predictor variables. The class variable that has to be predicted has to be binary or dichotomous. Logistic Regression is also a supervised machine learning algorithm which used data with predetermined classes to create a model and perform predictive analysis on unseen data.

### Decision Tree

It is a supervised machine learning algorithm. Due to its simple logic, effectiveness and interpretability it the most widely used classification algorithm. The model works by creating a tree-like structure by dividing the data-set into several smaller subsets based on different conditional logic. The main components of the decision tree are the decision nodes, leaf nodes and the branches. Nodes with multiple branches are the decision nodes, nodes with no branches are called the leaf nodes, and the top node is called the root node of the decision tree. The nodes are connected to each other via branches based which are different conditions. The root and decision nodes are created by computing the entropy and information gain for the data-set.