

# ADS2001

## Virtual Internship

# Nephrotex



MONASH  
University

By Isaac Wood, Louise Childs, Aiden  
Abraham, Yilin Han, Oscar Brazzale

## Table of Contents

Executive Summary .....	3
Brief Summary .....	5
Project Description and Details .....	5
Background Research .....	5
Dataset .....	6
Aim .....	7
Preprocessing and Data Wrangling .....	8
Understanding the Data .....	8
Cleaning the Data .....	9
Extracting Necessary Data .....	10
Exploratory Data Analysis .....	11
Method .....	11
Basic Statistics .....	11
Modelling .....	12
Linear Regression .....	12
Decision Trees.....	13
Random Forests .....	14
K-Nearest Neighbour .....	15
KMeans Clustering .....	16
Boosting .....	17
SMOTE .....	18
Natural Language Processing .....	19
Text-Data Pre-Processing .....	19
TF-IDF.....	19
New Features .....	20
Modelling with New Features .....	21
Results .....	22
Model Results.....	22
Conclusion.....	23
Appendix .....	24

# Executive Summary

Virtual Internships is a company that runs academic simulations for students in the US. The Nephrotex simulation is run based on a bioengineering design challenge. The data that was explored covered 15 implementations of the challenge with 5 groups of varying sizes from 4-7 students. The data that was collected consisted mainly of the chats between groups. This included between members and with mentors. The other part of the data was the outcome score. Students were marked on a scale of 0-8 based on an individual submission they made at the completion of the design challenge, this was the outcome score. The aim of the investigation was to attempt to use the given data to model the outcome score of the students.

The dataset that was worked with was reasonably clean and usable, there were small errors in the data frame that were removed due to the irrelevance to the overall project. From this new clean data two new datasets were formed, one with the contributions of individual users and their outcome scores, and another with the summed contributions of entire groups and the average score per user in that group.

Some initial testing and analysis found that the most common outcome score was 4 with just over 1/3<sup>rd</sup> of all students receiving that outcome score. It was also found to be incredibly difficult to get the higher scores with only 2.4% of students receiving a mark above a 6. This initial analysis led to the idea that modelling the data may become difficult with bias potentially coming into the model with the same scores being present so often and a lack of higher scores to potentially identify what traits lead to a higher score.

Some of the problems encountered when working with the data were, the inconsistencies in the number of people in each group, and the individual scoring within the same groups. This leads to difficulty associating chats to outcome score which is the initial problem that is explored. To help combat this problem, two different approaches were tested. First approach was to focus on the individuals who received high scores and analyse their chats to see the key ideas/words that were used and see if that had a direct effect on the outcome score or not. The other approach was to focus on the groups who had the highest average score across their members and focus on the content of the chats between the groups. Using this, a key set of words was found and this was used to help initial analysis of data, and later data modelling. This was found using NLP programming, more specifically TF-IDF.

Modelling was the next step in investigating the relationship between the chats and the outcome score. Many models were attempted with importance being placed on the decision tree classifier and the random forest modelling techniques. Due to the output of the model being discrete, these models were deemed most appropriate. Other factors included the binary classifications of the chats that seemed to fit the idea of the classifier models.

A variety of models were implemented in order to try and predict the final outcome scores of the students, and the overall groups, with varying accuracies. From these implemented models the most accurate for predicting individual user score was seen to be a Decision Tree Classifier, returning an accuracy of 0.311. For predicting the average score for a user in a group the most accurate model was seen to be a K Nearest Neighbours model which resulted in an accuracy of 0.370. Other unsuccessful models used included Linear Regression, Random Forest Classification, kMeans Clustering, AdaBoost and SMOTE.

The use of the NLP words was then applied to the models to see if they improved the accuracy. Unfortunately the use of the additional important words in the model did not lead to improved scores across any of the appropriate models.

All of these approaches promised their own limitations due to the nature of the task at hand, the chat data had little to no effect on the group outcome score as outcome score was scored on a group level and a third of the students achieved an outcome score of 4, below average groups were still just as active in the chats as the higher scoring teams and they still discuss the relevant terms. The virtual internship also takes place over 18 hours and many implementations occur at the same time so definitive conclusions are hard to determine. As we were faced with unbalanced data due to the frequency of an outcome score of 4, SMOTE, which is an improved overbalancing technique which utilises data augmentation could be used, this however reduces false negatives and increases true positives at the cost of accuracy, this wasn't suitable for our analysis due to the low accuracies to begin with, however a future in depth data analysis could investigate measures other than accuracies such as precision and recall when providing analytical advice. The natural language processing analysis also faced limitations as TF-IDF uses frequency of words to determine their importance and doesn't provide meaning to the words, future data analysis should delve into sentiment analysis so the words in the dataset can be provided meaning and used to draw concrete connections with the outcome score.

# Main Body

## Brief Summary

Throughout the weeks before the submission of our final report and our preparation for the presentation, every member of the group worked regularly and diligently. Members completed their work on time and contributed effectively to team discussions on Thursdays by offering their own thoughts to the group debate. Despite the fact that some of the group members lacked expertise in specific areas, other teammates stepped up to provide support and assistance. All members had developed valuable skills in both the technical and interpersonal components of excellent team discussion by the end of this project. Overall, everyone in the group contributed to the project's successful conclusion.

## Project Description and Details

### Background Research

Virtual internships are educational simulations designed to help students acquire the critical thinking, acting, communication, and argumentation abilities needed for a given career field. Our study focuses on Nephrotex, a virtual internship programme that is unique in that interns create teams in an online setting and communicate mostly using a chat application. This arrangement also improves communication between interns and their mentors, whose responsibility it is to answer queries and spark stimulating debates.

In Nephrotex's virtual environment, students work as interns for a biomedical engineering company and are required to collaborate in a team to design a prototype device that can help patients with kidney failure. During the internship, the students were required to conduct preliminary research, understand stakeholder needs, design and develop a prototype, as well as test and evaluate the prototype and rationalise their design decisions.

The design challenge also involves responding to stakeholders in the fictitious company. These stakeholders often have varying demands regarding the 5 key outputs of the project. The 5 key outputs are, Biocompatibility, Reliability, Marketability, Cost and Flux. These 5 outputs are important when analysing the importance of the chats between the groups.

## Dataset

In this project, we used a dataset containing chats from 15 implementations of the Nephrotex internship. These chats have been annotated to identify the presence or absence of key concepts related to specific design actions and arguments in engineering discussions. In addition, the dataset also includes a rating of each student's final design report during the internship, which is an important metric for assessing the effectiveness of the internship.

There are 16 factors used as explanatory variables in the investigation :

Features	Explanation
<b>userIDs</b>	A unique id for each student
<b>implementation</b>	A unique id for each implementation
<b>Line_ID</b>	A unique id for each chat utterance
<b>ChatGroup</b>	A string value indicates the team for the first half of the internship, while a numerical value indicates the team for the second half of the internship. Midway through the internship, individual members will change teams
<b>content</b>	The content of each chat utterance
<b>group_id</b>	A non-unique numerical id for the team the chatter was in
<b>RoleName</b>	Whether the chatter was a mentor or a player
<b>roomName</b>	A non-unique name for the internship activity the chatter was participating in.

<b>m_experimental_testing</b>	Talk about using experimental techniques to understand the technical features of design
<b>m_making_design_choices</b>	Talk about choosing a specification or characteristic for a design
<b>m_asking_questions</b>	Asking questions
<b>J_customer_consultant_requests</b>	Talk about justifying design choices by stating that they should meet or exceed stakeholder requests.
<b>J_performance_parameters_requirements</b>	Talk about justifying design choices by referring to performance parameters or experimental results
<b>j_communication</b>	Talk about justifying design choices by referring to facilitating communication among the engineers
<b>OutcomeScore</b>	A mark from 0 (low) to 8 (high) indicating the quality of the chatter's final design report
<b>wordCount</b>	The word count for the chat utterance

## Aim

The aim of this project is to use team-level discussion data from the internship to model and predict final report performance. Specifically, our project will be split into two sections. First, we will transform the data to the team-level statistic, clean it up, format it, and then apply machine learning models and create new features to produce the best possible forecasts. Extrapolate lessons from these prediction models about the relationship between discourse features and report scores in this setting, and look into the mentor's potential influence on these associations. This can provide valuable information to the supervisor, identify potential problems in a timely manner and improve the overall quality of the placement.

# Pre-Processing and Data Wrangling

## Understanding The Data

To understand the data, we firstly displayed our dataset by calling `df.shape()` to check the dimensionality of the DataFrame. It's shown that there are 19180 rows and 17 columns in the dataframe. It has been observed that the dataset size is extremely vast, which makes developing a model with a large number of observations inefficient. As a result, the data must be cleaned before it can be modelled. This will also improve the accuracy and efficiency of a machine learning model. Then we implemented the `df.dtypes()` to help us to get to know the data types in a dataframe.

For this image we can see that there are 5 columns containing objects( like strings) and others all contain the integer values. This is useful later on in data cleansing because it informs us about the type of data we're dealing with. This is due to the fact that certain operations can only be performed on specific data types (for example, you can't conduct math on strings or concatenate values). If a column that is supposed to be a number (int64 or float64) has a datatype object, it may indicate that some values are not numbers, possibly due to missing or malformed data. Also many machine learning models require the input data

Unnamed: 0	int64
userIDs	int64
implementation	object
Line_ID	int64
ChatGroup	object
content	object
group_id	int64
RoleName	object
roomName	object
m_experimental_testing	int64
m_making_design_choices	int64
m_asking_questions	int64
j_customer_consultants_requests	int64
j_performance_parameters_requirements	int64
j_communication	int64
OutcomeScore	int64
wordCount	int64
dtype: object	

to be in numeric format. Knowing your data type is useful at the pre-processing stage, when you may need to normalise numeric variables. After those, we displayed by calling `df.describe()` to provide a basic descriptive statistics for our dataset. The output shows that the data was found to be of high quality. One small problem was the overlap of group numbers across different implementations, this was easily solved by creating a unique group ID, using group number and implementation.

During the data processing phase, we encountered an issue with inconsistency in the number of members in each group and the scores of individuals inside the group. This made tying the discussion logs to the final scores more challenging, and it became a problem that we needed to debate and fix at first. To address this, we tested and implemented two distinct ways. In the first approach, we focused on high scoring individuals and analysed their chat logs in depth, aiming to understand whether the key ideas or words they used had a direct impact on the outcome score. In the other strategy, we focused on those groups with the highest average member scores and



delved into the content of the chats within the group. Through this approach, we identified a set of key words that played a key supporting role in both our initial data analysis and subsequent data modelling. We identified these keywords by using natural language processing (NLP) techniques, and more specifically, we used the word frequency-inverse document frequency (TF-IDF) method.

## Data Cleaning

Before deciding which way to approach the missing values within the dataset, we as data scientists must understand why there are missing values. Data scientists have grouped missing data into three categories which includes Missing at Random (MAR), Missing Completely at Random (MCAR) and Missing Not at Random (MNAR). (Swalin, 2018)

1. Missing at Random (MAR): Data may be missing due to external reasons that may not be able to be predicted. The values are just completely missing at random.
2. Missing Completely at Random (MCAR): Certain values of the data could be missing but have nothing to do with its hypothetical values. For example, people are given weights to weigh themselves for research but some of them are reluctant to weigh themselves. Therefore, the missing data cannot be concluded that it is missing due to observational error.
3. Missing not at Random (MNAR): Not missing at random but are ignorable. The missing value of the data can be determined by the value of interest.

In this case, there may be a number of factors causing missing values in the data. For example, users (mentors or players) may choose not to participate in certain chats, or their responses may be deleted by mistake or not saved correctly. There is no perfect way to handle missing values in data but what we can do is to minimise it if it does appear.

There are two main ways to handle missing values. One being deletion and second being imputation (Dixon, 2021). According to the number of entries, by calling the function `df.isnull().sum()`, the output shows that the RoleName got 3 missing values, then we deleted that column. After deleting the columns, we use `isnull().sum()` on the dataframe to check the number of missing values in each column again, so there are no missing values for the RoleName column.

Unnamed: 0	0	Unnamed: 0	0
userIDs	0	userIDs	0
implementation	0	implementation	0
Line_ID	0	Line_ID	0
ChatGroup	0	ChatGroup	0
content	0	content	0
group_id	0	group_id	0
RoleName	3	RoleName	0
roomName	0	roomName	0
m_experimental_testing	0	m_experimental_testing	0
m_making_design_choices	0	m_making_design_choices	0
m_asking_questions	0	m_asking_questions	0
j_customer_consultants_requests	0	j_customer_consultants_requests	0
j_performance_parameters_requirements	0	j_performance_parameters_requirements	0
j_communication	0	j_communication	0
OutcomeScore	0	OutcomeScore	0
wordCount	0	wordCount	0
dtype: int64		dtype: int64	

## Extracting Necessary Data

In order to eventually start predicting final outcome scores the data needed to be formed into two new data sets, grouped by users and by teams, so that we can attempt to predict the outcome of individual users and the average score of entire groups. This was done by grouping the data frame by either 'userID' or 'Group,' then finding the sum of this table to find the amount of times each contribution was made. However, for users as the 'OutcomeScore' column was now their actual outcome score multiplied by the amount of times they appeared in the chat a new column was needed to be added with this larger incorrect outcome score divided by the number of contributions.

For the group model this outcome score then also needed to be divided by the amount of people within the group in order to get the average score per student which ultimately ranged from 2 to 6.

# Exploratory Data Analysis

## Method

Given the project's aim and the datasets, some analytical techniques were more pertinent than others. As we are dealing with continuous data and building the relationship between variables, we chose to use techniques such as linear regression, decision tree, and random forest regression. The feature importance will assist to determine which features are directly influencing the model's performance, we will try different combinations of features and assess their impact on the model's performance. Furthermore, we use visualisations including count plots, scatter plots, and heatmap etc. to assist our approach. In addition, we have separated the user model from the group model, which allows us to better analyse individual and collective behaviour and performance separately, which can provide unique insights.

## Basic Statistics

The Nephrotex data was based on a virtual internship program that was run. This was an engineering design challenge that was given to groups of interns to work on. There were 15 different implementations of the design challenge with each implementation having 5 groups, and each group containing 4-7 members. At the end of the challenge each individual submitted a design and report that was then marked on a scale of 0-8 with 8 being the highest score. The spread of scores was extremely interesting, it appears to be extremely difficult to get a higher score with only 2.4% of interns scoring a 7 or 8. On the contrary, achieving a 4 was the most common outcome with 33.8% of interns receiving a score of 4.

When the different amounts of each feature were plotted against user score and average group score, as seen in Figures 1 and 2, it was clear that for the individual users communicating about the performance parameters and customer consultant requests had more of a relationship with a higher outcome score. For the group scores it can be seen that making design choices, experimental testing and word count contributed more to an improved final outcome score. However, in all of these graphs it can be seen that there is a lot of variance within the data. This is likely due to the very few users who achieved a higher score of 7 or 8 as discussed above.

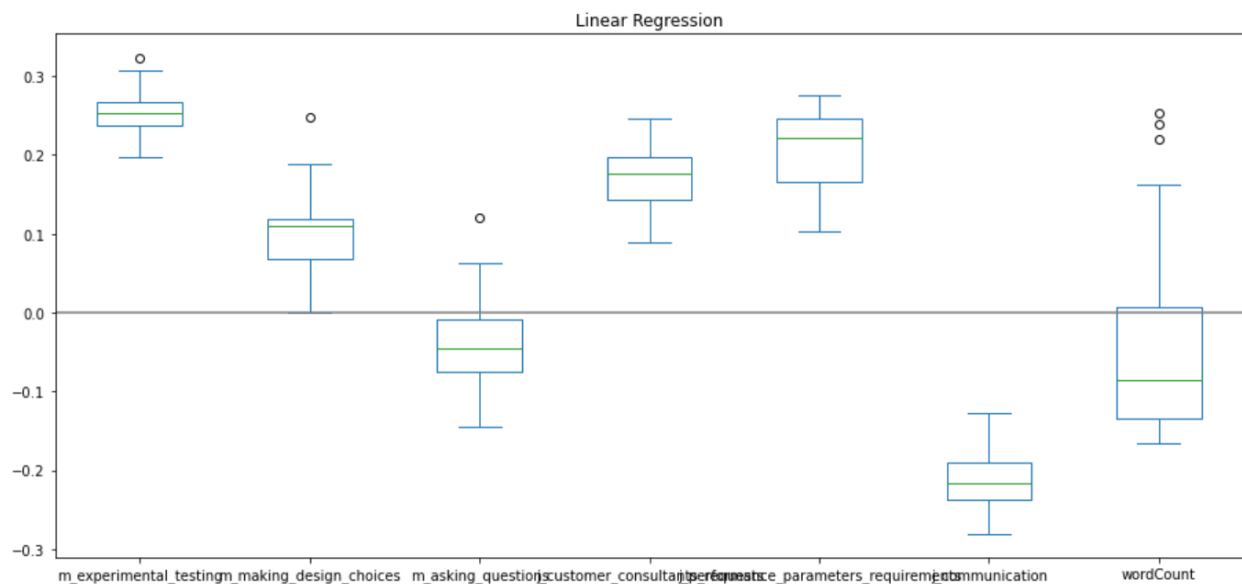
# Modelling

## Overview

The main goal of each of these models was to predict a final outcome score as accurately as possible. Two data sets were used in these predictions, the dataset grouped by individual users, in which we tried to predict individuals outcome scores, and the dataset grouped by teams in which we were attempting to predict the average score per user in each of these groups.

## Linear Regression

The first model utilised was a simple Linear Regression model as it is easy to implement and gave an initial baseline for performance. This however returned very low training scores of 0.108 and 0.230 and testing scores of -0.164 and -0.102 for the grouped by user and by team data respectively. This negative testing score means that our data doesn't follow this model at all. A potential cause of issue in Linear Regression models is often the high amount of variance within the data which causes the algorithm to model this erroneous data and not the intended outcome. This can be seen graphed below for the user model.

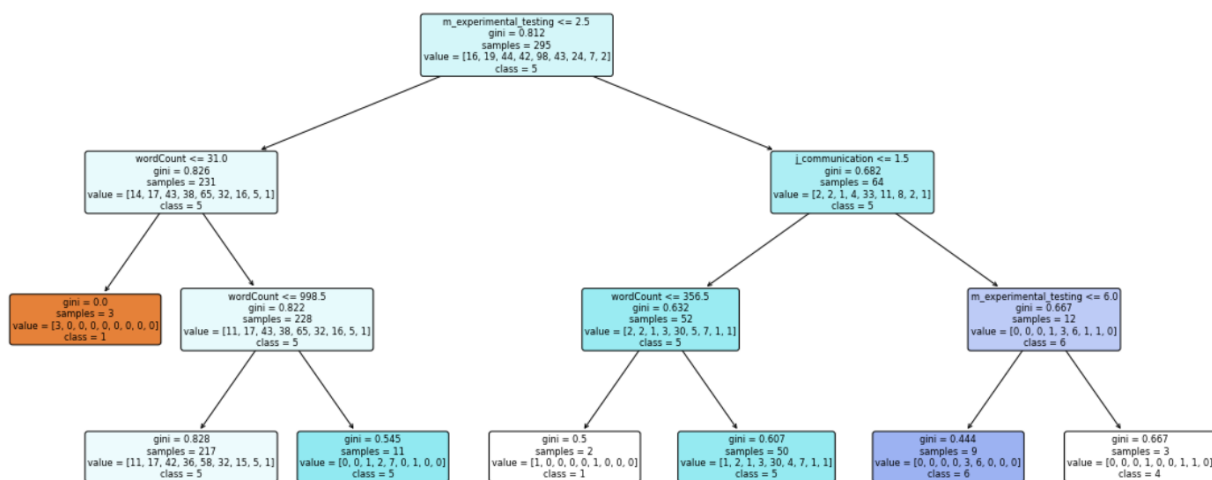


This variance can also lead to overfitting, meaning that the model is far too specific to the training data provided. We look to minimise the amount of variance in this model without increasing too much the bias in what is known as the Bias-Variance Tradeoff. Regularisation is designed to help with this by simplifying the data and there are two types which were utilised, Lasso and Ridge Regularisation. However, even with these

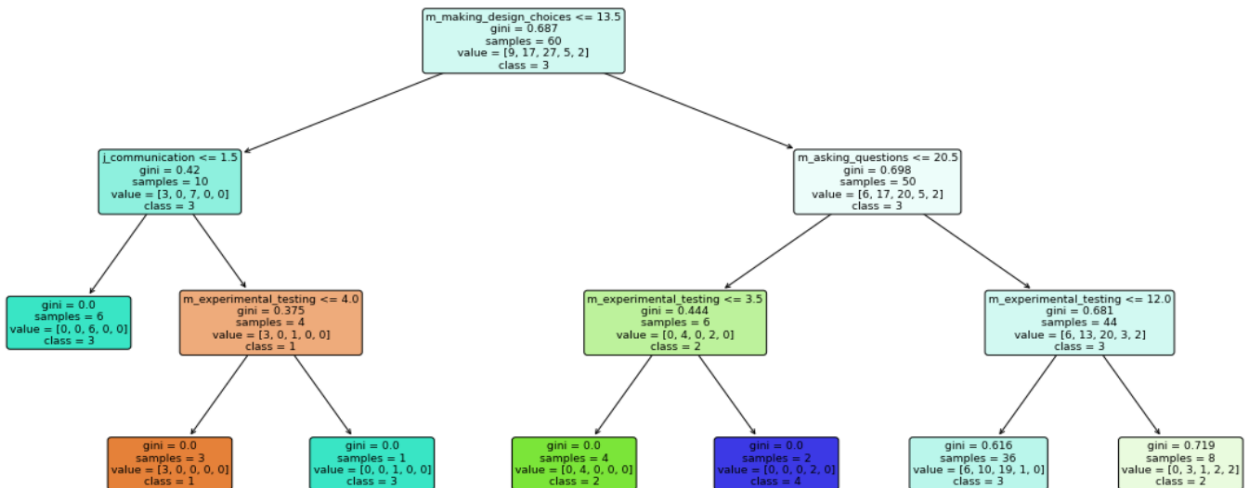
techniques implemented there was only a very slight improvement of testing score, but it still remained negative and therefore an unfit model for the data.

## Decision Trees

The next model implemented was a decision trees classifier model. This was deemed to be an appropriate model to use as each of the user scores was a discrete target value we were attempting to find. However, this did require the use of rounded versions of the average user score per group in order to make them discrete and provide easier comparison between the two models. Shown here is the user model with a depth of 3.



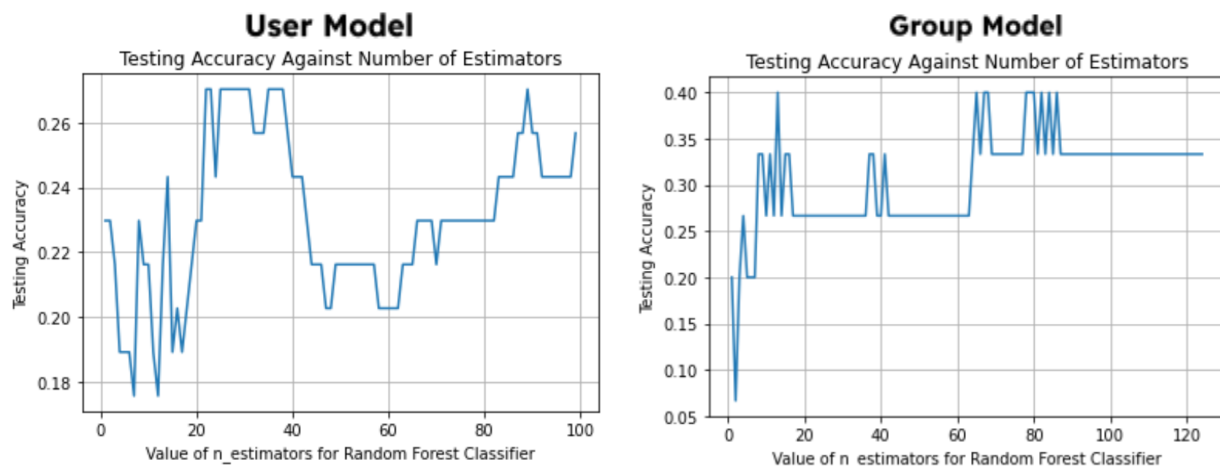
From this it can be seen the different requirements the model takes to classify input factors into different final outcome scores. For both the user and group based models, through manual testing and the gridsearchcv function, an optimal depth of 8 was found for both models. This provided enough depth to accurately classify scores while also being not too deep within the tree as after a certain point the model becomes highly specified to the provided data resulting in a very large bias. Ultimately, the user model obtained an accuracy of 0.311 and the group model, pictured below, had a final accuracy of 0.333.



Both of these models are improvements upon the initial linear regression, however still leave much to be desired.

## Random Forests

Following on from the decision tree modelling, we moved on to a Random Forest Classifier model. This was selected as it is an aggregation of multiple Decision Trees meaning that ideally, using these multiple models, it will provide us with a more accurate final model. Once again the rounded average scores were used for the group model as discrete values are required for classification as with Decision Trees. In order to find the best accuracy for the models we needed to calculate the optimal amount of estimators for the Random Forest Classifier to use. More estimators tends to result in a higher accuracy however it also makes it into a slower model. Below can be seen the graphs of number of estimators against model accuracy for both the user and group models.



Through the use of these visualisations in conjunction with `gridsearchcv` the optimal amount of estimators was found to be 30 and 100, for the user and group models respectively, which return accuracies of 0.270 and 0.333. These amounts of estimators were chosen as they were when the accuracy was highest but also had stopped varying as much, as seen in the plots.

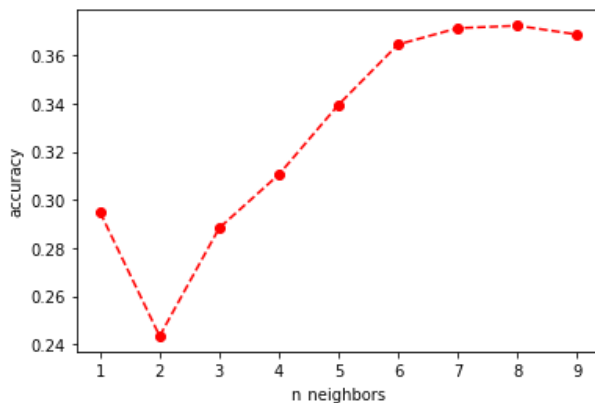
However, the accuracy scores obtained from this model showed no improvement whatsoever over the original Decision Tree model results with the user model actually decreasing in accuracy by 0.041 while the group model remained the same. While this result is unexpected, it is likely due to the fact that Random Forest models are a collected average of multiple Decision Trees. In the Decision Tree models we may have just been lucky to find an accurate random state, while on average the accuracy score of other Decision Tree models is much lower.

Another benefit of Random forests is that it can tell you which features are the most important for the model. For the user model 'word count', 'asking questions' and 'making

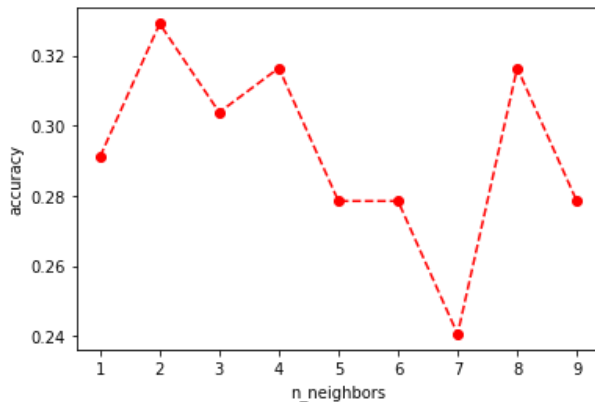
design choices' were deemed the most important whereas for the group model 'experimental testing' replaced 'word count.'

## K-Nearest Neighbour

- Group Model:



- User Model:



We implemented the K-Nearest Neighbour classification model. The K-Nearest Neighbour (KNN) algorithm is a data classification method used to estimate the likelihood that a data point is a member of one group or another based on which group the nearest data point belongs to (Joby, 2021).

To optimise our model, we varied the number neighbours parameter from 1 to 9 and evaluated the performance of each setting on both models. For each 'n\_neighbors' value, we trained the KNN classifier and then used it to predict the outcome of the test dataset.

Afterwards, we calculated the accuracy of each prediction, which measures the proportion of correct predictions out of the total number of instances. The results are displayed in a line graph, allowing us to see how the accuracy of the model varies as the number of neighbours grows.

The highest accuracy score we achieved was approximately 0.38 in the group model. For the user model, we achieved 0.28. While this score reflects the model's ability to correctly predict

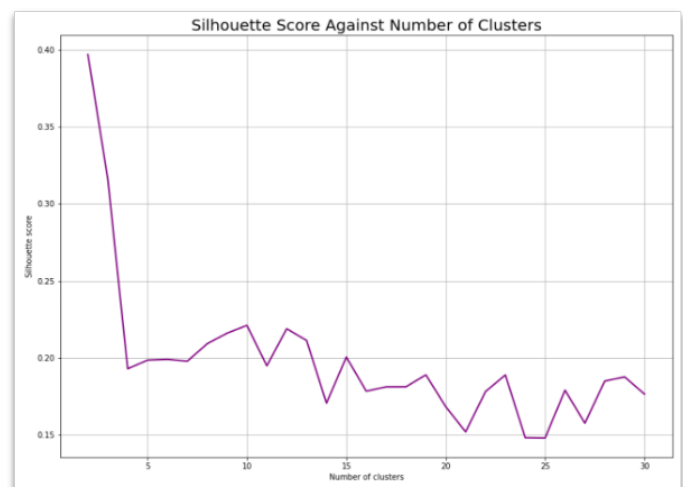
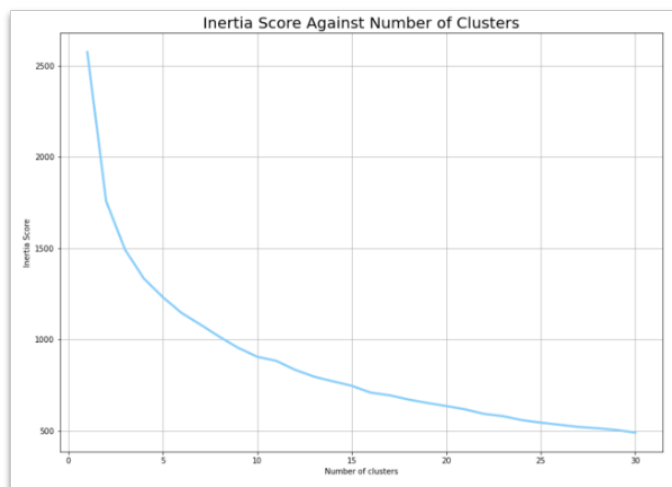
outcomes based on the given training data, its relatively low accuracy score suggests that this may be due to the fact that the words or phrases extracted from the chat logs do not adequately reflect the knowledge or skills that translate into reported performance, as well as if the chat data contains a lot of irrelevant or misleading information. KNN is a simple, distance\_based classification method (Aishearya Singh, 2018). If the relationship between chat content and reported scores is complex or nonlinear, KNN may not be the best model for this particular dataset.

# KMeans Clustering

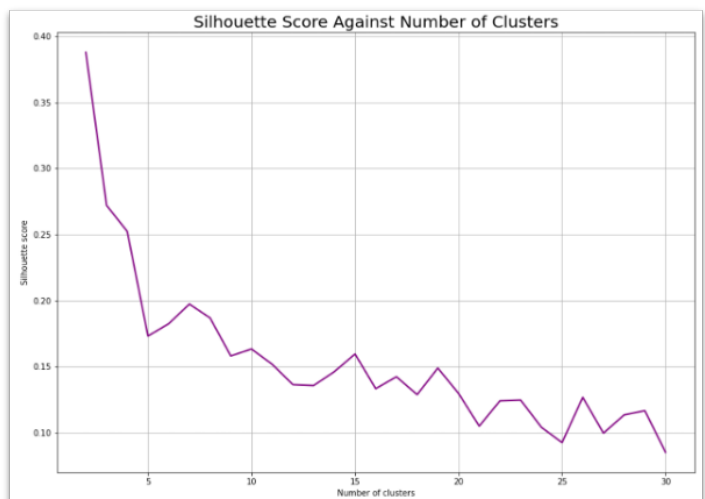
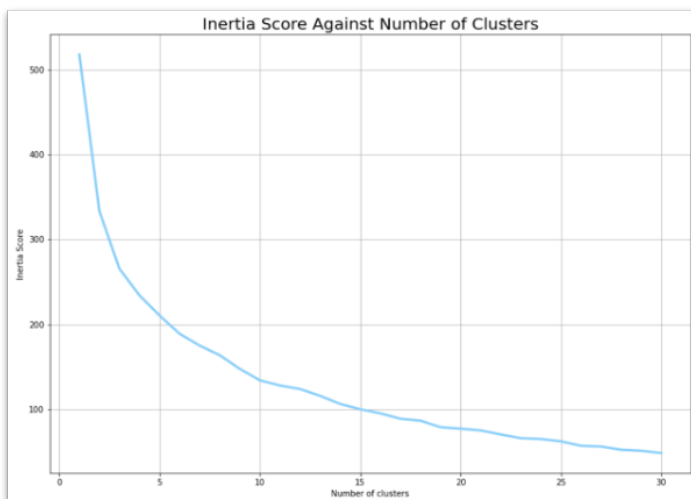
K-means clustering is utilised to identify potential structures or groups in our dataset that may not be initially apparent (LEDU, 2018). Clustering could reveal patterns related to how different team behaviours or chat patterns correlate with their final Outcomescores. We initially set the algorithm to  $k=3$ , which implies the data will be separated into three clusters. Each cluster will contain data points that are similar to one another in terms of the clustering features. The fitting approach is used to compute k-means clusters, while the prediction method is used to compute cluster centres and predict the clustering index for each data point.

The `kmeans.inertia_` parameter is used to assess the quality of our clusters. This value represents the total of the samples' squared distances to their nearest cluster centres. A lower inertia score indicates better clustering in this scenario since it indicates that the data points are closer to their respective centroids.

- User Model:



- Group Model:

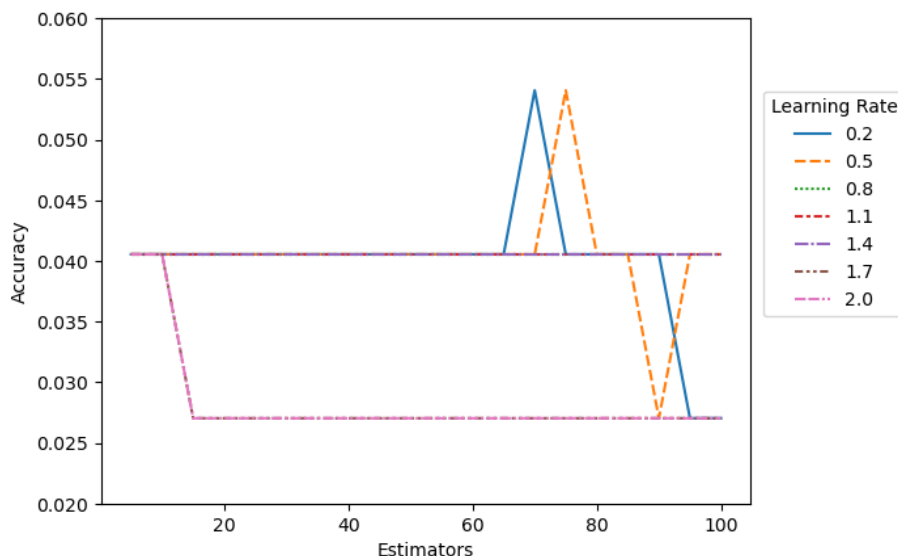




From those images, we can see that it loops over a range of possible cluster numbers (from 1 to 30), fitting a k-means model for each number of clusters and storing the resulting inertia score. In the User model, we set that `n_clusters=10` when initialising the KMeans model to have an optimal clustering solution. 905 is the inertia score, indicating the total within-cluster sum of squares. In the Group Model, we set `n_clusters=7` and got 175 inertia scores.

## Boosting

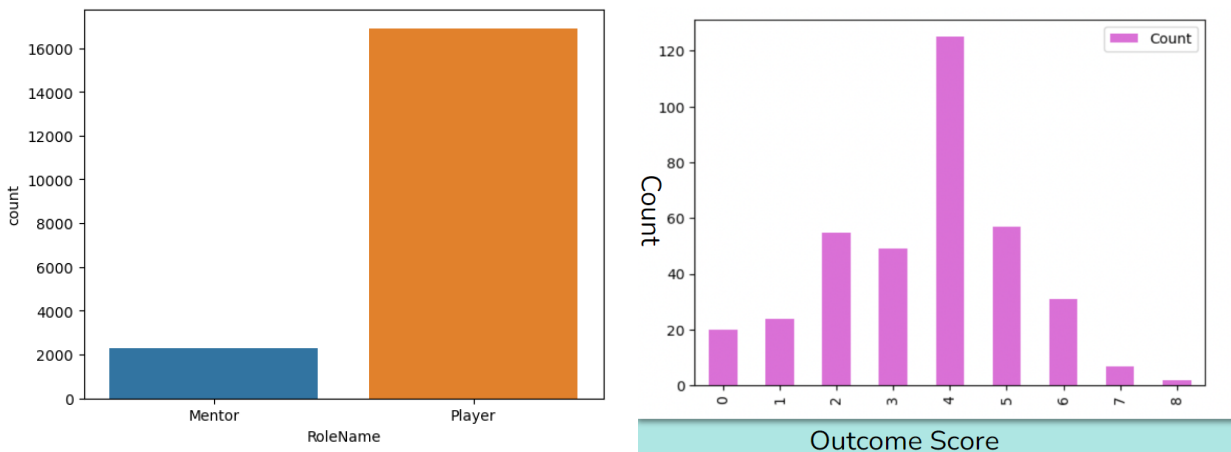
Boosting, similar to kMeans clustering is an iterative algorithm, this means an estimate is created, then used to generate an updated and theoretically improved estimate, this algorithm then continues until a stopping criterion is satisfied, or it reaches a maximum number of iterations. The boosting implemented with our data uses simple and fast regression and classifications models at each iteration to create the model. In particular AdaBoost (Adaptive Boosting) was used which is an iterative process that associates weights at each observation, the weights are then shifted after each iteration, so more emphasis is placed on classifications that are incorrectly predicted, these weighted targets then create the updated model. Learning rates the rate at which the weights are updated each iteration and consequently choosing an optimal learning rate is vital when optimising AdaBoost, a learning rate too high or too low decreases the accuracy of the model.



The graph above represents the optimal number of estimators and learning rate to achieve the highest accuracy, for the data provided the best learning rate is 0.5 and the corresponding optimal number of estimators is 75.

## SMOTE

Imbalanced data occurs when observed frequencies are clearly different across different values of a variable, in other words there are many of one variable type and few of another, two examples of this can be seen below with majority of the users being players over mentors and a third of the users having an outcome score of 4, and an uneven spread of scores from 1-8.



It is important to combat this imbalance as a model that predicts an outcome score of 4 may have a high accuracy, however this is irrelevant to our analysis due to the number of outcome scores of 4. One method to counter data imbalance is under sampling which is when a number of data points that are present too often are removed from the data, the disadvantage to this is that valuable data is lost. Oversampling on the other hand consists of making duplicates of the data values that are least present in the data, this however leads to the creation of many duplicate data points. Data augmentation is a method that works like oversampling, however rather than creating duplicate data points, minor deviations are made to the copied data points, the Synthetic Minority Oversampling Technique (SMOTE) uses data augmentation to derive an advanced version of oversampling, thus avoiding duplicates of data. The aim of the SMOTE algorithm is to increase false positives and decrease false negatives, however it does this at the cost of accuracy, as the accuracies of our models are already so low this is not a desirable approach when predicting outcome scores.

# NLP

## Overview

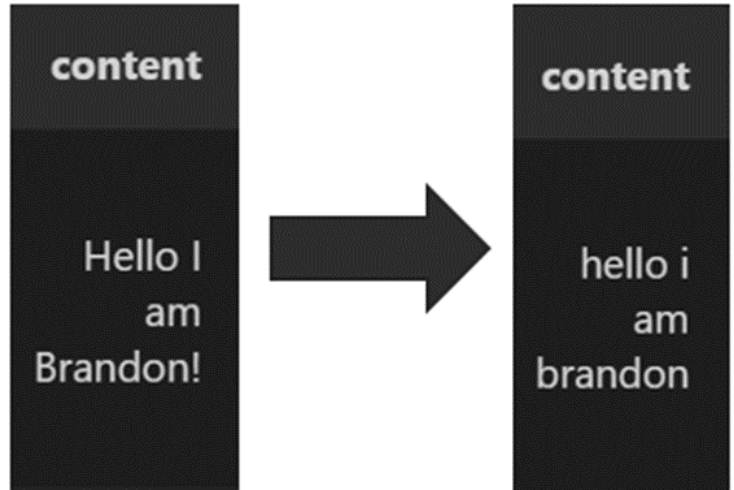
Natural Language Processing (NLP) is a field of programming concerned with extracting usable numeric data or other important information from human language or text.

When it comes to the Nephrotex Virtual Internships dataset, the actual chat or 'content' column is the primary source of information which determined other features such as 'm\_making\_design\_choices' or 'j\_communication', which are Boolean columns which denote the purpose or function of a chat message. As such, to introduce any new features or information to the dataset, they would have to be extracted from the actual chat transcripts using NLP techniques.

For this dataset, TF-IDF will be used to find the most important words in all the chat data and use them to establish new features which may have some stronger and more relevant correlations to the outcome score.

## Text-Data Pre-Processing

Before we apply any NLP techniques to the data, the text data needs to undergo some transformation to make the TF-IDF more effective. Most NLP techniques, including the one intended to be used on this dataset, involve tokenizing the text by words or splitting bodies of text by individual words. However, features of language things such as capital letters or standard features of punctuation, mean that the tokens 'Mean', 'mean.' and 'mean' are recognized as different words. To avoid any inaccuracies, the chat data was stripped of all punctuation and set to lower case only.



## TF-IDF

Term frequency-inverse document frequency (TF-IDF) is an NLP technique used to find the importance of each word in a document or series of documents, in relation to those documents. This is achieved by finding the number of times a certain term appears in

one document (the term frequency) and multiplying it by the proportion of documents in the set which contain that (inverse document frequency), this process usually returns a score between 0 and 1, where the higher the score the more relevant, or important, the word can be considered. Whilst the TF helps to recognize words with high usage, the IDF helps balance this frequency count by returning lower scores for those words which appear in more documents, to avoid confusion caused by standard lexical features of the English language, e.g., 'the', 'is', 'I', etc.

$$w_{x,y} = tf_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

## TF-IDF

Term  $x$  within document  $y$

$tf_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

$N$  = total number of documents

*The formula used to calculate TF-IDF*

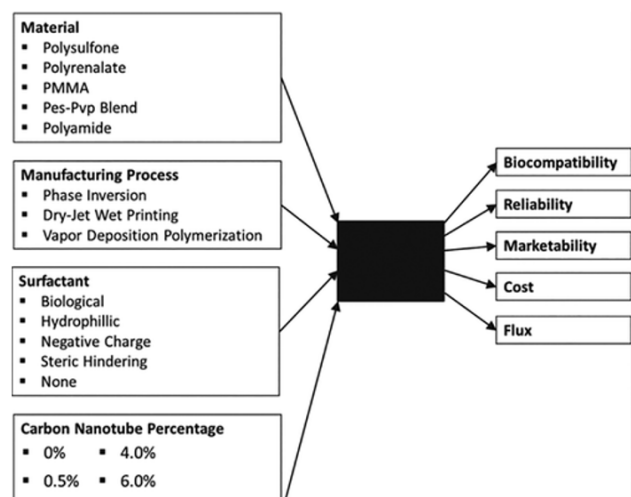
Implementing TF-IDF on the Virtual Internships data involved utilising the scikit learn package 'TfidfVectorizer'. When initialising the vectorizer, stop-words were set to 'English' so that basic structural words in the English language would be ignored by the vectorizer. The vectorizer was also initialised to return a list of 25 words with highest TF-IDF scores, using the 'max\_features' parameter, this was aimed to make initial analysis a bit easier and prevent over-analysis of the TF-IDF results in-order to avoid confusion.

```
Feature Names n ['agree' 'best' 'blood' 'cnt' 'cost' 'did' 'flux' 'good' 'hindering' 'im'
'just' 'like' 'make' 'marketability' 'negative' 'notebook' 'prototype'
'prototypes' 'reactivity' 'reliability' 'steric' 'surfactant' 'think'
'yeah' 'yes']
```

*List of the 25 most important words according to the TfidfVectorizer.*

## New Features

Though there are now 25 possible words from which new features could be designed it is important to note that the issue of relevancy had to be considered. Making features without any reason would cause inaccuracies in any later modelling and analyses.



*A diagram describing the inputs and outputs considered by participants in the Virtual Internship program.*

Upon some further research about Nephrotex, it was learned that some components of the virtual internship task could be categorised as either inputs or outputs. Using the diagram seen to the side and the list extracted from the TfidfVectorizer it was possible to construct two new, contextually relevant, features for the dataset. Both of which counted the following words in each chat message, producing a sum of the number of input and output terms contained in one chat entry. The hypothesis being the more a participant used input and output terms the more likely they were having relevant conversations which would lead to higher outcome scores.

#### Input:

- Surfactant
- Steric
- Hindering
- CNT (Carbon Nanotubes)

#### Output:

- Reliability
- Cost
- Flux
- Marketability

## Modelling with New Features

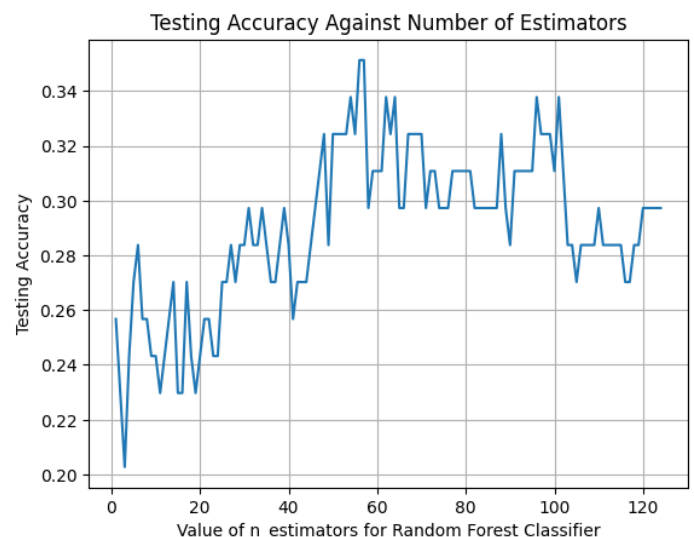
```
m_experimental_testing 0.078
m_making_design_choices 0.117
m_asking_questions 0.138
j_customer_consultants_requests 0.062
j_performance_parameters_requirements 0.093
j_communication 0.046
wordCount 0.195
input 0.136
output 0.134
```

To begin exploring the effects of the new features, a user-based model was constructed using Random Forest Classifier. This model produced a 0.297 score for accuracy at an optimal number of estimators of 69, which was an improvement from the previous user model's 0.27, however it was not substantial enough to be deemed relevant.

Though there was no evident improvement in accuracy due to addition of new features, when taking a look at the feature importances for the model, both 'input' and 'output' are the most important features after word count.

This was an indication that creating the new features were a step in the right direction when it came to constructing better models.

Taking the feature importances into consideration a new model using RFC was made using just 'output', 'input', 'wordCount', 'm\_asking\_questions' and 'm\_making\_design\_choices', the features with importance greater than 0.1.



This model was capable of achieving an accuracy of 0.351 when the number of estimators was 55, however as shown in the graph the accuracy is very volatile and lacks stability around any number of estimators.

## Results

Initially, a Linear Regression model returned a training and testing score of 0.108 and 0.164 for the user model and 0.230 and -0.102 for the group model. An attempt was made to reduce the variance within the data using Lasso and Ridge regularisation however this only slightly improved the scores with testing score still being negative indicating an unfit model. Following this, a Decision Tree Classifier was utilised which returned accuracies of 0.311 and 0.333 at a depth of 8 for the user and group models respectively. To improve upon this a Random Forests Classifier was employed however this returned a decreased accuracy score of 0.270 with 30 estimators for the user model, and an accuracy of 0.333 with 100 estimators for the group model. The Random Forests user model was further improved when using the new features, 'input' and 'output', whilst also removing some less important features to generate an accuracy of 0.351 with 55 estimators.

Both a K Nearest Neighbours (KNN) and kMeans clustering model was developed to analyse the data, the accuracy for the KNN model was 0.28 and 0.38 for the user and group model respectively, this group model was the highest accuracy achieved thus far. The kMeans clustering model doesn't use training and testing sets and instead only uses the data features to fit the model, it also doesn't calculate an accuracy score and instead determines the optimal number of clusters and the inertia score. The optimal number of clusters for the user and group models was 10 and 7 respectively and the user model obtained an inertia score of 905 while the group model obtained an inertia score of 175, the inertia score represents how well a dataset was clustered, and a good model has a low inertia score and a low number of clusters and thus the group model was the better model in both situations.

# Conclusion

Various models were applied to the virtual internships data in both a user and group sense, the highest accuracy score achieved was 0.38 in the group model in K Nearest Neighbours (KNN) with the next closest being a Random Forests Classifier produced using Natural Language Processing which returned an accuracy of 0.351. The approaches used to predict outcome score using chat data between groups faced limitations as the different topics discussed and the frequency of communication was shown to have little to no effect on the outcome score, groups with below average scores were still talking just as much and about the relevant things. A third of the users in the internship also achieved an outcome score of 4 and the outcome score is on a user level based on the report the individual submits, this made it difficult to predict a group outcome score as groups had to be given a score derived from the members averages. Natural Language Processing was utilised to provide an in depth analysis of the chat data between groups, specifically TF-IDF (Term frequency - Inverse document frequency) was used to determine the importance of certain words, however this method faced various limitations such as its lack of account for sentiment/meaning and instead is just frequency based, the words also had to have meaning assigned by us and we could be wrong. The next step in natural language processing would be conducting sentiment analysis to provide meaning for words and demonstrate a direct relation in the words to performance outcome. Mentors were also removed from the modelling as they were all given an outcome score of 4 and it would bias the results, in future implementations of the Nephrotex internship it would be beneficial to have the same set of mentors in each internship as it would provide opportunity to analyse the effect of mentors on the outcome score.

# Appendix

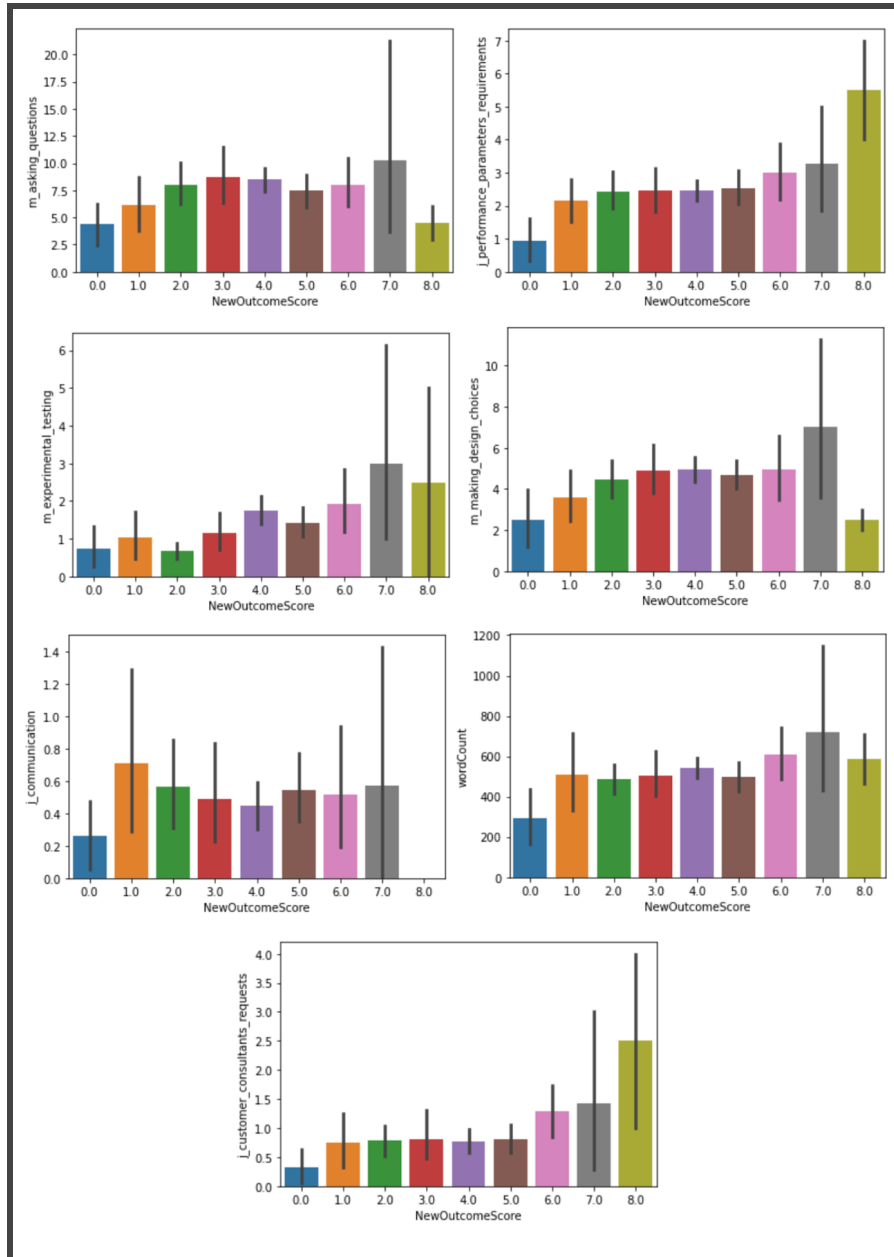
## References

- Nephrotex - International Society for Quantitative Ethnography.  
<https://www.virtualinterns.org/nephrotex/>
- Aishwarya Sigh. (2018, August 22). *KNN Algorithm: Introduction to K-Nearest Neighbours Algorithm in Regression*. Medium.  
<https://www.analyticsvidhya.com/blog/2018/08/k-nearest-neighbor-introduction-regression-python/>
- Swalin Alvira. (2018, Jan 31). *How to Handle Missing Data*. Towards Data Science. Medium.  
<https://towardsdatascience.com/how-to-handle-missing-data-8646b18db0d4>
- Education Ecosystem(LEDU). (2018, Sep 13). *Understanding K-means Clustering in Machine Learning*. Medium.  
<https://towardsdatascience.com/understanding-k-means-clustering-in-machine-learning-6a6e67336aa1>
- Shaffer, David & Arastoopour, Golnaz & Swiecki, Zachari & Ruis, Andrew & Chesler, Naomi. (2015). Teaching and Assessing Engineering Design Thinking with Virtual Internships and Epistemic Network Analysis.
- Golnaz Arastoopour, Chesler, N., D'Angelo, C. M., & Lepak, C. (2012). Nephrotex: Measuring first year students' ways of professional thinking in a virtual internship. ResearchGate; unknown.  
[https://www.researchgate.net/publication/283865973\\_Nephrotex\\_Measuring\\_first\\_year\\_students'\\_ways\\_of\\_professional\\_thinking\\_in\\_a\\_virtual\\_internship](https://www.researchgate.net/publication/283865973_Nephrotex_Measuring_first_year_students'_ways_of_professional_thinking_in_a_virtual_internship)
- Korstanje, J. (2021, August 30). SMOTE. Medium.  
<https://towardsdatascience.com/smote-fdce2f605729>

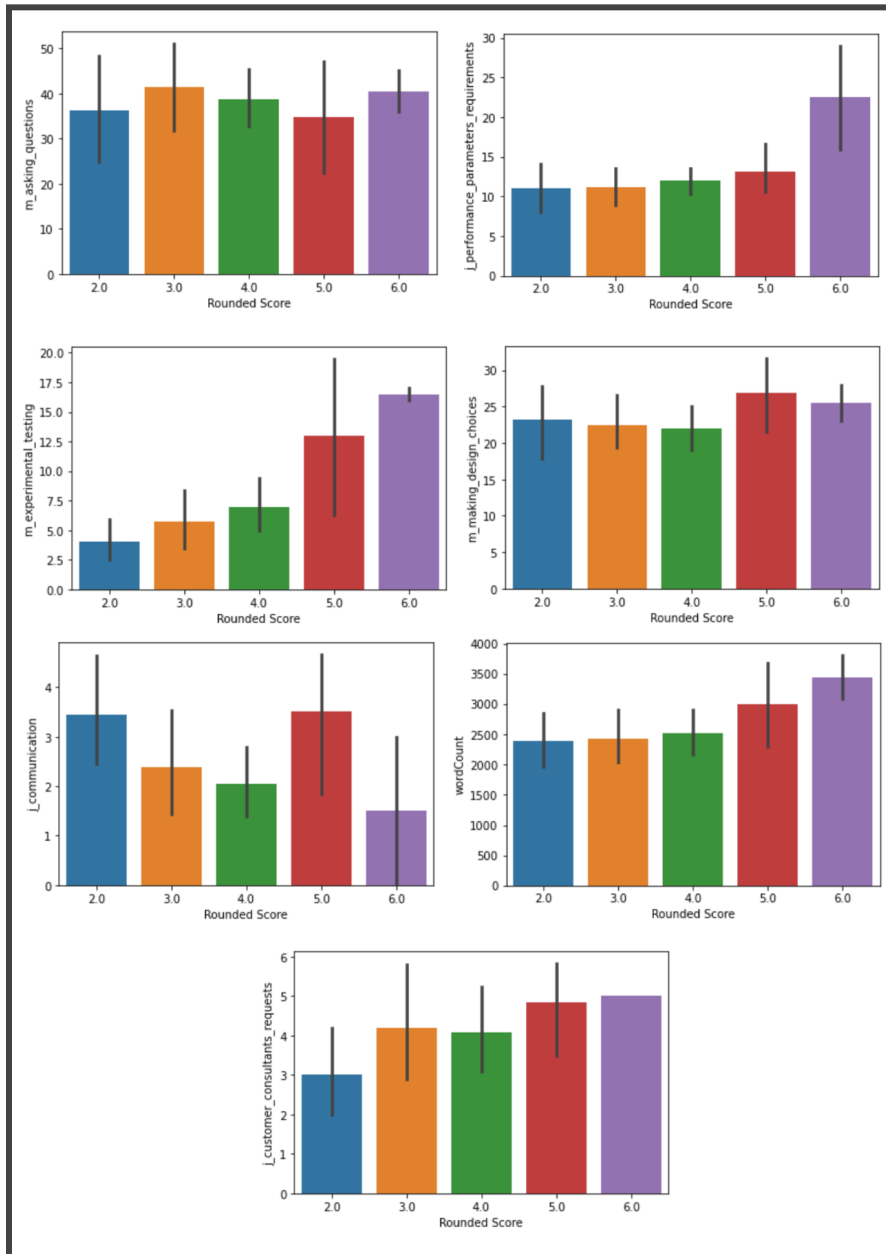


## Figures

Fig 1.



**Fig 2.**



# Virtual Internships

May 30, 2023

## 1 Virtual Internships

```
[221]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from matplotlib import ticker
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LinearRegression, LogisticRegression
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier, plot_tree
from sklearn.ensemble import RandomForestRegressor, BaggingClassifier, \
    →RandomForestClassifier
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.decomposition import PCA
from sklearn.svm import SVR
from sklearn.pipeline import Pipeline
from sklearn.metrics import confusion_matrix, accuracy_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn import metrics
from logitplots import plt_correlation_matrix

from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.ensemble import BaggingClassifier, RandomForestClassifier, \
    →VotingClassifier
from sklearn.model_selection import train_test_split, cross_validate, \
    →RepeatedKFold, GridSearchCV
from sklearn.metrics import accuracy_score, r2_score, silhouette_score
from sklearn.svm import SVC
from sklearn.linear_model import LogisticRegression, LinearRegression, Lasso, \
    →Ridge, RidgeCV, LassoCV
from sklearn.cluster import KMeans
from sklearn.ensemble import AdaBoostClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import GradientBoostingClassifier
```

```
# from logitplots import plt_correlation_matrix
import string
```

## 1.1 Data Cleaning

```
[222]: column_names = [
    "userIDs",
    "implementation",
    "Line_ID",
    "ChatGroup",
    "content",
    "group_id",
    "RoleName",
    "roomName",
    "m_experimental_testing",
    "m_making_design_choices",
    "m_asking_questions",
    "j_customer_consultants_requests",
    "j_performance_parameters_requirements",
    "j_communication",
    "OutcomeScore",
    "wordCount",
]

file_path = 'virtualInternshipData_ADS2001.csv'

try:
    with open(file_path, 'r', encoding='utf-8', errors='replace') as f:
        df = pd.read_csv(f)
        print(df.head())
except FileNotFoundError:
    print(f"File not found: {file_path}")
```

	Unnamed: 0	userIDs	implementation	Line_ID	ChatGroup	\
0	1	1	a	1	PRNLT	
1	2	1	a	2	PRNLT	
2	3	1	a	3	PRNLT	
3	4	1	a	4	PRNLT	
4	5	1	a	5	PRNLT	

	content	group_id	RoleName	\
0	Hello team. Welcome to Nephrotex!	2	Mentor	
1	I'm Maria Williams. I'll be your design adviso...	2	Mentor	
2	I'm here to help if you have any questions.	2	Mentor	
3	Please introduce yourselves with the name you ...	2	Mentor	
4	I just want to make sure everyone has found th...	2	Mentor	

		roomName	m_experimental_testing	\
0	Introduction and Workflow Tutorial with Entran...		0	
1	Introduction and Workflow Tutorial with Entran...		0	
2	Introduction and Workflow Tutorial with Entran...		0	
3	Introduction and Workflow Tutorial with Entran...		0	
4	Introduction and Workflow Tutorial with Entran...		0	

	m_making_design_choices	m_asking_questions	\
0	0	0	
1	0	0	
2	0	0	
3	0	0	
4	0	0	

	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0	0	
1	0	0	
2	0	0	
3	1	0	
4	0	0	

	j_communication	OutcomeScore	wordCount
0	0	4	5
1	0	4	11
2	0	4	9
3	0	4	51
4	0	4	39

```
[223]: df.columns
```

```
[223]: Index(['Unnamed: 0', 'userIDs', 'implementation', 'Line_ID', 'ChatGroup',
        'content', 'group_id', 'RoleName', 'roomName', 'm_experimental_testing',
        'm_making_design_choices', 'm_asking_questions',
        'j_customer_consultants_requests',
        'j_performance_parameters_requirements', 'j_communication',
        'OutcomeScore', 'wordCount'],
        dtype='object')
```

```
[224]: df.head()
```

```
[224]:   Unnamed: 0  userIDs implementation  Line_ID ChatGroup \
0          1        1             a         1    PRNLT
1          2        1             a         2    PRNLT
2          3        1             a         3    PRNLT
3          4        1             a         4    PRNLT
4          5        1             a         5    PRNLT
```

```
        content  group_id RoleName \
```

0	Hello team. Welcome to Nephrotex!	2	Mentor
1	I'm Maria Williams. I'll be your design adviso...	2	Mentor
2	I'm here to help if you have any questions.	2	Mentor
3	Please introduce yourselves with the name you ...	2	Mentor
4	I just want to make sure everyone has found th...	2	Mentor

	roomName	m_experimental_testing	\
0	Introduction and Workflow Tutorial with Entran...		0
1	Introduction and Workflow Tutorial with Entran...		0
2	Introduction and Workflow Tutorial with Entran...		0
3	Introduction and Workflow Tutorial with Entran...		0
4	Introduction and Workflow Tutorial with Entran...		0

	m_making_design_choices	m_asking_questions	\
0	0	0	
1	0	0	
2	0	0	
3	0	0	
4	0	0	

	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0	0	
1	0	0	
2	0	0	
3	1	0	
4	0	0	

	j_communication	OutcomeScore	wordCount
0	0	4	5
1	0	4	11
2	0	4	9
3	0	4	51
4	0	4	39

```
[225]: df.shape
```

```
[225]: (19180, 17)
```

```
[226]: df.dtypes
```

```
[226]: Unnamed: 0          int64
      userIDs          int64
      implementation    object
      Line_ID          int64
      ChatGroup         object
      content           object
      group_id         int64
```

RoleName	object
roomName	object
m_experimental_testing	int64
m_making_design_choices	int64
m_masking_questions	int64
j_customer_consultants_requests	int64
j_performance_parameters_requirements	int64
j_communication	int64
OutcomeScore	int64
wordCount	int64
dtype:	object

```
[227]: # make 'none' change to null
df.replace('None',np.nan,inplace=True)
# checking for the missing value
df.isnull().sum()
```

```
[227]: Unnamed: 0      0
      userIDs      0
      implementation  0
      Line_ID      0
      ChatGroup     0
      content       0
      group_id      0
      RoleName      3
      roomName      0
      m_experimental_testing  0
      m_making_design_choices  0
      m_masking_questions    0
      j_customer_consultants_requests  0
      j_performance_parameters_requirements  0
      j_communication        0
      OutcomeScore          0
      wordCount             0
      dtype: int64
```

```
[228]: # find the RoleName got the missing value the delete that column
df.dropna(axis=0,subset=['RoleName'],
          how='any',inplace=True)
# reset
df.reset_index(drop=True,inplace=True)
# check for the missing value again
df.isnull().sum()
```

```
[228]: Unnamed: 0      0
      userIDs      0
      implementation  0
```



```

Line_ID          0
ChatGroup        0
content          0
group_id        0
RoleName         0
roomName         0
m_experimental_testing 0
m_making_design_choices 0
m_asking_questions 0
j_customer_consultants_requests 0
j_performance_parameters_requirements 0
j_communication 0
OutcomeScore     0
wordCount        0
dtype: int64

```

```
[229]: # check if there are any True values for the test isnull on the dataframe
df.isnull().values.any()
```

```
[229]: False
```

```
[230]: df.shape
```

```
[230]: (19177, 17)
```

```
[231]: filtered_df = df[df['RoleName'].isnull()]

# Select the 'userIDs', 'implementation', and 'Line_ID' columns
output = filtered_df[['userIDs', 'implementation', 'Line_ID']]

print(output)
```

```

Empty DataFrame
Columns: [userIDs, implementation, Line_ID]
Index: []

```

```
[232]: data = df[['Unnamed: 0', 'content', 'RoleName']].drop_duplicates()
content = df['content']
data.shape
data
```

```
[232]:
```

	Unnamed: 0	content	RoleName
0	1	Hello team. Welcome to Nephrotex!	Mentor
1	2	I'm Maria Williams. I'll be your design adviso...	Mentor
2	3	I'm here to help if you have any questions.	Mentor
3	4	Please introduce yourselves with the name you ...	Mentor
4	5	I just want to make sure everyone has found th...	Mentor
...	...	...	...

19172	19176	yes	Player
19173	19177	sounds good	Player
19174	19178	Well, we are out of time for our meeting.	Mentor
19175	19179	Precisely	Player
19176	19180	Good discussion today! Don't forget to complet...	Mentor

[19177 rows x 3 columns]

```
[233]: #drop duplicate rows
data = df[['Unnamed: 0', 'content', 'RoleName']].drop_duplicates()
#drop duplicate columns
df = df.drop(columns = ["content", "RoleName"])
#merge dataests
clean_df = pd.merge(data, df, on='Unnamed: 0')
```

```
[234]: clean_df
```

```
[234]:      Unnamed: 0      content RoleName \
0           1      Hello team. Welcome to Nephrotex!  Mentor
1           2  I'm Maria Williams. I'll be your design adviso...  Mentor
2           3      I'm here to help if you have any questions.  Mentor
3           4  Please introduce yourselves with the name you ...  Mentor
4           5  I just want to make sure everyone has found th...  Mentor
...      ...      ...      ...
19172      19176      yes      Player
19173      19177      sounds good      Player
19174      19178      Well, we are out of time for our meeting.  Mentor
19175      19179      Precisely      Player
19176      19180  Good discussion today! Don't forget to complet...  Mentor
```

	userIDs	implementation	Line_ID	ChatGroup	group_id	\
0	1	a	1	PRNLT	2	
1	1	a	2	PRNLT	2	
2	1	a	3	PRNLT	2	
3	1	a	4	PRNLT	2	
4	1	a	5	PRNLT	2	
...	...	...	...	...	...	
19172	392	o	19179	PESPVP	6	
19173	388	o	19180	PESPVP	6	
19174	367	o	19181	PESPVP	6	
19175	393	o	19182	PESPVP	6	
19176	367	o	19183	PESPVP	6	

	roomName	\
0	Introduction and Workflow Tutorial with Entran...	
1	Introduction and Workflow Tutorial with Entran...	
2	Introduction and Workflow Tutorial with Entran...	

3	Introduction and Workflow Tutorial with Entran...
4	Introduction and Workflow Tutorial with Entran...
...	...
19172	Reflection team discussion of first batch results
19173	Reflection team discussion of first batch results
19174	Reflection team discussion of first batch results
19175	Reflection team discussion of first batch results
19176	Reflection team discussion of first batch results

	m_experimental_testing	m_making_design_choices	m_asking_questions	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	
...	...	...	...	
19172	0	0	0	
19173	0	0	0	
19174	0	0	0	
19175	0	0	0	
19176	0	0	0	

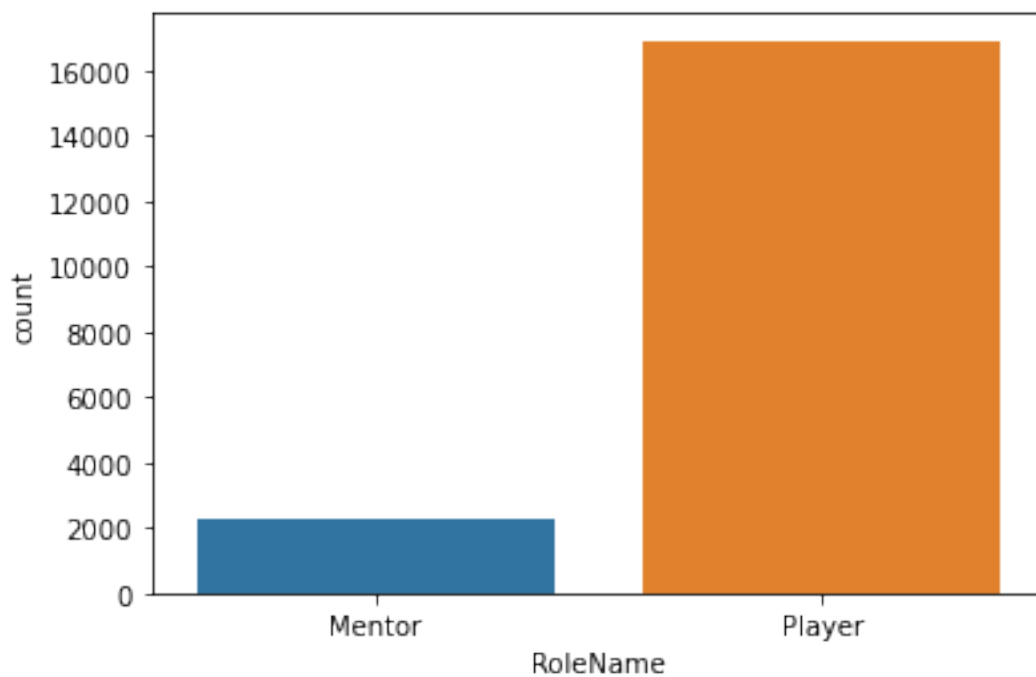
	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0	0	
1	0	0	
2	0	0	
3	1	0	
4	0	0	
...	...	...	
19172	0	0	
19173	0	0	
19174	0	0	
19175	0	0	
19176	0	0	

	j_communication	OutcomeScore	wordCount
0	0	4	5
1	0	4	11
2	0	4	9
3	0	4	51
4	0	4	39
...	...	...	...
19172	0	5	1
19173	0	8	2
19174	0	4	9
19175	0	4	1
19176	0	4	11

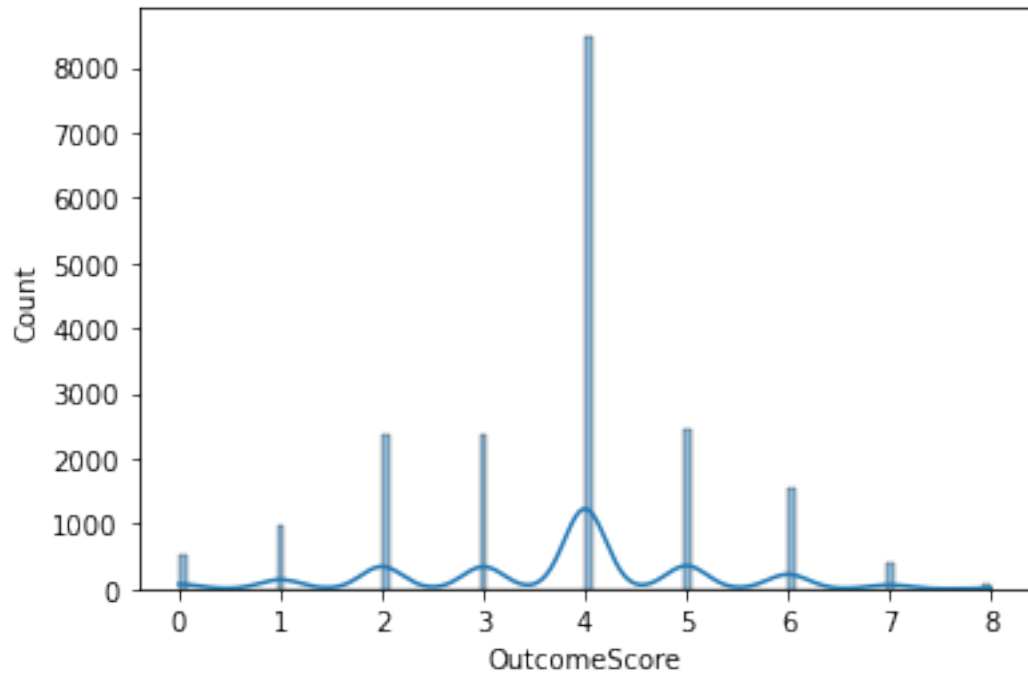
[19177 rows x 17 columns]

## 1.2 Visualization

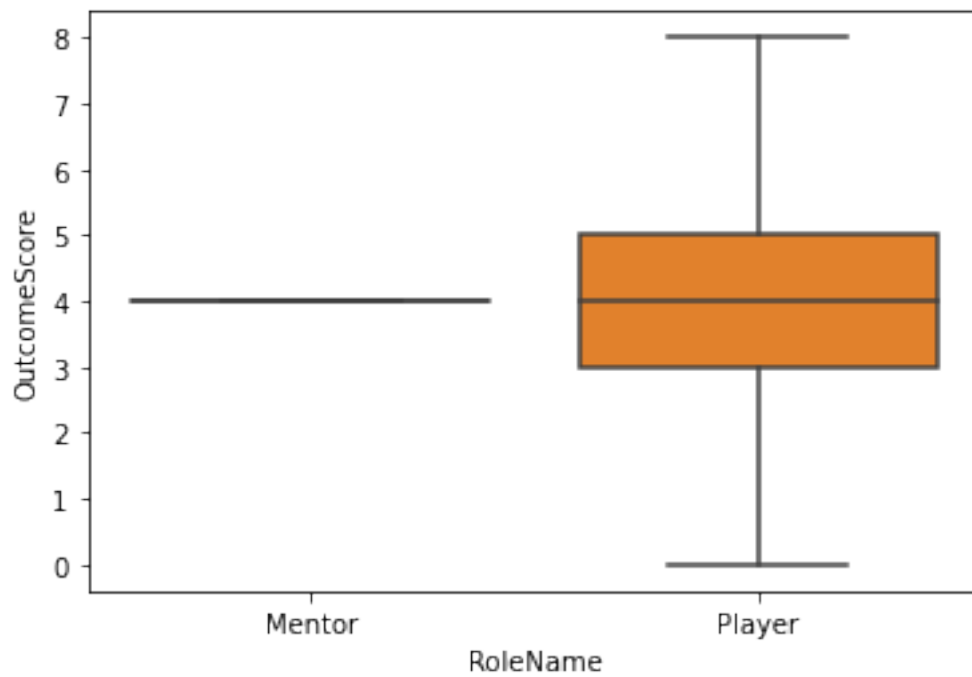
```
[235]: sns.countplot(x='RoleName', data=clean_df)  
plt.show()
```



```
[236]: # Histogram for 'OutcomeScore'  
sns.histplot(clean_df['OutcomeScore'], kde=True)  
plt.show()
```



```
[237]: sns.boxplot(x='RoleName', y='OutcomeScore', data=clean_df)  
plt.show()
```



```
[238]: df1=clean_df[['RoleName','content','OutcomeScore']]
for Mentor, Player in df1.groupby('RoleName'):
    print(Mentor)
    print(Player)
```

Mentor

	RoleName	content \
0	Mentor	Hello team. Welcome to Nephrotex!
1	Mentor	I'm Maria Williams. I'll be your design adviso...
2	Mentor	I'm here to help if you have any questions.
3	Mentor	Please introduce yourselves with the name you ...
4	Mentor	I just want to make sure everyone has found th...
...	...	...
19148	Mentor	Soon, you will be switching teams and sharing ...
19156	Mentor	You will bring knowledge about a unique materi...
19163	Mentor	Which attributes does your material maximize a...
19174	Mentor	Well, we are out of time for our meeting.
19176	Mentor	Good discussion today! Don't forget to complet...

	OutcomeScore
0	4
1	4
2	4
3	4
4	4
...	...
19148	4
19156	4
19163	4
19174	4
19176	4

[2275 rows x 3 columns]

Player

	RoleName	content	OutcomeScore
5	Player	Hello I am Brandon!	4
6	Player	I am Zelin	4
7	Player	Hi	4
8	Player	i am jack	4
9	Player	Hey! I'm Rachel!	2
...	...	...	...
19170	Player	exactly!	5
19171	Player	sounds good	7
19172	Player	yes	5
19173	Player	sounds good	8
19175	Player	Precisely	4

[16902 rows x 3 columns]

```
[239]: grouped = clean_df.groupby('RoleName')['wordCount'].agg(['mean', 'sum'])
print(grouped)
```

	mean	sum
RoleName		
Mentor	21.596484	49132
Player	11.264939	190400

```
[240]: grouped = clean_df.groupby(['RoleName', 'implementation'])['OutcomeScore'].
        →agg(['mean', 'sum'])
print(grouped)
```

		mean	sum
RoleName	implementation		
Mentor	a	4.000000	836
	b	4.000000	840
	c	4.000000	1000
	d	4.000000	668
	e	4.000000	688
	f	4.000000	656
	g	4.000000	280
	h	4.000000	564
	i	4.000000	524
	j	4.000000	472
	k	4.000000	248
	l	4.000000	700
	m	4.000000	284
	n	4.000000	792
Player	o	4.000000	548
	a	2.877704	3459
	b	3.583333	3870
	c	3.239863	4714
	d	3.315261	4953
	e	3.936170	4070
	f	4.109278	3986
	g	3.815860	2839
	h	4.207263	4750
	i	4.003709	5397
	j	3.861582	2734
	k	3.725888	2936
	l	3.098865	1912
	m	3.406037	3498
	n	4.093835	7504
	o	4.101833	6042

```
[241]: grouped = clean_df.groupby('RoleName')['OutcomeScore'].agg(['mean', 'sum'])
print(grouped)
```

	mean	sum
--	------	-----

RoleName		
Mentor	4.00000	9100
Player	3.70749	62664

## 1.2.1 Grouping by Users

```
[242]: #Dropping mentors
clean_df = clean_df[clean_df.RoleName != 'Mentor']
#Adding column for each group
clean_df['Group'] = clean_df['group_id'].astype(str) + clean_df['implementation']
```

C:\Users\owner\AppData\Local\Temp\ipykernel\_19800\2187175854.py:4:

SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
clean_df['Group'] = clean_df['group_id'].astype(str) +
clean_df['implementation']
```

```
[243]: #Grouping by users
user_group = clean_df.groupby(['userIDs'], as_index = False)
#Getting the sum of the contributions
user_sum = user_group.sum()
```

```
[244]: user_group.describe()
```

```
[244]: Unnamed: 0 \
      count      mean      std      min      25%      50%
0      65.0    189.015385  97.980785      6.0    104.00    207.0
1      22.0    127.090909  110.428132      7.0     41.50     71.0
2      31.0    208.193548   91.107416      9.0    146.00    220.0
3      45.0    208.555556   72.142396     10.0    167.00    213.0
4      37.0    182.162162  111.009888     11.0     84.00    141.0
..      ...      ...      ...      ...      ...      ...
364     46.0  19011.108696   90.436774  18867.0  18934.75  19009.5
365     48.0  19023.895833   83.428206  18868.0  18968.25  19032.5
366     51.0  19027.352941   88.813923  18869.0  18960.50  19013.0
367     36.0  19023.333333  101.943682  18870.0  18927.75  19001.5
368     72.0  19020.611111   82.448713  18873.0  18956.00  19030.0

      userIDs      ... OutcomeScore      wordCount \
      75%      max      count      mean      ...      75%      max      count
0      270.00    345.0     65.0      2.0      ...      4.0     4.0     65.0
1      214.25    327.0     22.0      3.0      ...      4.0     4.0     22.0
2      284.50    331.0     31.0      4.0      ...      4.0     4.0     31.0
3      260.00    323.0     45.0      5.0      ...      2.0     2.0     45.0
```



4	288.00	343.0	37.0	6.0	...	2.0	2.0	37.0
..	...	...	...	...	...	...	...	...
364	19082.25	19175.0	46.0	389.0	...	7.0	7.0	46.0
365	19078.25	19173.0	48.0	390.0	...	4.0	4.0	48.0
366	19109.50	19174.0	51.0	391.0	...	5.0	5.0	51.0
367	19122.25	19176.0	36.0	392.0	...	5.0	5.0	36.0
368	19075.50	19179.0	72.0	393.0	...	4.0	4.0	72.0

	mean	std	min	25%	50%	75%	max
0	10.830769	9.032181	1.0	4.0	9.0	14.00	35.0
1	7.136364	6.613806	1.0	3.0	5.5	9.00	27.0
2	11.258065	8.229896	1.0	5.0	10.0	15.00	31.0
3	7.600000	5.532548	1.0	4.0	7.0	11.00	22.0
4	10.783784	8.888617	1.0	5.0	8.0	15.00	43.0
..	...	...	...	...	...	...	...
364	16.413043	11.960260	1.0	5.0	17.0	21.75	53.0
365	8.979167	8.926388	1.0	2.0	6.0	13.50	32.0
366	11.862745	12.935254	1.0	3.0	9.0	15.50	70.0
367	12.527778	11.319985	1.0	2.0	9.0	20.00	37.0
368	12.055556	11.056690	1.0	2.0	10.0	17.25	43.0

[369 rows x 96 columns]

```
[245]: user_group.count()
```

```
[245]:
```

	userIDs	Unnamed: 0	content	RoleName	implementation	Line_ID	\
0	2	65	65	65	65	65	
1	3	22	22	22	22	22	
2	4	31	31	31	31	31	
3	5	45	45	45	45	45	
4	6	37	37	37	37	37	
..	...	...	...	...	...	...	
364	389	46	46	46	46	46	
365	390	48	48	48	48	48	
366	391	51	51	51	51	51	
367	392	36	36	36	36	36	
368	393	72	72	72	72	72	

	ChatGroup	group_id	roomName	m_experimental_testing	\
0	65	65	65	65	
1	22	22	22	22	
2	31	31	31	31	
3	45	45	45	45	
4	37	37	37	37	
..	...	...	...	...	
364	46	46	46	46	

365	48	48	48	48
366	51	51	51	51
367	36	36	36	36
368	72	72	72	72

	m_making_design_choices	m_asking_questions	\
0	65	65	
1	22	22	
2	31	31	
3	45	45	
4	37	37	
..	...	...	
364	46	46	
365	48	48	
366	51	51	
367	36	36	
368	72	72	

	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	65	65	
1	22	22	
2	31	31	
3	45	45	
4	37	37	
..	...	...	
364	46	46	
365	48	48	
366	51	51	
367	36	36	
368	72	72	

	j_communication	OutcomeScore	wordCount	Group
0	65	65	65	65
1	22	22	22	22
2	31	31	31	31
3	45	45	45	45
4	37	37	37	37
..	...	...	...	...
364	46	46	46	46
365	48	48	48	48
366	51	51	51	51
367	36	36	36	36
368	72	72	72	72

[369 rows x 18 columns]

```
[246]: #Adds news column with actual outcome score for each student
user_sum['NewOutcomeScore'] = user_sum['OutcomeScore'] / user_group.
      ↳count()['OutcomeScore']
user_data = user_sum
user_data
```

```
[246]:
```

	userIDs	Unnamed: 0	Line_ID	group_id	m_experimental_testing	\
0	2	12286	12286	130	2	
1	3	2796	2796	44	1	
2	4	6454	6454	62	2	
3	5	9385	9385	90	0	
4	6	6740	6740	74	0	
..	...	...	...	...	...	
364	389	874511	874649	276	2	
365	390	913147	913291	288	0	
366	391	970395	970548	306	2	
367	392	684840	684948	216	2	
368	393	1369484	1369700	432	5	

	m_making_design_choices	m_masking_questions	\
0	4	17	
1	4	3	
2	2	3	
3	0	6	
4	2	7	
..	...	...	
364	8	4	
365	5	6	
366	3	11	
367	4	4	
368	6	8	

	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0	4	
1	0	1	
2	1	5	
3	0	2	
4	1	3	
..	...	...	
364	1	7	
365	1	2	
366	1	4	
367	0	2	
368	1	7	

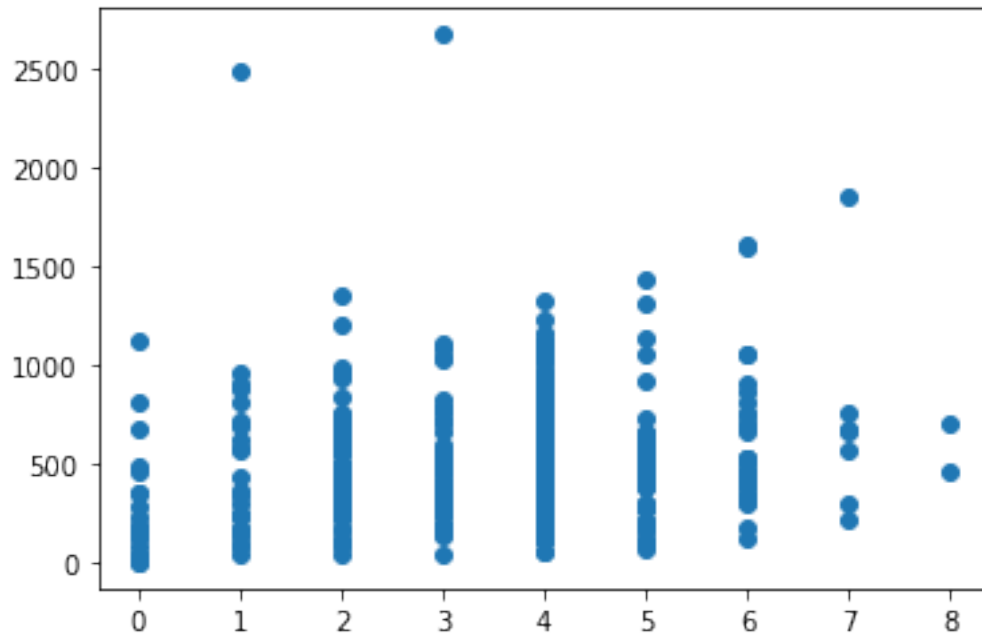
	j_communication	OutcomeScore	wordCount	NewOutcomeScore
0	0	260	704	4.0

1	0	88	157	4.0
2	0	124	349	4.0
3	0	90	342	2.0
4	0	74	399	2.0
...	...	...	...	...
364	0	322	755	7.0
365	0	192	431	4.0
366	0	255	605	5.0
367	0	180	451	5.0
368	0	288	868	4.0

[369 rows x 13 columns]

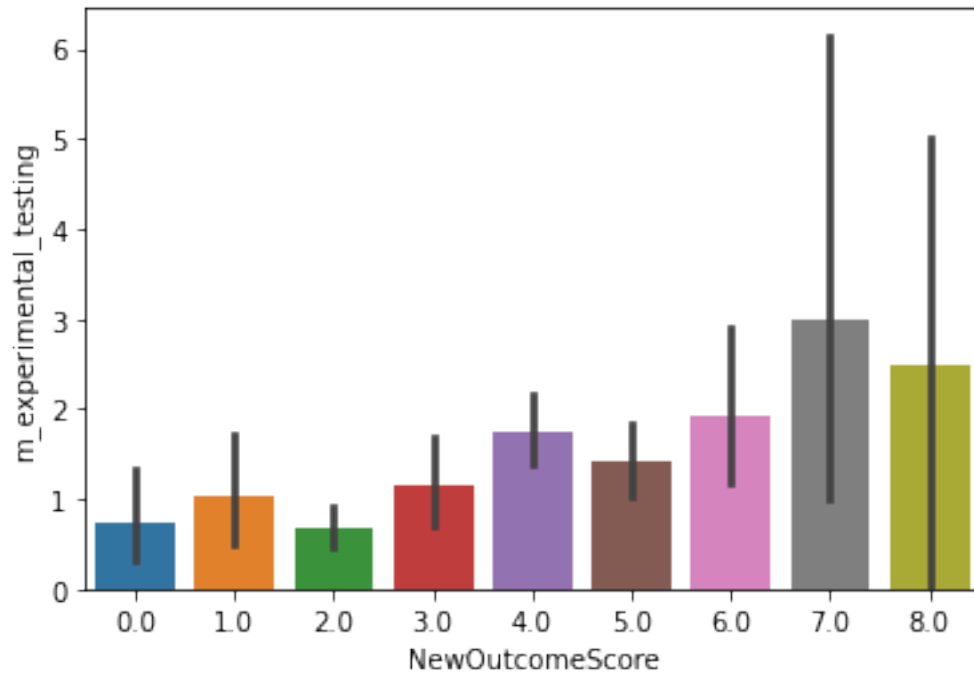
```
[247]: plt.scatter(user_data[['NewOutcomeScore']], user_data[['wordCount']])
```

```
[247]: <matplotlib.collections.PathCollection at 0x2c6368dcac0>
```



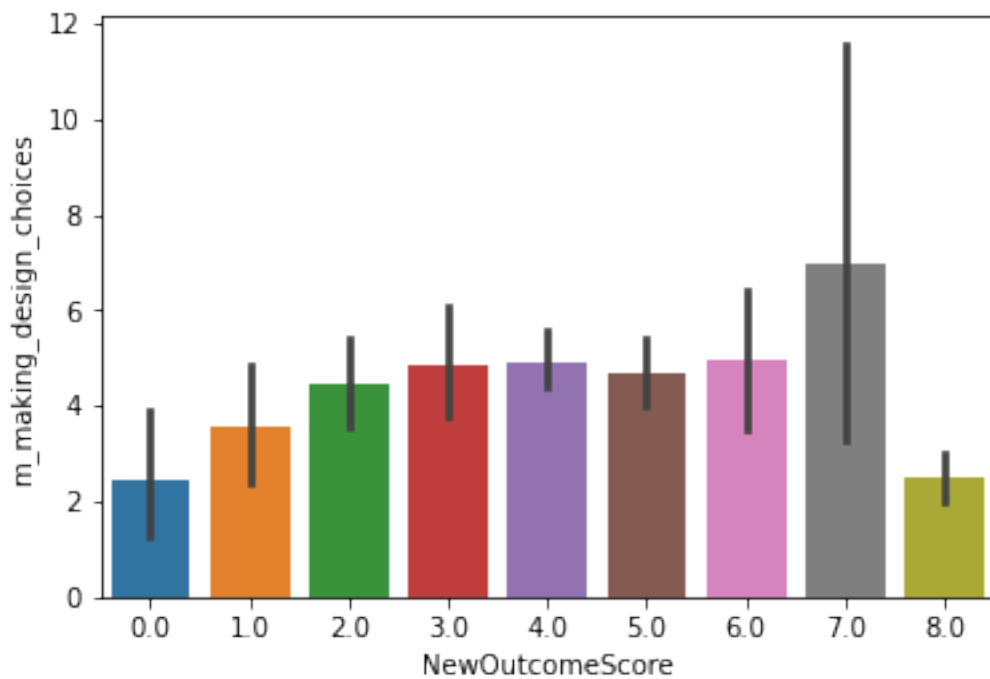
```
[248]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = 'm_experimental_testing')
```

```
[248]: <AxesSubplot:xlabel='NewOutcomeScore', ylabel='m_experimental_testing'>
```



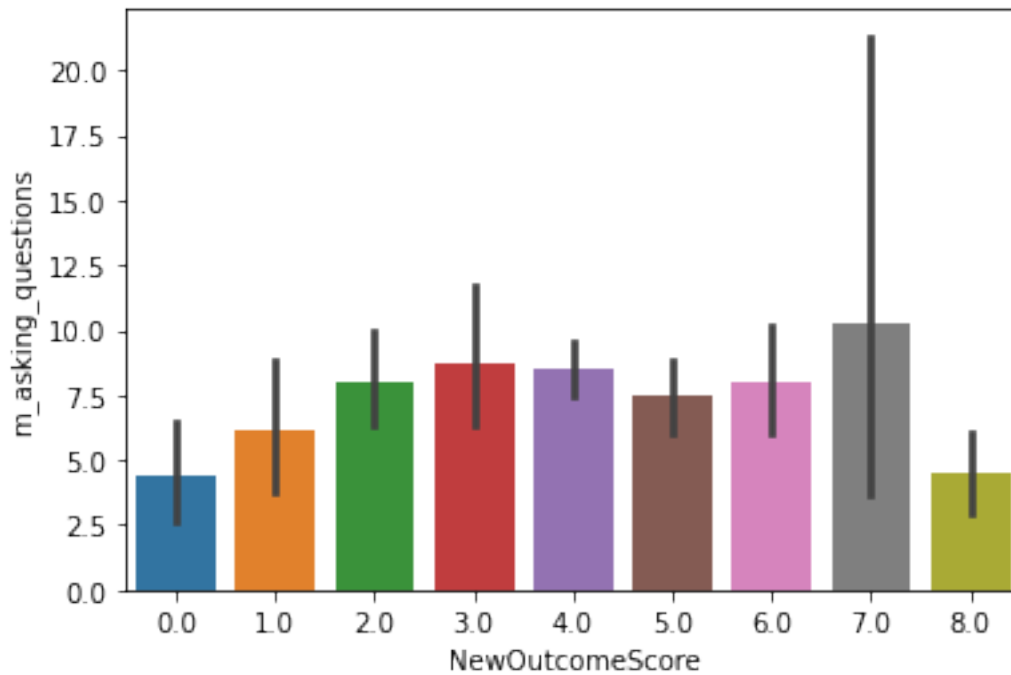
```
[249]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = 'm_making_design_choices')
```

```
[249]: <AxesSubplot:xlabel='NewOutcomeScore', ylabel='m_making_design_choices'>
```



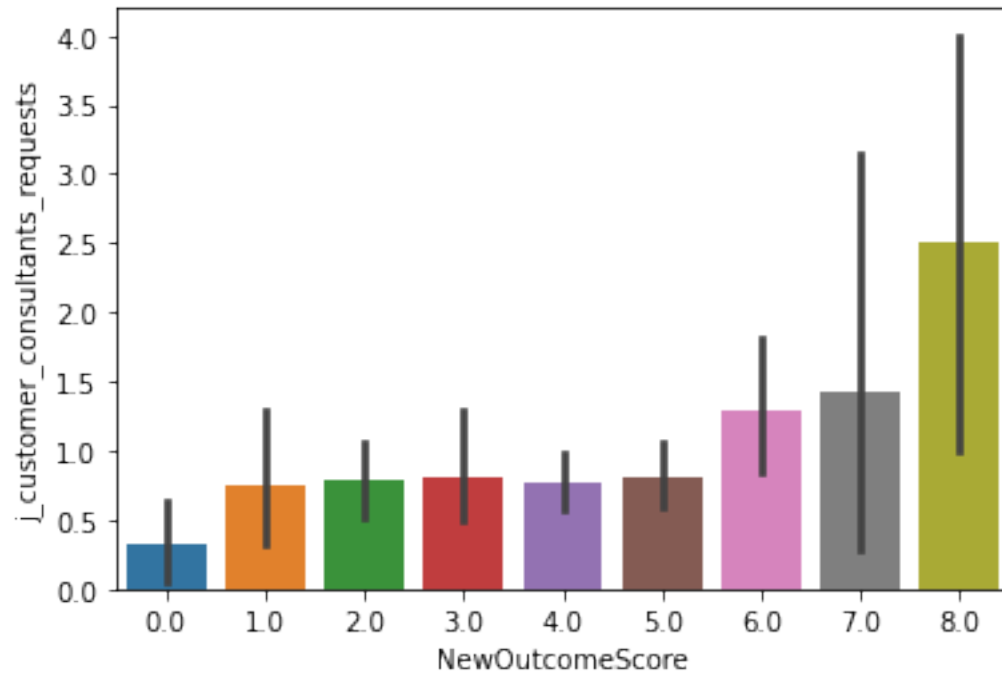
```
[250]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = 'm_asking_questions')
```

```
[250]: <AxesSubplot:xlabel='NewOutcomeScore', ylabel='m_asking_questions'>
```



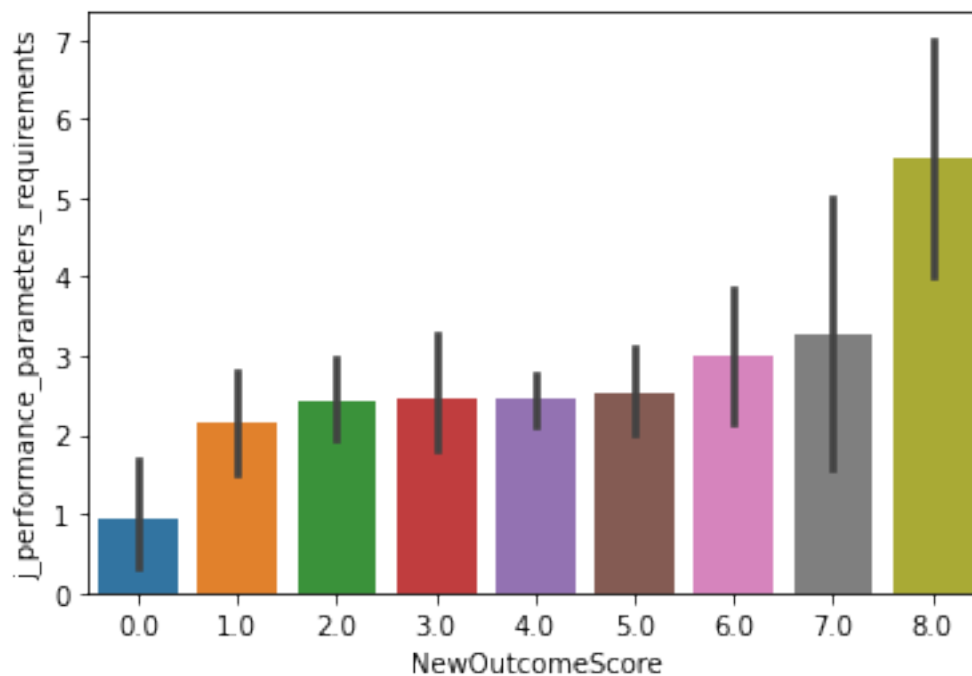
```
[251]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = 'j_customer_consultants_requests')
```

```
[251]: <AxesSubplot:xlabel='NewOutcomeScore', ylabel='j_customer_consultants_requests'>
```



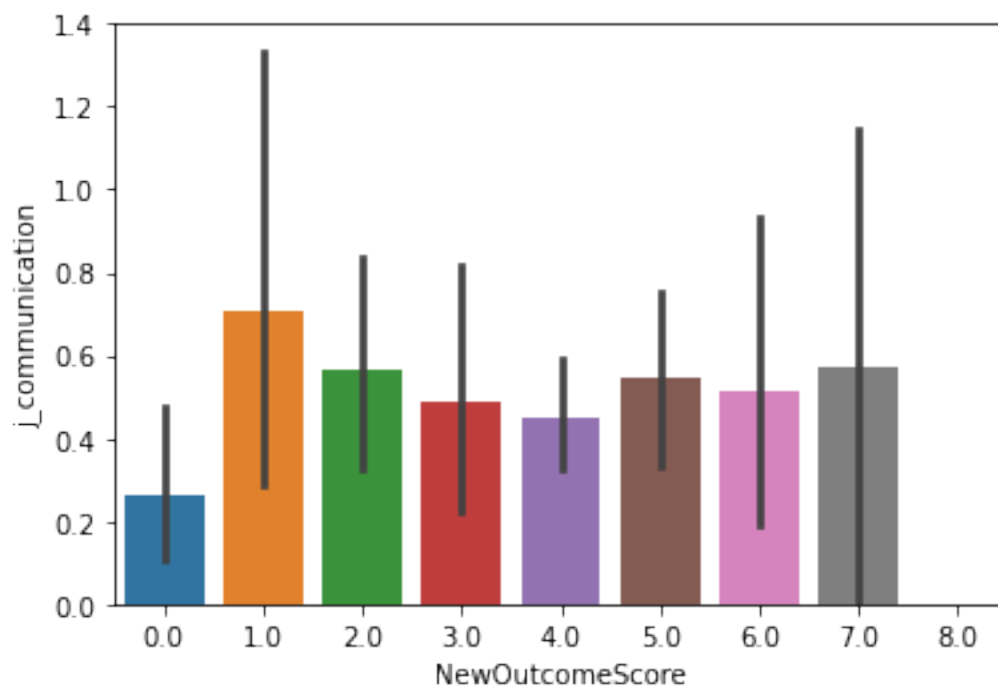
```
[252]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y =  
        ↪ 'j_performance_parameters_requirements')
```

```
[252]: <AxesSubplot:xlabel='NewOutcomeScore',  
        ylabel='j_performance_parameters_requirements'>
```



```
[253]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = 'j_communication')
```

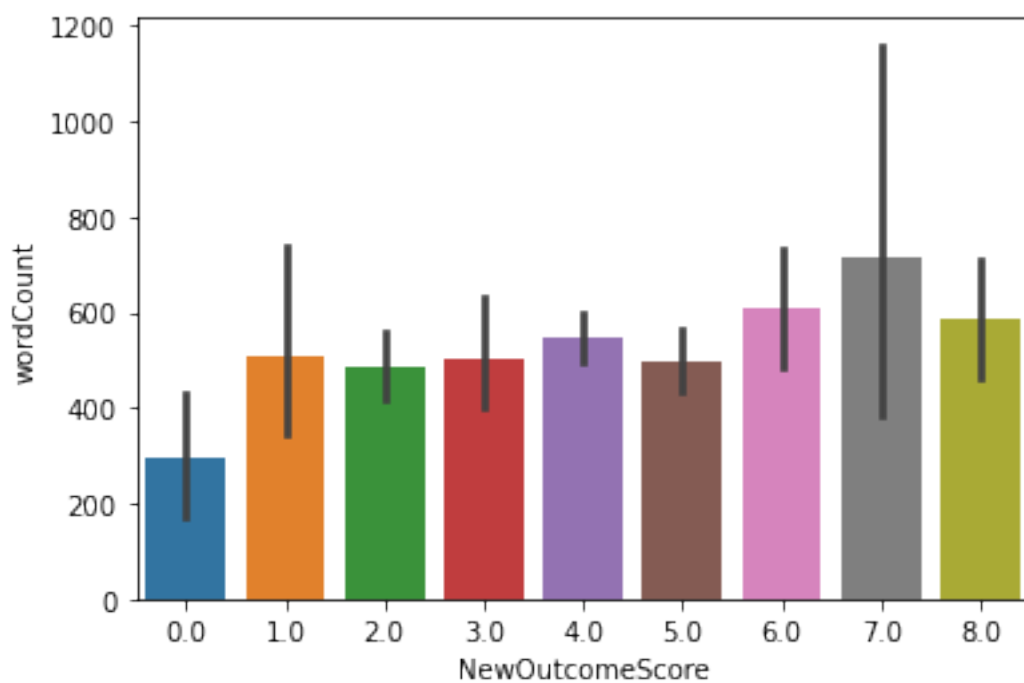
```
[253]: <AxesSubplot:xlabel='NewOutcomeScore', ylabel='j_communication'>
```



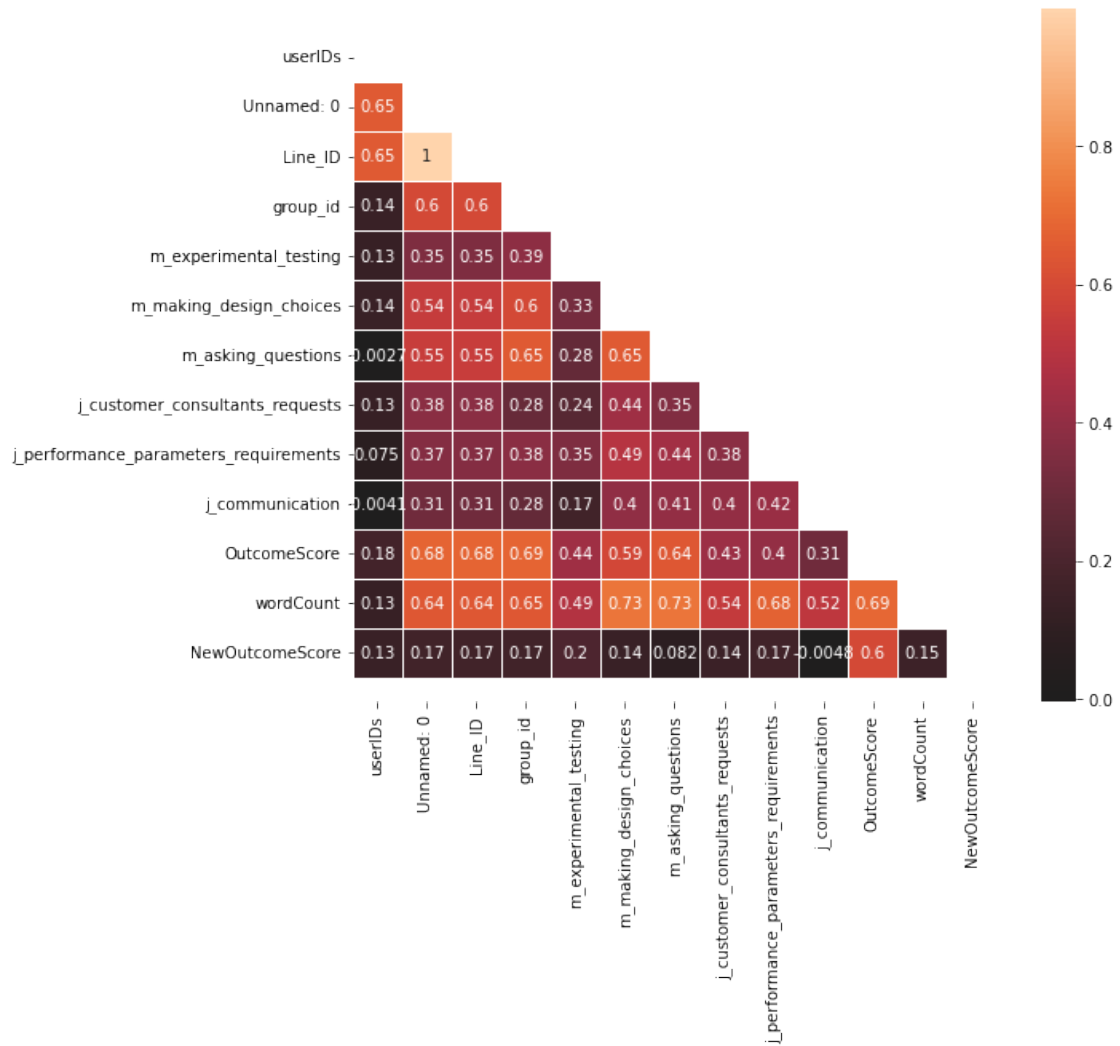


```
[254]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = 'wordCount')
```

```
[254]: <AxesSubplot:xlabel='NewOutcomeScore', ylabel='wordCount'>
```



```
[255]: corr = user_data.corr()  
plt_correlation_matrix(corr)
```



## 2 Modelling

[256]: user\_data

```
[256]:
```

	userIDs	Unnamed: 0	Line_ID	group_id	m_experimental_testing	\
0	2	12286	12286	130	2	
1	3	2796	2796	44	1	
2	4	6454	6454	62	2	
3	5	9385	9385	90	0	
4	6	6740	6740	74	0	
..	...	...	...	...	...	
364	389	874511	874649	276	2	
365	390	913147	913291	288	0	

366	391	970395	970548	306	2
367	392	684840	684948	216	2
368	393	1369484	1369700	432	5

	m_making_design_choices	m_asking_questions	\
0	4	17	
1	4	3	
2	2	3	
3	0	6	
4	2	7	
..	...	...	
364	8	4	
365	5	6	
366	3	11	
367	4	4	
368	6	8	

	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0	4	
1	0	1	
2	1	5	
3	0	2	
4	1	3	
..	...	...	
364	1	7	
365	1	2	
366	1	4	
367	0	2	
368	1	7	

	j_communication	OutcomeScore	wordCount	NewOutcomeScore
0	0	260	704	4.0
1	0	88	157	4.0
2	0	124	349	4.0
3	0	90	342	2.0
4	0	74	399	2.0
..	...	...	...	...
364	0	322	755	7.0
365	0	192	431	4.0
366	0	255	605	5.0
367	0	180	451	5.0
368	0	288	868	4.0

[369 rows x 13 columns]

## 2.1 Linear Regression

These seem to correspond to higher scores from our EDA.

```
[257]: XX = user_data.drop(['NewOutcomeScore', 'OutcomeScore', 'group_id', 'Line_ID',  
    ↳ 'Unnamed: 0', 'userIDs'], axis = 1)  
Y = user_data['NewOutcomeScore']  
  
#Normalise for regularisation later  
X = (XX-XX.mean())/XX.std()  
  
feature_names = X.columns  
  
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,  
    ↳ random_state=42)
```

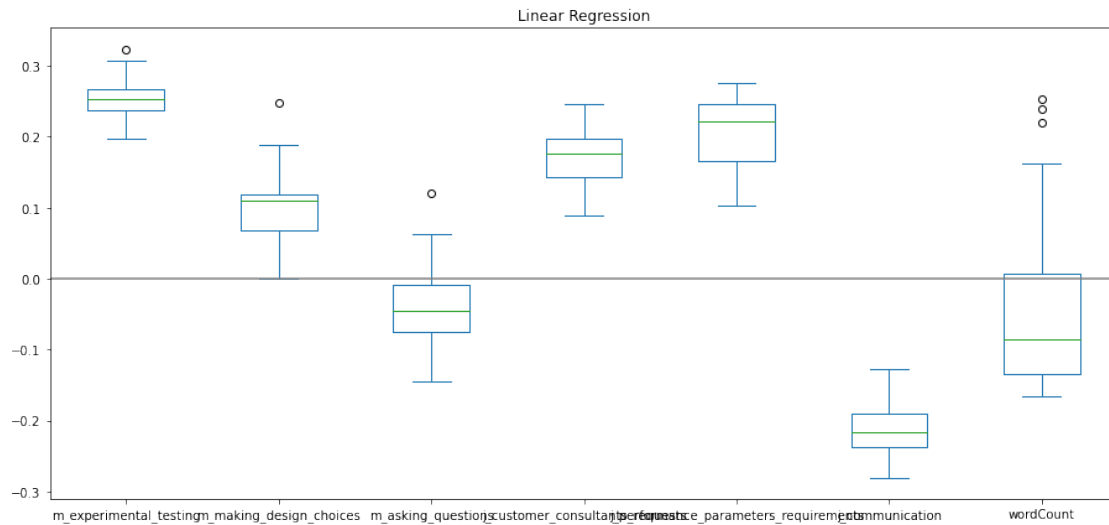
```
[258]: linear = LinearRegression() # instantiate the linear regression model  
linear.fit(X_train, Y_train) # fit the data to the model  
training_score = linear.score(X_train, Y_train) # calculate rsq for the training  
    ↳ set  
  
# use the independent variables for the testing set to predict the target  
    ↳ variable  
preds_linear = linear.predict(X_test)  
  
# calculate the correlation of the predicted and actual target variables  
rsquared_linear = r2_score(Y_test, preds_linear)  
  
# print the training and testing scores  
print(f'Training score is {training_score:.3f}')  
print(f'Testing score is {rsquared_linear:.3f}')
```

Training score is 0.108

Testing score is -0.164

```
[259]: scores = cross_validate(  
    linear, X, Y, cv=RepeatedKFold(n_splits=5, n_repeats=5, random_state=42),  
    return_estimator=True  
)  
# take the results for each simulation (estimator), extract the coefficients for  
    ↳ each run  
# and add them to a dataframe with columns being the feature names  
coefs = pd.DataFrame([est.coef_ for est in scores['estimator']], columns =  
    ↳ feature_names)  
# plot the descriptive statistics of the coefficients in a box and whisker plot to  
    ↳ show variability  
ax = coefs.plot(kind='box', figsize=(20, 7))  
plt.title('Linear Regression')
```

```
plt.axhline(y=0, color='.5')
plt.subplots_adjust(left=.3)
```



Most features seem to have quite a high variance but wordCount is the highest. Also m\_asking\_questions and wordCount seem to have less of an effect than the other features.

### 2.1.1 Ridge Regularisation

```
[260]: clf = RidgeCV(alphas = np.logspace(-2, 2, num=21)).fit(X, Y)
print(f"Optimal Alpha: {clf.alpha_}")
```

Optimal Alpha: 100.0

```
[261]: ridge = Ridge(alpha = 100, random_state = 1234) # instantiate Ridge
        ↪ regularisation with the current alfa
ridge.fit(X_train, Y_train)

print(f"Ridge Training Score: {ridge.score(X_train, Y_train)}")
preds_linear = ridge.predict(X_test)
print(f"Ridge Testing Score: {r2_score(Y_test, preds_linear)}")
```

Ridge Training Score: 0.10056014528221568

Ridge Testing Score: -0.08706108291429526

This is slightly better than the original but really there is not much difference.

## 2.1.2 Lasso Regularisation

```
[262]: clf = LassoCV(alphas = np.logspace(-2, 2, num=21)).fit(X, Y)
print(f"Optimal Alpha: {clf.alpha_}")
```

Optimal Alpha: 0.025118864315095794

```
[263]: lasso = Lasso(alpha = 0.025, random_state = 1234) # instantiate Ridge
      ↪regularisation with the current alfa
lasso.fit(X_train, Y_train)

print(f"Ridge Training Score: {lasso.score(X_train, Y_train)}")
preds_linear = ridge.predict(X_test)
print(f"Ridge Testing Score: {r2_score(Y_test, preds_linear)}")
```

Ridge Training Score: 0.10516819324125048

Ridge Testing Score: -0.08706108291429526

Basically the same as ridge regularisation.

## 2.2 SVM Modelling

```
[264]: user_data
```

```
[264]:
```

	userIDs	Unnamed: 0	Line_ID	group_id	m_experimental_testing	\
0	2	12286	12286	130	2	
1	3	2796	2796	44	1	
2	4	6454	6454	62	2	
3	5	9385	9385	90	0	
4	6	6740	6740	74	0	
..	...	...	...	...	...	
364	389	874511	874649	276	2	
365	390	913147	913291	288	0	
366	391	970395	970548	306	2	
367	392	684840	684948	216	2	
368	393	1369484	1369700	432	5	

	m_making_design_choices	m_asking_questions	\
0	4	17	
1	4	3	
2	2	3	
3	0	6	
4	2	7	
..	...	...	
364	8	4	
365	5	6	
366	3	11	
367	4	4	
368	6	8	

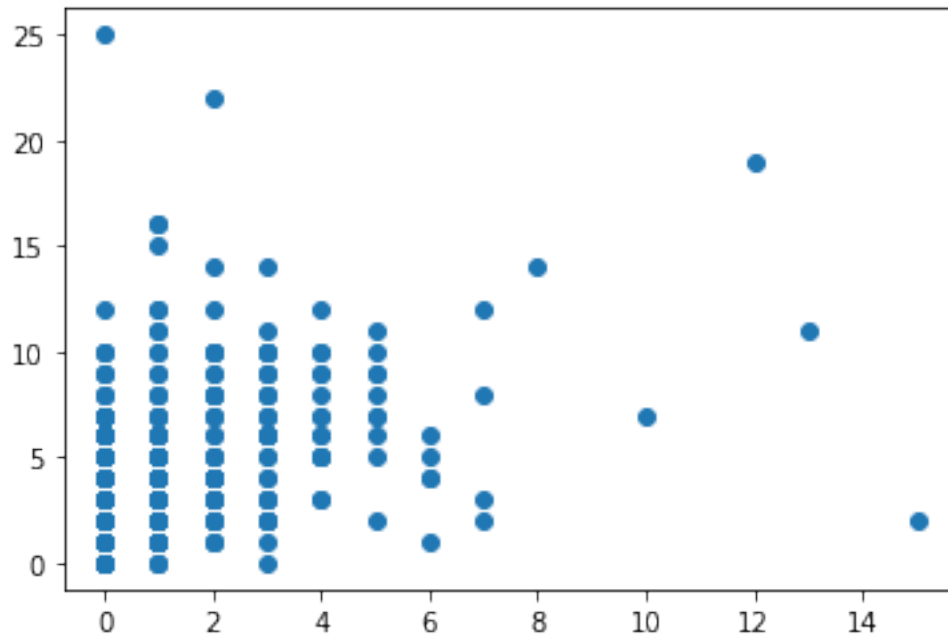
	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0		4
1	0		1
2	1		5
3	0		2
4	1		3
..	...		...
364	1		7
365	1		2
366	1		4
367	0		2
368	1		7

	j_communication	OutcomeScore	wordCount	NewOutcomeScore
0	0	260	704	4.0
1	0	88	157	4.0
2	0	124	349	4.0
3	0	90	342	2.0
4	0	74	399	2.0
..	...	...	...	...
364	0	322	755	7.0
365	0	192	431	4.0
366	0	255	605	5.0
367	0	180	451	5.0
368	0	288	868	4.0

[369 rows x 13 columns]

```
[265]: plt.scatter(user_data[['m_experimental_testing']],
    ↪user_data[['m_making_design_choices']])
```

```
[265]: <matplotlib.collections.PathCollection at 0x2c636eedac0>
```



## 2.3 Decision Trees

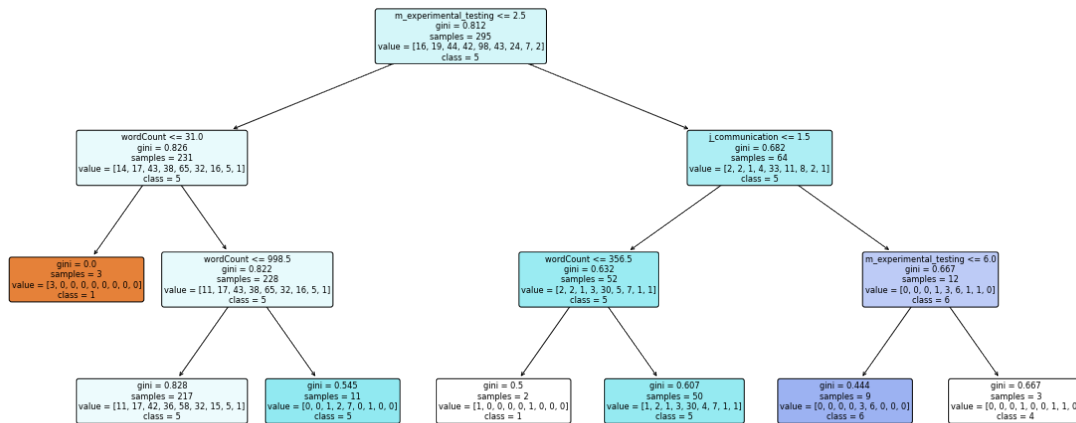
```
[266]: X = user_data.drop(['NewOutcomeScore', 'OutcomeScore', 'group_id', 'Line_ID',
    → 'Unnamed: 0', 'userIDs'], axis = 1)
Y = user_data['NewOutcomeScore']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,
    → random_state=42)
```

```
[267]: df_dtc = DecisionTreeClassifier(random_state = 42, max_depth = 3)
df_dtc = df_dtc.fit(X_train, Y_train)

fig, ax = plt.subplots(figsize=(20,9)) # initialise the plots and axes
# plot the decision tree for the model df_dtc
plot_tree(df_dtc,
    filled = True, # colour the nodes according to the classification
    rounded = True, # make the nodes have rounded corners
    class_names = ['1', '2', '3', '4', '5', '6', '7', '8'], # use these names for
    → targets
    feature_names = X.columns # use these names for features
);
```





```
[268]: for i in range(1, 21):
    df_dtc = DecisionTreeClassifier(random_state = 42, max_depth = i)
    df_dtc = df_dtc.fit(X_train, Y_train)

    depth = df_dtc.tree_.max_depth

    Y_pred = df_dtc.predict(X_test)
    accuracy = np.round(accuracy_score(Y_test, Y_pred), 3)
    print(f"Depth of {depth} gives accuracy: {accuracy}")
```

```

Depth of 1 gives accuracy: 0.365
Depth of 2 gives accuracy: 0.378
Depth of 3 gives accuracy: 0.365
Depth of 4 gives accuracy: 0.365
Depth of 5 gives accuracy: 0.27
Depth of 6 gives accuracy: 0.27
Depth of 7 gives accuracy: 0.23
Depth of 8 gives accuracy: 0.311
Depth of 9 gives accuracy: 0.297
Depth of 10 gives accuracy: 0.243
Depth of 11 gives accuracy: 0.243
Depth of 12 gives accuracy: 0.27
Depth of 13 gives accuracy: 0.203
Depth of 14 gives accuracy: 0.23
Depth of 15 gives accuracy: 0.23
Depth of 16 gives accuracy: 0.243
Depth of 17 gives accuracy: 0.27
Depth of 18 gives accuracy: 0.257
Depth of 19 gives accuracy: 0.216
Depth of 20 gives accuracy: 0.243
  
```

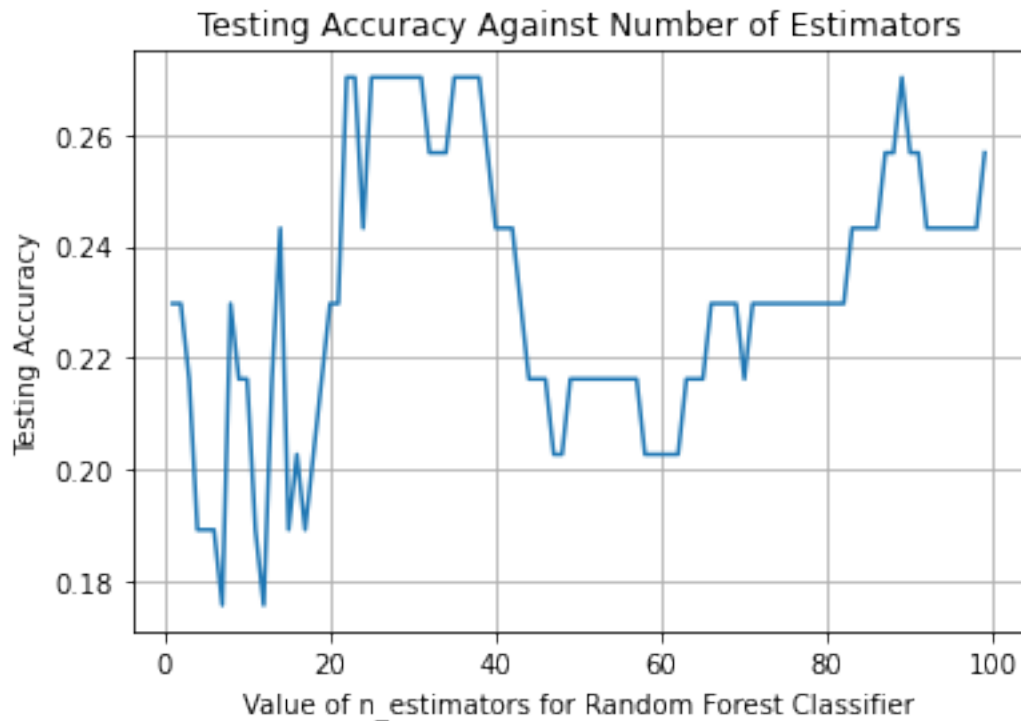
From this and other manual testing it can be seen that the best accuracy occurs at a depth of 2, giving 0.378 accuracy. But this should be improved upon by random forests as it is a group of decision trees.

## 2.4 Random Forests

```
[269]: X = user_data.drop(['NewOutcomeScore', 'OutcomeScore', 'group_id', 'Line_ID',  
→ 'Unnamed: 0', 'userIDs'], axis = 1)  
Y = user_data['NewOutcomeScore']  
  
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,  
→ random_state=42)
```

### 2.4.1 Finding best parameters

```
[270]: scores = []  
for i in range(1, 100):  
    rfc = RandomForestClassifier(n_estimators = i, random_state = 42)  
    rfc.fit(X_train, Y_train)  
    Y_pred = rfc.predict(X_test)  
    scores.append(accuracy_score(Y_test, Y_pred))  
  
plt.plot(range(1, 100), scores)  
plt.xlabel('Value of n_estimators for Random Forest Classifier')  
plt.ylabel('Testing Accuracy')  
plt.title('Testing Accuracy Against Number of Estimators')  
plt.grid(True);
```



We can see from this we can see that the accuracy hits a peak at around 30 estimators and then drops off and varies.

```
[271]: rfc = RandomForestClassifier(n_estimators = 30, random_state = 42).fit(X_train, Y_train)
Y_pred = rfc.predict(X_test)
accuracy = np.round(accuracy_score(Y_test, Y_pred), 3)

print(f"Accuracy: {accuracy}")
```

Accuracy: 0.27

We can then find the importance of each feature in the model.

```
[272]: for name, score in zip(X.columns, rfc.feature_importances_):
        print(name, np.round(score, 3))
```

```
m_experimental_testing 0.098
m_making_design_choices 0.171
m_asking_questions 0.197
j_customer_consultants_requests 0.08
j_performance_parameters_requirements 0.127
j_communication 0.059
wordCount 0.267
```

This is an improvement on some of the previous models but really is still not a very good predictor

at all of student grades.

## 2.5 KMeans Clustering

```
[273]: XX = user_data.drop(['NewOutcomeScore', 'OutcomeScore', 'group_id', 'Line_ID', 'OutcomeScore', 'userIDs'], axis = 1)
      → 'Unnamed: 0', 'userIDs'], axis = 1)
      Y = user_data['NewOutcomeScore']

      #Normalise for regularisation later
      X = (XX-XX.mean())/XX.std()
```

This is just a basic model using 8 clusters and it gives an inertia score of around 1010.

```
[274]: kmeans = KMeans(n_clusters = 3, random_state = 42)
      kmeans.fit(X)
      y_kmeans = kmeans.predict(X)
      print(kmeans.inertia_)
```

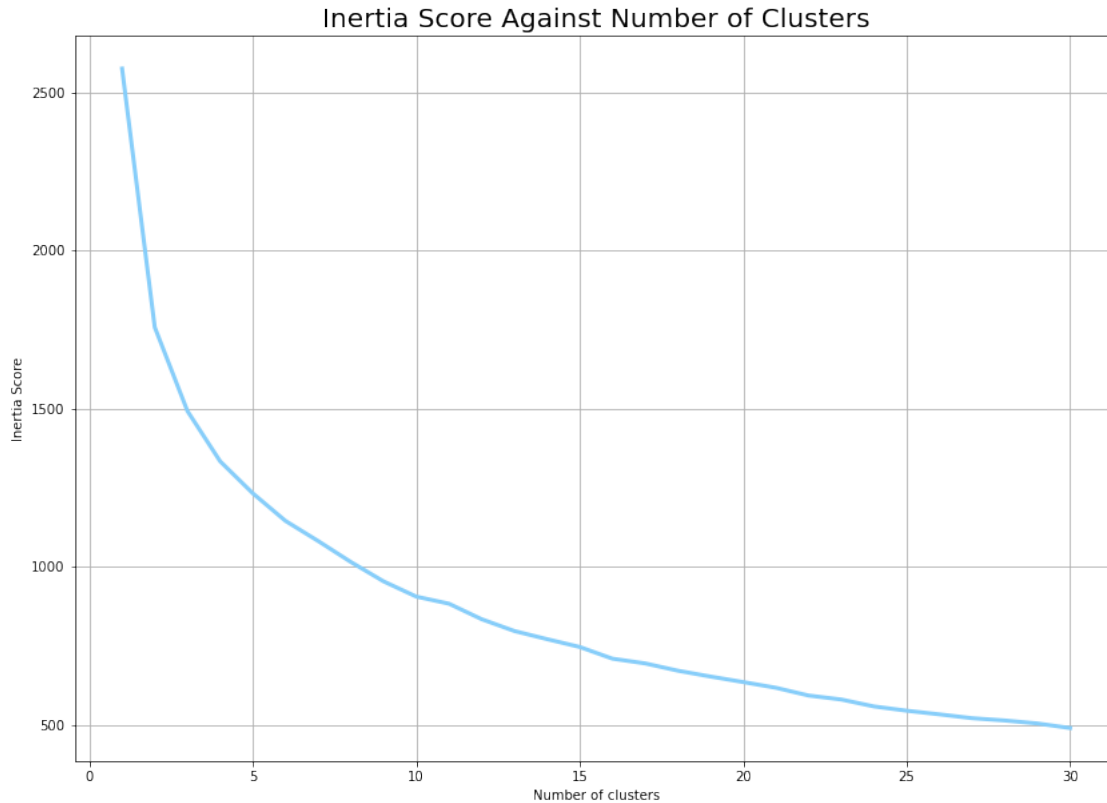
1492.2549855427221

We can graph the inertia score against the number of clusters and see that it decreases the more clusters we add. This is good but having too many clusters means it is very specified for this data.

```
[275]: nclusters = np.arange(1,31)
      inertia_score = []
      for ncl in nclusters: #Loop over the number of clusters
          kmeans = KMeans(n_clusters = ncl, random_state = 42)
          kmeans.fit(X)
          inertia_score.append(kmeans.inertia_)

      fig, ax = plt.subplots(figsize=(14,10))
      plt.plot(nclusters, inertia_score, color = 'lightskyblue', linewidth = 3)
      plt.grid(True)
      plt.xlabel('Number of clusters')
      plt.ylabel('Inertia Score')
      plt.title('Inertia Score Against Number of Clusters', size = 20);
```

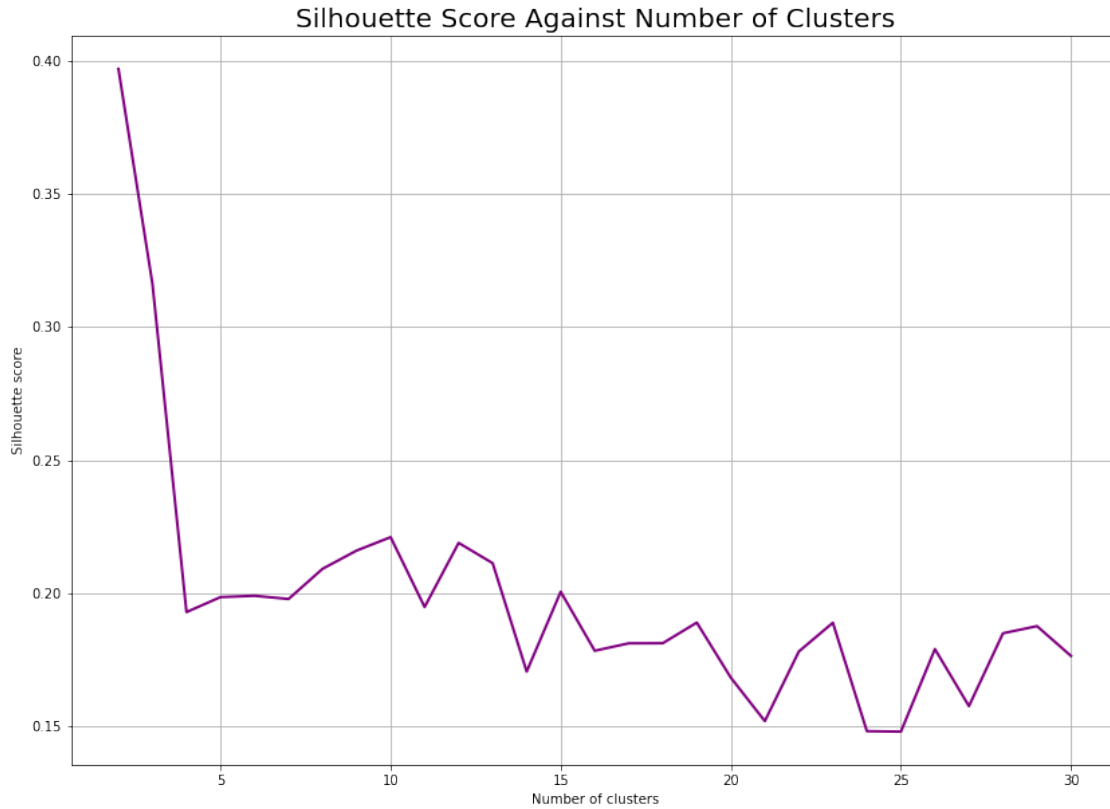
```
C:\Users\owner\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:881:
UserWarning: KMeans is known to have a memory leak on Windows with MKL, when
there are less chunks than available threads. You can avoid it by setting the
environment variable OMP_NUM_THREADS=2.
  warnings.warn(
```



We can graph the silhouette score against the number of clusters to see which amount of clusters gives the maximum value of the silhouette score. It can be seen that it happens around 10 clusters before degrading. Since the silhouette score isn't very high it means that a lot of the points are close to the borders of other clusters, reducing that accuracy of our model.

```
[276]: nclusters = np.arange(2,31) #n_clusters from 2 to 30
sil_score = [] #List for the silhouette scores
for cluster in nclusters: # loop over the number of clusters
    kmeans = KMeans(n_clusters = cluster, random_state = 42).fit(X)
    sil_score.append(silhouette_score(X, kmeans.labels_))

# plot the silhouette score
fig, ax = plt.subplots(figsize=(14, 10))
plt.plot(nclusters, sil_score, color = 'purple', linewidth = 2)
plt.grid(True)
plt.xlabel('Number of clusters')
plt.ylabel('Silhouette score')
plt.title('Silhouette Score Against Number of Clusters', size = 20);
```



KMeans model with 10 clusters.

```
[277]: kmeans = KMeans(n_clusters = 10, random_state = 42)
kmeans.fit(X)
y_kmeans = kmeans.predict(X)
inertia_score = kmeans.inertia_
print(f"Inertia Score: {np.round(inertia_score,3)}")
```

Inertia Score: 905.084

## 2.6 K Nearest Neighbours

```
[278]: user_sum['NewOutcomeScore'] = user_sum['OutcomeScore'] / user_group.
        ↳count()['OutcomeScore']
user_data = user_sum
user_data
```

```
[278]:
```

	userIDs	Unnamed: 0	Line_ID	group_id	m_experimental_testing	\
0	2	12286	12286	130	2	
1	3	2796	2796	44	1	
2	4	6454	6454	62	2	
3	5	9385	9385	90	0	

4	6	6740	6740	74	0
..	...	...	...	...	...
364	389	874511	874649	276	2
365	390	913147	913291	288	0
366	391	970395	970548	306	2
367	392	684840	684948	216	2
368	393	1369484	1369700	432	5

	m_making_design_choices	m_asking_questions	\
0	4	17	
1	4	3	
2	2	3	
3	0	6	
4	2	7	
..	...	...	
364	8	4	
365	5	6	
366	3	11	
367	4	4	
368	6	8	

	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0	4	
1	0	1	
2	1	5	
3	0	2	
4	1	3	
..	...	...	
364	1	7	
365	1	2	
366	1	4	
367	0	2	
368	1	7	

	j_communication	OutcomeScore	wordCount	NewOutcomeScore
0	0	260	704	4.0
1	0	88	157	4.0
2	0	124	349	4.0
3	0	90	342	2.0
4	0	74	399	2.0
..	...	...	...	...
364	0	322	755	7.0
365	0	192	431	4.0
366	0	255	605	5.0
367	0	180	451	5.0
368	0	288	868	4.0

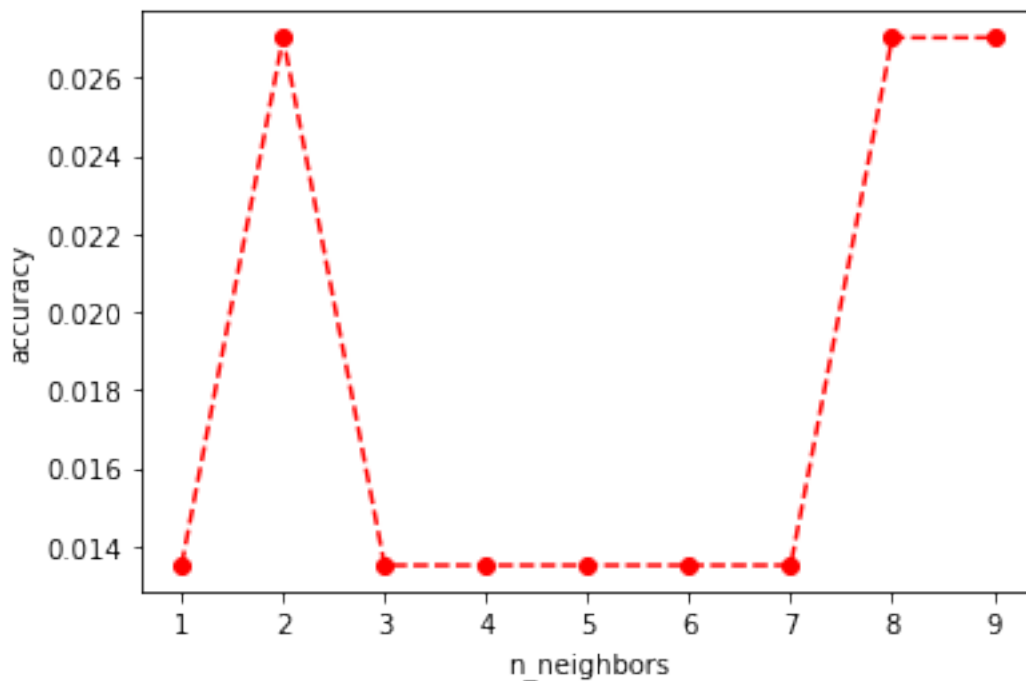
[369 rows x 13 columns]

```
[279]: X = user_data[['m_experimental_testing', 'm_making_design_choices',  
→ 'm_asking_questions', 'j_customer_consultants_requests',  
→ 'j_performance_parameters_requirements', 'j_communication', 'wordCount']]  
Y = user_data['OutcomeScore']
```

```
[280]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2,  
→ random_state=42)
```

```
[281]: # Check Accuracy  
Accuracy = []  
for i in range(1,10):  
    knn = KNeighborsClassifier(n_neighbors = i).fit(X_train, Y_train)  
    Y_pred = knn.predict(X_test)  
    Accuracy.append(metrics.accuracy_score(Y_test, Y_pred))  
plt.plot(range(1,10), Accuracy, color = 'red', linestyle = 'dashed', marker =  
→ 'o')  
plt.xlabel('n_neighbors')  
plt.ylabel('accuracy')
```

```
[281]: Text(0, 0.5, 'accuracy')
```





```
[282]: knn = KNeighborsClassifier(n_neighbors = 9).fit(X_train, Y_train)
Y_pred = knn.predict(X_test)
metrics.accuracy_score(Y_test, Y_pred)
```

```
[282]: 0.02702702702702703
```

## 2.7 Boosting

```
[283]: X = user_data[['m_experimental_testing', 'm_making_design_choices',
    → 'm_asking_questions', 'j_customer_consultants_requests',
    → 'j_performance_parameters_requirements', 'j_communication', 'wordCount']]
Y = user_data['OutcomeScore']

X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2,
    → random_state=42)
```

```
[284]: ada_clf = AdaBoostClassifier(DecisionTreeClassifier(max_depth=5),
n_estimators=10, learning_rate=0.2)
ada_clf.fit(X_train, y_train)
y_pred = ada_clf.predict(X_test) # calculate the predicted values
# print the accuracy of the classifier
print('Accuracy {0}'.format(np.round(accuracy_score(y_test, y_pred),3)))
```

Accuracy 0.027

```
[285]: lrates = np.arange(0.2, 2.3, 0.3)
nests = np.arange(5, 105, 5)
results = np.zeros((len(nests), len(lrates)))
for i in range(len(lrates)):
    for j in range(len(nests)):
        print(i, j, end='\r')
        ada_clf = AdaBoostClassifier(n_estimators=nests[j],
    → learning_rate=lrates[i])
        ada_clf.fit(X_train, y_train)
        y_pred = ada_clf.predict(X_test) # calculate the predicted values
        # evaluate the model and collect the results
        results[j,i] = accuracy_score(y_test, y_pred)

results_df = pd.DataFrame(results, columns=np.round(lrates, 1), index=nests)
results_df
```

6 197

```
[285]:
```

	0.2	0.5	0.8	1.1	1.4	1.7	2.0
5	0.040541	0.040541	0.040541	0.027027	0.040541	0.027027	0.040541
10	0.040541	0.040541	0.040541	0.040541	0.027027	0.027027	0.094595
15	0.040541	0.040541	0.040541	0.040541	0.040541	0.027027	0.027027
20	0.040541	0.040541	0.040541	0.027027	0.040541	0.027027	0.027027

25	0.040541	0.040541	0.040541	0.040541	0.040541	0.027027	0.027027
30	0.040541	0.040541	0.040541	0.040541	0.040541	0.027027	0.013514
35	0.040541	0.040541	0.040541	0.027027	0.040541	0.027027	0.013514
40	0.040541	0.040541	0.040541	0.027027	0.027027	0.027027	0.000000
45	0.040541	0.040541	0.054054	0.040541	0.027027	0.027027	0.000000
50	0.040541	0.040541	0.040541	0.027027	0.040541	0.027027	0.013514
55	0.040541	0.040541	0.040541	0.040541	0.027027	0.027027	0.000000
60	0.040541	0.054054	0.027027	0.040541	0.040541	0.027027	0.013514
65	0.040541	0.040541	0.040541	0.027027	0.027027	0.027027	0.013514
70	0.054054	0.040541	0.054054	0.027027	0.040541	0.027027	0.013514
75	0.040541	0.027027	0.054054	0.040541	0.040541	0.027027	0.000000
80	0.040541	0.040541	0.054054	0.040541	0.040541	0.027027	0.000000
85	0.040541	0.040541	0.054054	0.040541	0.027027	0.027027	0.027027
90	0.040541	0.027027	0.054054	0.040541	0.040541	0.027027	0.000000
95	0.027027	0.040541	0.054054	0.027027	0.040541	0.027027	0.000000
100	0.027027	0.040541	0.040541	0.040541	0.027027	0.027027	0.000000

```
[286]: results_df.agg(['idxmax', 'max'], axis=0)
```

```
[286]:
```

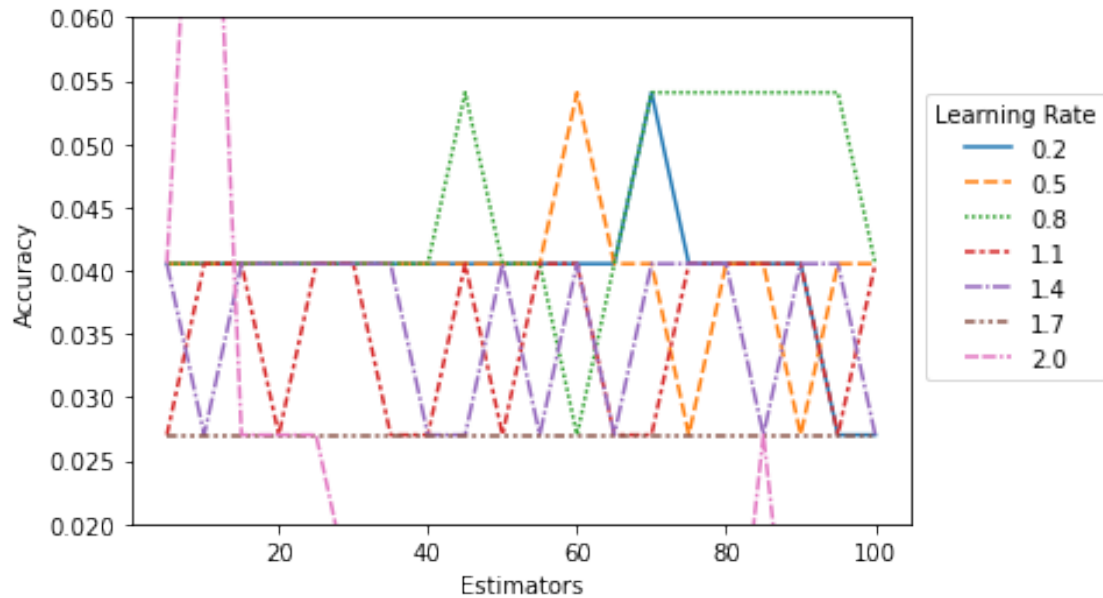
	0.2	0.5	0.8	1.1	1.4	1.7 \
idxmax	70.000000	60.000000	45.000000	10.000000	5.000000	5.000000
max	0.054054	0.054054	0.054054	0.040541	0.040541	0.027027

	2.0
idxmax	10.000000
max	0.094595

```
[287]: sns.lineplot(data=results_df)
plt.ylim((0.02,.06))
plt.xlabel('Estimators')
plt.ylabel('Accuracy')
plt.legend(bbox_to_anchor=(1.02, 0.85), loc='upper left', borderaxespad=0,
           title='Learning Rate')
plt.title('')
```

```
[287]: Text(0.5, 1.0, '')
```



### 3 Grouped by Team

First we need to find the number of people in each team to calculate the average score for each team.

This gets rid of the duplicate users for each group and then the count for any column is the amount of users in that group.

```
[288]: new_df = clean_df
#Dropping NAs and mentors
new_df = new_df.dropna()
new_df = new_df[new_df.RoleName != 'Mentor']
#Adding column for each group
new_df['Group'] = new_df['implementation'] + new_df['group_id'].astype(str)

without_user_dups = new_df.drop_duplicates(subset=['userIDs'])
dupless_group_nums = without_user_dups.groupby(['Group'], as_index = False)

sum_table = dupless_group_nums.sum()
count_table = dupless_group_nums.count()
sum_table['outcome_per_student'] = sum_table['OutcomeScore'] / \
    ↳count_table['userIDs']
sum_table['group_id'] = count_table['group_id']
sum_table[['outcome_per_student']]
```

```
[288]: outcome_per_student
0      3.571429
1      3.666667
2      2.571429
3      2.500000
4      2.166667
..      ...
70     2.200000
71     4.400000
72     4.800000
73     4.600000
74     5.500000
```

```
[75 rows x 1 columns]
```

```
[289]: #Grouping by teams
team_group = new_df.groupby(['Group'], as_index = False)
#Getting the sum of the contributions
team_sum = team_group.sum()
#Adding the outcome per student
team_sum[['Average Score']] = sum_table[['outcome_per_student']]
```

```
[290]: team_sum
```

```
[290]:   Group  Unnamed: 0  userIDs  Line_ID  group_id  m_experimental_testing \
0     a2      55756     1515     55756      598             10
1     a3      79134     1893     79134      510             6
2     a4     200807     4879     200807     1096            25
3     a5     193183     4902     193183      960             2
4     a6     340114     8465     340114     1602             2
..    ...         ...         ...         ...         ...         ...
70    o2     6216397    129357    6217447      700            10
71    o3     5305801    109842    5306680      879             6
72    o4     3676635     75978    3677235      800            11
73    o5     6314090    130116    6315104     1690            16
74    o6     5554645    114095    5555521     1752            16
```

```
   m_making_design_choices  m_asking_questions \
0                        21                    54
1                        16                    25
2                        25                    59
3                        32                    24
4                        34                    40
..                        ...                  ...
70                       31                    41
71                       26                    29
72                       21                    18
```

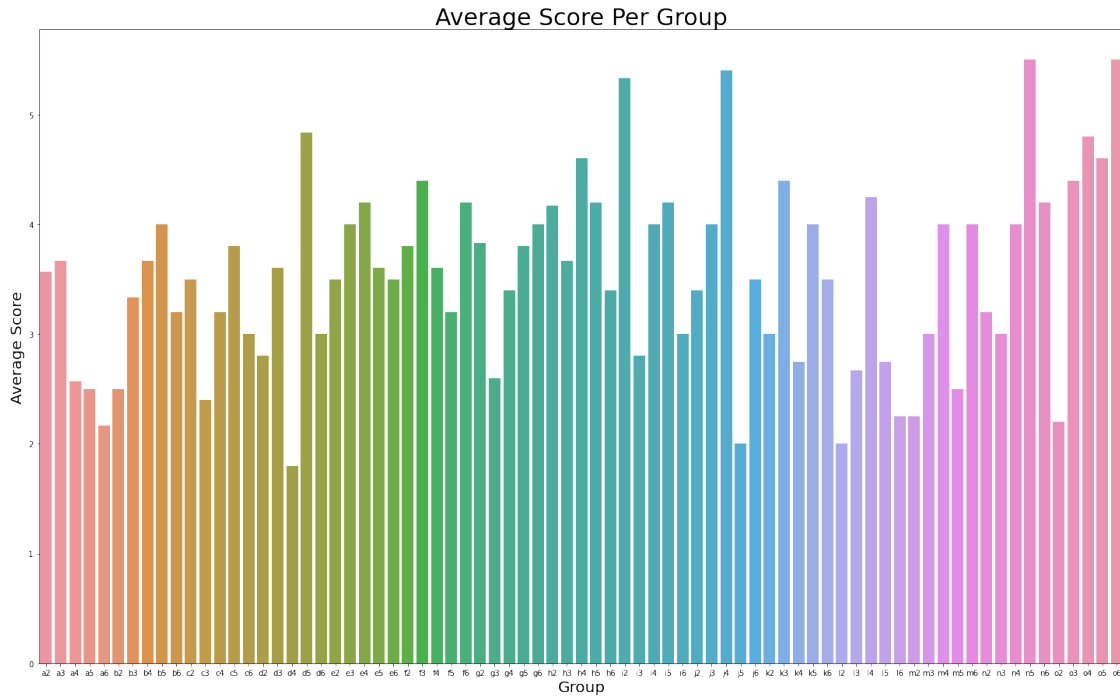
73	36	59
74	28	36

	j_customer_consultants_requests	j_performance_parameters_requirements \
0	2	21
1	3	11
2	2	18
3	6	14
4	1	21
..	...	...
70	6	18
71	4	19
72	6	21
73	6	11
74	5	29

	j_communication	OutcomeScore	wordCount	Average Score
0	1	1076	3007	3.571429
1	2	619	2455	3.666667
2	3	628	3234	2.571429
3	3	515	2470	2.500000
4	4	621	2827	2.166667
..	...	...	...	...
70	6	789	3956	2.200000
71	0	1192	3237	4.400000
72	3	972	2785	4.800000
73	5	1540	4484	4.600000
74	0	1549	3818	5.500000

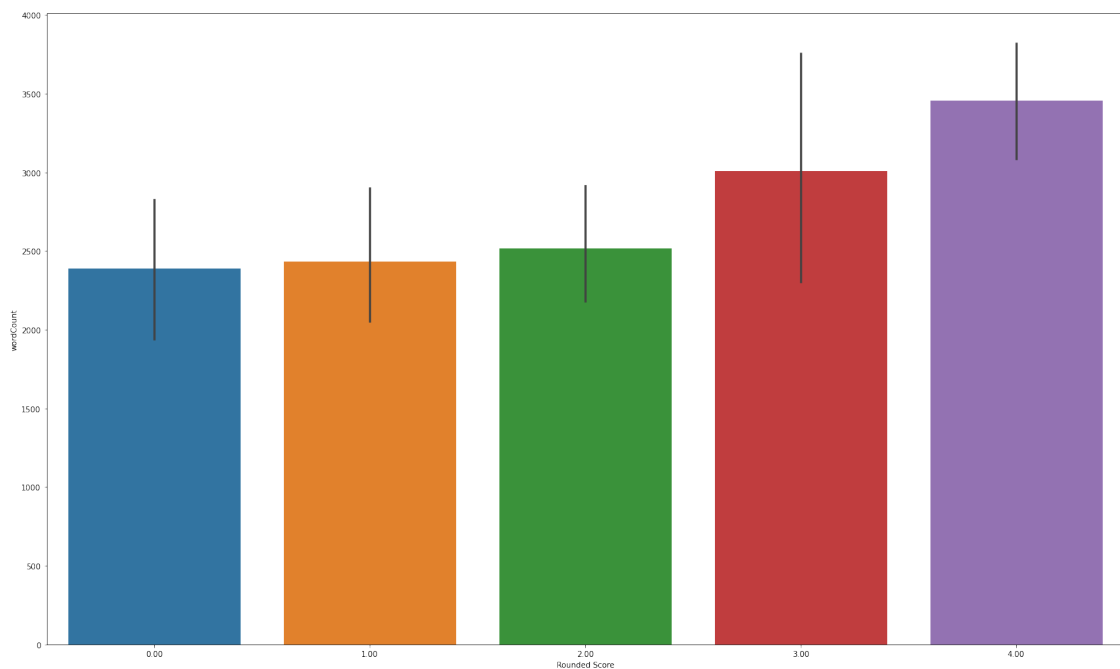
[75 rows x 14 columns]

```
[291]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Group', y = 'Average Score')
plt.title('Average Score Per Group', fontsize = 30)
plt.xlabel('Group', fontsize = 20)
plt.ylabel('Average Score', fontsize = 20);
```

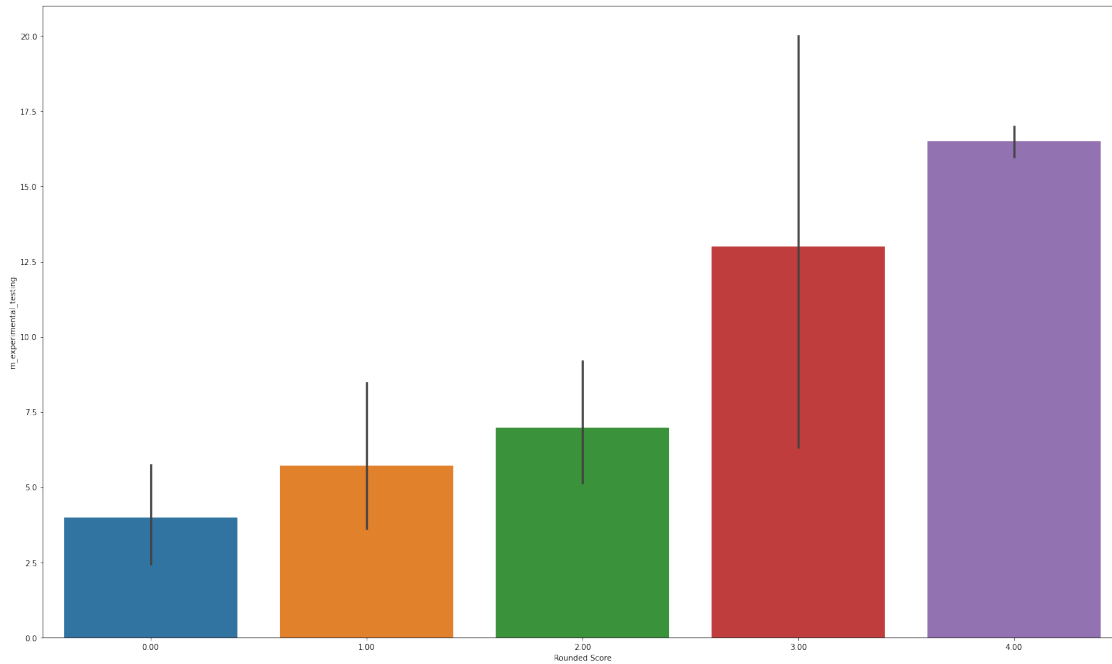


```
[292]: team_sum[['Rounded Score']] = np.round(team_sum[['Average Score']], 0)
```

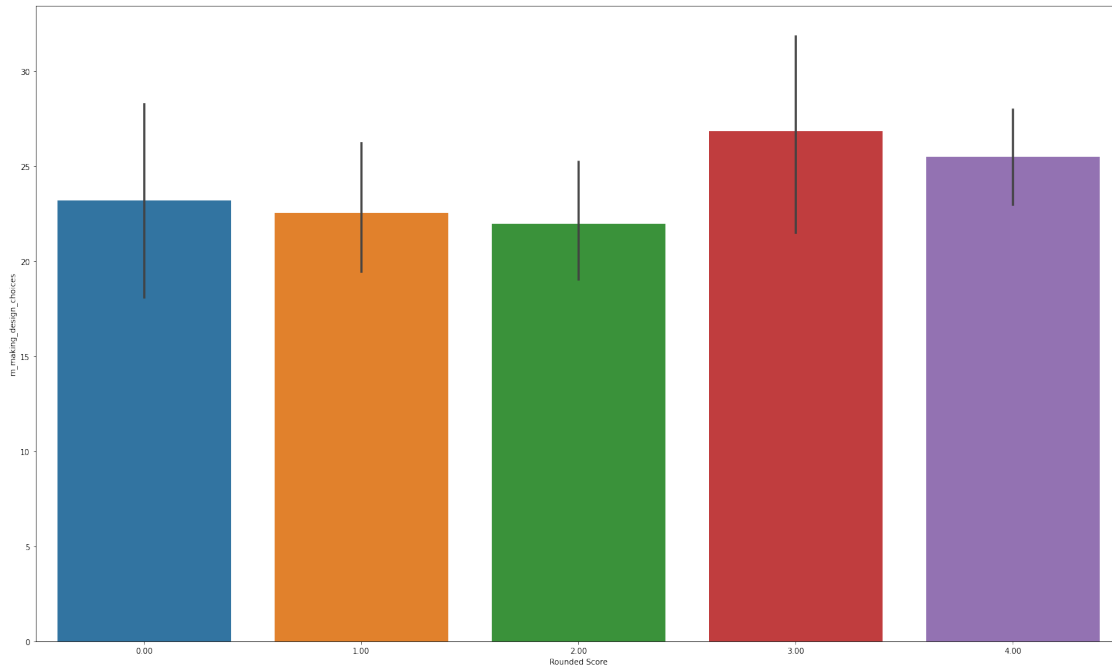
```
[293]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Rounded Score', y = 'wordCount')
ax.xaxis.set_major_formatter(ticker.StrMethodFormatter("{x:.2f}"));
```



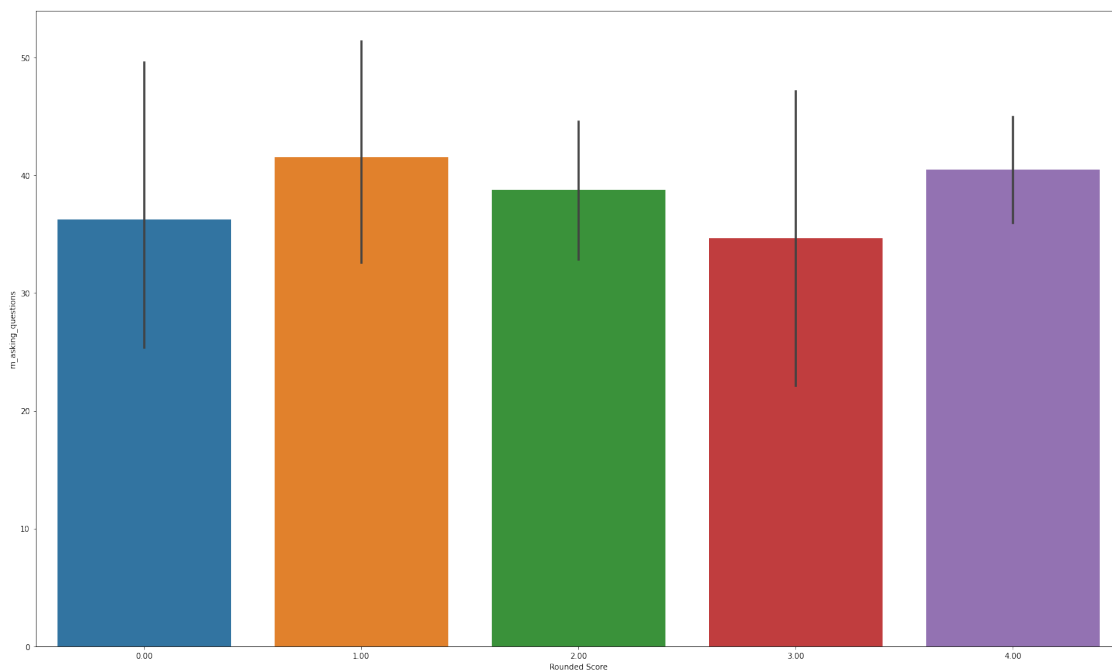
```
[294]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Rounded Score', y = 'm_experimental_testing')
ax.xaxis.set_major_formatter(ticker.StrMethodFormatter("{x:.2f}"));
```



```
[295]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Rounded Score', y = 'm_making_design_choices')
ax.xaxis.set_major_formatter(ticker.StrMethodFormatter("{x:.2f}"));
```

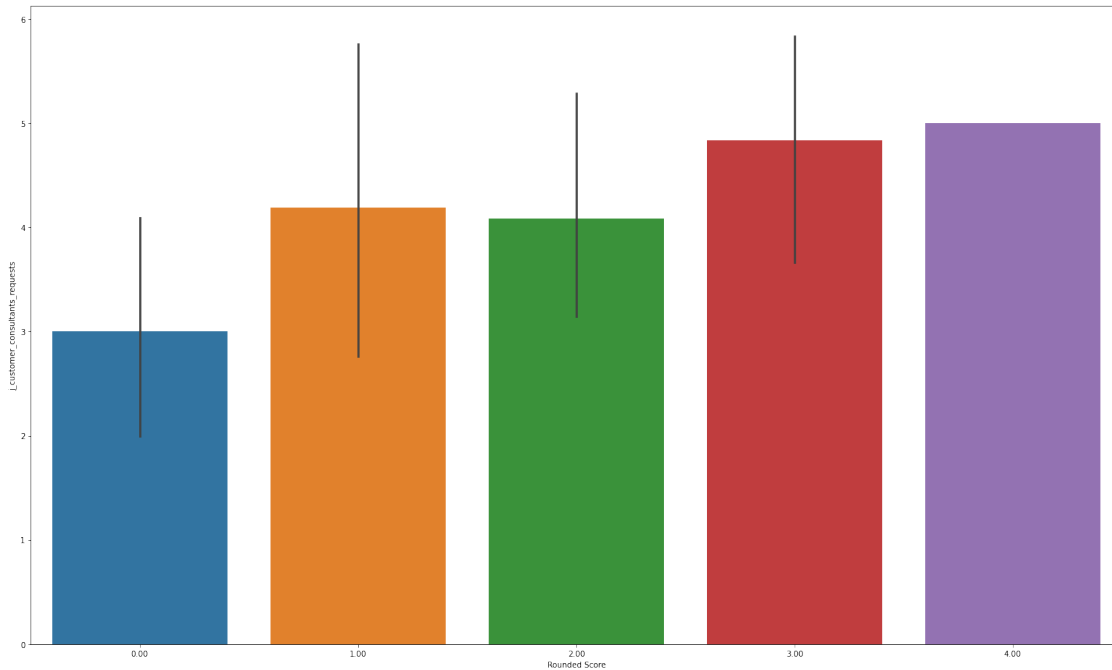


```
[296]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Rounded Score', y = 'm_masking_questions')
ax.xaxis.set_major_formatter(ticker.StrMethodFormatter("{x:.2f}"));
```

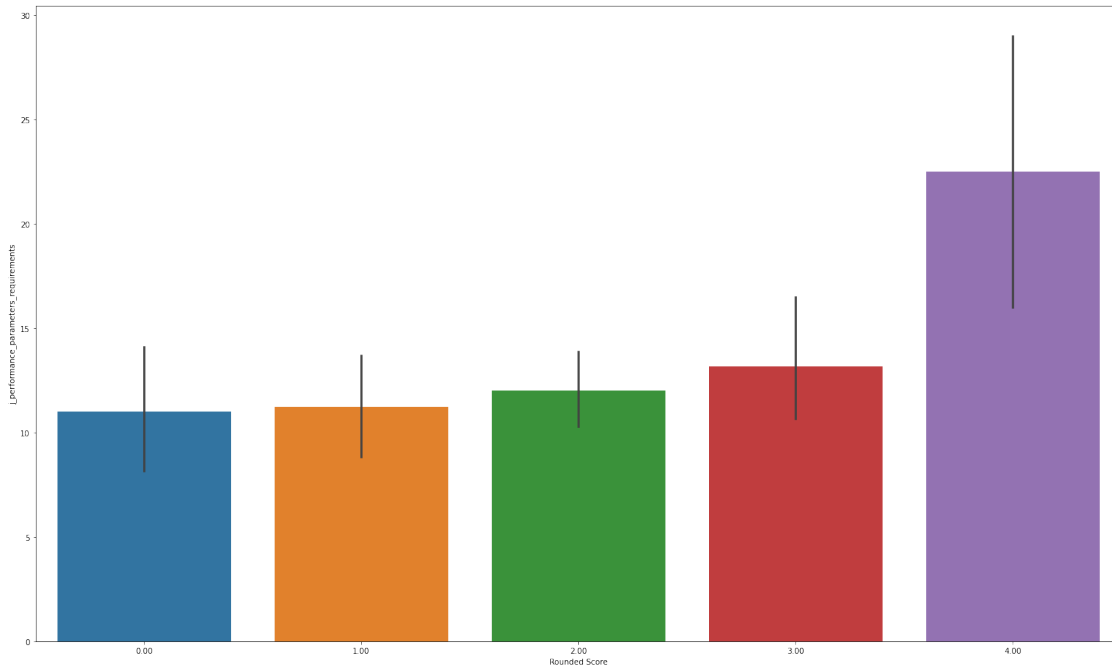




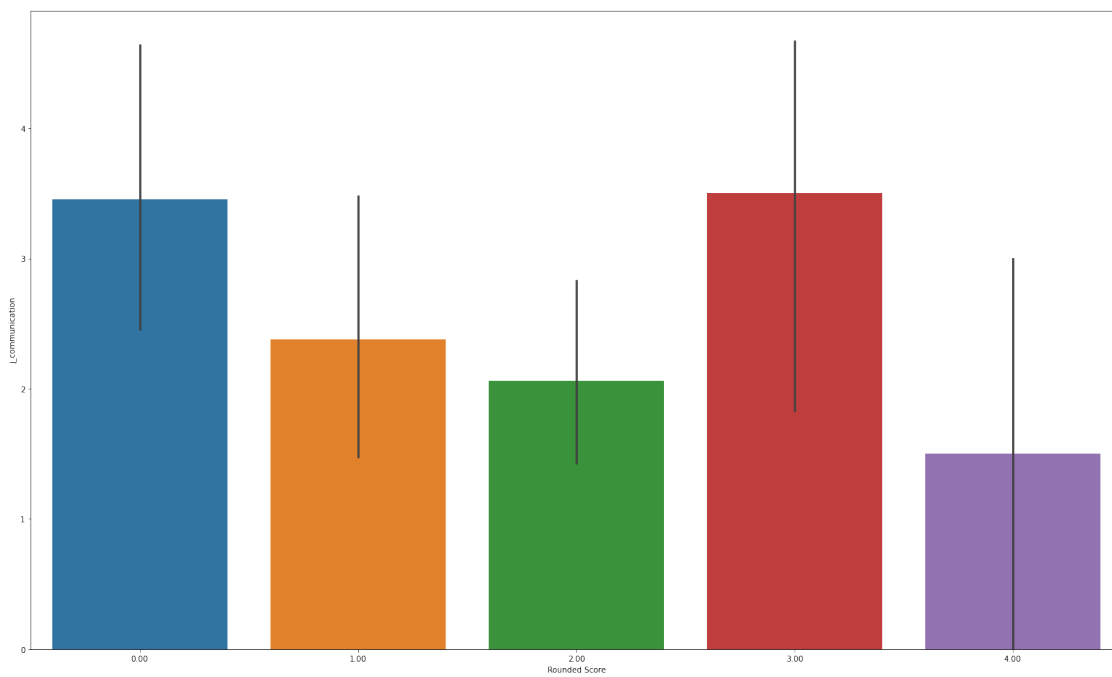
```
[297]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Rounded Score', y =
↳ 'j_customer_consultants_requests')
ax.xaxis.set_major_formatter(ticker.StrMethodFormatter("{x:.2f}"));
```



```
[298]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Rounded Score', y =
↳ 'j_performance_parameters_requirements')
ax.xaxis.set_major_formatter(ticker.StrMethodFormatter("{x:.2f}"));
```

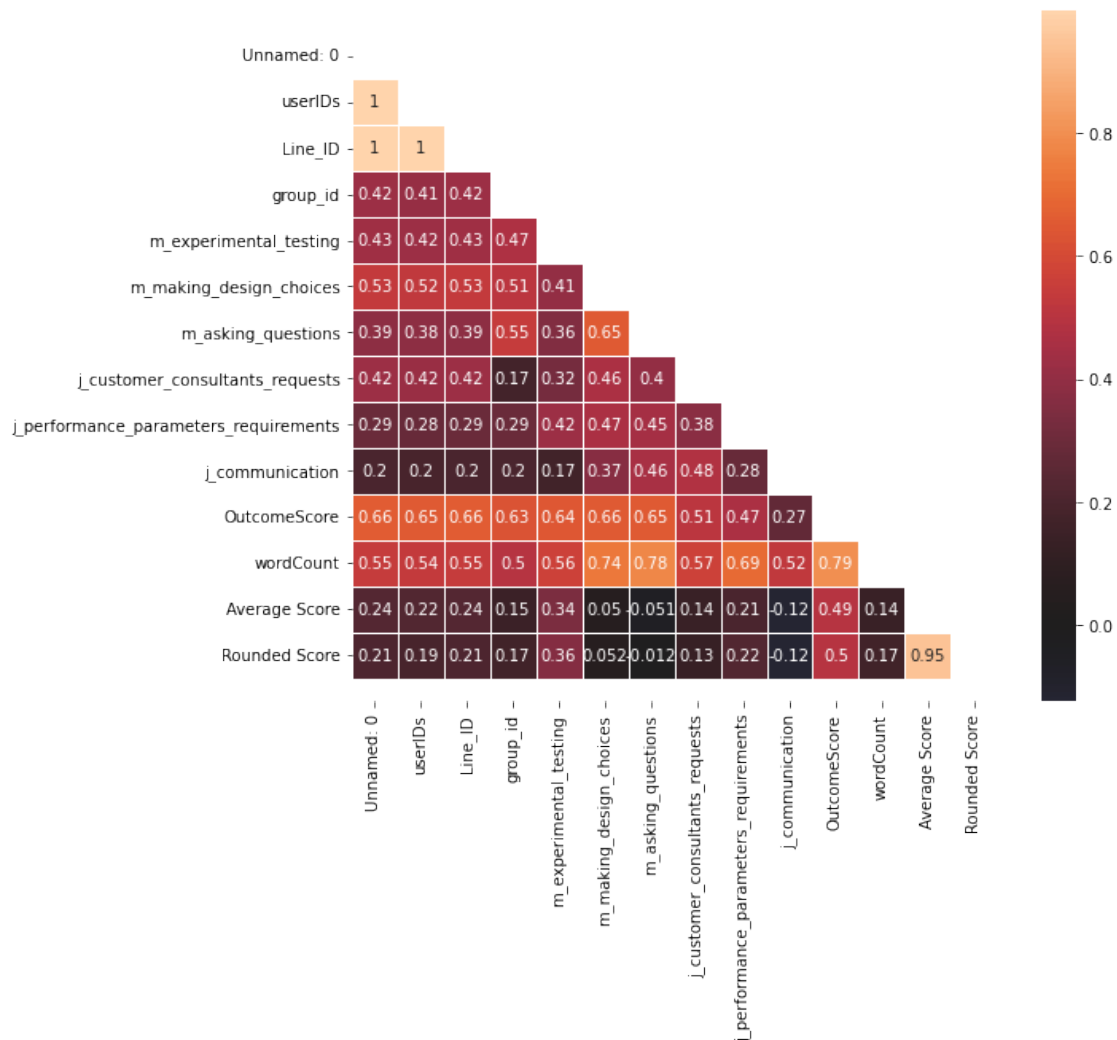


```
[299]: fig, ax = plt.subplots(figsize=(25, 15))
sns.barplot(data = team_sum, x = 'Rounded Score', y = 'j_communication')
ax.xaxis.set_major_formatter(ticker.StrMethodFormatter("{x:.2f}"));
```



It seems that 'wordCount', 'm\_experimental\_testing', 'j\_customer\_consultants\_requests' and 'j\_performance\_parameters\_requirements' have a reasonable correlation with the final average outcome score for each group.

```
[300]: corr = team_sum.corr()
plt_correlation_matrix(corr)
```



## 4 Modelling

### 4.1 Linear Regression

Trying to predict the Average Scores.

```
[301]: team_sum
```

```

[301]:
  Group  Unnamed: 0  userIDs  Line_ID  group_id  m_experimental_testing \
0    a2      55756    1515    55756      598             10
1    a3      79134    1893    79134      510             6
2    a4     200807    4879    200807     1096            25
3    a5     193183    4902    193183      960             2
4    a6     340114    8465    340114     1602             2
..    ...         ...         ...         ...         ...
70   o2     6216397  129357  6217447      700            10
71   o3     5305801  109842  5306680      879             6
72   o4     3676635   75978  3677235      800            11
73   o5     6314090  130116  6315104     1690            16
74   o6     5554645  114095  5555521     1752            16

      m_making_design_choices  m_asking_questions \
0                             21                    54
1                             16                    25
2                             25                    59
3                             32                    24
4                             34                    40
..                             ...                   ...
70                            31                    41
71                            26                    29
72                            21                    18
73                            36                    59
74                            28                    36

      j_customer_consultants_requests  j_performance_parameters_requirements \
0                                     2                                     21
1                                     3                                     11
2                                     2                                     18
3                                     6                                     14
4                                     1                                     21
..                                     ...                                   ...
70                                    6                                     18
71                                    4                                     19
72                                    6                                     21
73                                    6                                     11
74                                    5                                     29

      j_communication  OutcomeScore  wordCount  Average Score  Rounded Score
0                    1           1076       3007      3.571429         4.0
1                    2           619       2455      3.666667         4.0
2                    3           628       3234      2.571429         3.0
3                    3           515       2470      2.500000         2.0
4                    4           621       2827      2.166667         2.0
..                    ...           ...         ...         ...
70                    6           789       3956      2.200000         2.0

```

71	0	1192	3237	4.400000	4.0
72	3	972	2785	4.800000	5.0
73	5	1540	4484	4.600000	5.0
74	0	1549	3818	5.500000	6.0

[75 rows x 15 columns]

```
[302]: XX = team_sum.drop(['Average Score', 'Rounded Score',
    ↳ 'OutcomeScore', 'group_id', 'Line_ID', 'Unnamed: 0', 'userIDs', 'Group'], axis =
    ↳ 1)
Y = team_sum['Average Score']

#Normalise for regularisation later
X = (XX-XX.mean())/XX.std()

feature_names = X.columns

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,
    ↳ random_state=42)
```

```
[303]: linear = LinearRegression() # instantatiatate the linear regression model
linear.fit(X_train, Y_train) # fit the data to the model
training_score = linear.score(X_train, Y_train) # calculate rsq for the training
    ↳ set

# use the independent variables for the testing set to predict the target
    ↳ variable
preds_linear = linear.predict(X_test)

# calculate the correlation of the predicted and actual target variables
rsquared_linear = r2_score(Y_test, preds_linear)

# print the training and testing scores
print(f'Training score is {training_score:.3f}')
print(f'Testing score is {rsquared_linear:.3f}')
```

Training score is 0.230

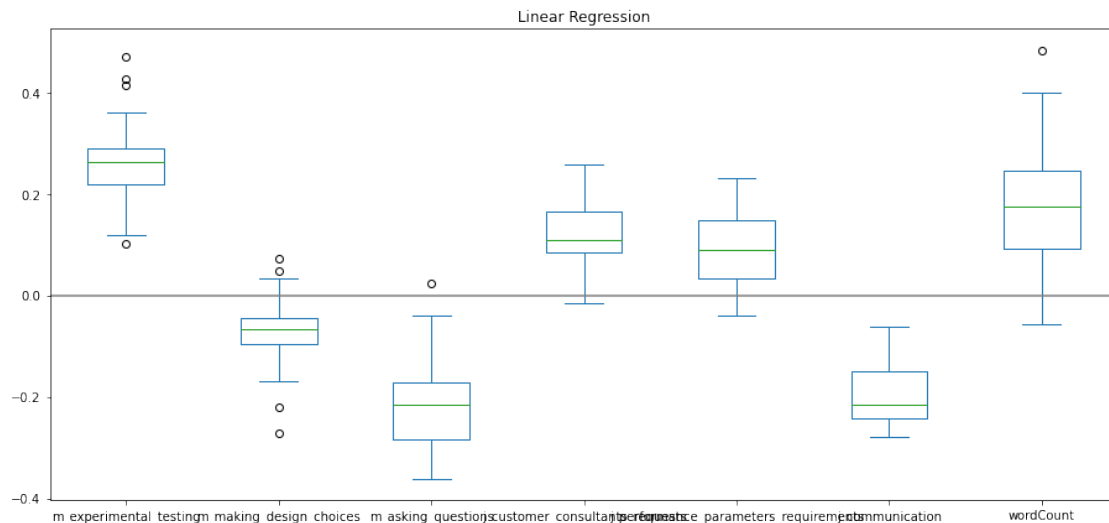
Testing score is -0.102

```
[304]: scores = cross_validate(
    linear, X, Y, cv=RepeatedKfold(n_splits=5, n_repeats=5, random_state=42),
    return_estimator=True
)
# take the results for each simulation (estimator), extract the coefficients for
    ↳ each run
# and add them to a dataframe with columns being the feature names
```

```

coefs = pd.DataFrame([est.coef_ for est in scores['estimator']], columns =
    →feature_names)
# plot the descriptive statics of the coefficients in a box and whisker plot to
    →show variability
ax = coefs.plot(kind='box',figsize=(20, 7))
plt.title('Linear Regression')
plt.axhline(y=0, color='.5')
plt.subplots_adjust(left=.3)

```



## 4.2 Decision Trees

### 4.2.1 Regressor

First using decision tree regressor to predict average group score.

```

[305]: X = team_sum.drop(['Average Score', 'Rounded Score',
    →'OutcomeScore', 'group_id', 'Line_ID', 'Unnamed: 0', 'userIDs', 'Group'], axis =
    →1)
Y = team_sum['Average Score']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size = 0.8,
    →random_state = 42)

```

```

[306]: df_dtr = DecisionTreeRegressor(random_state = 42, max_depth = 3)
df_dtr = df_dtr.fit(X_train, Y_train)

```

```

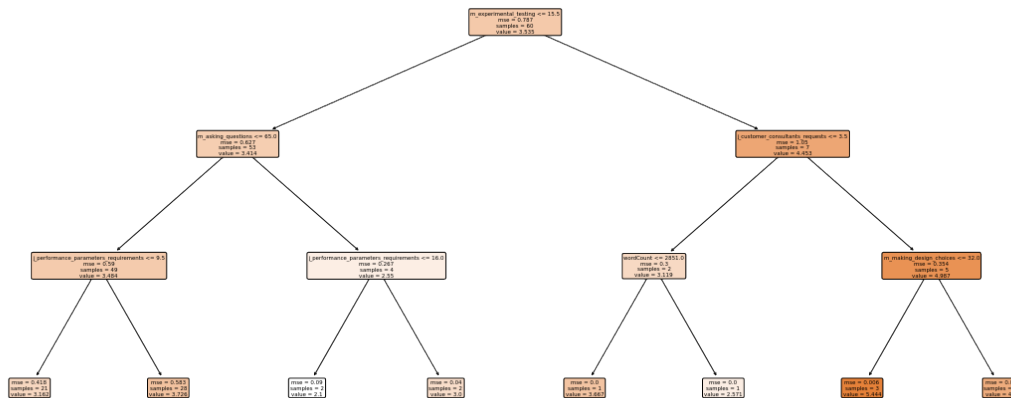
[307]: fig, ax = plt.subplots(figsize=(20,9)) # initialise the plots and axes
# plot the decision tree for the model df_dtr
plot_tree(df_dtr,
    filled = True, # colour the nodes according to the classification

```

```

rounded = True, # make the nodes have rounded corners
class_names = ['1', '2', '3', '4', '5', '6', '7', '8'], # use these names for targets
feature_names = X.columns # use these names for features
);

```



```

[308]: for i in range(1, 14):
    df_dtr = DecisionTreeRegressor(random_state = 42, max_depth = i)
    df_dtr = df_dtr.fit(X_train, Y_train)

    depth = df_dtr.tree_.max_depth

    Y_pred = df_dtr.predict(X_test)
    r2 = r2_score(Y_test, Y_pred)

    print(f"Depth of {depth} gives r2 score: {np.round(r2, 3)}")

```

```

Depth of 1 gives r2 score: 0.244
Depth of 2 gives r2 score: -0.146
Depth of 3 gives r2 score: -0.498
Depth of 4 gives r2 score: -0.333
Depth of 5 gives r2 score: -0.572
Depth of 6 gives r2 score: -0.764
Depth of 7 gives r2 score: -0.639
Depth of 8 gives r2 score: -0.717
Depth of 9 gives r2 score: -0.469
Depth of 10 gives r2 score: -0.4
Depth of 11 gives r2 score: -1.085
Depth of 12 gives r2 score: -0.95

```

Depth of 13 gives r2 score: -0.993

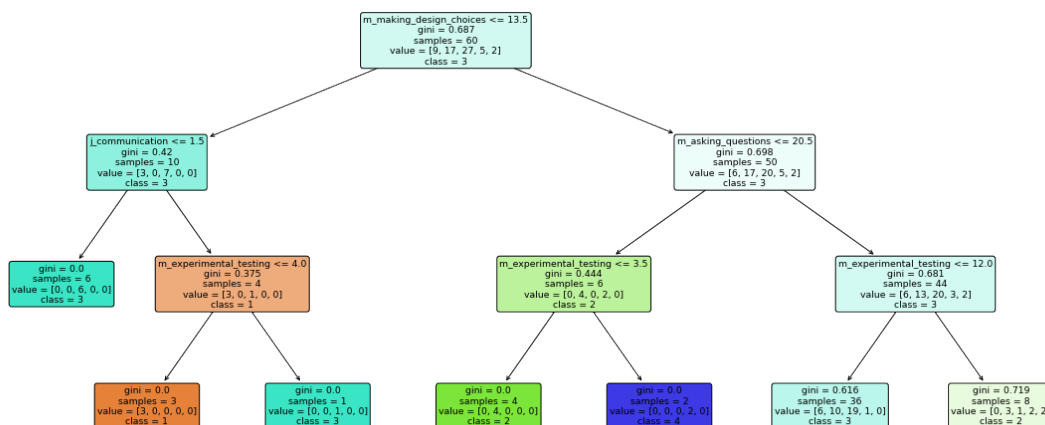
At depth of 13 it becomes specialised only to this dataset.

## 4.2.2 Classifier

Attempting to predict the rounded versions of the average score per group as decision tree classifier needs to be distinct.

```
[309]: X = team_sum.drop(['Average Score', 'Rounded Score',  
    ↳ 'OutcomeScore', 'group_id', 'Line_ID', 'Unnamed: 0', 'userIDs', 'Group'], axis =  
    ↳ 1)  
Y = team_sum['Rounded Score']  
  
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size = 0.8,  
    ↳ random_state = 42)
```

```
[310]: df_dtc = DecisionTreeClassifier(random_state = 42, max_depth = 3)  
df_dtc = df_dtc.fit(X_train, Y_train)  
  
fig, ax = plt.subplots(figsize=(20,9)) # initialise the plots and axes  
# plot the decision tree for the model df_dtc  
plot_tree(df_dtc,  
    filled = True, # colour the nodes according to the classification  
    rounded = True, # make the nodes have rounded corners  
    class_names = ['1', '2', '3', '4', '5', '6', '7', '8'], # use these names for  
    ↳ targets  
    feature_names = X.columns # use these names for features  
);
```





```
[311]: for i in range(1, 21):
        df_dtc = DecisionTreeClassifier(random_state = 42, max_depth = i)
        df_dtc = df_dtc.fit(X_train, Y_train)

        depth = df_dtc.tree_.max_depth

        Y_pred = df_dtc.predict(X_test)
        accuracy = np.round(accuracy_score(Y_test, Y_pred), 3)
        print(f"Depth of {depth} gives accuracy: {accuracy}")
```

```
Depth of 1 gives accuracy: 0.533
Depth of 2 gives accuracy: 0.4
Depth of 3 gives accuracy: 0.333
Depth of 4 gives accuracy: 0.2
Depth of 5 gives accuracy: 0.133
Depth of 6 gives accuracy: 0.267
Depth of 7 gives accuracy: 0.333
Depth of 8 gives accuracy: 0.333
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
Depth of 9 gives accuracy: 0.4
```

## 4.3 Random Forests

### 4.3.1 Regressor

```
[312]: X = team_sum.drop(['Average Score', 'Rounded Score',
    ↳ 'OutcomeScore', 'group_id', 'Line_ID', 'Unnamed: 0', 'userIDs', 'Group'], axis =
    ↳ 1)
        Y = team_sum['Rounded Score']

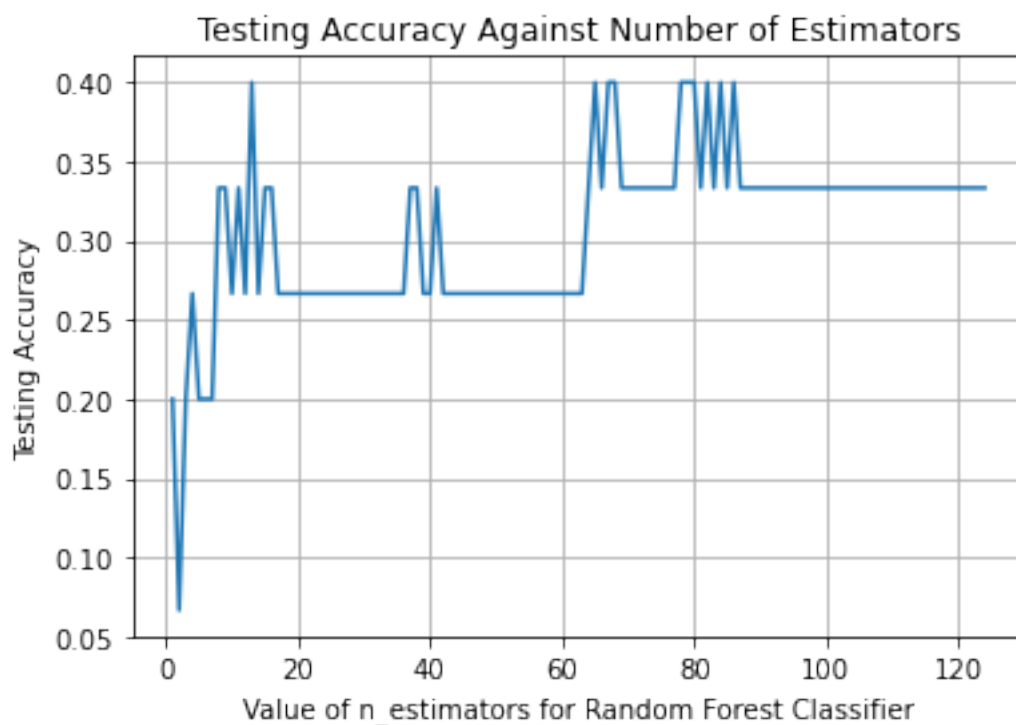
        X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size = 0.8,
    ↳ random_state = 42)
```

### 4.3.2 Finding Best Parameters

For Classifier

```
[313]: scores = []
for i in range(1, 125):
    rfc = RandomForestClassifier(n_estimators = i, random_state = 42)
    rfc.fit(X_train, Y_train)
    Y_pred = rfc.predict(X_test)
    scores.append(accuracy_score(Y_test, Y_pred))

plt.plot(range(1, 125), scores)
plt.xlabel('Value of n_estimators for Random Forest Classifier')
plt.ylabel('Testing Accuracy')
plt.title('Testing Accuracy Against Number of Estimators')
plt.grid(True);
```



```
[314]: rfc = RandomForestClassifier(n_estimators = 125, random_state = 42).fit(X_train, Y_train)
Y_pred = rfc.predict(X_test)
accuracy = np.round(accuracy_score(Y_test, Y_pred), 3)

print(f"Accuracy: {accuracy}")
```

Accuracy: 0.333

## 4.4 K Nearest Neighbours

```
[315]: user_sum['NewOutcomeScore'] = user_sum['OutcomeScore'] / user_group.
      ↪count()['OutcomeScore']
user_data = user_sum
user_data
```

```
[315]:
```

	userIDs	Unnamed: 0	Line_ID	group_id	m_experimental_testing	\
0	2	12286	12286	130	2	
1	3	2796	2796	44	1	
2	4	6454	6454	62	2	
3	5	9385	9385	90	0	
4	6	6740	6740	74	0	
..	...	...	...	...	...	
364	389	874511	874649	276	2	
365	390	913147	913291	288	0	
366	391	970395	970548	306	2	
367	392	684840	684948	216	2	
368	393	1369484	1369700	432	5	

	m_making_design_choices	m_asking_questions	\
0	4	17	
1	4	3	
2	2	3	
3	0	6	
4	2	7	
..	...	...	
364	8	4	
365	5	6	
366	3	11	
367	4	4	
368	6	8	

	j_customer_consultants_requests	j_performance_parameters_requirements	\
0	0	4	
1	0	1	
2	1	5	
3	0	2	
4	1	3	
..	...	...	
364	1	7	
365	1	2	
366	1	4	
367	0	2	
368	1	7	

	j_communication	OutcomeScore	wordCount	NewOutcomeScore
--	-----------------	--------------	-----------	-----------------

0	0	260	704	4.0
1	0	88	157	4.0
2	0	124	349	4.0
3	0	90	342	2.0
4	0	74	399	2.0
..	...	...	...	...
364	0	322	755	7.0
365	0	192	431	4.0
366	0	255	605	5.0
367	0	180	451	5.0
368	0	288	868	4.0

[369 rows x 13 columns]

```
[316]: XX = user_data.drop(['NewOutcomeScore', 'OutcomeScore', 'group_id', 'Line_ID',
    ↳ 'Unnamed: 0', 'userIDs'], axis = 1)
Y = user_data['NewOutcomeScore']

#Normalise for regularisation later
X = (XX-XX.mean())/XX.std()

feature_names = X.columns

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,
    ↳ random_state=42)
```

```
[317]: # normalize the data
nX =(X-X.mean())/X.std() # create nX, a normalised version of X
nX.describe() # show the descriptive statistics of nX
```

```
[317]:      m_experimental_testing  m_making_design_choices  m_asking_questions  \
count      3.690000e+02      3.690000e+02      3.690000e+02
mean      -8.544806e-17      2.542381e-17     -1.513111e-16
std       1.000000e+00      1.000000e+00      1.000000e+00
min      -7.088055e-01     -1.285083e+00     -1.131795e+00
25%      -7.088055e-01     -7.301439e-01     -7.021324e-01
50%      -2.048571e-01     -1.752045e-01     -2.724693e-01
75%       2.990913e-01      3.797350e-01      4.436359e-01
max       6.850420e+00      5.651660e+00      6.888582e+00

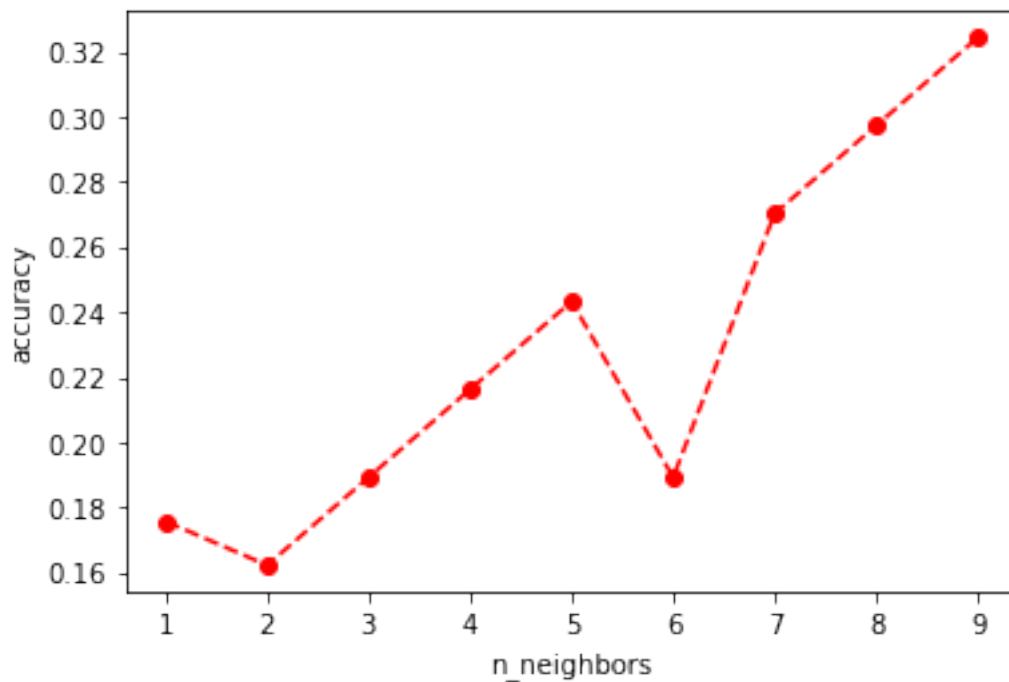
      j_customer_consultants_requests  j_performance_parameters_requirements  \
count      3.690000e+02      3.690000e+02
mean       2.226464e-16     -1.949660e-16
std       1.000000e+00      1.000000e+00
min      -6.893246e-01     -1.166287e+00
25%      -6.893246e-01     -6.886402e-01
50%      -6.893246e-01     -2.109931e-01
```

75%	1.501499e-01	2.666539e-01
max	6.865946e+00	5.520771e+00

	j_communication	wordCount
count	3.690000e+02	3.690000e+02
mean	4.392752e-17	1.022970e-17
std	1.000000e+00	1.000000e+00
min	-5.378988e-01	-1.541381e+00
25%	-5.378988e-01	-6.614284e-01
50%	-5.378988e-01	-1.526124e-01
75%	5.408221e-01	4.699390e-01
max	5.934427e+00	6.456010e+00

```
[318]: # Check Accuracy
Accuracy = []
for i in range(1,10):
    knn = KNeighborsClassifier(n_neighbors = i).fit(X_train, Y_train)
    Y_pred = knn.predict(X_test)
    Accuracy.append(metrics.accuracy_score(Y_test, Y_pred))
plt.plot(range(1,10), Accuracy, color = 'red', linestyle = 'dashed', marker = 'o')
plt.xlabel('n_neighbors')
plt.ylabel('accuracy')
```

```
[318]: Text(0, 0.5, 'accuracy')
```



```
[319]: knn = KNeighborsClassifier(n_neighbors = 9).fit(X_train, Y_train)
Y_pred = knn.predict(X_test)
metrics.accuracy_score(Y_test, Y_pred)
```

```
[319]: 0.32432432432432434
```

## 5 Importance of Certain Words

```
[320]: new_df = clean_df
#Dropping NAs and users
new_df = new_df.dropna()
new_df = new_df[new_df.RoleName == 'Mentor']
#Adding column for each group
new_df['Group'] = new_df['group_id'].astype(str) + new_df['implementation']
```

```
[321]: without_dups = new_df.drop_duplicates(subset=['userIDs'])
group_nums = without_dups.sort_values(by = ['Group'])
#group_nums = without_dups.groupby(['Group'], as_index = False)
```

```
[322]: group_nums
```

```
[322]: Empty DataFrame
Columns: [Unnamed: 0, content, RoleName, userIDs, implementation, Line_ID,
ChatGroup, group_id, roomName, m_experimental_testing, m_making_design_choices,
m_masking_questions, j_customer_consultants_requests,
j_performance_parameters_requirements, j_communication, OutcomeScore, wordCount,
Group]
Index: []
```

```
[323]: without_dups
```

```
[323]: Empty DataFrame
Columns: [Unnamed: 0, content, RoleName, userIDs, implementation, Line_ID,
ChatGroup, group_id, roomName, m_experimental_testing, m_making_design_choices,
m_masking_questions, j_customer_consultants_requests,
j_performance_parameters_requirements, j_communication, OutcomeScore, wordCount,
Group]
Index: []
```

## 6 Finding The Sum of Certain Words

### 6.1 By Team

#### 6.1.1 Sorting

First we need to find the number of people in each team to calculate the average score for each team.

This gets rid of the duplicate users for each group and then the count for any column is the amount of users in that group.

```
[324]: new_df = clean_df
        #Dropping NAs and mentors
        new_df = new_df.dropna()
        new_df = new_df[new_df.RoleName != 'Mentor']
        #Adding column for each group
        new_df['Group'] = new_df['implementation'] + new_df['group_id'].astype(str)

        without_user_dups = new_df.drop_duplicates(subset=['userIDs'])
        dupless_group_nums = without_user_dups.groupby(['Group'], as_index = False)

        sum_table = dupless_group_nums.sum()
        count_table = dupless_group_nums.count()
        sum_table['outcome_per_student'] = sum_table['OutcomeScore'] / \
        ↪count_table['userIDs']
        sum_table['group_id'] = count_table['group_id']
        sum_table[['outcome_per_student']]
```

```
[324]:      outcome_per_student
0          3.571429
1          3.666667
2          2.571429
3          2.500000
4          2.166667
..          ...
70         2.200000
71         4.400000
72         4.800000
73         4.600000
74         5.500000
```

[75 rows x 1 columns]

```
[325]: #Grouping by teams
        team_group = new_df.groupby(['Group'], as_index = False)
        #Getting the sum of the contributions
        team_sum = team_group.sum()
        #Adding the outcome per student
```

```
team_sum[['Average Score']] = sum_table[['outcome_per_student']]
```

```
[326]: team_sum
```

```
[326]:
```

	Group	Unnamed: 0	userIDs	Line_ID	group_id	m_experimental_testing \
0	a2	55756	1515	55756	598	10
1	a3	79134	1893	79134	510	6
2	a4	200807	4879	200807	1096	25
3	a5	193183	4902	193183	960	2
4	a6	340114	8465	340114	1602	2
..	...	...	...	...	...	...
70	o2	6216397	129357	6217447	700	10
71	o3	5305801	109842	5306680	879	6
72	o4	3676635	75978	3677235	800	11
73	o5	6314090	130116	6315104	1690	16
74	o6	5554645	114095	5555521	1752	16

	m_making_design_choices	m_asking_questions \
0	21	54
1	16	25
2	25	59
3	32	24
4	34	40
..	...	...
70	31	41
71	26	29
72	21	18
73	36	59
74	28	36

	j_customer_consultants_requests	j_performance_parameters_requirements \
0	2	21
1	3	11
2	2	18
3	6	14
4	1	21
..	...	...
70	6	18
71	4	19
72	6	21
73	6	11
74	5	29

	j_communication	OutcomeScore	wordCount	Average Score
0	1	1076	3007	3.571429
1	2	619	2455	3.666667
2	3	628	3234	2.571429



3	3	515	2470	2.500000
4	4	621	2827	2.166667
..	...	...	...	...
70	6	789	3956	2.200000
71	0	1192	3237	4.400000
72	3	972	2785	4.800000
73	5	1540	4484	4.600000
74	0	1549	3818	5.500000

[75 rows x 14 columns]

## 6.2 By User

### 6.2.1 Sorting

We need to find the amount of times certain words were said. We will be looking at some of the more common words found using nlp: 'surfactant' and 'steric.'

```
[327]: new_df = clean_df
#Dropping NAs and mentors
new_df = new_df.dropna()
new_df = new_df[new_df.RoleName != 'Mentor']

#cleaning text/language data
new_df['content'] = new_df['content'].str.lower()
new_df['content'] = new_df['content'].str.translate(str.maketrans('', '', string.
    ↳punctuation))
#Adding column for each group
new_df['Group'] = new_df['implementation'] + new_df['group_id'].astype(str)
```

```
[328]: new_df
```

```
[328]:      Unnamed: 0      content RoleName  userIDs implementation \
5           6  hello i am brandon  Player         2             a
6           7      i am zelin  Player         3             a
7           8           hi  Player         3             a
8           9      i am jack  Player         4             a
9          10  hey im rachel  Player         5             a
...         ...           ...           ...           ...           ...
19170      19174      exactly  Player        391             o
19171      19175  sounds good  Player        389             o
19172      19176          yes  Player        392             o
19173      19177  sounds good  Player        388             o
19175      19179  precisely  Player        393             o

      Line_ID ChatGroup  group_id \
5           6    PRNLT          2
6           7    PRNLT          2
```

7	8	PRNLT	2
8	9	PRNLT	2
9	10	PRNLT	2
...	...	...	...
19170	19177	PESPVP	6
19171	19178	PESPVP	6
19172	19179	PESPVP	6
19173	19180	PESPVP	6
19175	19182	PESPVP	6

	roomName \
5	Introduction and Workflow Tutorial with Entran...
6	Introduction and Workflow Tutorial with Entran...
7	Introduction and Workflow Tutorial with Entran...
8	Introduction and Workflow Tutorial with Entran...
9	Introduction and Workflow Tutorial with Entran...
...	...
19170	Reflection team discussion of first batch results
19171	Reflection team discussion of first batch results
19172	Reflection team discussion of first batch results
19173	Reflection team discussion of first batch results
19175	Reflection team discussion of first batch results

	m_experimental_testing	m_making_design_choices	m_asking_questions \
5	0	0	0
6	0	0	0
7	0	0	0
8	0	0	0
9	0	0	0
...	...	...	...
19170	0	0	0
19171	0	0	0
19172	0	0	0
19173	0	0	0
19175	0	0	0

	j_customer_consultants_requests	j_performance_parameters_requirements \
5	0	0
6	0	0
7	0	0
8	0	0
9	0	0
...	...	...
19170	0	0
19171	0	0
19172	0	0
19173	0	0

19175

0

0

	j_communication	OutcomeScore	wordCount	Group
5	0	4	4	a2
6	0	4	3	a2
7	0	4	1	a2
8	0	4	3	a2
9	0	2	3	a2
...	...	...	...	...
19170	0	5	1	o6
19171	0	7	2	o6
19172	0	5	1	o6
19173	0	8	2	o6
19175	0	4	1	o6

[16902 rows x 18 columns]

We can iterate through 16902 rows.

```
[329]: #print(new_df[['content']].iloc[0])

for i in new_df[['content']].iloc[0]:
    print(i)
```

hello i am brandon

```
[330]: from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
doc_corpus = []
for i in new_df.index:
    doc_corpus.append(new_df['content'].loc[i])

# print(doc_corpus)

vec=TfidfVectorizer(stop_words='english', max_features= 25)
matrix=vec.fit_transform(doc_corpus)
# print("Feature Names n",vec.get_feature_names_out())
# print("Sparse Matrix n",matrix.shape,"n",matrix.toarray())

# print(matrix)
```

```
[331]: word_count = []
tfidf_terms = ['surfactant','steric','hindering','cnt']

for i in range(len(new_df)):
    word_sum = 0
    for sentence in new_df[['content']].iloc[i]:
        split_sentence = sentence.split()
        for word in split_sentence:
```

```

        if word in tfidf_terms:
            word_sum += 1
    word_count.append(word_sum)

new_df['input'] = word_count

```

```

[332]: tfidf_terms = ['marketability', 'reliability', 'cost', 'flux']
word_count = []
for i in range(len(new_df)):
    word_sum = 0
    for sentence in new_df[['content']].iloc[i]:
        split_sentence = sentence.split()
        for word in split_sentence:
            if word in tfidf_terms:
                word_sum += 1
    word_count.append(word_sum)

new_df['output'] = word_count

```

Now we group by user and then we can see how many times each student said a certain word.

```

[333]: #Grouping by users
user_group = new_df.groupby(['userIDs'], as_index = False)
#Getting the sum of the contributions
user_sum = user_group.sum()

```

```

[334]: user_sum['NewOutcomeScore'] = user_sum['OutcomeScore'] / user_group.
        ↪count()['OutcomeScore']
user_data = user_sum
user_data

```

```

[334]:
   userIDs  Unnamed: 0  Line_ID  group_id  m_experimental_testing  \
0         2      12286    12286      130                2
1         3       2796     2796       44                1
2         4       6454     6454       62                2
3         5       9385     9385       90                0
4         6       6740     6740       74                0
...      ...        ...      ...      ...                ...
364       389     874511   874649      276                2
365       390     913147   913291      288                0
366       391     970395   970548      306                2
367       392     684840   684948      216                2
368       393    1369484  1369700      432                5

   m_making_design_choices  m_masking_questions  \
0                          4                   17
1                          4                    3
2                          2                    3

```

```

3           0           6
4           2           7
..         ...         ...
364         8           4
365         5           6
366         3          11
367         4           4
368         6           8

      j_customer_consultants_requests  j_performance_parameters_requirements \
0                                     0                                     4
1                                     0                                     1
2                                     1                                     5
3                                     0                                     2
4                                     1                                     3
..                                   ...                                   ...
364                                   1                                     7
365                                   1                                     2
366                                   1                                     4
367                                   0                                     2
368                                   1                                     7

      j_communication  OutcomeScore  wordCount  input  output  NewOutcomeScore
0                   0             260       704     4     11             4.0
1                   0              88       157     2      9             4.0
2                   0             124       349     6     14             4.0
3                   0              90       342     3      7             2.0
4                   0              74       399     6     16             2.0
..                 ...             ...       ...     ...     ...             ...
364                  0             322       755    33      9             7.0
365                  0             192       431     6      6             4.0
366                  0             255       605    19     13             5.0
367                  0             180       451     9      6             5.0
368                  0             288       868    27     11             4.0

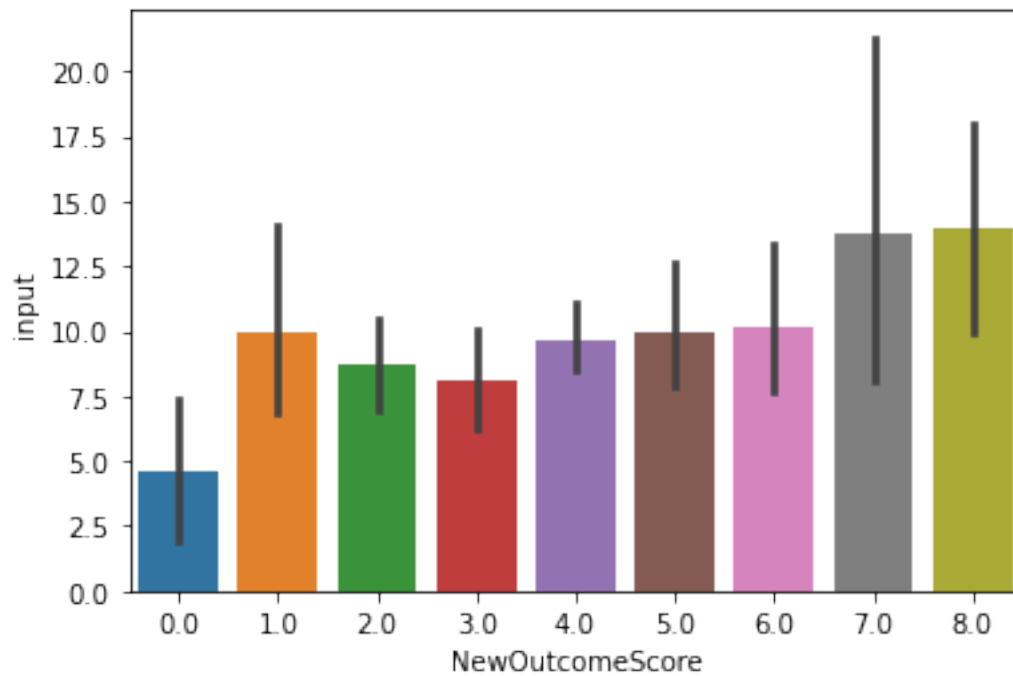
```

```
[369 rows x 15 columns]
```

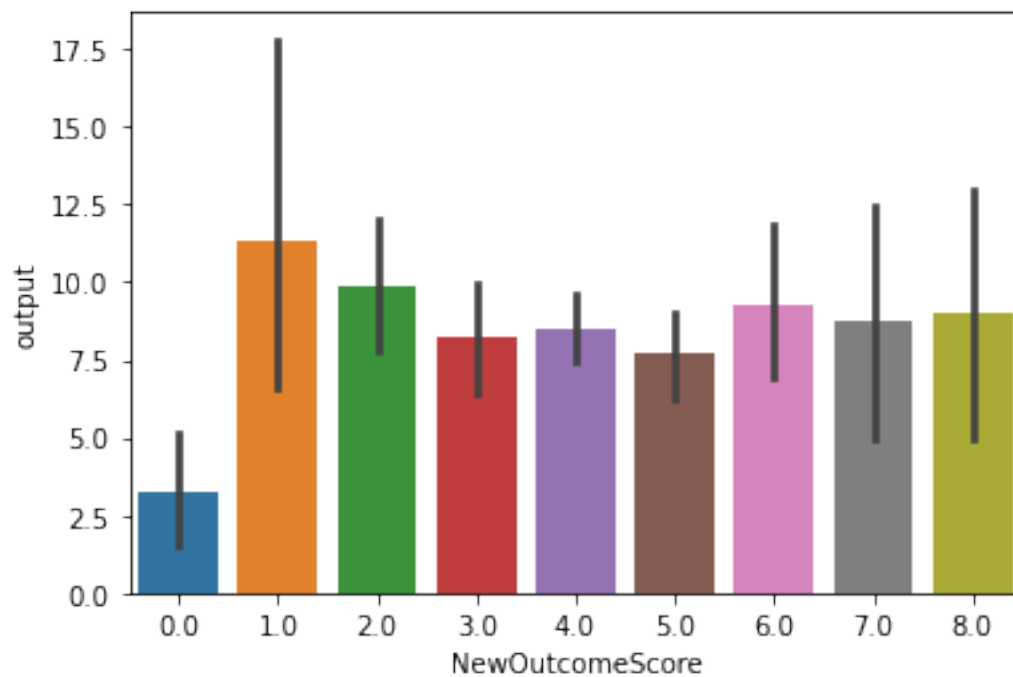
## 6.2.2 Plotting

We can plot the amount of times input and output terms were said on average per each outcome score and there seems to be a slight correlation with the higher marks.

```
[335]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = "input");
```



```
[336]: sns.barplot(data = user_data, x = 'NewOutcomeScore', y = "output");
```



Using just new features

```
[337]: X = user_data[['input', 'output']]
Y = user_data['NewOutcomeScore']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,
→random_state=42)

rfc = RandomForestClassifier(n_estimators = 20, max_samples = 0.5, random_state=
→42).fit(X_train, Y_train)
Y_pred = rfc.predict(X_test)
accuracy = np.round(accuracy_score(Y_test, Y_pred), 3)

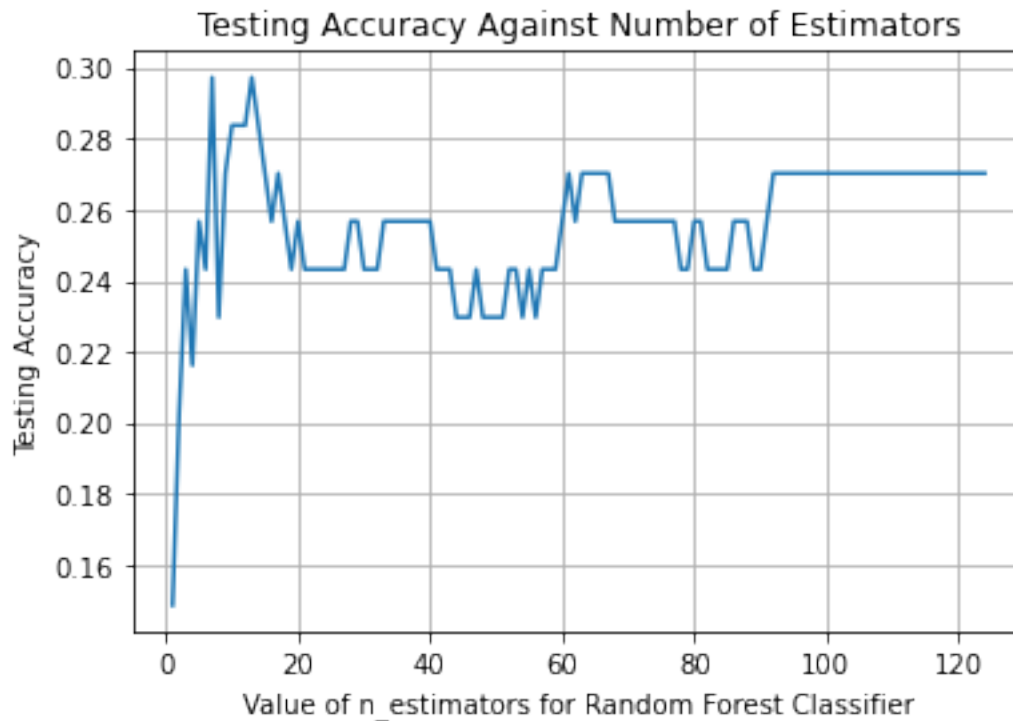
print(f"Accuracy: {accuracy}")
```

Accuracy: 0.297

```
[338]: scores = []
for i in range(1, 125):
    rfc = RandomForestClassifier(n_estimators = i, random_state = 42)
    rfc.fit(X_train, Y_train)
    Y_pred = rfc.predict(X_test)
    scores.append(accuracy_score(Y_test, Y_pred))

plt.plot(range(1, 125), scores)
plt.xlabel('Value of n_estimators for Random Forest Classifier')
plt.ylabel('Testing Accuracy')
plt.title('Testing Accuracy Against Number of Estimators')
plt.grid(True); GridSearchCV
```

```
[338]: sklearn.model_selection._search.GridSearchCV
```



All Variables

```
[339]: X = user_data.drop(['NewOutcomeScore', 'OutcomeScore', 'group_id', 'Line_ID', '
    ↳ 'Unnamed: 0', 'userIDs'], axis = 1)
Y = user_data['NewOutcomeScore']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,
    ↳ random_state=42)

rfc = RandomForestClassifier(n_estimators = 24, max_samples = 0.5, random_state=
    ↳ 42).fit(X_train, Y_train)
Y_pred = rfc.predict(X_test)
accuracy = np.round(accuracy_score(Y_test, Y_pred), 3)

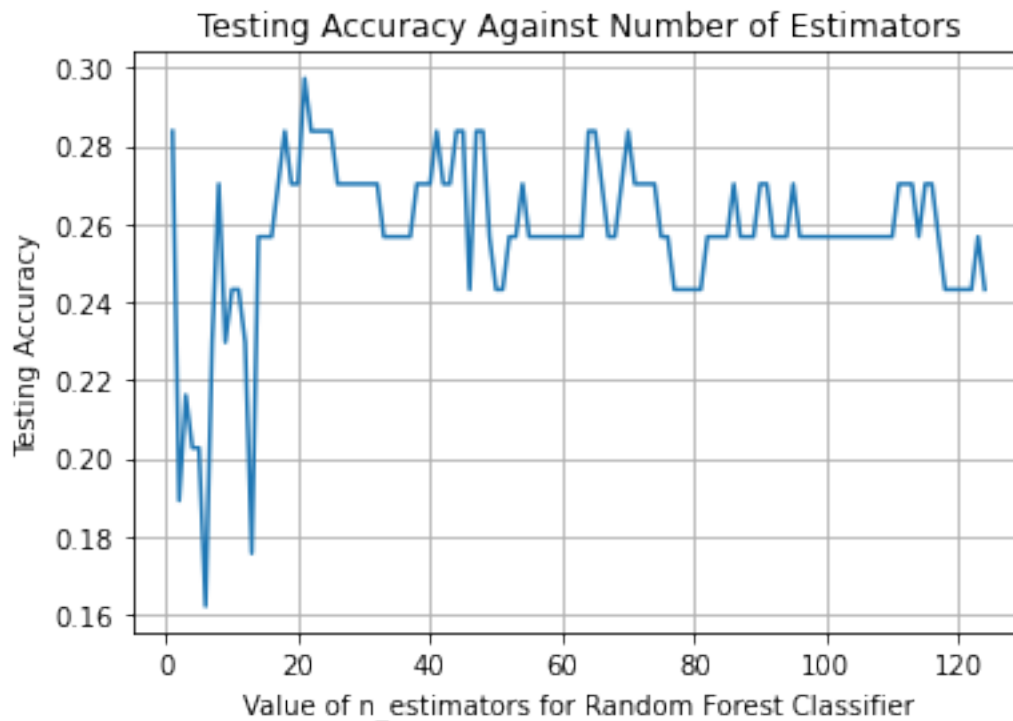
print(f"Accuracy: {accuracy}")
```

Accuracy: 0.311

```
[340]: scores = []
for i in range(1, 125):
    rfc = RandomForestClassifier(n_estimators = i, random_state = 42)
    rfc.fit(X_train, Y_train)
    Y_pred = rfc.predict(X_test)
    scores.append(accuracy_score(Y_test, Y_pred))
```



```
plt.plot(range(1, 125), scores)
plt.xlabel('Value of n_estimators for Random Forest Classifier')
plt.ylabel('Testing Accuracy')
plt.title('Testing Accuracy Against Number of Estimators')
plt.grid(True);
```



## Feature Importance

```
[341]: for name, score in zip(X.columns, rfc.feature_importances_):
        print(name, np.round(score, 3))
```

```
m_experimental_testing 0.078
m_making_design_choices 0.117
m_asking_questions 0.138
j_customer_consultants_requests 0.062
j_performance_parameters_requirements 0.093
j_communication 0.046
wordCount 0.195
input 0.136
output 0.134
```

Five most important features.

```
[342]: X = user_data[['wordCount', 'm_asking_questions', 'm_making_design_choices',
    → 'input', 'output']]
Y = user_data['NewOutcomeScore']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, train_size=0.8,
    → random_state=42)

rfc = RandomForestClassifier(n_estimators = 56, max_samples = 0.5, random_state=
    → 42).fit(X_train, Y_train)
Y_pred = rfc.predict(X_test)
accuracy = np.round(accuracy_score(Y_test, Y_pred), 3)

print(f"Accuracy: {accuracy}")
```

Accuracy: 0.297

```
[343]: scores = []
for i in range(1, 125):
    rfc = RandomForestClassifier(n_estimators = i, random_state = 42)
    rfc.fit(X_train, Y_train)
    Y_pred = rfc.predict(X_test)
    scores.append(accuracy_score(Y_test, Y_pred))

plt.plot(range(1, 125), scores)
plt.xlabel('Value of n_estimators for Random Forest Classifier')
plt.ylabel('Testing Accuracy')
plt.title('Testing Accuracy Against Number of Estimators')
plt.grid(True); GridSearchCV
```

[343]: sklearn.model\_selection.\_search.GridSearchCV

