



**ITAHARI**  
INTERNATIONAL  
COLLEGE



**Module Code & Module Title**  
**CU6051NT Artificial Intelligence**

**Assessment Weightage & Type**  
**20% Individual Coursework**

**Year and Semester**  
**2020-21 Autumn Year Long**

**Student Name: Girija Tamang**

**London Met ID: 18030995**

**Title: Red Wine Quality Prediction**

**Assignment Due Date: 17 Jan 2021**

**Assignment Submission Date: 17 Jan 2021**

**Module Tutor: Weenit Maharjan**

*I confirm that I understand my coursework needs to be submitted online via Google Classroom under the relevant module page before the deadline for my assignment to be accepted and marked. I am fully aware that late submissions will be treated as non-submission and a mark of zero will be awarded.*

## Abstract

This report is about “Red Wine Quality Prediction” using random forest algorithm. The report includes literature review of similar projects and briefly explains the AI concepts and terminologies used in the proposed solution. The flowchart, pseudocode, and the explanation of used algorithms with work analysis, and the future work is presented in the report. The main aim of this project is to predict the quality of red wine and help wine industries in checking the quality of wines for making good wine products in the future.

## Table of Contents

1. Introduction.....	1
1.1 Introduction to the AI concepts used. ....	1
1.2 Introduction of the chosen problem domain .....	2
2. Background .....	3
3. Solution .....	6
3.1 Explanation of the proposed solution .....	6
3.2 Explanation of the AI algorithm used. ....	8
3.3 Pseudocode of the solution .....	10
3.4 Diagrammatical representations of the solution .....	11
4. Conclusion .....	12
4.1 Analysis of the work done. ....	12
4.2 How the solution addresses real world problems. ....	12
4.3 Further work.....	13
5. References .....	14

## Tables of Figures

Figure 1: Article on Wine Quality and Taste Classification Using Machine Learning Model .....	3
Figure 2: Research Paper on The Classification of White Wine and Red Wine According to Their Physicochemical Qualities. ....	4
Figure 3: Article on Wine Quality Prediction using Machine Learning Algorithms. ....	5
Figure 4: Working of Random Forest Algorithm. ....	9
Figure 5: Flowchart of red wine quality prediction system. ....	11

## 1. Introduction

### 1.1 Introduction to the AI concepts used.

#### **Machine Learning**

Machine Learning at its most basic level is the practise of using algorithms to parse data, learn from it and then make a determination or prediction about something in the world. Machine learning is the branch of Artificial intelligence where a system is trained with a set of datasets or patterns using various mathematical models and algorithms to make a machine capable of making decisions or predictions independently without being explicitly programmed for performing the given task (Rouse, 2019).

#### **Random Forest**

Random forest is a popular supervised learning algorithm used for both classification and regression problems in machine learning. Random Forest is a classifier that includes a number of decision trees on different subsets of the dataset specified and takes the average to increase the predictive accuracy of the dataset. It is an ensemble technique that is better than a single decision tree because by averaging the result, it decreases the over-fitting.

Here are some points illustrating why the algorithm of the Random Forest is used:

- As compared to other algorithms, it takes less training time.
- Even for the large dataset it manages effectively, it predicts highly accurate results.
- It can also preserve precision when a significant proportion of the data is missing.
- Random forest has less variance than a single decision tree (Tutorials Point, 2019).

## 1.2 Introduction of the chosen problem domain

### Red Wine Quality Prediction

Product quality has been one of the essential components of any single industry in the recent years. Since the last decade, the wine industry has been rising well in the market. As wine demand has risen in recent years, the consumption of wine has also increased. With increasing demand, the quality checking of wine has been the major problem faced by wine industries. Wine quality is generally measured by professional tasters in the wine industry who render their decision based on several sensory criteria, such as color, taste, and odor, which is very complex and time-consuming since wine demands have been increasing worldwide.

Understanding the criteria of wine quality testing in industries can be a challenging activity for a laboratory with a wide variety of analyses and residues to track. Different kinds of machine learning algorithms should be implemented by wine industries for analyzing taste and other properties in wine. Machine learning makes it more effective to test or predict any kind of thing effectively, to find the quality of wine in a short time without the need for any human expertise.

## 2. Background

Some significant work has been done in the field of wine quality prediction. In the past decade, academic papers on the very subject have proliferated. To offer a clear understanding of wine quality prediction techniques and their implementation, this proposal analyses the various papers written.

### Wine Quality and Taste Classification Using Machine Learning Model

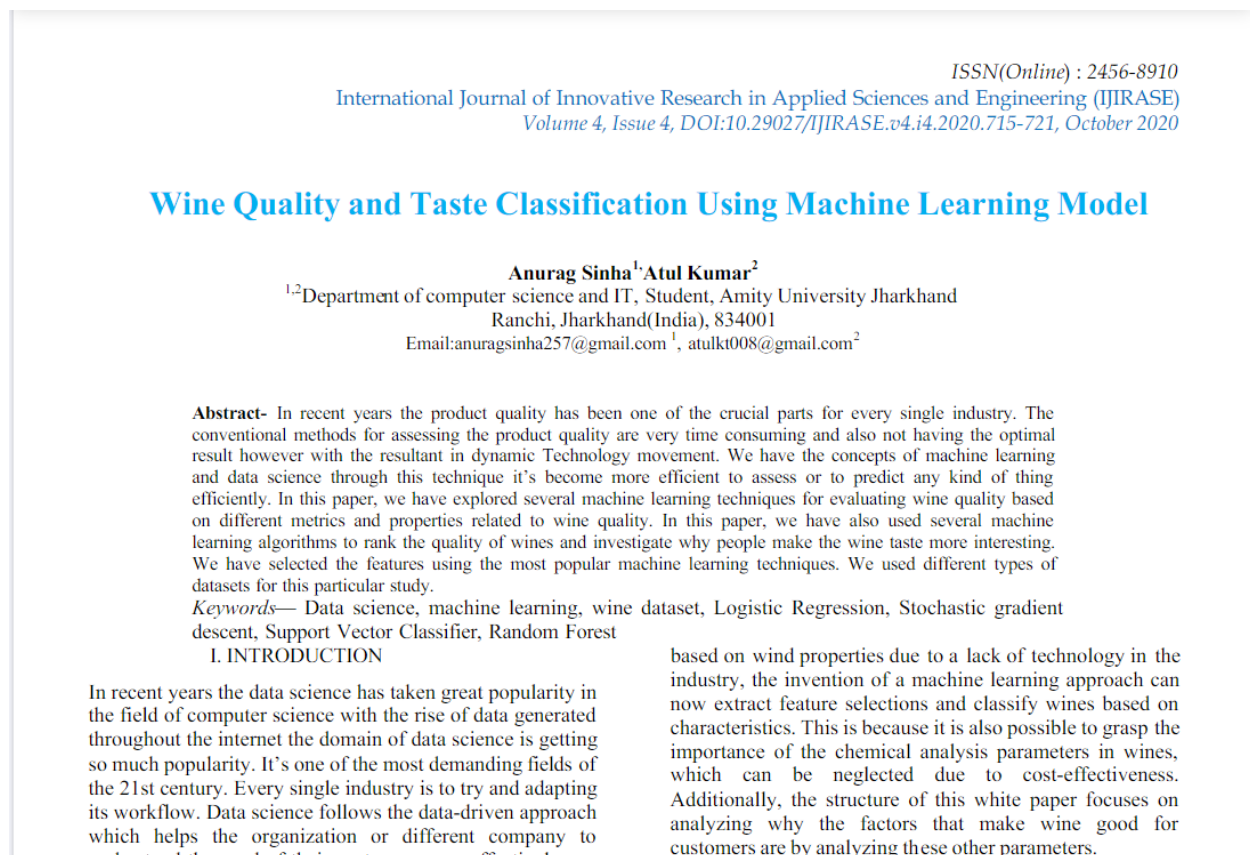


Figure 1: Article on Wine Quality and Taste Classification Using Machine Learning Model

Research on wine quality and taste classification using machine learning has been done by Anurag Sinha and Atul Kumar. In this paper, several machine learning techniques were explored to determine wine quality based on different parameters and properties related to wine quality. In this paper, they have used logistic regression, Stochastic descent of the gradient, support Vector classifier and Random forest machine learning algorithms for classification of wines.

Logistic regression and Random Forest provided 86% and 87.33% accuracy in predicting quality of wines. High quality of wine is usually associated with low levels of volatile acidity (Sinha & Kumar, 2010). They have used various datasets of wines for doing this research which was actually a good idea for assuring wine quality and taste classification. In this paper, they have described the used dataset clearly but there was no proper description of the algorithm used. This paper helps me to understand and select a dataset to predict the quality of red wine.

## The Classification of White Wine and Red Wine According to Their Physicochemical Qualities



Figure 2: Research Paper on The Classification of White Wine and Red Wine According to Their Physicochemical Qualities.

Research of this article was performed by Yesim Er and Ayten Atasoy in the International Journal of Intelligent Systems and Applications in Engineering. Predicting the quality of wine based on physicochemical data was the main objective of this research paper. In this study, two different large data sets taken from the UC Irvine Machine Learning Repository were used. They have used both red and white dataset for classification of wine. In this report they have used k-nearest-neighborhood, random forests, and support vector machines classifier for evaluating the datasets of both red and white wine.



Using the Random Forests Algorithm, the cases were successfully categorized as red wine and white wine with an accuracy of 99.5229 percent (Er & Atasoy, 2016). In this article they have properly described the dataset and all the used machine learning algorithms. This article helps me to choose a random forest algorithm for quality prediction of red wine.

## Wine Quality Prediction using Machine Learning Algorithms

International Journal of Computer Applications Technology and Research  
Volume 8–Issue 09, 385-388, 2019, ISSN:-2319–8656

### Wine Quality Prediction using Machine Learning Algorithms

Devika Pawar<sup>[1]</sup>  
M.Sc. (Big Data Analytics)  
MIT-WPU  
Pune, India

Aakanksha Mahajan<sup>[2]</sup>  
M.Sc. (Big Data Analytics)  
MIT-WPU  
Pune, India

Sachin Bhoithe<sup>[3]</sup>  
Faculty of Science  
MIT-WPU  
Pune, India

**Abstract:** Wine classification is a difficult task since taste is the least understood of the human senses. A good wine quality prediction can be very useful in the certification phase, since currently the sensory analysis is performed by human tasters, being clearly a subjective approach. An automatic predictive system can be integrated into a decision support system, helping the speed and quality of the performance. Furthermore, a feature selection process can help to analyze the impact of the analytical tests. If it is concluded that several input variables are highly relevant to predict the wine quality, since in the production process some variables can be controlled, this information can be used to improve the wine quality. Classification models used here are 1) Random Forest 2) Stochastic Gradient Descent 3) SVC 4) Logistic Regression.

**Keywords:** Machine Learning, Classification, Random Forest, SVM, Prediction.

#### I. INTRODUCTION

The aim of this project is to predict the quality of wine on a scale of 0–10 given a set of features as inputs. The dataset used is Wine Quality Data set from UCI Machine Learning Repository. Input variables are fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulphur dioxide, total sulphur dioxide, density, pH,

that the significant difference between the two is small. Then this paper uses the Cronbach Alpha coefficient method to analyze the credibility of the two groups of data.<sup>[1]</sup>

Paulo Cortez, Juliana Teixeira, António Cerdeira, Fernando Almeida, Telmo Matos, José Reis wrote a paper on wine Quality assesment using Data Mining techniques. In this

Figure 3: Article on Wine Quality Prediction using Machine Learning Algorithms.

Devika Pawar, Aakanksha Mahajan and Sachi Bhoithe wrote a paper on wine quality prediction using machine learning algorithms. Random Forest, Stochastic Gradient Descent, Logistic Regression are the classification models they have used for predicting the wine quality. They were able to achieve optimum precision using random forests of 88 percent where Stochastic gradient, SVC, and logistic regression were able to provide 81 %, 85%, and 86% accuracy, respectively.

This article is well prepared for learning basic machine learning concepts with various algorithms (Pawar, et al., 2019). This paper presented a proper explanation of the data set and past work done to predict the quality of wine. There was a lack of a proper explanation of how machine learning algorithms operate when predicting the quality of wine.

### 3. Solution

#### 3.1 Explanation of the proposed solution

The checking and predicting of wine quality have been the major problem faced by wine industries. After doing a lot of research, I found that there are many machine learning algorithms to determine wine quality based on various wine quality parameters and properties. I have chosen a random forest algorithm for red wine quality prediction. Here are the steps that are followed while testing quality of wines:

1. Data Collection of Red Wine from public datasets.
2. Data preparation for building models.
3. Feature selection
4. Implementing machine learning techniques
5. Comparison of performance.
6. Interpretation of results

#### Data set Information

The dataset is related to red variants of the Portuguese "Vinho Verde" wine. Only physicochemical (inputs) and sensory (output) variables are accessible due to privacy and logistical problems (example: there is no data about grape types, wine brand, wine selling price, etc.). The classes are ordered and not balanced. This dataset can be viewed as classification or regression tasks (P. Cortez, 2019).

**Attribute Information: Input variables (based on physicochemical tests):**

1. fixed acidity
2. volatile acidity
3. citric acid
4. residual sugar
5. chlorides
6. free sulfur dioxide
7. total sulfur dioxide
8. density
9. pH
10. sulphates
11. alcohol

Output variable (based on sensory data):

12. quality (score between 0 and 10).

After finding the dataset, I have performed initial visual analysis of data present in the dataset. So, I got a concise summary of the data frame like information from the data set, whether there are null values available or not in the data frame, shows the data types and detailed information of the columns. I will visualize the relationship between quality and the other columns using subplot. By using label encoding, I will categories the label data of quality of good or bad. I will assign 1 to good and 0 to bad quality of wine. After that I will separate the dataset as response variables and feature variables. I will split the data to both testing and training data sets and then I will standardize the data. Our training and testing data will be ready after standardization. I will perform classification using a random forest classifier machine learning algorithm and see how the developed model will perform.

### 3.2 Explanation of the AI algorithm used.

The random forest classifier is used for predicting the quality of red wine after preparing training and testing data from the red wine dataset. Random Forest is a popular algorithm for machine learning based on the principle of ensemble learning, which is a method of combining multiple classifiers to solve a complex problem and improve the model's output. Random Forest is a classifier that comprises a number of decision trees and takes the average to boost the predictive accuracy of that dataset on different subsets of the given dataset. The random forest takes the prediction from each tree and is based on the majority votes of predictions rather than depending on a decision tree and predicts the final output (JavaTpoint, 2019).

In two phases, Random Forest operates, first by combining the N decision tree to create the random forest, and secondly by making predictions for each tree generated in the first phase. The Working process can be explained in the below steps:

Step-1: Choose a random K data point from the training set.

Step-2: Create decision trees associated with the data points selected (Subsets)

Step-3: Choose the number N for the decision trees you want to create.

Step-4: Repeat step 1 and step 2.

Step-5: For new data points, find the predictions of each decision tree, and assign the new data points to the category that wins the majority votes (Sharma, 2019).

The following diagram will illustrate the working of a random forest algorithm.

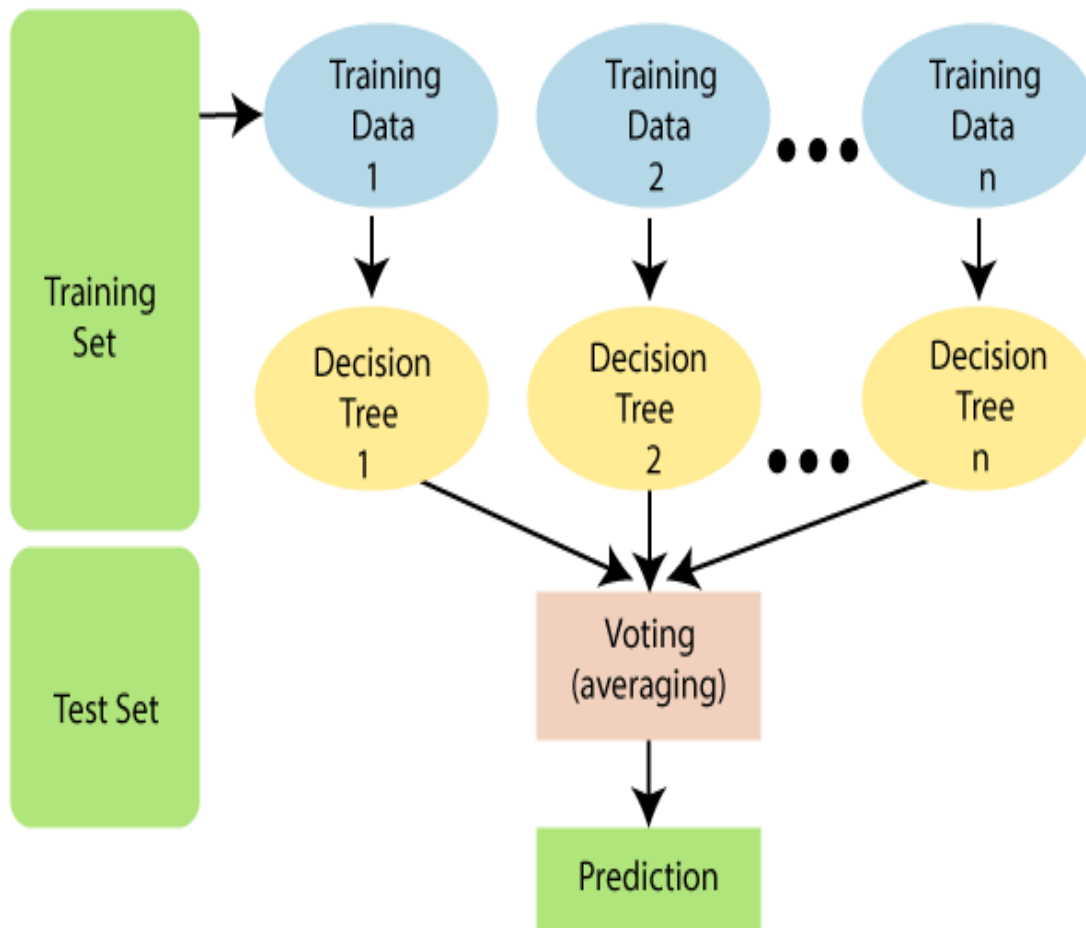


Figure 4: Working of Random Forest Algorithm.

### 3.3 Pseudocode of the solution

Step 1: Import Necessary Libraries

```
import pandas as pd

import seaborn as sns

import matplotlib.pyplot as plt

from sklearn.ensemble import RandomForestClassifier
```

Step 2: Load the data set for getting the data of wine.

```
df = pd.read_csv('filename.csv')
```

Step 3: Check how the data set looks or distributed by using dot info attribute and to know how the columns of data are distributed in the dataset, do some plotting.

Step 4: Divide wine as good and bad by giving the limit for the quality and assign a label to the quality variable using label encoding.

Step 5: Separate the data as response variables and feature variables.

Example: `X = wine.drop('quality', axis = 1)`

```
y = wine['quality']
```

Step 5: Split the data into training and testing sets and apply standard scaling to get optimized results. The training and testing data will be ready to perform with machine learning algorithms.

Step 6: Now use a random forest classifier for predicting the quality of red wine.

Example: `RF = RandomForestClassifier(n_estimators=100)`

```
RF.fit(X_train, Y_train)
```

```
pred_RF = RF.predict(X_test)
```

Step 7: Calculate the accuracy score with the target variable and with the predicted target variable using random forest classifier.

### 3.4 Diagrammatical representations of the solution

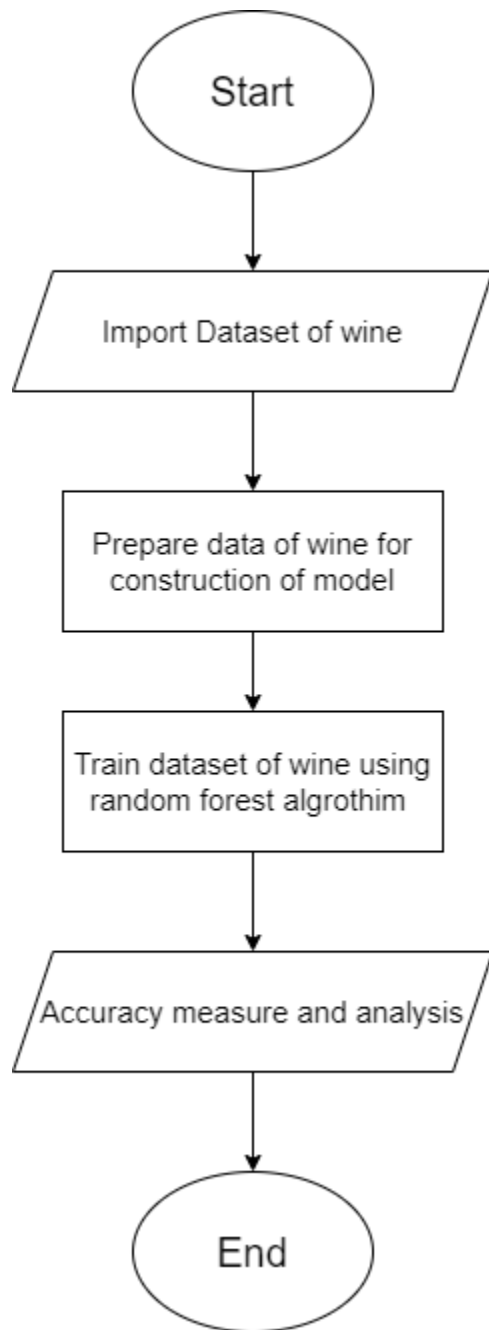


Figure 5: Flowchart of red wine quality prediction system.

## 4. Conclusion

### 4.1 Analysis of the work done.

In this report, the red wine quality prediction system is proposed which may help wine industries for checking and predicting wine quality. The random forest algorithm is used for classification of red wine. A lot of research was done on red wine quality prediction using machine learning algorithms. This report offered a good understanding of the value of quality prediction attributes using functions selected on the algorithm, which was time consuming and costly when performed in the conventional way. In the red wine data set, the Random Forest classifier can predict the quality of the wine better than any other classifier. This report has the solution and dataset information for predicting the red wine quality with the proper description of the algorithm that is used. This project may face some limitation such as:

1. Since the prediction is entirely data-based, accuracy is highly dependent on it.
2. Model was not checked with real-time data on red wines and this project has not focused on white wine.
3. The system may be ideal for factories and not for local citizens.

### 4.2 How the solution addresses real world problems.

This paper discussed the use of machine learning techniques to predict the quality of Red Wine. The feature selection algorithm provided a clear idea about the importance of the attributes for prediction of quality, which was time consuming and expensive when done in the traditional way. Random forest machine learning algorithms should be implemented by wine industries for analyzing taste and other properties in wine. It makes it more effective to test, predict or find the quality of wine in a short time without the need for any human expertise.

The wine industry should invest in new technology such as data mining to evaluate taste and other properties in wine in order to stay competitive in the future. Machine learning will help companies to predict the quality of the various types of wines on the basis of certain characteristics in a short period of time, and it will be beneficial for them to produce good wine products in coming days.



### 4.3 Further work

In the future, this device could be improved in several respects. To some degree, the prediction made by this system is reliable, but the prediction of systems can be made more precise by gathering more real-time data and other training attributes that influence it. In the future, I will try other machine learning methods for a better comparison of outcomes. I will try to add other machine learning algorithms for predicting the quality of the different types of wines based on certain attributes which will be helpful for industries to make good products of wine in the future.

## 5. References

Er, Y. & Atasoy, A., 2016. The Classification of White Wine and Red Wine According to Their Physicochemical Qualities. *Intelligent Systems and Applications in Engineering*, I(2147), pp. 23-26.

JavaTpoint, 2019. *Random Forest Algorithm*. [Online] Available at: <https://www.javatpoint.com/machine-learning-random-forest-algorithm> [Accessed 06 January 2021].

P. Cortez, A. C. F. A. T. M. a. J. R., 2019. *Wine Quality Data Set*. [Online] Available at: <https://archive.ics.uci.edu/ml/datasets/wine+quality> [Accessed 13 January 2021].

Pawar, D., Mahajan, A. & Bhoithe, S., 2019. Wine Quality Prediction using Machine Learning. *International Journal of Computer Applications Technology and Research*, 8(09), pp. 385-388.

Rouse, M., 2019. *Machine Learning*. [Online] Available at: <https://searchenterpriseai.techtarget.com/definition/machine-learning-ML> [Accessed 5 January 2021].

Sharma, N., 2019. Quality Prediction of Red Wine based on Different. *Department of Computer Science & Engineering*, 4(5), pp. 20-34.

Sinha, A. & Kumar, A., 2010. Wine Quality and Taste Classification Using Machine Learning Model. *International Journal of Innovative Research in Applied Sciences and Engineering (IJIRASE)*, 4(4), pp. 715-721.

Tutorials Point, 2019. *Classification Algorithms - Random Forest*. [Online] Available at: [https://www.tutorialspoint.com/machine\\_learning\\_with\\_python/machine\\_learning\\_with\\_python\\_classification\\_algorithms\\_random\\_forest.htm](https://www.tutorialspoint.com/machine_learning_with_python/machine_learning_with_python_classification_algorithms_random_forest.htm) [Accessed 5 January 2021].