# WINE QUALITY PREDICTION

**Submitted To -**   Mr. Mayur Dev Sewak
General Manager,
Operations Eisystems Services

&

Ms. Mallika Srivastava
Trainer, Data Science & Analytics
Domain Eisystems Services

**Submitted By**   – ABHILASH BURAGOHAIN

**Course**   – Data Science / Machine Learning

# CONTENTS

# ABSTRACT

The quality of a wine is important for the consumers as well as the wine industry. The traditional (expert) way of measuring wine quality is time-consuming. Nowadays, machine learning models are important tools to replace human tasks. In this case, there are several features to predict the wine quality but the entire features will not be relevant for better prediction. So, our thesis work is focusing on what wine features are important to get the promising result. For the purpose of classification model and evaluation of the relevant features, we used three algorithms namely support vector machine (SVM), naïve Bayes (NB), and artificial neural network (ANN). In this study, we used two wine quality datasets red wine and white wine. To evaluate the feature importance we used the Pearson coefficient correlation and performance measurement matrices such as accuracy, recall, precision, and f1 score for comparison of the machine learning algorithm. A grid search algorithm was applied to improve the model accuracy. Finally, we achieved the artificial neural network (ANN) algorithm has better prediction results than the Support Vector Machine (SVM) algorithm and the Naïve Bayes (NB) algorithm for both red wine and white wine datasets.

# KEYWORDS

# INTRODUCTION

The quality of the wine is a very important part for the consumers as well as the manufacturing industries. Industries are increasing their sales using product quality certification. Nowadays, all over the world wine is a regularly used beverage and the industries are using the certification of product quality to increases their value in the market. Previously, testing of product quality will be done at the end of the production, this is time taking process and it requires a lot of resources such as the need for various human experts for the assessment of product quality which makes this process very expensive. Every human has their own opinion about the test, so identifying the quality of the wine based on humans experts it is a challenging task.

There are several features to predict the wine quality but the entire features will not be relevant for better prediction.

Our goal for this project is to conduct an analysis on Wine Quality data from the University of California, Irvine Machine Learning Repository. The data are the results of a chemical analysis of wines grown in the same region in Italy, but derived from three different cultivars. The data has 13 attributes which are wine type, fixed acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates, alcohol, and quality. We are interested to see if these attributes has an impact on the quality of wine. Therefore, our goal is to use machine learning to see if we can build an accurate model to predict the wine quality. Below we can see the features that were used to build this model.

# PROJECT SUMMARY AND OBJECTIVES

My analysis will use Red Wine Quality Data Set, available on the UCI machine learning repository. I obtained the red wine samples from the north of Portugal model red wine quality based physiochemical tests. The dataset contains total 12 variables, which were recorded for 1599 observation. This data will allow us to create different regression models to determine how different independent variables help predict our dependent variable, quality. Knowing how each variable will impact the red wine quality will help producers, distributors and businesses in the red wine industry better assess their production, distribution, and pricing strategy.

Wine classification is a difficult task since taste is the least understood of 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

# LANGUAGE USED

Python is an OOPs (Object Oriented Programming) based, high level, interpreted programming language. It is a robust, highly useful language focused on rapid application development (RAD). Python helps in easy writing and execution of codes. Python can implement the same logic with as much as 1/5th code as compared to other OOPs languages.

Python provides a huge list of benefits to all. The usage of Python is such that it cannot be limited to only one activity. Its growing popularity has allowed it to enter into some of the most popular and complex processes like Artificial Intelligence (AI), Machine Learning (ML), natural language processing, data science etc. Python has a lot of libraries for every need of this project. For JIA, libraries used are speech recognition to recognize voice, Pyttsx for text to speech, selenium for web automation etc.

Libraries Used:

1. **NumPy**: It is a general-purpose array processing package. It provides a high-performance multidimensional array object and tools for working with these arrays. It is the fundamental package for scientific computing with Python. It is an open source software.

2. **Pandas**: It is a Python library used for working with data sets.It has functions for analyzing, cleaning, exploring, and manipulating data.The name "Pandas" has a reference to both "Panel Data", and "Python Data Analysis" and was created by Wes McKinney in 2008. Pandas allows us to analyse big data and make conclusions based on statistical theories.Pandas can clean messy data sets, and make them readable and relevant.Relevant data is very important in data science.

3. **Matplotlib**: Matplotlib is a cross-platform, data visualization and graphical plotting library for Python and its numerical extension NumPy. As such, it offers a viable open source alternative to MATLAB. Developers can also use matplotlib's APIs (Application Programming Interfaces) to embed plots in GUI applications. A Python matplotlib script is structured so that a few lines of code are all that is required in most instances to generate a visual data plot.

4. **Seaborn**: Seaborn is a Python data visualization library based on matplotlib.

It provides a high-level interface for drawing attractive and informative statistical graphics.

5. **Sklearn.model selection**: The classes in the sklearn.feature_selection module can be used for feature selection/dimensionality reduction on sample sets, either to improve estimators' accuracy scores or to boost their performance on very high-dimensional datasets.

6. **Sklearn.ensemble**: The sklearn.ensemble module includes two averaging algorithms based on randomized decision trees: the RandomForest algorithm and the Extra-Trees method. Both algorithms are perturb-and-combine techniques [B1998] specifically designed for trees. This means a diverse set of classifiers is created by introducing randomness in the classifier construction. The prediction of the ensemble is given as the averaged prediction of the individual classifiers.

7. **Sklearn.metrics**: The sklearn.metrics module implements functions assessing prediction error for specific purposes. These metrics are detailed in sections on Classification metrics, Multilabel ranking metrics, Regression metrics and Clustering metrics.

# SOFTWARE USED

## ANACONDA NAVIGATOR

**Anaconda** comes with over 250 packages automatically installed and over 7500 additional packages can be installed from PyPi as well as the conda package and environment manager. It also includes a GUI, Anaconda Navigator as a graphical alternative to the command-line interface (CLI).

Anaconda Navigator is a desktop graphical user interface (GUI) included in Anaconda distribution that allows users to launch applications and manage conda packages, environments and channels without using command-line commands. Navigator can search for packages on Anaconda Cloud or in a local Anaconda Repository, install them in an environment, run the packages and update them. It is available for Windows, macOS and Linux.

The following applications are available by default in Navigator:

- JupyterLab
- Jupyter Notebook
- QtConsole
- Spyder
- Glue
- Orange
- RStudio
- Visual Studio Code

## JUPYTER NOTEBOOK:

Jupyter Notebook (formerly IPython Notebook) is a web-based interactive computational environment for creating notebook documents. Jupyter Notebook is built using several open-source libraries, including IPython, ZeroMQ, Tornado, jQuery, Bootstrap, and MathJax. A Jupyter Notebook document is a browser-based REPL containing an ordered list of input/output cells which can contain code, text (using Markdown), mathematics, plots and rich media. Underneath the interface, a notebook is a JSON document, following a versioned schema, usually ending with the ".ipynb" extension

# DIAGRAMS



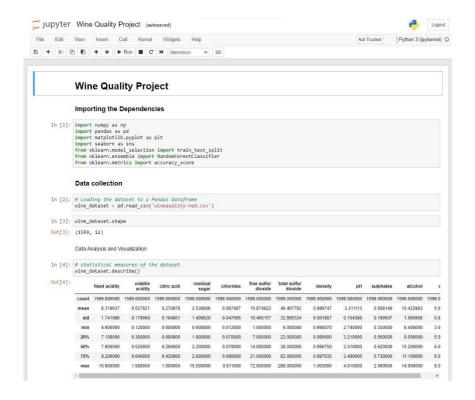**Figure 1:** Proposed Wine Quality Classification model

The overall flow structure of the proposed Wine quality classification scheme is depicted in figure 1. The wine quality dataset taken from Kaggle. Subsequently, the machine learning models where used which includes regression analysis and classification techniques to presume the better accuracy of the model. To classify given input data into good quality and bad quality, all data were fed into three different classifier.
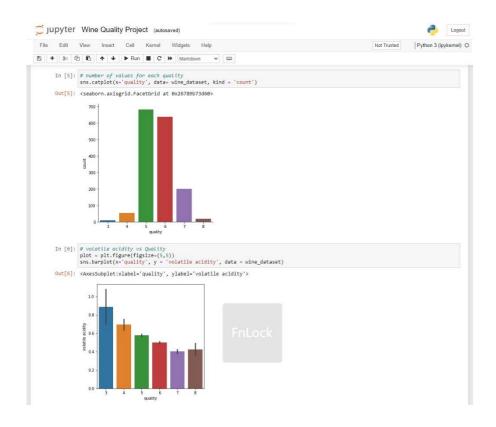
**DATA FLOW DIAGRAM**

# SCREENSHOTS OF THE CODE

IMG 1



IMG 2

IMG 3



IMG 4

IMG 5



IMG 6

IMG 7



IMG 8

# REFERENCE

• Websites referred

♣ www.kaggle.com

♣ www.pythonprogramming.net

♣ https://scikit-learn.org/

♣ www.geekforgeeks.com

♣ www.google.co.in


• Books referred

♣ Python Programming - Kiran Gurbani

♣ Learning Python - Mark Lutz


• YouTube Channels referred

♣ codewithharry

♣ edureka!

# *THANK YOU*