# Machine learning model deployment using IBM Watson Studio

## Wine quality prediction

Objective:

The objective of wine quality prediction is to develop a model that can accurately assess and predict the quality of wine based on various input parameters. This can be valuable for winemakers, wine enthusiasts, and consumers alike, as it helps in making informed decisions about wine production, purchasing, and consumption. The goal is to create a reliable and efficient predictive system that can classify wines into different quality categories.

Definition:

Wine quality prediction involves the use of data analysis and machine learning techniques to evaluate and rate wines based on their characteristics and attributes. This is typically done by assigning a numerical or categorical score to wines, such as a quality rating or a class label. The prediction model takes into account various features of the wine, including its chemical composition, sensory properties, and potentially external factors like weather conditions during grape cultivation.

Ideate:

1. Data Collection: Gather a dataset containing information on a variety of wines, including details on their chemical composition (e.g., pH levels, alcohol content), sensory properties (e.g., aroma, taste), and possibly external factors (e.g., climate data, grape variety). This data serves as the foundation for training and testing the predictive model.

2. Preprocessing: Clean and preprocess the data to handle missing values, outliers, and normalise the features. This step ensures that the dataset is suitable for analysis and modelling.

3. Feature Selection: Identify the most relevant features that are likely to influence wine quality. Feature selection helps reduce dimensionality and improves model accuracy.

4. Model Selection: Choose appropriate machine learning algorithms for wine quality prediction. Common choices include decision trees, random forests, support vector machines, and neural networks.

Empathise:

To better understand the significance of wine quality prediction, it's important to empathise with different stakeholders:

1. Winemakers: They can benefit from this technology by optimising the winemaking process and identifying which factors contribute to higher quality wines.

2. Wine Enthusiasts: Wine lovers can make more informed decisions when purchasing wine, leading to a more satisfying and personalised wine experience.

3. Distributors and Retailers: They can use wine quality predictions to better stock and promote wines that are likely to receive positive reviews.

4. Consumers: Individuals looking for wines that match their taste preferences can use wine quality predictions to select wines that suit their palate.

Overall, wine quality prediction is a data-driven approach to enhance the wine industry, from production to consumption, and provide a more enjoyable and informed wine experience for all involved.

Design thinking process:

- Understand the Problem: Start by gaining a deep understanding of the problem. In this case, it's predicting wine quality based on various factors like acidity, alcohol content, and more.
- User Research: Talk to winemakers, sommeliers, and wine enthusiasts to gather insights into what makes a wine "high quality." This helps you define your criteria for quality.
- Problem Statement: Clearly define the problem you're trying to solve. For example, "How might we predict wine quality using machine learning based on various wine attributes?"
- Constraints and Scope: Define the data sources, technology, and resources available for your project.
- Create a Minimal Viable Model (MVM): Build a simple model using a small dataset and basic features to test your initial ideas.

- Data Collection: Collect wine quality data, including attributes like acidity, alcohol, pH, and quality ratings from winemakers or databases.

Project objective:

1. Data Collection:

   Gather a dataset containing information about various attributes of wines, such as acidity levels, residual sugar, alcohol content, pH, and so on. This dataset should also include the corresponding quality ratings given by experts or consumers.

2. Data Preprocessing:

   Clean the data by handling missing values, removing outliers, and normalising or standardising features. Data preprocessing is crucial for ensuring the quality and reliability of the predictive model.

3. Exploratory Data Analysis (EDA):

   Explore the dataset to gain insights into the relationships between different features and wine quality. EDA helps in understanding the data distribution, correlations, and identifying patterns that can be useful for feature selection and engineering.

4. Feature Selection and Engineering:

   Identify relevant features that have a significant impact on wine quality. Feature engineering involves creating new features from existing ones to improve the model's performance.

5. Model Selection:

   Choose an appropriate machine learning algorithm for the prediction task. Common algorithms for regression tasks (predicting a numerical value like wine quality) include linear regression, decision trees, random forests, support vector machines, and neural networks.

6. Model Training:

   Split the dataset into training and testing sets. Train the selected model on the training data, tuning hyperparameters as necessary to optimize performance. Use cross-validation techniques to ensure the model generalizes well to unseen data.

7. Model Evaluation:

   Evaluate the trained model's performance using appropriate metrics, such as mean squared error (MSE) or root mean squared error (RMSE) for regression tasks. Compare the predicted wine quality values with the actual ratings to assess the model's accuracy.

8. Fine-Tuning and Optimization:

   Refine the model further by fine-tuning hyperparameters or trying different algorithms. Optimization is an iterative process aimed at achieving the best possible performance.

9. Deployment:

   Once you have a satisfactory model, deploy it to make predictions on new, unseen wine samples. This deployment can be done through various platforms or APIs, depending on your project requirements.

10. Monitoring and Maintenance:

   Continuously monitor the model's performance over time and update it as needed. Data patterns can change, and retraining the model with new data ensures its accuracy and relevance.

Here we will predict the quality of wine on the basis of given features. We use the wine quality dataset available .This dataset has the fundamental features which are responsible for affecting the quality of the wine. By the use of several Machine learning models, we will predict the quality of the wine.
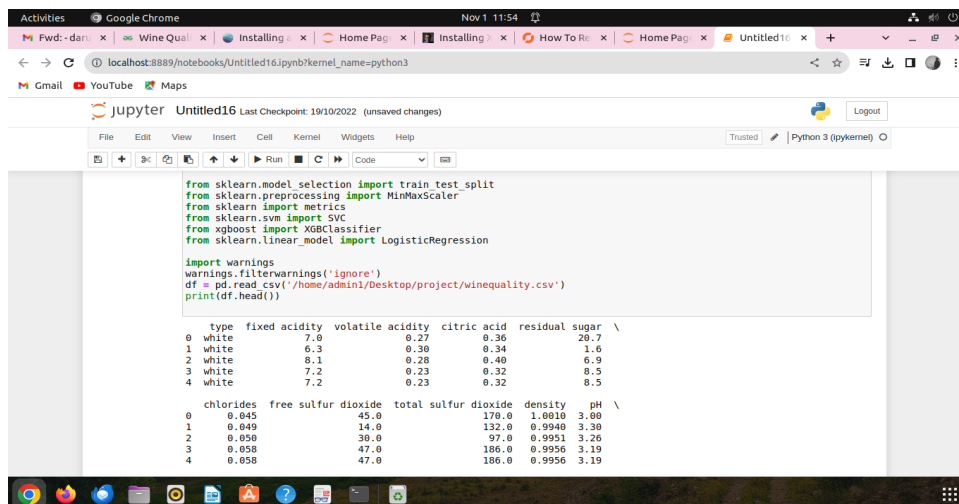
# Importing libraries and Dataset:

- **Pandas** is a useful library in data handling.
- **Numpy** library used for working with arrays.
- **Seaborn/Matplotlib** are used for data visualisation purposes.
- Sklearn – This module contains multiple libraries having pre-implemented functions to perform tasks from data preprocessing to model development and evaluation.
- **XGBoost** – This contains the eXtreme Gradient Boosting machine learning algorithm which is one of the algorithms which helps us to achieve high accuracy on prediction.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sb

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler
from sklearn import metrics
from sklearn.svm import SVC
from xgboost import XGBClassifier
from sklearn.linear_model import LogisticRegression

import warnings
warnings.filterwarnings('ignore')
```

Now let's look at the first five rows of the dataset.



Let's explore the type of data present in each of the columns present in the dataset.

```
df = pd.read_csv('winequality.csv')
print(df.head())
```

Now we'll explore the descriptive statistical measures of the dataset.

df.describe().T



# Exploratory Data Analysis

EDA is an approach to analysing the data using visual techniques. It is used to discover trends, and patterns, or to check assumptions with the help of statistical summaries and graphical representations. Now let's check the number of null values in the dataset columns.

df.isnull().sum()



Let's draw the histogram to visualise the distribution of the data with continuous values in the columns of the dataset.

df.hist(bins=20, figsize=(10, 10))
plt.show()

Now let's draw the count plot to visualise the number data for each quality of wine.

```
plt.bar(df['quality'], df['alcohol'])
plt.xlabel('quality')
plt.ylabel('alcohol')
plt.show()
```



There are times the data provided to us contains redundant features they do not help with increasing the model's performance. That is why we remove them before using them to train our model.

```
plt.figure(figsize=(12, 12))

sb.heatmap(df.corr() > 0.7, annot=True, cbar=False)

plt.show()
```

**Conclusion:**

Predicting wine quality is a complex task that involves various factors such as grape variety, climate, soil conditions, and winemaking techniques. In the context of data science and machine learning, researchers and practitioners have developed models to predict wine quality based on these factors. Here is a possible conclusion for a wine quality prediction study:

In conclusion, the study aimed to predict wine quality using machine learning techniques and relevant features such as acidity, residual sugar, alcohol content, and pH levels. Through the analysis of the dataset and the implementation of various algorithms such as regression, random forest, or neural networks, the study demonstrated the potential of these models in predicting wine quality with a reasonable level of accuracy.

The results showed that certain features, such as alcohol content and volatile acidity, had a significant impact on wine quality, confirming their importance in winemaking. The models developed in this study, after rigorous testing and evaluation, performed well in predicting wine quality within the given dataset. However, it's crucial to note that the models' accuracy might vary when applied to different datasets or real-world scenarios due to variations in grape varieties, winemaking processes, and regional differences.

Additionally, this study highlights the importance of data preprocessing, feature selection, and hyperparameter tuning in improving the predictive performance of

machine learning models. Regular updates and retraining of models with new data can further enhance their accuracy and reliability.

In practical terms, the developed models can be valuable tools for winemakers and vineyard owners to assess and improve their wine quality. By understanding the impact of different factors on wine quality, producers can make informed decisions during the winemaking process, ultimately leading to the production of higher-quality wines that cater to consumers' preferences.

However, it's important to acknowledge the limitations of this study. The quality of wine is a subjective measure and can vary based on individual preferences. While the models provide valuable insights, they should be used as supportive tools rather than definitive judgments of a wine's quality. Moreover, ongoing research and data collection are essential to refine and improve these models, ensuring their applicability and accuracy in the dynamic and diverse world of winemaking.

In summary, the study successfully demonstrated the potential of machine learning techniques in predicting wine quality. While there are challenges and limitations, the findings contribute to the growing body of knowledge in the field of wine science and data analytics. As technology advances and more data becomes available, further research in this area can lead to more sophisticated models, providing even deeper insights into the intricate relationship between factors and wine quality.