

Wine Quality Prediction

By Machine Learning

Project Report



SELF CERTIFICATE

This is to certify that the dissertation/project report entitled “**Wine Quality Prediction by Machine Learning**” has been carried out by me and my group, as an authentic piece of work for the partial fulfilment of the requirements for the award of the BTECH degree under the guidance of **Partha Koley Sir**. I hereby declare that I am fully aware of the guidelines stated in the “ B Tech Project & Project Report.”

A Project Report

SUBMITTED BY

<i>NAME</i>	<i>ROLL NO.</i>
Sitaram Murmu	34600122034
Radhesh Saha	34600122031
Namita Mahata	34600122028
Sandipan Das	34600122056

ACKNOWLEDGEMENT

We feel immense pleasure to introduce “**Wine Quality Prediction by Machine Learning**” as our major project.

We would like to express our special thanks to our teacher **Mr. PARTHA KOLEY** Sir who has been a constant source of knowledge and inspiration to us, and who gave us the opportunity to do this project. We would also like to express our gratitude to our beloved parents for their review and many helpful comments and enlightening us and guiding us throughout the finalization of this project within the limited time frame.

Last but not the least, we thank all our teachers and as well as friends who have given us that much strength to keep moving on forward every time. We are greatly thankful to one and all and have no words to express our gratitude to them.

Finally, we would like to thank members of Euphoria Genx family for their moral support and encouragement.

INDEX

SL NO.	<i>CONTANT</i>	PAGE NO.
1.	Scope on project	
2.	Abstract	
3.	Introduction	
4.	Literature Survey	
5.	Data Set & Importing Libraries	
6.	Heatmap & Displot	
7.	Histogram & Pairplot	
8.	Logistic Regression	
9.	KNN	
10.	GaussianNB	
11.	Codes	
12.	Conclusion	
13.	Reference	

SCOPE OF PROJECT

- ❖ Wine quality prediction.
- ❖ Machine learning.
- ❖ Feature engineering.
- ❖ Random forest.
- ❖ Physicochemical properties.
- ❖ Wine datasets.

ABSTRACT

- Wine quality prediction has garnered significant attention in recent years due to the growing demand for efficient and accurate methods to assess wine quality, pricing, and consumer preferences. This study aims to explore the effectiveness of machine learning techniques in predicting wine quality based on physicochemical and sensory attributes. Wine quality is traditionally assessed by human experts through sensory evaluation, which is subjective, time-consuming, and costly. However, by leveraging large datasets containing chemical properties such as pH, alcohol content, residual sugar, acidity, and sulfur dioxide levels, machine learning models can offer an automated and consistent approach to wine quality prediction.
- In this research, various machine learning algorithms—including decision trees, support vector machines (SVM), random forests, and neural networks—are applied to publicly available datasets, such as the Wine Quality Dataset, to predict quality scores assigned by wine experts. Feature engineering is conducted to optimize model performance, and the models are evaluated using accuracy, precision, and F1-score. This study also explores how external factors, such as region of production, grape variety, and vintage, impact prediction accuracy.
- The results indicate that machine learning models, particularly ensemble methods like random forests, outperform traditional statistical models in wine quality prediction. The study concludes by discussing the implications of these findings for winemakers, retailers, and consumers, highlighting the potential of predictive models in enhancing production processes, optimizing pricing strategies, and personalizing wine recommendations. Furthermore, this approach could reduce reliance on subjective human judgment and provide a scalable solution for evaluating and predicting wine quality.

INTRODUCTION

Wine has long been a product of cultural significance and economic value, appreciated for its complex flavors, varied styles, and regional characteristics. The traditional evaluation of wine quality, pricing, and consumer preference relies heavily on the expertise of sommeliers and wine critics. However, these evaluations, often based on sensory perception and subjective judgment, introduce variability and may not be scalable for large-scale production or consistent quality control.

In recent years, advancements in data science and machine learning have opened up new possibilities for automating the process of wine evaluation and prediction. By leveraging data on the physicochemical properties of wine—such as pH, alcohol content, residual sugar, acidity, and sulfates along with other relevant factors like vintage, grape variety, and region of production, machine learning models can be trained to predict wine quality, pricing, and consumer preferences with a higher degree of accuracy and objectivity.

Machine learning techniques offer a promising alternative to traditional wine evaluation methods, enabling wine producers, retailers, and consumers to make more informed decisions. For producers, predictive models can help optimize production processes, improve product consistency, and target consumer preferences more effectively. For retailers and distributors, wine pricing models can be used to adjust pricing strategies based on market trends and historical sales data. For consumers, personalized wine recommendations based on previous preferences or wine ratings can enhance the buying experience.

This paper explores the application of various machine learning models to wine prediction, with a focus on predicting wine quality based on physicochemical attributes. The primary objectives of the study are to assess the performance of different algorithms in predicting quality ratings and to investigate the factors that most strongly influence these predictions. Additionally, the paper discusses the potential for machine learning to revolutionize the wine industry by providing scalable, data-driven solutions that enhance quality control, pricing, and consumer satisfaction.

❖ LITERATURE SURVEY

Wine productivity prediction is an area of research that intersects agriculture, viticulture, and data science, focusing on estimating grape yield, wine quality, and related outcomes based on environmental, agricultural, and chemical factors. This review summarizes key methodologies, variables, and trends found in the literature.

1. Environmental Factors in Wine Productivity

Numerous studies highlight the role of environmental factors—such as soil composition, temperature, rainfall, and sunlight exposure—in influencing grape productivity and wine quality. These factors can vary significantly across regions and even within the same vineyard.

Climate Influence: A major focus in many studies is how climate affects wine yield. For instance, Jones et al. (2005) examine the impact of temperature variations and seasonal rainfall patterns on grapevine phenology and productivity. Warmer regions typically lead to earlier ripening, but can stress the vines, reducing yield.

Soil and Topography: Soil health and topography are crucial. Studies such as White (2003) demonstrate how soil organic matter, drainage, and mineral content influence vine productivity, as these factors directly affect root development and water retention.

2. Agricultural Practices and Vineyard Management

Effective vineyard management practices play a key role in wine productivity. Many researchers investigate the relationship between vineyard interventions and yield.

Pruning and Canopy Management: Techniques such as pruning, trellising, and canopy management are shown to influence both grape quantity and quality. For example, Intrieri et al. (2008) assess how different pruning methods lead to variations in grape production and sugar content.

Irrigation Systems: Precision agriculture techniques, particularly irrigation, are central to managing water stress. Studies like Romero et al. (2010) show that regulated deficit irrigation can optimize yield while maintaining grape quality, a crucial factor in regions facing water scarcity.

3. Phenological and Biochemical Data

The biological characteristics of grapevines—such as vine age, grape variety, and disease resistance—affect productivity.

Vine Age and Grape Varieties: Vine age and grape variety have been linked to differing productivity outcomes. Older vines often produce lower yields but higher quality grapes, as seen in studies like Ramos et al. (2018).

Biochemical Markers: The presence of certain biochemical compounds in grapes (e.g., phenolics, sugars, and acids) are used to predict both wine yield and quality. Bellincontro et al. (2004) illustrate how monitoring these markers can provide insights into the optimal harvest period, which directly impacts wine productivity.

4. Remote Sensing and Data-Driven Approaches

The rise of precision viticulture has been driven by advances in remote sensing technologies, artificial intelligence, and machine learning.

Remote Sensing: Satellite imagery, drones, and ground sensors have been employed to predict vine health and yield. Johnson et al. (2013) highlight how remote sensing data, combined with machine learning algorithms, can predict vineyard productivity by monitoring vegetative growth, soil moisture, and temperature.

Machine Learning and AI: Data-driven models have become increasingly prominent in wine productivity prediction. Machine learning techniques, such as artificial neural networks and random forest models, are used to analyze historical vineyard data to predict future yields. Tardaguila et al. (2020) use a combination of remote sensing and machine learning to create predictive models with high accuracy.

5. Climate Change and Its Impact on Wine Productivity

The growing awareness of climate change has brought attention to the long-term sustainability of wine production.

Shifting Growing Conditions: Studies by Mozell & Thach (2014) discuss how climate change will affect vineyard productivity, particularly in regions where temperature increases may stress traditional growing practices.

Adaptation Strategies: Adaptive measures, such as changing grape varieties, implementing more efficient water management techniques, and using new technologies to monitor environmental factors, are proposed as ways to mitigate the negative impacts of climate change on wine production.

6. Yield and Quality Trade-Off

A central challenge in wine production is balancing yield with quality. High yields often correlate with diluted grape quality, while lower yields can enhance concentration and flavor.

Yield-Quality Models: Research such as Poni et al. (2009) focuses on developing models that optimize this balance, showing that regulated management of vine vigor can help achieve an optimal yield-quality ratio.

❖ DATA SET:

fixed acids	volatile acids	citric acid	residual sugar	chlorides	free sulfur	total sulfur	density	pH	sulphates	alcohol	quality
7.4	0.7	0	1.9	0.076	11	34	0.9978	3.51	0.56	9.4	5
7.8	0.88	0	2.6	0.098	25	67	0.9968	3.2	0.68	9.8	5
7.8	0.76	0.04	2.3	0.092	15	54	0.997	3.26	0.65	9.8	5
11.2	0.28	0.56	1.9	0.075	17	60	0.998	3.16	0.58	9.8	6
7.4	0.7	0	1.9	0.076	11	34	0.9978	3.51	0.56	9.4	5
7.4	0.66	0	1.8	0.075	13	40	0.9978	3.51	0.56	9.4	5
7.9	0.6	0.06	1.6	0.069	15	59	0.9964	3.3	0.46	9.4	5
7.3	0.65	0	1.2	0.065	15	21	0.9946	3.39	0.47	10	7
7.8	0.58	0.02	2	0.073	9	18	0.9968	3.36	0.57	9.5	7
7.5	0.5	0.36	6.1	0.071	17	102	0.9978	3.35	0.8	10.5	5
6.7	0.58	0.08	1.8	0.097	15	65	0.9959	3.28	0.54	9.2	5
7.5	0.5	0.36	6.1	0.071	17	102	0.9978	3.35	0.8	10.5	5
5.6	0.615	0	1.6	0.089	16	59	0.9943	3.58	0.52	9.9	5
7.8	0.61	0.29	1.6	0.114	9	29	0.9974	3.26	1.56	9.1	5
8.9	0.62	0.18	3.8	0.176	52	145	0.9986	3.16	0.88	9.2	5
8.9	0.62	0.19	3.9	0.17	51	148	0.9986	3.17	0.93	9.2	5
8.5	0.28	0.56	1.8	0.092	35	103	0.9969	3.3	0.75	10.5	7
8.1	0.56	0.28	1.7	0.368	16	56	0.9968	3.11	1.28	9.3	5
7.4	0.59	0.08	4.4	0.086	6	29	0.9974	3.38	0.5	9	4
7.9	0.32	0.51	1.8	0.341	17	56	0.9969	3.04	1.08	9.2	6
8.9	0.22	0.48	1.8	0.077	29	60	0.9968	3.39	0.53	9.4	6
7.6	0.39	0.31	2.3	0.082	23	71	0.9982	3.52	0.65	9.7	5
7.9	0.43	0.21	1.6	0.106	10	37	0.9966	3.17	0.91	9.5	5
8.5	0.49	0.11	2.3	0.084	9	67	0.9968	3.17	0.53	9.4	5
6.9	0.4	0.14	2.4	0.085	21	40	0.9968	3.43	0.63	9.7	6
6.3	0.39	0.16	1.4	0.08	11	23	0.9955	3.34	0.56	9.3	5
7.6	0.41	0.24	1.8	0.08	4	11	0.9962	3.28	0.59	9.5	5
7.9	0.43	0.21	1.6	0.106	10	37	0.9966	3.17	0.91	9.5	5
7.1	0.71	0	1.9	0.08	14	35	0.9972	3.47	0.55	9.4	5
7.8	0.645	0	2	0.082	8	16	0.9964	3.38	0.59	9.8	6
6.7	0.675	0.07	2.4	0.089	17	82	0.9958	3.35	0.54	10.1	5

❖ IMPORTING LIBRARIES:

1. Import numpy as np
2. Import matplotlib.pyplot as plt
3. Import pandas as pd
4. Import seaborn as sns

Explain

1. Numpy:

- NumPy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays.
- The ancestor of NumPy, Numeric, was originally created by Jim Hugunin. NumPy targets the CPython reference implementation of Python, which is a non-optimizing bytecode interpreter.
- Mathematical algorithms written for this version of Python often run much slower than compiled equivalents. Using NumPy in Python gives functionality comparable to MATLAB since they are both interpreted, and they both allow the user to write fast programs as long as most operations work on arrays or matrices instead of scalars.

2. Matplotlib:

- Matplotlib is a comprehensive library in Python used for creating static, interactive, and animated visualizations.
- It provides a wide array of tools for generating plots, histograms, scatterplots, and more, making it a versatile tool for data visualization and analysis.
- Matplotlib's flexibility allows users to customize every aspect of their visualizations, from colors and line styles to annotations and axis formatting. Its integration with NumPy, SciPy, and Pandas makes it an essential tool in the data science and scientific computing communities.
- Whether for exploratory data analysis or publication-quality figures, Matplotlib empowers users to create informative and visually appealing graphics with ease.

3. Pandas:

- In computer programming, pandas are a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series.
- It is free software released under the three clause BSD license. “Panel data”, an econometrics term for multidimensional, structured data sets.

4. Seaborn:

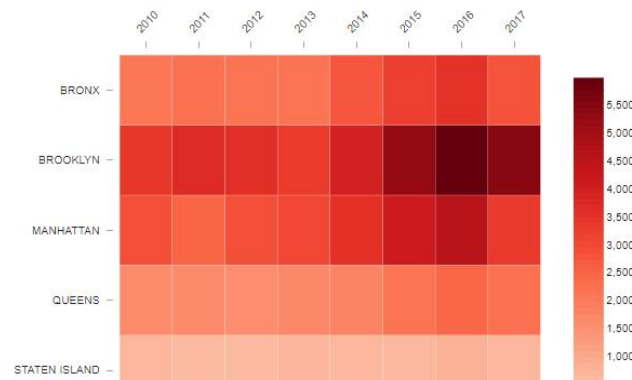
- Seaborn is a Python library for creating statistical graphics. It builds upon Matplotlib and closely integrates with Panda's data structures.
- We use it for Attractive Styles and Colour Palettes, Dataset-Oriented APIs, Categories of Plots (like Relational plots, Categorical plots, Distribution plots, Regression plots, and Matrix plots), Multi-plot grids, and Integration with Pandas.

5. Sklearn:

- Sklearn (scikit-learn) is a Python library that provides a wide range of unsupervised and supervised machine learning algorithms.
- It is built on top of SciPy and is one of the most widely used machine learning libraries. Sklearn offers tools for classification, regression, clustering, and dimensionality reduction, making it a powerful resource for data scientists and machine learning practitioners.
- There are lots of libraries under Sklearn e.g. model_selection, feature_selection, linear-model, ensemble, tree, svm, neural_network, metrics, neighbour etc. To work with linear regression, KNN Algorithm and SVM algorithm we need help of this kind of libraries.

❖ HEATMAP :

A **heatmap** is a great way to visualize the relationships or correlations between various features (attributes) of a dataset, particularly in wine prediction. For example, you can create a heatmap to show how the physicochemical properties (such as pH, alcohol content, residual sugar, etc.) are correlated with the wine quality rating. This helps in understanding which features most strongly affect the prediction of wine quality.



HEATMAP

Fig : 2

❖ DISTPLOT :

A **distplot** (distribution plot) is useful to visualize the distribution of data and observe trends in a specific variable, such as wine quality ratings or any other feature in a wine prediction model. By using distplots, you can quickly assess the distribution of features such as alcohol content, pH levels, or wine quality, and check if the data follows a normal distribution or is skewed, which can be important when developing predictive models.

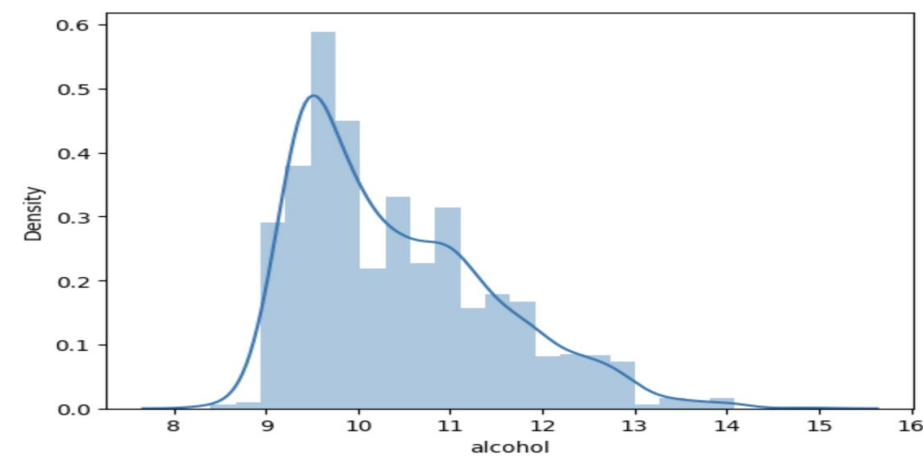


Fig: 3

❖ HISTOGRAM :

A **histogram** is another useful tool for visualizing the distribution of a single variable, such as wine quality, alcohol content, or any other feature within a wine prediction dataset. Unlike a distplot, which includes both a histogram and a density estimate, a histogram simply displays the count or frequency of data points that fall into each bin or range.

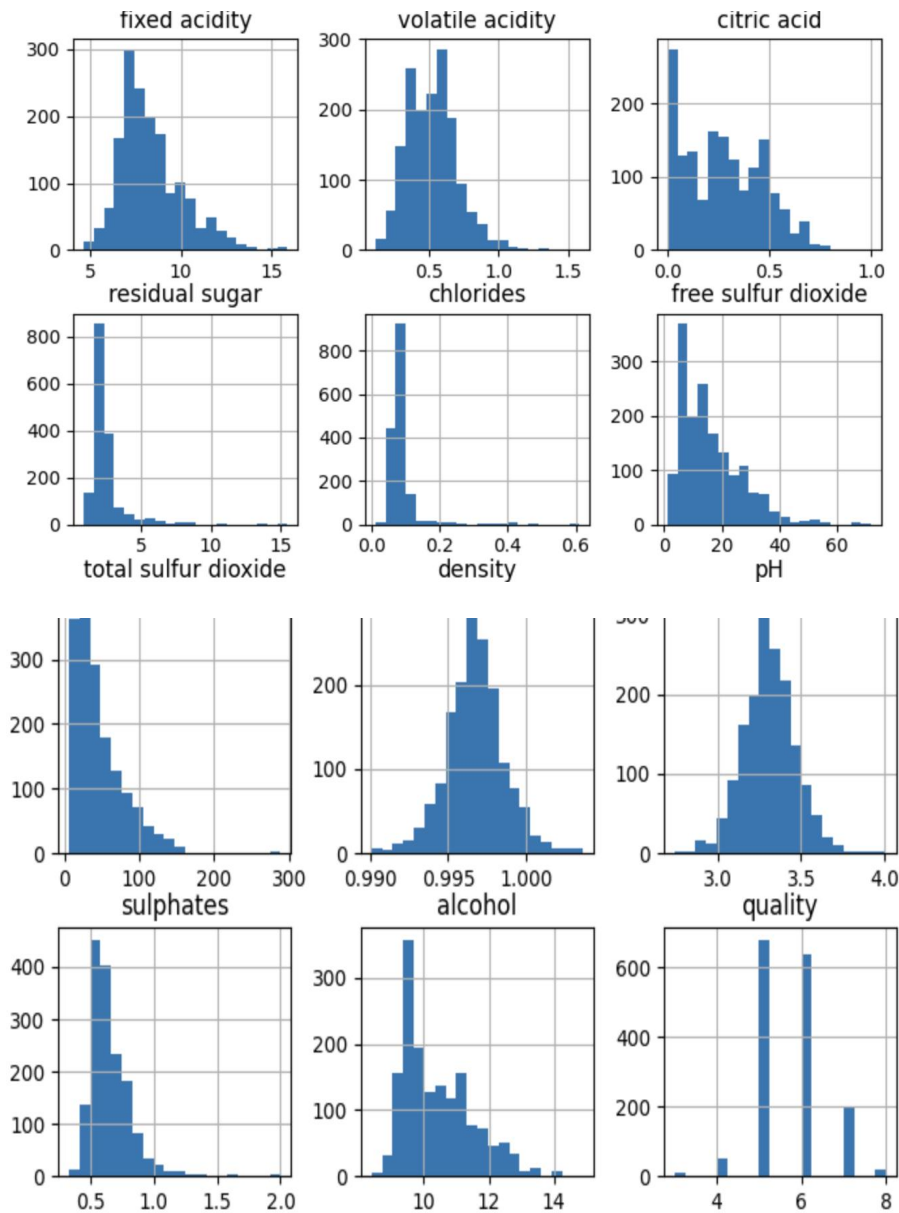
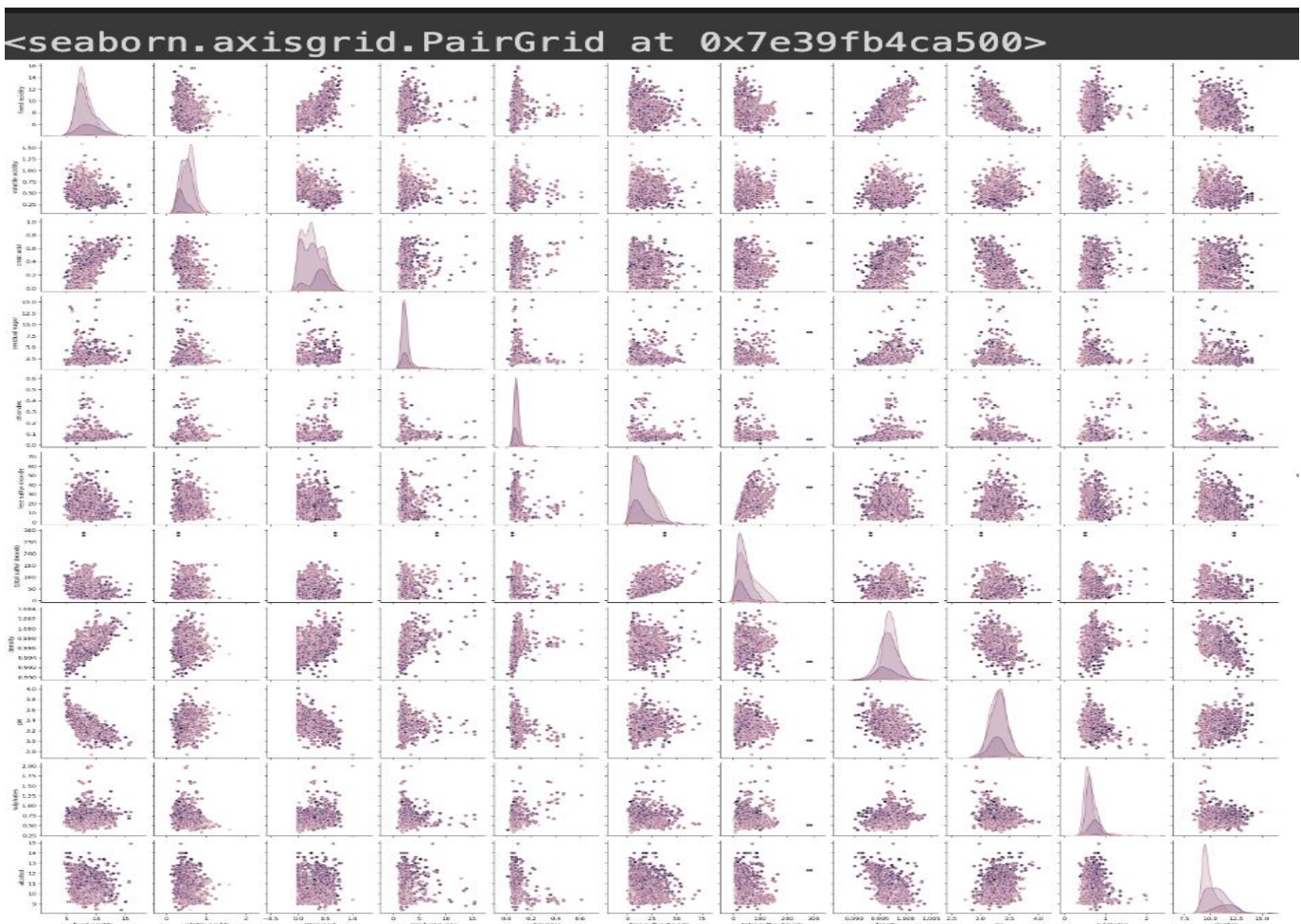


Fig: 4

A histogram can provide insights into how wine quality is distributed across samples. If the quality ratings are skewed or concentrated in certain ranges, this can influence the design of your predictive models.

❖ PAIR PLOT :

A **pairplot** is a powerful visualization tool that displays pairwise relationships between features in a dataset, especially useful when working with wine prediction datasets. It allows you to visualize correlations, patterns, or clusters between variables, such as how features like alcohol content, acidity, or residual sugar relate to wine quality. It also helps you spot outliers and understand the distribution of each feature.



Pair Plot Diagram:

Fig : 5

❖ Logistic Regression

Logistic regression is a statistical method used for binary classification problems, where the outcome variable is categorical and typically takes on two values (e.g., 0 and 1, true and false, yes and no). It predicts the probability that a given input point belongs to a particular category.

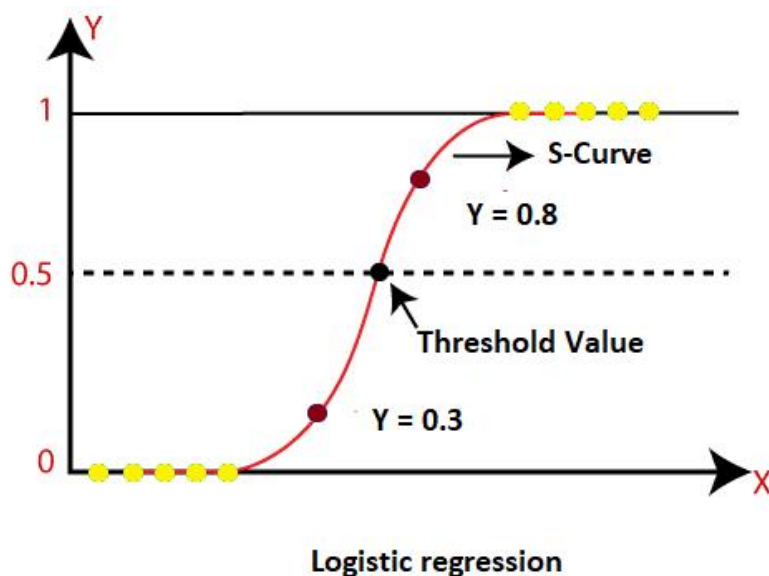


Fig: 6

Key Concepts of Logistic Regression:

1. Logistic Function (Sigmoid Function):-

The core of logistic regression is the logistic function, which maps any real-valued number into a value between 0 and 1. The formula is:

$$P(Y=1 | X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

Here, $P(Y=1 | X)$ is the probability that the dependent variable (Y) equals 1 given predictors (X) , and $(\beta_0, \beta_1, \dots, \beta_n)$ are the coefficients of the model.

2. Binary Outcome:-

Logistic regression is primarily used when the target variable is binary (e.g., spam vs. not spam). However, it can be extended to multiclass problems using techniques like one-vs-all or softmax regression.

3.Odds and Log-Odds:-

Odds represent the ratio of the probability of the event occurring to the probability of it not occurring.

Log-Odds (logit) is the natural logarithm of the odds and is a linear combination of the input features.

4.Training the Model:-

The coefficients β are estimated using the method of maximum likelihood estimation (MLE), which finds the parameters that maximize the likelihood of the observed data.

5.Threshold for Classification:-

A threshold (commonly 0.5) is used to convert probabilities into class labels. If $P(Y=1 | X)$ is greater than or equal to the threshold, the output is predicted as 1; otherwise, it is predicted as 0.

❖ KNN :

K-Nearest Neighbors (KNN) is a simple, versatile, and effective algorithm used for classification and regression tasks in machine learning. It is a type of instance-based learning, where the model makes predictions based on the closest training examples in the feature space.

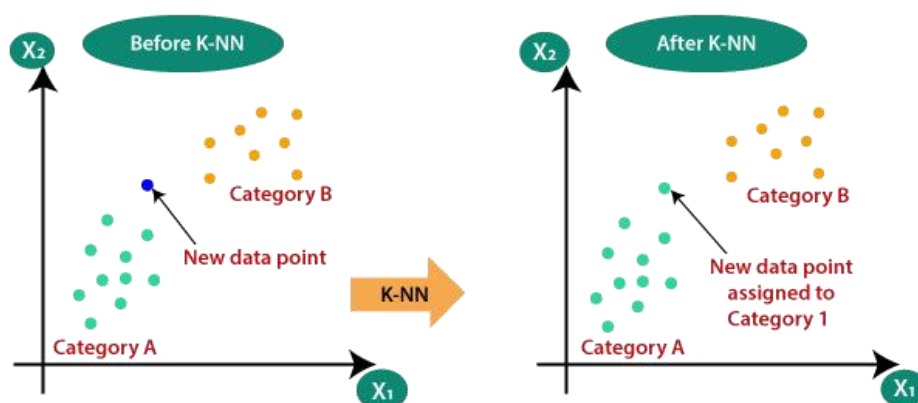


Fig: 7

◆ How does K-NN work?

The K-NN working can be explained on the basis of the below algorithm:

- **Step-1:** Select the number K of the neighbors
- **Step-2:** Calculate the Euclidean distance of **K number of neighbors**
- **Step-3:** Take the K nearest neighbors as per the calculated Euclidean distance.
- **Step-4:** Among these k neighbors, count the number of the data points in each category.
- **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.
- **Step-6:** Our model is ready.

Suppose we have a new data point and we need to put it in the required category. Consider the below image:

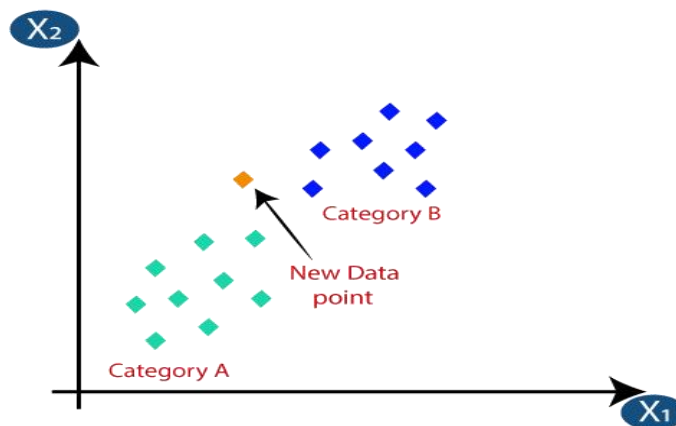


Fig: 8

- Firstly, we will choose the number of neighbors, so we will choose the $k=5$.
- Next, we will calculate the Euclidean distance between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as:

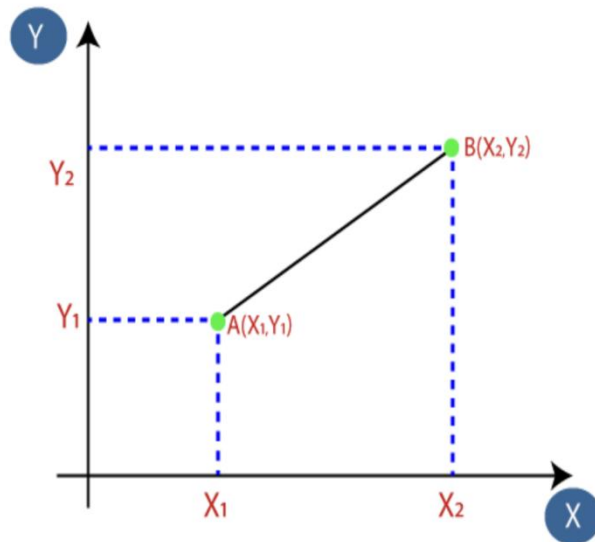


Fig : 9

By calculating the Euclidean distance we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:

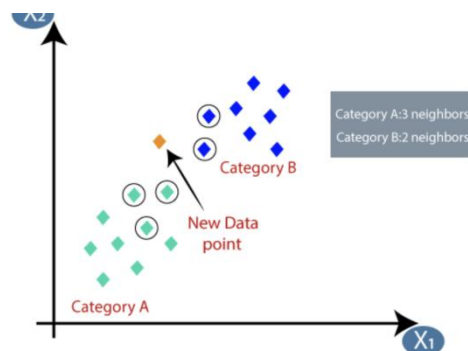


Fig: 10

- As we can see the 3 nearest neighbors are from category A, hence this new data point must belong to category A.

❖ GAUSSIAN NAIVE BAYES :

Gaussian Naive Bayes is a variant of the Naive Bayes algorithm for classification that assumes the features follow a Gaussian (normal) distribution. It is particularly useful for problems with continuous data where we can assume that the underlying distributions of features are Gaussian.

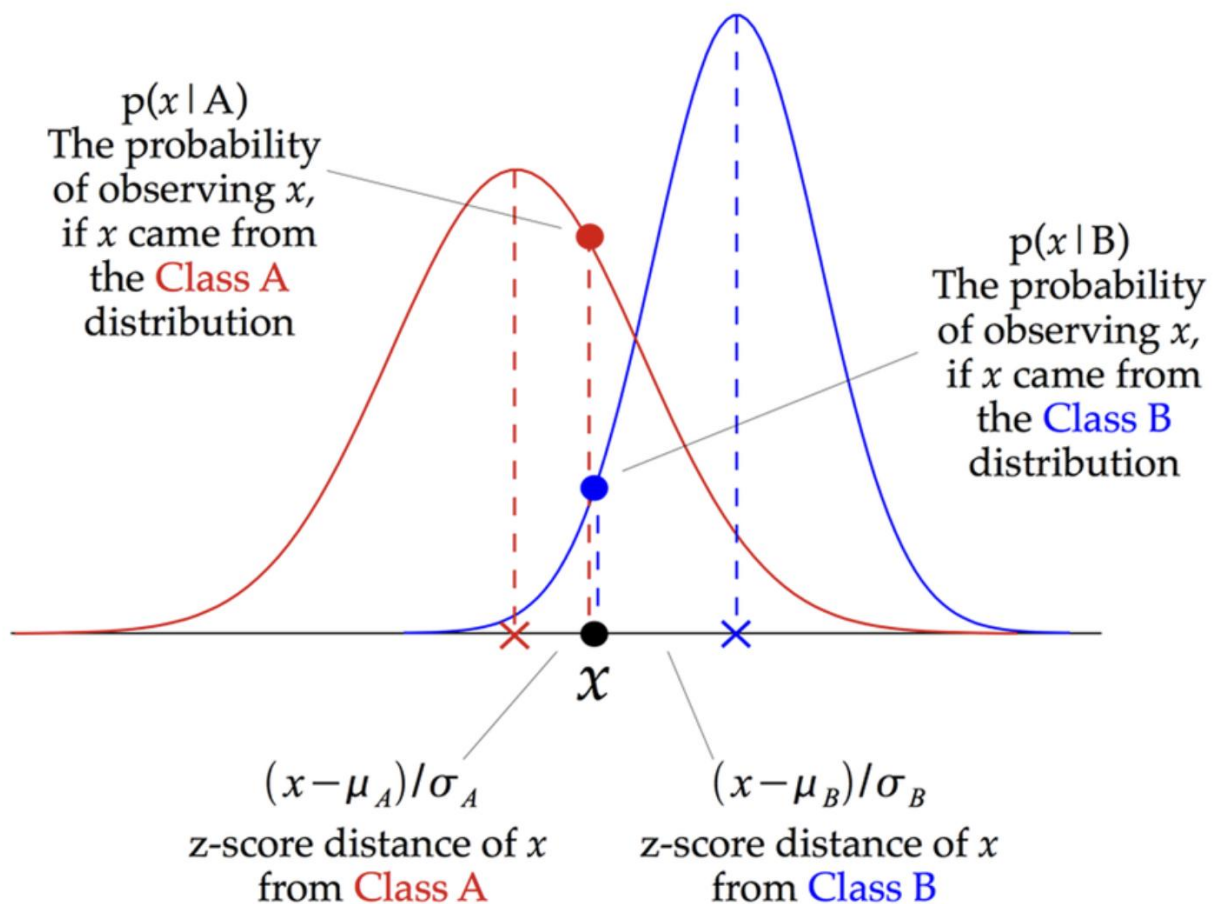


Fig : 11

Key Concepts of Gaussian Naive Bayes :-

1. Naive Bayes Classifier :-

Naive Bayes classifiers are based on Bayes' theorem, which describes the probability of a class given the features. The "naive" assumption refers to the idea that the features are independent of each other, given the class label. This simplifies calculations significantly, even though this assumption may not hold in reality.

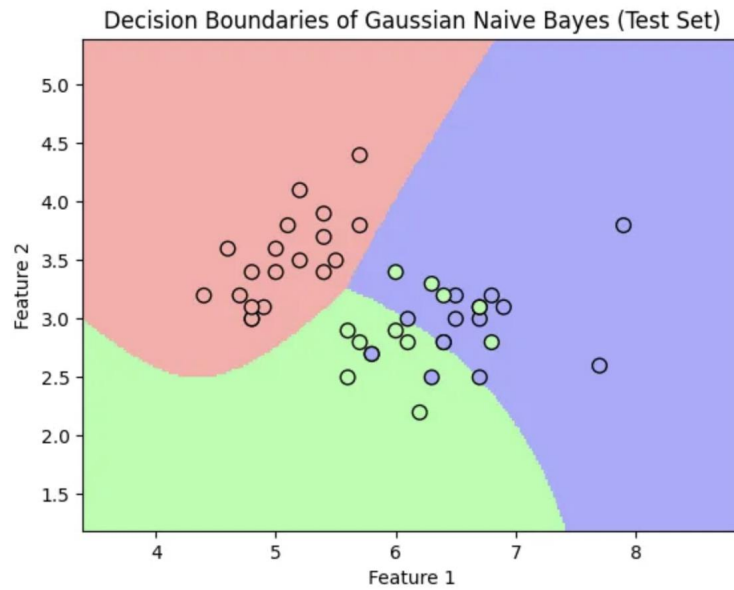


Fig: 12

2. Gaussian Distribution :-

In Gaussian Naive Bayes, the features are assumed to be normally distributed within each class. The probability density function of a Gaussian distribution is defined as:

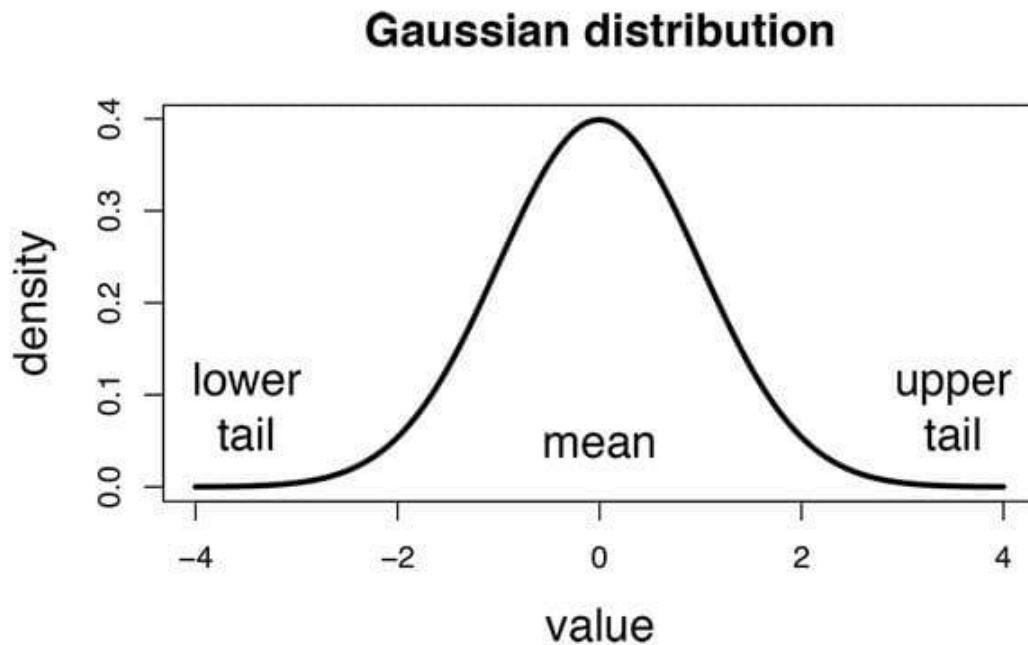


Fig : 13

\[

$$P(x | y) = \frac{1}{\sqrt{2 \pi \sigma^2}} e^{-\frac{(x - \mu)^2}{2 \sigma^2}}$$

\]

where:

- $P(x | y)$ is the probability of feature (x) given class (y) .
- (μ) is the mean of the feature values for class (y) .
- (σ^2) is the variance of the feature values for class (y) .

3. Bayes' Theorem :-

The classifier uses Bayes' theorem to calculate the posterior probability for each class given the observed feature values:

\[

$$P(y | X) = \frac{P(X | y) P(y)}{P(X)}$$

\]

where $P(y)$ is the prior probability of class (y) , and $P(X | y)$ is the likelihood of observing features (X) given class (y) .

4. Classification :-

For a given instance, Gaussian Naive Bayes calculates the posterior probability for each class using Bayes' theorem and then assigns the class with the highest probability.

✧ Importing Libraries :

CODE

```
[1] import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns

from warnings import filterwarnings
filterwarnings(action='ignore')
```

✧ Loading Data Set :

```
data=pd.read_csv("/content/winequality-red.csv")
print(data)
```

```
fixed acidity volatile acidity citric acid residual sugar chlorides \
0      7.4      0.700      0.00      1.9      0.076
1      7.8      0.880      0.00      2.6      0.098
2      7.8      0.760      0.04      2.3      0.092
3     11.2      0.280      0.56      1.9      0.075
4      7.4      0.700      0.00      1.9      0.076
...
1594    6.2      0.600      0.08      2.0      0.090
1595    5.9      0.550      0.10      2.2      0.062
1596    6.3      0.510      0.13      2.3      0.076
1597    5.9      0.645      0.12      2.0      0.075
1598    6.0      0.310      0.47      3.6      0.067

free sulfur dioxide total sulfur dioxide density pH sulphates \
0      11.0      34.0 0.99780 3.51 0.56
1      25.0      67.0 0.99680 3.20 0.68
2      15.0      54.0 0.99700 3.26 0.65
3      17.0      60.0 0.99800 3.16 0.58
4      11.0      34.0 0.99780 3.51 0.56
...
1594    32.0      44.0 0.99490 3.45 0.58
1595    39.0      51.0 0.99512 3.52 0.76
1596    29.0      40.0 0.99574 3.42 0.75
1597    32.0      44.0 0.99547 3.57 0.71
1598    18.0      42.0 0.99549 3.39 0.66

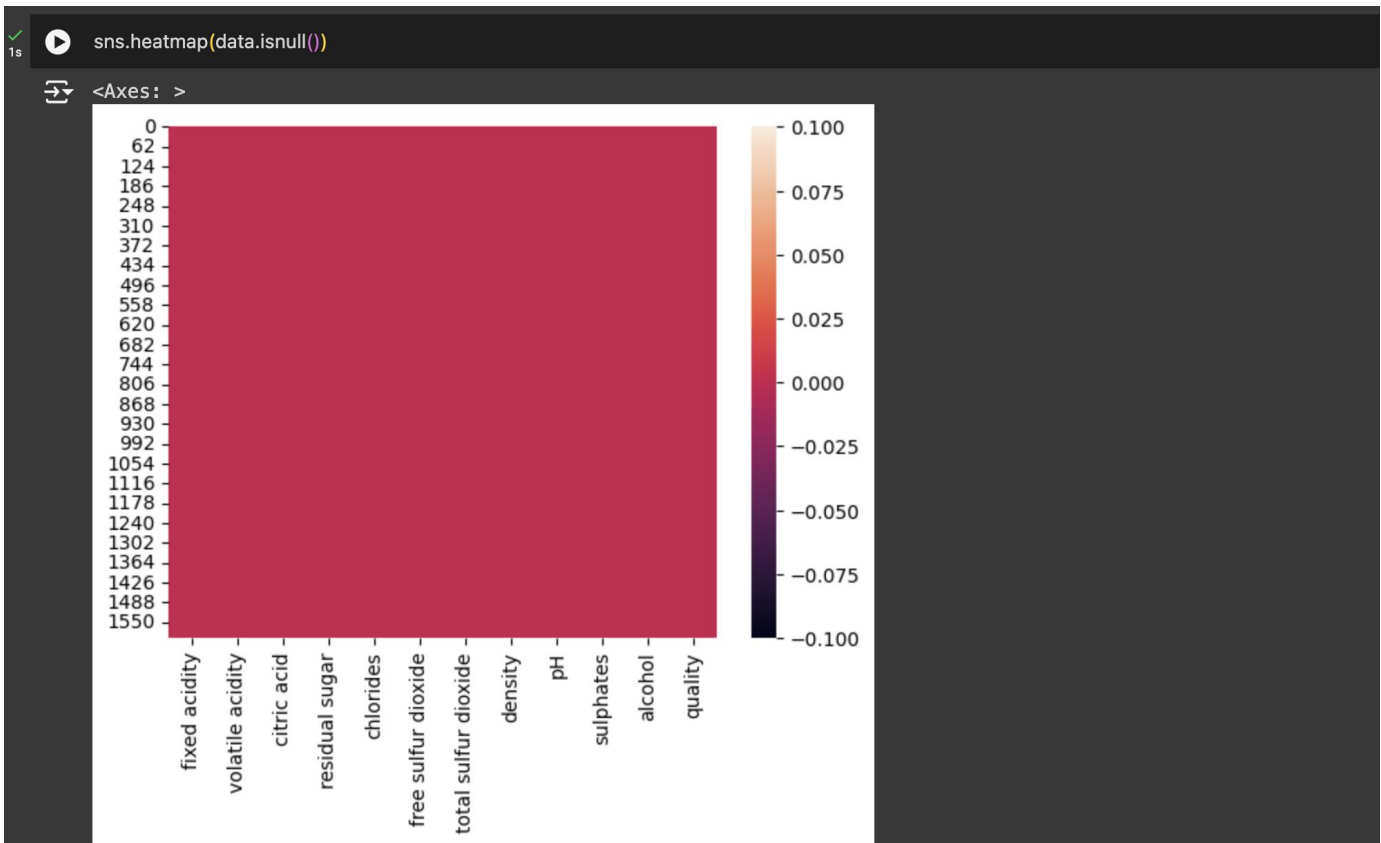
alcohol quality
0      9.4      5
1      9.8      5
2      9.8      5
3      9.8      6
4      9.4      5
...
1594    10.5      5
1595    11.2      6
1596    11.0      6
1597    10.2      5
1598    11.0      6

[1599 rows x 12 columns]
```

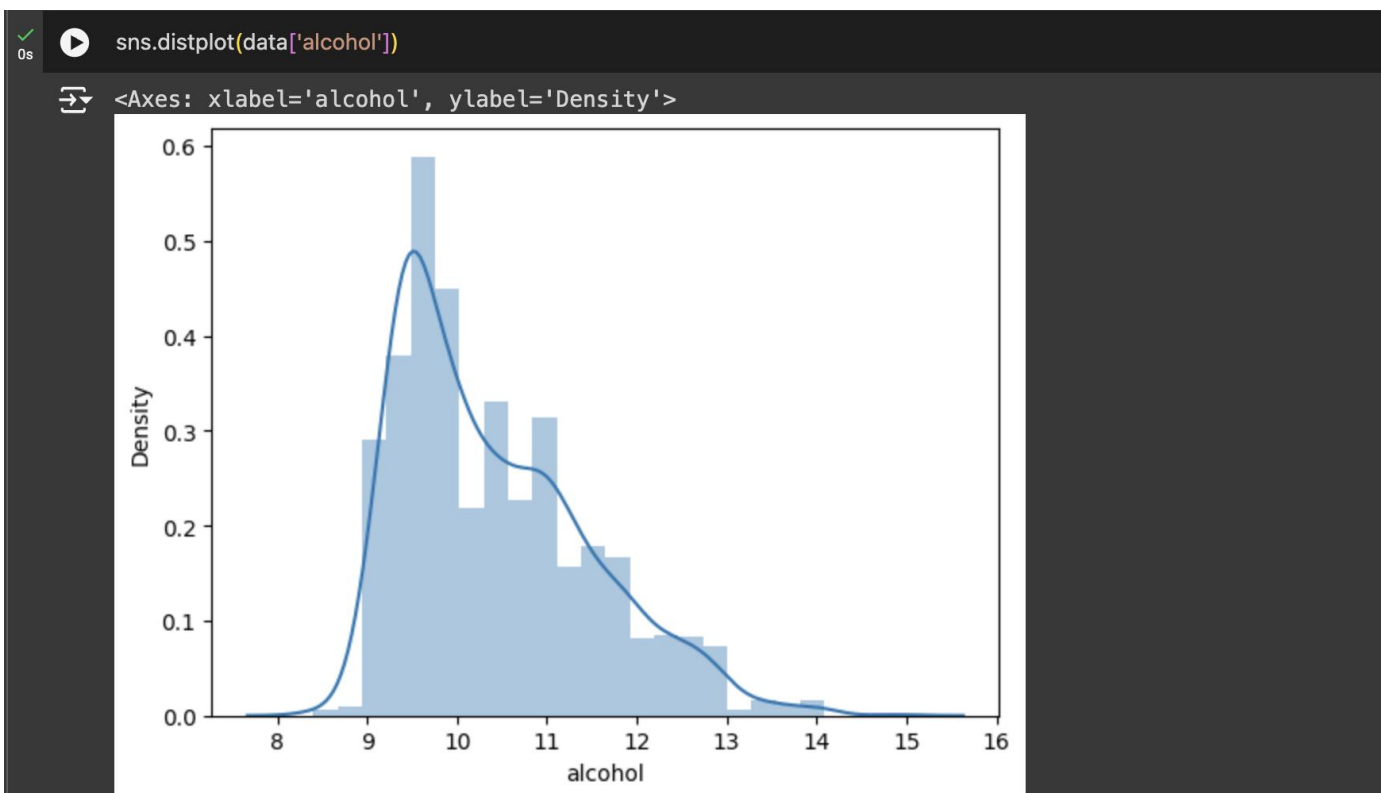
```
data.describe().T
```

	count	mean	std	min	25%	50%	75%	max
fixed acidity	1599.0	8.319637	1.741096	4.60000	7.1000	7.90000	9.200000	15.90000
volatile acidity	1599.0	0.527821	0.179060	0.12000	0.3900	0.52000	0.640000	1.58000
citric acid	1599.0	0.270976	0.194801	0.00000	0.0900	0.26000	0.420000	1.00000
residual sugar	1599.0	2.538806	1.409928	0.90000	1.9000	2.20000	2.600000	15.50000
chlorides	1599.0	0.087467	0.047065	0.01200	0.0700	0.07900	0.090000	0.61100
free sulfur dioxide	1599.0	15.874922	10.460157	1.00000	7.0000	14.00000	21.000000	72.00000
total sulfur dioxide	1599.0	46.467792	32.895324	6.00000	22.0000	38.00000	62.000000	289.00000
density	1599.0	0.996747	0.001887	0.99007	0.9956	0.99675	0.997835	1.00369
pH	1599.0	3.311113	0.154386	2.74000	3.2100	3.31000	3.400000	4.01000
sulphates	1599.0	0.658149	0.169507	0.33000	0.5500	0.62000	0.730000	2.00000
alcohol	1599.0	10.422983	1.065668	8.40000	9.5000	10.20000	11.100000	14.90000
quality	1599.0	5.636023	0.807569	3.00000	5.0000	6.00000	6.000000	8.00000

✧ Heatmap :



✧ Displot :

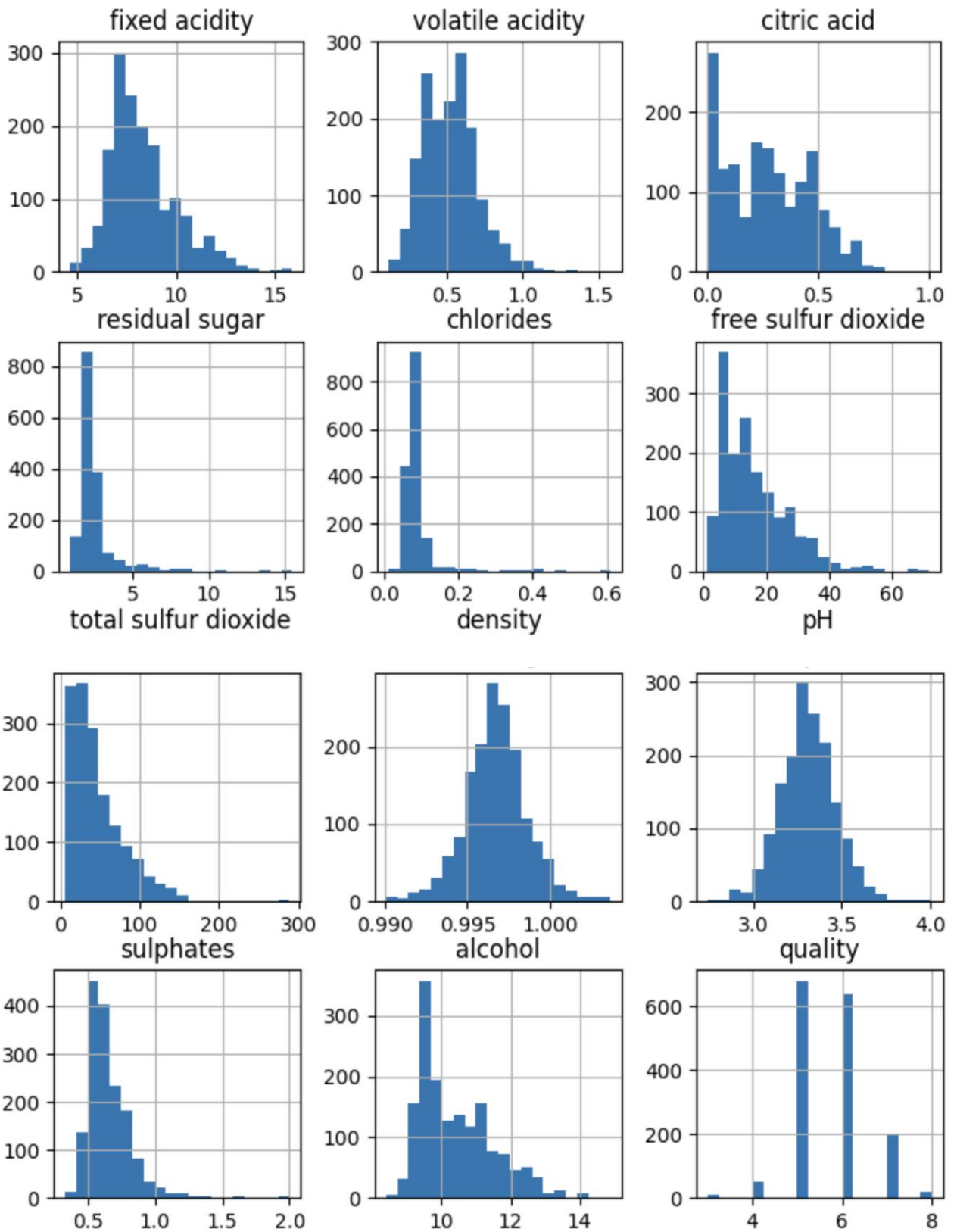


✧ Histogram :

3s

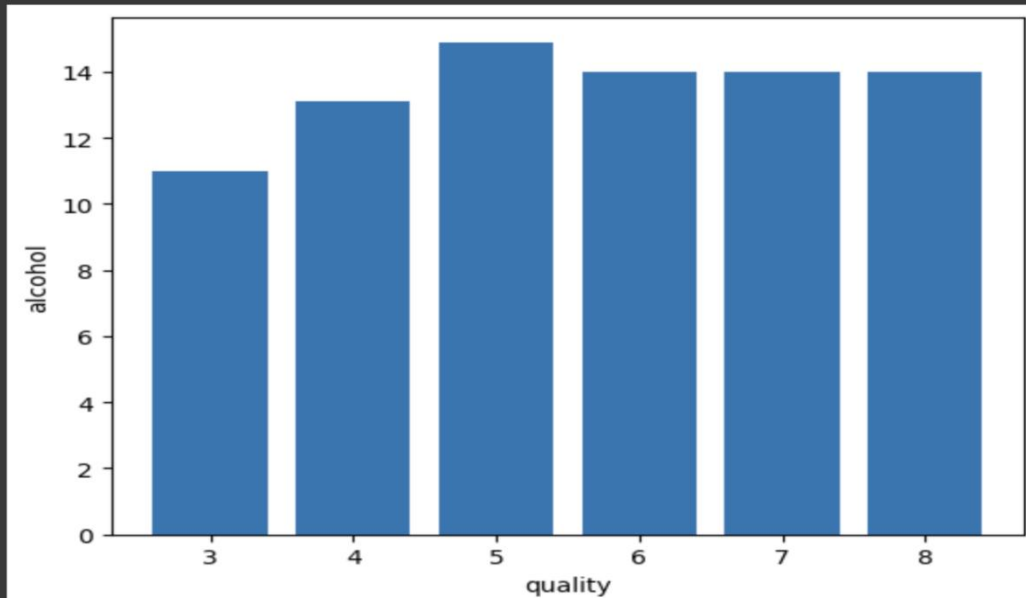


```
data.hist(figsize=(8,10),bins=20)  
plt.show()
```



✧ Barplot :

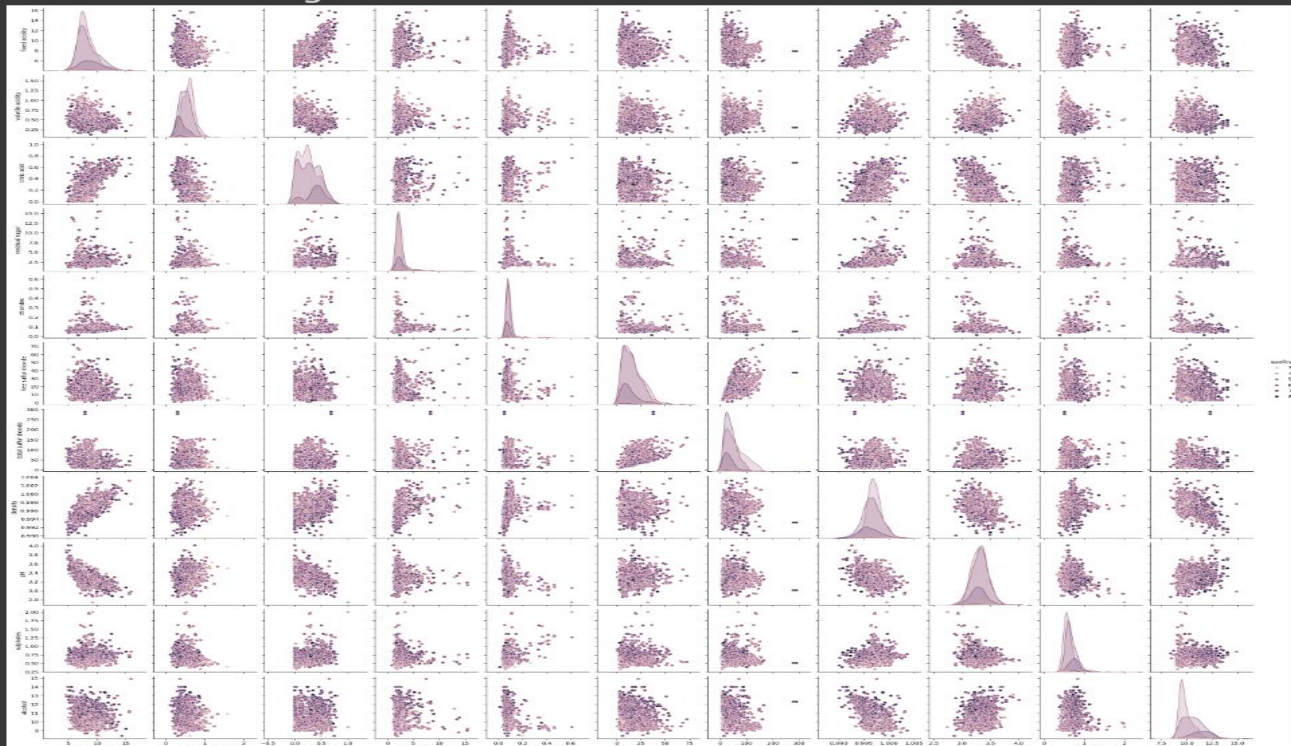
```
[7] plt.bar(data['quality'], data['alcohol'])
plt.xlabel('quality')
plt.ylabel('alcohol')
plt.show()
```



✧ Pairplot :

```
sns.pairplot(data, hue="quality")
```

<seaborn.axisgrid.PairGrid at 0x7e39fb4ca500>



```
✓ 0s [9] X=data.iloc[:,:].values  
      print("Row and Columns:",X.shape)
```

⇒ Row and Columns: (1599, 12)

```
✓ 0s [9] X = data.drop('quality',axis=1)  
      Y=data['quality']  
      print(Y)
```

⇒

0	5
1	5
2	5
3	6
4	5
..	
1594	5
1595	6
1596	6
1597	5
1598	6

Name: quality, Length: 1599, dtype: int64

✧ *Splitting Data Set :*

```
✓ 0s [23] from sklearn.model_selection import train_test_split  
      X_train, X_test, Y_train, Y_test = train_test_split(X,Y,test_size=0.4)  
      print(X_train.shape)
```

⇒ (959, 11)

✧ *Logistic Regression :*

```
✓ 0s [13] from sklearn.linear_model import LogisticRegression  
      model = LogisticRegression()  
      model.fit(X_train,Y_train)  
      Y_predict = model.predict(X_test)  
  
      from sklearn.metrics import accuracy_score,confusion_matrix  
      print("Accuracy Score:",accuracy_score(Y_test,Y_predict))
```

⇒ Accuracy Score: 0.559375

✧ *Classification Report :*

```
from sklearn.metrics import classification_report
CR=classification_report(Y_test,Y_predict)
print(CR)
```

```

      precision  recall  f1-score  support
3      0.00      0.00      0.00         4
4      0.00      0.00      0.00        23
5      0.65      0.70      0.67       273
6      0.49      0.68      0.57       247
7      0.00      0.00      0.00        90
8      0.00      0.00      0.00         3

 accuracy              0.56       640
 macro avg           0.19      0.23      0.21       640
 weighted avg        0.46      0.56      0.51       640
```

✧ *Confusion Matrix :*

```
✓ [15] confusion_mat = confusion_matrix(Y_test,Y_predict)
0s print(confusion_mat)
```

```

[[ 0  0  4  0  0  0]
 [ 0  0 12 11  0  0]
 [ 0  0 191 82  0  0]
 [ 0  0  78 167  2  0]
 [ 0  0  10  80  0  0]
 [ 0  0  0  3  0  0]]
```

✧ *KNN :*

```
✓ from sklearn.neighbors import KNeighborsClassifier
0s model = KNeighborsClassifier(n_neighbors=7)
model.fit(X_train,Y_train)
y_predict = model.predict(X_test)

from sklearn.metrics import accuracy_score
print("Accuracy Score: ",accuracy_score(Y_test,y_predict))
```

```

Accuracy Score: 0.4765625
```

✧ *GaussianNB* :

```
✓ [17] from sklearn.naive_bayes import GaussianNB  
0s      model_NB = GaussianNB()  
      model_NB.fit(X_train,Y_train)  
      y_pred_NB = model_NB.predict(X_test)  
  
      from sklearn.metrics import accuracy_score  
      print("Accuracy Score: ",accuracy_score(Y_test,y_pred_NB))
```

⇄ Accuracy Score: 0.525

CONCLUSION :

This literature survey highlights the evolution of crop prediction methodologies from traditional statistical models to modern machine learning and remote sensing-based approaches. The incorporation of large, real-time datasets through IoT devices and satellite imagery has significantly enhanced the accuracy and timeliness of crop yield predictions. As agriculture moves toward data-driven decision-making, the continued development of more integrated and scalable predictive models will be essential to address the challenges posed by climate change, resource constraints, and growing global food demand.

References

- ✓ • Jones, G.V., White, M.A., Cooper, O.R., & Storchmann, K. (2005). Climate change and global wine quality. *Climatic Change*, 73(3), 319–343.
- ✓ <https://link.springer.com/article/10.1007/s10584-005-4704-2>
- ✓ • Intrieri, C., Poni, S., & Silvestroni, O. (2008). Canopy management and vine productivity. *Precision Agriculture*, 9, 305–322.
- ✓ • Johnson, L.F., Roczen, D.E., Youkhana, S.K., Nemani, R.R., & Bosch, D.F. (2013). Mapping vineyard leaf area with multispectral satellite imagery. *Remote Sensing of Environment*, 71, 39–46.
- ✓ • Tardaguila, J., Diago, M.P., & Fernández-Novales, J. (2020). Advanced sensor technologies and data-driven models in viticulture: Current status and future directions. *Computers and Electronics in Agriculture*, 168, 105097.

THANK YOU

