# Red Wine Quality Analysis: Presentation of Findings

- Bhavya Reddy Patlolla

# Agenda

- Exploring Data
- Visualizations
- Identifying Additional Data/Variables
- Summary of Analysis

# Exploring Data

1. Loaded 'Red Wine Quality' data using pd.read_csv(), initiating the analysis process.

2. Utilized descriptive statistics (mean, median, etc.) for better comprehension of data distribution and central tendencies.

3. Employed a correlation matrix to pinpoint features strongly correlating with 'quality' and each other.

4. Descriptive stats and correlation matrix identified key features, aiding in hypothesis generation for wine quality factors.

5. These statistical measures provided a robust foundation for subsequent analysis, model building, and decision-making in the project.

# Visualizations - Scatter Plot

1. Worked on a scatter plot visual to explore relationships between variables.

2. Utilized the plot to reveal different features during data analysis.

3. Chose to represent alcohol quality in the scatter plot to uncover linear or non-linear patterns.

```
In [12]: #Visualizations
         #scatter plot visual
         plt.figure(figsize=(8, 6))
         sns.scatterplot(x='alcohol', y='quality', data=wine, alpha=0.5)
         plt.title("Scatter Plot: Alcohol vs. Quality")
         plt.show()
```

# Violin plot Visual

- Employed a violin plot to depict the distribution of sulphate quality levels.

- Chose this plot for its clear display of detailed feature distributions.

- Gained insights into the range, central tendency, and shape of feature distributions.

- Utilized the plot to assess variations across wine quality levels.

- Investigated for significant differences in feature distributions between different quality categories.

```
[14]: #Violin plot visual
plt.figure(figsize=(10, 6))
sns.violinplot(x='quality', y='sulphates', data=wine, palette="Blues")
plt.title("Sulphates Distribution by Quality")
plt.show()
```



Sulphates Distribution by Quality

# Histogram Visual

o Created a histogram visual to display the distribution of alcohol content in the dataset.

o Used the histogram to represent the frequency of alcohol content occurrences.

o Utilized the plot to gain insights into the shape of the feature's distribution.

o Examined central tendencies of the alcohol content distribution through the histogram.

o Investigated the presence of outliers to understand characteristics and their association with wine quality.

```
In [48]: # histogram visual
sns.countplot(data=wine, x='quality', palette='coolwarm')
plt.title('Number of wines in each quality category')
plt.show()
```

# Regression Plot Visual

```
In [19]: #regression plot visual
         sns.regplot(x='alcohol', y='quality', data=wine, line_kws={"color": "red"})

Out[19]: <Axes: xlabel='alcohol', ylabel='quality'>
```



1. Explored the relationship between alcohol content and wine quality using the red wine quality dataset.

2. Utilized a regression plot to visually represent the relationship between the two variables.

3. Determined whether a positive or negative correlation exists between alcohol content and wine quality.

4. Evaluated the strength and direction of the feature and wine quality relationship.

5. Incorporated a regression line to illustrate the trend and facilitate predictions, aiding in understanding how changes in the predictor variable impact the target variable.

# Heat Map Visual

- Visualized the correlation matrix as a heatmap for enhanced representation.

- Used color coding to distinguish and signify the strength and direction of correlations.

- Utilized the heatmap to identify both strong and weak correlations among features.

- Applied the heatmap as a valuable tool for guiding feature selection.

- Leveraged the plot to uncover patterns and inform subsequent analysis and modeling decisions in the exploration of the red wine quality dataset.

# Swarm Plot Visual

- Utilized a Swarm Plot to showcase the distribution of alcohol content in various quality categories within the wine dataset.

- Achieved a detailed analysis of variations in alcohol content within each quality category through the plot's point-by-point correspondence with individual wine samples.
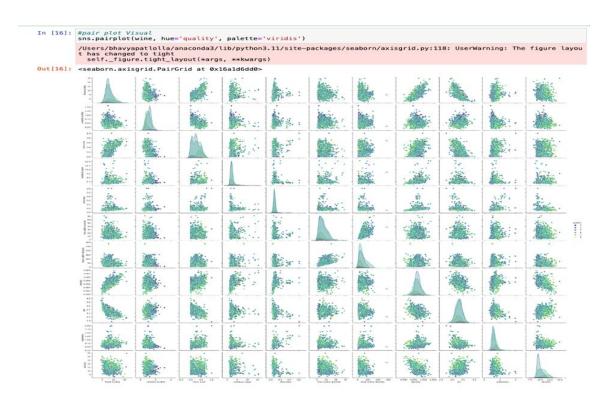
- Found the Swarm Plot particularly effective in identifying patterns and trends in the relationship between alcohol content and wine quality.

- The visualization is especially user-friendly, aiding in the easy comprehension of the dataset's complexities.

- Enhanced the visual presentation by using a "viridis" color scheme, providing a clear display of alcohol levels for wines across different quality classifications.

# Pair Plot Visual

o Employed a pair plot to gain an in-depth understanding of connections between various features in the wine dataset.

o Diagonal subplots were utilized to display the distribution of individual components within the dataset.

o Each subplot presented a scatterplot, illustrating the interaction between pairs of variables in the dataset.

o Enhanced visual clarity by implementing the 'viridis' color scheme, facilitating the distinction between different groups and aiding in pattern recognition.

o Identified two key benefits of using the pair plot: investigating correlations between variables and gaining insights into how different dataset attributes interact with each other.

# Pie Chart Visual

```
In [17]: #Pie chart Visual
         wine['quality'].value_counts().plot.pie(autopct='%1.1f%%', startangle=90, figsize=(8, 8), colormap='viridis')

Out[17]: <Axes: ylabel='count'>
```



- o Utilized a pie chart to represent the distribution of wine quality categories throughout the dataset.

- o Each wedge in the pie chart corresponds to a different quality level.

- o Emphasized the 'high-quality' category by subtly separating it using the explode parameter.

- o Applied the 'Viridis' color palette for visual appeal and clarity in distinguishing categories.

- o Displayed percentage labels using the 'autopct' parameter, offering a concise and easily understandable summary of the categorical distribution of wine quality levels in the dataset.

# Identifying Additional Data/Variables



- o Identified key variables such as Grape Variety, Winery Information, Geographic Origin, Vintage Year, Wine Aging, Climatic Data, pH Levels, Yeast Strain, Sensory Evaluation, and Chemical Analysis.

- o Emphasized the importance of understanding grape variety, winery details, geographic origin, vintage year, wine aging, climatic data, pH levels, yeast strain, sensory evaluation, chemical analysis, food pairings, market and pricing information, and consumer reviews.

- o Highlighted how these variables can significantly influence wine quality and contribute to a more thorough analysis.

- o Acknowledged that the availability of such data may vary and suggested careful consideration of including additional variables based on project objectives.

- o Concluded that incorporating these additional data points could lead to a more comprehensive assessment of factors influencing red wine quality, ultimately aligning with the study's objectives.

# Summary of Analysis

- Conducted a comprehensive analysis of the red wine quality dataset, revealing valuable insights into its characteristics and potential influencing factors.

- Utilized statistical measures and visualizations for an in-depth understanding, highlighting key features and relationships contributing to wine quality.

- Identified the importance of considering additional variables for future research, enhancing the assessment of factors influencing red wine quality.

- Proposed integrating diverse data on grape variety, winery details, geographic origin, vintage year, wine aging, climatic data, pH levels, yeast strain, sensory evaluation, chemical analysis, food pairings, market dynamics, and consumer reviews.

- Established the analysis as a foundation for informed decision-making in red wine quality, suggesting potential solutions and future research directions.

# Conclusion

- The analysis lays the foundation for informed decision-making in red wine quality.

- Proposed solutions include incorporating identified variables and conducting further research to validate their impact.

- Advocated careful consideration of data availability, alignment with objectives, and a balanced approach to including diverse data points.

- The study serves as a model for approaching wine quality analysis, encouraging researchers to explore the complexities of winemaking for a comprehensive understanding.

Thank you