# Practical 4: EDA

## EDA

```
[14]: import pandas as pd
      import numpy as np
      import matplotlib.pyplot as plt
      import seaborn as sns
      import warnings as wr
      wr.filterwarnings('ignore')
```

```
[2]: df =pd.read_csv("winequality-red.csv")
```

```
[3]: df
```

```
[3]:
```

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7.4 | 0.700 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.99780 | 3.51 | 0.56 | 9.4 | 5 |
| 1 | 7.8 | 0.880 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.99680 | 3.20 | 0.68 | 9.8 | 5 |
| 2 | 7.8 | 0.760 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.99700 | 3.26 | 0.65 | 9.8 | 5 |
| 3 | 11.2 | 0.280 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.99800 | 3.16 | 0.58 | 9.8 | 6 |

```
[4]: df.shape
```

```
[4]: (1599, 12)
```

```
[5]: df.describe
```

```
[5]: <bound method NDFrame.describe of      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
```

```
[6]: df.isnull().sum()
```

```
[6]: fixed acidity           0
     volatile acidity        0
     citric acid             0
     residual sugar          0
     chlorides               0
     free sulfur dioxide     0
     total sulfur dioxide    0
     density                 0
     pH                      0
     sulphates               0
     alcohol                 0
     quality                 0
     dtype: int64
```

```
[7]: df.columns.tolist()
```
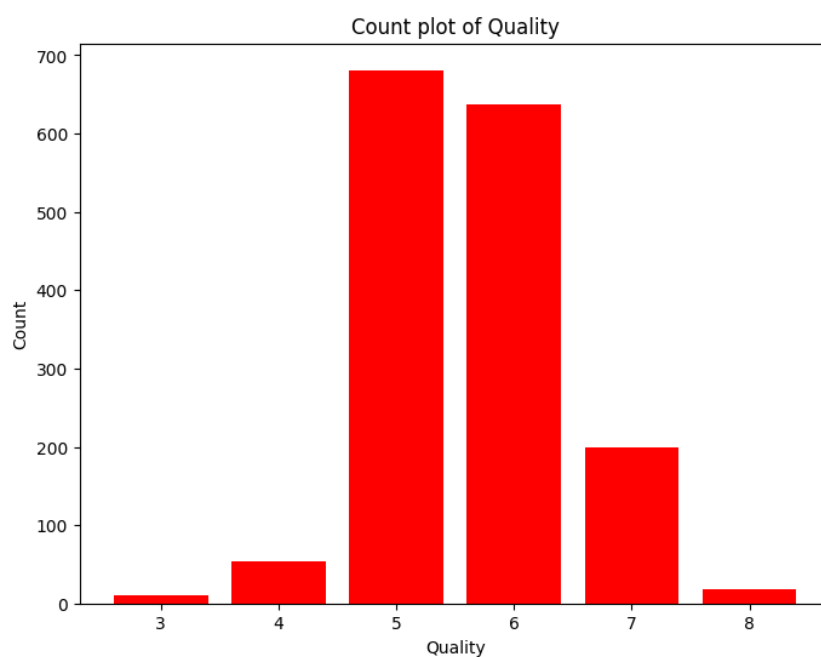
```
[7]: ['fixed acidity',
      'volatile acidity',
      'citric acid',
      'residual sugar',
      'chlorides',
      'free sulfur dioxide',
      'total sulfur dioxide',
      'density',
      'pH',
      'sulphates',
      'alcohol',
      'quality']
```

```
[8]: df.nunique()
```

```
[8]: fixed acidity          96
     volatile acidity      143
     citric acid            80
     residual sugar         91
     chlorides             153
     free sulfur dioxide    60
     total sulfur dioxide  144
     density               436
     pH                     89
     sulphates              96
     alcohol                65
     quality                 6
     dtype: int64
```

```
[9]: qc =df['quality'].value_counts()

     plt.figure(figsize=(8,6))
     plt.bar(qc.index,qc,color="red")
     plt.title("Count plot of Quality")
     plt.xlabel("Quality")
     plt.ylabel("Count")
     plt.show()
```
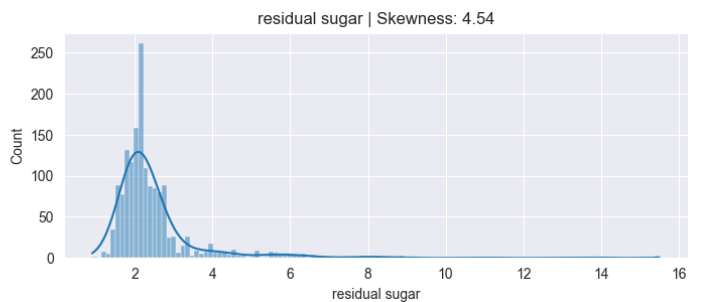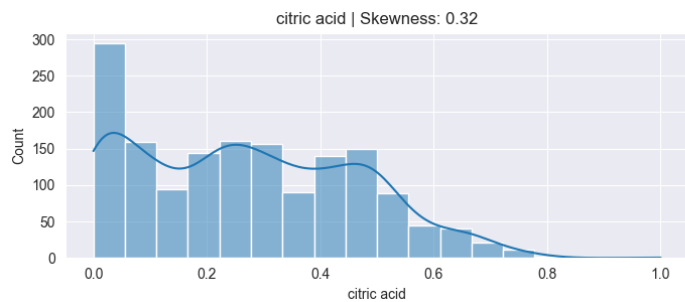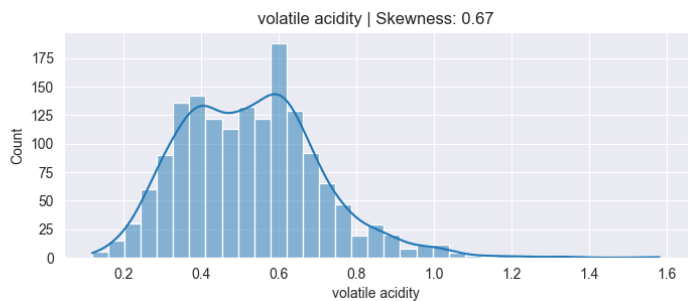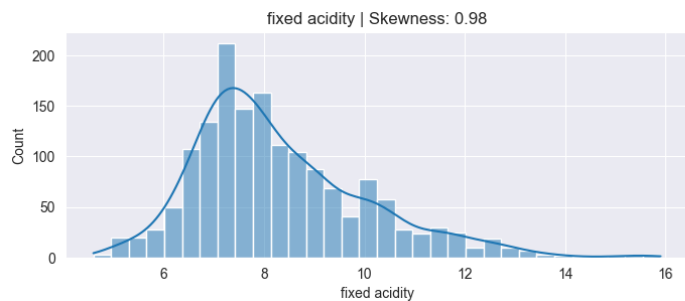


## Kernel Density Plots

```
[11]: sns.set_style("darkgrid")
      numerical_columns = df.select_dtypes(include=["int64", "float64"]).columns

      plt.figure(figsize=(14, len(numerical_columns) * 3))
      for idx, feature in enumerate(numerical_columns, 1):
       plt.subplot(len(numerical_columns), 2, idx)
       sns.histplot(df[feature], kde=True)
       plt.title(f"{feature} | Skewness: {round(df[feature].skew(), 2)}")

      plt.tight_layout()
      plt.show()
```
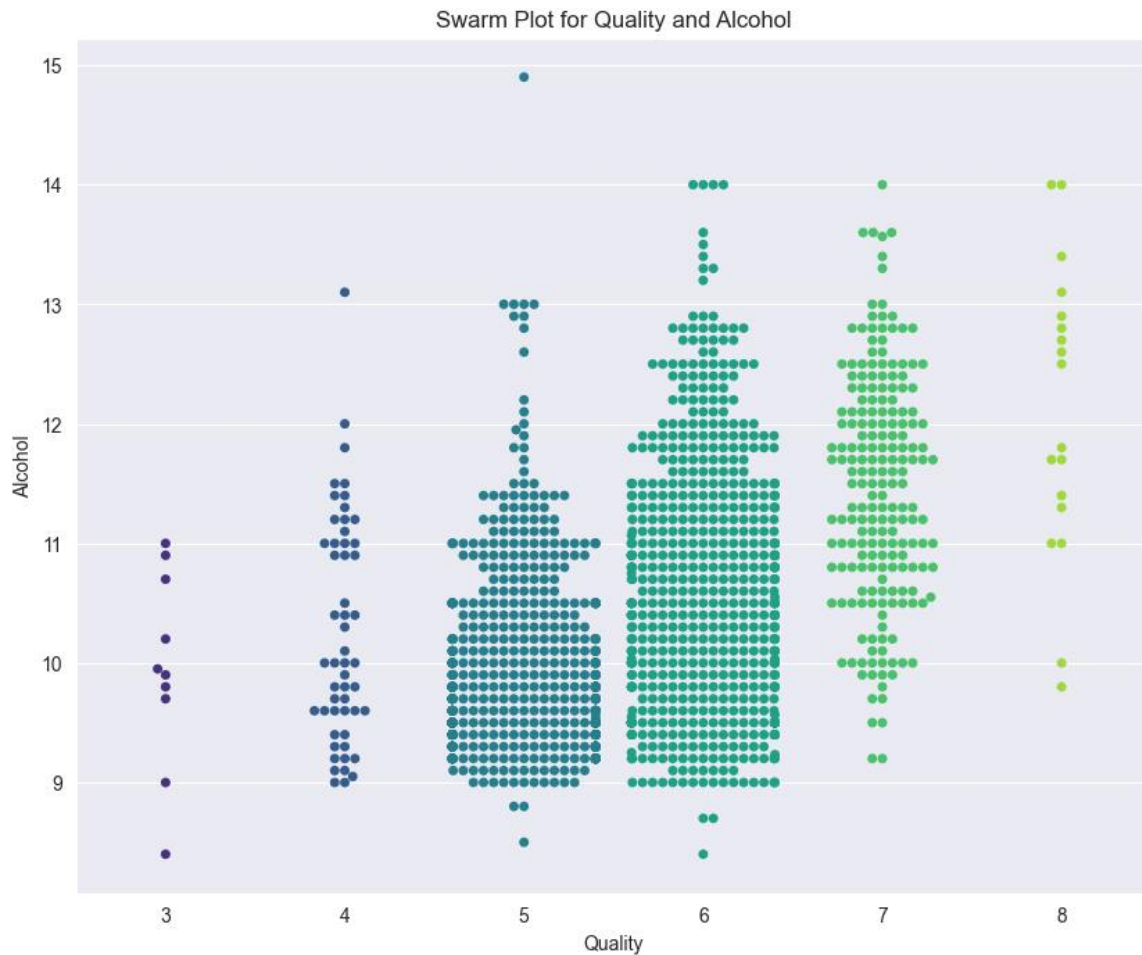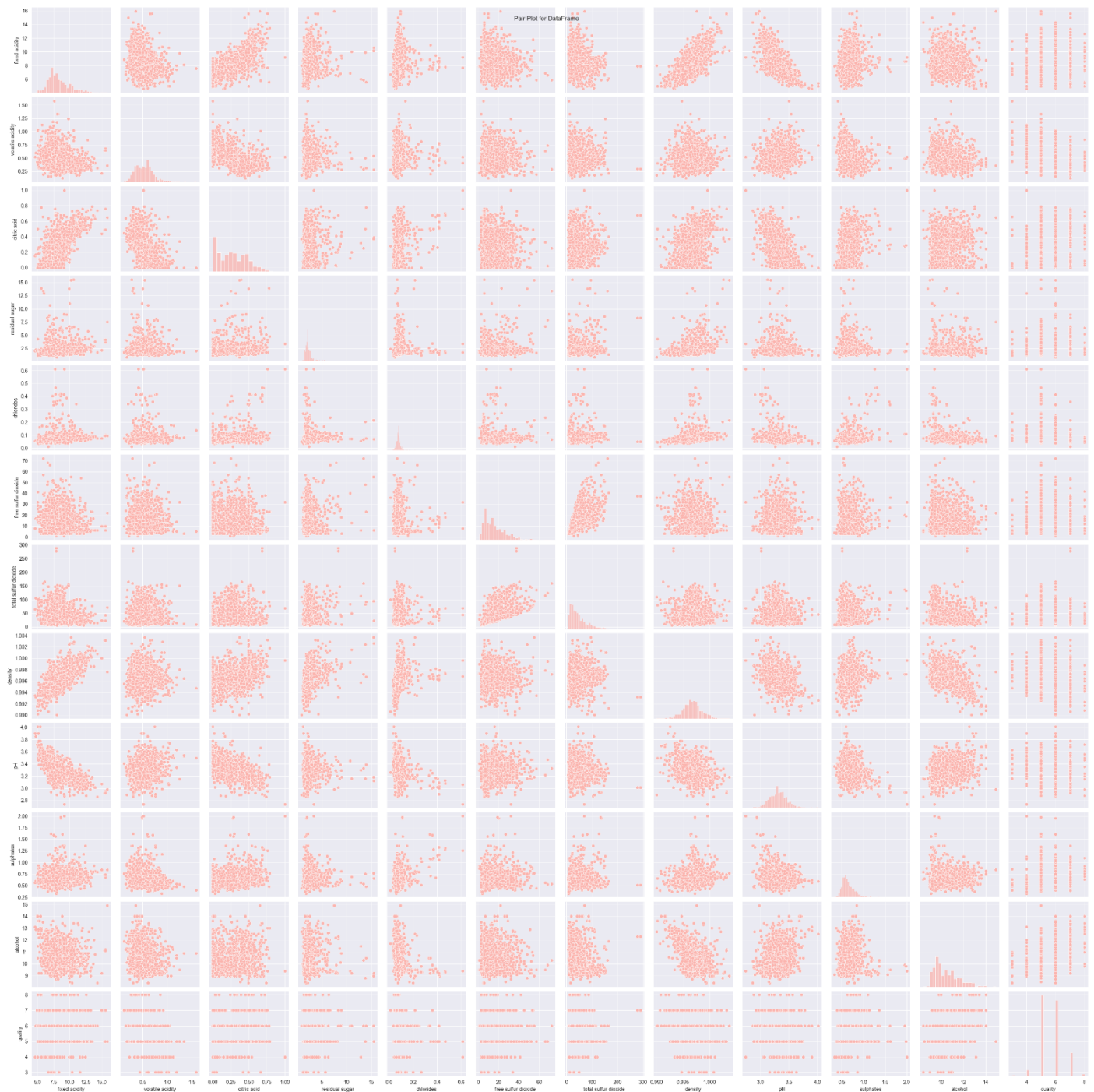
## Swarm plot

```
[15]: plt.figure(figsize=(10, 8))

      sns.swarmplot(x="quality", y="alcohol", data=df, palette='viridis')
      plt.title('Swarm Plot for Quality and Alcohol')
      plt.xlabel('Quality')
      plt.ylabel('Alcohol')
      plt.show()
```
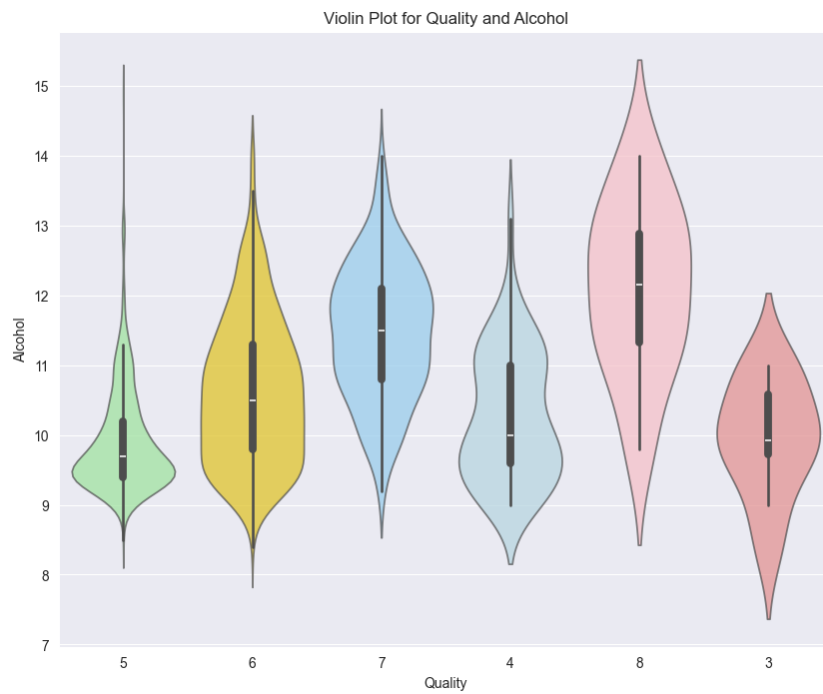


Swarm Plot for Quality and Alcohol

## Pair Plots

```
[17]: sns.set_palette("Pastel1")
      plt.figure(figsize=(10, 6))
      sns.pairplot(df)
      plt.suptitle('Pair Plot for DataFrame')
      plt.show()
```
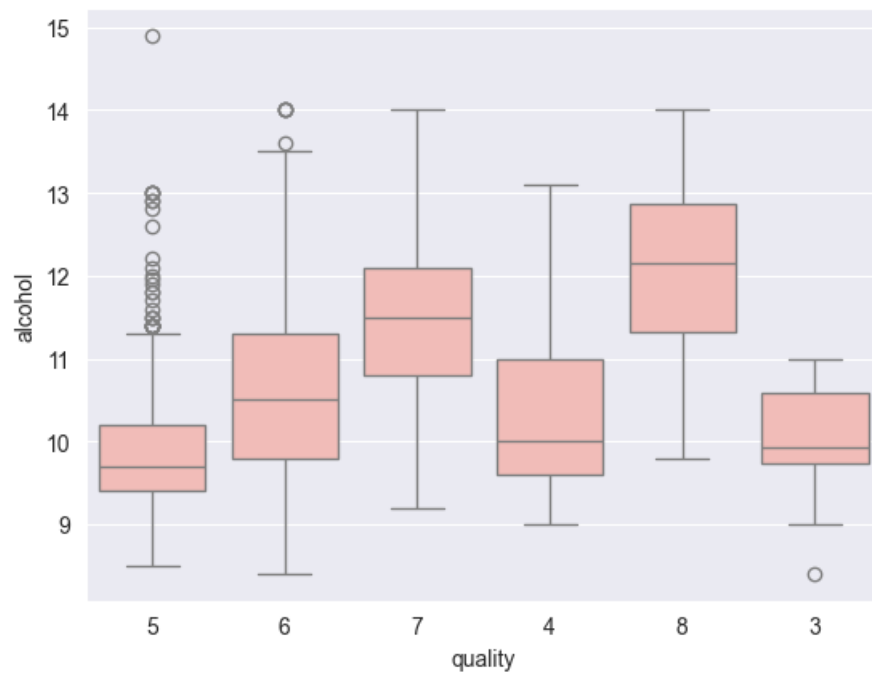
Pair Plot for DataFrame

## Violin plot

```python
df['quality'] = df['quality'].astype(str) # Convert 'quality' to categorical
plt.figure(figsize=(10, 8))
sns.violinplot(x="quality", y="alcohol", data=df, palette={
    '3': 'lightcoral', '4': 'lightblue', '5': 'lightgreen', '6': 'gold',
    '7': 'lightskyblue', '8': 'lightpink'},
    alpha=0.7)
plt.title('Violin Plot for Quality and Alcohol')
plt.xlabel('Quality')
plt.ylabel('Alcohol')
plt.show()
```

Violin Plot for Quality and Alcohol

## Box Plot

```
[20]: sns.boxplot(x='quality', y='alcohol', data=df)
```

## Correlation Matrix

```python
plt.figure(figsize=(15, 10))
sns.heatmap(df.corr(), annot=True, fmt='.2f', cmap='Blues', linewidths=2)
plt.title('Correlation Heatmap')
plt.show()
```



Correlation Heatmap