



```
1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4
5  # 데이터 불러오기
6  data = pd.read_csv('winequality-red.csv')
7
8  # 데이터 확인
9  print(data.head())
10
11 # 각 column의 데이터 타입 및 결측치 확인
12 print(data.info())
13
14 # 각 column의 기술통계량 확인
15 print(data.describe())
16
17 # 각 column의 분포 시각화
18 plt.figure(figsize=(12, 10))
19 for i, col in enumerate(data.columns):
20     plt.subplot(3, 4, i + 1)
21     sns.histplot(data[col], kde=True)
22     plt.title(col)
23 plt.tight_layout()
24 plt.show()
25
26 # 이상치 확인 (boxplot으로 확인)
27 plt.figure(figsize=(12, 10))
28 for i, col in enumerate(data.columns):
29     plt.subplot(3, 4, i + 1)
30     sns.boxplot(x=data[col])
31     plt.title(col)
32 plt.tight_layout()
33 plt.show()
34
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides \
0	7.4	0.70	0.00	1.9	0.076
1	7.8	0.88	0.00	2.6	0.098
2	7.8	0.76	0.04	2.3	0.092
3	11.2	0.28	0.56	1.9	0.075
4	7.4	0.70	0.00	1.9	0.076

	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates \
0	11.0	34.0	0.9978	3.51	0.56
1	25.0	67.0	0.9968	3.20	0.68
2	15.0	54.0	0.9970	3.26	0.65
3	17.0	60.0	0.9980	3.16	0.58
4	11.0	34.0	0.9978	3.51	0.56

	alcohol	quality
0	9.4	5
1	9.8	5
2	9.8	5
3	9.8	6
4	9.4	5

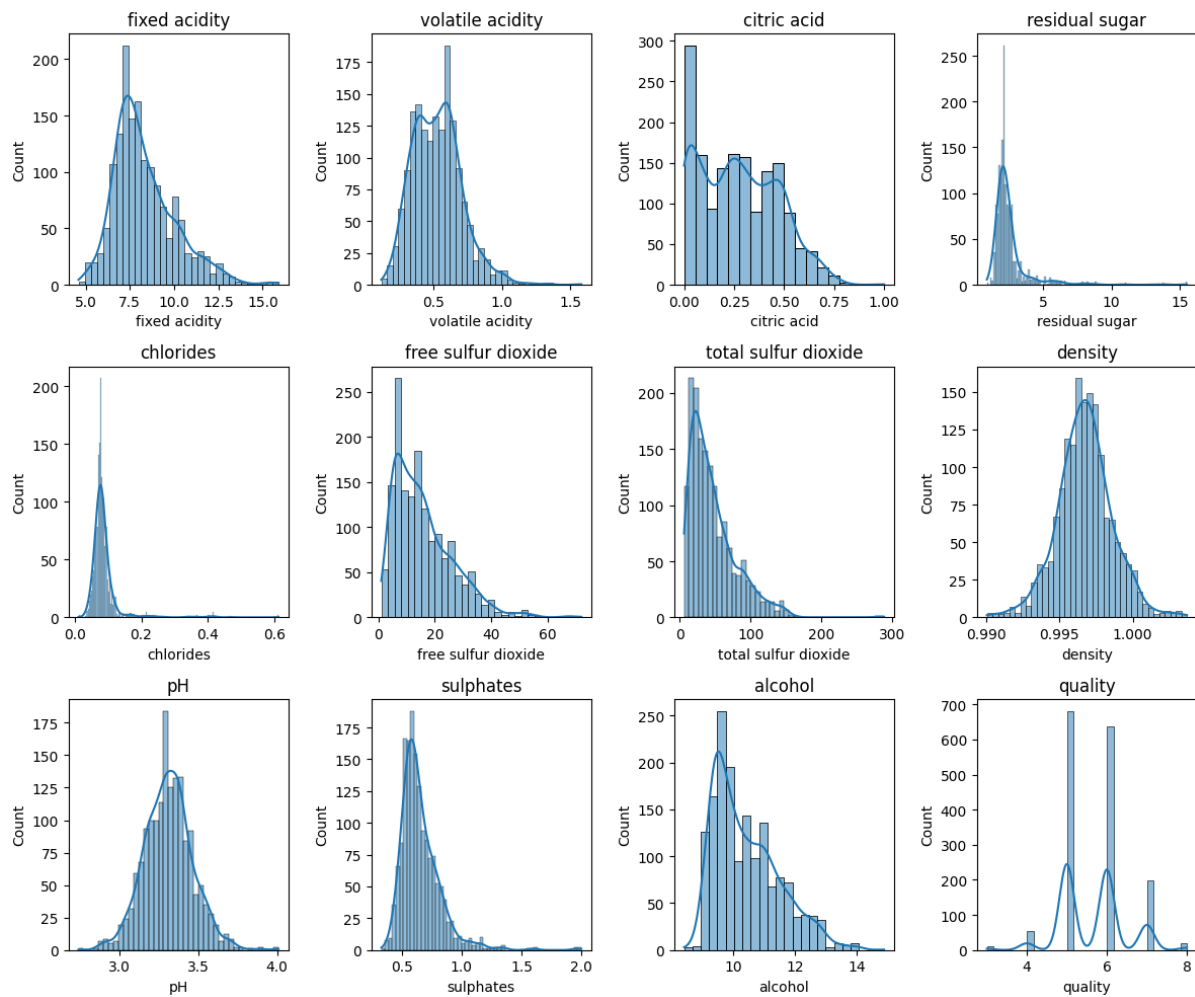
```
<class 'pandas.core.frame.DataFrame'>
```

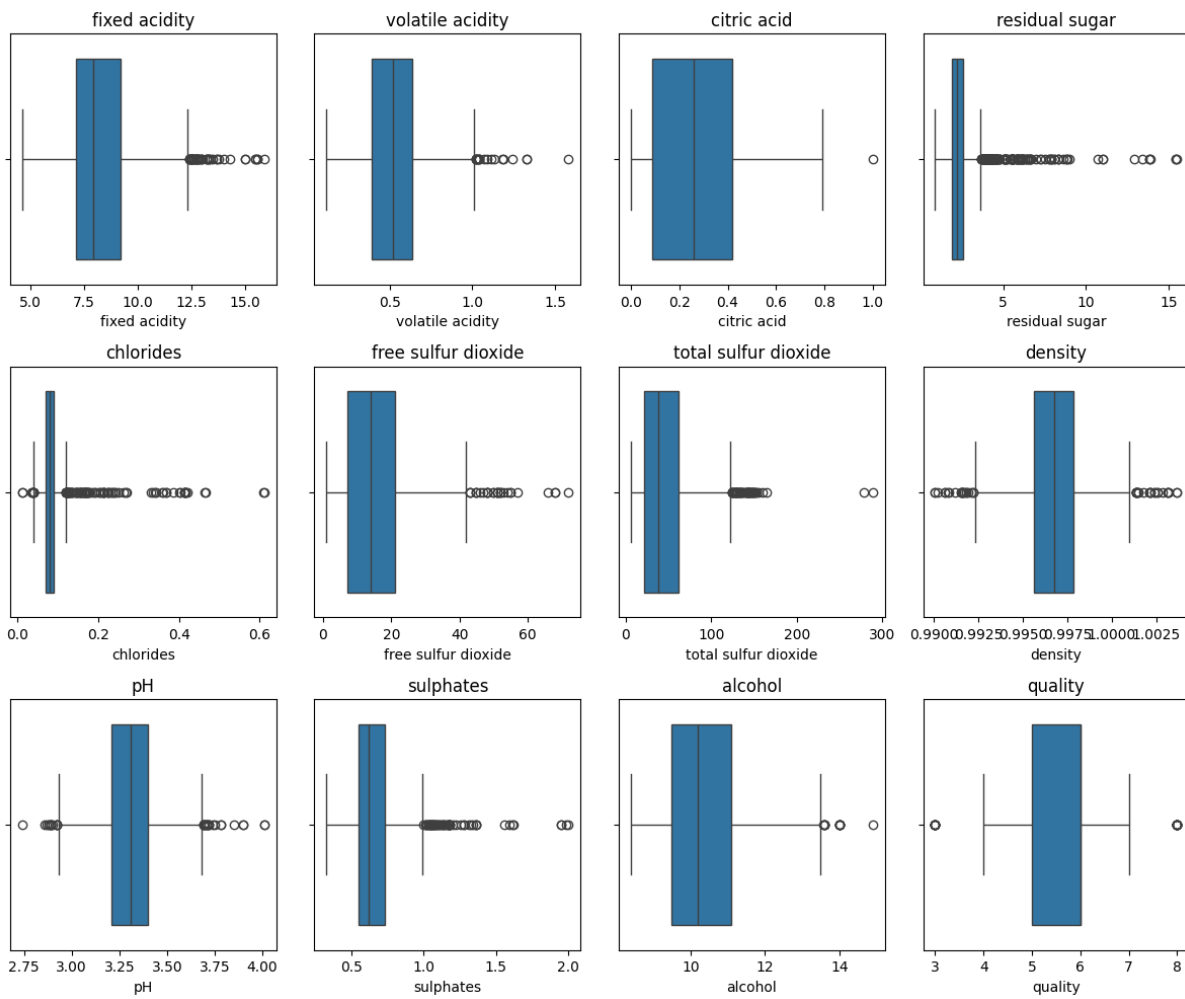
```
RangeIndex: 1599 entries, 0 to 1598
```


```
Data columns (total 12 columns):
```

#	Column	Non-Null Count	Dtype
...			
25%	3.210000	0.550000	9.500000
50%	3.310000	0.620000	10.200000
75%	3.400000	0.730000	11.100000
max	4.010000	2.000000	14.900000

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings...](#)







```

1  # 결측치 확인 및 처리
2  print(data.isnull().sum())
3
4  # Z-score 이용한 이상치 제거
5  from scipy.stats import zscore
6
7  z_scores = zscore(data)
8  abs_z_scores = abs(z_scores)
9  filtered_entries = (abs_z_scores < 3).all(axis=1)
10 data = data[filtered_entries]
11
12 # Min-Max 데이터 스케일링
13 from sklearn.preprocessing import MinMaxScaler
14
15 scaler = MinMaxScaler()
16 data_scaled = scaler.fit_transform(data)

```

```

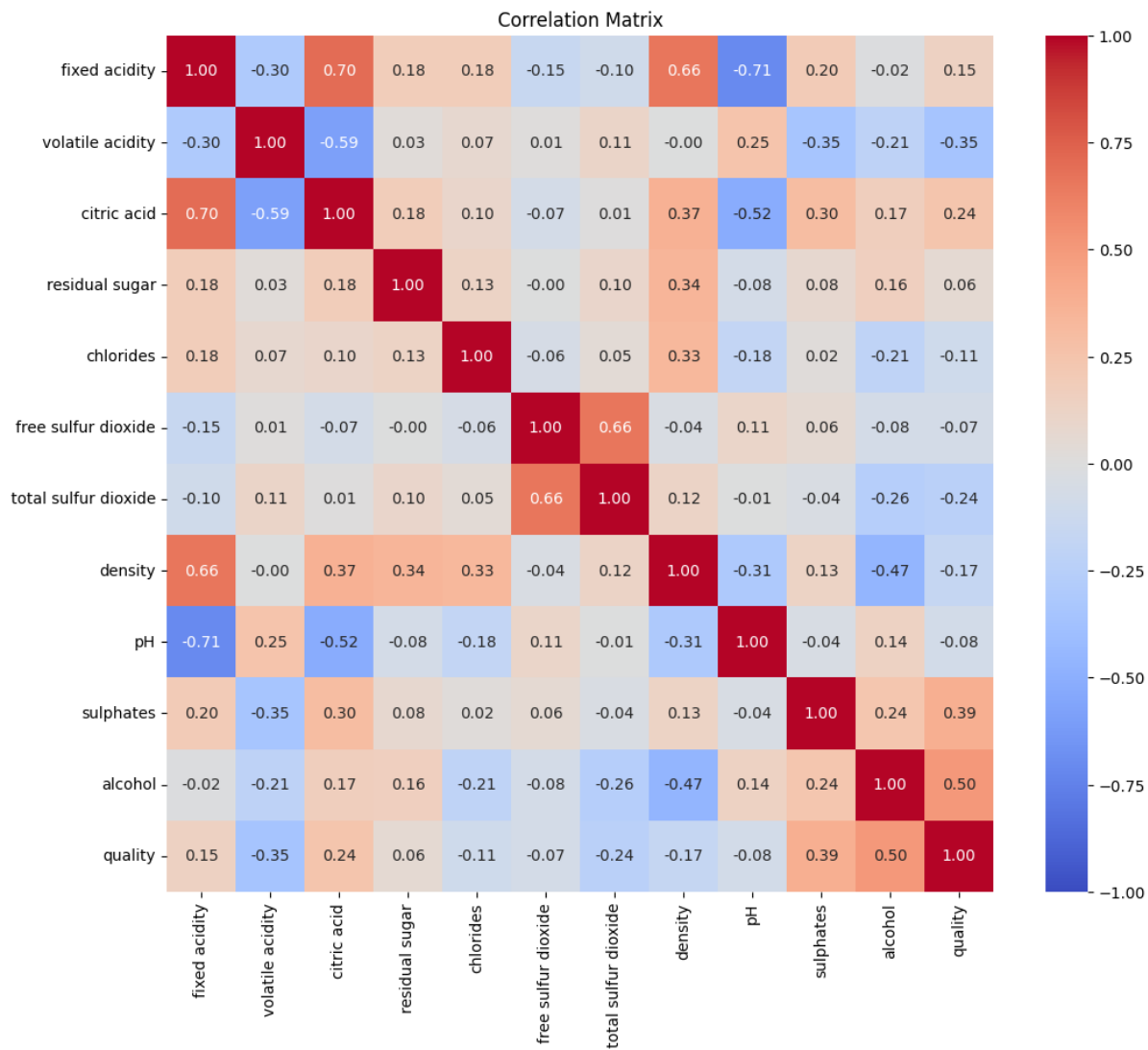
fixed acidity      0
volatile acidity  0
citric acid        0
residual sugar     0
chlorides          0
free sulfur dioxide 0
total sulfur dioxide 0
density           0
pH                0
sulphates         0
alcohol           0
quality           0
dtype: int64

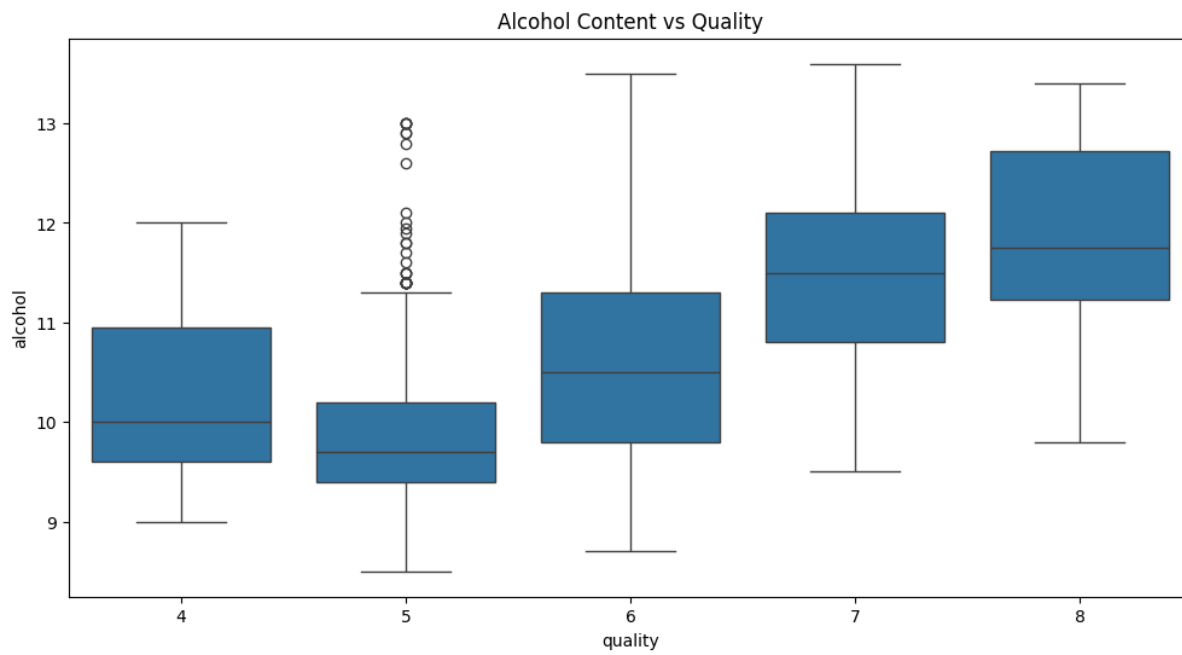
```

```

1 # 상관 분석
2 correlation_matrix = data.corr()
3 plt.figure(figsize=(12, 10))
4 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", vmin=-1, vmax=1)
5 plt.title('Correlation Matrix')
6 plt.show()
7
8 # 품질과 다른 변수들 간의 관계 시각화 (예시)
9 plt.figure(figsize=(12, 6))
10 sns.boxplot(x='quality', y='alcohol', data=data)
11 plt.title('Alcohol Content vs Quality')
12 plt.show()
13

```





```
1 from sklearn.model_selection import train_test_split
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.metrics import mean_squared_error, r2_score
4
5 # 데이터 준비
6 X = data.drop('quality', axis=1)
7 y = data['quality']
8
9 # 데이터 분할
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```

1  from sklearn.metrics import roc_auc_score
2  from itertools import combinations
3
4  # 모든 컬럼의 부분집합 생성
5  all_columns = X.columns
6  all_subsets = []
7
8  for i in range(1, len(all_columns) + 1):
9      subsets_i = list(combinations(all_columns, i))
10     all_subsets.extend(subsets_i)
11
12 # 모델 평가를 위한 함수 정의
13 def evaluate_model_auc(subset, X_train, X_test, y_train, y_test):
14     # 모델 선택 (예시로 랜덤 포레스트 분류 모델 사용)
15     model = RandomForestClassifier(random_state=42)
16
17     # 선택된 변수들로 데이터셋 재구성
18     X_train_subset = X_train[list(subset)]
19     X_test_subset = X_test[list(subset)]
20
21     # 모델 학습
22     model.fit(X_train_subset, y_train)
23
24     # 예측 확률 계산
25     y_pred_proba = model.predict_proba(X_test_subset)
26
27     # AUC 계산
28     auc = roc_auc_score(y_test, y_pred_proba, multi_class='ovo')
29
30     return auc
31
32 # 각 부분집합에 대해 AUC 계산
33 results = []
34 for subset in all_subsets:
35     auc = evaluate_model_auc(subset, X_train, X_test, y_train, y_test)
36     results.append((subset, auc))

```




```
1 # 결과를 AUC 기준으로 정렬
2 sorted_results = sorted(results, key=lambda x: x[1], reverse=True)
3
4 # 상위 10개 부분집합에 대한 결과 출력
5 for subset, auc in sorted_results[:10]:
6     print(f"Subset: {subset}")
7     print(f"AUC: {auc:.4f}")
8     print()
```

```
Subset: ('volatile acidity', 'citric acid', 'residual sugar', 'chlorides', 'free sulfur dioxide', 'sulphates', 'alcohol')
AUC: 0.8531

Subset: ('volatile acidity', 'citric acid', 'residual sugar', 'chlorides', 'density', 'sulphates', 'alcohol')
AUC: 0.8526

Subset: ('volatile acidity', 'residual sugar', 'free sulfur dioxide', 'density', 'pH', 'sulphates')
AUC: 0.8495

Subset: ('citric acid', 'residual sugar', 'chlorides', 'total sulfur dioxide', 'density', 'sulphates', 'alcohol')
AUC: 0.8492

Subset: ('fixed acidity', 'residual sugar', 'chlorides', 'total sulfur dioxide', 'pH', 'sulphates', 'alcohol')
AUC: 0.8479

Subset: ('volatile acidity', 'residual sugar', 'free sulfur dioxide', 'density', 'pH', 'sulphates', 'alcohol')
AUC: 0.8479

Subset: ('fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar', 'total sulfur dioxide', 'density', 'sulphates', 'alcohol')
AUC: 0.8472

Subset: ('citric acid', 'residual sugar', 'chlorides', 'total sulfur dioxide', 'pH', 'sulphates', 'alcohol')
AUC: 0.8463

Subset: ('volatile acidity', 'residual sugar', 'density', 'sulphates', 'alcohol')
...

Subset: ('fixed acidity', 'volatile acidity', 'residual sugar', 'free sulfur dioxide', 'density', 'sulphates', 'alcohol')
AUC: 0.8461
```

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#)...