

## Deskripsi Dataset: Kualitas Anggur Merah "Vinho Verde"

### 1. Gambaran Umum Dataset

Dataset yang digunakan dalam analisis ini adalah data kualitas anggur merah (*red wine*) dari varian "Vinho Verde" yang berasal dari Portugal. Dataset ini menyediakan informasi kuantitatif mengenai atribut fisikokimia yang didapat dari pengujian laboratorium dan penilaian sensorik berupa skor kualitas.

Penting untuk dicatat bahwa dataset ini tidak mencakup informasi mengenai jenis anggur, merek, atau harga jual, sehingga fokus analisis murni pada hubungan antara karakteristik fisikokimia dengan kualitas anggur yang dihasilkan. Permasalahan ini dapat didekati sebagai tugas regresi (memprediksi skor kualitas secara langsung) atau klasifikasi (mengelompokkan anggur ke dalam kategori kualitas tertentu), di mana pendekatan klasifikasi menjadi fokus dalam laporan ini.

### 2. Sumber dan Sitasi

Dataset ini merupakan koleksi publik yang tersedia di **UCI Machine Learning Repository** dan telah dibagikan kembali di platform **Kaggle** untuk kemudahan akses. Berikut tautan dataset dari Kaggle yang kami gunakan: <https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009>

### 3. Deskripsi Variabel (Atribut)

Dataset ini terdiri dari 11 variabel input (fitur) dan 1 variabel output (target). Semua fitur merupakan variabel numerik dengan tipe data `float64`.

#### Variabel Input (Fitur Fisikokimia):

1. **fixed acidity** (Keasaman Tetap):

- Mengukur konsentrasi asam-asam non-volatil (tidak mudah menguap) dalam anggur, seperti asam tartarat dan malat. Berkontribusi pada rasa "tajam" pada anggur. (Tipe Data: `float64`)

2. **volatile acidity** (Keasaman Volatil):

- Mengukur jumlah asam asetat dalam anggur. Tingkat yang terlalu tinggi dapat menyebabkan aroma dan rasa seperti cuka yang tidak diinginkan. (Tipe Data: `float64`)

3. **citric acid** (Asam Sitrat)

- Dalam jumlah kecil, asam sitrat dapat menambah kesegaran dan rasa "segar" pada anggur. (Tipe Data: `float64`)

4. **residual sugar** (Sisa Gula):

- Jumlah gula yang tersisa setelah proses fermentasi berhenti. Ini menentukan tingkat kemanisan anggur. (Tipe Data: `float64`)

5. **chlorides** (Klorida):

- Jumlah garam (khususnya natrium klorida) dalam anggur. (Tipe Data: `float64`)

6. **free sulfur dioxide** (Belerang Dioksida Bebas):

- Bagian dari  $\text{SO}_2$  yang berada dalam bentuk bebas dan berfungsi sebagai antimikroba serta antioksidan untuk mencegah pembusukan. (Tipe Data: `float64`)

7. **total sulfur dioxide** (Total Belerang Dioksida):

- Gabungan dari belerang dioksida dalam bentuk bebas dan terikat. Berfungsi sebagai pengawet. (Tipe Data: `float64`)

8. **density** (Kepadatan):

- Kepadatan anggur yang umumnya mendekati kepadatan air, dipengaruhi oleh kandungan alkohol dan gula. (Tipe Data: `float64`)

9. **pH**:

- Mengukur tingkat keasaman atau kebasaan pada skala pH, yang memengaruhi rasa dan stabilitas anggur. Sebagian besar anggur memiliki pH antara 3 dan 4. (Tipe Data: `float64`)

10. **sulphates** (Sulfat):

- Garam sulfat yang dapat berkontribusi pada efektivitas  $\text{SO}_2$  dan memengaruhi rasa anggur. (Tipe Data: `float64`)

11. **alcohol** (Alkohol):

- Persentase kandungan alkohol dalam anggur berdasarkan volume. (Tipe Data: `float64`)

#### Variabel Output (Target Sensorik):

12. **quality** (Kualitas):

- Skor kualitas yang diberikan oleh ahli berdasarkan penilaian sensorik, dengan skala asli antara 0 hingga 10. Untuk tugas klasifikasi ini, variabel ini telah ditransformasikan menjadi kategori biner:

- **1 (Baik):** Jika skor asli  $\geq 7$

- **0 (Buruk):** Jika skor asli < 7

#### 4. Karakteristik Statistik dan Pembersihan Data

Analisis eksplorasi data awal mengungkapkan beberapa karakteristik penting:

- **Ukuran Data Awal:** Dataset ini pada awalnya berisi **1599 baris** data.
- **Data Duplikat:** Teridentifikasi adanya **240 baris** data yang duplikat. Baris-baris ini telah dihapus untuk memastikan integritas data dalam pemodelan.
- **Ukuran Data Final:** Setelah pembersihan, dataset yang digunakan untuk analisis terdiri dari **1359 baris dan 12 kolom**.
- **Nilai Hilang:** **Tidak ditemukan** adanya nilai yang hilang (*missing values*) pada setiap kolom, sehingga tidak diperlukan proses imputasi.
- **Distribusi dan Skala:** Ringkasan statistik menunjukkan bahwa rentang nilai dan skala antar fitur sangat bervariasi. Sebagai contoh, total sulfur dioxide memiliki nilai maksimum 289.0, sedangkan chlorides memiliki nilai maksimum 0.611.
- **Keseimbangan Kelas:** Setelah variabel *quality* ditransformasikan menjadi kategori biner, teridentifikasi adanya **ketidakseimbangan kelas** (*class imbalance*). Terdapat **1175 sampel** untuk anggur berkualitas "Buruk" (kelas 0) dan hanya **184 sampel** untuk anggur berkualitas "Baik" (kelas 1). Kondisi ini menjadi pertimbangan krusial dalam proses pemodelan, di mana strategi seperti *stratified sampling* dan penyesuaian bobot kelas (*class weighting*) diterapkan untuk mencegah model menjadi bias terhadap kelas mayoritas.

## ✓ 1. Import Library

```

1  # Manipulasi dan Analisis Data
2  import pandas as pd
3  import numpy as np
4
5  # Visualisasi Data
6  import seaborn as sns
7  import matplotlib.pyplot as plt
8
9  # Pemuatan Dataset dari Kaggle
10 import kagglehub
11 from kagglehub import KaggleDatasetAdapter
12
13 # Pemodelan dan Evaluasi
14 from sklearn.model_selection import train_test_split, GridSearchCV
15 from sklearn.ensemble import RandomForestClassifier
16 from sklearn.metrics import (accuracy_score, precision_score, recall_score,
17                             classification_report, confusion_matrix,
18                             ConfusionMatrixDisplay,
19                             roc_auc_score, roc_curve)
20
21 # Visualisasi Decision Tree
22 from sklearn.tree import export_graphviz
23 import graphviz

```

### Penjelasan Library yang digunakan

#### 1. pandas (pd)

- Digunakan untuk manipulasi dan analisis data.
- Fungsi dalam proyek:
  - Membaca dataset
  - Mengolah data dalam bentuk *DataFrame*
  - Membersihkan data (menghapus duplikat, memeriksa missing values)
  - Memanipulasi kolom data

#### 2. numpy (np)

- Digunakan untuk operasi numerik.
- Meskipun tidak eksplisit, menjadi dasar operasi di library lain seperti pandas dan scikit-learn.

#### 3. seaborn (sns)

- Library visualisasi data tingkat tinggi.
- Fungsi: Membuat *histplot*, *boxplot*, dan *heatmap* untuk eksplorasi data. [1]

#### 4. matplotlib.pyplot (plt)

- Library visualisasi data dasar.
- Fungsi:
  - Membuat dan mengatur plot
  - Menambahkan judul, label, grid
  - Menampilkan plot [1]

**5. kagglehub**

- Library untuk berinteraksi dengan Kaggle Hub.
- Digunakan khusus untuk memuat dataset langsung dari Kaggle.

**6. KaggleDatasetAdapter**

- Adapter dari kagglehub untuk memuat dataset sebagai DataFrame pandas.

**7. sklearn.model\_selection.train\_test\_split**

- Membagi dataset menjadi set **pelatihan** dan **pengujian**.

**8. sklearn.model\_selection.GridSearchCV**

- Melakukan **hyperparameter tuning** dengan cross-validation.

**9. sklearn.ensemble.RandomForestClassifier**

- Model ML berbasis ensemble (gabungan decision tree).
- Digunakan untuk klasifikasi kualitas anggur.

**10. sklearn.metrics**

- Modul metrik evaluasi performa model klasifikasi:
  - `accuracy_score`, `precision_score`, `recall_score`
  - `classification_report` (laporan presisi, recall, f1-score per kelas)
  - `confusion_matrix` + `ConfusionMatrixDisplay` (visualisasi)
  - `roc_auc_score` (AUC-ROC)
  - `roc_curve` (menghitung FPR/TPR untuk plot ROC)

**11. sklearn.tree.export\_graphviz**

- Menghasilkan representasi DOT dari decision tree.

**12. graphviz**

- Merender output `export_graphviz` menjadi gambar decision tree.

## ✓ 2. Load data dan Exploratory Data Analysis (EDA)

Pada tahap ini, kita akan memuat dataset dan melakukan beberapa langkah eksplorasi data awal untuk memahami karakteristik data yang kita miliki.

### ✓ 2.1 Memuat Dataset

Dataset diunduh langsung dari Kaggle Hub menggunakan pustaka kagglehub.

```
1 df = kagglehub.load_dataset(
2     KaggleDatasetAdapter.PANDAS,
3     "uciml/red-wine-quality-cortez-et-al-2009",
4     "winequality-red.csv"
5 )
6
7 display(df.head())
```

 <ipython-input-162-92150c9e3da5>:1: DeprecationWarning: load\_dataset is deprecated and will be removed in a future version.  
df = kagglehub.load\_dataset(

|   | fixed<br>acidity | volatile<br>acidity | citric<br>acid | residual<br>sugar | chlorides | free<br>sulfur<br>dioxide | total<br>sulfur<br>dioxide | density | pH   | sulphates | alcohol | quality |
|---|------------------|---------------------|----------------|-------------------|-----------|---------------------------|----------------------------|---------|------|-----------|---------|---------|
| 0 | 7.4              | 0.70                | 0.00           | 1.9               | 0.076     | 11.0                      | 34.0                       | 0.9978  | 3.51 | 0.56      | 9.4     | 5       |
| 1 | 7.8              | 0.88                | 0.00           | 2.6               | 0.098     | 25.0                      | 67.0                       | 0.9968  | 3.20 | 0.68      | 9.8     | 5       |
| 2 | 7.8              | 0.76                | 0.04           | 2.3               | 0.092     | 15.0                      | 54.0                       | 0.9970  | 3.26 | 0.65      | 9.8     | 5       |
| 3 | 11.2             | 0.28                | 0.56           | 1.9               | 0.075     | 17.0                      | 60.0                       | 0.9980  | 3.16 | 0.58      | 9.8     | 6       |

### ✓ 2.2 Pemeriksaan Awal Data

Kita akan melihat informasi dasar dari dataset seperti tipe data setiap kolom, jumlah entri, dan ringkasan statistiknya.

```
1 df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   fixed acidity         1599 non-null   float64
 1   volatile acidity      1599 non-null   float64
 2   citric acid           1599 non-null   float64
 3   residual sugar        1599 non-null   float64
 4   chlorides             1599 non-null   float64
 5   free sulfur dioxide    1599 non-null   float64
 6   total sulfur dioxide   1599 non-null   float64
 7   density               1599 non-null   float64
 8   pH                   1599 non-null   float64
 9   sulphates             1599 non-null   float64
10   alcohol               1599 non-null   float64
11   quality               1599 non-null   int64   
dtypes: float64(11), int64(1)
memory usage: 150.0 KB

```

Metode `df.info()` memberikan kita ringkasan singkat dan penting mengenai DataFrame (tabel data) yang baru saja kita muat. Ini adalah langkah pertama yang krusial dalam setiap proyek analisis data.

Dari output di atas, kita dapat menyimpulkan beberapa hal:

- **Struktur Data:**

- Data kita memiliki **1599 baris** (*entries*) dan **12 kolom**. Ini memberikan gambaran awal tentang ukuran dataset.

- **Kolom dan Tipe Data:**

- **Non-Null Count** : Semua kolom menunjukkan 1599 non-null. Ini adalah kabar baik! Artinya, **tidak ada data yang hilang (missing values)** dalam dataset kita. Dengan demikian, kita tidak perlu melakukan penanganan data kosong seperti imputasi.
- **Dtype** : Menunjukkan tipe data setiap kolom.
  - Terdapat **11 kolom** dengan tipe `float64` (angka desimal), yang cocok untuk atribut fisiko-kimia seperti *fixed acidity*, *alcohol*, dll.
  - Ada **1 kolom** dengan tipe `int64` (bilangan bulat), yaitu kolom target kita, *quality*.

- **Penggunaan Memori:**

- DataFrame ini menggunakan sekitar **150.0 KB** memori.

Secara keseluruhan, data ini terlihat "bersih" dari nilai yang hilang dan siap untuk tahap eksplorasi lebih lanjut.

```
1 display(df.describe())
```

|              | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides   | free sulfur dioxide | total sulfur dioxide | density     | pH          | sulphates   |
|--------------|---------------|------------------|-------------|----------------|-------------|---------------------|----------------------|-------------|-------------|-------------|
| <b>count</b> | 1599.000000   | 1599.000000      | 1599.000000 | 1599.000000    | 1599.000000 | 1599.000000         | 1599.000000          | 1599.000000 | 1599.000000 | 1599.000000 |
| <b>mean</b>  | 8.319637      | 0.527821         | 0.270976    | 2.538806       | 0.087467    | 15.874922           | 46.467792            | 0.996747    | 3.311113    | 0.658149    |
| <b>std</b>   | 1.741096      | 0.179060         | 0.194801    | 1.409928       | 0.047065    | 10.460157           | 32.895324            | 0.001887    | 0.154386    | 0.169507    |
| <b>min</b>   | 4.600000      | 0.120000         | 0.000000    | 0.900000       | 0.012000    | 1.000000            | 6.000000             | 0.990070    | 2.740000    | 0.330000    |
| <b>25%</b>   | 7.100000      | 0.390000         | 0.090000    | 1.900000       | 0.070000    | 7.000000            | 22.000000            | 0.995600    | 3.210000    | 0.550000    |
| <b>50%</b>   | 7.900000      | 0.520000         | 0.260000    | 2.200000       | 0.079000    | 14.000000           | 38.000000            | 0.996750    | 3.310000    | 0.620000    |
| <b>75%</b>   | 8.300000      | 0.640000         | 0.400000    | 2.600000       | 0.080000    | 21.000000           | 62.000000            | 0.997825    | 3.400000    | 0.720000    |
| <b>max</b>   | 16.960000     | 1.910000         | 0.450000    | 15.000000      | 0.190000    | 170.000000          | 170.000000           | 1.050000    | 4.010000    | 0.950000    |

Metode `df.describe()` memberikan ringkasan statistik yang sangat berguna untuk kolom-kolom numerik. Dari sini, kita bisa mendapatkan pemahaman awal tentang pusat data, sebaran, dan potensi adanya nilai-nilai aneh (*outlier*).

Berikut adalah penjabaran dari setiap metrik yang ditampilkan:

- **count** : Jumlah baris data yang tidak kosong. Angka 1599 di semua kolom mengonfirmasi bahwa tidak ada nilai yang hilang.
- **mean** : Nilai rata-rata dari setiap kolom.
- **std** (Standar Deviasi) : Mengukur seberapa tersebar data dari nilai rata-ratanya. Nilai **std** yang besar, seperti pada *total sulfur dioxide* (32.89), menunjukkan sebaran data yang luas.
- **min** : Nilai terkecil dalam kolom.
- **25%** : Kuartil pertama (Q1). 25% dari data memiliki nilai di bawah angka ini.
- **50%** : Median atau kuartil kedua (Q2). Ini adalah nilai tengah dari data, yang seringkali lebih representatif daripada **mean** jika ada outlier.
- **75%** : Kuartil ketiga (Q3). 75% dari data memiliki nilai di bawah angka ini.
- **max** : Nilai terbesar dalam kolom.

Pengamatan Kunci dari Statistik:

### 1. Potensi Outlier

- Terlihat perbedaan yang sangat signifikan antara nilai 75% dan max pada beberapa fitur. Contohnya, pada **total sulfur dioxide**, nilai kuartil ke-3 adalah 62, namun nilai maksimumnya mencapai 289. Hal yang sama terlihat pada **residual sugar** dan **chlorides**. Ini adalah indikasi kuat adanya *outlier* yang perlu kita perhatikan.

### 2. Distribusi Data

- Untuk beberapa kolom seperti **residual sugar**, nilai mean (2.53) lebih tinggi dari median/50% (2.2). Ini mengindikasikan distribusi data yang miring ke kanan (*right-skewed*), yang berarti ada beberapa nilai yang sangat tinggi yang "menarik" rata-ratanya.

### 3. Variabel Target quality

- Kolom **quality** memiliki nilai rata-rata (mean) 5.63. Skoranya berkisar dari 3 (min) hingga 8 (max), menunjukkan tidak ada anggur dengan kualitas sangat buruk atau sangat baik dalam dataset ini. Sebagian besar data terkumpul di tengah.

### 4. Skala Fitur

- Nilai pada kolom **total sulfur dioxide** jauh lebih besar dibandingkan **chlorides**. Ini menunjukkan bahwa penskalaan fitur (*feature scaling*) mungkin akan diperlukan sebelum melatih beberapa jenis model machine learning untuk memastikan setiap fitur memberikan kontribusi yang seimbang.

## 2.3 Data Cleaning: Duplikat dan Missing Value

```
1 # Memeriksa data duplikat
2 print(f'Jumlah data duplikat : {df.duplicated().sum()}')
3 print(f'Jumlah baris sebelum menghapus duplikat: {len(df)}')
4
5 # Menghapus data duplikasi
6 df.drop_duplicates(inplace=True)
7
8 # Memverifikasi setelah penanganan
9 print(f'Jumlah baris setelah menghapus duplikat : {len(df)}')
```

➡ Jumlah data duplikat : 240  
 Jumlah baris sebelum menghapus duplikat: 1599  
 Jumlah baris setelah menghapus duplikat : 1359

```
1 # Memeriksa nilai yang hilang (missing values)
2 df.isna().sum()
```

➡

|                             | 0 |
|-----------------------------|---|
| <b>fixed acidity</b>        | 0 |
| <b>volatile acidity</b>     | 0 |
| <b>citric acid</b>          | 0 |
| <b>residual sugar</b>       | 0 |
| <b>chlorides</b>            | 0 |
| <b>free sulfur dioxide</b>  | 0 |
| <b>total sulfur dioxide</b> | 0 |
| <b>density</b>              | 0 |
| <b>pH</b>                   | 0 |
| <b>sulphates</b>            | 0 |
| <b>alcohol</b>              | 0 |
| <b>quality</b>              | 0 |

dtype: int64

Hasil di atas menunjukkan bahwa data sudah bersih dari nilai yang hilang.

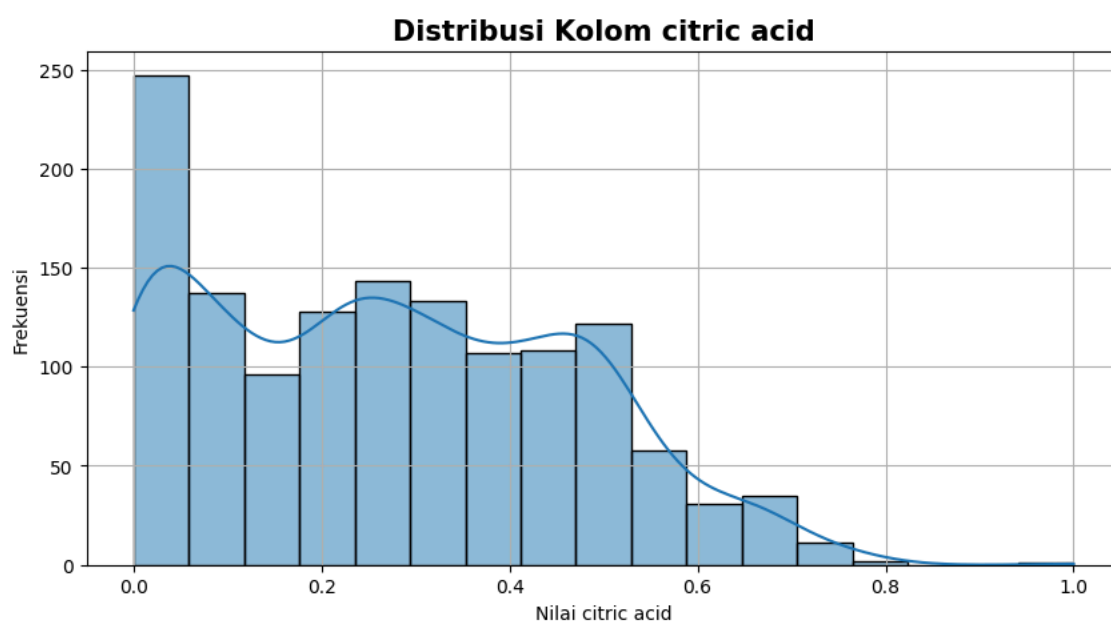
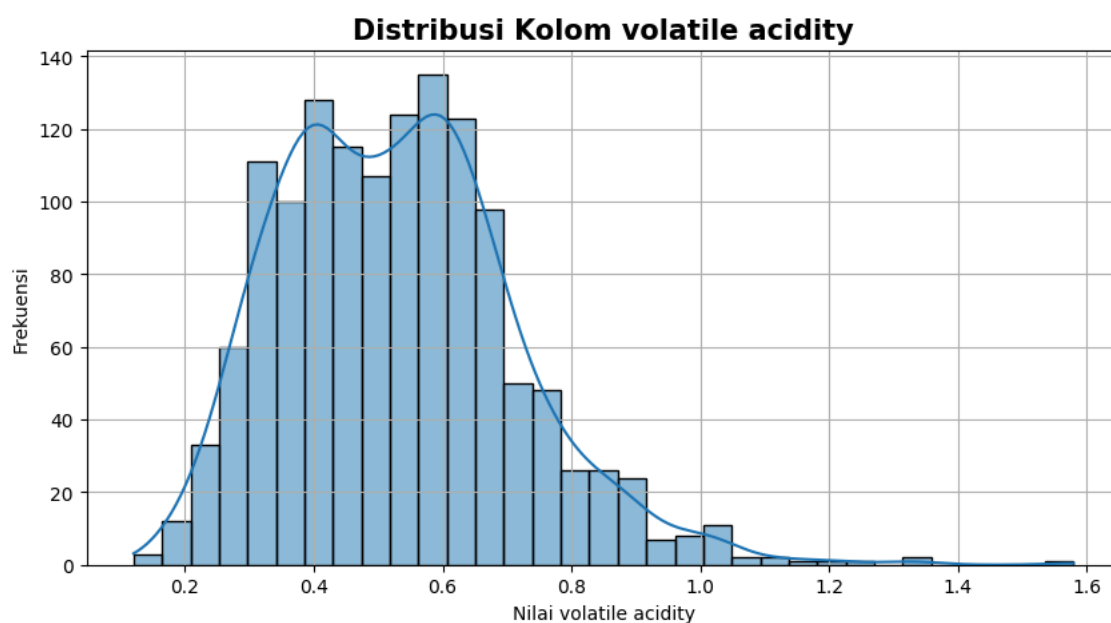
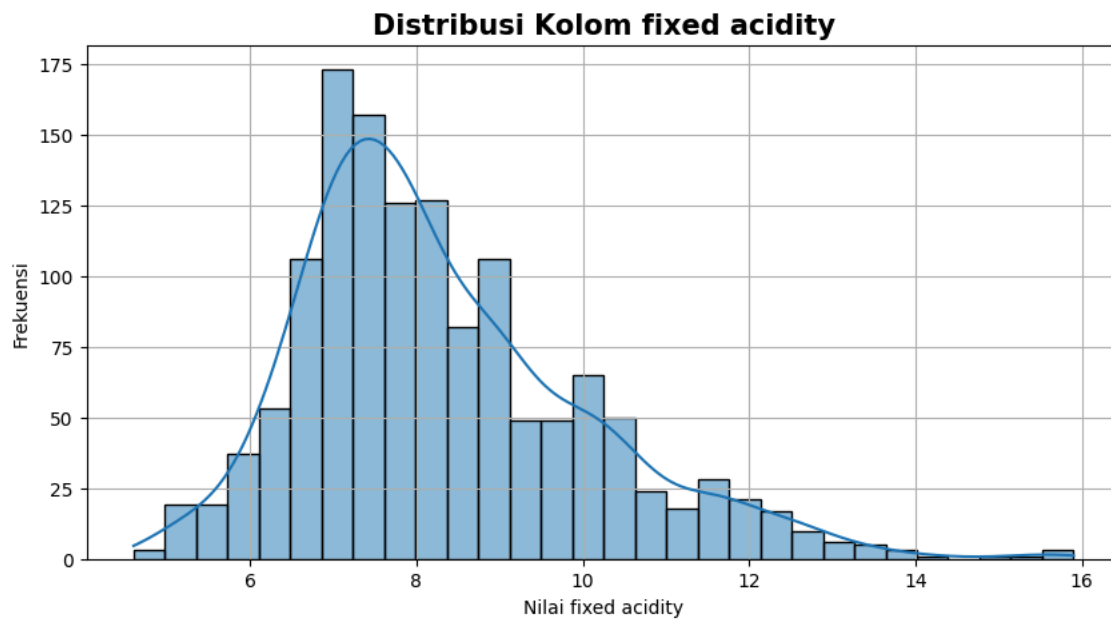
## 2.4 Visualisasi Distribusi Fitur

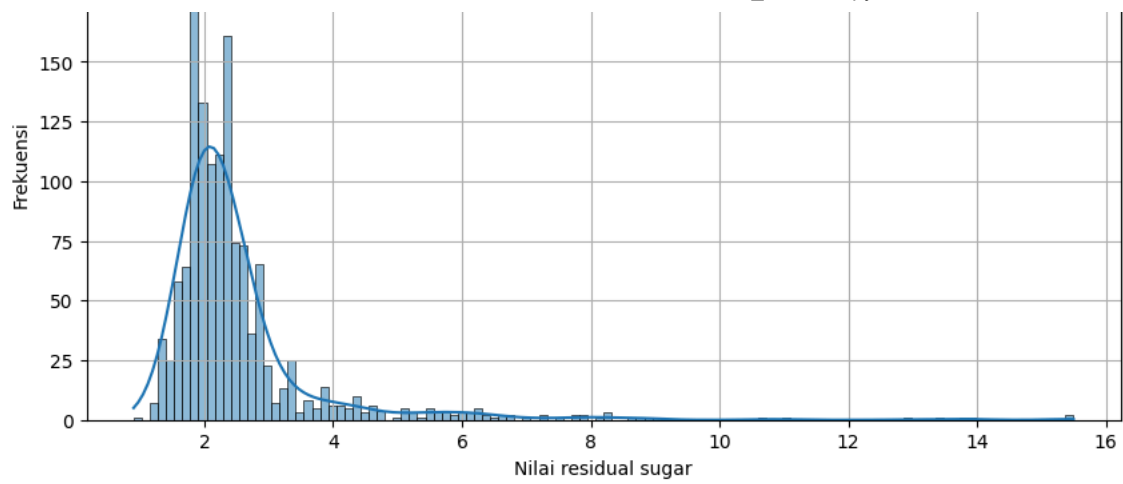
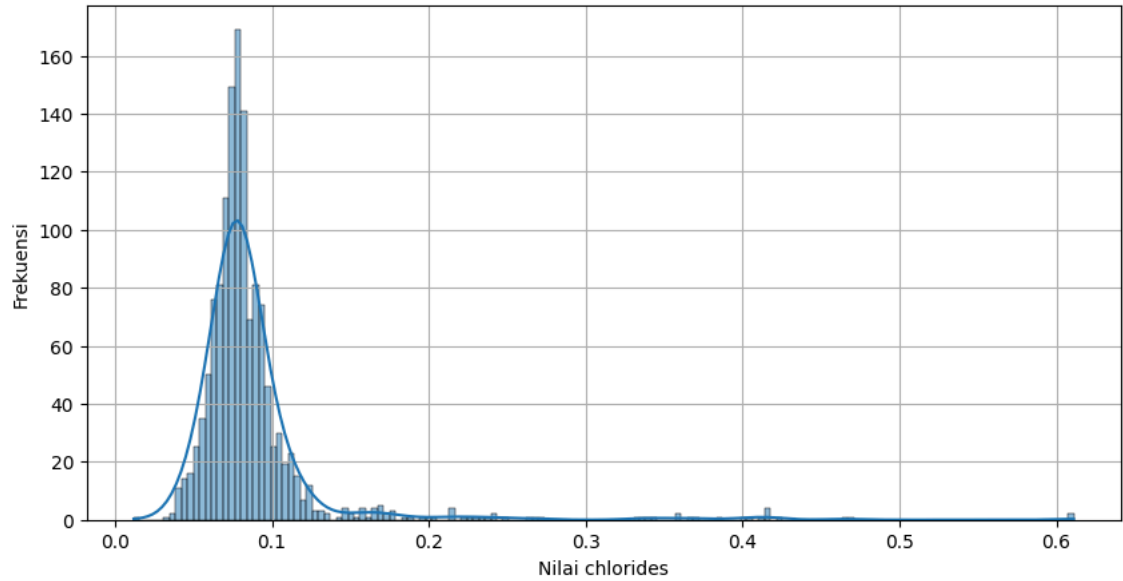
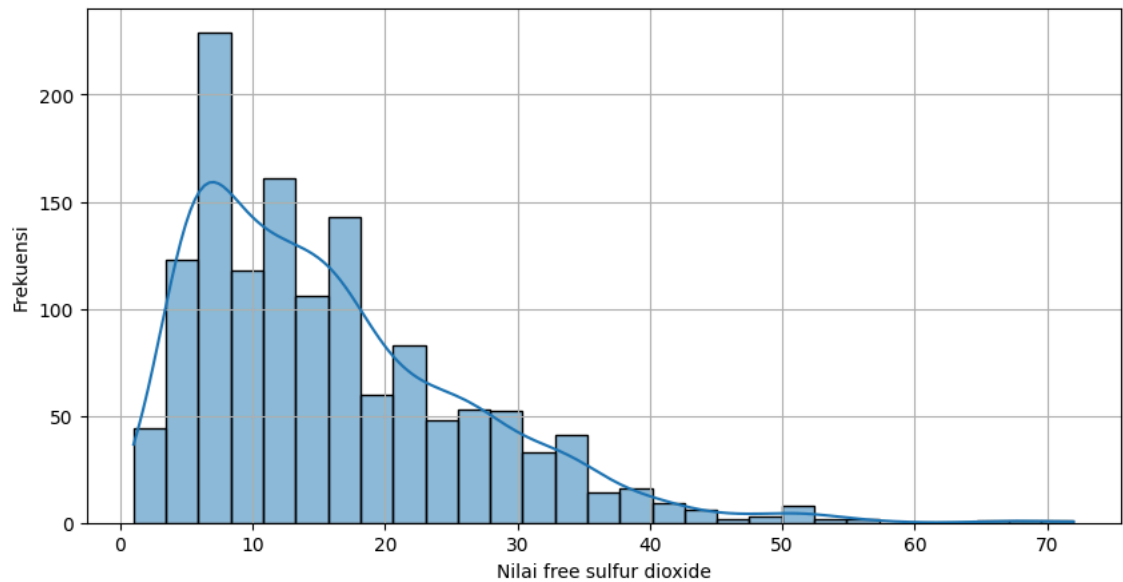
```
1 print("\nVisualisasi Distribusi Setiap Fitur:")
2 for column in df.columns:
3     plt.figure(figsize=(10, 5))
4     sns.histplot(data=df, x=column, kde=True)
5     plt.title(f'Distribusi Kolom {column}', fontweight='bold', fontsize=15)
```

```
6 plt.xlabel(f'Nilai {column}')
7 plt.ylabel(f'Frekuensi')
8 plt.grid(True)
```

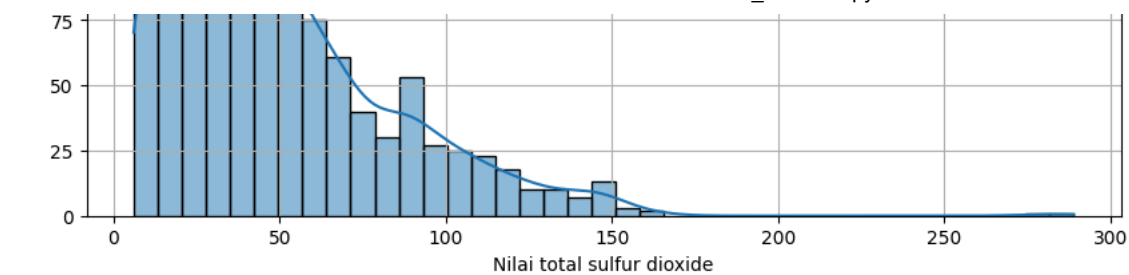
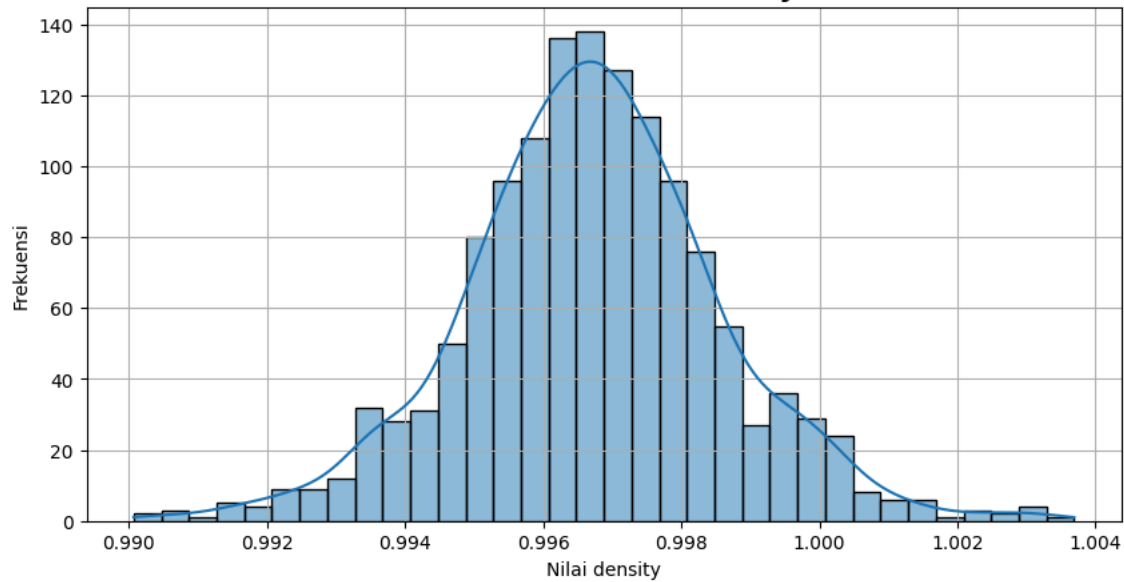
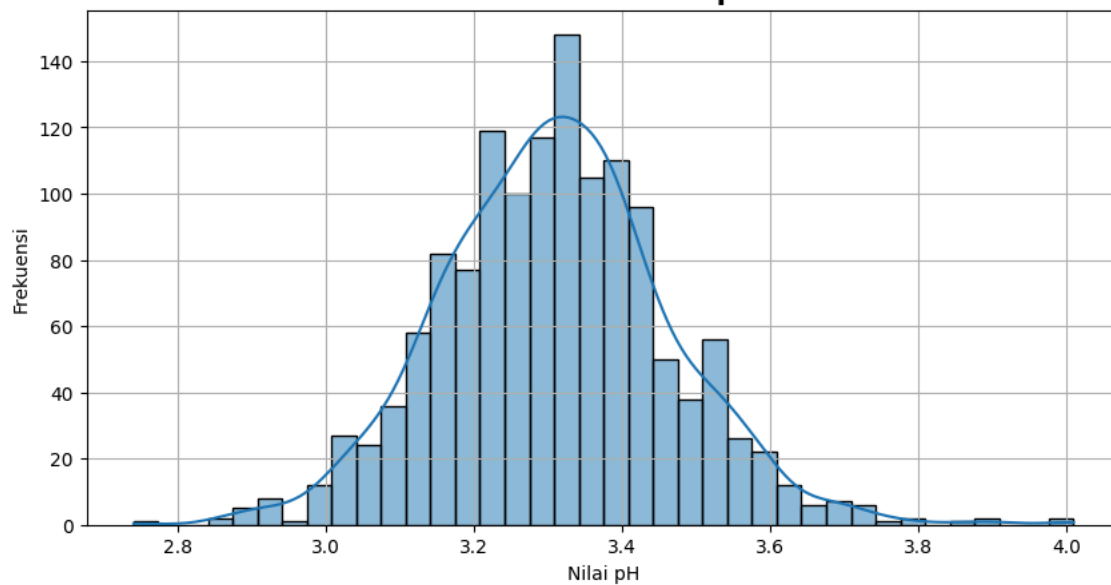
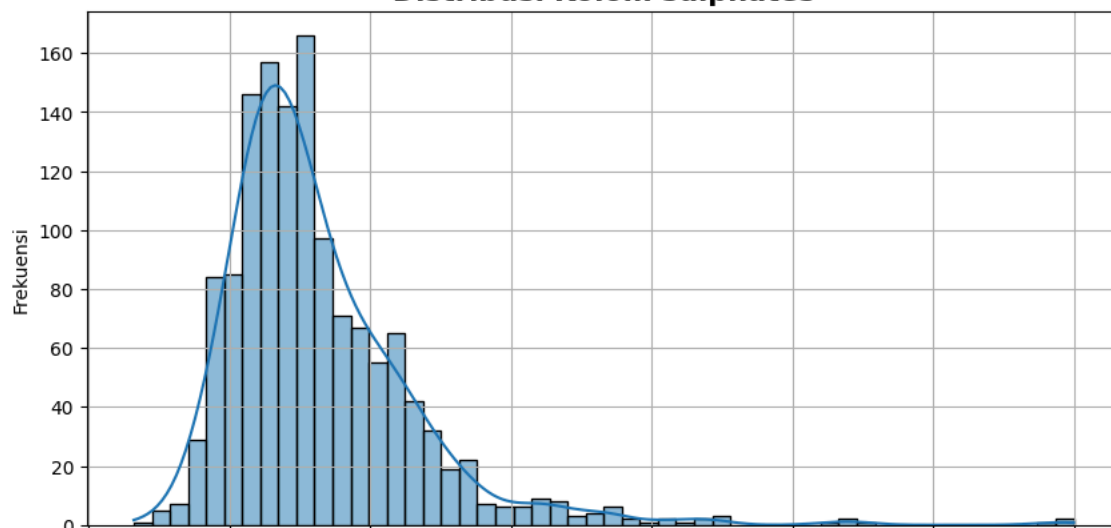


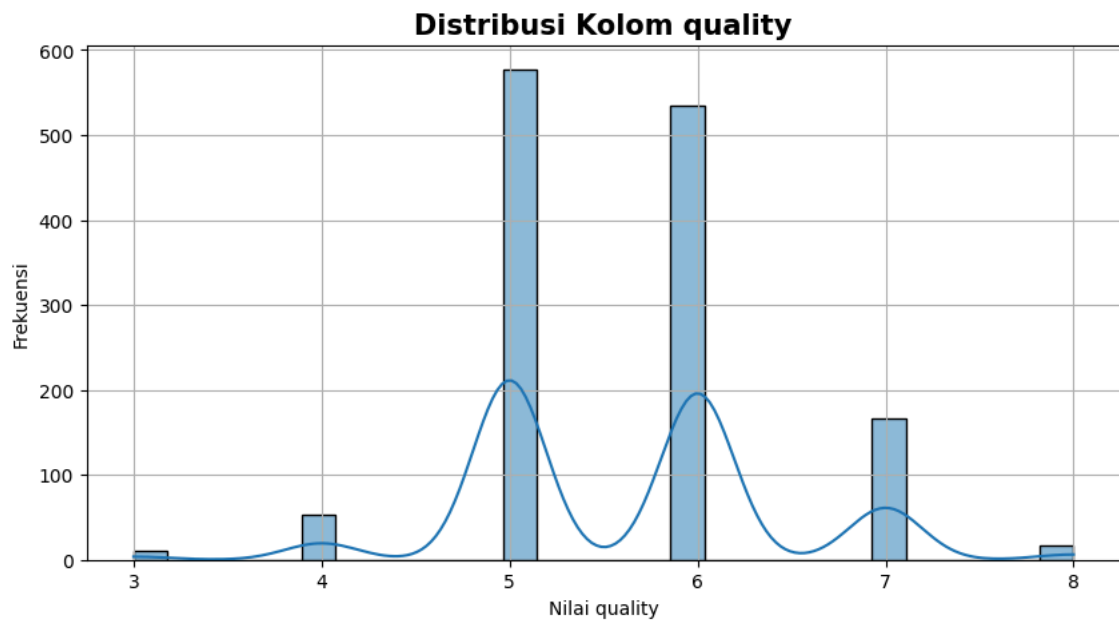
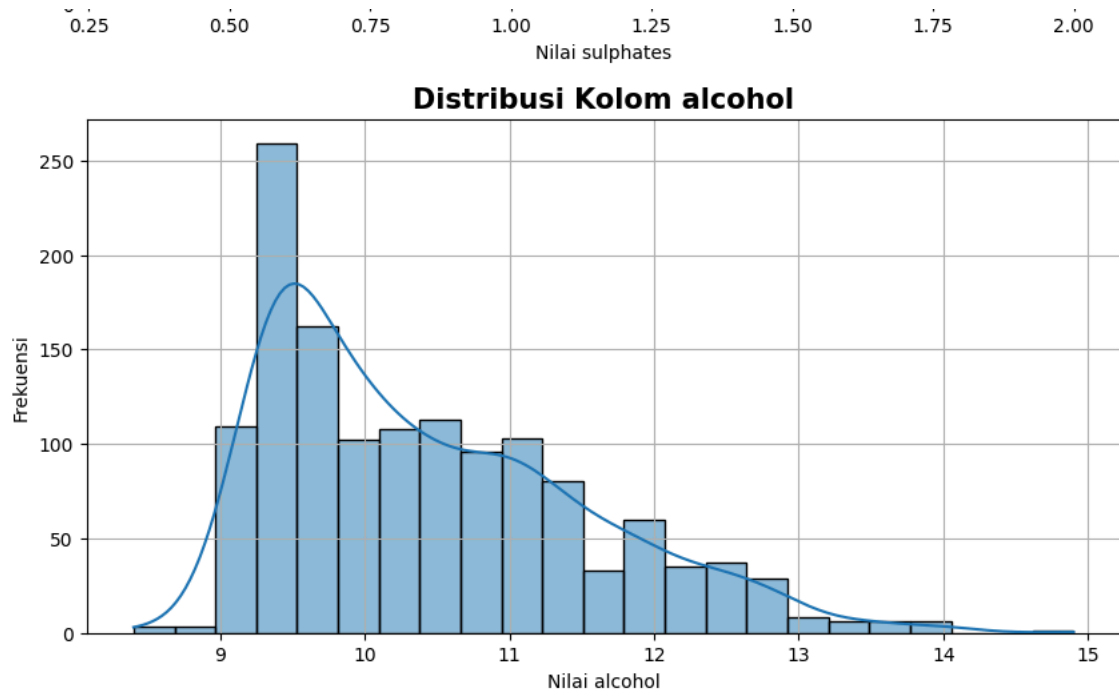
Visualisasi Distribusi Setiap Fitur:



**Distribusi Kolom chlorides****Distribusi Kolom free sulfur dioxide****Distribusi Kolom total sulfur dioxide**



**Distribusi Kolom density****Distribusi Kolom pH****Distribusi Kolom sulphates**



Langkah ini penting untuk memahami karakteristik sebaran data pada setiap kolom secara visual. Dengan menggunakan **histogram** yang dilengkapi dengan kurva **Kernel Density Estimate (KDE)**, kita dapat melihat bentuk distribusi dan mengidentifikasi potensi adanya *outlier* atau kemiringan data (*skewness*).

Hasil Visualisasi & Pengamatan Kunci:

- **Distribusi Normal:**
  - Fitur *density*, *pH*, dan *fixed acidity* menunjukkan distribusi yang mendekati kurva normal (lonceng), meskipun tidak sempurna. Ini berarti sebagian besar nilai terpusat di sekitar rata-rata.
- **Distribusi Miring ke Kanan (*Right-Skewed*):**
  - Banyak fitur yang distribusinya miring ke kanan. Ini terlihat dari "ekor" panjang di sisi kanan grafik.
  - Fitur-fitur seperti *residual sugar*, *chlorides*, *free sulfur dioxide*, dan *total sulfur dioxide* menunjukkan kemiringan yang sangat jelas. Ini mengindikasikan bahwa sebagian besar anggur memiliki nilai rendah untuk atribut ini, namun ada beberapa sampel dengan nilai yang sangat tinggi (*outlier*).
  - *alcohol* dan *sulphates* juga menunjukkan kemiringan ke kanan, menandakan mayoritas anggur memiliki kandungan alkohol dan sulfat yang moderat, dengan beberapa pengecualian yang nilainya lebih tinggi.
- **Distribusi Bimodal:**
  - *volatile acidity* menunjukkan dua puncak (bimodal), yang bisa berarti ada dua kelompok anggur yang berbeda dalam dataset berdasarkan keasaman volatilnya.
- **Distribusi Target quality:**
  - Grafik untuk kolom *quality* sangat jelas menunjukkan bahwa sebagian besar anggur dinilai dengan skor **5 dan 6**.
  - Sangat sedikit anggur yang memiliki skor kualitas rendah (3) atau sangat tinggi (8). Ini adalah konfirmasi visual dari **ketidakseimbangan data**, di mana model mungkin akan lebih mudah mempelajari karakteristik anggur berkualitas rata-rata dibandingkan anggur yang sangat baik atau sangat buruk.

Secara umum, visualisasi ini memperkuat temuan dari analisis statistik sebelumnya dan memberikan wawasan penting tentang perlunya penskalaan fitur (karena rentang nilai yang berbeda) dan penanganan outlier/data miring untuk beberapa model machine learning.

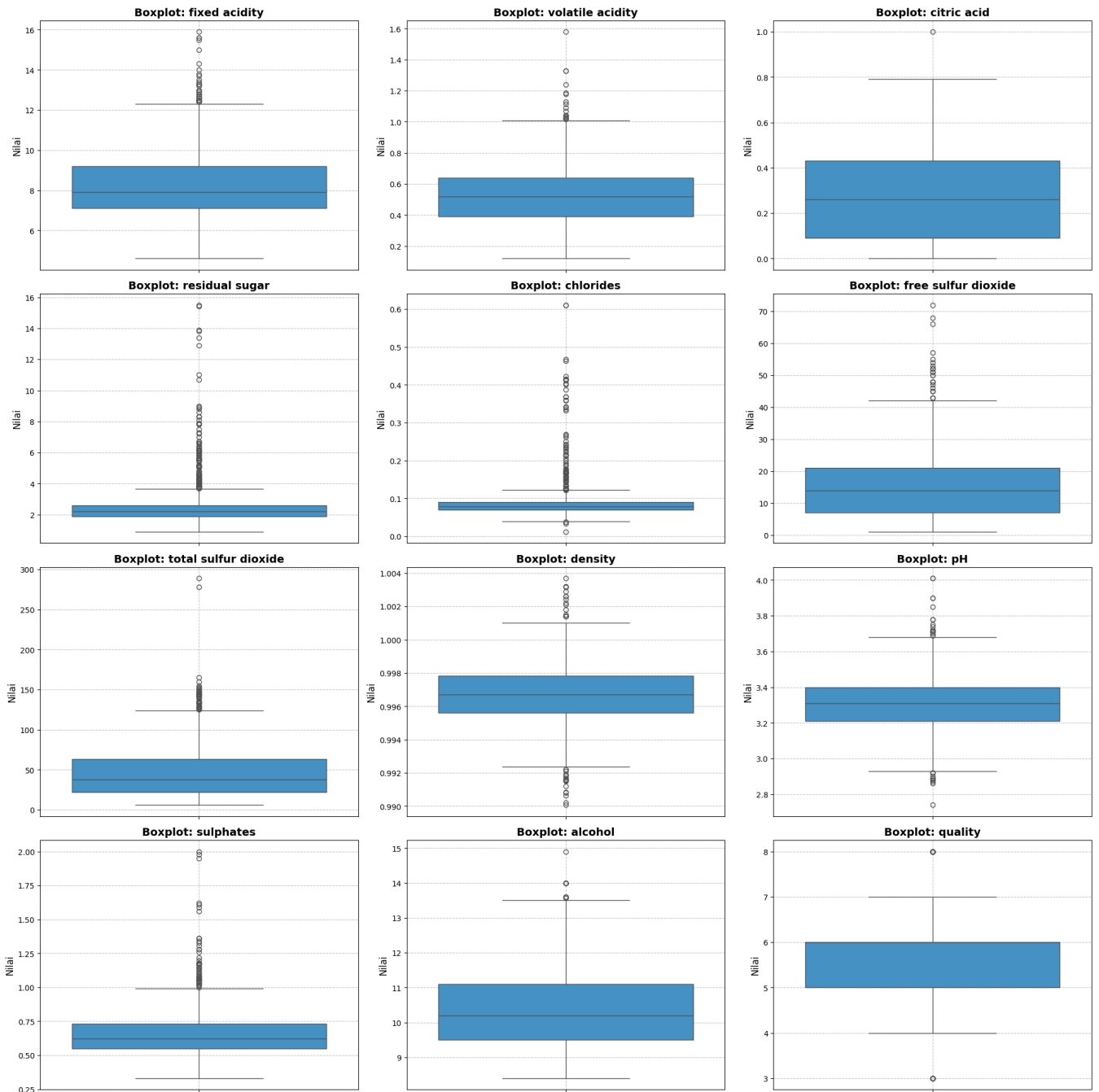
## ✓ 2.5 Identifikasi Outliers

Boxplot digunakan untuk melihat sebaran data dan mengidentifikasi adanya outliers (pencilan) pada setiap fitur.

```
1 # Visualisasi Boxplot untuk setiap fitur secara individual
2 print("\nBoxplot Individual untuk Setiap Fitur:")
3 columns = df.columns
4 num_columns = len(columns)
5 num_rows = (num_columns + 2) // 3 # Atur 3 kolom per baris
6
7 plt.figure(figsize=(20, 5 * num_rows)) # Sesuaikan ukuran figure berdasarkan jumlah baris
8
9 for i, column in enumerate(columns):
10     plt.subplot(num_rows, 3, i + 1)
11     sns.boxplot(y=df[column], color='#3498db') # Menggunakan seaborn untuk tampilan yang lebih baik dan konsisten
12     plt.title(f'Boxplot: {column}', fontsize=14, fontweight='bold')
13     plt.ylabel('Nilai', fontsize=12)
14     plt.grid(True, linestyle='--', alpha=0.7)
15
16 plt.tight_layout() # Menyesuaikan layout agar tidak tumpang tindih
17 plt.show()
```



Boxplot Individual untuk Setiap Fitur:



**Boxplot** adalah alat visual yang sangat efektif untuk memahami sebaran data dan mengidentifikasi keberadaan **outliers** (pencilan) pada setiap fitur.

Setiap "kotak" dalam grafik ini merepresentasikan *Interquartile Range (IQR)*, yaitu rentang antara kuartil pertama (Q1 atau 25%) dan kuartil ketiga (Q3 atau 75%). Garis di tengah kotak adalah median (50%). Garis vertikal (disebut *whiskers*) memanjang dari kotak untuk menunjukkan rentang data, biasanya hingga 1.5 kali IQR. **Titik-titik di luar whiskers** adalah outlier.

Hasil Visualisasi & Pengamatan Kunci:

- **Banyaknya Outlier:** Hampir semua fitur fisikokimia, kecuali citric acid dan pH, menunjukkan adanya outlier. Outlier paling ekstrem terlihat pada fitur total sulfur dioxide, free sulfur dioxide, dan residual sugar.
- **Konfirmasi Kemiringan Data:** Adanya banyak outlier di sisi atas (*upper outliers*) pada fitur seperti residual sugar, chlorides, dan sulphates mengonfirmasi kembali bahwa distribusi data pada fitur-fitur ini sangat miring ke kanan (*right-skewed*).
- **Fitur dengan Sebaran Normal:** Fitur pH dan density menunjukkan distribusi yang paling simetris (mendekati normal) dengan sedikit atau tanpa outlier yang signifikan. Ini menandakan bahwa sebagian besar nilai untuk fitur ini terpusat di tengah dengan sebaran yang wajar.
- **Variabel Target quality:** Kolom quality juga memiliki outlier, yaitu pada nilai 3 dan 8, yang menegaskan bahwa anggur dengan kualitas sangat rendah atau sangat tinggi adalah kasus yang jarang terjadi dalam dataset ini.

Tindak Lanjut:

Kehadiran outlier ini penting untuk diketahui. Meskipun Random Forest secara umum cukup tahan (*robust*) terhadap outlier, keberadaannya bisa memengaruhi metrik statistik seperti rata-rata dan standar deviasi. Untuk model lain yang lebih sensitif, penanganan outlier mungkin diperlukan. Namun, untuk kasus ini, kita akan melanjutkan tanpa menghapus outlier terlebih dahulu, mengingat ketahanan alami dari model Random Forest.

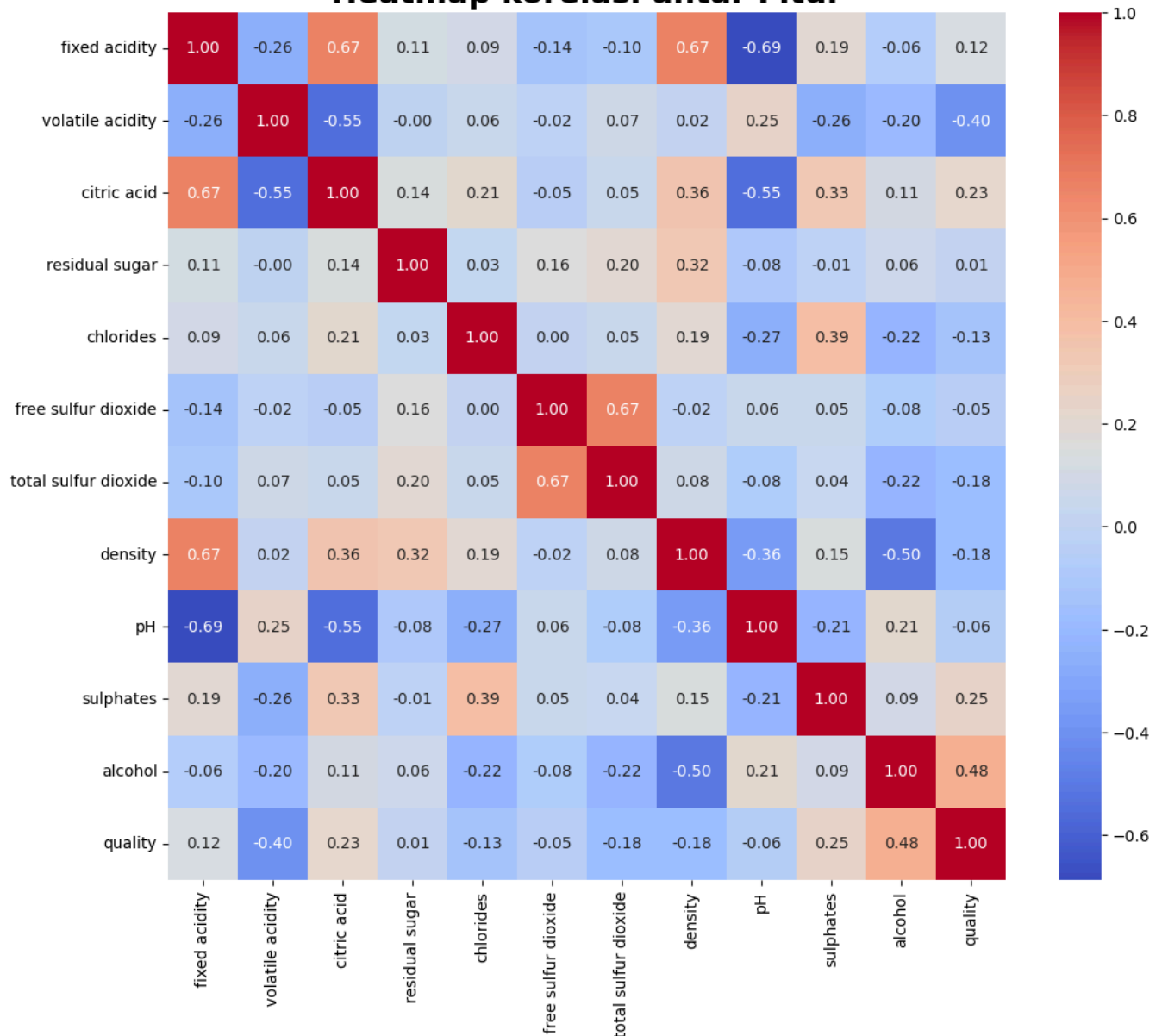
## ✓ 2.6 Analisis Korelasi

Heatmap korelasi menunjukkan hubungan linear antar variabel. Korelasi yang tinggi antara fitur-fitur independen (multikolinearitas) dapat memengaruhi interpretasi model.

```
1 plt.figure(figsize=(12, 10))
2 sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm', fmt='.2f')
3 plt.title('Heatmap korelasi antar Fitur', fontsize=20, fontweight='bold')
4 plt.show()
```



## Heatmap korelasi antar Fitur



**Heatmap korelasi** adalah visualisasi yang sangat berguna untuk melihat **hubungan linear antar variabel** dalam dataset. Setiap sel dalam heatmap menunjukkan koefisien korelasi antara dua variabel.

- **Skala Warna:**

- Warna **merah tua** menunjukkan korelasi positif yang kuat (mendekati +1). Artinya, jika satu variabel meningkat, variabel lainnya cenderung meningkat juga.
- Warna **biru tua** menunjukkan korelasi negatif yang kuat (mendekati -1). Artinya, jika satu variabel meningkat, variabel lainnya cenderung menurun.
- Warna **netral** (sekitar putih atau abu-abu muda) menunjukkan korelasi yang lemah atau tidak ada korelasi linear (mendekati 0).

- **Nilai dalam Sel:** Angka di setiap sel adalah nilai koefisien korelasi Pearson, yang berkisar antara -1 hingga +1. Format `fmt=' .2f'` memastikan nilai ditampilkan dengan dua angka desimal.

Hasil Visualisasi & Pengamatan Kunci:

- **Korelasi dengan Target ( quality ):**

- **alcohol** memiliki korelasi positif tertinggi dengan **quality (0.48)**. Ini mengindikasikan bahwa anggur dengan kandungan alkohol yang lebih tinggi cenderung memiliki skor kualitas yang lebih baik.
- **volatile acidity** menunjukkan korelasi negatif tertinggi dengan **quality (-0.40)**. Ini berarti semakin tinggi keasaman volatil, kualitas anggur cenderung semakin rendah.
- **sulphates** juga menunjukkan korelasi positif yang cukup baik dengan **quality (0.25)**.
- **citric acid** memiliki korelasi positif moderat dengan **quality (0.23)**.
- Fitur lain seperti **residual sugar**, **free sulfur dioxide**, dan **pH** menunjukkan korelasi yang sangat lemah dengan **quality**.

- **Korelasi Antar Fitur (Potensi Multikolinearitas):**

- Terdapat korelasi positif yang cukup kuat antara **fixed acidity** dan **citric acid (0.67)**. Ini masuk akal karena asam sitrat adalah salah satu komponen dari keasaman tetap.
- **fixed acidity** juga berkorelasi positif dengan **density (0.67)** dan negatif dengan **pH (-0.69)**. Ini juga logis karena peningkatan asam akan meningkatkan densitas dan menurunkan pH.
- **free sulfur dioxide** dan **total sulfur dioxide** memiliki korelasi positif yang cukup tinggi (**0.67**), yang memang diharapkan.

Implikasi untuk Pemodelan:

- Fitur-fitur dengan korelasi yang lebih tinggi terhadap **quality** (baik positif maupun negatif) kemungkinan besar akan menjadi prediktor yang lebih penting dalam model kita.
- Adanya korelasi yang cukup tinggi antar beberapa fitur input (seperti **fixed acidity** dan **citric acid**) menunjukkan potensi adanya **multikolinearitas**. Meskipun Random Forest relatif tahan terhadap multikolinearitas dibandingkan model regresi linear, ini adalah sesuatu yang perlu diingat, terutama jika kita ingin menginterpretasikan kontribusi masing-masing fitur secara individual.

Analisis korelasi ini memberikan panduan awal yang baik tentang fitur mana yang mungkin paling berpengaruh dan bagaimana fitur-fitur tersebut saling terkait.

## ✓ 3. Feature Engineering dan Data Splitting

Pada bagian ini, kita akan mempersiapkan data untuk pemodelan.

### ✓ 3.1 Transformasi Variabel Target

Sesuai dengan deskripsi dataset, kita akan mengubah masalah ini dari regresi menjadi klasifikasi. Kolom target **quality** akan diubah menjadi variabel biner:

- **1 (Baik):** Jika **quality**  $\geq 7$
- **0 (Buruk):** Jika **quality**  $< 7$

```
1 # Mengubah kolom 'quality' menjadi kategori biner (0 = bad, 1 = good)
2 df['quality'] = df['quality'].apply(lambda x: 1 if x >= 7 else 0)
3
4 print("Hasil transformasi kolom 'quality':")
5 df['quality'].head(10)
```

Hasil transformasi kolom 'quality':

|    | quality |
|----|---------|
| 0  | 0       |
| 1  | 0       |
| 2  | 0       |
| 3  | 0       |
| 5  | 0       |
| 6  | 0       |
| 7  | 1       |
| 8  | 1       |
| 9  | 0       |
| 10 | 0       |

dtype: int64

### ✓ 3.2 Pemeriksaan Keseimbangan Kelas (Class Balance)

Setelah kolom target **quality** ditransformasi menjadi kategori biner (0 untuk "Buruk" dan 1 untuk "Baik"), langkah krusial selanjutnya adalah memeriksa **distribusi atau keseimbangan kelas**. Ketidakseimbangan kelas terjadi ketika jumlah sampel pada satu kelas jauh lebih banyak daripada kelas lainnya. Hal ini dapat menyebabkan model machine learning menjadi bias dan lebih cenderung memprediksi kelas mayoritas, sehingga performanya pada kelas minoritas menjadi buruk.

#### ✓ 3.2.1. Memeriksa distribusi kelas setelah transformasi

Kode di bawah ini menghitung jumlah kemunculan dan persentase dari setiap kelas pada variabel target **quality**.

```

1 # Menghitung jumlah kemunculan setiap nilai di kolom 'quality'
2 jumlah_kualitas = df['quality'].value_counts()
3
4 # Menghitung persentase kemunculan setiap nilai
5 persentase_kualitas = df['quality'].value_counts(normalize=True) * 100
6
7 # Menggabungkan jumlah dan persentase ke dalam satu DataFrame
8 # menggunakan pandas.concat
9 df_kualitas_distribusi = pd.concat([jumlah_kualitas, persentase_kualitas], axis=1)
10
11 # Memberi nama kolom
12 df_kualitas_distribusi.columns = ['Jumlah', 'Persentase (%)']
13
14 # Memformat kolom 'Persentase (%)' menjadi 2 angka di belakang koma
15 df_kualitas_distribusi['Persentase (%)'] = df_kualitas_distribusi['Persentase (%)'].round(2)
16
17 # Menampilkan hasilnya
18 print("Distribusi Kualitas Anggur (0: Buruk, 1: Baik) dengan Persentase:")
19 display(df_kualitas_distribusi)

```

↗ Distribusi Kualitas Anggur (0: Buruk, 1: Baik) dengan Persentase:

|         | Jumlah | Persentase (%) |  |
|---------|--------|----------------|--|
| quality |        |                |  |
| 0       | 1175   | 86.46          |  |
| 1       | 184    | 13.54          |  |

Next steps:

[Generate code with df\\_kualitas\\_distribusi](#)

[View recommended plots](#)

[New interactive sheet](#)

Dari tabel di atas, kita dapat melihat:

- **Kelas 0 (Buruk):** Terdapat **1175 sampel**, yang merupakan **86.46%** dari total dataset.
- **Kelas 1 (Baik):** Terdapat **184 sampel**, yang hanya merupakan **13.54%** dari total dataset.

Ini menunjukkan adanya **ketidakseimbangan kelas yang signifikan**. Kelas "Buruk" (0) adalah kelas mayoritas, sedangkan kelas "Baik" (1) adalah kelas minoritas.

### ✓ 3.2.2. Visualisasi Distribusi Kelas

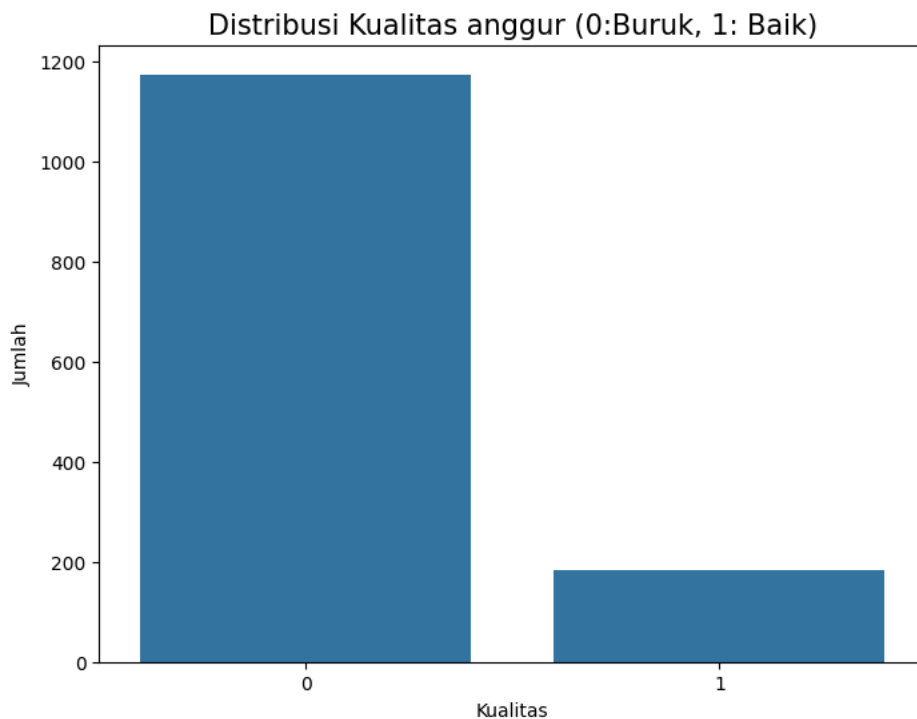
Untuk memperjelas ketidakseimbangan ini, kita akan memvisualisasikannya menggunakan *countplot* dari *seaborn*.

```

1 plt.figure(figsize=(8, 6))
2 sns.countplot(data=df, x="quality")
3 plt.title('Distribusi Kualitas anggur (0:Buruk, 1: Baik)', fontsize=15)
4 plt.xlabel('Kualitas')
5 plt.ylabel('Jumlah')
6 plt.show()

```





Grafik batang di atas secara visual mengonfirmasi temuan dari tabel. Batang untuk kualitas "Buruk" (0) jauh lebih tinggi dibandingkan batang untuk kualitas "Baik" (1).

Implikasi untuk Pemodelan:

Ketidakseimbangan kelas ini adalah masalah umum dalam tugas klasifikasi. Jika tidak ditangani, model akan cenderung "malas" dan lebih sering memprediksi kelas mayoritas karena itu akan memberikan akurasi yang tinggi secara keseluruhan, meskipun performanya buruk dalam mengidentifikasi kelas minoritas (yang seringkali justru lebih penting).

Oleh karena itu, dalam tahap pemodelan selanjutnya, kita perlu menerapkan strategi untuk mengatasi ketidakseimbangan ini, seperti:

#### 1. Penggunaan Metrik Evaluasi yang Tepat:

- Selain akurasi, kita akan fokus pada *precision*, *recall*, dan *F1-score*, terutama untuk kelas minoritas.

#### 2. Teknik Resampling:

- Seperti *oversampling* pada kelas minoritas (misalnya, SMOTE) atau *undersampling* pada kelas mayoritas.

#### 3. Penyesuaian Bobot Kelas:

- Memberikan bobot yang lebih tinggi pada kelas minoritas saat melatih model (seperti `class_weight='balanced'` pada Random Forest).

Kesadaran akan ketidakseimbangan ini sejak awal adalah kunci untuk membangun model klasifikasi yang robust dan dapat diandalkan.

### ✓ 3.3 Pemisahan Data (Train-Test Split)

Sebelum kita melatih model, dataset perlu dibagi menjadi dua bagian utama:

- Set Pelatihan (Train Set):** Bagian data yang akan digunakan untuk "mengajari" model kita pola-pola yang ada.
- Set Pengujian (Test Set):** Bagian data yang akan digunakan untuk mengevaluasi seberapa baik model yang telah dilatih dapat melakukan prediksi pada data baru yang belum pernah dilihat sebelumnya.

Kita akan memisahkan dataset menjadi data fitur (X) dan data target (y), kemudian membaginya menjadi set pelatihan dan set pengujian. Penggunaan `stratify=y` sangat penting untuk memastikan proporsi kelas pada data latih dan data uji sama dengan proporsi pada data asli, terutama pada kasus dataset tidak seimbang.

Pertama, kita memisahkan kolom-kolom fitur (variabel input) dan kolom target (variabel output).

- `x`: Berisi semua kolom kecuali kolom `quality`. Ini adalah fitur-fitur yang akan digunakan model untuk belajar.
- `y`: Hanya berisi kolom `quality` (yang sudah biner 0 atau 1). Ini adalah apa yang ingin diprediksi oleh model.

Selanjutnya, kita menggunakan fungsi `train_test_split` dari `sklearn` untuk membagi data.

- `test_size=0.2`: Menentukan bahwa **20%** dari data akan digunakan sebagai set pengujian, dan sisanya **80%** sebagai set pelatihan.

- `random_state=42`: Mengatur *seed* untuk generator angka acak. Ini memastikan bahwa jika kita menjalankan kode ini lagi, pembagian datanya akan selalu sama. Ini penting untuk reproduktifitas hasil.
- `stratify=y`: Ini adalah parameter yang **sangat penting**, terutama karena kita tahu dataset kita tidak seimbang. Dengan menggunakan `stratify=y`, kita memastikan bahwa proporsi kelas (0 dan 1) pada variabel target `y` akan **dijaga sama** baik di set pelatihan maupun di set pengujian. Artinya, jika kelas "Baik" (1) adalah 13.54% di dataset asli, maka di set pelatihan dan set pengujian juga akan sekitar 13.54%. Ini mencegah situasi di mana salah satu set (misalnya, set uji) secara kebetulan memiliki sangat sedikit atau tidak ada sampel dari kelas minoritas.

```
1 # Memisahkan fitur (X) dan target (y)
2 X = df.drop('quality', axis=1)
3 y = df['quality']
4
5 # Membagi data menjadi 80% data latih dan 20% data uji
6 # Menggunakan stratify=y untuk menjaga proporsi kelas
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
8
9 print(f'Ukuran X_train: {X_train.shape}')
10 print(f'Ukuran X_test: {X_test.shape}')
11 print(f'Ukuran y_train: {y_train.shape}')
12 print(f'Ukuran y_test: {y_test.shape}')
```

```
↗ Ukuran X_train: (1087, 11)
   Ukuran X_test: (272, 11)
   Ukuran y_train: (1087,)
   Ukuran y_test: (272,)
```

Verifikasi proporsi kelas setelah stratifikasi

```
1 print("\nProporsi kelas pada y_train:")
2 print(y_train.value_counts(normalize=True) * 100)
3 print("\nProporsi kelas pada y_test:")
4 print(y_test.value_counts(normalize=True) * 100)
```

```
↗ \nProporsi kelas pada y_train:
   quality
0    86.476541
1    13.523459
Name: proportion, dtype: float64
\nProporsi kelas pada y_test:
   quality
0    86.397059
1    13.602941
Name: proportion, dtype: float64
```

Terlihat bahwa proporsi kelas "Baik" dan "Buruk" pada `y_train` dan `y_test` sangat mirip dengan proporsi pada dataset keseluruhan (86.46% untuk Buruk dan 13.54% untuk Baik). Ini menunjukkan bahwa parameter `stratify=y` bekerja dengan baik.

Dengan pembagian data yang tepat ini, kita siap untuk melanjutkan ke tahap pra-pemrosesan lebih lanjut dan pelatihan model.

Melihat nilai `x` dan `y`

```
1 print("Nilai X:")
2 display(X)
```

Nilai X:

|      | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH   | sulphates | alcohol |
|------|---------------|------------------|-------------|----------------|-----------|---------------------|----------------------|---------|------|-----------|---------|
| 0    | 7.4           | 0.700            | 0.00        | 1.9            | 0.076     | 11.0                | 34.0                 | 0.99780 | 3.51 | 0.56      | 9.4     |
| 1    | 7.8           | 0.880            | 0.00        | 2.6            | 0.098     | 25.0                | 67.0                 | 0.99680 | 3.20 | 0.68      | 9.8     |
| 2    | 7.8           | 0.760            | 0.04        | 2.3            | 0.092     | 15.0                | 54.0                 | 0.99700 | 3.26 | 0.65      | 9.8     |
| 3    | 11.2          | 0.280            | 0.56        | 1.9            | 0.075     | 17.0                | 60.0                 | 0.99800 | 3.16 | 0.58      | 9.8     |
| 5    | 7.4           | 0.660            | 0.00        | 1.8            | 0.075     | 13.0                | 40.0                 | 0.99780 | 3.51 | 0.56      | 9.4     |
| ...  | ...           | ...              | ...         | ...            | ...       | ...                 | ...                  | ...     | ...  | ...       | ...     |
| 1593 | 6.8           | 0.620            | 0.08        | 1.9            | 0.068     | 28.0                | 38.0                 | 0.99651 | 3.42 | 0.82      | 9.5     |
| 1594 | 6.2           | 0.600            | 0.08        | 2.0            | 0.090     | 32.0                | 44.0                 | 0.99490 | 3.45 | 0.58      | 10.5    |
| 1595 | 5.9           | 0.550            | 0.10        | 2.2            | 0.062     | 39.0                | 51.0                 | 0.99512 | 3.52 | 0.76      | 11.2    |
| 1597 | 5.9           | 0.645            | 0.12        | 2.0            | 0.075     | 32.0                | 44.0                 | 0.99547 | 3.57 | 0.71      | 10.2    |
| 1598 | 6.0           | 0.310            | 0.47        | 3.6            | 0.067     | 18.0                | 42.0                 | 0.99549 | 3.39 | 0.66      | 11.0    |

1359 rows × 11 columns

Next steps:

[Generate code with X](#)[View recommended plots](#)[New interactive sheet](#)

```
1 print("Nilai y:")
2 display(y)
```

Nilai y:

|      | quality |
|------|---------|
| 0    | 0       |
| 1    | 0       |
| 2    | 0       |
| 3    | 0       |
| 5    | 0       |
| ...  | ...     |
| 1593 | 0       |
| 1594 | 0       |
| 1595 | 0       |
| 1597 | 0       |
| 1598 | 0       |

1359 rows × 1 columns

dtype: int64

## 4. Pelatihan Model dan Tuning Hyperparameter

Setelah data dipersiapkan, langkah selanjutnya adalah melatih model klasifikasi kita, yaitu **Random Forest Classifier**. Untuk mendapatkan performa model yang optimal, kita tidak hanya melatih model dengan parameter default, tetapi juga melakukan **tuning hyperparameter**.

**Hyperparameter** adalah parameter yang nilainya kita tentukan *sebelum* proses pelatihan model dimulai (berbeda dengan parameter model yang dipelajari selama pelatihan). Contoh hyperparameter pada Random Forest adalah jumlah pohon (`n_estimators`) atau kedalaman maksimum setiap pohon (`max_depth`).

Kita akan menggunakan `GridSearchCV` dari `sklearn` untuk mencari kombinasi hyperparameter terbaik secara sistematis. `GridSearchCV` akan mencoba semua kemungkinan kombinasi hyperparameter yang kita definisikan dan mengevaluasinya menggunakan *cross-validation*.

### Penjelasan Kode:


1. `param_grid`: Kita mendefinisikan sebuah *dictionary* yang berisi daftar hyperparameter dan nilai-nilai yang ingin diuji untuk masing-masing hyperparameter tersebut.
  - `n_estimators`: Jumlah pohon yang akan dibangun.
  - `max_depth`: Seberapa dalam setiap pohon dapat tumbuh. `None` berarti pohon akan tumbuh sampai semua daun murni atau sampai mencapai `min_samples_split`.
  - `min_samples_split`: Jumlah sampel minimum yang harus ada di sebuah node agar node tersebut bisa dipecah lagi.

- `min_samples_leaf`: Jumlah sampel minimum yang harus ada di setiap *leaf node* (ujung pohon).
2. `RandomForestClassifier(...)`: Kita membuat instance dari model Random Forest.
    - `class_weight='balanced'`: Ini adalah poin penting. Karena kita memiliki masalah ketidakseimbangan kelas, parameter ini akan secara otomatis menyesuaikan bobot kelas sedemikian rupa sehingga kelas minoritas (anggur "Baik") mendapatkan perhatian lebih selama pelatihan.
    - `random_state=42`: Untuk konsistensi hasil.
  3. `GridSearchCV(...)`:
    - `estimator=model`: Model yang akan di-tuning.
    - `param_grid=param_grid`: Kisi-kisi hyperparameter yang akan dicoba.
    - `cv=5`: Melakukan 5-fold cross-validation.
    - `n_jobs=-1`: Menggunakan semua prosesor yang tersedia untuk mempercepat pencarian.
    - `verbose=1`: Akan menampilkan pesan selama proses fitting, menunjukkan kombinasi mana yang sedang diuji.
    - `scoring='recall_weighted'`: Karena data kita tidak seimbang, akurasi saja tidak cukup. Kita memilih `recall_weighted` sebagai metrik utama untuk dioptimalkan oleh `GridSearchCV`. *Recall weighted* menghitung rata-rata recall dari setiap kelas, dengan bobot berdasarkan jumlah sampel sebenarnya di setiap kelas. Ini memberikan gambaran yang lebih baik tentang kemampuan model mengenali semua kelas, termasuk kelas minoritas.
  4. `grid_search.fit(X_train, y_train)`: Memulai proses pencarian hyperparameter terbaik pada data latih.
  5. `best_rf = grid_search.best_estimator_`: Menyimpan model terbaik yang ditemukan oleh `GridSearchCV`.
  6. **Output**: Menampilkan kombinasi hyperparameter terbaik dan skor `recall_weighted` rata-rata tertinggi yang dicapai selama cross-validation.

```

1 # Menentukan hyperparameter yang akan diuji
2 param_grid = {
3     'n_estimators': [50, 100, 200, 500],
4     'max_depth': [None, 10, 15, 20],
5     'min_samples_split': [2, 5, 10],
6     'min_samples_leaf': [1, 2, 4]
7 }
8
9 # Membuat model RandomForest. `class_weight='balanced'` digunakan untuk mengatasi masalah class imbalance.
10 model = RandomForestClassifier(random_state=42, class_weight='balanced')
11
12 # Menggunakan Grid Search dengan 5-fold cross-validation
13 grid_search = GridSearchCV(estimator=model,
14                             param_grid=param_grid,
15                             cv=5,
16                             n_jobs=-1, # Menggunakan semua core CPU yang tersedia
17                             verbose=1
18                             )
19
20 # Melatih model dengan Grid Search
21 grid_search.fit(X_train, y_train)
22
23 # Simpan model terbaik
24 best_rf = grid_search.best_estimator_
25
26 # Menampilka hyperparameter dan skor terbaik
27 print('\nHyperparameter terbaik:', grid_search.best_params_)
28 print('akurasi terbaik dari cross-validation:', grid_search.best_score_)

```

 Fitting 5 folds for each of 144 candidates, totalling 720 fits

```

Hyperparameter terbaik: {'max_depth': 15, 'min_samples_leaf': 2, 'min_samples_split': 2, 'n_estimators': 500}
akurasi terbaik dari cross-validation: 0.8859299031835285

```

#### Hasil & Interpretasi:

- **Output Fitting 5 folds for each of 144 candidates, totalling 720 fits:**
  - Ini menunjukkan bahwa `GridSearchCV` telah menguji 144 kombinasi hyperparameter yang berbeda (4 nilai `n_estimators` \* 4 nilai `max_depth` \* 3 nilai `min_samples_split` \* 3 nilai `min_samples_leaf` = 144). Karena kita menggunakan 5-fold cross-validation, maka total ada  $144 * 5 = 720$  model yang dilatih dan dievaluasi.
- **Hyperparameter Terbaik:**
  - Output menunjukkan kombinasi hyperparameter yang menghasilkan skor `recall_weighted` tertinggi selama cross-validation. Dalam kasus ini:
    - `'max_depth': 15`
    - `'min_samples_leaf': 2`

- 'min\_samples\_split': 2
- 'n\_estimators': 500

- **Skor Recall Weighted Terbaik:**

- **0.8859.** Ini adalah skor `recall_weighted` rata-rata dari 5-fold cross-validation menggunakan hyperparameter terbaik tersebut. Angka ini memberikan estimasi yang cukup baik tentang seberapa baik model kita akan bekerja pada data yang belum pernah dilihat.

Dengan model terbaik yang telah ditemukan melalui `GridSearchCV` ini (`best_rf`), kita sekarang siap untuk melakukan prediksi pada data uji dan mengevaluasi kinerjanya secara lebih komprehensif.

## ✓ 5. Evaluasi Model

Setelah mendapatkan model terbaik dari proses *tuning hyperparameter* (`best_rf`), langkah selanjutnya adalah mengevaluasi kinerjanya pada **data uji** (`x_test`, `y_test`), yaitu data yang belum pernah dilihat oleh model selama proses pelatihan. Ini memberikan estimasi yang lebih realistis tentang seberapa baik model akan berkinerja di dunia nyata.

### ✓ 5.1 Laporan Klasifikasi & Confusion Matrix

Kita akan menggunakan dua alat utama untuk evaluasi:

1. **Laporan Klasifikasi** (`classification_report`): Memberikan rincian metrik seperti presisi, recall, dan F1-score untuk setiap kelas.
2. **Confusion Matrix**: Memvisualisasikan performa klasifikasi model dalam bentuk matriks yang menunjukkan prediksi yang benar dan salah.

```
1 # Melakukan prediksi pada data uji
2 y_pred = best_rf.predict(X_test)
```

```
1 # Menampilkan laporan klasifikasi
2 print("Laporan Klasifikasi:")
3 print(classification_report(y_test, y_pred))
```

```
↔ Laporan Klasifikasi:
              precision    recall  f1-score   support

     0       0.91      0.96      0.93       235
     1       0.60      0.41      0.48        37

 accuracy          0.88       272
 macro avg       0.76      0.68      0.71       272
 weighted avg    0.87      0.88      0.87       272
```

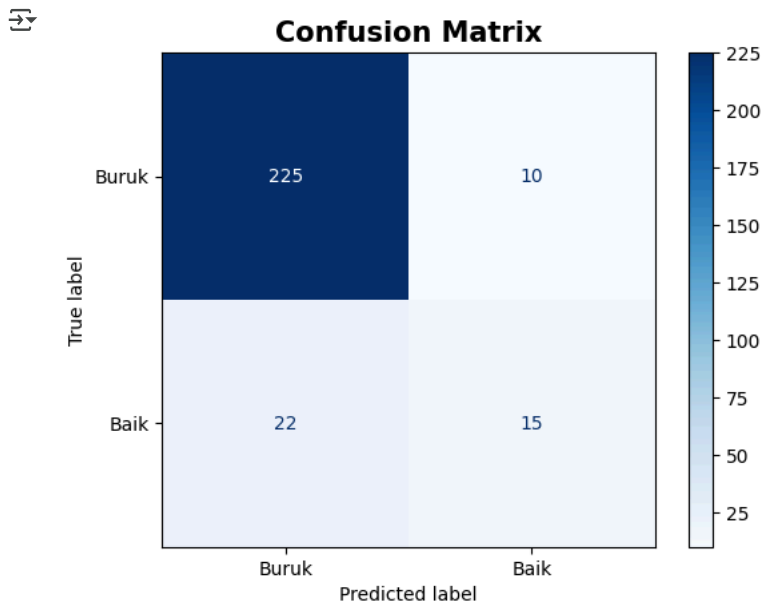
Laporan Klasifikasi memberikan metrik performa untuk masing-masing kelas ("Buruk (0)" dan "Baik (1)") serta rata-rata:

- **precision:**
  - **Untuk Kelas "Buruk (0)" (0.91):** Dari semua anggur yang diprediksi sebagai "Buruk", 91% di antaranya memang benar-benar "Buruk".
  - **Untuk Kelas "Baik (1)" (0.60):** Dari semua anggur yang diprediksi sebagai "Baik", 60% di antaranya memang benar-benar "Baik". Ini adalah peningkatan yang cukup baik dibandingkan model awal (sebelumnya 0.42).
- **recall:**
  - **Untuk Kelas "Buruk (0)" (0.96):** Model berhasil mengidentifikasi 96% dari semua anggur yang sebenarnya "Buruk".
  - **Untuk Kelas "Baik (1)" (0.41):** Model berhasil mengidentifikasi 41% dari semua anggur yang sebenarnya "Baik". Meskipun masih belum sempurna, ini adalah peningkatan signifikan dari model awal (sebelumnya 0.26), menunjukkan bahwa penggunaan `class_weight='balanced'` dan tuning hyperparameter (terutama dengan *scoring* `recall_weighted`) memberikan dampak positif.
- **f1-score:** Rata-rata harmonik dari precision dan recall.
  - **Untuk Kelas "Buruk (0)" (0.93):** Sangat baik.
  - **Untuk Kelas "Baik (1)" (0.48):** Cukup baik, dan merupakan peningkatan dari model awal (0.32). Ini menunjukkan keseimbangan yang lebih baik antara precision dan recall untuk kelas minoritas.
- **support:** Jumlah sampel sebenarnya untuk setiap kelas di data uji.
- **accuracy (0.88):** Akurasi keseluruhan model pada data uji adalah 88%.
- **macro avg:** Rata-rata metrik tanpa memperhitungkan proporsi kelas. F1-score 0.71 menunjukkan performa rata-rata yang lebih baik antar kelas.
- **weighted avg:** Rata-rata metrik dengan memperhitungkan proporsi kelas. F1-score 0.87 menunjukkan performa keseluruhan yang baik, tetapi kita tahu ini lebih dipengaruhi oleh kelas mayoritas.

```

1 # Membuat dan menampilkan confusion matrix
2 cm = confusion_matrix(y_test, y_pred)
3 ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=['Buruk', 'Baik']).plot(cmap='Blues')
4 plt.title('Confusion Matrix', fontsize=15, fontweight='bold')
5 plt.show()

```



Confusion matrix memberikan visualisasi yang lebih detail:

- **True Positives (TP)** untuk kelas "Baik" (pojok kanan bawah): **15**. Model memprediksi 15 anggur sebagai "Baik" dan prediksi tersebut benar.
- **True Negatives (TN)** untuk kelas "Buruk" (pojok kiri atas): **225**. Model memprediksi 225 anggur sebagai "Buruk" dan prediksi tersebut benar.
- **False Positives (FP)** untuk kelas "Baik" (Type I Error, pojok kanan atas): **10**. Model memprediksi 10 anggur sebagai "Baik", padahal sebenarnya "Buruk".
- **False Negatives (FN)** untuk kelas "Baik" (Type II Error, pojok kiri bawah): **22**. Model memprediksi 22 anggur sebagai "Buruk", padahal sebenarnya "Baik".

Dari confusion matrix:

- Precision untuk kelas "Baik" =  $TP / (TP + FP) = 15 / (15 + 10) = 15 / 25 = 0.60$ .
- Recall untuk kelas "Baik" =  $TP / (TP + FN) = 15 / (15 + 22) = 15 / 37 = 0.405$  (dibulatkan menjadi 0.41 di laporan).

Model yang telah di-tuning dengan GridSearchCV dan menggunakan `class_weight='balanced'` menunjukkan **peningkatan yang signifikan** dalam kemampuannya mengidentifikasi kelas minoritas ("Baik"), terutama pada metrik *recall* dan *F1-score* untuk kelas 1, dibandingkan dengan model awal. Meskipun *recall* untuk kelas "Baik" masih 0.41, ini sudah jauh lebih baik. Masih ada ruang untuk perbaikan lebih lanjut, mungkin dengan teknik *oversampling* seperti SMOTE jika *recall* yang lebih tinggi untuk kelas "Baik" sangat diinginkan, namun dengan risiko peningkatan *false positives*.

## ✓ 5.2 Kurva ROC (Receiver Operating Characteristic) dan AUC

**Kurva ROC** adalah alat evaluasi grafis yang sangat baik untuk masalah klasifikasi biner. Kurva ini memvisualisasikan kemampuan model dalam membedakan antara dua kelas (dalam kasus ini, "Baik" dan "Buruk") pada berbagai ambang batas (*threshold*) klasifikasi.

- **Sumbu X (False Positive Rate - FPR)**: Proporsi dari sampel negatif yang salah diklasifikasikan sebagai positif. (Anggur "Buruk" yang salah diprediksi sebagai "Baik").
  - $FPR = FP / (FP + TN)$
- **Sumbu Y (True Positive Rate - TPR atau Recall/Sensitivity)**: Proporsi dari sampel positif yang benar diklasifikasikan sebagai positif. (Anggur "Baik" yang benar diprediksi sebagai "Baik").
  - $TPR = TP / (TP + FN)$

**Nilai AUC (Area Under the Curve)** adalah ukuran numerik dari performa keseluruhan model.

- Nilai AUC berkisar antara 0 hingga 1.

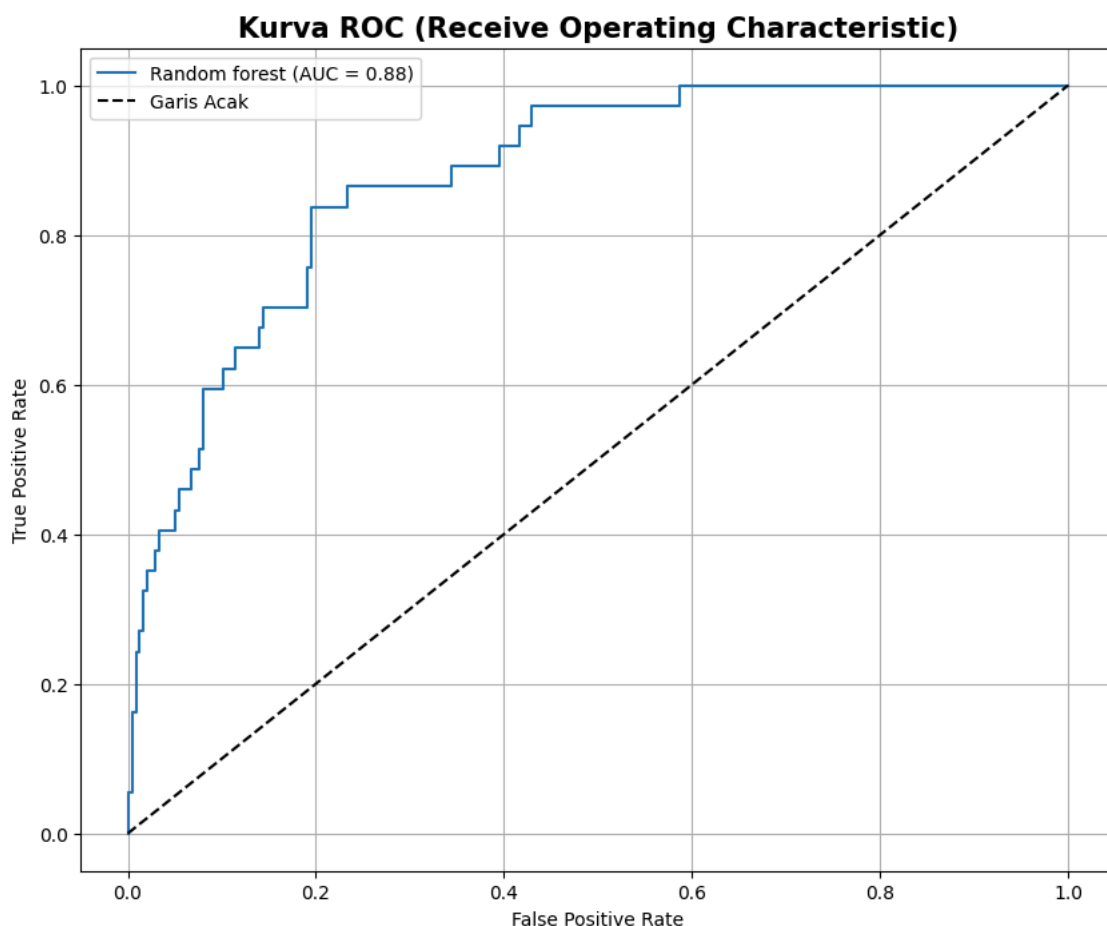
- AUC **0.5** menunjukkan model yang tidak lebih baik dari tebakan acak (diwakili oleh garis diagonal putus-putus pada grafik, disebut "Garis Acak").
- AUC **1.0** menunjukkan model yang sempurna dalam membedakan kelas.
- Semakin tinggi nilai AUC (mendekati 1), semakin baik kemampuan model dalam membedakan antara kelas positif dan negatif.

```
1 # Menghitung probabilitas prediksi
2 y_pred_proba = best_rf.predict_proba(X_test)[: , 1]
```

```
1 # Menghitung AUC Score
2 auc = roc_auc_score(y_test, y_pred_proba)
3 print(f"AUC score: {auc:.4f}")
```

AUC score: 0.8771

```
1 # Membuat plot ROC Curve
2 fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
3 plt.figure(figsize=(10, 8))
4 plt.plot(fpr, tpr, label=f'Random forest (AUC = {auc:.2f})')
5 plt.plot([0, 1], [0, 1], 'k--', label='Garis Acak')
6 plt.xlabel('False Positive Rate')
7 plt.ylabel('True Positive Rate')
8 plt.title('Kurva ROC (Receive Operating Characteristic)', fontsize=15,
9 fontweight='bold')
10 plt.legend(loc='best')
11 plt.grid()
12 plt.show()
```



- **AUC Score: 0.8771** (atau 0.88 jika dibulatkan) Nilai AUC sebesar **0.88** menunjukkan bahwa model kita memiliki kemampuan yang **sangat baik** dalam membedakan antara anggur berkualitas "Baik" dan "Buruk". Ini jauh lebih baik daripada tebakan acak (AUC = 0.5).
- **Kurva ROC:**
  - Kurva ROC untuk model Random Forest kita (garis biru) berada **jauh di atas** garis diagonal acak (garis hitam putus-putus).
  - Semakin kurva mendekati pojok kiri atas (di mana TPR = 1 dan FPR = 0), semakin baik performa model. Kurva kita menunjukkan tren yang baik ke arah tersebut.
  - Bentuk kurva yang naik tajam di awal dan kemudian mendatar menunjukkan bahwa model dapat mencapai *true positive rate* yang tinggi dengan *false positive rate* yang relatif rendah pada beberapa *threshold*.

Meskipun laporan klasifikasi sebelumnya menunjukkan tantangan dalam *recall* untuk kelas "Baik", nilai AUC yang tinggi (0.88) mengindikasikan bahwa model secara keseluruhan memiliki daya diskriminatif yang baik. Artinya, jika kita menyesuaikan *threshold* probabilitas untuk klasifikasi, kita berpotensi menemukan titik operasi di mana model lebih sensitif terhadap kelas "Baik", meskipun mungkin dengan mengorbankan sedikit presisi (meningkatkan *false positives*).

AUC adalah metrik yang berguna, terutama untuk dataset yang tidak seimbang, karena ia mengevaluasi kinerja model di semua *threshold* klasifikasi yang mungkin.

## 6. Analisis Kepentingan Fitur (Feature Importance)

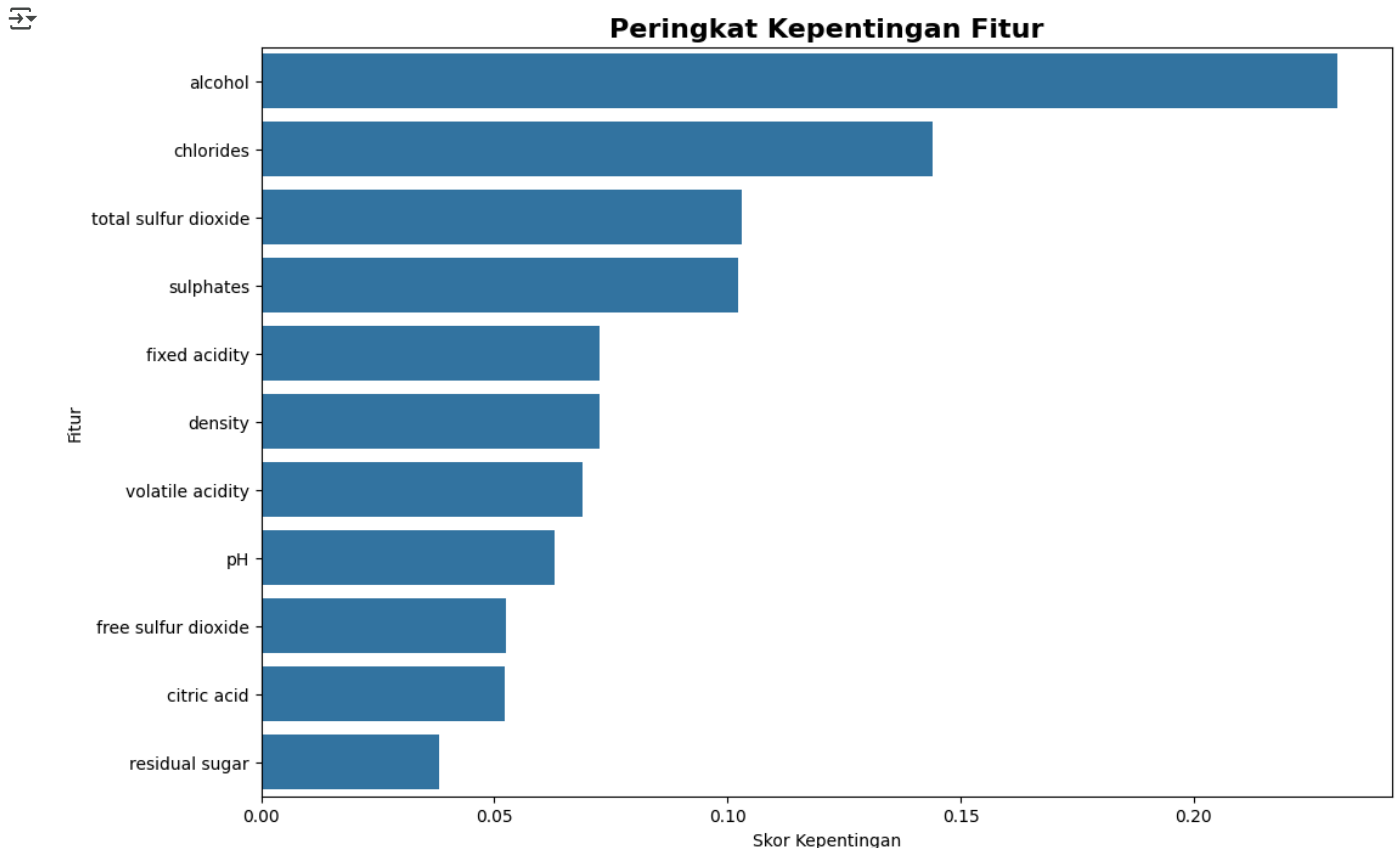
Salah satu keunggulan model berbasis pohon seperti Random Forest adalah kemampuannya untuk memberikan ukuran **kepentingan fitur** (*feature importance*). Ini memberi tahu kita seberapa besar kontribusi masing-masing fitur input dalam membuat prediksi.

Fitur yang lebih penting memiliki dampak yang lebih besar pada keputusan yang dibuat oleh pohon-pohon di dalam *forest*. Informasi ini sangat berguna untuk:

- Memahami faktor-faktor kunci yang memengaruhi kualitas anggur.
- Potensi untuk seleksi fitur (memilih fitur yang paling relevan untuk model yang lebih sederhana).

```
1 # Membuat series yang menyimpan feature importance dari model dan feature names dari training data
2 feature_importances = pd.Series(best_rf.feature_importances_, index=X.columns.sort_values(ascending=False))
```

```
1 # Membuat plot bar chart
2 plt.figure(figsize=(12, 8))
3 feature_importances_sorted = feature_importances.sort_values(ascending=False)
4 sns.barplot(x=feature_importances_sorted, y=feature_importances_sorted.index)
5 plt.title('Peringkat Kepentingan Fitur', fontsize=16, fontweight='bold')
6 plt.xlabel('Skor Kepentingan')
7 plt.ylabel('Fitur')
8 plt.show()
```



```
1 print('Fitur Paling Berpengaruh')
2 feature_importances.sort_values(ascending=False)
```



Fitur Paling Berpengaruh

|                      | 0        |
|----------------------|----------|
| alcohol              | 0.230887 |
| chlorides            | 0.144058 |
| total sulfur dioxide | 0.103051 |
| sulphates            | 0.102256 |
| fixed acidity        | 0.072671 |
| density              | 0.072492 |
| volatile acidity     | 0.069007 |
| pH                   | 0.062800 |
| free sulfur dioxide  | 0.052376 |
| citric acid          | 0.052307 |
| residual sugar       | 0.038095 |

dtype: float64

Grafik batang dan output numerik di atas menunjukkan peringkat kepentingan fitur berdasarkan kontribusinya dalam model Random Forest yang telah di-tuning:

1. **alcohol (Skor ~0.231)**: Fitur ini adalah yang **paling berpengaruh** dalam menentukan kualitas anggur. Semakin tinggi kandungan alkohol, semakin besar kemungkinannya anggur dianggap "Baik". Ini konsisten dengan temuan kita dari analisis korelasi.
2. **chlorides (Skor ~0.144)**: Klorida menjadi fitur terpenting kedua. Ini sedikit mengejutkan karena korelasinya dengan *quality* (-0.13) tidak sebesar fitur lain seperti *sulphates* atau *volatile acidity*. Namun, dalam model non-linear seperti Random Forest, interaksi antar fitur dapat membuat fitur dengan korelasi individual yang lebih rendah menjadi penting.
3. **total sulfur dioxide (Skor ~0.103)**: Jumlah total SO<sub>2</sub> juga memiliki kontribusi yang signifikan.
4. **sulphates (Skor ~0.102)**: Sulfat, yang juga berkorelasi positif dengan kualitas, berada di peringkat keempat.
5. **fixed acidity (Skor ~0.073)**
6. **density (Skor ~0.072)**
7. **volatile acidity (Skor ~0.069)**
8. **pH (Skor ~0.063)**
9. **free sulfur dioxide (Skor ~0.052)**
10. **citric acid (Skor ~0.052)**
11. **residual sugar (Skor ~0.038)**: Sisa gula memiliki pengaruh paling kecil dalam model ini.

#### Kesimpulan Feature Importance:

- Seperti pada analisis korelasi, *alcohol* kembali menjadi fitur yang paling penting, diikuti oleh *chlorides* dan *total sulfur dioxide*.
- Kepentingan fitur ini memberikan wawasan tentang atribut mana yang paling menjadi fokus model dalam melakukan klasifikasi. Ini bisa menjadi dasar untuk eksplorasi lebih lanjut atau rekayasa fitur di masa depan.
- Penting untuk diingat bahwa kepentingan fitur di sini dihitung berdasarkan bagaimana model Random Forest *menggunakan* fitur tersebut untuk mengurangi ketidakmurnian (*impurity*) pada *node-node* pohonnya, dan mungkin berbeda dengan korelasi linear sederhana.

## ✓ 7. Visualisasi Tree

Untuk memahami bagaimana model Random Forest membuat keputusan, kita bisa memvisualisasikan beberapa *decision tree* (pohon keputusan) pertama yang membentuk *forest* tersebut. Random Forest adalah kumpulan dari banyak pohon keputusan, dan prediksi akhirnya adalah hasil mayoritas (untuk klasifikasi) atau rata-rata (untuk regresi) dari prediksi semua pohon.

Dengan melihat struktur pohon individual, kita bisa mendapatkan intuisi tentang aturan-aturan (*rules*) yang dipelajari model berdasarkan fitur-fitur yang ada.

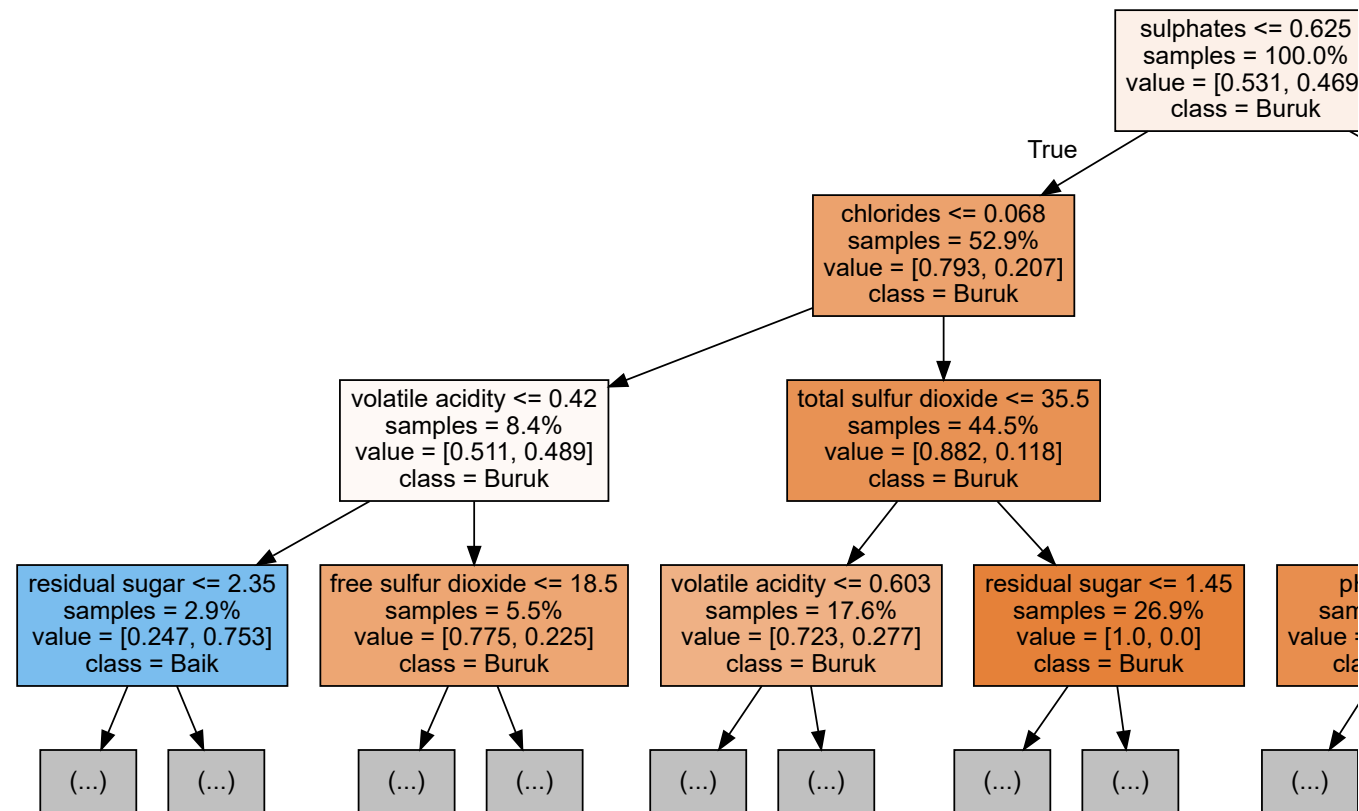
```

1 # Visualisasi 3 pohon keputusan pertama dari model Random Forest
2 for i in range(3):
3     tree = best_rf.estimators_[i]
4     dot_data = export_graphviz(tree,
5                                 feature_names=X_train.columns, # Nama-nama fitur
6                                 class_names=['Buruk', 'Baik'], # Nama kelas target
7                                 filled=True, # Memberi warna pada node berdasarkan kelas mayoritas
8                                 max_depth=3, # Membatasi kedalaman untuk visualisasi yang lebih sederhana
9                                 impurity=False, # Tidak menampilkan impurity
10                                proportion=True) # Menampilkan proporsi sampel, bukan jumlah absolut

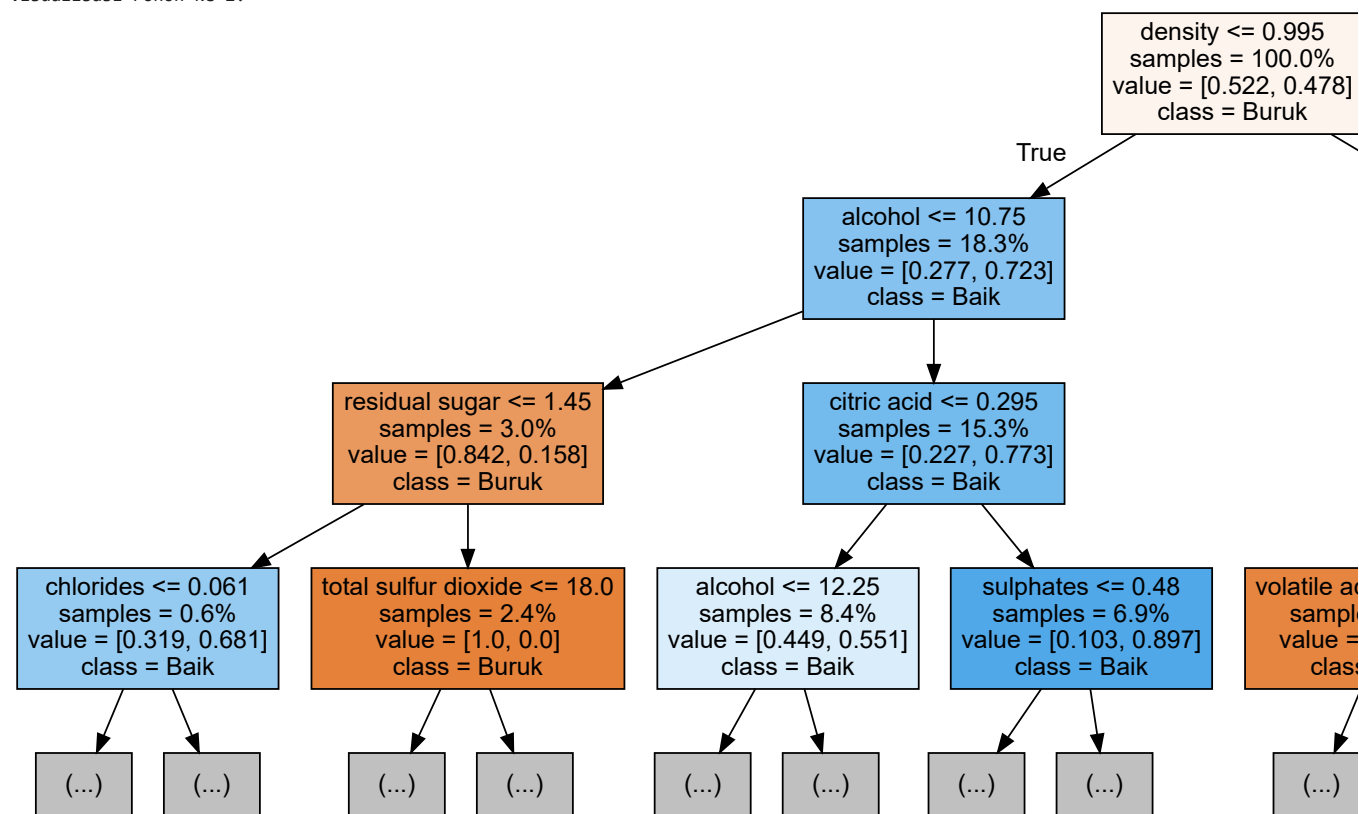
```

```
11
12     graph = graphviz.Source(dot_data)
13     print(f'Visualisasi Pohon ke-{{i+1}}:')
14     display(graph)
```

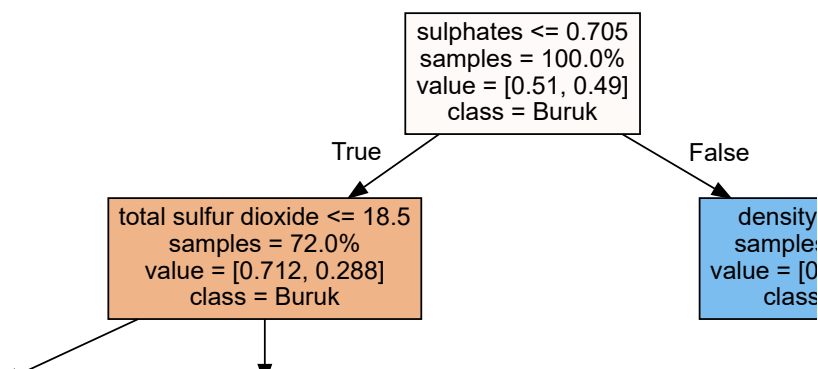
Visualisasi Pohon ke-1:

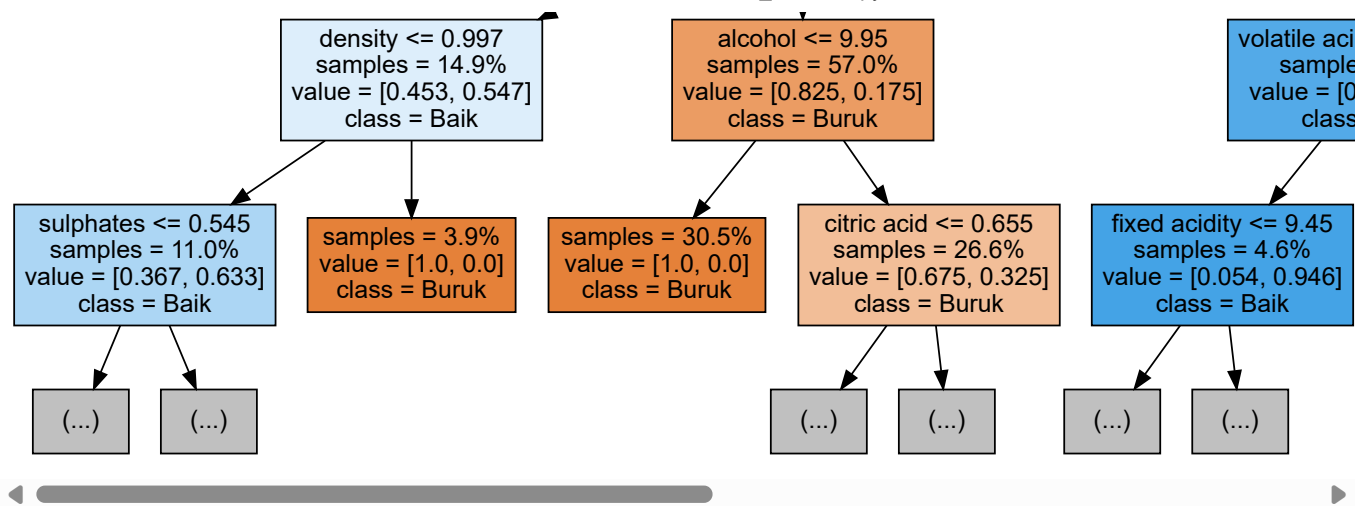


Visualisasi Pohon ke-2:



Visualisasi Pohon ke-3:





Dengan melihat beberapa pohon, kita bisa melihat bahwa fitur yang berbeda mungkin menjadi penting di pohon yang berbeda, atau pada kedalaman yang berbeda. Ini adalah salah satu kekuatan Random Forest – ia menggabungkan "pandangan" dari banyak pohon yang beragam.

```

1 # Hasil dari X_test dataframe yang digabung dengan prediksi dan target column yg asli(quality)
2 final_df = pd.concat([X_test.reset_index(drop=True),
3                       y_test.reset_index(drop=True)], axis=1)
4 final_df['quality_prediction'] = pd.Series(y_pred)
5
6 # Print dataframe X_test dengan y_pred dan quality
7 print(f"Hasil prediksi dari X_test :")
8 display(final_df)

```

➡ Hasil prediksi dari X\_test :

|     | fixed<br>acidity | volatile<br>acidity | citric<br>acid | residual<br>sugar | chlorides | free<br>sulfur<br>dioxide | total<br>sulfur<br>dioxide | density | pH   | sulphates | alcohol | quality | quality_prediction |
|-----|------------------|---------------------|----------------|-------------------|-----------|---------------------------|----------------------------|---------|------|-----------|---------|---------|--------------------|
| 0   | 8.0              | 1.18                | 0.21           | 1.9               | 0.083     | 14.0                      | 41.0                       | 0.99532 | 3.34 | 0.47      | 10.5    | 0       | 0                  |
| 1   | 7.5              | 0.55                | 0.24           | 2.0               | 0.078     | 10.0                      | 28.0                       | 0.99830 | 3.45 | 0.78      | 9.5     | 0       | 0                  |
| 2   | 7.2              | 0.34                | 0.24           | 2.0               | 0.071     | 30.0                      | 52.0                       | 0.99576 | 3.44 | 0.58      | 10.1    | 0       | 0                  |
| 3   | 9.0              | 0.60                | 0.29           | 2.0               | 0.069     | 32.0                      | 73.0                       | 0.99654 | 3.34 | 0.57      | 10.0    | 0       | 0                  |
| 4   | 7.1              | 0.66                | 0.00           | 3.9               | 0.086     | 17.0                      | 45.0                       | 0.99760 | 3.46 | 0.54      | 9.5     | 0       | 0                  |
| ... | ...              | ...                 | ...            | ...               | ...       | ...                       | ...                        | ...     | ...  | ...       | ...     | ...     | ...                |
| 267 | 7.4              | 0.64                | 0.07           | 1.8               | 0.100     | 8.0                       | 23.0                       | 0.99610 | 3.30 | 0.58      | 9.6     | 0       | 0                  |
| 268 | 10.4             | 0.28                | 0.54           | 2.7               | 0.105     | 5.0                       | 19.0                       | 0.99880 | 3.25 | 0.63      | 9.5     | 0       | 0                  |
| 269 | 5.4              | 0.42                | 0.27           | 2.0               | 0.092     | 23.0                      | 55.0                       | 0.99471 | 3.78 | 0.64      | 12.3    | 1       | 0                  |
| 270 | 9.6              | 0.60                | 0.50           | 2.3               | 0.079     | 28.0                      | 71.0                       | 0.99970 | 3.50 | 0.57      | 9.7     | 0       | 0                  |

Next steps:

[Generate code with final\\_df](#)

[View recommended plots](#)

[New interactive sheet](#)

### Interpretasi Tabel Hasil Prediksi:

Tabel `final_df` di atas menampilkan:

- Kolom-kolom fitur dari `x_test`.
- Kolom `quality`: Ini adalah label kualitas **sebenarnya** dari anggur tersebut (0 untuk "Buruk", 1 untuk "Baik").
- Kolom `quality_prediction`: Ini adalah label kualitas yang **diprediksi** oleh model Random Forest kita.

Dengan membandingkan kolom `quality` dan `quality_prediction`, kita bisa melihat secara langsung sampel mana yang berhasil diprediksi dengan benar dan mana yang salah.

### Contoh:

- Jika pada suatu baris `quality` adalah 1 dan `quality_prediction` juga 1, berarti model berhasil memprediksi anggur "Baik" dengan benar (True Positive).
- Jika `quality` adalah 0 dan `quality_prediction` adalah 1, berarti model salah memprediksi anggur "Buruk" sebagai "Baik" (False Positive).