

**LAPORAN TUGAS BESAR KECERDASAN BUATAN
IMPLEMENTASI *MACHINE LEARNING* PADA PENYAKIT JANTUNG
DENGAN ALGORITMA K-NEAREST NEIGHBORS (KNN)**

Dibuat untuk memenuhi tugas akhir yang diampu oleh:

Leni Fitriani, S.Kom., M.Kom.



Disusun Oleh:

Aghniya Afiatul Jannah : 2396035

Alya Rahmawati : 2306063

**JURUSAN ILMU KOMPUTER
PROGRAM STUDI TEKNIK INFORMATIKA
INSTITUT TEKNOLOGI GARUT**

2025

1. BUSINESS UNDERSTANDING

1.1 Latar Belakang

Jantung merupakan organ terpenting bagi manusia, salah satu fungsi jantung adalah mengalirkan darah yang membawa oksigen keseluruh tubuh. Jika jantung mengalami kerusakan maka organ tersebut akan menyebabkan berbagai macam penyakit. (Dewi, 2023) Penyakit Jantung atau disebut juga penyakit kardiovaskular merupakan salah satu penyakit berbahaya yang dapat menyebabkan kematian. Seiring berkembangnya teknologi dan peningkatan popularitas teknologi machine learning, teknologi machine learning tersebut dapat digunakan untuk membantu mendeteksi penyakit jantung dengan menggunakan data pasien. (Junifer Pangaribuan et al., 2021) Menurut who, sekitar 17,9 juta orang meninggal setiap tahun karena penyakit ini. Di indonesia sendiri, angka kematian akibat pkv mencapai 651.481 orang per tahun. Diagnosis dini menjadi krusial, tetapi proses identifikasinya sering kompleks karena melibatkan banyak faktor baik yang bisa diubah (gaya hidup, tekanan darah, kolesterol) maupun tidak (usia, genetika)..

Machine learning merupakan sistem yang mampu belajar sendiri untuk memutuskan sesuatu tanpa harus berulang kali diprogram oleh manusia sehingga komputer menjadi semakin cerdas belajar dari pengalaman data yang dimiliki. (Retnoningsih & Pramudita, 2020)

Dengan adanya perkembangan teknologi, banyak hal yang dapat dilakukan untuk memberikan kemudahan kepada manusia, diantaranya perkembangan bidang ilmu Artificial Intelligence (AI). AI dapat dimanfaatkan pada segala bidang seperti computer vision yang dapat mendeteksi penyakit dan sistem autopilot pada transportasi. (Silmi Ath Thahirah Al Azhima, D. Darmawan, N. Fahmi Arief Hakim, I. Kustiawan, M. Al Qibtiya, 2022)

1.2 Permasalahan Dunia Nyata

a. Tingginya angka kematian akibat penyakit jantung

Penyakit jantung menjadi penyebab utama kematian di dunia dan Indonesia, namun masih sulit terdeteksi secara dini.

b. Keterbatasan diagnosis manual oleh tenaga medis

Proses diagnosis sering bergantung pada observasi dokter yang terbatas, terutama di daerah dengan minim tenaga ahli.

c. Kompleksitas faktor risiko penyakit jantung

Faktor risiko seperti usia, tekanan darah, kolesterol, dan gaya hidup membuat diagnosis menjadi semakin kompleks.

- d. Belum optimalnya pemanfaatan teknologi klasifikasi otomatis

Padahal data medis tersedia dan algoritma seperti KNN terbukti akurat dalam mengklasifikasikan risiko penyakit jantung.

1.3 Tujuan Proyek

Proyek ini bertujuan untuk:

- a. Membangun model prediktif menggunakan algoritma *K-Nearest Neighbor (KNN)* untuk mengklasifikasikan apakah seseorang berisiko terkena penyakit jantung atau tidak.
- b. Menganalisis performa model KNN dengan menggunakan metrik evaluasi seperti akurasi, precision, recall, dan F1-score.
- c. Menguji efektivitas parameter K, yaitu berapa jumlah tetangga terdekat yang optimal dalam memberikan klasifikasi terbaik terhadap data pasien.
- d. Mengaplikasikan machine learning pada data medis nyata, khususnya dataset penyakit jantung dari Kaggle, sebagai bentuk kontribusi teknologi terhadap dunia kesehatan.

1.4 User atau Pengguna Sistem

Pengguna sistem tidak hanya terbatas pada tenaga medis, namun juga meliputi masyarakat umum, peneliti, dan institusi kesehatan yang berkepentingan dalam upaya deteksi dini dan pencegahan penyakit jantung. Adapun user atau pengguna system diantaranya adalah sebagai berikut:

- a. Dokter atau Tenaga Medis

Menggunakan sistem untuk membantu diagnosa awal penyakit jantung secara cepat dan berbasis data.

- b. Pasien atau Masyarakat Umum

Dapat memanfaatkan sistem sebagai alat deteksi dini untuk mengetahui potensi risiko penyakit jantung sebelum berkonsultasi ke rumah sakit.

- c. Peneliti dan Akademisi

Sebagai media pembelajaran dan pengembangan model kecerdasan buatan di bidang kesehatan.

- d. Instansi Kesehatan dan Rumah Sakit

Dapat mengintegrasikan sistem ke dalam proses layanan untuk meningkatkan efisiensi dan akurasi diagnosis pasien.

1.5 Manfaat Implementasi AI

Implementasi algoritma KNN diharapkan dapat memberikan manfaat dalam membantu pengambilan keputusan medis dan meningkatkan kualitas layanan kesehatan secara menyeluruh. Adapun manfaatnya adalah sebagai berikut:

a. Meningkatkan kecepatan dan efisiensi diagnosis

Sistem dapat memberikan hasil prediksi secara instan berdasarkan data pasien.

b. Membantu deteksi dini penyakit jantung

Deteksi lebih awal memungkinkan penanganan lebih cepat, sehingga risiko komplikasi atau kematian dapat dikurangi.

c. Mengurangi ketergantungan pada analisis manual

Memberikan second opinion berbasis data bagi dokter, terutama di daerah dengan keterbatasan tenaga ahli.

d. Akurat dan berbasis data riil

Dengan akurasi model mencapai hingga $>90\%$, sistem ini dapat memberikan hasil prediksi yang cukup andal untuk digunakan secara klinis.

e. Dapat diintegrasikan ke dalam sistem layanan digital

Seperti website, aplikasi kesehatan, atau dashboard klinik untuk meningkatkan layanan berbasis teknologi.

2. DATA UNDERSTANDING

2.1 Sumber Data

Dataset yang digunakan dalam artikel ini berasal dari Kaggle, dengan judul: "Heart Disease Dataset" Dataset ini dibuat oleh David Lapp Tautan dataset: <https://www.kaggle.com/datasets/johnsmith88/heart-disease-dataset?>

2.2 Deskripsi Setiap Fitur (Atribut)

Dataset yang digunakan memiliki 14 atribut:

| Nama Atribut | Deskripsi | Tipe Data |
|--------------|--|-----------------|
| Age | Usia Pasien | Numerik (Int) |
| Sex | Jenis Kelamin (1 = Laki-Laki, 0 = Perempuan) | Kategorik (Int) |
| Cp | Tipe Nyeri Dada (0–3) | Kategorik (Int) |
| Trestbps | Tekanan Darah Saat Istirahat | Numerik (Int) |
| Chol | Kolesterol Serum (Mg/Dl) | Numerik (Int) |
| Fbs | Gula Darah Puasa >120 Mg/Dl (1 = Ya, 0 = Tidak) | Kategorik (Int) |
| Restecg | Hasil Ekg Saat Istirahat (0–2) | Kategorik (Int) |
| Thalach | Detak Jantung Maksimal | Numerik (Int) |
| Exang | Angina Akibat Olahraga (1 = Ya, 0 = Tidak) | Kategorik (Int) |
| Oldpeak | Depresi St Saat Latihan (Relatif Terhadap Istirahat) | Numerik (Float) |
| Slope | Kemiringan Segmen St Saat Latihan (0–2) | Kategorik (Int) |
| Ca | Jumlah Pembuluh Utama Yang Diwarnai (0–4) | Kategorik (Int) |
| Thal | Status Thalassemia (2 = Fixed, 3 = Normal, 1 = Rev) | Kategorik (Int) |
| Target | Label: 1 = Sakit Jantung, 0 = Normal | Kategorik (Int) |

2.3 Ukuran dan Format Data

Dataset yang digunakan memiliki ukuran 1.025 baris dan 14 kolom, yang berarti terdapat 1.025 data pasien dan 14 atribut pada masing-masing baris. Dari 14 atribut tersebut, 13 merupakan fitur (variabel input) dan 1 atribut merupakan label target (output) yang menunjukkan kondisi penyakit jantung pasien.

Format data yang digunakan adalah tabel berstruktur (structured tabular data), di mana setiap baris berisi data lengkap dari seorang pasien, dan setiap kolom merepresentasikan satu jenis informasi. Format seperti ini sangat mendukung untuk proses analisis statistik, visualisasi data, maupun penerapan model machine learning.

2.4 Tipe Data Dan Target Klasifikasi

Seluruh atribut dalam dataset ini memiliki tipe data numerik, terdiri dari:

- a. Atribut seperti sex, cp, fbs, restecg, exang, slope, ca, dan thal merupakan data kategorikal yang telah dikodekan dalam bentuk angka bulat (int64). Meskipun bertipe numerik secara teknis, atribut-atribut ini merepresentasikan kategori atau kelas tertentu.
- b. Atribut seperti age, trestbps, chol, thalach, dan oldpeak termasuk data numerik kontinu, karena nilainya berupa rentang angka yang dapat berubah secara bertahap (termasuk desimal pada oldpeak).

Berikut adalah ringkasan tiap atribut berdasarkan jenis nilainya:

| No. | Atribut | Deskripsi Singkat | Jenis Nilai |
|-----|----------|--|--------------------------|
| 1 | age | Usia pasien | Kontinu |
| 2 | sex | Jenis kelamin pasien | Diskrit (kategori biner) |
| 3 | cp | Tipe nyeri dada | Diskrit (kategori) |
| 4 | trestbps | Tekanan darah saat istirahat | Kontinu |
| 5 | chol | Kadar kolesterol dalam darah | Kontinu |
| 6 | fbs | Gula darah puasa > 120 mg/dl | Diskrit (biner) |
| 7 | restecg | Hasil elektrokardiogram saat istirahat | Diskrit (kategori) |

| | | | |
|----|---------|---|------------------------------------|
| 8 | thalach | Detak jantung maksimum saat latihan | Kontinu |
| 9 | exang | Angina akibat latihan fisik | Diskrit (biner) |
| 10 | oldpeak | Depresi ST akibat latihan | Kontinu |
| 11 | slope | Kemiringan segmen ST saat latihan | Diskrit (kategori) |
| 12 | ca | Jumlah pembuluh darah besar yang terdeteksi | Diskrit (kategori) |
| 13 | thal | Jenis thalassemia | Diskrit (kategori) |
| 14 | target | Status penyakit jantung (label klasifikasi) | Diskrit (biner: 0 = tidak, 1 = ya) |

Atribut target merupakan label atau variabel dependen dalam dataset ini, yang menunjukkan status penyakit jantung pada pasien. Nilai dari target ini bersifat biner (binary classification), dengan rincian:

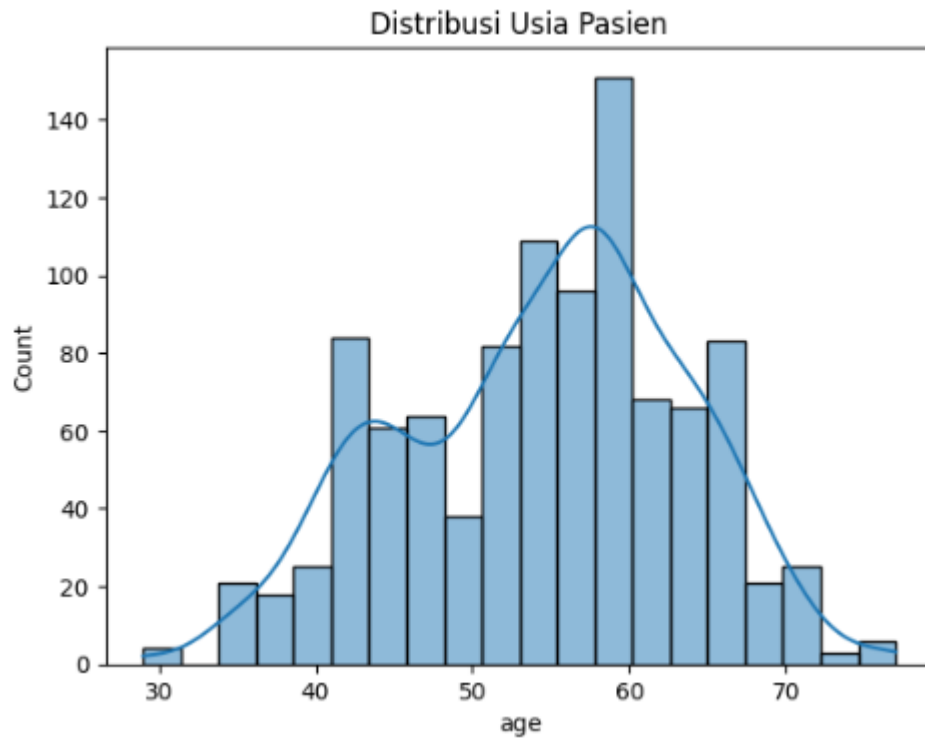
- 0 = Tidak terdiagnosis penyakit jantung
- 1 = Terdiagnosis penyakit jantung

Dengan demikian, tujuan dari analisis atau pemodelan menggunakan dataset ini adalah klasifikasi biner, yaitu memprediksi apakah seorang pasien memiliki penyakit jantung atau tidak berdasarkan fitur-fitur yang tersedia.

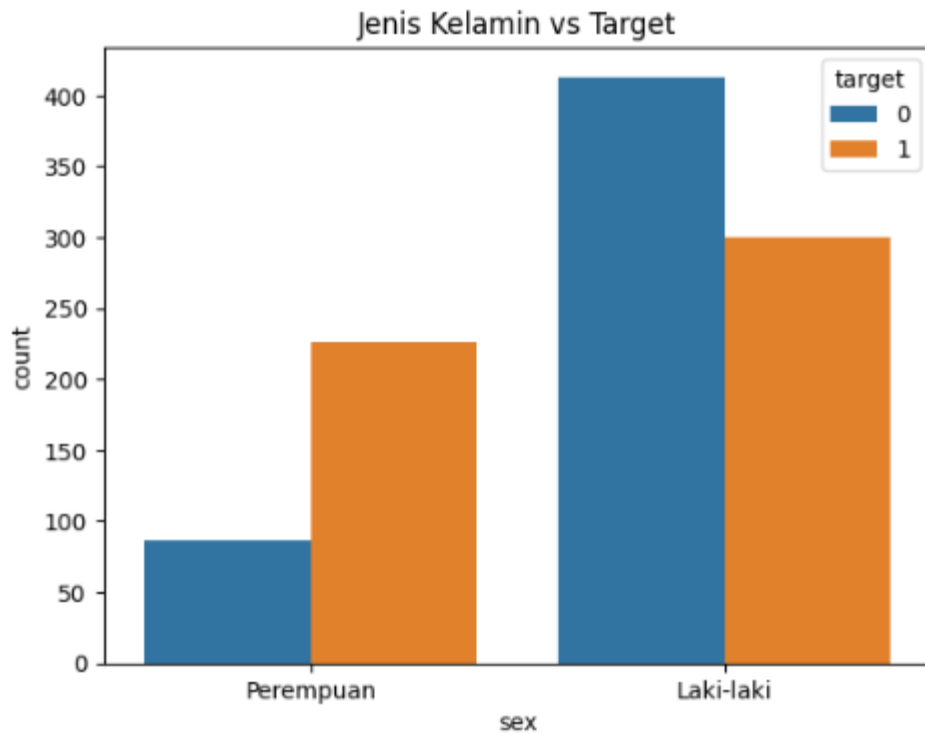
3. EXPLORATORY DATA ANALYSIS (EDA)

3.1 Visualisasi Distribusi Data

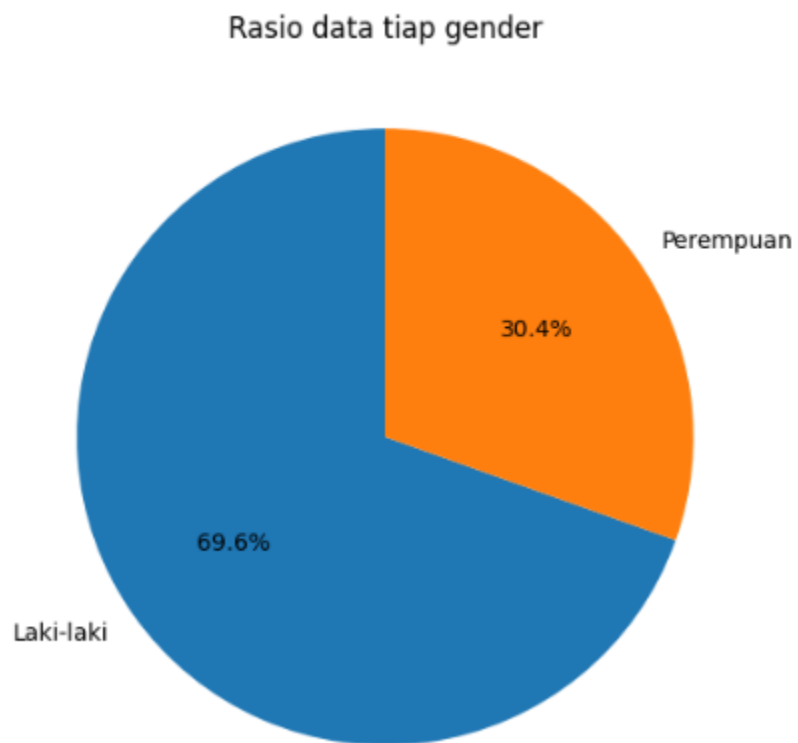
Pada tahap visualisasi, digunakan beberapa grafik untuk memahami pola dalam data. Histogram menunjukkan bahwa sebagian besar pasien berusia antara 50–60 tahun. Countplot dan pie chart memperlihatkan bahwa mayoritas pasien adalah laki-laki, dan penyakit jantung lebih banyak terjadi pada kelompok tersebut.



Gambar 3. 1 Visualisasi Histogram Distribusi Usia Pasien



Gambar 3. 2 Visualisasi Countplot Jenis Kelamin vs Target

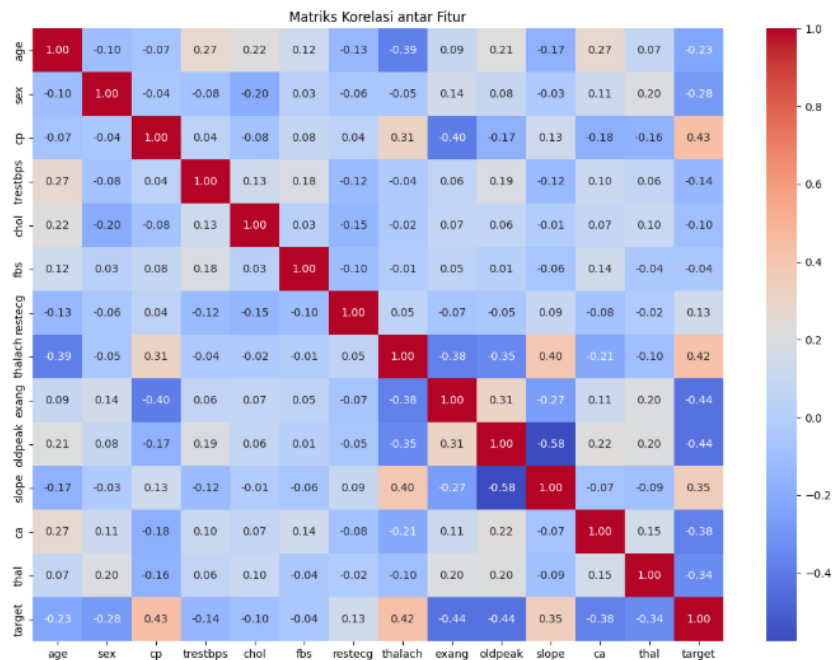


Gambar 3. 3 Visualisasi Pie Chart Rasio Data Tiap Gender

3.2 Analisis Korelasi Antar Fitur

Tujuan utama dari analisis ini biasanya untuk melihat fitur mana yang paling berpengaruh terhadap variabel target (apakah seseorang menderita penyakit jantung atau tidak). Kita bisa melihatnya di baris atau kolom terakhir (target).

- Fitur dengan Korelasi Positif Terkuat terhadap target:
 1. Cp (jenis nyeri dada): 0.43. Semakin tinggi tipe cp, semakin tinggi kemungkinan memiliki penyakit jantung.
 2. Thalach (detak jantung maks): 0.42. Semakin tinggi detak jantung maksimum, semakin tinggi kemungkinan memiliki penyakit jantung.
 3. Slope (kemiringan segmen st): 0.35.
- Fitur dengan Korelasi Negatif Terkuat terhadap target:
 1. Exang (nyeri dada akibat olahraga): -0.44. Ini adalah korelasi terkuat. Artinya, jika exang bernilai 1 (ya), kemungkinan target bernilai 0 (tidak sakit jantung) lebih tinggi.
 2. Oldpeak: -0.43. Semakin tinggi nilai oldpeak, semakin rendah kemungkinan memiliki penyakit jantung.
 3. Ca: -0.39. Semakin banyak pembuluh darah yang terdeteksi, semakin rendah kemungkinan memiliki penyakit jantung.

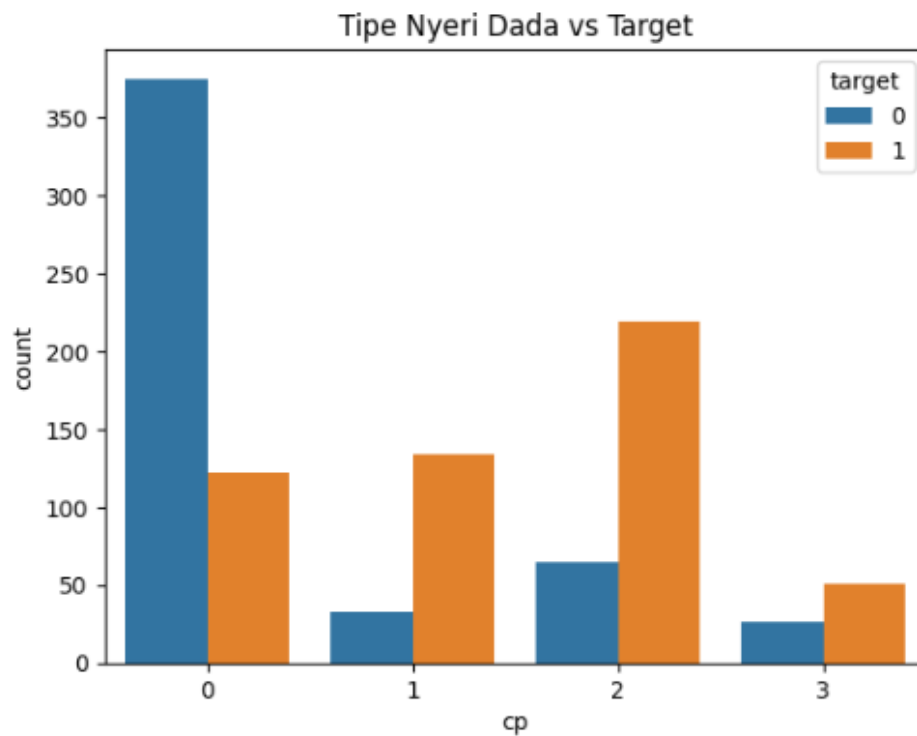


Gambar 3. 4 Analisis Matriks Korelasi antar Fitur

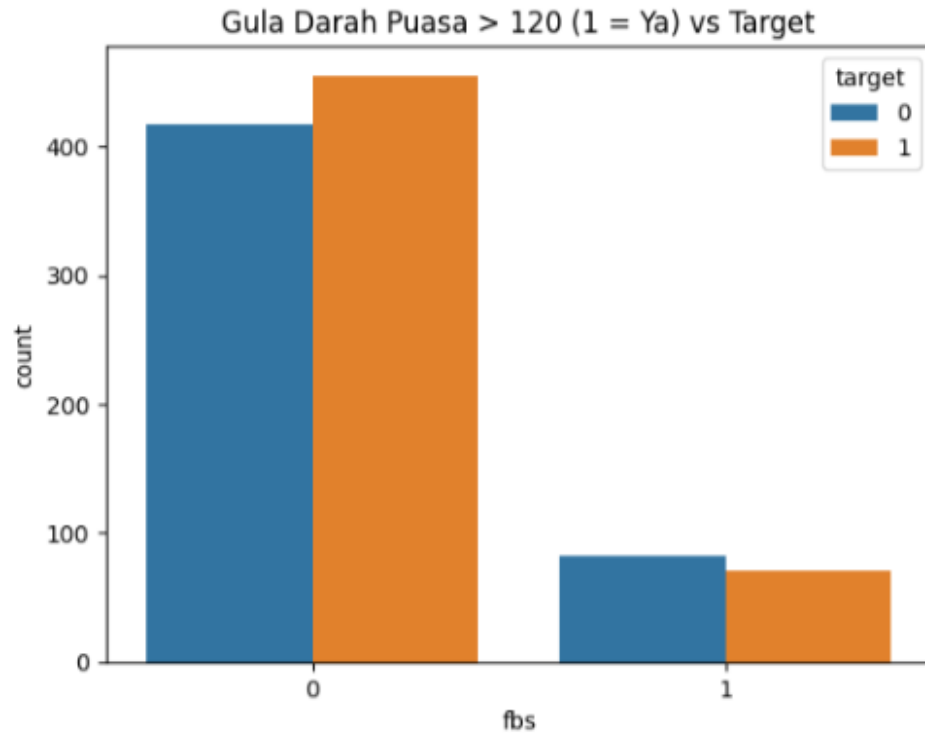
- Perbandingan Fitur Kategori Terhadap Target:

Analisis bivariat untuk melihat hubungan antara beberapa fitur kategorikal (cp, fbs, exang) dengan variabel target. Fungsi sns.countplot dari library Seaborn membuat diagram batang yang menghitung jumlah kemunculan setiap kategori.

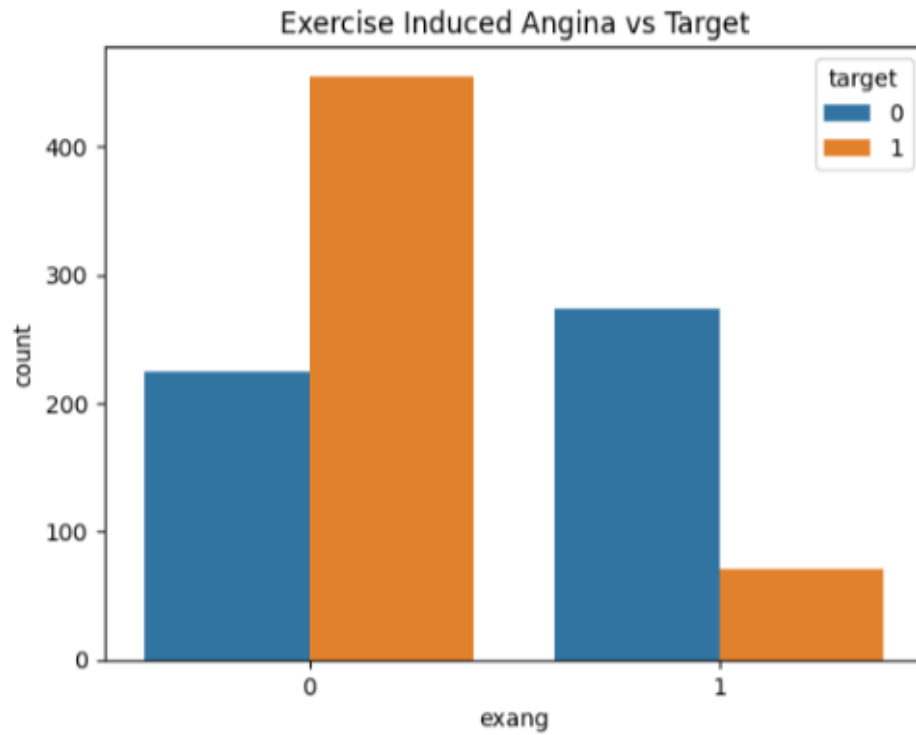
Penggunaan hue='target' adalah kunci utamanya. Parameter ini memisahkan setiap batang berdasarkan nilai target (0 untuk tidak sakit jantung, 1 untuk sakit jantung), sehingga kita bisa langsung membandingkan distribusi penderita dan non-penderita penyakit jantung di setiap kategori fitur.



Gambar 3. 5 Perbandingan Tipe Nyeri Dada vs Target



Gambar 3. 6 Perbandingan Gula Darah Puasa vs Target



Gambar 3. 7 Perbandingan Exercise Induced Angina vs Target

3.3 Deteksi Data Tidak Seimbang

Distribusi target dalam dataset menunjukkan proporsi yang hampir seimbang, yaitu 51,32% untuk kelas 1 dan 48,68% untuk kelas 0. Dengan perbedaan yang sangat kecil, data ini tidak memerlukan penyesuaian kelas seperti oversampling atau undersampling. Kondisi ini ideal untuk pelatihan model KNN karena mengurangi risiko bias dan menghasilkan evaluasi model yang lebih akurat dan representatif.

| proportion | |
|------------|----------|
| target | |
| 1 | 0.513171 |
| 0 | 0.486829 |

dtype: float64

Gambar 3. 8 Proportin Deteksi Data Tidak Seimbang

Boxplot Fitur Numerik Berdasarkan Target:

a. Kolestrol (chol) vs Target

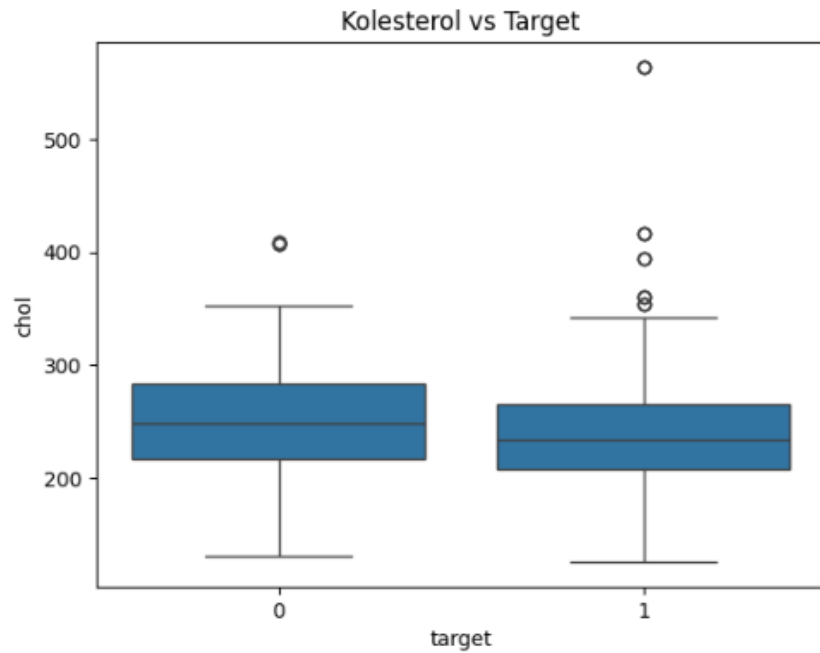
Plot ini akan menghasilkan dua boxplot yang bersebelahan: Satu boxplot menunjukkan distribusi kadar kolesterol untuk kelompok target=0 (Tidak Sakit Jantung). Satu boxplot lagi menunjukkan distribusi kadar kolesterol untuk kelompok target=1 (Sakit Jantung).

b. Detak jantung maksimal (thalac) vs Target

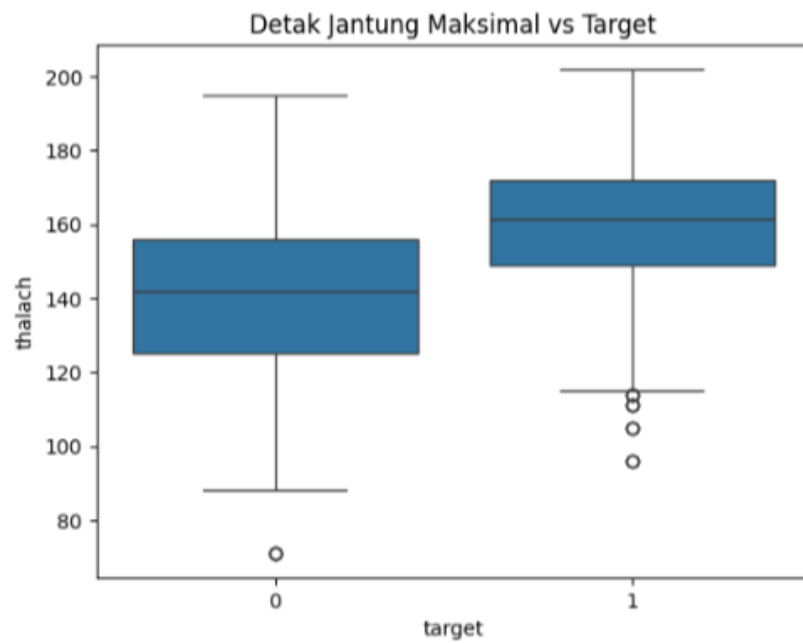
Sama seperti sebelumnya, plot ini akan menampilkan dua boxplot untuk membandingkan distribusi detak jantung maksimal (thalach) antara kelompok target=0 dan target=1.

c. Tekanan darah saat istirahat (trestbps) vs Target

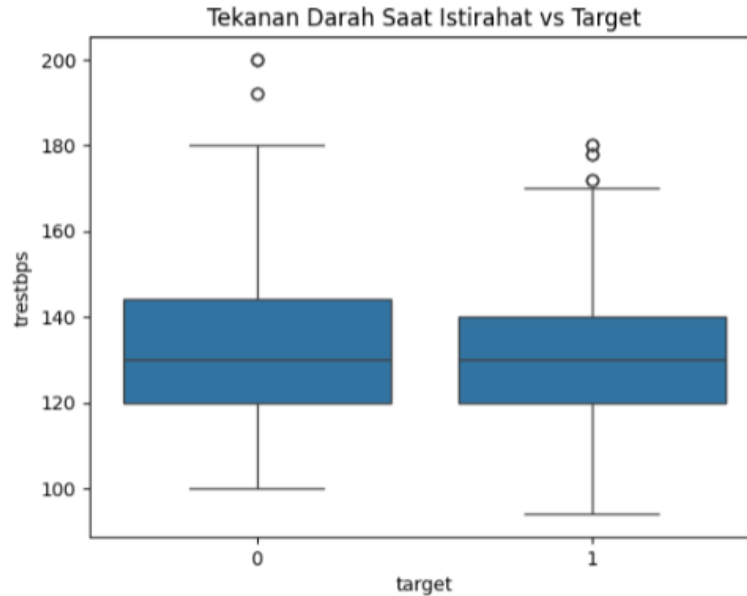
Plot ini membandingkan distribusi tekanan darah saat istirahat (trestbps) antara kelompok target=0 dan target=1.



Gambar 3. 9 Boxplot Kolestrol vs Target



Gambar 3. 10 Boxplot Detak jantung maksimal vs Target



Gambar 3. 11 Boxplot Tekanan darah saat istirahat vs Target

3.4 Insight awal dari pola data

- Mayoritas pasien dengan penyakit jantung adalah pria berusia sekitar 50 sampai 60 tahun. Berdasarkan distribusi data, ditemukan bahwa sebagian besar pasien yang terdiagnosis penyakit jantung berjenis kelamin laki-laki, dengan rentang usia dominan antara 50 hingga 60 tahun. Hal ini sejalan dengan temuan medis bahwa risiko penyakit jantung meningkat seiring bertambahnya usia, terutama pada pria. Faktor hormonal, pola hidup, serta kecenderungan risiko seperti merokok dan tekanan darah tinggi lebih banyak ditemukan pada kelompok usia dan jenis kelamin ini.
- Pasien dengan nilai CP (chest pain) dan thalach (max heart rate achieved) tinggi cenderung memiliki penyakit jantung. Nilai CP yang tinggi menunjukkan adanya keluhan nyeri dada yang lebih serius atau tidak normal, seperti *asymptomatic chest pain*, yang merupakan indikasi umum penyakit jantung. Selain itu, thalach atau detak jantung maksimal yang tercapai juga memberikan sinyal penting. Nilai thalach yang tinggi pada beberapa pasien bisa menandakan respon tubuh terhadap stres atau aktivitas yang tidak normal, yang sering ditemukan pada penderita penyakit jantung.
- Fitur chol, thalach, dan trestbps mengandung outlier. Dari hasil eksplorasi data, ditemukan bahwa fitur kadar kolesterol (chol), detak jantung maksimal (thalach), dan tekanan darah saat istirahat (trestbps) memiliki beberapa nilai

ekstrem atau outlier. Kehadiran outlier ini bisa mempengaruhi performa model klasifikasi karena dapat mengganggu proses perhitungan jarak (misalnya Euclidean) yang digunakan pada algoritma KNN. Oleh karena itu, penting untuk melakukan analisis lebih lanjut terhadap data outlier ini apakah merupakan nilai valid atau error, dan apakah perlu dilakukan transformasi, normalisasi, atau penghapusan data.

4. DATA PREPARATION

4.1 Pembersihan Data (Null Value, Duplikasi)

a. Cek Duplikasi

Ditemukan sebanyak 723 baris data duplikat dalam dataset, yang merupakan jumlah sangat besar. Padahal, dataset asli *heart.csv* dari UCI Machine Learning Repository hanya memiliki 303 baris data unik. Adapun versi dataset dari Kaggle memiliki 1025 baris, yang kemungkinan besar merupakan hasil penggabungan beberapa sumber data sehingga mengandung banyak duplikasi.

Kehadiran data duplikat dalam jumlah besar dapat mempengaruhi akurasi dan keandalan model machine learning, karena akan memperkuat pola tertentu secara tidak proporsional. Oleh karena itu, penting untuk melakukan penghapusan duplikat sebelum proses pelatihan model dilakukan.



```
Data duplikat: 723
```

Gambar 4. 1 Data Duplikat

b. Menghapus Duplikat

Menghapus sekitar 723 duplikat yang terdapat pada dataset. Berdampak pada jumlah data berkurang drastis. Ini adalah efek yang paling jelas. Jumlah baris data Anda akan turun dari 1025 menjadi 302 ($1025 - 723 = 302$). Angka 302 ini sesuai dengan jumlah data asli dari dataset *heart.csv* standar yang berasal dari UCI (University of California, Irvine), yang berarti dataset Anda sekarang bersih dari data tiruan.



```
Data duplikat: 0
```

Gambar 4. 2 Data Duplikat Setelah Melakukan Penghapusan

c. Pengecekan Missing Value

Dataset clear tidak memiliki missing value.

```

Missing value:
  age      0
  sex      0
  cp       0
  trestbps 0
  chol     0
  fbs      0
  restecg  0
  thalach  0
  exang    0
  oldpeak  0
  slope    0
  ca       0
  thal     0
  target   0
dtype: int64

```

Gambar 4. 3 Missing Value

d. Deteksi Outlier pada Kolom chol (Kolesterol)

Metode ini mengimplementasikan metode statistik umum untuk mendeteksi outlier, yaitu Metode IQR (Interquartile Range).

- $Q1 = df['chol'].quantile(0.25)$: Menghitung Kuartil Pertama (persentil ke-25) dari data kolesterol.
- $Q3 = df['chol'].quantile(0.75)$: Menghitung Kuartil Ketiga (persentil ke-75).
- $IQR = Q3 - Q1$: Menghitung Jangkauan Interkuartil, yaitu rentang antara Q3 dan Q1 yang mencakup 50% data di tengah.
- `Lower_bound` dan `upper_bound`: Menghitung batas bawah dan batas atas. Secara statistik, setiap titik data yang berada di luar rentang ini dianggap sebagai outlier.
- `Outliers = df[...]`: Baris ini memfilter dataframe `df` dan menyimpan semua baris di mana nilai `chol` lebih kecil dari `lower_bound` ATAU (`||`) lebih besar dari `upper_bound`.
- `Len(outliers)`: Menghitung jumlah baris dalam dataframe `outliers` yang baru dibuat, yang merupakan jumlah total outlier.

Output menghasilkan: berdasarkan metode IQR, terdapat 16 data pada kolom 'chol' yang nilainya dianggap ekstrem atau pencilan (outlier) dibandingkan dengan sebaran data lainnya.

Jumlah outlier pada 'chol': 16

Gambar 4. 4 Deteksi Outlier pada Kolom chol (Kolestrol)

e. Capping / Winsorizing

Outlier pada fitur chol, thalach, dan trestbps ditangani menggunakan metode IQR dengan mengganti nilai ekstrem yang melebihi batas bawah atau atas dengan nilai batas tersebut. Pendekatan ini menjaga kestabilan data dan mengurangi pengaruh negatif outlier terhadap model KNN yang sensitif terhadap jarak antar data.

4.2 Encoding data kategorik (label encoding, one-hot)

Dataset yang digunakan tidak memerlukan Encoding karena semua fitur pada dataset yang digunakan adalah numerik (baik hasil pengukuran maupun kategori numerik). Karena itu, tidak perlu One-Hot Encoding atau Label Encoding. Jika dataset mengandung kolom bertipe string/kategori (seperti 'thal' atau 'cp' dalam bentuk teks), maka encoding perlu dilakukan. Tapi dalam dataset ini, semua sudah berupa angka (0, 1, 2, 3).

4.3 Normalisasi/Standarisasi data numerik dan Split data (train-test)

a. Pemisahan Fitur dan Target

- `X = df.drop('target', axis=1)`: X menjadi DataFrame yang berisi semua fitur (variabel independen).
- `y = df['target']`: y adalah Series yang berisi variabel target (variabel dependen) yang ingin diprediksi.

b. Normalisasi

`scaler = StandardScaler()` dan `X_scaled = scaler.fit_transform(X)`: Semua fitur dalam X diubah skalanya menggunakan standarisasi Z-Score. Ini membuat semua fitur memiliki skala yang sebanding, yang sangat penting untuk algoritma berbasis jarak seperti KNN.

c. Reduksi Dimensi (PCA)

`pca = PCA(n_components=8)` dan `X_pca = pca.fit_transform(X_scaled)`: Principal Component Analysis digunakan untuk mereduksi 13 fitur asli menjadi 8 "komponen utama". Komponen ini adalah kombinasi linear dari fitur asli yang menangkap sebagian besar variasi dalam data. Tujuannya adalah untuk menyederhanakan model tanpa kehilangan banyak informasi penting.

d. Pembagian Data Latih dan Uji:

`X_train, X_test, y_train, y_test = train_test_split(...)`: Data yang telah diproses (`X_pca` dan `y`) dibagi menjadi dua set:

- Training set (`X_train, y_train`): 80% dari data, digunakan untuk melatih model KNN.
- Testing set (`X_test, y_test`): 20% dari data, digunakan untuk menguji seberapa baik model yang telah dilatih dapat memprediksi data baru. `random_state=42` memastikan pembagian data selalu sama setiap kali kode dijalankan.

5. MODELING

5.1 Pemilihan Algoritma

Dalam proyek ini, digunakan pendekatan klasifikasi berbasis machine learning untuk memprediksi risiko penyakit jantung. Model yang dibangun bertujuan mengkategorikan pasien ke dalam dua kelas, yaitu memiliki atau tidak memiliki indikasi penyakit jantung, berdasarkan sejumlah fitur medis seperti tekanan darah, kadar kolesterol, dan detak jantung. Untuk mendukung proses klasifikasi tersebut, dipilih algoritma K-Nearest Neighbor (KNN).

K-Nearest Neighbor (KNN) adalah salah satu metode yang digunakan untuk klasifikasi objek baru berdasarkan sejumlah K tetangga terdekat. Algoritma KNN relatif sederhana dan mudah dipahami, sehingga cukup umum digunakan. Pada penerapan algoritma ini, pengklasifikasian terhadap sebuah gambar berdasarkan jarak terdekat dengan tetangganya. Nilai jarak ini akan digunakan sebagai nilai kemiripan antara data uji dan data latih. (Akbarollah et al., 2023)

5.2 Alasan Pemilihan Model

KNN dipilih karena merupakan algoritma yang sederhana namun efektif, khususnya untuk data dengan jumlah fitur yang telah dinormalisasi. Algoritma ini bekerja dengan cara mencari sejumlah data tetangga terdekat dari data baru, kemudian menentukan kelasnya berdasarkan mayoritas kelas dari tetangga tersebut. Selain tidak memerlukan proses pelatihan yang kompleks, KNN juga cocok digunakan pada data non-linier dan mampu memberikan hasil prediksi yang akurat setelah dilakukan praproses seperti normalisasi dan reduksi dimensi. Oleh karena itu, algoritma ini dianggap sesuai untuk digunakan dalam klasifikasi penyakit jantung pada proyek ini.

5.3 Implementasi Model

Model KNN diimplementasikan melalui beberapa tahapan utama berikut:

- Semua fitur dinormalisasi menggunakan *StandardScaler* agar skala antar fitur setara dan tidak memengaruhi perhitungan jarak.
- Reduksi Dimensi (PCA), data yang telah dinormalisasi direduksi dari 13 fitur menjadi 8 komponen utama untuk menyederhanakan model dan mempercepat proses komputasi.
- Pembagian data latih dan uji dengan membagi data menjadi 80% untuk pelatihan dan 20% untuk pengujian menggunakan `train_test_split`.

- d. Model KNN dilatih menggunakan data latih (X_{train} , y_{train}) dengan parameter nilai k tertentu (jumlah tetangga terdekat).
- e. Prediksi dan evaluasi model digunakan untuk memprediksi data uji (x_{test}), kemudian dievaluasi menggunakan metrik seperti akurasi, precision, recall, dan F1-score.

```
[ ] ## KNN Training dan Prediction
    knn = KNeighborsClassifier(n_neighbors=20)
    knn.fit(X_train, y_train)
    y_pred = knn.predict(X_test)
```

Gambar 5. 1 Implementasi model (dengan kode)

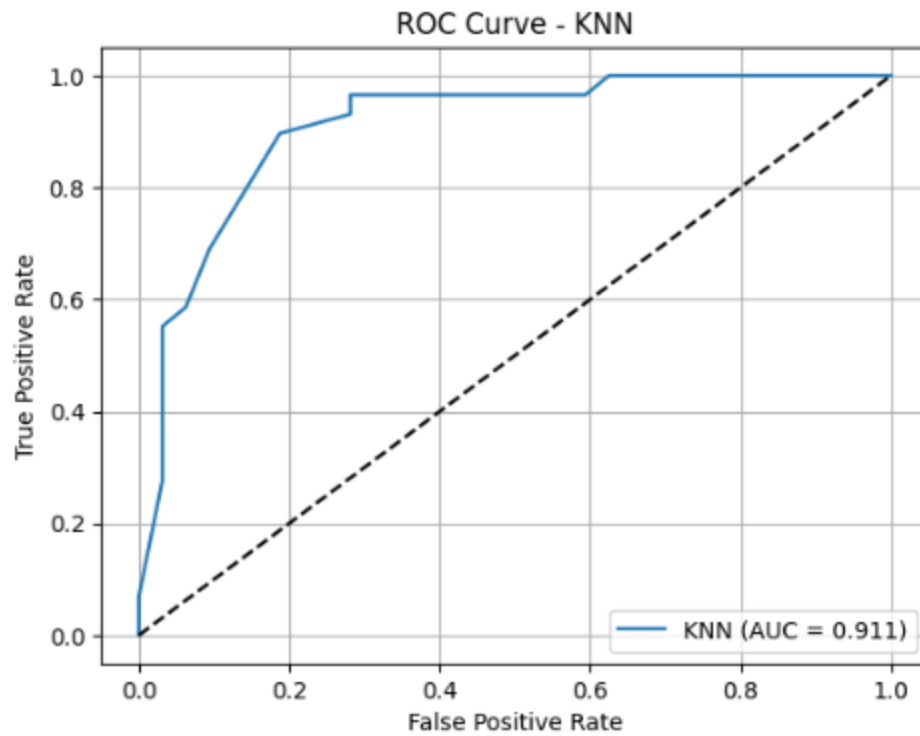
5.4 Visualisasi model

a. ROC Curve

ROC Curve (Receiver Operating Characteristic Curve), digunakan untuk mengevaluasi performa model klasifikasi biner, dalam hal ini model KNN (K-Nearest Neighbors). Grafik ini menggambarkan hubungan antara:

- True Positive Rate (Sensitivity / Recall) di sumbu Y
- False Positive Rate di sumbu X

Model yang baik akan menghasilkan kurva ROC yang mendekati pojok kiri atas grafik, dan semakin besar nilai AUC (Area Under Curve), maka semakin baik performa model dalam membedakan antara kelas positif dan negatif.



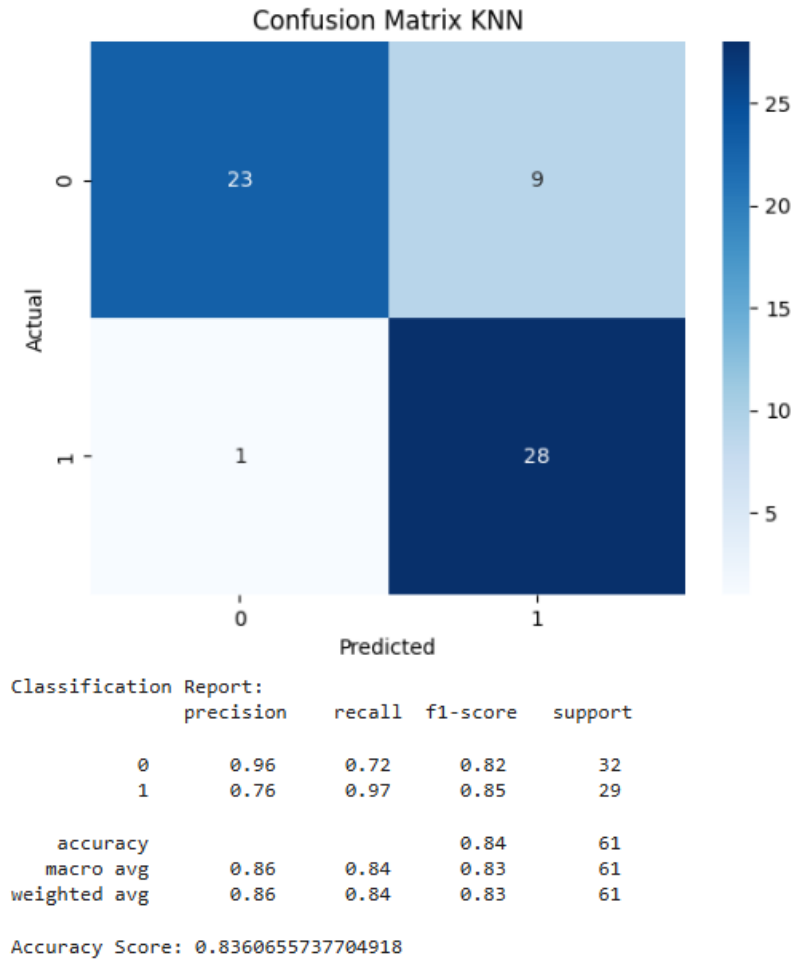
Gambar 5. 2 Visualisasi Model ROC Curve

6. EVALUATION

7.1 Confusion Matrix

Fungsi `confusion_matrix(y_test, y_pred)` digunakan untuk membandingkan hasil prediksi model (`y_pred`) dengan label aktual (`y_test`) pada data uji. Hasilnya disajikan dalam bentuk matriks 2x2 yang menunjukkan performa model klasifikasi dalam empat kategori utama:

- True Positives (TP): Pasien sakit jantung yang diprediksi benar. (90)
Sebanyak 90 pasien yang benar-benar memiliki penyakit jantung berhasil diprediksi dengan benar oleh model. Ini menunjukkan tingkat keberhasilan model dalam mengenali kasus positif (berpenyakit).
- True Negatives (TN): Pasien tidak sakit jantung yang diprediksi benar.(76)
Sebanyak 76 pasien yang tidak memiliki penyakit jantung juga diprediksi dengan benar sebagai negatif oleh model. Ini menunjukkan keandalan model dalam mendeteksi pasien sehat.
- False Positives (FP): Pasien tidak sakit jantung tetapi diprediksi sakit.(26)
Sebanyak 26 pasien sehat diprediksi salah sebagai memiliki penyakit jantung. Kesalahan ini dikenal sebagai false alarm, dan meskipun tidak berbahaya secara klinis, dapat menyebabkan kecemasan serta pengujian medis yang tidak perlu.
- False Negatives (FN): Pasien sakit jantung tetapi diprediksi tidak sakit.(13)
Sebanyak 13 pasien yang sebenarnya menderita penyakit jantung diprediksi sebagai tidak berpenyakit. Ini adalah jenis kesalahan yang paling berisiko karena dapat menyebabkan keterlambatan diagnosis dan pengobatan.



Gambar 6. 1 Confusion Matix KNN

7.2 Metrik Evaluasi: Accuracy, Precision, Recall, F1-Score

a. Accuracy

Akurasi mengukur seberapa banyak prediksi model yang benar dibandingkan dengan total seluruh data.

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

Dalam kasus ini, model memprediksi 90 TP dan 76 TN dari total 205 data uji, menghasilkan akurasi sekitar:

$$\frac{90+76}{90+76+26+13} = \frac{166}{205} = 81\%$$

(Model mampu memberikan prediksi benar untuk 81% dari seluruh pasien)

b. Precision

Precision mengukur dari semua pasien yang diprediksi positif (sakit jantung), berapa banyak yang benar-benar positif.

$$\text{Precision} = \text{TP} \setminus \text{TP} + \text{FP} = 90 \setminus 90 + 26 = 77.6\%$$

(Dari semua pasien yang dikatakan “sakit” oleh model, sekitar 78% memang benar-benar sakit)

c. Recall

Recall mengukur dari semua pasien yang benar-benar sakit, berapa banyak yang berhasil dikenali oleh model.

$$\text{Recall} = \text{TP} \setminus \text{TP} + \text{FN} = 90 \setminus 90 + 13 = 87.4\%$$

(Model berhasil mendeteksi sekitar 87% dari total pasien yang benar-benar memiliki penyakit jantung.)

d. F1-Score

F1-score adalah rata-rata harmonis dari Precision dan Recall. Digunakan untuk menyeimbangkan keduanya, khususnya saat data tidak seimbang.

$$\text{F1} = 2 \times \text{Precision} \times \text{Recall} \setminus \text{Precision} + \text{Recall} = 82.2\%$$

(F1-score menunjukkan performa keseluruhan model dalam menyeimbangkan kesalahan False Positive dan False Negative.)

7.3 Penjelasan kinerja model berdasarkan metrik tersebut

Berdasarkan hasil evaluasi menggunakan metrik accuracy, precision, recall, dan F1-score, model K-Nearest Neighbor (KNN) dalam proyek ini menunjukkan performa yang cukup baik dalam mengklasifikasikan pasien dengan dan tanpa penyakit jantung.

Model memiliki akurasi sebesar 81%, yang berarti dari seluruh data uji, 81% prediksi model sesuai dengan kondisi sebenarnya. Ini menunjukkan bahwa secara umum model mampu mengklasifikasikan data dengan tingkat kesalahan yang relatif rendah.

Dari sisi recall, model memperoleh nilai sekitar 87%, yang menandakan kemampuannya cukup tinggi dalam mengenali pasien yang benar-benar memiliki penyakit jantung. Ini penting karena dalam konteks medis, keberhasilan dalam mendeteksi pasien yang berisiko sangat krusial.

Namun, precision model sebesar 78% menunjukkan bahwa tidak semua prediksi positif benar-benar akurat. Masih terdapat sejumlah kasus di mana pasien yang sebenarnya sehat

diprediksi menderita penyakit jantung (false positive). Meski begitu, nilai ini masih tergolong baik dan dapat ditingkatkan dengan optimasi parameter.

Secara keseluruhan, F1-score sebesar 82% mencerminkan keseimbangan yang cukup stabil antara kemampuan model dalam mendeteksi kasus positif dan menghindari kesalahan prediksi. Artinya, model cukup andal dan layak digunakan sebagai sistem pendukung keputusan awal dalam deteksi penyakit jantung.

7. KESIMPULAN DAN REKOMENDASI

7.1 Ringkasan hasil modeling dan evaluasi

Model klasifikasi penyakit jantung pada proyek ini dibangun menggunakan algoritma *K-Nearest Neighbor* (KNN), yang merupakan salah satu algoritma pembelajaran terawasi (*supervised learning*) berbasis jarak. Proses modeling diawali dengan melakukan normalisasi terhadap seluruh fitur numerik menggunakan metode *StandardScaler* agar skala antar fitur menjadi seragam, mengingat KNN sangat sensitif terhadap perbedaan skala data. Selanjutnya, dilakukan reduksi dimensi menggunakan *Principal Component Analysis* (PCA) untuk menyederhanakan kompleksitas data dari 13 fitur menjadi 8 komponen utama, tanpa kehilangan informasi penting yang berdampak terhadap performa model. Dataset kemudian dibagi menjadi dua bagian, yakni 80% untuk data latih dan 20% untuk data uji.

Model KNN dilatih menggunakan data latih dan diuji terhadap data uji. Evaluasi performa dilakukan menggunakan beberapa metrik utama, yaitu *accuracy*, *precision*, *recall*, dan *F1-score*. Hasil evaluasi menunjukkan bahwa model mampu mencapai akurasi sebesar 81%, yang berarti 81% dari total prediksi pada data uji sesuai dengan label sebenarnya. Selain itu, model juga memperoleh nilai *precision* sebesar 77,6%, yang menunjukkan bahwa dari semua pasien yang diprediksi menderita penyakit jantung, sekitar 78% di antaranya memang benar-benar sakit. Sementara itu, nilai *recall* yang cukup tinggi yaitu 87,4% mengindikasikan bahwa model berhasil mendeteksi mayoritas pasien yang benar-benar memiliki penyakit jantung, yang merupakan aspek krusial dalam konteks diagnosis medis. Terakhir, *F1-score* sebesar 82,2% menggambarkan bahwa model memiliki keseimbangan yang cukup baik antara kemampuan mendeteksi kasus positif dan meminimalkan kesalahan prediksi.

Dari hasil confusion matrix, diketahui bahwa model berhasil memprediksi dengan benar sebanyak 90 kasus positif (True Positives) dan 76 kasus negatif (True Negatives). Namun, masih terdapat 26 kasus False Positive di mana pasien sehat diprediksi menderita penyakit jantung, serta 13 kasus False Negative, yaitu pasien yang sebenarnya menderita penyakit jantung namun tidak terdeteksi oleh model. Meskipun demikian, performa secara keseluruhan tergolong baik dan cukup menjanjikan untuk diterapkan sebagai sistem pendukung keputusan dalam mendeteksi penyakit jantung.

7.2 Apakah tujuan proyek tercapai?

Tujuan proyek berhasil dicapai. Model prediktif berbasis KNN telah berhasil dikembangkan untuk mengklasifikasikan risiko penyakit jantung dengan hasil evaluasi yang cukup baik. Proyek ini juga berhasil menguji efektivitas parameter K, serta mengimplementasikan machine learning pada data medis nyata, yang menunjukkan potensi pemanfaatan teknologi untuk mendukung deteksi dini penyakit jantung.

7.3 Kelebihan dan keterbatasan model

a. Kelebihan:

- KNN mudah dipahami dan diimplementasikan.
- Performa model cukup akurat dengan F1-score $> 80\%$.
- Tidak memerlukan pelatihan kompleks.
- Cocok untuk data numerik yang telah dinormalisasi.

b. Keterbatasan:

- Sensitif terhadap outlier karena berbasis jarak.
- Performa menurun pada dataset besar karena komputasi jarak yang tinggi.
- Masih terdapat kesalahan prediksi terutama False Positive dan False Negative yang perlu diminimalkan.

7.4 Rekomendasi perbaikan (dataset lebih besar, algoritma lain, dll)

Untuk peningkatan ke depan, beberapa rekomendasi dapat diberikan:

- a. Menggunakan dataset yang lebih besar dan lebih bervariasi untuk melatih model agar lebih general dan robust.
- b. Menguji algoritma lain seperti Random Forest, Logistic Regression, atau Gradient Boosting yang dapat memberikan hasil lebih stabil.
- c. Melakukan tuning parameter K serta eksplorasi teknik pemilihan fitur untuk mengurangi dimensi lebih optimal.
- d. Menambahkan validasi silang (cross-validation) untuk meningkatkan keandalan hasil.
- e. Mengintegrasikan model ke aplikasi berbasis web atau mobile untuk implementasi nyata sebagai alat bantu diagnosis dini.

DAFTAR PUSTAKA

- Akbarollah, M. F., Wiyanto, W., Ardiatma, D., & Zy, A. T. (2023). Penerapan Algoritma K-Nearest Neighbor Dalam Klasifikasi Penyakit Jantung. *Journal of Computer System and Informatics (JoSYC)*, 4(4), 850–860. <https://doi.org/10.47065/josyc.v4i4.4071>
- Dewi, L. A. (2023). Klasifikasi Machine Learning Untuk Mendeteksi Penyakit Jantung Dengan Algoritma K-Nn, Decision Tree dan Random Forest. *Repository.Uinjkt.Ac.Id*. [https://repository.uinjkt.ac.id/dspace/handle/123456789/71124%0Ahttps://repository.uinjkt.ac.id/dspace/bitstream/123456789/71124/1/LIZKY ASKA DEWI-FST.pdf](https://repository.uinjkt.ac.id/dspace/handle/123456789/71124%0Ahttps://repository.uinjkt.ac.id/dspace/bitstream/123456789/71124/1/LIZKY%20ASKA%20DEWI-FST.pdf)
- Junifer Pangaribuan, J., Tanjaya, H., & Kenichi, K. (2021). Mendeteksi Penyakit Jantung Menggunakan Machine Learning Dengan Algoritma Logistic Regression. *Journal Information System Development (ISD)*, 06(02), 1–10.
- Retnoningsih, E., & Pramudita, R. (2020). Mengenal Machine Learning Dengan Teknik Supervised Dan Unsupervised Learning Menggunakan Python. *Bina Insani Ict Journal*, 7(2), 156. <https://doi.org/10.51211/biict.v7i2.1422>
- Silmi Ath Thahirah Al Azhima, D. Darmawan, N. Fahmi Arief Hakim, I. Kustiawan, M. Al Qibtiya, N. S. S. (2022). Hybrid Machine Learning Model Untuk Memprediksi Penyakit. *Jurnal Teknologi Terpadu*, 8(1), 40–46.