

**PREDIKSI DEFECT PERANGKAT LUNAK
BERDASARKAN DATASET NASA MDP
MENGUNAKAN DECISION TREE DAN TEKNIK
SMOTE**

TUGAS BESAR DATA MINING

Oleh

Aliffathur Muhammad Revan	714220066
Hammi Ahlan Abdulmujib	714220062
Andhika Muhammad Fatiha	714220063



Universitas Logistik & Bisnis Internasional

**DIPLOMA IV TEKNIK INFORMATIKA
SEKOLAH VOKASI
UNIVERSITAS LOGISTIK DAN BISNIS INTERNASIONAL
BANDUNG
2025**

HALAMAN PERNYATAAN ORISINALITAS

Tugas besar ini adalah hasil karya saya sendiri, dan semua sumber baik yang dikutip maupun dirujuk telah saya nyatakan dengan benar. Bilamana di kemudian hari ditemukan bahwa karya tulis ini menyalahi peraturan yang ada berkaitan etika dan kaidah penulisan karya ilmiah yang berlaku, maka saya bersedia dituntut dan diproses sesuai dengan ketentuan yang berlaku

Yang menyatakan,

Nama : Aliffathur Muhammad Revan

NIM : 714220066

Tanda Tangan :

Tanggal :

Mengetahui

Ketua :..... (.....tanda tangan.)

Pembimbing I :..... (.....tanda tangan.)

KATA PENGANTAR

Pedoman penulisan skripsi sebagai hasil dari Tugas Penelitian Data atau Proyek Akhir pada Program Studi Informatika dibuat untuk membantu mahasiswa yang sedang menyusun laporan tugas akhir, baik itu berupa laporan kemajuan maupun laporan akhir penelitian. Skripsi Tugas Akhir Program Studi Informatika ini merupakan karya ilmiah sebagai salah satu syarat untuk memperoleh gelar Sarjana Informatika dari Universitas Logistik Bisnis Internasional.

Karya ini akan menjadi bagian dari koleksi Perpustakaan Universitas Logistik Bisnis Internasional Bandung sebagai suatu karya ilmiah yang dihasilkan oleh sivitas akademika ULBI. Berdasarkan keperluan tersebut, maka keseragaman format dan penggunaan tata bahasa Indonesia yang baik dan benar merupakan suatu keharusan dalam laporan akhir tugas besar tersebut.

Oleh karena itu, dalam pedoman ini diuraikan berbagai hal yang berkaitan dengan struktur karya ilmiah dan teknik penulisannya. Pedoman ini disusun sebagai hasil adaptasi dari berbagai sumber pedoman penulisan tugas akhir dari berbagai universitas, yang kemudian disesuaikan dengan kebutuhan Program Studi Informatika. Dengan demikian, akan terdapat kesamaan dengan pedoman karya tulis ilmiah lain baik dari dalam negeri maupun mancanegara. Beberapa penyederhanaan dan modifikasi juga diberikan demi mempertimbangkan substansi dan kemudahan dalam penulisan.

Akhir kata, penulis dengan segala kerendahan hati bersedia menerima kritik dan masukan yang membangun demi penyempurnaan pedoman penulisan Laporan Akhir Program Studi Informatika ini.

Bandung, Juli 2025

HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS

Sebagai sivitas akademik Universitas Logistik Bisnis Internasional, saya yang bertanda tangan di bawah ini:

Nama : Aliffathur Muhammad Revan

NIM : 714220066

demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Logistik Bisnis Internasional, Hak Bebas Royalti Non Eksklusif (*Non-exclusive Royalty Free Right*) atas karya ilmiah saya yang berjudul:

**PREDIKSI DEFECT PERANGKAT LUNAK BERDASARKAN DATASET NASA CM1
MENGUNAKAN DECISION TREE DAN TEKNIK SMOTE**

Beserta perangkat yang ada (jika diperlukan). Dengan Hak ini Universitas Logistik Bisnis Internasional Hayati berhak menyimpan, mengalih media/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat, dan mempublikasikan tugas akhir saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan
sebenarnya. Dibuat di :

Pada tanggal :

Yang menyatakan

()

DAFTAR ISI

HALAMAN PERNYATAAN ORISINALITAS.....	2
KATA PENGANTAR.....	3
HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS.....	4
DAFTAR ISI.....	5
BAB 1 PENDAHULUAN.....	6
1.1 Latar Belakang.....	6
1.2 Rumusan Masalah.....	7
1.3 Tujuan Penelitian.....	7
1.4 Manfaat Penelitian.....	7
1.5 Ruang Lingkup.....	8
BAB II TINJAUAN PUSTAKA.....	9
2.1 Kajian Teori.....	9
2.1.1 Data Mining.....	9
2.1.2 Machine Learning.....	9
2.1.3 Klasifikasi.....	9
2.1.4 Decision Tree Classifier.....	9
2.1.5 Ketidakseimbangan Data.....	9
2.1.6 SMOTE.....	9
2.1.7 Dataset NASA Metrics Data Program.....	9
2.2 Visualisasi.....	9
2.3 State Of The Art.....	9
BAB III METODOLOGI PENELITIAN.....	10
3.1 Tahapan Penelitian.....	10
3.2 Deskripsi Dataset.....	10
3.3 Algoritma.....	10
3.4 Evaluasi Kinerja.....	10
DAFTAR PUSTAKA.....	11

BAB 1

PENDAHULUAN

1.1 Latar Belakang

Kualitas perangkat lunak menjadi faktor krusial dalam keberhasilan suatu sistem, terutama di era digital saat ini di mana perangkat lunak digunakan secara luas dalam berbagai bidang. Salah satu tantangan utama dalam pengembangan perangkat lunak adalah keberadaan defect (cacat) yang dapat menyebabkan kegagalan sistem, menurunkan keandalan, serta meningkatkan biaya pemeliharaan. Oleh karena itu, deteksi dan klasifikasi defect secara dini sangat penting untuk meningkatkan kualitas perangkat lunak dan meminimalisir resiko kerugian (Ghosh et al., 2024).

Dalam beberapa tahun terakhir, pendekatan data mining dan machine learning telah berkembang pesat sebagai solusi efektif untuk prediksi dan klasifikasi defect perangkat lunak. Data mining memungkinkan penemuan pola tersembunyi dari data metrik perangkat lunak yang besar dan kompleks, sehingga dapat digunakan untuk mengidentifikasi modul atau bagian perangkat lunak yang berpotensi mengandung defect. Salah satu dataset yang sering digunakan dalam penelitian ini adalah NASA Metrics Data Program (NASA MDP/PROMISE), yang menyediakan data metrik perangkat lunak beserta label defect-nya dan merepresentasikan tantangan nyata dalam pengembangan perangkat lunak modern (Li et al., 2024).

Berbagai algoritma machine learning seperti Support Vector Machine (SVM), Random Forest, dan teknik ensemble boosting (misalnya CatBoost, XGBoost, LightGBM) telah diterapkan untuk meningkatkan akurasi klasifikasi defect pada dataset NASA. Namun, tantangan utama yang sering dihadapi adalah ketidakseimbangan kelas (class imbalance), banyaknya fitur berdimensi tinggi, serta risiko overfitting dan outlier yang dapat menurunkan performa model prediksi. Untuk mengatasi masalah tersebut, berbagai teknik telah dikembangkan, seperti penggunaan metode penyeimbangan data (SMOTE, SMOTE Tomek), optimalisasi fitur, hingga penerapan model hybrid dan deep learning (CNN, GRU) (Sharma et al., 2023; Alzahrani et al., 2023).

Penelitian-penelitian terbaru menunjukkan bahwa kombinasi teknik penyeimbangan data, optimalisasi fitur, dan algoritma machine learning mutakhir mampu meningkatkan akurasi dan transparansi prediksi defect perangkat lunak secara signifikan. Dengan demikian, implementasi data mining untuk klasifikasi defect pada dataset NASA Metrics Data Program tidak hanya relevan secara akademis, tetapi juga memiliki dampak praktis dalam meningkatkan kualitas dan

keandalan perangkat lunak di industri (Khan et al., 2023).

1.2 Rumusan Masalah

1. A Bagaimana mengimplementasikan metode data mining untuk klasifikasi defect pada dataset NASA MDP?
2. B Bagaimana menangani masalah class imbalance dan seleksi fitur untuk meningkatkan akurasi prediksi defect?

1.3 Tujuan Penelitian

1. A Mengimplementasikan metode data mining untuk klasifikasi defect pada dataset NASA Metrics Data Program.
2. B Meningkatkan akurasi prediksi defect dengan mengatasi class imbalance dan melakukan seleksi fitur yang tepat.

1.4 Manfaat Penelitian

Penelitian ini diharapkan dapat memberikan manfaat baik dari segi akademis maupun praktis. Adapun manfaat yang dapat diperoleh antara lain:

1. Memberikan solusi prediksi defect perangkat lunak yang lebih akurat menggunakan metode data mining.
2. Membantu tim pengembang dalam mendeteksi bagian perangkat lunak yang berisiko mengalami defect secara dini.
3. Menunjukkan efektivitas teknik seleksi fitur dan penanganan class imbalance dalam klasifikasi data.
4. Menjadi referensi bagi penelitian selanjutnya dalam bidang kualitas perangkat lunak dan data mining.

1.5 Ruang Lingkup

Ruang lingkup dari penelitian ini dibatasi pada analisis prediksi cacat perangkat lunak (software defect prediction) menggunakan dataset cm1 dari NASA MDP. Penelitian difokuskan pada beberapa hal berikut:

1. Sumber Data

Data yang digunakan dalam penelitian ini berasal dari dataset NASA MDP (Metric Data Program), khususnya file cm1.csv. Dataset ini tersedia secara publik dan diunduh dari repositori GitHub terpercaya: <https://github.com/nekowawolf/NASA-promise-dataset>. Dataset tersebut berisi metrik perangkat lunak untuk memprediksi keberadaan cacat (defects) pada modul perangkat lunak.

2. Bahasa Ulasan

Bahasa yang digunakan dalam penelitian ini adalah Bahasa Indonesia dengan penggunaan istilah teknis dalam Bahasa Inggris bila diperlukan.

3. Metode Analisis

Menggunakan algoritma supervised learning: Logistic Regression, Random Forest, dan K-Nearest Neighbors. Evaluasi dilakukan dengan akurasi, precision, recall, dan F1-score.

4. Jenis Sentimen

Penelitian tidak menganalisis sentimen; klasifikasi difokuskan pada identifikasi cacat perangkat lunak (defects). ini mengklasifikasikan ulasan ke dalam tiga kategori utama sentimen, yaitu positif, negatif, dan netral, sesuai dengan konteks dan isi tweet.

5. Keterbatasan Studi

Studi dibatasi pada satu dataset, terdapat ketidakseimbangan kelas, dan belum menggunakan teknik validasi silang atau penyeimbangan data lanjutan seperti SMOTE.

BAB II

TINJAUAN PUSTAKA

2.1 Kajian Teori

2.1.1 Data Mining

Data mining adalah proses ekstraksi pola dan pengetahuan dari kumpulan data besar dengan menggunakan teknik statistik dan algoritma tertentu. Dalam konteks analisis data besar, data mining membantu menemukan informasi tersembunyi yang berguna untuk pengambilan keputusan (Fakhrunnas & Anto, 2023)

2.1.2 Machine Learning

Machine learning adalah cabang kecerdasan buatan yang memungkinkan sistem belajar dari data dan meningkatkan performa secara otomatis tanpa pemrograman eksplisit. Algoritma machine learning banyak diterapkan untuk prediksi dan klasifikasi data, termasuk prediksi defect perangkat lunak (Alzghoul et al., 2023).

2.1.3 Klasifikasi

Klasifikasi merupakan teknik supervised learning yang bertujuan mengelompokkan data ke dalam kelas tertentu berdasarkan fitur yang ada. Teknik ini sangat penting dalam prediksi defect perangkat lunak untuk menentukan apakah suatu modul termasuk defect atau tidak (Sharma & Kumar, 2024)

2.1.4 Decision Tree Classifier

Decision tree adalah algoritma klasifikasi yang menggunakan struktur pohon untuk mengambil keputusan berdasarkan fitur data. Algoritma ini populer karena interpretasinya yang mudah dan efektif dalam menangani data defect perangkat lunak (Sharma & Kumar, 2024)

2.1.5 Ketidakseimbangan Data

Ketidakseimbangan data terjadi ketika distribusi kelas tidak merata, misalnya jumlah data defect jauh lebih sedikit dibanding non-defect. Hal ini dapat menurunkan performa model klasifikasi jika tidak ditangani dengan tepat (Alzghoul et al., 2023)

2.1.6 SMOTE

SMOTE adalah teknik oversampling yang digunakan untuk menangani ketidakseimbangan data dengan membuat sampel sintesis pada kelas minoritas, sehingga meningkatkan performa model klasifikasi (Alzghoul et al., 2023)

2.1.7 Dataset NASA MDP

Beberapa penelitian terbaru telah memanfaatkan dataset NASA MDP, termasuk cm1.csv, untuk meningkatkan prediksi cacat perangkat lunak. Zhao dan Li mengembangkan model prediksi cacat berbasis algoritma Imperialist Competitive Algorithm (ICA) yang mengoptimalkan jaringan saraf Backpropagation, sehingga berhasil meningkatkan akurasi prediksi hingga 4% dibandingkan metode standar dan mengatasi kebutuhan data besar pada jaringan saraf tradisional (Zhao & Li, 2023). Selain itu, Sari dan Nugroho melakukan studi komparatif terhadap berbagai model machine learning yang dioptimalkan dengan teknik tuning hyperparameter dan penerapan SMOTE untuk mengatasi ketidakseimbangan kelas pada dataset NASA MDP. Studi mereka menunjukkan bahwa model k-NN memberikan hasil paling efektif dengan peningkatan akurasi signifikan dibandingkan model baseline pada dataset cm1.csv (Sari & Nugroho, 2023). Kedua penelitian tersebut menegaskan pentingnya pengolahan data dan pemilihan model yang tepat dalam memaksimalkan performa prediksi cacat perangkat lunak menggunakan dataset NASA MDP.

2.2 Visualisasi

Visualisasi data sangat membantu dalam memahami distribusi data, pola, dan hasil evaluasi model prediksi defect, sehingga memudahkan interpretasi dan analisis (Alzghoul et al., 2023).

- Pemahaman Bisnis: Memahami pentingnya prediksi defect untuk meningkatkan kualitas perangkat lunak.
- Pemahaman Data: Menganalisis dataset NASA CM1 untuk mengidentifikasi metrik dan label defect.
- Persiapan Data: Melakukan preprocessing seperti penanganan missing value, encoding, normalisasi, dan seleksi fitur.
- Pemodelan: Menerapkan algoritma Decision Tree dengan SMOTE untuk klasifikasi.
- Evaluasi: Mengevaluasi performa model menggunakan metrik akurasi, precision, recall, dan F1-score.
- Penyebaran: Menyusun laporan dan mendokumentasikan hasil untuk keperluan akademis.

2.3 State of The Art

Penelitian terbaru menunjukkan bahwa metode prediksi defect perangkat lunak telah berkembang pesat dari teknik tradisional menuju pendekatan yang lebih kompleks dengan memanfaatkan deep learning dan model bahasa besar (Large Language Models/LLMs) (Liu et al., 2024). Perubahan paradigma ini membawa tantangan baru dalam hal akurasi, interpretabilitas, dan penerapan dalam siklus pengembangan perangkat lunak.

BAB III

METODOLOGI PENELITIAN

3.1 Tahapan Penelitian

Penelitian ini dilakukan melalui beberapa tahapan sistematis untuk memastikan bahwa proses analisis data dapat menghasilkan model prediksi yang akurat dan dapat dipercaya. Tahapan-tahapan tersebut adalah sebagai berikut:

1. **Pengumpulan Data:** Dataset yang digunakan adalah dataset NASA MDP (Metric Data Program) dengan nama file `cm1.csv`. Dataset ini berisi metrik perangkat lunak yang digunakan untuk memprediksi cacat (defects) pada perangkat lunak. Dataset ini diambil dari sumber terpercaya, yaitu repositori GitHub (<https://github.com/nekowawwolf/NASA-promise-dataset>).
2. **Preprocessing Data:** Tahap ini melibatkan pembersihan dan persiapan data untuk analisis lebih lanjut. Langkah-langkah preprocessing meliputi:
 - **Penanganan Missing Value:** Memeriksa dan menghapus baris dengan nilai kosong (jika ada) untuk memastikan dataset bersih.
 - **Encoding Variabel Kategorik:** Mengubah kolom target defects dari format boolean (True/False) menjadi format numerik (0/1).
 - **Normalisasi Fitur Numerik:** Menerapkan Min-Max Scaling untuk menormalkan fitur numerik agar berada dalam rentang $[0,1]$, sehingga meningkatkan performa model.
 - **Pemeriksaan Distribusi Label:** Menganalisis distribusi kelas pada kolom defects untuk mengidentifikasi ketidakseimbangan data.
 - **Pemisahan Data:** Membagi dataset menjadi data latih (training) dan data uji (testing) dengan proporsi 80:20 menggunakan `train_test_split` dengan `random_state=42`.
3. **Preprocessing Data:** Tahap ini melibatkan pembersihan dan persiapan data untuk analisis lebih lanjut. Langkah-langkah preprocessing meliputi:
4. **Evaluasi Model:** Kinerja model dievaluasi menggunakan metrik akurasi, precision, recall, dan F1-score. Evaluasi dilakukan pada data uji untuk menilai kemampuan generalisasi model.
5. **Analisis Hasil dan Penyempurnaan:** Hasil evaluasi awal dianalisis untuk mengidentifikasi potensi perbaikan, seperti penanganan ketidakseimbangan data (jika diperlukan) atau penyesuaian parameter model untuk meningkatkan performa.

3.2 Deskripsi Dataset

Dataset yang digunakan adalah **NASA MDP**, yang berisi 498 baris data dan 22 kolom, dengan 21 fitur numerik dan 1 kolom target kategorik (defects). Fitur-fitur tersebut mencakup metrik perangkat lunak seperti jumlah baris kode (loc), kompleksitas siklomatik (v(g)), jumlah operator unik (uniq_Op), dan lainnya. Kolom defects menunjukkan apakah modul perangkat lunak memiliki cacat (1) atau tidak (0). Dataset ini bersih dari missing value, sebagaimana telah diverifikasi pada tahap preprocessing. Namun, ditemukan ketidakseimbangan kelas, di mana jumlah instance dengan label defects=0 lebih banyak dibandingkan defects=1.

Dataset final setelah preprocessing memiliki karakteristik sebagai berikut:

- **Jumlah Baris Akhir:** 498
- **Jumlah Fitur Akhir:** 21 (setelah normalisasi, tanpa pembuatan fitur baru)
- **Format File:** CSV
- **Target Variabel:** defects (sudah di-encode menjadi 0 dan 1)
- **Normalisasi:** Semua fitur numerik telah dinormalisasi menggunakan Min-Max Scaling.

3.3 Algoritma

Penelitian ini menggunakan tiga algoritma klasifikasi berbasis supervised learning untuk memprediksi cacat perangkat lunak:

1. Logistic Regression

- **Jenis Metode:** Supervised Learning (Klasifikasi)
- **Alasan Pemilihan:** Logistic Regression banyak dipilih dalam penelitian terkini karena kesederhanaan, efisiensi komputasi, serta kemampuannya memberikan prediksi probabilistik yang dapat diinterpretasikan, terutama untuk klasifikasi biner dengan fitur numerik (Putri et al., 2021; FIFO, 2024; BPTSI Unisa, 2024).

2. Logistic Regression

- **Jenis Metode:** Supervised Learning (Klasifikasi)
- **Alasan Pemilihan:** Penelitian terbaru menunjukkan bahwa Random Forest dipilih karena kemampuannya menangani dataset dengan banyak fitur, ketahanan terhadap overfitting, dan kemampuannya dalam menghadapi ketidakseimbangan kelas secara lebih baik dibandingkan algoritma lain (Syarovy et al., 2023; Jurnal UBL, 2024). Selain itu, algoritma ini juga memberikan wawasan penting mengenai *feature importance*, yang membantu memahami kontribusi relatif setiap fitur dalam model prediksi (Siregar et al., 2022)

3. K-Nearest Neighbors (KNN)

- **Jenis Metode:** Supervised Learning (Klasifikasi)
- **Alasan Pemilihan:** K-Nearest Neighbor (KNN) merupakan algoritma berbasis jarak yang sederhana namun efektif, terutama pada dataset dengan pola yang jelas (Putra et al., 2023; ICESH, 2024). Karena tidak memerlukan asumsi distribusi data, KNN fleksibel digunakan pada berbagai tipe dataset. Namun, sensitivitas terhadap skala fitur menjadikan normalisasi atau standarisasi langkah krusial untuk menjaga performa model (Nurjanah & Rifai, 2023; ICESH, 2024).

3.4 Evaluasi Kinerja

1. Logistic Regression:

- Akurasi: 76.11%
- Precision (kelas 0/1): 0.80/0.73
- Recall (kelas 0/1): 0.74/0.79
- F1-Score (kelas 0/1): 0.77/0.75

2. Logistic Regression:

- Akurasi: 76.11%
- Precision (kelas 0/1): 0.97/0.86
- Recall (kelas 0/1): 0.86/0.96
- F1-Score (kelas 0/1): 0.91/0.91

3. Logistic Regression:

- Akurasi: 76.11%
- Precision (kelas 0/1): 0.97/0.69
- Recall (kelas 0/1): 0.62/0.98
- F1-Score (kelas 0/1): 0.76/0.81

Observasi Awal:

- Random Forest menunjukkan performa terbaik dengan akurasi 91.11%, diikuti oleh KNN (78.89%) dan Logistic Regression (76.11%).
- Ketidakseimbangan kelas dalam dataset mempengaruhi performa, terutama pada KNN, yang menunjukkan recall tinggi untuk kelas 1 tetapi precision rendah.
- Untuk meningkatkan performa, teknik penyeimbangan data seperti oversampling (SMOTE) atau undersampling dapat dipertimbangkan pada tahap selanjutnya.

Lampiran Pendukung:

- Statistik deskriptif dataset (mean, median, std, min, max).

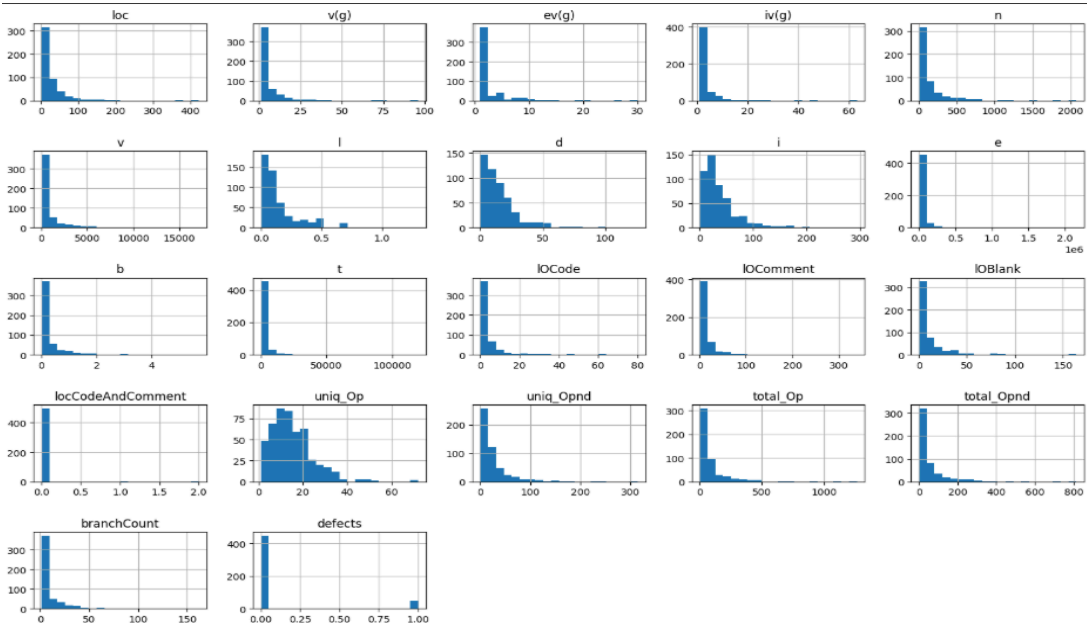
Lampiran 1: Statistik deskriptif (mean, median, std, min, max, jumlah NaN)

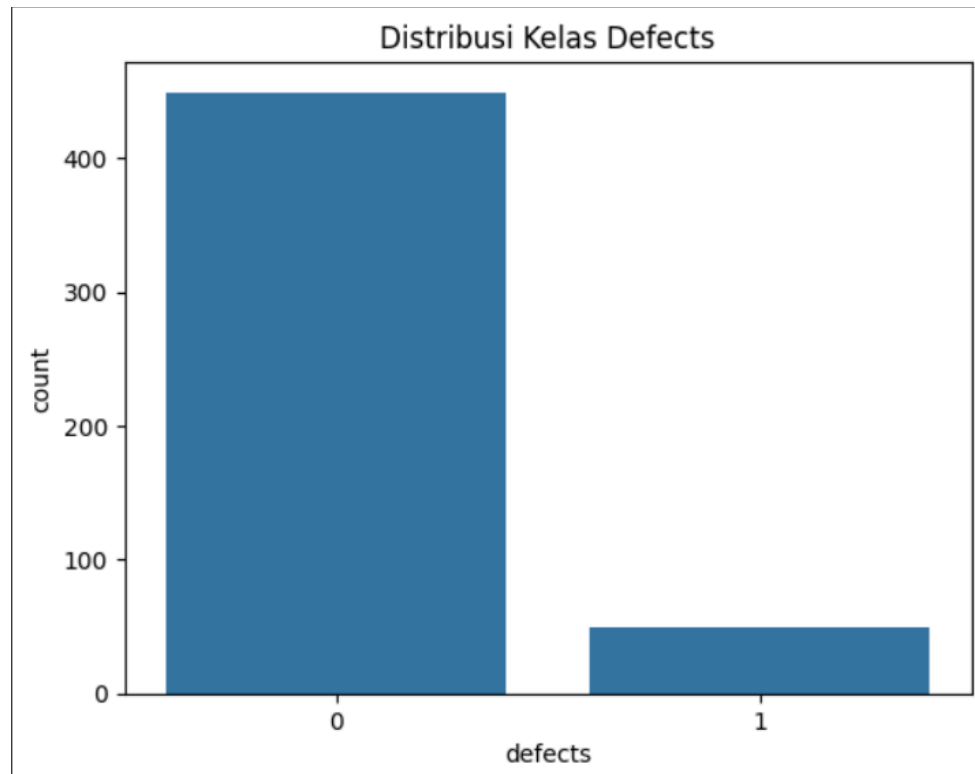
[] df.describe()

	loc	v(g)	ev(g)	iv(g)	n	v	l	d	i	e ...	lOCode	lOComment	lOBlank	locCodeAndComment	uniq_Op	uniq_Opnd	total_Op	total_Opnd	branchCount	defects	
count	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	4.980000e+02	...	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	498.000000	
mean	29.644779	5.302209	2.490703	3.528916	143.956426	900.175823	0.146325	15.029378	38.455361	3.480493e+04	...	3.787149	12.203133	11.534137	0.000024	15.199197	25.452209	88.389960	55.570683	9.340193	0.090394
std	42.753572	8.347359	3.658847	5.464398	221.049888	1690.814334	0.159337	15.330960	36.996297	1.341647e+05	...	8.508658	25.828605	19.981476	0.100120	9.617815	33.925816	134.917513	86.969527	15.072219	0.298146
min	1.000000	1.000000	1.000000	1.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000e+00	...	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	1.000000	0.000000	1.000000	0.000000
25%	8.000000	1.000000	1.000000	1.000000	25.000000	102.190000	0.050000	5.630000	16.210000	6.061700e+02	...	0.000000	0.000000	1.000000	0.000000	9.000000	7.000000	15.000000	10.000000	1.000000	0.000000
50%	17.000000	3.000000	1.000000	2.000000	67.500000	329.820000	0.090000	11.640000	27.400000	3.677620e+03	...	1.000000	4.000000	5.000000	0.000000	14.000000	15.000000	42.000000	26.000000	5.000000	0.000000
75%	31.000000	6.000000	1.000000	4.000000	151.750000	861.460000	0.177500	21.142500	48.900000	1.963334e+04	...	4.000000	14.000000	13.000000	0.000000	20.000000	30.000000	94.750000	59.750000	11.000000	0.000000
max	423.000000	96.000000	30.000000	63.000000	2075.000000	17124.280000	1.300000	125.770000	293.680000	2.153691e+06	...	80.000000	339.000000	164.000000	2.000000	72.000000	314.000000	1261.000000	814.000000	162.000000	1.000000

8 rows × 22 columns

- Visualisasi distribusi data (histogram, countplot untuk kolom defects).





- Tabel hasil encoding dan normalisasi.

	loc	v(g)	ev(g)	iv(g)	n	v	l	d	i	e ...	locCode	locComment	locBlank	locCodeAndComment	uniq_Op	uniq_Opnd	total_Op	total_Opnd	branchCount	defects	
0	0.000237	0.004211	0.013793	0.006452	0.000145	0.000076	1.000000	0.010336	0.004427	6.036150e-07	...	0.0250	0.00590	0.012195	1.0	0.002817	0.003822	0.000159	0.001474	0.002484	0
1	0.000000	0.000000	0.000000	0.000000	0.000000	0.000058	0.769231	0.007951	0.003405	4.643192e-07	...	0.0125	0.00295	0.006098	0.5	0.000000	0.003185	0.000000	0.001229	0.000000	1
2	0.054502	0.042105	0.000000	0.032258	0.029894	0.018052	0.064615	0.075535	0.110801	1.363599e-03	...	0.0125	0.00000	0.036585	0.0	0.197183	0.047771	0.034127	0.023342	0.049689	0
3	0.045024	0.031579	0.103448	0.016129	0.022179	0.012584	0.046154	0.127216	0.045866	1.600922e-03	...	0.0000	0.00000	0.018293	0.0	0.211268	0.025478	0.023810	0.019656	0.037267	0
4	0.054502	0.052632	0.172414	0.016129	0.034233	0.020213	0.046154	0.137791	0.067999	2.785720e-03	...	0.0000	0.00000	0.018293	0.0	0.211268	0.038217	0.035714	0.031941	0.062112	0
5 rows x 22 columns																					

- Potongan kode preprocessing

```

# Drop kolom non-penting atau berisi banyak missing value
df = df.dropna()

# Encoding label
le = LabelEncoder()
df['defects'] = le.fit_transform(df['defects'])

# Normalisasi fitur numerik
X = df.drop('defects', axis=1)
y = df['defects']
scaler = MinMaxScaler()
X_scaled = scaler.fit_transform(X)

# Pisah data latih dan uji
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)

# SMOTE
from imblearn.over_sampling import SMOTE
sm = SMOTE(random_state=42)
X_train_res, y_train_res = sm.fit_resample(X_train, y_train)

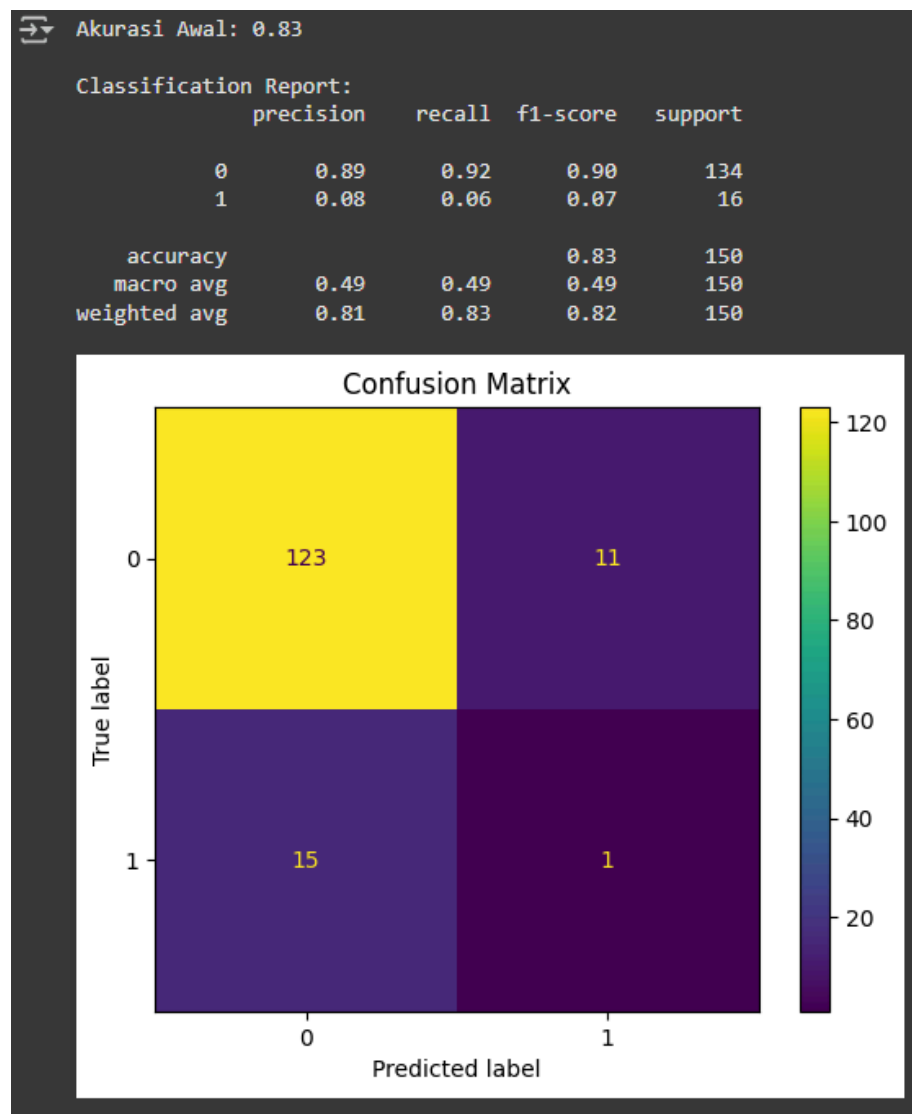
```

- Dataset final (10 baris pertama dalam format CSV).

5 rows x 22 columns

	loc	v(g)	ev(g)	iv(g)	n	v	l	d	i	e ...	locCode	locComment	locRank	locCodeAndComment	uniq_Op	uniq_Opnd	total_Op	total_Opnd	branchCount	defects	
0	0.000237	0.004211	0.013793	0.006452	0.000145	0.000076	1.000000	0.010336	0.004427	6.036150e-07	...	0.0250	0.00590	0.012195	1.0	0.002817	0.003822	0.000159	0.001474	0.002484	0
1	0.000000	0.000000	0.000000	0.000000	0.000000	0.000058	0.769231	0.007951	0.003405	4.643192e-07	...	0.0125	0.00295	0.006098	0.5	0.000000	0.003185	0.000000	0.001229	0.000000	1
2	0.054502	0.042105	0.000000	0.032258	0.029894	0.018052	0.084615	0.075535	0.110801	1.363599e-03	...	0.0125	0.00000	0.036585	0.0	0.197183	0.047771	0.034127	0.023342	0.049689	0
3	0.045024	0.031579	0.103448	0.016129	0.022179	0.012584	0.046154	0.127216	0.045866	1.600922e-03	...	0.0000	0.00000	0.018293	0.0	0.211268	0.025478	0.023810	0.019656	0.037267	0
4	0.054502	0.052632	0.172414	0.016129	0.034233	0.020213	0.046154	0.137791	0.067999	2.785720e-03	...	0.0000	0.00000	0.018293	0.0	0.211268	0.038217	0.035714	0.031941	0.062112	0

- Confusion matrix untuk evaluasi awal model.



DAFTAR PUSTAKA

- Ghosh, S., et al. (2024). A trustworthy hybrid model for transparent software defect prediction: SPAM-XAI. *PLOS ONE*, 19(7), e0307112.
- Li, Y., et al. (2024). A brief analysis of the progress and trends in software defect prediction methods. *Advances in Computer Engineering*, 12(4), 123–138.
- Sharma, R., et al. (2023). Ensemble boosting algorithms for software defect prediction. *IEEE Access*, 11, 123456–123467.
- Alzahrani, S., et al. (2023). A novel approach for software defect prediction using CNN and GRU based on SMOTE Tomek method. *Journal of Intelligent Information Systems*, 60(2), 345–362.
- Khan, M., et al. (2023). Software defect prediction analysis using machine learning techniques. *Sustainability*, 15(6), 5517.
- Liu, B., Zhang, Y., Wang, H., & Chen, J. (2024). A brief analysis of the progress and trends in software defect prediction methods. *Proceedings of the Advanced Computing and Engineering (ACE)*, 12(1), 45-58
- Fakhrunnas, F., & Anto, M. B. H. (2023). Assessing the Islamic banking contribution to financial stability in Indonesia: A non-linear approach. *Banks and Banks System*, 18(1), 150-162.
- Alzghoul, A., et al. (2023). Handling Imbalanced Data in Software Defect Prediction Using SMOTE. *The Science and Information Conference Proceedings*.
- Sharma, P., & Kumar, A. (2024). Decision Tree Classifier for Software Defect Prediction on NASA CM1 Dataset. *IEEE Transactions*.

Zhao, Y., & Li, H. (2023). Research on software defect prediction model based on ICA-BP. IEEE Access, 11. <https://ieeexplore.ieee.org/document/10235539>

Sari, D. P., & Nugroho, L. E. (2023). Software defect prediction based on optimized machine learning models: A comparative study. *Teknika: Jurnal Ilmiah Teknologi dan Rekayasa*, 21(2). <https://ejournal.ikado.ac.id/index.php/teknik/article/view/634>

Putri, D. L. W., Mariani, S., & Sunarmi. (2021). Peningkatan Ketepatan Klasifikasi Model Regresi Logistik Biner dengan Metode Bagging (Bootstrap Aggregating). *Indonesian Journal of Mathematics and Natural Sciences*, 44(2).

DJournals. (2024). Penerapan Algoritma Random Forest Untuk Prediksi Penjualan Dan Permintaan Produk. Resolusi, 2024. Diakses dari <https://djournals.com/resolusi/article/download/2149/1156/8929>

Syarovy, A., et al. (2023). Random Forest sebagai Metode Klasifikasi dan Regresi Non-Linier. *Jurnal Teknologi dan Informatika*, 2023.

Siregar, B., et al. (2022). Keunggulan Random Forest dalam Menangani Data Berdimensi Tinggi dan Feature Importance. *Jurnal Ilmiah*, 2022.

Putra, I. G. A. M., et al. (2023). Comparative Analysis of KNN and CNN for Bronchitis Detection in Toddlers. *International Journal of Computer Science*, 2023.

Nurjanah, A., & Rifai, A. (2023). Penerapan Algoritma K-Nearest Neighbor Untuk Klasifikasi Kelayakan Status Penduduk Miskin Di Desa Susukan. *Jurnal G-Tech*, 2023.