

NAMA: Syaifa Maulana
KELAS: 05 TPLE 004
NIM: 231011401802
MATA KULIAH: Machine Learning

Berikut adalah deskripsi lengkap untuk dataset iris yang digunakan dalam analisis klasifikasi:

1. Informasi Umum

- Nama Dataset: Iris Flower Dataset
- Sumber: Dikumpulkan oleh Edgar Anderson (1935) dan dipopulerkan oleh Ronald Fisher (1936)
- Jenis: Dataset klasik untuk klasifikasi multivariat
- Domain: Botani/Biologi
- Ukuran: 150 sampel

2. Visualisasi Karakteristik Dataset

2.1 Perbandingan Fitur per Spesies

- Iris-setosa:
 - Memiliki petal yang paling kecil (panjang $< 2\text{cm}$, lebar $< 0.5\text{cm}$)
 - Sepal relatif pendek dan lebar
 - Paling mudah dibedakan
- Iris-versicolor:
 - Ukuran sedang untuk semua fitur
 - Berada di antara setosa dan virginica
- Iris-virginica:
 - Memiliki petal terbesar (panjang $> 5\text{cm}$, lebar $> 1.5\text{cm}$)
 - Sepal paling panjang

2.2 Korelasi Antar Fitur

- petal length dan petal width memiliki korelasi positif tinggi (~ 0.96)
- sepal length dan petal length juga berkorelasi positif (~ 0.87)
- sepal width memiliki korelasi negatif dengan fitur lainnya

3. Karakteristik yang Membuat Dataset Ideal untuk Klasifikasi

3.1 Kelebihan Dataset

Seimbang: 50 sampel untuk setiap kelas

Bersih: Tidak ada missing values atau outlier signifikan

Terpisah dengan jelas: Iris-setosa mudah dibedakan dari dua kelas lainnya

Dimensi optimal: Cukup fitur untuk pembelajaran, tidak terlalu kompleks

Well-documented: Banyak referensi dan benchmark

4. Aplikasi dalam Machine Learning

4.1 Use Cases

- Benchmark algoritma klasifikasi
- Edukasi machine learning untuk pemula
- Testing teknik preprocessing dan feature engineering
- Demonstrasi evaluasi model multiclass

4.2 Expected Performance

- Accuracy tinggi: >95% untuk kebanyakan algoritma
- Precision/Recall: Sangat baik untuk setosa, sedikit lebih rendah untuk versicolor-virginica
- Model sederhana: Decision tree dengan depth 3 sudah cukup akurat

5. Insight dari EDA

6.1 Feature Importance

Berdasarkan analisis:

1. Petal length → Fitur paling informatif
2. Petal width → Fitur kedua terpenting
3. Sepal length → Membantu membedakan versicolor-virginica
4. Sepal width → Fitur paling tidak informatif

6.2 Visual Patterns

python

Pola yang teramati:

- Setosa: Petal kecil + Sepal lebar
- Versicolor: Ukuran sedang di semua fitur
- Virginica: Petal besar + Sepal panjang

7. Rekomendasi untuk Modeling

7.1 Preprocessing

- Scaling: Diperlukan untuk SVM dan KNN
- No encoding: Target sudah dalam format numerik
- No missing value handling: Dataset sudah bersih

7.2 Model Selection

- Linear models: Cocok karena sebagian data linearly separable
- Tree-based models: Dapat menangkap decision boundaries kompleks
- Distance-based: Membutuhkan feature scaling

8. Kesimpulan Dataset

Dataset Iris merupakan dataset yang sangat well-structured dengan:

- Kualitas data tinggi
- Pattern yang jelas
- Keseimbangan kelas sempurna
- Relevansi untuk multiple algorithms

Dataset ini sering disebut sebagai "Hello World" dalam machine learning classification karena karakteristiknya yang ideal untuk pembelajaran dan benchmarking.

KESIMPULAN EVALUASI

Kelebihan Kode Asli:

Struktur sangat baik dan terorganisir
Visualisasi komprehensif dan informatif
Implementasi algoritma lengkap
Evaluasi metrik menyeluruh
Dokumentasi yang jelas

Rekomendasi Perbaikan:

Tambah data validation (missing values, duplicates)

Implementasi cross-validation untuk evaluasi lebih robust
Detailed per-class analysis untuk insight lebih dalam
Feature importance analysis untuk understanding model
Hyperparameter tuning untuk optimasi performa

Tingkat Kematangan kode:

Completeness: 90%
Code Quality: 85%
Best Practices: 80%
Documentation: 95%
Overall Score: 87.5%