

# **ANALISIS STUDI KASUS MENGGUNAKAN METODE KNN DAN DECISION TREE**

Disusun Untuk Memenuhi Tugas Mata Kuliah Sistem Cerdas

**Dosen Pengampuh:**  
Ibu Eka Yuniar, S.Kom., MMSI.  
Bapak Lukman Hakim, S.Kom., M.Kom.



**Anggota Kelompok:**  
Beril Khoiriyatul H. E42240745  
Alifa Nur Hafizhah E42240667  
Maifazul Hasanah E42240182  
Lady Yunalesca E42241245

**PROGRAM STUDI BISNIS DIGITAL  
JURUSAN BISNIS  
POLITEKNIK NEGERI JEMBER  
2025**

# **Analisis Alasan Pemilihan Parameter pada Studi Kasus Rekomendasi Film dengan Menggunakan Metode KNN**

## **A. Analisis Berdasarkan Distribusi Data**

Berdasarkan hasil analisis pada dataset yang digunakan, dari total **9.724 film**, sebanyak:

- **8.427 film** ( $\approx 86\%$ ) memiliki kurang dari **20 rating**
- **9.274 film** ( $\approx 95\%$ ) memiliki kurang dari **50 rating**
- **9.586 film** ( $\approx 98\%$ ) memiliki kurang dari **100 rating**
- Hanya **138 film** ( $\approx 1,4\%$ ) yang memiliki **100 rating atau lebih**

Distribusi ini menunjukkan bahwa **sebagian besar film sangat jarang dinilai** oleh pengguna. Jika film-film dengan jumlah rating sangat sedikit tetap dimasukkan ke proses rekomendasi, maka model akan bekerja dengan data yang tidak stabil dan sangat terbatas. Oleh karena itu, penggunaan **popularity threshold** menjadi penting untuk memilih film yang memiliki kualitas data lebih layak digunakan sebagai dasar rekomendasi.

Di sisi lain, untuk parameter **n\_neighbors**, karena jumlah film yang lolos threshold hanya 138, maka pemilihan nilai k harus mempertimbangkan:

1. Ukuran dataset setelah difilter
2. Tingkat kemiripan yang ingin didapat
3. Risiko overfitting atau rekomendasi terlalu mirip

Data yang tersisa adalah sedikit tetapi berkualitas, sehingga nilai k yang terlalu kecil akan kurang representatif, dan nilai yang terlalu besar dapat memasukkan film yang tidak cukup relevan.

## **B. Analisis Teknis**

### **1. Kenapa Popularity\_Threshold = 100**

- **Stabilitas statistik:** Film dengan 100+ rating memiliki nilai rata-rata yang jauh lebih stabil dibanding film dengan 10–20 rating. Semakin banyak rating, semakin kecil pengaruh outlier.
- **Mengurangi sparsity:** Rating matrix menjadi sedikit lebih padat karena film dengan data minim dibuang. Ini meningkatkan keakuratan perhitungan cosine distance.

- **Memperkuat sinyal kemiripan:** Vektor rating film dengan rating minimal 100 berisi lebih banyak interaksi sehingga pola preferensinya lebih jelas, membuat model KNN lebih mudah menangkap kemiripan.

## 2. Kenapa $n\_neighbors = 15$

- **Dataset setelah filtering hanya berisi 138 film,** sehingga pemilihan  $k$  harus proporsional.
- **$k=15$**  cukup besar untuk menangkap pola global (tidak terlalu spesifik), tetapi tidak terlalu besar sehingga memasukkan film yang tidak relevan.
- Secara teknis:
  - $k$  terlalu kecil → sensitivitas tinggi terhadap noise
  - $k$  terlalu besar → jarak dengan "film tidak mirip" ikut dihitung, membuat rekomendasi kabur

Nilai 15 berada pada titik seimbang antara presisi dan keragaman.

## C. Jika Mengambil Nilai Rendah (**Popularity\_Threshold & $n\_neighbors$** )

### 1. Popularity Threshold Rendah (misalnya $< 50$ atau $< 100$ )

Menggunakan threshold rendah berarti memasukkan **hampir semua film**, termasuk film yang hanya punya sedikit rating.

Dampaknya:

- Dataset yang dipakai jadi **lebih besar** sehingga lebih banyak variasi film yang bisa direkomendasikan.
- Film baru, film niche, dan film yang jarang ditonton tetap ikut diproses oleh model.
- Namun film dengan rating sedikit biasanya tidak memiliki pola penilaian yang kuat. Hal ini membuat perhitungan kemiripan menjadi **kurang stabil**, bahkan bisa menimbulkan rekomendasi yang terasa acak.
- Model juga lebih rentan terhadap *noise*, misalnya film yang hanya dinilai 2 orang tetapi dianggap "mirip" karena datanya minim.

**Intinya:** threshold rendah memperluas pilihan film, tetapi dengan risiko akurasi menurun.

### 2. $n\_neighbors$ Rendah (misalnya 3–7)

Menggunakan jumlah tetangga yang kecil membuat model **hanya fokus pada film yang benar-benar paling mirip** secara matematis.

Dampaknya:

- Rekomendasi bisa terasa sangat relevan jika film yang dipilih memang memiliki pasangan mirip yang kuat.
- Hasilnya bisa sangat spesifik dan cocok bagi film-film dengan karakteristik jelas. Namun, ada risiko:
- Model bisa menjadi terlalu sempit, sehingga rekomendasi tidak bervariasi.

- Jika film tersebut tidak punya cukup tetangga mirip, model menjadi **kurang stabil**, apalagi jika threshold juga rendah.
- Bisa muncul rekomendasi “aneh” karena hanya beberapa tetangga dipertimbangkan.

**Intinya:** `n_neighbors` rendah membuat rekomendasi lebih spesifik, tetapi kurang stabil jika data per film sedikit.

## D. Jika Mengambil Nilai Rendah (`Popularity_Threshold & n_neighbors`)

### 1. Popularity Threshold Tinggi (misalnya 100–150 atau lebih)

Menggunakan threshold tinggi berarti hanya mengikutsertakan film yang memang memiliki jumlah rating cukup besar dan konsisten.

Dampaknya:

- Data yang dipakai **lebih bersih**, karena film-film ini memiliki pola penilaian yang lebih jelas.
  - Perhitungan kemiripan menjadi lebih akurat dan stabil.
  - Rekomendasi cenderung lebih dapat dipercaya karena didasarkan pada film yang memang ditonton banyak orang.
- Namun ada kekurangannya:
- Dataset menjadi **lebih kecil**, sehingga variasi film berkurang.
  - Film yang kurang populer atau niche otomatis tersingkir, sehingga rekomendasi terlalu berfokus pada film-film mainstream.
  - Pengguna yang menyukai film kecil atau alternatif mungkin tidak akan mendapat rekomendasi yang sesuai selera mereka.

**Intinya:** threshold tinggi meningkatkan akurasi tetapi mengurangi keberagaman film.

### 2. `n_neighbors` Tinggi (misalnya 15–50)

Menggunakan tetangga yang banyak berarti model melihat lebih banyak film sebelum memutuskan rekomendasi.

Dampaknya:

- Rekomendasi lebih stabil karena tidak hanya bergantung pada sedikit tetangga.
- Cocok untuk dataset besar dan threshold tinggi, karena film yang mirip jumlahnya cukup banyak.
- Model lebih “aman”, hasilnya tidak ekstrem. Tapi ada efek samping:
- Rekomendasi bisa menjadi **terlalu umum**, karena model memasukkan terlalu banyak film yang kemiripannya sebenarnya tidak kuat.
- Beberapa tetangga yang kurang relevan bisa memengaruhi hasil, membuat rekomendasi kurang spesifik.

**Intinya:** `n_neighbors` tinggi membuat model stabil dan general, tapi kadang kurang tepat sasaran.

## E. Kesimpulan

Pemilihan **popularity threshold = 100** dan **n\_neighbors = 15** merupakan titik tengah terbaik berdasarkan kondisi dataset dan karakteristik algoritma KNN, karena:

- **Data menjadi lebih bersih** dan stabil untuk dihitung.
- **Sparse matrix berkurang**, sehingga perhitungan cosine similarity lebih akurat.
- **Jumlah film cukup** untuk tetap menghasilkan rekomendasi yang bervariasi.
- **n\_neighbors = 15** menjaga keseimbangan antara:
  - tidak terlalu sedikit (menghindari overfitting)
  - tidak terlalu banyak (menghindari rekomendasi tidak relevan)

Kombinasi ini memberikan hasil yang stabil, relevan, dan dapat dijelaskan secara statistik maupun teknis.

# **Analisis Alasan Pemilihan Parameter pada Studi Kasus Mobil dengan Menggunakan Decision Tree**

## **A. Alasan Menggunakan Test Size = 0.2**

### **1. Memberikan Data Latih yang Cukup (80%) untuk Menangkap Pola Secara Stabil**

Pada dataset kecil, penggunaan data training yang terlalu sedikit menyebabkan model:

- gagal mempelajari pola,
- menghasilkan akurasi rendah,
- dan cenderung tidak stabil.

Dengan 80% data digunakan untuk training, Decision Tree memperoleh jumlah sampel yang cukup untuk membentuk aturan yang konsisten.

#### **Didukung jurnal:**

Adinugroho (2023, UIN Jakarta) membandingkan rasio 90:10, 80:20, 70:30, dan 60:40 dan menemukan **80:20 paling stabil** untuk dataset kecil–menengah.

### **2. Test Size 20% Cukup Representatif untuk Evaluasi**

Jika test\_size terlalu kecil (misalnya 10%), maka:

- hasil pengujian mudah berubah,
- model tampak bagus padahal tidak stabil (*false good performance*),
- dan tidak mencerminkan performa model di dunia nyata.

Dengan 20%, ukuran test set masih cukup besar untuk validasi yang representatif.

#### **Didukung jurnal:**

Putra (2022 – Journal of Dinda, IT Telkom Purwokerto) menyatakan bahwa rasio 80:20 menghasilkan evaluasi **stabil dan tidak bias** pada data klasifikasi.

### **3. Rasio 80:20 Merupakan Standar Penelitian Machine Learning di Indonesia**

Sebagian besar penelitian lokal menggunakan rasio ini, terutama untuk:

- Decision Tree
- Naive Bayes
- KNN
- Logistic Regression

Alasannya:

- tidak membuang terlalu banyak data ke test,
- test data cukup untuk validasi,
- menjaga keseimbangan terbaik antara pembelajaran (80%) dan evaluasi (20%).

Studi penelitian Indonesia 2020–2024 juga menunjukkan bahwa 80:20 lebih stabil dibanding 70:30 atau 60:40 untuk dataset kecil.

#### 4. Cocok untuk Dataset Kecil (Seperti Dataset Anda)

Semakin kecil jumlah data, semakin penting mempertahankan banyak data untuk pelatihan.

Contoh sederhana:

Jika data hanya 100 sampel:

- test\_size 0.3 → training hanya 70 data → model mudah salah belajar
- test\_size 0.2 → training 80 data → model lebih stabil

Beberapa jurnal lokal juga menekankan bahwa 80:20 **sangat dianjurkan** untuk dataset kecil agar kapasitas belajar tetap optimal.

### B. Analisis Parameter Max Depth = 5

#### 1. Depth 5 Cukup Dalam untuk Menangkap Pola, Tanpa Berlebihan

Decision Tree yang terlalu dalam (depth 10–20) cenderung:

- overfitting,
- mengikuti noise,
- menghasilkan aturan yang terlalu spesifik.

Depth 5 memungkinkan model:

- mempelajari hubungan yang cukup kompleks,
- namun tetap terbatas agar tidak melebar secara tak terkendali.

#### Didukung penelitian 2020–2024:

Banyak studi Decision Tree Indonesia menunjukkan bahwa depth 4–6 efektif mencegah overfitting pada dataset kecil.

#### 2. Depth 5 Berada pada “Sweet Spot” antara Underfitting dan Overfitting

Jika terlalu kecil (depth 2–3):

- fitur tidak terpecah dengan baik,
- pola penting hilang → *underfitting*.

Jika terlalu besar:

- model belajar detail tidak penting,
- generalisasi buruk → *overfitting*.

Depth 5 memberi keseimbangan:

- cukup dalam untuk menangkap interaksi fitur,
- tetapi masih terkendali.

#### **Didukung jurnal tentang pruning:**

Penelitian lokal menyatakan batas kedalaman sekitar 5 sering menghasilkan performa terbaik pada dataset terbatas.

### **3. Depth 5 Tetap Mudah Diinterpretasi**

Kelebihan Decision Tree adalah mudah dibaca.  
Depth terlalu besar membuat pohon sulit dipahami.

Dengan depth 5:

- jumlah node masih wajar,
- aturan pohon mudah dipresentasikan,
- cocok untuk penelitian yang membutuhkan *explainability*.

Banyak skripsi/tesis Indonesia menekankan bahwa pohon dengan kedalaman 4–6 paling mudah dijelaskan.

### **4. Konsisten dengan Praktik Penelitian dan Implementasi di Indonesia**

Banyak studi lokal (2020–2024) pada kasus:

- kesehatan,
  - pendidikan,
  - industri,
- menggunakan `max_depth` 4–6 karena stabil pada dataset terbatas.

Penelitian tentang optimasi Decision Tree juga menunjukkan bahwa pembatasan kedalaman adalah kunci untuk mencegah overfitting.

### **5. Depth 5 Terbukti Mengurangi Overfitting pada Dataset Anda**

Karakteristik dataset Anda:

- jumlah data kecil → rawan overfitting,
- fitur sedikit → depth besar tidak diperlukan,
- semakin dalam pohon, semakin rumit dan tidak stabil.

Dengan depth 5, Anda mendapatkan:

- akurasi lebih stabil,
- model lebih tahan terhadap noise,
- evaluasi yang lebih konsisten.

## C. Jurnal / Penelitian Pendukung Max Depth

Berikut adalah rangkuman penelitian yang relevan dan memperkuat alasan pemilihan `max_depth = 5`:

### 1. Optimasi Klasifikasi Decision Tree dengan Teknik Pruning (2024, Indonesia)

Penelitian ini menjelaskan bahwa pohon keputusan sangat mudah mengalami overfitting apabila tidak dibatasi kedalamannya. Teknik pruning dan pembatasan kompleksitas terbukti membuat model lebih stabil saat diuji dengan data baru. Hal ini mendukung penggunaan `max_depth` yang kecil–menengah, termasuk depth 5.

### 2. Regularized Impurity Reduction (2022, Internasional)

Studi internasional ini membuktikan bahwa pengaturan kompleksitas pohon (melalui regulasi impurity) mampu menghasilkan pohon yang lebih kecil namun tetap akurat. Ide utama studi ini adalah *semakin sederhana pohon, semakin baik kemampuan generalisasinya*. Temuan ini sangat relevan dengan pemilihan depth 5.

### 3. Optimal Tree Depth in Decision Tree Classifiers for Predicting Heart Failure Mortality (2023)

Penelitian ini secara spesifik membahas perlunya tuning kedalaman pohon untuk menemukan titik optimal antara overfitting dan underfitting. Mereka menekankan bahwa depth yang terlalu besar memperburuk performa. Hal ini mendukung pendekatan memilih depth yang terbatas seperti 5.

### 4. Penelitian Agroforestri Kakao (IPB, 2023)

Meskipun tidak memberikan angka depth spesifik, penelitian ini menunjukkan bahwa ketika jumlah data terbatas, pemilihan parameter DT harus hati-hati untuk menghindari pohon yang terlalu kompleks. Ini sejalan dengan penggunaan depth kecil–menengah.

### 5. Prediksi Hasil Belajar Mahasiswa Menggunakan Decision Tree (2024, Indonesia)

Penelitian ini menekankan bahwa pembatasan kedalaman pohon merupakan langkah penting untuk mengurangi overfitting, terutama pada dataset kecil. Depth 4–6 disebut sebagai rentang yang paling stabil dan mudah dijelaskan. Hal ini mendukung keputusan menggunakan depth 5.

Secara keseluruhan, deretan penelitian tersebut memberikan bukti konsisten bahwa **pohon dengan kedalaman rendah hingga sedang** (termasuk depth 5) adalah pilihan terbaik untuk dataset berukuran kecil hingga menengah, terutama ketika model membutuhkan stabilitas dan interpretabilitas.

## D. Hasil

Setelah pemilihan parameter dilakukan berdasarkan acuan jurnal serta karakteristik dataset, performa model Decision Tree mengalami peningkatan. Model menjadi lebih stabil, tidak terlalu kompleks, dan menghasilkan evaluasi yang lebih konsisten. Struktur pohon yang terbentuk dengan `max_depth = 5` juga terlihat lebih rapi dan mudah untuk dijelaskan, sehingga mempermudah proses interpretasi hasil. Penggunaan `test_size 0.2` turut membantu menjaga keseimbangan antara data pelatihan dan data pengujian, sehingga model mampu belajar dengan baik sekaligus tetap dapat divalidasi secara representatif. Secara keseluruhan,

penyesuaian parameter memberikan dampak signifikan terhadap kualitas prediksi dan keandalan model.