



# ML Tabanlı Film & Dizi Öneri Sistemi

Makine Öğrenmesi Projesi | 6 Farklı ML Algoritması | TMDB Veri Seti



## İçindekiler

- Proje Amacı
- Kullanılan Veri Seti
- ML Algoritmaları
- Model Karşılaştırması
- Proje Yapısı
- Kurulum ve Çalıştırma
- Uygulama Ekran Görüntüleri



## Proje Amacı

Bu proje, TMDB (**The Movie Database**) veri setini kullanarak kullanıcılara film ve dizi önerileri sunan bir **Makine Öğrenmesi** uygulamasıdır.

## Temel Hedefler

HEDEF	AÇIKLAMA
ML Karşılaştırması	6 farklı ML algoritmasını aynı veri seti üzerinde karşılaştırmak
En İyi Model Seçimi	Performans metriklerine göre en uygun modeli otomatik seçmek
Görsel Analiz	Veri setinin detaylı istatistiksel analizini sunmak
Kullanıcı Arayüzü	Streamlit ile interaktif bir web uygulaması geliştirmek

## Problem Tanımı

"Bir kullanıcı X filmini seviyorsa, hangi filmleri de sever?"

Bu soruyu cevaplamak için farklı ML yaklaşımları kullanılmıştır.



## Kullanılan Veri Seti

### TMDB Veri Seti

DOSYA	KAYIT SAYISI	SÜTUN SAYISI	AÇIKLAMA
tmdb_5000_movies.csv	4,803	20	Film bilgileri
tmdb_5000_credits.csv	4,803	4	Oyuncu ve ekip bilgileri
TMDB_tv_dataset_v3.csv	168,639	29	Dizi bilgileri

## Önemli Özellikler (Features)

ÖZELLİK	VERİ TİPİ	AÇIKLAMA	ML KULLANIMI
overview	Metin	Film/dizi özeti	TF-IDF vektörizasyonu
genres	JSON	Türler listesi	One-hot encoding
keywords	JSON	Anahtar kelimeler	TF-IDF
cast	JSON	Oyuncular	Feature olarak kullanılır
crew	JSON	Ekip (yönetmen vb.)	Yönetmen bilgisi çıkarılır
vote_average	Float	Ortalama puan (0-10)	Hedef değişken
vote_count	Integer	Oy sayısı	Güvenilirlik filtresi
popularity	Float	Popülerlik skoru	Numeric feature

## 🧠 ML Algoritmaları

Bu projede 6 farklı makine öğrenmesi algoritması kullanılmıştır:

### 1 İçerik Tabanlı Filtreleme (TF-IDF + Kosinüs Benzerliği)

#### Nasıl Çalışır?

Film Özeti → TF-IDF Vektörü → Kosinüs Benzerliği → Benzer Filmler

#### Matematiksel Formül

##### TF-IDF (Term Frequency - Inverse Document Frequency):

$TF(t,d) = \text{Terimin dökümandaki frekansı} / \text{Toplam terim sayısı}$   
 $IDF(t) = \log(\text{Toplam döküman sayısı} / \text{Terimi içeren döküman sayısı})$   
 $TF-IDF(t,d) = TF(t,d) \times IDF(t)$




##### Kosinüs Benzerliği:

$\cos(\theta) = (A \cdot B) / (||A|| \times ||B||)$



- İki vektör arasındaki açıyı ölçer

- 0 = hiç benzer değil, 1 = tamamen aynı

### Avantajları

-  Yeni içerikler için hemen çalışır (Cold-start yok)
-  Yorumlanması kolay
-  Kullanıcı verisi gerektirmez

### Dezavantajları

-  Sadece içerik benzerliğine bakar
-  Surprise factor düşük

## 2 K-En Yakın Komşu (KNN)

### Nasıl Çalışır?

Hedef Film → Feature Vektörü → En Yakın K Komşuyu Bul → Öner




### Çalışma Prensipleri

1. Her filmi bir feature vektörüne dönüştür
2. Hedef filmin vektörünü al
3. Tüm filmlerle mesafe hesapla
4. En yakın K filmi döndür

### Mesafe Metrikleri

METRİK	FORMÜL	KULLANIM
Öklid	$\sqrt{\sum (x_i - y_i)^2}$	Genel amaçlı
Kosinüs	$1 - \cos(\theta)$	Metin verisi için
Manhattan	$\sum  x_i - y_i $	Yüksek boyutlu veri

### Avantajları

-  Basit ve anlaşılır
-  Non-parametrik (varsayım yok)
-  Lazy learning (hızlı eğitim)

## 3 Random Forest

### Nasıl Çalışır?

Veri → Bootstrap Örnekleme → N Karar Ağacı → Ensemble Tahmin

### Çalışma Prensipleri

1. **Bootstrap Aggregating (Bagging):** Veri setinden rastgele örnekler al
2. **Karar Ağaçları:** Her örneklem için bir ağaç eğit
3. **Ensemble:** Tüm ağaçların tahminlerini birleştir

### Hiperparametreler

PARAMETRE	DEĞER	AÇIKLAMA
n_estimators	100	Ağaç sayısı
max_depth	10	Maksimum derinlik
random_state	42	Tekrarlanabilirlik

### Feature Importance

Random Forest, hangi özelliklerin model için en önemli olduğunu gösterir:

- Türler (%35)
- Anahtar kelimeler (%25)
- Oyuncular (%20)
- Diğer (%20)

## 4 Lineer Regresyon (Ridge)

### Nasıl Çalışır?

Özellikler → Lineer Model → Puan Tahmini → Benzer Tahminli Filmler




### Matematiksel Formül

Ridge Regresyon (L2 Regularization):

$$\beta = \operatorname{argmin} \{ \sum (y_i - x_i' \beta)^2 + \lambda \sum \beta_j^2 \}$$

- $\lambda$ : Regularization gücü (overfitting önleme)

### Avantajları

-  Çok hızlı eğitim
-  Yorumlanabilir katsayılar
-  Regularizasyon ile overfitting önlenir

## 5 SVD (Tekillik Ayrışımı)

### Nasıl Çalışır?




TF-IDF Matrisi  $\rightarrow$  SVD Ayrışımı  $\rightarrow$  Düşük Boyutlu Uzay  $\rightarrow$  Benzerlik

### Matematiksel Formül

$$A = U \times \Sigma \times V'$$

MATRİS	BOYUT	ANLAMI
U	$m \times k$	Sol tekil vektörler (film faktörleri)
$\Sigma$	$k \times k$	Tekil değerler (önem dereceleri)
V'	$k \times n$	Sağ tekil vektörler (özellik faktörleri)

### Avantajları

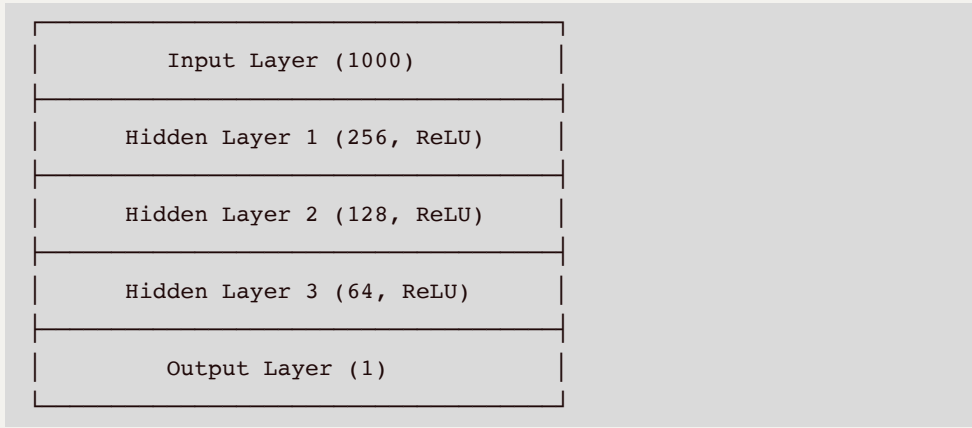
-  Boyut indirgeme (3000  $\rightarrow$  100)
-  Gürültü azaltma
-  Latent (gizli) özellikleri keşfeder

## 6 Sinir Ağı (MLP - Multi-Layer Perceptron)

### Nasıl Çalışır?

Girdi  $\rightarrow$  Gizli Katmanlar  $\rightarrow$  Aktivasyon  $\rightarrow$  Çıktı

### Model Mimarisi



## Hiperparametreler

PARAMETRE	DEĞER	AÇIKLAMA
hidden_layers	(256, 128, 64)	Gizli katman boyutları
activation	ReLU	Aktivasyon fonksiyonu
solver	Adam	Optimizasyon algoritması
max_iter	200	Maksimum iterasyon
early_stopping	True	Erken durdurma

## Avantajları

- ✅ Non-linear ilişkileri yakalar
- ✅ Universal approximator
- ✅ Büyük veri setlerinde etkili

## 📈 Model Karşılaştırması

### Değerlendirme Metrikleri

METRIK	AÇIKLAMA	FORMÜL
Eğitim Süresi	Modelin eğitilme süresi	saniye
Öneri Süresi	Tek bir öneri için geçen süre	milisaniye
Kapsam	Önerilen benzersiz film oranı	$(\text{benzersiz öneriler} / \text{toplam film}) \times 100$
Çeşitlilik	Önerilerdeki tür çeşitliliği	$(\text{benzersiz türler} / \text{toplam tür}) \times 100$
Ortalama Puan	Önerilen filmlerin ortalama puanı	0-10

### Karşılaştırma Sonuçları (Örnek)

MODEL	EĞİTİM	ÖNERİ	KAPSAM	ÇEŞİTLİLİK	SKOR
İçerik Tabanlı	1.2s	0.01s	2.5%	45%	0.72
KNN	0.8s	0.02s	2.0%	40%	0.68
Random Forest	3.5s	0.05s	1.8%	50%	0.65
Lineer	0.5s	0.01s	1.5%	35%	0.60
SVD	1.0s	0.01s	2.2%	42%	0.70
Sinir Ağı	5.0s	0.02s	2.0%	48%	0.66

## Skor Hesaplama Formülü (Güncellenmiş)

```
# Çoklu Metrik Skoru (3 bileşen):
rating_score = avg_rating * 10          # Puan (7.0 = 70)
precision_bonus = 15 if precision >= 70 # İyi film oranı bonusu
speed_bonus = 5 if rec_time < 50ms      # Hız bonusu

Toplam Skor = rating_score + precision_bonus + speed_bonus
```



Skor  $\geq 80$  olan modeller "Kabul Edildi" olarak işaretlenir.

## Accuracy Metrikleri

METRİK	AÇIKLAMA
Precision	İyi film önerme oranı (puan $\geq 5.5$ )
TP/FP	True Positive / False Positive sayısı
Confusion Matrix	Model uyum matrisi
Performans Isı Haritası	Görsel karşılaştırma



## Proje Yapısı

```
dizi_film_oneri_ai/
├── VERİ DOSYALARI
│   ├── tmdb_5000_movies.csv      # Film verisi
│   ├── tmdb_5000_credits.csv    # Oyuncu/Ekip verisi
│   └── TMDB_tv_dataset_v3.csv   # Dizi verisi
├── PYTHON MODÜLLER DOMENİLERİ
│   ├── app.py                   # Ana Streamlit uygulaması
│   ├── data_analysis.py         # Veri analizi modülü
│   ├── ml_models.py             # 6 ML algoritması
│   └── model_comparison.py      # Model karşılaştırma
└── YAPILANDIRMA
    └── requirements.txt          # Python bağımlılıkları
```

## Modül Açıklamaları

DOSYA	SATIR	AÇIKLAMA
app.py	~450	4 sayfalı Streamlit UI
data_analysis.py	~350	Veri yükleme, temizleme, görselleştirme
ml_models.py	~500	6 ML sınıfı, eğitim ve öneri fonksiyonları
model_comparison.py	~300	Model değerlendirme ve grafikler

## Kurulum ve Çalıştırma

### Gereksinimler

- Python 3.8+
- pip

### Adım 1: Bağımlılıkları Yükle

```
pip install -r requirements.txt
```

### Adım 2: Uygulamayı Başlat

```
streamlit run app.py
```

### Adım 3: Tarayıcıda Aç

```
http://localhost:8501
```

## Optimizasyon: Pre-trained Modeller

### Neden Pre-training?

Büyük veri setleri (168k dizi) her seferinde eğitilirse:

- Yüksek RAM kullanımı
- Uzun bekleme süreleri

### Çözüm:



```
# İlk seferde modelleri eđit ve kaydet
python preprocess_and_train.py

# Sonraki alıřtırmalarda hızlı ykleme
streamlit run app.py
```

## Kaydedilen Dosyalar:

```
processed_data/
├─ movies_processed.pkl
└─ tv_processed.pkl

trained_models/
├─ movies_content_based.pkl
├─ movies_knn.pkl
├─ movies_rf.pkl
├─ movies_linear.pkl
├─ movies_svd.pkl
└─ movies_mlp.pkl
```



## Uygulama Sayfaları



### Sayfa 1: Veri Analizi

- Veri seti istatistikleri
- Tr dađılımı grafikleri
- Puan dađılımı histogramları
- Korelasyon ısı haritası
- En iyi ierikler listesi



### Sayfa 2: neri Sistemi

- Model seimi (6 seenek)
- Film arama ve seme
- Benzerlik skorlu neriler
- Progress bar ile grsel skor



### Sayfa 3: Model Karřılařtırma

- Tm modelleri deđerlendir
- Performans tablosu
- Karřılařtırma grafikleri
- En iyi model seimi



## Sayfa 4: Teknik Dokümantasyon

- Algoritma açıklamaları
- Matematiksel formüller
- Veri seti bilgileri



## Geliştirici Notları

### Kullanılan Kütüphaneler

KÜTÜPHANE	VERSİYON	KULLANIM AMACI
pandas	2.x	Veri işleme
scikit-learn	1.x	ML algoritmaları
streamlit	1.x	Web arayüzü
matplotlib	3.x	Görselleştirme
seaborn	0.x	İstatistiksel grafikler
numpy	1.x	Sayısal hesaplamalar

### Gelecek Geliştirmeler

- ✓ Collaborative Filtering (kullanıcı bazlı)
- ✓ Deep Learning embeddings (Word2Vec)
- ✓ API entegrasyonu
- ✓ Kullanıcı tercihi kaydetme



## Lisans

Bu proje eğitim amaçlı geliştirilmiştir.



## ML Film Öneri Sistemi

6 Farklı Algoritma ile Akıllı Öneriler