

## پروژه درس هوش مصنوعی

### دسته بندی مشتریان بانک

دستیاران آموزشی: بهداد نادری فرد، سجاد مهرپیما، نفیسه احمدی

این پروژه به پیاده سازی و مقایسه الگوریتم های یادگیری ماشین (K-Nearest Neighbors (KNN) - ، درخت تصمیم (Decision Tree)، جنگل تصادفی (Random Forest) و Ada Boost می پردازد. هدف اصلی، ارزیابی عملکرد این الگوریتم ها در یک مسئله طبقه بندی دوتایی است.

### ساختار پروژه

پروژه از چندین بخش اصلی تشکیل شده است:

#### 1. پیاده سازی KNN

کلاس KNNClassifier با سه متد اصلی:

- `__init__(k)`: مقداردهی اولیه با پارامتر `k`
- `fit(X, y)`: آموزش مدل با داده های ورودی
- `predict(X)`: پیش بینی برچسب های نمونه های جدید

#### 2. پیش پردازش داده

تابع `load_and_preprocess_data` مراحل زیر را انجام می دهد:

- بارگذاری داده ها
- نمونه برداری متوازن از هر کلاس
- کدگذاری متغیرهای کیفی
- تقسیم داده ها به سه مجموعه آموزش، اعتبارسنجی و آزمون
- نرمال سازی داده ها

#### 3. آموزش مدل ها

دو تابع اصلی برای آموزش مدل ها:

`train_decision_tree`

- آموزش درخت تصمیم با عمق `maximum depth` برابر `d`

`train_knn`

- آموزش مدل KNN با پارامتر `k`

train\_adaboost

- آموزش مدل جنگل تصادفی با پارامترهای:
- n\_estimators: تعداد درختان در جنگل
- max\_depth: حداکثر عمق هر درخت

- آموزش مدل AdaBoost با پارامترهای:
- n\_estimators: تعداد طبقه‌بندهای پایه

**\*\* ستون لیبل‌ها، ستون bad\_loans است.**

#### 4. مقایسه مدل‌ها

تابع compare\_models برای مقایسه بصری عملکرد دو مدل با استفاده از نمودار میله‌ای

#### 5. تابع اصلی(main)

در تابع main، وظایف زیر انجام می‌شود:

- بارگذاری و پردازش داده‌ها
- بهینه‌سازی ابرپارامترها :
  - یافتن بهترین مقدار k برای KNN
  - یافتن بهترین عمق برای درخت تصمیم
- آموزش مدل‌های نهایی
- گزارش دقت روی مجموعه آزمون
- نمایش ساختار درخت تصمیم

#### نکات پیاده‌سازی

##### الف. پیاده‌سازی KNN

1. محاسبه فاصله اقلیدسی بین نقاط
2. یافتن k همسایه نزدیک
3. رأی‌گیری اکثریت برای تعیین برچسب

##### ب. تنظیم ابرپارامترها

1. برای KNN:
  - آزمایش مقادیر مختلف
  - استفاده از مجموعه اعتبارسنجی برای انتخاب بهترین k
2. برای درخت تصمیم :
  - آزمایش عمق‌های مختلف

- جلوگیری از over-fitting با محدود کردن عمق

3. برای Random Forest:

- بهینه‌سازی تعداد درختان ( $n\_estimators$ )
- تنظیم حداکثر عمق درختان برای جلوگیری از over-fitting

4. برای AdaBoost:

- تنظیم تعداد طبقه‌بندهای پایه ( $n\_estimators$ )

**\*\* توجه کنید برای یافتن بهترین هایپر پارامترها برای Random Forest از GridSearch استفاده کنید.**

#### خروجی‌های مورد انتظار

1. نمودار مقایسه‌ای دقت مدل‌ها
2. نمایش گرافیکی ساختار درخت تصمیم
3. گزارش دقت نهایی روی مجموعه آزمون

#### معیارهای ارزیابی

- دقت (Accuracy) به عنوان معیار اصلی ارزیابی
- مقایسه زمان اجرا و پیچیدگی محاسباتی