

پیاده‌سازی KNN:

```

1. class KNNClassifier:
2.     def __init__(self, k):
3.         self.k = k
4.
5.     def fit(self, X, y):
6.         self.X_train = X
7.         self.y_train = np.array(y)
8.
9.     def predict(self, X):
10.        predictions = []
11.        for x in X:
12.            distances = np.sqrt(np.sum((self.X_train - x)**2, axis=1)) # calculate
euclid. dist.
13.
14.            # get k nearest neighbors
15.            k_indices = np.argsort(distances)[:self.k]
16.            k_nearest_labels = self.y_train[k_indices]
17.
18.            prediction = np.bincount(k_nearest_labels).argmax()
19.            predictions.append(prediction)
20.        return np.array(predictions)
21.

```

در تابع fit، داده‌های تست در مدل ذخیره می‌شوند. برای انجام پیش‌بینی توسط مدل، فاصله نقطه مد نظر تا تمامی نقاط موجود در داده‌های train محاسبه می‌گردد و K داده نزدیک به آن بررسی می‌شوند و بیشترین لیبل تکراری به عنوان نتیجه گزارش می‌شود.

پیش‌پردازش و آماده‌سازی داده‌ها:

```

1. def load_and_preprocess_data(path):
2.     df = pd.read_csv(path)
3.     cols = ["grade", "term", "home_ownership", "emp_length", "bad_loans"]
4.     df = df[cols].copy() # use .copy method to prevent pd to place a view instead of a
complete dataframe
5.
6.     min_class_count = min(df['bad_loans'].value_counts()) # 1 is less, use it as the
count of samples
7.     df_balanced = pd.concat([
8.         df[df['bad_loans'] == 0].sample(min_class_count),
9.         df[df['bad_loans'] == 1].sample(min_class_count)
10.    ])
11.    # df_balanced now have the same amount of records with values 0 and 1
12.
13.    # print(df_balanced.head())
14.
15.    le = LabelEncoder()
16.    for col in ['grade', 'term', 'home_ownership', 'emp_length']:
17.        df_balanced[col] = le.fit_transform(df_balanced[col].astype(str))
18.
19.    # print(df_balanced.head())
20.
21.    X = df_balanced.drop('bad_loans', axis=1)
22.    y = df_balanced['bad_loans']
23.
24.    X_train, X_non_train, y_train, y_non_train = train_test_split(X, y, test_size=0.3,
random_state=42) # leaves 30% of data for test and valid.
25.    X_val, X_test, y_val, y_test = train_test_split(X_non_train, y_non_train, test_size=0.5,
random_state=42) # 15% for valid. and 15% test
26.
27.

```

```

28.     scaler = StandardScaler()
29.     X_train_scaled = scaler.fit_transform(X_train)    # calculate mean and variance here to
use later for valid. and test
30.     X_val_scaled = scaler.transform(X_val)
31.     X_test_scaled = scaler.transform(X_test)
32.
33.     return X_train_scaled, X_val_scaled, X_test_scaled, y_train, y_val, y_test, X.columns
34.

```

در این بخش ابتدا ستون‌های مورد نیاز جدا می‌شوند و بقیه drop می‌شوند. در ادامه، به تعداد مساوی sample از انواع bad_loan گرفته می‌شود. (به تعداد لیبل‌ی که تعداد کمتری دارد.) سپس برای متغیرهای کیفی یک لیبل عددی در نظر گرفته می‌شود.

برای فیچرهای موجود، ستون drop bad_loans می‌شود و مقدار بدست آمده در X ریخته می‌شود. همچنین خود این ستون بعنوان target در متغیر y نگهداری می‌شود.

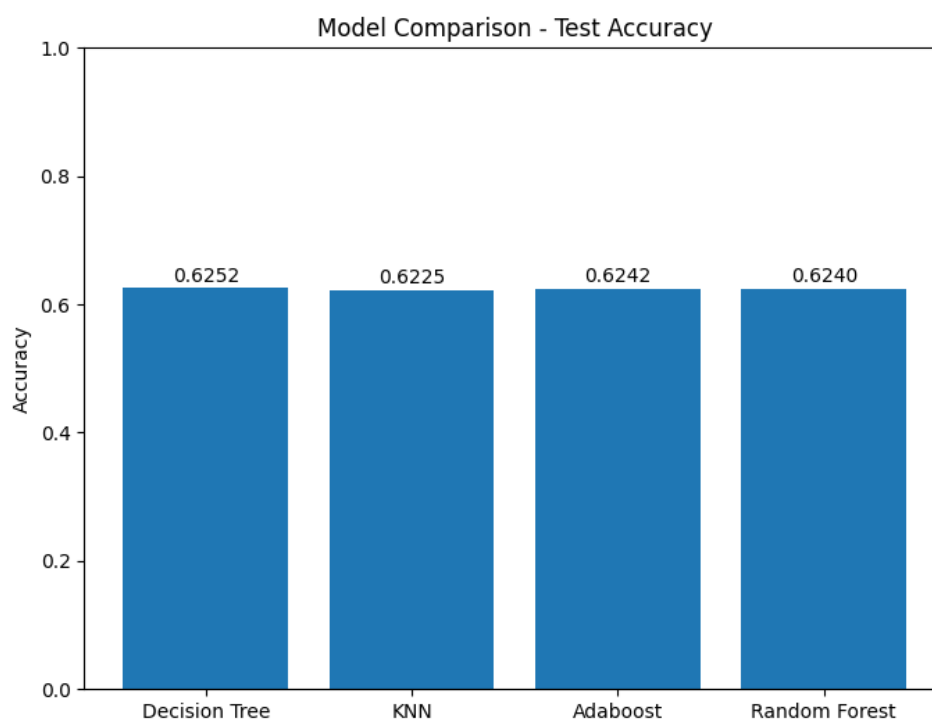
۷۰ درصد از داده موجود به عنوان داده آموزشی در نظر گرفته می‌شود. نیمی از داده باقی‌مانده برای تست و نیم دیگر برای VALIDATION به کار می‌رود.

سپس با استفاده از رابطه $z = \frac{x - \mu}{\sigma}$ داده‌ها SCALE می‌شوند تا الگوریتم‌ها نتیجه منطقی‌تری ارائه دهند. (داده‌های VALIDATION و تست نیز با میانگین و انحراف معیار داده‌های آموزشی SCALE می‌شوند.)

TUNE نمودن ابرپارامترها

برای بهینه‌سازی الگوریتم‌های KNN و درخت تصمیم از یک لیست از مقادیری که به تجربه PROMISING به نظر می‌رسند استفاده گردید. برای الگوریتم‌های جنگل تصادفی و ADABOOST نیز از GRIDSEARCH استفاده شد و بر اساس دقت، این ابرپارامترها محاسبه شدند.

نمودار دقت مدل‌ها



مقایسه الگوریتم‌ها از نظر پیچیدگی زمانی:

| | زمان یادگیری | زمان پیش‌بینی |
|----------|------------------------------|----------------------|
| KNN | $O(1) - (\text{LAZY})$ | $O(N * D)$ |
| DT | $O(N * D * \text{LOGN})$ | $O(\text{LOGN})$ |
| RF | $O(T * N * D * \text{LOGN})$ | $O(T * \text{LOGN})$ |
| ADABOOST | $O(T * N * D)$ | $O(T * D)$ |