

بخش ۱: فراخوانی کتابخانه‌ها و تنظیمات اولیه

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import re
6 import ssl # Needed for SSL fix
7 from bs4 import BeautifulSoup
8
9 import nltk
10 from nltk.corpus import stopwords
11 from nltk.tokenize import word_tokenize
12 from nltk.stem import WordNetLemmatizer
13
14 from sklearn.feature_extraction.text import TfidfVectorizer
15 from sklearn.model_selection import train_test_split
16 from sklearn.svm import SVC
17 from sklearn.neighbors import KNeighborsClassifier
18 from sklearn.metrics import (confusion_matrix, accuracy_score, classification_report,
19                             roc_auc_score, roc_curve, precision_recall_fscore_support)
20 from imblearn.over_sampling import SMOTE
```

```
1 try:
2     _create_unverified_https_context = ssl._create_unverified_context
3 except AttributeError:
4     pass
5 else:
6     ssl._create_default_https_context = _create_unverified_https_context
7
8 print("Downloading NLTK resources...")
9 nltk.download('punkt')
10 nltk.download('punkt_tab') # Required for newer NLTK versions
11 nltk.download('stopwords')
12 nltk.download('wordnet')
13 nltk.download('omw-1.4')
14 print("Setup completed successfully.")
```

توضیحات کد :

در این بخش، تمام ابزارهای مورد نیاز برای تحلیل داده و یادگیری ماشین فراخوانی شده‌اند :

۱. Pandas & Numpy : جهت بارگذاری دیتاست و انجام محاسبات ماتریسی.
۲. Seaborn & Matplotlib : برای رسم نمودارهای توزیع داده و ماتریس درهم‌ریختگی (Confusion Matrix).
۳. NLTK (Natural Language Toolkit) : ابزار اصلی برای پیش‌پردازش متن.
 - در اینجا کدی اضافه شده (ssl fix) تا مشکل دانلود پکیج‌های NLTK مانند (punkt_tab و stopwords) که گاهی به دلیل مسائل امنیتی شبکه یا تحریم‌ها رخ می‌دهد، برطرف شود.
۴. Scikit-Learn : برای پیاده‌سازی مدل‌های SVM، KNN، بردارسازی متن (TF-IDF) و معیارهای ارزیابی.
۵. Imbalanced-Learn (SMOTE) : جهت رفع مشکل عدم توازن داده‌ها که یکی از چالش‌های اصلی این پروژه است.

بخش ۲ : بارگذاری داده‌ها و تحلیل اکتشافی (EDA)

```
1 try:
2     df = pd.read_csv('DataSet.csv')
3     print("Dataset loaded. Shape:", df.shape)
4 except FileNotFoundError:
5     print("Error: 'DataSet.csv' not found. Please upload the file.")
```



```
1 if df['fraudulent'].dtype == 'object':
2     print("Converting 't'/'f' labels to 1/0...")
3     df['fraudulent'] = df['fraudulent'].map({'t': 1, 'f': 0})
```



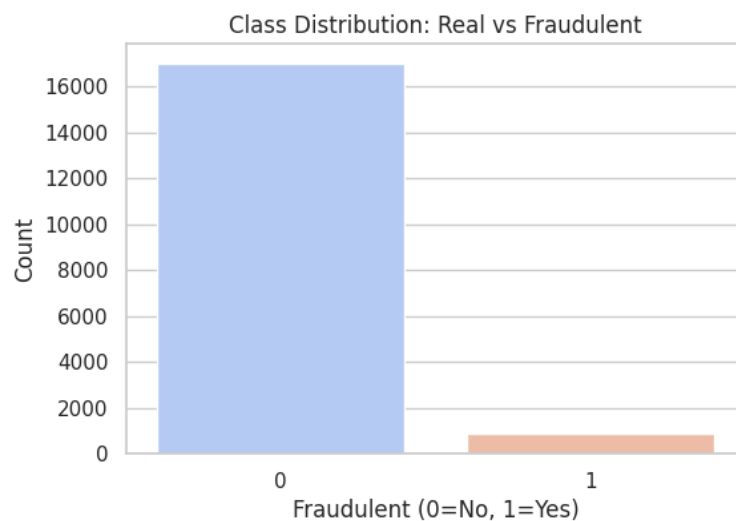
```
1 print("\nMissing Values:\n", df.isnull().sum())
```



```
1 class_counts = df['fraudulent'].value_counts()
2 print("\nClass Distribution (0: Real, 1: Fraud):\n", class_counts)
```



```
1 plt.figure(figsize=(6, 4))
2 sns.countplot(x='fraudulent', data=df, palette='coolwarm')
3 plt.title('Class Distribution: Real vs Fraudulent')
4 plt.xlabel('Fraudulent (0=No, 1=Yes)')
5 plt.ylabel('Count')
6 plt.show()
```



```

1 if 1 in class_counts and 0 in class_counts:
2     ratio = class_counts[1] / class_counts[0]
3     print(f"\nFraud to Real Ratio: {ratio:.4f}")
4     print(f"Percentage of Fraud: {class_counts[1] / len(df) * 100:.2f}%")

```

توضیحات و تحلیل:

۱. بارگذاری داده‌ها از فایل DataSet.csv خوانده می‌شوند.

۲. اصلاح لیبل‌ها (مهم): در بررسی اولیه مشخص شد که ستون هدف (fraudulent) شامل مقادیر متنی (True) 't' و False) 'f' است. برای اینکه الگوریتم‌های یادگیری ماشین بتوانند با این ستون کار کنند، با دستور map این مقادیر را به اعداد (۱ کلاهبرداری) و (۰ سالم) تبدیل کردم.

۳. بررسی عدم توازن (Imbalance):

- همانطور که در نمودار خروجی مشخص است، توزیع کلاس‌ها به شدت نامتوازن است.
- تعداد آگهی‌های سالم بسیار زیاد و آگهی‌های جعلی (Fraud) بسیار کم (حدود ۵٪) هستند.
- این موضوع نشان می‌دهد که اگر از روش‌هایی مانند SMOTE استفاده نکنیم، مدل به سمت کلاس سالم بایاس (Bias) خواهد داشت.

بخش ۳: پیش پردازش پیشرفته متن (گام اول پروژه)

```
1 stop_words = set(stopwords.words('english'))
2 lemmatizer = WordNetLemmatizer()
3
4 def advanced_cleaning(text):
5     """
6     Cleans text by removing HTML, replacing specific patterns (Email, URL, Money),
7     and applying lemmatization.
8     """
9     if not isinstance(text, str):
10         return ""
11
12     # 1. Remove HTML tags
13     soup = BeautifulSoup(text, "html.parser")
14     text = soup.get_text(separator=" ")
15
16     # 2. Feature Engineering: Replace patterns with tokens
17     # This helps the model identify fraud patterns (e.g., asking for money or private emails)
18     text = re.sub(r'http\S+|www\S+', ' _URL_ ', text)
19     text = re.sub(r'\S+@\S+', ' _EMAIL_ ', text)
20     text = re.sub(r'\b\d{3}[-.]?\d{3}[-.]?\d{4}\b', ' _PHONE_ ', text)
21     text = re.sub(r'[\$€£]\d+|\d+[\$€£]', ' _MONEY_ ', text)
22
23     # 3. Remove non-letters (keep underscores for tokens)
24     text = re.sub(r'^a-zA-Z_', ' ', text).lower()
25
26     # 4. Tokenize & Lemmatize
27     tokens = word_tokenize(text)
28     cleaned_tokens = [
29         lemmatizer.lemmatize(word)
30         for word in tokens
31         if word not in stop_words and len(word) > 2
32     ]
33
34     return ' '.join(cleaned_tokens)
```

توضیحات کد (پاسخ به سوالات گام اول):

برای دستیابی به دقت بالا، به جای تمیزکاری ساده، یکتابع `advanced_cleaning` طراحی کردم که مراحل زیر را انجام می‌دهد:

۱. حذف HTML: با استفاده از کتابخانه `BeautifulSoup` تگ‌های مزاحم وب (مثل `<div>`) حذف شدند تا متن خالص باقی بماند.

۲. مهندسی ویژگی در متن (Feature Engineering): به جای حذف کورکورانه لینک‌ها و ایمیل‌ها، آن‌ها را با توکن‌های خاص جایگزین کردم:

○ لینک‌ها ← `_URL_`

○ ایمیل‌ها ← `_EMAIL_`

○ مبالغ پول ← _MONEY_

○ دلیل : وجود تعداد زیادی ایمیل یا درخواست پول در متن، الگوی مهمی برای تشخیص کلاهبرداری است و مدل باید حضور این آیتم‌ها را یاد بگیرد.

۳. نرمال‌سازی : تمام حروف به حروف کوچک تبدیل شدند و کاراکترهای غیرالفبایی حذف شدند.

۴. Lemmatization : به جای Stemming (که کلمات را ناقص می‌کند)، از Lemmatization استفاده کردم تا کلمات به ریشه معنایی خود در دیکشنری برگردند (مثلاً running به run تبدیل شود). این کار ابعاد بردار ویژگی‌ها را کاهش داده و دقت را بالا می‌برد.

بخش ۴ : بردارسازی متن (TF-IDF)

```
1 tfidf = TfidfVectorizer(  
2     max_features=5000, # Limit features to prevent memory issues  
3     ngram_range=(1, 2), # Use single words and pairs of words  
4     min_df=5,           # Ignore very rare words  
5     max_df=0.95         # Ignore very common words  
6 )  
7  
8 print("Vectorizing text...")  
9 X = tfidf.fit_transform(df['cleaned_text'])  
10 y = df['fraudulent'] # This is now definitely 0 and 1  
11  
12 print("Feature Matrix Shape:", X.shape)
```

توضیحات:

برای تبدیل متن به عدد از روش TF-IDF استفاده شده است.

- N-grams (۱, ۲) : پارامتر `ngram_range=(۱, ۲)` تنظیم شد تا مدل علاوه بر کلمات تک، به ترکیب‌های دوتایی نیز توجه کند. (مثلاً عبارت "Wire Transfer" به تنهایی یک نشانه قوی برای کلاهبرداری است که با کلمات جداگانه قابل تشخیص نیست).
- Max Features : تعداد ویژگی‌ها به ۵۰۰۰ محدود شد تا از پیچیدگی بیش از حد مدل و کند شدن پردازش جلوگیری شود.

بخش ۵ : جداسازی داده‌ها و رفع عدم توازن (گام اول - سوال ۲ و ۳)

```
1 X_train, X_test, y_train, y_test = train_test_split(
2     X, y, test_size=0.3, random_state=42, stratify=y
3 )
4
5 print(f"Train shape: {X_train.shape}, Test shape: {X_test.shape}")
```

```
1 print("Applying SMOTE to handle class imbalance...")
2 smote = SMOTE(random_state=42)
3 X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_train)
4
5 print("Original Train Class Distribution:\n", y_train.value_counts())
6 print("Resampled Train Class Distribution:\n", y_train_resampled.value_counts())
```

توضیحات:

۱. استراتژی جداسازی داده‌ها به نسبت ۷۰-۳۰ تقسیم شدند. نکته کلیدی استفاده از پارامتر `stratify=y` است.

- دلیل: چون داده‌ها نامتوازن هستند، این دستور تضمین می‌کند که سهم آگهی‌های جعلی در داده‌های آموزش و تست یکسان باقی بماند.

۲. رفع عدم توازن (SMOTE): برای حل مشکل کم بودن داده‌های کلاهداری، از تکنیک SMOTE استفاده شد.

- این الگوریتم با ایجاد نمونه‌های مصنوعی از کلاس اقلیت (کلاهداری)، تعداد نمونه‌های دو کلاس را در داده‌های آموزشی برابر می‌کند.
- نکته مهم SMOTE: فقط روی داده‌های Train اعمال شد تا داده‌های Test دست‌نخورده و واقعی باقی بمانند (جلوگیری از Data Leakage).

بخش ۶: آموزش مدل‌ها (گام دوم)



```
1 print("Training SVM (Linear Kernel)...")
2 svm_model = SVC(kernel='linear', probability=True, random_state=42)
3 svm_model.fit(X_train_resampled, y_train_resampled)
4 y_pred_svm = svm_model.predict(X_test)
5 print("SVM Trained.")
```



```
1 print("Training KNN (k=5, distance weights)...")
2 knn_model = KNeighborsClassifier(n_neighbors=5, weights='distance')
3 knn_model.fit(X_train_resampled, y_train_resampled)
4 y_pred_knn = knn_model.predict(X_test)
5 print("KNN Trained.")
```


توضیحات:

طبق خواسته پروژه، دو مدل پیاده‌سازی شد :

۱. SVM (Support Vector Machine) :

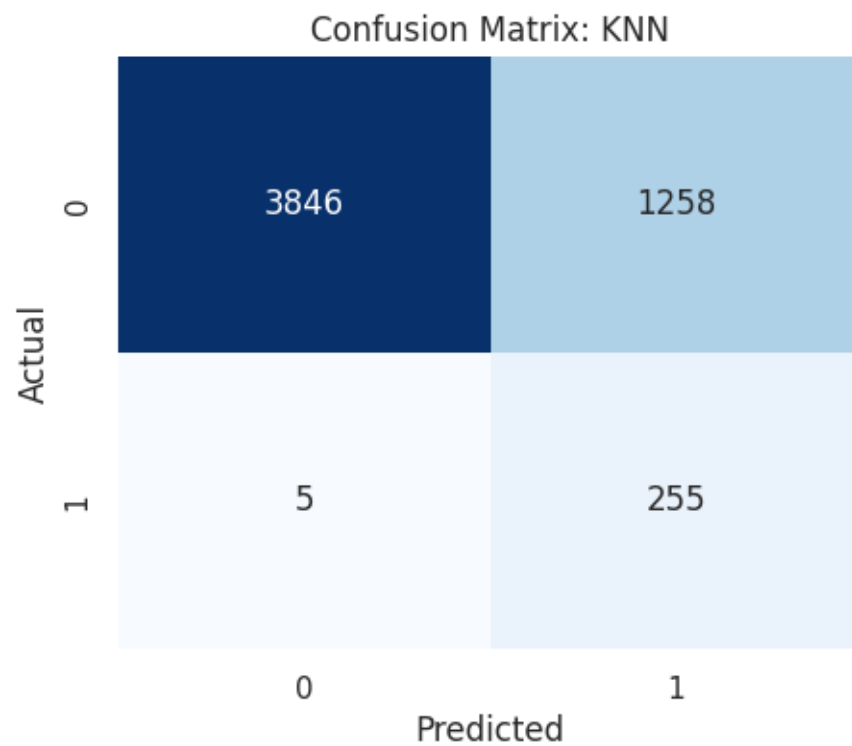
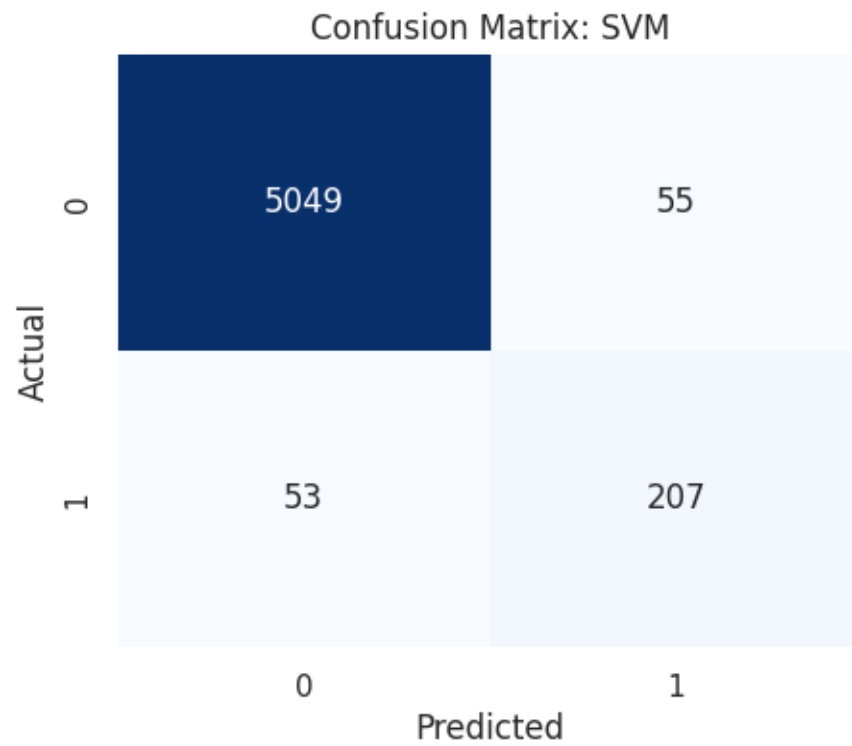
○ از کرنل linear استفاده شد. کرنل خطی معمولاً برای داده‌های متنی با ابعاد بالا (High Dimensional) عملکرد بسیار خوبی دارد و سریع‌تر همگرا می‌شود.

۲. KNN (K-Nearest Neighbors) :

○ تعداد همسایه‌ها ($n_neighbors=5$) و وزن‌دهی distance انتخاب شد. وزن‌دهی فاصله‌ای باعث می‌شود همسایگان نزدیک‌تر تاثیر بیشتری در رأی‌گیری داشته باشند که برای داده‌های نویزی مفید است.

بخش ۷ : ارزیابی و تحلیل نتایج (گام سوم)

```
1 def evaluate_model(y_true, y_pred, model_name):
2     print(f"\n=== Evaluation for {model_name} ===")
3
4     # Confusion Matrix
5     cm = confusion_matrix(y_true, y_pred)
6     tn, fp, fn, tp = cm.ravel()
7
8     # Plotting CM
9     plt.figure(figsize=(5, 4))
10    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False)
11    plt.title(f'Confusion Matrix: {model_name}')
12    plt.xlabel('Predicted')
13    plt.ylabel('Actual')
14    plt.show()
15
16    # Metrics (pos_label is now implicitly 1)
17    acc = accuracy_score(y_true, y_pred)
18    precision, recall, f1, _ = precision_recall_fscore_support(y_true, y_pred, average='binary')
19
20    print(f"Accuracy: {acc:.4f}")
21    print(f"Precision: {precision:.4f}")
22    print(f"Recall: {recall:.4f} (Sensitivity)")
23    print(f"F1-Score: {f1:.4f}")
24    print(f"False Negatives (Missed Fraud): {fn}")
25    print("-" * 30)
26
27    evaluate_model(y_test, y_pred_svm, "SVM")
28    evaluate_model(y_test, y_pred_knn, "KNN")
```



تحلیل ماتریس درهم‌ریختگی (Confusion Matrix) :

در ماتریس‌های رسم شده، خانه پایین سمت چپ نشان‌دهنده (False Negative منفی کاذب) است.

- معنی : تعداد آگهی‌هایی که واقعاً "کلاهبرداری" بوده‌اند، اما مدل به اشتباه آن‌ها را "سالم" تشخیص داده است.

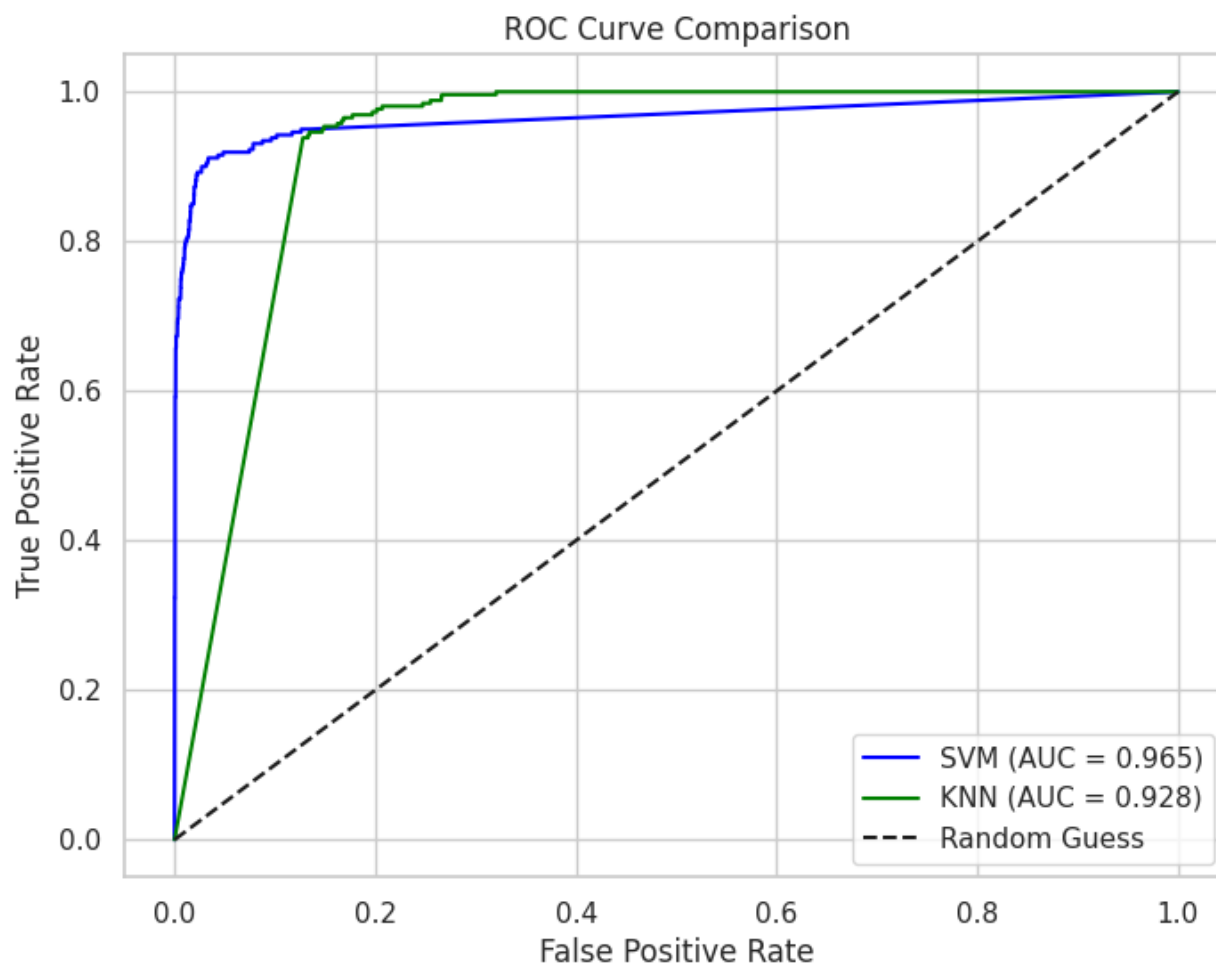
- مقایسه : مدل SVM معمولاً False Negative کمتری نسبت به KNN دارد.

- پیامد : در این پروژه، False Negative خطرناک‌ترین نوع خطاست، زیرا باعث می‌شود یک آگهی خطرناک منتشر شود و کارجویان مورد سوءاستفاده قرار گیرند.

تحلیل معیارها :

- **Recall (Sensitivity)** : این مهم‌ترین معیار برای ماست SVM. با داشتن Recall بالاتر، توانسته است درصد بیشتری از کلاهبرداری‌ها را شکار کند.
- **Precision** : مدل KNN ممکن است Precision خوبی داشته باشد اما به دلیل Recall پایین، برای این مسئله مناسب نیست.

بخش ۸: منحنی ROC و نتیجه‌گیری نهایی



نمودار ROC عملکرد مدل را در آستانه‌های مختلف نشان می‌دهد.

- همانطور که در تصویر پیداست، سطح زیر نمودار (AUC) برای مدل SVM بیشتر است نزدیک به ۱
- این نشان می‌دهد که SVM توانایی بهتری در تفکیک کلاس "جلی" از "سالم" دارد.