

Московский государственный технический университет им. Н.Э. Баумана
Кафедра «Системы обработки информации и управления»



Лабораторная работа №3
по дисциплине
«Методы машинного обучения»
на тему

Выполнила:
студент Ли Лююй
группы ИУ5И-21М

Москва — 2024 г.

1. Цель лабораторной работы

Цель лабораторной работы: изучение продвинутых способов предварительной обработки данных для дальнейшего формирования моделей.

2. Задание

1. Выбрать один или несколько наборов данных (датасетов) для решения следующих задач. Каждая задача может быть решена на отдельном датасете, или несколько задач могут быть решены на одном датасете. Просьба не использовать датасет, на котором данная задача решалась в лекции.
2. Для выбранного датасета (датасетов) на основе материалов лекций решить следующие задачи:
 - 2.1 масштабирование признаков (не менее чем тремя способами);
 - 2.2 обработку выбросов для числовых признаков (по одному способу для удаления выбросов и для замены выбросов);
 - 2.3 обработку по крайней мере одного нестандартного признака (который не является числовым или категориальным);
 - 2.4 отбор признаков;
 - 2.5 один метод из группы методов фильтрации (filter methods);
 - 2.6 один метод из группы методов обертывания (wrapper methods);
 - 2.7 один метод из группы методов вложений (embedded methods).

3. Ход выполнения работы

```
import pandas as pd
```

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler
```

```
from sklearn.feature_selection import SelectKBest, f_classif
```

```
from sklearn.ensemble import RandomForestRegressor
```

```
from sklearn.feature_selection import RFE
```

```
from sklearn.feature_selection import SelectFromModel
```

```
# Загрузка данных
```

```
df = pd.read_csv('processed_dataset.csv')
```

Масштабирование признаков

Метод 1: Стандартизация

```
scaler_standard = StandardScaler()
```

```
df_standard_scaled = pd.DataFrame(scaler_standard.fit_transform(df), columns=df.columns)
```

Метод 2: Масштабирование Min-Max

```
scaler_minmax = MinMaxScaler()
```

```
df_minmax_scaled = pd.DataFrame(scaler_minmax.fit_transform(df), columns=df.columns)
```

Метод 3: Масштабирование, устойчивое к выбросам

```
scaler_robust = RobustScaler()
```

```
df_robust_scaled = pd.DataFrame(scaler_robust.fit_transform(df), columns=df.columns)
```

Обработка выбросов в числовых признаках

Удаление выбросов: Использование метода IQR

```
Q1 = df.quantile(0.25)
```

```
Q3 = df.quantile(0.75)
```

```
IQR = Q3 - Q1
```

```
df_no_outliers = df[~((df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))).any(axis=1)]
```

Замена выбросов: Использование случайного леса для регрессии

```
def replace_outliers_with_predictions(df, column):
```

```
    model = RandomForestRegressor()
```

```
    train_data = df[df[column].notnull()]
```

```
    test_data = df[df[column].isnull()]
```

```
if not test_data.empty:

    model.fit(train_data.drop(column, axis=1), train_data[column])

    predicted_values = model.predict(test_data.drop(column, axis=1))

    df.loc[df[column].isnull(), column] = predicted_values

return df
```

columns_with_outliers = ['Age', 'Income'] # Здесь предполагается, что в этих столбцах есть выбросы

for column in columns_with_outliers:

```
    df = replace_outliers_with_predictions(df, column)
```

Обработка нестандартных признаков

Предположим, что 'Occupation' является нечисловым некатегориальным признаком, мы преобразуем его в категориальный

if 'Occupation' in df.columns:

```
    df['Occupation'] = df['Occupation'].astype('category')
```

```
    df = pd.get_dummies(df, columns=['Occupation'])
```

Выбор признаков

Фильтрационные методы: Использование SelectKBest и f_classif

```
selector_kbest = SelectKBest(score_func=f_classif, k=5)
```

```
X = df.drop('Income', axis=1)
```

```
y = df['Income']
```

```
X_kbest = selector_kbest.fit_transform(X, y)
```

Обертка: Использование рекурсивного исключения признаков (RFE)

```
estimator = RandomForestRegressor()

selector_rfe = RFE(estimator, n_features_to_select=5, step=1)

X_rfe = selector_rfe.fit_transform(X, y)
```

Встроенные методы: Использование выбора признаков на основе регуляризации

```
selector_embedded = SelectFromModel(RandomForestRegressor())

X_embedded = selector_embedded.fit_transform(X, y)
```

Экспорт обработанных данных в CSV файл

```
processed_data_path = 'processed_data.csv'

df.to_csv(processed_data_path, index=False)
```

```
print(f"Предварительная обработка данных завершена, файл '{processed_data_path}' сохранен.")
```

Результат:

```
1  UserID, Age, Income, Gender_Male, Gender_Non-Binary, Gender_Prefer not to say, Occupation_Docto
2  -0.8183965640695533, 0.9426371625691136, -1.0165262175896297, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.
3  -1.0608844349049766, -0.8186156129590085, -1.335040195829793, 0.0, 0.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.
4  -0.0476315460569581, -0.5984590160179932, -0.9186034021429064, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
5  -0.4373441956138882, 0.5023239686870831, -1.6181307756651704, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
6  1.1561475270188928, 0.2821673717460678, 1.2196817379278146, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0
7  -0.8790185317784002, 1.3095648241374724, -0.0012814499117903, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.
8  -1.1474872459176275, 1.3095648241374724, -1.503386083224787, 1.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.
9  -0.6278703798417208, -0.2315313544496344, -0.8009773806456743, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
10 0.1602152003734046, -0.0113747575086192, -0.8481379577085054, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
11 -0.5499278499303347, 1.456335888764816, -0.2229319248584839, 0.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.
12 -1.2687311813353392, -1.6992420007230695, 0.672949549390788, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0
13 -1.0349035916011813, -0.8920011452726803, -1.1041525189067318, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
14 0.2554782924873209, 1.2361792918238006, -1.2656764360347754, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0
15 1.312032586841665, -0.671844548331665, 1.1626059991105335, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0
16 -1.5718410198796182, -0.0113747575086192, -0.1921271274418816, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
17 -1.0002624671961209, 1.7498780180195028, -0.9283490739528356, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
18 0.6018895365379254, 0.7958600979417701, 0.2571907154819833, 0.0, 0.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0
19 0.723133471955637, 0.5023239686870831, 1.6626860804182908, 0.0, 0.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0
20 0.9569610616897952, -0.3049168867633062, -0.5389035539308398, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0
21 0.7837554396644929, -1.552470936095726, -0.2198387333709846, 0.0, 1.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0
22 -0.8616979695758789, -1.6992420007230695, -1.4152936845602533, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
23 0.0043301405506325, 1.0894082271964571, 0.6718902372375348, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0
24 -1.4419368033606417, -0.378302419076978, -1.573851527659188, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0
25 -0.8010760018670231, 0.4289384363734113, -0.0688231928032125, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.
26 -1.0695447160062417, 0.8692516302554418, 0.0046083256602972, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0
27 0.194856324778465, 1.2361792918238006, 0.968116287732516, 1.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0
28 0.0129904216518976, 1.7498780180195028, -0.5313612513996772, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0
29 -1.554520457677088, 0.0620107748050525, 0.4390957984386179, 1.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0
```

