



Московский государственный технический университет им. Н.Э. Баумана

Кафедра «Системы обработки информации и управления» – ИУ5

Факультет «Радиотехнический» – РТ5

Отчёт по лабораторной работе №4 по курсу Технологии машинного обучения

5

(количество листов)

Исполнитель

студент группы РТ5-616

Нижаметдинов М. Ш.

“ ____ ” _____ 2023 г.

Проверил

Преподаватель кафедры ИУ5

Гапанюк Ю. Е.

“ ____ ” _____ 2023 г.

Москва, 2023 г.

Задание

1. Выберите набор данных (датасет) для решения задачи классификации или регрессии.
2. В случае необходимости проведите удаление или заполнение пропусков и кодирование категориальных признаков.
3. С использованием метода `train_test_split` разделите выборку на обучающую и тестовую.
4. Обучите следующие модели:
 - о одну из линейных моделей (линейную или полиномиальную регрессию при решении задачи регрессии, логистическую регрессию при решении задачи классификации);
 - о SVM;
 - о дерево решений.
5. Оцените качество моделей с помощью двух подходящих для задачи метрик. Сравните качество полученных моделей.
6. Постройте график, показывающий важность признаков в дереве решений.
7. Визуализируйте дерево решений или выведите правила дерева решений в текстовом виде.

Набор данных

https://scikit-learn.org/stable/datasets/toy_dataset.html#wine-recognition-dataset

Исходный текст проекта

Задание

1. Выберите набор данных (датасет) для решения задачи классификации или регрессии.
2. В случае необходимости проведите удаление или заполнение пропусков и кодирование категориальных признаков.
3. С использованием метода `train_test_split` разделите выборку на обучающую и тестовую.
4. Обучите следующие модели:
 1. одну из линейных моделей (линейную или полиномиальную регрессию при решении задачи регрессии, логистическую регрессию при решении задачи классификации);
 2. SVM;
 3. дерево решений.
5. Оцените качество моделей с помощью двух подходящих для задачи метрик. Сравните качество полученных моделей.
6. Постройте график, показывающий важность признаков в дереве решений.
7. Визуализируйте дерево решений или выведите правила дерева решений в текстовом виде.

Ход работы

Выбор и загрузка датасета

```

# %matplotlib inline
# sns.set(style="ticks")

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.datasets import *
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn import svm, tree
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
from operator import itemgetter

def make_dataframe(ds_function):
    ds = ds_function()
    df = pd.DataFrame(data= np.c_[ds['data'], ds['target']],
                      columns= list(ds['feature_names']) + ['target'])
    return df

wine = load_wine()

df = make_dataframe(load_wine)

# Первые 5 строк датасета
df.head()

df.dtypes

# Проверим наличие пустых значений
# Цикл по колонкам датасета
for col in df.columns:
    # Количество пустых значений - все значения заполнены
    temp_null_count = df[df[col].isnull()].shape[0]
    print('{} - {}'.format(col, temp_null_count))

### Разделение на тестовую и обучающую выборки

y = df['target']
x = df.drop('target', axis = 1)

scaler = MinMaxScaler()
scaled_data = scaler.fit_transform(x)

```

```
x_train, x_test, y_train, y_test = train_test_split(scaled_data, y, test_size = 0.4, random_state = 0)
```

```
print(f'Обучающая выборка:\n{x_train, y_train}')
```

```
print(f'Тестовая выборка:\n{x_test, y_test}')
```

```
### Логическая регрессия
```

```
lr = LogisticRegression(random_state=0)
```

```
lr_prediction = lr.fit(x_train, y_train).predict(x_test)
```

```
### SVM
```

```
svc = svm.SVC(random_state=0)
```

```
svc_prediction = svc.fit(x_train, y_train).predict(x_test)
```

```
### Дерево решений
```

```
dt = DecisionTreeClassifier(random_state=0)
```

```
dt_prediction = dt.fit(x_train, y_train).predict(x_test)
```

```
### Оценка качества решений
```

```
print("Logistic regression: ", accuracy_score(y_test, lr_prediction))
```

```
print("SVM: ", accuracy_score(y_test, svc_prediction))
```

```
print("Decision tree: ", accuracy_score(y_test, dt_prediction))
```

```
print("Logistic regression: ", accuracy_score(y_test, lr_prediction))
```

```
cm = confusion_matrix(y_test, lr_prediction, labels=np.unique(df.target), normalize='true')
```

```
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=np.unique(df.target))
```

```
disp.plot()
```

```
print("SVM: ", accuracy_score(y_test, svc_prediction))
```

```
cm = confusion_matrix(y_test, svc_prediction, labels=np.unique(df.target), normalize='true')
```

```
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=np.unique(df.target))
```

```
disp.plot()
```

```
print("Decision tree: ", accuracy_score(y_test, dt_prediction))
```

```
cm = confusion_matrix(y_test, dt_prediction, labels=np.unique(df.target), normalize='true')
```

```
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=np.unique(df.target))
```

```
disp.plot()
```

```
### Важность признаков
```

```
list(zip(x.columns.values, dt.feature_importances_))
```

```
def draw_feature_importances(tree_model, X_dataset, figsize=(18,5)):
    # Sorting the values of the importance of features in descending order
    list_to_sort = list(zip(X_dataset.columns.values, tree_model.feature_importances_))
    sorted_list = sorted(list_to_sort, key=itemgetter(1), reverse = True)
    # Features names
    labels = [x for x,_ in sorted_list]
    # Features importance
    data = [x for _,x in sorted_list]
    # Graph output
    fig, ax = plt.subplots(figsize=figsize)
    ind = np.arange(len(labels))
    plt.bar(ind, data)
    plt.xticks(ind, labels, rotation='vertical')
    # Values output
    for a,b in zip(ind, data):
        plt.text(a-0.05, b+0.01, str(round(b,3)))
    plt.show()
    return labels, data
```

```
dt_fl, dt_fd = draw_feature_importances(dt, x)
```

```
### Визуализация дерева решений
```

```
tree.plot_tree(dt)
```