

Министерство науки и высшего образования Российской Федерации Федеральное государственное бюджетное образовательное учреждение высшего образования

«Московский государственный технический университет имени Н.Э. Баумана (национальный исследовательский университет)» (МГТУ им. Н.Э. Баумана)

Факультет «Информатика и системы управления» Кафедра «Системы обработки информации и управления»

Отчет по лабораторной работе №4 «Линейные модели, SVM и деревья решений»

по дисциплине «Технологии машинного обучения»

Выполнил: студент группы ИУ5Ц-84Б Тихонова Д.Д. подпись, дата

Проверил: к.т.н., доц., Ю.Е. Гапанюк подпись, дата

1. Цель лабораторной работы

Изучение линейных моделей, SVM и деревьев решений.

2. Описание задания

Выберите набор данных (датасет) для решения задачи классификации или регрессии.

В случае необходимости проведите удаление или заполнение пропусков и кодирование категориальных признаков.

С использованием метода train_test_split разделите выборку на обучающую и тестовую.

Обучите следующие модели:

одну из линейных моделей (линейную или полиномиальную регрессию при решении задачи регрессии, логистическую регрессию при решении задачи классификации);

SVM;

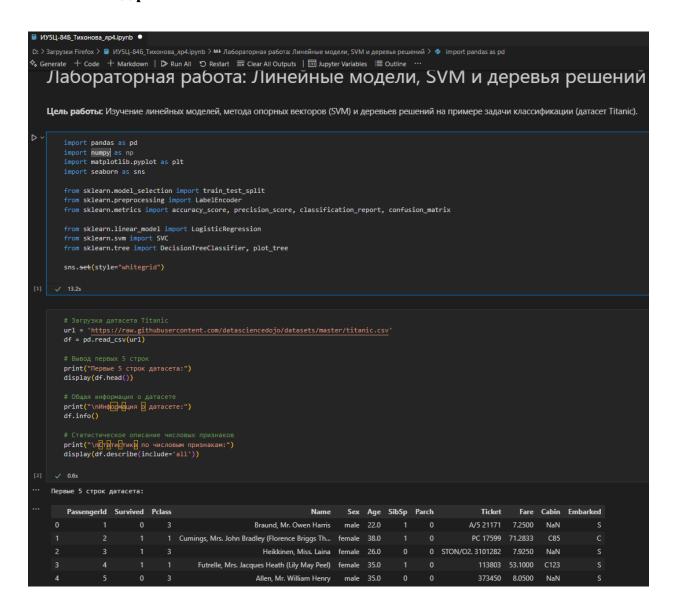
дерево решений.

Оцените качество моделей с помощью двух подходящих для задачи метрик. Сравните качество полученных моделей.

Постройте график, показывающий важность признаков в дереве решений.

Визуализируйте дерево решений или выведите правила дерева решений в текстовом виде.

3. Ход работы



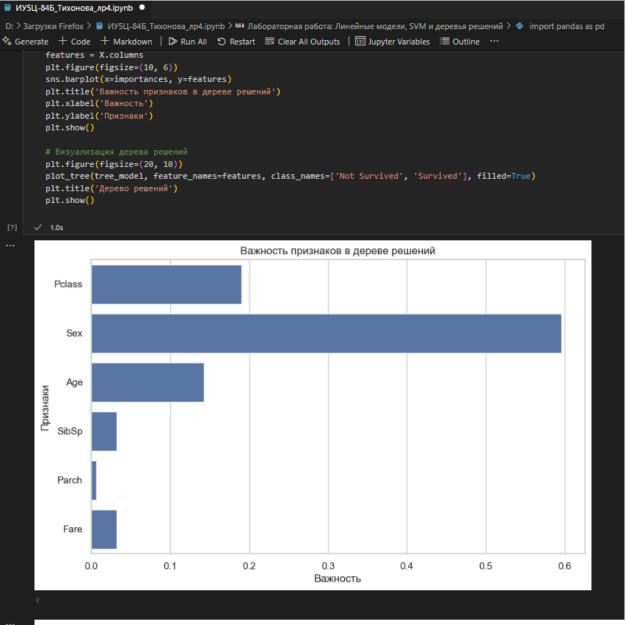
```
■ ИУ5Ц-84Б_Тихонова_лр4.ipynb •
🍫 Generate 🕂 Code 🕂 Markdown | ⊳ Run All 🖰 Restart 🚍 Clear All Outputs | 📾 Jupyter Variables 🗏 Outline 👑
    Информация о датасете:
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 891 entries, 0 to 890 Data columns (total 12 columns):
                    Non-Null Count Dtype
        PassengerId 891 non-null
                                    int64
                                     int64
        Name
                     891 non-null
                                    obiect
                     891 non-null
                                    object
                     714 non-null
         SibSp
                     891 non-null
                                    int64
                                    int64
         Parch
                     891 non-null
                                    object
         Fare
                     891 non-null
                                    float64
     10
        Cabin
                     204 non-null
                                    obiect
     11 Embarked
                     889 non-null
    dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
    Статистика по числовым признакам:
            Passengerld Survived
                                       Pclass
                                                       Name Sex
                                                                         Age
                                                                                   SibSp
                                                                                              Parch Ticket
                                                                                                                 Fare Cabin Embarked
                                                                                                     891 891.000000 204
      count 891.000000 891.000000 891.000000 891 891 891 714.000000 891.000000 891.000000
                                                                                                                                   889
                                                         891
                                                                                                       681
                   NaN
                             NaN
                                        NaN
                                                                         NaN
                                                                                    NaN
                                                                                              NaN
                                                                                                                 NaN
     unique
                   NaN
                             NaN
                                        NaN Dooley, Mr. Patrick male
                                                                         NaN
                                                                                    NaN
                                                                                              NaN 347082
                                                                                                                 NaN
       top
                   NaN
                              NaN
                                        NaN
                                                                         NaN
                                                                                    NaN
                                                                                              NaN
                                                                                                                 NaN
                                                                                                                                   644
      mean 446.000000 0.383838
                                                                    29.699118 0.523008
                                   2.308642
                                                        NaN NaN
                                                                                          0.381594
                                                                                                     NaN 32.204208
                                                                                                                        NaN
                                                                                                                                  NaN
        std 257.353842 0.486592
                                                                    14.526497
                                                                                         0.806057 NaN 49.693429
                                   0.836071
                                                        NaN NaN
                                                                                1.102743
                                                                                                                                  NaN
                                                                                                                        NaN
             1.000000 0.000000 1.000000
                                                        NaN NaN
                                                                   0.420000 0.000000 0.000000 NaN 0.000000 NaN
                                                                                                                                  NaN
       25% 223.500000 0.000000 2.000000
                                                        NaN NaN 20.125000
                                                                                0.000000 0.000000 NaN
                                                                                                            7.910400
       50% 446.000000 0.000000 3.000000
                                                        NaN NaN 28.000000
                                                                                0.000000 0.000000 NaN 14.454200
                                                                                                                        NaN
                                                                                                                                  NaN
                                                        NaN NaN 38.000000
                                                                                                     NaN 31.000000
       75% 668.500000 1.000000
                                    3.000000
                                                                                1.000000 0.000000
                                                                                                                       NaN
                                                                                                                                  NaN
       max 891.000000 1.000000
                                    3.000000
                                                        NaN NaN
                                                                    80.000000
                                                                                8.000000
                                                                                          6.000000
                                                                                                     NaN 512.329200 NaN
                                                                                                                                  NaN
       # Удалим ненужные признаки и обработаем пропуски

df = df[['Survived', 'Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare']]

df.dropna(inplace=True)
       df['Sex'] = LabelEncoder().fit_transform(df['Sex'])
       # Проверим после обработки
print("Данные после предобработки:")
       display(df.head())
```

```
В ИУ5Ц-84Б_Тихонова_лр4.ipynb
D: > Загрузки Firefox > 🚦 ИУ5Ц-845_Тихонова_лр4.ipynb > 👫 Лабораторная работа: Линейные модели, SVM и деревья решений > 💠 import pandas as pd
🗞 Generate + Code + Markdown | ▶ Run All 🤝 Restart 🚍 Clear All Outputs | 🔯 Jupyter Variables 🗏 Outline …
… Данные после предобработки:
         Survived Pclass Sex Age SibSp Parch Fare
               0 3 1 22.0
                                     1 0 7.2500
                          0 38.0
                                              0 71.2833
                                            0 7.9250
                    3 0 26.0
                                             0 53.1000
              0 3 1 35.0 0 0 8.0500
       X = df.drop('Survived', axis=1)
        y = df['Survived']
        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
        print(f"Pазмер обучающей выборки: {X_train.shape}")
print(f"Pазмер тестовой выборки: {X_test.shape}")
 … Размер обучающей выборки: (571, 6)
     Размер тестовой выборки: (143, 6)
D ~
        log_reg = LogisticRegression(max_iter=1000)
        log_reg.fit(X_train, y_train)
        y_pred_log = log_reg.predict(X_test)
        svm_model = SVC(kernel='linear')
        svm_model.fit(X_train, y_train)
        y_pred_svm = svm_model.predict(X_test)
        tree_model = DecisionTreeClassifier(max_depth=4, random_state=42)
        tree_model.fit(X_train, y_train)
        y_pred_tree = tree_model.predict(X_test)
        def evaluate_model(name, y_true, y_pred):
            print(f"\n[ {name}")
            print(f"Accuracy: {accuracy_score(y_true, y_pred):.2f}")
```

```
В ИУ5Ц-84Б_Тихонова_лр4.ipynb
D: > Загрузки Firefox > 📳 ИУ5Ц-84Б_Тихонова_лр4.ipynb > 👫 Лабораторная работа: Линейные модели, SVM и деревья решений > 💠 import pandas as pd
🗫 Generate 🕂 Code 🕂 Markdown | ⊳ Run All 😏 Restart 🚃 Clear All Outputs | 📼 Jupyter Variables 🗮 Outline …
        def evaluate_model(name, y_true, y_pred):
           print(f"\n | {name}")
            print(f"Accuracy: {accuracy_score(y_true, y_pred):.2f}")
           print("Confusion Matrix:")
            print(confusion_matrix(y_true, y_pred))
            print("Classification Report:")
            print(classification_report(y_true, y_pred))
        evaluate_model("Логистическая регрессия", y_test, y_pred_log)
        evaluate_model("SVM", y_test, y_pred_svm)
evaluate_model("Дерево решений", y_test, y_pred_tree)
     Ⅲ Логистическая регрессия
     Accuracy: 0.75
     Precision: 0.69
     Confusion Matrix:
     [[71 16]
     [20 36]]
     Classification Report:
                   precision
                              recall f1-score support
                0
                        0.78
                                  0.82
                                            0.80
                                                        87
                                0.64
                        0.69
                                           0.67
                                                        56
         accuracy
                                            0.75
        macro avg
                     0.74 0.73
                                         0.73
                                                       143
     weighted avg
                      0.75 0.75
                                           0.75
                                                       143
     II SVM
     Accuracy: 0.73
     Precision: 0.67
     Confusion Matrix:
     [[69 18]
     [20 36]]
     Classification Report:
        accuracy
                                            0.76
                      0.74 0.74
                                                     143
        macro avg
                                         0.74
     weighted avg
                                                       143
                      0.76
                                 0.76
                                           0.76
     Output is truncated. View as a scrollable element or open in a text editor. Adjust cell output settings...
                                                                                                                            💠 G
        # Важность признаков
         importances = tree_model.feature_importances_
        features = X.columns
         plt.figure(figsize=(10.
```



Дерево решений

Sex <= 0.5 gini = 0.484 samples = 571 value = [337, 234] class = Not Survived

