

# 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 ИУ5 «Системы обработки информации и управления»

Отчёт по рубежному контролю №2

Вариант 11

Выполнила: Студентка группы ИУ5-65Б Е. И. Мащенко

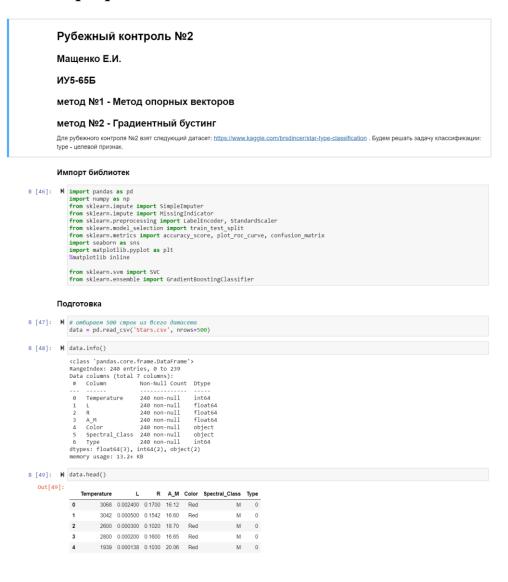
# Задание

Для заданного набора данных (по Вашему варианту) постройте модели классификации или регрессии (в зависимости от конкретной задачи, рассматриваемой в наборе данных). Для построения моделей используйте методы 1 и 2 (по варианту для Вашей группы). Оцените качество моделей на основе подходящих метрик качества (не менее двух метрик). Какие метрики качества Вы использовали и почему? Какие выводы Вы можете сделать о качестве построенных моделей? Для построения моделей необходимо выполнить требуемую предобработку данных: заполнение пропусков, кодирование категориальных признаков, и т.д.

Набор данных: https://www.kaggle.com/brsdincer/star-type-classification.

Метод №1: Метод опорных векторов. Метод №2: Градиентный бустинг.

# Текст программы



```
B [50]: № # Оцениваем баланс классов целевого признака data['Type'].value_counts()/data['Type'].shape[0]*100
    Out[50]: 5
                          16.666667
                          16.666667
                          16.666667
                          16.666667
                          16.666667
                          16,666667
                   Name: Type, dtype: float64
В [51]: 🔰 # Проверяем процент пропусков в данных для всех колонок
                  (data.isnull().sum()/data.shape[0]*100).sort_values(ascending=False)
    Out[51]: Type
                  Spectral_Class
Color
                                              0.0
                                              0.0
                   АМ
                                              0.0
                                              0.0
                                              0.0
                   Temperature
                                              0.0
                   dtype: float64
B [52]: М # Проберяем категориальные признаки на уникальность col_obj = data.dtypes[data.dtypes==object].index.values.tolist()
                   for i in enumerate(col_obj):
    uniq_obj = data[i[1]].unique()
    print(f'{i[0]+1}. {i[1]}: {uniq_obj} | KON-BO: {len(uniq_obj)}')

    Color: ['Red' 'Blue White' 'White' 'Yellowish White' 'Blue white' 'Pale yellow orange' 'Blue' 'Blue-white' 'Whitish' 'Yellow-white' 'Orange' 'White-Yellow' 'white' 'yellowish' 'Yellowish' 'Orange-Red'

                   'slue-white'] | KOM-BO: 17
2. Spectral_class: ['M' 'B' 'A' 'F' 'O' 'K' 'G'] | KOM-BO: 7
В [53]: 🔰 # Оцениваем важность признаков для целевого
                  dataLE = data.copy()
                   le = LabelEncoder(
                  le = LabelEncoder()
col_obj = dataLE.dtypes[dataLE.dtypes==object].index.values.tolist()
for i in col_obj:
    dataLE[i] = le.fit_transform(dataLE[i])
(dataLE.corr()['Type']*100).sort_values(ascending=False)
    Out[53]: Type
                                              100.000000
                                                67.684495
                                                66.097527
                   Temperature
                                               41.112879
                   Spectral_Class
                                                -4.913076
                   Color
                                              -30.499279
                                              -95.527558
                   Name: Type, dtype: float64
В [54]: № #По результатам корреляционного анализа удаляем столбцы, которые имеют меньшую значимость по отношению к целевому признаку del_data = (datalE.corr()['Type']*100).sort_values(ascending=False) del_col = del_data[(del_data < 10) & (del_data > -10) | (del_data.isnull())].index.values.tolist() data.drop(columns=del_col, inplace=True)
                  dataLE.drop(columns=del_col, inplace=True)
B [55]: ► data.info()
                   <class 'pandas.core.frame.DataFrame'>
                  RangeIndex: 240 entries, 0 to 239
Data columns (total 6 columns):
                                              Non-Null Count Dtype
                          Temperature 240 non-null
                                              240 non-null
                                                                       float64
                                              240 non-null
                    3 A M
                                              240 non-null
                                                                      float64
                          Color
                                              240 non-null
                                                                      object
                   5 Type 240 non-null int64 dtypes: float64(3), int64(2), object(1) memory usage: 11.4+ KB
```

#### Метод опорных векторов

```
В [61]: М ргint_metrics(data_X_train, data_y_train, data_x_test, data_y_test, SVC())

Точность: 0.9583333333333334
Матрица ошибок:
[[10 0 0 0 0 0]
[ 3 10 0 0 0 0 0]
[ 0 0 8 0 0 0]
[ 0 0 0 8 0 0 0]
[ 0 0 0 16 0 0]
[ 0 0 0 0 13 0]
[ 0 0 0 0 0 12]]
```

### Градиентный бустинг

```
B [62]: И рrint_metrics(data_X_train, data_y_train, data_X_test, data_y_test, GradientBoostingClassifier(random_state=RANDOM_STATE))

Точность: 1.0
Матрица ошибок:
[[10 0 0 0 0 0 0]
[ 0 13 0 0 0 0 0]
[ 0 0 8 0 0 0]
[ 0 0 0 16 0 0]
[ 0 0 0 0 16 0 0]
[ 0 0 0 0 0 13 0]
[ 0 0 0 0 0 0 12]]
```

## Вывод

В данной работе для оценки моделей были использованы следующие метрики, подходящие для задачи классификации:

1) accuracy 2) confusion matrix

По результатам оценивания можно сделать следующий вывод: модель GradientBoostingClassifier по обеим метрикам обладает немного большей предсказательной способностью, чем Support Vector Machine.