

МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ  
УНИВЕРСИТЕТ  
им. Н.Э. Баумана

---

Факультет «Информатика и система управления»  
Кафедра ИУ-5 «Системы обработки информации и  
управления»

ОТЧЕТ

**Лабораторная работа №5 по курсу  
«Методы машинного обучения»**

«Изучение линейных моделей, SVM и деревьев решений»

Исполнитель - студент группы ИУ5-21М:

Кауров Максим \_\_\_\_\_

Москва – 2020 год



## Лабораторная работа №5

### Цель работы: Изучить линейные модели, SVM и деревья решений.

#### Задание

Требуется выполнить следующие действия:

- Взглянуть на исходные данные (dataset) для решения задачи классификации или регрессии.
- В случае необходимости провести операции или заполнить пропуски в исходные категориальные признаки.
- С использованием метода `train_test_split` разделить выборку на обучающую и тестовую.
- Применить один из линейных моделей, SVM или дерево решений. Сравните качество модели с помощью тестов подбора для задачи метрик.
- Сравните качество полученных моделей.
- Примените для каждой модели модор одного гиперпараметра с использованием `GridSearchCV` и кросс-валидации.
- Получите метрику для наиболее оптимальных значений гиперпараметров. Сравните качество полученных моделей с качеством моделей, полученных в пункте 4.

Подключаем библиотеки:

```
In [1]: from google.colab import files
from sklearn import datasets
import graphviz
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import mean_absolute_error
from sklearn.metrics import median_absolute_error, r2_score
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import ShuffleSplit
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeRegressor
from sklearn.tree import export_graphviz, plot_tree

# Enable inline plots
%matplotlib inline

# Set plots formats to save high resolution PNG
from IPython.display import set_matplotlib_formats
set_matplotlib_formats('retina')
```

In [3]:

```
dataset = datasets.load_svmlight_file('data/train')
X, y = dataset[0], dataset[1]

# Choose File No files chosen

Upload widget is only available when the cell has been executed in the current browser session. Please return this cell to enable.

Saving datasets_418778_1848669_c8b20.csv to file explorer
Saving datasets_418778_1848669_c8b20.csv to file explorer
```

#### Предварительная подготовка данных

### Информация о датасете

В качестве датасета возьмем датасет с информацией о играх баскетбольных команд колледжей.

Данные о содержании (из источника, Kaggle)

RK (Only in c8b20): The ranking of the team at the end of the regular season according to bartonovk

TEAM: The Division I college basketball school

CONF: The Athletic Conference in which the school participates in (A10 = Atlantic 10, ACC = Atlantic Coast Conference, AE = America East, Amer = American, ASun = ASUN, Big10 = Big Ten, Big12 = Big 12, BigE = Big East, BigS = Big Sky, BSW = Big West, CAA = Colonial Athletic Association, CUSA = Conference USA, Horz = Horizon League, Ivy = Ivy League, MAC = Metro Atlantic Athletic Conference, MEAC = Mid-American Conference, MEAC = Mid-Eastern Athletic Conference, MVC = Missouri Valley Conference, MWC = Mountain West, NEC = Northeast Conference, OVC = Ohio Valley Conference, Pac12 = Pac-12, Pac = Pacific League, SB = Sun Belt, SC = Southern Conference, SEC = South Eastern Conference, Sun = Southern Conference, Sum = Summit League, SWAC = Southwestern Athletic Conference, WAC = Western Athletic Conference, WCC = West Coast Conference)

G: Number of games played

W: Number of games won

ADJOE: Adjusted Offensive Efficiency (An estimate of the offensive efficiency (points scored per 100 possessions) a team would have against the average Division I defense)

ADJDE: Adjusted Defensive Efficiency (An estimate of the defensive efficiency (points allowed per 100 possessions) a team would have against the average Division I offense)

BARTNAG: Power Rating (Chance of beating an average Division I team)

EFQ: O: Effective Field Goal Percentage Shot

EFQ: D: Effective Field Goal Percentage Allowed

TOR: Turnover Percentage Allowed (Steal Rate)

TORD: Turnover Percentage Committed (Steal Rate)

ORB: Offensive Rebound Percentage

DRB: Defensive Rebound Percentage

FTFR: Free Throw Rate (FTR when the given team shoots Free Throws)

FTDR: Free Throw Rate Allowed

2P\_D: Two-Point Shooting Percentage

2P\_D: Two-Point Shooting Percentage Allowed

3P\_D: Three-Point Shooting Percentage

3P\_D: Three-Point Shooting Percentage Allowed

ADJ\_T: Adjusted Tempo (An estimate of the tempo (possessions per 40 minutes) a team would have against the team that wants to play at an average Division I tempo)

WAB: Wins Above Bubble (The bubble refers to the cut-off between making the NCAA March Madness Tournament and not making it)

POSTSEASON: Round where the given team was eliminated or where their season ended (R6B = First Four, R64 = Round of 64, R32 = Round of 32, S16 = Sweet Sixteen, E8 = Elite Eight, F4 = Final Four, 2ND = Runner-up, Champion = Winner of the NCAA March Madness Tournament for that given year)

SEED: Seed in the NCAA March Madness Tournament

YEAR: Season

In [6]: data = pd.read\_csv("datasets\_418778\_1848669\_c8b20.csv")

Проверим типы данных в датасете:

In [6]:

data.dtypes

```
RK      int64
TEAM    object
CONF    object
G        int64
W        int64
ADJOE    float64
ADJDE    float64
BARTNAG  float64
EFQ_D    float64
EFQ_O    float64
TOR       int64
TORD      int64
ORB       int64
DRB       int64
FTFR      int64
FTDR      int64
2P_D      int64
2P_D      int64
3P_D      int64
3P_D      int64
ADJ_T     float64
WAB       int64
dtype: object
```

Посмотрим как это выглядит:

In [7]:

data.head()

```
Out[7]:
```

	RK	TEAM	CONF	G	W	ADJOE	ADJDE	BARTNAG	EFQ_O	EFQ_D	TOR	TORD	ORB	DRB	FTFR	FTDR	2P_D	2P_D	3P_D	3P_D	ADJ_T	WAB
0	1	Kansas	B12	30	28	116.1	87.7	0.9636	53.7	45.7	18.7	18.6	32.6	26.4	38.8	13.2	54.9	42.4	34.1	36.5	67.4	106
1	2	Butler	B12	30	26	114.5	88.4	0.9613	49.4	45.2	17.8	23.7	35.8	29.8	30.8	30.8	47.0	44.4	36.1	31.1	66.2	81
2	3	Georgia	WCC	33	31	121.3	94.3	0.9472	57.5	47.6	15.3	18.4	33.6	22.7	38.8	21.8	57.4	47.4	38.6	32.0	72.0	71
3	4	Duquesne	AD	31	29	115.5	93.4	0.9445	59.7	46.6	18.0	18.9	26.4	26.6	33.9	30.9	62.3	45.1	37.1	33.0	67.5	66
4	5	Marquette	SE	33	31	114.8	91.5	0.9336	52.6	45.5	18.1	15.9	32.6	26.0	39.8	29.3	52.9	43.4	34.8	28.7	69.3	52

Нам не очень интересны столбцы TEAM и CONF, так что избавимся от них

In [6]:

```
data.drop(['TEAM', 'CONF'], axis='columns', inplace=True)
```

In [11]:

data.dtypes

```
Out[11]:
```

	RK	G	W	ADJOE	ADJDE	BARTNAG	EFQ_O	EFQ_D	TOR	TORD	ORB	DRB
0	1	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
1	2	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
2	3	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
3	4	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
4	5	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
5	6	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
6	7	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
7	8	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
8	9	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
9	10	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
10	11	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
11	12	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
12	13	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
13	14	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
14	15	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
15	16	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
16	17	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
17	18	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
18	19	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
19	20	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
20	21	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
21	22	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
22	23	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
23	24	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
24	25	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
25	26	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
26	27	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
27	28	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
28	29	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
29	30	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
30	31	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
31	32	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
32	33	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
33	34	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
34	35	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
35	36	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
36	37	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
37	38	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
38	39	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
39	40	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
40	41	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
41	42	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
42	43	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
43	44	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
44	45	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
45	46	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
46	47	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
47	48	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
48	49	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
49	50	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
50	51	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
51	52	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
52	53	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
53	54	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
54	55	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
55	56	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
56	57	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
57	58	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
58	59	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
59	60	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
60	61	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
61	62	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
62	63	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
63	64	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
64	65	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
65	66	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
66	67	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
67	68	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
68	69	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
69	70	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
70	71	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
71	72	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
72	73	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
73	74	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
74	75	int64	int64	float64	float64	float64	float64	float64	int64	int64	int64	int64
75</												