

ИУ5-62Б Васильченко Д.Д.

Рубежный контроль №1 (вариант 6)

Задание

Для заданного набора данных проведите корреляционный анализ. В случае наличия пропусков в данных удалите строки или колонки, содержащие пропуски. Сделайте выводы о возможности построения моделей машинного обучения и о возможном вкладе признаков в модель.

Набор данных

<https://www.kaggle.com/mohansacharya/graduate-admissions> (файл Admission_Predict.csv)

Дополнительное требование

Для произвольной колонки данных построить гистограмму.

Решение

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: data = pd.read_csv("Admission_Predict.csv")
data.head()
```

```
Out[2]:
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
0	1	337	118	4	4.5	4.5	9.65	1	0.92
1	2	324	107	4	4.0	4.5	8.87	1	0.76
2	3	316	104	3	3.0	3.5	8.00	1	0.72
3	4	322	110	3	3.5	2.5	8.67	1	0.80
4	5	314	103	2	2.0	3.0	8.21	0	0.65

```
In [3]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype  
---  --
0   Serial No.            400 non-null   int64   
1   GRE Score              400 non-null   int64   
2   TOEFL Score            400 non-null   int64   
3   University Rating      400 non-null   int64   
4   SOP                    400 non-null   float64  
5   LOR                    400 non-null   float64  
6   CGPA                   400 non-null   float64  
7   Research               400 non-null   int64   
8   Chance of Admit        400 non-null   float64  
dtypes: float64(4), int64(5)
memory usage: 28.2 KB
```

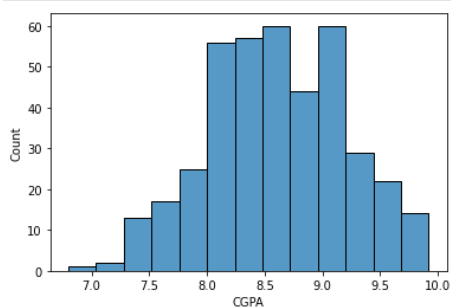
```
In [4]: print('Количество пропущенных значений')
data.isnull().sum()
```

```
Out[4]:
```

Количество пропущенных значений
Serial No. 0
GRE Score 0
TOEFL Score 0
University Rating 0
SOP 0
LOR 0
CGPA 0
Research 0
Chance of Admit 0
dtype: int64

Пропуски в данных не обнаружены.

```
In [5]: sns.histplot(data['CGPA']);
```



Корреляционный анализ

In [6]:

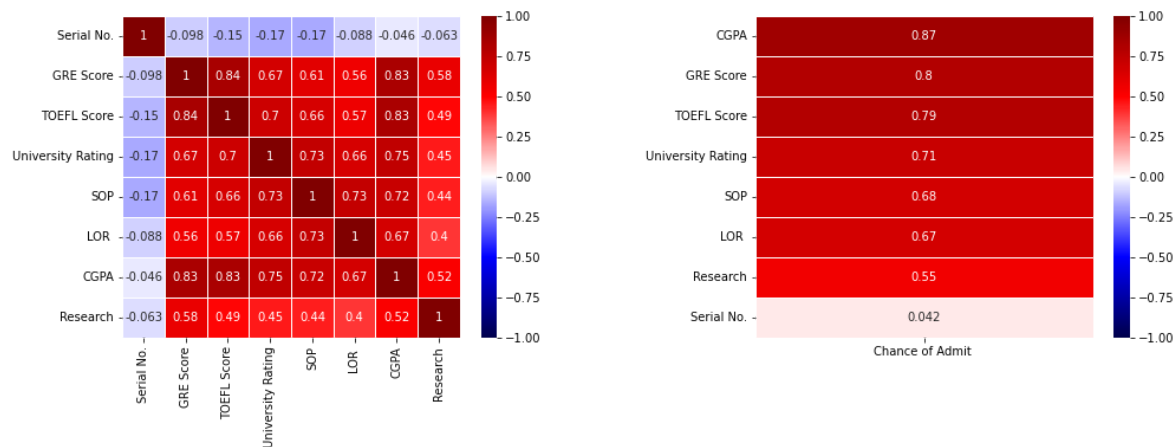
```
data.corr()
```

Out[6]:

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research	Chance of Admit
Serial No.	1.000000	-0.097526	-0.147932	-0.169948	-0.166932	-0.088221	-0.045608	-0.063138	0.042336
GRE Score	-0.097526	1.000000	0.835977	0.668976	0.612831	0.557555	0.833060	0.580391	0.802610
TOEFL Score	-0.147932	0.835977	1.000000	0.695590	0.657981	0.567721	0.828417	0.489858	0.791594
University Rating	-0.169948	0.668976	0.695590	1.000000	0.734523	0.660123	0.746479	0.447783	0.711250
SOP	-0.166932	0.612831	0.657981	0.734523	1.000000	0.729593	0.718144	0.444029	0.675732
LOR	-0.088221	0.557555	0.567721	0.660123	0.729593	1.000000	0.670211	0.396859	0.669889
CGPA	-0.045608	0.833060	0.828417	0.746479	0.718144	0.670211	1.000000	0.521654	0.873289
Research	-0.063138	0.580391	0.489858	0.447783	0.444029	0.396859	0.521654	1.000000	0.553202
Chance of Admit	0.042336	0.802610	0.791594	0.711250	0.675732	0.669889	0.873289	0.553202	1.000000

In [7]:

```
_, axes = plt.subplots(1, 2, figsize=(16, 5))
sns.heatmap(data.drop('Chance of Admit ', axis=1).corr(), annot=True, vmin=-1, vmax=1, cmap='seismic', linewidth=1, ax=axes[0])
sns.heatmap(pd.DataFrame(data.corr()['Chance of Admit ']).sort_values(ascending=False)[1:]),
            annot=True, vmin=-1, vmax=1, cmap='seismic', linewidth=1, ax=axes[1])
plt.subplots_adjust(wspace=0.5)
plt.show()
```



Выше представлены матрица корреляций признаков между собой и матрица корреляции между признаками и прогнозируемой величиной.

Из значений второй матрицы видим, что признак Serial No. не оказывает никакого влияния на прогнозируемую величину Chance of Admit. Также видно, что остальные признаки имеют положительную связь с прогнозируемым.

Из значений первой матрицы видим крайне высокую корреляцию между следующими парами признаков:

- TOEFL Score и GRE Score
- CGPA и GRE Score
- CGPA и TOEFL Score

Так как одновременное использование этих пар признаков в моделях машинного обучения привело бы к мультиколлинеарности, следует оставить только один признак из этого множества. Вторая матрица демонстрирует, что наибольшая связь наблюдается между прогнозируемой величиной и признаком CGPA, поэтому логичнее оставить именно его, так как его вклад в модель обучения будет наибольшим.

Таким образом, в результате корреляционного анализа было принято решение в первую очередь пробовать использовать в моделях машинного обучения для прогноза величины Chance of Admit все признаки, кроме Serial No., GRE Score, TOEFL Score.