

**Московский государственный технический  
университет им. Н. Э. Баумана**

Факультет «Радиотехнический»  
Кафедра ИУ5 «Системы обработки информации и управления»

Курс «Технологии машинного обучения»

Отчет по лабораторной работе №5  
«Ансамбли моделей машинного обучения»

Выполнил:

студент группы РТ5-61Б

Проверил:

доцент каф. ИУ5

Алиев Тимур

Подпись и дата:

Гапанюк Ю.Е.

Подпись и дата:

Москва, 2022 г.

## Описание задания

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 ансамблевые модели:
  - одну из моделей группы бэггинга (бэггинг или случайный лес или сверхслучайные деревья);
  - одну из моделей группы бустинга;
  - одну из моделей группы стекинга.

# Jupyter notebook

## Лабораторная работа №5. Алиев Тимур РТ5-61Б.

### Подготовка данных

```
Ввод [1]: 1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 import numpy as np
5 from sklearn.preprocessing import PolynomialFeatures, MinMaxScaler, StandardScaler
6 from sklearn.linear_model import LinearRegression, Lasso, Ridge
7 from sklearn.tree import DecisionTreeRegressor, export_graphviz, export_text
8 from sklearn.svm import SVR
9 from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
10 from sklearn.model_selection import train_test_split, GridSearchCV
11 from IPython.display import Image
12 from IPython.core.display import HTML
13 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
14 from heamy.dataset import Dataset
15 from heamy.estimator import Regressor
16 from heamy.pipeline import ModelsPipeline
17 from sklearn.neural_network import MLPRegressor
18 from gmdhpy import gmdh
19 from warnings import simplefilter
20
21 simplefilter('ignore')
```

```
Ввод [2]: 1 data = pd.read_csv('WineQT.csv')
2 data.head()
```

```
Out[2]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality	Id
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5	0
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5	1
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5	2
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6	3
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5	4

```
Ввод [3]: 1 # Корреляция с алкоголем по модулю - top
2 best_params = data.corr()['alcohol'].map(abs).sort_values(ascending=False)[1:]
3 best_params = best_params[best_params.values > 0.35]
4 best_params
```

```
Out[3]: density      0.494727
quality    0.484866
Name: alcohol, dtype: float64
```

### Разделение выборки на обучающую и тестовую

```
Ввод [4]: 1 x_train, x_test, y_train, y_test = train_test_split(data[best_params.index], data['alcohol'], test_size=0.3, random_state=3)
2
```

### Масштабирование данных

```
Ввод [5]: 1 scaler = StandardScaler().fit(x_train)
2 x_train_scaled = pd.DataFrame(scaler.transform(x_train), columns=x_train.columns)
3 x_test_scaled = pd.DataFrame(scaler.transform(x_test), columns=x_train.columns)
```

```
Ввод [6]: 1 # функция для вывода метрики
2 def print_metrics(y_test, y_pred):
3     print(f"R^2: {r2_score(y_test, y_pred)}")
4     print(f"MSE: {mean_squared_error(y_test, y_pred)}")
5     print(f"MAE: {mean_absolute_error(y_test, y_pred)}")
```

## 1) Случайный лес

Ввод [7]: 1 print\_metrics(y\_test, RandomForestRegressor(random\_state=17).fit(x\_train, y\_train).predict(x\_test))

R<sup>2</sup>: 0.436356891520067  
MSE: 0.6606961915609506  
MAE: 0.6317937744229445

Ввод [8]: 1 # Подбор гиперпараметров  
2 rf = RandomForestRegressor(random\_state=17)  
3 params = {'n\_estimators': [100, 1000], 'criterion': ['squared\_error', 'absolute\_error', 'poisson'],  
4 'max\_features': ['auto', 'sqrt'], 'min\_samples\_leaf': [1, 3, 5]}  
5 grid\_cv = GridSearchCV(estimator=rf, cv=5, param\_grid=params, n\_jobs=-1, scoring='neg\_mean\_absolute\_error')  
6 grid\_cv.fit(x\_train, y\_train)  
7 print(grid\_cv.best\_params\_)  
  
{'criterion': 'absolute\_error', 'max\_features': 'sqrt', 'min\_samples\_leaf': 5, 'n\_estimators': 1000}

Ввод [9]: 1 best\_rf = grid\_cv.best\_estimator\_  
2 best\_rf.fit(x\_train, y\_train)  
3 y\_pred\_rf = best\_rf.predict(x\_test)  
4 print\_metrics(y\_test, y\_pred\_rf)

R<sup>2</sup>: 0.4900907703255626  
MSE: 0.5977099356296611  
MAE: 0.601562973760933

## 2) Градиентный бустинг

Ввод [10]: 1 print\_metrics(y\_test, GradientBoostingRegressor(random\_state=17).fit(x\_train, y\_train).predict(x\_test))

R<sup>2</sup>: 0.5187256098228614  
MSE: 0.564144494808706  
MAE: 0.592595580465614

Ввод [11]: 1 # Подбор гиперпараметров  
2 gb = GradientBoostingRegressor(random\_state=17)  
3 params = {'loss': ['squared\_error', 'absolute\_error', 'huber'], 'n\_estimators': [10, 50, 100, 200],  
4 'criterion': ['friedman\_mse', 'squared\_error', 'mse', 'mae'], 'min\_samples\_leaf': [1, 3, 5]}  
5 grid\_cv = GridSearchCV(estimator=gb, cv=5, param\_grid=params, n\_jobs=-1, scoring='r2')  
6 grid\_cv.fit(x\_train, y\_train)  
7 print(grid\_cv.best\_params\_)  
  
{'criterion': 'friedman\_mse', 'loss': 'squared\_error', 'min\_samples\_leaf': 1, 'n\_estimators': 50}

Ввод [12]: 1 best\_gb = grid\_cv.best\_estimator\_  
2 best\_gb.fit(x\_train, y\_train)  
3 y\_pred\_gb = best\_gb.predict(x\_test)  
4 print\_metrics(y\_test, y\_pred\_gb)

R<sup>2</sup>: 0.5265950960016097  
MSE: 0.5549199705137814  
MAE: 0.589858793540792

## 3) Стекинг

Ввод [13]: 1 dataset = Dataset(x\_train, y\_train, x\_test)

Ввод [14]: 1 model\_lr = Regressor(dataset=dataset, estimator=LinearRegression, name='lr')  
2 model\_rf = Regressor(dataset=dataset, estimator=RandomForestRegressor,  
3 parameters={'criterion': 'absolute\_error', 'n\_estimators': 1000, 'random\_state': 17}, name='rf')  
4 model\_gb = Regressor(dataset=dataset, estimator=GradientBoostingRegressor,  
5 parameters={'loss': 'huber', 'random\_state': 17}, name='gb')

Ввод [15]: 1 pipeline = ModelsPipeline(model\_lr, model\_rf)  
2 stack\_ds = pipeline.stack(k=10, seed=1)  
3 stacker = Regressor(dataset=stack\_ds, estimator=GradientBoostingRegressor)  
4 results = stacker.validate(k=10, scorer=mean\_absolute\_error)

Metric: mean\_absolute\_error  
Folds accuracy: [0.6194110924984411, 0.661961907062875, 0.6253281354764146, 0.647977261569514, 0.46334402329637425, 0.5428994509989188, 0.5361135714228367, 0.47879256948668053, 0.4813988120354393, 0.634546756726363]  
Mean accuracy: 0.5691773580573857  
Standard Deviation: 0.07323033099064555  
Variance: 0.005362681376999502

Ввод [16]: 1 y\_pred\_stack = stacker.predict()  
2 print\_metrics(y\_test, y\_pred\_stack)

R<sup>2</sup>: 0.4654915465854482  
MSE: 0.6265448724823064  
MAE: 0.6213646378838801

## 4) Многослойный перцептрон ¶

Ввод [17]: 

```
1 print_metrics(y_test, MLPRegressor(random_state=17).fit(x_train, y_train).predict(x_test))
```

R<sup>2</sup>: 0.30581687622485887  
MSE: 0.8137137476247516  
MAE: 0.7275229510423318

Ввод [18]: 

```
1 # Подбор гиперпараметров
2 mlp = MLPRegressor(random_state=17)
3 params = {'solver': ['lbfgs', 'sgd', 'adam'], 'hidden_layer_sizes': [(100,), (50, 30,), (100, 40,)],
4           'alpha': [1e-4, 3e-4, 5e-4], 'max_iter': [500, 1000]}
5 grid_cv = GridSearchCV(estimator=mlp, cv=5, param_grid=params, n_jobs=-1, scoring='r2')
6 grid_cv.fit(x_train, y_train)
7 print(grid_cv.best_params_)

{'alpha': 0.0005, 'hidden_layer_sizes': (100,), 'max_iter': 500, 'solver': 'lbfgs'}
```

Ввод [19]: 

```
1 best_mlp = grid_cv.best_estimator_
2 best_mlp.fit(x_train, y_train)
3 y_pred_mlp = best_mlp.predict(x_test)
4 print_metrics(y_test, y_pred_mlp)
```

R<sup>2</sup>: 0.4791641439073404  
MSE: 0.6105180057586421  
MAE: 0.6211604340563389

## 5) Метод группового учёта аргументов

Ввод [20]: 

```
1 gm = gmdh.Regressor(n_jobs=-1)
2 gm.fit(np.array(x_train_scaled), np.array(y_train))
3 y_pred_gm = gm.predict(np.array(x_test_scaled))
4 print()
5 print_metrics(y_test, y_pred_gm)
```

train layer0 in 0.00 sec  
train layer1 in 0.00 sec  
train layer2 in 0.00 sec  
train layer3 in 0.00 sec  
train layer4 in 0.00 sec  
train layer5 in 0.00 sec  
train layer6 in 0.00 sec

R<sup>2</sup>: 0.4630314373942027  
MSE: 0.6294285851526578  
MAE: 0.6345983648345276

## Сравнение моделей

Ввод [21]: 

```
1 print("Случайный лес")
2 print_metrics(y_test, y_pred_rf)
3
4 print("\nГрадиентный бустинг")
5 print_metrics(y_test, y_pred_gb)
6
7 print("\nСтекинг")
8 print_metrics(y_test, y_pred_stack)
9
10 print("\nМногослойный перцептрон")
11 print_metrics(y_test, y_pred_mlp)
12
13 print("\nМетод группового учёта аргументов")
14 print_metrics(y_test, y_pred_gm)
```

Случайный лес  
R<sup>2</sup>: 0.4900907703255626  
MSE: 0.5977099356296611  
MAE: 0.601562973760933

Градиентный бустинг  
R<sup>2</sup>: 0.5265950960016097  
MSE: 0.5549199705137814  
MAE: 0.589858793540792

Стекинг  
R<sup>2</sup>: 0.4654915465854482  
MSE: 0.6265448724823064  
MAE: 0.6213646378838801

Многослойный перцептрон  
R<sup>2</sup>: 0.4791641439073404  
MSE: 0.6105180057586421  
MAE: 0.6211604340563389

Метод группового учёта аргументов  
R<sup>2</sup>: 0.4630314373942027  
MSE: 0.6294285851526578  
MAE: 0.6345983648345276