

## ▼ Ансамбли моделей машинного обучения

Цель лабораторной работы: изучение ансамблей моделей машинного обучения. Задание:

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 две ансамблевые модели. Оцените качество моделей с помощью одной из подходящих для задачи метрик. Сравните качество полученных моделей.
5. Произведите для каждой модели подбор значений одного гиперпараметра. В зависимости от используемой библиотеки можно применять функцию `GridSearchCV`, использовать перебор параметров в цикле, или использовать другие методы.
6. Повторите пункт 4 для найденных оптимальных значений гиперпараметров. Сравните качество полученных моделей с качеством моделей, полученных в пункте 4.

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 sns.set(style="ticks")
7 data = pd.read_csv('Data/lab_4/heart.csv', sep=",")
8 data.head(5)
```



	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope
0	63	1	3	145	233	1	0	150	0	2.3	0
1	37	1	2	130	250	0	1	187	0	3.5	0
2	41	0	1	130	204	0	0	172	0	1.4	2
3	56	1	1	120	236	0	1	178	0	0.8	2
4	57	0	0	120	354	0	1	163	1	0.6	2

```
1 data.shape
```



(303, 14)

```
1 # Проверка на пустые значения
2 data.isnull().sum()
```



```

age          0
sex          0
cp          0
trestbps    0
chol        0
fbs         0


```

## ▼ Feature Scaling

```

1 from sklearn.preprocessing import MinMaxScaler
2 import warnings
3 warnings.filterwarnings('ignore')
4
5 # Create the scaler object with a range of 0-1
6 scaler = MinMaxScaler(feature_range=(0, 1))
7 # Fit on data, transform data
8 scaler.fit_transform(data)
9

```

 array([[0.70833333, 1. , 1. , ..., 0. , 0.33333333,
 1. ],
 [0.16666667, 1. , 0.66666667, ..., 0. , 0.66666667,
 1. ],
 [0.25 , 0. , 0.33333333, ..., 0. , 0.66666667,
 1. ],
 ...,
 [0.8125 , 1. , 0. , ..., 0.5 , 1. ,
 0. ],
 [0.58333333, 1. , 0. , ..., 0.25 , 1. ,
 0. ],
 [0.58333333, 0. , 0.33333333, ..., 0.25 , 0.66666667,
 0. ]])

```

1 from sklearn import svm
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import roc_curve, auc
6 import pylab as pl

```

```

1 # Пустых значений нет
2 # Перейдем к разделению выборки на обучающую и тестовую.
3 X = data.drop('target',axis = 1).values
4 y = data['target'].values

```

## ▼ Ансамблевые модели

```

1 from sklearn.model_selection import train_test_split
2 kfold = 5 #количество подвыборок для валидации
3 itog_val = {} #список для записи результатов кросс валидации разных алгоритмов

```

```

1 ROctrainTRN, ROCTestTRN, ROctrainTRG, ROCTestTRG = train_test_split(X, y, test_

```

```

1 model_rfc = RandomForestClassifier(n_estimators = 70) #в параметре передаем кол-во д
2 model_knc = KNeighborsClassifier(n_neighbors = 18) #в параметре передаем кол-во соседей
3 model_lr = LogisticRegression(penalty='l1', tol=0.01)
4 model_svc = svm.SVC() #по умолчанию kernel='rbf'

```

```
1 from sklearn.model_selection import cross_val_score
```

SVM - метод опорных векторов(SVC): Суть работы "Машин" Опорных Векторов проста: алгоритм создает линию или гиперплоскость, которая разделяет данные на классы. Метод k-ближайших соседей(KNeighborsClassifier) Random forest(RandomForestClassifier) Логистическая регрессия (LogisticRegression)

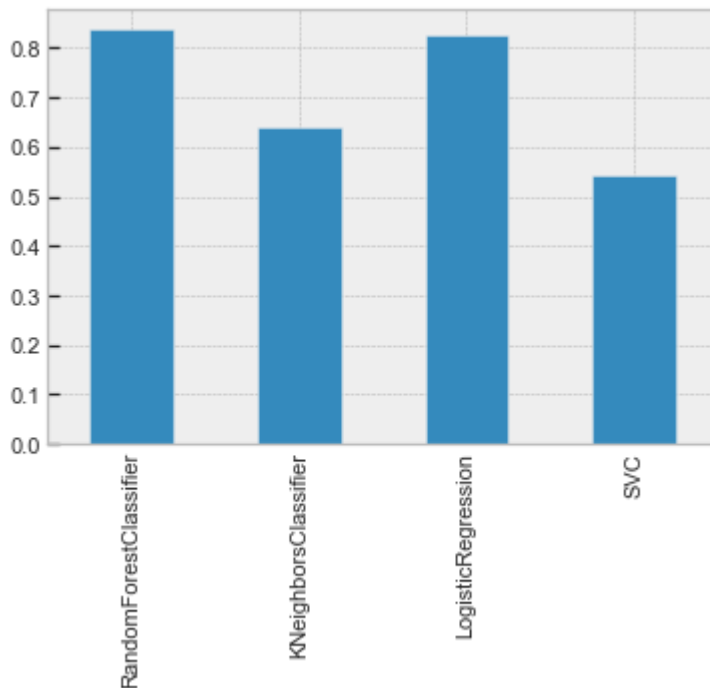
```
1 from sklearn.ensemble import ExtraTreesClassifier, ExtraTreesRegressor
```

```
1 scores = cross_val_score(model_rfc, X, y, cv = kfold)
2 itog_val['RandomForestClassifier'] = scores.mean()
3 scores = cross_val_score(model_knc, X, y, cv = kfold)
4 itog_val['KNeighborsClassifier'] = scores.mean()
5 scores = cross_val_score(model_lr, X, y, cv = kfold)
6 itog_val['LogisticRegression'] = scores.mean()
7 scores = cross_val_score(model_svc, X, y, cv = kfold)
8 itog_val['SVC'] = scores.mean()
```

```
1 import matplotlib.pyplot as plt
2 plt.style.use('bmh')
3 data.from_dict(data = itog_val, orient='index').plot(kind='bar', legend=False)
```



<matplotlib.axes.\_subplots.AxesSubplot at 0x1a1e343240>



```
1 pl.clf()
2 plt.figure(figsize=(8,6))
3 #SVC
4 model_svc.probability = True
5 probas = model_svc.fit(ROctrainTRN, ROctrainTRG).predict_proba(ROctestTRN)
6 fpr, tpr, thresholds = roc_curve(ROctestTRG, probas[:, 1])
7 roc_auc = auc(fpr, tpr)
8 pl.plot(fpr, tpr, label='%s ROC (area = %0.2f)' % ('SVC', roc_auc))
9 #RandomForestClassifier
10 probas = model_rfc.fit(ROctrainTRN, ROctrainTRG).predict_proba(ROctestTRN)
11 fpr, tpr, thresholds = roc_curve(ROctestTRG, probas[:, 1])
12 roc_auc = auc(fpr, tpr)
13 pl.plot(fpr, tpr, label='%s ROC (area = %0.2f)' % ('RandomForest', roc_auc))
14 #KNeighborsClassifier
15 probas = model_knc.fit(ROctrainTRN, ROctrainTRG).predict_proba(ROctestTRN)
16 fpr, tpr, thresholds = roc_curve(ROctestTRG, probas[:, 1])
```

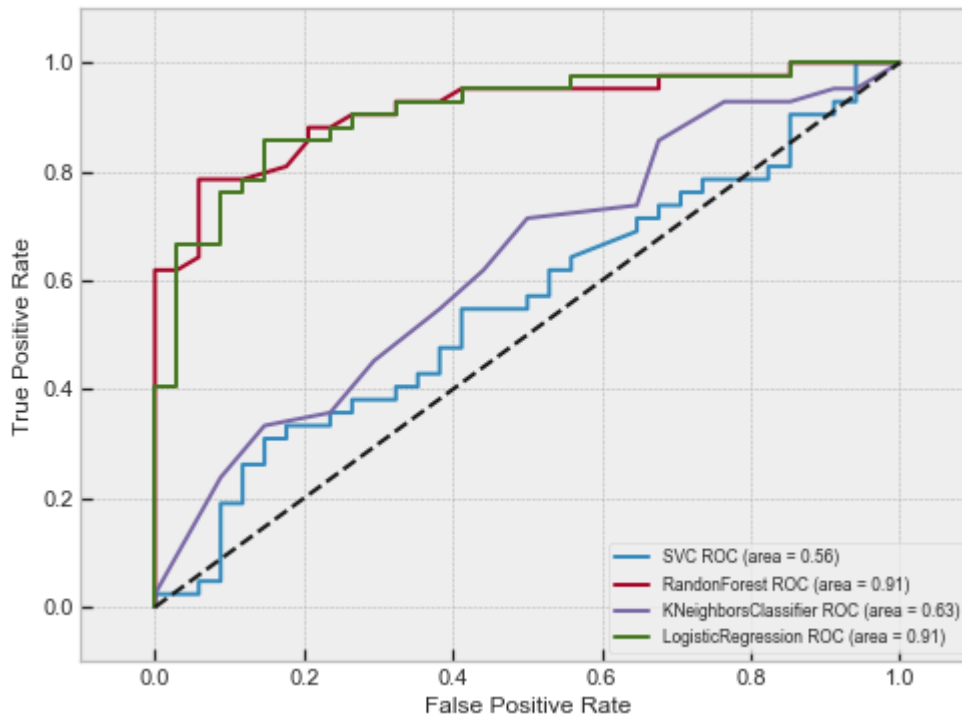
```

17 roc_auc = auc(fpr, tpr)
18 pl.plot(fpr, tpr, label='%s ROC (area = %0.2f)' % ('KNeighborsClassifier',roc_a
19 #LogisticRegression
20 probas = model_lr.fit(ROctrainTRN, ROctrainTRG).predict_proba(ROctestTRN)
21 fpr, tpr, thresholds = roc_curve(ROctestTRG, probas[:, 1])
22 roc_auc = auc(fpr, tpr)
23 pl.plot(fpr, tpr, label='%s ROC (area = %0.2f)' % ('LogisticRegression',roc_auc
24 pl.plot([0, 1], [0, 1], 'k--')
25 pl.xlim([-0.1, 1.1])
26 pl.ylim([-0.1, 1.1])
27 pl.xlabel('False Positive Rate')
28 pl.ylabel('True Positive Rate')
29 pl.legend(loc=0, fontsize='small')
30 pl.show()

```



<Figure size 432x288 with 0 Axes>



```

1 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
2 from sklearn.metrics import accuracy_score
3 from sklearn.metrics import balanced_accuracy_score
4 from sklearn.metrics import precision_score, recall_score, f1_score

```

```

1 from sklearn.model_selection import train_test_split
2 # Функция train_test_split разделила исходную выборку таким образом,
3 #чтобы в обучающей и тестовой частях сохранились пропорции классов.
4 X_train, X_test, y_train, y_test = train_test_split(
5     X, y, test_size=0.35, random_state=1)

```

```

1 from sklearn.preprocessing import MinMaxScaler
2 import warnings
3 warnings.filterwarnings('ignore')
4
5 # Create the scaler object with a range of 0-1
6 scaler = MinMaxScaler(feature_range=(0, 1))
7 # Fit on data, transform data
8 scaler.fit_transform(X)
9 scaler.fit_transform(X_train)
10 scaler.fit_transform(X_test)

```




```
array([[0.77777778, 0.          , 0.          , ..., 0.          , 0.75          ,
        1.          ],
       [0.61111111, 1.          , 0.33333333, ..., 1.          , 0.          ,
        1.          ],
       [0.38888889, 1.          , 0.          , ..., 1.          , 0.5          ,
        1.          ],
       ...,
       [0.52777778, 1.          , 0.          , ..., 0.          , 0.          ,
        1.          ],
       ...])
```

```
1 # n_estimators = 10 (default)
2 rfc = RandomForestClassifier().fit(X_train, y_train)
3 predicted_rfc = rfc.predict(X_test)
```


```
1 accuracy_score(y_test, predicted_rfc)
```

 0.7570093457943925


```
1 balanced_accuracy_score(y_test, predicted_rfc)
```

 0.7547368421052632

```
1 (precision_score(y_test, predicted_rfc, average='weighted'),
2  recall_score(y_test, predicted_rfc, average='weighted')).
```


 (0.7567717408522097, 0.7570093457943925)

```
1 f1_score(y_test, predicted_rfc, average='weighted')
```


 0.7566245963419206

```
1 # n_estimators = 50 (default)
2 abc = AdaBoostClassifier().fit(X_train, y_train)
3 predicted_abc = abc.predict(X_test)
```


```
1 accuracy_score(y_test, predicted_abc)
```

 0.7289719626168224


```
1 balanced_accuracy_score(y_test, predicted_abc)
```

 0.7284210526315789

```
1 (precision_score(y_test, predicted_abc, average='weighted'),
2  recall_score(y_test, predicted_abc, average='weighted')).
```

 (0.7293842770753162, 0.7289719626168224)

```
1 f1_score(y_test, predicted_abc, average='weighted')
```

 0.7291144464706996

```
1 rfc_n_range = np.array(range(5,100,5))
2 rfc_tuned_parameters = [{'n_estimators': rfc_n_range}]
```

```
3 rfc_tuned_parameters
```

```
[{'n_estimators': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 90, 95])}]
```

```
1 import warnings
2 from sklearn.model_selection import GridSearchCV
3 warnings.filterwarnings('ignore')
4
5 gs_rfc = GridSearchCV(RandomForestClassifier(), rfc_tuned_parameters, cv=5,
6                     scoring='accuracy')
7 gs_rfc.fit(X_train, y_train)
```

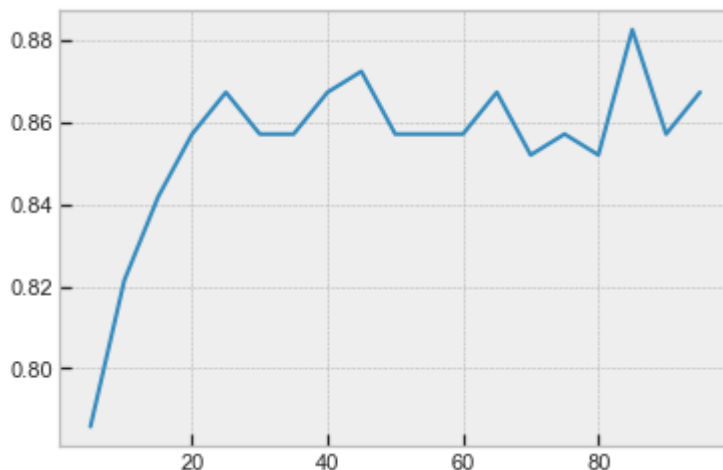
```
[GridSearchCV(cv=5, error_score='raise-deprecating',
              estimator=RandomForestClassifier(bootstrap=True, class_weight=None, c
              max_depth=None, max_features='auto', max_leaf_nodes=None,
              min_impurity_decrease=0.0, min_impurity_split=None,
              min_samples_leaf=1, min_samples_split=2,
              min_weight_fraction_leaf=0.0, n_estimators='warn', n_jobs=None,
              oob_score=False, random_state=None, verbose=0,
              warm_start=False),
              fit_params=None, iid='warn', n_jobs=None,
              param_grid=[{'n_estimators': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45
              90, 95])}],
              pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
              scoring='accuracy', verbose=0)
```

```
1 gs_rfc.best_params_
```

```
{'n_estimators': 85}
```

```
1 plt.plot(rfc_n_range, gs_rfc.cv_results_['mean_test_score'])
```

```
[<matplotlib.lines.Line2D at 0x1ale9302e8>]
```



```
1 abc_n_range = np.array(range(5,100,5))
2 abc_tuned_parameters = [{'n_estimators': abc_n_range}]
3 abc_tuned_parameters
```

```
[{'n_estimators': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 90, 95])}]
```

```
1 gs_abc = GridSearchCV(AdaBoostClassifier(), abc_tuned_parameters, cv=5,
2                     scoring='accuracy')
3 gs_abc.fit(X_train, y_train)
```

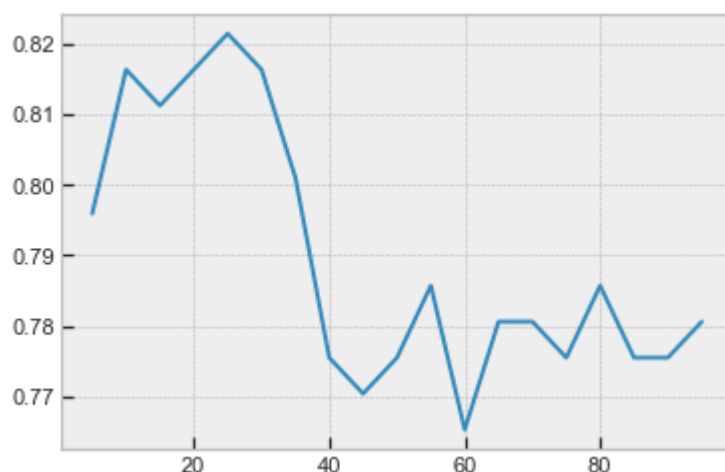
```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=AdaBoostClassifier(algorithm='SAMME.R', base_estimator=None,
             learning_rate=1.0, n_estimators=50, random_state=None),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid=[{'n_estimators': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45,
             50, 55, 60, 65, 70, 75, 80, 85, 90, 95])}],
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='accuracy', verbose=0)
```

```
1 gs_abc.best_params_
```

```
{'n_estimators': 25}
```

```
1 plt.plot(abc_n_range, gs_abc.cv_results_['mean_test_score'])
```

```
[<matplotlib.lines.Line2D at 0x1a1e985d68>]
```



```
1 rfc_optimized = RandomForestClassifier(n_estimators=gs_rfc.best_params_['n_estimators'])
2 predicted_rfc_opt = rfc_optimized.predict(X_test)
```

```
1 accuracy_score(y_test, predicted_rfc_opt)
```

```
0.7663551401869159
```

```
1 balanced_accuracy_score(y_test, predicted_rfc_opt)
```

```
0.7635087719298246
```

```
1 (precision_score(y_test, predicted_rfc_opt, average='weighted'),
2  recall_score(y_test, predicted_rfc_opt, average='weighted')).
```


```
(0.7663352555179958, 0.7663551401869159)
```

```
1 f1_score(y_test, predicted_rfc_opt, average='weighted')
```


```
0.765737522265126
```

```
1 abc_optimized = RandomForestClassifier(n_estimators=gs_abc.best_params_['n_estimators'])
2 predicted_abc_opt = abc_optimized.predict(X_test)
```


```
1 | accuracy_score(y_test, predicted_abc_opt)
```

 0.7476635514018691


```
1 | balanced_accuracy_score(y_test, predicted_abc_opt)
```

 0.7471929824561403

```
1 | .(precision_score(y_test, predicted_abc_opt, average='weighted'),  
2 | recall_score(y_test, predicted_abc_opt, average='weighted')).
```

 (0.748059504175502, 0.7476635514018691)

```
1 | f1_score(y_test, predicted_abc_opt, average='weighted')
```

 0.7477962087830651