

**Московский государственный технический
университет им. Н. Э. Баумана
Факультет «Информатика и системы управления»**

Кафедра «Системы обработки информации и управления»
Курс «Технологии машинного обучения»

Отчет по лабораторной работе №4

Подготовка обучающей и тестовой выборки, кросс-валидация и подбор
гиперпараметров на примере метода ближайших соседей

Группа: ИУ5-62Б

Студент: Селедкина А.С.

Преподаватель: Гапанюк Ю.Е.

Москва, 2020 г.

Цель лабораторной работы: изучение сложных способов подготовки выборки и подбора гиперпараметров на примере метода ближайших соседей.

Описание задания

1. Выберите набор данных (датасет) для решения задачи классификации или регрессии.
2. С использованием метода `train_test_split` разделите выборку на обучающую и тестовую.
3. Обучите модель ближайших соседей для произвольно заданного гиперпараметра `K`. Оцените качество модели с помощью подходящих для задачи метрик.
4. Постройте модель и оцените качество модели с использованием кросс-валидации.
5. Произведите подбор гиперпараметра `K` с использованием `GridSearchCV` и кросс-валидации.

Текст программы и примеры выполнения

Будем использовать датасет по определению наличия сердечного заболевания у пациента: <https://www.kaggle.com/ronitf/heart-disease-uci>.

```
data = pd.read_csv('data/heart.csv')
data.head()
```

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

```
data.shape
```

```
(303, 14)
```

```
data.columns
```

```
Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',  
      'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],  
      dtype='object')
```

```
data.dtypes
```

```
age          int64  
sex          int64  
cp           int64  
trestbps     int64  
chol         int64  
fbs          int64  
restecg      int64  
thalach      int64  
exang        int64  
oldpeak      float64  
slope        int64  
ca           int64  
thal         int64  
target       int64  
dtype: object
```

```
data.isnull().sum()
```

```
age          0  
sex          0  
cp           0  
trestbps     0  
chol         0  
fbs          0  
restecg      0  
thalach      0  
exang        0  
oldpeak      0  
slope        0  
ca           0  
thal         0  
target       0  
dtype: int64
```

Разделение выборки

```
x = data[['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',  
         'exang', 'oldpeak', 'slope', 'ca', 'thal']]  
y = data['target']
```

```
# С использованием метода train_test_split разделим выборку на обучающую и тестовую  
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=1)  
print("x_train:", x_train.shape)  
print("x_test:", x_test.shape)  
print("y_train:", y_train.shape)  
print("y_test:", y_test.shape)
```

```
x_train: (227, 13)  
x_test: (76, 13)  
y_train: (227,)  
y_test: (76,)
```

Обучение модели

```
class_model = KNeighborsClassifier(n_neighbors=5)
```

```
class_model.fit(x_train, y_train)
```

```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',  
                    metric_params=None, n_jobs=None, n_neighbors=5, p=2,  
                    weights='uniform')
```

```
y_predicted = class_model.predict(x_test)
```

Оценка качества модели

```
accuracy_score(y_test, y_predicted)
```

```
0.5921052631578947
```

```
precision_score(y_test, y_predicted)
```

```
0.6086956521739131
```

```
recall_score(y_test, y_predicted)
```

```
0.6829268292682927
```

```
f1_score(y_test, y_predicted)
```

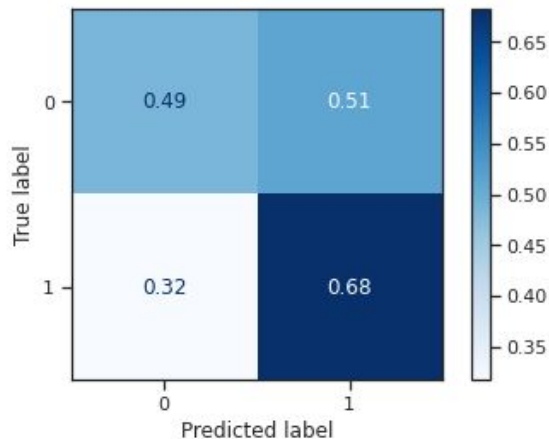
```
0.6436781609195402
```

```
roc_auc_score(y_test, y_predicted)
```

```
0.5843205574912893
```

```
plot_confusion_matrix(class_model, x_test, y_test,
                      display_labels=['0', '1'],
                      cmap=plt.cm.Blues, normalize='true')
```

```
<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7f5d4606eee0>
```



```
scores = cross_val_score(class_model,
                          x_train, y_train, cv=3,
                          scoring='f1_weighted')
scores, np.mean(scores)
```

```
(array([0.57356459, 0.67688787, 0.67542522]), 0.6419592283654437)
```

Подбор гиперпараметра

```
n_range = np.array(range(1,100,2))
tuned_parameters = [{'n_neighbors': n_range}]
tuned_parameters
```

```
[{'n_neighbors': array([ 1,  3,  5,  7,  9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31, 33,
                        35, 37, 39, 41, 43, 45, 47, 49, 51, 53, 55, 57, 59, 61, 63, 65, 67,
                        69, 71, 73, 75, 77, 79, 81, 83, 85, 87, 89, 91, 93, 95, 97, 99])}]
```

```
%%time
clf_gs = GridSearchCV(KNeighborsClassifier(), tuned_parameters, cv=5, scoring='roc_auc')
clf_gs.fit(x_train, y_train)
```

```
CPU times: user 1.79 s, sys: 4.26 ms, total: 1.79 s
Wall time: 1.79 s
```

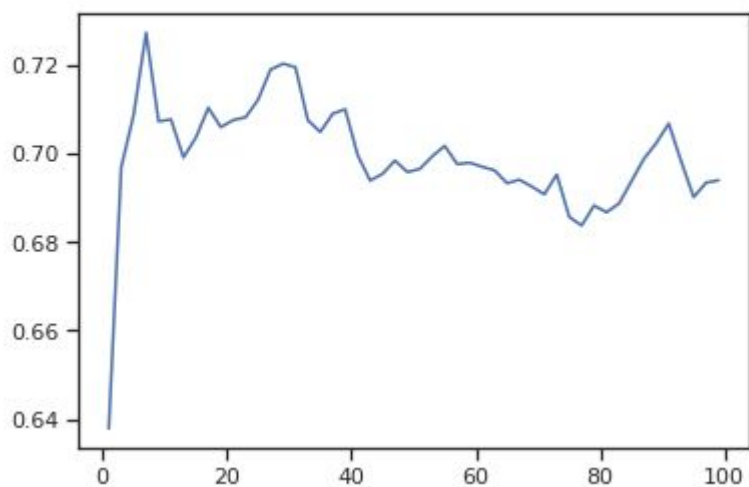
```
GridSearchCV(cv=5, error_score=nan,
             estimator=KNeighborsClassifier(algorithm='auto', leaf_size=30,
                                           metric='minkowski',
                                           metric_params=None, n_jobs=None,
                                           n_neighbors=5, p=2,
                                           weights='uniform'),
             iid='deprecated', n_jobs=None,
             param_grid=[{'n_neighbors': array([ 1,  3,  5,  7,  9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31, 33,
                                                35, 37, 39, 41, 43, 45, 47, 49, 51, 53, 55, 57, 59, 61, 63, 65, 67,
                                                69, 71, 73, 75, 77, 79, 81, 83, 85, 87, 89, 91, 93, 95, 97, 99])}]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
             scoring='roc_auc', verbose=0)
```

```
# Лучшее значение параметра
clf_gs.best_params_
```

```
{'n_neighbors': 7}
```

```
# Изменение качества на тестовой выборке в зависимости от K-соседей  
plt.plot(n_range, clf_gs.cv_results_['mean_test_score'])
```

```
[<matplotlib.lines.Line2D at 0x7f5d438f5c10>]
```



Самым оптимальным значением количества ближайших соседей оказалось 7.