

Лабораторная работа №4  
по дисциплине  
«Методы машинного обучения»  
на тему  
«Подготовка обучающей и тестовой выборки,  
кросс-валидация и подбор гиперпараметров  
на примере метода ближайших соседей.»

Выполнил:  
студент группы ИУ5-23М  
Богомолов Д.Н.

# Подготовка обучающей и тестовой выборки, кросс-вали гиперпараметров на примере метода ближайших сосед

## Цель лабораторной работы:

изучение сложных способов подготовки выборки и подбора гиперпараметров на примере

## Задание:

1. Выберите набор данных (датасет) для решения задачи классификации или регрессии.
2. В случае необходимости проведите удаление или заполнение пропусков и кодирован
3. С использованием метода `train_test_split` разделите выборку на обучающую и тестову
4. Обучите модель ближайших соседей для произвольно заданного гиперпараметра  $K$ .  
подходящих для задачи метрик.
5. Постройте модель и оцените качество модели с использованием кросс-валидации. П  
различными стратегиями кросс-валидации.
6. Произведите подбор гиперпараметра  $K$  с использованием `GridSearchCV` и кросс-вали
7. Повторите пункт 4 для найденного оптимального значения гиперпараметра  $K$ . Сравни  
качеством модели, полученной в пункте 4.
8. Постройте кривые обучения и валидации.

```
1 from google.colab import drive #Монтирование Google Drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call `drive.m`

## Задание 1. Выберите набор данных (датасет) для решения зада регрессии.

Выбранный датасет - HEART.CSV Параметры датасета:

- age - age in years
- sex (1 = male; 0 = female)
- cp - chest pain type
- trestbps - resting blood pressure (in mm Hg on admission to the hospital)
- chol - serum cholestoral in mg/dl fbs (fasting blood sugar > 120 mg/dl) (1 = true; 0 = false)
- restecg - resting electrocardiographic results
- thalach - maximum heart rate achieved
- exang - exercise induced angina (1 = yes; 0 = no) oldpeak ST depression induced by exercise
- slope - the slope of the peak exercise ST segment
- ca - number of major vessels (0-3) colored by flourosopy

- ♦ thal: 3 = normal; 6 = fixed defect; 7 = reversible defect target 1 or 0

```

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 heart = pd.read_csv('/content/drive/My Drive/MMO_Datasets/heart.csv', sep=',')
8 heart.head(10)
9

```

```

/usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning
import pandas.util.testing as tm

```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	tha
0	63	1	3	145	233	1	0	150	0	2.3	0	0	
1	37	1	2	130	250	0	1	187	0	3.5	0	0	
2	41	0	1	130	204	0	0	172	0	1.4	2	0	
3	56	1	1	120	236	0	1	178	0	0.8	2	0	
4	57	0	0	120	354	0	1	163	1	0.6	2	0	
5	57	1	0	140	192	0	1	148	0	0.4	1	0	
6	56	0	1	140	294	0	0	153	0	1.3	1	0	
7	44	1	1	120	263	0	1	173	0	0.0	2	0	
8	52	1	2	172	199	1	1	162	0	0.5	2	0	
9	57	1	2	150	168	0	1	174	0	1.6	2	0	

Задание 2. В случае необходимости проведите удаление или за кодирование категориальных признаков.

```

1 heart.shape

(303, 14)

1 heart.dtypes

```

age	int64
sex	int64
cp	int64
trestbps	int64
chol	int64
fbs	int64
restecg	int64
thalach	int64

### Проверка на пустые значения

```
slope      int64
1 heart.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

Пустых значений не обнаружено

```
1 X = heart.drop('target',axis = 1).values
2 y = heart['target'].values
```

## Задание 3. С использованием метода `train_test_split` разделите тестовую

Функция `train_test_split` используется для того, чтобы разделить исходную выборку так, что сохранились пропорции классов.

Разделим выборку следующим образом:

- Размер обучающей выборки (65%)
- Размер тестовой выборки (35%)

```
1 from sklearn.model_selection import train_test_split
2 heart_X_train, heart_X_test, heart_y_train, heart_y_test = train_test_split(
3     X, y, test_size=0.35, random_state=1)
```

```
1 print('heart_X_train: {} heart_y_train: {}'.format(heart_X_train.shape, heart_y_train.shape))
```

```

heart_X_train: (196, 13) heart_y_train: (196,)

1 print('heart_X_test: {} heart_y_test: {}'.format(heart_X_test.shape, heart_y_test.sh

heart_X_test: (107, 13) heart_y_test: (107,)

1 np.unique(heart_y_train)

array([0, 1])

1 np.unique(heart_y_test)

array([0, 1])

```

## Задание 4. Обучение модели ближайших соседей для произвол гиперпараметра K.

Обучите модель ближайших соседей для произвольно заданного гиперпараметр  
помощью трех подходящих для задачи метрик.

```

1 neighbors = np.arange(1,14)
2 len(neighbors)

13

```

### Обучение при различном количестве соседей

```

1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.model_selection import GridSearchCV
3 from sklearn.model_selection import learning_curve, validation_curve
4 from sklearn.model_selection import KFold, RepeatedKFold, LeaveOneOut, LeavePOut, Shu
5 from sklearn.model_selection import cross_val_score, cross_validate
6 from sklearn.metrics import roc_curve, confusion_matrix, roc_auc_score, accuracy_score

```

**Возвращает новый массив заданной формы и типа без инициализации записей.**

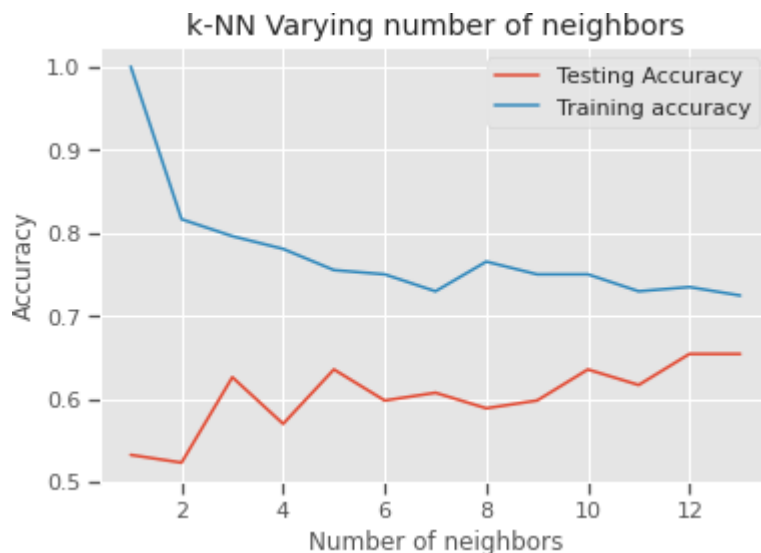
```

1 train_accuracy = np.empty(len(neighbors))
2 test_accuracy = np.empty(len(neighbors))
3
4 for i,k in enumerate(neighbors):
5     # Настройка классификатора Knn с K соседями
6     knn = KNeighborsClassifier(n_neighbors=k)
7
8     # Обучение модели
9     knn.fit(heart_X_train, heart_y_train)
10    # Вычленение точности на тренировочном наборе
11    train_accuracy[i] = knn.score(heart_X_train, heart_y_train)
12    # Вычисление точности на тестовом наборе
13    test_accuracy[i] = knn.score(heart_X_test, heart_y_test)

```

## Построить график для метода ближайших соседей

```
1 plt.style.use('ggplot')
2
3 plt.title('k-NN Varying number of neighbors')
4 plt.plot(neighbors, test_accuracy, label='Testing Accuracy')
5 plt.plot(neighbors, train_accuracy, label='Training accuracy')
6 plt.legend()
7 plt.xlabel('Number of neighbors')
8 plt.ylabel('Accuracy')
9 plt.show()
```



## ▼ Изучение работы KNeighborsClassifier

Установка a knn classifier with k neighbors, где k = 12

```
1 knn = KNeighborsClassifier(n_neighbors=12)
2 #Fit the model
3 knn.fit(heart_X_train,heart_y_train)

KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                    metric_params=None, n_jobs=None, n_neighbors=12, p=2,
                    weights='uniform')
```

Определение точности

```
1 knn.score(heart_X_test,heart_y_test)

0.6542056074766355

1 #import classification_report
2 from sklearn.metrics import classification_report
```

```

3
4 y_pred = knn.predict(heart_X_test)
5 print(classification_report(heart_y_test,y_pred))

```

	precision	recall	f1-score	support
0	0.62	0.66	0.64	50
1	0.69	0.65	0.67	57
accuracy			0.65	107
macro avg	0.65	0.65	0.65	107
weighted avg	0.66	0.65	0.65	107

## ▼ Оценка точности

```

1 from sklearn.metrics import roc_curve,confusion_matrix, roc_auc_score, accuracy_score

1 cl1_1 = KNeighborsClassifier(n_neighbors=12)
2 cl1_1.fit(heart_X_train, heart_y_train)
3 target1_1 = cl1_1.predict(heart_X_test)
4 accuracy_score(heart_y_test, target1_1)

0.6542056074766355

```

## ▼ Confusion Matrix - матрица ошибок

Матрица ошибок (confusion matrix) – таблица с 4 различными комбинациями прогнозируем Прогнозируемые значения описываются как положительные и отрицательные, а фактичес матрица ошибок используется для оценки точности моделей в задачах классификации. Но образов можно рассматривать как частный случай этой проблемы, поэтому confusion matr предсказаний.

```

1 y_pred = knn.predict(heart_X_test)
2 confusion_matrix(heart_y_test,y_pred)
3 pd.crosstab(heart_y_test, y_pred, rownames=['True'], colnames=['Predicted'], margins=

```

Predicted	0	1	All
True			
0	33	17	50
1	20	37	57
All	53	54	107

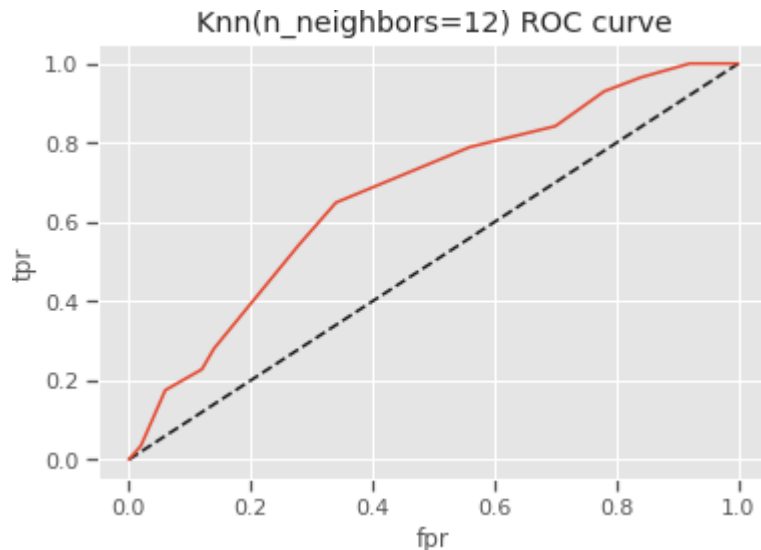
## ▼ ROC(Reciever Operating Charecteristic)-кривая

```

1 import matplotlib.pyplot as plt
2 plt.style.use('ggplot')

1 y_pred_proba = knn.predict_proba(heart_X_test)[: ,1]
2 fpr, tpr, thresholds = roc_curve(heart_y_test, y_pred_proba)
3 plt.plot([0,1],[0,1], 'k--')
4 plt.plot(fpr,tpr, label='Knn')
5 plt.xlabel('fpr')
6 plt.ylabel('tpr')
7 plt.title('Knn(n_neighbors=12) ROC curve')
8 plt.show()

```



```

1 y_pred_proba

array([0.5      , 0.66666667, 0.41666667, 0.66666667, 0.16666667,
       0.5      , 0.16666667, 0.58333333, 0.5      , 0.66666667,
       0.66666667, 0.5      , 0.58333333, 0.25     , 0.66666667,
       0.91666667, 0.91666667, 0.41666667, 0.91666667, 0.41666667,
       0.5      , 0.41666667, 0.41666667, 0.5      , 0.66666667,
       0.83333333, 0.5      , 1.        , 0.41666667, 0.66666667,
       0.91666667, 0.83333333, 0.66666667, 0.58333333, 0.41666667,
       1.        , 0.33333333, 0.66666667, 0.83333333, 0.25     ,
       0.66666667, 0.33333333, 0.66666667, 0.25     , 0.66666667,
       0.66666667, 0.5      , 0.75     , 0.91666667, 0.16666667,
       0.83333333, 0.83333333, 0.33333333, 0.33333333, 0.75     ,
       0.66666667, 0.66666667, 0.58333333, 1.        , 0.08333333,
       0.33333333, 0.91666667, 0.58333333, 0.66666667, 0.83333333,
       0.58333333, 0.5      , 0.75     , 0.16666667, 0.33333333,
       0.33333333, 0.75     , 0.08333333, 0.66666667, 0.16666667,
       0.33333333, 0.33333333, 0.5      , 0.41666667, 0.91666667,
       0.91666667, 0.66666667, 0.08333333, 0.5      , 0.5      ,
       0.5      , 0.25     , 0.5      , 0.91666667, 0.5      ,
       0.16666667, 0.5      , 0.66666667, 0.25     , 0.91666667,
       0.08333333, 0.5      , 0.58333333, 0.58333333, 0.58333333,
       0.66666667, 0.41666667, 0.5      , 0.66666667, 0.5      ,
       0.41666667, 0.66666667])

```

```

1 roc_auc_score(heart_y_test,y_pred_proba)

```



0.6740350877192982

## Задание 6. Произведите подбор гиперпараметра K с использованием валидации.

↳ 3 cells hidden

```
[ ] 1 import warnings
    2 warnings.filterwarnings('ignore')
    3
    4 param_grid = {'n_neighbors': np.arange(1,14)}
    5 knn = KNeighborsClassifier()
    6 knn_cv = GridSearchCV(knn, param_grid, cv=5)
    7 knn_cv.fit(X, y)
```

```
↳ 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,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13])},
               pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
               scoring=None, verbose=0)
```

```
[ ] 1 knn_cv.best_score_
```

```
↳ 0.6601092896174864
```

```
[ ] 1 knn_cv.best_params_
```

```
↳ {'n_neighbors': 12}
```

## Задание 5. Постройте модель и оцените качество модели с использованием валидации. Проведите эксперименты с тремя различными стратегиями K-fold. Перекрестная проверка — это процедура повторной выборки, используемая для оценки качества модели на ограниченной выборке данных.

```
1 scores = cross_val_score(KNeighborsClassifier(n_neighbors=12),
2                           X, y,
3                           cv=KFold(n_splits=10))
```

Значение метрики ассигасы для 10 фолдов Кол-во фолдов - кол-во разбиений

```
1 scores
array([0.5483871 , 0.61290323, 0.64516129, 0.6         , 0.5         ,
       0.66666667, 0.4         , 0.63333333, 0.66666667, 0.56666667])
```

Усредненное значение метрики ассигасы для 10 фолдов

```
1 np.mean(scores)
0.5839784946236558
```

```
1 import warnings
2 warnings.filterwarnings('ignore')
```

```

3  scoring = {'precision': 'precision_weighted',
4             'recall': 'recall_weighted',
5             'f1': 'f1_weighted'}
6
7  scores = cross_validate(KNeighborsClassifier(n_neighbors=4),
8                          X, y, scoring=scoring,
9                          cv=KFold(n_splits=5), return_train_score=True)
10 scores

```

```

{'fit_time': array([0.00094604, 0.00057364, 0.00054836, 0.00054193, 0.00052929]),
 'score_time': array([0.01519489, 0.00420713, 0.00418401, 0.0041275 , 0.00414848]),
 'test_f1': array([0.53012048, 0.65934066, 0.62005786, 0.62068966, 0.75      ]),
 'test_precision': array([1.      , 1.      , 0.74327869, 1.      , 1.      ]),
 'test_recall': array([0.36065574, 0.49180328, 0.60655738, 0.45      , 0.6      ]),
 'train_f1': array([0.76729516, 0.75673435, 0.74007707, 0.78961321, 0.79856665]),
 'train_precision': array([0.79272947, 0.78911533, 0.76025871, 0.79711127, 0.80984164]),
 'train_recall': array([0.7768595 , 0.76859504, 0.74380165, 0.78600823, 0.79423868])}

```

## ▼ LOOCV Leave One Out cross-validation

k равен общему размеру наблюдений в датасете. В тестовую выборку помещается единств  
фолдов в этом случае определяется автоматически и равняется количеству элементов. Да  
KFold. Существует правило, что вместо Leave One Out лучше использовать KFold на 5 или 1

```
1 # Эквивалент KFold(n_splits=n)
2 loo = LeaveOneOut()
3 loo.get_n_splits(X)
4
5 for train_index, test_index in loo.split(X):
6     heart_X_train, heart_X_test = X[train_index], X[test_index]
7     heart_y_train, heart_y_test = y[train_index], y[test_index]
```

## ▼ Repeated K-Fold

k-fold перекрестная валидация повторяется n раз, и датасет перетасовывается на каждой  
данных.

```
1 scores2 = cross_val_score(KNeighborsClassifier(n_neighbors=12),
2                             X, y,
3                             cv=RepeatedKFold(n_splits=5, n_repeats=2))

1 scores2

array([0.72131148, 0.6557377 , 0.6557377 , 0.6          , 0.61666667,
       0.60655738, 0.6557377 , 0.73770492, 0.68333333, 0.56666667])
```

## ▼ ShuPe Split

Генерируется N случайных перемешиваний данных, в каждом перемешивании заданная до

```
1 X = range(12)
2 # Эквивалент KFold(n_splits=n)
3 kf = ShuffleSplit(n_splits=5, test_size=0.35)
4 for train, test in kf.split(X):
5     print("%s %s" % (train, test))

[ 3  2  9  5  8 10  1] [ 4  6  7  0 11]
[11  7  9  6  0 10  2] [1 3 5 8 4]
[ 1  8  3  5  9 11  6] [ 7 10  2  4  0]
[ 4  8  3 11  0 10  6] [1 2 5 9 7]
[ 9  2 10  7  5  4  1] [ 8  6  0  3 11]
```

Задание 7. Повторите пункт 4 для найденного оптимального  
▼ значения

Сравните качество полученной модели с качеством модели

```

1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.35, random_stat

1 knn = KNeighborsClassifier(n_neighbors=12)
2 knn.fit(X_train,y_train)
3 knn.score(X_test,y_test)

0.6542056074766355

```

## ▼ Задание 8. Постройте кривые обучения и валидации.

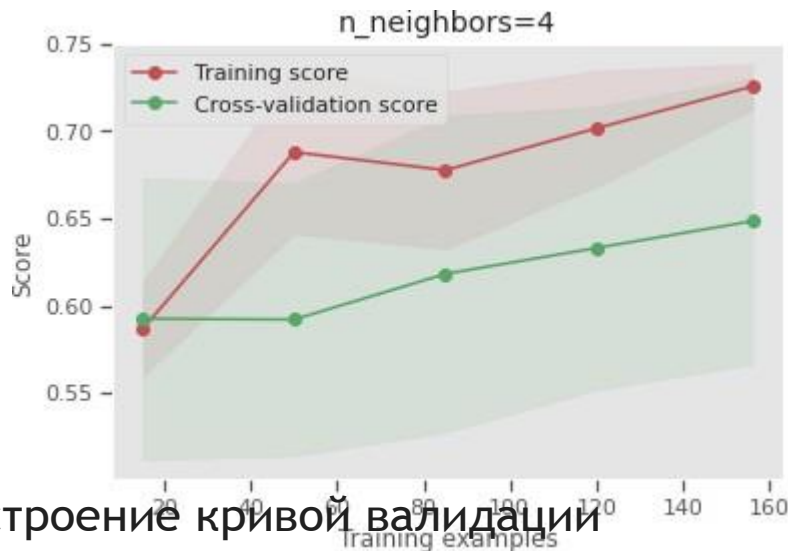
### ▼ Построение кривых обучения

```

1 def plot_learning_curve(estimator, title, X, y, ylim=None, cv=None,
2                         n_jobs=None, train_sizes=np.linspace(.1, 1.0, 5)):
3
4     plt.figure()
5     plt.title(title)
6     if ylim is not None:
7         plt.ylim(*ylim)
8     plt.xlabel("Training examples")
9     plt.ylabel("Score")
10    train_sizes, train_scores, test_scores = learning_curve(
11        estimator, X, y, cv=cv, n_jobs=n_jobs, train_sizes=train_sizes)
12    train_scores_mean = np.mean(train_scores, axis=1)
13    train_scores_std = np.std(train_scores, axis=1)
14    test_scores_mean = np.mean(test_scores, axis=1)
15    test_scores_std = np.std(test_scores, axis=1)
16    plt.grid()
17
18    plt.fill_between(train_sizes, train_scores_mean - train_scores_std,
19                    train_scores_mean + train_scores_std, alpha=0.1,
20                    color="r")
21    plt.fill_between(train_sizes, test_scores_mean - test_scores_std,
22                    test_scores_mean + test_scores_std, alpha=0.1, color="g")
23    plt.plot(train_sizes, train_scores_mean, 'o-', color="r",
24            label="Training score")
25    plt.plot(train_sizes, test_scores_mean, 'o-', color="g",
26            label="Cross-validation score")
27
28    plt.legend(loc="best")
29    return plt

1 plot_learning_curve(KNeighborsClassifier(n_neighbors=12), 'n_neighbors=4',
2                     X_train, y_train, cv=5)

```



## ▼ Построение кривой валидации

```

1 def plot_validation_curve(estimator, title, X, y,
2                           param_name, param_range, cv,
3                           scoring="accuracy"):
4
5     train_scores, test_scores = validation_curve(
6         estimator, X, y, param_name=param_name, param_range=param_range,
7         cv=cv, scoring=scoring, n_jobs=1)
8     train_scores_mean = np.mean(train_scores, axis=1)
9     train_scores_std = np.std(train_scores, axis=1)
10    test_scores_mean = np.mean(test_scores, axis=1)
11    test_scores_std = np.std(test_scores, axis=1)
12
13    plt.title(title)
14    plt.xlabel(param_name)
15    plt.ylabel("Score")
16    plt.ylim(0.0, 1.0)
17    lw = 4
18    plt.plot(param_range, train_scores_mean, label="Training score",
19             color="darkred", lw=lw)
20    plt.fill_between(param_range, train_scores_mean - train_scores_std,
21                    train_scores_mean + train_scores_std, alpha=0.2,
22                    color="darkred", lw=lw)
23    plt.plot(param_range, test_scores_mean, label="Cross-validation score",
24             color="green", lw=lw)
25    plt.fill_between(param_range, test_scores_mean - test_scores_std,
26                    test_scores_mean + test_scores_std, alpha=0.2,
27                    color="green", lw=lw)
28    plt.legend(loc="best")
29    return plt

1 n_range = np.array(range(5,55,5))
2 plot_validation_curve(KNeighborsClassifier(n_neighbors=4), 'knn',
3                       X_train, y_train,
4                       param_name='n_neighbors', param_range=n_range,
5                       cv=5, scoring="accuracy")

```

