	Подготовка обучающей и тестовой выборки. Кросс-валидация и подбор гиперпараметров на примере метода ближайших соседей. 1. Цель работы Изучение способов подготовки выборки и подбора гиперпараметров на примере метода ближайших соседей.
In [2]	2. Подготовка данных и библиотек from operator import itemgetter import matplotlib. pyplot as plt import matplotlib. ticker as ticker
	<pre>import numpy as np import pandas as pd import pandas_profiling import math from simple_kNN import * from sklearn.datasets import *</pre>
	from typing import Dict, Tuple from scipy import stats from sklearn.model_selection import train_test_split from sklearn.model_selection import GridSearchCV, RandomizedSearchCV from sklearn.model_selection import cross_val_score from sklearn.model_selection import cross_val_score from sklearn.neighbors import KNeighborsRegressor, KNeighborsClassifier from sklearn.metrics import accuracy_score, balanced_accuracy_score
	from sklearn.metrics import plot_confusion_matrix from sklearn.metrics import precision_score, recall_score, f1_score from sklearn.metrics import confusion_matrix from sklearn.metrics import mean_absolute_error, mean_squared_error from sklearn.metrics import roc_curve, roc_auc_score, classification_report import seaborn as sns import matplotlib. pyplot as plt
In [3]	*matplotlib inline sns.set(style="ticks") # Загрузка датасета из скалерна data = load_boston() # Конвертирование массива в датафрейм
	boston_df = pd.DataFrame(data=np.c_[data['data']], columns=data.feature_names) boston_df['target'] = data['target'] # Cmonδιμω data.feature_names
	array(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD',
	# Column Non-Null Count Dtype CRIM 506 non-null float64 I ZN 506 non-null float64 I NDUS 506 non-null float64 CHAS 506 non-null float64 NOX 506 non-null float64 RM 506 non-null float64
In [6]	8 RAD 506 non-null float64 9 TAX 506 non-null float64 10 PTRATIO 506 non-null float64 11 B 506 non-null float64 12 LSTAT 506 non-null float64 13 target 506 non-null float64 dtypes: float64(14) memory usage: 55.5 KB
Out[6]	
In [7]	25% 0.082045 0.00000 5.19000 0.00000 0.44900 5.88500 45.02500 2.100175 4.00000 279.00000 17.02500 d.95000 17.02500 d.95000 0.00000 9.69000 0.00000 0.53800 6.20850 77.50000 320.00000 19.05000 391.44000 11.36000 21.20000 d.95000 18.10000 0.00000 18.10000 0.00000 0.62400 6.62350 94.07500 5.188425 24.00000 666.00000 20.20000 396.22500 16.95500 25.00000 d.95000 27.74000 1.00000 0.87100 0.87100 8.78000 100.00000 12.126500 24.00000 711.00000 375.97750 6.95000 17.02500 25.00000 d.95000 17.02500 d.95000 100.00000 27.74000 1.00000 0.87100 8.78000 100.00000 12.126500 24.00000 711.00000 396.90000 37.97000 50.00000 d.95000 d.95
	mask = np.zeros_like(boston_df.corr(), dtype=np.bool) mask[np.triu_indices_from(mask)] = True sns.heatmap(boston_df.corr(), cmap='YlGnBu', mask=mask, annot=True, fmt='.3f') <pre> </pre> <pre> </pre> <pre> </pre> <pre> <pre> </pre> <pre> <pre> </pre> <pre> <pre> <pre> </pre> <pre> <pre> <pre> <pre> </pre> <pre> <pre> <pre> <pre> <pre> </pre> <pre> <pre> <pre> <pre> <pre> </pre> <pre> <pr< th=""></pr<></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre>
	ZN - 0.200 INDUS - 0.407
	RM - 0.219 0.312 0.392 0.091 0.302 - 0.392 0.091 0.302 0.22 AGE - 0.353 0.570 0.645 0.087 0.731 0.240 - 0.200 0.769 0.205 0.748 - 0.090 0.205 0.748 RAD - 0.626 0.312 0.595 0.007 0.611 0.210 0.456 0.495
	TAX = 0.583
In [8]	LSTAT - 0.456
Out[8]	3. Разделение датафрейма на обучающую и тестовую выборку boston_x_train, boston_x_test, boston_y_train, boston_y_test = train_test_split(data.data, data.target, test_size=0.25, random_state=1)
In [10]	print(f'boston_X_train\tboston_y_train\nCтроки : {boston_x_train.shape[0]}\tCтроки : {boston_y_train.shape[0]}\nCтолбцы : {boston_x_train.shape[1]}\tCтолбцы : 1') boston_X_train
	boston_X_test boston_y_test Строки: 127 Строки: 127 Столбцы: 13 Столбцы: 1 пр.unique(boston_y_train) array([5. , 5.6, 7. , 7.2, 7.4, 7.5, 8.3, 8.4, 8.5, 8.7, 8.8, 9.5, 9.6, 9.7, 10.2, 10.4, 10.5, 10.8, 10.9, 11.3, 11.7, 11.8
	11.9, 12., 12.1, 12.3, 12.5, 12.7, 12.8, 13.1, 13.2, 13.3, 13.4, 13.5, 13.6, 13.8, 13.9, 14. , 14.1, 14.2, 14.3, 14.4, 14.5, 14.8, 14.9, 15., 15.1, 15.2, 15.4, 15.6, 15.7, 16.8, 15.7, 16. , 16.1, 16.2, 16.3, 16.4, 16.5, 16.6, 16.7, 16.8, 17. , 17.1, 17.2, 17.3, 17.4, 17.5, 17.6, 17.7, 17.8, 17.9, 18.2, 18.3, 18.4, 18.5, 18.6, 18.7, 18.8, 18.9, 19. , 19.1, 19.2, 19.3, 19.4, 19.5, 19.6, 19.7, 19.8, 19.9, 20. , 20.1, 20.2, 20.3, 20.4, 20.5, 20.6, 20.7, 20.8, 20.9, 21. , 21.1, 21.2, 21.4, 21.5, 21.6, 20.7, 21.8, 21.9, 22. , 22.2, 22.3, 22.4, 22.5, 22.6, 22.7, 22.8, 22.9, 23. , 23.1, 23.2, 23.3,
In [13]	23.4, 23.6, 23.7, 23.8, 23.9, 24., 24.1, 24.2, 24.3, 24.4, 24.5, 24.6, 24.7, 24.8, 25., 25.1, 25.2, 26.2, 26.4, 26.6, 26.7, 27.1, 27.5, 28., 28.1, 28.4, 28.5, 28.6, 28.7, 29., 29.1, 29.4, 29.6, 29.8, 29.9, 30.1, 30.3, 30.5, 30.7, 30.8, 31., 31.1, 31.5, 31.6, 31.7, 32., 32.2, 32.4, 32.5, 32.7, 32.9, 33., 33.1, 33.2, 33.3, 33.3, 33.4, 33.8, 34.6, 34.7, 34.9, 35.1, 35.2, 35.4, 36.1, 36.2, 36.5, 37., 37.6, 37.9, 38.7, 39.8, 41.3, 42.3, 42.8, 43.8, 44., 45.4, 46., 46.7, 48.5, 48.8, 50.]) print(f'Paamep boston_y_test : {len(boston_y_test)}')
	Print(boston_y_test) Pasmep boston_y_test: 127 [28.2 23.9 16.6 22. 20.8 23. 27.9 14.5 21.5 22.6 23.7 31.2 19.3 19.4 19.4 27.9 13.9 50. 24.1 14.6 16.2 15.6 23.8 25. 23.5 8.3 13.5 17.5 43.1 11.5 24.1 18.5 50. 12.6 19.8 24.5 14.9 36.2 11.9 19.1 22.6 20.7 30.1 13.3 14.6 8.4 50. 12.7 25. 18.6 29.8 22.2 28.7 23.8 8.1 22.2 6.3 22.1 17.5 48.3 16.7 26.6 8.5 14.5 23.7 37.2 41.7 16.5 21.7 22.7 23. 10.5 21.9 21. 20.4 21.8 50. 22. 23.3 37.3 18. 19.2 34.9 13.4 22.9 22.5 13. 24.6 18.3 18.1 23.9 50. 13.6 22.9 10.9 18.9 22.4 22.9
In [16]	44.8 21.7 10.2 15.4 25.3 23.3 7.2 21.2 11.7 27. 29.6 26.5 43.5 23.6 11. 33.4 36. 36.4 19. 20.2 34.9 50. 19.3 14.9 26.6 19.9 24.8 21.2 23.9] 4. Построение базовой модели на основе ближайших соседей # 2 ближайших соседа c11_1 = KNeighborsRegressor(n_neighbors=2)
Out[16]	<pre>cl1_1.fit(boston_x_train, boston_y_train) target1_1 = cl1_1.predict(boston_x_test) len(target1_1), target1_1 (127, array([26.4 , 28.4 , 18. , 19.15, 33.05, 17.3 , 29.65, 15.1 , 22.35,</pre>
	23.95, 12.6, 20.95, 18.8, 22.2, 21.6, 30.25, 11.05, 13.55, 17.3, 42.65, 13.7, 33.4, 19.05, 41.95, 15.1, 21.1, 21.05, 15.7, 26.05, 11.15, 22.9, 22.2, 21.4, 22.2, 19, 17.8, 12.95, 50. , 13.95, 21.6, 14.35, 21.15, 22.65, 18.6, 17.2, 12.1, 21.3, 15.55, 20.45, 18. , 28. , 15.7, 28.45, 14.1, 14.2, 18.9, 43.1, 28.3, 22. , 24.65, 32.8, 26.65, 93.5, 15.7, 17.4, 19.75, 20.9, 20.95, 22.75, 45.65, 29.65, 17.45, 20.65, 31.25, 13. , 24.3, 24.8, 15.7, 23.65, 18.6, 17.45, 22.5, 42.65, 17.65, 20.65, 20.65, 22.85, 26.83, 22.8, 22.8, 42.65, 17.65, 20.65, 20.65, 21.3, 16.25, 15.1, 24.95, 26.43, 16.25, 15.1, 24.95, 26.43, 16.25, 25.1, 24.95, 26.43, 16.25, 26.53, 20.85, 2
In [17]	38. , 22.5 , 22.85, 22.35, 15.2 , 36.9 , 39.45, 30.15, 20.05, 14.35, 29.65, 28.3 , 20.85, 13.05, 24.95, 15.3 , 21.55, 18.6 , 25.6]))
Out[17]	len(target1_2), target1_2 (127, array([31.08, 29.37, 21.51, 26.99, 24.53, 16.73, 35.35, 16.9 , 21.33, 25.37, 21.75, 29.37, 18.03, 26.44, 20.25, 20.77, 18.89, 26.81, 26.44, 20.73, 22.08, 17.12, 20.11, 26.21, 26.86, 12.29, 17.34, 20.79, 29.31, 12.67, 29.09, 20.94, 36.9 , 15.86, 21.16, 21.3 , 18.79, 22.76, 18.41, 21.04, 25. , 21.71, 25.25, 18.19, 17.54,
	12. 22, 36. 12, 12. 78, 17. 74, 16. 66, 22. 33, 23. , 25. 9, 18. 15, 17. 94, 16. 66, 22. 34, 23. 25. 9, 18. 15, 17. 94, 17. 92, 173, 20. 77, 22. 92, 21. 18, 26. 27, 23. 34, 26. 78, 26. 58, 9. 67, 20. 92, 23. 18, 19. 91, 21. 68, 20. 37, 23. 43, 27. 43, 35. 35, 16. 86, 18. 56, 25. 83, 12. 87, 23. 85, 28. 28, 17. 93, 27. 74, 20. 16, 16. 86, 23. 15, 29. 99, 15. 59, 22. 68, 12. 83, 21. 24, 24. 44, 24. 44, 28. 14, 20. 89, 13. 22, 14. 73, 22. 27, 24. 21, 10. 81, 20. 89, 12. 29, 24. 21, 10. 81, 20. 89, 12. 29, 24. 21, 23. 30. 72, 28. 83, 25. , 12. 29, 22. 7, 25. 13, 23. 74, 20. 49, 26. 9]))
	5. Метрики качества классификации cl1_1.score(boston_x_train, boston_y_train) e.8566879344566602
Out[19]	cl1_2.score(boston_x_train, boston_y_train) e.5342300718439038 6. Произведите подбор гиперпараметра К с использованием GridSearchCV
ın [21]	param_grid = {'n_neighbors' : np.arange(1, 25)} grid = GridSearchCV(KNeighborsRegressor(), param_grid, cv=3) grid.fit(boston_x_train, boston_y_train) print(f'Подобранный параметр: {grid.best_params_}') print(f'Оценка при подобранном параметрe: {grid.best_score_}') Подобранный параметр: {'n_neighbors': 7} Оценка при подобранном параметрe: 0.3376584947980437
In [22]	Подобранный параметр: (n_netgnotes : // Оценка при подобранном параметре: (0.3376584947980437 7. Кросс валидация DEFAULT_FOLDS = 3 DEFAULT_RANGES = np.arange(1, 25)
Out[23]	<pre>cv1 = cross_val_score(KNeighborsRegressor(n_neighbors=5), data.data, data.target, cv=DEFAULT_FOLDS, scoring='max_error') cv1 array([-31.72, -31.66, -31.66]) from sklearn.model_selection import *</pre>
	<pre>def knn_cross_strtegies(mass_strategies, mass_matrics, k_range, cv, x, y): k_strategies_scores = {} for strat in mass_strategies: strategia = strat(cv) # n_splits k_scores = knn_cross_val_score(mass_matrics, k_range, strategia, x, y)</pre>
	<pre>k_scores = knn_cross_val_score(mass_matrics, k_range, strategia, x, y)</pre>
	<pre>k_scores = {} for metric in mass_matrics: for k in k_range: knn = KNeighborsRegressor(n_neighbors=k) scores = cross_val_score(knn, x, y, cv=cv, scoring=metric) if metric not in k_scores:</pre>
	<pre>k_scores[metric].append(scores.mean()) return k_scores def plot_knn_cross_val_score(k_scores, cv):</pre>
	<pre>columns=len(k_scores.items()) fig = plt.figure(figsize=(10*columns, 10)) rows=1 index = 1 for stratName, metrics in k_scores.items(): ax = fig.add_subplot(rows, columns, index) for name, scores in metrics.items():</pre>
Tr. F.	<pre>ax.plot(np.arange(1, len(scores)+1), scores) ax.set_title(stratName) ax.legend(metrics.keys(), loc='lower right') index+=1</pre>
In [26] In [27]	strategies = [KFold, RepeatedKFold, ShuffleSplit]
In [28]	warnings.warn("Pass {} as keyword args. From version 0.25 " 7. Сравнение моделей
	-20 - -20 -
	-40- -30-
	-60 - -40 - -50 -
	-805060 -