о 1 2 3 4	from sklearn.model_selection import RepeatedKFold from sklearn.model_selection import cross_val_score from sklearn.metrics import mean_absolute_error, mean_squared_error, median_absolute_error, r2_score from sklearn.model_selection import GridSearchCV, RandomizedSearchCV pd.options.display.float_format = '{:.0f}'.format Вагрузка данных рочитаем имеющийся набора данных в переменную
1 2 3 4	рочитаем имеющийся набора данных в переменную data = pd.read_csv('houses_to_rent_v2.csv') ыведем основую информацию и статистические показатели по датасету data.head()
	São Paulo 320 4 4 0 20 acept not furnished 1200 4960 1750 63 7973 Porto Alegre 80 1 1 6 acept not furnished 1000 2800 0 41 3841 Porto Alegre 51 2 1 0 2 acept not furnished 270 1112 22 17 1421
< c Ra Da	São Paulo 25 1 1 0 1 not acept not furnished 0 800 25 11 836
	0 city 10692 non-null object 1 area 10692 non-null int64 2 rooms 10692 non-null int64 3 bathroom 10692 non-null int64 4 parking spaces 10692 non-null int64 5 floor 10692 non-null object 6 animal 10692 non-null object 7 furniture 10692 non-null object 8 hoa (R\$) 10692 non-null int64 9 rent amount (R\$) 10692 non-null int64 10 property tax (R\$) 10692 non-null int64
d- me	11 fire insurance (R\$) 10692 non-null int64 12 total (R\$) 10692 non-null int64 ltypes: int64(9), object(4) lemory usage: 1.1+ MB data.describe() area rooms bathroom parking spaces hoa (R\$) rent amount (R\$) property tax (R\$) fire insurance (R\$) total (R\$)
m	mean 149 3 2 2 1174 3896 367 53 5490 std 537 1 1 2 15592 3409 3108 48 16485 min 11 1 1 0 0 450 0 3 499 25% 56 2 1 0 1530 38 21 2062 50% 90 2 2 1 560 2661 125 36 3582
	75% 182 3 3 2 1238 5000 375 68 6768 max 46335 13 10 12 1117000 45000 313700 677 1120000 sns.heatmap(data.corr(), annot=True, fmt='.1f') AxesSubplot:>
	area - 10
	total (R\$) - 0.2 0.6 0.7 0.6 0.0 1.0 0.1 0.2 0.1 1.0 0.3 0.2 0.1 1.0 0.3 0.2 0.3 1.0 0.3 0.2 0.3 1.0 0.3 0.2 0.3 0.3 0.3 0.3 0.3 0.3 0.3 0.3 0.3 0.3
Π _i	огласно полученной документации мы имеем датасет из 13 колонок и 10692 строк: ризнаки сity - город, в котором расположена квартира аrea - площадь квартиры гоотs - количество комнат в квартире bathroom - количество ванных комнат
	 parking spaces - количество парковочных мест floor - этаж, на котором расположена квартира animal - разрешение на проживание с домашним животным furniture - мебель hoa - налог на товарищество собственников жилья rent amount - размер арендной платы property tax - налог на имущество fire insurance - страховка от пожара
Ц •	елевая переменная • total - суммарная стоимость аренды данных остутствуют явные пропуски, однако о качестве данных все же известно мало, поэтому потребуется предобработка. Тредобработка данных
(data.columns index(['city', 'area', 'rooms', 'bathroom', 'parking spaces', 'floor',
	data = data.rename(columns={'hoa (R\$)':'hoa',
В (хр	1081
3 4 5 6 7 8 9 10 11	.1 303 .2 257
14 1! 10 11 18 19 20 21 21 21 21 21	.5 147 .6 109 .7 96 .8 75 .9 53 .0 44 .1 42 .5 25 .25 25 .25 24
20 22 22 23 33 30 51 40 31	20
\(\frac{1}{2}\)	ак и есть, значений, заполненных "заглушкой" в виде - почти 25% датасета. Видимо при сборе данных некоторые пользователи не указали этаж и на этом месте появилась заглушка, или же опрос проводили несколько раз, и поле с номером этажа было н data['floor'].value_counts().plot(kind='bar') AxesSubplot:>
19	2000 - 150
	0 '
	2000 - 800 - 600 -
	data.query('floor != "-"')['floor'].median()
((data = data.query('floor != "-"').reset_index(drop=True) data = data.astype({'floor':'int32'}) data = data.astype({'floor':'int32'}) data.info() calass 'pandas.core.frame.DataFrame'>
Da	tangeIndex: 8231 entries, 0 to 8230 tata columns (total 13 columns): # Column Non-Null Count Dtype **Column Non-Null count object **Column Non-Null count object **Column Non-Null count object **Column Non-Null count object **Column Non-Null objec
d ⁻ me	7 furniture 8231 non-null object 8231 non-null int64 9 rent_amount 8231 non-null int64 10 property_tax 8231 non-null int64 11 fire_insurance 8231 non-null int64 12 total 8231 non-null int64 12 total 8231 non-null int64 14 ltypes: int32(1), int64(9), object(3) 12 ltypes: int32(1), int64(9), object(3) 14 ltypes: Non-null int64 15 ltypes: Non-null int64 16 ltypes: Non-null int64 17 ltypes: Non-null int64 17 ltypes: Non-null int64 18
(акодируем категориальные признаки с помощью техники One-Hot Encoding data_ohe = pd.get_dummies(data, columns=['city', 'animal', 'furniture'], drop_first=1) data_ohe area rooms bathroom parking_spaces floor hoa rent_amount property_tax fire_insurance total city_Campinas city_Porto Alegre city_Rio de Janeiro city_São Paulo animal_not acept furniture_not furnished
82 82 Р	228 285 4 4 4 4 17 3100 15000 973 191 19260 0 0 0 1 0 1 0 1 1 1 0 1 1 1 0 1 1 1 0 1
() ()	print(features_train.shape, features_test.shape) (6173, 13) (2058, 13) print(target_train.shape, target_test.shape) (6173,) (2058,) Обучение моделей
r 9 r	Model_DT = DecisionTreeRegressor(random_state=1, max_depth=10) %%time model_DT.fit(features_train, target_train) Mall time: 19 ms
ţ	predictions_DT = model_DT.predict(features_test) print('Decision Tree:', mean_absolute_error(target_test, predictions_DT)) print('Decision Tree: 332.5153566973952
r (Model = XGBRegressor(objective='reg:squarederror') cv = RepeatedKFold(n_splits=10, n_repeats=3, random_state=1) n_scores = cross_val_score(model, features, target, scoring='neg_mean_absolute_error', cv=cv, n_jobs=-1, error_score='raise')
MA r	print('MAE (Средняя Абсолютная Ошибка): %.3f (%.3f)' % (пр.mean(n_scores), пр.std(n_scores))) МАЕ (Средняя Абсолютная Ошибка): -181.386 (126.665) model = XGBRegressor(objective='reg:squarederror') %%time model.fit(features_train, target_train)
	GGRegressor(base_score=0.5, booster='gbtree', callbacks=None,
ŗ	reg_lambda=1,) predictions_GB = model.predict(features_test) print('Gradient boosting:', mean_absolute_error(target_test, predictions_GB)) Gradient boosting: 137.5728101368896
Γν r t	loдбор гиперпараметров иперпараметры модели дерева решений n_range = np.array(range(1, 41, 2)) tuned_parametrs = [{'max_depth' : n_range}] tuned_parametrs
t t	<pre>{'max_depth': array([1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31, 33,</pre>
1	inidSearchCV(cv=5, estimator=DecisionTreeRegressor(random_state=1),
ŗ	param_test = { 'gamma':[0, 0.3, 0.4, 0.5], 'subsample':[0, 0.3, 0.5, 0.7, 0.9], 'reg_lambda':[1,5,10], 'max_depth':[1, 3, 5], 'colsample_bytree': [0.2, 0.4, 0.8, 1]} %*time
t t Wa	<pre>tuned_GB = GridSearchCV(XGBRegressor(objective='reg:squarederror'), param_grid = param_test, scoring='neg_mean_squared_error', cv=5) tuned_GB.fit(features, target) vall time: 6min 42s iridSearchCV(cv=5,</pre>
	colsample_byree=None, early_stopping_rounds=None, enable_categorical=False, eval_metric=None, gamma=None, gpu_id=None, grow_policy=None, importance_type=None, interaction_constraints=None, learning_rate=None, max_bin=None, max_cat min_child_weight=None, missing=nan, monotone_constraints=None, n_estimators=100, n_jobs=None, num_parallel_tree=None, predictor=None, random_state=None,
	reg_alpha=None, reg_lambda=None,), param_grid={'colsample_bytree': [0.2, 0.4, 0.8, 1],
{	'colsample_bytree': 1, 'gamma': 0, 'max_depth': 3, 'reg_lambda': 1, 'subsample': 0.5} Обучение оптимальных моделей дерево решений
{ О Д	tuned DT.best estimator fit(features train target train)
{ О Д	tuned_DT.best_estimatorfit(features_train, target_train) becisionTreeRegressor(max_depth=27, random_state=1) predictions_DT_tuned_train = tuned_DT.best_estimatorpredict(features_train) predictions_DT_tuned_test = tuned_DT.best_estimatorpredict(features_test) радиентный бустинг
{ О Д г г г	predictions_DT_tuned_train = tuned_DT.best_estimatorpredict(features_train) predictions_DT_tuned_test = tuned_DT.best_estimatorpredict(features_test) радиентный бустинг tuned_GB.best_estimatorfit(features_train, target_train) GGRegressor(base_score=0.5, booster='gbtree', callbacks=None,
О Д , г	predictions_DT_tuned_train = tuned_DT.best_estimatorpredict(features_train) predictions_DT_tuned_test = tuned_DT.best_estimatorpredict(features_test)
{ O Д I I I I I I I I I I I I I I I I I I	ecisionTreeRegressor(wax_depth=27, random_state=1) predictions_Di_tuned_resis = tuned_Di_best_estimatorpredict(Features_test) predictions_Di_tuned_resis = tuned_Di_best_estimatorpredict(Features_test) pagueththsis Gyctust tuned_GSt.best_estimatorfit(Features_test) tuned_GSt.best_estimatorfit(Features_test) pagueththsis Gyctust paguethtsis Gyctust pagueththsis Gyctust pagueththsis Gyctust paguethtsis Gyctust pagueththsis Gyctust
{ O Д i i i i i i i i i i i i i i i i i i	predictions_Of_tuned_train = tuned_Of_hest_extinatorpredict(footume_trait) predictions_Of_tuned_train = tuned_Of_hest_extinatorpredict(footume_trait) page(Bit best_extinatorfit(footume_trait) predictions_Of_tuned_train = tuned_Of_hest_extinatorpredict(footume_trait) Deleta kauectrae nonyvee+hext модевей использовансь стандартные метрики качестве nonyvee+hext модевей использовансь стандартные metallic non- non-point = 0 non-point =
{ O Д I I I I I I I I I I I I I I I I I I	productions DI tamod train = tamod DI base_estimator_product(features_train) productions_DI tamod_tamod_tamod_tamod_tamod_train) productions_DI tamod_
{ O Д	prelictions_IL_core_Il_core_I = ane_IL_I bestschedus_prelict_feature_Lords) prelictions_IL_core_Il_core = ane_IL_I bestschedus_prelict_feature_Lords) prelictions_IL_core_Il_core = ane_IL_I bestschedus_prelict_feature_Lords) prelictions_IL_core_Il_core_I
{ O 4	Account reclamation that any interface and any interface profit profit or any interface and any interf
{ O Д	The second of the content of the con
{ OД 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	AND CONTROL OF THE PROPERTY OF
{ OД 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	STANDARD STA