

Министерство образования и науки
федеральное государственное автономное образовательное учреждение
высшего образования
«Национальный исследовательский университет ИТМО»
Факультет инфокоммуникационных технологий

Отчет по дисциплине: **«Современные инструменты анализа данных»**

Лабораторная работа 3

Выполнил:

—

Проверила:

—

Санкт-Петербург

2023

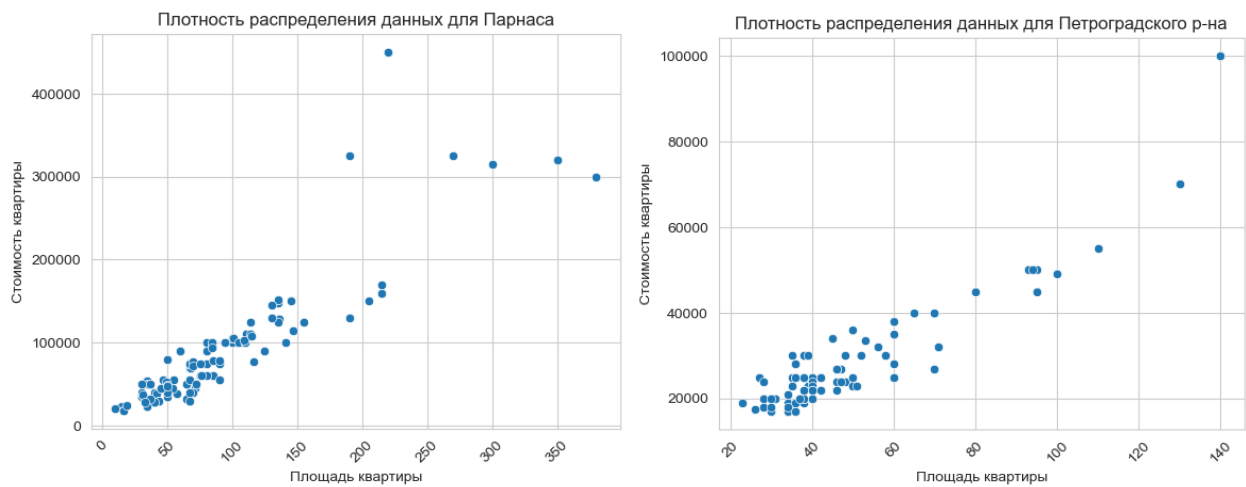
Работа делалась с использованием **python** и последующих модулей для работы с данными и их визуализации: **pandas, numpy, seaborn, matplotlib, statsmodels, sklearn**

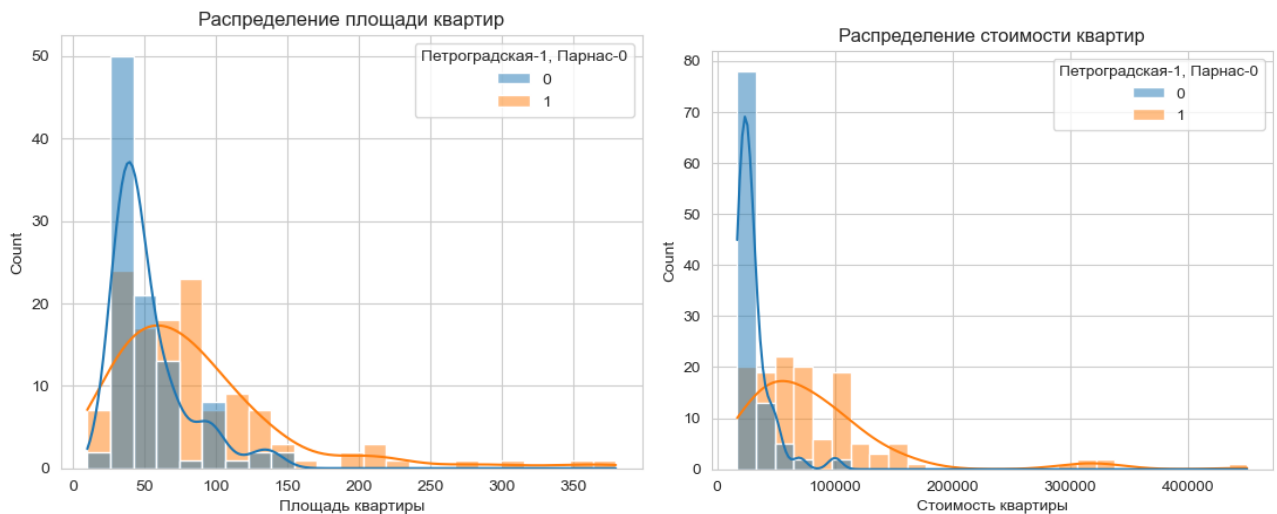
Часть 1. Анализ многомерных данных. Парная и множественная регрессия

- 1. Указать описательные статистики, графики плотности и гистограммы для каждого района

	Парнас - стоимость кв	Парнас - площадь кв	Петроградская - стоимость кв	Петроградская - площадь кв
count	126.000000	126.000000	100.000000	100.0000
mean	83115.079365	84.742063	29880.000000	52.3300
std	68314.044316	62.512390	14634.282429	25.6598
min	18000.000000	10.000000	17000.000000	23.0000
25%	45000.000000	43.000000	22000.000000	36.0000
50%	70000.000000	70.000000	25000.000000	42.0000
75%	100000.000000	104.500000	32000.000000	60.0000
max	450000.000000	380.000000	100000.000000	140.0000

Статистическое описание данных в зависимости от района





2. Построить модель **парной регрессии** для квартир площадью от 20 до 110 кв.м включительно стоимости от площади. Оценить характеристики построенной модели.

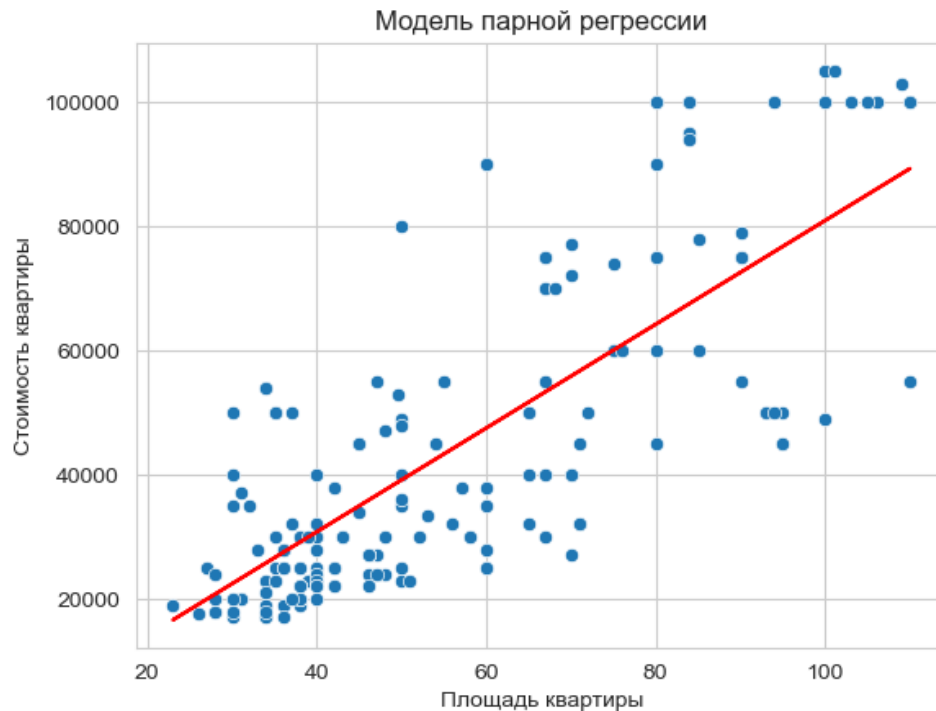
	Стоимость квартиры	Площадь квартиры	Петроградская-1, Парнас-0
Стоимость квартиры	1.000000	0.887266	0.455191
Площадь квартиры	0.887266	1.000000	0.309301
Петроградская-1, Парнас-0	0.455191	0.309301	1.000000

Матрица корреляции признаков

OLS Regression Results						
=====						
Dep. Variable:	Стоимость квартиры	R-squared:	0.591			
Model:	OLS	Adj. R-squared:	0.589			
Method:	Least Squares	F-statistic:	269.9			
Date:	Wed, 11 Oct 2023	Prob (F-statistic):	4.01e-38			
Time:	18:58:08	Log-Likelihood:	-2097.0			
No. Observations:	189	AIC:	4198.			
Df Residuals:	187	BIC:	4204.			
Df Model:	1					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

const	-2637.4656	3091.642	-0.853	0.395	-8736.443	3461.512
Площадь квартиры	835.3907	50.851	16.428	0.000	735.075	935.706
=====						
Omnibus:	4.598	Durbin-Watson:	0.691			
Prob(Omnibus):	0.100	Jarque-Bera (JB):	4.660			
Skew:	0.378	Prob(JB):	0.0973			
Kurtosis:	2.863	Cond. No.	161.			

Результаты построенной парной линейной регрессии



3. Построить модель **множественной регрессии** для квартир площадью от 20 до 110 кв.м включительно стоимости от площади и района. Оценить характеристики построенной модели.

OLS Regression Results						
=====						
Dep. Variable:	Стоимость квартиры	R-squared:	0.801			
Model:	OLS	Adj. R-squared:	0.799			
Method:	Least Squares	F-statistic:	374.7			
Date:	Wed, 11 Oct 2023	Prob (F-statistic):	5.80e-66			
Time:	18:58:12	Log-Likelihood:	-2028.8			
No. Observations:	189	AIC:	4064.			
Df Residuals:	186	BIC:	4073.			
Df Model:	2					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

const	-4750.7798	2165.986	-2.193	0.030	-9023.838	-477.722
Площадь квартиры	661.4266	37.641	17.572	0.000	587.169	735.684
Петроградская-1, Парнас-0	2.42e+04	1725.202	14.030	0.000	2.08e+04	2.76e+04
=====						
Omnibus:	3.103	Durbin-Watson:	1.088			
Prob(Omnibus):	0.212	Jarque-Bera (JB):	3.437			
Skew:	0.021	Prob(JB):	0.179			
Kurtosis:	3.659	Cond. No.	162.			
=====						

Результаты построенной множественной линейной регрессии

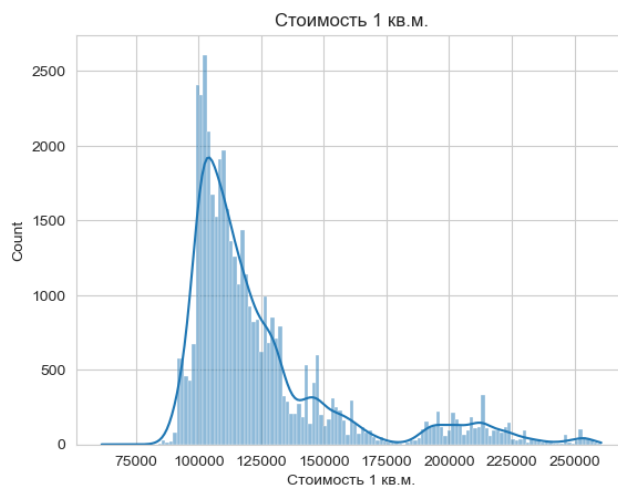
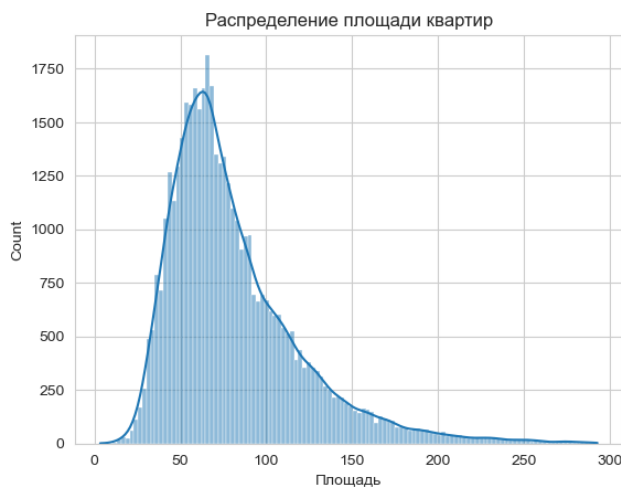
Часть 2. Анализ многомерных данных. Множественная регрессия

1. Скорректировать выборку случайным образом, чтобы осталось минимум 5000 записей.

Из выборки были удалены выбросы по всем непрерывным признакам с использованием интервала $[\text{mean} - 3 * \text{std}; \text{mean} + 3 * \text{std}]$ и выбраны случайным образом 5000 записей.

	Площадь	Стоимость 1 кв.м.	Стоимость
count	44035.000000	44035.000000	4.403500e+04
mean	81.370508	124625.932105	1.004983e+07
std	38.888211	32352.222625	5.480365e+06
min	3.300000	61266.670000	4.095822e+05
25%	54.900000	103689.685000	6.588206e+06
50%	71.600000	113406.300000	8.395278e+06
75%	98.400000	130333.160000	1.150167e+07
max	292.700000	260715.100000	3.383975e+07

Статистика по непрерывным данным после преобразования

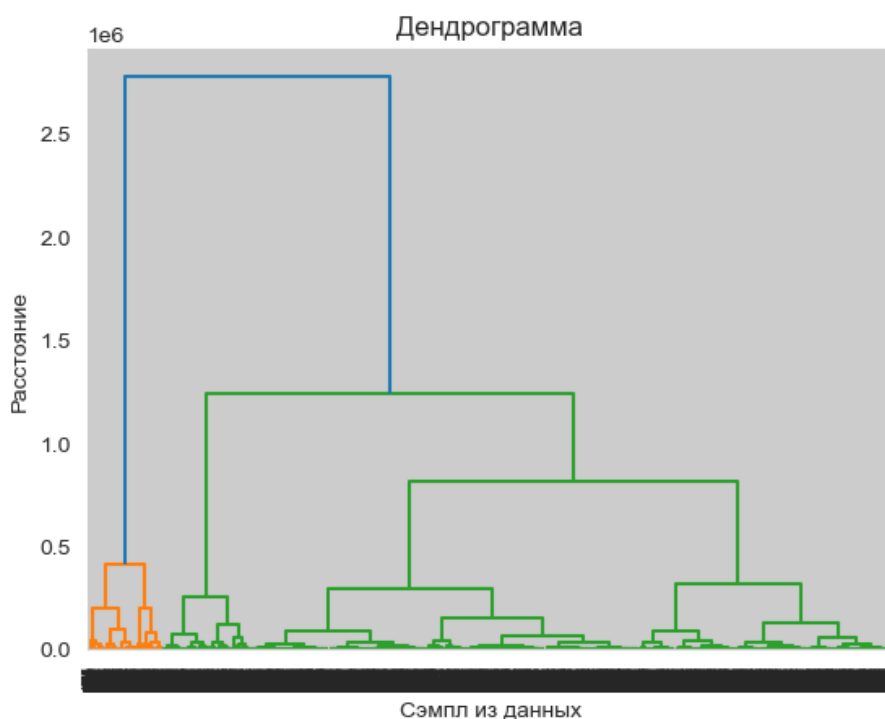


Выше представлены результаты множественной линейной регрессии. Категориальные признаки были закодированы с использованием one-hot encoding.

Часть 3. Анализ многомерных данных. Кластеризация

Для выбранных данных провести кластеризацию каждым из методов. Переменные подобрать самостоятельно из количественных. Оценить характеристики построенной модели и выводы по результатам.

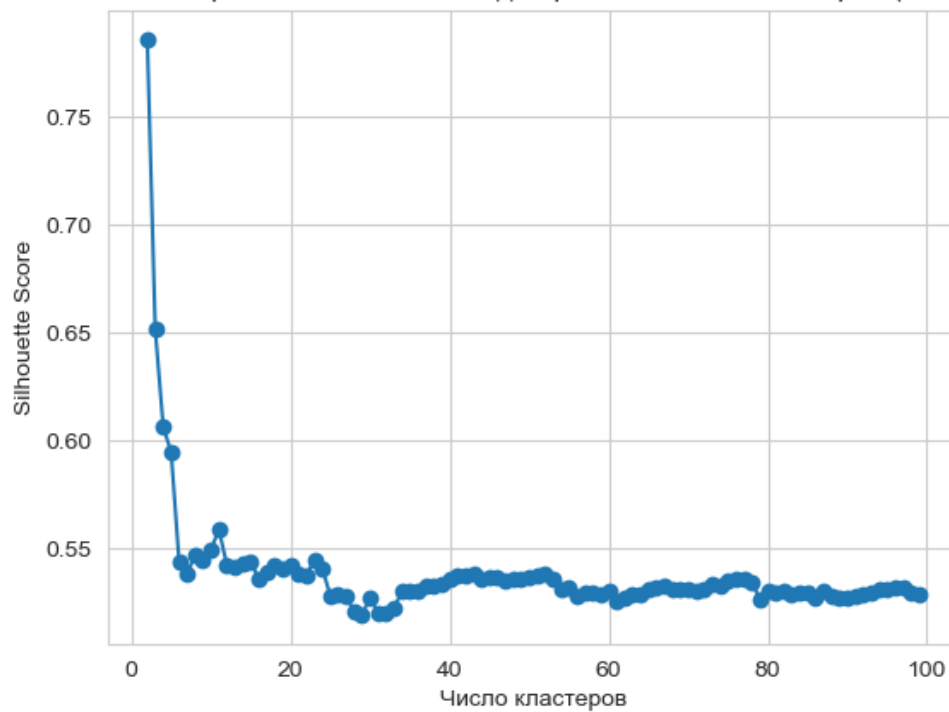
- метод к-средних
- древовидная классификация



Дендрограмма на основе древовидной классификации данных

Для продолжения поиска оптимального количества кластеров (хотя из полученной дендрограммы можно сделать вывод о диапазоне 2-5 кластеров) была проведена оценка разбиения на разное кол-во кластером с учетом silhouette оценки.

Изменение метрики Silhouette Score для разного числа кластеров (от 2 до 100)



Наилучшее значение получается на **двух кластерах**.

Кластеризация K-Means с понижением размерности признаков

