

# Демонстрация возможностей нотбука с Кластеризацией

## Первоначальный отбор признаков для предобработки

Необходимо ли отобразить признаки из всех? (y/n)

y

**Признаки, участвующие в кластеризации:**

- ☒ Client ID
- ☒ Genre
- ☒ Age
- ☒ Annual Income (k\$)
- ☒ Spending Score (1-100)
- ☐ Unknown feature
- ☐ Unknown feature 2
- ☒ category

## Выбор индекса

**Установите название столбца с индексами (id)**

Column:

## Обработка дат при наличии

**Обработка дат**

```
cb = date_preproc_widget(df)
```

В признаках присутствуют даты?(y/n)

y

Укажите столбцы с датой:

- ☐ Genre
- ☒ date\_test

```
df = date_preproc(df, cb, dayfirst = True)
```

Столбец date\_test

Что делать с датой?

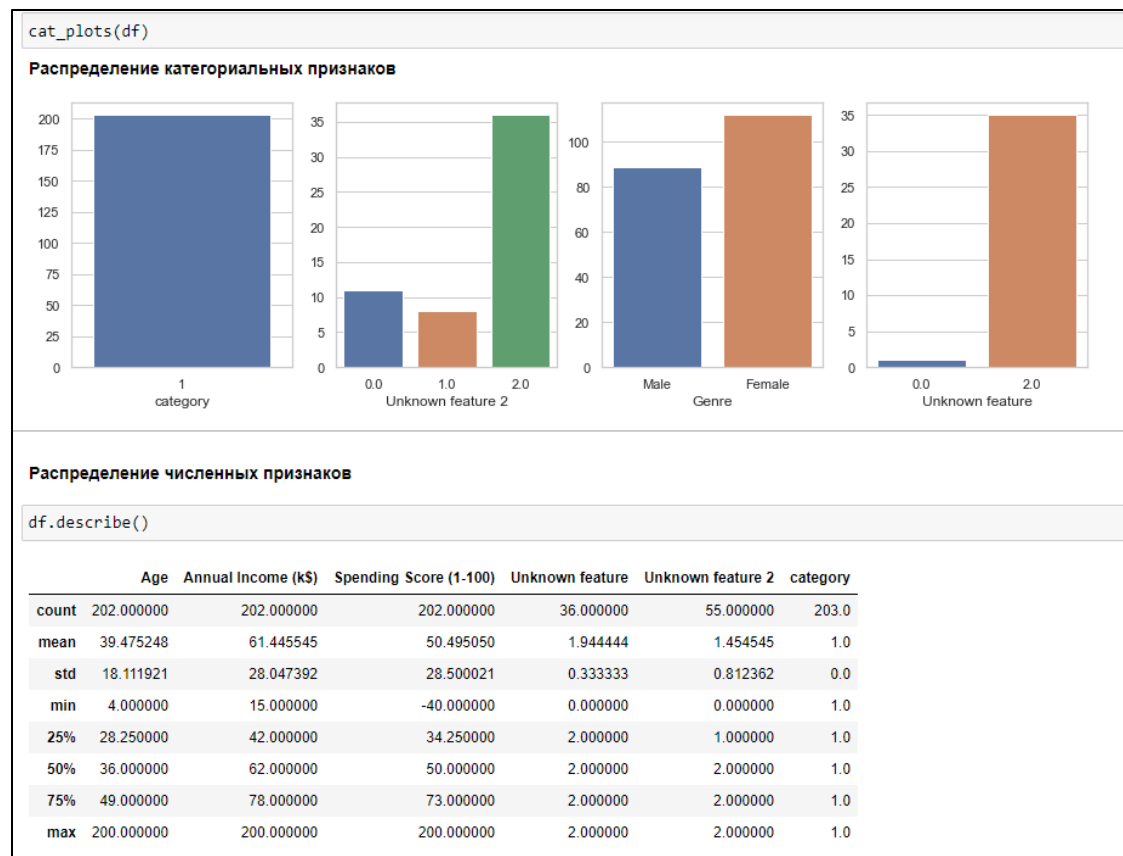
- Преобразовать
- Удалить

1

- Количество дней до текущей даты
- Количество дней до другой даты

1

## Вывод распределений признаков



## Удаление объектов с большим числом пропусков

**Пропуски**

```
df = drop_nan_rows(df)
```

Количество пустых значений в объектах:

	Число NA	Число объектов
0	6	1
1	2	147
2	1	20
3	0	35

Удалить объекты с большим числом NA?(y/n) y

**Выберите отсечение, по которому удалить объекты**

Удалить объекты, у которых число признаков с NA>=6

Количество оставшихся объектов: 202

## Удаление признаков с низкой вариативностью

### Удаление низковариативных признаков

```
df = drop_unvariative_cols(df)
```

Признаки с низким коэффициентом вариативности  $\frac{std}{mean}$  :

**category** : 0.000

Удалить признаки с низкой вариативностью? (y/n):

y

Удалить признаки, с коэфф. вариативности не больше: 0

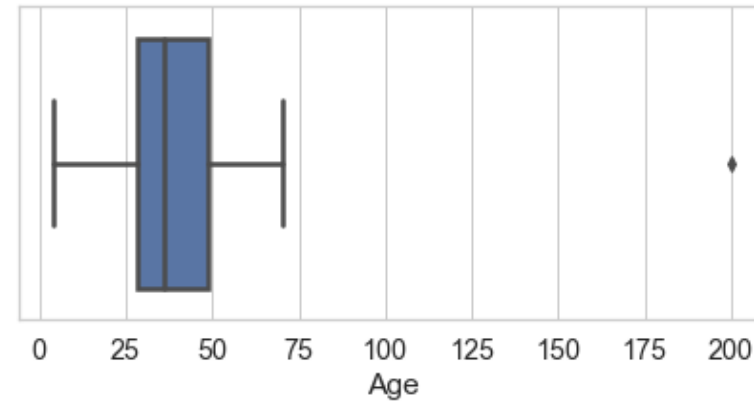
## Обработка аномалий (для каждого признака)

Выберите метод обработки аномалий для признаков

Age

Распределение признака "Age".

Аномальных значений: 1



Метод обработки аномальных значений:

1. Удалить
2. Ограничить

## Удаление признаков с большим числом пропусков

Количество пропусков в признаках:

	%NA	Тип
Unknown feature	82.09%	float64
Unknown feature 2	72.64%	float64
Genre	0.5%	object

Удалить признаки с большим процентом NA? (y/n) y

Удалить признаки, с % NA >= 50

## Выбор заполнения пропусков для остальных признаков

Выберите метод заполнения пропусков для каждого признака

Признак **Genre**

Распределение:

Female 0.56

Male 0.44

Name: Genre, dtype: float64

Заполнить NA:

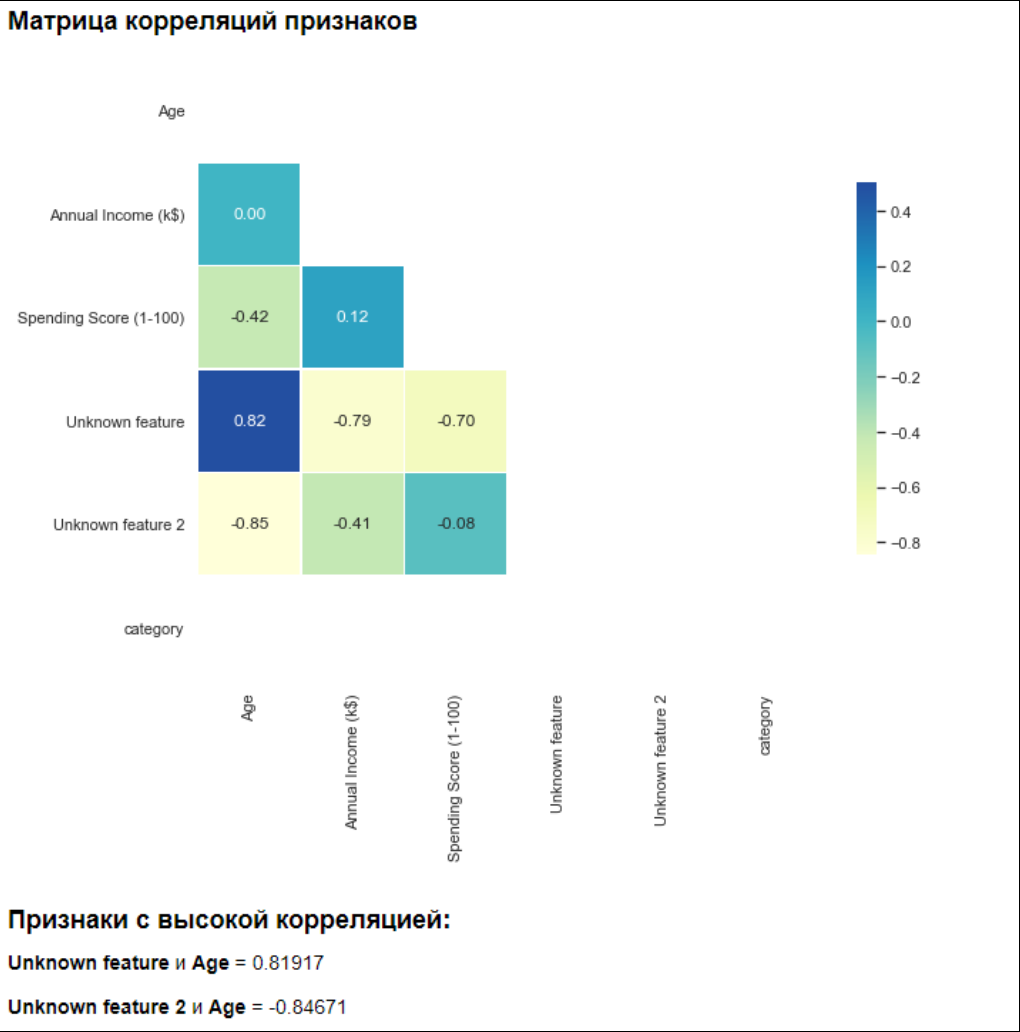
1.Female (Мода)

2.Male

3.Случайно (с весами из распределения)

4.Ручной ввод

# Построение матрицы корреляции



## Кодирование и масштабирование признаков

**Кодирование признаков**

```
df = encoding(df)
```

**Genre**, уникальных значений: 2

Способ кодирования:

- 1. LabelEncoder
- 2. OneHot
- 3. Удалить этот признак

Выбор: 1

**Масштабирование**

```
dfc = df.copy()
df = scaling(df)
```

Способ масштабирования данных:

- 1. Стандартизация
- 2. Нормализация [0, 1]
- 3. Не масштабировать

Выбор: 2

# Выбор метода для рассмотрения

Кластеризация

Выбор метода

method = choice\_method\_widget()

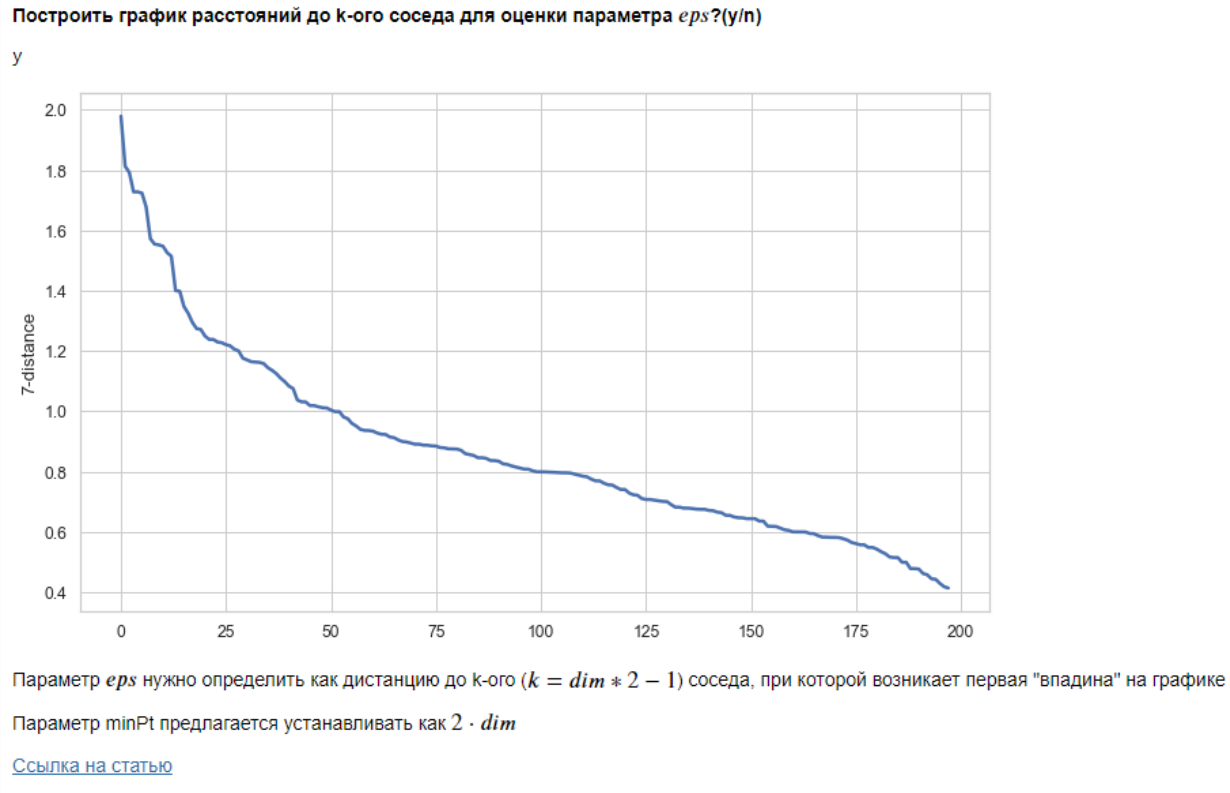
Выберите метод кластеризации

Метод: 

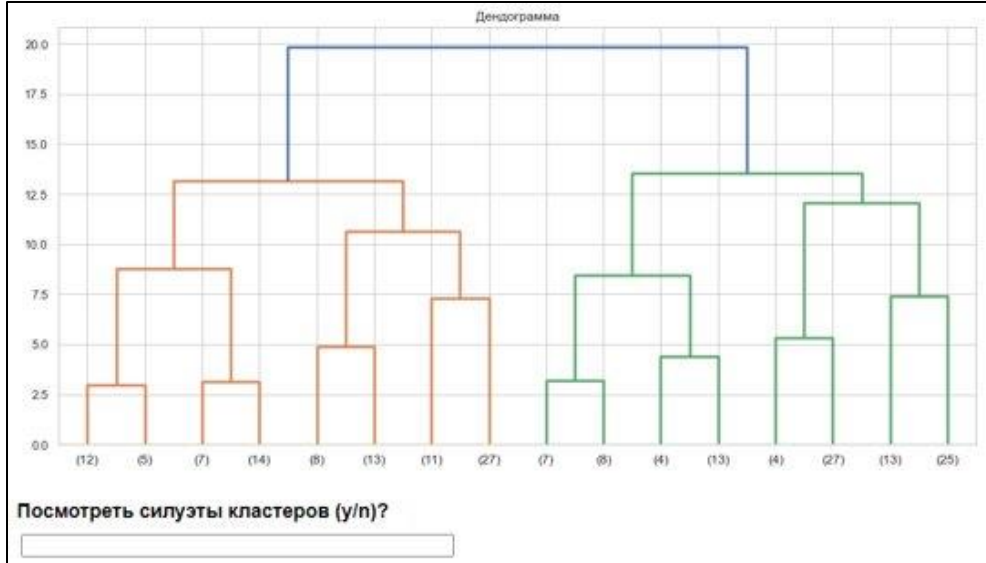
K-means  
Agglomerative  
DBSCAN

# Подбор параметров для методов кластеризации

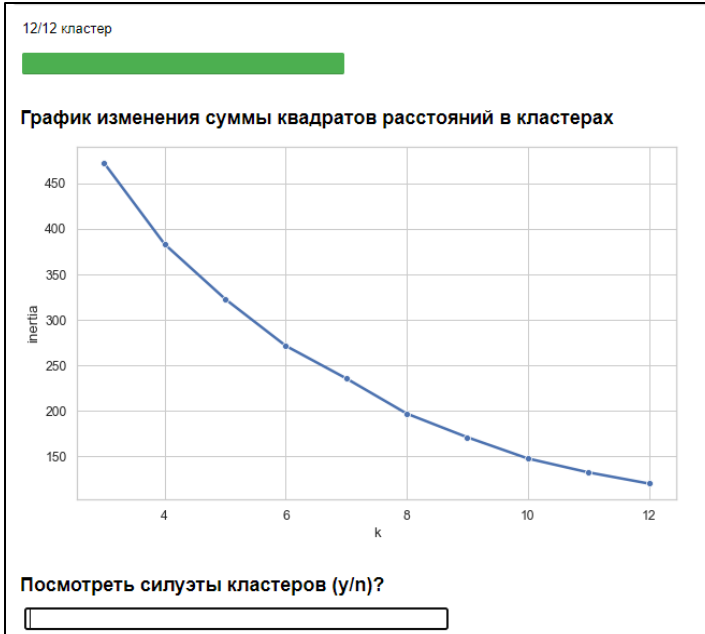
dbscan



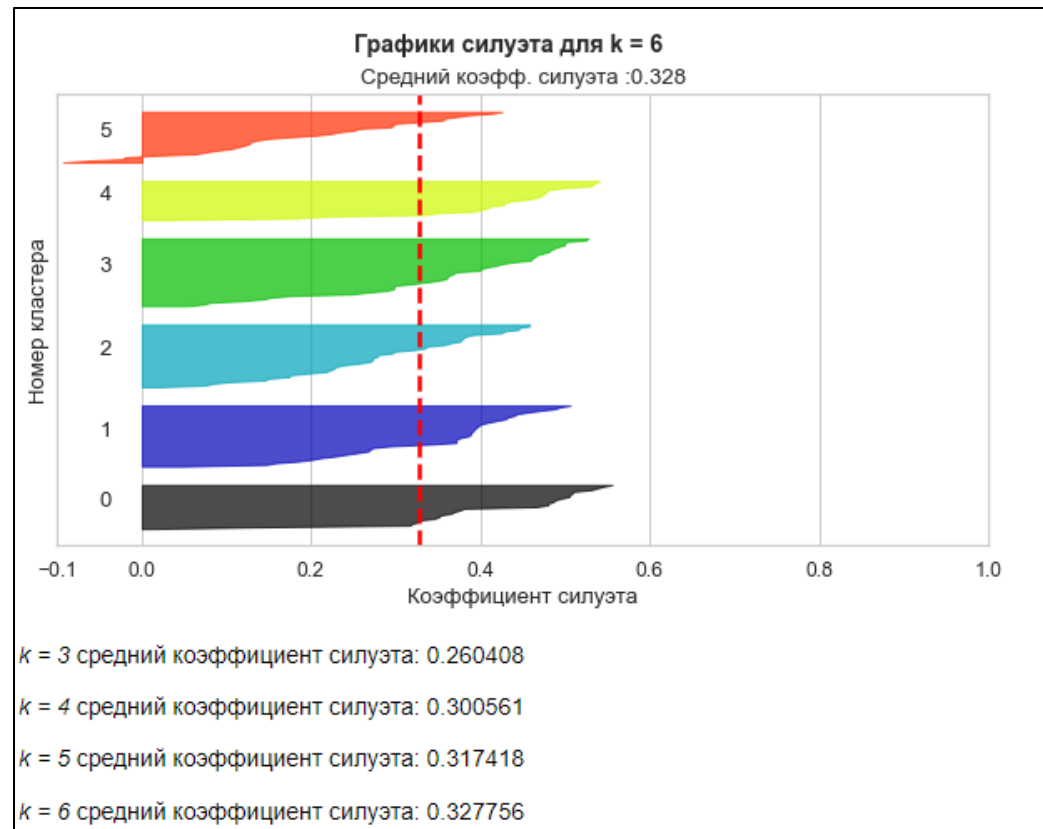
agglomerative



kmeans



## Построение силуэтов для различного количества кластеров

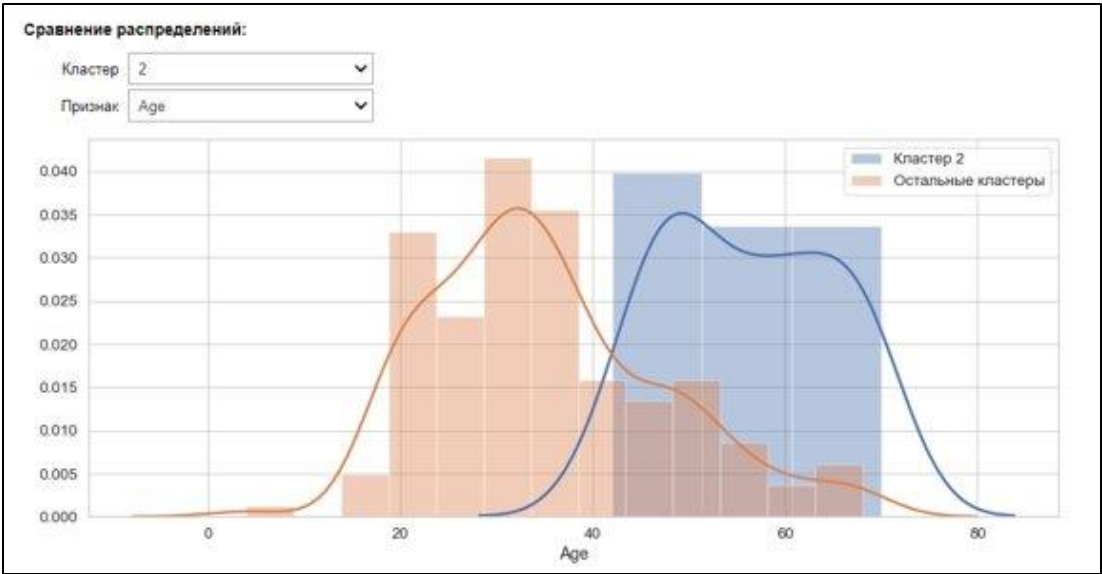




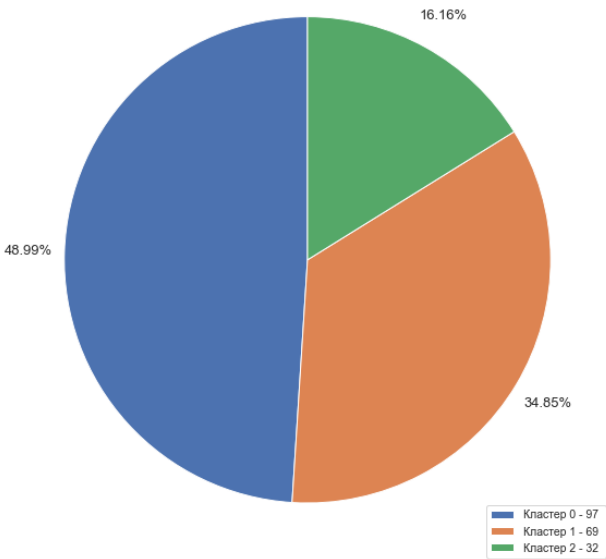
Построение характеристик кластеров

Средние значения характеристик в кластерах

	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
cluster				
0	0.391753	28.577320	60.185567	68.000000
1	0.449275	52.173913	46.333333	40.971014
2	0.531250	41.750000	87.593750	16.125000



Размеры кластеров



Выберите кластер для выделения на графике:

Кластер: 1

Кластер 0   Кластер 1   Кластер 2

