

Кодирование категориальных переменных

[Материалы](#) > [Анализ и обработка данных](#)

Алгоритмы машинного обучения, как мы знаем, не умеют работать с категориальными данными, выраженными с помощью строковых значений. Для этого строки необходимо **закодировать** (encode) числами. Сегодня мы рассмотрим основные способы такой кодировки.

Откроем блокнот к этому занятию

Подготовим простые учебные данные кредитного скоринга.

```
1 scoring = {
2     'Name' : ['Иван', 'Николай', 'Алексей', 'Александра', 'Евгений', 'Елена'],
3     'Age' : [35, 43, 21, 34, 24, 27],
4     'City' : ['Москва', 'Нижний Новгород', 'Санкт-Петербург', 'Владивосток', 'Москва',
5     'Experience' : [7, 13, 2, 8, 4, 12],
6     'Salary' : [95, 135, 73, 100, 78, 110],
7     'Credit_score' : ['Good', 'Good', 'Bad', 'Medium', 'Medium', 'Good'],
8     'Outcome' : [1, 1, 0, 1, 0, 1]
9 }
10
11 df = pd.DataFrame(scoring)
12 df
```

	Name	Age	City	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	7	95	Good	Вернул
1	Николай	43	Нижний Новгород	13	135	Good	Вернул
2	Алексей	21	Санкт-Петербург	2	73	Bad	Не вернул
3	Александра	34	Владивосток	8	100	Medium	Вернул
4	Евгений	24	Москва	4	78	Medium	Не вернул
5	Елена	27	Екатеринбург	12	110	Good	Вернул

Про категориальные переменные

Вначале в целом повторим, как выявлять и исследовать категориальные переменные.

Методы .info(), .unique(), value_counts()

Начать исследование категориальных переменных можно с изучения типа данных. Для этого

подойдут метод **.info()** или атрибут **dtypes**.

```
1 df.info()
```

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 6 entries, 0 to 5
3 Data columns (total 7 columns):
4  #   Column          Non-Null Count  Dtype
5  ---  ---
6  0   Name             6 non-null     object
7  1   Age              6 non-null     int64
8  2   City             6 non-null     object
9  3   Experience       6 non-null     int64
10  4   Salary           6 non-null     int64
11  5   Credit_score     6 non-null     object
12  6   Outcome          6 non-null     object
13 dtypes: int64(3), object(4)
14 memory usage: 464.0+ bytes
```

```
1 df.dtypes
```

```
1 Name             object
2 Age              int64
3 City             object
4 Experience       int64
5 Salary           int64
6 Credit_score     object
7 Outcome          object
8 dtype: object
```

При этом категориальные признаки часто могут «прятаться» в типах `int` и `float`. В этом случае для их выявления можно изучить распределение данных.

Отдельные категории можно посмотреть с помощью метода **.unique()**.

```
1 df.City.unique()
```

```
1 array(['Москва', 'Нижний Новгород', 'Санкт-Петербург', 'Владивосток',
2       'Екатеринбург'], dtype=object)
```

С помощью методов **.value_counts()** библиотеки Pandas и **np.unique()** библиотеки Numpy можно посмотреть и количество объектов в каждой категории.

```
1 # метод .value_counts() сортирует категории по количеству объектов
2 # в убывающем порядке
3 df.City.value_counts()
```

```
1 Москва          2
2 Нижний Новгород  1
3 Санкт-Петербург  1
4 Владивосток     1
5 Екатеринбург     1
6 Name: City, dtype: int64
```

```
1 np.unique(df.City, return_counts = True)
```

```
1 (array(['Владивосток', 'Екатеринбург', 'Москва', 'Нижний Новгород',
2       'Санкт-Петербург'], dtype=object), array([1, 1, 2, 1, 1]))
```

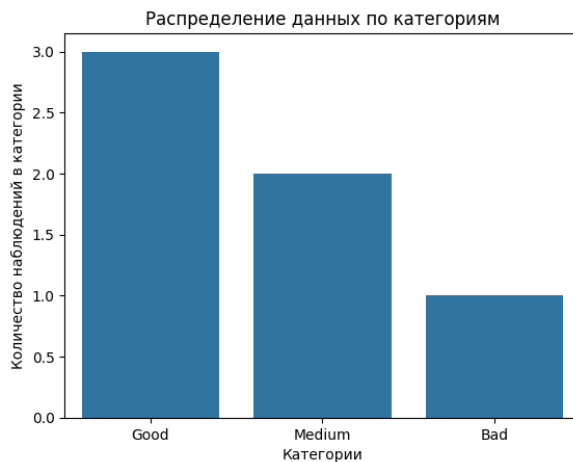
Последовательное применение методов **.value_counts()** и **.count()** выведет общее количество уникальных категорий.

```
1 # посмотрим на общее количество уникальных категорий
2 df.City.value_counts().count()
```

```
1 5
```

Выведем категории на графике.

```
1 score_counts = df.Credit_score.value_counts()
2 sns.barplot(x = score_counts.index, y = score_counts.values)
3 plt.title('Распределение данных по категориям')
4 plt.ylabel('Количество наблюдений в категории')
5 plt.xlabel('Категории');
```



Тип данных category

Хорошая практика — перевести категориальную переменную в тип данных category. Зачастую (например, если много категорий) это ускоряет работу с категориями и уменьшает использование памяти.

Можно воспользоваться уже знакомым нам **методом .astype()**.

```
1 df = df.astype({'City' : 'category', 'Outcome' : 'category'})
```

Функция pd.Categorical() позволяет прописать, в частности, сами категории, а также указать, есть ли в переданных категориях порядок или нет.

```
1 df.Credit_score = pd.Categorical(df.Credit_score,
2                                 categories = ['Bad', 'Medium', 'Good'],
3                                 ordered = True)
```

Воспользуемся атрибутами categories и dtype.

```
1 df.Credit_score.cat.categories
```

```
1 Index(['Bad', 'Medium', 'Good'], dtype='object')
```

```
1 df.Credit_score.dtype
```

```
1 CategoricalDtype(categories=['Bad', 'Medium', 'Good'], ordered=True, categories_dtype=c
```

Атрибут codes выводит коды каждой из категорий (мы воспользуемся этим в дальнейшем при кодировании).

```
1 df.Credit_score.cat.codes
```

```
1 0 2
2 1 2
3 2 0
4 3 1
5 4 1
```

```
6 | 5    2
7 | dtype: int8
```

Категории можно переименовать.

```
1 | df.Outcome = df.Outcome.cat.rename_categories(new_categories =
2 |                                             {'Вернул': 'Yes', 'Не вернул': 'No'})
3 |
4 | df
```

	Name	Age	City	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	7	95	Good	Yes
1	Николай	43	Нижний Новгород	13	135	Good	Yes
2	Алексей	21	Санкт-Петербург	2	73	Bad	No
3	Александра	34	Владивосток	8	100	Medium	Yes
4	Евгений	24	Москва	4	78	Medium	No
5	Елена	27	Екатеринбург	12	110	Good	Yes

Убедимся, что нужные нам признаки преобразованы в тип category.

```
1 | df.info()
```

```
1 | <class 'pandas.core.frame.DataFrame'>
2 | RangeIndex: 6 entries, 0 to 5
3 | Data columns (total 7 columns):
4 | #   Column      Non-Null Count  Dtype
5 | ---  ---
6 | 0    Name        6 non-null     object
7 | 1    Age         6 non-null     int64
8 | 2    City        6 non-null     category
9 | 3    Experience   6 non-null     int64
10 | 4    Salary      6 non-null     int64
11 | 5    Credit_score 6 non-null     category
12 | 6    Outcome     6 non-null     category
13 | dtypes: category(3), int64(3), object(1)
14 | memory usage: 806.0+ bytes
```

Кардинальность данных

Большое количество уникальных категорий в столбце называется **высокой кардинальностью** (high cardinality) признака. В частности, потенциально (если бы у нас было больше данных) признак City мог бы обладать высокой кардинальностью.

Высокая кардинальность, среди прочего, может привести к разреженности данных. Одним из способов преодоления этой сложности могло бы стать создание нового признака, например, региона, группирующего несколько городов.

```
1 | region = np.where(((df.City == 'Екатеринбург') | (df.City == 'Владивосток')), 0, 1)
2 |
3 | df.insert(loc = 3, column= 'Region', value = region)
4 |
5 | df
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	1	7	95	Good	Yes

1	Николай	43	Нижний Новгород	1	13	135	Good	Yes
2	Алексей	21	Санкт-Петербург	1	2	73	Bad	No
3	Александра	34	Владивосток	0	8	100	Medium	Yes
4	Евгений	24	Москва	1	4	78	Medium	No
5	Елена	27	Екатеринбург	0	12	110	Good	Yes

Дополнительным полезным свойством нового признака будет то, что на основе изначальных данных алгоритм бы не увидел разницы между Москвой и Владивостоком или Москвой и Екатеринбургом (а вполне вероятно, что в данных она есть). В новом же признаке, по сути делящем города по принадлежности к европейской и азиатской частям России, такую разницу выявить получится.

Базовые методы кодирования

Кодирование через cat.codes

Как уже было сказано выше, кодировать категориальную переменную можно через атрибут `cat.codes`.

```
1 df_cat = df.copy()
2 df_cat.Credit_score.cat.codes
```

```
1 0 2
2 1 2
3 2 0
4 3 1
5 4 1
6 5 2
7 dtype: int8
```

```
1 df_cat.Credit_score = df_cat.Credit_score.astype('category').cat.codes
2 df_cat
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	1	7	95	2	Yes
1	Николай	43	Нижний Новгород	1	13	135	2	Yes
2	Алексей	21	Санкт-Петербург	1	2	73	0	No
3	Александра	34	Владивосток	0	8	100	1	Yes
4	Евгений	24	Москва	1	4	78	1	No
5	Елена	27	Екатеринбург	0	12	110	2	Yes

Mapping

Этот способ мы уже применяли на прошлых занятиях. Суть его заключается в том, чтобы передать схему кодирования в виде словаря в **функцию `map()`** и применить к соответствующему столбцу.

```
1 df_map = df.copy()
2
3 # ключами будут старые значения признака
```

```
4 # значениями словаря - новые значения признака
5 map_dict = {'Bad' : 0,
6             'Medium' : 1,
7             'Good' : 2}
8
9 df_map['Credit_score'] = df_map['Credit_score'].map(map_dict)
10 df_map
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	1	7	95	2	Yes
1	Николай	43	Нижний Новгород	1	13	135	2	Yes
2	Алексей	21	Санкт-Петербург	1	2	73	0	No
3	Александра	34	Владивосток	0	8	100	1	Yes
4	Евгений	24	Москва	1	4	78	1	No
5	Елена	27	Екатеринбург	0	12	110	2	Yes

Словарь в функцию **map()** можно передать и так.

```
1 # сделаем еще одну копию датафрейма
2 df_map = df.copy()
3
4 df_map.Credit_score = df_map.Credit_score.map(dict(Bad = 0, Medium = 1, Good = 2))
5 df_map
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	1	7	95	2	Yes
1	Николай	43	Нижний Новгород	1	13	135	2	Yes
2	Алексей	21	Санкт-Петербург	1	2	73	0	No
3	Александра	34	Владивосток	0	8	100	1	Yes
4	Евгений	24	Москва	1	4	78	1	No
5	Елена	27	Екатеринбург	0	12	110	2	Yes

Label Encoder

Рассмотрим класс `LabelEncoder` библиотеки `sklearn`. Этот класс преобразует n категорий в числа от 1 до n . Применим его к целевой переменной (бинарная категориальная переменная).

На вход `LabelEncoder` принимает только одномерные массивы (например, `Series`)

```
1 from sklearn.preprocessing import LabelEncoder
2
3 labelencoder = LabelEncoder()
4
5 df_le = df.copy()
6
7 df_le.loc[:, 'Outcome'] = labelencoder.fit_transform(df_le.loc[:, 'Outcome'])
8 df_le
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	1	7	95	Good	1
1	Николай	43	Нижний Новгород	1	13	135	Good	1
2	Алексей	21	Санкт-Петербург	1	2	73	Bad	0
3	Александра	34	Владивосток	0	8	100	Medium	1
4	Евгений	24	Москва	1	4	78	Medium	1
5	Елена	27	Екатеринбург	0	12	110	Good	1

1	Николай	43	Пижмий повород	1	13	135	Good	1
2	Алексей	21	Санкт-Петербург	1	2	73	Bad	0
3	Александра	34	Владивосток	0	8	100	Medium	1
4	Евгений	24	Москва	1	4	78	Medium	0
5	Елена	27	Екатеринбург	0	12	110	Good	1

Для категорий, в которых больше двух классов, но нет внутренней иерархии (номинальные данные), этот encoder подходит хуже, потому что построенная на основе преобразованных данных модель может подумать, что между категориями есть иерархия, когда в действительности ее нет.

```
1 # применим LabelEncoder к номинальной переменной City
2 df_le.loc[:, 'City'] = labelencoder.fit_transform(df_le.loc[:, 'City'])
3 df_le
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	2	1	7	95	Good	1
1	Николай	43	3	1	13	135	Good	1
2	Алексей	21	4	1	2	73	Bad	0
3	Александра	34	0	0	8	100	Medium	1
4	Евгений	24	2	1	4	78	Medium	0
5	Елена	27	1	0	12	110	Good	1

При этом даже для порядковых категориальных данных этот способ вряд ли подойдет, потому что LabelEncoder не видит порядка в данных.

```
1 # применим LabelEncoder к номинальной переменной Credit_score
2 df_le.loc[:, 'Credit_score'] = labelencoder.fit_transform(df_le.loc[:, 'Credit_score'])
3 df_le
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	2	1	7	95	1	1
1	Николай	43	3	1	13	135	1	1
2	Алексей	21	4	1	2	73	0	0
3	Александра	34	0	0	8	100	2	1
4	Евгений	24	2	1	4	78	2	0
5	Елена	27	1	0	12	110	1	1

```
1 labelencoder.classes_
```

```
1 array(['Bad', 'Good', 'Medium'], dtype=object)
```

Как мы видим, на второе место в иерархии категорий LabelEncoder поместил Good, что конечно является ошибкой. Таким образом, можно сказать, что LabelEncoder лучше всего справляется с бинарными категориальными данными.

Ordinal Encoder

С порядковыми категориальными данными справится OrdinalEncoder, которому при создании

объекта класса можно передать иерархию категорий.

На вход `OrdinalEncoder` принимает только двумерные массивы.

```
1 from sklearn.preprocessing import OrdinalEncoder
2
3 ordinalencoder = OrdinalEncoder(categories = [['Bad', 'Medium', 'Good']])
4
5 df_oe = df.copy()
6
7 # используем метод .to_frame() для преобразования Series в датафрейм
8 df_oe.loc[:, 'Credit_score'] = ordinalencoder.fit_transform(df_oe.loc[:, 'Credit_score'])
9 df_oe
```

	Name	Age	City	Region	Experience	Salary	Credit_score	Outcome
0	Иван	35	Москва	1	7	95	2.0	Yes
1	Николай	43	Нижний Новгород	1	13	135	2.0	Yes
2	Алексей	21	Санкт-Петербург	1	2	73	0.0	No
3	Александра	34	Владивосток	0	8	100	1.0	Yes
4	Евгений	24	Москва	1	4	78	1.0	No
5	Елена	27	Екатеринбург	0	12	110	2.0	Yes

Убедимся, что иерархия категорий не нарушена.

```
1 ordinalencoder.categories_
```

```
1 [array(['Bad', 'Medium', 'Good'], dtype=object)]
```

OneHotEncoding

Как уже было сказано, номинальные данные нельзя заменять числами $1, 2, 3, \dots$, так как алгоритм ML на этапе обучения подумает, что речь идет о порядковых данных. Нужно использовать one-hot encoder. С этим инструментом мы уже познакомились, когда рассматривали [основы нейронных сетей](#).

Класс OneHotEncoder

Вначале применим класс `OneHotEncoder` библиотеки `sklearn`.

```
1 df_onehot = df.copy()
2
3 from sklearn.preprocessing import OneHotEncoder
4
5 # создадим объект класса OneHotEncoder
6 # параметр sparse = True выдал бы результат в сжатом формате
7 onehotencoder = OneHotEncoder(sparse_output = False)
8
9 encoded_df = pd.DataFrame(onehotencoder.fit_transform(df_onehot[['City']]))
10 encoded_df
```

0 1 2 3 4

0	0.0	0.0	1.0	0.0	0.0
1	0.0	0.0	0.0	1.0	0.0
2	0.0	0.0	0.0	0.0	1.0
3	1.0	0.0	0.0	0.0	0.0
4	0.0	0.0	1.0	0.0	0.0
5	0.0	1.0	0.0	0.0	0.0

Выведем новые признаки с помощью метода `.get_feature_names_out()`.

```
1 onehotencoder.get_feature_names_out()

1 array(['City_Владивосток', 'City_Екатеринбург', 'City_Москва',
2        'City_Нижний Новгород', 'City_Санкт-Петербург'], dtype=object)
```

Используем вывод этого метода, чтобы добавить названия столбцов.

```
1 encoded_df.columns = onehotencoder.get_feature_names_out()
2 encoded_df
```

	City_Владивосток	City_Екатеринбург	City_Москва	City_Нижний Новгород	City_Санкт-Петербург
0	0.0	0.0	1.0	0.0	0.0
1	0.0	0.0	0.0	1.0	0.0
2	0.0	0.0	0.0	0.0	1.0
3	1.0	0.0	0.0	0.0	0.0
4	0.0	0.0	1.0	0.0	0.0
5	0.0	1.0	0.0	0.0	0.0

Присоединим новые признаки к исходному датафрейму, удалив, разумеется, признак City.

```
1 df_onehot = df_onehot.join(encoded_df)
2 df_onehot.drop('City', axis = 1, inplace = True)
```

Обратите внимание, на самом деле нам не нужен первый признак (в данном случае, Владивосток). Если его убрать, при «срабатывании» этого признака (наблюдение с индексом три) все остальные признаки будут иметь нули (так мы поймем, что речь идет именно об этом отсутствующем признаке).

```
1 df_onehot = df.copy()
2
3 # чтобы удалить первый признак, используем параметр drop = 'first'
4 onehot_first = OneHotEncoder(drop = 'first', sparse = False)
5
6 encoded_df = pd.DataFrame(onehot_first.fit_transform(df_onehot[['City']]))
7 encoded_df.columns = onehot_first.get_feature_names_out()
8
9 df_onehot = df_onehot.join(encoded_df)
10 df_onehot.drop('Outcome', axis = 1, inplace = True)
11 df_onehot
```

	Name	Age	City	Region	Experience	Salary	Credit_score	City_Екатеринбург	City_Москва	City_Нижний Новгород	City_Санкт-Петербург
0	Иван	35	Москва	1	7	95	Good	0.0	1.0	0.0	0.0
1	Николай	43	Нижний Новгород	1	13	135	Good	0.0	0.0	1.0	0.0

2	Алексей	21	Санкт-Петербург	1	2	73	Bad	0.0	0.0	0.0	1.0
3	Александра	34	Владивосток	0	8	100	Medium	0.0	0.0	0.0	0.0
4	Евгений	24	Москва	1	4	78	Medium	0.0	1.0	0.0	0.0
5	Елена	27	Екатеринбург	0	12	110	Good	1.0	0.0	0.0	0.0

Функция `pd.get_dummies()`

Еще один способ — использовать функцию `pd.get_dummies()` библиотеки Pandas. Применим функцию к столбцу `City`.

```
1 df_dum = df.copy()
2 pd.get_dummies(df_dum, columns = ['City'])
```

City_Владивосток	City_Екатеринбург	City_Москва	City_Нижний Новгород	City_Санкт-Петербург
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1
1	0	0	0	0
0	0	1	0	0
0	1	0	0	0

Уменьшить ширину новых столбцов можно через параметры `prefix` и `prefix_sep`.

```
1 pd.get_dummies(df_dum, columns = ['City'], prefix = '', prefix_sep = '')
```

Владивосток	Екатеринбург	Москва	Нижний Новгород	Санкт-Петербург
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1
1	0	0	0	0
0	0	1	0	0
0	1	0	0	0

Опять же, можно не использовать первую dummy-переменную.

```
1 pd.get_dummies(df_dum, columns = ['City'],
2                 prefix = '',
3                 prefix_sep = '',
4                 drop_first = True)
```

Екатеринбург	Москва	Нижний Новгород	Санкт-Петербург
0	1	0	0
0	0	1	0
0	0	0	1

0	0	0	0
0	1	0	0
1	0	0	0

Библиотека `category_encoders`

Рассмотрим еще один способ выполнить one-hot encoding через соответствующий инструмент [в](#) очень полезной библиотеки `category_encoders`.

```
1 # установим библиотеку
2 !pip install category_encoders
```

Импортируем библиотеку и применим класс `OneHotEncoder`.

```
1 df_catenc = df.copy()
2
3 import category_encoders as ce
4
5 # в параметр cols передадим столбцы, которые нужно преобразовать
6 ohe_encoder = ce.OneHotEncoder(cols = ['City'])
7 # в метод .fit_transform() мы передадим весь датафрейм целиком
8 df_catenc = ohe_encoder.fit_transform(df_catenc)
9 df_catenc
```

	Name	Age	City_1	City_2	City_3	City_4	City_5	Region	Experience	Salary
0	Иван	35	1	0	0	0	0	1	7	95
1	Николай	43	0	1	0	0	0	1	13	135
2	Алексей	21	0	0	1	0	0	1	2	73
3	Александра	34	0	0	0	1	0	0	8	100
4	Евгений	24	1	0	0	0	0	1	4	78
5	Елена	27	0	0	0	0	1	0	12	110

Что очень удобно, класс `OneHotEncoder` библиотеки `category_encoders` вставил новые столбцы сразу в изначальный датафрейм и удалил исходный признак.

Сравнение инструментов

Создадим два очень простых датасета из одного признака: один обучающий, второй — тестовый. В первом датасете в этом признаке (назовем его `recom`) будет три категории: `yes`, `no`, `maybe`. Во втором, только две, `yes` и `no`.

```
1 train = pd.DataFrame({'recom' : ['yes', 'no', 'maybe']})
2 train
```

	recom
0	yes
1	no
2	maybe

```
1 test = pd.DataFrame({'recom' : ['yes', 'no', 'yes']})
2 test
```

	recom
0	yes
1	no
2	yes

Теперь применим каждый из приведенных выше инструментов к этим датасетам (напомню, что обучать кодировщик мы должны на обучающей выборке, чтобы избежать утечки данных).

pd.get_dummies()

Функция **pd.get_dummies()** не «запоминает» категории при обучении.

```
1 pd.get_dummies(train)
```

	recom_maybe	recom_no	recom_yes
0	False	False	True
1	False	True	False
2	True	False	False

```
1 pd.get_dummies(test)
```

	recom_no	recom_yes
0	False	True
1	True	False
2	False	True

При попытке обучить модель будет ошибка.

OHE sklearn

Посмотрим, как с этим справится класс OneHotEncoder библиотеки sklearn.

```
1 ohe = OneHotEncoder()
2 ohe_model = ohe.fit(train)
3 ohe_model.categories_
```

```
1 [array(['maybe', 'no', 'yes'], dtype=object)]
```

```
1 train_arr = ohe_model.transform(train).toarray()
2 pd.DataFrame(train_arr, columns = ['maybe', 'no', 'yes'])
```

	maybe	no	yes
0	0.0	0.0	1.0
1	0.0	1.0	0.0
2	1.0	0.0	0.0

```

1 test_arr = ohe_model.transform(test).toarray()
2 pd.DataFrame(test_arr, columns = ['maybe', 'no', 'yes'])

```

	maybe	no	yes
0	0.0	0.0	1.0
1	0.0	1.0	0.0
2	0.0	0.0	1.0

Мы видим, что этот кодировщик учел отсутствующую в тестовой выборке категорию. Впрочем, в обратном случае, когда категория отсутствует в обучающей выборке, OneHotEncoder не будет иметь возможности правильно закодировать датасеты.

```

1 ohe = OneHotEncoder()
2 ohe_model = ohe.fit(test)
3 ohe_model.categories_

1 [array(['no', 'yes'], dtype=object)]

```

ONE category_encoders

Попробуем инструмент из библиотеки category_encoders.

```

1 ohe_encoder = ce.OneHotEncoder()
2 ohe_encoder.fit(train)

1 OneHotEncoder(cols=['recom'])

1 # категория maybe стоит на последнем месте
2 ohe_encoder.transform(test)

```

	recom_1	recom_2	recom_3
0	1	0	0
1	0	1	0
2	1	0	0

```

1 # убедимся в этом, добавив названия столбцов
2 test_df = ohe_encoder.transform(test)
3 test_df.columns = ohe_encoder.category_mapping[0]['mapping'].index[:3]
4 test_df

```

	yes	no	maybe
0	1	0	0
1	0	1	0
2	1	0	0

Binning

Некоторые (в частности, мультимодальные) количественные распределения не поддаются трансформации и приведению, например, к нормальному распределению.

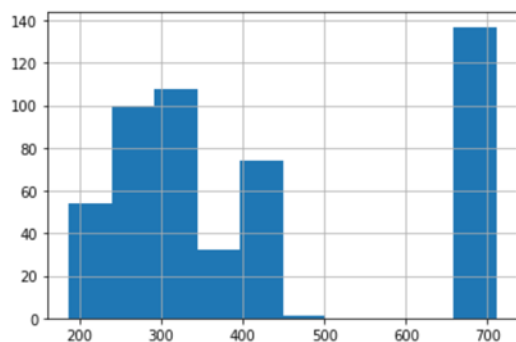
Для того чтобы извлечь ценную информацию из таких признаков можно попробовать сделать переменные категориальными, разбив данные на интервалы, которые и будут классами нового признака. Такой подход называется **binning** или **bucketing**.

Вновь обратимся к датасету о недвижимости в Бостоне и, в частности, рассмотрим переменную TAX.

boston.csv

[Скачать](#)

```
1 boston = pd.read_csv('/content/boston.csv')
2 boston.TAX.hist();
```



Как мы видим, распределение вряд ли можно трансформировать, используя какое-либо преобразование. Применим binning.

На равные интервалы

Подход **binning на равные интервалы** (binning with equally spaced boundaries) предполагает, что мы берем диапазон от минимального до максимального значений и делим его на нужное нам количество *равных* частей (если мы хотим получить три интервала, то нам нужно четыре границы).

```
1 min_value = boston.TAX.min()
2 max_value = boston.TAX.max()
3
4 bins = np.linspace(min_value, max_value, 4)
5 bins
```

```
1 array([187.          , 361.66666667, 536.33333333, 711.          ])
```

Создадим названия категорий.

```
1 labels = ['low', 'medium', 'high']
```

Применим **функцию pd.cut()**. В параметр bins мы передадим интервалы, в labels — названия категорий.

```
1 boston['TAX_binned'] = pd.cut(boston.TAX,
2                               bins = bins,
```

```
3 labels = labels,  
4 # уточним, что первый интервал должен включать  
5 # нижнюю границу (значение 187)  
6 include_lowest = True)
```

Посмотрим на результат.

```
1 boston[['TAX', 'TAX_binned']].sample(5, random_state = 42)
```

	TAX	TAX_binned
173	296.0	low
274	254.0	low
491	711.0	high
72	305.0	low
452	666.0	high

Границы и количество элементов в них можно получить с помощью метода `.value_counts()`.

```
1 boston.TAX.value_counts(bins = 3, sort = False)
```

```
1 (186.475, 361.667]    273  
2 (361.667, 536.333]    96  
3 (536.333, 711.0]     137  
4 Name: TAX, dtype: int64
```

Результат этого метода позволяет выявить недостаток подхода binning на равные интервалы. Количество объектов внутри интервалов сильно различается. Преодолеть эту особенность можно с помощью деления по квантилям.

По квантилям

Binning по квантилям (quantile binning) позволяет разделить наблюдения не по значениям признака, а по количеству объектов в интервале. Например, выберем разделение на три части.

```
1 # для наглядности вначале найдем интересующие нас квантили  
2 np.quantile(boston.TAX, q = [1/3, 2/3])
```

```
1 array([300., 403.])
```

Применим функцию `pd.qcut()`.

```
1 boston['TAX_qbinned'], boundaries = pd.qcut(boston.TAX,  
2                                             q = 3,  
3                                             # precision определяет округление  
4                                             precision = 1,  
5                                             labels = labels,  
6                                             retbins = True)  
7  
8 boundaries
```

```
1 array([187., 300., 403., 711.])
```

	TAX	TAX_qbinned
173	296.0	low

274	254.0	low
491	711.0	high
72	305.0	medium
452	666.0	high

```
1 boston.TAX_qbinned.value_counts()
```

```
1 low      172
2 high     168
3 medium   166
4 Name: TAX_qbinned, dtype: int64
```

Как вы видите, в данном случае количество объектов примерно одинаковое. Наглядную иллюстрацию двух подходов можно [посмотреть здесь](#).

KBinsDiscretizer

Эти же задачи можно решить с помощью [класса KBinsDiscretizer](#) библиотеки sklearn.

Рассмотрим три основных параметра класса:

- параметр **strategy** определяет, как будут делиться интервалы:
 - на равные части (uniform);
 - по квантилям (quantile); или
 - так, чтобы значения в каждом кластере относились к центроиду (kmeans);
- параметр **encode** определяет, как закодировать интервалы:
 - ordinal, т.е. числами от 1 до n интервалов; или
 - one-hot encoding;
- количество интервалов **n_bins**.

Применим каждую из стратегий. Так как в категориях заложен порядок, выберем ordinal кодировку.

```
1 from sklearn.preprocessing import KBinsDiscretizer
```

strategy = uniform

```
1 est = KBinsDiscretizer(n_bins = 3, encode = 'ordinal',
2                        strategy = 'uniform', subsample = None)
3
4 est.fit(boston[['TAX']])
5 est.bin_edges_
```

```
1 array([array([187.          , 361.66666667, 536.33333333, 711.          ])],
2       dtype=object)
```

```
1 np.unique(est.transform(boston[['TAX']]), return_counts = True)
```

```
1 (array([0., 1., 2.]), array([273, 96, 137]))
```


strategy = quantile

```

1 est = KBinsDiscretizer(n_bins = 3, encode = 'ordinal', strategy = 'quantile')
2 est.fit(boston[['TAX']])
3 est.bin_edges_

1 array([array([187., 300., 403., 711.]), dtype=object)

1 np.unique(est.transform(boston[['TAX']]), return_counts = True)

1 (array([0., 1., 2.]), array([165, 143, 198]))

```

strategy = kmeans

```

1 est = KBinsDiscretizer(n_bins = 3, encode = 'ordinal',
2                         strategy = 'kmeans', subsample = None)
3
4 est.fit(boston[['TAX']])
5 est.bin_edges_

1 array([array([187.          , 338.7198937 , 535.07350433, 711.          ])],
2        dtype=object)

1 np.unique(est.transform(boston[['TAX']]), return_counts = True)

1 (array([0., 1., 2.]), array([262, 107, 137]))

```

Еще одно сравнение стратегий разделения на интервалы можно посмотреть [здесь](#).

С помощью статистических показателей

Дополнительно замечу, что интервалы можно заполнить каким-либо статистическим показателем. Например, медианой. Для наглядности снова создадим только три интервала и найдем медианное значение внутри каждого из них.

Воспользуемся [функцией `binned_statistic\(\)`](#) модуля `scipy.stats`. Функция возвращает медианы каждого из интервалов, границы интервалов, а также к какому из интервалов относится каждое из наблюдений.

Нас будут интересовать метрики и границы интервалов.

```

1 from scipy.stats import binned_statistic
2
3 medians, bin_edges, _ = binned_statistic(boston.TAX,
4                                          np.arange(0, len(boston)),
5                                          statistic = 'median',
6                                          bins = 3)
7
8 medians, bin_edges

1 (array([216. , 147.5, 424. ]),
2  array([187.          , 361.66666667, 536.33333333, 711.          ]))

```

Подставим эти значения в функцию `pd.cut()`.

```

1 boston['TAX_binned_median'] = pd.cut(boston.TAX,
2                                     bins = bin_edges,
3                                     labels = medians,
4                                     include_lowest = True)

```

```
5 |
6 | boston['TAX_binned_median'].value_counts()

1 | 216.0    273
2 | 424.0    137
3 | 147.5     96
4 | Name: TAX_binned_median, dtype: int64
```

Алгоритм Дженкса

Алгоритм естественных границ Дженкса (Jenks natural breaks optimization) делит данные на группы (кластеры) таким образом, чтобы минимизировать отклонение наблюдений от среднего каждого класса (дисперсию внутри классов) и максимизировать отклонение среднего каждого класса от среднего других классов (дисперсию между классами).

Этот алгоритм также можно использовать для определения границ интервалов. Установим библиотеку `jenkspy`.

```
1 | !pip install jenkspy
```

Найдем оптимальные границы. Количество интервалов (`n_classes`) нужно по-прежнему указывать вручную.

```
1 | import jenkspy
2 |
3 | breaks = jenkspy.jenks_breaks(boston.TAX, n_classes = 3)
4 | breaks

1 | [187.0, 337.0, 469.0, 711.0]
```

Подставим интервалы в функцию `pd.cut()`.

```
1 | boston['TAX_binned_jenks'] = pd.cut(boston.TAX,
2 |                                   bins = breaks,
3 |                                   labels = labels,
4 |                                   include_lowest = True)
5 |
6 | boston['TAX_binned_jenks'].value_counts()

1 | low      262
2 | high     137
3 | medium   107
4 | Name: TAX_binned_jenks, dtype: int64
```

Подведем итог

Сегодня мы рассмотрели базовые методы кодирования категориальных переменных, а также стратегии binning/bucketing.