

**Московский государственный технический
университет им. Н. Э. Баумана**

Отчёт по лабораторной работе №2 по курсу «Технологии машинного
обучения».

«Обработка пропусков в данных, кодирование категориальных признаков,
масштабирование данных».

Выполнил:
Головацкий А. Д.
студент группы ИУ5-61Б

Проверил:
Галанюк Ю. Е.

Подпись и дата:

Подпись и дата:

Москва, 2022 г.

Обработка пропусков в данных, кодирование категориальных признаков, масштабирование данных.

Выбрать набор данных (датасет), содержащий категориальные признаки и пропуски в данных. Для выполнения следующих пунктов можно использовать несколько различных наборов данных (один для обработки пропусков, другой для категориальных признаков и т.д.)

Для выбранного датасета (датасетов) на основе материалов лекции решить следующие задачи: обработку пропусков в данных; кодирование категориальных признаков; масштабирование данных.

Выбор и загрузка данных

В качестве датасета будем использовать набор данных, содержащий данные по продажам автомобилей. А также его модифицированную версию, содержащую пропуски.

Набор данных имеет следующие атрибуты:

- 1) manufacturer - марка
- 2) model - модель
- 3) sales_in_thousands - продажи в тысячах
- 4) year_resale_value - годовой объем продаж
- 5) vehicle_type - тип автомобиля
- 6) price_in_thousands - цена в тысячах
- 7) engine_size - объем двигателя
- 8) horsepower - лошадиные силы
- 9) wheelbase - колесная база
- 10) width - ширина
- 11) length - длина
- 12) curb_weight - масса
- 13) fuel_capacity - топливный бак

14) fuel_efficiency - расход топлива

15) latest_Launch - начало производства модели

16) power_perf_factor - мощностной коэффициент

Импорт библиотек

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(style="ticks")
```

Загрузка данных

```
data = pd.read_csv('C:\\Users\\Andrew\\Anaconda Projects\\datasets\\
car_sales.csv')
data_mod = pd.read_csv('C:\\Users\\Andrew\\Anaconda Projects\\
datasets\\car_sales_passes.csv')
```

Первичный анализ данных

Выведем первые 5 строк датасета:

```
data.head()
```

	Manufacturer	Model	Sales_in_thousands	__year_resale_value
0	Acura	Integra	16.919	16.360
1	Acura	TL	39.384	19.875
2	Acura	CL	14.114	18.225
3	Acura	RL	8.588	29.725
4	Audi	A4	20.397	22.255

	Price_in_thousands	Engine_size	Horsepower	Wheelbase	Width
0	21.50	1.8	140.0	101.2	67.3
1	28.40	3.2	225.0	108.1	70.3
2	NaN	3.2	225.0	106.9	70.6
3	42.00	3.5	210.0	114.6	71.4
4	23.99	1.8	150.0	102.6	68.2

	Curb_weight	Fuel_capacity	Fuel_efficiency	Latest_Launch	\
0	2.639	13.2	28.0	2/2/2012	
1	3.517	17.2	25.0	6/3/2011	
2	3.470	17.2	26.0	1/4/2012	
3	3.850	18.0	22.0	3/10/2011	
4	2.998	16.4	27.0	10/8/2011	

	Power_perf_factor
0	58.280150
1	91.370778
2	NaN
3	91.389779
4	62.777639

Определим размер датасета:

```
data.shape
```

```
(157, 16)
```

В датасете 157 строк и 16 столбцов. Определим типы столбцов:

```
data.dtypes
```

```
Manufacturer      object
Model             object
Sales_in_thousands  float64
__year_resale_value float64
Vehicle_type      object
Price_in_thousands float64
Engine_size       float64
Horsepower        float64
Wheelbase         float64
Width             float64
Length            float64
Curb_weight       float64
Fuel_capacity      float64
Fuel_efficiency    float64
Latest_Launch     object
Power_perf_factor  float64
dtype: object
```

Проверим наличие пропусков:

```
data.isnull().sum()
```

```
Manufacturer      0
Model             0
Sales_in_thousands 0
__year_resale_value 36
Vehicle_type      0
```

```

Price_in_thousands    2
Engine_size            1
Horsepower             1
Wheelbase              1
Width                  1
Length                 1
Curb_weight            2
Fuel_capacity          1
Fuel_efficiency        3
Latest_Launch          0
Power_perf_factor      2
dtype: int64

```

Видим, что пропуски наблюдаются в некоторых столбцах.

Обработка пропусков данных

Удалим колонки, содержащие пустые значения:

```

data_new_1 = data.dropna(axis=1, how='any')
(data.shape, data_new_1.shape)

((157, 16), (157, 5))

```

Выведем первые строки датасета на экран:

```

data_new_1

```

	Manufacturer	Model	Sales_in_thousands	Vehicle_type
Latest_Launch				
0	Acura	Integra	16.919	Passenger
2/2/2012				
1	Acura	TL	39.384	Passenger
6/3/2011				
2	Acura	CL	14.114	Passenger
1/4/2012				
3	Acura	RL	8.588	Passenger
3/10/2011				
4	Audi	A4	20.397	Passenger
10/8/2011				
..
..				
152	Volvo	V40	3.545	Passenger
9/21/2011				
153	Volvo	S70	15.245	Passenger
11/24/2012				
154	Volvo	V70	17.531	Passenger
6/25/2011				
155	Volvo	C70	3.493	Passenger
4/26/2011				
156	Volvo	S80	18.969	Passenger

11/14/2011

[157 rows x 5 columns]

Удалим строки, содержащие пустые значения:

```
data_new_2 = data.dropna(axis=0, how='any')  
(data.shape, data_new_2.shape)
```

```
((157, 16), (117, 16))
```

```
data_new_2.head()
```

	Manufacturer	Model	Sales_in_thousands	__year_resale_value
0	Acura	Integra	16.919	16.360
1	Acura	TL	39.384	19.875
3	Acura	RL	8.588	29.725
4	Audi	A4	20.397	22.255
5	Audi	A6	18.780	23.555

	Price_in_thousands	Engine_size	Horsepower	Wheelbase	Width
0	21.50	1.8	140.0	101.2	67.3
1	28.40	3.2	225.0	108.1	70.3
3	42.00	3.5	210.0	114.6	71.4
4	23.99	1.8	150.0	102.6	68.2
5	33.95	2.8	200.0	108.7	76.1

	Curb_weight	Fuel_capacity	Fuel_efficiency	Latest_Launch
0	2.639	13.2	28.0	2/2/2012
1	3.517	17.2	25.0	6/3/2011
3	3.850	18.0	22.0	3/10/2011
4	2.998	16.4	27.0	10/8/2011
5	3.561	18.5	22.0	8/9/2011

	Power_perf_factor
0	58.280150
1	91.370778
3	91.389779

```
4          62.777639
5          84.565105
```

Заполним все пропущенные значения нулями:

```
data_new_3 = data.fillna(0)
```

Выведем на экран:

```
data_new_3.head()
```

	Manufacturer	Model	Sales_in_thousands	__year_resale_value
0	Acura	Integra	16.919	16.360
1	Acura	TL	39.384	19.875
2	Acura	CL	14.114	18.225
3	Acura	RL	8.588	29.725
4	Audi	A4	20.397	22.255

	Price_in_thousands	Engine_size	Horsepower	Wheelbase	Width
0	21.50	1.8	140.0	101.2	67.3
1	28.40	3.2	225.0	108.1	70.3
2	0.00	3.2	225.0	106.9	70.6
3	42.00	3.5	210.0	114.6	71.4
4	23.99	1.8	150.0	102.6	68.2

	Curb_weight	Fuel_capacity	Fuel_efficiency	Latest_Launch
0	2.639	13.2	28.0	2/2/2012
1	3.517	17.2	25.0	6/3/2011
2	3.470	17.2	26.0	1/4/2012
3	3.850	18.0	22.0	3/10/2011
4	2.998	16.4	27.0	10/8/2011

	Power_perf_factor
0	58.280150
1	91.370778
2	0.000000
3	91.389779
4	62.777639

Импьютация данных

Обработка пропусков в числовых данных

Выберем числовые столбцы с пропущенными значениями и посчитаем количество пустых значений:

```
num_cols = []
for col in data.columns:
    temp_null_count = data[data[col].isnull()].shape[0]
    dt = str(data[col].dtype)
    if temp_null_count>0 and (dt=='float64' or dt=='int64'):
        num_cols.append(col)
        temp_perc = round((temp_null_count / data.shape[0]) * 100.0,
2)
        print('Столбец {}. Тип данных {}. Количество пустых значений {}, {}%.'.format(col, dt, temp_null_count, temp_perc))
```

Столбец __year_resale_value. Тип данных float64. Количество пустых значений 36, 22.93%.

Столбец Price_in_thousands. Тип данных float64. Количество пустых значений 2, 1.27%.

Столбец Engine_size. Тип данных float64. Количество пустых значений 1, 0.64%.

Столбец Horsepower. Тип данных float64. Количество пустых значений 1, 0.64%.

Столбец Wheelbase. Тип данных float64. Количество пустых значений 1, 0.64%.

Столбец Width. Тип данных float64. Количество пустых значений 1, 0.64%.

Столбец Length. Тип данных float64. Количество пустых значений 1, 0.64%.

Столбец Curb_weight. Тип данных float64. Количество пустых значений 2, 1.27%.

Столбец Fuel_capacity. Тип данных float64. Количество пустых значений 1, 0.64%.

Столбец Fuel_efficiency. Тип данных float64. Количество пустых значений 3, 1.91%.

Столбец Power_perf_factor. Тип данных float64. Количество пустых значений 2, 1.27%.

Отфильтруем по столбцам:

```
data_num = data[num_cols]
data_num
```

	__year_resale_value	Price_in_thousands	Engine_size	Horsepower
0	16.360	21.50	1.8	140.0
1	19.875	28.40	3.2	225.0

2	18.225	NaN	3.2	225.0
3	29.725	42.00	3.5	210.0
4	22.255	23.99	1.8	150.0
..
152	NaN	24.40	1.9	160.0
153	NaN	27.50	2.4	168.0
154	NaN	28.80	2.4	168.0
155	NaN	45.50	2.3	236.0
156	NaN	36.00	2.9	201.0

	Wheelbase	Width	Length	Curb_weight	Fuel_capacity
Fuel_efficiency \					
0	101.2	67.3	172.4	2.639	13.2
28.0					
1	108.1	70.3	192.9	3.517	17.2
25.0					
2	106.9	70.6	192.0	3.470	17.2
26.0					
3	114.6	71.4	196.6	3.850	18.0
22.0					
4	102.6	68.2	178.0	2.998	16.4
27.0					
..
...					
152	100.5	67.6	176.6	3.042	15.8
25.0					
153	104.9	69.3	185.9	3.208	17.9
25.0					
154	104.9	69.3	186.2	3.259	17.9
25.0					
155	104.9	71.5	185.7	3.601	18.5
23.0					
156	109.9	72.1	189.8	3.600	21.1
24.0					

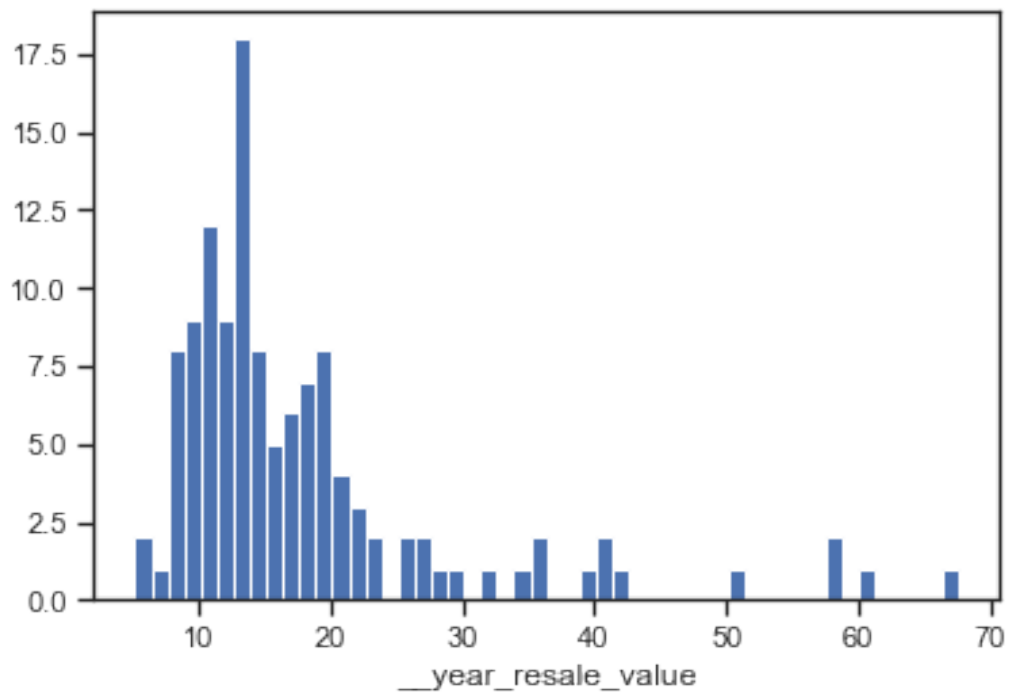
	Power_perf_factor
0	58.280150
1	91.370778
2	NaN

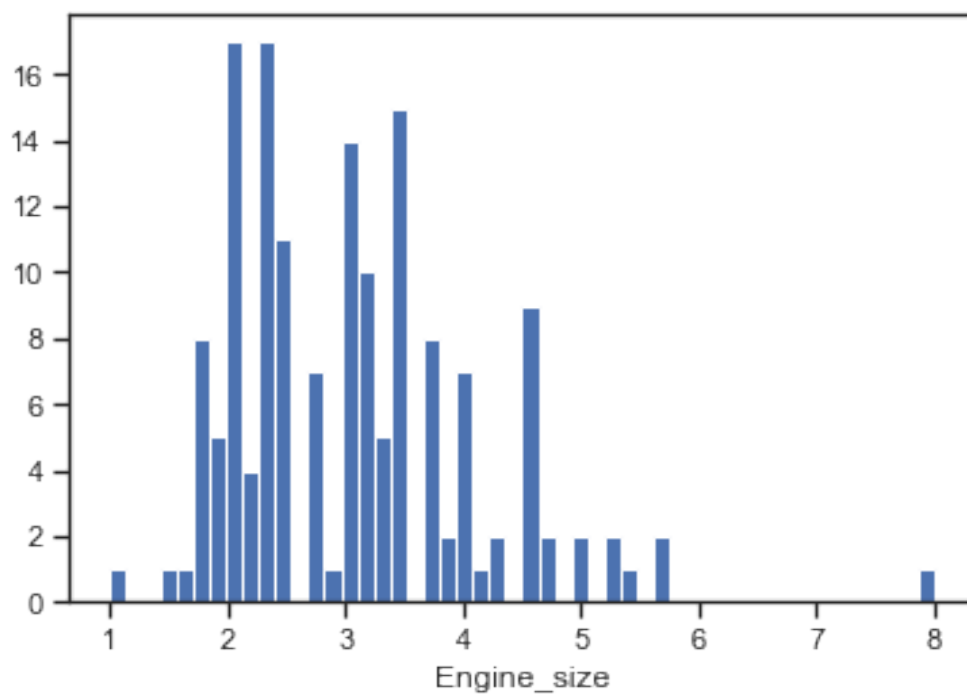
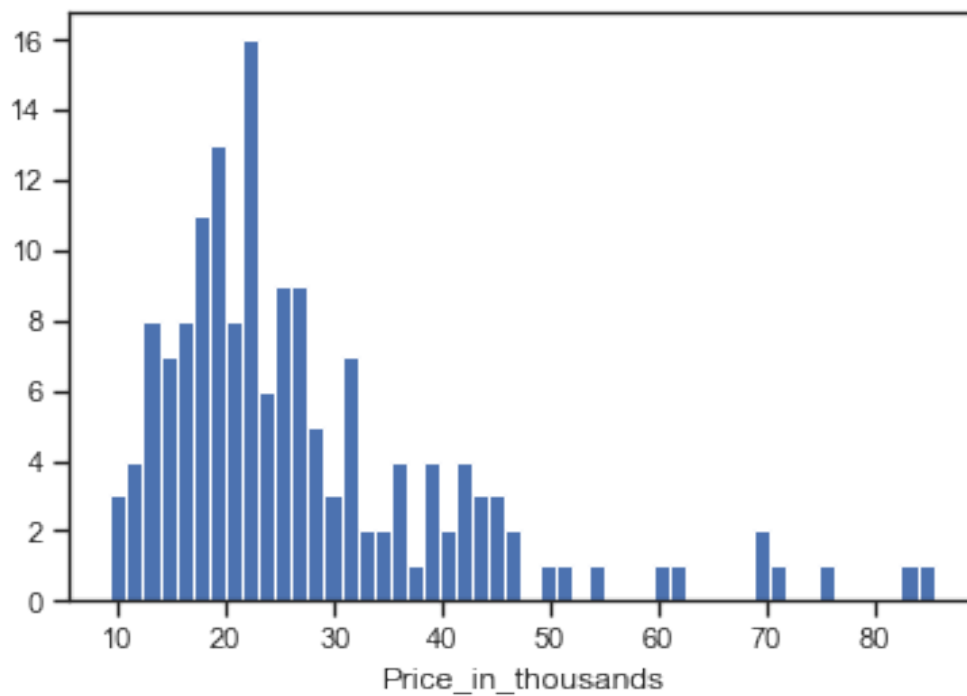
```
3          91.389779
4          62.777639
..          ...
152        66.498812
153        70.654495
154        71.155978
155        101.623357
156        85.735655
```

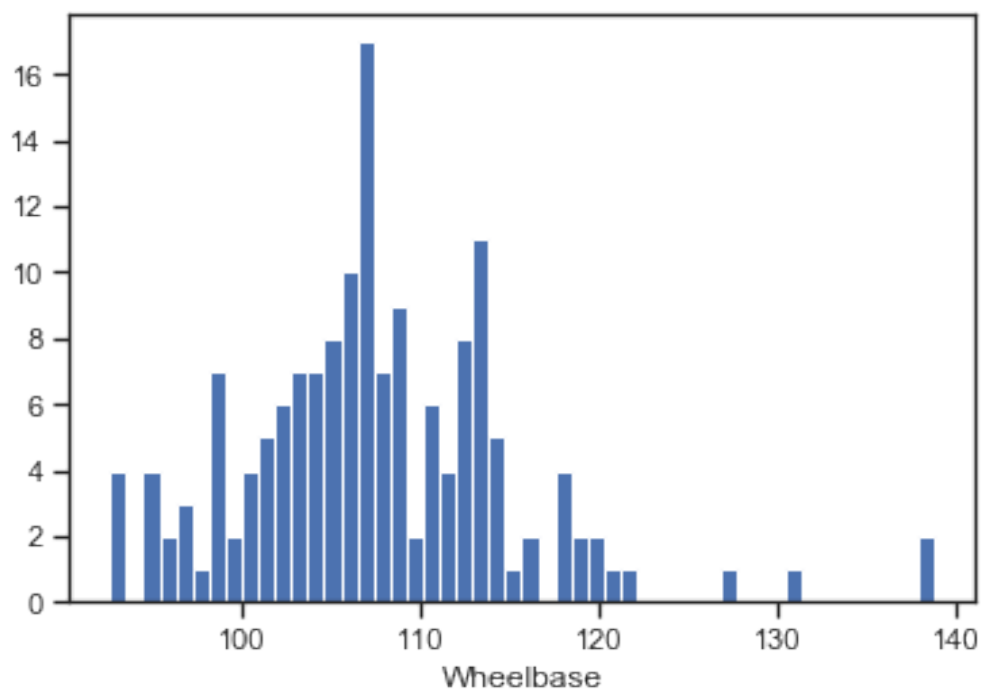
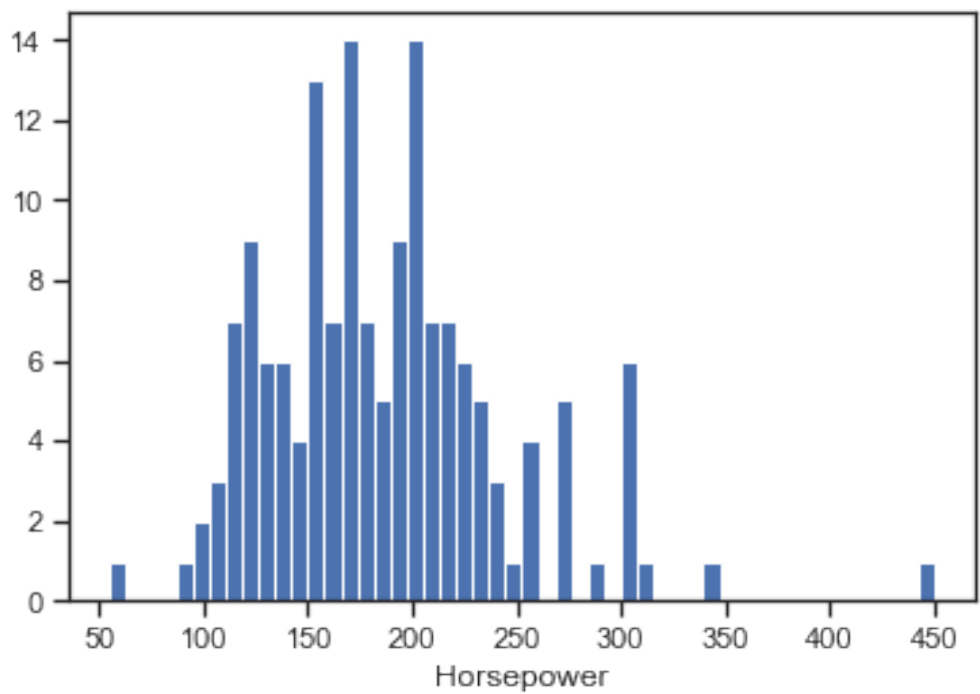
[157 rows x 11 columns]

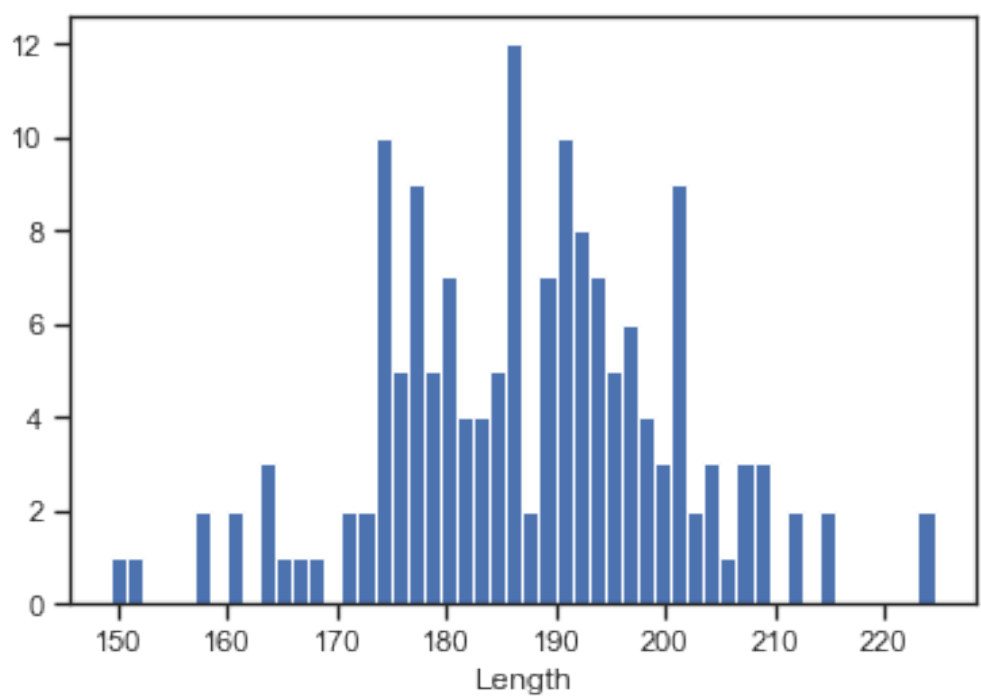
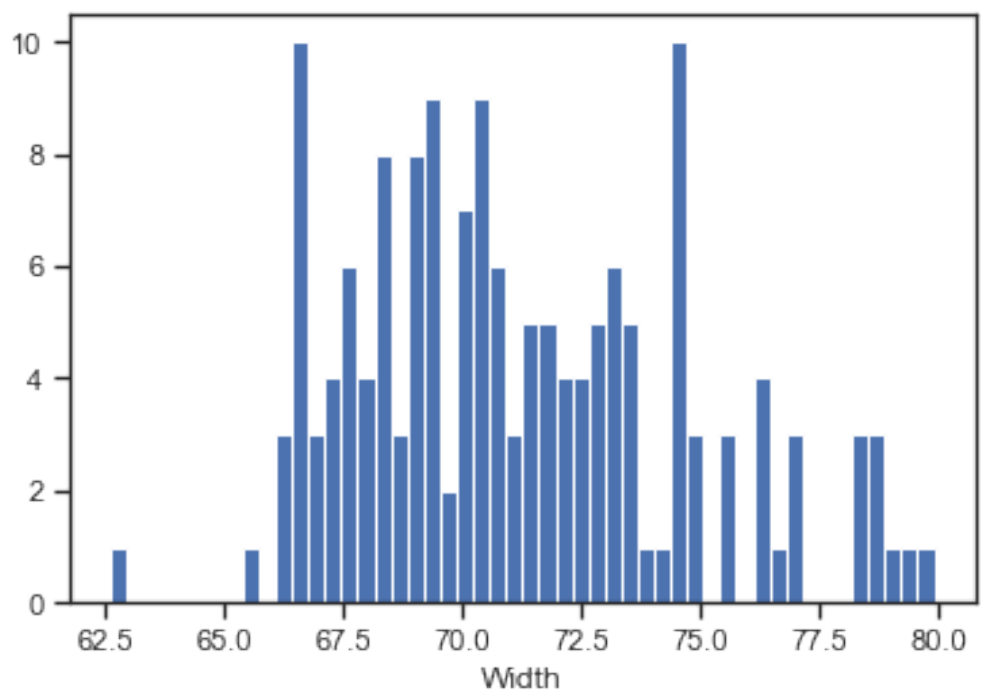
Гистограмма по признакам:

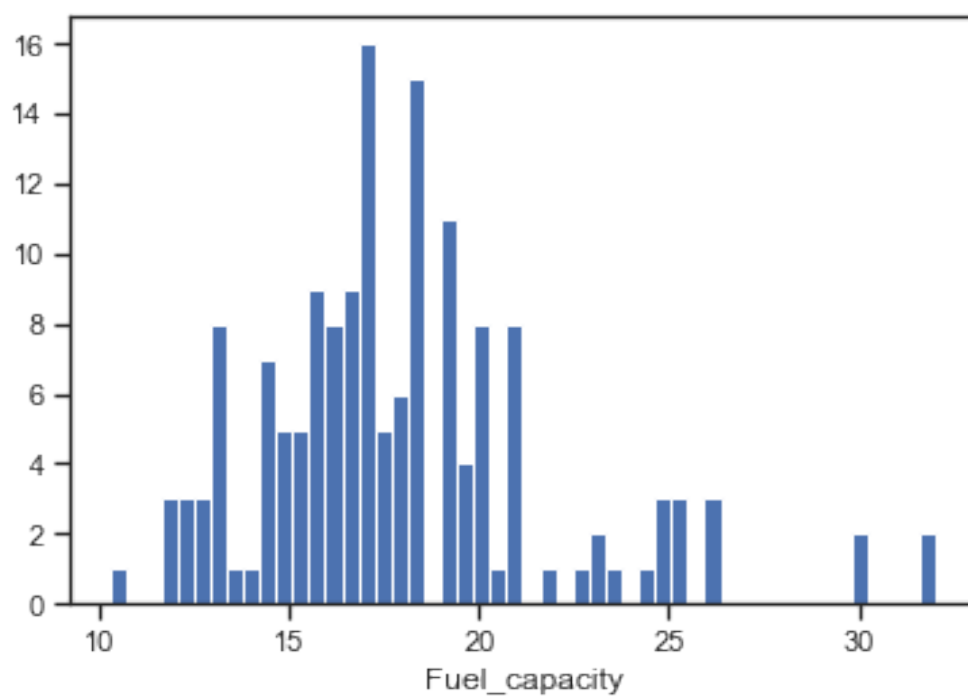
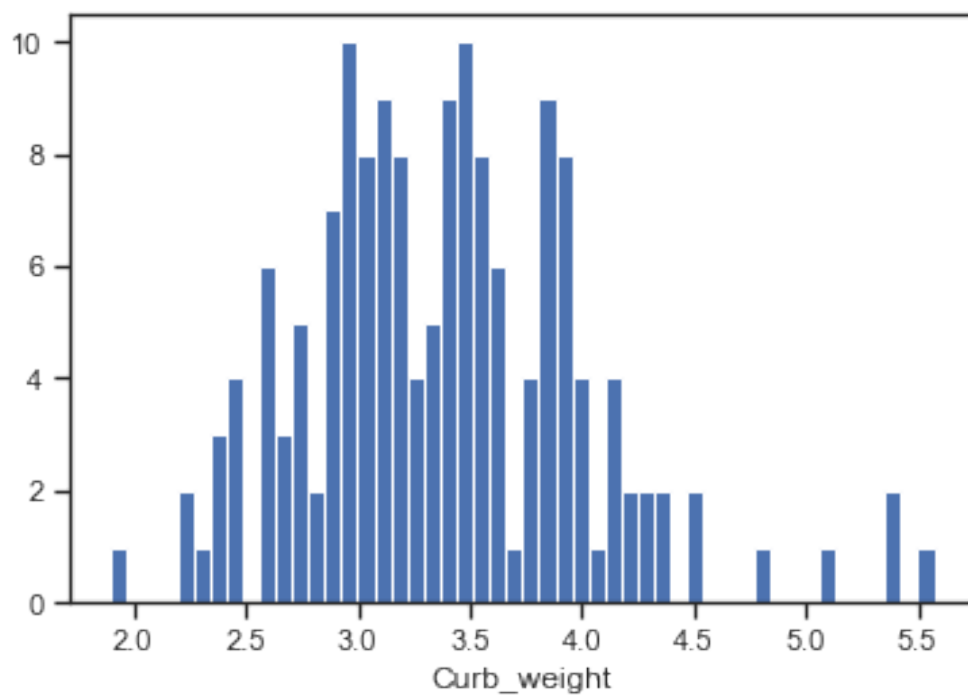
```
for col in data_num:
    plt.hist(data[col], 50)
    plt.xlabel(col)
    plt.show()
```

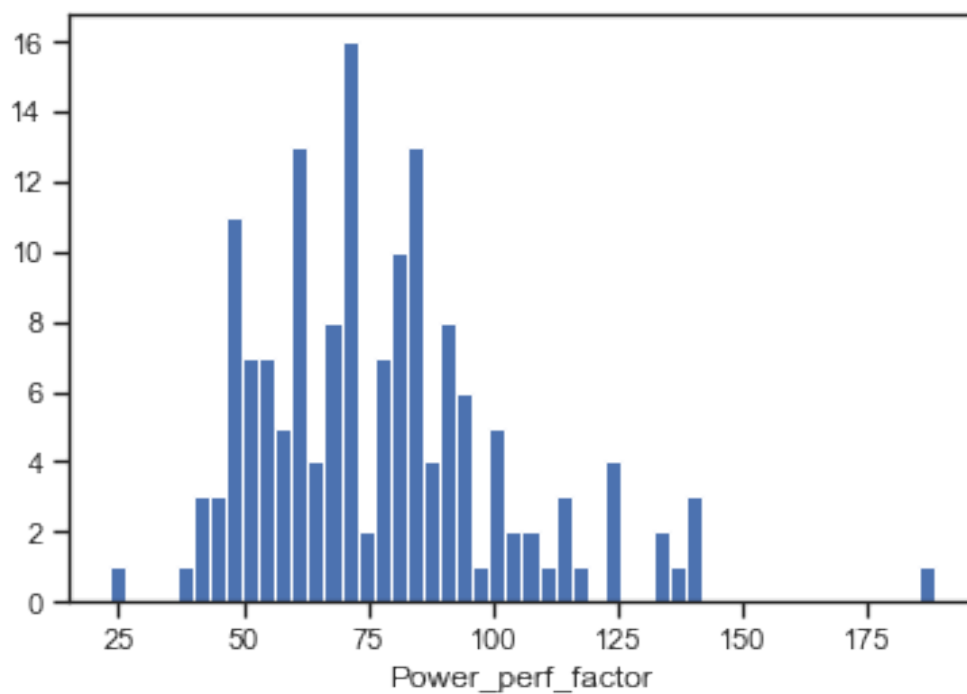
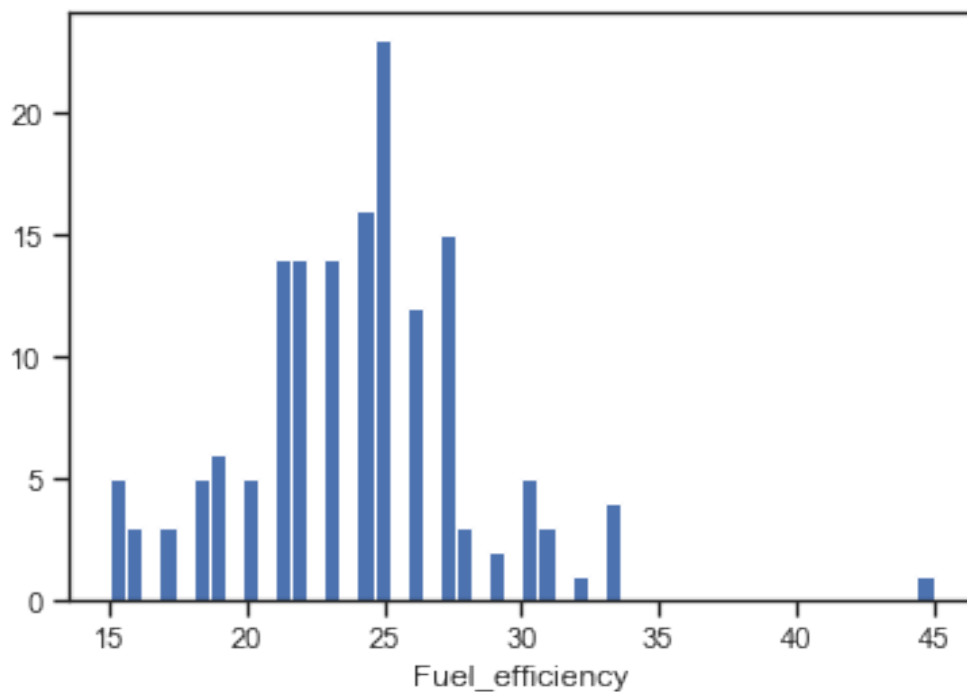












Будем использовать встроенные средства импутации библиотеки scikit-learn

```
data_num_pit = data_num[['Price_in_thousands']]
```

```
from sklearn.impute import SimpleImputer
from sklearn.impute import MissingIndicator
```

Фильтр для проверки заполнения пустых значений:

```
indicator = MissingIndicator()  
mask_missing_values_only = indicator.fit_transform(data_num_pit)  
mask_missing_values_only
```

[illegible]

[illegible]

[illegible]

```
[False],
[False],
[False],
[False],
[False],
[False],
[False],
[False],
[False],
[False],
[False],
[False]])
```

Проведем импьютацию различными показателями центра распределения:

```
strategies=['mean', 'median', 'most_frequent']
```

Создадим функцию, позволяющую задавать столбец и вид импьютации:

```
def test_num_impute_col(dataset, column, strategy_param):
    temp_data = dataset[[column]]

    indicator = MissingIndicator()
    mask_missing_values_only = indicator.fit_transform(temp_data)

    imp_num = SimpleImputer(strategy=strategy_param)
    data_num_imp = imp_num.fit_transform(temp_data)

    filled_data = data_num_imp[mask_missing_values_only]

    return column, strategy_param, filled_data.size, filled_data[0],
    filled_data[filled_data.size-1]
```

Проверим работу функции по продажам автомобилей:

```
data[['__year_resale_value']].describe()
```

	__year_resale_value
count	121.000000
mean	18.072975
std	11.453384
min	5.160000
25%	11.260000
50%	14.180000
75%	19.875000
max	67.550000

```
test_num_impute_col(data, '__year_resale_value', strategies[0])
```

```
('__year_resale_value', 'mean', 36, 18.07297520661157,
18.07297520661157)
```

```
test_num_impute_col(data, '__year_resale_value', strategies[1])
('__year_resale_value', 'median', 36, 14.18, 14.18)
test_num_impute_col(data, '__year_resale_value', strategies[2])
('__year_resale_value', 'most_frequent', 36, 7.75, 7.75)
```

Обработка пропусков в категориальных данных

Проверим категориальный признак:

```
cat_cols = []
for col in data.columns:
    temp_null_count = data_mod[data_mod[col].isnull()].shape[0]
    dt = str(data_mod[col].dtype)
    if temp_null_count > 0 and (dt == 'object'):
        cat_cols.append(col)
        temp_perc = round((temp_null_count / data.shape[0]) * 100.0,
2)
        print('Столбец {}. Тип данных {}. Количество пустых значений
{}, {}%.'.format(col, dt, temp_null_count, temp_perc))
```

Столбец Manufacturer. Тип данных object. Количество пустых значений 15, 9.55%.

```
cat_temp_data = data_mod[['Manufacturer']]
cat_temp_data.head()
```

```
Manufacturer
0      Acura
1      Acura
2      Acura
3      Acura
4      Audi
```

```
cat_temp_data['Manufacturer'].unique()
```

```
array(['Acura', 'Audi', 'BMW', 'Buick', 'Cadillac', 'Chevrolet', nan,
      'Dodge', 'Ford', 'Honda', 'Hyundai', 'Infiniti', 'Jaguar',
      'Jeep',
      'Lexus', 'Mitsubishi', 'Mercury', 'Mercedes-B', 'Nissan',
      'Oldsmobile', 'Plymouth', 'Pontiac', 'Porsche', 'Saab',
      'Subaru',
      'Toyota', 'Volkswagen', 'Volvo'], dtype=object)
```

```
cat_temp_data[cat_temp_data['Manufacturer'].isnull()].shape
(15, 1)
```

Импьютация наиболее частыми значениями:

```
imp2 = SimpleImputer(missing_values=np.nan, strategy='most_frequent')
data_imp2 = imp2.fit_transform(cat_temp_data)
data_imp2
```

[illegible]

['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Honda'],
['Honda'],
['Honda'],
['Honda'],
['Honda'],
['Hyundai'],
['Hyundai'],
['Hyundai'],
['Infiniti'],
['Jaguar'],
['Jeep'],
['Jeep'],
['Jeep'],
['Lexus'],
['Lexus'],
['Lexus'],
['Lexus'],
['Lexus'],
['Lexus'],
['Dodge'],
['Dodge'],
['Dodge'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],

[illegible]

```

        ['Volkswagen'],
        ['Volkswagen'],
        ['Volkswagen'],
        ['Volkswagen'],
        ['Volkswagen'],
        ['Volvo'],
        ['Volvo'],
        ['Volvo'],
        ['Volvo'],
        ['Volvo'],
        ['Volvo']], dtype=object)

np.unique(data_imp2)

array(['Acura', 'Audi', 'BMW', 'Buick', 'Cadillac', 'Chevrolet',
       'Dodge',
       'Ford', 'Honda', 'Hyundai', 'Infiniti', 'Jaguar', 'Jeep',
       'Lexus',
       'Mercedes-B', 'Mercury', 'Mitsubishi', 'Nissan', 'Oldsmobile',
       'Plymouth', 'Pontiac', 'Porsche', 'Saab', 'Subaru', 'Toyota',
       'Volkswagen', 'Volvo'], dtype=object)

```

Наблюдаем отсутствие пустых значений.

Импьютация константой:

```

imp3 = SimpleImputer(missing_values=np.nan, strategy='constant',
fill_value='NULL')
data_imp3 = imp3.fit_transform(cat_temp_data)
data_imp3

array([[ 'Acura'],
       [ 'Acura'],
       [ 'Acura'],
       [ 'Acura'],
       [ 'Audi'],
       [ 'Audi'],
       [ 'Audi'],
       [ 'BMW'],
       [ 'BMW'],
       [ 'BMW'],
       [ 'Buick'],
       [ 'Buick'],
       [ 'Buick'],
       [ 'Buick'],
       [ 'Cadillac'],
       [ 'Cadillac'],
       [ 'Cadillac'],
       [ 'Cadillac'],
       [ 'Cadillac'],
       [ 'Chevrolet'],

```



```
['Chevrolet'],
['Chevrolet'],
['Chevrolet'],
['Chevrolet'],
['Chevrolet'],
['Chevrolet'],
['Chevrolet'],
['Chevrolet'],
['NULL'],
['NULL'],
['NULL'],
['NULL'],
['NULL'],
['NULL'],
['NULL'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Dodge'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Ford'],
['Honda'],
['Honda'],
['Honda'],
['Honda'],
['Honda'],
['Hyundai'],
['Hyundai'],
['Hyundai'],
['Infiniti'],
['Jaguar'],
['Jeep'],
['Jeep'],
['Jeep'],
```

```
['Lexus'],
['Lexus'],
['Lexus'],
['Lexus'],
['Lexus'],
['Lexus'],
['Lexus'],
['NULL'],
['NULL'],
['NULL'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mitsubishi'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercury'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Mercedes-B'],
['Nissan'],
['Nissan'],
['Nissan'],
['Nissan'],
['Nissan'],
['Nissan'],
['Nissan'],
['Oldsmobile'],
['Oldsmobile'],
['Oldsmobile'],
['Oldsmobile'],
['Oldsmobile'],
['Oldsmobile'],
['Plymouth'],
['Plymouth'],
['Plymouth'],
['Plymouth'],
['Pontiac'],
['Pontiac'],
```

```

['Pontiac'],
['Pontiac'],
['Pontiac'],
['Pontiac'],
['Porsche'],
['Porsche'],
['Porsche'],
['Saab'],
['Saab'],
['NULL'],
['NULL'],
['NULL'],
['NULL'],
['NULL'],
['Subaru'],
['Subaru'],
['Toyota'],
['Toyota'],
['Toyota'],
['Toyota'],
['Toyota'],
['Toyota'],
['Toyota'],
['Toyota'],
['Toyota'],
['Volkswagen'],
['Volkswagen'],
['Volkswagen'],
['Volkswagen'],
['Volkswagen'],
['Volkswagen'],
['Volvo'],
['Volvo'],
['Volvo'],
['Volvo'],
['Volvo'],
['Volvo']], dtype=object)

np.unique(data_imp3)

array(['Acura', 'Audi', 'BMW', 'Buick', 'Cadillac', 'Chevrolet',
'Dodge',
      'Ford', 'Honda', 'Hyundai', 'Infiniti', 'Jaguar', 'Jeep',
'Lexus',
      'Mercedes-B', 'Mercury', 'Mitsubishi', 'NULL', 'Nissan',
      'Oldsmobile', 'Plymouth', 'Pontiac', 'Porsche', 'Saab',
'Subaru',
      'Toyota', 'Volkswagen', 'Volvo'], dtype=object)

data_imp3[data_imp3==0].size

```

0

Значения были заменены на "NULL".

Преобразование категориальных признаков в числовые

```
cat_enc = pd.DataFrame({'cat1':data_imp2.T[0]})
cat_enc
```

```
   cat1
0  Acura
1  Acura
2  Acura
3  Acura
4   Audi
..    ..
152 Volvo
153 Volvo
154 Volvo
155 Volvo
156 Volvo
```

[157 rows x 1 columns]

Кодирование категорий целочисленными значениями

LabelEncoder

```
from sklearn.preprocessing import LabelEncoder
```

```
cat_enc['cat1'].unique()
```

```
array(['Acura', 'Audi', 'BMW', 'Buick', 'Cadillac', 'Chevrolet',
       'Dodge',
       'Ford', 'Honda', 'Hyundai', 'Infiniti', 'Jaguar', 'Jeep',
       'Lexus',
       'Mitsubishi', 'Mercury', 'Mercedes-B', 'Nissan', 'Oldsmobile',
       'Plymouth', 'Pontiac', 'Porsche', 'Saab', 'Subaru', 'Toyota',
       'Volkswagen', 'Volvo'], dtype=object)
```

```
le = LabelEncoder()
```

```
cat_enc_le = le.fit_transform(cat_enc['cat1'])
```

```
le.classes_
```

```
array(['Acura', 'Audi', 'BMW', 'Buick', 'Cadillac', 'Chevrolet',
       'Dodge',
       'Ford', 'Honda', 'Hyundai', 'Infiniti', 'Jaguar', 'Jeep',
       'Lexus',
       'Mercedes-B', 'Mercury', 'Mitsubishi', 'Nissan', 'Oldsmobile',
       'Plymouth', 'Pontiac', 'Porsche', 'Saab', 'Subaru', 'Toyota',
       'Volkswagen', 'Volvo'], dtype=object)
```

```
cat_enc_le
array([ 0,  0,  0,  0,  1,  1,  1,  2,  2,  2,  3,  3,  3,  3,  4,  4,
        4,  4,  4,  5,  5,  5,  5,  5,  5,  5,  5,  6,  6,  6,  6,  6,
        6,  6,  6,  6,  6,  6,  6,  6,  6,  6,  6,  6,  7,  7,  7,  7,
        7,  7,  7,  7,  7,  7,  8,  8,  8,  8,  8,  9,  9,  9, 10, 11,
       12, 12, 12, 13, 13, 13, 13, 13, 13, 13,  6,  6,  6, 16, 16, 16, 16, 16,
       16, 16, 15, 15, 15, 15, 15, 15, 15, 14, 14, 14, 14, 14, 14, 14, 14,
       17, 17, 17, 17, 17, 17, 18, 18, 18, 18, 18, 18, 19, 19, 19, 19,
       20, 20, 20, 20, 20, 21, 21, 21, 22, 22,  6,  6,  6,  6,  6, 23,
       23, 24, 24, 24, 24, 24, 24, 24, 24, 24, 25, 25, 25, 25, 25, 25, 26,
       26, 26, 26, 26])
```

```
np.unique(cat_enc_le)
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15,
       16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26])
```

```
le.inverse_transform([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11,
                       12, 13, 14, 15, 16,
                       17, 18, 19, 20, 21, 22, 23, 24, 25, 26])
```

```
array(['Acura', 'Audi', 'BMW', 'Buick', 'Cadillac', 'Chevrolet',
       'Dodge',
       'Ford', 'Honda', 'Hyundai', 'Infiniti', 'Jaguar', 'Jeep',
       'Lexus',
       'Mercedes-B', 'Mercury', 'Mitsubishi', 'Nissan', 'Oldsmobile',
       'Plymouth', 'Pontiac', 'Porsche', 'Saab', 'Subaru', 'Toyota',
       'Volkswagen', 'Volvo'], dtype=object)
```

OrdinalEncoder

```
from sklearn.preprocessing import OrdinalEncoder
```

```
data_oe = data_mod[['Manufacturer', 'Model']]
data_oe.head()
```

	Manufacturer	Model
0	Acura	Integra
1	Acura	TL
2	Acura	CL
3	Acura	RL
4	Audi	A4

```

imp4 = SimpleImputer(missing_values=np.nan, strategy='constant',
fill_value='NULL')
data_oe_filled = imp4.fit_transform(data_oe)
data_oe_filled
array([[ 'Acura', 'Integra'],
      [ 'Acura', 'TL'],
      [ 'Acura', 'CL'],
      [ 'Acura', 'RL'],
      [ 'Audi', 'A4'],
      [ 'Audi', 'A6'],
      [ 'Audi', 'A8'],
      [ 'BMW', '323i'],
      [ 'BMW', '328i'],
      [ 'BMW', '528i'],
      [ 'Buick', 'Century'],
      [ 'Buick', 'Regal'],
      [ 'Buick', 'Park Avenue'],
      [ 'Buick', 'LeSabre'],
      [ 'Cadillac', 'DeVille'],
      [ 'Cadillac', 'Seville'],
      [ 'Cadillac', 'Eldorado'],
      [ 'Cadillac', 'Catera'],
      [ 'Cadillac', 'Escalade'],
      [ 'Chevrolet', 'Cavalier'],
      [ 'Chevrolet', 'Malibu'],
      [ 'Chevrolet', 'Lumina'],
      [ 'Chevrolet', 'Monte Carlo'],
      [ 'Chevrolet', 'Camaro'],
      [ 'Chevrolet', 'Corvette'],
      [ 'Chevrolet', 'Prizm'],
      [ 'Chevrolet', 'Metro'],
      [ 'Chevrolet', 'Impala'],
      [ 'NULL', 'Sebring Coupe'],
      [ 'NULL', 'Sebring Conv.'],
      [ 'NULL', 'Concorde'],
      [ 'NULL', 'Cirrus'],
      [ 'NULL', 'LHS'],
      [ 'NULL', 'Town & Country'],
      [ 'NULL', '300M'],
      [ 'Dodge', 'Neon'],
      [ 'Dodge', 'Avenger'],
      [ 'Dodge', 'Stratus'],
      [ 'Dodge', 'Intrepid'],
      [ 'Dodge', 'Viper'],
      [ 'Dodge', 'Ram Pickup'],
      [ 'Dodge', 'Ram Wagon'],
      [ 'Dodge', 'Ram Van'],
      [ 'Dodge', 'Dakota'],
      [ 'Dodge', 'Durango'],

```

```
['Dodge', 'Caravan'],
['Ford', 'Escort'],
['Ford', 'Mustang'],
['Ford', 'Contour'],
['Ford', 'Taurus'],
['Ford', 'Focus'],
['Ford', 'Crown Victoria'],
['Ford', 'Explorer'],
['Ford', 'Windstar'],
['Ford', 'Expedition'],
['Ford', 'Ranger'],
['Ford', 'F-Series'],
['Honda', 'Civic'],
['Honda', 'Accord'],
['Honda', 'CR-V'],
['Honda', 'Passport'],
['Honda', 'Odyssey'],
['Hyundai', 'Accent'],
['Hyundai', 'Elantra'],
['Hyundai', 'Sonata'],
['Infiniti', 'I30'],
['Jaguar', 'S-Type'],
['Jeep', 'Wrangler'],
['Jeep', 'Cherokee'],
['Jeep', 'Grand Cherokee'],
['Lexus', 'ES300'],
['Lexus', 'GS300'],
['Lexus', 'GS400'],
['Lexus', 'LS400'],
['Lexus', 'LX470'],
['Lexus', 'RX300'],
['NULL', 'Continental'],
['NULL', 'Town car'],
['NULL', 'Navigator'],
['Mitsubishi', 'Mirage'],
['Mitsubishi', 'Eclipse'],
['Mitsubishi', 'Galant'],
['Mitsubishi', 'Diamante'],
['Mitsubishi', '3000GT'],
['Mitsubishi', 'Montero'],
['Mitsubishi', 'Montero Sport'],
['Mercury', 'Mystique'],
['Mercury', 'Cougar'],
['Mercury', 'Sable'],
['Mercury', 'Grand Marquis'],
['Mercury', 'Mountaineer'],
['Mercury', 'Villager'],
['Mercedes-B', 'C-Class'],
['Mercedes-B', 'E-Class'],
['Mercedes-B', 'S-Class'],
```

```
['Mercedes-B', 'SL-Class'],
['Mercedes-B', 'SLK'],
['Mercedes-B', 'SLK230'],
['Mercedes-B', 'CLK Coupe'],
['Mercedes-B', 'CL500'],
['Mercedes-B', 'M-Class'],
['Nissan', 'Sentra'],
['Nissan', 'Altima'],
['Nissan', 'Maxima'],
['Nissan', 'Quest'],
['Nissan', 'Pathfinder'],
['Nissan', 'Xterra'],
['Nissan', 'Frontier'],
['Oldsmobile', 'Cutlass'],
['Oldsmobile', 'Intrigue'],
['Oldsmobile', 'Alero'],
['Oldsmobile', 'Aurora'],
['Oldsmobile', 'Bravada'],
['Oldsmobile', 'Silhouette'],
['Plymouth', 'Neon'],
['Plymouth', 'Breeze'],
['Plymouth', 'Voyager'],
['Plymouth', 'Prowler'],
['Pontiac', 'Sunfire'],
['Pontiac', 'Grand Am'],
['Pontiac', 'Firebird'],
['Pontiac', 'Grand Prix'],
['Pontiac', 'Bonneville'],
['Pontiac', 'Montana'],
['Porsche', 'Boxter'],
['Porsche', 'Carrera Coupe'],
['Porsche', 'Carrera Cabrio'],
['Saab', '5-Sep'],
['Saab', '3-Sep'],
['NULL', 'SL'],
['NULL', 'SC'],
['NULL', 'SW'],
['NULL', 'LW'],
['NULL', 'LS'],
['Subaru', 'Outback'],
['Subaru', 'Forester'],
['Toyota', 'Corolla'],
['Toyota', 'Camry'],
['Toyota', 'Avalon'],
['Toyota', 'Celica'],
['Toyota', 'Tacoma'],
['Toyota', 'Sienna'],
['Toyota', 'RAV4'],
['Toyota', '4Runner'],
['Toyota', 'Land Cruiser'],
```



```

['Volkswagen', 'Golf'],
['Volkswagen', 'Jetta'],
['Volkswagen', 'Passat'],
['Volkswagen', 'Cabrio'],
['Volkswagen', 'GTI'],
['Volkswagen', 'Beetle'],
['Volvo', 'S40'],
['Volvo', 'V40'],
['Volvo', 'S70'],
['Volvo', 'V70'],
['Volvo', 'C70'],
['Volvo', 'S80']], dtype=object)

oe = OrdinalEncoder()
cat_enc_oe = oe.fit_transform(data_oe_filled)
cat_enc_oe
array([[ 0.,  79.],
       [ 0., 143.],
       [ 0.,  25.],
       [ 0., 115.],
       [ 1.,   8.],
       [ 1.,   9.],
       [ 1.,  10.],
       [ 2.,   3.],
       [ 2.,   4.],
       [ 2.,   7.],
       [ 3.,  38.],
       [ 3., 121.],
       [ 3., 107.],
       [ 3.,  89.],
       [ 4.,  51.],
       [ 4., 137.],
       [ 4.,  58.],
       [ 4.,  35.],
       [ 4.,  59.],
       [ 5.,  36.],
       [ 5.,  92.],
       [ 5.,  90.],
       [ 5.,  97.],
       [ 5.,  30.],
       [ 5.,  46.],
       [ 5., 111.],
       [ 5.,  94.],
       [ 5.,  78.],
       [17., 135.],
       [17., 134.],
       [17.,  42.],
       [17.,  40.],
       [17.,  83.],
       [17., 146.]])

```

```
[ 17.,  2.],
[  6., 104.],
[  6.,  17.],
[  6., 141.],
[  6.,  80.],
[  6., 151.],
[  6., 117.],
[  6., 119.],
[  6., 118.],
[  6.,  50.],
[  6.,  53.],
[  6.,  32.],
[  7.,  60.],
[  7., 101.],
[  7.,  44.],
[  7., 145.],
[  7.,  65.],
[  7.,  48.],
[  7.,  62.],
[  7., 153.],
[  7.,  61.],
[  7., 120.],
[  7.,  63.],
[  8.,  41.],
[  8.,  12.],
[  8.,  28.],
[  8., 109.],
[  8., 105.],
[  9.,  11.],
[  9.,  57.],
[  9., 140.],
[ 10.,  77.],
[ 11., 123.],
[ 12., 154.],
[ 12.,  39.],
[ 12.,  74.],
[ 13.,  55.],
[ 13.,  68.],
[ 13.,  69.],
[ 13.,  85.],
[ 13.,  87.],
[ 13., 116.],
[ 17.,  43.],
[ 17., 147.],
[ 17., 103.],
[ 16.,  95.],
[ 16.,  56.],
[ 16.,  71.],
[ 16.,  52.],
[ 16.,   1.],
```

```
[ 16.,  98.],  
[ 16.,  99.],  
[ 15., 102.],  
[ 15.,  47.],  
[ 15., 133.],  
[ 15.,  75.],  
[ 15., 100.],  
[ 15., 150.],  
[ 14.,  23.],  
[ 14.,  54.],  
[ 14., 122.],  
[ 14., 129.],  
[ 14., 130.],  
[ 14., 131.],  
[ 14.,  27.],  
[ 14.,  26.],  
[ 14.,  91.],  
[ 18., 136.],  
[ 18.,  14.],  
[ 18.,  93.],  
[ 18., 113.],  
[ 18., 110.],  
[ 18., 155.],  
[ 18.,  67.],  
[ 19.,  49.],  
[ 19.,  81.],  
[ 19.,  13.],  
[ 19.,  15.],  
[ 19.,  21.],  
[ 19., 139.],  
[ 20., 104.],  
[ 20.,  22.],  
[ 20., 152.],  
[ 20., 112.],  
[ 21., 142.],  
[ 21.,  73.],  
[ 21.,  64.],  
[ 21.,  76.],  
[ 21.,  19.],  
[ 21.,  96.],  
[ 22.,  20.],  
[ 22.,  34.],  
[ 22.,  33.],  
[ 23.,   6.],  
[ 23.,   0.],  
[ 17., 128.],  
[ 17., 127.],  
[ 17., 132.],  
[ 17.,  86.],  
[ 17.,  84.],
```

```
[ 24., 106.],
[ 24.,  66.],
[ 25.,  45.],
[ 25.,  31.],
[ 25.,  16.],
[ 25.,  37.],
[ 25., 144.],
[ 25., 138.],
[ 25., 114.],
[ 25.,   5.],
[ 25.,  88.],
[ 26.,  72.],
[ 26.,  82.],
[ 26., 108.],
[ 26.,  29.],
[ 26.,  70.],
[ 26.,  18.],
[ 27., 124.],
[ 27., 148.],
[ 27., 125.],
[ 27., 149.],
[ 27.,  24.],
[ 27., 126.]]])
```

Уникальные значения столбца "Производитель":

```
np.unique(cat_enc_oe[:, 0])
array([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9., 10., 11.,
        12., 13., 14., 15., 16., 17., 18., 19., 20., 21., 22., 23., 24.,
        25., 26., 27.]])
```

Уникальные значения столбца "Модель":

```
np.unique(cat_enc_oe[:, 1])
array([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9.,
        10., 11., 12., 13., 14., 15., 16., 17., 18., 19., 20.,
        21., 22., 23., 24., 25., 26., 27., 28., 29., 30., 31.,
        32., 33., 34., 35., 36., 37., 38., 39., 40., 41., 42.,
        43., 44., 45., 46., 47., 48., 49., 50., 51., 52., 53.,
        54., 55., 56., 57., 58., 59., 60., 61., 62., 63., 64.,
        65., 66., 67., 68., 69., 70., 71., 72., 73., 74., 75.,
```

```

76.,
    77., 78., 79., 80., 81., 82., 83., 84., 85., 86.,
87.,
    88., 89., 90., 91., 92., 93., 94., 95., 96., 97.,
98.,
    99., 100., 101., 102., 103., 104., 105., 106., 107., 108.,
109.,
    110., 111., 112., 113., 114., 115., 116., 117., 118., 119.,
120.,
    121., 122., 123., 124., 125., 126., 127., 128., 129., 130.,
131.,
    132., 133., 134., 135., 136., 137., 138., 139., 140., 141.,
142.,
    143., 144., 145., 146., 147., 148., 149., 150., 151., 152.,
153.,
    154., 155.])

```

Все значения:

oe.categories_

```

[array(['Acura', 'Audi', 'BMW', 'Buick', 'Cadillac', 'Chevrolet',
'Dodge',
      'Ford', 'Honda', 'Hyundai', 'Infiniti', 'Jaguar', 'Jeep',
'Lexus',
      'Mercedes-B', 'Mercury', 'Mitsubishi', 'NULL', 'Nissan',
'Oldsmobile', 'Plymouth', 'Pontiac', 'Porsche', 'Saab',
'Subaru',
      'Toyota', 'Volkswagen', 'Volvo'], dtype=object),
array(['3-Sep', '3000GT', '300M', '323i', '328i', '4Runner', '5-Sep',
'528i', 'A4', 'A6', 'A8', 'Accent', 'Accord', 'Alero',
'Altima',
      'Aurora', 'Avalon', 'Avenger', 'Beetle', 'Bonneville',
'Boxster',
      'Bravada', 'Breeze', 'C-Class', 'C70', 'CL', 'CL500', 'CLK
Coupe',
      'CR-V', 'Cabrio', 'Camaro', 'Camry', 'Caravan', 'Carrera
Cabrio',
      'Carrera Coupe', 'Catera', 'Cavalier', 'Celica', 'Century',
'Cherokee', 'Cirrus', 'Civic', 'Concorde', 'Continental',
'Contour', 'Corolla', 'Corvette', 'Cougar', 'Crown Victoria',
'Cutlass', 'Dakota', 'DeVille', 'Diamante', 'Durango', 'E-
Class',
      'ES300', 'Eclipse', 'Elantra', 'Eldorado', 'Escalade',
'Escort',
      'Expedition', 'Explorer', 'F-Series', 'Firebird', 'Focus',
'Forester', 'Frontier', 'GS300', 'GS400', 'GTI', 'Galant',
'Golf',
      'Grand Am', 'Grand Cherokee', 'Grand Marquis', 'Grand Prix',
'I30',
      'Impala', 'Integra', 'Intrepid', 'Intrigue', 'Jetta', 'LHS',

```

```

'LS',
    'LS400', 'LW', 'LX470', 'Land Cruiser', 'LeSabre', 'Lumina',
    'M-Class', 'Malibu', 'Maxima', 'Metro', 'Mirage', 'Montana',
    'Monte Carlo', 'Montero', 'Montero Sport', 'Mountaineer',
    'Mustang', 'Mystique', 'Navigator', 'Neon', 'Odyssey',
'Outback',
    'Park Avenue', 'Passat', 'Passport', 'Pathfinder', 'Prizm',
    'Prowler', 'Quest', 'RAV4', 'RL', 'RX300', 'Ram Pickup', 'Ram
Van',
    'Ram Wagon', 'Ranger', 'Regal', 'S-Class', 'S-Type', 'S40',
'S70',
    'S80', 'SC', 'SL', 'SL-Class', 'SLK', 'SLK230', 'SW', 'Sable',
    'Sebring Conv.', 'Sebring Coupe', 'Sentra', 'Seville',
'Sienna',
    'Silhouette', 'Sonata', 'Stratus', 'Sunfire', 'TL', 'Tacoma',
    'Taurus', 'Town & Country', 'Town car', 'V40', 'V70',
'Villager',
    'Viper', 'Voyager', 'Windstar', 'Wrangler', 'Xterra'],
dtype=object)]

```

```
oe.inverse_transform(cat_enc_oe)
```

```

array([[ 'Acura', 'Integra'],
       [ 'Acura', 'TL'],
       [ 'Acura', 'CL'],
       [ 'Acura', 'RL'],
       [ 'Audi', 'A4'],
       [ 'Audi', 'A6'],
       [ 'Audi', 'A8'],
       [ 'BMW', '323i'],
       [ 'BMW', '328i'],
       [ 'BMW', '528i'],
       [ 'Buick', 'Century'],
       [ 'Buick', 'Regal'],
       [ 'Buick', 'Park Avenue'],
       [ 'Buick', 'LeSabre'],
       [ 'Cadillac', 'DeVille'],
       [ 'Cadillac', 'Seville'],
       [ 'Cadillac', 'Eldorado'],
       [ 'Cadillac', 'Catera'],
       [ 'Cadillac', 'Escalade'],
       [ 'Chevrolet', 'Cavalier'],
       [ 'Chevrolet', 'Malibu'],
       [ 'Chevrolet', 'Lumina'],
       [ 'Chevrolet', 'Monte Carlo'],
       [ 'Chevrolet', 'Camaro'],
       [ 'Chevrolet', 'Corvette'],
       [ 'Chevrolet', 'Prizm'],
       [ 'Chevrolet', 'Metro'],
       [ 'Chevrolet', 'Impala'],
       [ 'NULL', 'Sebring Coupe'],

```

['NULL', 'Sebring Conv.'],
['NULL', 'Concorde'],
['NULL', 'Cirrus'],
['NULL', 'LHS'],
['NULL', 'Town & Country'],
['NULL', '300M'],
['Dodge', 'Neon'],
['Dodge', 'Avenger'],
['Dodge', 'Stratus'],
['Dodge', 'Intrepid'],
['Dodge', 'Viper'],
['Dodge', 'Ram Pickup'],
['Dodge', 'Ram Wagon'],
['Dodge', 'Ram Van'],
['Dodge', 'Dakota'],
['Dodge', 'Durango'],
['Dodge', 'Caravan'],
['Ford', 'Escort'],
['Ford', 'Mustang'],
['Ford', 'Contour'],
['Ford', 'Taurus'],
['Ford', 'Focus'],
['Ford', 'Crown Victoria'],
['Ford', 'Explorer'],
['Ford', 'Windstar'],
['Ford', 'Expedition'],
['Ford', 'Ranger'],
['Ford', 'F-Series'],
['Honda', 'Civic'],
['Honda', 'Accord'],
['Honda', 'CR-V'],
['Honda', 'Passport'],
['Honda', 'Odyssey'],
['Hyundai', 'Accent'],
['Hyundai', 'Elantra'],
['Hyundai', 'Sonata'],
['Infiniti', 'I30'],
['Jaguar', 'S-Type'],
['Jeep', 'Wrangler'],
['Jeep', 'Cherokee'],
['Jeep', 'Grand Cherokee'],
['Lexus', 'ES300'],
['Lexus', 'GS300'],
['Lexus', 'GS400'],
['Lexus', 'LS400'],
['Lexus', 'LX470'],
['Lexus', 'RX300'],
['NULL', 'Continental'],
['NULL', 'Town car'],
['NULL', 'Navigator'],

['Mitsubishi', 'Mirage'],
['Mitsubishi', 'Eclipse'],
['Mitsubishi', 'Galant'],
['Mitsubishi', 'Diamante'],
['Mitsubishi', '3000GT'],
['Mitsubishi', 'Montero'],
['Mitsubishi', 'Montero Sport'],
['Mercury', 'Mystique'],
['Mercury', 'Cougar'],
['Mercury', 'Sable'],
['Mercury', 'Grand Marquis'],
['Mercury', 'Mountaineer'],
['Mercury', 'Villager'],
['Mercedes-B', 'C-Class'],
['Mercedes-B', 'E-Class'],
['Mercedes-B', 'S-Class'],
['Mercedes-B', 'SL-Class'],
['Mercedes-B', 'SLK'],
['Mercedes-B', 'SLK230'],
['Mercedes-B', 'CLK Coupe'],
['Mercedes-B', 'CL500'],
['Mercedes-B', 'M-Class'],
['Nissan', 'Sentra'],
['Nissan', 'Altima'],
['Nissan', 'Maxima'],
['Nissan', 'Quest'],
['Nissan', 'Pathfinder'],
['Nissan', 'Xterra'],
['Nissan', 'Frontier'],
['Oldsmobile', 'Cutlass'],
['Oldsmobile', 'Intrigue'],
['Oldsmobile', 'Alero'],
['Oldsmobile', 'Aurora'],
['Oldsmobile', 'Bravada'],
['Oldsmobile', 'Silhouette'],
['Plymouth', 'Neon'],
['Plymouth', 'Breeze'],
['Plymouth', 'Voyager'],
['Plymouth', 'Prowler'],
['Pontiac', 'Sunfire'],
['Pontiac', 'Grand Am'],
['Pontiac', 'Firebird'],
['Pontiac', 'Grand Prix'],
['Pontiac', 'Bonneville'],
['Pontiac', 'Montana'],
['Porsche', 'Boxter'],
['Porsche', 'Carrera Coupe'],
['Porsche', 'Carrera Cabrio'],
['Saab', '5-Sep'],
['Saab', '3-Sep'],


```

['NULL', 'SL'],
['NULL', 'SC'],
['NULL', 'SW'],
['NULL', 'LW'],
['NULL', 'LS'],
['Subaru', 'Outback'],
['Subaru', 'Forester'],
['Toyota', 'Corolla'],
['Toyota', 'Camry'],
['Toyota', 'Avalon'],
['Toyota', 'Celica'],
['Toyota', 'Tacoma'],
['Toyota', 'Sienna'],
['Toyota', 'RAV4'],
['Toyota', '4Runner'],
['Toyota', 'Land Cruiser'],
['Volkswagen', 'Golf'],
['Volkswagen', 'Jetta'],
['Volkswagen', 'Passat'],
['Volkswagen', 'Cabrio'],
['Volkswagen', 'GTI'],
['Volkswagen', 'Beetle'],
['Volvo', 'S40'],
['Volvo', 'V40'],
['Volvo', 'S70'],
['Volvo', 'V70'],
['Volvo', 'C70'],
['Volvo', 'S80']], dtype=object)

```

Кодирование шкал порядка

Для кодирования шкал порядка воспользуемся функцией map:

```
sizes = ['small', 'medium', 'large', 'small', 'medium', 'large',
'small', 'medium', 'large']
```

```
pd_sizes = pd.DataFrame(data={'sizes':sizes})
pd_sizes
```

```

      sizes
0    small
1  medium
2    large
3    small
4  medium
5    large
6    small
7  medium
8    large

```

```
pd_sizes['sizes_codes'] = pd_sizes['sizes'].map({'small':1,
'medium':2, 'large':3})
pd_sizes
```

	sizes	sizes_codes
0	small	1
1	medium	2
2	large	3
3	small	1
4	medium	2
5	large	3
6	small	1
7	medium	2
8	large	3

```
pd_sizes['sizes_decoded'] = pd_sizes['sizes_codes'].map({1:'small',
2:'medium', 3:'large'})
pd_sizes
```

	sizes	sizes_codes	sizes_decoded
0	small	1	small
1	medium	2	medium
2	large	3	large
3	small	1	small
4	medium	2	medium
5	large	3	large
6	small	1	small
7	medium	2	medium
8	large	3	large

Кодирование категорий наборами бинарных значений - one-hot encoding

Каждое уникальное значение признака становится новым отдельным признаком:

```
from sklearn.preprocessing import OneHotEncoder
```

```
ohe = OneHotEncoder()
cat_enc_ohe = ohe.fit_transform(cat_enc[['cat1']])
cat_enc_ohe
```

```
<157x27 sparse matrix of type '<class 'numpy.float64'>'
  with 157 stored elements in Compressed Sparse Row format>
```

```
cat_enc.shape
```

```
(157, 1)
```

```
cat_enc_ohe.shape
```

```
(157, 27)
```

```
cat_enc_ohe.todense()[0:10]
```

```
matrix([[1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
[0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
0.,
        0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.]])
```

```
cat_enc.head(10)
```

```
   cat1
0  Acura
1  Acura
2  Acura
3  Acura
4   Audi
5   Audi
6   Audi
7   BMW
8   BMW
9   BMW
```

```
pd.get_dummies(cat_enc)
```

```
   cat1_Acura  cat1_Audi  cat1_BMW  cat1_Buick  cat1_Cadillac  \
0            1          0          0            0            0
1            1          0          0            0            0
2            1          0          0            0            0
```

3	1	0	0	0	0
4	0	1	0	0	0
..
152	0	0	0	0	0
153	0	0	0	0	0
154	0	0	0	0	0
155	0	0	0	0	0
156	0	0	0	0	0

	cat1_Chevrolet	cat1_Dodge	cat1_Ford	cat1_Honda
cat1_Hyundai	...			
0	0	0	0	0
0	...			
1	0	0	0	0
0	...			
2	0	0	0	0
0	...			
3	0	0	0	0
0	...			
4	0	0	0	0
0	...			
..
...				...
152	0	0	0	0
0	...			
153	0	0	0	0
0	...			
154	0	0	0	0
0	...			
155	0	0	0	0
0	...			
156	0	0	0	0
0	...			

	cat1_Nissan	cat1_Oldsmobile	cat1_Plymouth	cat1_Pontiac
cat1_Porsche	\			
0	0	0	0	0
0				
1	0	0	0	0
0				
2	0	0	0	0
0				
3	0	0	0	0
0				
4	0	0	0	0
0				
..
...				
152	0	0	0	0
0				

153	0	0	0	0	
0					
154	0	0	0	0	
0					
155	0	0	0	0	
0					
156	0	0	0	0	
0					
	cat1_Saab	cat1_Subaru	cat1_Toyota	cat1_Volkswagen	cat1_Volvo
0	0	0	0	0	0
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0
..
152	0	0	0	0	1
153	0	0	0	0	1
154	0	0	0	0	1
155	0	0	0	0	1
156	0	0	0	0	1

[157 rows x 27 columns]

```
pd.get_dummies(cat_temp_data, dummy_na=True)
```

	Manufacturer_Acura	Manufacturer_Audi	Manufacturer_BMW	\
0	1	0	0	
1	1	0	0	
2	1	0	0	
3	1	0	0	
4	0	1	0	
..	
152	0	0	0	
153	0	0	0	
154	0	0	0	
155	0	0	0	

156	0	0	0
	Manufacturer_Buick	Manufacturer_Cadillac	Manufacturer_Chevrolet
\			
0	0	0	0
1	0	0	0
2	0	0	0
3	0	0	0
4	0	0	0
..
152	0	0	0
153	0	0	0
154	0	0	0
155	0	0	0
156	0	0	0

	Manufacturer_Dodge	Manufacturer_Ford	Manufacturer_Honda	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	
..	
152	0	0	0	
153	0	0	0	
154	0	0	0	
155	0	0	0	
156	0	0	0	

	Manufacturer_Hyundai	...	Manufacturer_Oldsmobile	\
0	0	...	0	
1	0	...	0	
2	0	...	0	
3	0	...	0	
4	0	...	0	
..	
152	0	...	0	
153	0	...	0	

154	0	...	0
155	0	...	0
156	0	...	0

	Manufacturer_Plymouth	Manufacturer_Pontiac	Manufacturer_Porsche
\			
0	0	0	0
1	0	0	0
2	0	0	0
3	0	0	0
4	0	0	0
..
152	0	0	0
153	0	0	0
154	0	0	0
155	0	0	0
156	0	0	0

	Manufacturer_Saab	Manufacturer_Subaru	Manufacturer_Toyota	\
0	0	0	0	
1	0	0	0	
2	0	0	0	
3	0	0	0	
4	0	0	0	
..	
152	0	0	0	
153	0	0	0	
154	0	0	0	
155	0	0	0	
156	0	0	0	

	Manufacturer_Volkswagen	Manufacturer_Volvo	Manufacturer_nan
0	0	0	0
1	0	0	0
2	0	0	0
3	0	0	0
4	0	0	0
..

152	0	1	0
153	0	1	0
154	0	1	0
155	0	1	0
156	0	1	0

[157 rows x 28 columns]

Масштабирование данных

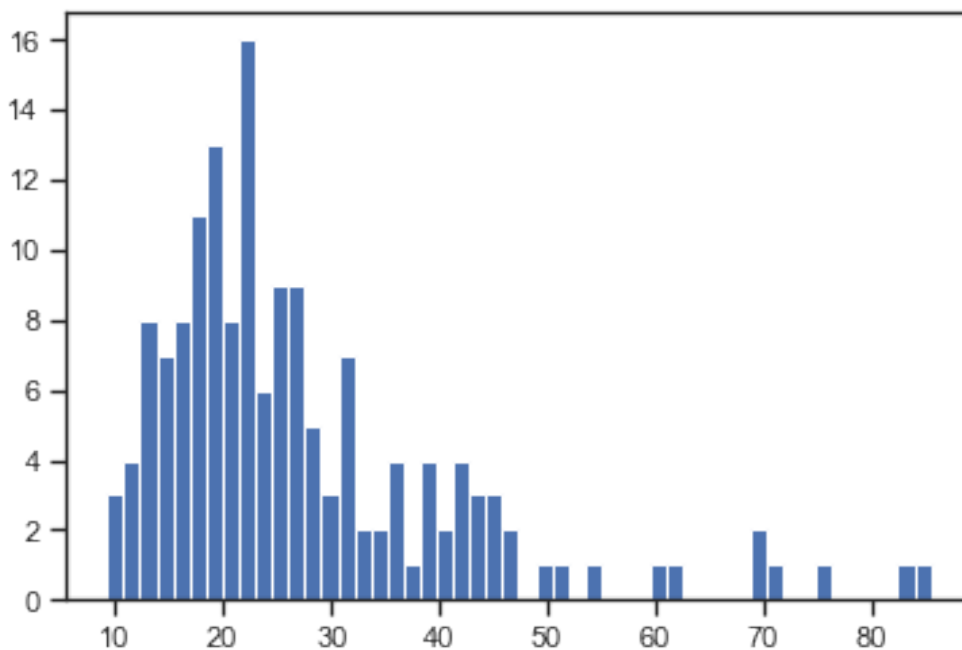
Масштабирование предполагает изменение диапазона измерения величины.

```
from sklearn.preprocessing import MinMaxScaler, StandardScaler,
Normalizer
```

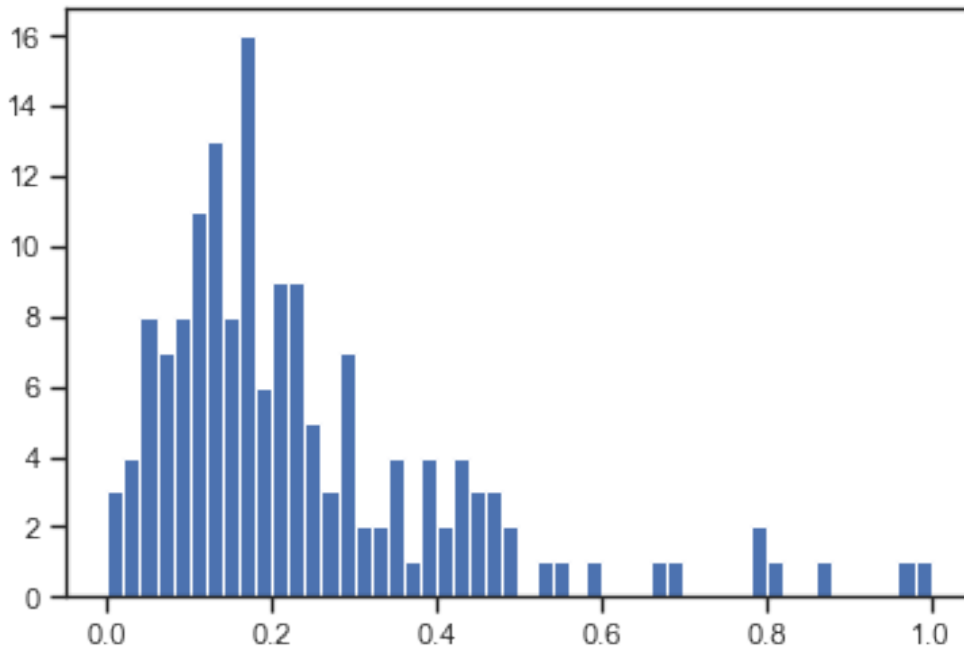
MinMax масштабирование

```
sc1 = MinMaxScaler()
sc1_data = sc1.fit_transform(data[['Price_in_thousands']])

plt.hist(data['Price_in_thousands'], 50)
plt.show()
```



```
plt.hist(sc1_data, 50)
plt.show()
```

Масштабирование данных на основе Z-оценки

```
sc2 = StandardScaler()  
sc2_data = sc2.fit_transform(data[['Price_in_thousands']])  
  
plt.hist(sc2_data, 50)  
plt.show()
```

