

```

1 import numpy as np
2 import pandas as pd
3

1 !unzip /content/drive/MyDrive/Colab_data/MMO/archive.zip

Archive: /content/drive/MyDrive/Colab_data/MMO/archive.zip
  inflating: autoru_total.csv


```

```
1 data = pd.read_csv('/content/autoru_total.csv', sep=',')
```


```
1 data.shape
```

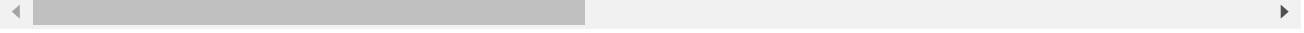
```
(36906, 17)
```

```
1 data.head()
```



	Model	Year	Mileage	V_engine	EngineType	HorsePower	Tax	State	
0	Subaru Forester IV	2013	83800.0	2.5	Бензин	171.0	8379.0	Не требует ремонта	вла,
1	Opel Zafira B Рестайлинг	2014	97265.0	1.8	Бензин	140.0	4900.0	Не требует ремонта	вл,
2	Kia Rio IV	2017	48000.0	1.6	Бензин	123.0	3075.0	Не требует ремонта	вла,
3	Skoda Octavia II (A5) Рестайлинг	2011	292000.0	1.4	Бензин	122.0	4148.0	Не требует ремонта	вла,
4	Audi A6 IV (C7) Рестайлинг	2015	106205.0	1.8	Бензин	190.0	9500.0	Не требует ремонта	вла,





```

1 hcols_with_na = [c for c in data.columns if data[c].isnull().sum() > 0]
2 hcols_with_na

```

```
['Tax', 'State', 'Owners', 'Customs', 'Accidents']
```

```
1 [(c, data[c].isnull().sum()) for c in hcols_with_na]
```

```
[('Tax', 476),
```

```

('State', 4),
('Owners', 2),
('Customs', 1),
('Accidents', 28527)]

```

```

1 # Доля (процент) пропусков
2 [(c, data[c].isnull().mean()) for c in hcols_with_na]
3

```

```

[('Tax', 0.012897631821384056),
 ('State', 0.00010838346068389963),
 ('Owners', 5.419173034194982e-05),
 ('Customs', 2.709586517097491e-05),
 ('Accidents', 0.7729637457324012)]

```

## Удаление пропущенных значений

```

1
2 data_drop = data['Customs'].dropna()
3 data_drop.shape

```

```

(36905,)

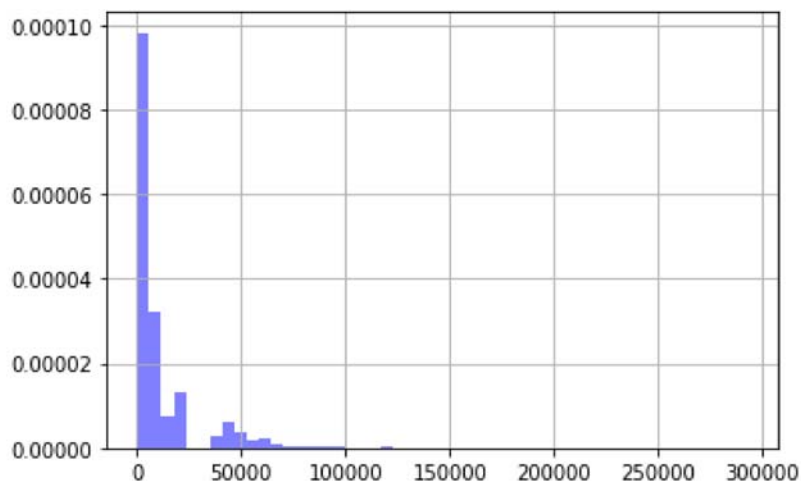
```

## Заполнение показателями центра распределения

```

1 import matplotlib.pyplot as plt
2 %matplotlib inline
3
4 fig = plt.figure()
5 ax = fig.add_subplot(111)
6 data['Tax'].hist(bins=50, ax=ax, color='blue', density=True, alpha=0.5)
7 plt.show()

```



```

1 from sklearn.impute import SimpleImputer
2 from sklearn.impute import MissingIndicator
3
4 indicator = MissingIndicator()

```

```

5 mask_missing_values_only = indicator.fit_transform(data[['Tax']].values)
6
7 imputer = SimpleImputer(strategy='mean')
8 all_data = imputer.fit_transform(data[['Tax']].values)
9 all_data

array([[8379.],
       [4900.],
       [3075.],
       ...,
       [4760.],
       [1176.],
       [9500.]])

```

## Заполнение наиболее распространенным значением категории

```

1 indicator = MissingIndicator()
2 mask_missing_values_only = indicator.fit_transform(data[['State']].values)
3
4 imputer = SimpleImputer(strategy='most_frequent')
5 all_data = imputer.fit_transform(data[['State']].values)
6 filled_data = all_data[mask_missing_values_only]
7
8 print(filled_data)
9 all_data

['Не требует ремонта' 'Не требует ремонта' 'Не требует ремонта'
 'Не требует ремонта']
array(['Не требует ремонта',
       'Не требует ремонта',
       'Не требует ремонта',
       ...,
       'Не требует ремонта',
       'Не требует ремонта',
       'Не требует ремонта'], dtype=object)

```

## Введение отдельного значения категории для пропущенных значений

```

1 indicator = MissingIndicator()
2 mask_missing_values_only = indicator.fit_transform(data[['Accidents']].values)
3
4 imputer = SimpleImputer(strategy='constant', fill_value='NA')
5 all_data = imputer.fit_transform(data[['Accidents']].values)
6 all_data

array(['ДТП не найдены',
       'NA',
       'NA',
       ...,
       'NA',
       'NA',
       'ДТП не найдены'], dtype=object)

```

## Добавление флагов пропусков

```
1 hdata_mis = data[['Tax']].copy()
2 indicator = MissingIndicator()
3 Tax_missing = indicator.fit_transform(hdata_mis[['Tax']])
4 Tax_missing
```

```
array([[False],
       [False],
       [False],
       ...,
       [False],
       [False],
       [False]])
```

```
1 hdata_mis['flag'] = Tax_missing
2 hdata_mis.head()
```

	Tax	flag
0	8379.0	False
1	4900.0	False
2	3075.0	False
3	4148.0	False
4	9500.0	False

## ▼ Кодирование категориальных признаков

```
1 from sklearn.preprocessing import LabelEncoder
2
3 le = LabelEncoder()
4 cat_enc_le = le.fit_transform(data['Model'])
```

```
1 data['Model'].unique()
```

```
array(['Subaru Forester IV', 'Opel Zafira B Рестайлинг', 'Kia Rio IV',
       ..., 'Mitsubishi Chariot II', 'Porsche 911 Carrera VI (997)',
       'BMW 3 серии Gran Turismo 335i VI (F3x)'], dtype=object)
```

```
1 np.unique(cat_enc_le)
```

```
array([ 0, 1, 2, ..., 3160, 3161, 3162])
```

```
1 le.inverse_transform(np.unique(cat_enc_le))
```

```
array(['Audi 100 III (C3)', 'Audi 100 III (C3) Рестайлинг',
```

```
'Audi 100 IV (C4)', ..., 'Volkswagen Type 1',  
'Volkswagen Type 2 T2', 'Volkswagen Vento'], dtype=object)
```

## Кодирование категорий наборами бинарных значений - one-hot encoding

```
1 from sklearn.preprocessing import OneHotEncoder  
2  
3 ohe = OneHotEncoder()  
4 cat_enc_ohe = ohe.fit_transform(data[['Model']])  
5 cat_enc_ohe  
  
<36906x3163 sparse matrix of type '<class 'numpy.float64'>'  
  with 36906 stored elements in Compressed Sparse Row format>
```

```
1 cat_enc_ohe.todense()[0:10]  
  
matrix([[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., ..., 0., 0., 0.],  
        [0., 0., 0., ..., 0., 0., 0.]])
```

```
1 !pip install category_encoders
```

Collecting category\_encoders

Downloading category\_encoders-2.4.0-py2.py3-none-any.whl (86 kB)

|██| 86 kB 3.5 MB/s

Requirement already satisfied: statsmodels>=0.9.0 in /usr/local/lib/python3.7/dist-p

Requirement already satisfied: pandas>=0.21.1 in /usr/local/lib/python3.7/dist-packa

Requirement already satisfied: patsy>=0.5.1 in /usr/local/lib/python3.7/dist-package

Requirement already satisfied: numpy>=1.14.0 in /usr/local/lib/python3.7/dist-packag

Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.7/dist-package

Requirement already satisfied: scikit-learn>=0.20.0 in /usr/local/lib/python3.7/dist

Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/dist-package

Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/di

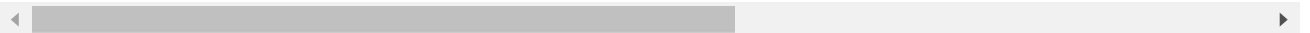
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from p

Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist

Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-package

Installing collected packages: category-encoders

Successfully installed category-encoders-2.4.0



## Count (frequency) encoding

```
1 from category_encoders.count import CountEncoder as ce_CountEncoder  
2  
3 ce_CountEncoder1 = ce_CountEncoder()  
4 data_COUNT_ENC = ce_CountEncoder1.fit_transform(data[data.columns.difference(['Model'])])  
5  
6 data_COUNT_ENC
```

	Accidents	CarBodyType	Color	Customs	Drive	EngineType	HorsePower	Mileage
0	8379	13303	3941	36905	14584	30981	171.0	83800
1	28527	946	9789	36905	20236	30981	140.0	97265
2	28527	11985	3594	36905	20236	30981	123.0	48000
3	28527	1870	3941	36905	20236	30981	122.0	292000
4	28527	11985	9789	36905	20236	30981	190.0	106205
...	...	...	...	...	...	...	...	...
36901	28527	11985	7264	36905	20236	30981	102.0	133789
36902	28527	13303	9789	36905	14584	5138	249.0	103965
36903	28527	11985	7264	36905	2086	30981	136.0	108000
36904	28527	11985	9789	36905	20236	30981	98.0	259000
36905	8379	11985	5173	36905	20236	30981	190.0	265000

36906 rows × 16 columns



## Target (Mean) encoding

```

1 from category_encoders.target_encoder import TargetEncoder as ce_TargetEncoder
2
3 ce_TargetEncoder1 = ce_TargetEncoder()
4 data_MEAN_ENC = ce_TargetEncoder1.fit_transform(data[data.columns.difference(['Tax'])]),
5 data_MEAN_ENC

```

	Accidents	CarBodyType	Color	Customs	Drive	EngineT
0	12194.545894	17928.761822	7368.860160	12357.519229	22441.953945	12086.564
1	12405.048242	3726.181818	18208.519405	12357.519229	4326.840010	12086.564
2	12405.048242	8045.025446	11227.700526	12357.519229	4326.840010	12086.564

## ▼ Нормализация числовых признаков

3	12405.048242	8045.025446	18208.519405	12357.519229	4326.840010	12086.564
---	--------------	-------------	--------------	--------------	-------------	-----------

```

1 import scipy.stats as stats
2
3 def diagnostic_plots(df, variable):
4     plt.figure(figsize=(15,6))
5     # гистограмма
6     plt.subplot(1, 2, 1)
7     df[variable].hist(bins=30)
8     ## Q-Q plot
9     plt.subplot(1, 2, 2)
10    stats.probplot(df[variable], dist="norm", plot=plt)
11    plt.show()

```



```

1 norm_data = data[['Tax']].dropna()
2 norm_data

```

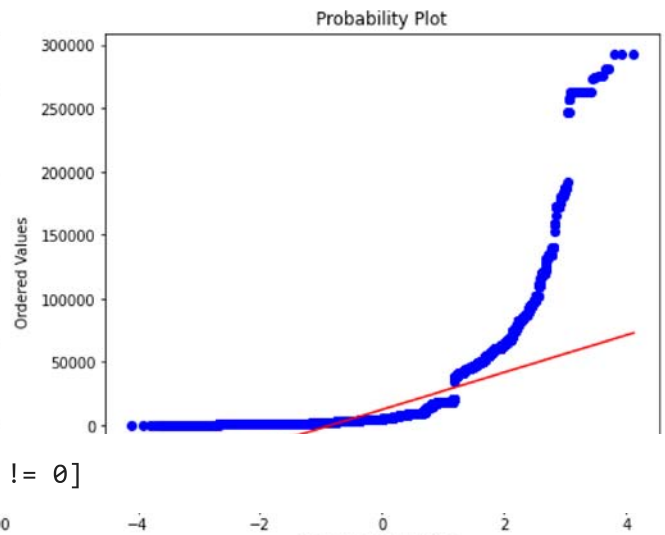
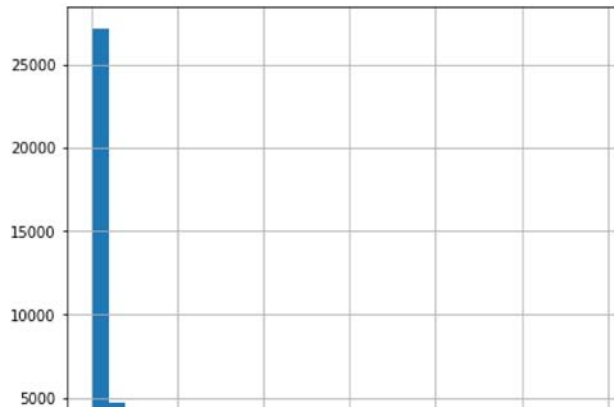
	Tax
0	8379.0
1	4900.0
2	3075.0
3	4148.0
4	9500.0
...	...
36901	2550.0
36902	18675.0
36903	4760.0
36904	1176.0
36905	9500.0

36430 rows × 1 columns

```

1 diagnostic_plots(norm_data, 'Tax')

```



```
1 norm_data = norm_data.loc[norm_data['Tax'] != 0]
```

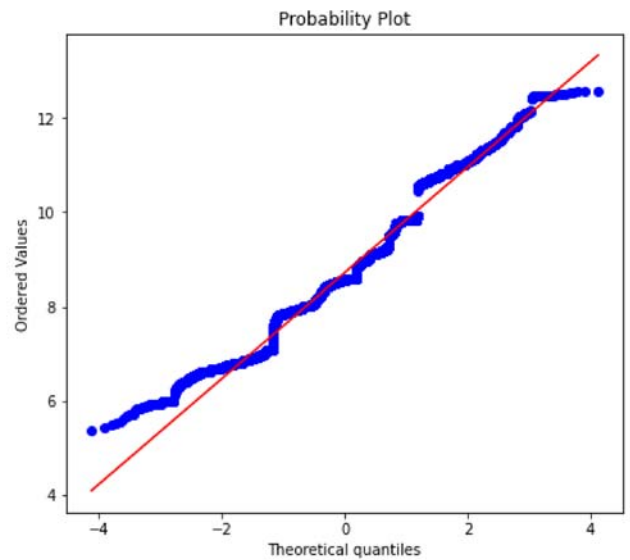
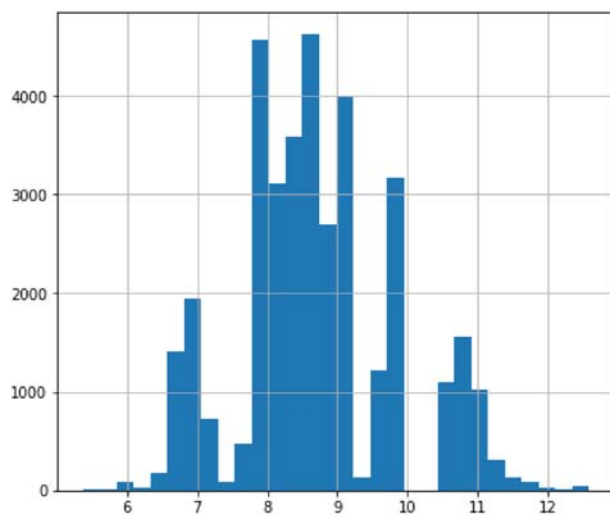
```
1 norm_data['Tax_log'] = np.log(norm_data['Tax'])
```

```
2
```

```
3
```

```
4 norm_data['Tax_log'].min()
```

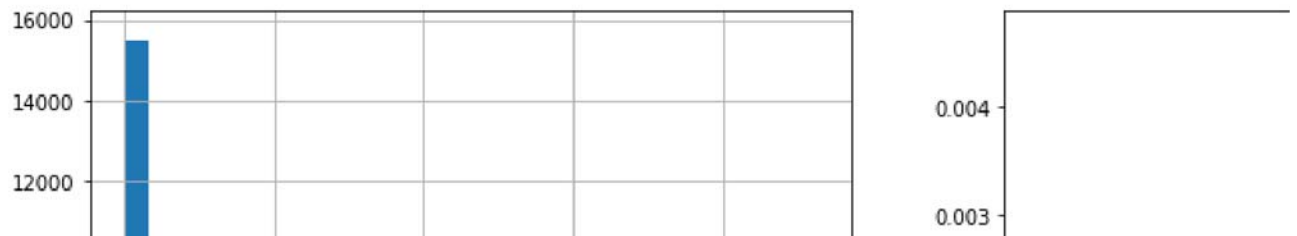
```
5 diagnostic_plots(norm_data, 'Tax_log')
```



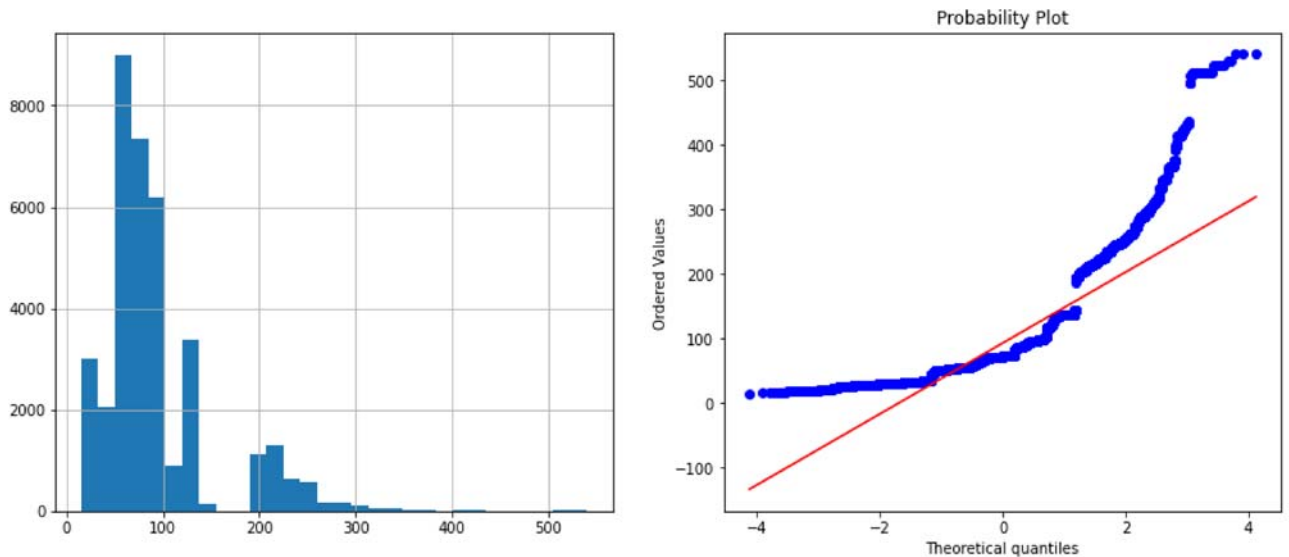
```
1 norm_data['reciprocal'] = 1 / (norm_data['Tax'])
```

```
2 diagnostic_plots(norm_data, 'reciprocal')
```



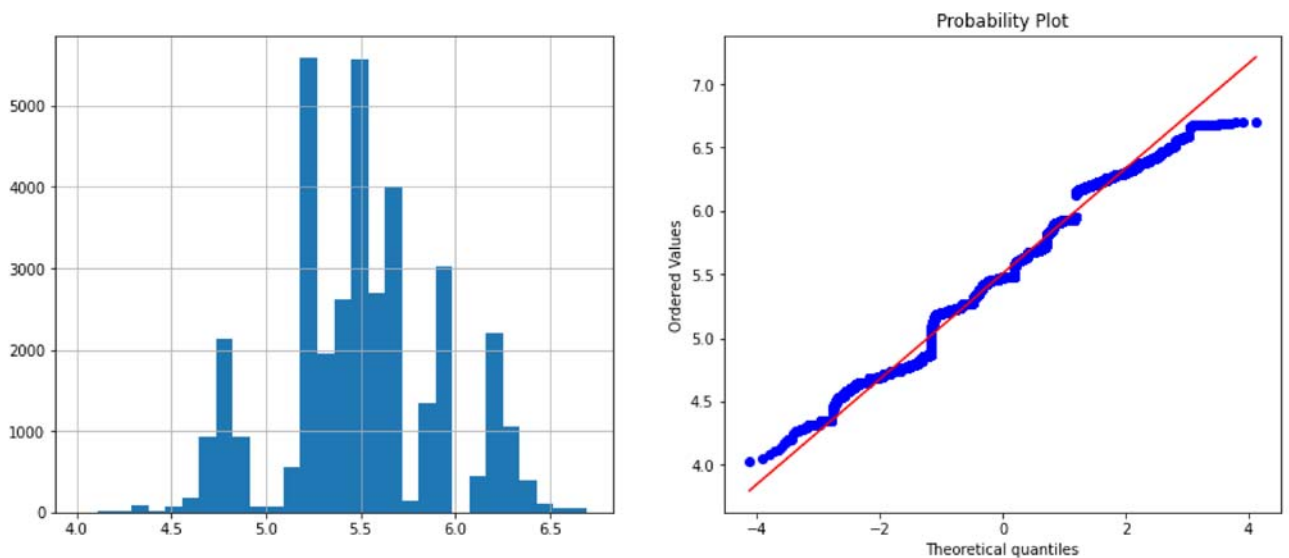


```
1 norm_data['sqr'] = norm_data['Tax']**(1/2)
2 diagnostic_plots(norm_data, 'sqr')
```

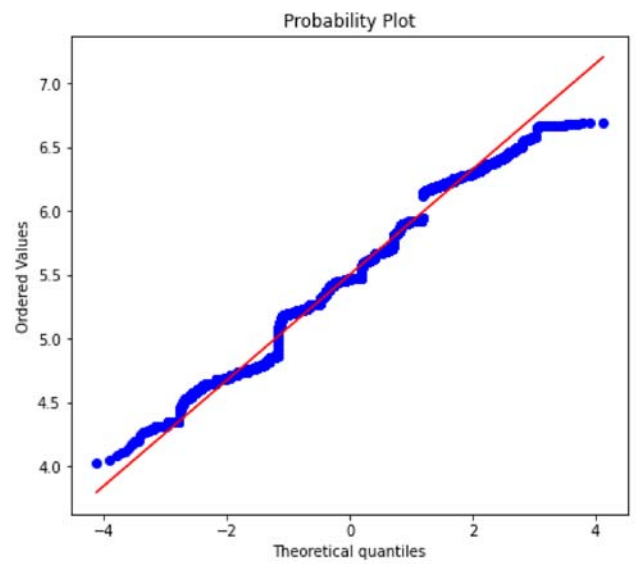
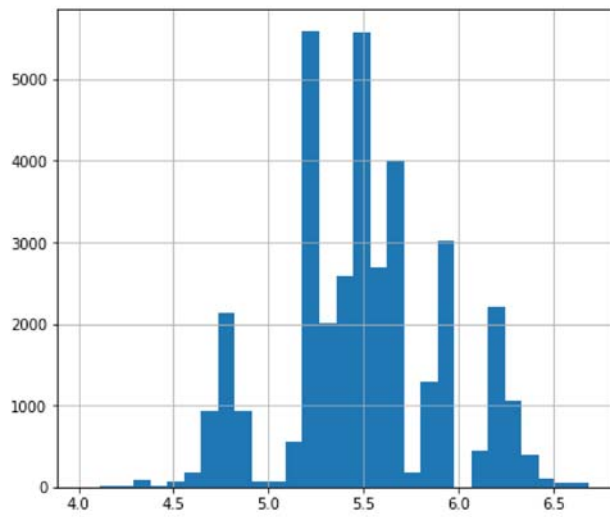


```
1 norm_data['boxcox'], param = stats.boxcox(norm_data['Tax'])
2 print('Оптимальное значение  $\lambda$  = {}'.format(param))
3 diagnostic_plots(norm_data, 'boxcox')
```

Оптимальное значение  $\lambda$  = -0.11353789281093771



```
1 norm_data['Tax'] = norm_data['Tax'].astype('float')
2 norm_data['yeojohnson'], param = stats.yeojohnson(norm_data['Tax'])
3 diagnostic_plots(norm_data, 'yeojohnson')
```



✓ 1 сек. выполнено в 23:31

● ✕