

Преобразование данных. Часть 2

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

Во второй части рассмотрим нелинейные преобразования количественных данных.

Продолжим работу в том же блокноте

Нелинейные преобразования

Нелинейные преобразования, как уже было сказано, меняют структуру распределения.

```
1 # вновь подгрузим полный датасет boston
2 boston = pd.read_csv('/content/boston.csv')
```

Логарифмическое преобразование

Смысл преобразования

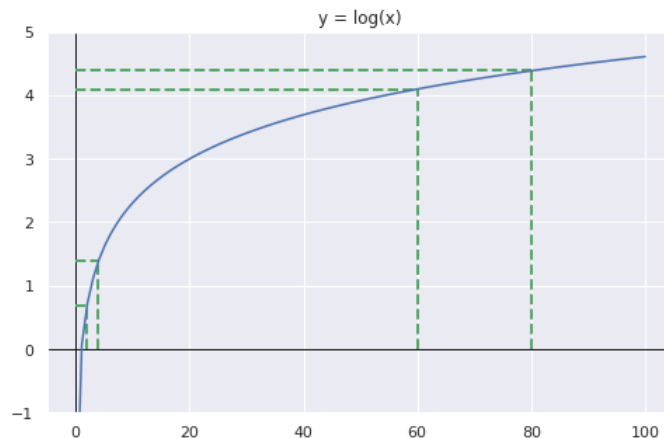
Рассмотрим график логарифмической функции.

```
1 # построим график логарифмической функции
2 x = np.linspace(0.05, 100, 100)
3 y = np.log(x)
4
5 ax = plt.axes()
6
7 plt.xlim([-5, 105])
8 plt.ylim([-1, 5])
9
10 ax.hlines(y = 0, xmin = -5, xmax = 105, linewidth = 1, color = 'k')
11 ax.vlines(x = 0, ymin = -1, ymax = 5, linewidth = 1, color = 'k')
12
13 plt.plot(x, y)
14
15 # возьмем произвольные промежутки между малыми
16 ax.vlines(x = 2, ymin = 0, ymax = np.log(2), linewidth = 2, color = 'g', linestyle =
17 ax.vlines(x = 4, ymin = 0, ymax = np.log(4), linewidth = 2, color = 'g', linestyle =
18 ax.hlines(y = np.log(2), xmin = 0, xmax = 2, linewidth = 2, color = 'g', linestyle =
19 ax.hlines(y = np.log(4), xmin = 0, xmax = 4, linewidth = 2, color = 'g', linestyle =
20
21 # и большими значениями
22 ax.vlines(x = 60, ymin = 0, ymax = np.log(60), linewidth = 2, color = 'g', linestyle
23 ax.vlines(x = 80, ymin = 0, ymax = np.log(80), linewidth = 2, color = 'g', linestyle
24 ax.hlines(y = np.log(60), xmin = 0, xmax = 60, linewidth = 2, color = 'g', linestyle
25 ax.hlines(y = np.log(80), xmin = 0, xmax = 80, linewidth = 2, color = 'g', linestyle
```

```

26
27 plt.title('y = log(x)')
28
29 plt.show()

```



По оси x располагаются значения до трансформации, по оси y — после. Ценность логарифмического преобразования в том, что:

- расстояние между небольшими значениями увеличивается;
- расстояние между большими значениями наоборот уменьшается.

И, таким образом, это преобразование делает скошенное распределение более симметричным и приближенным к нормальному. Замечу, что, как видно из графика, в общем случае преобразование возможно только для положительных исходных значений.

Скошенное вправо распределение

В силу описанной выше особенности логарифмическое преобразование чаще применяют к **скошенным вправо** (right-skewed) распределениям. В этих распределениях большая часть наблюдений находится как раз в диапазоне меньших значений.

```

1 fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
2
3 sns.histplot(x = boston.LSTAT, bins = 15, ax = ax[0])
4 ax[0].set_title('Скошенное вправо распределение')
5
6 sns.histplot(x = np.log(boston.LSTAT),
7             bins = 15, color = 'green',
8             ax = ax[1])
9 ax[1].set_title('Log transformation')
10
11 plt.tight_layout()
12 plt.show()

```





Количественно оценим изменения в скошенности и островершинности (то есть сосредоточенности значений вокруг среднего) распределения.

Коэффициент асимметрии

Первое свойство измеряется **коэффициентом асимметрии** (skewness), который в нормальном распределении должен быть равен нулю. При этом,

- положительные значения говорят о скошенности вправо (positively или right-skewed);
- отрицательные, о скошенности влево (negatively или left-skewed).

```
1 # импортируем необходимые функции
2 from scipy.stats import skew, kurtosis

1 # рассчитаем асимметричность до и после преобразования
2 skew(boston.LSTAT), skew(np.log(boston.LSTAT))

1 (0.9037707431346133, -0.3192822699479382)
```

Выраженная скошенность вправо превратилась в меньшую скошенность влево.

Коэффициент эксцесса

Коэффициент эксцесса (kurtosis) измеряет островершинность распределения. Можно сказать, что эксцесс показывает сосредоточенность (плотность) значений вокруг среднего.

По *формуле Фишера* (Fisher's definition) для нормального распределения значение этого коэффициента также равно нулю. Одновременно,

- положительные значения говорят о большей сосредоточенности значений около среднего (острая вершина);
- отрицательные — о наличии более «тяжелых хвостов» (плоская вершина).

```
1 # рассчитаем коэффициент эксцесса до и после преобразования
2 kurtosis(boston.LSTAT), kurtosis(np.log(boston.LSTAT))

1 (0.476544755729746, -0.4390590293275558)
```

Вершина сменилась с более острой на более плоскую.

График нормальной вероятности

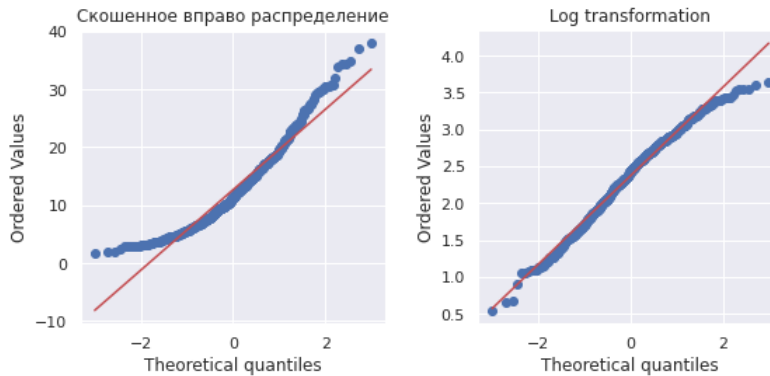
Наконец с помощью **графика нормальной вероятности** (normal probability plot) мы можем визуально оценить, приблизилось ли распределение к нормальному.

```
1 # построим графики нормальной вероятности
2 from scipy.stats import probplot
3
```

```

4 fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
5
6 probplot(boston.LSTAT, dist = 'norm', plot = ax[0])
7 ax[0].set_title('Скошенное вправо распределение')
8
9 probplot(np.log(boston.LSTAT), dist = 'norm', plot = ax[1])
10 ax[1].set_title('Log transformation')
11
12 plt.tight_layout()
13 plt.show()

```

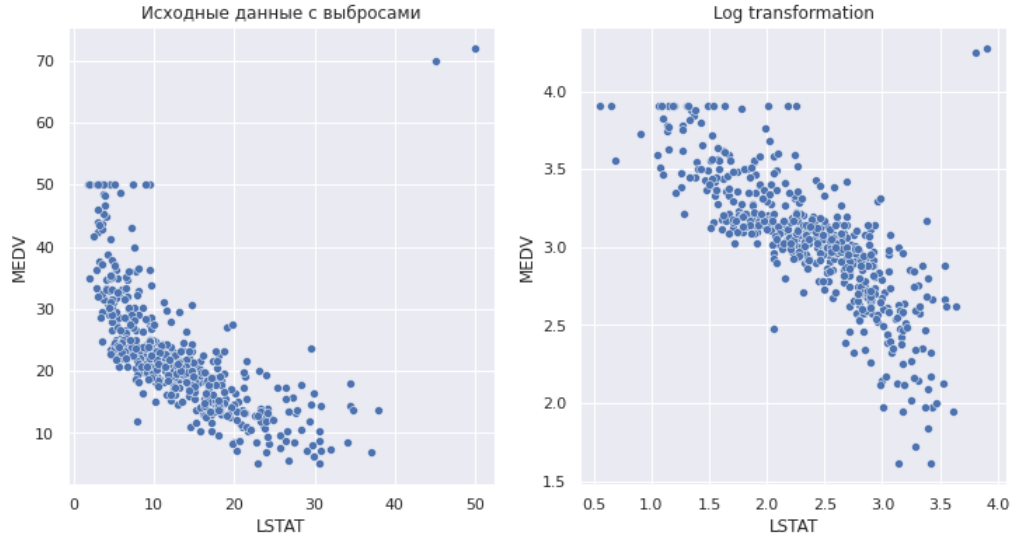


Разумеется, логарифмическое преобразование снижает эффект выбросов «справа».

```

1 fig, ax = plt.subplots(1, 2, figsize = (12,6))
2
3 sns.scatterplot(x = boston_outlier.LSTAT, y = boston_outlier.MEDV, ax = ax[0]).set(title='Исходные данные с выбросами')
4 sns.scatterplot(x = np.log(boston_outlier.LSTAT), y = np.log(boston_outlier.MEDV), ax = ax[1]).set(title='Log transformation')

```



Скошенное влево распределение

Логарифмическое преобразование, скорее всего, не подойдет для скошенного влево распределения, потому что здесь наоборот большая часть наблюдений находится в правой части диапазона. Применив лог-преобразование, мы только увеличим скошенность.

```

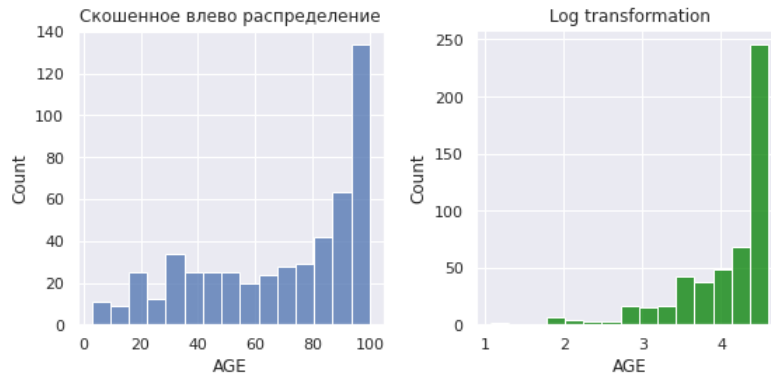
1 fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
2
3 sns.histplot(x = boston.AGE, bins = 15, ax = ax[0])
4 ax[0].set_title('Скошенное влево распределение')
5

```

```

6 | sns.histplot(x = np.log(boston.AGE),
7 |             bins = 15, color = 'green',
8 |             ax = ax[1])
9 | ax[1].set_title('Log transformation')
10 |
11 | plt.tight_layout()
12 | plt.show()

```



Посмотрим на коэффициент асимметрии, коэффициент эксцесса и график нормальной вероятности.

```
1 | skew(boston.AGE), skew(np.log(boston.AGE))
```

```
1 | (-0.5971855948016143, -1.6706835909283215)
```

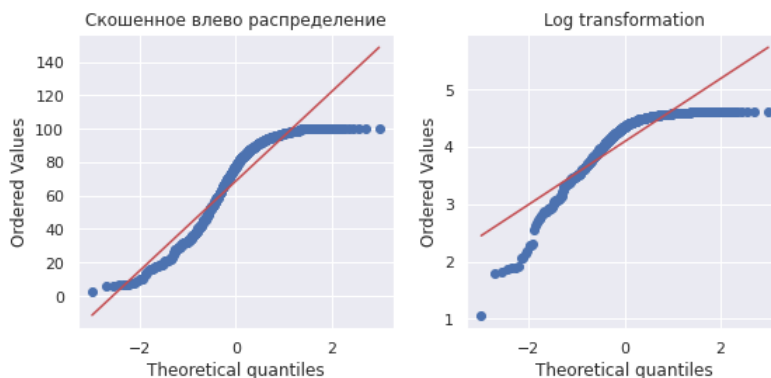
```
1 | kurtosis(boston.AGE), kurtosis(np.log(boston.AGE))
```

```
1 | (-0.97001392664039, 2.907332087827127)
```

```

1 | fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
2 |
3 | probplot(boston.AGE, dist = 'norm', plot = ax[0])
4 | ax[0].set_title('Скошенное влево распределение')
5 |
6 | probplot(np.log(boston.AGE), dist = 'norm', plot = ax[1])
7 | ax[1].set_title('Log transformation')
8 |
9 | plt.tight_layout()
10 | plt.show()

```



Как мы видим, показатели только ухудшились.

Большие значения

Дополнительно замечу, что даже если распределение скошено вправо, но в нем присутствуют

только большие (удаленные от нуля) значения, то логарифмическое преобразование может не приблизить распределение к нормальному, поскольку эффект расширения и сжатия на больших значениях менее заметен.

Логарифм нуля и отрицательных значений

Так как логарифм нуля и отрицательных значений неопределен, при наличии нулевых значений мы можем **добавить к значениям константу** (например, $\log(x + 0.001)$).

```
1 # в переменной ZN есть нулевые значения
2 # добавим небольшую константу
3 np.log(boston.ZN + 0.0001)[:5]
```

```
1 0    2.890377
2 1   -9.210340
3 2   -9.210340
4 3   -9.210340
5 4   -9.210340
6 Name: ZN, dtype: float64
```

Так как в данном случае мы все-таки произвольным образом меняем исходные данные, то можно использовать **преобразование обратного гиперболического синуса** (inverse hyperbolic sine (IHS) transformation).

$$IHS(x) = \ln \left(x + \sqrt{x^2 + 1} \right)$$

```
1 np.log(boston.ZN + np.sqrt(boston.ZN ** 2 + 1))[:5]
```

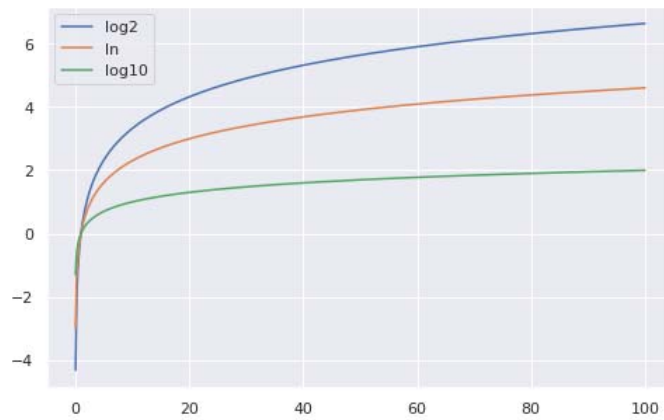
```
1 0    3.58429
2 1    0.00000
3 2    0.00000
4 3    0.00000
5 4    0.00000
6 Name: ZN, dtype: float64
```

Дополнительным преимуществом такого подхода является то, что мы можем преобразовывать любые действительные числа (а не только неотрицательные).

Основание логарифма

Чем меньше основание логарифма ($2 < e < 10$), тем больше диапазон преобразованных значений. Это можно наблюдать на графиках логарифмических функций, где функция с бóльшим основанием становится «плоской» быстрее.

```
1 x = np.linspace(0.05, 100, 500)
2 y_2 = np.log2(x)
3 y_ln = np.log(x)
4 y_10 = np.log10(x)
5
6 plt.plot(x, y_2, label = 'log2')
7 plt.plot(x, y_ln, label = 'ln')
8 plt.plot(x, y_10, label = 'log10')
9
10 plt.legend()
11
12 plt.show()
```



При этом разумеется в целом логарифмическое преобразование с любым основанием действует примерно одинаково.

Линейная взаимосвязь

Многие процессы удобно моделировать с помощью линейных моделей. И хотя, как мы увидим в следующем разделе, можно использовать линейную модель с полиномиальными коэффициентам, зачастую проще наоборот «выпрямить» (straighten) сами данные.

```
1 # визуально оценим "выпрямление" данных
2 fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
3
4 sns.scatterplot(x = boston.LSTAT, y = boston.MEDV, ax = ax[0])
5 ax[0].set_title('Изначальное распределение')
6
7 sns.scatterplot(x = np.log(boston.LSTAT), y = boston.MEDV, ax = ax[1])
8 ax[1].set_title('Log transformation')
9
10 plt.tight_layout()
11
12 plt.show()
```



В качестве метрики линейности взаимосвязи можно использовать коэффициент линейной корреляции Пирсона (максимизация коэффициента корреляции означает максимизацию линейности).

```
1 # посмотрим, как изменится корреляция, если преобразовать
2 # одну, вторую или сразу обе переменные
3 boston['LSTAT_log'] = np.log(boston['LSTAT'])
4 boston['MEDV_log'] = np.log(boston['MEDV'])
5
6 boston[['LSTAT', 'LSTAT_log', 'MEDV', 'MEDV_log']].corr()
```

	LSTAT	LSTAT_log	MEDV	MEDV_log
LSTAT	1.000000	0.944031	-0.737663	-0.805034
LSTAT_log	0.944031	1.000000	-0.815442	-0.822960
MEDV	-0.737663	-0.815442	1.000000	0.953155
MEDV_log	-0.805034	-0.822960	0.953155	1.000000

Как мы видим, в данном случае коэффициент корреляции будет наибольшим в том случае, когда мы преобразовываем обе переменные (−0,82).

Если вы преобразовываете целевую переменную, важно выполнить обратное преобразования после формирования прогноза. Для операции взятия натурального логарифма обратным преобразованием будет возведение числа Эйлера в степень преобразованного числа.

$$y = \ln(x) \rightarrow x = e^y$$

```
1 # сравним исходный датасет и лог-преобразование + обратную операцию
2 # (округлим значения, чтобы ошибка округления не мешала сравнению)
3 boston.MEDV.round(2).equals(np.exp(np.log(boston.MEDV)).round(2))
```

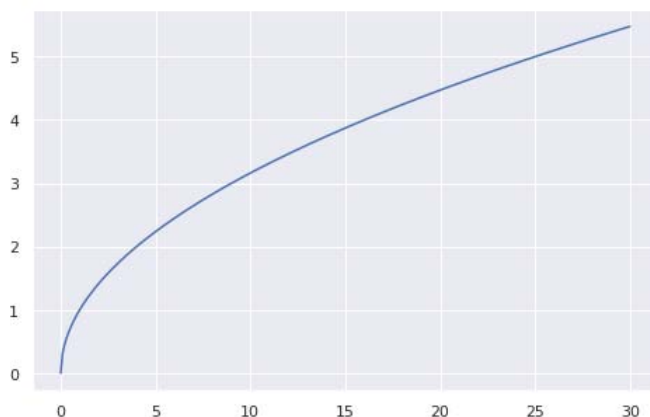
```
1 True
```

Примечание. Логарифмическое преобразование также легко интерпретировать. Например, снижение на −0,162 функции натурального логарифма свидетельствует о снижении на 15% в исходных данных вне зависимости от их масштаба.

Преобразование квадратного корня

Рассмотрим преобразование с помощью квадратного корня (square root transformation).

```
1 x = np.linspace(0, 30, 300)
2 y = np.sqrt(x)
3
4 plt.plot(x, y);
```



В целом, как видно из формы кривой, такое преобразование действует аналогично логарифмическому, однако менее агрессивно. С другой стороны, его без изменений можно применять к нулевым значениям.

Лестница степеней Тьюки

Помимо преобразования квадратного корня (то есть возведение в $\frac{1}{2}$ степени) можно пробовать и другие показатели. Для того чтобы понять, какое преобразование окажется наиболее удачным, используют **лестницу степеней Тьюки** (Tukey's Ladder of Powers).

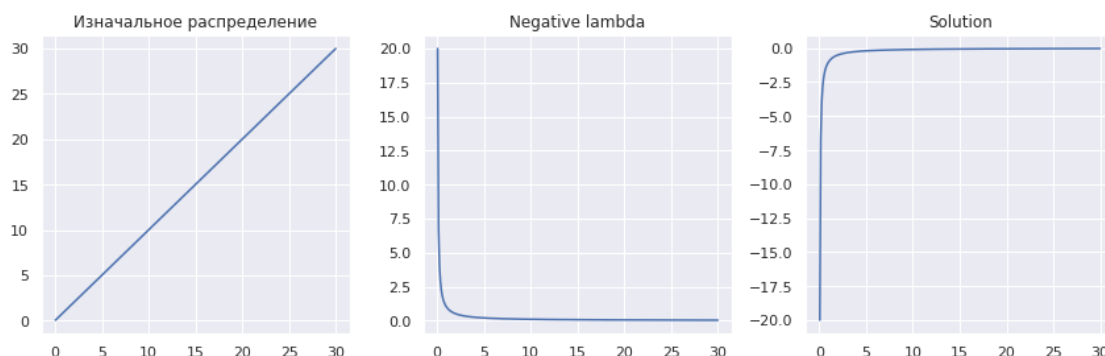
λ		-2	-1	$-\frac{1}{2}$	0	$\frac{1}{2}$	1	2
y		$\frac{1}{x^2}$	$\frac{1}{x}$	$\frac{1}{\sqrt{x}}$	$\log x$	\sqrt{x}	x	x^2

По сути, мы по очереди применяем каждое из значений λ и смотрим, что получится. Здесь важно вставить несколько примечаний:

- вместо возведения в нулевую степень применяется логарифмическое преобразование;
- если переменная принимает отрицательные значения, то после, например, возведения в квадрат и обратного преобразования, такое значение потеряет смысл;
- если параметр λ принимает отрицательные значения, то зависимость переворачивается; например, если $y = x$ — это возрастающая зависимость, то $\frac{1}{x}$ — убывающая; преодолеть это можно, прописав что $-(x^\lambda)$, если $\lambda < 0$.

```

1  x = np.linspace(0.05, 30, 300)
2
3  y0 = x
4  y1 = x ** (-1)
5  y2 = -(x ** (-1))
6
7  fig, ax = plt.subplots(nrows = 1, ncols = 3, figsize = (12,4))
8
9  ax[0].plot(x, y0);
10 ax[0].set_title('Изначальное распределение')
11
12 ax[1].plot(x, y1);
13 ax[1].set_title('Negative lambda')
14
15 ax[2].plot(x, y2);
16 ax[2].set_title('Solution')
17
18 plt.tight_layout()
19
20 plt.show()
```



Тогда,

$$y = \begin{cases} x^\lambda & \text{if } \lambda > 0 \\ \log x & \text{if } \lambda = 0 \\ -(x^\lambda) & \text{if } \lambda < 0 \end{cases}$$

λ		-2	-1	$-\frac{1}{2}$	0	$\frac{1}{2}$	1	2
y		$-\frac{1}{x^2}$	$-\frac{1}{x}$	$-\frac{1}{\sqrt{x}}$	$\log x$	\sqrt{x}	x	x^2

Этот инструмент удобно использовать, когда нужно «выправить» нелинейную взаимосвязь между переменными, например, в задаче линейной регрессии. Напишем функцию, которая на вход будет принимать признак и целевую переменную и выдавать оптимальную λ и соответствующий ей коэффициент корреляции.

```

1 def tukey(x, y):
2
3     x, y = x.to_numpy(), y.to_numpy()
4
5     # в lambdas поместим возможные степени
6     lambdas = [-2, -1, -1/2, 0, 1/2, 1, 2]
7     # в corrs будем записывать получающиеся корреляции
8     corrs = []
9
10    # в цикле последовательно применим каждую lambda
11    for l in lambdas:
12
13        if l < 0:
14            # рассчитаем коэффициент корреляции Пирсона и добавим результат в corrs
15            corrs.append(np.corrcoef(x ** l, y ** l)[0][1])
16
17        elif l == 0:
18            corrs.append(np.corrcoef(np.log(x + np.sqrt(x**2 + 1)),
19                                    np.log(y + np.sqrt(y**2 + 1)))[0][1])
20
21        else:
22            corrs.append(np.corrcoef(-(x ** l), -(y ** l))[0][1])
23
24    # теперь найдем индекс наибольшего значения корреляции
25    idx = np.argmax(np.abs(corrs))
26
27    # выведем оптимальную lambda и соответствующую корреляцию
28    return lambdas[idx], np.round(corrs[idx], 3)

```

```
1 tukey(boston.LSTAT, boston.MEDV)
```

```
1 (0, -0.824)
```

Оптимальным будет логарифмическое преобразование. Применим функцию к нескольким признакам с положительными значениями.

```

1 for col in boston[['CRIM', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'LSTAT']]
2     print(str(col) + '\t' + str(tukey(boston[col], boston.MEDV)))

```

```

1 CRIM    (0, -0.593)
2 NOX     (-0.5, -0.526)
3 RM      (2, 0.724)
4 AGE     (0.5, -0.402)
5 DIS     (-1, 0.489)
6 RAD     (0, -0.44)
7 TAX     (-0.5, -0.558)
8 PTRATIO (0.5, -0.509)
9 LSTAT   (0, -0.824)

```

Посмотрим на корреляцию изначальных переменных.

```

1 boston[['CRIM', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'LSTAT', 'MEDV']].c

```

1	CRIM	-0.39
2	NOX	-0.43
3	RM	0.70
4	AGE	-0.38
5	DIS	0.25
6	RAD	-0.38
7	TAX	-0.47
8	PTRATIO	-0.51
9	LSTAT	-0.74
10	MEDV	1.00
11	Name: MEDV, dtype: float64	

Для дальнейшей работы отберем RM, PTRATIO и LSTAT. Выполним необходимые преобразования.

```

1 boston_transformed = {}
2
3 boston_transformed['RM'] = boston.RM ** 2
4 boston_transformed['PTRATIO'] = np.sqrt(boston.PTRATIO)
5 boston_transformed['LSTAT'] = np.log(boston.LSTAT)
6 boston_transformed['MEDV'] = boston.MEDV
7
8 boston_transformed = pd.DataFrame(boston_transformed,
9                                   columns = ['RM', 'PTRATIO', 'LSTAT', 'MEDV'])
10
11 boston_transformed.head()

```

	RM	PTRATIO	LSTAT	MEDV
0	43.230625	3.911521	1.605430	24.0
1	41.229241	4.219005	2.212660	21.6
2	51.624225	4.219005	1.393766	34.7
3	48.972004	4.324350	1.078410	33.4
4	51.079609	4.324350	1.673351	36.2

Построим модель линейной регрессии на исходных и преобразованных данных и рассчитаем коэффициент детерминации.

```

1 from sklearn.linear_model import LinearRegression

```

```

1 model = LinearRegression()
2 model.fit(boston[['RM', 'PTRATIO', 'LSTAT']], boston.MEDV)
3 model.score(boston[['RM', 'PTRATIO', 'LSTAT']], boston.MEDV)

```

```

1 0.6786241601613112

```

```

1 model = LinearRegression()
2 model.fit(boston_transformed[['RM', 'PTRATIO', 'LSTAT']], boston_transformed.MEDV)
3 model.score(boston_transformed[['RM', 'PTRATIO', 'LSTAT']], boston_transformed.MEDV)

```

```

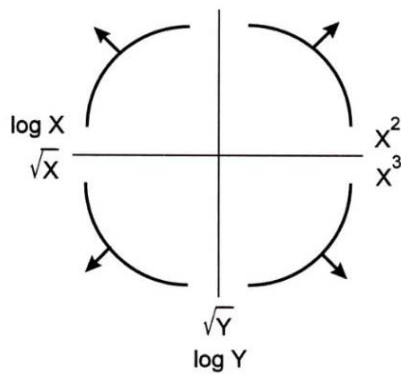
1 0.7446785206677597

```

Еще одним вариантом было бы применение различных преобразований к признаку и целевой переменной. В частности, Тьюки и Мостеллер предложили следующее правило выпуклости (**Tukey and Mosteller's Bulging Rule**).

$$Y^3$$

$$Y^2$$



В данном случае преобразования выбираются в зависимости от четырех приведенных на схеме форм зависимости данных.

Теперь рассмотрим более сложные, но близкие по смыслу к лестнице Тьюки степенные преобразования (power transformations) Бокса-Кокса и Йео-Джонсона.

Преобразование Бокса-Кокса

Приведем формулу **преобразования Бокса-Кокса** (Box-Cox transformation).

$$y_i^{(\lambda)} = \begin{cases} \frac{y_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln y_i & \text{if } \lambda = 0 \end{cases}$$

Очевидно, что преобразуемые значения y_i могут быть только положительными. Параметр λ выбирается методом наибольшего (максимального) правдоподобия (maximum likelihood method). Его рассмотрение выходит за рамки сегодняшнего занятия.

```
1 from sklearn.preprocessing import PowerTransformer
2
3 pt = PowerTransformer(method = 'box-cox')
4
5 # найдем оптимальный параметр лямбда
6 pt.fit(boston[['LSTAT']])
7
8 pt.lambdas_
```

```
1 array([0.22776737])
```

```
1 # преобразуем данные
2 bc_pt = pt.transform(boston[['LSTAT']])
3
4 # метод .transform() возвращает двумерный массив
5 bc_pt.shape
```

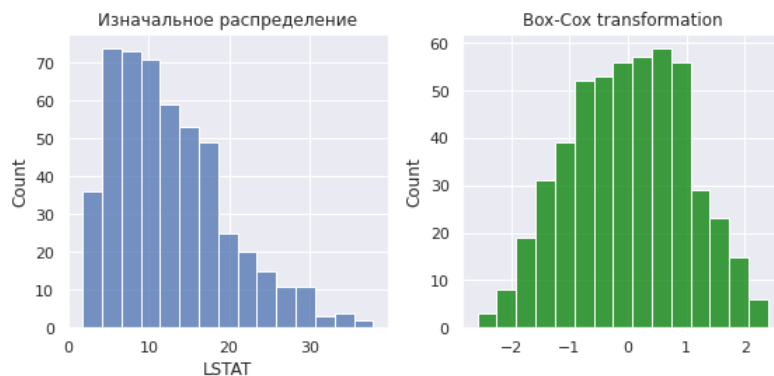
```
1 (506, 1)
```

```
1 # сравним изначальное распределение и распределение после преобразования Бокса-Кокса
2 fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
3
4 sns.histplot(x = boston.LSTAT, bins = 15, ax = ax[0])
5 ax[0].set_title('Изначальное распределение')
6
7 # так как на выходе метод .transform() выдает двумерный массив,
8 # его необходимо преобразовать в одномерный
9 sns.histplot(x = bc_pt.flatten(),
10             bins = 15, color = 'green',
```

```

11         ax = ax[1])
12     ax[1].set_title('Box-Cox transformation')
13
14     plt.tight_layout()
15     plt.show()

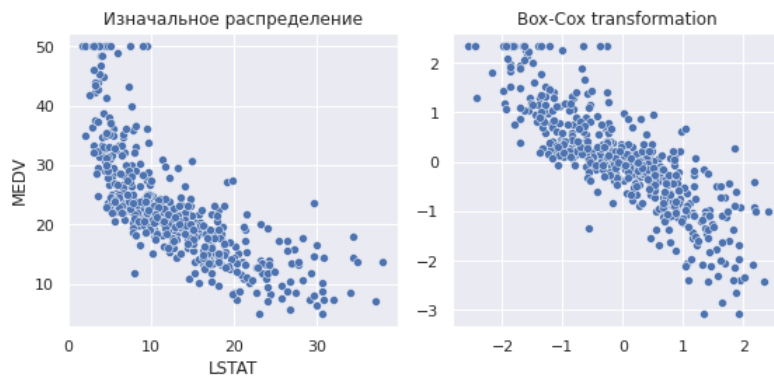
```



```

1     # оценим изменение взаимосвязи после преобразования Бокса-Кокса
2     fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
3
4     sns.scatterplot(x = boston.LSTAT, y = boston.MEDV, ax = ax[0])
5     ax[0].set_title('Изначальное распределение')
6
7     # можно использовать функцию power_transform(),
8     # она действует аналогично классу, но без estimator
9     sns.scatterplot(x = power_transform(boston[['LSTAT']], method = 'box-cox').flatten(),
10                    y = power_transform(boston[['MEDV']], method = 'box-cox').flatten(),
11                    ax = ax[1])
12     ax[1].set_title('Box-Cox transformation')
13
14     plt.tight_layout()
15
16     plt.show()

```



```

1     # посмотрим на достигнутый коэффициент корреляции
2     pd.DataFrame(power_transform(boston[['LSTAT', 'MEDV']], method = 'box-cox'),
3                  columns = [['LSTAT', 'MEDV']]).corr()

```

	LSTAT	MEDV
LSTAT	1.000000	-0.830424
MEDV	-0.830424	1.000000

```

1     # сравним корреляцию признаков с целевой переменной
2     # после преобразования Бокса-Кокса
3     MEDV_bc = power_transform(boston[['MEDV']], method = 'box-cox').flatten()

```

```

4
5 for col in boston[['CRIM', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'LSTAT']]
6     col_bc = power_transform(boston[[col]], method = 'box-cox').flatten()
7     print(col + '\t' + str(np.round(np.corrcoef(col_bc, MEDV_bc)[0][1], 3)))

```

```

1 CRIM    -0.528
2 NOX     -0.5
3 RM      0.64
4 AGE     -0.452
5 DIS      0.392
6 RAD     -0.403
7 TAX     -0.538
8 PTRATIO -0.522
9 LSTAT   -0.83

```

```

1 # возьмем признаки RM, PTRATIO, LSTAT и целевую переменную MEDV
2 # и применим преобразование
3 pt = PowerTransformer(method = 'box-cox')
4 boston_bc = pt.fit_transform(boston[['RM', 'PTRATIO', 'LSTAT', 'MEDV']])
5 boston_bc = pd.DataFrame(boston_bc, columns = ['RM', 'PTRATIO', 'LSTAT', 'MEDV'])
6
7 # построим линейную регрессию
8 # в данном случае показатель чуть хуже, чем при лестнице Тьюки
9 model = LinearRegression()
10 model.fit(boston_bc[['RM', 'PTRATIO', 'LSTAT']], boston_bc.MEDV)
11 model.score(boston_bc[['RM', 'PTRATIO', 'LSTAT']], boston_bc.MEDV)

```

```
1 0.7331845210304999
```

```

1 # посмотрим на лямбды
2 pt.lambdas_

```

```
1 array([0.44895975, 4.35021552, 0.22776736, 0.2166209 ])
```

```

1 # выполним обратное преобразование
2 pd.DataFrame(pt.inverse_transform(boston_bc),
3             columns = ['RM', 'PTRATIO', 'LSTAT', 'MEDV']).head()

```

	RM	PTRATIO	LSTAT	MEDV
0	6.575	15.3	4.98	24.0
1	6.421	17.8	9.14	21.6
2	7.185	17.8	4.03	34.7
3	6.998	18.7	2.94	33.4
4	7.147	18.7	5.33	36.2

Преобразование Йео-Джонсона

Преобразование Йео-Джонсона (Yeo-Johnson transformation) позволяет работать с нулевыми и отрицательными значениями.

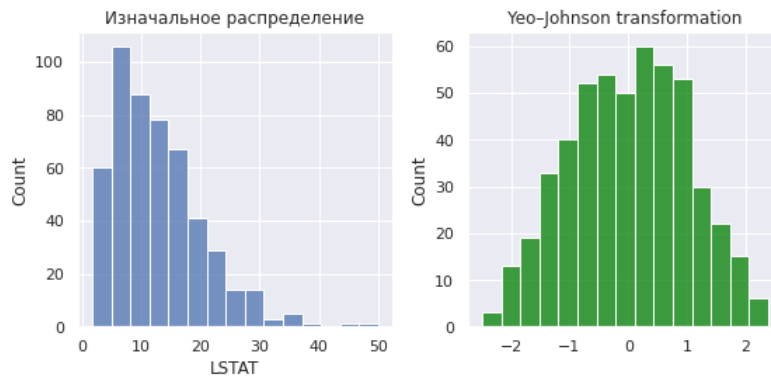
$$y = \begin{cases} ((y+1)^\lambda - 1)/\lambda & \text{if } \lambda \neq 0, y \geq 0 \\ \log(y+1) & \text{if } \lambda = 0, y \geq 0 \\ -[(-y+1)^{(2-\lambda)} - 1]/(2-\lambda) & \text{if } \lambda \neq 2, y < 0 \\ -\log(-y+1) & \text{if } \lambda = 2, y < 0 \end{cases}$$

Попробуем преобразование Йео-Джонсона.

```

1 fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
2
3 sns.histplot(x = boston_outlier.LSTAT, bins = 15, ax = ax[0])
4 ax[0].set_title('Изначальное распределение')
5
6 sns.histplot(x = power_transform(boston[['LSTAT']], method = 'yeo-johnson').flatten(),
7             bins = 15, color = 'green',
8             ax = ax[1])
9 ax[1].set_title('Yeo-Johnson transformation')
10
11 plt.tight_layout()
12 plt.show()

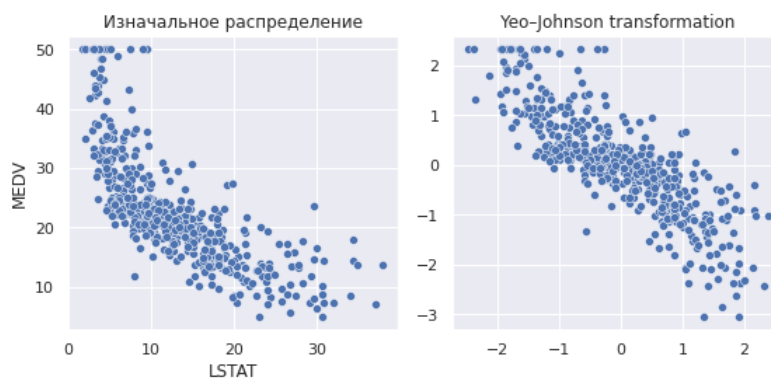
```



```

1 # посмотрим, как изменится линейность взаимосвязи
2 fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
3
4 sns.scatterplot(x = boston.LSTAT, y = boston.MEDV, ax = ax[0])
5 ax[0].set_title('Изначальное распределение')
6
7 sns.scatterplot(x = power_transform(boston[['LSTAT']], method = 'yeo-johnson').flatten(),
8               y = power_transform(boston[['MEDV']], method = 'yeo-johnson').flatten(),
9               ax = ax[1])
10 ax[1].set_title('Yeo-Johnson transformation')
11
12 plt.tight_layout()
13
14 plt.show()

```



```

1 # возьмем те же признаки и целевую переменную, преобразуем их
2 # преобразование Йео-Джонсона является методом по умолчанию
3 pt = PowerTransformer()
4 boston_yj = pt.fit_transform(boston[['RM', 'PTRATIO', 'LSTAT', 'MEDV']])
5 boston_yj = pd.DataFrame(boston_yj, columns = ['RM', 'PTRATIO', 'LSTAT', 'MEDV'])
6
7 # построим модель
8 model = LinearRegression()
9 model.fit(boston_yj.iloc[:, :3], boston_yj.iloc[:, -1])

```

```
10 | model.score(boston_yj.iloc[:, :3], boston_yj.iloc[:, -1])
```

```
1 | 0.7333775808517045
```

QuantileTransformer

При **квантильном преобразовании** (quantile transformation) исходным данным присваивается квантильный ранг в целевом (равномерном или нормальном) распределении. Этот ранг и есть новая преобразованная оценка.

Особенность такого преобразования заключается в том, что новое распределение никак не связано с исходным. Рассмотрим пример преобразования данных с выбросами.

```
1 | from sklearn.preprocessing import QuantileTransformer
2 |
3 | # приведем переменные с выбросами (!) к нормальному распределению
4 | # с помощью квантиль-функции
5 | qt = QuantileTransformer(n_quantiles = len(boston_outlier),
6 |                          output_distribution = 'normal',
7 |                          random_state = 42)
8 |
9 | # для каждого из столбцов вычислим квантили нормального распределения,
10 | # соответствующие заданному выше количеству квантилей (n_quantiles)
11 | # и преобразуем (map) данные к нормальному распределению
12 | boston_qt = pd.DataFrame(qt.fit_transform(boston_outlier),
13 |                          columns = boston_outlier.columns)
14 |
15 | # посмотрим на значения, на основе которых будут рассчитаны квантили
16 | qt.quantiles_[-5:]
```

```
1 | array([[34.77, 50. ],
2 |        [36.98, 50. ],
3 |        [37.97, 50. ],
4 |        [45.  , 70. ],
5 |        [50.  , 72. ]])
```

```
1 | # посмотрим на соответствующие им квантили нормального распределения
2 | qt.references_[-5:]
```

```
1 | array([0.99211045, 0.99408284, 0.99605523, 0.99802761, 1.          ])
```

```
1 | # рассчитаем предпоследнее значение с помощью библиотеки scipy
2 | from scipy.stats import norm
3 |
4 | norm.ppf(0.99802761, loc = 0, scale = 1)
```

```
1 | 2.8825440308212347
```

```
1 | # сравним с преобразованными значениями
2 | boston_qt.LSTAT.sort_values()[-5:]
```

```
1 | 373    2.413985
2 | 414    2.517047
3 | 374    2.656761
4 | 506    2.882545
5 | 507    5.199338
6 | Name: LSTAT, dtype: float64
```

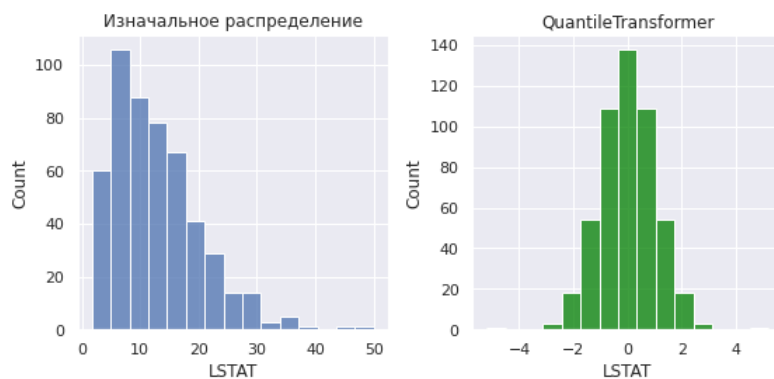
```
1 | # выведем результат
2 | fig, ax = plt.subplots(nrows = 1, ncols = 2, figsize = (8,4))
3 |
4 | sns.histplot(x = boston_outlier.LSTAT, bins = 15, ax = ax[0])
5 | ax[0].set_title('Изначальное распределение')
6 |
```



```

7 | sns.histplot(x = boston_qt.LSTAT,
8 |             bins = 15, color = 'green',
9 |             ax = ax[1])
10 | ax[1].set_title('QuantileTransformer')
11 |
12 | plt.tight_layout()
13 | plt.show()

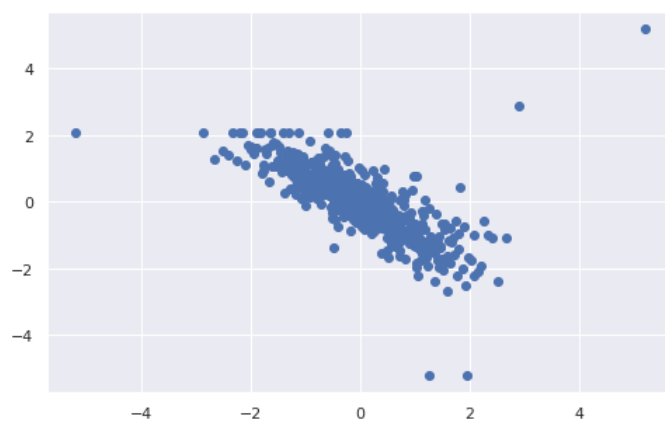
```



```

1 | # посмотрим, выправилась ли взаимосвязь
2 | plt.scatter(boston_qt.LSTAT, boston_qt.MEDV);

```



Исходя из преобразованных значений (посмотрите на значение с индексом 507) и точечной диаграммы мы видим, что эффект выбросов сохранился.

```

1 | max(boston_qt.LSTAT), max(boston_qt.MEDV)

```

```

1 | (5.19933758270342, 5.19933758270342)

```

```

1 | # сравним исходную корреляцию
2 | boston_outlier[['LSTAT', 'MEDV']].corr().iloc[0,1]

```

```

1 | -0.5772033139947359

```

```

1 | # с корреляцией после преобразования
2 | boston_qt.corr().iloc[0,1]

```

```

1 | -0.7037287662365327

```

```

1 | # посмотрим на корреляцию после преобразования Йео-Джонсона
2 | boston_yj = pd.DataFrame(power_transform(boston_outlier, method = 'yeo-johnson'),
3 |                          columns = boston_outlier.columns)
4 |
5 | boston_yj.corr().iloc[0,1]

```

```

1 | -0.7661930913306837

```

Подведем итог

На сегодняшнем занятии мы рассмотрели линейные и нелинейные способы преобразования количественных данных.

Ответы на вопросы

Вопрос. Не очень понял применение формулы евклидова расстояния на занятии по кластерному анализу.

Ответ. На занятии по кластеризации мы измеряли с помощью евклидова расстояния дистанцию между двумя векторами, здесь же мы использовали формулу для расчета длины одного и того же вектора.
