

# Regression Intro: Transforming Target

- Applying transformations to make target variable more normally distributed for regression
- Applying inverse transformations to be able to use these in a regression context

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

## Loading in Boston Data

```
In [2]: #loading boston dataset from sklearn
from sklearn.datasets import load_boston
boston = load_boston()
boston_data = pd.DataFrame(boston.data, columns=boston.feature_names) # Creates dataframe with the feature names as columns
boston_data['MEDV'] = boston.target # Appends the MEDV or target onto the dataframe
boston_description = boston['DESCR'] # Description of the columns
```

## Determining Normality

Making our target variable normally distributed often will lead to better results

If our target is not normally distributed, we can apply a transformation to it and then fit our regression to predict the transformed values.

How can we tell if our target is normally distributed? There are two ways:

- Visually
- Using a statistical test

```
In [3]: boston_data.MEDV.hist()
```



Does not look normal due to that right tail. Let's try to verify statistically.

```
In [6]: from scipy.stats.stats import normaltest # D'Apostolno K-S Test
```

Without getting into Bayesian vs. frequentist debates, for the purposes of this lesson, the following will suffice:

- This is a statistical test that tests whether a distribution is normally distributed or not. It isn't perfect, but suffice it to say:
  - These test outputs a "p-value". The higher this p-value is the closer the distribution is to normal.
  - Frequentist statisticians would say that you accept that the distribution is normal (more specifically: fail to reject the null hypothesis that it is normal) if  $p > 0.05$ .

```
In [7]: normaltest(boston_data.MEDV.values)
Out[7]: NormaltestResult(statistic=90.9746373709967, pvalue=1.7583188871696995e-20)
```

p-value extremely low. Our y variable we've been dealing with this whole time was not normally distributed!

Linear Regression assumes a normally distributed residuals which can be aided by transforming y variable. Let's try some common transformations to try and get y to be normally distributed:

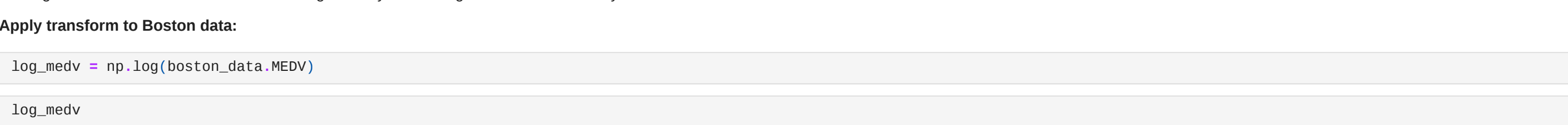
- Log
- Square root
- Box cox

The log transform can transform data that is significantly skewed right to be more normally distributed:

Apply transform to Boston data:

```
In [12]: log_medv = np.log(boston_data.MEDV)
```

```
In [13]: log_medv
```



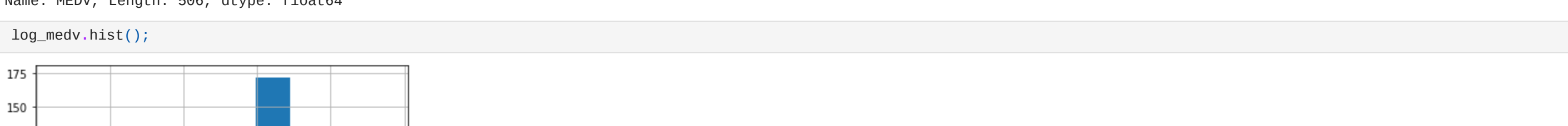
```
In [15]: normaltest(log_medv)
Out[15]: NormaltestResult(statistic=17.21981096640697, pvalue=0.68918245472768345307)
```

Conclusion: closer, but still not normal.

Slightly skewed right.

Apply the square root transformation to the Boston data target and test whether the result is normally distributed.

```
In [20]: # Instructor Solution
sqrt_medv = np.sqrt(boston_data.MEDV)
plt.hist(sqrt_medv)
```



```
In [21]: normaltest(sqrt_medv)
Out[21]: NormaltestResult(statistic=20.48799082683967, pvalue=3.558645781429252e-05)
```

## Box Cox

The box cox transformation is a parameterized transformation that tries to get distributions "as close to a normal distribution as possible".

It is defined as:

$$\text{boxcox}(y) = \frac{y^{\lambda} - 1}{\lambda}$$

You can think of  $\lambda$  as a generalization of the square root function: the square root function uses the exponent of 0.5, but box cox lets its exponent vary so it can find the best one.

```
In [22]: from scipy.stats import boxcox
```

```
In [25]: bc_result = boxcox(boston_data.MEDV)
boxcox_medv = bc_result[0]
lam = bc_result[1]
```

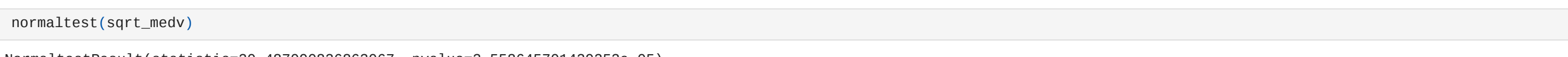
```
In [24]: lam
```

```
Out[24]: 0.2166299012915364
```

```
In [30]: boston_data['MEDV'].hist()
```



```
In [31]: plt.hist(boxcox_medv)
```



```
In [32]: normaltest(boxcox_medv)
Out[32]: NormaltestResult(statistic=4.513528775533945, pvalue=0.1046886692917602)
```

Significantly more normally distributed (according to p value) than the other two distributions - above 0.05, even!

Now that we have a normally distributed y-variable, let's try a regression!

## Testing regression:

```
In [33]: from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import (StandardScaler,
PolynomialFeatures)
```

```
In [34]: lr = LinearRegression()
```

Reload clean version of boston\_data:

```
In [37]: from sklearn.datasets import load_boston
boston = load_boston()
boston_data = pd.DataFrame(boston.data, columns=boston.feature_names) # Creates dataframe with the feature names as columns
boston_data['MEDV'] = boston.target # Appends the MEDV or target onto the dataframe
boston_description = boston['DESCR'] # Description of the columns
```

Same steps as before.

```
In [38]: Create X and y
y_col = "MEDV"
X = boston_data.drop(y_col, axis=1)
y = boston_data[y_col]
```

```
In [48]: X
```



```
Out[48]: 0 0.0362 18.0 2.31 0.0 0.538 6.575 65.2 4.090 1.0 296.0 15.3 396.9 9.67
1 0.0273 0.0 7.07 0.0 0.469 6.421 7.89 4.9671 2.0 242.0 17.8 396.90 9.14
2 0.0279 0.0 7.07 0.0 0.469 7.185 6.11 4.9671 2.0 242.0 17.8 392.83 4.03
3 0.0237 0.0 2.18 0.0 0.458 6.998 45.8 6.0622 3.0 222.0 18.7 394.63 2.94
4 0.0905 0.0 2.18 0.0 0.458 7.347 54.2 6.0622 3.0 222.0 18.7 396.90 5.33
```

```
... ..
501 0.0623 0.0 11.93 0.0 0.573 6.993 69.1 2.4786 1.0 273.0 21.0 391.99 9.67
502 0.0457 0.0 11.93 0.0 0.573 6.130 76.7 2.2875 1.0 273.0 21.0 396.90 9.08
503 0.0978 0.0 11.93 0.0 0.573 6.976 91.0 2.1675 1.0 273.0 21.0 396.90 5.64
504 0.1095 0.0 11.93 0.0 0.573 6.794 89.3 2.3889 1.0 273.0 21.0 393.45 6.48
505 0.0474 0.0 11.93 0.0 0.573 6.030 80.8 2.5050 1.0 273.0 21.0 396.90 7.88
```

506 rows x 13 columns

```
In [49]: y
```



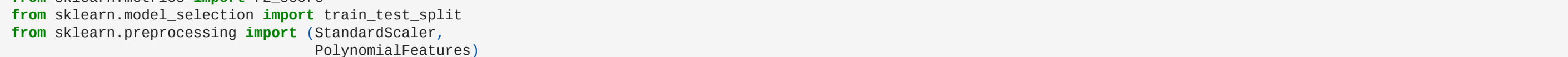
```
Out[49]: 0 24.0
1 21.6
2 34.7
3 33.4
4 38.2
501 22.4
502 29.6
503 23.9
504 22.0
505 11.9
```

lam = MEDV, Length: 506, dtype: float64

## Create Polynomial Features

```
In [46]: pf = PolynomialFeatures(degree=2, include_bias=False)
X_pf = pf.fit_transform(X)
```

```
In [47]: X_pf
```



```
Out[47]: array([[6.32096960e-03, 1.88090908e+01, 2.31090908e+00, ...,
1.57529610e+05, 1.97656200e+03, 2.48004000e+04, ...,
2.72309696e-02, 0.86909090e+00, 7.07090909e+00, ...,
1.57529610e+05, 3.62766600e+03, 8.35393000e+01],
1.54315409e+05, 1.58310409e+03, 1.62409090e+01],
... ..
1.57529610e-02, 0.86909090e+00, 1.19309090e+01, ...,
1.57529610e+05, 2.23851600e+03, 3.18890900e+01],
1.69580900e-01, 0.86909090e+00, 1.19309090e+01, ...,
1.54802020e+02, 2.54855600e+03, 4.18864000e+01],
4.74180900e-02, 0.86909090e+00, 1.19309090e+01, ...,
1.57529610e+05, 3.12757200e+03, 6.28844000e+01]])
```

```
In [51]: X.shape
Out[51]: (506, 13)
```

```
In [52]: X_pf.shape
Out[52]: (506, 184)
```

## Train test split

```
In [45]: X_train, X_test, y_train, y_test = train_test_split(X_pf, y, test_size=0.3,
random_state=72018)
```

```
In [54]: X_train.shape # X_train is 70% of OUR X_pf
Out[54]: (354, 184)
```

```
In [59]: X_test.shape
Out[55]: (152, 184)
```

## Fit StandardScaler on X\_train as before

We're going to use the StandardScaler, which is another transform object. We set that equal to s. Again, we used the fit transform method and we pass in x\_train, and that will calculate, given our x\_train, those 354 rows and 104 columns. For each one of those columns, what is the mean value, what is the standard deviation, and subtract the mean and divide by that standard deviation. Then, we set the standardized version of x\_train equal to x\_train\_s. That should have the same shape as before. Just now, the variables rather than being on different scales, they're now all on the same scale.

```
In [56]: s = StandardScaler()
X_train_s = s.fit_transform(X_train)
```

```
In [57]: X_train_s.shape
Out[57]: (354, 184)
```

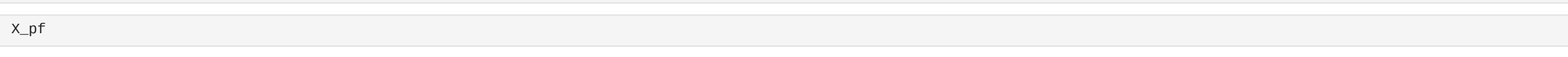
## What transformation do we need to apply next?

Apply the appropriate transformation.

So we're going to use the boxcox that we defined earlier. We pass in y\_train, which is the subset of just 354 rows, it's just one column because the y is just one column. We only want to test on our y\_train, we're only learning our parameters from our x\_train and y\_train, so we're only going to do our transformations there first. From that, remember when we run boxcox, we get two outputs, both the actual transformed values, which is going to be index 0 of our output here, as well as the Lambda value that we have learned that parameter that we learned given our data in order to come up with the normal distribution. So we run this and we have our new transform version of the outcome variable, as well as lam2, which is just going to be for that training set what the Lambda value, you see it's slightly different than before because we're only working on a subset of the data in order to come up with a normal distribution.

```
In [59]: # Instructor Solution
bc_result2 = boxcox(y_train)
y_train_bc = bc_result2[0]
lam2 = bc_result2[1]
```

```
In [61]: y_train_bc
```



```
Out[61]: array([4.8625668 , 6.0620177 , 4.51310411, 4.81957899, 4.95725428,
5.77277497, 5.41197888, 5.75540824, 3.72665146, 4.94686892,
5.7659624 , 6.0671112 , 3.9225236 , 4.68262566, 4.85108878,
6.82748198, 3.86236699, 5.8643167 , 4.31536775, 6.91872348,
4.93530901, 4.64209542, 5.84297841, 5.42960436, 5.44832381,
3.98464825, 4.88299811, 4.19919987, 3.59736813, 5.2309719 ,
5.30478958, 5.85928444, 4.78820348, 4.9844728 , 6.78622385,
3.9269121 , 3.97763516, 2.98361841, 6.8642437 , 4.40872355,
5.5174939 , 3.79092744, 5.61456675, 3.84824493, 5.35252888,
3.36764987, 4.97580975, 5.14213658, 4.37539217, 5.07041151,
6.31926686, 4.90476787, 5.25799412, 5.69133946, 4.21362253,
5.232357 , 5.49371542, 4.38659851, 3.95393734, 4.18951812,
2.75465748, 4.06367514, 3.01030741, 4.5542641 , 4.8099537 ,
5.74517281, 3.26239318, 3.5862887 , 6.17678314, 4.23677257,
6.1805343 , 4.41206646, 4.841458 , 5.7757157 , 3.93757579,
5.8709227 , 5.2224684 , 4.59501124, 4.43662547, 3.8263134 ,
5.54227886, 4.39650786, 4.48825286, 4.2752375 , 5.0827772 ,
4.78697889, 4.86974534, 5.83759891, 4.04151359, 5.51858119,
3.1923357 , 4.86203543, 6.02962036, 4.7099583 , 6.75824002,
4.97861144, 4.95471759, 4.80828669, 4.58268801, 4.9495399 ,
5.7085544 , 5.06741314, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.81537118, 4.93454864, 4.6227336 , 4.7645648 , 4.8113669 ,
4.39336941, 3.74218874, 4.6145823 , 5.1092456 , 5.03817894,
4.79851144, 4.95471759, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
6.02962036, 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 4.5902127 , 4.2882772 ,
4.63012198 , 3.9518773 , 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 4.5902127 , 4.2882772 ,
4.63012198 , 3.9518773 , 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 4.5902127 , 4.2882772 ,
4.63012198 , 3.9518773 , 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 4.5902127 , 4.2882772 ,
4.63012198 , 3.9518773 , 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 4.5902127 , 4.2882772 ,
4.63012198 , 3.9518773 , 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 4.5902127 , 4.2882772 ,
4.63012198 , 3.9518773 , 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 4.5902127 , 4.2882772 ,
4.63012198 , 3.9518773 , 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
4.9846805 , 4.47712986, 4.87522696, 4.5902127 , 4.8611435 ,
4.7268886 , 3.19821381, 3.87049243, 4.95584843, 5.22849234,
4.2252288 , 4.80821357, 4.82227987, 4.80828669, 4.58268801, 4.9495399 ,
3.47670383 , 4.57216244, 4.34144687, 3.84861819, 4.97792636,
4.10355832 , 4.46902222, 3.58276147, 6.54991774, 4.3295266 ,
5.82948262 , 4.6073081 , 3.84861819, 4.69828951, 4.77608569,
4.8916183 , 4.841458 , 4.74293114, 5.08895195, 5.17780898,
3.9805418 , 5.89317686, 4.87522696, 
```