

# Machine Learning - Übung 1

Estimate the price of real estate using linear regression with different features, Module *ml*@FHNW.

## Excercise 1

### Load data

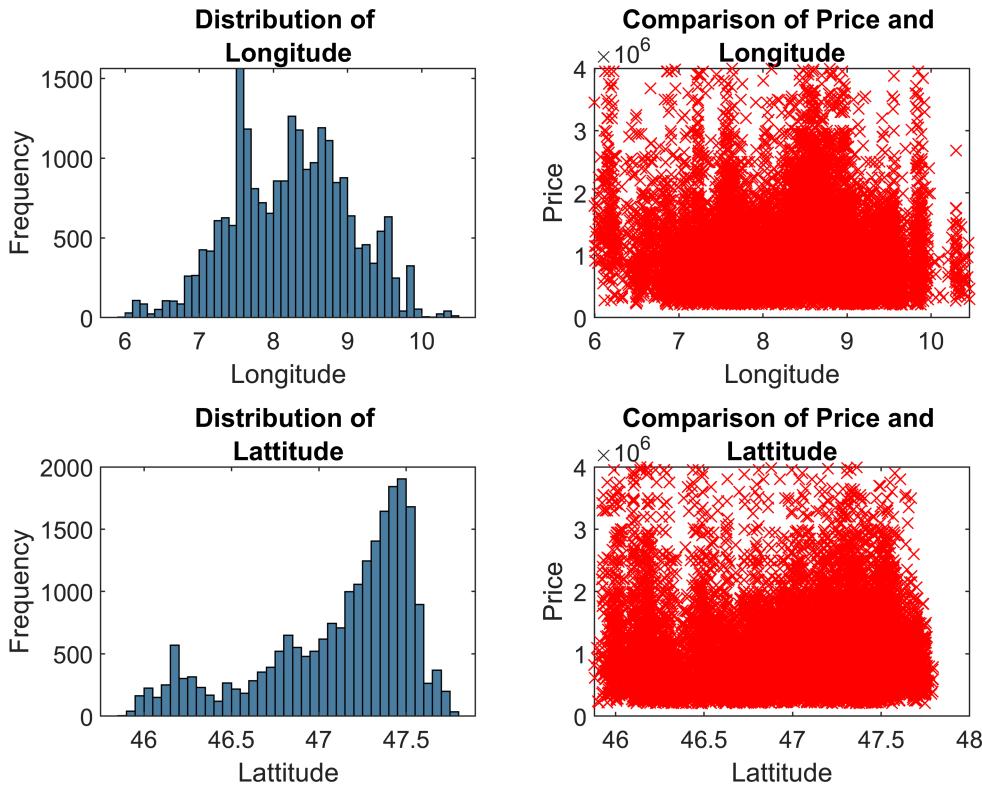
The data used for the upcoming tasks will be read from a given file `house_data.csv`. The description of the data and their origin is not part of this assignment. Below the first five rows will be displayed:

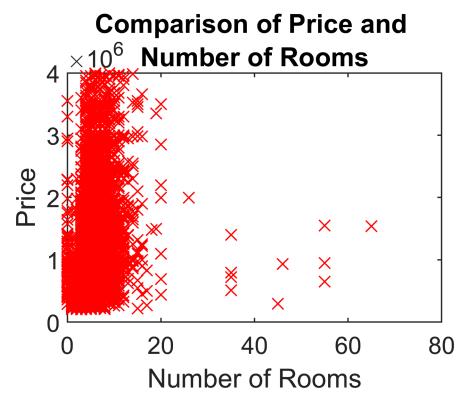
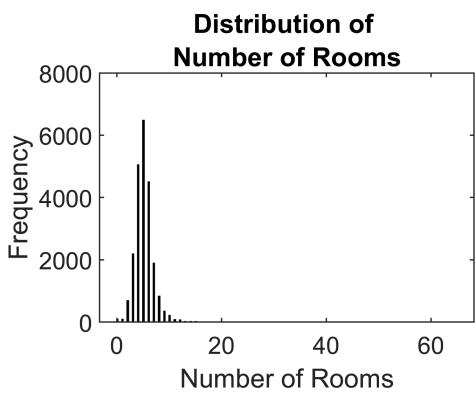
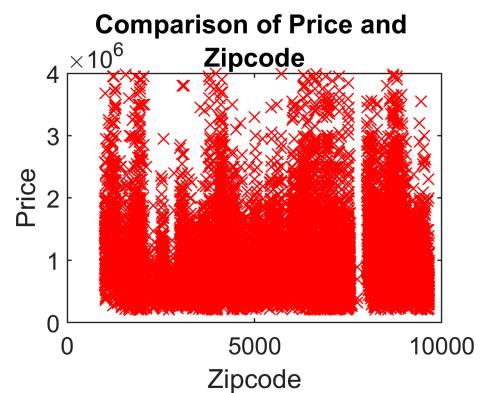
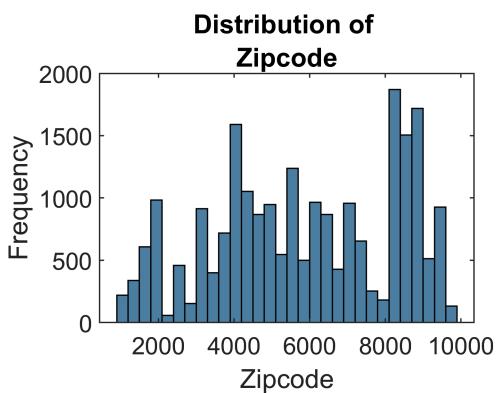
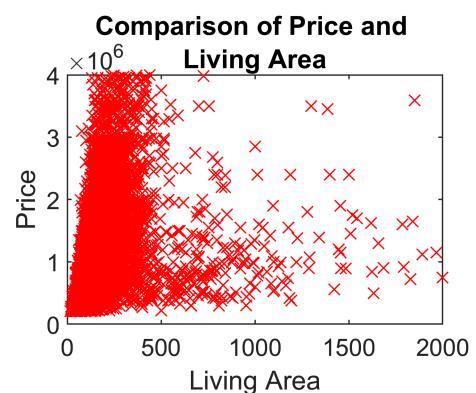
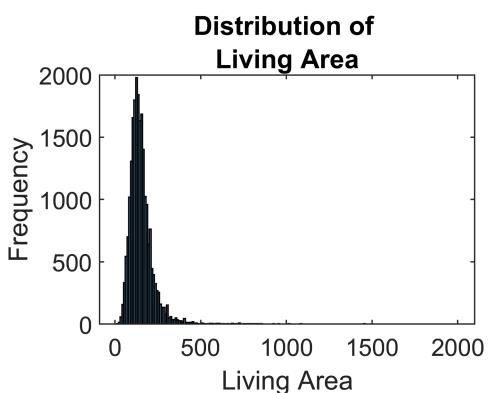
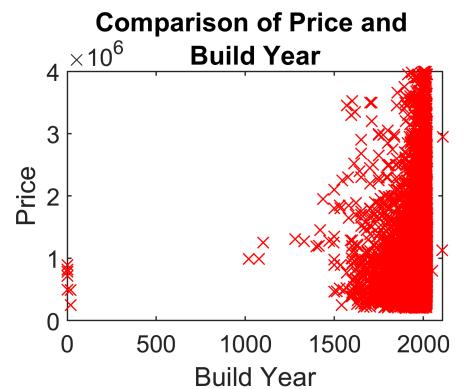
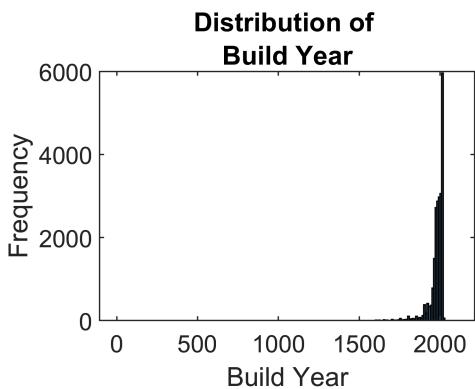
```
ans = 5x19 table
```

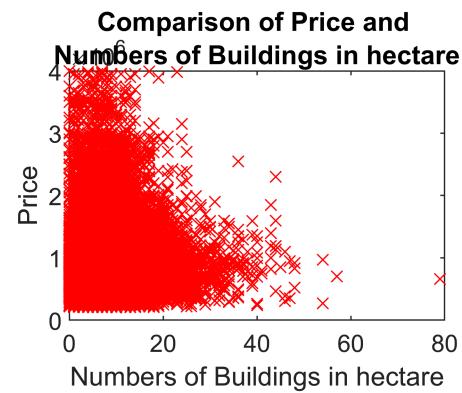
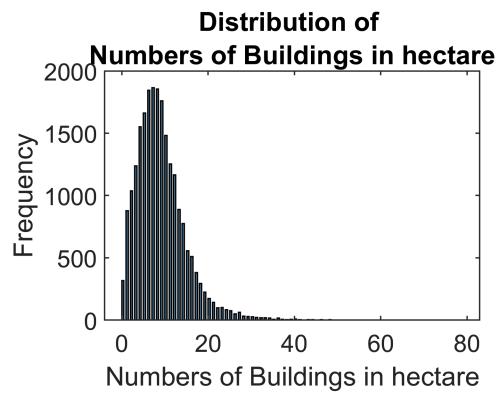
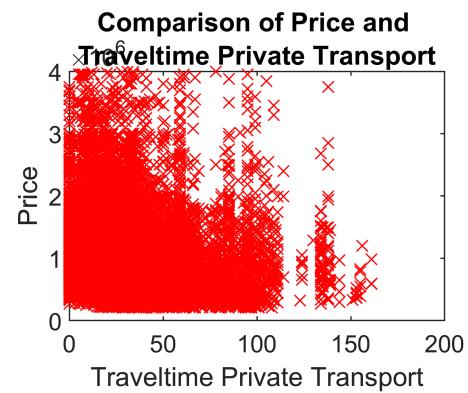
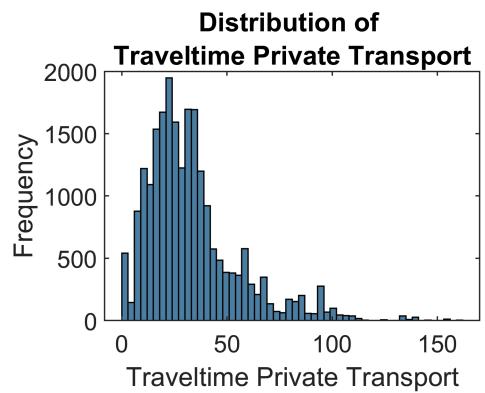
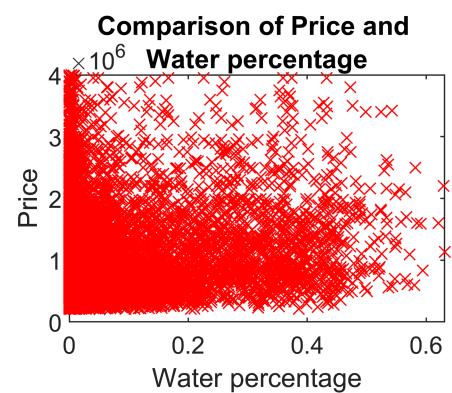
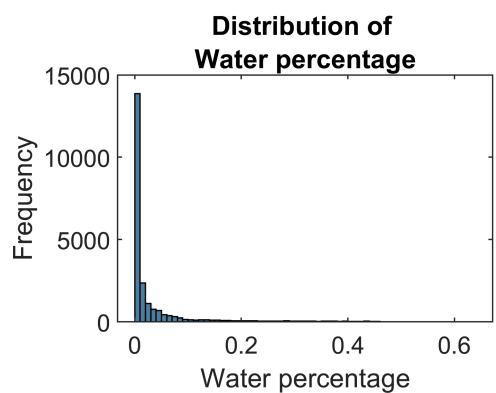
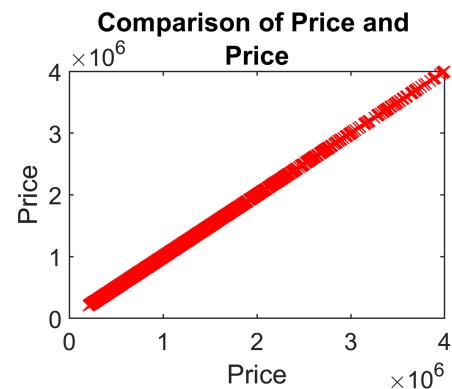
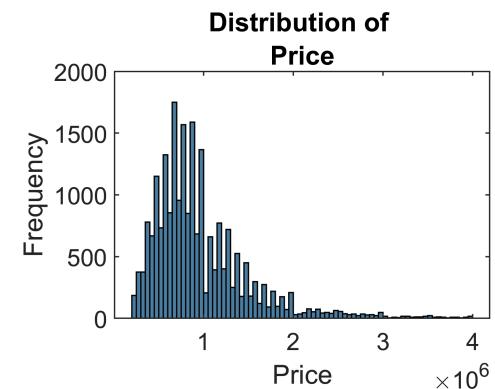
	build_year	lat	living_area	long	municipality_name	zipcode
1	1990	47.0098...	110	8.48378...	'Vitznau'	6354
2	2017	46.8721...	120	9.88018...	'Klosters-Serneus'	7250
3	2010	46.5191...	107	6.52588...	'Préverenges'	1028
4	2018	47.5213...	103	8.53697...	'Bülach'	8180
5	2007	46.5102...	95	9.85242...	'Celerina/Schlarigna'	7505

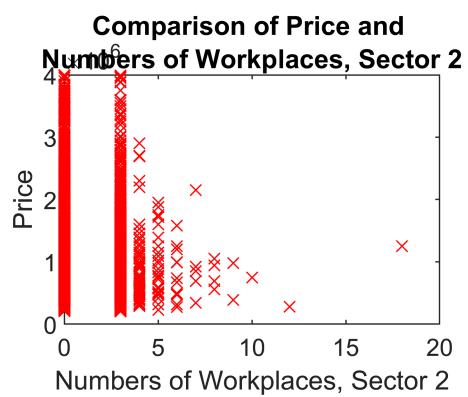
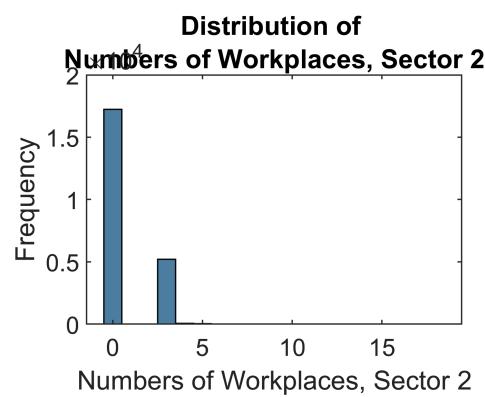
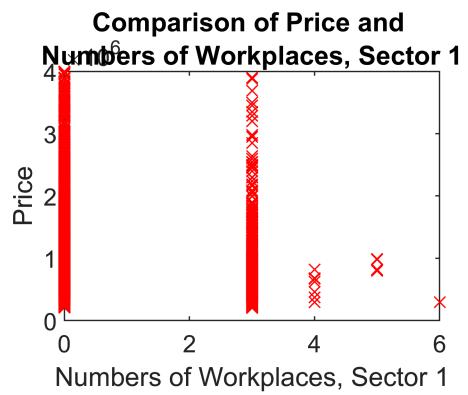
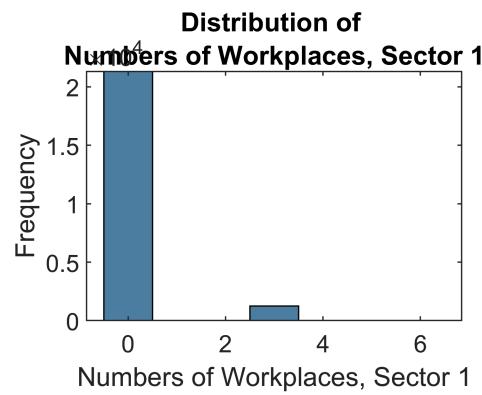
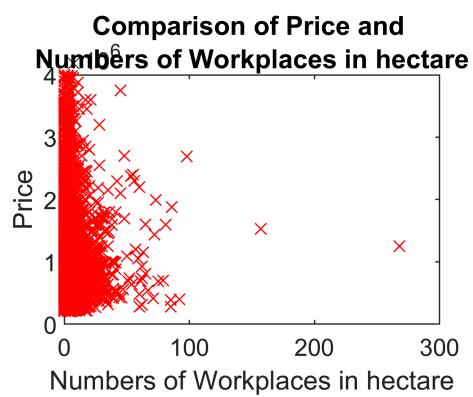
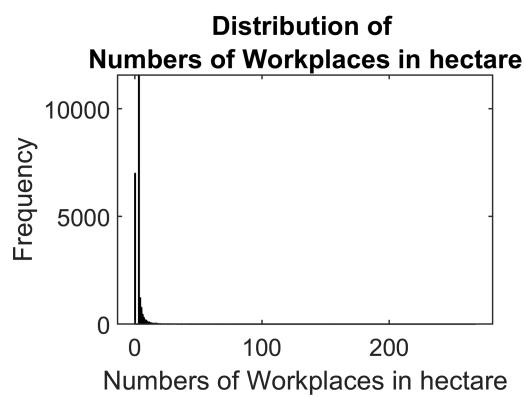
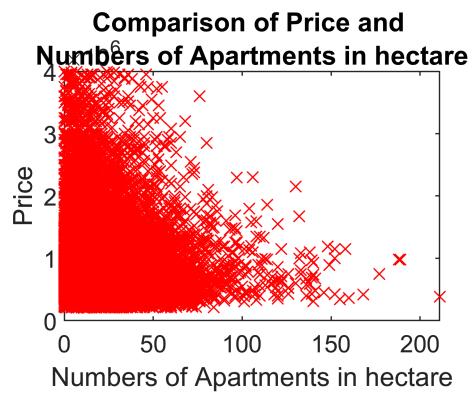
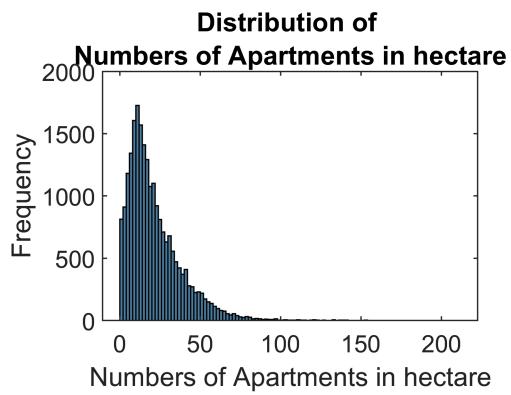
### Visualize data

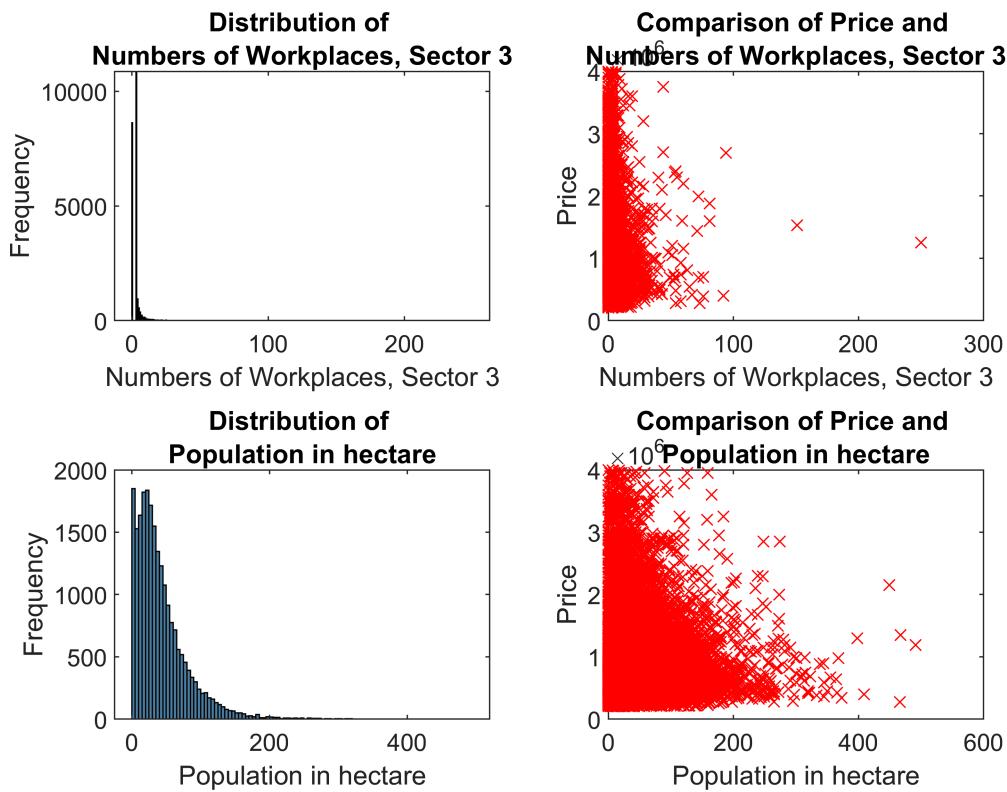
The loaded data is visualized in a next step. For all numerical features a histogram plot of the values is created to get a rough idea of the distribution of the values. In addition, a scatter plot is used to display the correlation between the respective feature and the feature price.



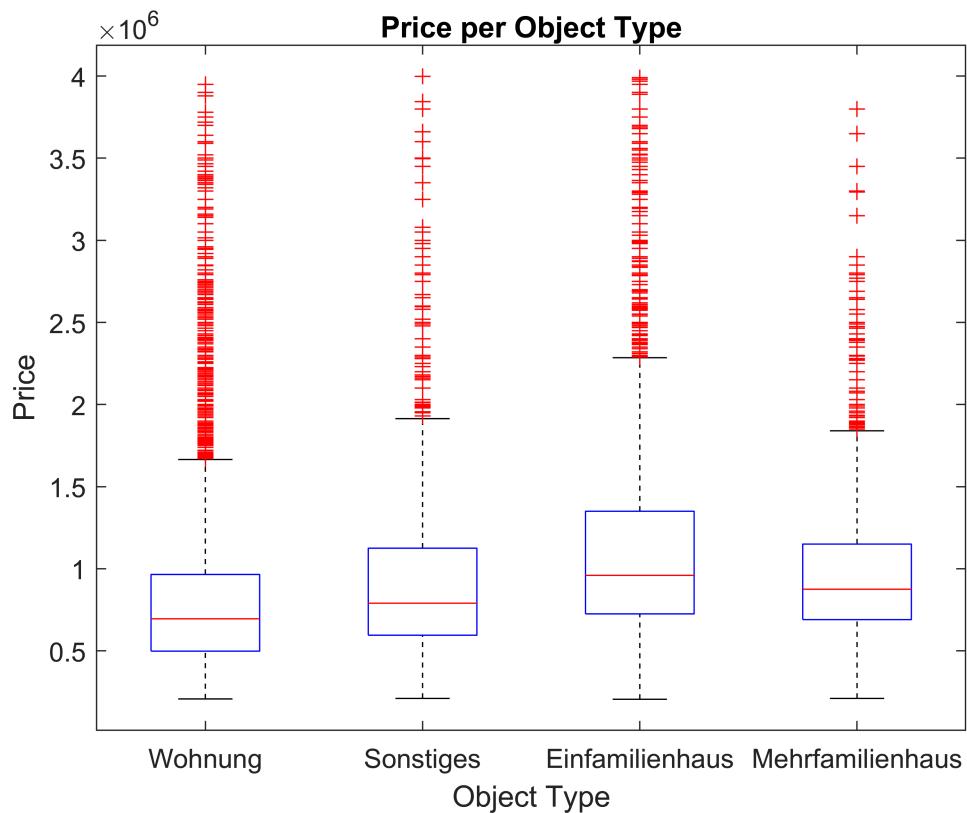








The only categorical feature object\_type is represented by a boxplot. In the boxplot we can see the distribution of the achieved prices per object type.



## Excercise 2

### Normalize feature living area

To standardize the range of the feature, mean normalization is applied to the feature living area. The normalization is done according to [https://en.wikipedia.org/wiki/Feature\\_scaling#Mean\\_normalization](https://en.wikipedia.org/wiki/Feature_scaling#Mean_normalization).

### Predict price using gradient descent

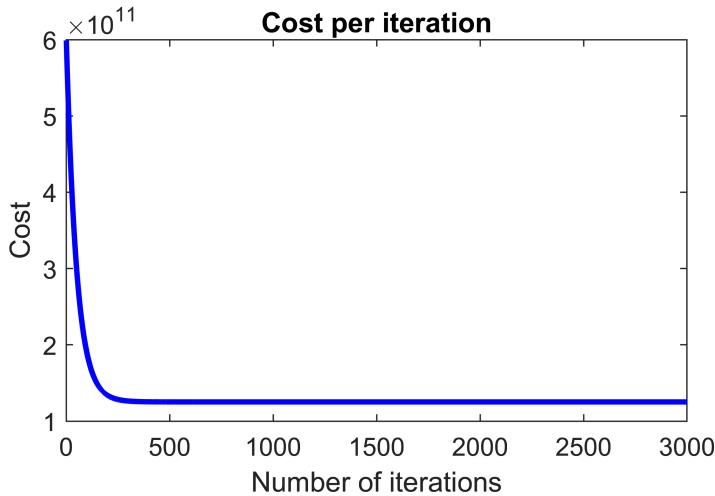
In the next step gradient descent is used, to predict the price using the given and normalized feature living area.

Theta found by gradient descent:

960568.080018  
210590.105394

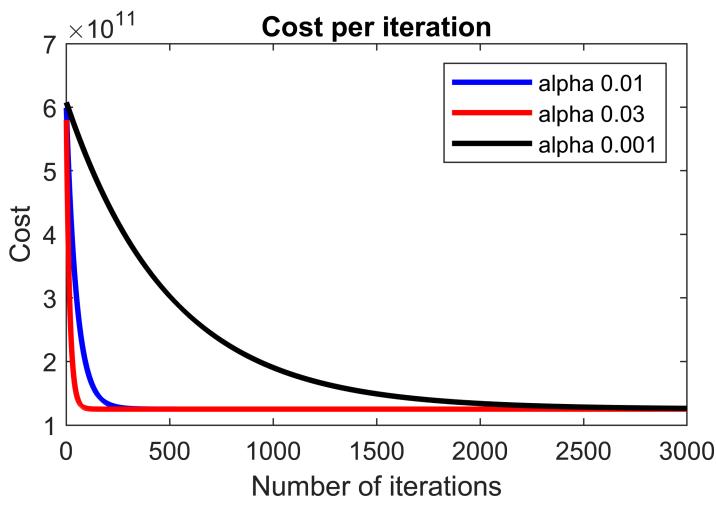
### Visualize costs per iteration

To ensure that the gradient descent procedure has been implemented correctly, the development of costs per iteration is plotted:



### Vary alpha

To show the effect of using different parameters, the value of the learning rate alpha is changed in the following executions. The larger the alpha is chosen, the faster the costs settle down. However, it must be taken care that no underfitting occurs.



## Predict price using normal equation

To verify gradient descent results normal quidation is used, to predict the price using the given and normalized feature living area.

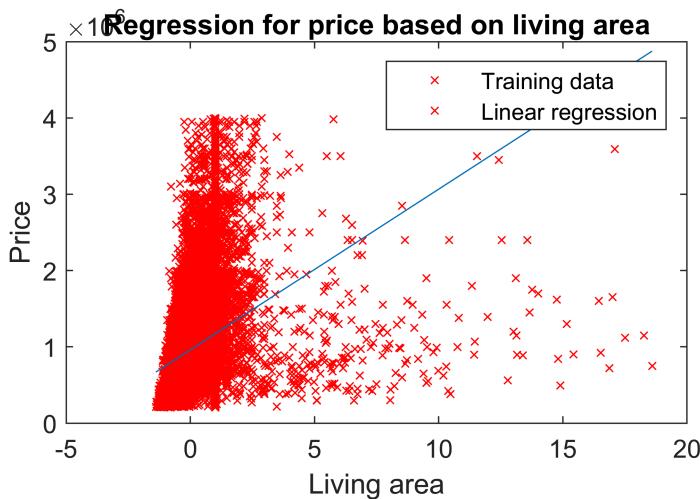
Theta computed from the normal equation:

960568.080018  
210590.105394

While the learning rate and number of iterations must be determined for the gradient descent procedure, this is omitted for the calculation using normal equation. If a large number of features are examined, the calculation using normal equiation can become more expensive. The results of both gradient descent and normal equation are consistent. In both cases we get 960568.080018 and 210590.105394 as optimal elements for theta.

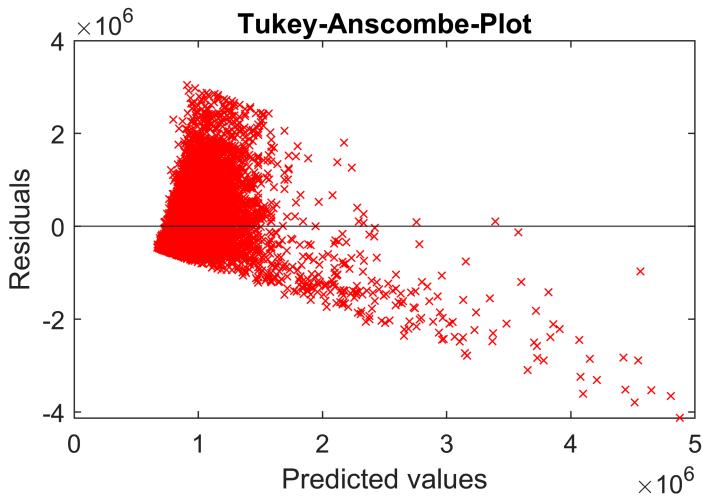
## Visualize found regression

In the following plot, the found regression is added to our former plot, where we displayed data points for price on the y axis and living area on the x axis. Already in this plot it can be seen that the linear regression fits only conditionally to the data. Especially when predicting the price of objects with a large living space, the predictions seem inaccurate.



## Visualize residuals

Residuals are deviations from the expected value. In a tukey-anscombe-plot those residuals can be visualized. The tukey-anscombe plot is further described under <http://stat.ethz.ch/~stahel/courses/regression/regresanal.pdf>. In the plot a trend can be seen that the predicted values are less accurate the bigger the values are. In addition, the forecasts for large values tend to be too low.

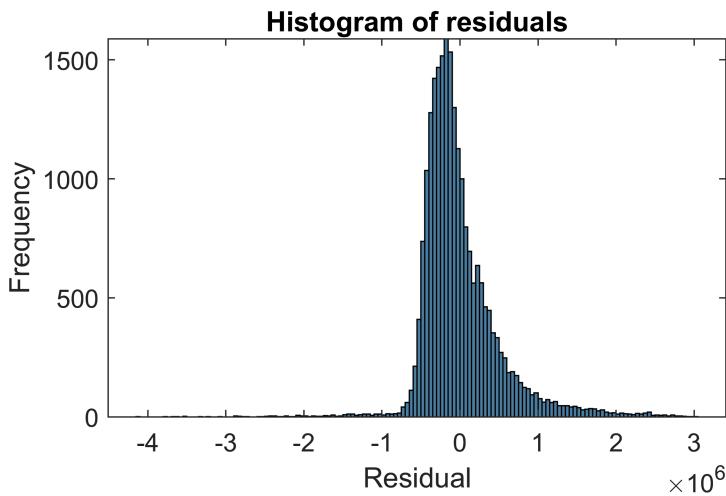


The high deviations of the model to the effective prices already assumed in the plot are underlined by the values obtained for mean value and standard deviation:

Mean: 0.000000

Standard deviation: 500349.786811

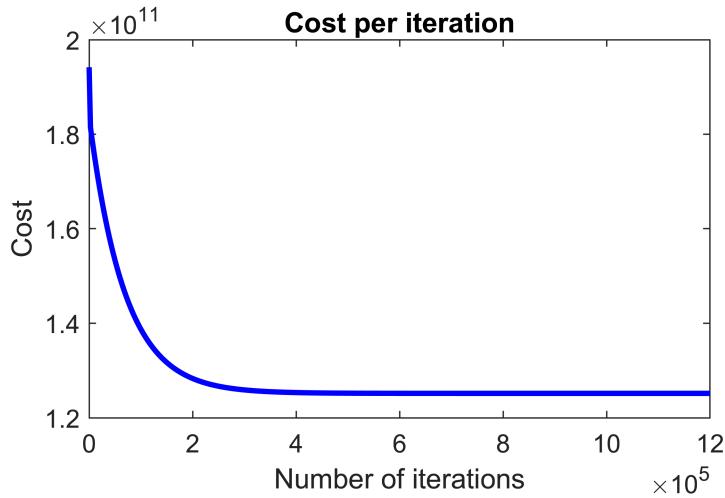
Furthermore the distribution of the residuals can be visualized in a histogram.



The histogram shows the frequencies of the respective residuals. Also in this diagram it can be seen that the differences are mostly in a range from -1'000'000 to 1'000'000. Our assumption of a linear model seems to be inaccurate.

## Predict price using non-normalized feature *living area*

In the following the linear regression is calculated with the non-normalized feature living area. For the gradient descent procedure to work, the number of iterations must be very high and the alpha very low. Due to the high number of iterations and the low learning rate, the complete process takes a relatively long time. The following example shows a run using 1'200'000 iterations and a learning rate of 0.000025:



The result can be checked again with normal equation. Various attempts show that the values obtained are only approximately the same.

Theta found by gradient descent:

632584.058974  
2123.445110

Theta computed from the normal equations:

632683.237456  
2122.990475

## Excercise 3

### Predict price using gradient descent

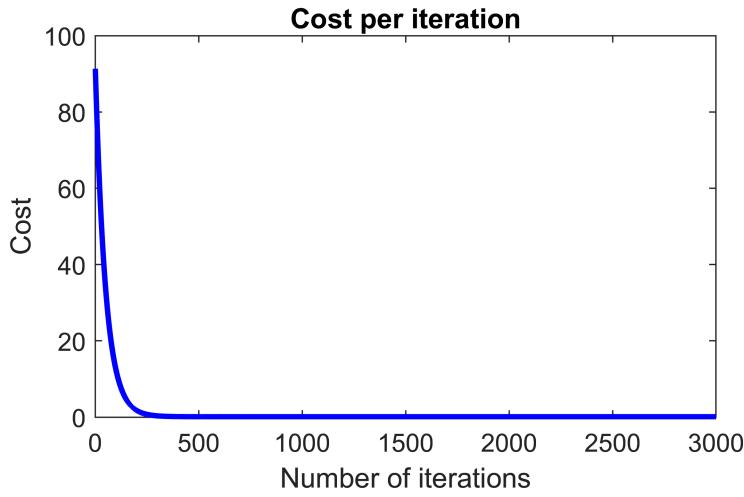
In the next step gradient descent is used, to predict the price using the given and normalized feature living area. Instead of linear, price is assumed to be a logarithmic-function and therefore transformed to  $\log(\text{price})$ .

Theta found by gradient descent:

13.641726  
0.201981

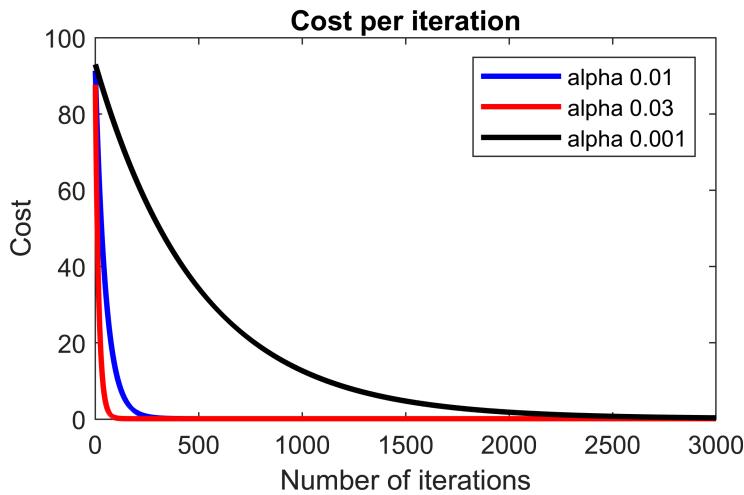
### Visualize costs per iteration

To ensure that the gradient descent procedure has been implemented correctly, the development of costs per iteration is plotted. In comparison to the previous model, it can be seen that the model is more quickly stabilising at a certain value.



## Vary alpha

To show the effect of using a different parameters, the value of alpha ist changed in the following executions:



## Predict price using normal equiation

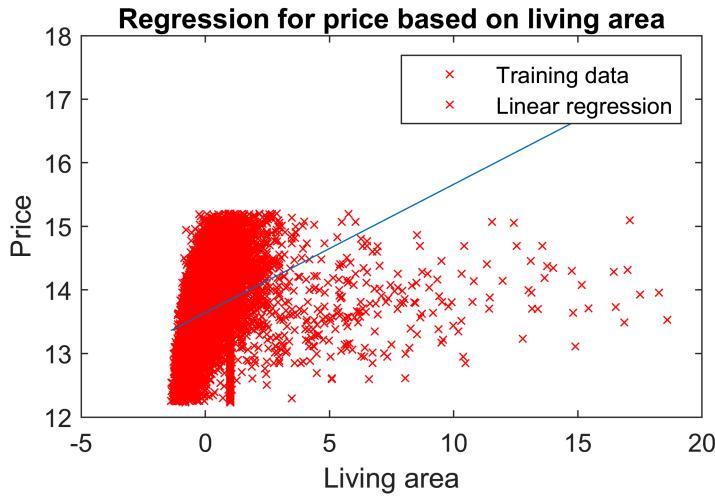
To verify gradient descent results normal quidation is used, to predict the price using the given and normalized feature living area. Again we receive the same results for gradient descent and normal equiation.

Theta computed from the normal equations:

```
13.641726
0.201981
```

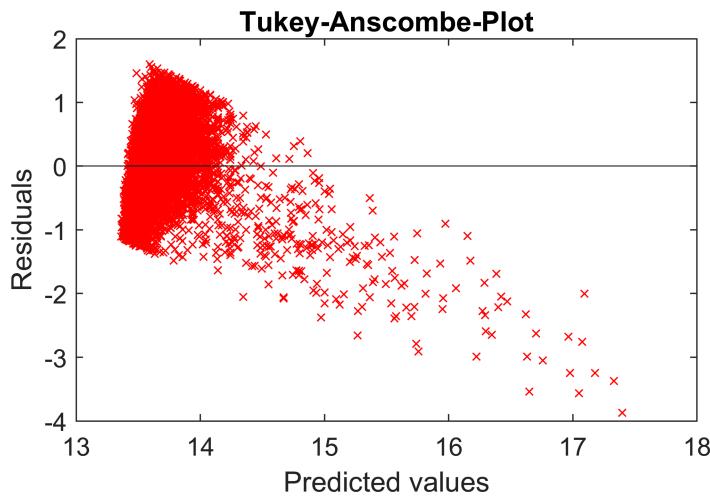
## Visualize found regression

In the following plot, the found regression is added to our former plot, where we displayed data points for  $\log(\text{price})$  and living area. In general, plots show that price developments often behave approximately logarithmically and that the model therefore delivers better results. In the next plot, these findings are still difficult to see, but at the latest the calculation of the standard deviation provides corresponding results.

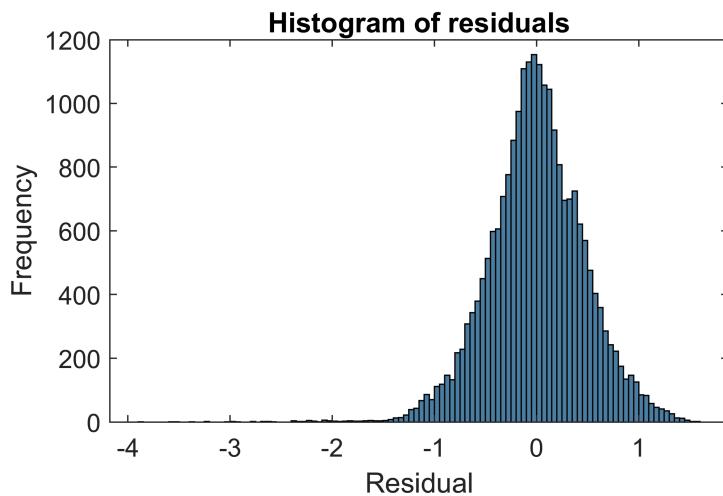


## Visualize residuals

Residuals are deviations from the expected value. In a tukey-anscombe-plot those residuals can be visualized.



Furthermore the distribution of the residuals can be visualized in a histogram. The deviations on the x-axis allow the conclusion that the model can now predict values more accurately.



The above assumptions are confirmed by a significantly lower standard deviation (compared with ~500'349 using  $y = \text{price}$ ).

Mean: 0.000000

Standard deviation: 0.469710

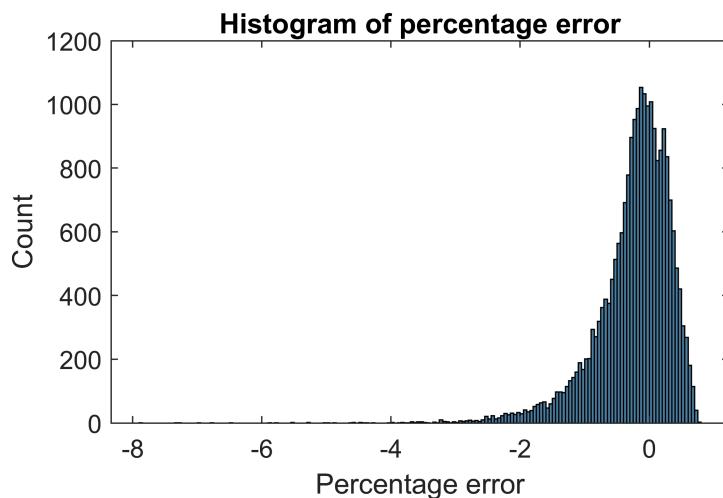
## Excercise 4

### Calculate mean absolute percentage error

In the following section the mean absolute percentage error (MAPE) for  $y = \text{price}$  is calculated and the distribution of de percentage error displayed in a histogram. The MAPE of the first linear model is relatively high. This coincides with the high standard deviation of its residuals.

Computed MAPE from  $y = \text{price}$ :

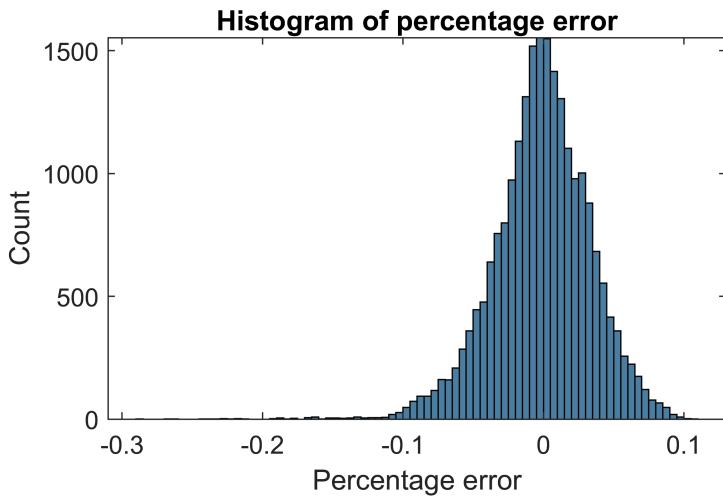
0.438663



The same can be done for  $y = \log(\text{price})$ . The comparison shows that the model makes much more accurate predictions and the percentage error is much lower.

Computed MAPE from  $y = \log(\text{price})$ :

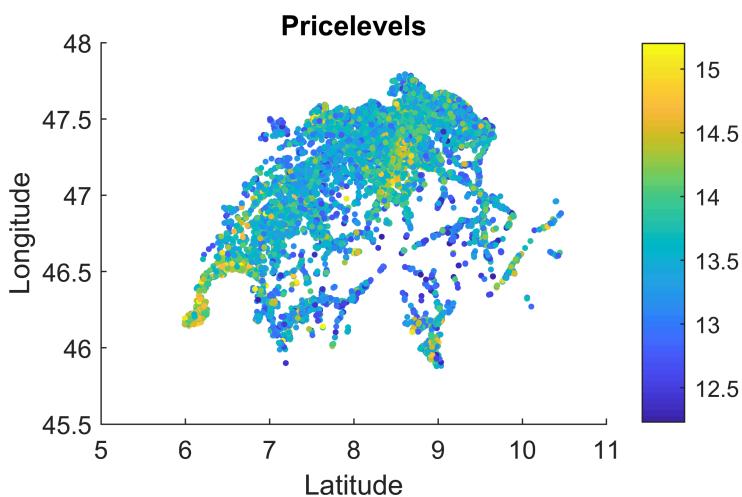
0.026266



## Excercise 5

### Visualize price in relation to longitude/latitude

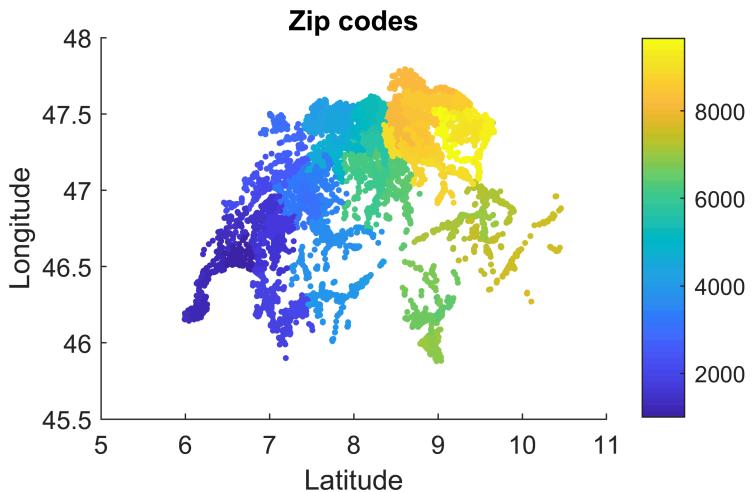
The output variable  $y = \log(\text{price})$  can be visualized on a map using longitude and latitude values. In the plot the yellow data points represent the most expensive objects. These buildings are located mainly in the regions of Geneva and Zurich.



## Excercise 6

### Visualize zip code in relation to longitude/latitude

A similar visualization can be done using  $y = \text{zip}$ . The plot shows the distribution of the zip code ranges in Switzerland.



## Excercise 7

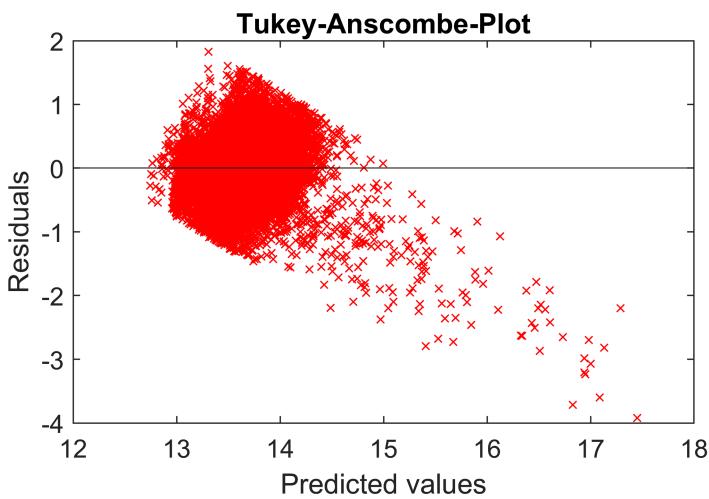
### Predict price using additional one-hot-encoded feature zipcode

For the following prediction we additionally use the first two digits of the feature zipcode to predict an objects price. This should reduce outliers in high-price areas and allow to predict prices more accurately overall.

Theta computed from the normal equations (partly shown):

```
13.418721
0.208330
0.000000
-0.000000
-0.000000
```

In the tukey-anscombe plot only tiny differences are visible:

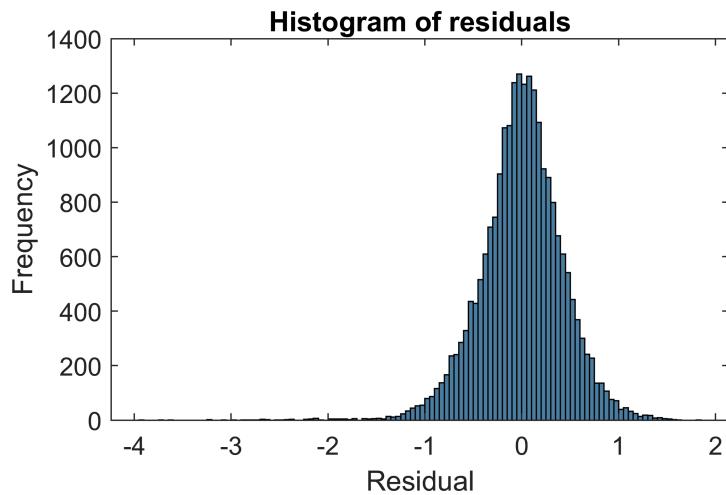


However, the standard deviations show that the predicted values could be improved further (compared to a std. deviation of 0.469710):

Mean: -0.000000

Standard deviation: 0.425935

The improved model is also reflected in the histogram. The deviations from the mean value are again smaller.



## Excercise 8

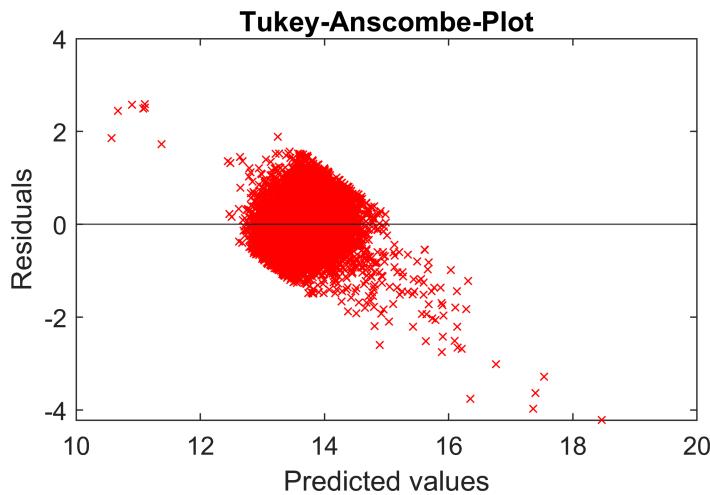
### Predict price using further additional features

Based on the previous findings, we consider other features to determine the price.

Theta computed from the normal equations (partly shown):

```
10.841215  
0.137625  
0.001143  
0.076836  
-0.001264
```

The tukey-anscombe plot shows that the model could be improved significantly:

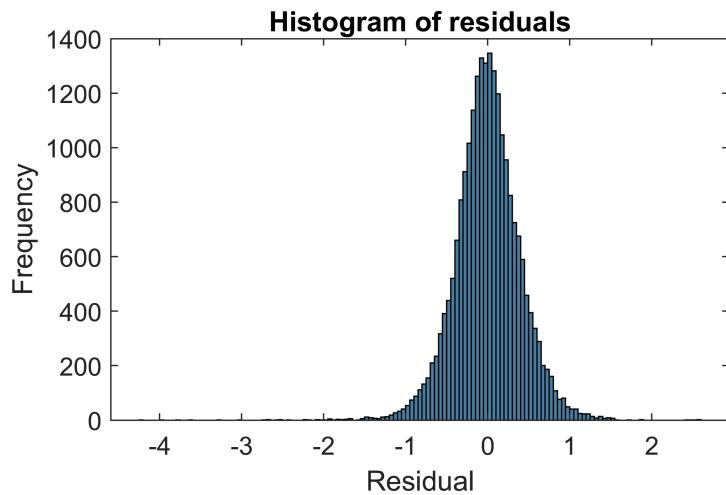


With the help of the additional features the standard deviation is even smaller.

Mean: 0.000000

Standard deviation: 0.399064

The histogram of the residuals of this model (so far the most accurate) looks as follows:



## Excercise 9

### Predict price using mape instead of residual-sum-of-squares

Instead of the residual sum of squares (RSS), the mean absolute percentage error (MAPE) can also be used as a cost function. The formula of the mape is defined as follows:

$$\text{MAPE} := \frac{1}{N} \sum_{i=1}^N \frac{|\hat{y}_i - y_i|}{y_i}$$

In order to perform the gradient descent procedure with the mean absolute percentage error, the resulting costs must be minimized using the derivatives. The derivatives of the MAPE function are as follows:

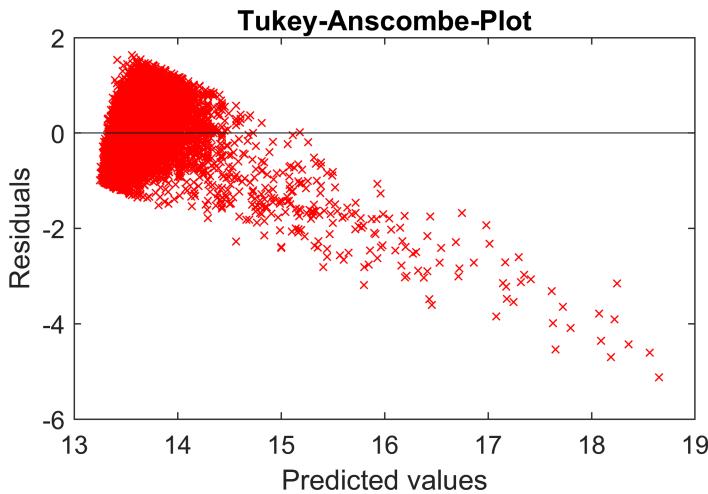
$$\frac{\partial \text{MAPE}}{\partial \hat{y}_i} = \begin{cases} -\frac{1}{N}, & \text{if } \hat{y}_i < y_i \\ \text{undefined,} & \text{if } \hat{y}_i = y_i \\ \frac{1}{N}, & \text{if } \hat{y}_i > y_i \end{cases}$$

Within the gradient descent function, the code part in which the costs are calculated must be exchanged. In the function `gradientDescentMape` this adjustment has been made and the MAPE is used instead of RSS.

Theta found by gradient descent:

```
13.624230
0.270425
```

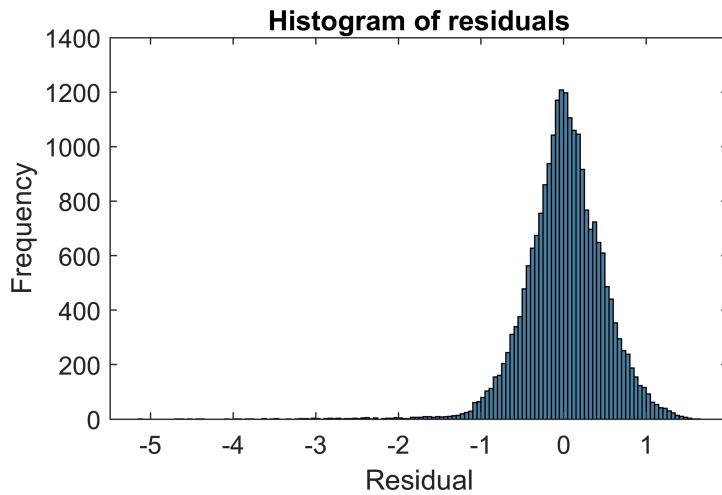
The obtained theta almost matches the results from the previous exercises and also the plot with the residuals looks plausible.



Mean: 0.017496

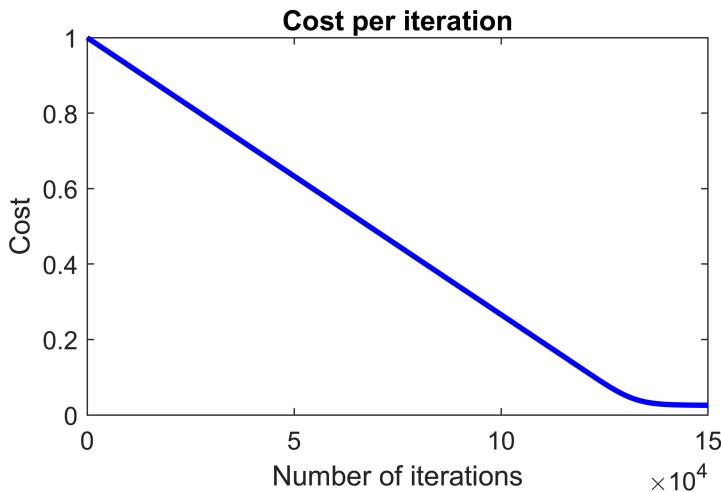
Standard deviation: 0.474671

The obtained standard deviation of ~0.5 is also reflected in the diagramm:

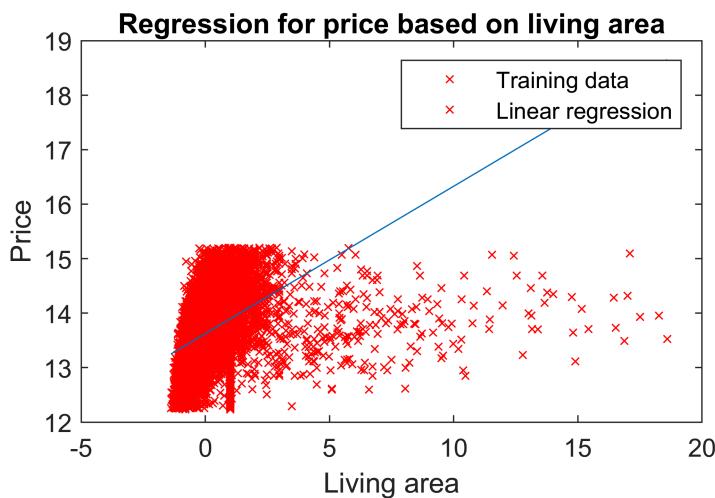


### Visualize costs per iteration

To ensure that the gradient descent procedure has been implemented correctly, the development of costs per iteration is plotted. The development of costs behaves as expected: By correcting the values using the derivations described above, costs can be minimized. The plot shows that the costs settle at a minimum.



In the next plot, the regression line is displayed together with the data points from the training data set. The plot shows that the function obtained corresponds roughly to the model in which RSS was used.



If MAPE is used as a cost function, but the model is then evaluated with RSS, it must be noted that two different aspects were considered. Accordingly, the predictions could be less accurate.