

# Homework2\_report\_Jan\_Wojtas

March 11, 2024

## 1 Dataset

	no2	cars	temp	windspeed	tempdiff	winddir	hour	day
0	22.199949	2206.008383	6.4	3.5	-0.3	56	14	196
1	80.500227	1044.997738	-7.2	1.7	1.2	74	23	143
2	77.200042	1840.992302	-1.3	2.6	-0.1	65	11	115
3	46.200009	333.999668	-3.1	1.8	0.3	78	2	55
4	88.399826	3323.986186	1.0	1.2	1.5	215	18	47

First we analyze the distributions and pairplots of features present in the dataset. For example we can notice that the NO2 distribution has heavy tail. Pairplots for features (excluding target variable) do not show any particular strong correlation. Another note is that treating day as numerical predictor might not be that sensible, however any one hot encoding for this variable makes even less sense.



	no2	cars	temp	windspeed	tempdiff	\
count	251.000000	251.000000	251.000000	251.000000	251.000000	
mean	50.510360	1598.581544	0.707171	3.008765	0.333068	
std	39.490877	1158.345338	5.626803	1.733783	0.877805	
min	3.900013	75.000141	-18.600000	0.300000	-3.500000	
25%	24.799909	456.498585	-2.800000	1.700000	-0.100000	
50%	45.399904	1444.993352	1.500000	2.800000	0.200000	
75%	64.950166	2442.996788	4.800000	4.150000	0.800000	
max	324.099316	4224.009186	12.200000	9.900000	4.300000	

	winddir	hour	day
count	251.000000	251.000000	251.000000
mean	151.633466	12.139442	119.358566

std	87.117949	6.965664	52.133664
min	2.000000	1.000000	32.000000
25%	77.000000	6.000000	78.000000
50%	135.000000	12.000000	119.000000
75%	224.000000	18.000000	165.000000
max	359.000000	24.000000	212.000000

## 2 Feature collinearity

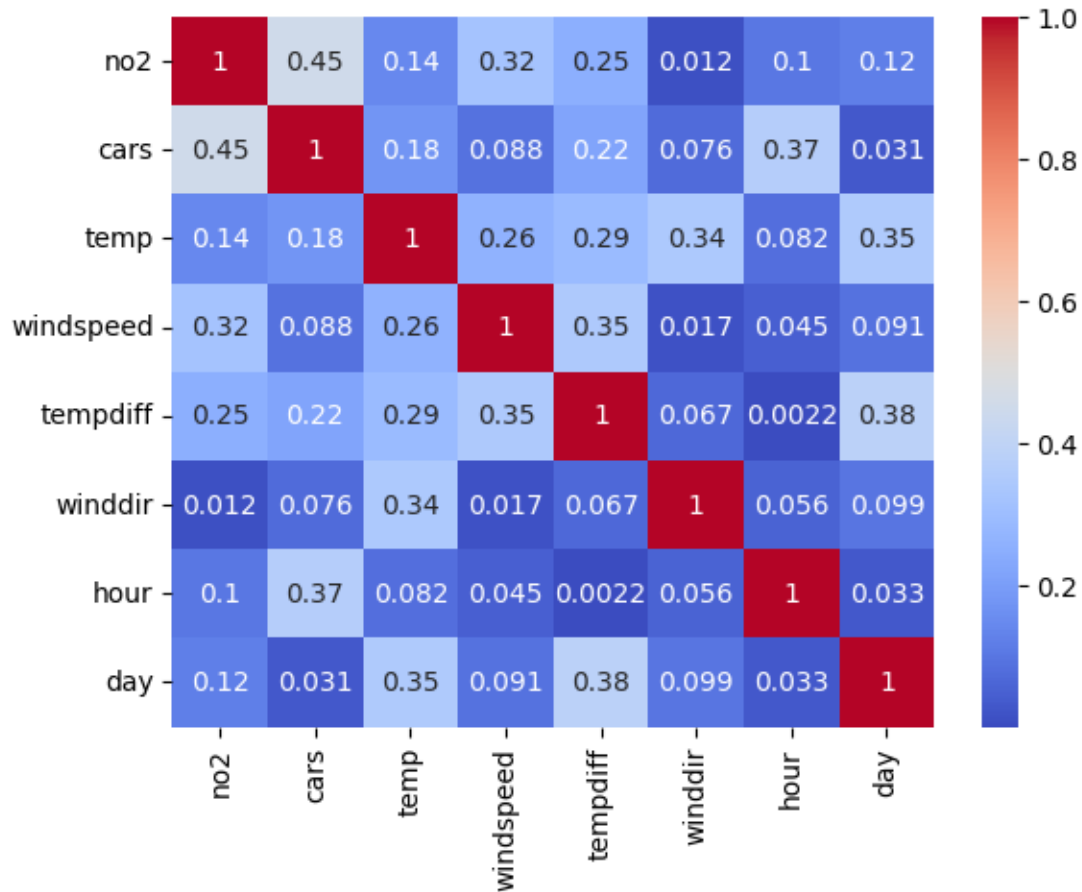
We verify if there are any collinearities between predictors and thus if any of the predictors should be removed from our training set. We use Variance Inflation Factor (VIF) and correlation matrix for validation of our findings.

### 2.1 VIF

	Feature	VIF
0	const	24.241779
1	cars	1.243074
2	temp	1.512238
3	windspeed	1.192568
4	tempdiff	1.418973
5	winddir	1.229814
6	hour	1.167743
7	day	1.351740

All of the features have  $VIF \sim 1$ , suggesting no major collinearities. The variance of the regression coefficients is not inflated. We validate that observation by looking at the correlation matrix. As we can see, the predictors are not highly correlated, which is a good sign. The biggest observable correlation is between `no2` and `cars` ( $\sim 0.45$ ), which might be a suggestion that the number of cars is a valuable predictor.

## 2.2 Correlation matrix



## 3 Linear regression on full data

### 3.1 Summary

Dep. Variable:	no2	R-squared:	0.411
Model:	OLS	Adj. R-squared:	0.394
Method:	Least Squares	F-statistic:	24.21
Date:	Mon, 11 Mar 2024	Prob (F-statistic):	6.77e-25
Time:	18:00:57	Log-Likelihood:	-1211.9
No. Observations:	251	AIC:	2440.
Df Residuals:	243	BIC:	2468.
Df Model:	7		
Covariance Type:	nonrobust		

	coef	std err	t	P>  t	[0.025	0.975]
<b>Intercept</b>	18.6681	4.777	3.908	0.000	9.258	28.078
<b>cars</b>	0.0198	0.002	10.603	0.000	0.016	0.024
<b>const</b>	18.6681	4.777	3.908	0.000	9.258	28.078
<b>day</b>	0.0131	0.043	0.303	0.762	-0.072	0.099
<b>hour</b>	-0.5144	0.302	-1.705	0.089	-1.109	0.080
<b>temp</b>	-0.6995	0.425	-1.646	0.101	-1.537	0.138
<b>tempdiff</b>	12.0540	2.639	4.568	0.000	6.857	17.252
<b>winddir</b>	-0.0025	0.025	-0.100	0.920	-0.051	0.046
<b>windspeed</b>	-5.6553	1.225	-4.618	0.000	-8.068	-3.243
<b>Omnibus:</b>	210.841	<b>Durbin-Watson:</b>	2.228			
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	4141.571			
<b>Skew:</b>	3.223	<b>Prob(JB):</b>	0.00			
<b>Kurtosis:</b>	21.827	<b>Cond. No.</b>	1.11e+18			

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 7.93e-28. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

The  $R^2$  and adjusted  $R^2$  indexes are relatively small and might signal not great fit of linear regression model. P-value for F-statistic is small and we can assume that not all of the predictors are insignificant. However, looking at the t-tests for predictor significance, we can observe high p-value for **winddir** and **day** columns, suggesting the removal of these predictors from model. Of course, removing variables can only lead to increase of SSE on the training data, but can improve generalisation on test data and overall interpretability of model.

Also p-values for **hour** and **temp** are higher than 0.05 which could indicate that these variables are also insignificant and removing these variables should at least be considered.

### 3.2 MSE for full model

MSE: 915.0655311951564

## 4 Model after excluding winddir and day predictors

### 4.1 Summary

<b>Dep. Variable:</b>	no2	<b>R-squared:</b>	0.411
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.399
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	34.14
<b>Date:</b>	Mon, 11 Mar 2024	<b>Prob (F-statistic):</b>	2.08e-26
<b>Time:</b>	04:22:17	<b>Log-Likelihood:</b>	-1212.0
<b>No. Observations:</b>	251	<b>AIC:</b>	2436.
<b>Df Residuals:</b>	245	<b>BIC:</b>	2457.
<b>Df Model:</b>	5		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P>  t	[0.025	0.975]
<b>Intercept</b>	19.4244	2.872	6.763	0.000	13.768	25.081
<b>cars</b>	0.0198	0.002	10.654	0.000	0.016	0.023
<b>const</b>	19.4244	2.872	6.763	0.000	13.768	25.081
<b>hour</b>	-0.5170	0.300	-1.721	0.086	-1.109	0.075
<b>temp</b>	-0.6792	0.368	-1.844	0.066	-1.405	0.046
<b>tempdiff</b>	11.7343	2.461	4.768	0.000	6.887	16.582
<b>windspeed</b>	-5.6911	1.212	-4.696	0.000	-8.078	-3.304
<b>Omnibus:</b>	210.498	<b>Durbin-Watson:</b>	2.225			
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	4120.157			
<b>Skew:</b>	3.217	<b>Prob(JB):</b>	0.00			
<b>Kurtosis:</b>	21.777	<b>Cond. No.</b>	9.70e+18			

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 1.04e-29. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

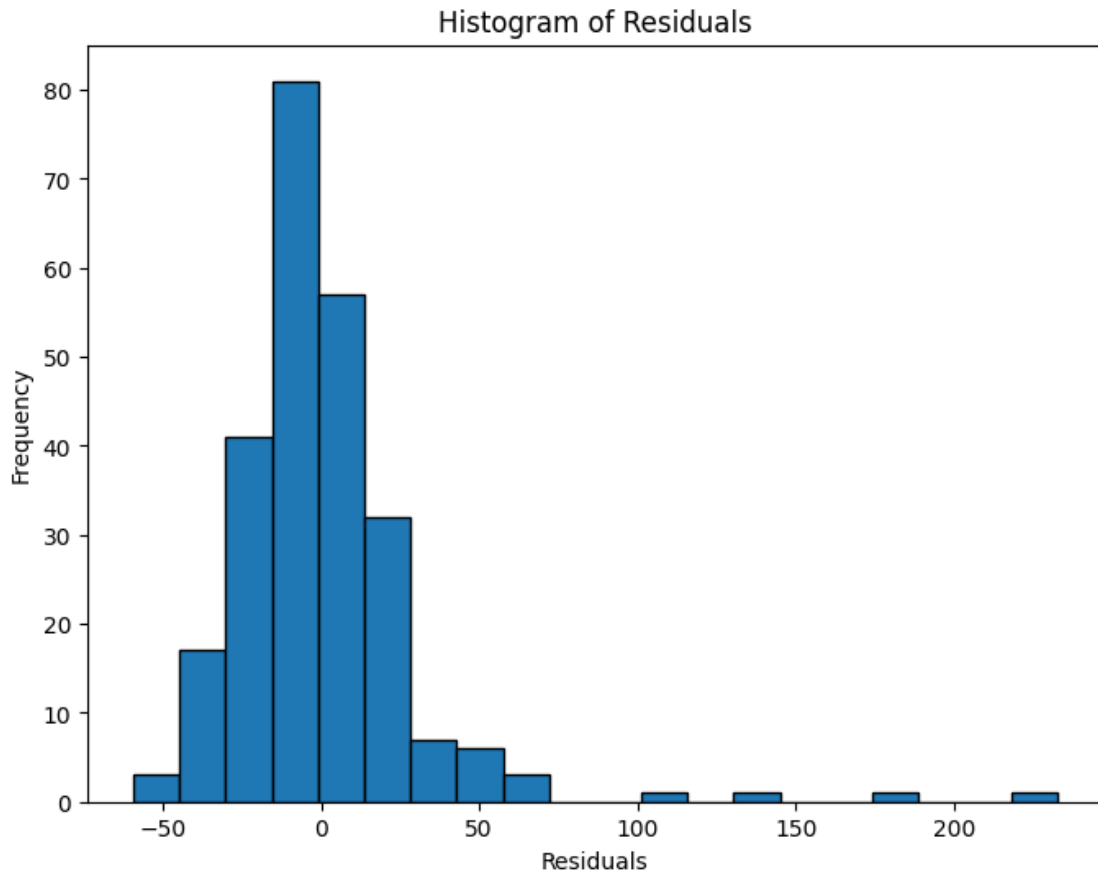
## 4.2 MSE

915.5139090872674

After excluding 2 day and winddir variables from the set of predictors we barely increase MSE on training, while reduce complexity of model and potentially improve generalization.

## 5 Diagnostics

### 5.1 Distribution of residuals

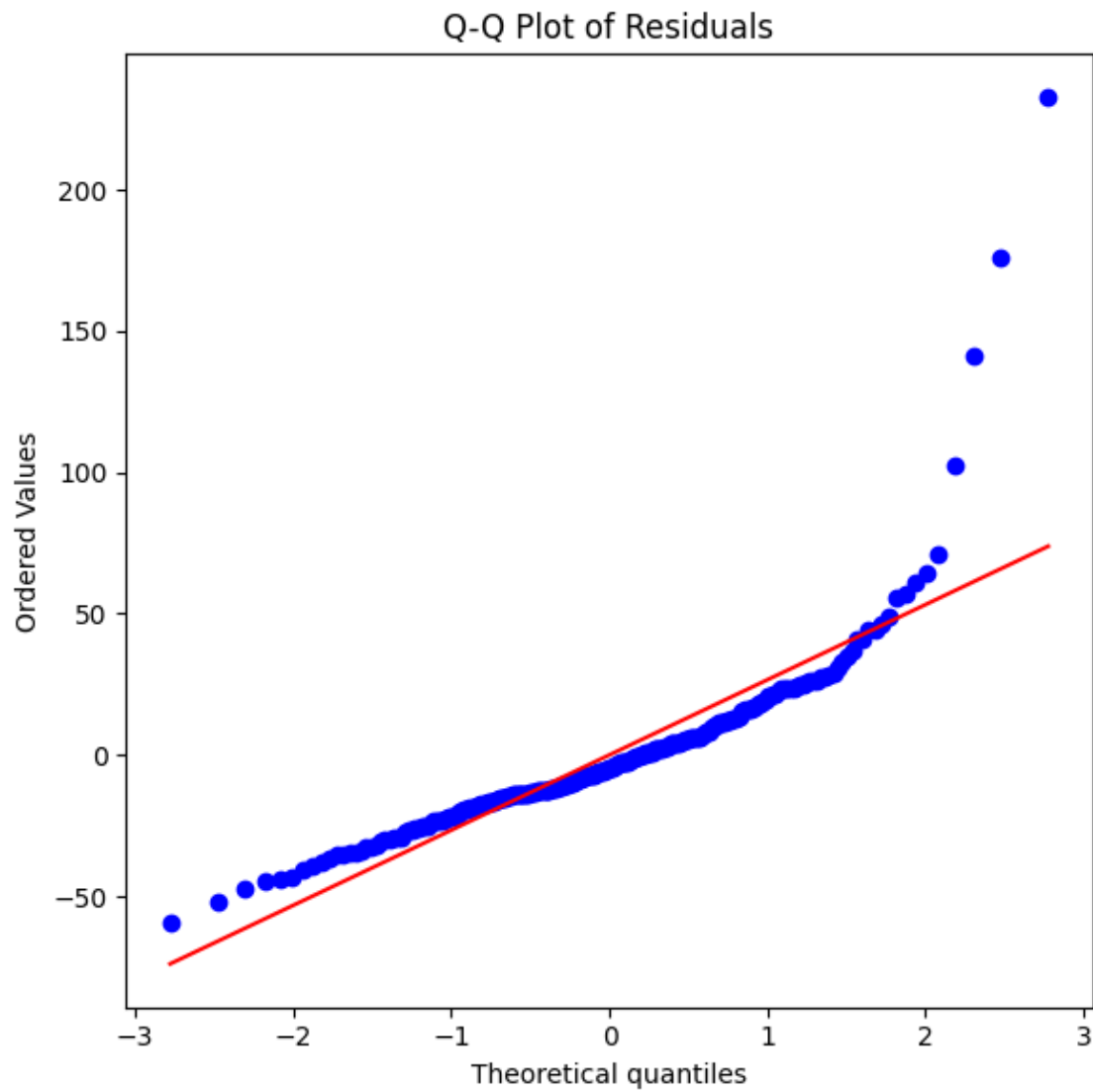


Shapiro-Wilk Test Statistic: 0.7688661343173065

Shapiro-Wilk p-value: 1.6377012765726963e-18

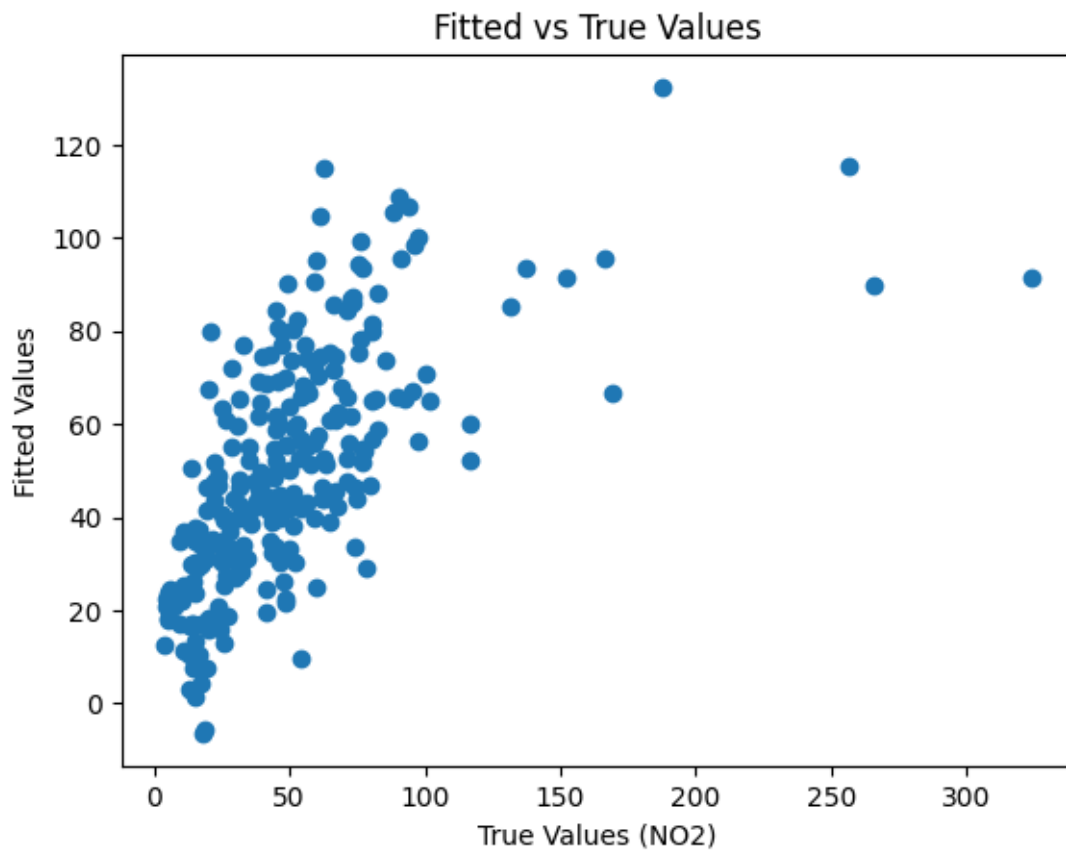
The residuals of the linear regression model should optimally have normal distribution. The histogram shows some concern and Shapiro-Wilk test was conducted to statistically verify the hypothesis of normality of residuals. The p-value is very low, suggesting to discard the null hypothesis about the normality of the residuals.

## 5.2 QQ-plot

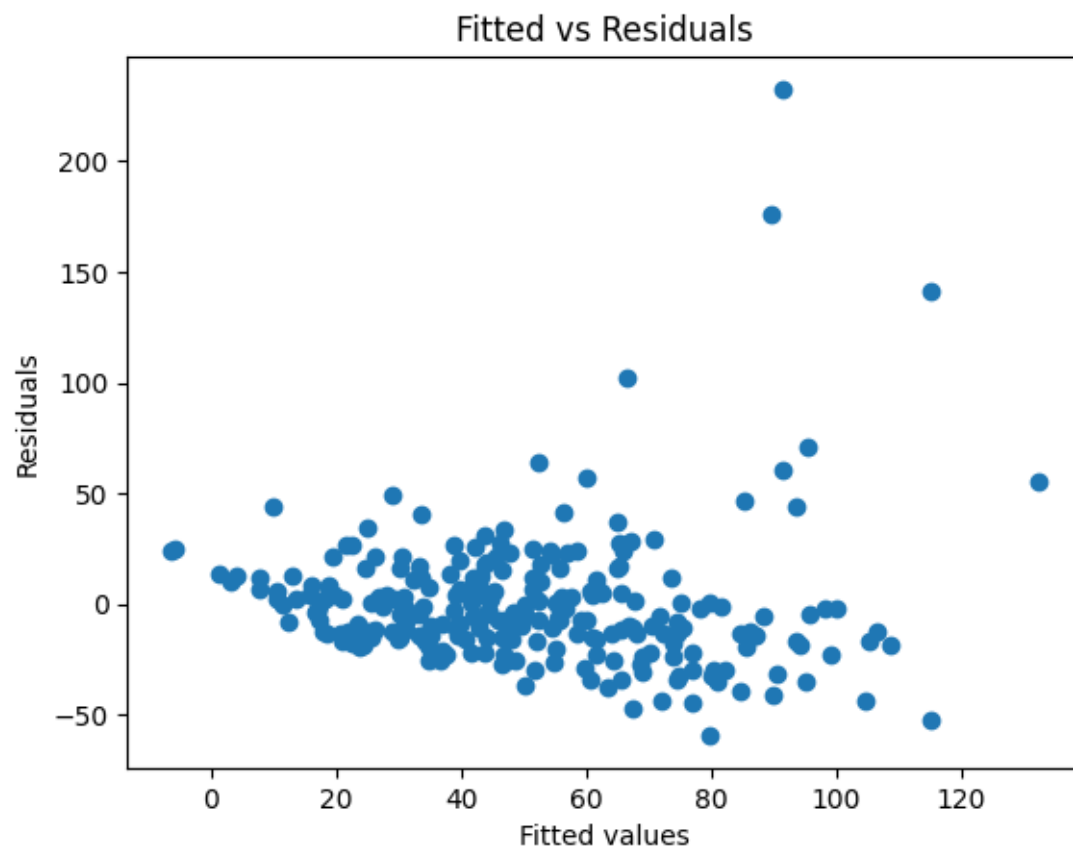


The QQ-plot reinforces the hypotheses that the residuals do not fulfill normality assumption.

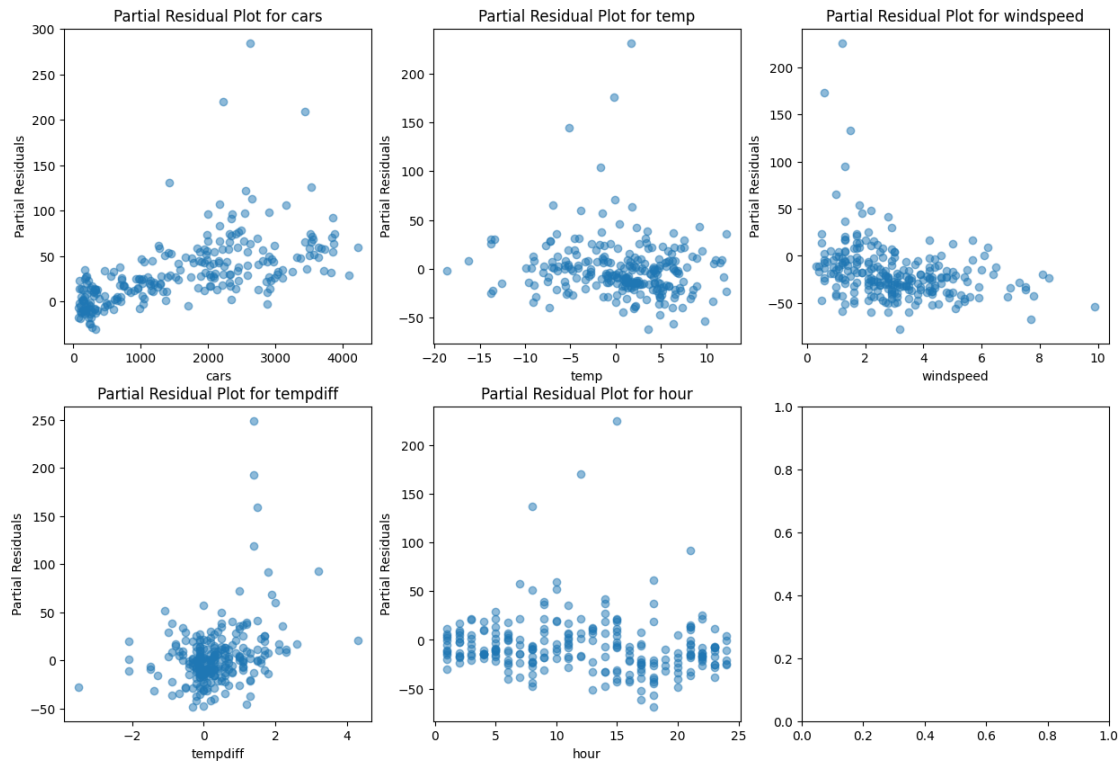




The desirable shape for the points from this plot would be a straight line. Unfortunately we can see some anomalies for larger values of the NO<sub>2</sub>.



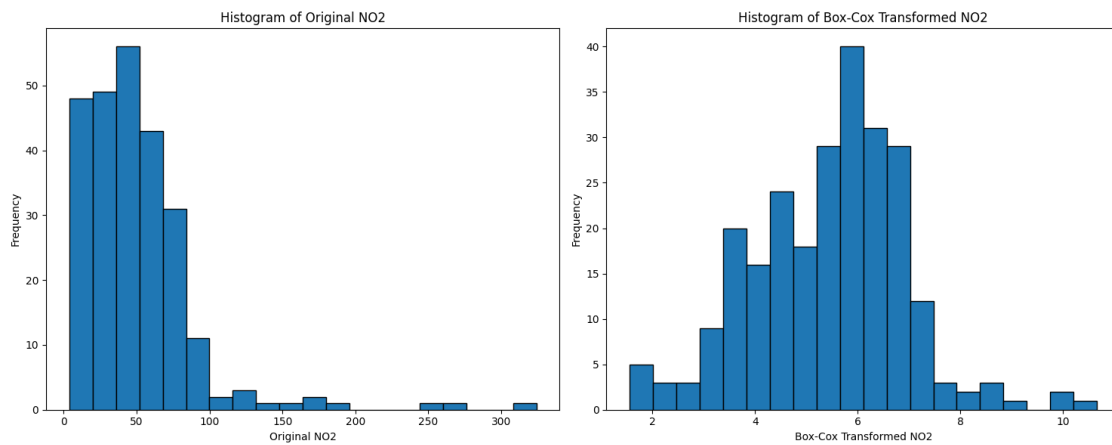
### 5.3 Partial residual plots



### 5.4 Box-cox transformation

One possibility to stabilize variance and normalize distributions of `N02` and `cars` variables is to apply Box-Cox transformation. After the transformation, I verify if it gives some improvement for the model.

Lambda value used for transformation: 0.1939379324703674



As we can see, box-cox transformation normalized values of NO2 variable. We do the same for `cars` variable and fit a model with transformed ones.

## 6 Model after transformation

### 6.1 Summary

<b>Dep. Variable:</b>	no2_boxcox	<b>R-squared:</b>	0.553
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.543
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	60.52
<b>Date:</b>	Mon, 11 Mar 2024	<b>Prob (F-statistic):</b>	6.96e-41
<b>Time:</b>	18:07:30	<b>Log-Likelihood:</b>	-355.77
<b>No. Observations:</b>	251	<b>AIC:</b>	723.5
<b>Df Residuals:</b>	245	<b>BIC:</b>	744.7
<b>Df Model:</b>	5		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P>  t	[0.025	0.975]
<b>Intercept</b>	1.9732	0.110	17.930	0.000	1.756	2.190
<b>cars_boxcox</b>	0.0467	0.003	14.171	0.000	0.040	0.053
<b>const</b>	1.9732	0.110	17.930	0.000	1.756	2.190
<b>hour</b>	-0.0308	0.011	-2.922	0.004	-0.052	-0.010
<b>temp</b>	-0.0457	0.012	-3.756	0.000	-0.070	-0.022
<b>tempdiff</b>	0.4094	0.082	4.991	0.000	0.248	0.571
<b>windspeed</b>	-0.2521	0.040	-6.299	0.000	-0.331	-0.173

<b>Omnibus:</b>	1.314	<b>Durbin-Watson:</b>	2.188
<b>Prob(Omnibus):</b>	0.518	<b>Jarque-Bera (JB):</b>	1.010
<b>Skew:</b>	0.119	<b>Prob(JB):</b>	0.604
<b>Kurtosis:</b>	3.200	<b>Cond. No.</b>	3.03e+17

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 1e-29. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

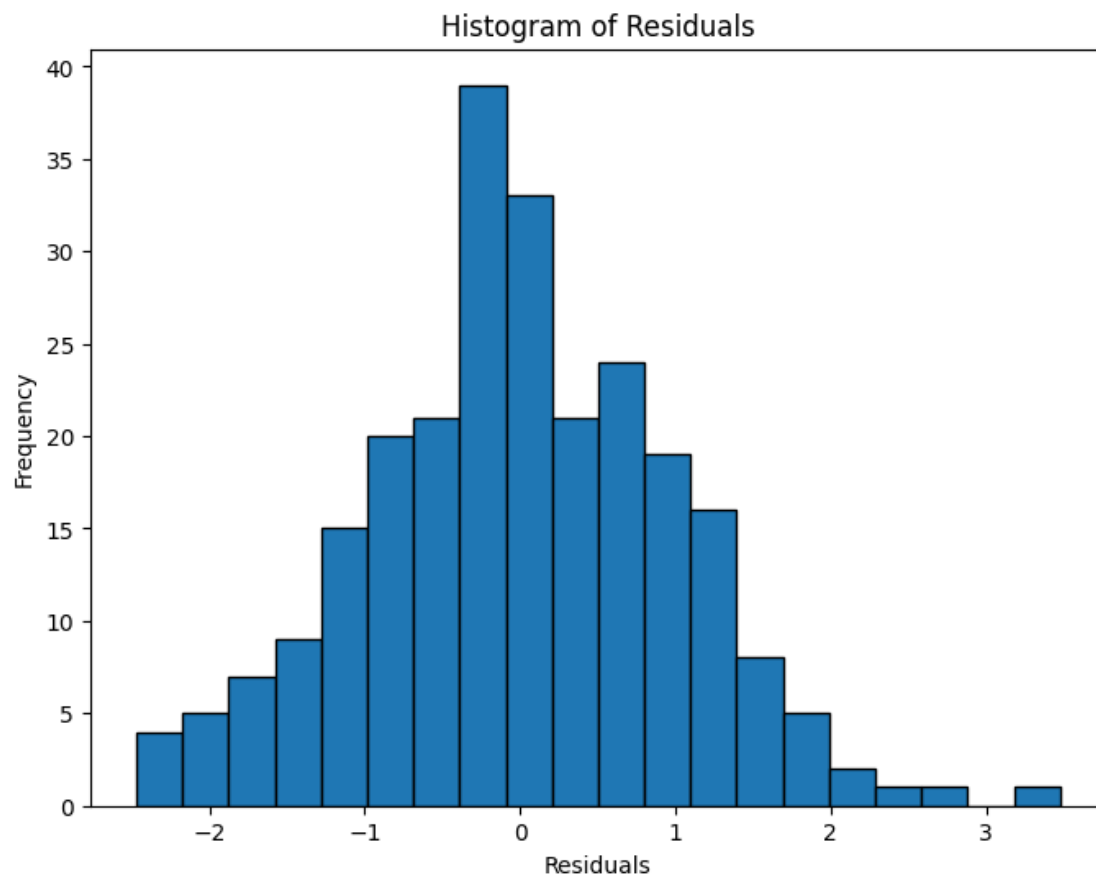
In order to calculate MSE, we need to find the inverse function of our Box-Cox transformation.

### 6.2 MSE

MSE: 862.3729217742655

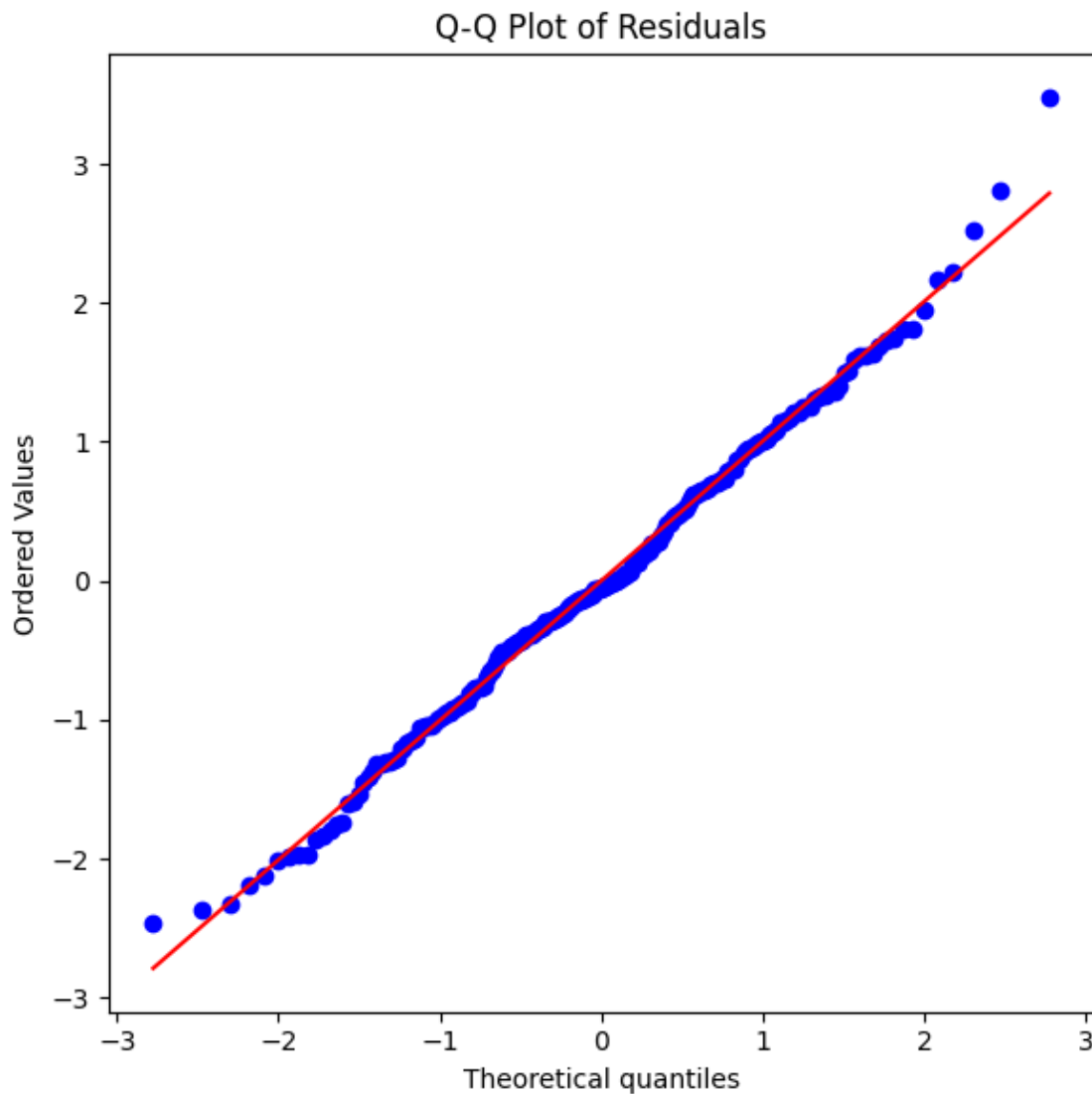
The MSE is significantly reduced and  $R^2$  index increased by over 0.1, which is a solid improvement.

### 6.3 Residuals



Shapiro-Wilk Test Statistic: 0.9948681045940815

Shapiro-Wilk p-value: 0.5631110423306341



Also the residuals improved and they pass Shapiro-Wilk test for normality of distribution. Quantiles on QQ-plot almost perfectly resemble the normal distribution.

## 7 Outlier detection

At last we analyse the influence of outlier removal on the results. We use Cook's distance for outlier detection and set threshold based on sample size (typically  $\frac{4}{N}$  is considered a reasonable threshold)

<b>Dep. Variable:</b>	no2_boxcox	<b>R-squared:</b>	0.565
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.556
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	61.32
<b>Date:</b>	Mon, 11 Mar 2024	<b>Prob (F-statistic):</b>	9.12e-41
<b>Time:</b>	18:10:10	<b>Log-Likelihood:</b>	-320.16
<b>No. Observations:</b>	242	<b>AIC:</b>	652.3
<b>Df Residuals:</b>	236	<b>BIC:</b>	673.2
<b>Df Model:</b>	5		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P>  t	[0.025	0.975]
<b>Intercept</b>	1.9831	0.101	19.704	0.000	1.785	2.181
<b>cars_boxcox</b>	0.0445	0.003	14.496	0.000	0.038	0.051
<b>const</b>	1.9831	0.101	19.704	0.000	1.785	2.181
<b>hour</b>	-0.0280	0.010	-2.876	0.004	-0.047	-0.009
<b>temp</b>	-0.0478	0.011	-4.175	0.000	-0.070	-0.025
<b>tempdiff</b>	0.3677	0.077	4.757	0.000	0.215	0.520
<b>windspeed</b>	-0.2364	0.037	-6.364	0.000	-0.310	-0.163

<b>Omnibus:</b>	1.507	<b>Durbin-Watson:</b>	2.083
<b>Prob(Omnibus):</b>	0.471	<b>Jarque-Bera (JB):</b>	1.578
<b>Skew:</b>	-0.150	<b>Prob(JB):</b>	0.454
<b>Kurtosis:</b>	2.742	<b>Cond. No.</b>	4.70e+16

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 3.99e-28. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

MSE: 890.2473820578655

We can notice that the fit of the model is better, but it should be intuitive as we removed the annoying observations. But the MSE calculated on the whole data grows.