

Execute a multiple linear regression (MLR) for BME680 gas readings of a single sensor

Read the CSV file created by the CCU Historian, ensure that field separator is ';' and decimal separator is '.', if necessary edit CSV file before reading it by the provided script 'csv_convert_historian.bsh'

The provided script 'get_new_history.bsh' is search for the CCU Historian's CSV in the directory '\${HOME}/Downloads'. The conversion script 'csv_convert_historian.bsh' is invoked inside 'csv_convert_historian.bsh'

```
In [1]: 1 import warnings
2 warnings.simplefilter(action='ignore', category=FutureWarning)
3
4 import pandas as pd
5 from datetime import datetime
6
7 import numpy as np
8
9
10 dateparse = lambda x: pd.datetime.strptime(x, '%d.%m.%Y %H:%M:%S,%f')
11
12 df0 = pd.read_csv("historian.csv", sep=';', thousands=".", decimal=",", skiprows = [0,1,2],c
13
14 df0.head(9)
15 #df0['Datum']
16 #type(df0['sensor 1'][0])
```

Out[1]:

	Datum	Mode	raw_gas_resistance	relative_humidity	temperature
0	2020-12-31 16:26:56.439	2	147800	36.9	23.278
1	2020-12-31 16:30:32.862	2	147800	36.9	22.800
2	2020-12-31 16:30:32.866	2	147800	35.6	22.800
3	2020-12-31 16:30:32.873	2	157780	35.6	22.800
4	2020-12-31 16:35:04.477	2	157780	35.6	22.600
5	2020-12-31 16:35:04.486	2	157780	36.5	22.600
6	2020-12-31 16:35:04.492	2	161480	36.5	22.600
7	2020-12-31 16:39:36.091	2	161480	36.5	22.700
8	2020-12-31 16:39:36.099	2	161480	36.8	22.700

```
In [2]: 1 # keep every 3rd row (CCU historian is tracking every change of a datapoint separately)
2 # each three consecutive entries in history.csv are identical; therefore we take every third
3 df = df0[(df0.index % 3 == 0)]
4
5 df.head(9)
```

Out[2]:

	Datum	Mode	raw_gas_resistance	relative_humidity	temperature
0	2020-12-31 16:26:56.439	2	147800	36.9	23.278
3	2020-12-31 16:30:32.873	2	157780	35.6	22.800
6	2020-12-31 16:35:04.492	2	161480	36.5	22.600
9	2020-12-31 16:39:36.106	2	162460	36.8	22.700
12	2020-12-31 16:44:07.740	2	160520	37.2	22.800
15	2020-12-31 16:48:39.353	2	158000	37.5	22.900
18	2020-12-31 16:53:10.973	2	155280	37.8	22.900
21	2020-12-31 16:57:42.592	2	153440	38.1	23.000
24	2020-12-31 17:02:14.215	2	152540	38.1	23.000

```
In [3]: 1 df.head(-1)
```

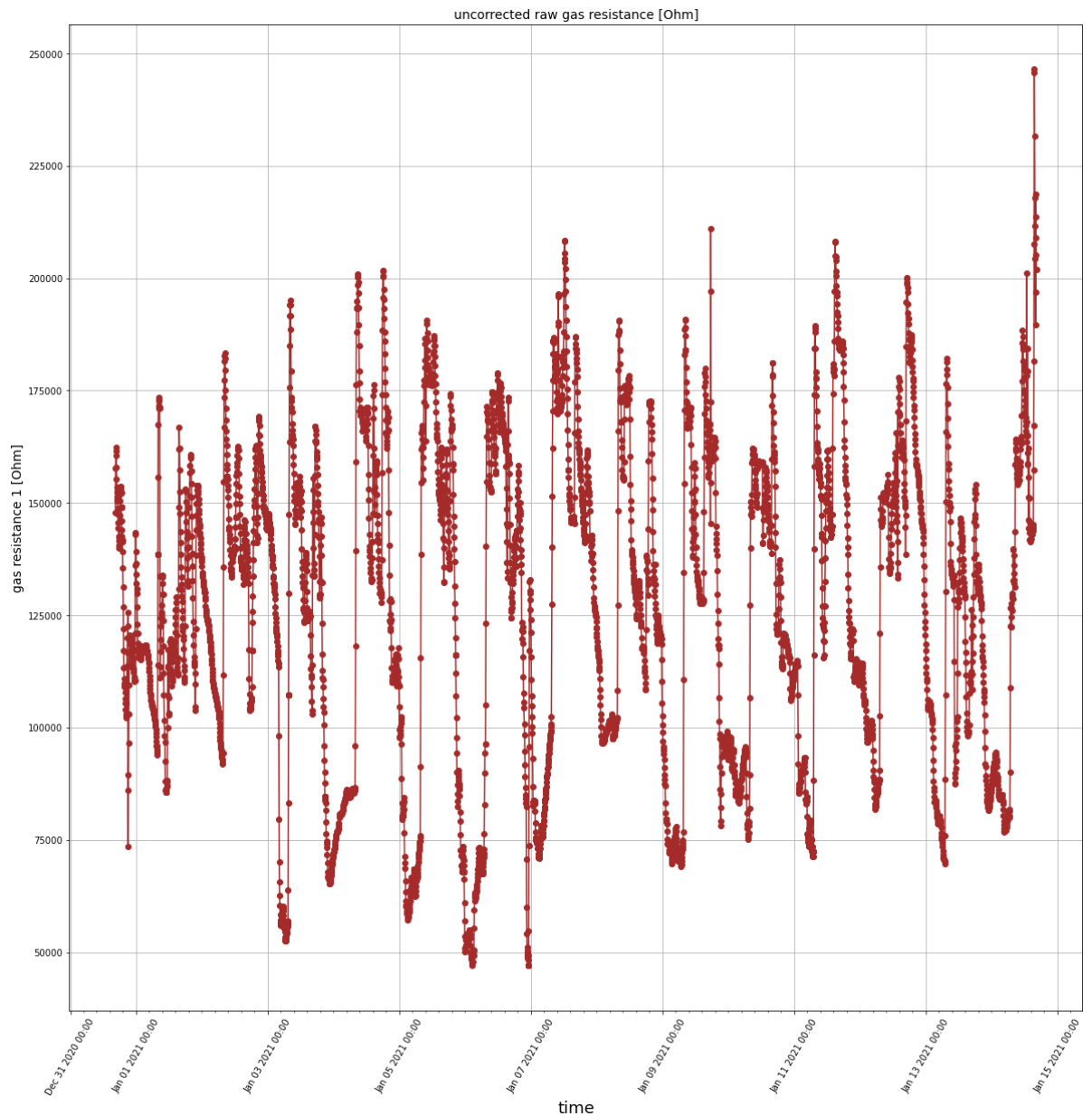
Out [3]:

	Datum	Mode	raw_gas_resistance	relative_humidity	temperature
0	2020-12-31 16:26:56.439	2	147800	36.9	23.278
3	2020-12-31 16:30:32.873	2	157780	35.6	22.800
6	2020-12-31 16:35:04.492	2	161480	36.5	22.600
9	2020-12-31 16:39:36.106	2	162460	36.8	22.700
12	2020-12-31 16:44:07.740	2	160520	37.2	22.800
...
14007	2021-01-14 16:08:06.854	2	189640	27.8	26.100
14010	2021-01-14 16:09:02.681	2	213600	28.0	25.000
14013	2021-01-14 16:11:22.754	2	218780	26.9	25.500
14016	2021-01-14 16:15:55.504	2	209040	26.8	26.400
14019	2021-01-14 16:20:29.140	2	205100	27.1	26.100

4674 rows × 5 columns

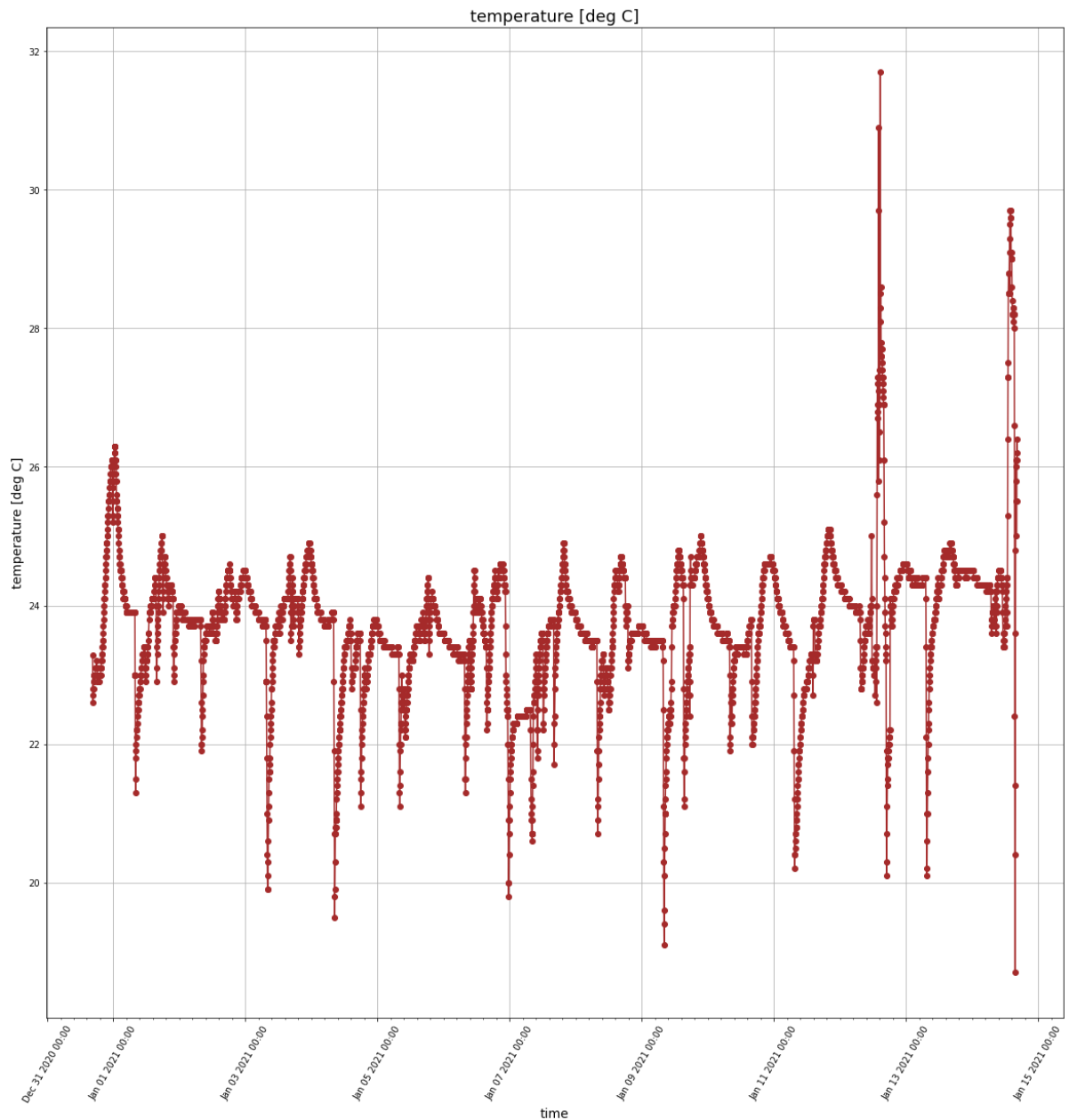
Time series diagram of the measured raw gas resistance

```
In [4]: 1 import matplotlib.pyplot as plt
2 import matplotlib.dates as mdates
3 from matplotlib.dates import DateFormatter
4 from matplotlib.ticker import (MultipleLocator, FormatStrFormatter,
5                               AutoMinorLocator)
6
7 fig, ax = plt.subplots(figsize=(20, 20))
8 plt.xticks(rotation=60)
9 ax.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
10
11 ax.xaxis.set_minor_locator(AutoMinorLocator())
12
13 ax.plot_date(df['Datum'], df['raw_gas_resistance'], linestyle='solid', color='brown')
14 plt.title('uncorrected raw gas resistance [Ohm]', fontsize=14)
15 plt.xlabel('time', fontsize=18)
16 plt.ylabel('gas resistance 1 [Ohm]', fontsize=14)
17 plt.grid(True)
18 plt.show()
```



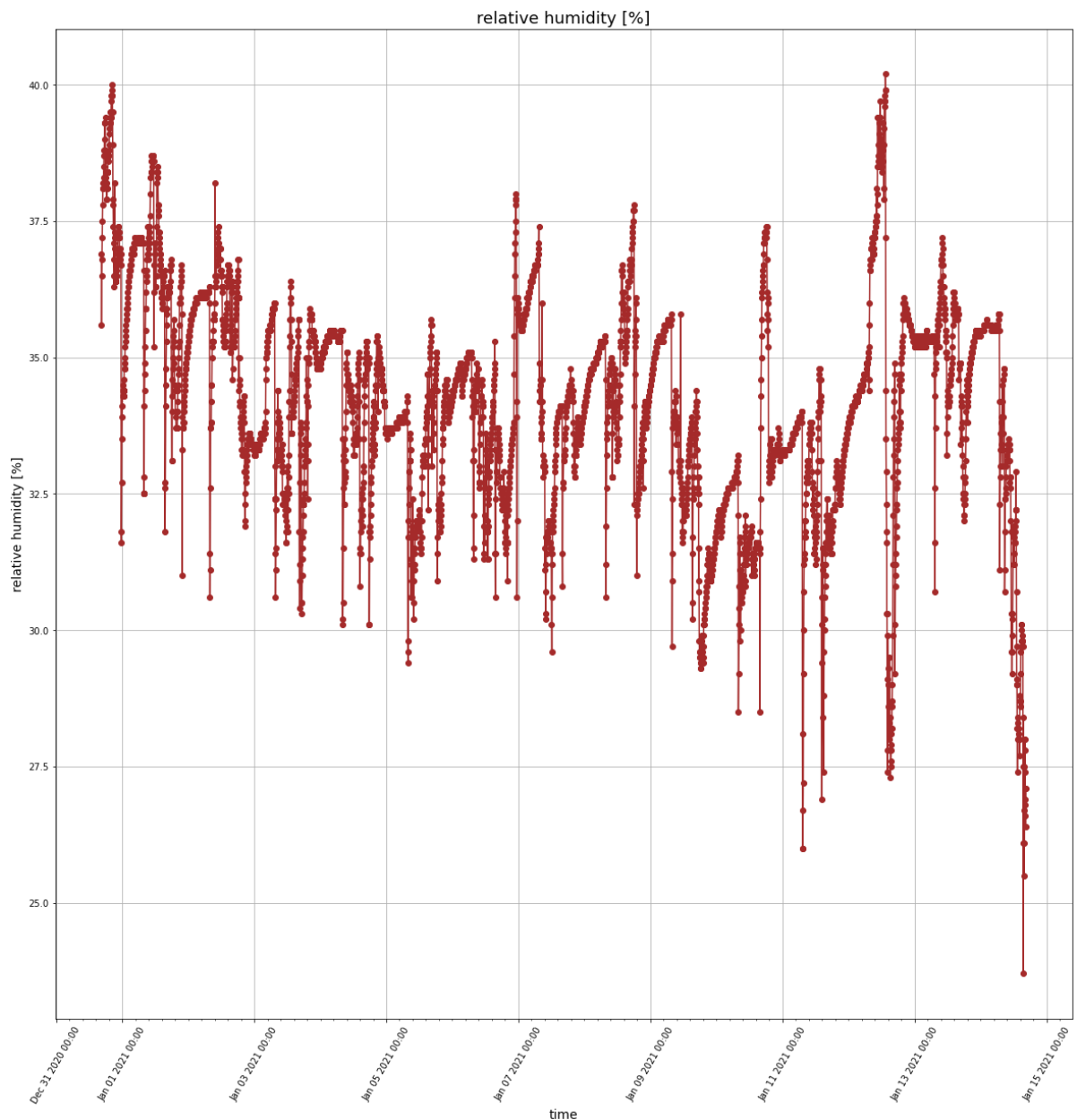
Time series diagram of the measured temperature

```
In [5]: 1 import matplotlib.pyplot as plt
2 import matplotlib.dates as mdates
3 from matplotlib.dates import DateFormatter
4 from matplotlib.ticker import (MultipleLocator, FormatStrFormatter,
5                               AutoMinorLocator)
6
7 fig, ax = plt.subplots(figsize=(20, 20))
8 plt.xticks(rotation=60)
9 ax.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
10
11 ax.xaxis.set_minor_locator(AutoMinorLocator())
12
13 ax.plot_date(df['Datum'], df['temperature'], linestyle='solid', color='brown')
14 plt.title('temperature [deg C]', fontsize=18)
15 plt.xlabel('time', fontsize=14)
16 plt.ylabel('temperature [deg C]', fontsize=14)
17 plt.grid(True)
18 plt.show()
```



Time series diagram of the measured relative humidity

```
In [6]: 1 import matplotlib.pyplot as plt
2 import matplotlib.dates as mdates
3 from matplotlib.dates import DateFormatter
4 from matplotlib.ticker import (MultipleLocator, FormatStrFormatter,
5                               AutoMinorLocator)
6
7 fig, ax = plt.subplots(figsize=(20, 20))
8 plt.xticks(rotation=60)
9 ax.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
10
11 ax.xaxis.set_minor_locator(AutoMinorLocator())
12
13 ax.plot_date(df['Datum'], df['relative_humidity'], linestyle='solid', color='brown')
14 plt.title('relative humidity [%]', fontsize=18)
15 plt.xlabel('time', fontsize=14)
16 plt.ylabel('relative humidity [%]', fontsize=14)
17 plt.grid(True)
18 plt.show()
```



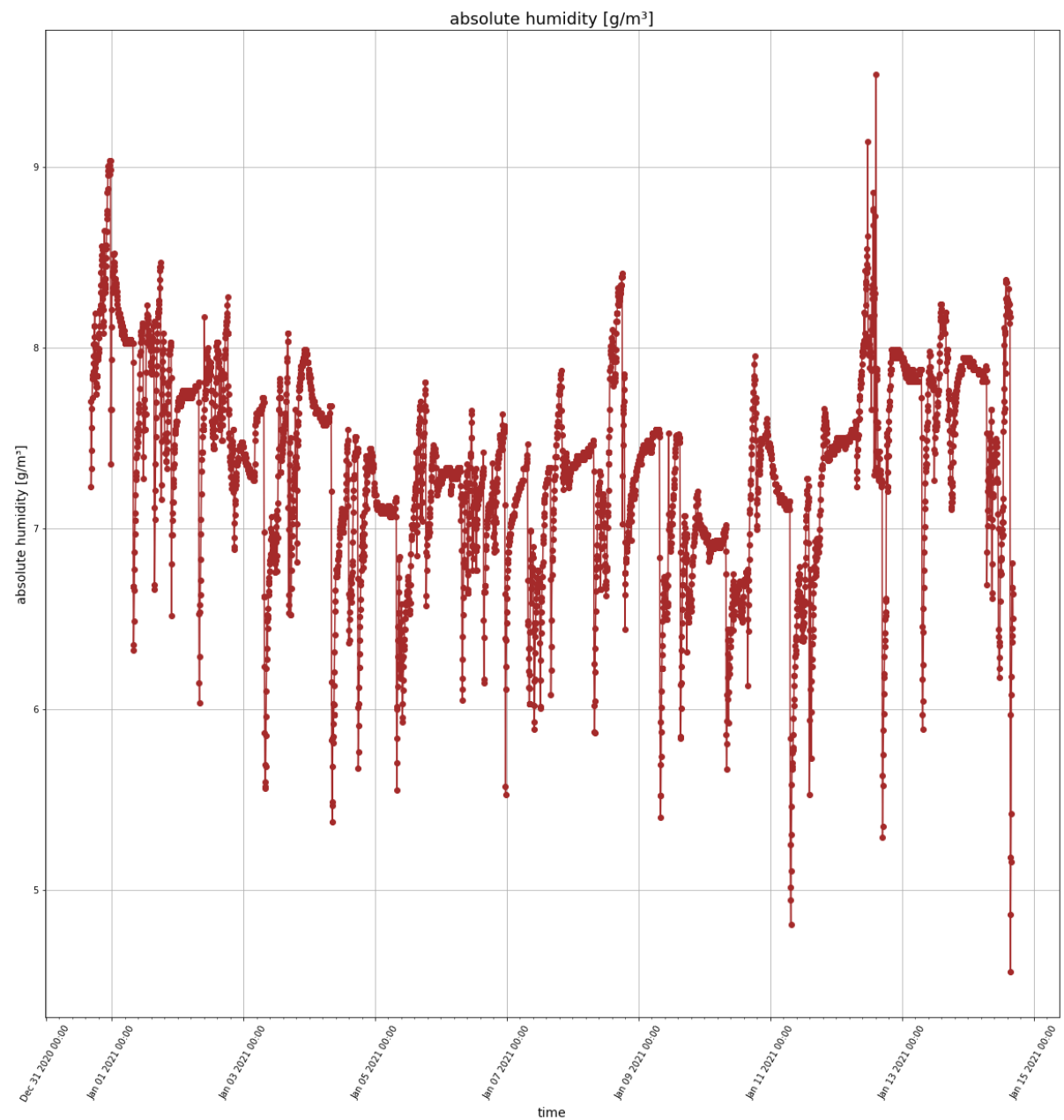
Calculate and plot the absolute humidity from temperature and relative humidity by an approximate formula (same formula as used inside the sensor)

In [7]:

```

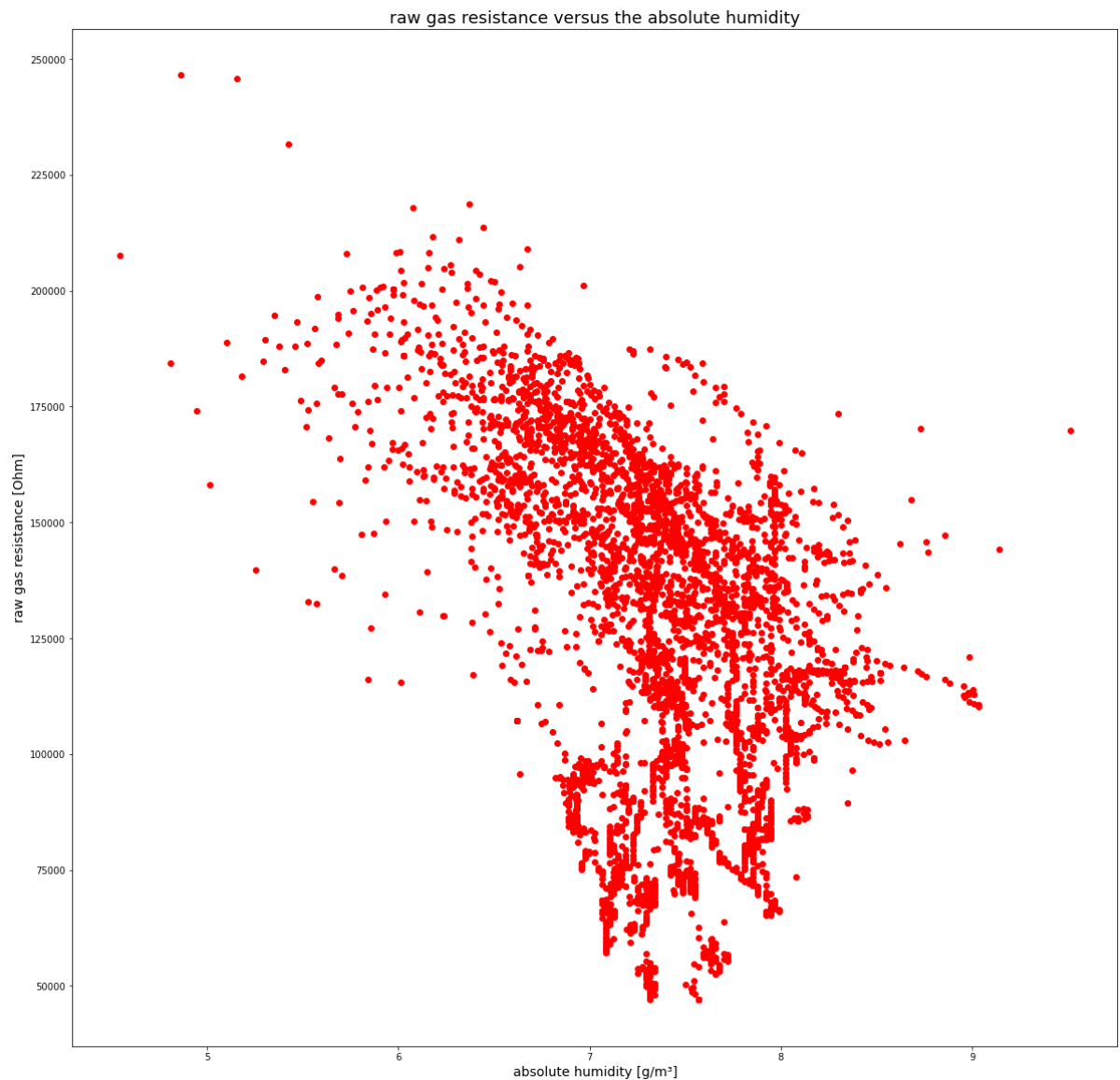
1 import numpy as np
2 # Create a function that calculates the absolute humidity from the two arguments 'temperatu
3 # see for details https://www.kompf.de/weather/vent.html or https://rechneronline.de/barome
4
5 a = 6.112
6 b = 17.67
7 c = 243.5
8
9 # Compute saturated water vapor pressure in hPa
10 # Param t - temperature in °C
11 def svp(t):
12     svp = a * np.exp((b*t)/(c+t))
13     return svp
14
15 # Compute actual water vapor pressure in hPa
16 # Param rh - relative humidity in %
17 # Param t - temperature in °C
18 def vp(rh, t):
19     vp = rh/100. * svp(t)
20     return vp
21
22 # Compute the absolute humidity in g/m³
23 # Param rh - relative humidity in %
24 # Param t - temperature in °C
25 def calculate_absolute_humidity(t, rh):
26     mw = 18.016 # kg/kmol (Molekulargewicht des Wasserdampfes)
27     rs = 8314.3 # J/(kmol*K) (universelle Gaskonstante)
28     ah = 10**5 * mw/rs * vp(rh, t)/(t + 273.15)
29     #return the absolute humidity in [g/m³]
30     return ah
31
32 pd.set_option('mode.chained_assignment', None)
33 # now apply the above defined formulas to get the pandas dataframe column 'absolute_humidity
34 df['absolute_humidity'] = calculate_absolute_humidity(df['temperature'], df['relative_humidi
35
36 import matplotlib.pyplot as plt
37 import matplotlib.dates as mdates
38 from matplotlib.dates import DateFormatter
39 from matplotlib.ticker import MultipleLocator, FormatStrFormatter,
40     AutoMinorLocator
41
42 fig, ax = plt.subplots(figsize=(20, 20))
43 plt.xticks(rotation=60)
44 ax.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
45
46 ax.xaxis.set_minor_locator(AutoMinorLocator())
47
48 ax.plot_date(df['Datum'], df['absolute_humidity'], linestyle='solid', color='brown')
49 plt.title('absolute humidity [g/m³]', fontsize=18)
50 plt.xlabel('time', fontsize=14)
51 plt.ylabel('absolute humidity [g/m³]', fontsize=14)
52 plt.grid(True)
53 plt.show()
54
55

```



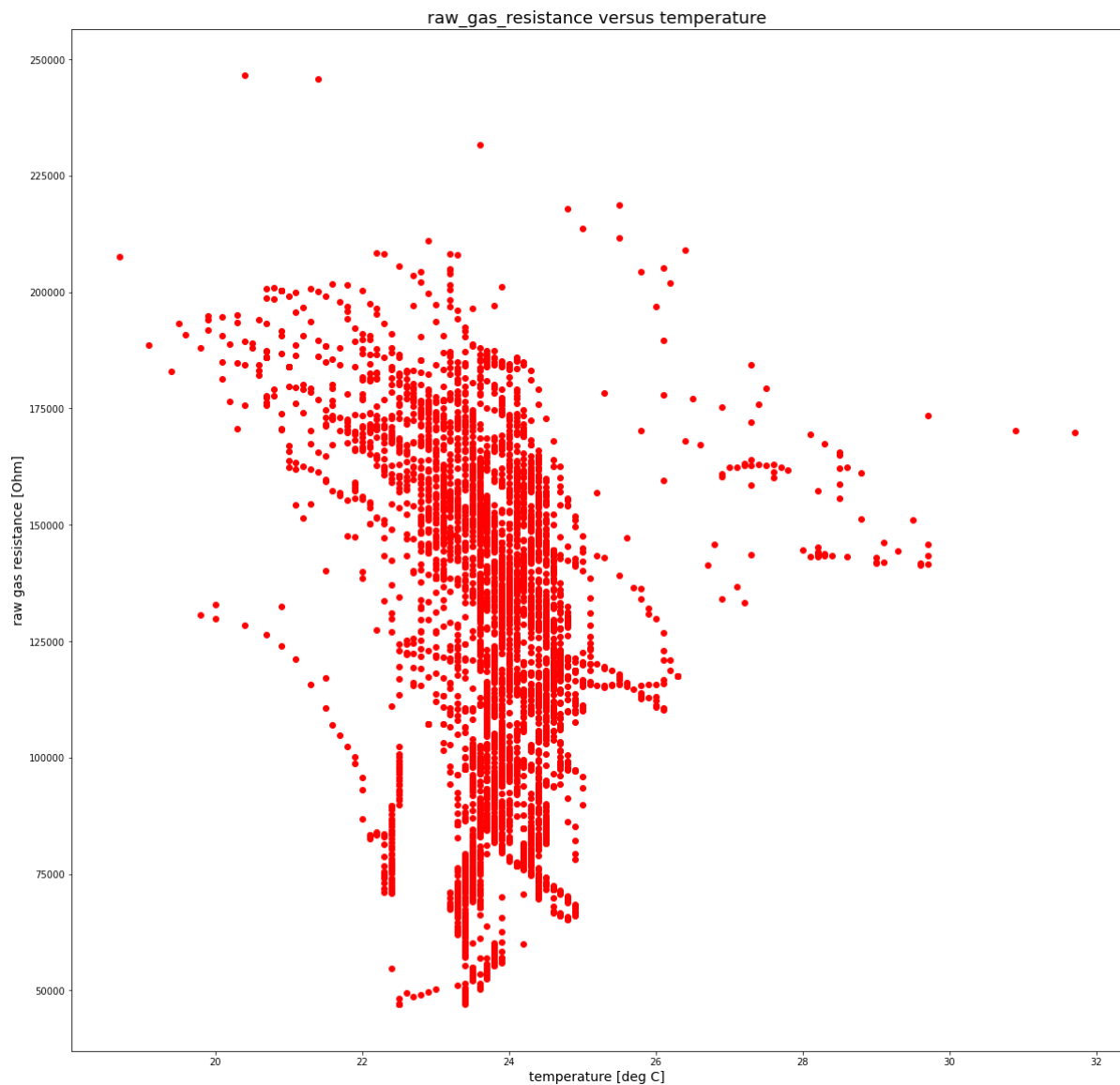
Scatter plot of raw gas resistance versus the absolute humidity, is the dependency somehow linear (should be for a multilinear regression)?

```
In [8]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5
6 plt.figure(figsize=(20,20))
7
8 plt.scatter(df['absolute_humidity'], df['raw_gas_resistance'], color='red')
9 plt.title('raw gas resistance versus the absolute humidity', fontsize=18)
10 plt.xlabel('absolute humidity [g/m³]', fontsize=14)
11 plt.ylabel('raw gas resistance [Ohm]', fontsize=14)
12 plt.show()
```



Scatter plot of raw gas resistance versus the temperature, is the dependency somehow linear (should be for a multilinear regression)?


```
In [9]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5
6 plt.figure(figsize=(20,20))
7 plt.scatter(df['temperature'], df['raw_gas_resistance'], color='red')
8 plt.title('raw_gas_resistance versus temperature', fontsize=18)
9 plt.xlabel('temperature [deg C]', fontsize=14)
10 plt.ylabel('raw gas resistance [Ohm]', fontsize=14)
11 plt.grid(False)
12 plt.show()
```



Execute a multiple linear regression of raw gas resistance on dependency of the absolute humidity and the temperature
use the prediction 'predictions1' of the mutiple linear regression to create a corrected gas resistance 'residuals' with eliminated
influence of the absolute humidity and the temperature
create a normalized scaled corrected gas resistance 'normalized_residuals'

```

In [10]: 1 import pandas as pd
2 from sklearn import linear_model
3 import statsmodels.api as sm
4 import matplotlib.pyplot as plt
5 import matplotlib.dates as mdates
6 from matplotlib.dates import DateFormatter
7 from matplotlib.ticker import (MultipleLocator, FormatStrFormatter,
8                               AutoMinorLocator)
9
10
11 X = df[['temperature','absolute_humidity']] # here we have 2 variables for multiple regressi
12 Y = df['raw_gas_resistance']
13
14 # with sklearn
15 regr = linear_model.LinearRegression()
16 regr.fit(X, Y)
17
18 print('Intercept: \n', regr.intercept_)
19 print('Coefficients: \n', regr.coef_)
20
21 X = sm.add_constant(X)
22 model = sm.OLS(Y, X).fit()
23 predictions = model.predict(X)
24
25 print_model = model.summary()
26 print(print_model)
27 print(model.rsquared)
28
29 residuals=df['raw_gas_resistance']-predictions
30 min_res=min(residuals)
31 max_res=max(residuals)
32
33 #clip min of residual to epsilon in order to avoid a log(0) trap
34 epsilon=0.0001
35
36 normalized_residuals=((residuals-min_res)/(max_res-min_res)).clip(epsilon,None)*100
37
38
39 fig, ax1 = plt.subplots(figsize=(20, 20))
40 plt.xticks(rotation=60)
41 ax1.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
42
43 ax1.xaxis.set_minor_locator(AutoMinorLocator())
44
45 lns1=ax1.plot_date(df['Datum'], predictions, linestyle='solid',color='brown', label='linear
46 lns2=ax1.plot_date(df['Datum'], df['raw_gas_resistance'], linestyle='solid',color='green', l
47 color = 'tab:red'
48 ax1.set_xlabel('time', fontsize=14)
49 ax1.set_ylabel('gas resistance 1 [Ohm]', color=color, fontsize=14)
50
51 ax2 = ax1.twinx() # instantiate a second axes that shares the same x-axis
52 ax2.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
53
54 ax2.xaxis.set_minor_locator(AutoMinorLocator())
55
56 color = 'tab:blue'
57 ax2.set_ylabel('normalized gas resistance', color=color, fontsize=14) # we already handled
58 lns3=ax2.plot_date(df['Datum'], normalized_residuals, linestyle='solid', color='red', label=
59
60 plt.title('compensated gas resistance', fontsize=18)
61
62 ax1.grid(True)
63 lns = lns1+lns2+lns3
64 labs = [l.get_label() for l in lns]
65 ax1.legend(lns, labs, loc="lower left")
66
67 fig.tight_layout() # otherwise the right y-label is slightly clipped
68 plt.show()
69

```

Intercept:
243521.11446688976

Coefficients:
[6603.48933636 -37355.38940408]

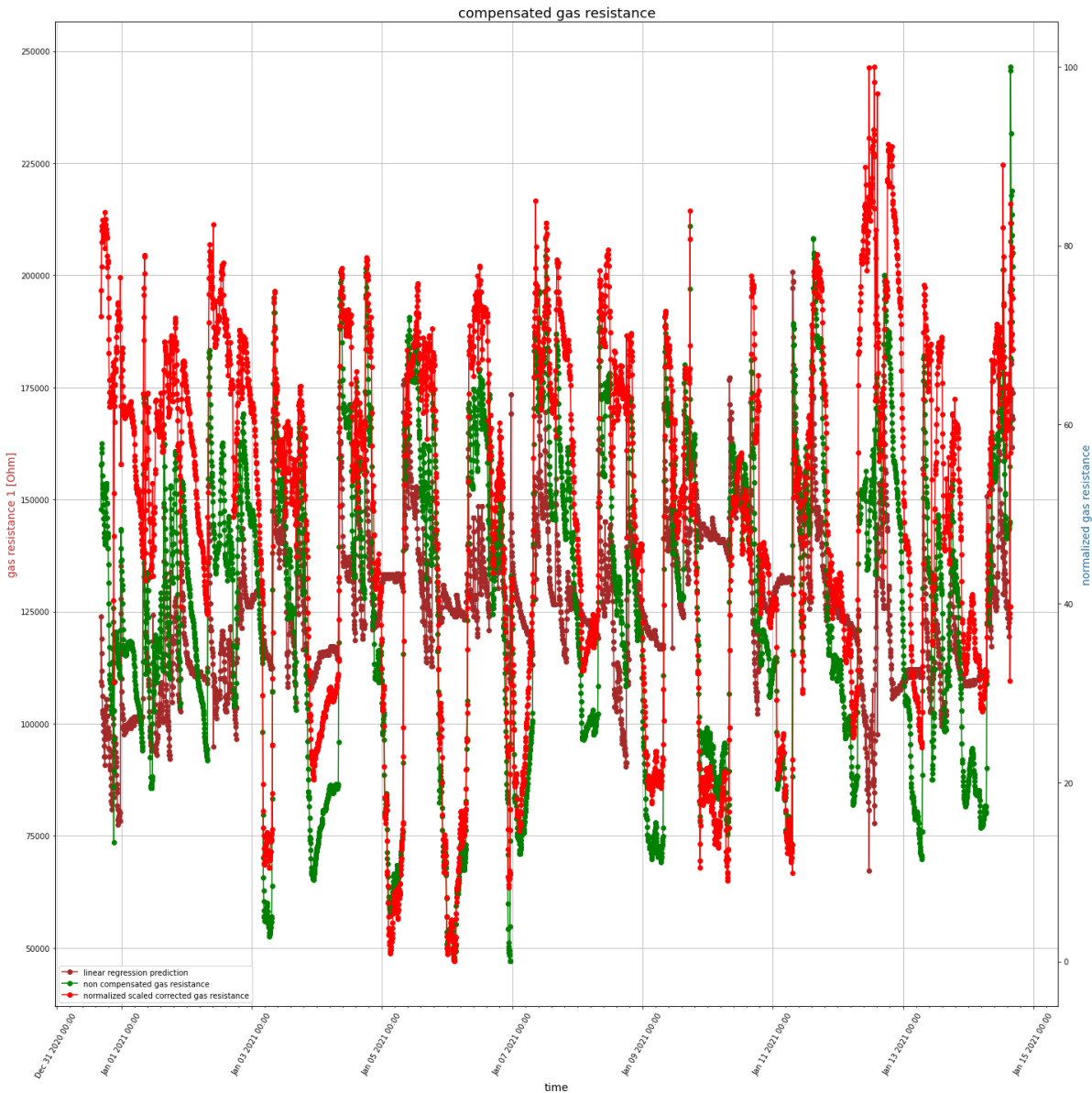
OLS Regression Results

Dep. Variable:	raw_gas_resistance	R-squared:	0.223
Model:	OLS	Adj. R-squared:	0.223
Method:	Least Squares	F-statistic:	671.6
Date:	Sun, 04 Jun 2023	Prob (F-statistic):	4.33e-257

Time:	13:30:56	Log-Likelihood:	-55087.			
No. Observations:	4675	AIC:	1.102e+05			
Df Residuals:	4672	BIC:	1.102e+05			
Df Model:	2					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

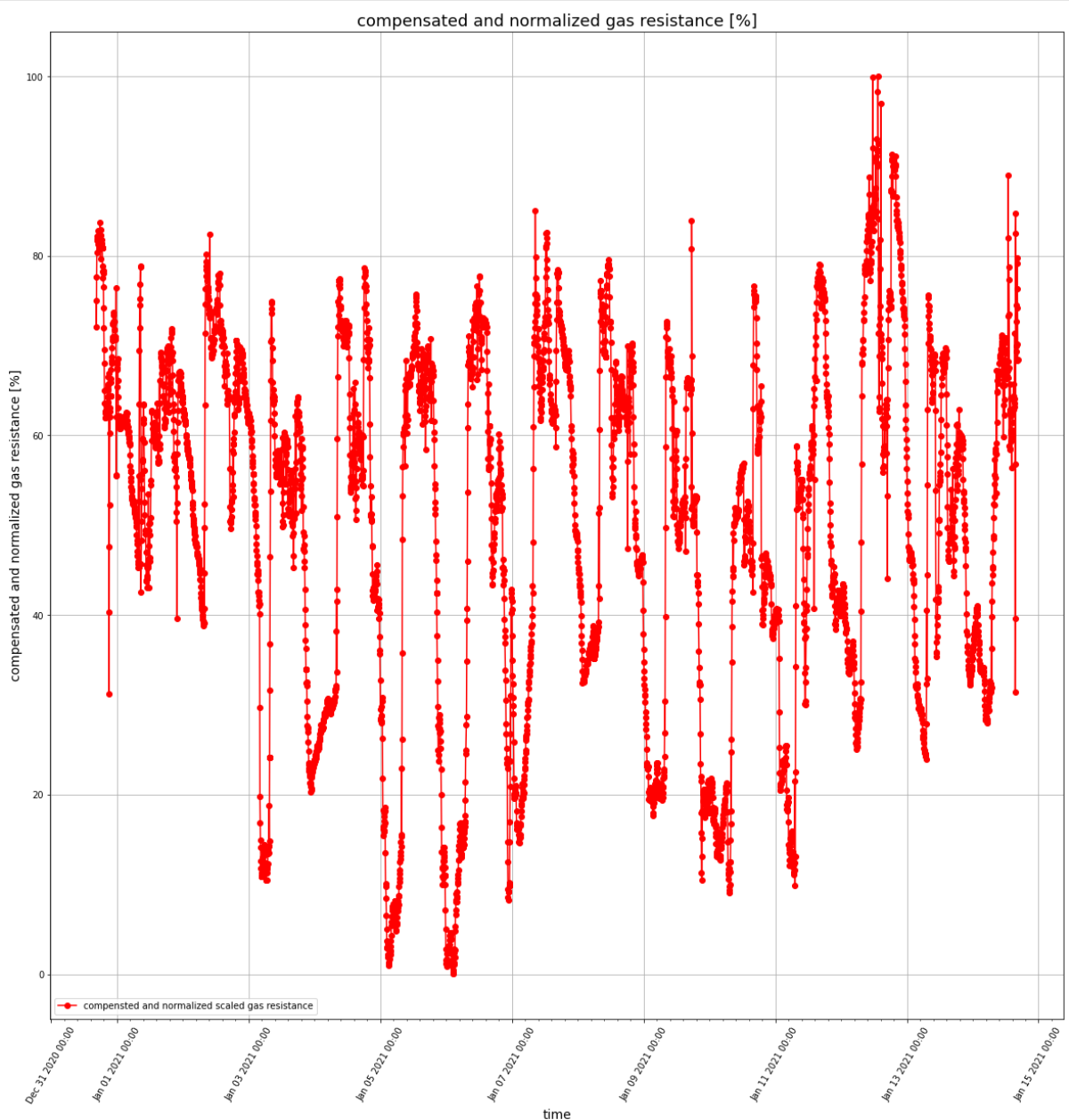
const	2.435e+05	1.11e+04	21.999	0.000	2.22e+05	2.65e+05
temperature	6603.4893	610.281	10.820	0.000	5407.051	7799.928
absolute_humidity	-3.736e+04	1111.429	-33.610	0.000	-3.95e+04	-3.52e+04
=====						
Omnibus:	395.538	Durbin-Watson:	0.015			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	277.025			
Skew:	-0.488	Prob(JB):	6.99e-61			
Kurtosis:	2.315	Cond. No.	594.			
=====						

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
0.2232967344642237



time series diagrams of normalized scaled gas resistance; y range is 0.0..100.0 [%]

```
In [11]: 1 import pandas as pd
2 from sklearn import linear_model
3 import statsmodels.api as sm
4 import matplotlib.pyplot as plt
5 from datetime import datetime
6 from matplotlib.dates import DateFormatter
7 from matplotlib.ticker import MultipleLocator, FormatStrFormatter,
8     AutoMinorLocator
9
10
11 fig, ax = plt.subplots(figsize=(20, 20))
12 plt.xticks(rotation=60)
13 ax.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
14
15 ax.xaxis.set_minor_locator(AutoMinorLocator())
16
17 plt.plot_date(df['Datum'], normalized_residuals, linestyle='solid', color='red', label='comp
18
19 plt.title('compensated and normalized gas resistance [%]', fontsize=18)
20 plt.xlabel('time', fontsize=14)
21 plt.ylabel('compensated and normalized gas resistance [%]', fontsize=14)
22 plt.grid(True)
23 plt.legend(loc="lower left")
24 plt.show()
```

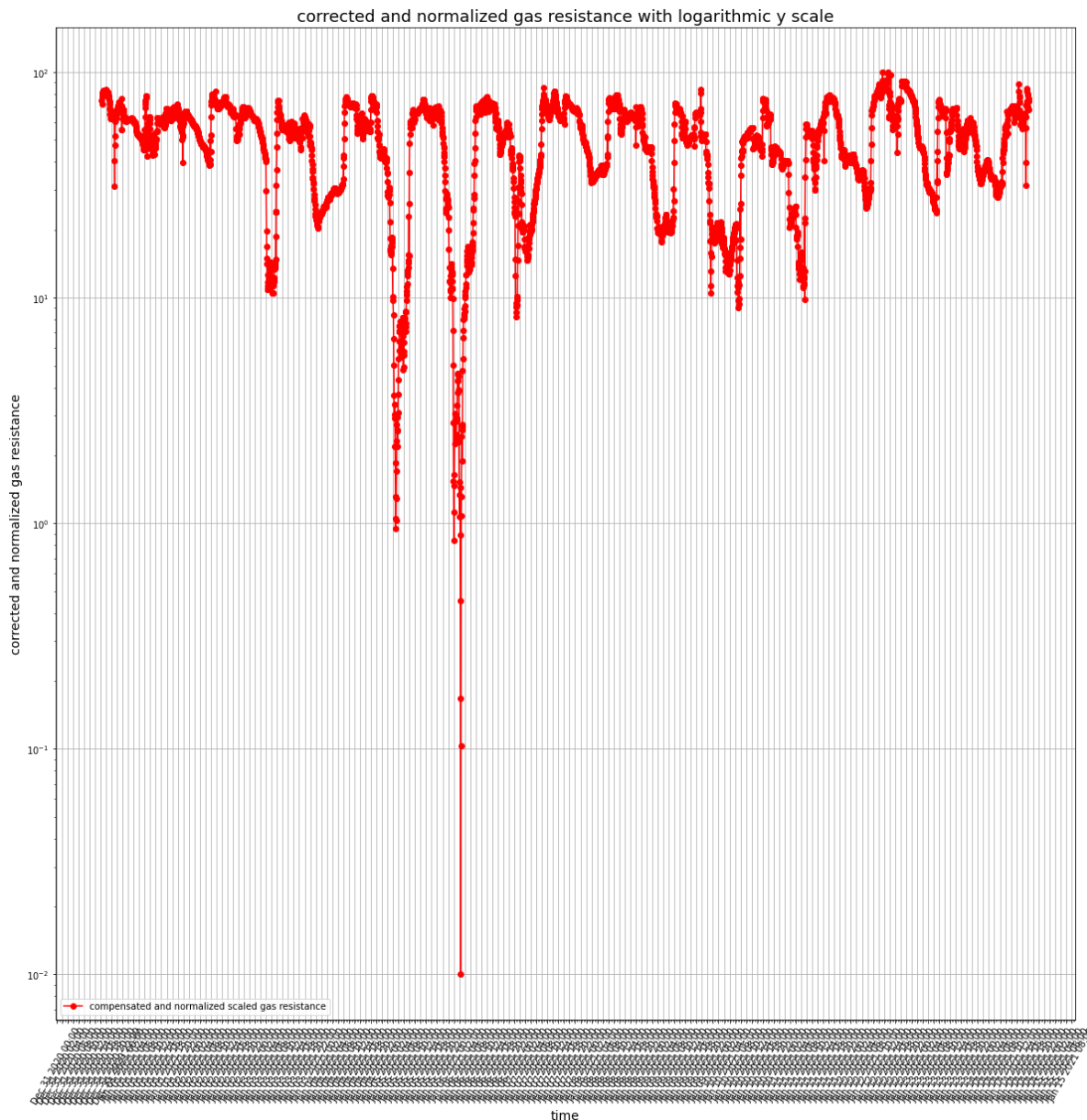


time series diagrams of compensated and normalized scaled gas resistance with logarithmic y scale 0.01 .. 100

```

In [12]: 1 import pandas as pd
2 from sklearn import linear_model
3 import statsmodels.api as sm
4 import matplotlib.pyplot as plt
5 from datetime import datetime
6 from matplotlib.dates import DateFormatter
7 from matplotlib.ticker import (MultipleLocator, FormatStrFormatter,
8                               AutoMinorLocator)
9
10
11 fig, ax = plt.subplots(figsize=(20, 20))
12 plt.yscale('log')
13 plt.xticks(rotation=60)
14 ax.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
15 hours = mdates.HourLocator(interval = 2)
16 ax.xaxis.set_major_locator(hours)
17 ax.xaxis.set_minor_locator(AutoMinorLocator())
18
19
20 plt.xticks(rotation=60)
21 plt.plot_date(df['Datum'], normalized_residuals, linestyle='solid', color='red', label='comp
22
23 plt.title('corrected and normalized gas resistance with logarithmic y scale', fontsize=18)
24 plt.xlabel('time', fontsize=14)
25 plt.ylabel('corrected and normalized gas resistance', fontsize=14)
26 plt.grid(True)
27 plt.legend(loc = "lower left")
28 plt.show()

```



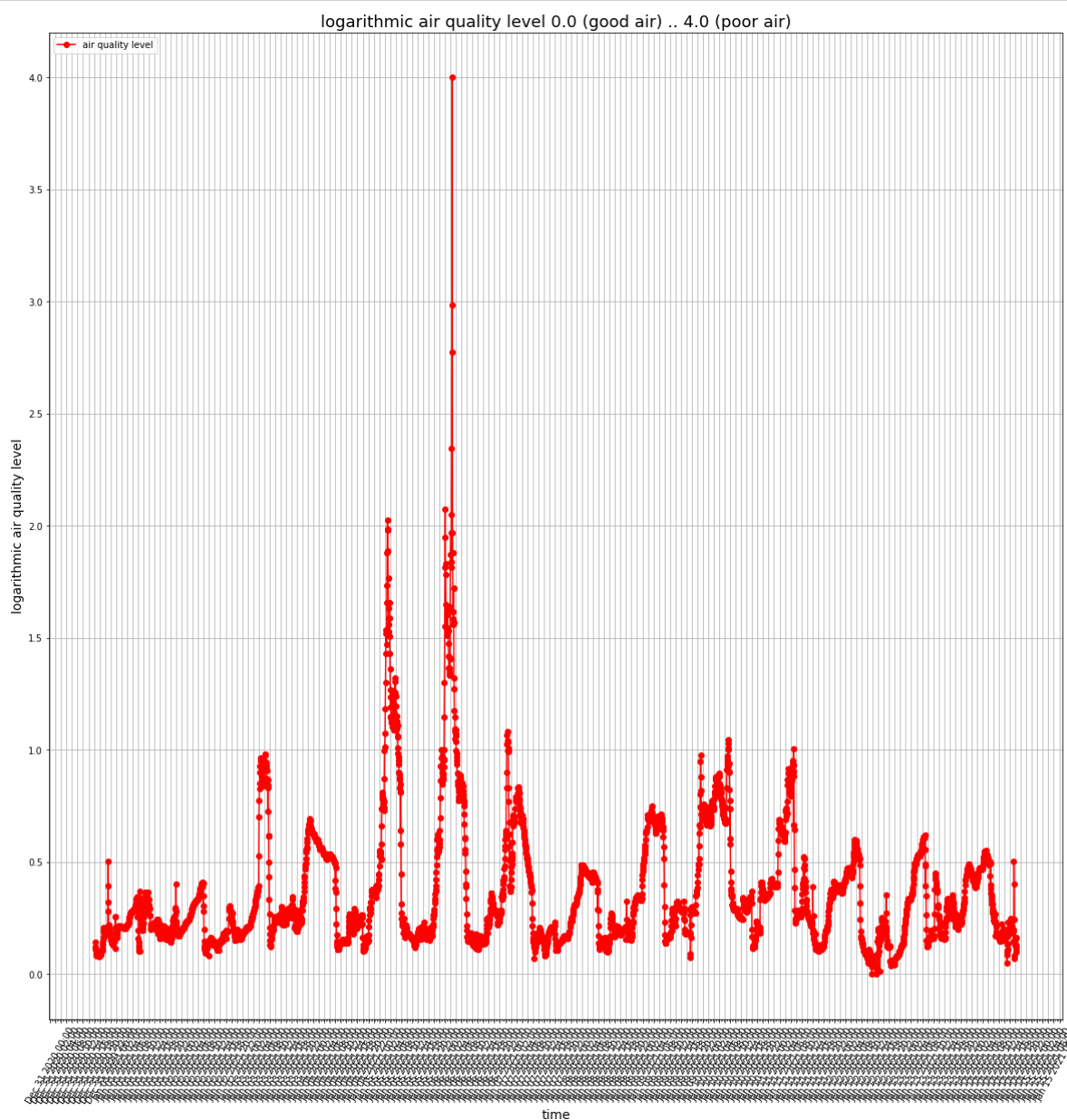
time series diagrams of the logarithmic air quality level the air quality level can vary between 0.0 (fresh air) and 4.0 (very poor air

quality)

```

In [13]: 1 import pandas as pd
2 from sklearn import linear_model
3 import statsmodels.api as sm
4 import matplotlib.pyplot as plt
5 from datetime import datetime
6 from matplotlib.dates import DateFormatter
7 from matplotlib.ticker import (MultipleLocator, FormatStrFormatter,
8                               AutoMinorLocator)
9
10 log_normalized_residuals = -(np.log10(normalized_residuals)-2)
11
12
13 fig, ax = plt.subplots(figsize=(20, 20))
14 plt.xticks(rotation=60)
15 ax.xaxis.set_major_formatter(DateFormatter('%b %d %Y %H:%M'))
16 hours = mdates.HourLocator(interval = 2)
17 ax.xaxis.set_major_locator(hours)
18 ax.xaxis.set_minor_locator(AutoMinorLocator())
19
20
21 plt.xticks(rotation=60)
22 plt.plot_date(df['Datum'], log_normalized_residuals, linestyle='solid', color='red', label='
23
24 plt.title('logarithmic air quality level 0.0 (good air) .. 4.0 (poor air)', fontsize=18)
25 plt.xlabel('time', fontsize=14)
26 plt.ylabel('logarithmic air quality level', fontsize=14)
27 plt.grid(True)
28 plt.legend(loc = "upper left")
29 plt.show()

```



Please check whether the R-squared (uncentered) of the multiple linear regression above is sufficiently good (should be > 0.7): R-squared (also called coefficient of determination) is the portion of variance in the dependent variables that can be explained by the independent variables. Hence, as a rule of thumb for interpreting the strength of a relationship based on its R-squared value is:

- if R-squared value < 0.3 this value is generally considered as None or very weak effect size
- if R-squared value $0.3 < r < 0.5$ this value is generally considered as weak or low effect size
- if R-squared value $0.5 < r < 0.7$ this value is generally considered as moderate effect size
- if R-squared value $0.7 < r < 1.0$ this value is generally considered as strong effect size

If R-squared value is < 0.3 , the collected history may be too short. Please try to collect datapoints for a longer timeframe!

```
In [14]: 1 print("\nR-squared (uncentered) of the multiple linear regression          = %11.2lf\n\n" % n
```

```
R-squared (uncentered) of the multiple linear regression          =          0.22
```

Please enter the following parameters of the multilinear regression into the Homematic/RaspberryMatic WebUI page 'Startseite > Einstellungen > Geräte > Geräte-/ Kanalparameter einstellen' of your concerning BME680 AQ sensor device which was the source of the history.csv file above

```
In [15]: 1 print("\nNumber of captured data points used for the MLR                      = %11d" % len(>
2
3 print("\nPlease enter the WebUI device parameter 'WEATHER|mlr_alpha'              = %11.3lf" % re
4 print("Please enter the WebUI device parameter 'WEATHER|mlr_beta'                = %11.3lf" % reg
5 print("Please enter the WebUI device parameter 'WEATHER|mlr_delta'                = %11.3lf" % reg
6
7 import datetime
8 now = datetime.datetime.now()
9 print("\n\nPlease check whether current date and time are correct: ")
10 print(str(now))
11
12
```

```
Number of captured data points used for the MLR                      =          4675
```

```
Please enter the WebUI device parameter 'WEATHER|mlr_alpha'          =        6603.489
Please enter the WebUI device parameter 'WEATHER|mlr_beta'            =       -37355.389
Please enter the WebUI device parameter 'WEATHER|mlr_delta'           =       243521.114
```

```
Please check whether current date and time are correct:
2023-06-04 13:31:02.235573
```

Congratulations, you are done!

Please repeat the multilinear regression update of the WebUI device parameters on a regular basis every month or similar as appropriate ..