

# [EDA] Exploratory Data Analysis

Note: this report summarizes the output that can be found on the Jupyter Notebook (EDA/EDA\_ExploratoryDataAnalysis.ipynb) which also includes the source code in-line (html can be consulted). This report has been created with the sole purpose of meeting the requirement of 5 pages for the visual report.

The report is an analysis of the file `dataset_processed.csv` which has been created by performing an ETL process (Extract Transform Load) of the original CSV file, running 10 minutes averages, on the cloud.

## Missing Data Analysis

First things first, ETL (Extract Transform Load) shows that the dataset contains gaps (identified by NaN values when pivoting on the output format). The following gaps are identified. The analysis also reveals that the gaps are per asset (all the variables for that asset are missing for those gaps).

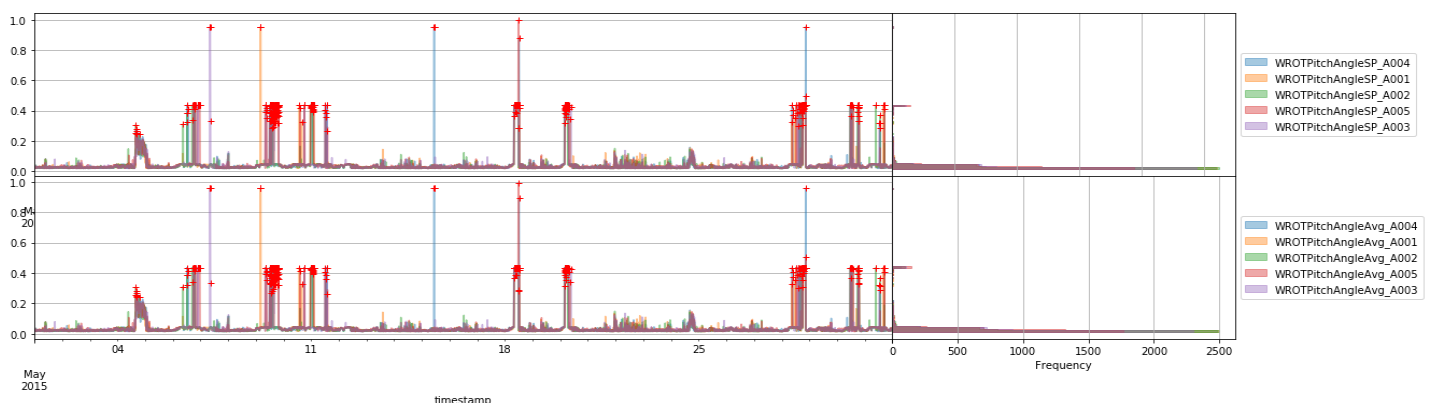
start	stop	duration	asset
2015-05-06 06:20:00	2015-05-06 13:30:00	0 days 07:10:00	A005
2015-05-07 06:20:00	2015-05-07 13:20:00	0 days 07:00:00	A005
2015-05-18 09:40:00	2015-05-18 12:20:00	0 days 02:40:00	A004
2015-05-21 11:10:00	2015-05-21 11:40:00	0 days 00:30:00	A002
2015-05-21 11:40:00	2015-05-21 12:10:00	0 days 00:30:00	A004
2015-05-25 09:50:00	2015-05-25 13:20:00	0 days 03:30:00	A003
2015-05-25 11:40:00	2015-05-25 12:00:00	0 days 00:20:00	A005
2015-05-26 09:00:00	2015-05-26 11:10:00	0 days 02:10:00	A001
2015-05-27 11:10:00	2015-05-27 11:30:00	0 days 00:20:00	A003

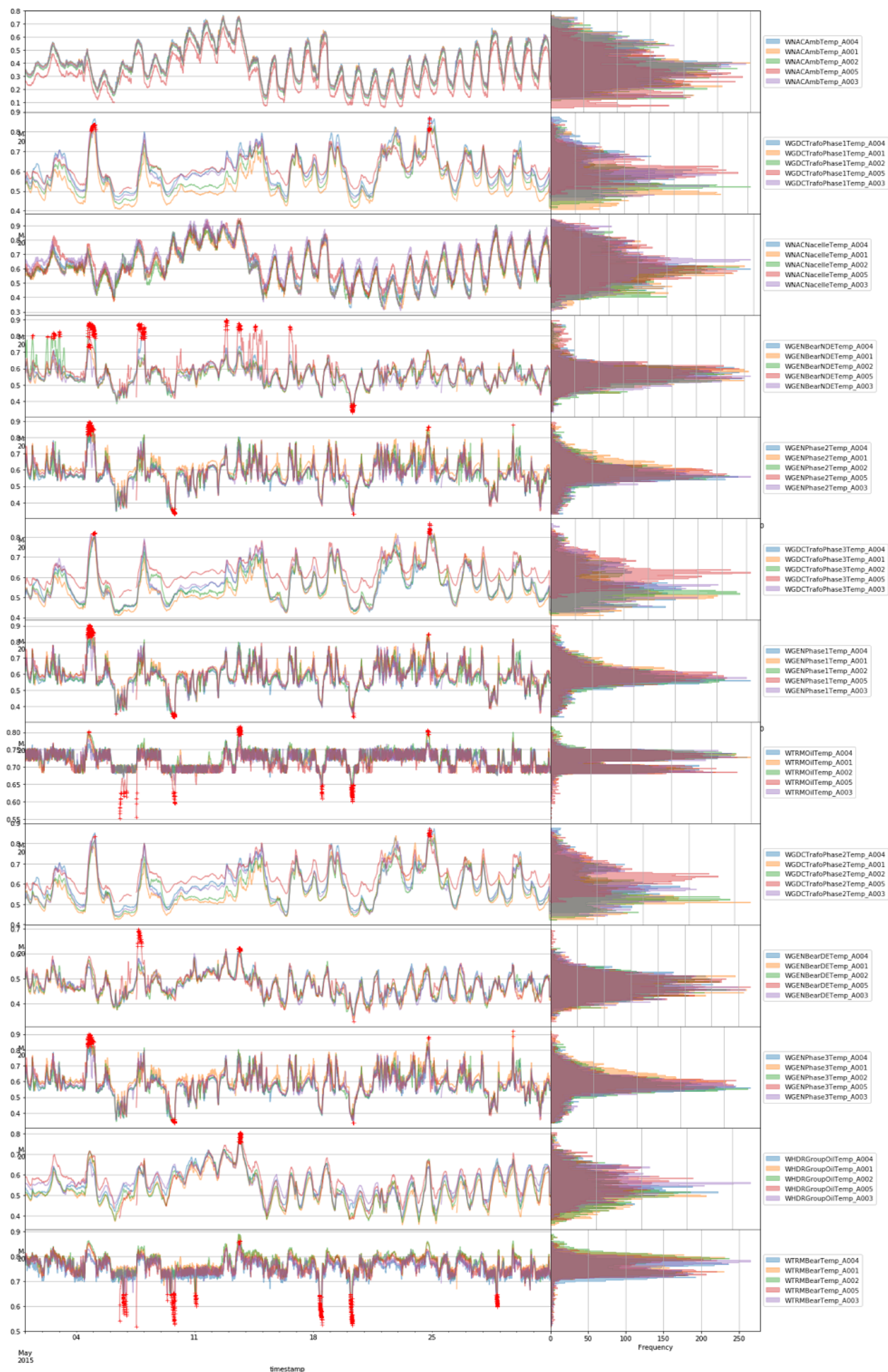
## Time series analysis

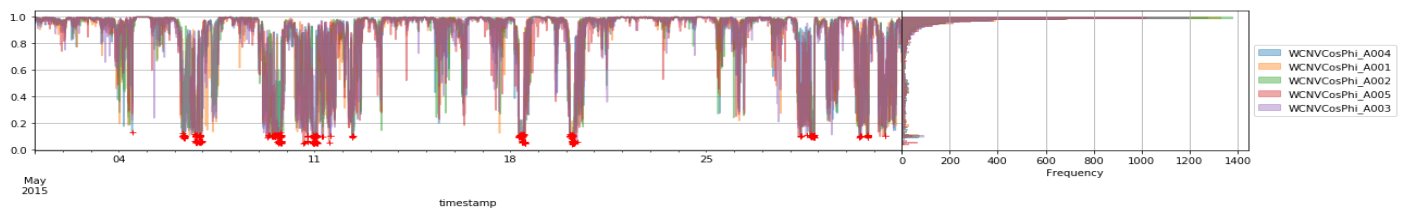
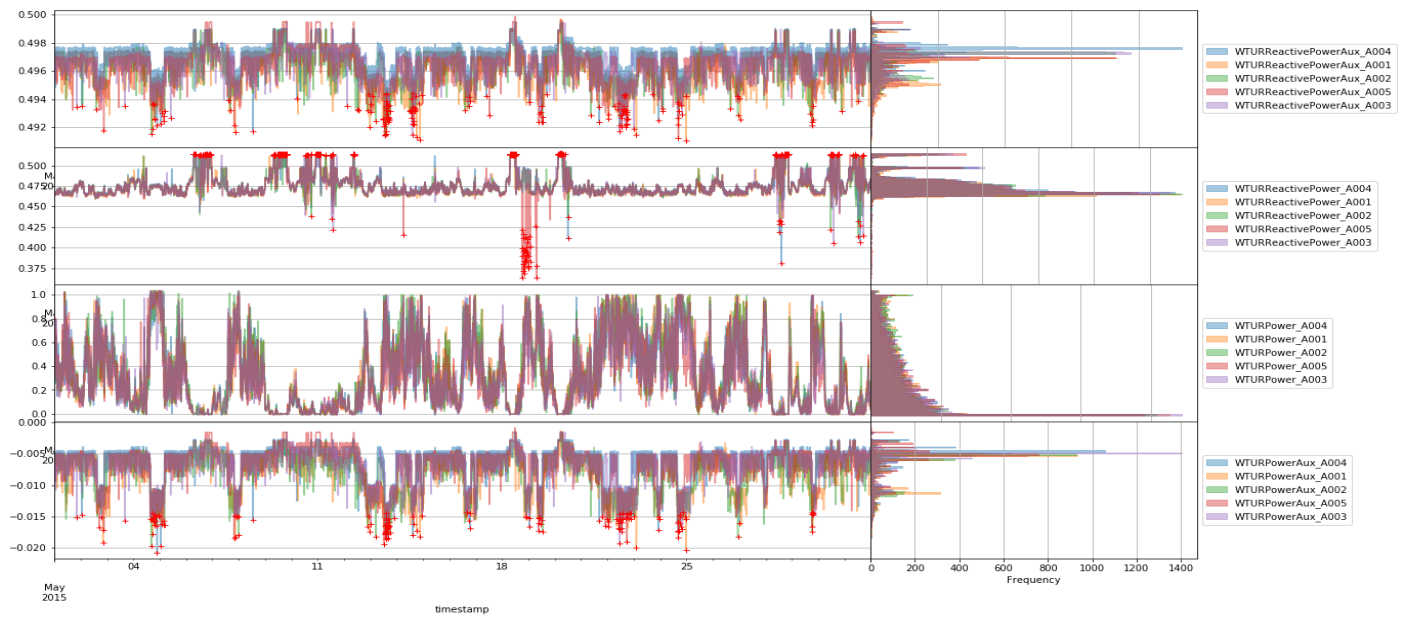
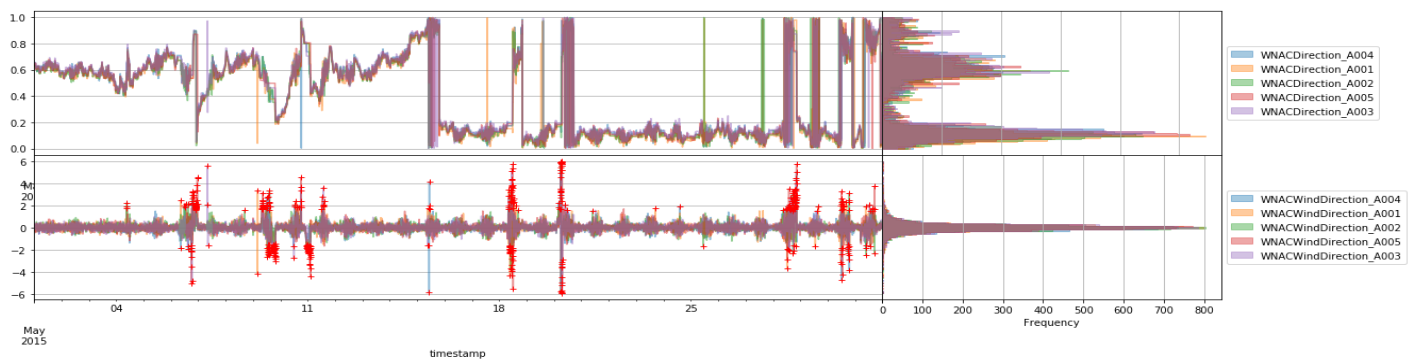
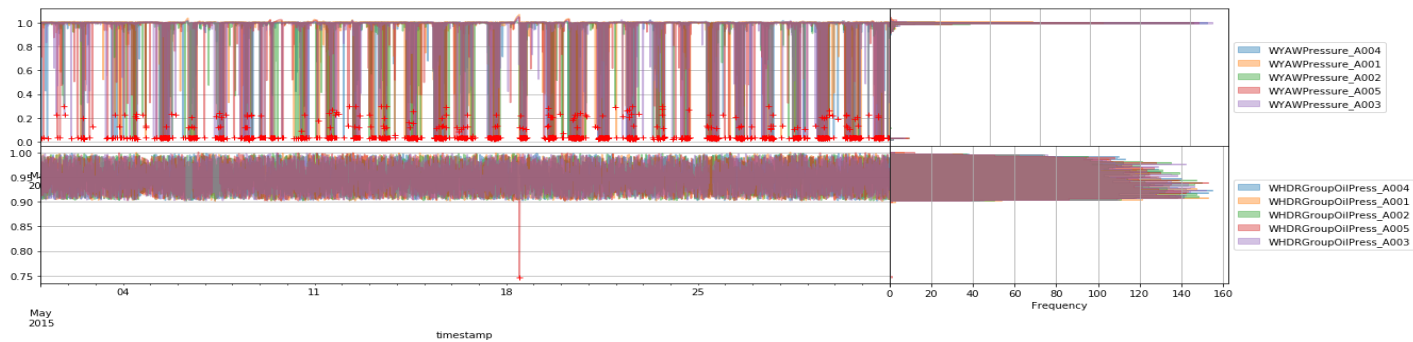
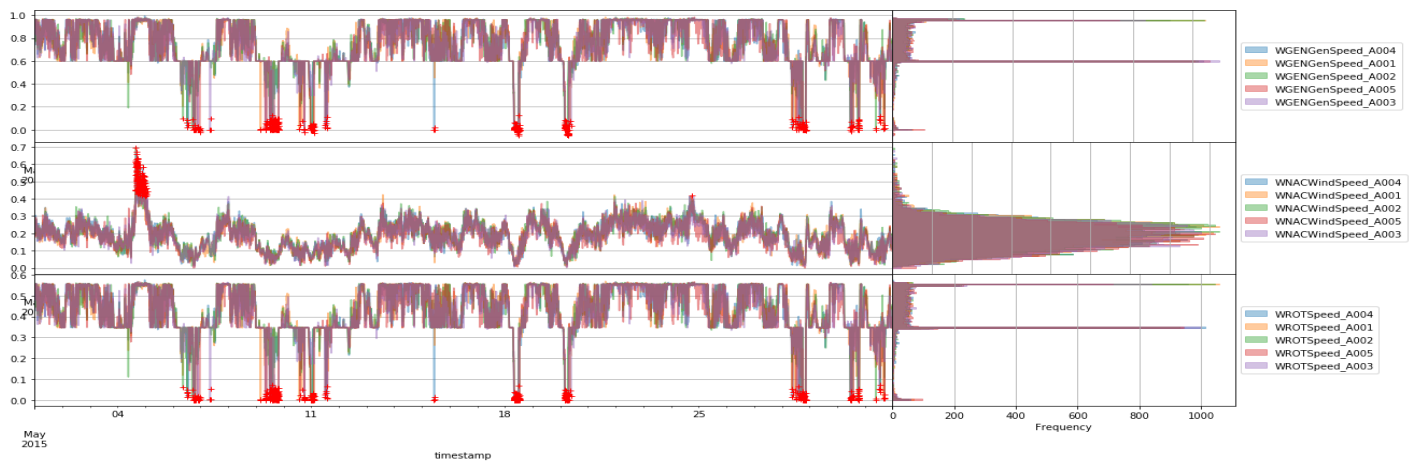
The data contains signals which come from sensors and therefore represent a physical quantity (temperature, speed, pressure). The statistical nature of these signals is to be analyzed, but a very good first approximation is to find outliers defined in a statistical sense as those points that are more than 3-sigmas away from the mean. The outliers are highlighted in red in the signal.

$$o_x(n) = 1 \leftarrow (x(n) > \mu_x + 3\sigma_x) \cup (x(n) < \mu_x - 3\sigma_x)$$

Given the high level of correlation among the same variable of different assets the following plot shows the 28 variables, where each color is an asset. The histogram is plotted on the side and the outliers are highlighted in red. The plots are grouped per variable nature (e.g. temperature, speed, angle, power)



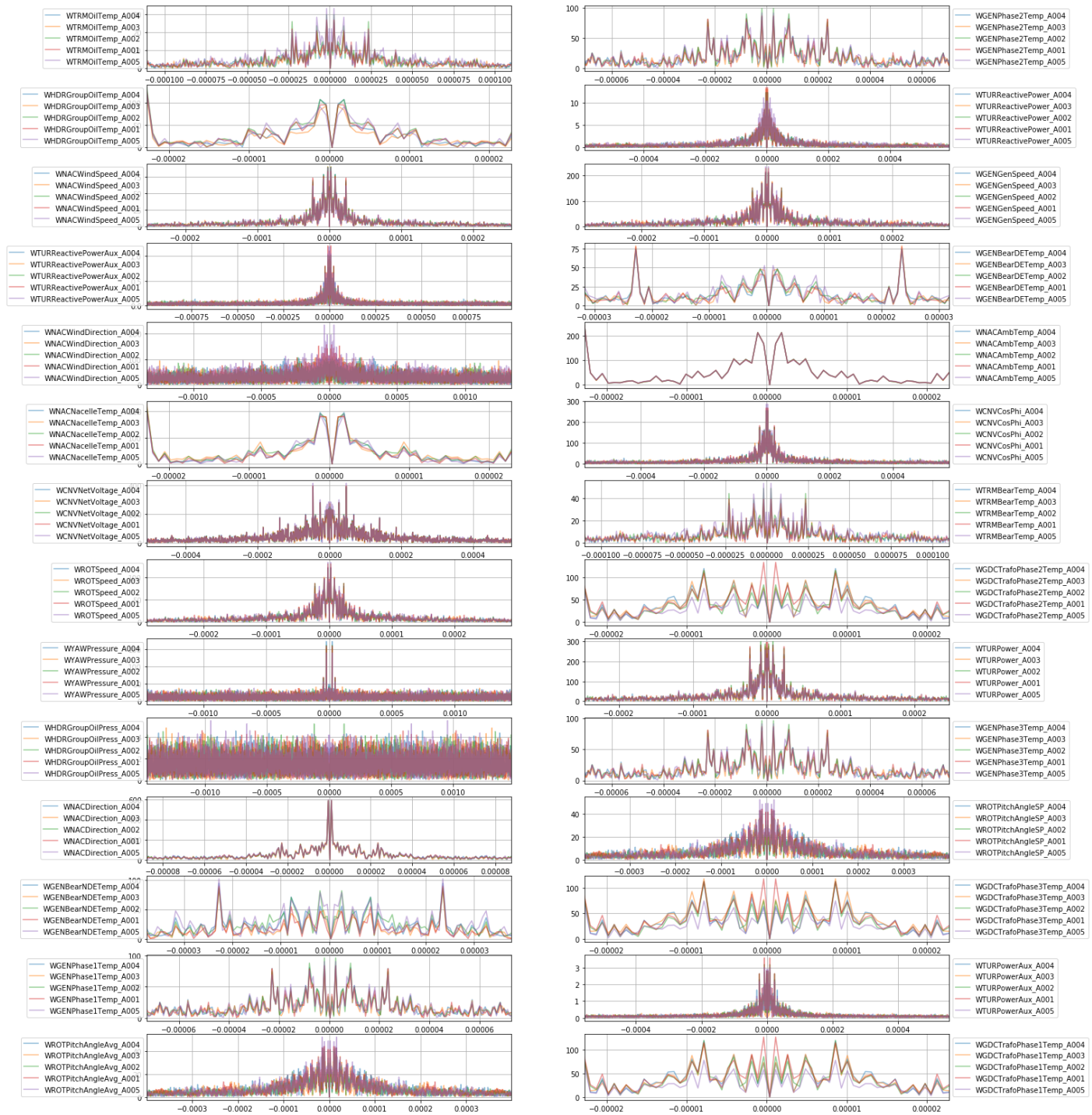






# Spectral analysis

These time signals are continuous and can also be inspected from the spectral point of view. Since almost all signals contain a bias, the best visualization is given by removing the DC component and adjusting the shown bandwidth so 95% of the energy of the signal is contained (otherwise most plots would be seen as a dirac delta at DC).

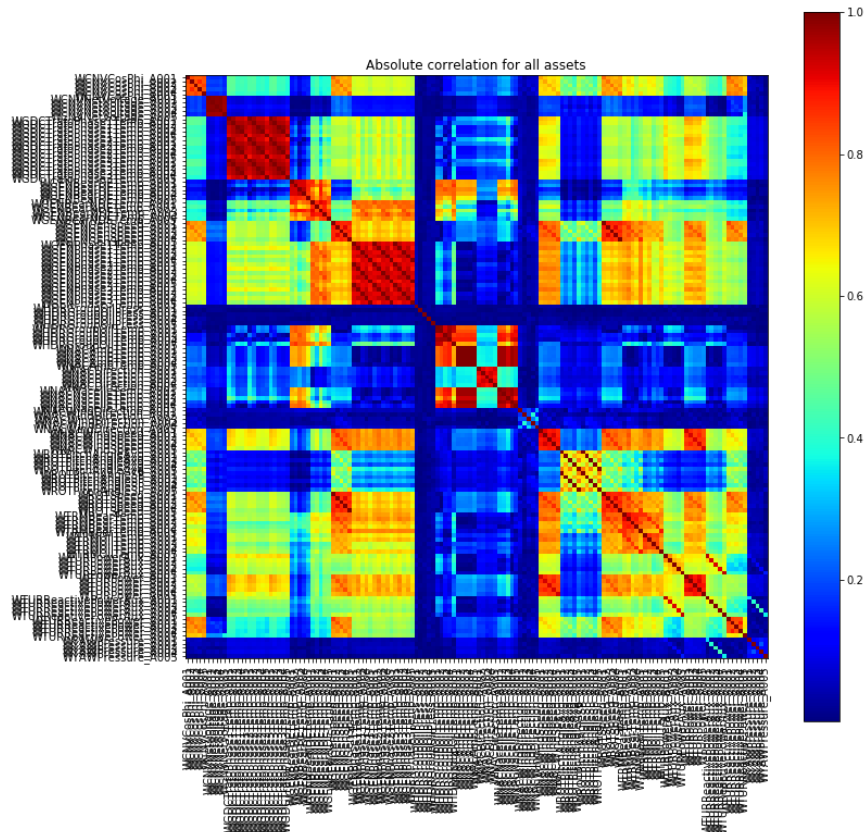


# Correlation analysis

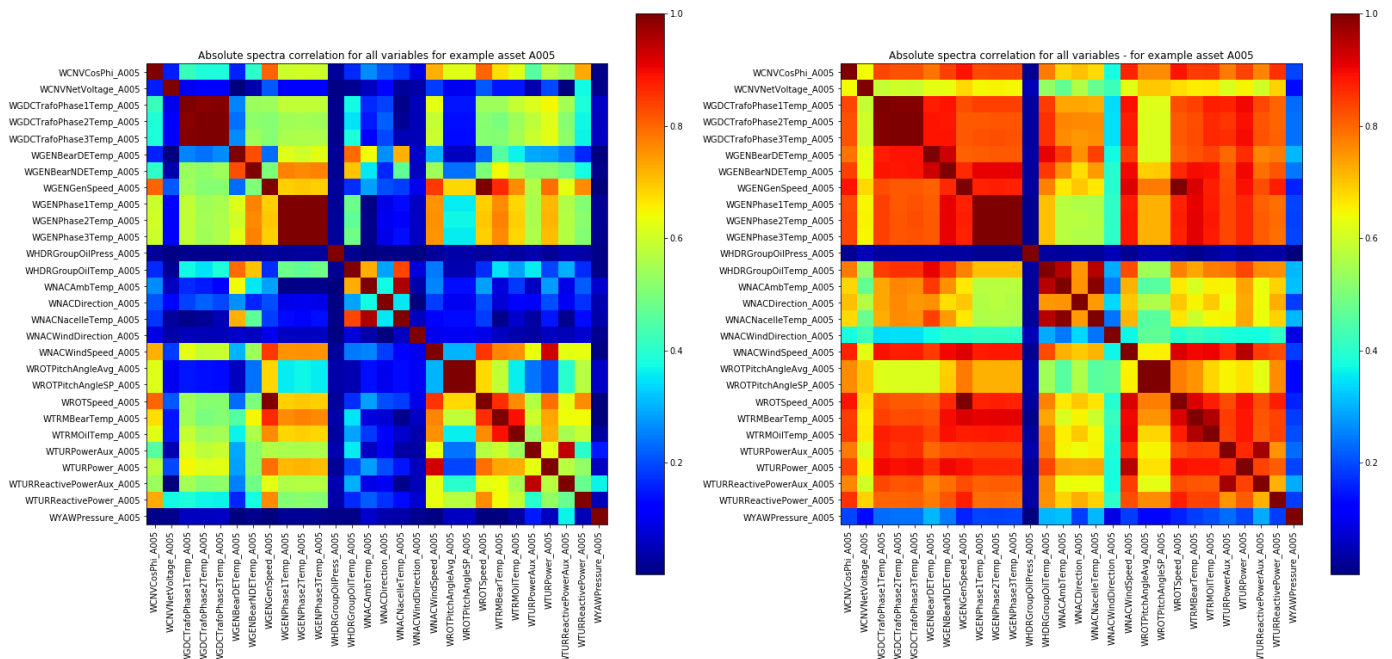
There is a high level of correlation among the signals. There are three types of correlations observed: correlation among assets, correlation among variables within an asset in the time domain, correlation among variables within an asset in the spectral domain.

## Correlation among assets

The following plot represents the correlation matrix of all variables for all assets (28\*5=140). The hypothesis to confirm is that the data is highly correlated (i.e. correlation shows a clustering of 5 elements around the same variable), which does.



## Correlation matrix among variables on the same asset on time and spectral domain



One can observe high correlation among temperatures, angles, power and reactive power, wind speed and power, speed and temperature, phase1 temp and speed. Others like pressure seem to be totally decorrelated to the rest.

### AOB:

An analysis on signal clustering has also been carried out on the notebook with inconclusive results.