# Seismic Time Series Analysis
# Statistics for Spatio-Temporal Data

Marco Ciotola, 848222

June 12, 2020

# Contents

# 1 Introduction

**Motivation and objective**

Thanks to the recent COVID-19 outbreak worldwide, we currently have data about many different phoenomena from periods with standard human activity, and from a period with reduced human activity.

Inspired by an article published during the initial phases of the Belgium-wide lockdown (1), we want to assess whether the lockdown period can be correlated with a visible change in background seismic activity.

It is important to notice that some industrial events, such as mine explosions, are already mapped by seismologic institutions since they interfere with seismic surveys (2). Such events are sometimes comparable with natural seismic events, depending on the distance from the seismograph.

The correct estimation of a background noise, correlated or not with human activity, could help institutions to better isolate less prominent earthquakes and similar events from the seismic data. This could result in better understanding of the frequency of earthquakes by studying their more precise evolution over time.

**Data source and description**

The Royal Observatory of Belgium (ROB) gathers and publishes (3, 4) data about vertical ground displacement from 9 stations through the Country, linking them to seismic events whenever possible (5). They also offer a data visualization tool for the same data (6).

The vertical ground displacement is represented by the minimum and maximum displacements of the ground in the vertical axis registered in a second, with *nm* as unit of measure. This can easily be transformed to total ground displacement for each second by taking their difference. There are 86400 measurements per day, around 30 million per year.

Out of the 9 stations, the following 6 have interesting positions, related to possible correlations with human activity:

- **Uccle**, **Sart-Tilman**, and **Ében-Émael** are near a city (Brussels, Liege, and Maastricht).

- **Ostenda** is near the city of Ostenda, that borders with the sea.

- **Membach** and **Dourbes** are in towns near a natural park.

Data collected from near a city and from near a natural park could be selected, in order to compare possible seasonality results and differences from before and after the lockdown. Considering these factors, we selected the **Uccle** and **Membach** stations.

It is important to notice that the daily data in each file refers to the UTC time instead of the local Europe/Brussels time, an optimal choice to publish data in a machine-readable way. This causes an issue fitting a fixed seasonality, when it depends on the legal time rather than on the absolute time.

## 1.1 Data cleanup and transformation

**Data gathering and cleanup**

The data is provided in CSV format, one file per day, with only two columns related to the min-max pair (4) and no column related to time. This causes two problems:

1. For merging days, we need to add a column with the second-precise datetime

2. In days with missing values there are placeholder values instead of a missing row, that would make the following rows misinterpreted for measurements at the wrong second. The placeholder values are one of these two min-max pairs: $(0,0)$, $(-1,1)$.

Data from each station have different domains, as distinct equipment can register a specific amplitude of the seismic movements. In fact, we found some values[1] outside the domains, and we consider them as missing values.

In particular, we find that out of $\approx$ 50 million seconds, from 27-10-2018 to 28-05-2020, only 0.67% (Uccle) and 0.37% (Membach) are missing values. We will see later how they are distributed for each station.

**Data aggregation and initial transformation**

Since we are not interested in a second-precise analysis, and it would be computationally expensive to deal with the whole dataset, we aggregate the data by hour. The aggregation results in around 14 thousand points, that can still be captured by interpretable models.

While performing this operation, we should consider the correct time-zone: as we anticipated, any seasonal analysis dealing with UTC dates would be out of phase for the half-year of legal time, with respect to the half-year of solar time. Employing the `lubridate` R package we converted times to local time-zones, and then we kept them as absolute instants, because all the R functions we use work on the UTC conversion of datetimes by default. This transformation made it possible to have each hour in the timeseries to represent the local time in Belgium (*Europe/Brussels*), without being influenced by the legal time.

The aggregation was performed considering the mean total ground displacement at each hour, for the following reasons:

- When an hour contains some missing seconds, these do not heavily affect the aggregation result, in contrast with summing the values

- Seismic events, such as quakes and mine explosions, normally last some minutes at most, and will affect the result only slightly, as opposed to the hourly maximum

- Usually, maximum value aggregation does not follow the Normal distribution, while some models that we are going to use assume the data follow a Normal distribution

# 2 Analysis and decomposition

**Missing values analysis**

After aggregating by hour, the percentages of missing values remain quite stable. As can be seen from Figures 1a and 1b, missing values at both stations are grouped in a few consecutive periods. To deal with them when R methods need complete timeseries, we define different completion strategies, depending on the seasonality, using the R function `na.aggregate`.

**Seasonality assumptions**

The already cited article (1) suggests the existence of a weekly seasonality. Looking at the data it seems confirmed, in addition to a daily seasonality (Figure 2).

Since one of the main focuses is on the seasonalities themselves, we aim to avoid searching for correlation with human activity while imposing human-based seasonality periods.

---

[1]Values out of domain for each station are less than 10 out of $\approx$ 50 million.
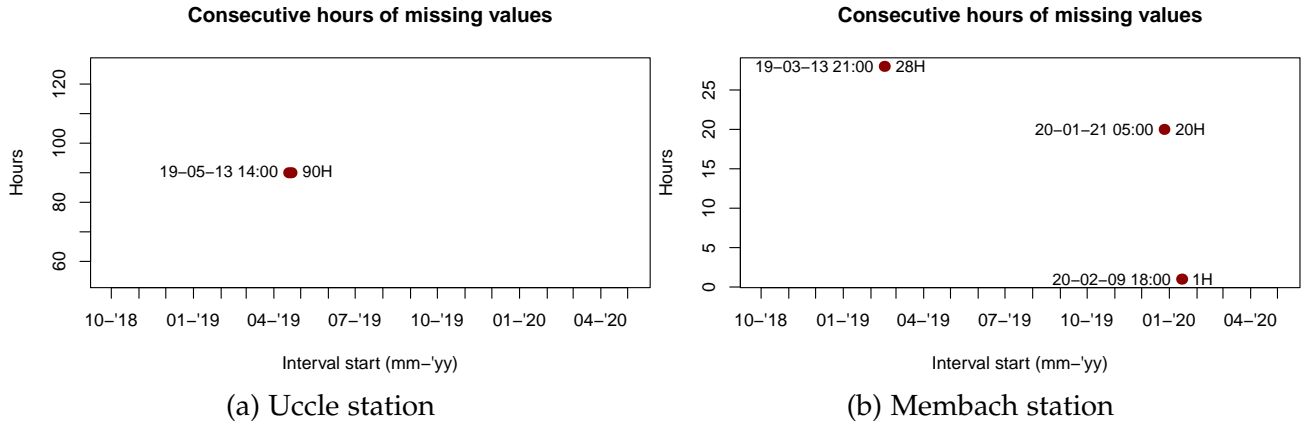
(a) Uccle station



(b) Membach station

Figure 1: Missing values representation

Thus, we use an analytical method to extract seasonalities that the data itself suggests to us. Still, we recognize that the already performed transformations impose an human-based time alteration to the data itself (Section 1.1).

The method we are going to apply is based on the Fourier Transform (FT), and generates a *periodgram* highlighting the frequencies that might correspond to a seasonality component (7, chapter 11). Applying the FT for any frequency, we transform our 1-dimensional data into 2-dimensional data mapped around a circle. Then, we can sum all the resulting vectors (that are our new data points) to obtain a vector for each frequency, whose absolute value represents the power level of the same frequency. An high power level represents a significant periodicity on that frequency in our original data, since it means that the corresponding 2-dimensional data behaves in a highly similar way during each cycle.
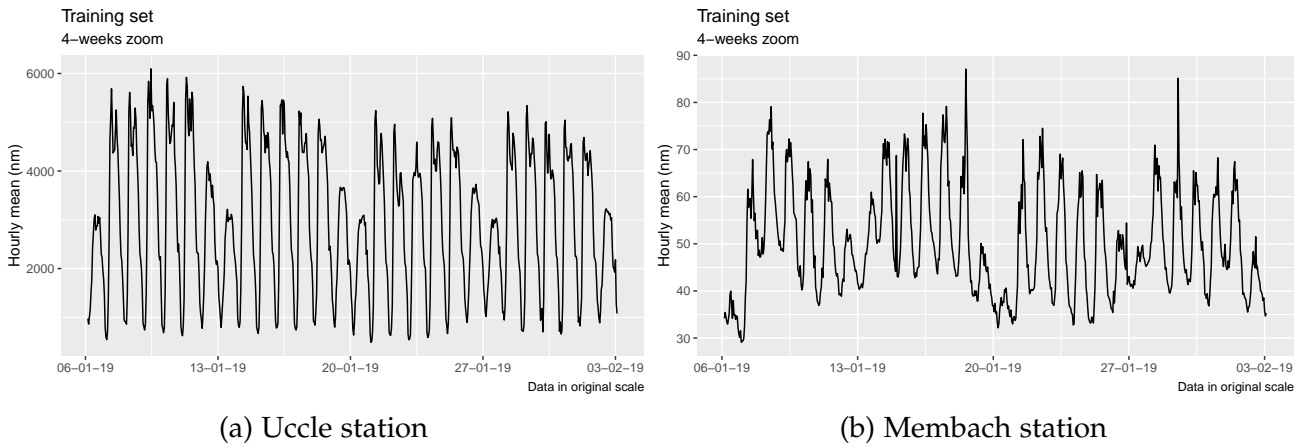


(a) Uccle station



(b) Membach station

Figure 2: Zoom from January $6^{th}$ to February $3^{rd}$ 2019

A common operation to increase the difference between significant and non-significant periodicities is to apply a Moving Average (MA) smoothing filter with a smaller window than the interesting ones. We will use such a filter with $p = 2$.

On this note we must specify that, when a significant periodicity is found, it influences smaller periodicities given by $f * \{2, 3, \dots\}$ the significant frequency (7, chapter 11). Also, if any periodicity is found that is longer than the original dataset, it must be considered as simple chance, rather than serendipity.

Since we apply this method on a dataset indexed by seconds, the transformation of the frequency $f$ into hourly periodicity $p$ is obtained by $p = (1/f)/3600$. We can analyse the

resulting periodgrams for Uccle and Membach stations in Figures 3a and 3b. An higher power value for a frequency means that the corresponding periodicity is meaningful for describing the dataset.



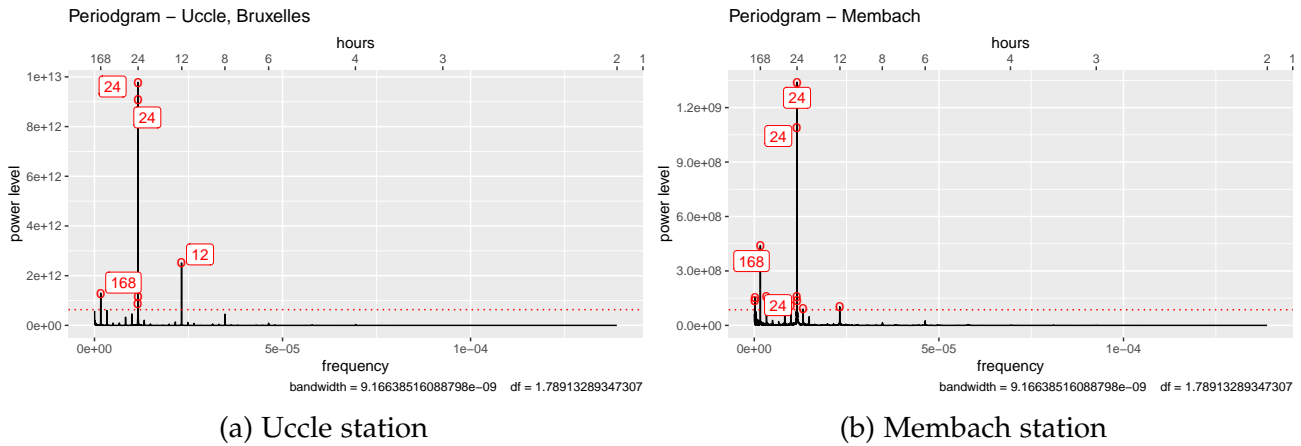(a) Uccle station

(b) Membach station

Figure 3: Periodgram plots

Both plots highlight multiple 24-hour periodicities. This is given by approximation factors, since we round the divisions result to the nearest natural number. It also highlights that using a single 24-hour seasonality could leave some information on the residuals. Still, in both stations it is the most significant periodicity.

Both stations highlight a significant 168-hour (1-week) periodicity, that also confirms our initial assumptions.

Both stations show a possibly significant 12-hour seasonality, but since it corresponds to double the frequency of the 24-hour seasonality, we can safely ignore it.

**Dataset division**

As best practices suggest, we divide the dataset into training and testing datasets with a ratio around 85%/15% − 75%/25%. Since one of our objectives is to assess any significant difference from before and after the Belgium lockdown (on March $14^{th}$ the *soft* lockdown, while on March $18^{th}$ the complete lockdown), we selected data from the start of our time-series until 20-01-2020 (64 weeks ≈ 82%) as training, while the test follows ending on 27-04-2020 (14 weeks ≈ 18%).

Having the test set span from before to after the lockdown enables us to assess significant changes in the forecast performance on the two different periods.

This initial division demostrated to be computationally heavy while fitting some models with weekly seasonality, so that it would be challenging to fit more than a couple of models. Given this issue, we choose the training/testing sets to hold 52/18 weeks: this is the greater train-to-test ratio we deem usable (≈ 75%/25%), while maintaining computational feasibility. The assessment of the lockdown influence will be discussed in Section 3.3.

Some R functions cannot handle timeseries with missing values, but as we have seen both stations' datasets include missing values. Also, please note that we are going to have a timeseries object indexed by weeks, so that it contains the 168 hours between each integer index. Since we know that the longest missing values series is 90 hours long, we can apply a simple completion strategy for serving those commands a filled timeseries: we will substitute missing values with the mean of the week they are in; in R, `na.aggregate(x, floor)`.

**Decomposition**

In order to decompose our dataset into trend and seasonality(-ties) components, we followed these steps:

1. Apply a smoothing filter to extract the trend component

2. Remove the trend from the dataset and extract the seasonal figure using the local mean method

3. Remove the trend and the seasonal components to obtain the residuals

While doing so, we tested different smoothing filters: simple filter (p=seasonality), Spencer's 15-point filter and MA smoothing filter (p=seasonality). Since they do not produce particularly different seasonalities, we choose to keep the results obtained by applying the MA smoothing filter, for its significance in terms of seasonally adjusted data.
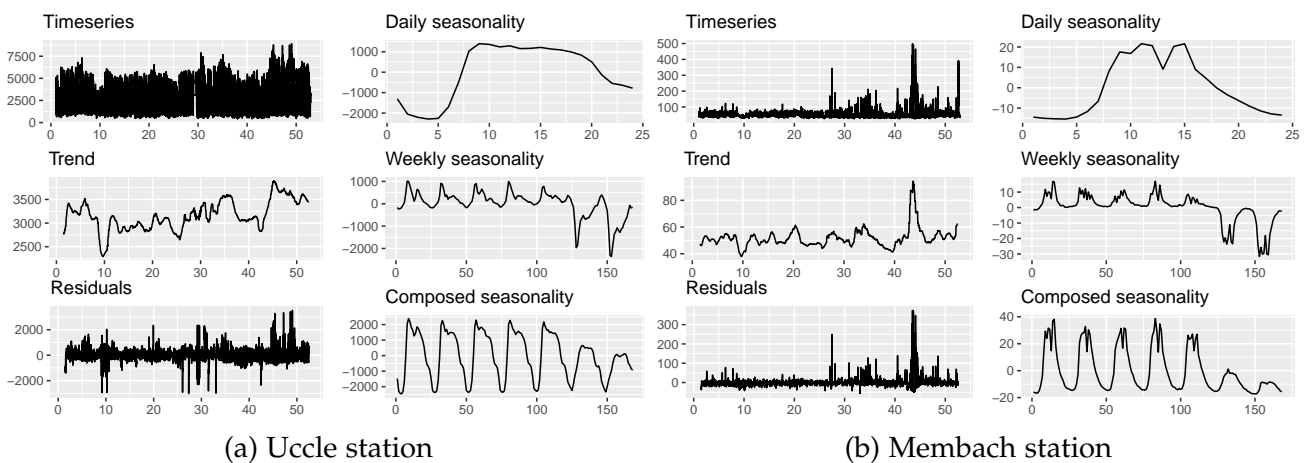


(a) Uccle station                    (b) Membach station

Figure 4: Timeseries decomposition
Only the second method results are shown for clarity. Week 10 includes January $1^{st}$.
x-axes on left columns represent weeks, on right columns represent hours.

Since the periodgrams suggested the presence of two significative periodicities, we also tested the extraction of both corresponding seasonalities (Figure 4):

1. Apply the MA smoothing filter (p=24) on the dataset to extract the trend component

2. Remove the trend from the dataset and extract the 24-hour seasonal figure using the local mean method

3. Remove the 24-hour seasonal components to obtain a deseasoned dataset

4. Apply the MA smoothing filter (p=168) on the deseasoned dataset to extract the updated trend component

5. Remove the trend from the already deseasoned dataset and extract the 168-hour seasonal figure using the local mean method

6. Repeat steps 1-5 using the last deseasoned dataset until convergence (2-3 cycles)

7. Removing the last fitted trend and the two seasonalities we obtain the final residuals

In Figure 5 we see both approaches applied to both stations. Those plots show interesting properties of our seasonalities, probably given by the method (local mean) and the fact that the weekly period is an exact multiple of the daily period:

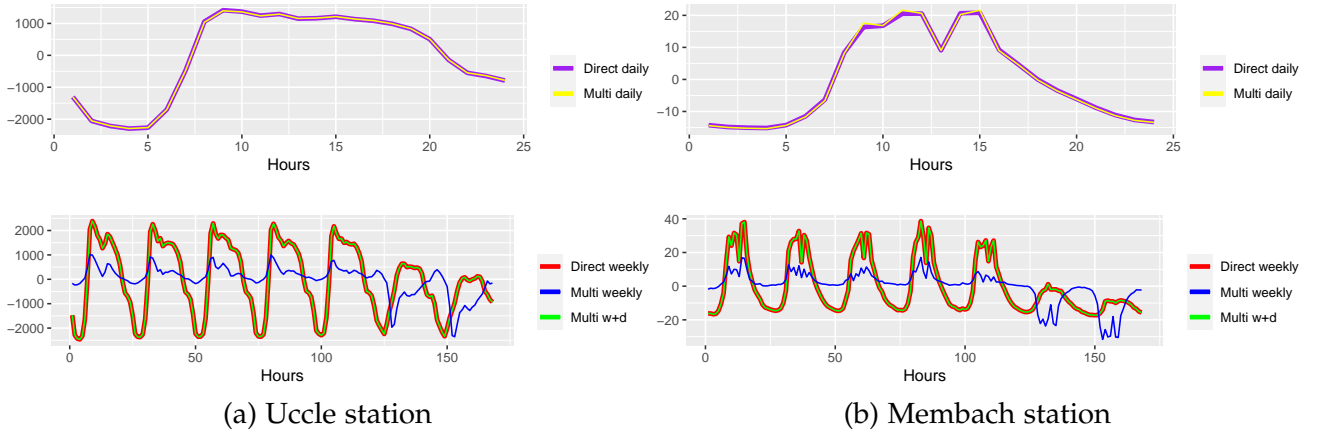• Both ways of extraction result in the same daily seasonality extracted

(a) Uccle station          (b) Membach station

Figure 5: Seasonality figures
*Direct* refers to single seasonality extraction
*Multi* refers to both seasonalities extracted at the same time

- The sum of the daily and weekly seasonalities extracted together is the same as the weekly seasonality extracted directly

Given in particular the second property, we will only use the 168-hours seasonality while fitting our models, for both the stations.

In Figure 4 we see the double decomposition results. While data from Uccle station have a clear significativity of both trend and seasonalities, data from Membach could seem not significant. Still, if we notice that the main timeseries lies between values of 0 and 100, we can find out that around half of the value will be composed by the trend, a 30% by the joined seasonality components, and the remaining value by the residuals. We will need to apply some transformation to flatten the spikes that start from the middle of that dataset.

From the same decomposition results, we notice that data around January $1^{st}$ show a significant drop in values. That is highlighted by both trends, where the overall minimum falls in that period, and by the residuals for Uccle data, showing a clear overestimation.

# 3   Methodology

Having two different datasets, we will operate on them separately and then we will compare the model interpretations.

For each station, we fit a model that generalises our training dataset in such a way that forecasting until March $1^{st}$ becomes effective, also assessing errors analytically. Obtaining such a reliable model will be the first step to verify any significant difference forecasting over the lockdown period.

The model we fit is the ARIMA[2] model, including the seasonality components. To identify the order of the different components described by ARIMA, we will analyse the autocorrelation (ACF) and partial autocorrelation (PACF) functions of our datasets. As a reminder:

- non-seasonal AR process of order $p$ is described by an exponentially or sigmoidally decaying ACF, and by a cut-off of significant PACF values at lag $p$

- non-seasonal MA process of order $q$ is described by an exponentially or sigmoidally decaying PACF, and by a cut-off of significant ACF values at lag $q$

---

[2]AutoRegressive Integrated Moving Average

- seasonal AR process of order $P$ is described by an exponential or sigmoidal decay at seasonal lags on ACF values, and by a cut-off of significant PACF values at each seasonal lag until the $P^{th}$

- seasonal MA process of order $Q$ is described by an exponential or sigmoidal decay at seasonal lags on PACF values, and by a cut-off of significant ACF values at each seasonal lag until the $Q^{th}$

- given how significance is assessed for the ACF/PACF spikes, one over 20 spikes (5%) could be *coincidentally significant* or *coincidentally **non**-significant*

The same model will be then used to forecast other 8 weeks from March $2^{nd}$, so that we can assess the forecast power over the lockdown period. Given the fact that both our series are quite stable over time (Figures 4a, 9a), we can expect that the predicted variability for further forecasts will not significantly increase, compromising the long-term forecasts.

Also, we will decompose the seasonalities from data before and after the lockdown, in order to assess any significant change regarding that aspect.

Per our hypotesis, we expect the Uccle station's data to show significant changes after the lockdown, for its position near Brussels, while we expect Membach's forecasts not to show any significant change, for its position is further from busy cities.

We will not directly compare the two models, since data from the two stations have different orders of magnitude, and that difference should be firstly assessed in order to meaningfully compare the models.

## 3.1 Uccle station

**ARIMA**

Our training set, despite the precautions taken during the initial data cleanup and aggregation, does not satisfy the Normality assumption, even after applying the *log* or *Box-Cox* transformations. Thus, we will estimate the models using the original scale.
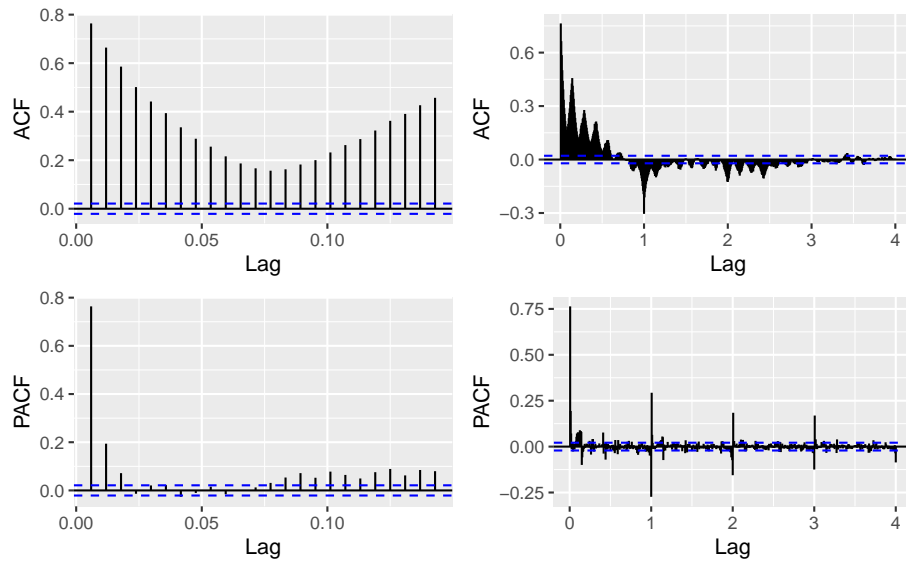


Figure 6: ACF and PACF plots after *diff* at lag 168

Since our dataset shows a relevant trend component, resulting in non-constant mean, we consider the usage of the *diff*[3] operator to smooth the mean. After applying the *diff*

---

[3]The *diff* operator applies the $d$-order differences. All *diff* operators we apply use $d = 1$.

operator with lag 168, in conjunction with the lag 1 *diff* or not, we come up with two possible models by checking the ACF and PACF plots:

- from using both *diff*s, the ACF and PACF on the long run hint to $ARIMA(0, 1, 0)(0, 1, 1)_{168}$ (plots omitted for brevity)
- from using just the lag 168 *diff*, the PACF on the first lags and the ACF/PACF combination on the long run hint to $ARIMA(2, 0, 0)(0, 1, 1)_{168}$ (Figure 6)

We also employ the *auto.arima* stepwise methodology[4], locking the *diff* values to 0 or 1, to get some hints on the possible models.

Following both the analytical ways and the *auto.arima* hints, and refining the models through further inspection of ACF/PACF plots of the residuals, we end up with three characteristic models. We selected them through the Akaike Information Criterion (AIC)[5], being the only three out of the fitted models with AIC < 120.000: $ARIMA(2, 0, 0)(0, 1, 1)_{168}$, $ARIMA(3, 1, 1)(0, 1, 1)_{168}$, and $ARIMA(4, 1, 2)(0, 1, 1)_{168}$.

Also checking the Mean Absolute Percentage Error (MAPE) on the training set, they actually have the lowest three values (MAPE < 6.2%). Nevertheless, MAPE is not a single valid criterion for selection, since a low error predicting the training set itself does not necessarily mean a good forecast power.
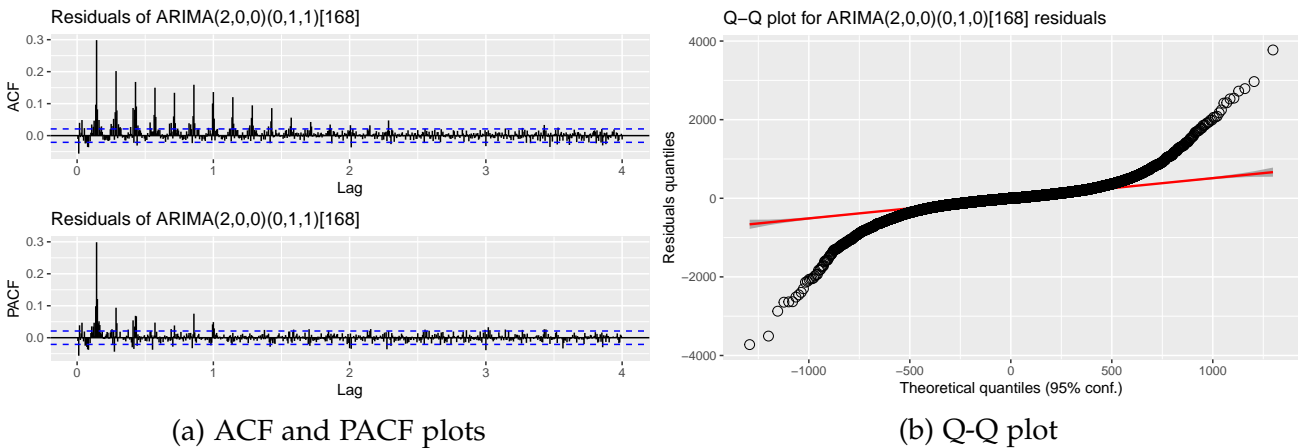


(a) ACF and PACF plots

(b) Q-Q plot

Figure 7: Plots for $ARIMA(2, 0, 0)(0, 1, 1)_{168}$

The ACF and PACF of the residuals for all three models are very similar (Figure 7a), even if the models themselves have great differences in meaning and in number of parameters: they highlight a clear seasonal ARMA process with periodicity 24. This shows how our previous conclusion that the 168 hours seasonality could include the 24 hours one is flawed. Unfortunately, how we are using the R function do not permit us to estimate multiple seasonalities[6].

Probably related to the 24 hours periodicity that remains on the residuals, together with the non-Normality of the original data, all Q-Q plots show clear divergence on both sides of the distribution (Figure 7b). The Box-Pierce test statistics accept the independence hypothesis between residuals.

---

[4]The *auto.arima* stepwise methodology fits different ARIMA models within a range of orders. Based on each result's information criterion (Akaike), it changes each time one of the *p*, *d*, *q*, *P*, *D*, or *Q* values, depending on the configured limits. These changes minimizes the information criteria of the final model.

[5]Actually, we could have used the corrected AIC (AICc) as it is usually more precise, but the values of AIC and AICc are equal to the first decimal for all our models ($\lesssim 10^{-7}$% difference).

[6]*Arima* method from the *forecast* R package, only specifying the seasonal and non-seasonal orders.

Then, we assess how the forecasts behave on the test set. The MAPE for all three models more than doubles to around 17% for $ARIMA(2,0,0)(0,1,1)_{168}$ and around 15% for the other two. This is expected, since normally the model is more fitted to the training data than to unknown future data.

Also, it must be noted that a January $1^{st}$ period lies in our test dataset, that was already noticed as low-valued period (Figure 4). We can notice a clear overestimation during those couple of weeks (around week 62 on Figure 8).
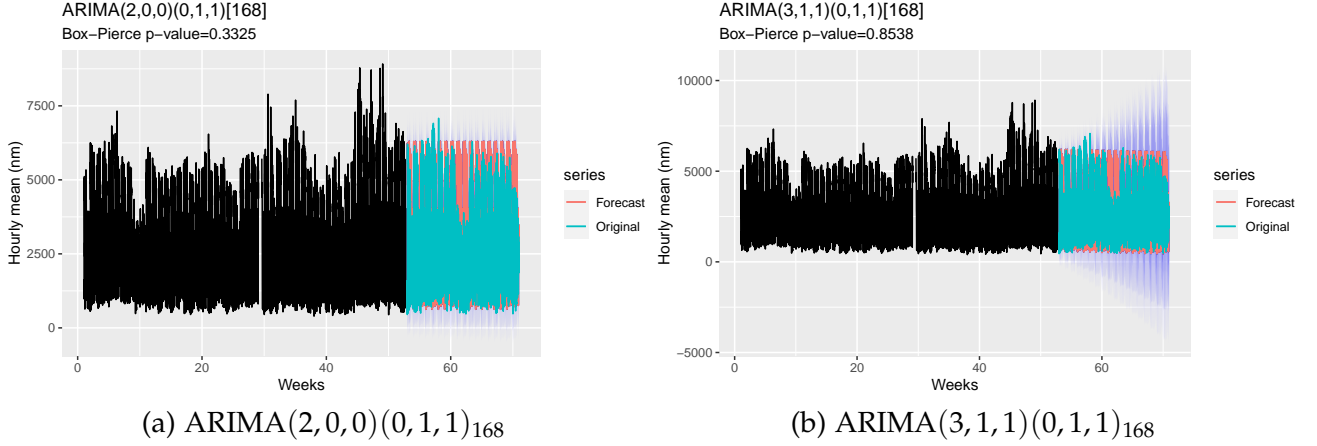


(a) $ARIMA(2,0,0)(0,1,1)_{168}$ 　　　　　　　 (b) $ARIMA(3,1,1)(0,1,1)_{168}$

Figure 8: ARIMA forecasts
$ARIMA(4,1,2)(0,1,1)_{168}$ behaves like $ARIMA(3,1,1)(0,1,1)_{168}$ with faster increasing variablity

We can see in Figure 8 the forecast behaviour for our models. We can appreciate the constant span of the confidence intervals in Figure 8a, in constrast with the more complex models that have increasing intervals further away from the training set. It must be noted that the increasing intervals include negative values, that are actually impossible to obtain given the data meaning (difference between maximum and minimum measurements).

The increasing variability also make the percentage of test values falling inside the confidence intervals to be higher than normal:

- $ARIMA(2,0,0)(0,1,1)_{168}$: $\approx 81\%$ of data is in the 80% CI and $\approx 91\%$ in the 95% CI
- $ARIMA(3,1,1)(0,1,1)_{168}$: $\approx 96\%$ of data is in the 80% CI and $\approx 98\%$ in the 95% CI
- $ARIMA(4,1,2)(0,1,1)_{168}$: $\approx 98\%$ of data is in the 80% CI and $\approx 99\%$ in the 95% CI

Given all these factors, we are confident choosing the $ARIMA(2,0,0)(0,1,1)_{168}$ model, being aware of the missed 24 hours seasonality.

## 3.2 Membach station

**ARIMA**

Similarly to data from Uccle, data from Membach does not respect the Normality assumption. We tested if the Q-Q plot showed significantly better approximation of the Normal distribution after applying the *log* or *Box-Cox* transformations, and the latter makes the data approach Normality with $\lambda \approx -1$ (Figure 9). This solves the high peaks problem we already noticed while decomposing to highlight the trend and seasonalities.

Since the $\lambda$ value approximately simplifies the Box-Cox transformation formula from $(y^\lambda - 1)/\lambda$ into $(y - 1)/y$, and since we only have positive values, all transformed values fall into the higher $0 - 1$ spectrum, in particular $0.95 - 1$ (Figure 9a). This will probably cause some approximation errors that could rig the final results.

(a) Transformed training set
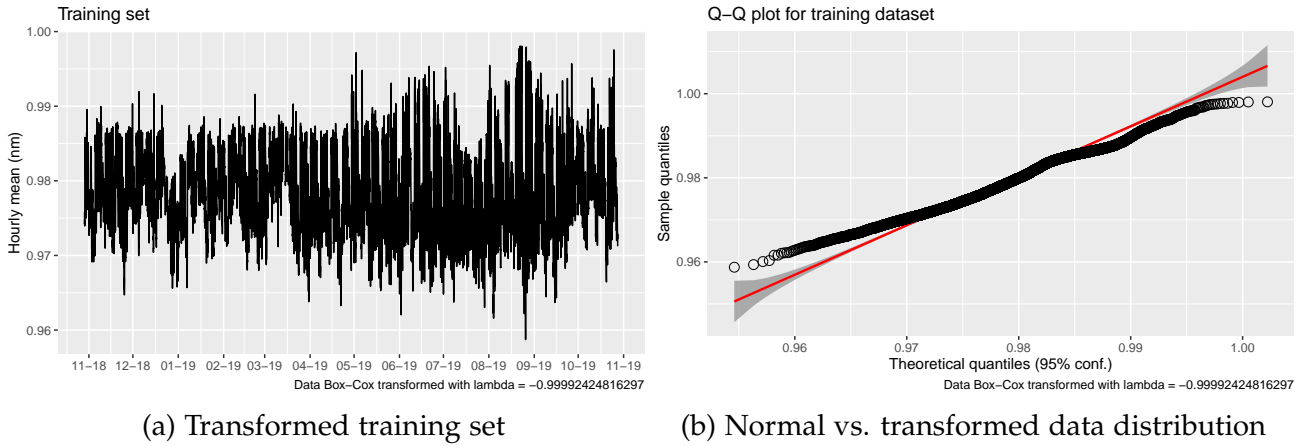


(b) Normal vs. transformed data distribution

Figure 9: Results of Box-Cox transformation

The transformed data still shows no constant mean nor variance, even if they show less variability around the peaks with respect to the original data, confirmed by the Q-Q plot (Figure 9b). Thus, we consider the usage of the *diff* operator, at lag 1 and 168. The analysis of ACF and PACF plots hints to two models: $\text{ARIMA}(0,1,0)(0,1,1)_{168}$ and $\text{ARIMA}(2,0,0)(0,1,1)_{168}$. All ACF and PACF plots for Membach model estimation are not shown for brevity.

The employment of *auto.arima* stepwise methodology[4] highlights the first model improved with one more non-seasonal MA parameter. This is confirmed by the ACF/PACF plots of its residuals, while their analysis for the second model do not lead to any improvement. The ACF and PACF plots of both models are similar to the ones from the chosen model for the Uccle station (Figure 7b), and we can notice the same 24 hours remaining periodicity.

The Box-Pierce test statistics for both models accept the independence hypothesis between residuals. The two models have similar AIC values, and similar MAPE on the training set ($\approx 6\%$). Thus, we can assess the forecasts for both models in Figure 10.



(a) $\text{ARIMA}(0,1,1)(0,1,1)_{168}$



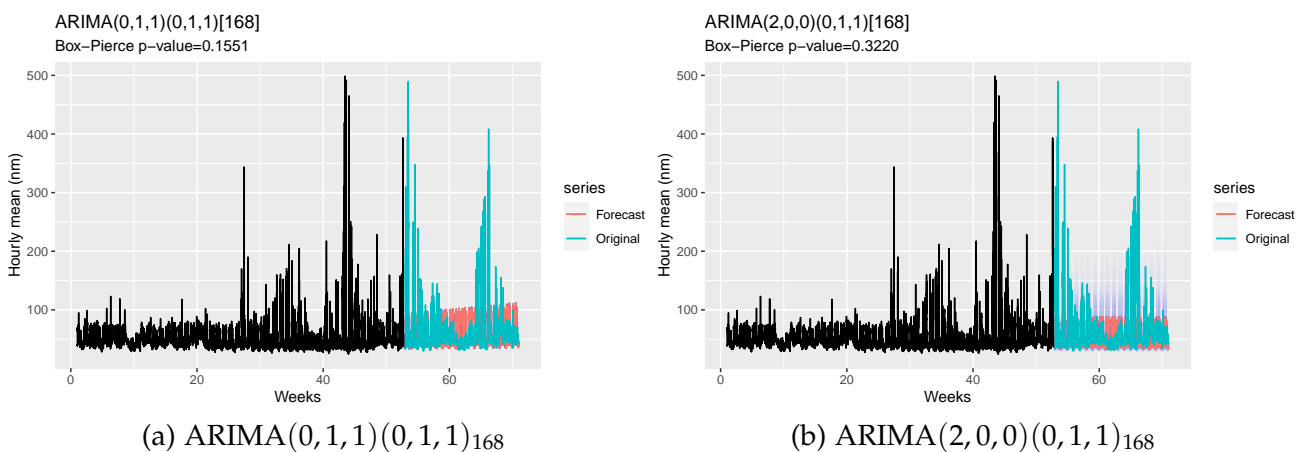(b) $\text{ARIMA}(2,0,0)(0,1,1)_{168}$

Figure 10: ARIMA forecasts

The first model catches a rapidly-increasing trend, while the latter predicts more stable values moving the most part of the increasing trend into slightly increasing variability. Visually, both could be acceptable, even if the latter seems to better behave on *normal* days losing precision on high-valued ones, while the first tries to better accomodate the high-valued days losing precision on the *normal* ones.

Also, it is important to notice that while the first model includes in the confidence intervals negative values, non existent in our domain, the latter directs all the variability to higher values, while maintaining a consistent lower bound. This appears to be a good feature to correctly model our data.

Analytically, we have quite a difference between the two models in terms of adherence of the confidence intervals:

- $\text{ARIMA}(0,1,1)(0,1,1)_{168}$: $\approx 98\%$ of data is in the 80% CI and $\approx 99\%$ in the 95% CI
- $\text{ARIMA}(2,0,0)(0,1,1)_{168}$: $\approx 73\%$ of data is in the 80% CI and $\approx 89\%$ in the 95% CI

As we see, the first model over-generalises the variability, while the latter seems more precise, even if losing some information. Considering the MAPE on the test set, the latter model fits the data better with a 16% error, instead of the 18% for the first model.

Thus, we can select the $\text{ARIMA}(2,0,0)(0,1,1)_{168}$ model on the Box-Cox($\lambda \approx -1$) transformed data, still being aware of the missed 24 hours seasonality, similarly to the Uccle station.

Since checking the ACF and PACF plots for the original data could lead to a model with the same order, we fitted it for comparison. The MAPEs are higher but comparable, 7% on the test set and 19% on the train set, the Box-Pierce test refutes the independence between residuals, and 88%/94% of test data fits inside the 80% and 95% CI. Finally, it shows a variability that tends to admit lower values than the training data shows, including negative ones, while having a lower acceptance for higher values. We omit plots related to this model, since such results confirm the choice of the model on the transformed data.

It is important to highlight that no model we assessed found any pattern for predicting the high-valued spikes.

## 3.3 Lockdown influence

Even if the two models have the same order, they cannot be directly compared using the Information Criteria, for the Membach model works after a transformation pass, and they are fitted from different data.

Still, we can assess how the lockdown period is forecasted for both stations, and then compare their efficency. Also, we can extract the seasonalities, using the same MA filters we already used, separately from the periods before and during the lockdown, and compare also these results.
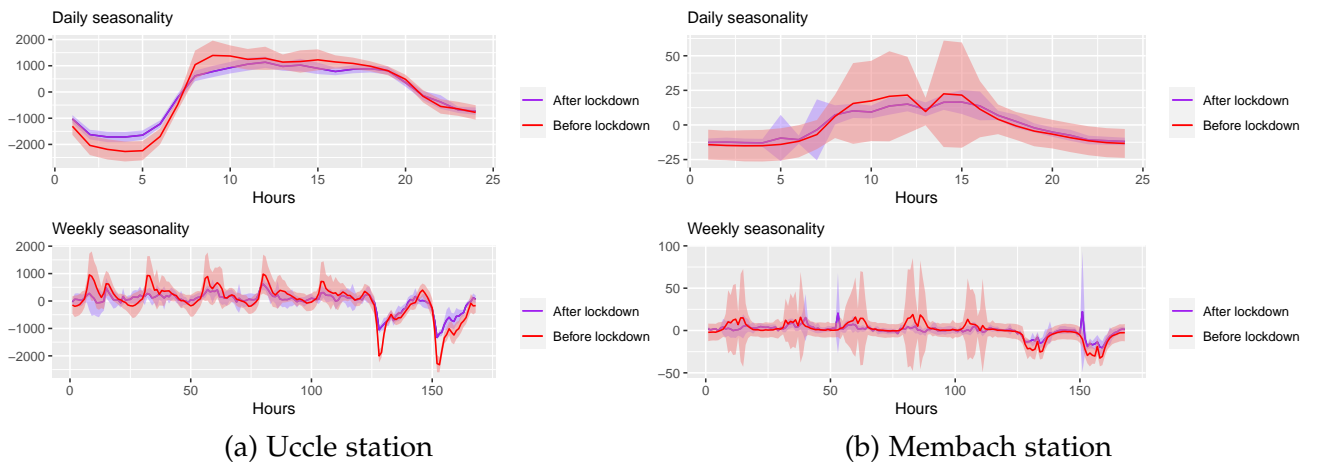


(a) Uccle station                    (b) Membach station

Figure 11: Seasonalities computed before and after the lockdown
Seasonalities extracted together with MA filters, confidence bands $\pm$sd

## Seasonalities divergence

We highlight the confidence bands for an estimated error of 67%. This interval estimation lies on the assumption that values for each single hour in the hourly/weekly seasonality is Normally distributed, a slightly different assumption that is not necessarily refuted by the already seen Q-Q plots.

As we can see in Figure 11, the seasonalities' shapes seem different, in particular comparing the weekly seasonality. The seasonalities from before and after the lockdown are still included in the highlighted confidence bands, apart from a few hours. This is true for both stations.



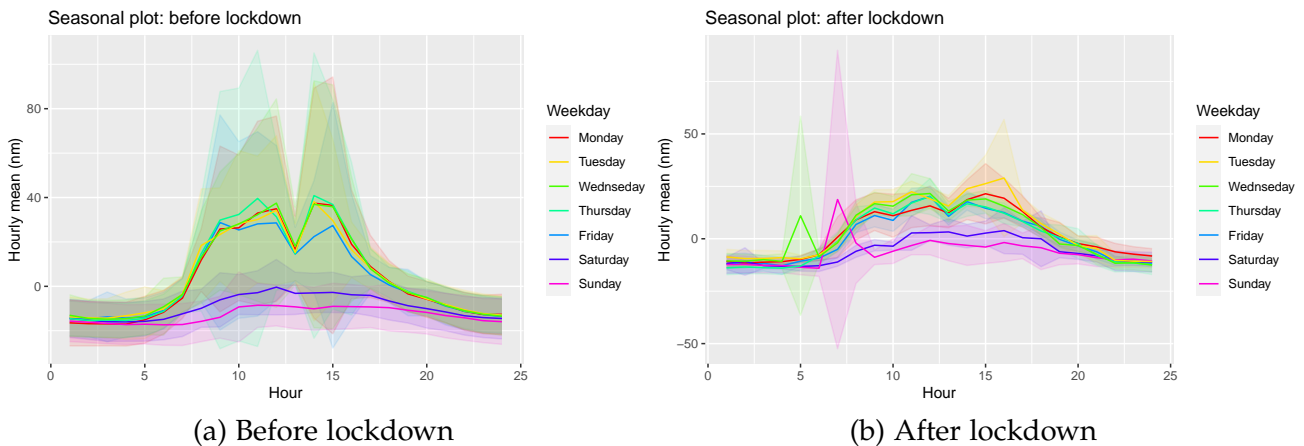(a) Before lockdown                    (b) After lockdown

Figure 12: Seasonal plots for Membach station

The differences in shape are more highlighted checking the seasonal plots (8) in Figure 12. It is possible to notice a couple of changes: the Monday to Friday shapes are heavily flattened almost on the weekend values, and all the values are flattened towards 0. Both flattenings can be noticed also checking the Uccle station, not showed for brevity.

## Forecast divergence

Since both our models forecast stable values with almost stable variance (Figures 8a, 10b), we can forecast another 8 weeks ahead, until April $26^{th}$, and then compare the efficency of the results with the original test set. These 8 weeks comprise 6 weeks after March $14^{th}$, first day of lockdown in Belgium.

As we can see in Figure 13, at both stations the models overestimate starting after the lockdown. Uccle data in particular show a greater overestimation.

We assess analytical values obtained from the new 8 weeks forecast:

- the MAPE for Uccle triples with respect to the original test set to $\approx 46\%$, and only 39% and 63% of the lockdown data falls inside the confidence intevals (80% and 95%).

- regarding the Membach station forecasts, the MAPE is reduced to 14%, while 82% and 96% of data falls inside the confidence intervals.

Thus, we notice a visible degradation in performance for Uccle station data, and a slight improvement for Membach station's. The slight improvement is most certainly driven by the single peak with lower values than the multiple previous peaks, and by the increased variability for data further from the end of the training set.
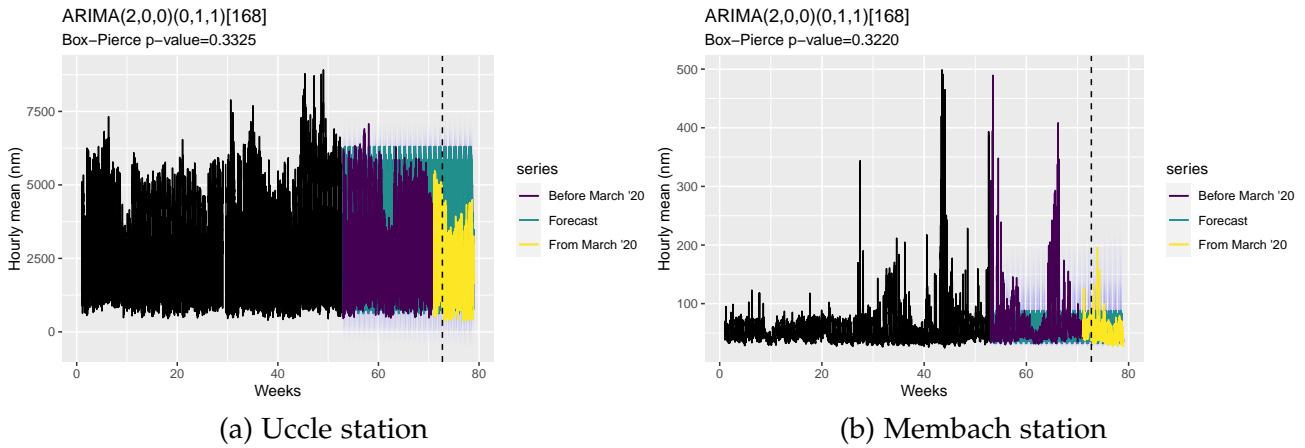
(a) Uccle station  (b) Membach station

Figure 13: Forecasts from March $2^{nd}$.
The forecasted values are on the background. March $14^{th}$ at midnight is highlighted.

# 4 Conclusions

The final models for both stations show interesting predictive efficiency, even if analysing the ACF and PACF plots of their residuals we notice a remaining daily seasonality that is not explained by the models (Figure 7b).

The initial assumptions on the seasonalities are confirmed by our study, given the periodgrams (Figure 3) and the efficiency and structure of the models (seasonal terms). Though, the origin of those assumptions (human activity) cannot be solely confirmed by our methodology.

Interestingly, both final models have the same order, $ARIMA(2, 0, 0)(0, 1, 1)_{168}$: still, they cannot be actually compared, given the Box-Cox transformation applied to the Membach data.

**Lockdown influence**  The assessment of the lockdown influence on seismic data show relevant information for further studies.

Extracting both daily and weekly seasonalities with the local mean methodology (Figure 11), it is not possible to appreciate any statistically significant difference from before and during the lockdown on both datasets.

Checking the seasonal plots of the complete weekly seasonality (Figure 12), it is possible to notice a flattening of the values towards 0, in particular regarding Monday to Friday values. This seem to confirm the initial hypotesis of lockdown influence.

The Uccle station, located near the Belgium capital city, show a relevant difference in forecasts after the lockdown, that are heavily out of the confidence bands, reaching over 45% of MAPE. This is the opposite for data from Membach station, located in a small town near a natural park, that show no significant difference in forecasts performance.

The forecast results hint to a significant correlation between the location of the seismic station and the influence of the lockdown.

## 4.1 Limitations and future works

The main problems identified during this analysis are the length and number of the datasets, the multiple seasonalities and the Normality assumption.

Even if after a first decomposition we chose to only consider the weekly seasonality,

our models show a remaining daily seasonality. This could be addressed with various approach. One initial test could be to forecast a different model on the residuals with a daily seasonality, so that the composition of the two models could include both seasonalities.

A second mean to address the multi-seasonalities is through a harmonic regression approach. Using that, we could model multiple seasonalities at once through Fourier terms. One difficulty we could encounter is that our weekly seasonality appear to be flat for most of the week (Figure 11), and such a behaviour is not easy to model with a few Fourier terms per seasonality, leading to computationally heavy modelling and forecasting.

Different models, such as STL (Seasonal and Trend decomposition using Loess) or TBATS (Trigonometric seasonality, Box-Cox transformation, ARMA errors, Trend and Seasonal components), are suggested in the literature (8, chapter 12) to address multiple and complex seasonalities, also when they change in time. These could better tackle the missing Normality assumption of our data, along with relaxing the initial transformation between timezones.

A future study should be done when more data is available, in order to assess the presence of higher order seasonalities, such as a yearly one.

The confidence bands used to assess any significant change in seasonalities lie on another Normality assumption, that data for each hour of the day/week follow a Normal distribution. This should be assessed in a future study integration.

On a more broad perspective, other analyses should be done on the lockdown impact, and more generally on isolating different factors that influence our datasets. Repeating the study on the other 7 stations will surely give interesting insights.

# References

[1] Elizabeth Gibney. Coronavirus lockdowns have changed the way Earth moves. *Nature (online)*, March 2020. ISSN 1476-4687. doi: 10.1038/d41586-020-00965-x. URL `https://www.nature.com/articles/d41586-020-00965-x`.

[2] Royal Observatory of Belgium. Other types of seismic events, April 2020. URL `http://seismologie.be/en/seismology/other-types-of-seismic-events`.

[3] Royal Observatory of Belgium. Data policy, April 2020. URL `http://seismologie.be/en/legal-notices/data-policy`.

[4] Royal Observatory of Belgium. Uccle station data from 02-04-2020, April 2020. URL `http://seismologie.be/data/csv/SEIS-2020-093-UCCS.csv`.

[5] Royal Observatory of Belgium. Uccle station data events from 02-04-2020, April 2020. URL `http://seismologie.be/data/csv/events_2020-04-02.json`.

[6] Royal Observatory of Belgium. Uccle station data visualization from 02-04-2020, April 2020. URL `http://seismologie.be/en/seismology/seismograms/uccs/20200402`.

[7] Peter J Brockwell, Richard A Davis, and Stephen E Fienberg. *Time series: theory and methods*. Springer Science & Business Media, 1991.

[8] Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*, chapter 2. OTexts, 2018.