

Assignment 1 Geo1001

Lisa Geers

September 21, 2020

1 Introduction

This assignment was made for the class Geo1001. A statistical analysis was done using a heat stress measurement dataset with five sensors [1]. Analysing was done in Visual Studio Code using Python 3.7.5. Plotting was done with the use of Matplotlib. The source code for this assignment can be found on GitHub: https://github.com/Lisageers/geo1001_hw1

2 A1

2.1 Mean statistics

In Table 1, the calculated mean statistics of all sensors are displayed. When zoomed-in, the individual values can be seen. The means of the sensors are quite similar for all variables. The means of the wind variables Direction - True, Wind Speed, Crosswind Speed and Headwind Speed differ the most between sensors. This is logical, because wind can differ greatly over short distances. This is in contrast with other variables like Temperature and Relative Humidity, which are less dynamic and thus differ less in means.

It stands out that the standard deviations of the variables Wind Direction True, Wind Direction Magnetic and Density Altitude are higher than the other variables, meaning that there are more values further from the mean. The differences in standard deviations between the sensors are comparable to the differences between means, differences are higher for the wind variables.

The variances of Wind Direction True, Wind Direction Magnetic, Density Altitude and Altitude are much higher than the other variables, which means that the values of these variables are very spread out. The differences between sensors are, like the means and standard deviations, the largest for the wind-related variables.

Table 1: Mean Statistics of all sensors

2.2 Histograms

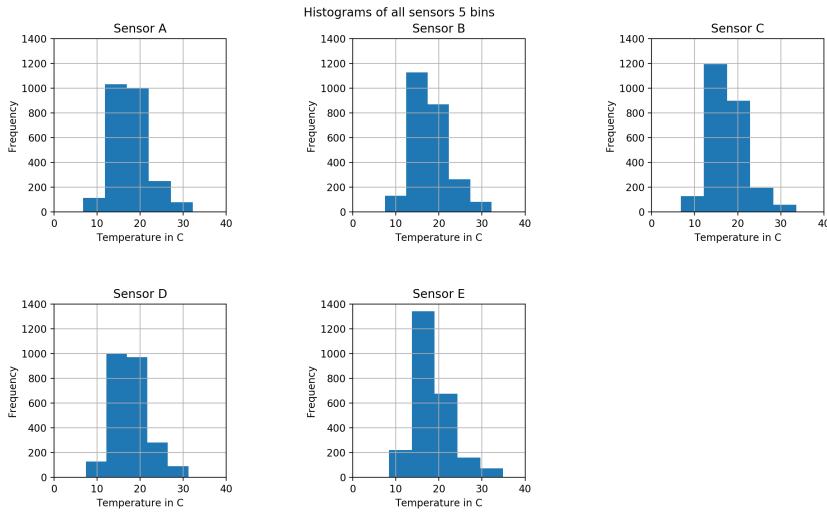


Figure 1: Histograms of all sensors with 5 bins

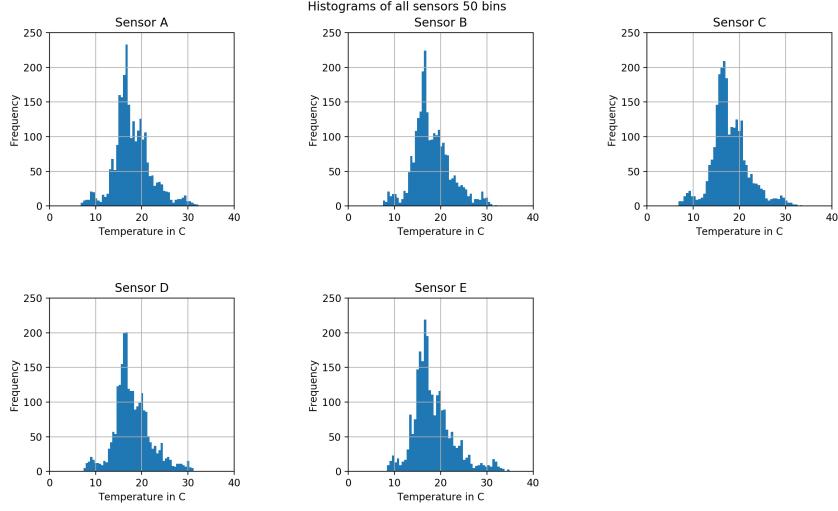


Figure 2: Histograms of all sensors with 50 bins

In Figures 1 and 2, the five histograms of the Temperature variable are displayed. As can be seen, there is a significant difference between the figures due to the bin sizes. Figure 2 with binsize 50 is much more detailed, which makes this figure more useful for analysis. This shows that the number of bins is important in order to be able to do the right analysis. The binsize calculated with Rice's rule is approximately in the middle between 5 and 50. Rice's rule $2 * \sqrt[3]{N}$ with $N = 2474$ gives 27 as a number of bins.

2.3 Frequency polygons

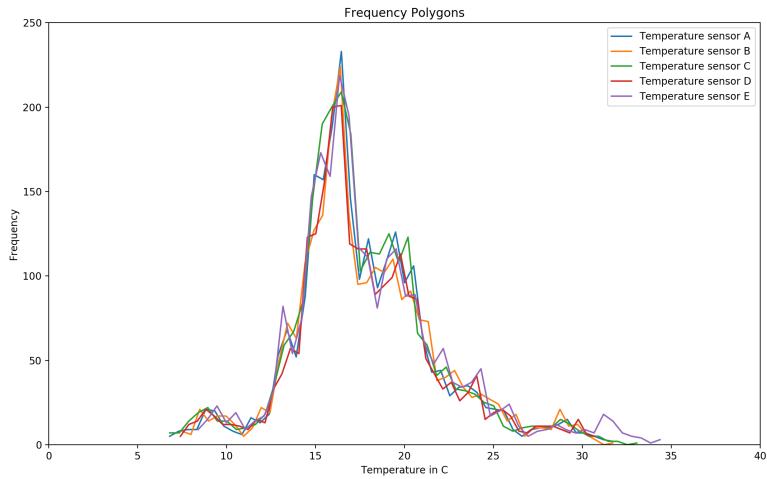


Figure 3: Frequency polygon of all sensors for the variable Temperature

2.4 Boxplots

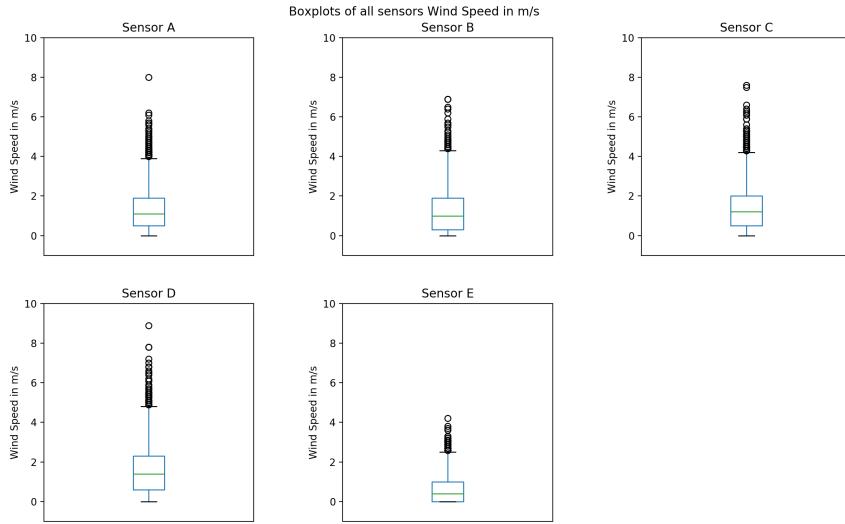


Figure 4: Boxplots of all sensors for the variable Wind Speed

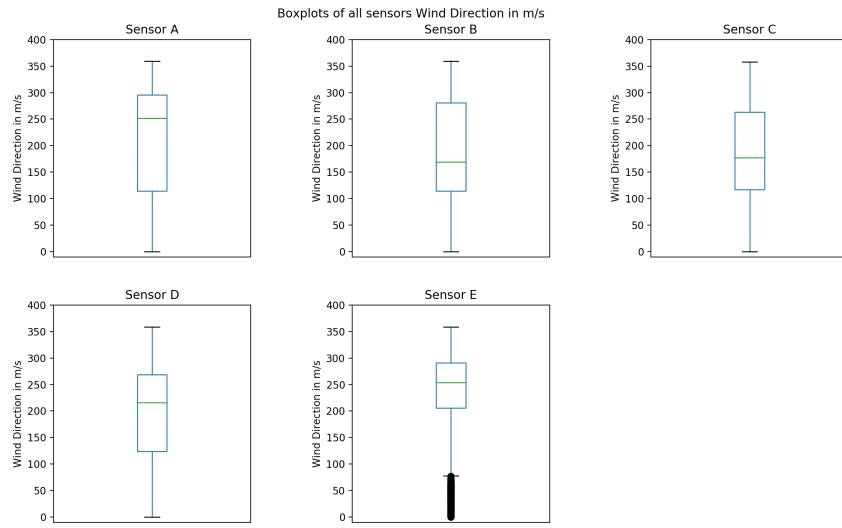


Figure 5: Boxplots of all sensors for the variable Wind Direction

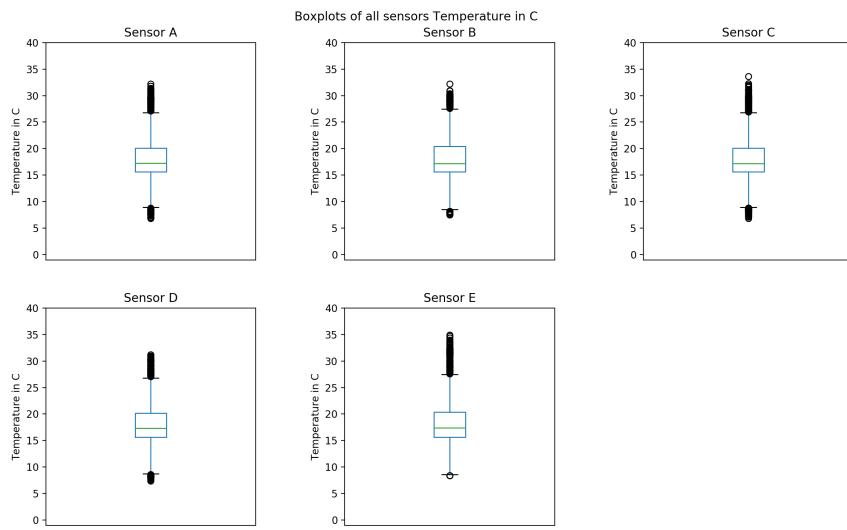


Figure 6: Boxplots of all sensors for the variable Temperature

3 A2

3.1 Functions Temperature

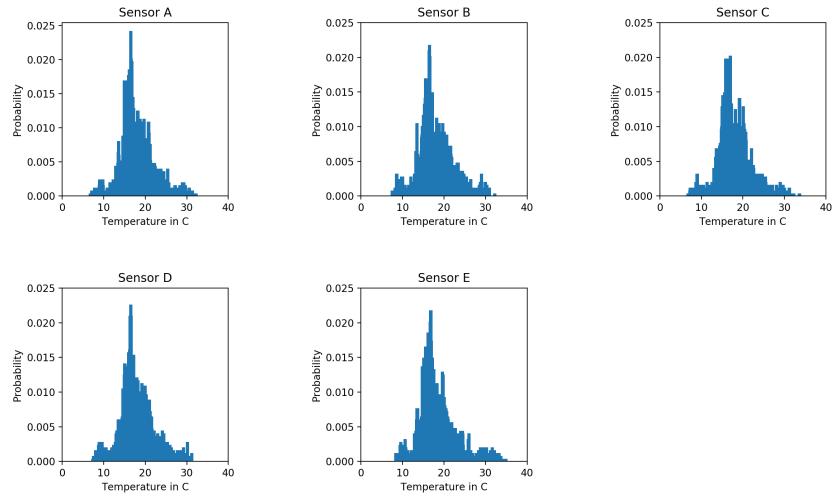


Figure 7: Probability Mass Functions of Temperature for all sensors

The probability mass functions (Figure 7) of the different sensors look similar and are relatively normal distributed. The left side of the curves look steeper, meaning that the distributions are slightly skewed positively. Sensor E is skewed the most compared to the other sensors.

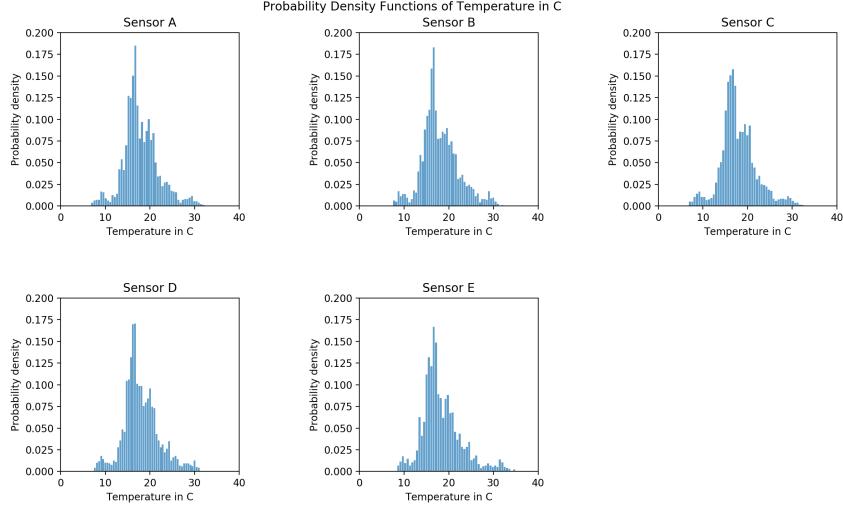


Figure 8: Probability Density Functions of Temperature for all sensors

The behaviours of the probability density functions (Figure 8) are similar to the probability mass functions. They are slightly positively skewed. In these plots the bimodal properties of the data can be seen more clearly. All plots show two peaks in the data, with the first peak being higher than the other.

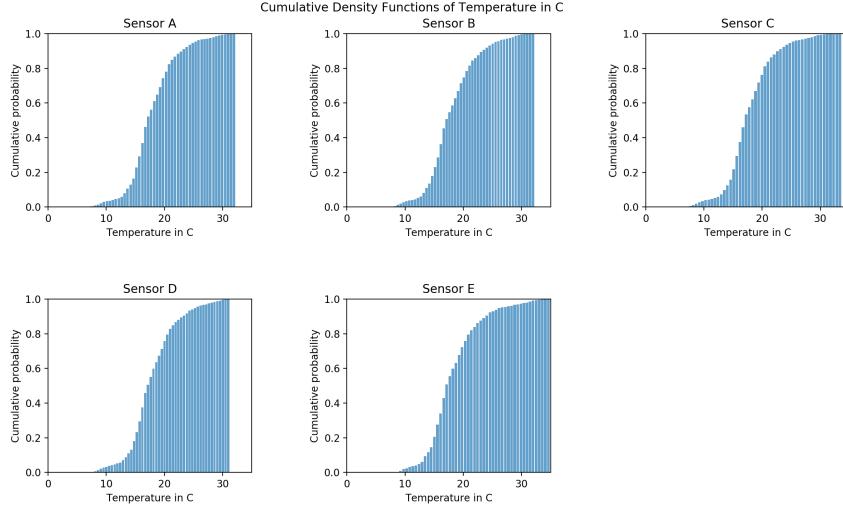


Figure 9: Cumulative Density Functions of Temperature for all sensors

Figure 9 shows the cumulative density functions of all sensors. The functions look similar for all five sensors. The cumulative probability increases the most in the middle (around $p=0.5$), which indicates a normal distribution. Sensor E seems to have higher temperature values, because the plot shows more lines on the right side.

3.2 Functions Wind Speed

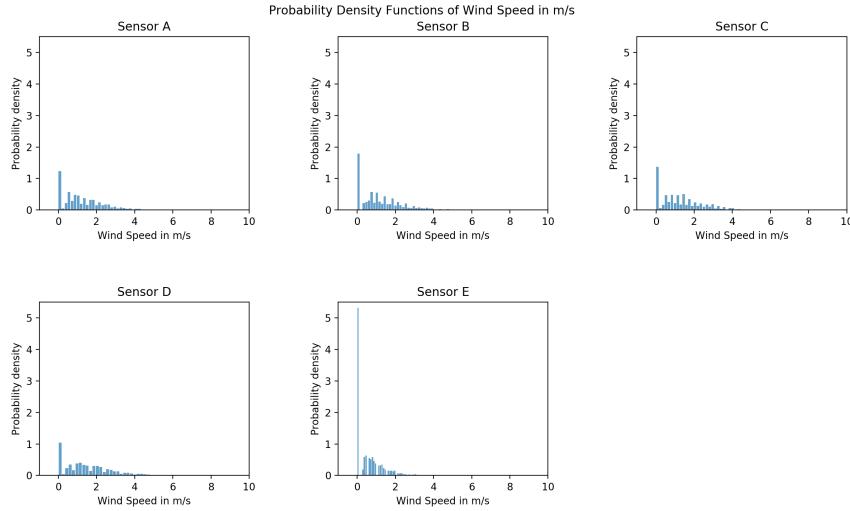


Figure 10: Probability Density Functions of Wind Speed for all sensors

The PDF plots of Wind Speed (Figure 10) are very different from the Temperature plots. They are not normally distributed, have a very high peak at zero and only a right tail. This makes sense, because having no wind occurs more frequently than, for example, a wind speed of 2 m/s. Also, low wind speeds are more common than high wind speeds, which is why the plots are descending to the right side.

There is also a big difference between the sensors. The plot of sensor E deviates the most. This is mostly due to the very high frequency of the wind speed zero. It is about three times as high as the other sensors.

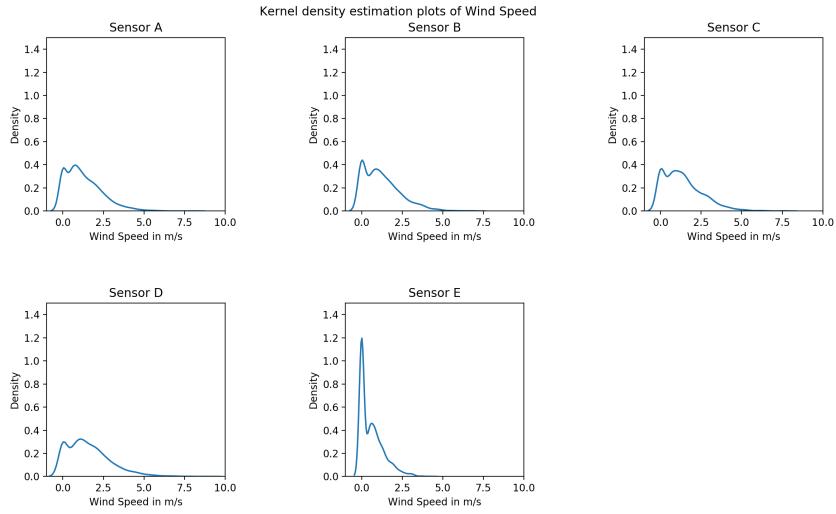


Figure 11: Kernel Density Estimation of Wind Speed for all sensors

In the kernel density plots (Figure 11) it can also be seen that the data is not normally distributed. The functions have two peaks: one for no wind and one for the wind speed with the highest frequency. The differences between sensors can also be seen clearly. The first peak of sensor E is much higher than the other sensors and the line stops at a lower value, because the wind speed of sensor E has a much higher frequency of zero.

4 A3

4.1 Correlation

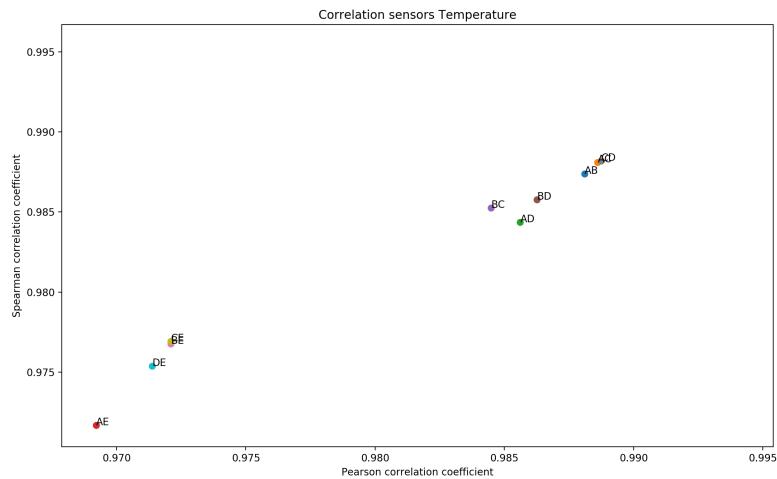


Figure 12: Spearman and Pearson correlation plot for all sensor combinations of Temperature

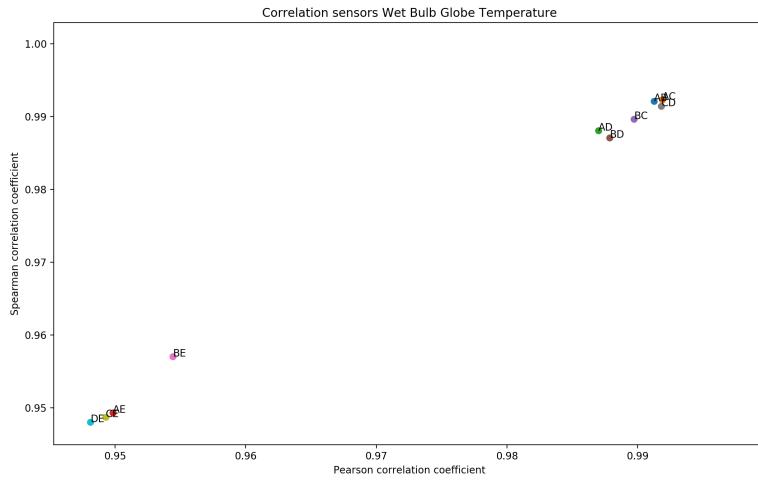


Figure 13: Spearman and Pearson correlation plot for all sensor combinations of Wet Bulb Globe Temperature

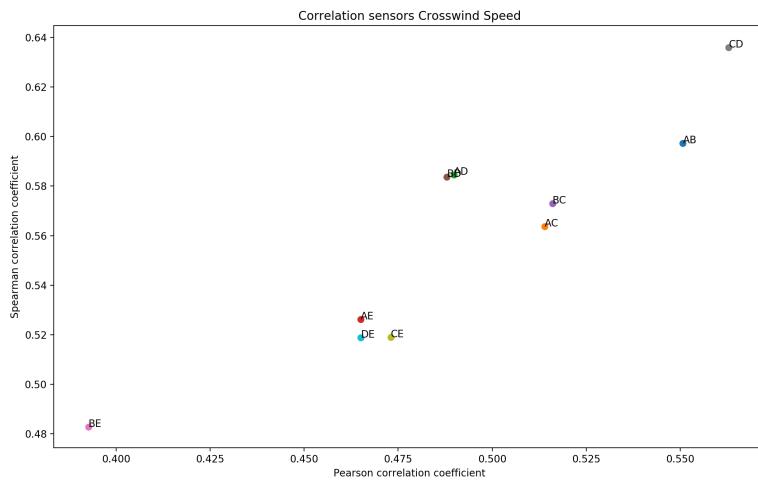


Figure 14: Spearman and Pearson correlation plot for all sensor combinations of Crosswind Speed

In the correlation plot of Temperature (Figure 12) it can be seen that all pairs are positively correlated. The correlation values are close to one, so there is a

strong positive relationship between all variables. The pairs with sensor E have a relatively lower value.

The correlation values of Wet Bulb Globe Temperature (Figure 13) are also high, so there is a positive relationship between all sensors. The pairs with sensor E are, once again, somewhat lower than the other pairs.

The correlation values of Crosswind Speed (Figure 14) are a lot lower than correlation of the other variables. However, there is still a moderate positive relationship between all variable pairs. The pairs with sensor E are the three lowest correlation values and the pair of sensor B and sensor E is clearly the lowest.

4.2 Sensor locations



Figure 15: Sensor locations by number

In Figure 15 the five potential locations of the sensors are shown. As was clear in the last section, sensor E differs the most from the other sensors. Therefore, its location would lay the furthest away, which is location 3. Furthermore, the pair of sensors C and D had the highest correlations for all variables, so they should lay the closest together, which are locations 1 and 2. This leaves locations 4 and 5 for sensors A and B. This seems correct, because the correlation between A and B are the second highest and locations 4 and 5 are the second closest together. The AE pair has a lower correlation for variables Temperature and WBGT compared to pair BE. Only for Crosswind Speed the correlation of AE

is higher than BE, but because wind is so variable, the temperature variables are valued more. So Sensor B is hypothesized to be situated at location 4 and sensor A at location 5, because B has a higher correlation with E. Because location 1 is the closest to sensor A, the sensor that is situated there should be more correlated to A than the sensor at location 2. Although is is very close, sensor C is more strongly correlated with sensor A than sensor D for variables Temperature and WBGT. For variable Crosswind Speed, pair AD is slightly more correlated, but as mentioned before, this variable is less valued due to its high variability. So, sensor C is estimated to be situated at location 1 and sensor D is situated at location 2. The hypothesized locations are shown in Figure 16.

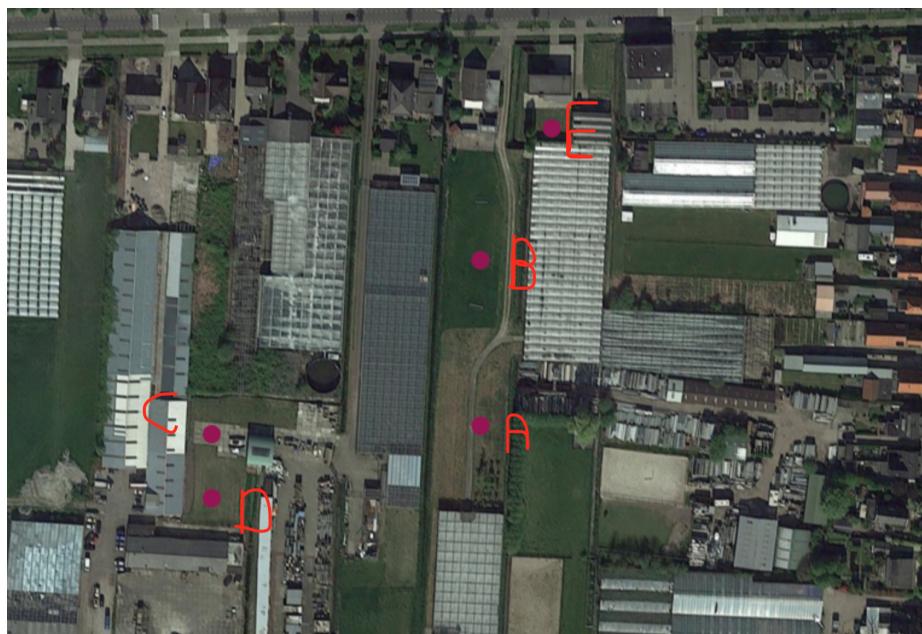


Figure 16: Hypothesized locations of sensors

5 A4

5.1 Cumulative Density Functions

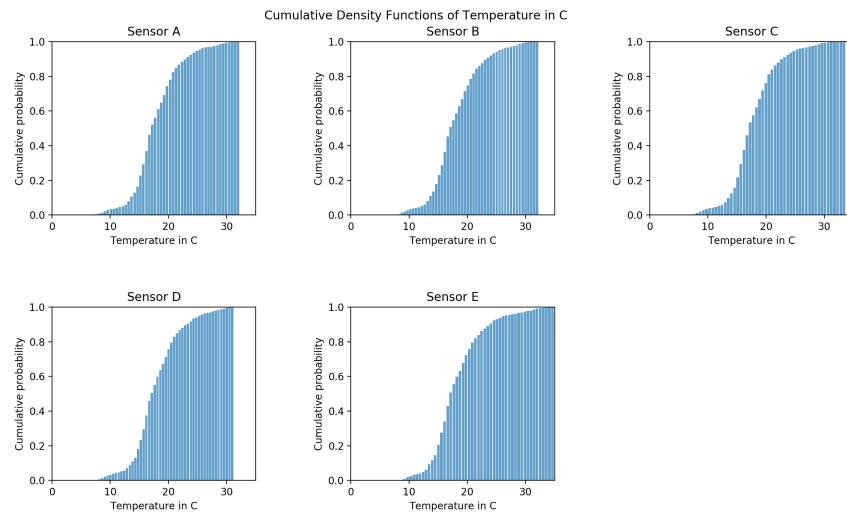


Figure 17: Cumulative Density Functions of Temperature for all sensors

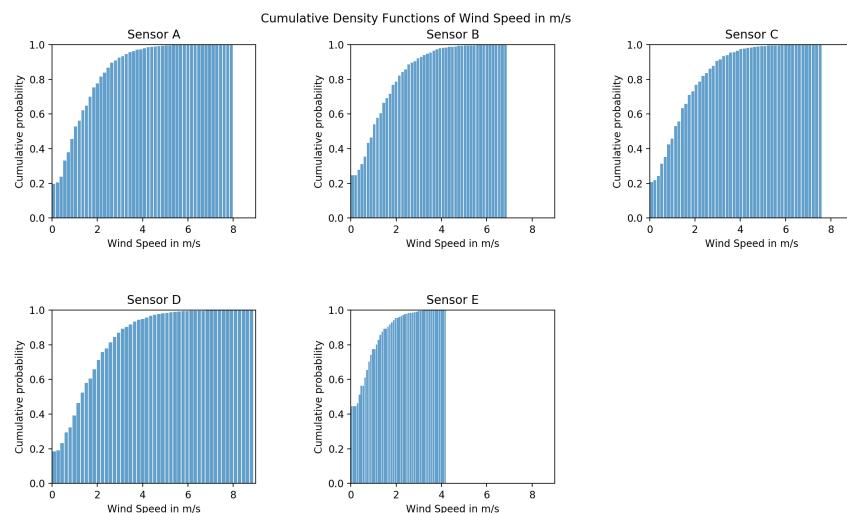


Figure 18: Cumulative Density Functions of Wind Speed for all sensors

5.2 Confidence Intervals

Table 2: Confidence intervals of Temperature for all sensors

Temperature	
A	(17.81214113267346, 18.126065652463858)
B	(17.90472689963894, 18.226129320070267)
C	(17.754926235060246, 18.071347006653575)
D	(17.83814660824381, 18.15457772482005)
E	(18.181933946027776, 18.525944841851015)

Table 3: Confidence intervals of Wind Speed for all sensors

Wind Speed	
A	(1.246227038990971, 1.3343868543854427)
B	(1.1971663346979249, 1.287082453670411)
C	(1.3243037885948932, 1.418622646328308)
D	(1.5296480419653757, 1.633650260379006)
E	(0.5680599051948441, 0.6244249432900044)

5.3 Hypothesis Test

Table 4: Confidence intervals of Wind Speed for all sensors

	Temperature	Wind Speed
p-value E, D	0.0027270117155346967	4.899592405994867e-212
p-value C, D	0.4657972008220813	4.610149126224334e-09
p-value B, C	0.18562772895626528	9.40075204600199e-05
p-value A, B	0.40185871871215073	0.13247973112544695

H0: $\mu_1 = \mu_2$

Ha: $\mu_2 \neq \mu_2$

The H0 is: the means of both sensors are the same.

The Ha is: the means of the sensors are not the same.

For the Temperature variable, only the pair ED has a p-value below 0.05, which rejects the null-hypothesis that the sensors are similar. All the other pairs have a p-value above the α of 0.05, so they are assumed to be similar.

For the variable Wind Speed pairs ED, CD and BC have a very low p-value, so the H0 is rejected and they are not assumed similar. So, only pair AB is assumed similar for Wind Speed.

References

- [1] Daniela Maiullari and Clara Garcia Sanchez. Measured Climate Data in Rijsenhout, 8 2020.