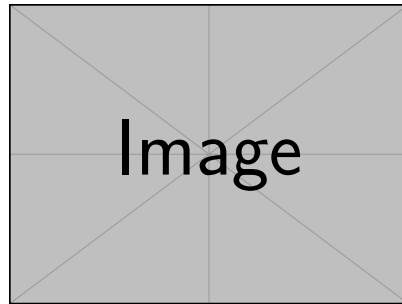


AUTUMN SEMESTER 2025



DATA-MINING IN ENVIRONMENTAL SCIENCE

---

LEARNING DIARY

---

<b>Date:</b>	November 22, 2025
<b>Author:</b>	Stephen Weybrecht
<b>Student number:</b>	2505376
<b>Supervisor:</b>	Mikko Kolehmainen

## Contents

1. Orientation	3
2. Introduction and basics	4
3. Environmental data and its pre-processing	10
4. Data visualization	13
5. Clustering in Python	14
6. Predictive modeling	15
7. Summary	16
8. Self-evaluation	17

# 1. Orientation

My personal background is that I am an exchange student in physics, studying at the University of Eastern Finland for the autumn semester. As such I have quite a strong background in programming, especially in Python, statistics and machine learning already. As I am quite interested in these topics and my home university offers a rather flexible study plan, I further deepened my knowledge by choosing multiple electives and a Bachelor thesis topic that were deeply connected with data analysis already. Although I expect that many things in the beginning of the course will be topics I already learned, I am very much looking forward to developing a deeper understanding of Data mining and seeing this done in a context I have no previous knowledge yet – namely Environmental science. Looking at the curriculum, there are also many topics I have had no prior experience in which is quite exciting to me. In summary, I expect this course to build nicely on my previous knowledge while additionally providing interesting insights to the field of Environmental science.

As a physics student, I am very used to writing scientific paper-like reports. This is the idea behind many reports of practicals I already needed to write as well as my Bachelor thesis. In these the expression of ones own opinion is actively discouraged. I would even go further and say that conciseness and scientific correctness are virtues hammered into us for years during our studies. Naturally a Learning diary such as this where the expression of a personal opinion and a critical reflection about the topics learned is not only encouraged but actively required is therefore quite a step out of my comfort zone. Still, I am looking forward to experiencing this new concept and seeing how it will shape my learning experience. At least at the time of starting this course this integrated approach of always putting learned things in ones own context, thinking critically and still performing quantitative task during the exercises seems like a very natural way to learn. It will be quite interesting to see how this will change during the course.

My future job prospects as a physicist will most likely revolve about programming and handling large amounts of data, regardless of whether I will pursue a career in industry or I will stay in academia. Jobs as a Data Scientist, Programmer or in the engineering direction are quite common when getting a Masters in physics and experimental physics in academia has mostly been computing, simulation or the analysis of huge amounts of data since many years already. Therefore, having a strong basis in programming, data analysis and visualization are skills one should have after the studies. I expect that this course will deepen my knowledge in Data Mining by not only introducing new concepts but also connecting those learned already on an even deeper level and will thus be a helpful resource for my future.

I will use Large Language Models in the following mainly for getting code suggestions for the exercises and help in the layout of this report (as LaTeX can be rather cumbersome at times). The text will mainly be written by myself, although sometimes AI is used for translation and paraphrasing purposes. All code of the exercises as well as the LaTeX files to create this report will be made available on a public GitHub repository [1].

## 2. Introduction and basics

### Lecture

The lectures this week dealt with introductory topics regarding the structure of this course, a math and statistics rehearsal as well as an introduction to the topic Data Mining. The math and statistics chapter covers many basic definitions like the axioms and basic properties of probabilities, the definitions of partial derivatives, vectors and matrices and their addition and multiplication properties. Although these are nice to have for completeness sake and should prove helpful for students which do not have a background in statistics yet, for me, they were already known and will therefore not be repeated here. Instead, I will focus on definitions of this chapter which I do not know by heart and which I think will prove useful for the following and will put them into context of what I have already learned.

The sample mean  $\hat{\mu}$  and sample variance  $\hat{\sigma}^2$  provide unbiased estimators for the population mean  $\mu$  and population variance  $\sigma^2$  given a certain sample  $v_i$  of size  $N$ , i.e.  $i \in (1, N)$ . They are defined as

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N v_i \quad (1)$$

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (v_i - \hat{\mu})^2 \quad (2)$$

I recall from a prior course on statistics that the sample variance with the  $N - 1$  term in the denominator should be used in the case of an *unknown* mean, i.e. if the mean is estimated by Equation 1, as it provides an unbiased estimator and therefore better convergence to the true but unknown sample variance for small  $N$ . If instead the mean is inferred through different means, the minus 1 term can be dropped. In `numpy` the sample variance can be simply calculated by setting the parameter `ddof=1`:

```
1 import numpy as np
2 # array is a sample array of data
3 sample_var = np.std(array, ddof=1)
```

Listing 1: Sample variance, calculated in numpy

If there are more than just one random variates a variance can be calculated for each one. Additionally, a so-called covariance between two different variated can also be calculated. Covariances and variances are summarized in a covariance matrix, whose elements are defined as follows:

$$\text{Cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^N (x_i - \hat{\mu}_x)(y_i - \hat{\mu}_y) = \hat{\sigma}_x \hat{\sigma}_y \rho_{x,y} \quad (3)$$

Here  $\hat{\sigma}_i$  are the standard deviations of the variates  $x$  and  $y$  according to Equation 2 and  $\rho_{x,y}$  is the degree of correlation between  $x$  and  $y$ . It holds that  $\rho_{x,y}$  is always between  $-1$  and  $1$  with  $-1$  indicating maximum negative correlation,  $1$  maximum positive correlation and  $0$  no correlation at all. The covariance matrix (or the reduced correlation matrix obtained when

removing all individual standard deviations) therefore indicates correlations between different parameters in a dataset which can be used for explorative data analysis. Another use of it is when fitting a model to a dataset. Here a large degree of correlation between model parameters indicates a surplus of model parameters.

The normal distribution is the most important continuous probability distribution as a lot of measured data follows it. This is due to the central limit theorem stating that means of random variables taken from arbitrary probability distributions will be distributed normally. As measured quantities are usually means over some finite measurement integration time due to a finite detector resolution many data are distributed normally. The normal or Gaussian probability distribution is given by

$$P(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (4)$$

with  $\mu$  its mean and  $\sigma$  its width, corresponding to the population standard deviation in  $x$ . It is a symmetric distribution and used to formulate confidence intervals as follows:

$$\begin{aligned} P(\mu - 1\sigma < x < \mu + 1\sigma) &\approx 68\% \\ P(\mu - 2\sigma < x < \mu + 2\sigma) &\approx 95.5\% \\ P(\mu - 3\sigma < x < \mu + 3\sigma) &\approx 99.7\% \end{aligned} \quad (5)$$

Further important definitions are those of the Jacobian and Hessian matrices, which are matrices of first order derivatives of vector valued functions  $\mathbf{f}(\mathbf{x})$  or second order derivatives of scalar valued functions  $f(\mathbf{x})$  respectively. They are used in optimization algorithms like data fitting or neural network learning. These methods are usually implemented already in various Python packages, however I still list the form of the Jacobian and Hessian here for completeness:

$$\mathbf{J} = \nabla_{\mathbf{x}} \cdot \mathbf{f}(\mathbf{x}) \quad \text{or element-wise:} \quad J_{ij} = \frac{\partial f_i}{\partial x_j} \quad (6)$$

$$\mathbf{H} = \nabla_{\mathbf{x}}^2 f(\mathbf{x}) \quad \text{or element-wise:} \quad H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j} \quad (7)$$

The last part of the lecture considered Data Mining. We discussed that Data Mining is a very structured process, where certain steps are done in order to hopefully get meaningful insights from the data at hand and be able to form a model that generalizes well. The first step is explorative data analysis (EDA), where plotting the data (e.g. in histograms or box plots) or calculating covariance matrices to get a rough idea about correlation are utilized. This step makes a lot of sense to me, especially once one deals with larger and larger datasets. However, this has been quite underutilized during my physics studies so far, as we have usually worked the other way around by first forming a hypothesis (for example a model) and then applying it to the data and doing statistics to see if it fits. It is quite interesting to me to do a different approach in this lecture and to see, where this may lead me. I am looking forward to applying this to real datasets in the following weeks during the tutorials. The next step is descriptive modelling, where one tries to find patterns in the data by, fundamentally, playing around with it. One can for example try to fit functions like multivariate Gaussians or use

at  
end:  
see if  
this  
really  
was  
inter-  
est-  
ing!

clustering methods on the data. Additionally, one can also try to use predictive modelling like classification neural networks to group the data and gather new insights. Hopefully these steps then lead to discovering patterns and rules in the data and finding a strong, i.e. simple model with high degree of generalization and prediction-power instead of remaining with weak models, which do offer some insights but fail to generalize and really understand the data. Weak models are also often classified by many parameters leading to the aforementioned overfitting. During my prior classes where I often used model fitting and machine learning I have often had contacts with overfitting due to fit functions with too many free parameters or overly large neural networks. I therefore know, that avoiding overfitting is one of the most difficult parts in data analysis and am definitely looking forward to learning more about this.

## Exercise

The goal of the first exercise was to get a handle on the provided `tool10.csv` dataset by using the aforementioned technique of EDA. First data from leap years is dropped for better comparability. After this the `.describe` method of `pandas-dataframes` is used to get a summary of the nitrogen dioxide and ozone concentration columns. By this we find mean concentrations and standard deviations of:

$$\text{NO}_2 : 38.4 \pm 23.2$$

$$\text{O}_3 : 37.0 \pm 22.0$$

This makes clear, that the nitrogen dioxide measurement has a larger variability than the ozone measurement. Furthermore, this simple analysis indicates already some erroneous data points. The minimum of the  $\text{NO}_2$  concentration is for example at  $-3$ . This is unphysical and could stem from an erroneous measurement e.g. an error when digitizing the measured concentration. Sometimes measurement devices also deliberately save unphysical values to indicate an error that happened during the measurement. In this case the value of  $-3$  could be understood as an error code. To see whether this is the case, one would need to take a look at the manufacturers data sheet of the specific detector in question. In principle such data points should be excluded from further analysis to not skew the results one gets from the physical data.

units?

units

unit

In a next step the columns corresponding to the particulate matter mass under  $10\text{ }\mu\text{m}$ , carbon monoxide concentration, temperature, humidity and wind speed have been analyzed additionally in a similar manner. From this it becomes clear, that the mass and carbon monoxide concentration show similar unphysical negative values as the nitrogen dioxide concentration. As also these unphysical values occur at perfect integer values, a look into the datasheet or the digitization software of the detectors would be helpful to understand these values and see, if they maybe really correspond to error codes. Of course, filtering this data out is simple, but it is worth to investigate where this erroneous data stems from, what could be wrong during measurement and also how much data is affected. Additionally, there are some interesting insights to be gained by looking at the mean, standard deviation, minimum and maximum values of the data besides faulty measurements. One can for example see, that there is a seemingly small number of very heavy small particles in the analyzed air, as the mean value of the mass is roughly 25 times smaller than the maximum value. Furthermore, one can see

that the carbon monoxide concentration in the air is much smaller than that of any other gas. The temperature and humidity show very large spreads over the dataset that are on the order of a forth of the whole data range. Lastly, one can see that the wind speed is a strictly positive quantity and is therefore not directional. This shows that already a simple analysis requiring only a few lines of code can give valuable insights on the data range, its variability and possible errors during measurement.

Next the correlation between the individual variates is calculated (compare Equation 3) and plotted. This is shown in Figure 1. The strongest positive correlation is present between the

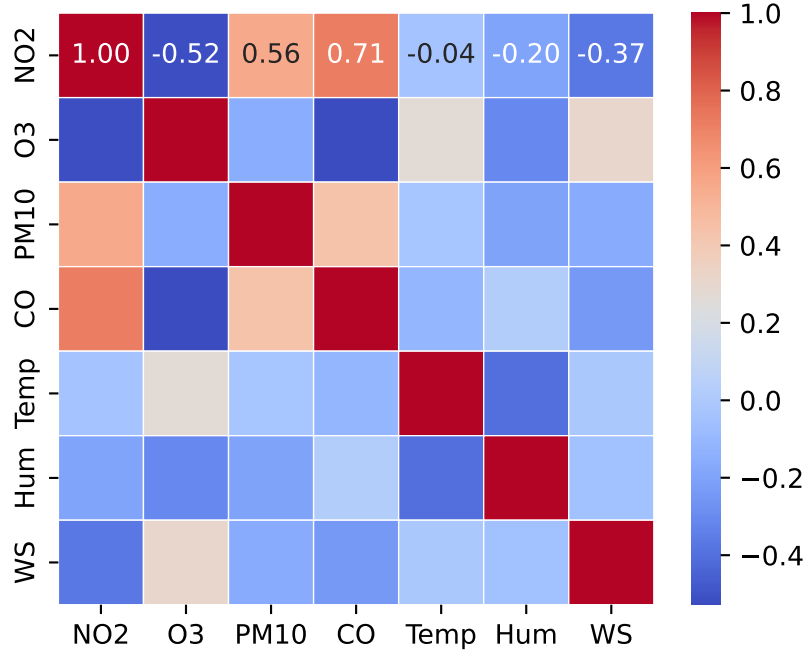


Figure 1: Cross correlation matrix calculated using the mentioned columns of the `todo.csv` dataset.

nitrogen dioxide and carbon monoxide concentration. This indicates, that the processing leading to an increase or decrease of these gases in the atmosphere are strongly linked such that an increase in one is connected to an increase in the other gas. At this point it is important to keep in mind that this analysis shows only correlation but not causality meaning that we cannot conclusively state which change leads to which but only that they are connected. Similarly, an anticorrelation between the concentrations of nitrogen dioxide and ozone can be seen as the correlation coefficient between those two quantities is negative. This indicates that a rise in one of the quantities is linked to a fall in the other. Additionally, a correlation between particle mass under ten microns and nitrogen dioxide concentration is present, suggesting a link between these two. There are also smaller degrees of anticorrelation between both the wind speed and nitrogen dioxide concentration as well as temperature and humidity. While the latter seems reasonable (higher temperatures seem to be connected with lower relative humidities), the prior is quite hard for me to understand. It could be interesting to investigate this further. It is quite fascinating to me, that such a simple analysis, requiring only a few

simple lines of code can already extract interesting research questions that can be expanded on.

In a last task the individual distributions of the quantities was looked at by plotting histograms of the data shown in Figure 2. While temperature as well as the gas concentrations and wind

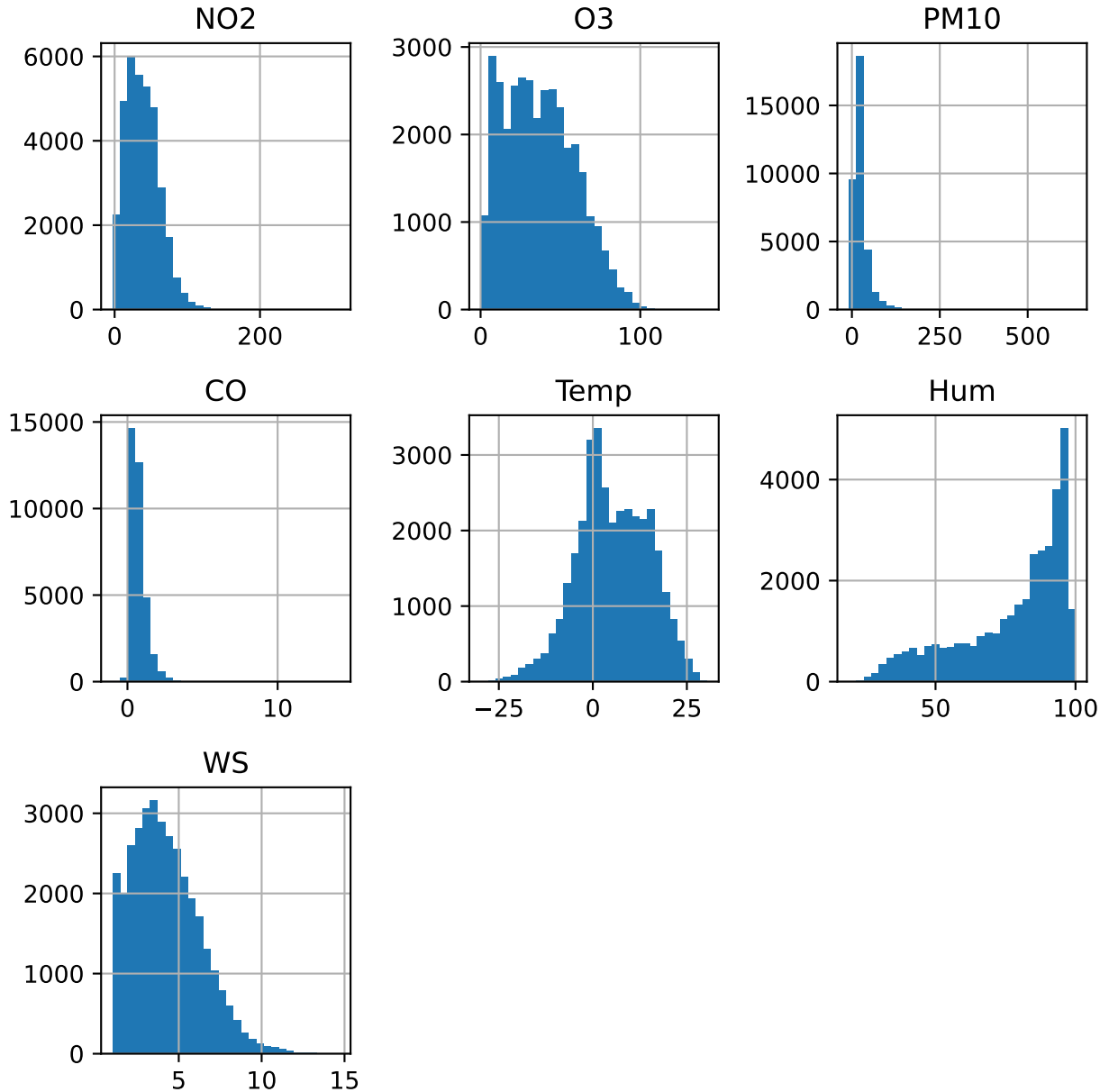


Figure 2: Histogram of the data columns of the `todo.csv` dataset.

speed qualitatively look roughly normal distributed (taking into account, that the distribution is cut due to unphysical values), the humidity and particle masses under ten microns seem quite asymmetric. The histogram of the particle masses roughly fits to the realization I already had when looking at the mean, minimum and maximum values of the taken data in a prior step. It seems there are a lot of particles with low mass which influence the mean heavily, but there also remain a few particles of very high mass. In fact, to me the distribution looks similar to an exponential when postulating that very small masses are underrepresented probably due



to the detector not being able to measure arbitrarily small particulates accurately. Of course, this would still need to be quantified, but it poses an interesting research idea for further analysis. The humidity diagram is heavily skewed towards high humidities.

### Reflection

Although the topics presented this week were mostly nothing new to me, I found applying explorative data analysis quite interesting. It was quite fascinating to see how much knowledge can already be gained from a complex dataset by just looking at statistical properties of the data and its correlation. This is especially true, as I do not have a background in environmental science and therefore discovering correlations between different gases in the atmosphere or temperature and humidity did always come as a surprise and led to the curiosity to read up on the phenomena happening in the atmosphere that could lead to these links. My learning itself was quite well-structured this week. I was successful in writing the learning diary in multiple sessions and doing small updates at a time which also lead to me thinking more about the presented topics and dataset and also re-discovering topics that I already learned in previous courses which fit nicely into this week's lecture. Having a stronger background in statistics I do however think that the meaning of the sample variance and mean as unbiased *estimators* of the true and unknown population quantities could have been motivated more thoroughly and clearly. I feel like this is one of the most important concepts to understand when starting to discuss statistics, as for people having no background so far these definitions must have come out of nowhere. I do however also understand the time constraints of this course, as it is a synthesis of two different courses which were held prior. Additionally, I think it would have been interesting to be required to do some analysis/coding myself for the exercises. The way the exercises are structured now, the entire code is already provided and only the description part is required. I feel like more interesting and maybe more rewarding insights could be gained if students are asked to perform some exploration on their own. This will maybe be the case in future exercises.

## 3. Environmental data and its pre-processing

### Lecture

The focus of this weeks lecture were data pre-processing steps necessary for typical datasets in environmental science. For this we discussed three important steps: Handling missing data, metrics, i.e. a way to compare data and transformations as a way to scale data to similar magnitudes. For each step we talked about the principle behind it and why it is necessary after which a few example algorithms were discussed. In the following I will focus on the idea behind each step and the most interesting and new algorithms to me and will compare their qualities.

The first topic we discussed was missing data. It is clear that efficiently handling these gaps is vital as firstly in any real dataset these are numerous and secondly most further analysis steps cannot handle missing values. These NaN values can stem from measurement errors, human mistakes or instrument failures. The most straightforward although very brute-force way is to just drop any entry of the dataframe, that contains at least 1 or some number of NaN values. This can be done rather easily as can be seen in Listing 2.

```
1 import pandas as pd
2 df = pd.read_csv(foo.csv)          # read some data
3 thresh = 1
4 df = df.dropna(thresh=thresh)      # drop every row with at least thresh NaNs
```

Listing 2: Example for dropping rows with missing values

If one finds out that mainly one variable is problematic, one could also delete this variable (column) from the dataset. Both options have the strong downside of cutting also viable data which generally leads to a loss of predictive power of a resulting model. Another way of dealing with missing values is imputation, meaning filling the missing value with another one based on some algorithm. A few of these are listed here, ordered simple to more complex:

- Imputation with mean, median or a random value  $\Rightarrow$  can alter the distribution of data significantly
- Imputation by linear interpolation  $\Rightarrow$  useful for time-series data
- Nearest neighbor imputation: For each  $N$  dimensional feature vector of the dataframe  $\mathbf{x}_i$  where the index  $i$  represents the row that has  $N_{\text{miss}}$  missing values, we calculate a distance to every other feature vector as follows:  $d_{ij} = \frac{N}{N - N_{\text{miss}}} \|\mathbf{x}_i - \mathbf{x}_j\|^2$ . The missing values are taken from the feature vector with the smallest distance  $\Rightarrow$ . The most interesting thing about this metric for me is that there is no introduction of new values, while also taking into account the reliability of data by making vectors with many NaNs have a large distance.

It is important to note that the choice of imputation method can heavily influence the model so taking care is needed. A good choice is to use univariate methods like interpolation for short gaps and multivariate approaches like the nearest neighbor imputation for longer gaps

in the data.

The next topic concerned metrics, a way to compare how similar or dissimilar two entries of a dataframe are. Metrics are a quite general mathematical concept allowing many different forms with different usecases, however here only a few examples will be named. If the data entries are numeric a simple measure of similarity is the dot product between vectors, while a measure of dissimilarity could be the distance between them. An interesting distance measure has already been shown above, but a very general and often used one is the so called Minkowski distance

$$d(x, y) = \left( \sum_{i=1}^N |x_i - y_i|^p \right)^{1/p} \quad (8)$$

which reduces to the normal Euclidean distance for  $p = 2$  but also has interesting cases for  $p = 1$  (Taxicab/Manhattan distance) or  $p \rightarrow \infty$  (Chebyshev distance). These are illustrated in Figure 3 for a 2-dimensional case. Additionally, there are also many different metrics used

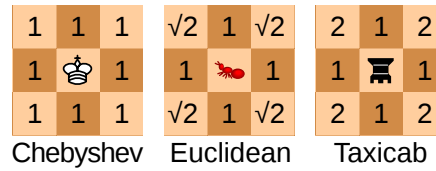


Figure 3: Different cases of the general Minkowski distance illustrated in 2D with the analogy of chess pieces. Taken from [2].

for class-like or binary data.

The last topic we covered were transformations which can affect single values, rows or columns of the dataset or even the dataset as a whole. Column transformations are especially important, as different data columns are usually very different in the magnitude of their entries. If no normalization is done, large entries will dominate small ones, which is usually not wanted. Additionally, columns could contain discontinuous data. Some useful and interesting transformations include:

- Discontinuous, cyclic data (like hours of the day) can be handled efficiently by expressing it by its sine and cosine components leading to a range of values between -1 and 1. The additional cosine is needed to distinguish between the rising and falling edge of the sine function. This seems very useful, as it normalizes the data as well.
- Logarithmic pre-scaling of columns can be useful if the data has a large spread over multiple orders of magnitude. This also leads to a rough normalization but alters the distribution.
- Variance scaling, i.e. shifting the data columns by its mean and normalizing by its variance is also a standard technique. It could prove especially useful if the data can be assumed to be Gaussian, as we then expect a standard normal distribution afterward. It does suppress outliers, however.

- Equalization forces the data to be in the range of  $(-1, 1)$  but can lead to pronounced outliers. This is a problem I already faced in a previous machine learning project of mine. The `sklearn` documentation provides a great comparison [3].

The main take-away from this lecture and my prior experience with machine learning and analysing datasets is the following: It is crucially important to think about all pre-processing steps that one performs on their dataset. They can have a huge impact on the statistics of the data after the processing, the model one gains and especially the machine learning efficiency and features that are learned by the neural network.

## Exercise

## Reflection

## 4. Data visualization

## 5. Clustering in Python

## 6. Predictive modeling

## 7. Summary



## 8. Self-evaluation

## References

- [1] Stephen Weybrecht. *Data Mining in Environmental Science: Repository*. <https://github.com/stewey0/uef-data-mining>. Source code. Accessed: November 22, 2025.
- [2] Cmglee. *Comparison of Chebyshev, Euclidean and taxicab/Manhattan distances for the hypotenuse of a 3-4-5 triangle on a chessboard [image]*. [https://commons.wikimedia.org/wiki/File:Minkowski\\_distance\\_examples.svg](https://commons.wikimedia.org/wiki/File:Minkowski_distance_examples.svg). Image cropped. Licensed under CC BY-SA 4.0. Accessed: November 22, 2025. n.d.
- [3] The scikit-learn developers. *Compare the effect of different scalers on data with outliers*. [https://scikit-learn.org/stable/auto\\_examples/preprocessing/plot\\_all\\_scaling.html](https://scikit-learn.org/stable/auto_examples/preprocessing/plot_all_scaling.html). sklearn documentation. Accessed: November 22, 2025.