



Data Science in Life Science

SS20

Quentin Quarantino



Copyright © Quentin Quarantino

Licensed under the Creative Commons Attribution-NonCommercial 3.0 Unported License (the “License”). You may not use this file except in compliance with the License. You may obtain a copy of the License at <http://creativecommons.org/licenses/by-nc/3.0>. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

First printing, August 2019

Contents

1	Overall Introduction	9
1.1	Background	9
1.2	Methods	9
1.3	Results	10
1.4	Discussion	10
 I Project 1: Relevant Research Topics for COVID-19		
2	Introduction	13
2.1	Background	13
2.2	Goal of the Project	13
2.3	Outcome	13
3	Topic 1: Spreading Models	15
3.1	Background	15
3.2	Data and Methods	15
3.3	Results and Discussion	17
4	Topic 2: Data-based Time Series Prediction	19
4.1	Background	19
4.2	Data and Methods	19
4.3	Results and Discussion	20

5	Topic 3: Risk Factor Analysis	21
5.1	Background	21
5.2	Data and Methods	21
5.3	Results	21
5.4	Discussion	22
6	Topic 4: Diagnostic	23
6.1	Background	23
6.2	Data and Methods	23
6.3	Results	23
6.4	Discussion	24
7	Topic 5: Origin Analysis	27
7.1	Background	27
7.2	Data and Methods	27
7.3	Results	27
7.4	Discussion	28

II Project 2: Literature Clustering and Visualization

8	Introduction	31
8.1	Background	31
8.2	Goal of the Project	31
8.3	Outcome	31
9	Part 1	33
9.1	Redo and Understand an Existing Clustering Approach	33
10	Part 2	37
10.1	Add Related Papers to the CORD-19 Data Set	37
10.2	Redo the Clustering	37
11	Part 3	43
11.1	Loading in the Student Data Set	43
11.2	Preprocessing, Clustering and Results	43
11.3	Creating a word cloud	48

III Project 3: System Dynamics Modeling using a SIR Model

12	Introduction	51
12.1	Background	51

12.2	Goal of the Project	51
12.3	Outcome	51
13	The Model	53
13.1	A simple SIR model	53
13.2	Extending the SIR model	55
13.3	Parameter Fitting	57
14	Scenario Studies	59
14.1	Implemented Measures	59
14.2	Methods	60
14.3	Results	60
14.4	Discussion	60

IV

Project 4: Agent-based Simulation

15	Introduction	65
15.1	Background	65
15.2	Goal of the Project	65
15.3	Outcome	65
16	Basic Principles	67
16.1	Agent-based Modeling for COVID-19 Spreading Simulations	67
17	A simple ABM	69
17.1	Model Description	69
17.2	Fitting to Real Data	71
18	Extending the Model	73
18.1	Introduction	73
18.2	Incubation and Exposed State	73
18.3	Chronic Conditions and Comorbidities	73
18.4	Central Locations	74
18.4.1	Supermarkets	74
18.4.2	School	75
19	Scenario Studies	77
19.1	Methods	77
19.2	Results	78
19.3	Discussion	84

20	Comparison of EBM and ABM	85
20.1	Fundamental differences	85
20.2	Case Study Covid-19	86

V

Project 5: Time-Series Prediction for COVID-19 Cases

21	Introduction	89
21.1	Background	89
21.2	Goal of the Project	89
21.3	Outcome	89
22	Model-based vs Data-based	91
22.1	Model-based Approaches	91
22.2	Data-based Approaches	91
23	Approaches	93
23.1	Prophet Library	93
23.2	LSTM Neural Networks	96
23.3	Classical Models	97
24	Comparison of Data-based TSP Approaches	99
24.1	Data	99
24.2	Results and Discussion	100
25	Model- vs. Data-based Time-Series Prediction	103
25.1	Fitting the Model	103
25.2	Results for fitted Model	104
25.3	Discussion	104
26	Towards COVID-19 Outbreak Prediction	107
26.1	Outlook and Evaluation	107

VI

Project 6: Time-series Prediction for COVID-19 Cases II

27	Introduction	111
27.1	Background	111
27.2	Goal of the Project	111
27.3	Outcome	111
28	Solution Approaches	113
28.1	Data	113

28.2	Visual Exploration	113
28.3	Time-Series Prediction via Prophet	114
28.4	Clustering	114
29	Results	115
29.1	Visual Exploration	115
29.2	Time-Series and Prophet Prediction	121
29.3	Clustering	122
30	Evaluation	129
30.1	Project Rating	129
30.2	Problems	129

VII

Project 7: Origin Analysis of COVID-19

31	Introduction	133
31.1	Background	133
31.2	Goal of the project	133
31.3	Outcomes	133
32	Methods for Phylogenetic Analysis	135
32.1	Data	135
32.2	Methods	136
32.2.1	Hierarchical Approaches	136
32.2.2	Non-Hierarchical Approaches	137
32.2.3	Phyldynamics	139
32.2.4	Phylogeography	139
33	Analysing the Spread of SARS-CoV-2	141
33.1	Results	141
33.1.1	Origin Analysis	141
33.1.2	Phyldynamics and Phylogeographics	142
33.1.3	TreeTime	143
33.2	Discussion	146
33.2.1	Origin Analysis	146
33.2.2	Phyldynamics and Phylogeographics	146
33.3	Conclusion	147

VIII

Project 8: Sequence-curve-based Phylogeny Analysis

34	Introduction	151
34.1	Background	151
34.2	Project Description	151

34.3	Outcomes	151
35	Solution Approach	153
35.1	Data	153
35.2	Methods	153
35.2.1	2D Graphical Representation of genomic Sequences	153
35.2.2	Building phylogenetic Trees based on 2D Curves	154
35.2.3	Metrics to compare phylogenetic Trees	155
36	Results and Discussion	157
37	Evaluation	163
37.1	Project Rating and Problems	163

IX

Project 9: Drug Repurposing

38	Introduction to Drug Repurposing	167
38.1	Background	167
38.2	Goal of the Project	167
38.3	Outcomes	167
39	Methods for Drug Repurposing	169
39.1	Computational Methods	169
39.2	Experimental Methods	170
39.3	Neural Networks that predict Drug-Target Interactions	170
39.3.1	DeepDTA	170
39.3.2	MT-DTI	171
40	Find potential Drugs for treating COVID-19	173
40.1	Methods	173
40.2	Results and Discussion	174
40.3	Conclusion	174

X

Appendix

41	Weekly Member Contribution	177
42	Bibliography	181
	Bibliography	181
	Articles	181
	Books	183
	Webpages	183



1. Overall Introduction

1.1 Background

Within the scope of the university course "Data Science in the Life Sciences" different aspects of the global COVID-19 pandemic were examined and analyzed. Each week of the semester, the focus was put on another research field to gain a complete overview of the possible analysis and application techniques.

1.2 Methods

In the first few weeks, the foundation was laid by summarizing the main research topics in the world of virology and epidemiology. Additionally, the already existing literature related to the novel virus was clustered according to their topics and the results were presented in an interactive two-dimensional t-SNE plot.

Furthermore, an equation-based model (EBM) as well as an agent-based model (ABM) were created to predict the spread of the disease. The basic EBM was extended by additional states (exposed, dead) and population characteristics (age, gender, smoking status, ICU bed capacity) and afterwards fitted to actual case data of Germany. The ABM was also refined by a set of extensions (incubation and exposed state, chronic health conditions, comorbidities and central locations). Finally, different scenarios (reducing social contacts, lockdown, wearing masks) were implemented and tested according to their effectiveness on the possible prevention of the spreading activity.

Based on the calculated EBM and ABM prediction results as well as actual German case data, a forecasting of the COVID-19 outbreak was performed using a classical approach (ARIMA), the *prophet* library and a long short-term memory network (LSTM). A comparison between the forecasting performances was evaluated and based on the *prophet* results an exploratory visual data analysis was made.

To identify the most recent common ancestor of the human SARS-CoV-2 genome, it was compared with six other virus sequences of the *Coronaviridae* family from different non-human hosts. In addition, a phylogenetic analysis with 14 other human SARS-CoV-2 sequences of different countries was performed with the aim to identify possible pathways of the global virus spread. In

the subsequent weekly project classical hierarchical methods (UPGMA, TreeTime) were compared to a graphical approach (pyrimidine-purine graph) for phylogenetic analysis.

Finally, an overview of the available methods for drug repurposing was compiled and an attempt was made to discover potential drugs for the treatment of COVID-19 using deepDTA and MTI-DTA.

1.3 Results

Both the EBM and the ABM could be significantly improved by the implemented extensions. The addition of central locations (schools and supermarkets) in the agent-based simulation had the largest impact on the results. By applying different intervention strategies, the combination of social distancing and wearing masks has been confirmed to be the most effective.

For the time-series forecasting LSTM turned out to be the clear winner and produced the lowest root mean square error (RMSE) for the prediction on the real-life data set ($RMSE = 0.007$) and the SIR data set ($RMSE = 0.0$) while *prophet* achieved the best prediction results on the ABM data set ($RMSE = 0.0155$). The further comparison of data-based and model-based approaches (LSTM vs SIR model) also confirmed LSTM to be the better option for the given data.

The *Rhinolophus* (horseshoe bat) was identified to be the most probable common ancestor among the six *Coronaviridae* sequences of different hosts with a sequence similarity of 94%. The phylogenetic analysis of the 14 human samples showed a clear spatial pattern. Thus leading to the identification of reasonable infection pathways. While TreeTime convinced with the best possibility to estimate the temporal transmission patterns, UPGMA showed the most accurate results in terms of sequence similarity calculation, especially on very similar sequences, where the graphical pyrimidine-purine approach tended to show inaccurate performances.

Discussing the drug repurposing analysis of Beck et al.[5], Atazanavir ($K_d = 94.94\text{ nM}$) was determined to be the best fitting chemical compound in terms of inhibitory potency against the SARS-CoV-2 3C-like proteinase.

1.4 Discussion

Overall, it can be said that the conducted analyses delivered plausible results. However, since the experiments were carried out with small sample sizes and had mainly representative character to understand the analysis principles, the results should be considered cautiously. Nevertheless, there has definitely been a huge increase in knowledge among the students of this group.

Project 1: Relevant Research Topics for COVID-19

2	Introduction	13
2.1	Background	
2.2	Goal of the Project	
2.3	Outcome	
3	Topic 1: Spreading Models	15
3.1	Background	
3.2	Data and Methods	
3.3	Results and Discussion	
4	Topic 2: Data-based Time Series Prediction	19
4.1	Background	
4.2	Data and Methods	
4.3	Results and Discussion	
5	Topic 3: Risk Factor Analysis	21
5.1	Background	
5.2	Data and Methods	
5.3	Results	
5.4	Discussion	
6	Topic 4: Diagnostic	23
6.1	Background	
6.2	Data and Methods	
6.3	Results	
6.4	Discussion	
7	Topic 5: Origin Analysis	27
7.1	Background	
7.2	Data and Methods	
7.3	Results	
7.4	Discussion	



2. Introduction

2.1 Background

The current pandemic is globalized, with severe consequences to social, health, and economic order. As of 11th of May 212 countries are affected, with a total of 4,215,274 confirmed cases and a death toll of 284,672 people [59]. Analyzing Covid-19 in every scientific angle is crucial towards forming new policies to combat this pandemic.

2.2 Goal of the Project

In this week's project, each group member was assigned one overarching topic pertaining to the current COVID-19 pandemic, highlighting the main concepts, methods and results within each.

2.3 Outcome

Within each topic a short introduction to the general concept is given. The understanding of these concepts is then deepened by real world code examples, showing a glimpse of what is possible in each topic in regards to COVID-19.



3. Topic1: Spreading Models

3.1 Background

Modeling the spread of infectious diseases is not only an essential tool in understanding the transmission rates and the trajectory of future cases but also has a significant influence on the appropriate guidelines to control the course of an pandemic. The approaches towards modeling the spread can range from computational models (e.g. agent-based) to mathematical modeling (e.g. compartmentalized models). In this short introduction we will focus on a SIR-Model which is part of the compartmentalized subgroup. These models follow a deterministic pattern where each subpopulation is divided into groups. In SIR-Models each letter stands for one group: $S = \text{susceptible}$, $I = \text{infectious}$ and $R = \text{recovered/death}$. Then it follows that for each time independent point t the rates for each subgroup can be calculated by differential ordinary functions (ODEs):

$$dS/dt = vN - \beta SI/N - \mu S \quad (3.1)$$

$$dI/dt = \beta SI/N - \gamma I - \mu I \quad (3.2)$$

$$dR/dt = \gamma I - \mu R \quad (3.3)$$

with γ denoting the time rate of death/recovery, β denoting the number of new infections one case causes per time point t , μ denoting general death rate and v denoting the birthrate.

3.2 Data and Methods

The code example is based on the kaggle notebook created by the user Lisphilar [32]. It uses python and a basic framework of libraries e.g *pandas*, *sklearn*, *datetime*. The main data used is from

the World Health Organization (WHO) showing novel corona infections by country. Furthermore supplementary data is used to include the age pyramid for each country. The WHO data set is preprocessed to include the variables: Date, Country, Province, Confirmed, Infected, Deaths and Recovered. A first visualization shows the global rate of infected, deaths and recovered people (Figure 3.1). Next, the growth factor is calculated, which is given by: G_n/G_{n-1} with $G = \text{confirmedcases}$. Countries with growth factor higher then one have an increasing number of cases. In contrast, a growth factor less then one shows a declining number of cases. The actual analysis is done for five countries: Italy, Japan, India, USA and New Zealand. Giving one case

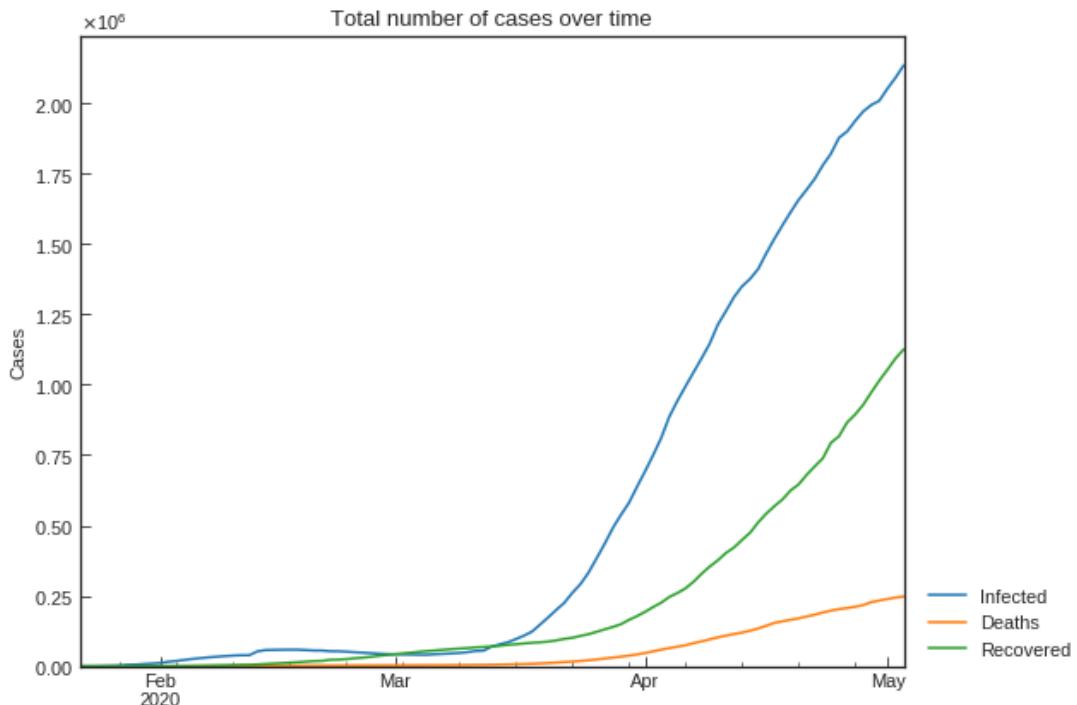


Figure 3.1: Global infections, deaths and recovered people for COVID-19. Infected people follow an exponential trajectory. Recovered people follow a delayed increase, which is due to the long incubation and illness period of 14 days.

as an example as a first step a S-R trend is plotted (Figure 3.2). It shows the trend of susceptible against recovered people. Five change points can be identified. Next the SIR-F model parameters are estimated for each change point. As a last step the changes in the p value are contrasted with measures taken by the country. While these results are interesting SIR modeling also has its limitations.

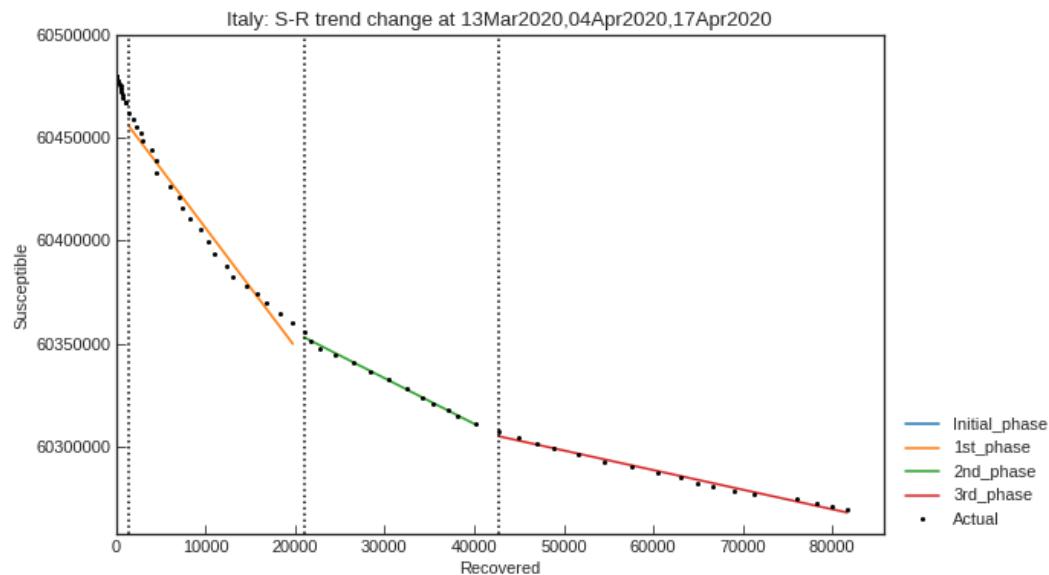
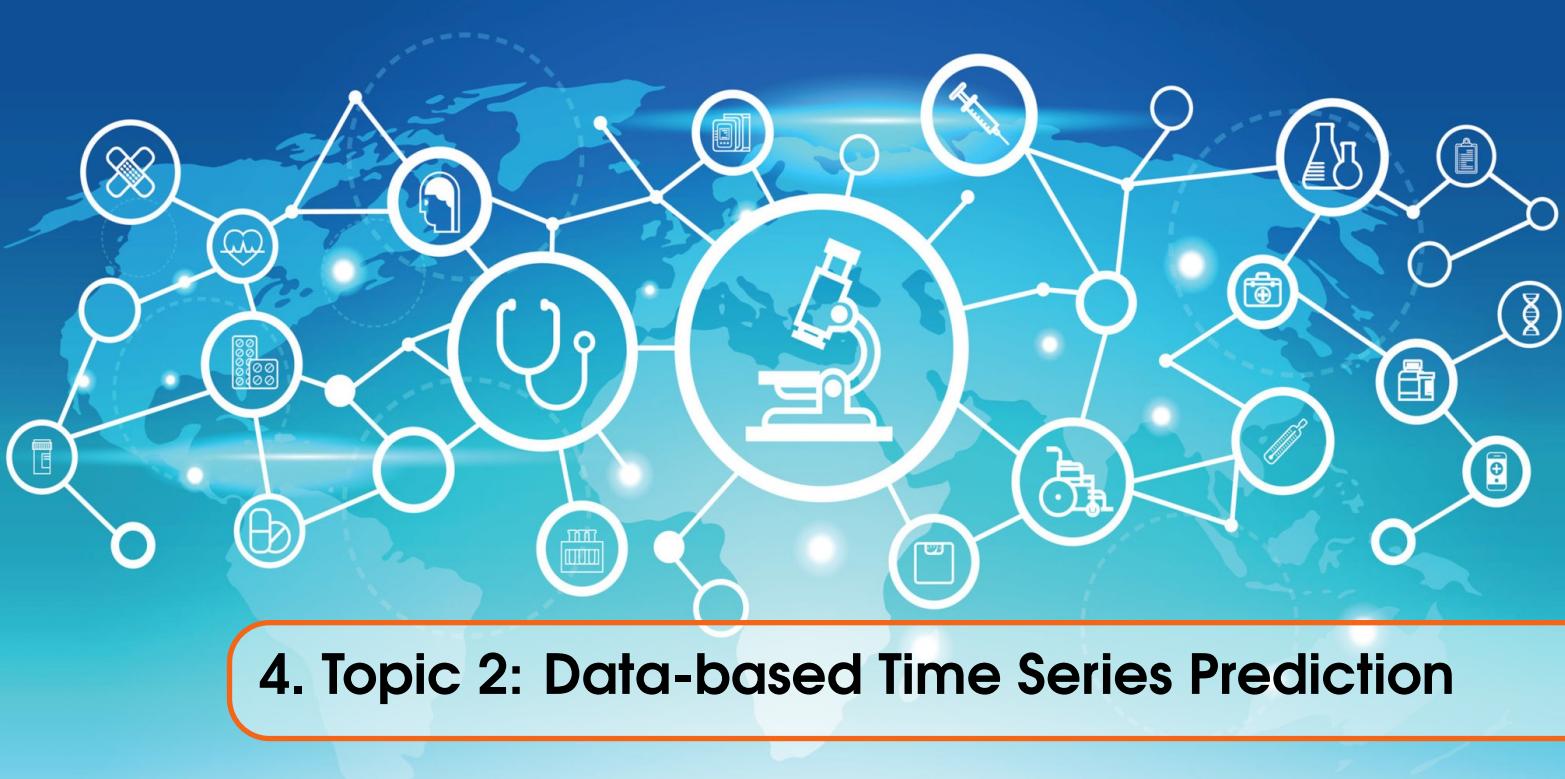


Figure 3.2: Trend of susceptible people versus recovered. Three distinct change points can be identified.

3.3 Results and Discussion

Even though the SIR-Model is one of the most basic infectious disease models available, it can show promising results with careful consideration for parameter selection and data processing steps. In the case of Italy three measures could be shown to reduce the p value: quarantine of persons contacted with positive patients, school closure and lock-down. It also has to be considered that the SIR model is based on very basic assumptions. For example the number of susceptible people needs to be defined before simulation.



4. Topic 2: Data-based Time Series Prediction

4.1 Background

Many governments around the world are building their political decisions around the number of current confirmed cases of people infected by COVID-19. Nonetheless, not only the current number of confirmed cases is of greater interest, but also how the virus spreads in the future. One way of forecasting the spread of the virus is by using data-based time series prediction.

4.2 Data and Methods

Machine learning models are calibrated using publicly available data sources like the WHO health report. Time series forecasting can be framed as a supervised learning problem. Other than model-based spreading simulation such as the SIR model, this approach does not simulate a population. The forecasting is performed using pythons *numpy* and *sklearn* libraries. At first, the data is downloaded. Since no data points are missing no preprocessing is performed with the exception of converting integers into date times and reorganizing dataframes. The data contains a wide range of countries with the number of infected people per day starting January 22. A support vector machine model is implemented to forecast the number of infected people. The parameters that have been set can be seen in Figure 4.1 in line two. The test and test training data sets are generated by splitting them without shuffling them, such that the time series is preserved. The model has been trained

```
[ ] 1 # svm_confirmed = svm_search.best_estimator_
2 svm_confirmed = SVR(shrinking=True, kernel='poly', gamma=0.01, epsilon=1, degree=5, C=0.1)
3 svm_confirmed.fit(X_train_confirmed, y_train_confirmed)
4 svm_pred = svm_confirmed.predict(future_forcast)
```

Figure 4.1: Initial parameter setup for the SVM Model.

using the first 75 days since January 22. In Figure 4.2 it can be seen that the model overestimates the number of confirmed infections by over 1.5 Million.

4.3 Results and Discussion

The results (Figure 4.2) are quite underwhelming. Pandemic infection curves usually increase at the beginning exponentially but then start to flatten e.g. because of restrictions in the society to decrease the spread of the disease. The model was trained using data at the pandemic's start where the number of cases rapidly grew. Based on this fact the estimated number of infected people massively exceeds the confirmed case numbers. Nonetheless, due to the lack of test capacities around the world, the estimated number of cases might be closer to the real number than it initially seemed. Still, the used model was quite simple and no testing was performed how the parameters were estimated. A more complex model may give a better insight to the spread of the virus.

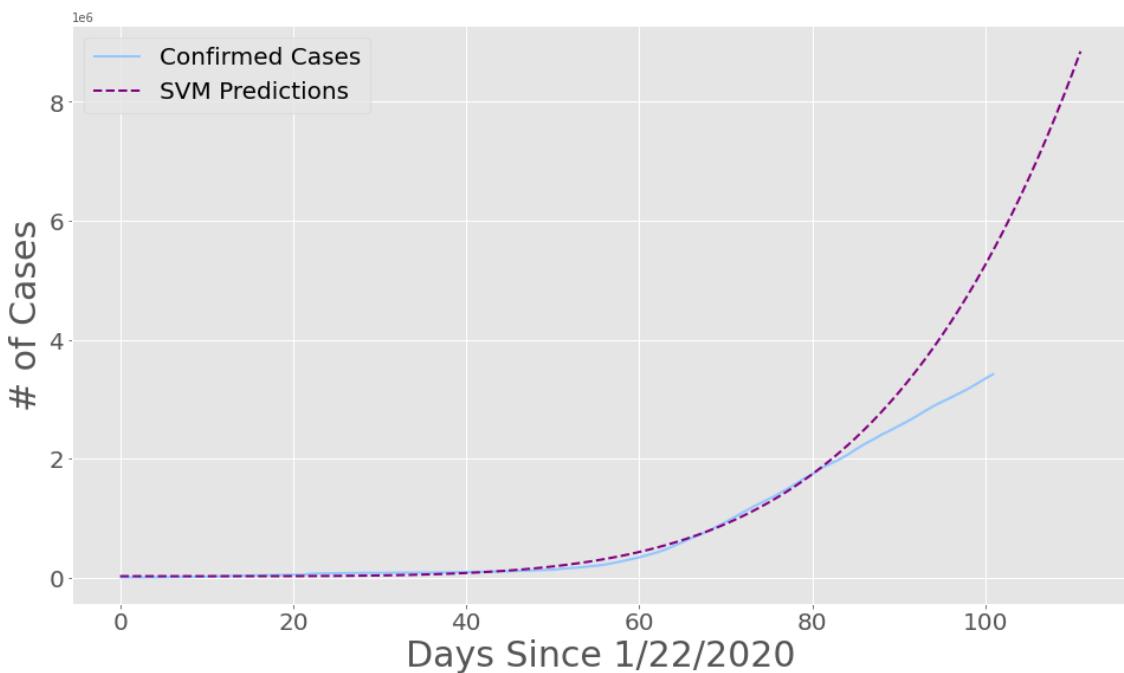


Figure 4.2: Comparison between the observed and the estimated number of infected people. The SVM prediction overshoots the number of cases in contrast to actual confirmed cases.



5. Topic 3: Risk Factor Analysis

5.1 Background

Many studies have shown the dangers the world is facing due to the recent 2019-ncov outbreak. The spread of the disease is related to the number of infections. Especially for undiagnosed infections it is crucial to estimate their numbers. In order to better assess the epidemic risk, there are two key parameters which come into the picture. The first is the basic reproduction number represented by R_0 and the second is the incubation period [42]. Initially, the cumulative number of cases were estimated in China outside the Hubei province after January 23. For that the time-dependent compartment model is used. We take that number as an input which contains the transmission dynamics to the global transportation network to generate probability distributions of the travellers who arrives outside China to their respective destinations. The usage of the Galton-Watson branching process makes the task easier by modeling the initial spread of the virus.

5.2 Data and Methods

The analysis is performed using the python libraries *numpy*, *matplotlib*, *kiwis solver*, *scipy* and *cycler*. The risk analysis is based on three key parameters which mainly contribute to the risk: the cumulative number of cases which arise due to the influence of the areas of China where the lockdown is not yet applied, the connectivity between China and its destination countries and finally the areas which are low connected to China with high local transmission potential of the virus. We computed the risk of the individual countries with the selected possible parameters like connectivity and Rloc, where Rloc is the local reproduction number of the infection. Finally, the collected data of the neighbouring regions of China was also considered. Based on the different variables, we generated a heat map.

5.3 Results

The generated heat map provides the information about the outbreak risks as functions of Θ and Rloc, when C=200,000. The white arrows in Figure 5.1 indicate the directions corresponding to the largest reductions in the risk.

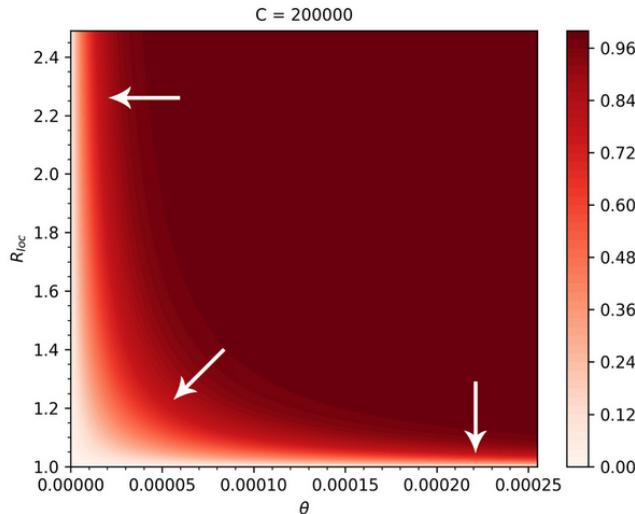


Figure 5.1: Heatmap of the outbreak risks as functions of Theta and Rloc [42]. The arrows represent the directions corresponding to the largest reductions in the risk.

5.4 Discussion

We can say that the three modeling approaches combined give us new information for assessing the risk of SARS-CoV-2 outbreaks in countries outside of China. However, these key parameters have a relatively high Rloc value. To tackle this kind of situation, suitable control measures need to be taken. Some of the most beneficial one include entry screening or travel restrictions. It is insufficient to only know the Rloc and the generation interval to effectively quantitate the risk estimation. However the idea of which types of control measures need to be implemented in order to lessen the risks play a significant role in risk assessment.



6. Topic 4: Diagnostic

6.1 Background

Combining medical diagnostics with computational approaches can help to detect COVID-19 infections automatically, which may lower error rates and speeding up the overall diagnostic process. One of such methods is based on RT-PCR-analysis, which are not very secure due to a high number of false positives.

6.2 Data and Methods

For our analysis the publicly available IEEE data set [12] was used, which contains 906 X-ray chest images provided by hospitals and physicians around the globe. The images show chest from healthy patients but also patients who have an bacterial or viral pneumonia. The metadata file was used to only filter those images who have been infected with corona and those from healthy patients. The X-ray detection method reproduced here is done by training a convolutional neural network in python with the help of *TensorFlow* and *Keras* using X-ray images (Figure 6.1) from the IEEE data set to predict whether a patient is infected with corona or not. As a first step X-ray data and python scripts were downloaded from the kaggle repository by Nabeel Sajid [46]. In the next step, a virtual environment was created to setup the package requirements for the analysis pipeline. The filtered data set contained only 70 infected and 28 normal X-ray images. The small sample size might be problematic since convolutional neural networks (CNN) require a wide range of samples to extract meaningful features. Thus, data augmentation was performed to increase the number of Covid-19 related images to 5088 and the number of non COVID-19 related X-ray images to 2424.

6.3 Results

The final analysis was performed on a train and test data set by splitting the augmented data 80/20. The model was trained for 100 epochs with a batch size of 128 images. The fitted model was then validated on the test data. In Figure 6.2 the training and validation accuracy and loss can be seen.

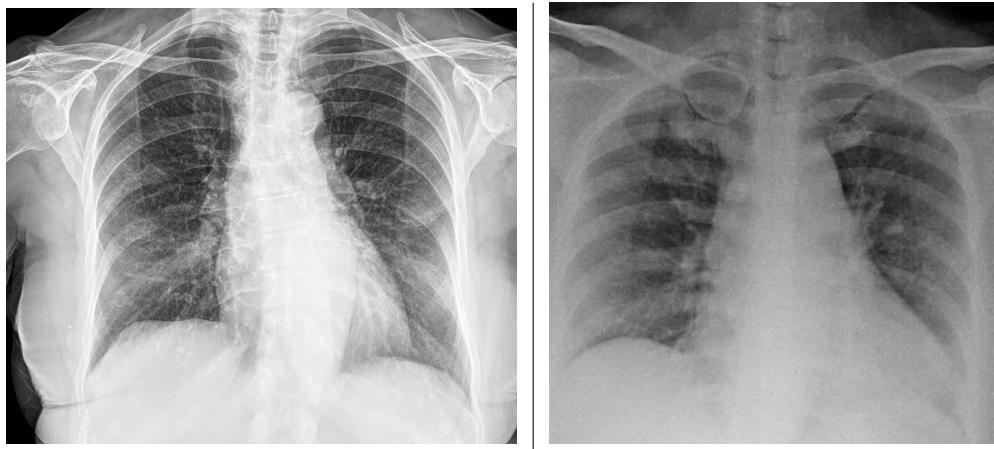


Figure 6.1: Exemplary X-ray images of a patient having an infection of SARS-CoV-2 (left) and a non-infected patient (right). A white obstruction within the left image implicates the corona infection.

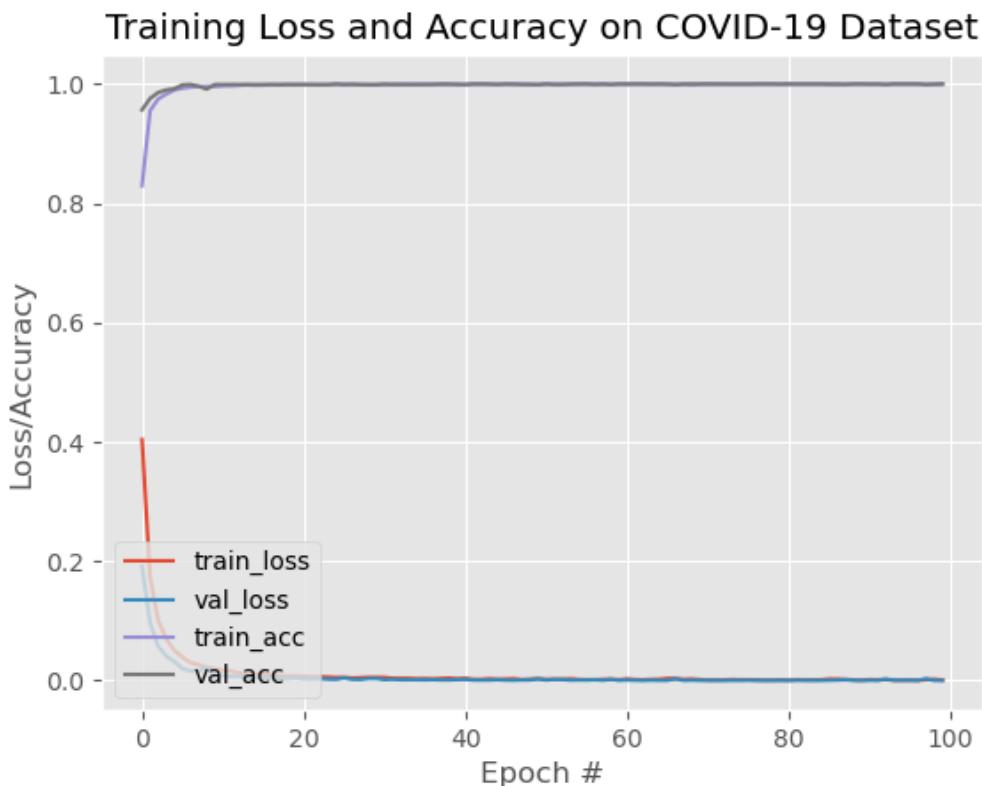
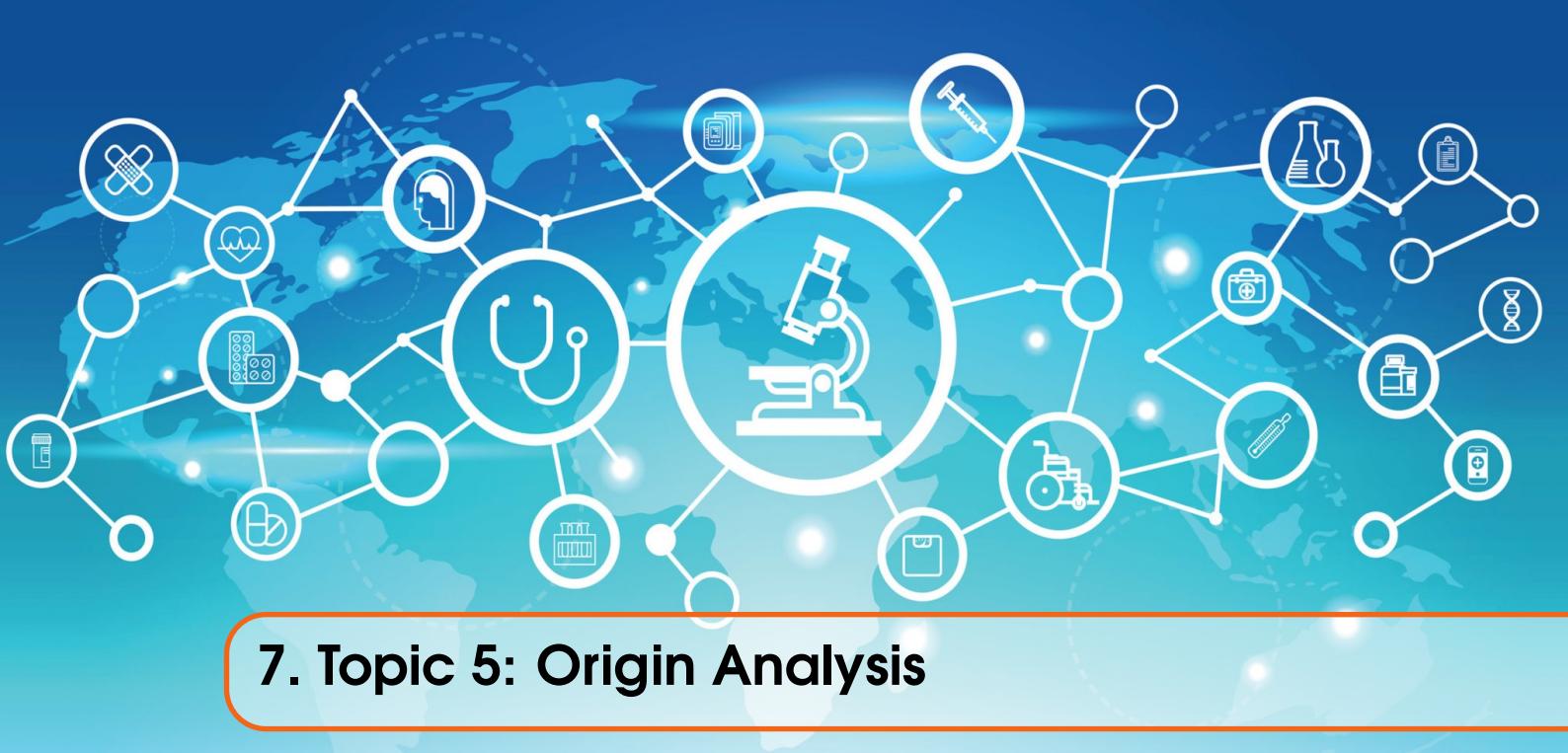


Figure 6.2: Plot of validation loss, training loss and accuracy. After 10 epochs a limit is reached

6.4 Discussion

The described CNN approach is a promising tool for COVID-19 detection in the lungs. However, it is very time-consuming. Augmentation of the images took 2.5 hours and training the model took another 40 hours. As seen in Figure 6.2 the accuracy of the model for both the train as well as the test data set is already maximized after just a few epochs. The same can be observed for the

loss that converges to zero also after just a few epochs. While on the first view this might look like a perfect result some things need to be taken into consideration. First, we just classified chest images from persons who are infected with the coronavirus and from those who are completely healthy. It would be devastating if the model wrongly classifies a patient as COVID-positive only caused by the presence of any other pneumonia. Furthermore, the model might be overfitted and not even suited for many COVID-19 cases since the samples used for training the model originated from only a few samples. Still, the shown analysis has a lot of potentials when also trained to other pneumonia with more available samples.



7. Topic 5: Origin Analysis

7.1 Background

Phylogenetic analysis aims to reconstruct phylogenies both for a group of species and also for the individuals within those species. For pathogens, we are interested in the evolution of the diseases at the genetic level, and what this can tell us about their past and present diversity [56]. An important fact about the *Coronaviridae* family is that its members tend to “jump” from one species to another. When the transmission occurs from a non-human host to a human host it is called zoonosis [18]. The determination of the most recent common ancestor of the human SARS-CoV-2 and the zoonotic transmission can provide important information about biological features, key mutations and properties of the virus. A detailed understanding of how an animal virus jumped species boundaries to infect humans will help in the prevention of future zoonotic events. [2].

7.2 Data and Methods

We will compare the genetic sequence of SARS-CoV-2 with other viruses of the *Coronaviridae* family in different hosts. The following analysis is based on a Github repository of Simon Burgermeister [11]. Six complete genomes were considered, whose names and hosts are listed in Table 7.1. The sequence data (fasta files) were downloaded from the NCBI Virus public library [22]. To compare the genetic sequences, a multiple sequence alignment needed to be performed. Clustal Omega is a software that uses seeded guide trees and HMM profile-profile techniques to generate alignments between multiple sequences. Unfortunately, my local computer was not able to compute the alignment due to RAM exceedance. Therefore, I submitted a request to the online version of Clustal Omega [36]. Based on the resulting alignment, a distance matrix was calculated with the *TreeConstruction* package from Biopython. Afterwards, the same package was used to create the phylogenetic tree base on the UPGMA algorithm.

7.3 Results

The resulting phylogenetic tree (Figure 7.1) shows that our human SARS-CoV-2 sequence is most similar to the SARS-like coronavirus sequence of the Rhinolophus (horseshoe bat) with a similarity

Accession number	Host	Description
MN996528	H. Sapiens	Human SARS-CoV-2
NC_019843	H. Sapiens	Human MERS-CoV
JQ065048	Anatidae	Ducks, geese and swans
MG772934	Rhinolophus	Horseshoe bats
NC_034972	Apodemus chevrieri	Rodent
KX38909	Gallus gallus	Chicken
MT084071	Manis javanica	Pangolin

Table 7.1: Considered *Coronaviridae* strains and hosts.

of 96%. The host with the next similar sequence is the *Manis javanica* (Pangolin) with a similarity of 0.89% between their genomes. The human MERS-Cov genome and the SARS-CoV-2 genome share only a sequence similarity of 0.74%.

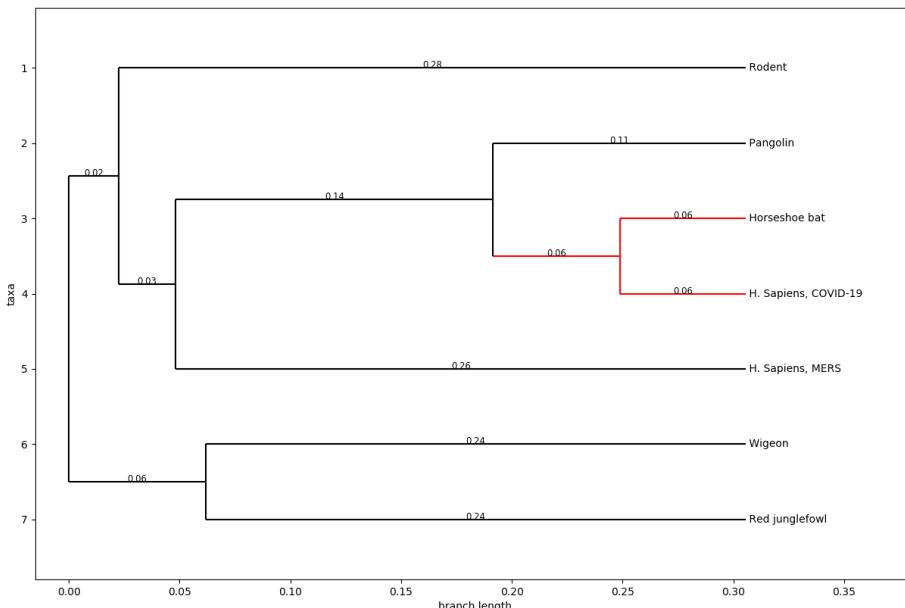
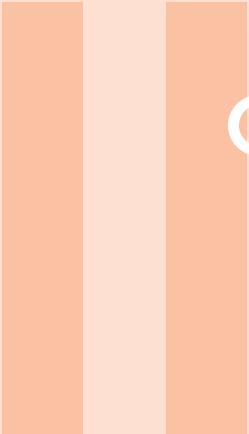


Figure 7.1: Phylogenetic tree of the origin detection analysis.

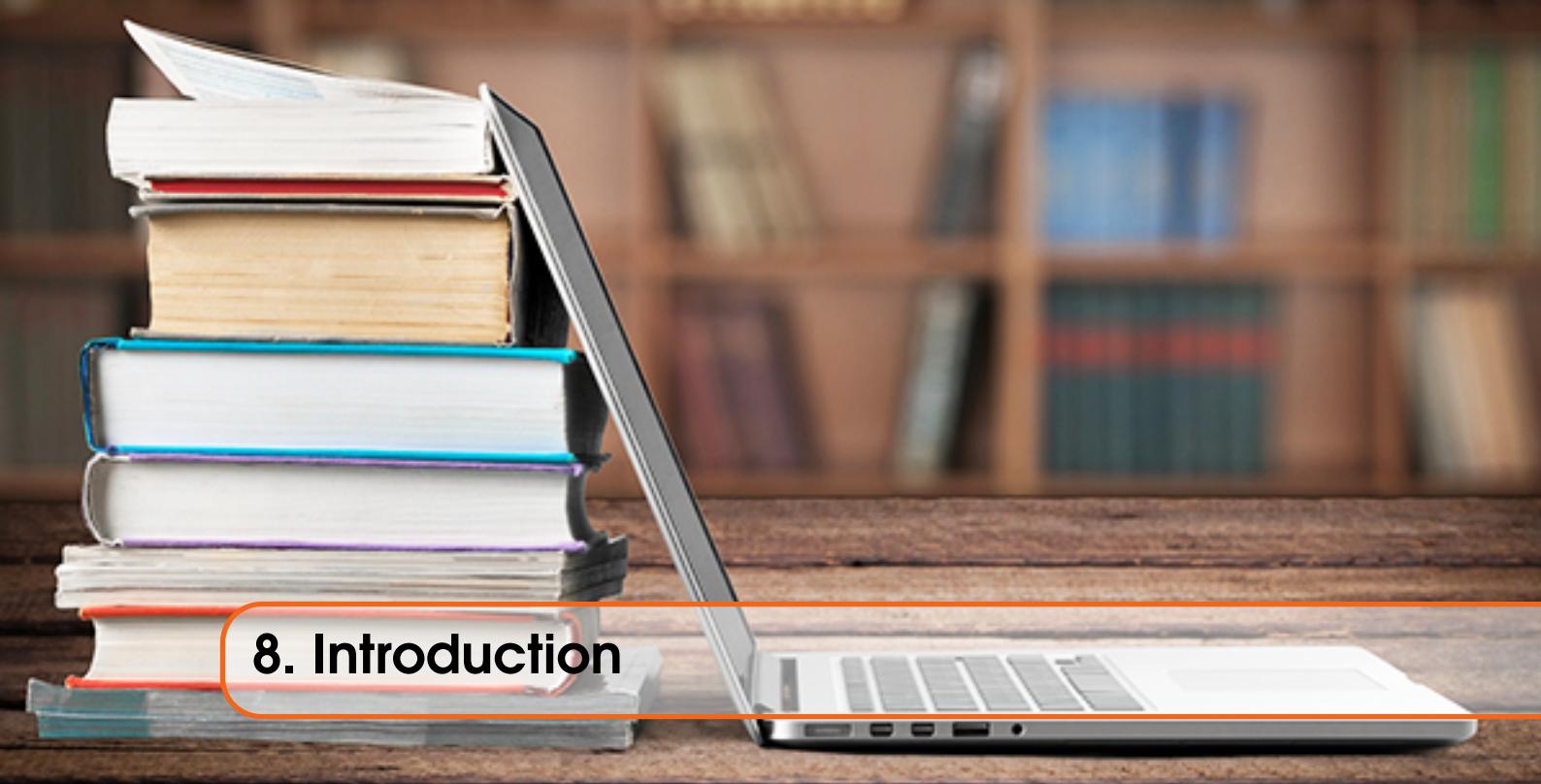
7.4 Discussion

As many early cases of COVID-19 were linked to the Huanan market in Wuhan [61], it is possible that an animal source was present at this location. Given the similarity of SARS-CoV-2 to bat SARS-CoV-like coronaviruses, it is likely that bats serve as reservoir hosts for its progenitor. Although the similarity of 96% to the coronavirus sequence hosted by the *Rhinolophus*, Andersen et al. [2] identified that its spike protein diverges in the receptor binding domain (RBD), which suggests that it may not bind efficiently to the human ACE2 receptor. Furthermore, it is assumed in this and other studies [29, 33] that an intermediate host was probably involved.



Project 2: Literature Clustering and Visualization

8	Introduction	31
8.1	Background	
8.2	Goal of the Project	
8.3	Outcome	
9	Part 1	33
9.1	Redo and Understand an Existing Clustering Approach	
10	Part 2	37
10.1	Add Related Papers to the CORD-19 Data Set	
10.2	Redo the Clustering	
11	Part 3	43
11.1	Loading in the Student Data Set	
11.2	Preprocessing, Clustering and Results	
11.3	Creating a word cloud	



8. Introduction

8.1 Background

The overwhelming amount of daily published papers correlated to the corona virus make it difficult, even for health professionals, to keep up with new information about the virus. One way of managing the flood of information is by clustering them according to their topics to simplify the search. Therefore, we have performed a cluster analysis for the CORD-19 data set, which contains roughly 60,000 articles.

8.2 Goal of the Project

After parsing the body of each article in the data set, the extracted information is transformed into a feature vector. We applied dimensionality reduction using PCA and performed k-means clustering. Subsequently, t-SNE is applied to project the original feature vector into two dimensions such that clusters become visible in the two dimensional space. Each course participant selected five scientific papers that cluster in the same group as the article they introduced in *Part I*. The submissions have been used to create a new data set. Furthermore K-means was performed on the data set to determine specific clusters. Finally, the selected articles of *Part I* were added to the CORD-19 data set. The clustering was redone to see if the papers will be assigned to the expected clusters. In addition, two methods for selecting the best k value with two distance metrics were compared: Silhouette Scoring vs Distortion and Euclidean vs Cosine Similarity.

8.3 Outcome

The five papers we selected in Part I were clustered into three different groups. Comparing both, the method of choosing k by elbow point or silhouette scoring and the distance metrics euclidean and cosine similarity, we determined that silhouette scoring and euclidean distance performed better. Ten clusters with unique topics were determined (Table 11.1).



9. Part 1

9.1 Redo and Understand an Existing Clustering Approach

The literature clustering pipeline started with the data import of the CORD-19 data set. The resulting metadata dataframe listed 59887 entries of coronavirus related publications, containing the body text, abstract, pub-id and several other features that describe the papers.

The metadata information is subsequently merged with the body text of the papers which are stored in separate JSON files. Due to partially missing information only 43331 entries of the metadata could be merged with the JSON files. To get an overview of the average text length of the abstracts and the body text information (on which the clustering will be performed) the overall and unique number of words were calculated. An average abstract length of 157 words and an average body text length of 4528 words (1376 unique) was found. Since the data was uploaded by many different sources, duplicates were present in the data set. Those needed to be filtered out such that. Consequently, only 30960 publications remained in the set. The subsequent calculation steps of the pipeline will require very high computing resources. Therefore, we randomly subsampled (seed=42) the data set to a maximum size of 10,000 instances. Unfortunately, we noticed afterward that both null values (1073) and non-English publications (242) were still present in the sampled data set. Those would massively reduce the interpretability of the clustering result. Thus we removed them. The final data set consisted then of 8685 entries.

In the next preprocessing step we detected and removed stop words. These are common words in the written text, that does not contribute to the content and act as noise in the clustering procedure. The *spacy* package was used to determine the stop words. Additionally, a predefined list of stop words was appended to the list, which contained frequently used words of scientific publication in general. The last step of the preprocessing was to vectorize the cleaned data. Hereby, the string formatted data is converted into a vector-based measure of how important each word is to the instance out of the literature as a whole using the *tf-idf* package. This method creates a very high feature space and since a clustering by k-means needs to be performed, a Principle Component Analysis (PCA) was applied to reduce the number of features by simultaneously keeping 95% of the data variance and immensely reducing the algorithm's runtime. For k-means the optimal k was determined by iterating through different values of k from two to 50. The resulting elbow plot

(Figure 9.1) shows the elbow point at $k=27$, which is subsequently used as the optimal number of clusters.

A t-Distributed Stochastic Neighbor Embedding (t-SNE) was used to reduce the high dimensional features vector to two dimensions. This step provides the possibility to represent the clustering result in a plain 2D coordinate system. The aim of the entire pipeline was to create an interactive bokeh plot. The location of each paper on the plot is determined by t-SNE while the label (color) is determined by k-means. Interestingly, the clusters determined by k-means also cluster together in the 2D representation given by our t-SNE analysis as it can be seen in Figure 9.2. Even though the clusters are now determined the information about the kind of papers within a cluster is still missing. To solve this task, a Latent Dirichlet Analysis (LDA) was performed to model the most important topics for each cluster. This information is also included in the final bokeh plot.

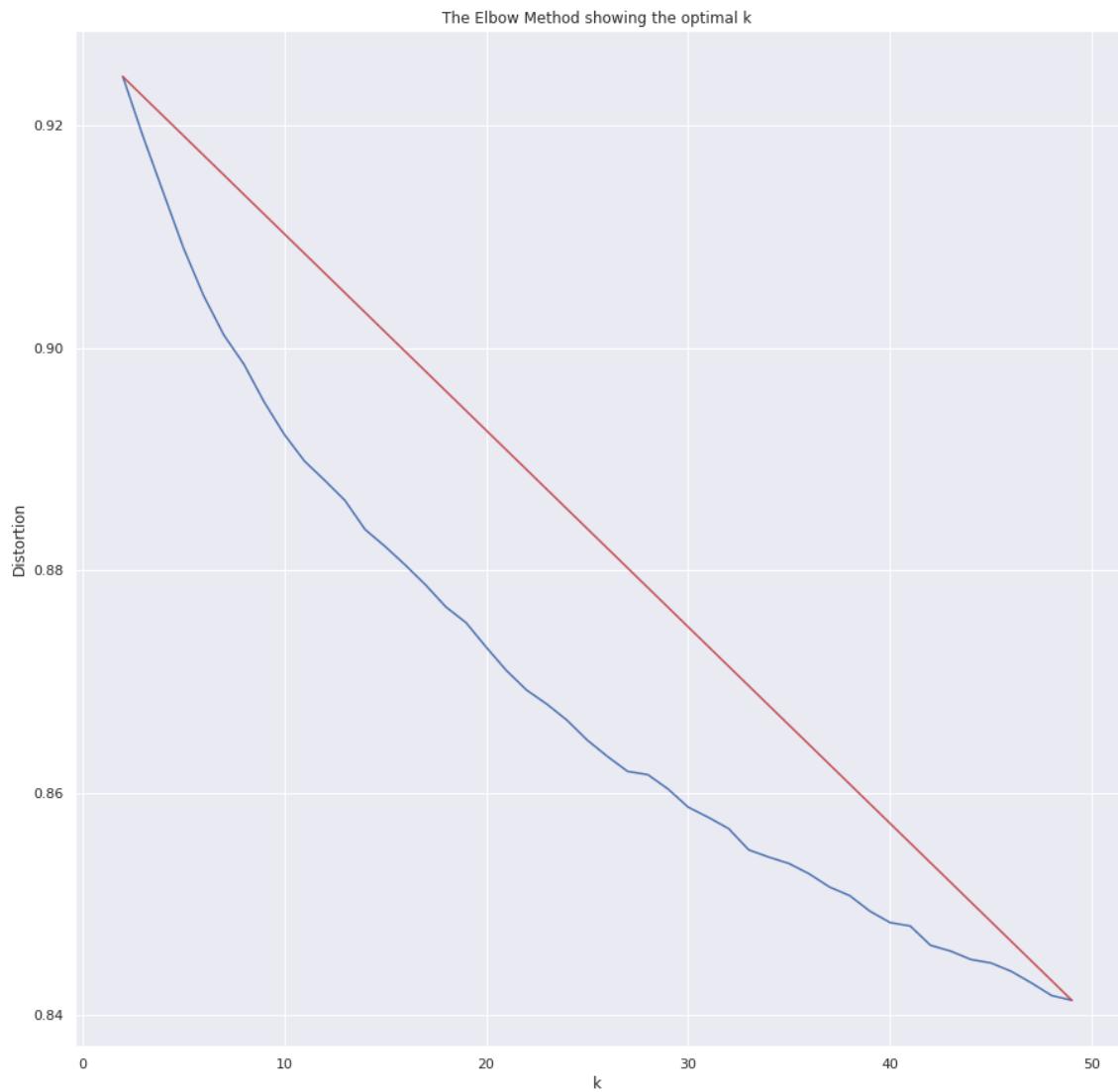


Figure 9.1: The Figure shows an elbow plot of the k-means clustering of the CORD-19 data set for k values from two to 50. On the y-axis the distortion and on the x-axis the number of clusters is represented. A clear elbow point cannot be identified but the incline decreases at $k=27$.

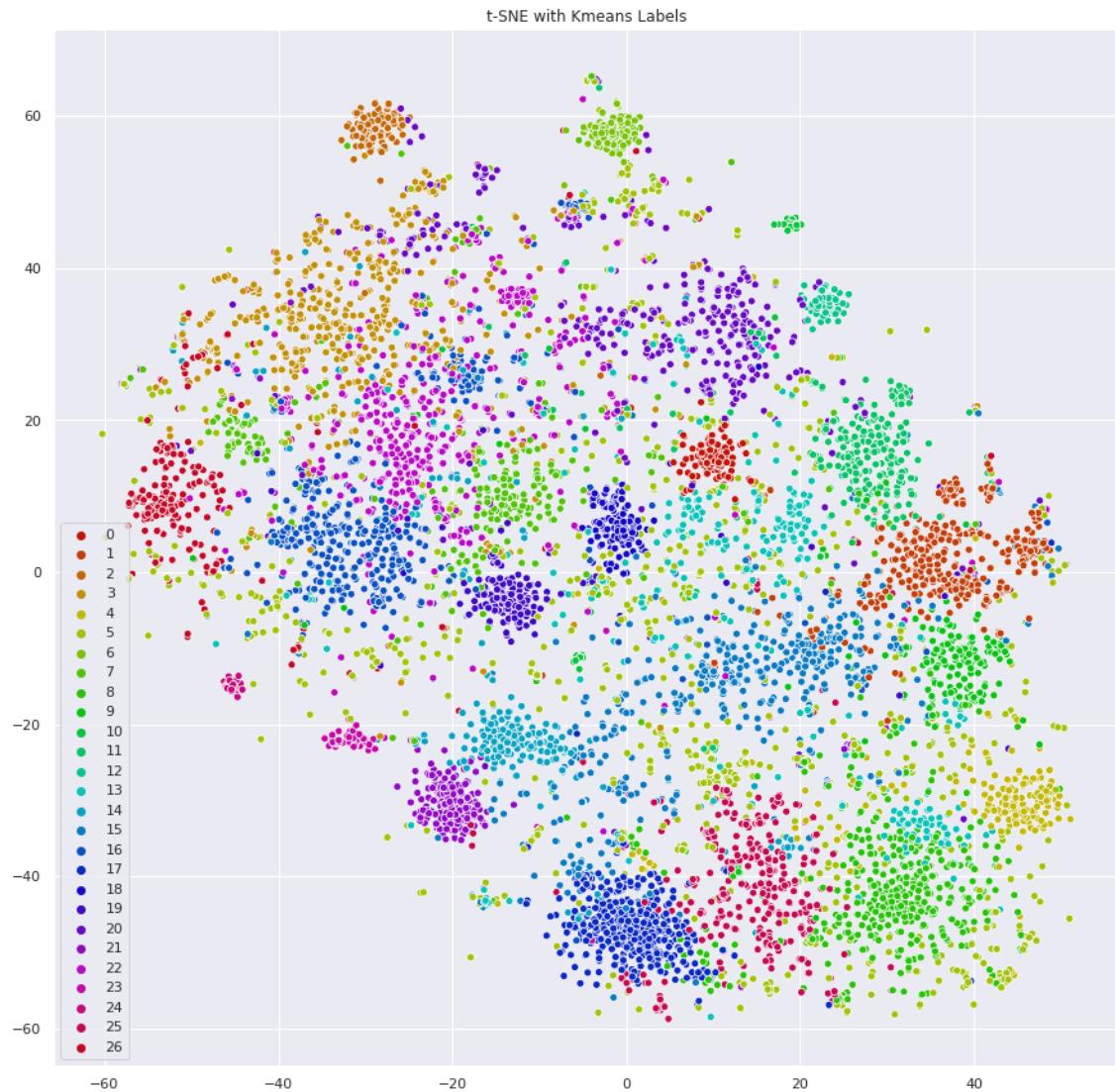
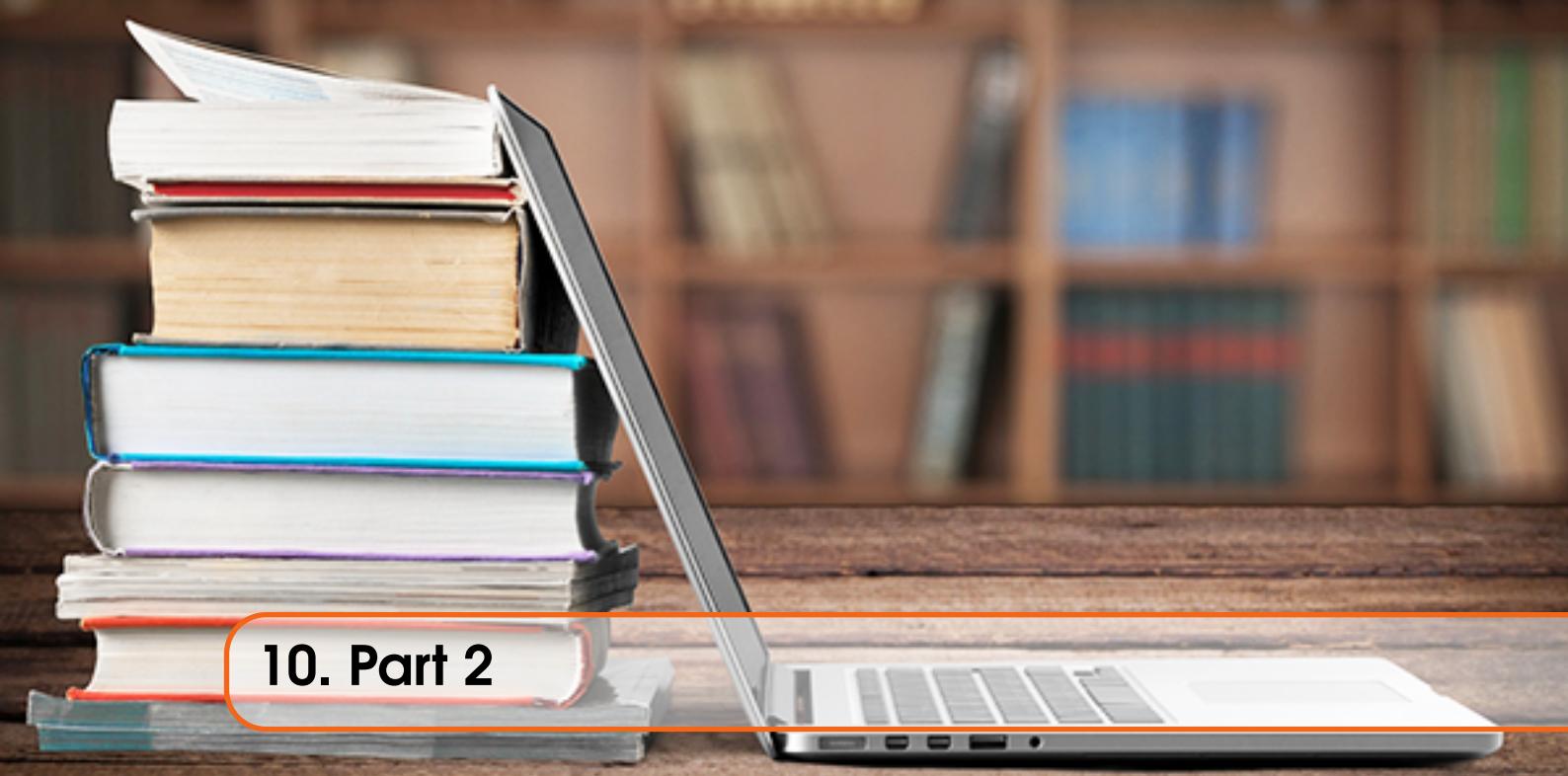


Figure 9.2: The plot shows the clustering result with the t-SNE positioning and the k-means labels. Various distinct clusters represented with different colors, can be identified indicating a good separation via k-means clustering.



10. Part 2

10.1 Add Related Papers to the CORD-19 Data Set

In this task we were supposed to add the metadata information of the selected papers from *Part I* to the CORD-19 data set and redo the clustering. We wanted to find out if this approach improves and facilitates the search of papers for a specific topic field. In order to solve this task, we tried to generate a dataframe using the *pyPDF2* package to extract the metadata from the pdf files of the articles. Unfortunately, it did not work for all pdfs, because they did not contain uniform metadata fields. Consequently, we decided to create a CSV file and added all relevant metadata fields manually (Table 10.1). The manually created table was then added to the CORD-19 data set (Figure 10.2).

Link	Title
https://www.ncbi.nlm.nih.gov/pubmed/32276116	Rapid and visual detection of 2019 novel coronavirus (SARS-CoV-2) by a reverse transcription loop-mediated isothermal amplification assay
https://www.nature.com/articles/s41598-018-37483-w	A method to identify respiratory virus infections in clinical samples using next-generation sequencing
https://www.tandfonline.com/doi/full/10.1586/14737159.2014.888313	Advances and challenges in biosensor-based diagnosis of infectious diseases
https://ann-clinmicrob.biomedcentral.com/articles/10.1186/1476-0711-7-18	Predicting the sensitivity and specificity of published real-time PCR assays
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3522074/	Application of Molecular Diagnostic Techniques for Viral Testing

Table 10.1: Papers to diagnostics with columns providing the link to the respective title.

10.2 Redo the Clustering

Then we applied the previously described clustering pipeline to our new data set. Finally, the cluster membership of each paper from *Part I* could be identified. The titles of papers from the cluster were saved as .csv-file respectively (see Figure 10.3).

We figured out that three papers (spreading models, database time-series prediction and risk factor analysis) from Part I belong to the same cluster 6 (Figure 10.4). The paper related to diagnostics was assigned to cluster 15 (Figure 10.5) and the origin analysis article was a member of cluster 9 (Figure 10.6).

A	B	C	D	E	F	G	H
paper_id	doi	abstract	body_text	authors	title	journal	abstract_summary
week_2_spreading_models	10 3390/ijerph 16234683	infectious diseases are an important cause of human death. The study of the pathogenesis, spread regularity and development trend of infectious diseases not only provides a theoretical basis for future research on infectious diseases but also has practical guiding significance for the prevention and control of their spread. In this paper, a controlled differential equation and an ordinary differential equation are combined to model the transmission process of the novel coronavirus. We developed a computational tool to assess the risks of novel coronavirus outbreaks outside of China. We estimate the dependence of the risk of a major outbreak in a country from imported cases on key parameters such as: (i) the evolution of the cumulative number of cases in mainland China. Based on the official data modeling, this paper studies the transmission process of the novel coronavirus disease 2019 (COVID-19). The error between the model and the official data curve is quite small. At the same time, it realized forward prediction and backward inference of the epidemic situation and the relevant analysis helps relevant countries to take effective measures to prevent and control the spread of COVID-19.	infectious diseases are diseases that can be transmitted from person to person, from person to animal or from animal to person after proto-microorganisms and parasites infect human beings or animals [1–3]. Infectivity, epidemic and uncertainty are the three main characteristics of infectious diseases. A thorough study of the second cluster of pneumonia cases in Wuhan, China was reported to the World Health Organization (WHO) on 31 December 2019. The cause of the pneumonia cases was identified as a novel betacoronavirus, the 2019 novel coronavirus (2019-nCoV, recently renamed as SARS-CoV-2). At the end of 2019, the new coronavirus (COVID-19) spread widely in China and a large number of people became infected. At present, the domestic outbreak has been effectively controlled while the new coronavirus is spreading rapidly in other areas. Currently, Europe has become the center of identifying who has the COVID-19 virus.	Bin Sheng Sun Gengxin Chen Chih-Cheng	Spread of Infectious Disease Modeling and Analysis of Different Factors on Risk	International Journal of Environmental Research and Public Health	something
week_2_risk	10 3390/jcm90 20571	We developed a computational tool to assess the risks of novel coronavirus outbreaks outside of China. We estimate the dependence of the risk of a major outbreak in a country from imported cases on key parameters such as: (i) the evolution of the cumulative number of cases in mainland China. Based on the official data modeling, this paper studies the transmission process of the novel coronavirus disease 2019 (COVID-19). The error between the model and the official data curve is quite small. At the same time, it realized forward prediction and backward inference of the epidemic situation and the relevant analysis helps relevant countries to take effective measures to prevent and control the spread of COVID-19.	A cluster of pneumonia cases in Wuhan, China was reported to the World Health Organization (WHO) on 31 December 2019. The cause of the pneumonia cases was identified as a novel betacoronavirus, the 2019 novel coronavirus (2019-nCoV, recently renamed as SARS-CoV-2). At the end of 2019, the new coronavirus (COVID-19) spread widely in China and a large number of people became infected. At present, the domestic outbreak has been effectively controlled while the new coronavirus is spreading rapidly in other areas. Currently, Europe has become the center of identifying who has the COVID-19 virus.	Boldog Péter Tekeli Tamás Vizi Zsolt Dénes Li Lixiang	Assessment of Novel Coronavirus COVID-19 Outbreaks	Journal of Clinical Medicine	something
week2_forecasting	https://doi.org/10.1016/j.idm.2020.03.002	https://doi.org/10.1016/j.idm.2020.03.002	Testing for COVID-19 has been unable to keep up	Yang Zihang Dang Zhongkai Meng Cui	Propagation analysis and prediction of the COVID-19	KeAi	something
				Hall	Finding		

Figure 10.1: csv-table showing columns like paper_id, doi, abstract etc

2.3 Appending metadata of papers from week2

Loading metadata of papers from week2 from csv file as dataframe, preprocessing (stripping whitespace, removing char, adding word counts of abstract, body text and unique words) and appending metadata to a general dataframe containing all the papers), dropping duplicates, data summary

```
In [39]: import pandas as pd
df_papers_week2 = pd.read_csv('C:/Users/Natalja/shared_folder/DSinLS20/week3/Week2_Papers.csv', sep=';', dtype={'paper_id': str})
for column in df_papers_week2:
    df_papers_week2[column] = df_papers_week2[column].apply(lambda x: x.strip())
df_papers_week2['title'] = df_papers_week2['title'].str.replace('<br>', ' ')
df_papers_week2.drop_duplicates(subset='title', keep='first', inplace=True)
df_papers_week2['abstract_word_count'] = df_papers_week2['abstract'].apply(lambda x: len(x.strip().split())) # word count in abstract
df_papers_week2['body_word_count'] = df_papers_week2['body_text'].apply(lambda x: len(x.strip().split())) # word count in body text
df_papers_week2['body_unique_words'] = df_papers_week2['body_text'].apply(lambda x:len(set(str(x).split()))) # number of unique words
df= df_papers_week2.append(df)
df.drop_duplicates(['abstract', 'body_text'], inplace=True)
df['abstract'].describe(include='all')

Out[39]: count    10005
unique    7265
top
freq      2794
Name: abstract, dtype: object
```

Figure 10.2: Appending metadata of week2 papers to a data frame containing metadata to CORD-19-research-challenge.

Identifying of cluster for each week2 paper

using helper function save_cluster_week2_titles clusters of each paper from week2 identified and paper titles of this cluster will be saved in extra .csv-file

```
In [113]: def save_cluster_week2_titles(title, table):
    cluster=table.loc[table['title']== title, 'y'].values
    #print(cluster)
    cluster_value=cluster[0]
    print('Cluster {}'.format(cluster_value))
    cluster=table.loc[table['y'] == cluster_value, 'title']
    nameId_week2=table.loc[table['title'] == title, 'paper_id'].values
    nameId_week2=nameId_week2[0]
    print(nameId_week2)
    savePath = 'C:/Users/Natalja/shared_folder/DSinLS20/week3/df_covid_cluster_topic_{}.csv'.format(nameId_week2)
    cluster.to_csv(savePath, index = False)
for index, row in df_papers_week2.iterrows():
    print ('Title of week 2 paper is {}'.format(row['title']))
    save_cluster_week2_titles(row['title'], df)

Title of week 2 paper is Spread of Infectious Disease Modeling and Analysis of Different Factors on Spread of Infectious Disease Based on Cellular Automata
Cluster 6
week_2_spreading_models
Title of week 2 paper is Risk Assessment of Novel Coronavirus COVID-19 Outbreaks Outside China
Cluster 6
week_2_risk
Title of week 2 paper is Propagation analysis and prediction of the COVID-19
Cluster 6
week2_forecasting
Title of week 2 paper is Finding COVID-19 from Chest X-rays using Deep Learning on a Small Dataset
Cluster 15
week2_diagnostics
Title of week 2 paper is Identification of a new coronavirus
Cluster 9
week2_phylogenetic_analysis
```

Figure 10.3: Code snipped for the clustering of the week2 papers.

```
title
Spread of Infectious Disease Modeling and Analysis of Different Factors on Spread of Infectious Disease Based on Cellular Automata
Risk Assessment of Novel Coronavirus COVID-19 Outbreaks Outside China
Propagation analysis and prediction of the COVID-19
Anthropological Perspectives on the Health<br>Transition
How HIV patients construct liveable<br>identities in a shame based culture: the case of Singapore
Estimating the economic impact of pandemic<br>influenza: An application of the computable general<br>equilibrium model to the
" Travelling to scientific meetings is a<br>mission, not a vacation"
Perspectives of public health laboratories in<br>emerging infectious diseases
D(2)EA: Depict the Epidemic Picture of<br>COVID-19
Suicide news reporting accuracy and<br>stereotyping in Hong Kong
Pandemic Risk Modelling
Reflections on travel-associated infections<br>in Europe
Chapter 27 Disaster Mitigation
Chapter 3 Emerging Infectious Diseases and the<br>International Traveler
Learning from recent outbreaks to strengthen<br>risk communication capacity for the next influenza<br>pandemic in the Western
Impact of the topology of metapopulations on<br>the resurgence of epidemics rendered by a new<br>multiscale hybrid modeling a
" After Malaria Is Controlled, What's Next?"*
A High-Resolution Human Contact Network for<br>Infectious Disease Transmission
Generality of the Final Size Formula for an<br>Epidemic of a Newly Invading Infectious Disease
Committed to Health: Key Factors to Improve<br>Users' Online Engagement through Facebook
Temporal patterns and geographic<br>heterogeneity of Zika virus (ZIKV) outbreaks in French<br>Polynesia and Central America
The challenges of implementing an integrated<br>One Health surveillance system in Australia
The legal determinants of health: harnessing<br>the power of law for global health and sustainable<br>development
Beyond the 'nanny state': Stewardship and<br>public health
Pandethics
International Organizations and Their<br>Approaches to Fostering Development
Using core competencies to build an evaluative<br>framework: outcome assessment of the University of Guelph<br>Master of Publ
China's distinctive engagement in global<br>health
A planetary vision for one health
```

Figure 10.4: Cluster 6: titles of related papers to the papers from week2 about risk factor analysis, spreading models and forecasting

```

|title
Finding COVID-19 from Chest X-rays using Deep Learning on a Small Dataset
COVID-19 and Dialysis Units: What Do We Know Now<br>and What Should We Do?
G6PD deficiency in COVID-19 pandemic: "a ghost<br>in the ghost"
COVID-19 pneumonia with hemoptysis: Acute<br>segmental pulmonary emboli associated with novel<br>coronavirus infection
" Maintenance Hemodialysis and Coronavirus<br>Disease 2019 (COVID-19): Saving Lives With Caution,<br>Care, and Courage"
Continuing education in oral cancer during<br>coronavirus disease 2019 (covid-19) outbreak
Inuit communities can beat COVID-19 and<br>tuberculosis
Tackling the COVID-19 Pandemic
Fellowship Training in Adult Cardiothoracic<br>Anesthesiology - navigating the new educational landscape due<br>to the corona
Pediatric Airway Management in Coronavirus<br>Disease 2019 Patients: Consensus Guidelines From the<br>Society for Pediatric A
" COVID-19, A Clinical Syndrome Manifesting as<br>Hypersensitivity Pneumonitis"
Editorial. Endonasal neurosurgery during the<br>COVID-19 pandemic: the Singapore perspective
Increased risk of ocular injury seen during<br>lockdown due to COVID-19
COVID-19 in pregnancy: early lessons
Clinical course and mortality risk of severe<br>COVID-19
Reply to "The use of traditional Chinese<br>medicines to treat SARS-CoV-2 may cause more harm than<br>good"
Knowledge and attitudes of medical staff in<br>Chinese psychiatric hospitals regarding COVID-19
The preventive strategies of GI physicians<br>during the COVID-19 pandemic
" Coronavirus disease (COVID-19) in a<br>paucisymptomatic patient: epidemiological and clinical<br>challenge in settings with
Perspectives from the Cancer and Aging<br>Research Group: Caring for the vulnerable older patient<br>with cancer and their car
SARS-CoV-2 infection in a patient on chronic<br>hydroxychloroquine therapy: Implications for prophylaxis
COVID-19 Diagnostic and Management Protocol<br>for Pediatric Patients
" Spinal anaesthesia for patients with<br>coronavirus disease 2019 and possible transmission rates<br>in anaesthetists: retros
" Epidemiology, causes, clinical<br>manifestation and diagnosis, prevention and control of<br>coronavirus disease (COVID-19) o
Concerns for activated breathing control<br>(ABC) with breast cancer in the era of COVID-19:<br>Maximizing infection control
Heart Failure Editorial Emergencies in the<br>COVID-19 Era
Ayurveda and COVID-19: where<br>psychoneuroimmunology and the meaning response meet
Pulmonary Pathology of Early-Phase 2019 Novel<br>Coronavirus (COVID-19) Pneumonia in Two Patients With Lung<br>Cancer
WFUMB Position Statement: How to perform a safe<br>ultrasound examination and clean equipment in the context<br>of COVID-19

```

Figure 10.5: Cluster 15: titles of related papers to the paper from week 2 about diagnostics

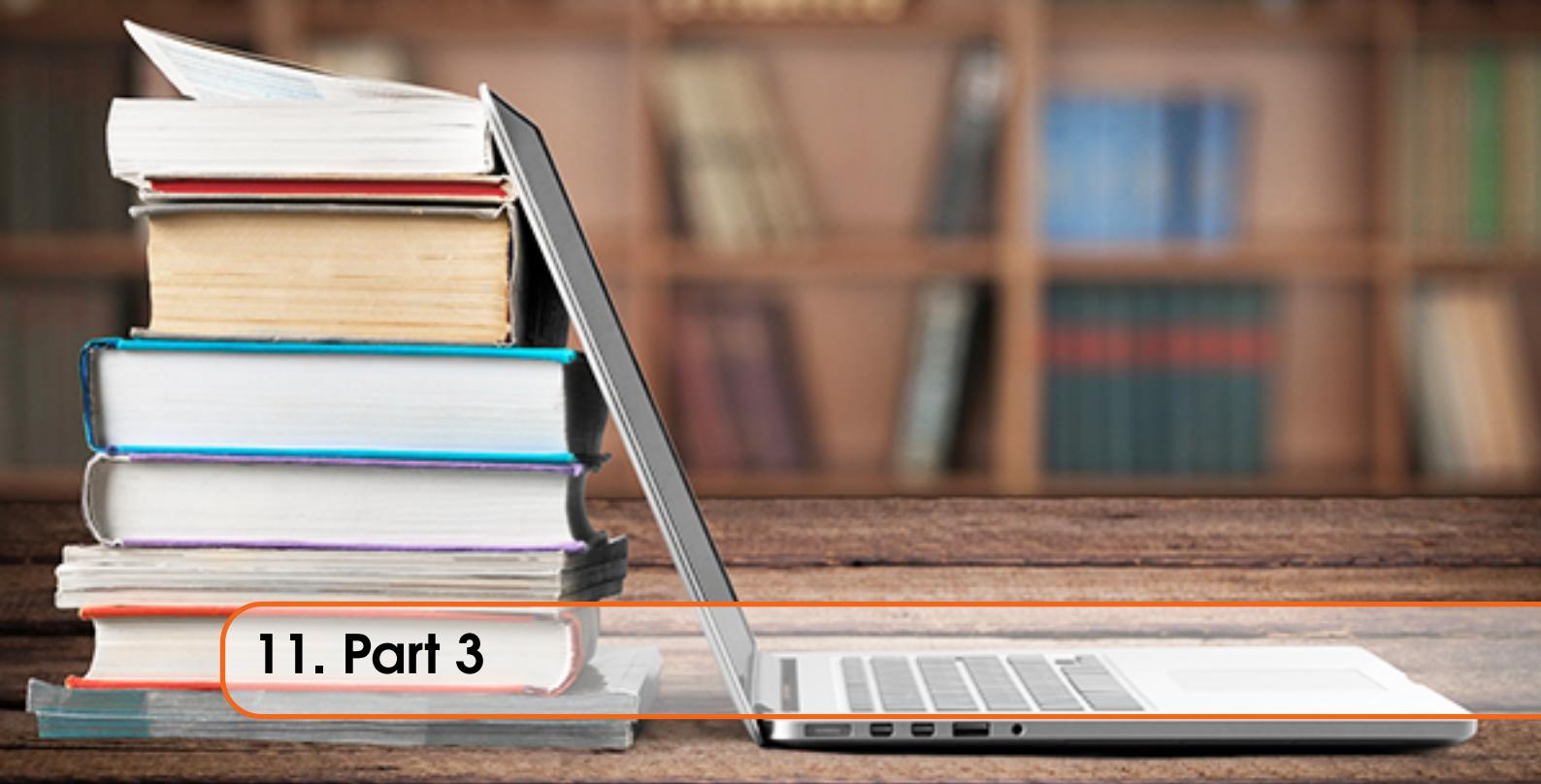
```

|title
Identification of a new coronavirus
High Resolution Analysis of Respiratory<br>Syncytial Virus Infection In Vivo
Detection of Novel SARS-like and Other<br>Coronaviruses in Bats from Kenya
Recombinant infectious bronchitis<br>coronavirus H120 with the spike protein S1 gene of the<br>nephropathogenic IBYZ strain
Nucleotide Sequence of the Inter-Structural<br>Gene Region of Feline Infectious Peritonitis Virus
Molecular characterization of bovine<br>noroviruses and neboviruses in Turkey: detection of<br>recombinant strains
" Detection and characterisation of canine<br>astrovirus, canine parvovirus and canine papillomavirus<br>in puppies using next
Identification and Characterization of<br>Severe Acute Respiratory Syndrome Coronavirus<br>Subgenomic RNAs
Coevolution of activating and inhibitory<br>receptors within mammalian carcinoembryonic antigen<br>families
CHAPTER 1 Remarks on the Classification of<br>Viruses
Canine kobuvirus infections in Korean dogs
Coronavirus Transcription: A Perspective
Identification and Analysis of Frameshift<br>Sites
" Codon usage in Alphabaculovirus and<br>Betabaculovirus hosted by the same insect species is weak,<br>selection dominated by
Genic amplification of the entire coding<br>region of the HEF RNA segment of influenza C virus
Comprehensive codon usage analysis of porcine<br>deltacoronavirus
The First Detection of Equine Coronavirus in<br>Adult Horses and Foals in Ireland
Sequences Promoting Recoding Are Singular<br>Genomic Elements
Recombination and Coronavirus Defective<br>Interfering RNAs
A recombinant infectious bronchitis virus<br>from a chicken with a spike gene closely related to<br>that of a turkey coronaviru
Single Stranded DNA Viruses Associated with<br>Capybara Faeces Sampled in Brazil
Genetic diversification of penaeid shrimp<br>infectious myonecrosis virus between Indonesia and<br>Brazil
" Discovery of novel virus sequences in an<br>isolated and threatened bat species, the New Zealand<br>lesser short-tailed bat
" Spliced Leader RNAs, Mitochondrial Gene<br>Frameshifts and Multi-Protein Phylogeny Expand Support<br>for the Genus Perkinsu
" Genomic Organization, Biology, and Diagnosis<br>of Taura Syndrome Virus and Yellowhead Virus of<br>Penaeid Shrimp"
" Polymorphisms and Tissue Expression of the<br>Feline Leukocyte Antigen Class I Loci FLAI-E, -H and -K"
WHO says coronavirus causes SARS
Conserved tertiary structure elements in the<br>5' untranslated region of human enteroviruses<br>and rhinoviruses
Standards for Sequencing Viral Genomes in the<br>Era of High-Throughput Sequencing

```

Figure 10.6: Cluster 9: titles of related papers to the paper from week 2 about origin analysis

As can be seen, the titles of the articles in the identified clusters match quite well the papers from Part I, which indicated that the literature clustering can be effectively used to organize newly published articles.



11. Part 3

11.1 Loading in the Student Data Set

After adding only the information about the five papers that we have chosen in *Part I* to the CORD-19 data set, we created a new data set containing all submitted papers by the course participants. It turned out that opening each of the 60,000 JSON files in the CORD-19 data set, to filter those JSON's that do not match one of the submitted articles is too time consuming. The runtime was dramatically decreased by first joining the course data set with the metadata. Title names have been stripped since traveling whitespaces result in mismatches. The included paper id ("sha") was then used to search in the file names of the JSON's for the papers of interest. This leads to a runtime reduction from several hours to less than three minutes. We were able to match 177 of the 195 submitted articles. And after removing duplicated entries our new data set contains 146 papers.

11.2 Preprocessing, Clustering and Results

Next, we proceeded with the preprocessing and final clustering step. For the preprocessing, we followed the previously described pipeline of the kaggle notebook closely. We removed common stopwords as they act as noise. Next, a vectorizer with a noise filer of 2^{12} is applied, counting words and scoring less frequent words higher. Afterwards, the dimension of the data set was reduced by PCA from over 4000 to 119. For the clustering we deviated from the kaggle notebook. First we tested a second method to determine the best k value. This was done by computing the silhouette score (Figure 11.1). Second we added another preprocessing step, to transform the data to a unit vector of 1, thus using the equivalent of cosine similarity as distance metric. After carefully comparing both methods with and without cosine similarity, one of elbow distortions and the other of silhouette scoring, we determined k by silhouette scoring and euclidean to be the best method (Figure 11.2). Thus we proceeded with k of 18. We choose 18, because it showed a noticeable bump in scoring, but is still a reasonable estimate based on the chosen paper data set by students. After running a t-SNE visualisation (Figure 11.3) and a Keyword extraction per cluster using LDA, with a minimum number constraint of 10 topics per cluster, we found 10 distinct clusters with an overarching topic. 8 clusters were removed due to insufficient number of data points. The results

can be found in the Table 11.1.

The assignment of topics to each unique cluster was surprisingly easy, indicating a meaningful result. Some problems could be identified, for one some clusters had not enough articles assigned to them. Thus we could stipulate that the chosen k value of 18 is too high, rerunning the clustering at a lower value might give better results. On the contrary some clusters with high assignment could be identified (cluster 11). This could mean two things: Not much variation in the chosen cluster topics from the students or a more granular clustering is needed.

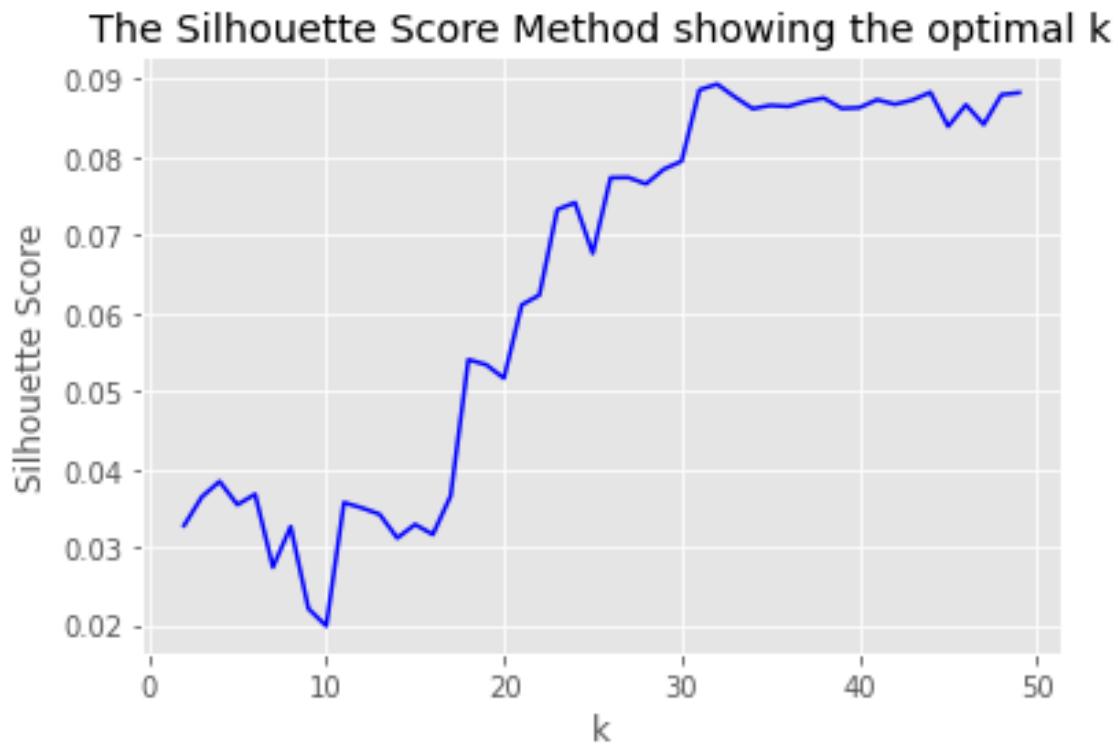


Figure 11.1: Silhouette scoring method for k-means with x axis indicating k value and y axis with Silhouette scores. At the cluster point 17 a significant bump in scoring can be seen.

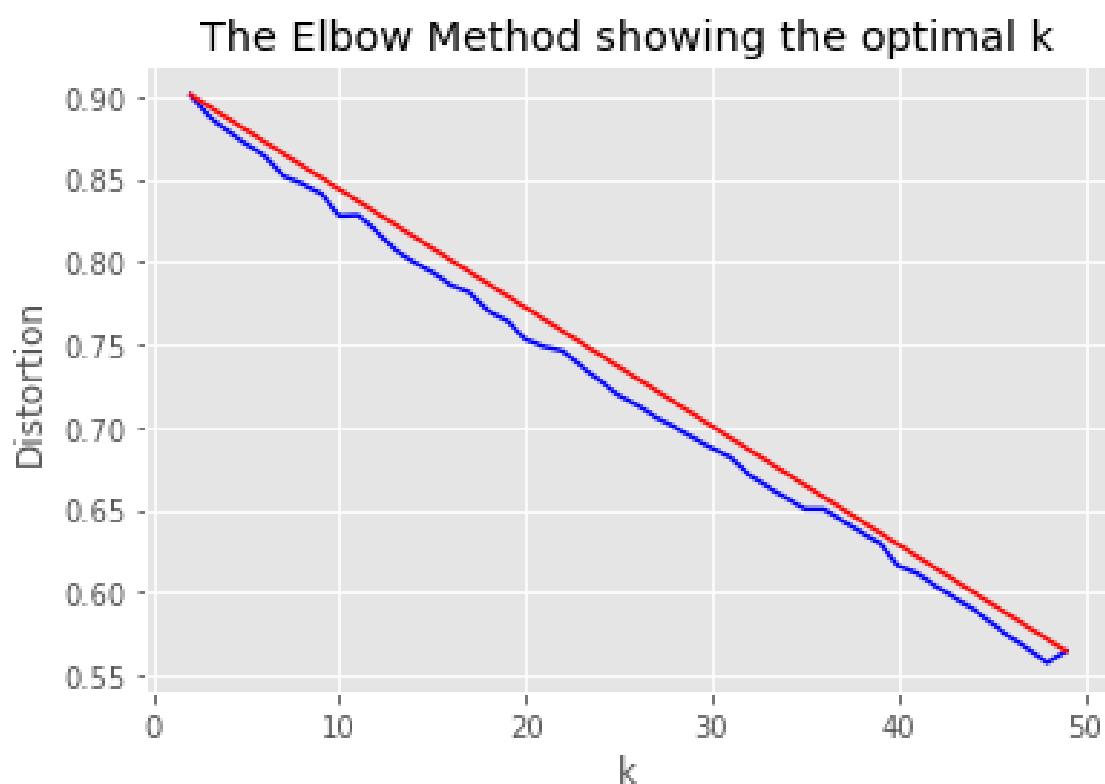


Figure 11.2: Distortion scoring method for k-means. x axis shows the k values and y axis showing Distortion. Due to the low difference in scoring an almost linear line is produced. Thus in this case the method is unsuited to determine the optimal k for clustering



Figure 11.3: 2 dimensional t-SNE visualisation of the data set. Due to the low presence of articles the projection seems to be sparse. Still some clusters can be determined like: Top middle cluster 14: virus detection or middle left cluster 15: compartmentalized models. Big clusters like 11: modeling of spread are not clumped together.

Cluster	No. Articles	Assigned Topic	Keywords
0	11	phylogenetics	'sars', 'set', 'gene', 'rate', 'orf', 'recombination', 'datum', 'frequency', 'region', 'distance'
1	13	study of initial outbreak	'symptom', 'hospital', 'sars-cov-', 'china', 'country', 'mortality', 'use', 'respiratory', 'risk', 'evidence'
2	8	network-modeling of spread	'community', 'degree', 'threshold', 'mix', 'heterogeneity', 'group', 'node', 'use', 'transmission', 'size'
5	6	network-modeling of spread	'mix', 'wave', 'estimate', 'outbreak', 'total', 'community', 'human', 'delay', 'overall', 'additional'
9	11	disease forecasting	'case', 'disease', 'outbreak', 'estimate', 'influenza', 'process', 'interval', 'forecast', 'day', 'peak'
11	31	modeling of spread	'parameter', 'sequence', 'disease', 'method', 'change', 'city', 'spread', 'network', 'value', 'interaction'
12	18	origin detection	'case', 'camel', 'sars-cov-', 'human', 'protein', 'isolate', 'healthcare', 'bat', 'viral', 'sars-cov'
14	12	virus detection	'sars', 'sample', 'lung', 'serum', 'finding', 'detection', 'virus', 'care', 'study', 'swab'
15	8	compartmentalized modeling	'rate', 'risk', 'model', 'outbreak', 'death', 'datum', 'day', 'virus', 'state', 'patient'
16	6	diagnostics	'pcr', 'rsv', 'sequence', 'pneumonia', 'age', 'child', 'rhinovirus', 'associate', 'young', 'presence'

Table 11.1: Results of the clustering with unique topic assignment based on the first 10 keywords. 8 clusters are omitted due to low assignment of articles.

11.3 Creating a word cloud

Finally, we generated a basic word cloud. Here, we took the extracted keywords from the final clustering and used the word cloud package. The package provides a basic understanding of the word cloud with the use of some simple python libraries like *numpy*, *pandas*, *matplotlib* and *pillow*. The Figure 11.4 shows the visualized wordcloud. Words like Sars, Virus and Cov are bold and large which shows that the frequency of usage of these words is high and denotes their importance of it during the recent times. While these wordclouds are easy to look at, not much useful information can be extracted..

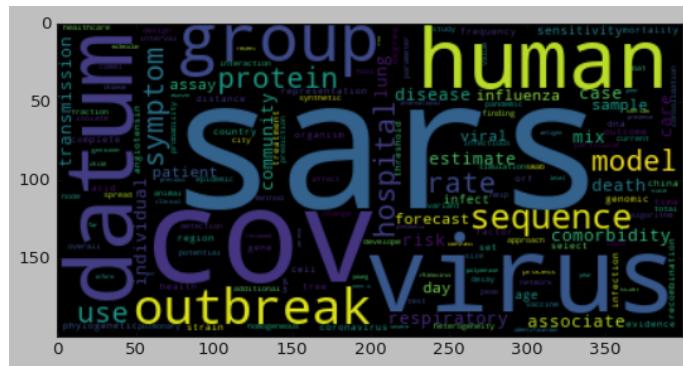
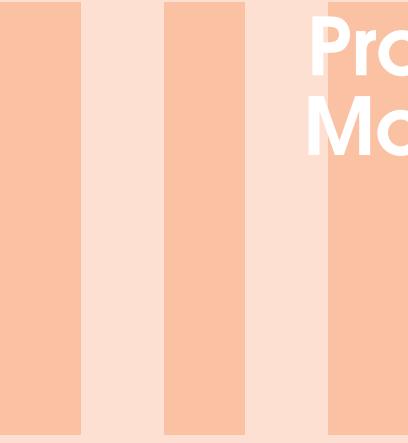


Figure 11.4: Word cloud of most frequent words. Most frequent words are clearly visible e.g. SARS and Human. More specific categories are less visible e.g. protein and estimate.



Project 3: System Dynamics Modeling using a SIR Model

12	Introduction	51
12.1	Background	
12.2	Goal of the Project	
12.3	Outcome	
13	The Model	53
13.1	A simple SIR model	
13.2	Extending the SIR model	
13.3	Parameter Fitting	
14	Scenario Studies	59
14.1	Implemented Measures	
14.2	Methods	
14.3	Results	
14.4	Discussion	



12. Introduction

12.1 Background

Dating back to the 1920's [27] for its first inception, a classical approach towards modeling the spread of diseases in epidemiology are SIR-models. SIR-Models are based on the idea of compartmentalization, where the dynamics of an epidemic are studied by dividing the populations into distinct subgroups. The name **SIR** is an abbreviation for its most simple form: **S** standing for **susceptible** (i.e. individuals not yet infected), **I** standing for **infectious** (i.e. infected and infectious individuals) and **R** standing for **recovered** (i.e. individuals which are not infected and infectious anymore). Each compartment can be understood as a state, with a flow from one state to another. By using an equation of the simple form $N = S + I + R$ the whole population N stays static, while the ratios between the states change. Each state is then modeled by ordinary differential equations, thus describing the fluctuations of each state at different time steps t . Furthermore, it is possible to freely add compartments by branching out from current ones, thus making the model adaptable to very different scenarios of an epidemic.

12.2 Goal of the Project

The objective of this week's project is the application of a SIR-model on current COVID-19 case data taken either from a city (e.g. Berlin) or national (e.g. Germany) scale. The model itself is extended beyond the simple case by integrating two new states (Exposed, Dead) to the model and by studying the impact of independent features (ICUbed-capacity, Age, Smoking, and Gender) on the epidemic. By fitting the model to actual case data, possible projections can be made. Furthermore, different scenarios such as lockdown, reducing social contacts, and wearing masks, are explored by simulating their effect on the fitted model. Each prevention method is simulated over different periods and in combination with and without wearing masks on top.

12.3 Outcome

A simple SIR model was implemented while exploring differences in the rate of infection and time to recovery. Extending the model with the independent features age, smoking, and gender

significantly altered the α values (i.e. rate of death) with smoking increasing the factor by more than double from 0.07 to 0.16. Two compartments were added, simulating the incubation period and extending the recovered individuals with death. A simulation with ICU bed capacity was performed, reaching the cap after 50 days. The fitting of the model to the data of Germany resulted in some real numbers as the range R₀ values were kept within the possible range and showed real declined behavior. Other predicted curves for susceptible number, exposed, dead and recovery number were also kept within the realistic bounds. The performed simulations suggest that both the duration as well as intensity of the restrictions plays an important role when fighting the outbreak of corona. The usage of masks is even more important for minor social reduction scenarios.



13. The Model

13.1 A simple SIR model

The simple SIR-model includes three subgroups: *Susceptible*, *Infectious* and *Recovered* cases.

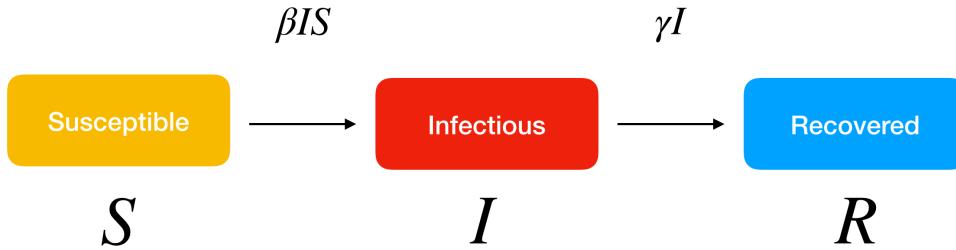


Figure 13.1: A flow diagram showing the state transitions between the subgroups. The whole population is constant, while the flow is unidirectional. (Taken from [50])

Because each compartment can be understood as a state, we can visualize their transitions as a flow diagram (Figure 13.1). The variables above the state transitions describe the rates of individuals switching between the different compartments. For susceptible people becoming infected we introduce the factor β denoting the rate of one person infecting another person and for infected people becoming recovered we introduce γ denoting the rate of infected people developing immunity any given day. Thus we can infer three differential equations for each different subgroup, with N denoting the whole population:

$$\begin{aligned} dS/dt &= -\beta * S * \frac{I}{N} \\ dI/dt &= \beta * S * \frac{I}{N} - \gamma * I \\ dR/dt &= \gamma * I \end{aligned}$$

By integrating these equations over the time point t using the `odeint` function from the `sklearn` package in python3, we can develop a model simulation of the developing compartments for any initial starting conditions. As an example, we can compare the impact of tripling the rate of infection by using a β value of 3.0 compared to 1.0 (Figure 13.2 and 13.3). The β value of 3.0 shows a drastic change. All three compartments are shifted to the left, while the curve of infectious people has a much higher and steeper initial incline, which in turn results in a fast drop of susceptible people. In contrast reducing the γ value from $\frac{1}{4}$ to $\frac{1}{8}$ (Figure 13.4), results in a much lower incline of recovered people and a much longer period of infectious people, as shown by the higher maximum of the yellow line. This shows the importance of both, the β and γ value. Thus, it is of high value to determine the ratio of $\frac{\beta}{\gamma}$ denoted as R_0 to study the dynamics of a developing epidemic.

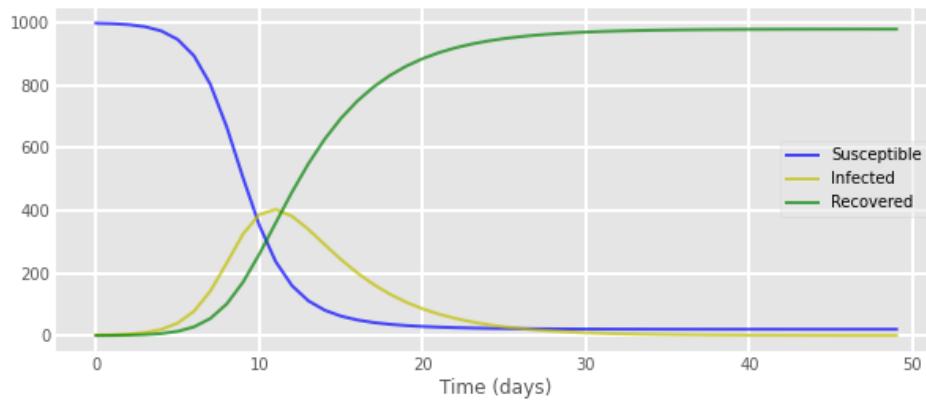


Figure 13.2: Basic SIR model simulation with starting values of $S:999, I:1, R:0, \beta:1.0$ and $\gamma:1/4$.

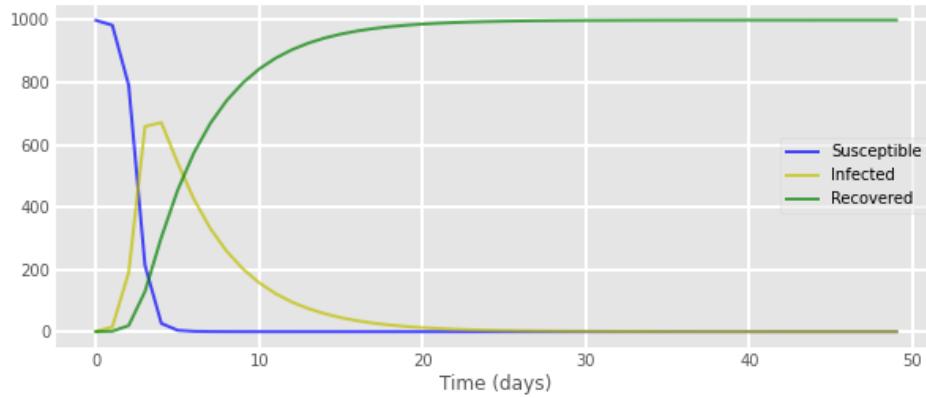


Figure 13.3: Basic SIR model simulation with starting values of $S:999, I:1, R:0, \beta:3.0$ and $\gamma:1/4$.

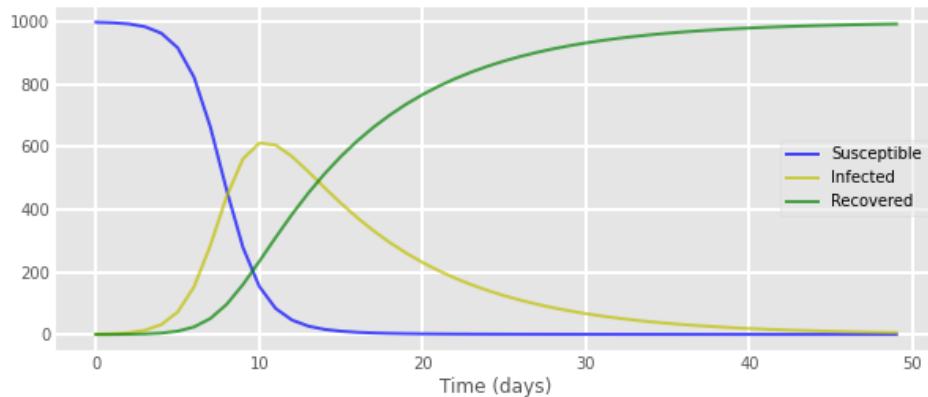


Figure 13.4: Basic SIR model simulation with starting values of $S:999$, $I:1$, $R:0$, $\beta:1.0$ and $\gamma:1/8$.

13.2 Extending the SIR model

The next step was to extend the basic SIR model with 2 new compartments (Figure 13.5).

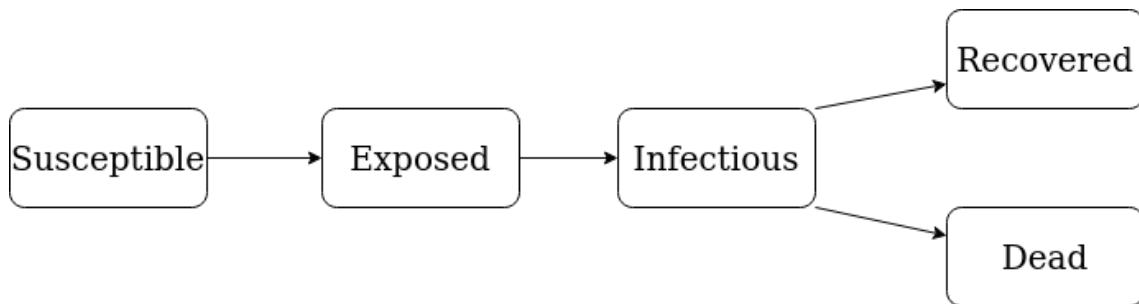


Figure 13.5: Flow diagram of the extended SIR-model. Two new compartments are added: Exposed and Dead.

Between susceptible and infectious people the exposed state is introduced. Exposed individuals carry the virus with an incubation period factor δ but are not infectious. Furthermore, the infectious group now branches out into recovered and dead. Thus, a death rate factor α is introduced to simulate the chance of death for infectious people, while also introducing a factor of ρ for the length of time until death. It follows that the differential equations had to be altered:

$$\begin{aligned} dS/dt &= -\beta * S * \frac{I}{N} \\ dE/dt &= \beta * S * \frac{I}{N} - \delta * E \\ dI/dt &= \delta * E - (1 - \alpha) * \gamma * I - \alpha * \rho * I \\ dR/dt &= (1 - \alpha) * \gamma * I \\ dD/dt &= \alpha * \rho * I \end{aligned}$$

At next, the influence of different population proportions on the fatality rate factor α is introduced to the model. Since the age of infected people has an impact on the severity of the disease and the death rate [63], we created four age groups: 0-29, 30-59, 60-89, and 89+. Based on the death rate calculations by age groups in Italy [13], we assigned differing α values to each age group and added the percentages of the age group distribution in Germany (Table 13.1). The resulting α value for the entire population was 0.07726. Another factor that has a considerable effect on COVID-19 outcomes is smoking behavior. We used an RKI report about the prevalence

of smoking in the adult population of Germany from 2013 [30] to integrate the proportion of daily smokers for each age group (Table 13.2). Abrams et al. [1] analyzed that current or former smokers have increased COVID-19-related mortality by 2.4 [95% CI 1.43–4.04]. Therefore, we multiplied the α values of the smoking people in each age group by this value. This affected the overall α value to be increased to 0.16389.

Age group	% in Germany	α
0-29	30.1	0.001
30-59	41.51	0.013
60-89	28.13	0.2267
89+	0.27	0.285

Table 13.1: Alpha values assigned to different age groups.

Age group	smokers	non-smokers
0-29	31.95	68.05
30-59	26.375	73.625
60-89	8.45	91.55
89+	-	100.0

Table 13.2: Proportion of smoking in the adult population of Germany.

The third and last population proportion, we added to our model was the gender information. Zhang et al. analyzed potential risk factors in a study of n=663 COVID-19 patients and identified that male patients have an odds ratio of 0.486 [95% CI 0.311–0.758] to unimproved from the disease. We integrated this information into our model combined with the proportion of males and females in Germany from 2018 [7],[6]. Since there are 50,7% females and 49,3% males, α was slightly reduced to 0.16383. The final simulation model can be seen in Figure 13.6.

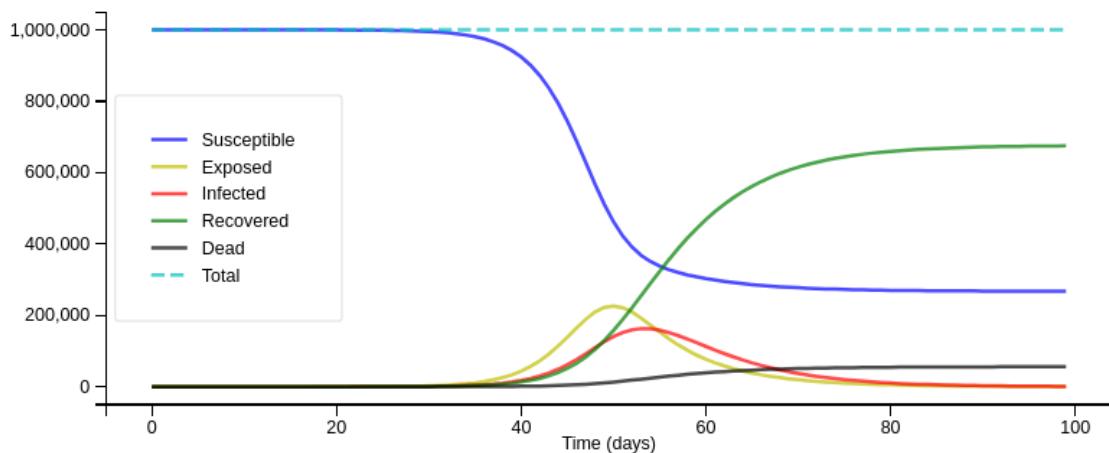


Figure 13.6: Resulting SIR model simulation with adjusted α value.

To get the information how many ICU beds are in use at each time point of the simulation, we created an equation based on two information: 17% of patients infected with COVID-19 need hospitalization in Germany [14] and 48% of the hospitalized patients need ventilation [47] and

therefore an ICU bed. The resulting equation is $\text{occupied ICU beds} = I * 0.17 * 0.48$. When the capacity of ICU beds is exceeded the death rate in our model is changed to 0.6 (Figure 13.7).

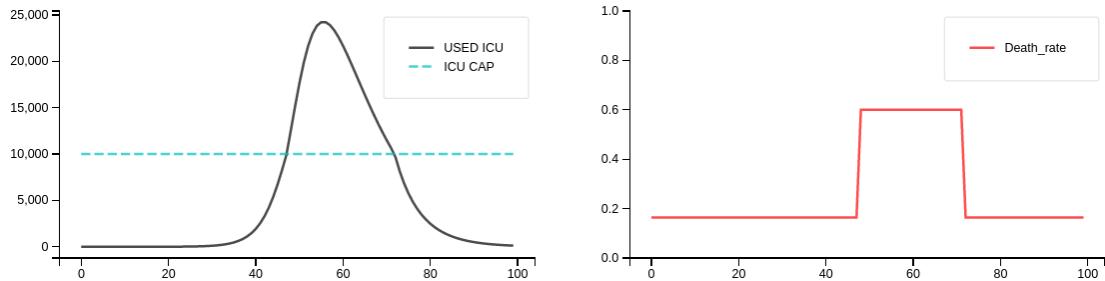


Figure 13.7: The left plot shows the occupied ICU beds (black) and the total amount of ICU beds (blue, dashed) while the right plot shows the corresponding death rate (red).

parameter	description	value
α	fatality rate	0.16
β	expected amount of people an infected person infects per day	1.25
γ	proportion of infected recovering per day	1/10 [47]
δ	incubation period (1/days)	1/5 [47]
<i>inf_to_dead_d</i>	days from infection until death	50 [47]

Table 13.3: Initial parameter setup for the SEIRD-Model described in Section 13.1.

13.3 Parameter Fitting

The next part deals with fitting the extended SIR models time-dependent R_0 values and resource-dependent death rates to real COVID-19 data of Germany [16], which have been obtained from the dashboard of John Hopkins University Center for Systems Science and Engineering [15]. The goal of our fitting is to come as close as possible to the reported numbers for Germany and make predictions about possible future developments. The data set recorded by RKI [45] contains ICU Beds capacity, age groups, R_0 values, beginning of lockdown, and decrease rate of R_0 (k) in Germany. We only considered the data within the range of 01.03.2020 – 30.04.2020 to fit our model.

Here, we initially loaded the data for the age groups, probabilities, created some lookup dictionaries for easy access of the data parameters. We mainly focused on the probabilities of infected to death. Finally, for curve fitting we first set the parameters we knew and assumed with an upper and lower bound values with unknown data.

The resulting simulations showed quite similar values to the reported data, e.g. a R_0 of 2.07 at the start of march 2020, a R_0 of 0.49 at the end of April, and a death rate of 0.16. It is quite comparable with the reported data for Germany as seen if Figure 29.3. With January 21 as the beginning of the outbreak and a 30-day shift in the outbreak, the model predicts that the moderate lockdown in Germany started almost 106 days after the first officially reported case, around the end of April. This also corresponds to the findings in Germany and suggests that the model can be used to predict the progress of the spread of the virus within Germany. Figure 13.9 shows the prediction model. R_0 is stabilized, until the number of cases increases to due lifting of corona measurements.

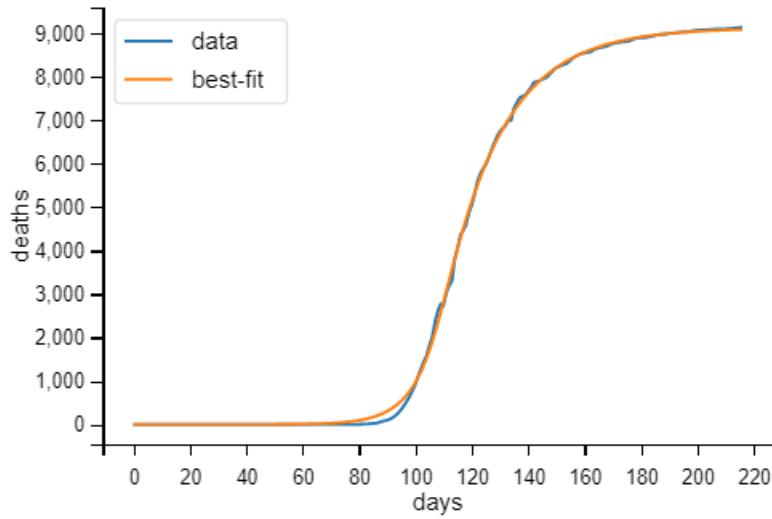


Figure 13.8: Model fit for Germany: Comparison of real and predicted accumulated deaths since 22.01.2020 with an outbreak shift of 30 days. The blue line corresponds to the observed fatalities and the orange line to the ones predicted by the model

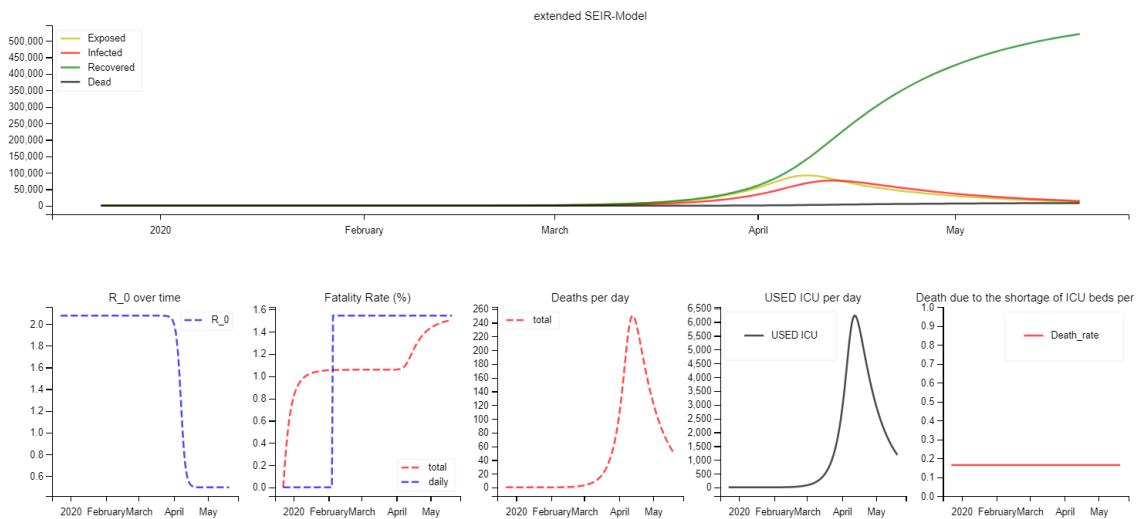


Figure 13.9: Prediction for Germany with fitted model parameters: reproduction number (R_0_{start} and R_0_{end} values), date of the steepest decrease in R_0 (x_0), rate at which R_0 decreases(k), and probability of going from infected to dead ($inf_to_dead_p$).

The main plot shows the infection, death and recovery curves of our SEIRD model using appropriate parameters. The plots below from left to right: R_0 plot, mortality rate plot, deaths per day plot, used ICU beds per day plot and deaths due to ICU bed shortage per day plot. The period for all plots is from January 2020 to the middle of May.



14. Scenario Studies

14.1 Implemented Measures

In the last step, the fitted model is used to simulate the spread of the virus with various prevention methods being implemented. Those methods affect the spread of the virus by reducing the expected amount of people an infected person infects per day (β). Recall in the equations in Section 13.2, where the parameter β is only present in the equations for the number of susceptible and the number of exposed people. The following prevention has been implemented to reduce β :

Reducing Social Contacts:

Maybe the most effective way of decreasing the spread of the virus is to limit social contacting. A short and lose restriction period increases the risk of an uncontrolled spread of the virus or the arise of a second wave which will lead to many deaths. But a too long and tight restriction can lead to economic and psychological incisions. Performing simulations can help to see the effect of social distancing with varying periods and intensities. For our simulation, we implemented social distancing by reducing β by 0/25/50/75%.

Lockdown:

The so-called lockdown is a special case of social reduction, in which social contacts are reduced to an absolute minimum by preventing the population from leaving their homes. But even for the case of a lockdown, peoples' social contacts cannot be stopped completely because part of the population like cashiers or hospital staff needs to go to work. Therefore we modelled the lockdown as a reduction of β by 90%.

Masks:

Several studies [53][57] state that wearing masks can effectively reduce the spread of the virus by 8%-16%. Each of the social distancing simulations has been performed with and without the usage of masks. Wearing a mask is modelled as an additional reduction of β by 12%.

14.2 Methods

After a fixed unrestricted time of 30 days the restriction period starts, followed by 100 simulation days with the length of the restriction varying within the simulation days by 0/25/50/100% of the simulation days. As before we initialized our model with one individual being infected while the rest of the population is susceptible.

14.3 Results

From Figure 14.1 we can make several observations:

- Without implemented restrictions (plot top left) the number of infected people is by far the highest and the number of deaths exceeds the linear growth
- Even minor restrictions for a small period of time significantly reduce the number of infected people
- The wearing of masks has has a greater effect in both absolute and relative terms when combined with less severe social restriction
- Increasing the time of social contact reduction has less impact than the intensity of social distancing
- The difference between 50% and 100% restricted simulated days is just minor when combined with masks and lockdown

The parameters that has not been fitted to the German data have been set as shown in Table 13.3.

14.4 Discussion

The simulations performed suggest that both the duration as well as the intensity of the restrictions play an important role when fighting the outbreak of the corona. The use of masks is even more important in scenarios with low social reduction. Nevertheless, the simulations seem to underestimate the true case. The number of people infected might be underestimated. As a reminder, the model was fitted to Germany with the data from the beginning of March to the end of April. At that time the German government already introduced several restrictions to keep the spread of the virus under control, such as advising people to stay at home, the closure of public locations and shifting of many jobs to the home office. Our data is therefore already being fitted to restrictions and then used to simulate restrictions, which doubles the effect of the restrictions implemented in the different scenarios. Surprisingly, without any simulated restrictions and with a model fitted to data from a period when restrictions were present in Germany, the model still simulates more than 10 million infected people within 130 days.

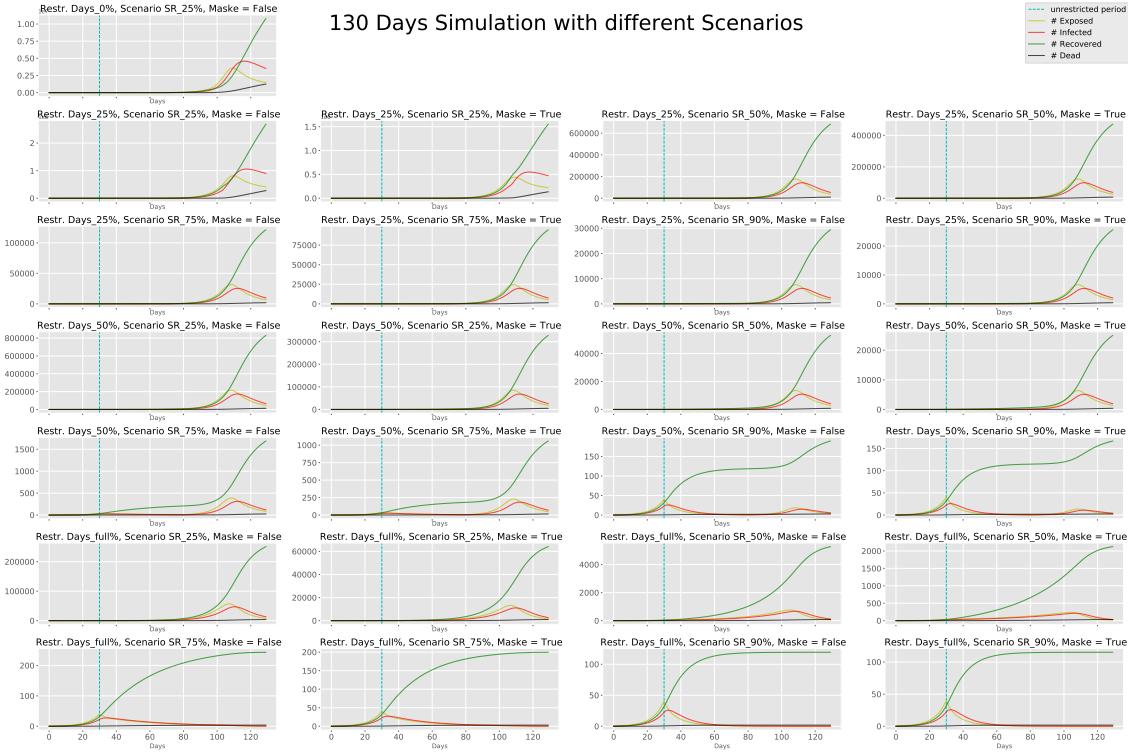


Figure 14.1: Each plot represents an independent scenario simulation and shows Exposed (yellow), Infected(red), Recovered(green) and Dead(black). The time without restrictions is indicated by the dashed blue vertical line. Note that the number of susceptible instances is not shown and the size of the population (y-axis) is not normalized to improve visibility. Each diagram . **Top left:** 0% restriction days, 0% social reduction, no mask; **second and third rows:** 25% restriction days, social reduction ranging from 25% to 90% and mask/without mask; **fourth and fifth rows:** 50% restriction days, social reduction ranging from 25% to 90% and mask/without mask; **sixth and seventh rows:** full restriction days, social reduction ranging from 25% to 90% and mask/without mask;

IV Project 4: Agent-based Simulation

15	Introduction	65
15.1	Background	
15.2	Goal of the Project	
15.3	Outcome	
16	Basic Principles	67
16.1	Agent-based Modeling for COVID-19 Spreading Simulations	
17	A simple ABM	69
17.1	Model Description	
17.2	Fitting to Real Data	
18	Extending the Model	73
18.1	Introduction	
18.2	Incubation and Exposed State	
18.3	Chronic Conditions and Comorbidities	
18.4	Central Locations	
19	Scenario Studies	77
19.1	Methods	
19.2	Results	
19.3	Discussion	
20	Comparison of EBM and ABM	85
20.1	Fundamental differences	
20.2	Case Study Covid-19	

15. Introduction

15.1 Background

Agent Based Modeling (ABM) has gained significance in the last 30 years due to ever increasing computational efficiency. It has a wide variety of applications including but not limited to biology, businesses, technology, social sciences and economics. The idea of ABM is based on simulating independent agents operating and interacting with each other within a micro scale computational model confined to a predetermined rule set. Especially within the field of epidemiology ABM's are characterized by their ability to capture heterogeneity of complex interactions between different agents [4]. The complexity of different agent behaviour can always be mirrored within the simulation by including new constraints or rules on the model. Thus, ABM's are very effective in modeling different epidemiological outcomes with different scenario rule sets.

15.2 Goal of the Project

The aim of this project is to use agent-based simulation to model the interactions of individuals within a population during the COVID-19 outbreak, so that one can determine how small changes in behavior and interaction can affect population level output. Different extensions (incubation and exposed state; chronic conditions and comorbidities; central locations) are implemented to refine the model. In the end, the variability of human behaviour can be shown with the purpose to understand the variability in the likely effectiveness of proposed interventions.

15.3 Outcome

The integration of central location had the largest impact of the added model extension. In the basic scenario without movement restrictions 90% of the population becomes infected by the virus after 21 days of simulation when supermarkets and schools are both opened. The ICU capacity was exceeded after 28 days. By applying different intervention strategies, the combination of social distancing and wearing masks has been confirmed to be the most effective.



16. Basic Principles

16.1 Agent-based Modeling for COVID-19 Spreading Simulations

Agent Based Models (ABM) can be implemented in very different fashions. For this week's project, we used the so-called simple billiard balls model (Silva) which is composed of a population of agents, within a loop where the agents run and interact. The agents are initialized with properties such as working place, age, or health conditions that drive their mobility patterns. As the name suggests the agents are represented as billiard balls that can transmit the virus when they get in touch with each other. The big advantage of this approach is its simplicity and modularity. On the other hand, the billiard balls model is very abstract and most likely produces wrong results. Nevertheless, all models are wrong, but some are useful (George P. Box) [10].

Taking this approach a step further the *Spatiotemporal Epidemic Model* introduced in April 2020 by Lorch [35] makes spatiotemporal predictions by making use of data from contact tracing technologies. By using Bayesian optimization the model estimates the risk of exposure based on the moving habits (e.g. going to a certain bar) of each individual, the percentage of symptomatic individuals, and the difference in transmission rate between asymptomatic and symptomatic individuals from historical longitudinal data.

Instead of modeling people as billiard balls, one can model a population as a network. The so-called *Network Based Model* (NBM) takes several assumptions into account, like social interactions, the probability to spread the disease, relationships between the individuals, immunity after infection, and many more. Based on those, a graph is built where each agent is represented as a node and the relationships between agents as an edge. The assumed baseline network structure is an input of the model but health policies, e.g. lockdowns, quarantine, etc., can be interpreted as (temporarily) changing the social network by eliminating edges.

17. A simple ABM

17.1 Model Description

The ABM should have agents (i.e. people) with the same characteristics like a national population. For this purpose, we selected the age distribution of USA, because it was easily accessible: 0-14 years: 18.62%, 15-24 years: 13.12%, 25-54 years: 39.29%, 55-64 years: 12.94%, 65 years and over: 16.03%), which were determined in 2018. In our model this was achieved by using the beta probability distribution with parameters $\alpha = 2$ and $\beta = 5$, so that the age is determined with $\sim \beta(2, 5)$. In python it can be done like this:

```
age = int(np.random.beta(2, 5, 1)) * 100
```

Our simple ABM is designed such that each agent must be in one of these states: Susceptible, Infected and Immune-Recovered, Dead. There are adjustable initial percentages of Infected (0.02) and Immune agents (0.01). The rest of the population has the status Susceptible. Once an agent is infected it takes one of the three sub-states that explain the severity of the infection: Asymptomatic, Hospitalization and Severe. Spread through contagion is determined by the interaction of the infected agents through proximity or contact. This means that the faster the agent moves, the greater the probability that he will approach an infected agent and become infected as well. A contagion distance defines the minimum distance (set to five units) that two agents must have for virus transmission to take place. The terrain where the agents are simulated is squared and bi-dimensional (100x100). Each agent is randomly created within this terrain. The initial horizontal and vertical position can be coded like this:

```
self.x = kwargs.get('x', 0), self.y = kwargs.get('y', 0)
```

During simulation, the mobility amplitudes can be defined for any possible agent status, in each iteration, each of the agents also moves randomly within the environment. Only the steps of its position are defined by the code:

```
x,y = np.random.normal(0, self.environment.amplitudes [self.status], 2)
self.x = int(self.x + x)
self.y = int(self.y + y)
```

The distance between two agents a_1 and a_2 is determined by

```
np.sqrt(((ai.x + self.positions[m][0]) - (aj.x + self.positions[n][0]))** 2  
+((ai.y + self.positions[m][1]) - (aj.y + self.positions[n][1]))** 2)
```

For **all** dead agents and **all** agents with the status Infected and with their severity of Hospitalization or Severe Infection, their movement will be set to zero.

Furthermore, the effects of mobility restrictions on the economy - especially on the income and wealth of individual agents - are simulated. The agents' income is simulated as a function of their mobility. Then the mobility of the agent is defined by the Euclidean distance from his previous position. The wealth of the agents is initialized according to the equal distribution of the society. This distribution is measured using quintiles, each quintile represents a social class: 20% most poor, poor class, working class, rich class and 20% richest. The minimum income defined by the first quintile of the poorest is used as the unit of expenditure and income of each stratum. During iteration, the wealth of each agent is reduced by its minimum fixed expenses, where the constant value in units of minimum income is proportional to its actual wealth, such as expenses = minimum_income [wealth quintile]. Furthermore, in each iteration, the wealth is increased by the daily income of the agent. The income is a random value which is proportional to his actual wealth and his mobility. Then, the final income is replaced by the minimum income [wealth quintile].

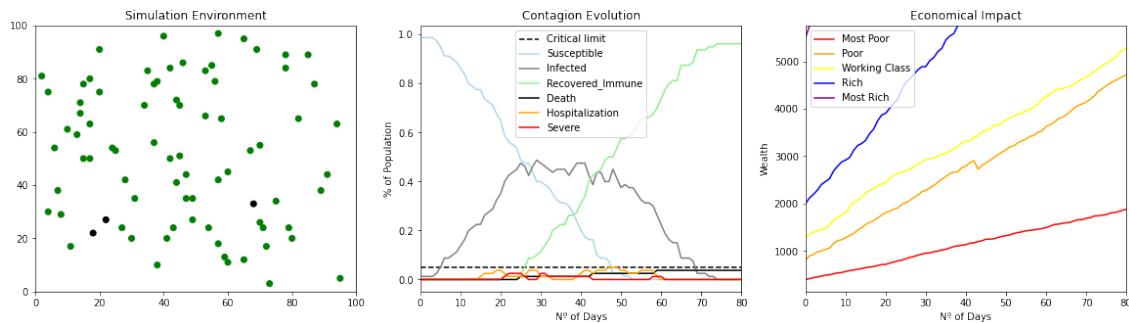


Figure 17.1: Simulation with COVID-19 ABS: First Run. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time; **Right:** A plot showing the economic changes of the agents with time

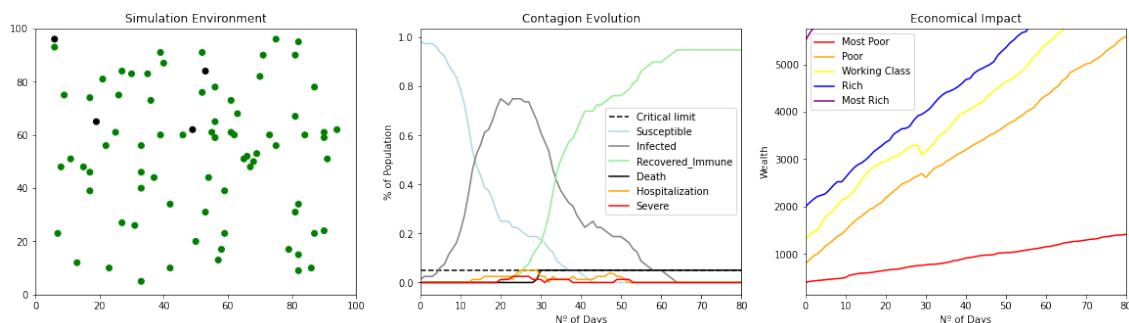


Figure 17.2: Simulation with COVID-19 ABS: Second run. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time; **Right:** A plot showing the economic changes of the agents with time

In order to obtain a reliable results for this weeks project, the simulation was executed 50 times simulations and the confidence intervals are displayed in Figure 17.3. Here is the plot for simple ABM 17.1. In the plot for the evolution of the contagion risk, you can see how much the critical

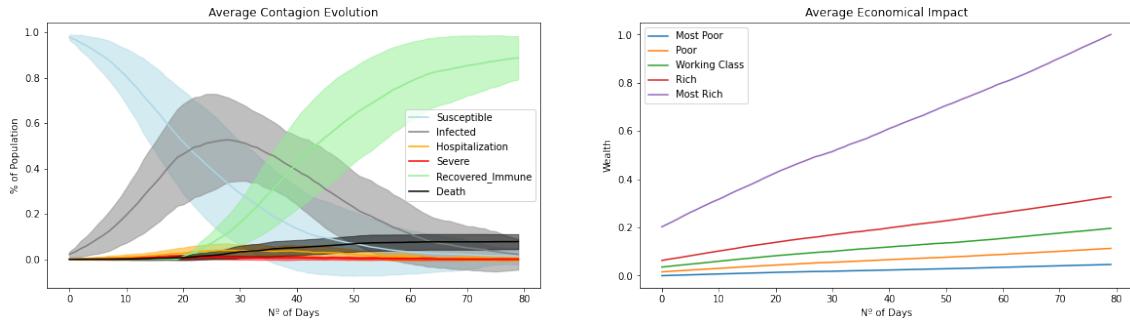


Figure 17.3: Simulation with COVID-19 ABS. Average results with confidence intervals for 50 executions. **Left:** A plot showing the health changes of the agents with time; **Right:** A plot showing the economic changes of the agents with time

limit of the health care system has been implemented and how many lives have been lost. This is the simulation of a catastrophic situation that will occur if nothing is done. The plot of the economic impact shows that it is not so bad, because the economy does not stop growing. If you compare the two diagrams for two runs, 17.1 and 17.2, you can see that the curves for the contagion status susceptible, infected and immune recovered already differ significantly. When comparing the plot for 50 executions 17.3, it is clear that curves from two plots for one run each are in the confidentiality area in the plot for 50 executions. In order to have reliable results, plot with confidentiality area should be used.

17.2 Fitting to Real Data

Unfortunately we could not fit the model to real world data. Since running the ABM is already a runtime expensive process we decided to set some parameters according to research papers in which the authors have already reported statistics for the current corona outbreak. Values like infection risk when being exposed, incubation time, infection duration could be set according to the data in WHO's Health Report [39]. Other values like the portion of immune individuals by the beginning of the outbreak are not detectable. Consequently we ran our simulations for those parameters with varying initial settings (see Figure 17.4). The scenario was build by setting the

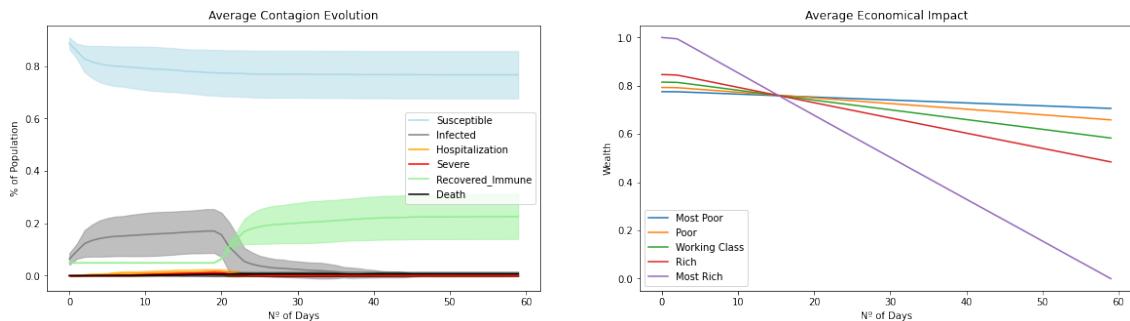


Figure 17.4: Simulation with COVID-19 ABS. Average results for 50 executions with confidence intervals fitting to real data. **Left:** A plot showing the health changes of the agents with time; **Right:** A plot showing the economic changes of the agents with time

percentage of initially infected persons to 2% and of recovered/immune people to 1%. To simulate a lockdown the mobility amplitudes of Susceptible, Recovered Immune were set to 0.5 and to 0 for the Infected when 5% of the population is infected. To adjust the time frames from 1.3.20-30.4.20,

iterations number was set to 60. Note that ABM approaches belong to the class of Monte Carlo algorithms whose results are generated using randomness and statistics. Therefore a bunch of runs is computed and the distribution within the results is evaluated. You can see that the infected number goes down by reaching the condition of 5% infected, which is reduced to about 0% after about 30 days, which is also the start of lockdown. The number of susceptible gradually decreases over the whole period. After about 30 days, the recovering immune increases abruptly and remains at the same level. The plot, which represents economic conditions, shows the decreasing trend, which also happened in reality.



18. Extending the Model

18.1 Introduction

One of the biggest advantages of the presented ABM approaches is its modularity. By adding properties to each agent and relationships between agents an entire society can be modelled. The simple ABM is a good starting point to understand the mechanisms of the model, but simplifies too much. The following extensions have been added to make the predictions more realistic:

18.2 Incubation and Exposed State

Two states have been introduced to capture the characteristics of virus spread. On the one hand, the *Exposed* state is used to count all individuals that have had contact with an infected person but have not infect themselves. The *Exposed* state can be used to measure how infectious the disease is by checking how many people who had contact with an infected person have infected themselves. A low ratio would indicate that the disease spreads only under tight conditions. Note, that the new state is interesting for our modeling but hardly applicable to real-world scenarios, as it would require a bunch of test kits and the contact chains of infected people needed to be tracked down. Neither of these requirements is available in any country at this time. On the other hand, we have added another Infection state to make our model predictions more realistic by implementing an *Incubation* state. People who become infected will go through an *Incubation* period during which they cannot infect other people. Concerning the 73th WHO's health report [39] for the current SARS-CoV-2 virus, most infected people went through an incubation period of 5-6 days. This implementation mostly delays the outbreak by preventing the spread of the newly infected people.

18.3 Chronic Conditions and Comorbidities

A second expansion to the model was performed by taking chronic health conditions of infected individuals and their effect on the fatality rate into account. Three common German health conditions were identified: obesity, hypertension and diabetes. For each risk factor the prevalence rates within the German population were taken from the literature (54% [38] for obesity, 33% [25] for hypertension and 13% [19] for diabetes). Next, each agent is initialized by random chance

with none, one or multiple conditions representing the true prevalence rates of the population. To accurately simulate the comorbidities for people dying of COVID-19 we used a report by Solis et al. [51]. Individuals under the age of 18 were assumed to be healthy.

comorbidities	death rate
0	5.58%
1	13.5%
2	23.2%
3	32.9%

Table 18.1: Death rates and comorbidities, the numbers were adjusted based on age due to inflated death values.

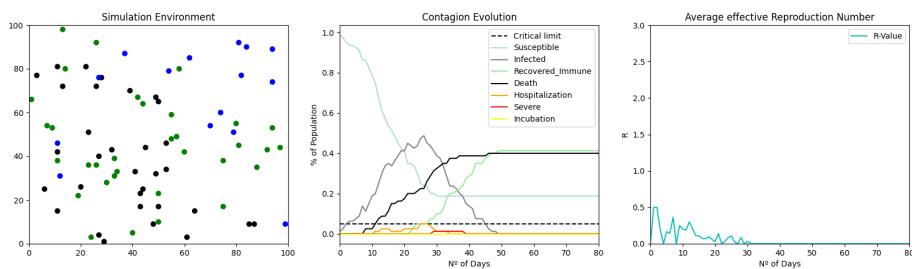


Figure 18.1: Simulation of chronic conditions effect on the base simulation. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time

By simulating the extension and comparing it to the base model a visual bump in the death rate can be seen (Figure 18.1). The dead agents are more than quadrupled, which suggests, that our values for the death rate are inflated. While we account for healthy individuals as well as people under the age of 18, the over-inflation might be due to the fact, that the report measured only hospitalized people with stays over 14 days. The value ranges are too specific for one subgroup of people. Thus the basic model, which uses death rates per age, is much more true in its simulation.

18.4 Central Locations

With the additions of central locations to the model, the agents are directed to predefined places at specific time points, instead of performing only arbitrary movement within the environment.

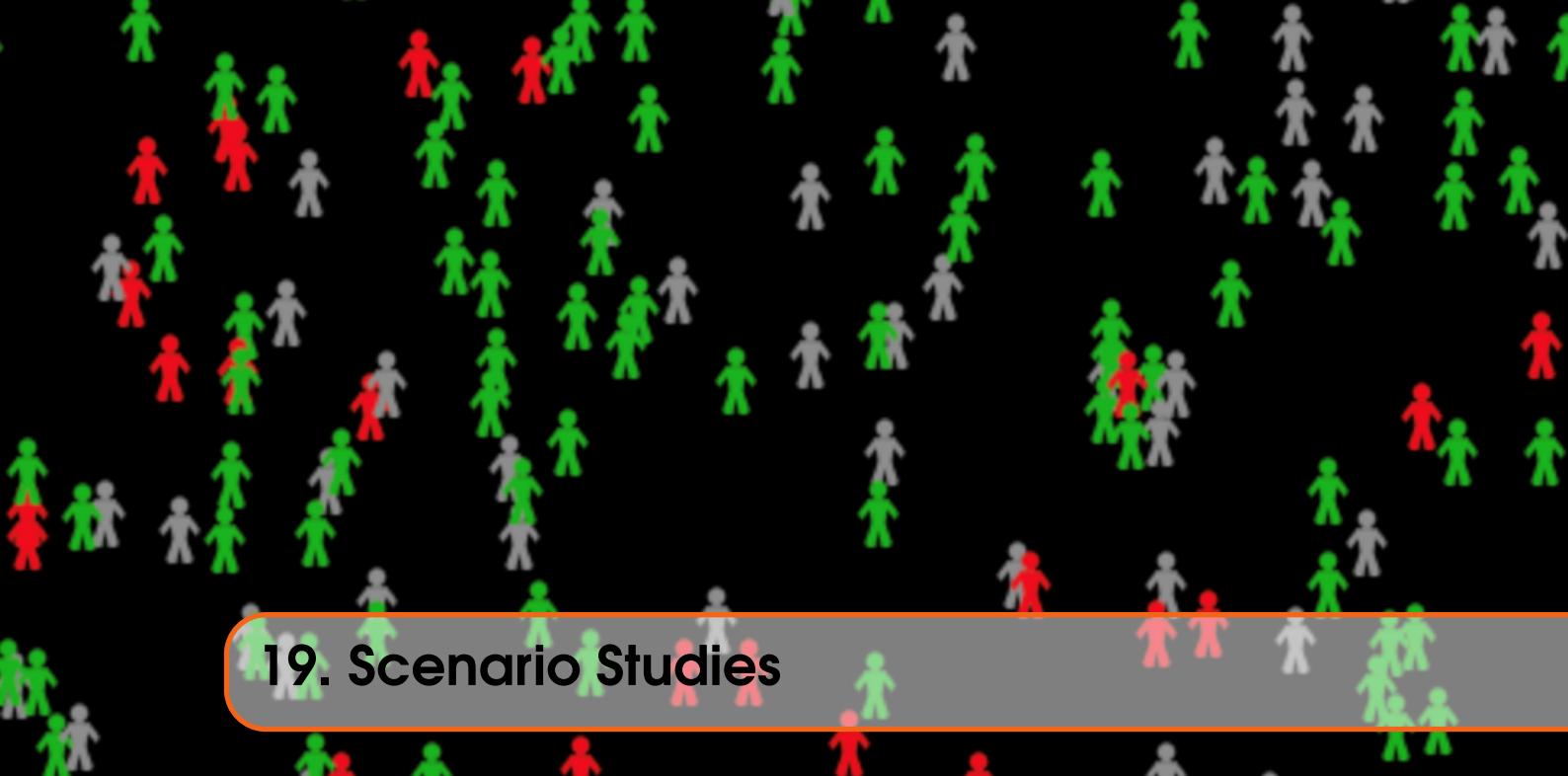
18.4.1 Supermarkets

Since all people have to provide themselves or their families with groceries, we decided to include supermarkets. Each agent with an age of 18+ (assuming that younger agents are supplied with food by older family members) has to go to the supermarket at least once every seven days. Based on the average retail sales area per 1,000 inhabitants in Germany [58] the size of the supermarket is adjusted in our model. For the number of 2438 inhabitants described later, a retail area of 3600 m^2 was created, divided into three different locations (2000 m^2 , 1000 m^2 and 600 m^2). To monitor the time (in days) an agent did not visit the supermarket, an attribute was added to the *agents* package. It is a simple counter that is reset when the agent is either placed in the supermarket after six days of absence or when the agent enters the supermarket area by a random movement. The counter is

initialized by random values between one and six, so that all agents have to visit the supermarket at different times. If an agent is placed in one of the supermarkets, his position within that area is randomly calculated, so that he is within contagious distance of some but not all others visitors of the supermarket. The probability that the agent goes to the biggest supermarket was set to 50%, while he is placed with 30% and 20% to the medium and small supermarket respectively. After a shopping day, the next position of the agent is calculated randomly within the entire environment.

18.4.2 School

Since only 18+ agents visit the supermarket, younger agents still perform only random movement in our model. Therefore, we created a school. Every agent with an age between 6 and 17 attend the school five of seven days a week. Here, in contrast to the supermarket all kids visit the school at the same days. We decided to place the school in an outer area of the environment, because it is less possible that people unintended come by a school (compared to the supermarket). The area of the school is 2000 m^2 and after five consecutive days in school, the new position is again calculated randomly. For the placement of the agents within the school the same concept as for the supermarket is applied, simulating interactions on the schoolyard only to some but differing school mates. The presence of adults (i.e. teachers) is not yet implemented.



19. Scenario Studies

19.1 Methods

The new extensions were tested by simulating different scenarios, where each scenario represents restrictions that have been implemented to decrease the spread of the virus. The following prevention have been implemented. For the reasons outlined in Section 4.4.2 we could not include the chronic expansion:

Reducing Social Contacts

Maybe the most effective way of decreasing the spread of the virus is to limit social contacting. A short and lose restriction period increases the risk of an uncontrolled spread of the virus or the arise of a second wave which will lead to many deaths. But a too long and tight restriction can lead to economic and psychological incisions. Performing simulations can help to see the effect of social distancing with varying periods and intensities. For our simulation we implemented social distancing by reducing the movement of the individuals by 60% after 10% of the population is infected. This change is reverted when only 5% of the population is still susceptible.

Lockdown

The so-called lockdown is a special case of social reduction, in which social contacts are reduced to an absolute minimum by restricting the population leaving their homes. But even in the case of a lockdown, peoples' social contacts cannot be stopped completely because part of the population, such as cashiers or hospital staff, have to go to work. Therefore, we modelled the lockdown as a reduction by reducing the travel amplitudes of the agents by 90%. The travelling reduction is revoked when a individual needs to go to a supermarket.

Masks

Several studies [53][57] state that wearing masks can effectively reduce the spread of the virus by 8% – 16%. Each of the social distancing simulations have been performed with and without the usage of masks. Wearing a mask is modelled by reducing the initial

contagion_rate by 16%.

Each simulation was executed with some fixed parameters that did not change for the different scenario simulations (Table 19.1). The simulation environment was initialized to represent 1 km² by setting width and height to 1000. The population size of 2438 was determined to reflect the population density of a German city. We chose Hamburg's population density [21] and adjusted our population such that each point in our graph still represents 1 m² while being displayed as 0.5 m² for a better visualization. The maximal distance between agents for contagion was defined to be 5 meters. This value does not correspond with the reality, but since each agent spends the entire day at the same position, the model would not provide meaningful results with a maximal contagion distance of 1.5 or 2 meters, which corresponds with the reality [48]. Each scenario was performed for a span of 80 days.

parameter	description	value
<i>w</i>	Width of the environment	1000
<i>h</i>	Height of the environment	1000
<i>pop_size</i>	Population Size	2348
<i>crit_limit</i>	Maximum percentage of population, which the Healthcare System can handle simultaneously	0.05
<i>dist</i>	maximal distance between agents for contagion	5
δ	Percentage of infected in initial population	0.02
β	Percentage of immune in initial population	0.01
<i>M</i>	Mobility ranges for agents by the beginning of simulation	5
<i>i_risk</i>	Prob of being exposed when being in contact with infected agent	0.9

Table 19.1: Initial parameter setup for the performed simulation, which are chosen to match Hamburg population density.

19.2 Results

The basic scenario with no restrictions as same as the implemented intervention scenarios were analyzed with and without the presence of central locations to estimate the risk of visiting supermarkets and opening schools. For all resulting figures, the left plot shows the agents represented as billiard balls while their color shows the agent's state after 80 simulation days. The centered plot displays the course of the health state portions for the entire population on each day. The right plot indicates the effective reproductive Number (R_t). It is the average number of secondary cases per infectious case in a population made up of both susceptible and non-susceptible hosts. If R_t exceeds 1.0, the virus will spread exponentially.

At first, we will compare the different model implementations for the basic scenario where no movement restrictions are applied. In Figure 19.1 it can be seen that when both schools and supermarkets are opened more than 90% of the population becomes infected by the virus after 21 days of simulation. The threshold for available ICUs is reached after 28 days, which cause an immense increase in the death rate. After 45 days, the entire population was infected with Covid: 95% recovered while 5% died. When only supermarkets are opened and schools are closed (Figure 19.2) the speed of the spread is slightly reduced. Nevertheless, the ICU capacity is exceeded after 29 days with the effect that more than 6% of the population died. At the end of the simulation, 20% of the people remain susceptible. Analyzing the model where supermarkets are excluded, but the school is opened, the infection curve is visibly flattened. The ICU capacity is never exceeded since

most of the infected people are kids (6-17 years), which have a lower probability of a severe course of the disease. Due to the fact that only young people are visiting the school nearly 50% of the entire population stays susceptible until the end of the simulation.

By examining the plots for the lockdown (Figure 19.4, 19.5, 19.6) and the contact reduction scenario (Figure 19.7, 19.8, 19.9) we unfortunately noticed that our code contains an implementation error. After visiting the supermarket, the next position of the agents is calculated randomly. Instead, we should store the position from where the agents were moved to the supermarket and reassign that to that location on the next day. The random position calculation after the supermarket shopping conflicts with the concept of lockdown and social distancing. Although the movement of the agents is reduced, most of the population becomes infected after 30 days of simulation and the results do not differ as much as they should compared to the basic scenario.

However, wearing masks (combined with social distancing) has an remarkable impact on the spread of the virus. In this approach the contagion rate is reduced by 16% which lowers the impact of contacts in central locations. In the model where only schools are opened (Figure 19.12), the young population (25-30%) becomes infected after 12 days (five days in school + weekend + five days in school) and recover very quickly. When the school is closed and only the supermarket is opened the infection curve is shifted to the right, because the people visit the supermarket at different days. Here, the peak of infected people is visible after 50 days with 30% of the population infected. When both institution are opened simultaneously, 40% of the population is infected after 30 days, but the curve rises slowly enough that the ICU limit is never exceeded.

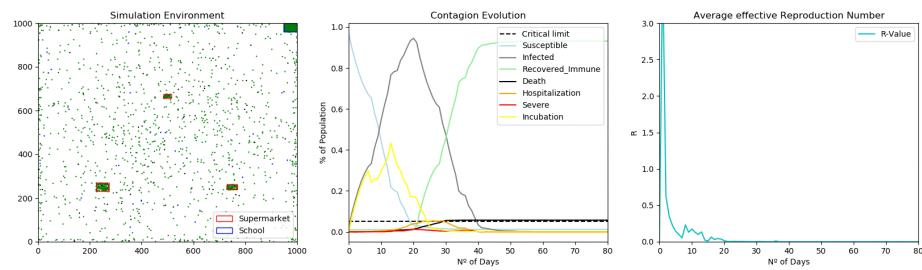


Figure 19.1: Basic Scenario with no restriction and the presence of three different sized **supermarkets** in the center and a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time

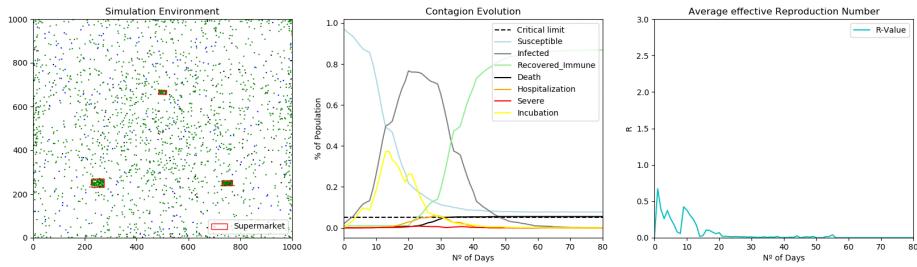


Figure 19.2: **Basic Scenario** with no restriction and the presence of three different sized **supermarkets** in the center of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state **Incubation**; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time

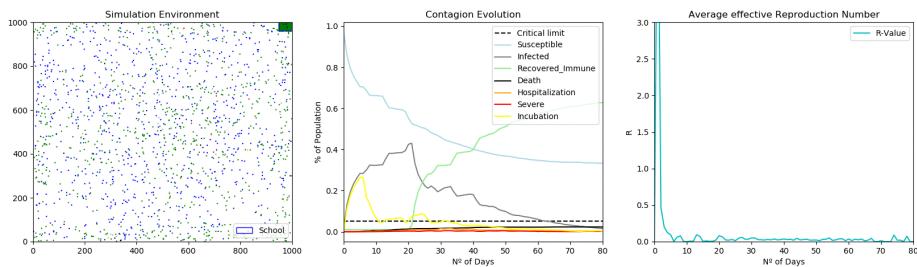


Figure 19.3: **Basic Scenario** with no restriction and the presence of a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state **Incubation**; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time

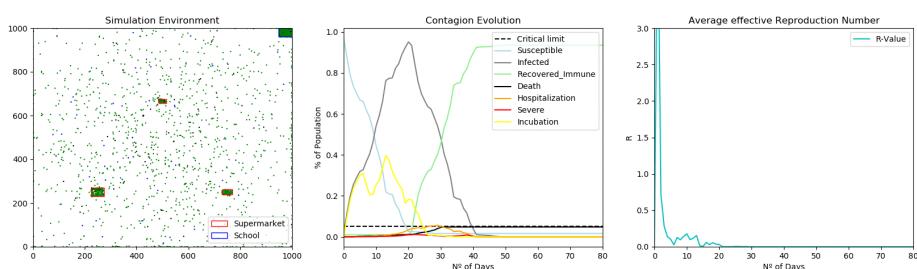


Figure 19.4: **Lockdown Scenario** with the presence of three different sized **supermarkets** in the center and a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state **Incubation**; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time

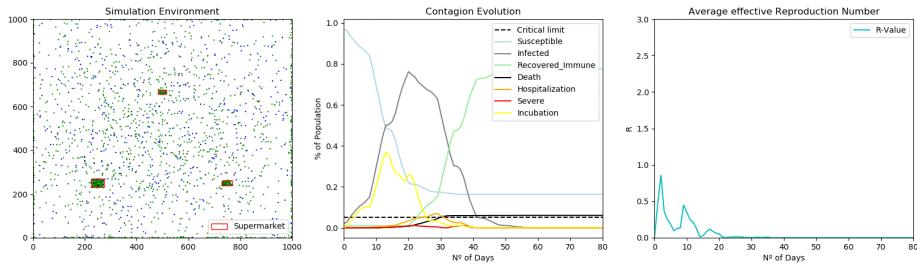


Figure 19.5: Lockdown Scenario with the presence of three different sized **supermarkets** in the center of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time

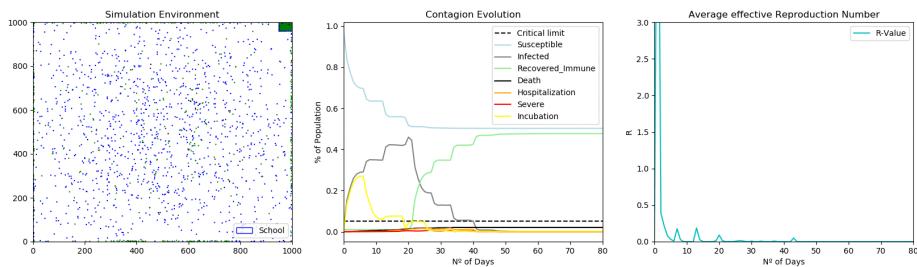


Figure 19.6: Lockdown Scenario with the presence of a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time

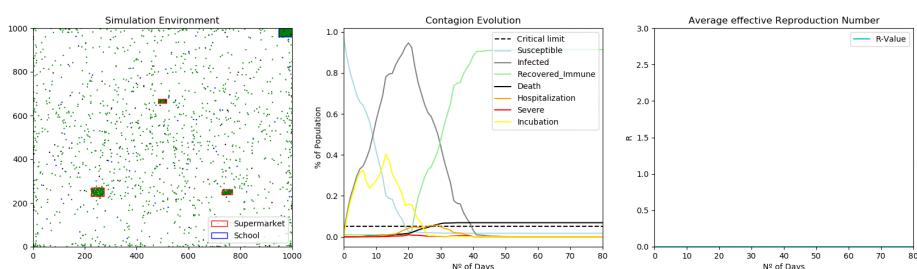


Figure 19.7: Social contact reduction Scenario with the presence of three different sized **supermarkets** in the center and a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time.

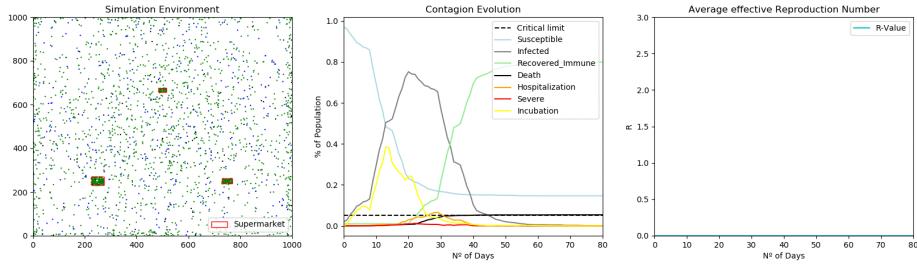


Figure 19.8: **Social contact reduction Scenario** with the presence of three different sized **supermarkets** in the center of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time.

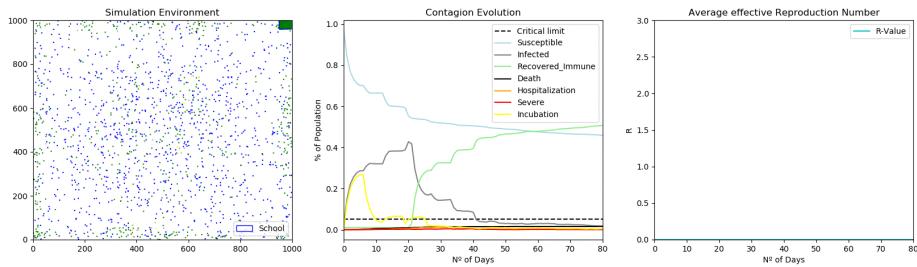


Figure 19.9: **Social contact reduction Scenario** with the presence of a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time.

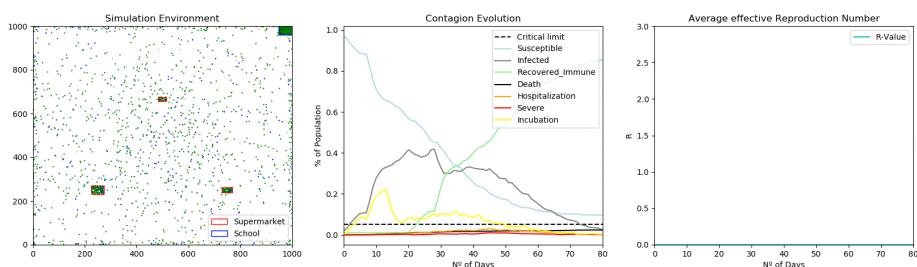


Figure 19.10: **Social contact reduction combined with wearing masks Scenario** with the presence of three different sized **supermarkets** in the center and a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time.

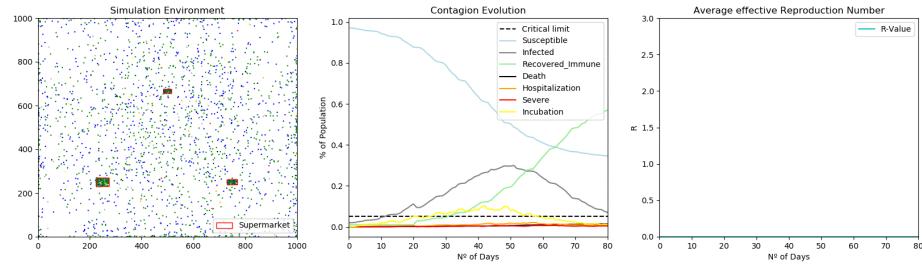


Figure 19.11: Social contact reduction combined with wearing masks Scenario with the presence of three different sized **supermarkets** in the center of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time.

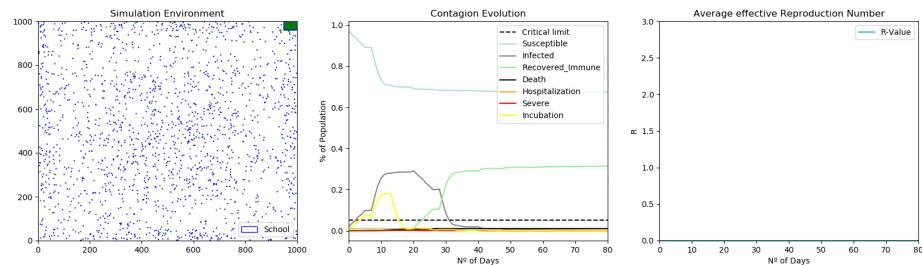


Figure 19.12: Social contact reduction combined with wearing masks Scenario with the presence of a **school** in an outer region of the environment. **Left:** A plot showing the interactions of the agents within the terrain; **Middle:** A plot showing the health changes of the agents with time including the new state Incubation; **Right:** A plot showing the changes of R_e (average effective reproduction number) with time.

19.3 Discussion

The results quickly demonstrate the need to integrate central institutions into the model, because of their large impact in agent-based systems. With different age groups (supermarket: 18+, school: 6-17) assigned to the visit of different central locations another effect becomes visible. If older people have to visit crowded places, the amount of available ICUs can be reached very fast, which need to be avoided at all costs. The analysis of the basic scenario clearly shows, that interventions are absolutely necessary. Unfortunately, the results of the lockdown and social distancing were difficult to analyze due to the implementation error. Nevertheless, the impact of wearing masks on the simulation results clearly demonstrates the effectiveness of intervention strategies to reduce the spread of the virus.



20. Comparison of EBM and ABM

20.1 Fundamental differences

In the final task we draw comparisons between an ABM and EBM for the recent COVID-19 outbreak. Both models are based on a common concept: modeling entities and observables within a complex system over a temporal timeline. While the basic concept is the same, their level of attention to the relationships between entities and the abstraction level of the system itself is different. Beginning with the modeling of entities, a typical approach with ABM would be going bottom up. Each agent itself is understood as an individual, with inherent properties and rule sets for interactions defined by observations. Using these definitions the model is build and simulated, showing us the macroscopic view of the individual interactions of the agents. In contrast, EBM could be defined as a top down approach. The system itself is modelled with complex equations, each entity not understood as an individual but as a compartment of subgroups.

These fundamental differences can already be observed in the computational power needed for simulation. The ABM runtime grows exponentially with higher population numbers, which is due to the nature of individual agent modeling. EBM models can cope with big data sets because only simple equations need to be solved. The individual nature of ABM allows the observer to follow singular agent interactions, thus giving a unique microscopic insight into the epidemiological effects of the disease. This could lead to new insights into which factors are responsible for the spreading of the disease as well as the infection rate. The EBM model does not allow this microscopic view and is very rigid in its simulation practices as no deviations from the set of equations are possible. On the other hand, ABM can capture the inherent stochasticity of real world systems. For example, agents might make decisions quite similar to real world individuals. This stochasticity can also have a negative influence, introducing noise into the system or leading to implausible outcomes. Alone with our simulations we had multiple outcomes where the epidemic did not start off at all or some subgroups exploding. The effect of introducing a new rule set for interactions can have dramatic influences on the ABM system at a whole, thus the only way of fine tuning ABM was to subtly tune the parameters with ongoing runs. This is a very time expensive endeavour. Expanding the model in population and rule sets might make it impossible.

Event-Based (Discrete) Modeling	Agent-Based Modeling
Macrospecifications reveal microstructures (top-down view)	Microspecifications generate macrostructure (bottom-up view)
Externally observable phenomenon (events)	Autonomous decision making entities (agents)
Programmed response to discrete events	Programmed functionality of agents
Events adhere to system-level observable information	Agents adhere to behavioral rules (boundedly rational)
System of interest changes state in response to events	Agents function independently and flexibly
Event impacts the entire entity	Agents interact as distinct parts of simulation
Simplicity in modeling inputs, state, and outputs	Simplicity in modeling rules
Internal behavior is unknown	Events emerge
Easy to test	Difficult to validate

Figure 20.1: Comparison of ABM and EBM characteristics. The Figure is taken from [51].

20.2 Case Study Covid-19

Due to COVID-19, there is a deep and continuous spread of infections between infectious and susceptible people. If we show any negligence in the control measures, the outbreak will start to grow rapidly within no time. The infection is supposed to grow if the people infected by the infection is greater than one. However, the only solution for this is people developing immunity. It is possible that the curve of infection starts up again after a continues period of low transmissions. This is called a second wave. This happens mainly due to the negligence of the people in not following the suitable control measures like social distancing, mobility etc. By considering some suitable scenarios and applying the simulation with suitable parameters it is possible to simulate a second wave. The measure of immunity within a population changes the possibility of a second wave happening. Both models take immune people into account. In ABM it is easy to create a change point within the system to model an upcoming second wave. By simply increasing certain parameters for infection rates, travel restrictions, social distancing and exposition rates a second wave can be simulated. In contrast, EBM does not have an easy way to simulate a second wave. A change point has to be defined and the equations have to be introduced beforehand. Overall, ABM is more suitable for modeling a second wave then EBM.

Project 5: Time-Series Prediction for COVID-19 Cases

21	Introduction	89
21.1	Background	
21.2	Goal of the Project	
21.3	Outcome	
22	Model-based vs Data-based	91
22.1	Model-based Approaches	
22.2	Data-based Approaches	
23	Approaches	93
23.1	Prophet Library	
23.2	LSTM Neural Networks	
23.3	Classical Models	
24	Comparison of Data-based TSP Approaches	99
24.1	Data	
24.2	Results and Discussion	
25	Model- vs. Data-based Time-Series Prediction	103
25.1	Fitting the Model	
25.2	Results for fitted Model	
25.3	Discussion	
26	Towards COVID-19 Outbreak Prediction	107
26.1	Outlook and Evaluation	



21. Introduction

21.1 Background

A time series is defined as a set of observations arranged in chronological order. The time interval between each data point remains constant, such that a sequence of discrete-time data is generated [28]. One aim is to extract meaningful statistics and interesting characteristics from the time series data structure. Thus, investigating the stationary and seasonality is part of the standard protocol when dealing with time-series data. It is said to be stationary if it's mean and variance does not change over time, while seasonality can be identified depending on the data, which refers to periodic fluctuations of the values. The second major intention is to perform forecasting. Here, a model is applied to the data to predict future values based on past observations.

21.2 Goal of the Project

This week's project aims to forecast the COVID-19 outbreak using a classical approach, the Prophet library, and a machine learning technique. The methods are applied to three different time series data sets: two data sets each generated by a SIR and an ABM model and another data set containing the confirmed COVID-19 cases for Germany from the beginning of March until the end of April 2020. The results are evaluated by comparing the forecasting plots with the actual data as reference and analyzed using the root mean square (RMSE) values. Additionally, a comparison between data-based and model-based prediction methods is performed by investigating the performance of the best SIR model of the previous week's project on the same data.

21.3 Outcome

LSTM turned out to be the clear winner and produced the lowest RMSE values for the prediction on the real-life data set ($\text{RMSE}=0.007$) and the SIR data set ($\text{RMSE}=0.0$) while Prophet achieved the best prediction results on the ABM data set ($\text{RMSE}=0.0155$). The comparison of data-based and model-based approaches (LSTM vs SIR model) also confirmed LSTM to be the better option for the given data.



22. Model-based vs Data-based

22.1 Model-based Approaches

Data- and model-based approaches are different ways to target the same problem. In the context of COVID-19, they can be used to predict the number of infections that will take place in the future. Data-driven predictors convince with their simplicity of implementation. If a sufficient amount of experimental data is available, they fit the existing curve using mathematical calculations and make meaningful future predictions. However, if not enough or only bad quality data is available (e.g. at the beginning of an epidemic), a data-driven will most likely struggle to make a reliable prediction. With rising complexity of a model more complex scenarios can be modelled.

22.2 Data-based Approaches

This is where the advantages of model-based approaches come to the fore. High prediction precision can be achieved by modeling any kind of scenario, that has not yet taken place in reality (e.g. behavioral constraints like social distancing or wearing masks). The dynamic of the states can be estimated and predicted at each time point, which makes models more versatile in comparison to data-driven approaches. Nevertheless, applying model-based approaches consume in general more computational power and are more time-consuming. Summing up, at different times in an epidemic either data-driven or model-based approaches can be the better choice.



23. Approaches

23.1 Prophet Library

The Prophet model [55] is composed of the three components *trend*, *seasonality*, and *holiday effects*, where each of its components is added together with time as its regressor:

$$y(t) = \text{trend}(t) + \text{seasonality}(t) + \text{holidays}(t) + \text{error}(t)$$

where:

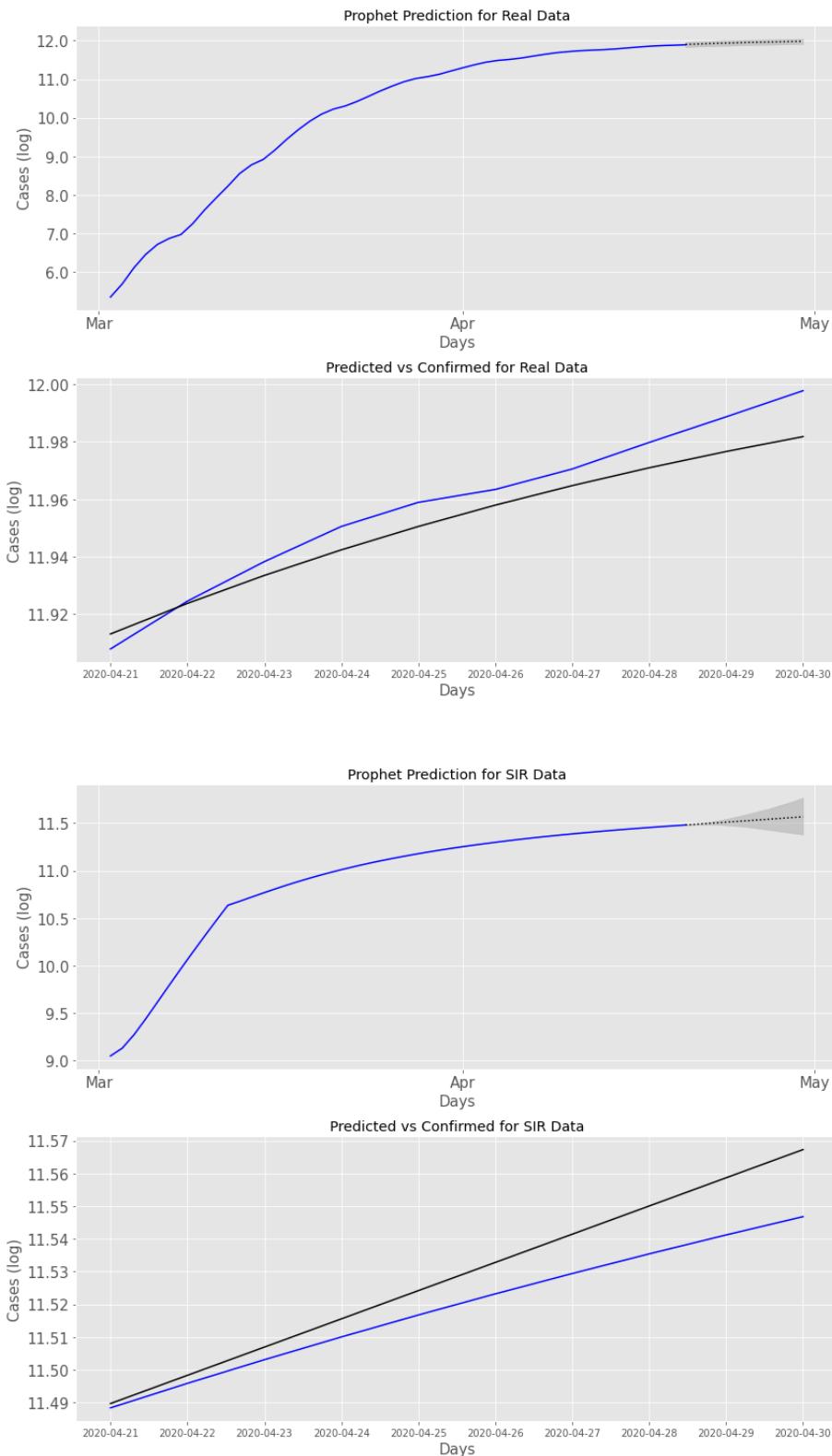
Trend models *non-periodic* changes. This component is the most important one for our analysis since our data has no seasonality in our data like it is the case for i.e. the amount of rain in a certain country over many years.

Seasonality models *periodic* changes (weekly, monthly, ...). Since our data sets include only 50 days of records, there is no seasonal trend present yet. This component would be more useful when we fight COVID-19 for the next upcoming years and use this data to predict years that lie even more in the feature. Over the years there might be an increase over the colder months of the year and decrease of new cases over warmer months of the year.

Holiday Effects allow the model to include short time intervals with an abnormal trend. When modeling corona in Berlin with no restrictions, this could be a public event like the Karneval der Kulturen where many people from different households interact closely with each other, which would probably lead into a spike in the number of confirmed cases. Since Germany restricted most public events early in the pandemic this component is also not of greater interest for our analysis.

Prophet forecasts by fitting a curve to the input data that is then prolonged for future values with the three mentioned components. Although, Prophet aims to produce a strong forecast without much hyper-parameter tuning. For our analysis, the prediction did not fit the actual trend of the test data. We achieved better results by changing the *growth* parameter from linear to logistic. Since the three data sets do not contain many breakpoints (points in the data, where the growth changes)

the nonlinearity of a logistic function can nestle closer to the slow decrease in the growth of the confirmed cases of the test data.



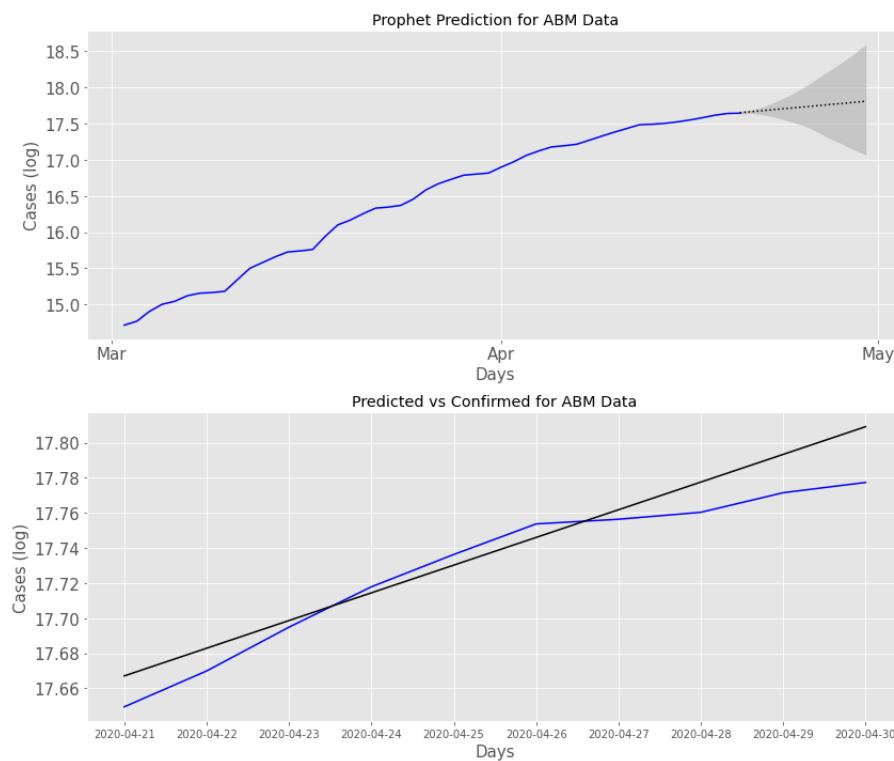


Figure 23.1: In each pair of figures, the results of the prophet’s time series prediction is displayed. The top plot shows in blue the recorded data and black the forecast. The dark gray area represents the confidence intervals. The bottom plot zooms into the forecast depicts the test data to the predicted data. The results were achieved as introduced in Section 23.

23.2 LSTM Neural Networks

Long Short-Term Memory Networks (LSTM) is a prominent deep learning method for time series prediction. LSTM is part of Recurrent Neural Networks (RNN). In contrast to common feed-forward networks where each layer is only connected to the subsequent layer, RNN has layers connected to itself or even previous layers. This is more closely modeled on the neural connections exhibited by the neocortex and allows the network a greater prediction accuracy for time-coded data.

The LSTM network is employed by using the python library *keras*. For the prediction, an LSTM of the length of the training data with a simple 1-dense hidden-layer is created. The model itself is optimized with adam and loss by the mean_squared_error. The data needed to be preprocessed. The first step of preprocessing is called shifting. In shifting the time series data is shifted by a constant value, dividing the data in a training value and expected value. Subsequently, because LSTM needs to evaluate local differences in the data, each time step is subtracted by the shifted value. As in all neural networks, a common scalar is used normalizing the data between -1 and 1 with a mean of 0. We plotted the predicted and actual time series data for the last 10 days. For all three test sets the RMSE is calculated (Figure 23.2, 23.3, 23.4).

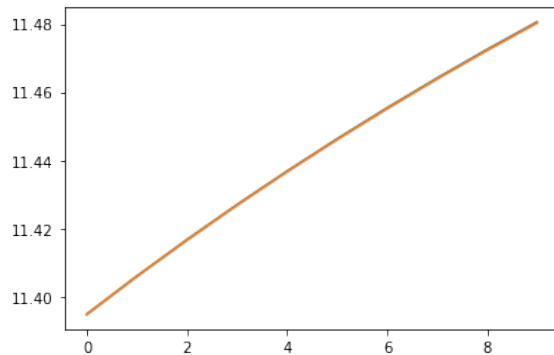


Figure 23.2: Predicted (orange) vs expected (blue) on SIR simulated data for the LSTM model as introduced in Section 23.1. (y-axis = log of infected people, x-axis = days)

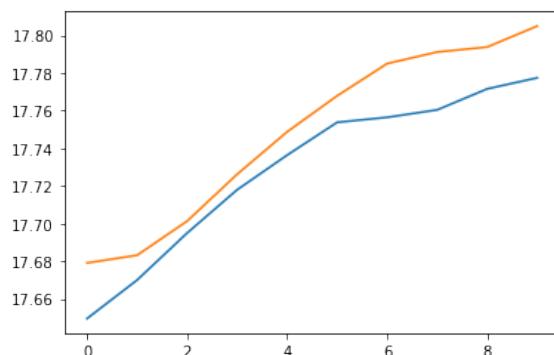


Figure 23.3: Predicted (orange) vs expected (blue) on ABM simulated data for the LSTM model as introduced in Section 23.1.(y-axis = log of infected people, x-axis = days)

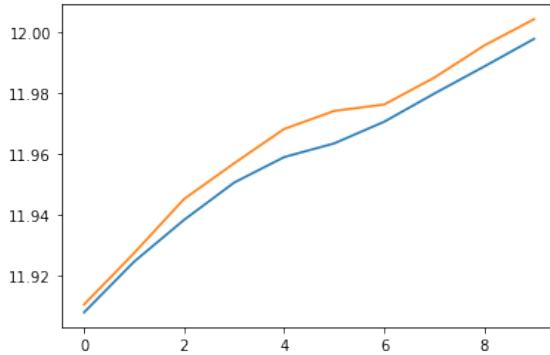


Figure 23.4: Predicted (orange) vs expected (blue) on RKI actual data for the LSTM model as introduced in Section 23.1. (y-axis = log of infected people, x-axis = days)

23.3 Classical Models

While there are various classical models they all focus on various linear relationships. From these models, we chose the AutoRegressive Integrated Moving Average method (ARIMA), which captures the auto-correlation in the data series by modeling it. One benefit of ARIMA models is that it tries to reduce the overall residual sum of squares (RSS) which leads to a small variance. The basic ARIMA approach processes the data in three main steps. The data has to be: (1) differenced to remove seasonality, (2) regressing of variables with its past values and (3) random movements are removed by using a moving average. Those three steps are represented by the parameters p (number of lag observations included in the model), d (number of times a raw observation is differenced), and q (window size) that lead to the classical presentation of the model: ARIMA(p,d,q).

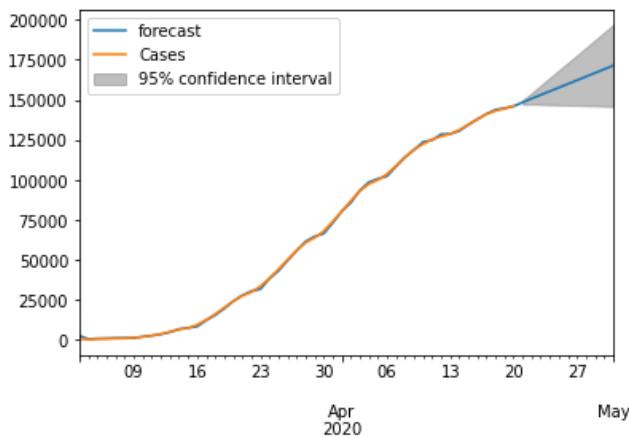


Figure 23.5: Cases (orange) vs forecast (blue) on RKI data for the ARIMA model as introduced in Section 23.2 (y-axis = confirmed cases, x-axis = time period). The gray area represents the 95% confidence interval for the 10 days forecast.

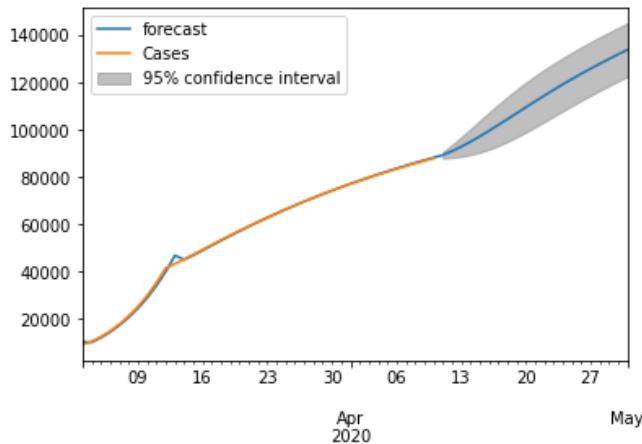


Figure 23.6: Cases (orange) vs forecast (blue) on SIR simulated data for the ARIMA model as introduced in Section 23.2 (y-axis = confirmed cases, x-axis = time period). The gray area represents the 95% confidence interval for the 10 days forecast.

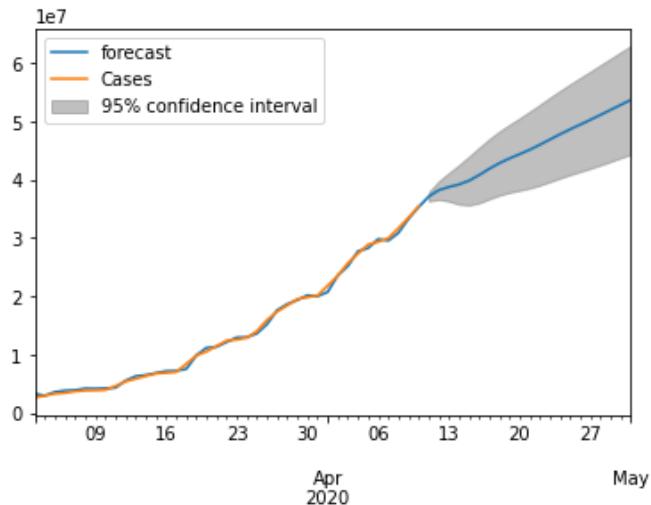


Figure 23.7: Cases (orange) vs forecast (blue) on ABM simulated data for the ARIMA model as introduced in Section 23.2 (y-axis = confirmed cases, x-axis = time period). The gray area represents the 95% confidence interval for the 10 days forecast.

At first, we imported the data sets, which contain 60 observations (days) of the confirmed number of corona cases. In the next step the data was processed to be stationary by applying the three processing steps. Afterward, we plotted Autocorrelation (ACF) and Partial Autocorrelation (PACF) to understand the parameters (p, d, q) of the ARIMA model. They basically indicate the correlation between two-time instances as well as the degree of association. Finally the model was trained and then evaluated using the RMSE. Figures 23.5, 23.6 and 23.7 represent the ARIMA model plotted for the RKI, SIR and ABM data sets respectively.



24. Comparison of Data-based TSP Approaches

24.1 Data

To evaluate the forecasting approaches three different data sets were created. Each data set contains the accumulated confirmed cases for the COVID-19 pandemic for 60 days.

Actual Data for Germany:

The actual COVID-19 case numbers for Germany were downloaded from a git repository hosted by the RKI [45]. The analysis of this project have been performed on the reported COVID-19 cases from 02.03.2020 till 30.04.2020.

SIR Data Set:

The SIR data was generated using the setup and the parameters as described in 17.2. The initial population size was set to the population size of Germany and the number of initially exposed people was set to 8000 to model an ongoing spread. The parameters have been chosen to make the data set more comparable to the real data from RKI. Note, that exposed means for the case of the SIR model, that an exposed individual will be infected at some point. In contrast to the ABM model, where exposed is defined as individuals that had contact with an infected person but will not necessarily get infected as well.

ABM Data Set:

To generate a third data set, an ABM approach was used that simulates central locations (supermarket and school) and models a scenario where the social contact is reduced and masks are worn. The number of agents was set to the population density of Hamburg (2438 inhabitants per square kilometer), but the results where upscaled concerning the entire German population, which explains the high number of infections in this data set.

In Figure 24.1 one can see that the curves of the RKI data set and the SIR data set are more flattened compared to the ABM data. For testing, all three data sets are split into a training and a test data set by the last 10 days.

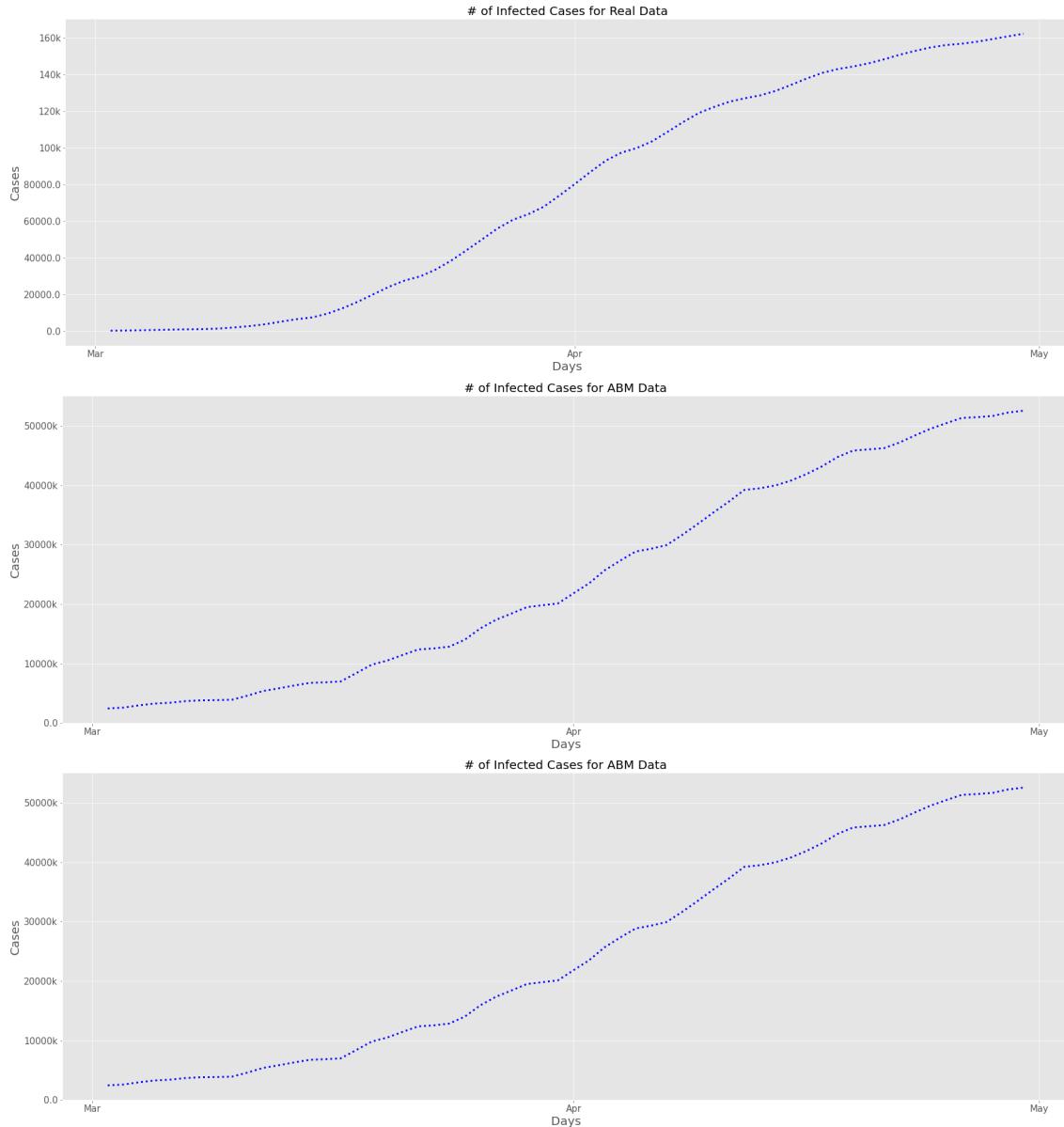


Figure 24.1: Number of confirmed cases for the three data sets in section 24.

24.2 Results and Discussion

Prophet :

In Figure 23.1 it can be seen that the fit of the Prophet model is rather suboptimal, especially for the ABM data. The fluctuating trend of the AMB and the RKI data set is not captured at all and the only reason for the small RMSE values is the low amplitudes of the fluctuations. The best fit we expected for the SIR data whose trend is almost linear because it fits the general trend of the Prophet predictions. However, the predicted period is overestimating the true values. After ten days the residual is highest for the SIR data set among the results for the three data sets.

LSTM :

LSTM is the "*winner*" and produced the best RMSE values for the SIR and RKI data set

(Table 24.1). For the SIR simulated data, it is due to the almost linear trajectory (Figure 23.2), where neural networks have an effortless fit. For the RKI data, it could also reproduce a delayed flattened of the curve at day four (Figure 23.4). The high ABM value is due to the high fluctuations within the data where the influence of the local differences might be too high. Still a delayed flattening of the curve was predicted (Figure 23.3)

ARIMA :

It can be seen in Figure 23.5, 23.6 and 23.7 that the fit is bad for the SIR data set when compared with the RKI and ABM data (Figure 23.6). By interpreting the RMSE of all three predictions, it can be seen that the RMSE is very high for the SIR model showing there are high variations in the predictions when compared with the results from the other data sets. The obtained RMSE differs significantly with 0.022, 0.147, and 0.028 for logarithmized RKI, SIR, and ABM data sets respectively.

Data set	Prophet	LSTM	Classical
SIR	0.011	0.000	0.147
ABM	0.0155	0.021	0.028
RKI	0.008	0.007	0.022

Table 24.1: Comparison of RMSE values for three data sets and three models as described in Chapter 22.2.

63.772

44.291

48.991

44.870

31.012

26.417

27.00

3.877

21.21

20.556

25. Model- vs. Data-based Time-Series Prediction

25.1 Fitting the Model

For the fitting of the data-based prediction model a SIR model was used that was fitted to Germany's reported cases in a previous step (see Section 13.1). The data used for the adjustment was generated for the period 03.02.2020-30.04.2020, using both the other SIR model and the ABM model.

The obtained data and the parameters are already known, and it is necessary to define initial estimates with lower and upper limits for the unknown ones in order to support the curve fitter. For the fit, a function is needed that takes exactly one x-value as the first argument (the tag) and all parameters to be fitted. It returns the confirmed cases predicted by the model for that x-value and parameters so that the curve fitter can compare the model prediction with the exact data. To perform the fit, it is necessary to initialize a curve-fitting model, set the parameters according to the initial, minimum, and maximum, specify a fitting method (e.g. *leastsq*). One of the important parameters is the *outbreak_shift*. The case data starts on 02.03.2020, so our model assumes that the virus started to spread on this day. For this reason, this parameter was set to zero. Others were the parameters *R0_start* (initial value of R_0 for the period), *R0_end* (final value of R_0), *k* (rate of decline of R_0), *x0* (start time of the decline of R_0), *inf_to_rec_d* (daily probability of transition from infected to recovered state), which influence R_0 and the infection rate (β), respectively, and *inf_to_dead_p* (rate of transition from infected to dead state). In Figure 25.1 one can see a plot of the fit for the extracted data from the sir model and the ABM model respectively.

The parameters predicted by the fitter are shown in Figure 25.2. Most of the parameters are nearly the same for both fitted data. Only the parameter *R_0_start* and *x0* differ significantly between them. The *R_0_start* of 2.76 for SIR data is more expected. The high value for *R_0_start* of 10.68 for the ABM may be possibly explained as an effect of upscaling the ABM data. Also, the numbers of confirmed cases from SIR data and from ABM data differ significantly, which may also be due to the upscaling of ABM data. The parameter *x0* (begin of R_0 decline) for the SIR data was fitted to 17, and for ABM data to 9. The difference might be caused by the different scenario setting of the two compared models. The possible problem is that the data were created with a modified SIR model that was adopted for scenarios with a parameter that changes the rate of transition from susceptible to the exposed state due to the restriction and its duration. For this reason, the data may

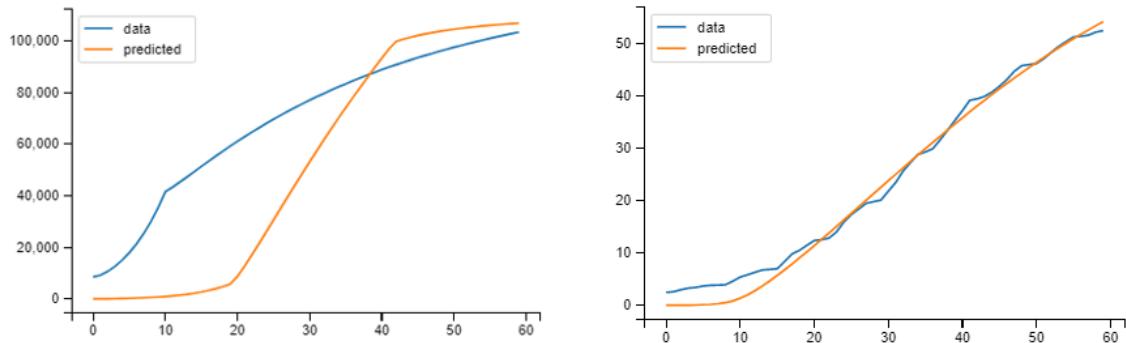


Figure 25.1: Fitting of the SIR model to the ABM (right) and SIR (left) data set. Predicted (orange) vs expected (blue). The y-axis shows confirmed cases for SIR data, confirmed cases $\times 10^6$ for ABM data and the x-axis shows the days. The SIR model is described in Chapter 24.2 and the data sets in Section 24.

not fit together well. In Figure 25.2 it can be seen in the right plot that the fit to the data generated with the ABM is much tighter.

25.2 Results for fitted Model

By predicting for 10 days it can be seen that the confirmed cases for the ABM data are higher than for SIR data (see Figure 25.3). The curves of expected and predicted for SIR data are not optimally adjusted ($rmse=0.052$). The curves of expected and predicted cases for ABM data are better adjusted ($rmse=0.015$). For the full comparison see Table 25.1.

25.3 Discussion

In comparison to the SIR model-based prediction, LSTM was able to better predict SIR data ($rmse=0$) but ABM data the prediction of SIR model ($rmse=0.015$) was better than LSTM ($rmse=0.021$). Taking all together, both approaches were able to deal with non-stationary data as the number of confirmed cases are. However, the SIR model can not always make optimal forecasting even for data extracted from another SIR-model with a different parameter setup.

```
{'R_0_end': 0.3413676232485514,          {'R_0_end': 0.6903763998947229,
'R_0_start': 2.7584555587301667,           'R_0_start': 10.67957469641348,
'inf_to_dead_p': 0.16000000000000003,        'inf_to_dead_p': 0.16000000000000003,
'inf_to_rec_d': 0.32371753521779184,         'inf_to_rec_d': 0.20168712000334932,
'k': 22.681536008041032,                      'k': 19.9999962543517,
'x0': 17.06903340839481}                     'x0': 9.177538495703322}
```

Figure 25.2: Fitted parameters for SIR simulated data (left) and ABM simulated data (right) for the SIR model in Chapter 24.2.

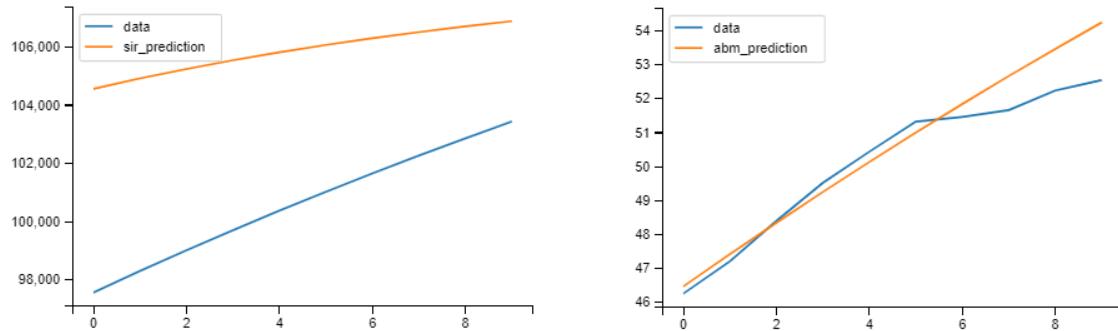


Figure 25.3: Prediction (orange) for SIR simulated data (left) and ABM simulated data (right) for the SIR model (see Section 24.2) of confirmed cases in contrast to the actual number confirmed cases (blue). The y-axis shows confirmed cases for SIR data, confirmed cases $\ast 10^6$ for ABM data and the x-axis shows the days.

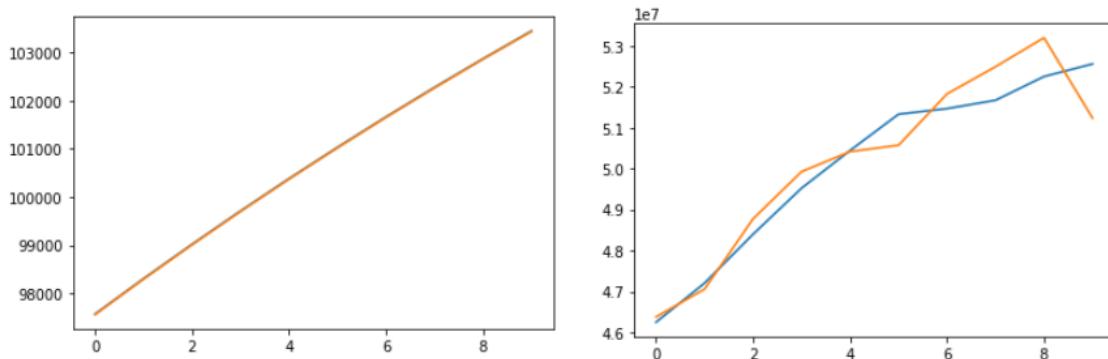


Figure 25.4: Prediction (orange) for SIR simulated data (left) and ABM simulated data (right) for a LSTM model (see Section 23.1) of confirmed cases in contrast to the actual number of confirmed cases (blue). The y-axis shows confirmed cases for SIR data, confirmed cases $\ast 10^6$ for ABM data and the x-axis shows the days.

Data set	SIR-model	LSTM
SIR	0.052	0.000
ABM	0.015	0.021

Table 25.1: Comparison of RMSE values for the three time-series prediction models as described in Chapter 22.2 and 24.2.

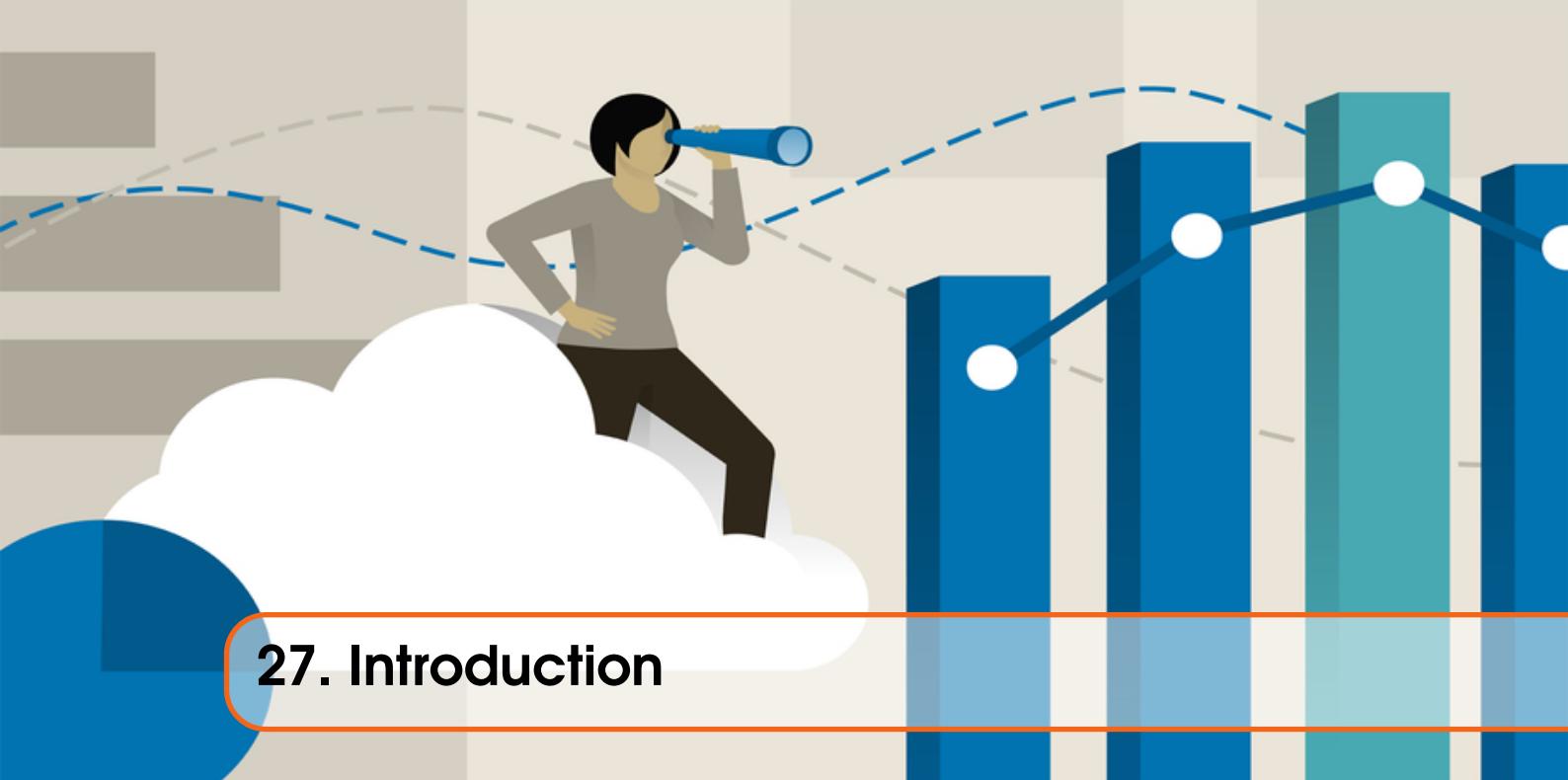
26. Towards COVID-19 Outbreak Prediction

26.1 Outlook and Evaluation

The results of the forecasting using the distinct time series approaches ARIMA, Prophet, and LSTM as well as the model-based approach using a SIR model were evaluated by comparing the forecasting plots with the actual data as reference and analyzing the root mean square values. ARIMA performs well on forecasting stationary data of short periods. Prophet is more advanced than ARIMA and offers also the possibility to identify trends and seasonality. Unfortunately, our data sets only contained 60 days of time series data. Since seasonality is already proven to exist for other viruses (e.g. influenza [34]), it will be an interesting experiment to analyze the existence of seasonality for COVID-19 on long-term time-series data. Such a study would distinguish Prophet's strengths. However, the data we used was non-stationary. That could be an explanation of why ARIMA underperformed and showed the biggest RMSE values for all three data sets compared to the other methods. LSTM was the clear winner as it is optimized for dealing with non-stationary data and our results confirmed: LSTM produced the lowest RMSE values for the RKI and the SIR data set. The results further demonstrate that, given data of confirmed COVID-19 cases, LSTM can learn and scale to more or less accurately estimate the amount of the people that will become infected in the future. Naturally, the prediction is only a tendency and need to be scrutinized very critically. Nevertheless, it is possible to predict a course of the outbreak. And the more data becomes available, the better the results of the data-driven forecasting methods will be.

Project 6: Time-series Prediction for COVID-19 Cases II

27	Introduction	111
27.1	Background	
27.2	Goal of the Project	
27.3	Outcome	
28	Solution Approaches	113
28.1	Data	
28.2	Visual Exploration	
28.3	Time-Series Prediction via Prophet	
28.4	Clustering	
29	Results	115
29.1	Visual Exploration	
29.2	Time-Series and Prophet Prediction	
29.3	Clustering	
30	Evaluation	129
30.1	Project Rating	
30.2	Problems	



27. Introduction

27.1 Background

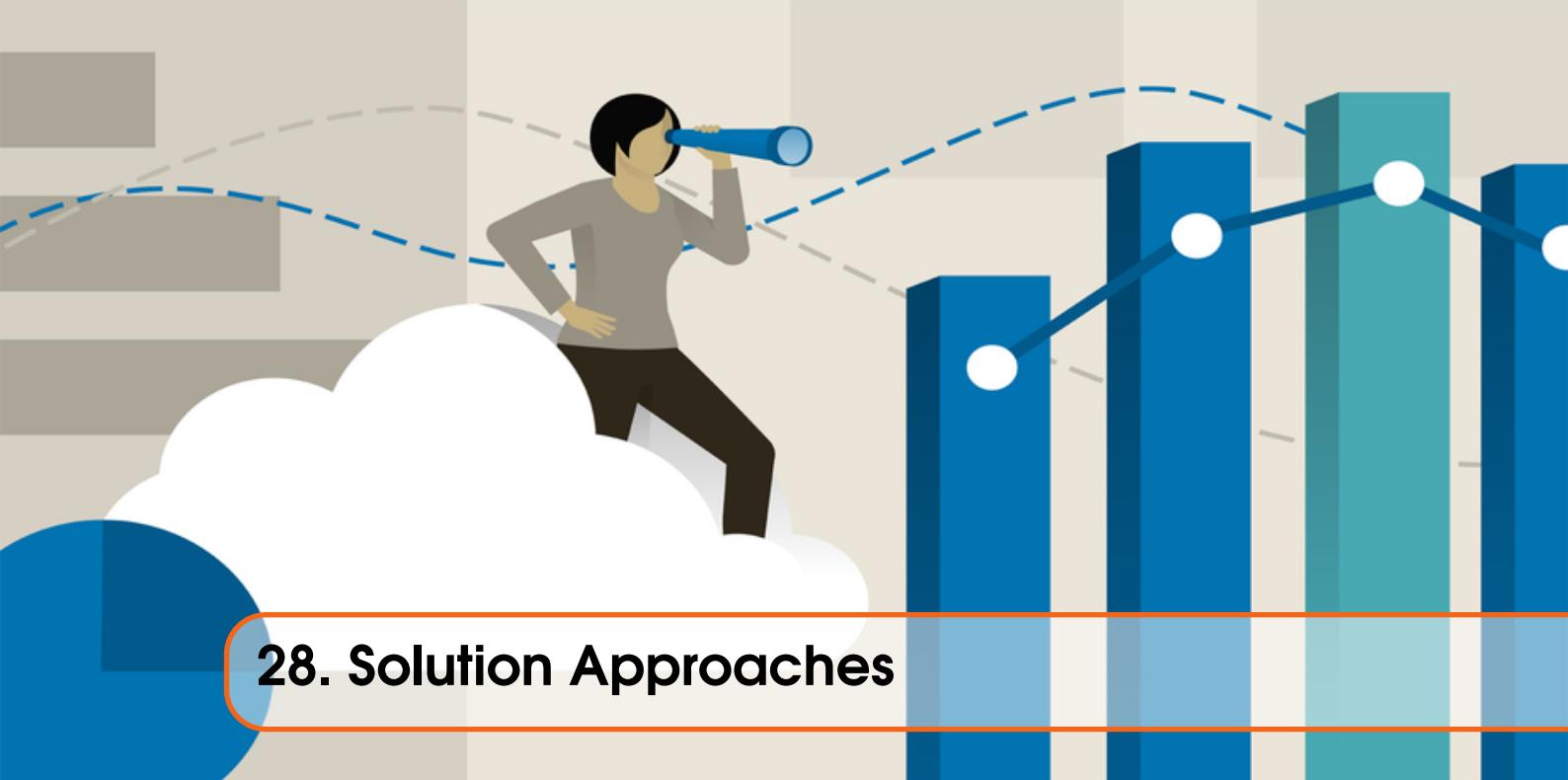
After exploring several time-series prediction methods in the last week's project, we will focus in this week's project on their visualization. Recall, forecasting is done by developing models that capture the characteristics of present data to make future predictions. The underlying mathematics and statistics of those models can be very overwhelming for people from other fields such that an easily interpretable presentation is an important but sometimes neglected part of the forecasting pipeline.

27.2 Goal of the Project

The aim of this weeks project is to combine time-series prediction methods as introduced in the last week combined with various visualization techniques. The time-series prediction is performed with Facebook's prophet library (Section 23) on the data given by RKI [45] for the federal states of Germany. Additionally, clustering is performed to group the federal states of Germany by the confirmed number of cases, deaths, and recovered individuals. The visual data exploration involves tree-maps, age and federal state dependent bar plots and GIFs that illustrate the change of the confirmed number of cases over time.

27.3 Outcome

Visualizations with bar charts and tree maps for the number of infected cases and deaths in Germany and its federal states and distribution of the same for different ages and gender. The resulting plots represent a detailed exploration of the outbreak. A clear distinction between case trajectory curves of different federal states for Germany could be drawn. The time-series prediction via prophet showed an increase in overall cases for each federal state. A **GitHub Repository** has been created that contains the three mentioned GIFs.



28. Solution Approaches

28.1 Data

The data set is provided by the Robert Koch Institute (RKI) site [45] and for the analysis, Germany and its federal states are considered. The time-span is defined from the 28th of January till the 15th of June. The starting date is based on the first occurrence of the coronavirus within the German borders in the state of Bavaria.

28.2 Visual Exploration

In order to understand the severe disease and to analyse the outbreak closely, a user friendly data visualization model with the help of different plots related to the various criteria was used. The various criteria are based on parameters like number of confirmed cases and deaths. When performing the visualizations we had some key thoughts in mind:

Which states in Germany are mostly affected?

How the confirmed and death cases are distributed all over the Germany?

Which countries have the most deaths?

Comparison of the cases with respect to age groups for various federal states

As a first visual overview each per state case trajectory was plotted (Figure 29.10). Two national counter measures were integrated, first being the general closure of schools and public spaces beginning with the 16th of March and the second being the introduction of mandatory masks beginning with the 27th of April. Next, we began plotting the data in the form of bar charts, here we obtained the visualizations for the number of infected cases and deaths in Germany and its federal states and distribution of the same for different ages and gender. The resulted plots represent a detailed exploration of the outbreak. Also represented the visualization of the number of confirmed, death and recovered with the tree map. Additionally, three GIFs are generated for the number of confirmed, death, and recovered cases. The time-span is defined from the 3th of March till the 19th and the code is taken from the **GitHub Repository et al. Chang Chia-huan** and extended by also

plotting the number of recovered individuals.

28.3 Time-Series Prediction via Prophet

While masks are still mandatory a lot of countermeasures were eased up in the month of June, thus we performed a time series prediction for each individual state to follow up on the trajectories of cases (Figure 29.11). For the time series prediction we used a 14 day future prediction from the facebook prophet library. For the parameters seasonality was disabled because the explored time scale is to small to account for trends due to seasonal changes. All other parameters have been remained unchanged.

28.4 Clustering

Clustering can be done in several ways. To get a quick insight into how data can be clustered, one way is to form hierarchical clusters. A hierarchy can be formed in two ways: start from the top and split or start from the bottom and merge. It was decided to perform the latter. First, the raw data containing confirmed cases, deaths, for each day for the rough period from March to 15th of June were imported for all German states. The data had to be made consistent such that the dates were synchronized for all federal states. Clustering should be done using time series data such as confirmed cases, deaths and recovered cases. To do this, you should first define the linkage function. The linkage function takes the distance information and groups pairs of objects into clusters based on their similarity. The parameters of the linkage function like metric and method may be adjusted and are set to "euclidian" and "ward" by default correspondingly. The keyword 'ward' causes the linkage function to use the algorithm to minimize the ward variance, and the keyword 'Euclidean' causes the Euclidean method to be used as a distance measure. These newly formed clusters are then linked together to form larger clusters. This process is repeated until all objects in the original data set are linked together in a hierarchical tree (dendrogram). So a dendrogram is a plot of clusters by hierarchical clustering, where the length of the bars represents the distance to the nearest cluster center. To find similarities between time series, it is worth using k-means to cluster them, since k-means is a kind of unsupervised learning and one of the most popular methods to combine unmarked data to k-clusters. The process starts with k centroids, which are randomly initialized. These centroids are used to assign points to the nearest cluster. The mean value of all points within the cluster is then used to update the position of the centroids. These steps are repeated until the centroid values stabilize. Before performing k-means clustering, it is necessary to identify the corresponding optimal number of clusters k. The elbow method as one way to estimate the value of k. To give equal importance to all features, it is needed to scale the continuous features, which will be done using the StandardScaler. For each k-value k-mean values will be initialised and the inertia attribute will be used to identify the sum of the squared distances of the samples to the nearest cluster center. As k increases, the sum of the squared distances tends toward zero. By plotting k values against the sum of the squared distances a plot should have a form of the arm and the optimal k should be then at the elbow.

29. Results

29.1 Visual Exploration

As a first overview it can be seen that the number of deaths for Germany is quite low compared to the confirmed cases, meaning that most infected individuals outlived the disease. Also, the confirmed cases and the recovered go hand in hand which is seen in Figure 29.1. The three states

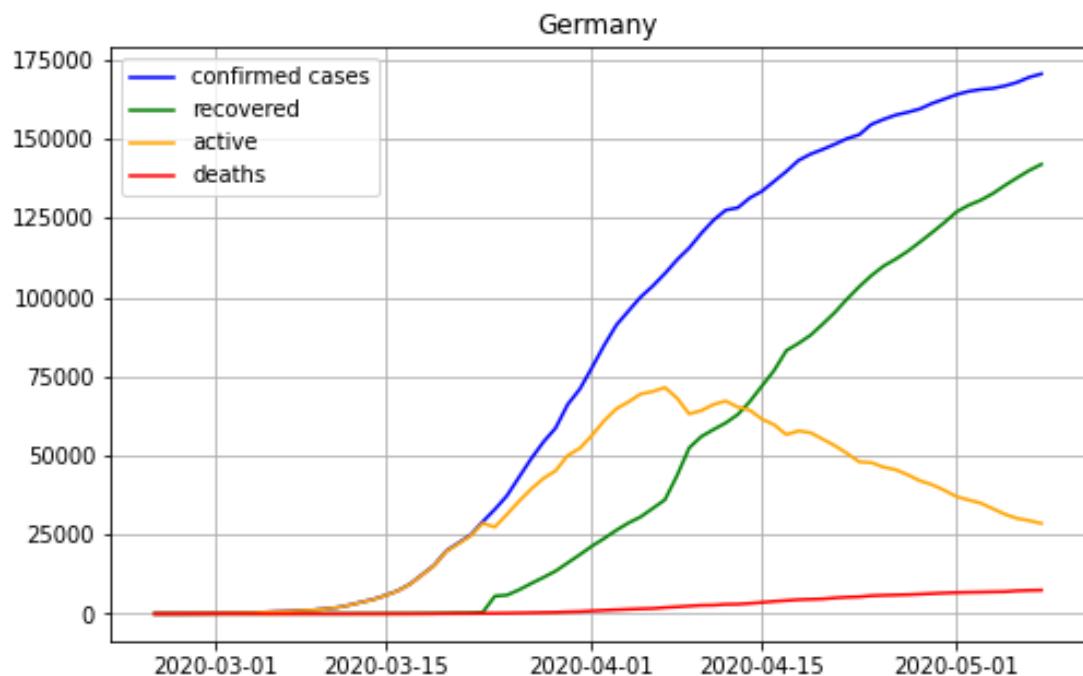


Figure 29.1: A trajectory plotting for entire Germany. The number of confirmed cases, active, recovered and deaths were visualized as dashed lines with different color codes. The number of confirmed cases and recovered cases go hand in hand

visualized Bayern, Baden-Württemberg, and Saarland (Figure 29.2) are on the top with the highest number of cases, slightly followed by Hamburg and so on. The deaths are minimum as most of the preliminary measures were taken place soon after the spread. The next plots show the distribution

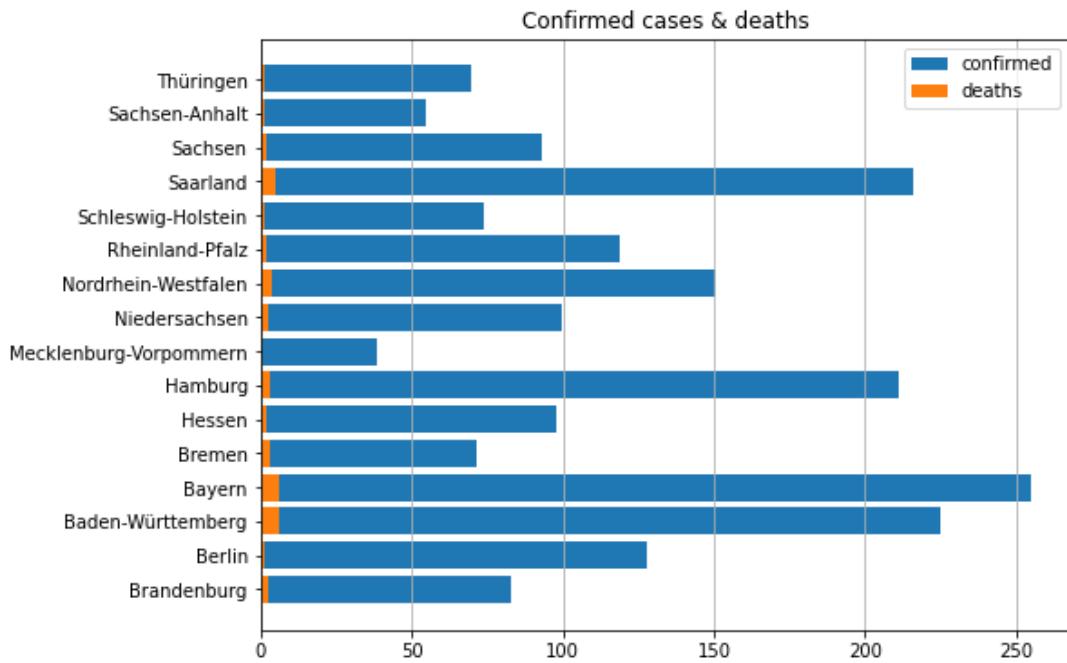


Figure 29.2: Bar chart representing confirmed cases and fatalities per 100k population. x axis: number of cases, y axis: federal states of Germany. Blue bars: number of confirmed cases, orange: number of deaths.

of confirmed cases and deaths for different ages and gender for Germany and its federal states which are shown (Figure 29.3, 29.4 and 29.5). It can be seen that the risk is increasing notably when the individuals were older than 40. Another spike can be observed for patients older than 80, where the death rate rises to a concerning possibility of 25%. When the federal states are considered then the results vary consistently. However some states like Mecklenburg-Vorpommern, Rheinland-Pfalz and Thueringen has shown more deaths with the men signifying the dominance in the male death rate.

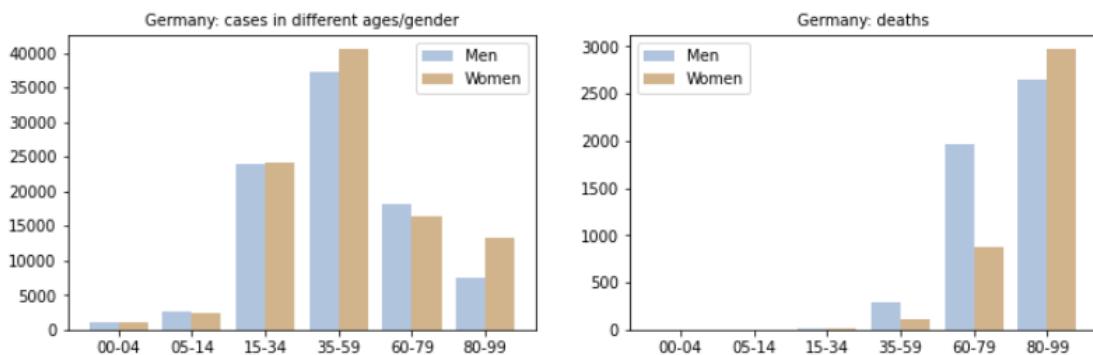


Figure 29.3: Bar charts showing the distribution of confirmed cases and deaths for different ages for men and women in Germany. x axis: age groups, y axis: number of cases.

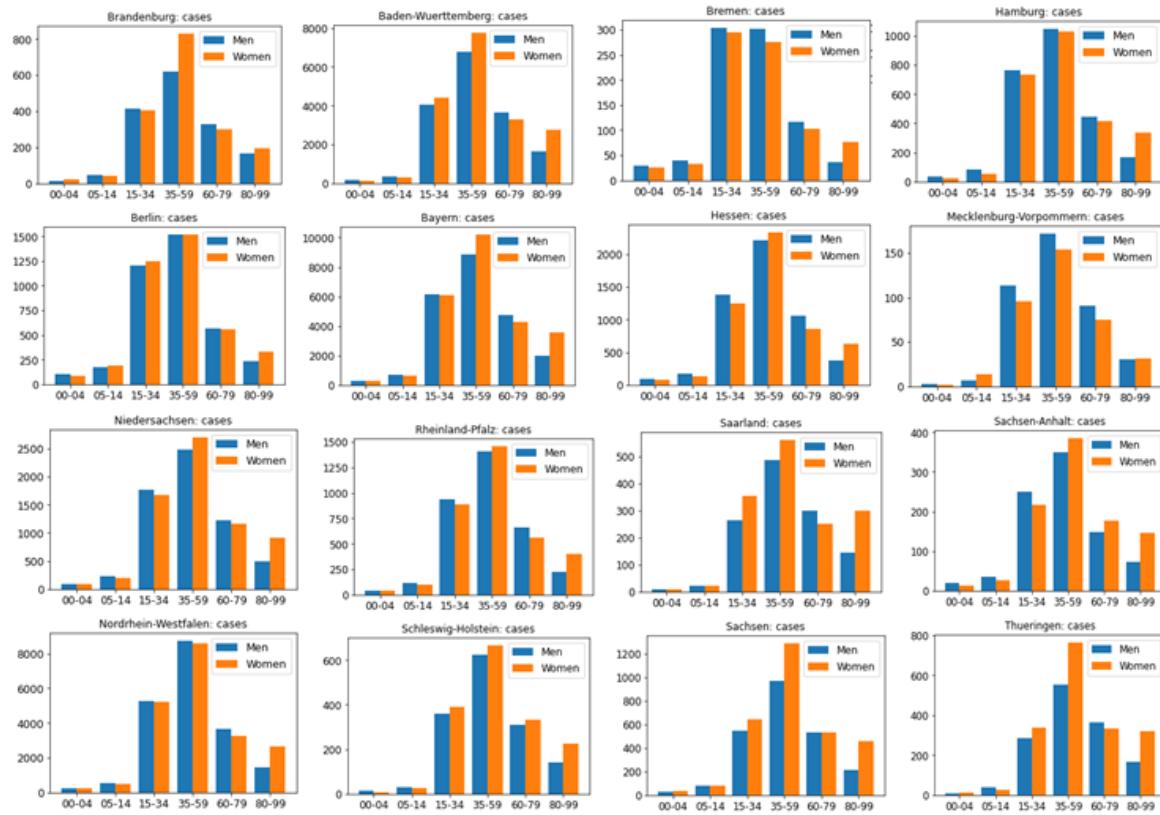


Figure 29.4: Bar charts showing the distribution of confirmed cases for different ages and gender for all the federal states of Germany. x axis: age groups, y axis: number of cases. Blue bars: number of cases in men, orange: number of cases in women

The treemap gives a brief understandings of distribution of number of confirmed, deaths and recovered cases for the federal states of Germany which is shown in the Figure 29.6.

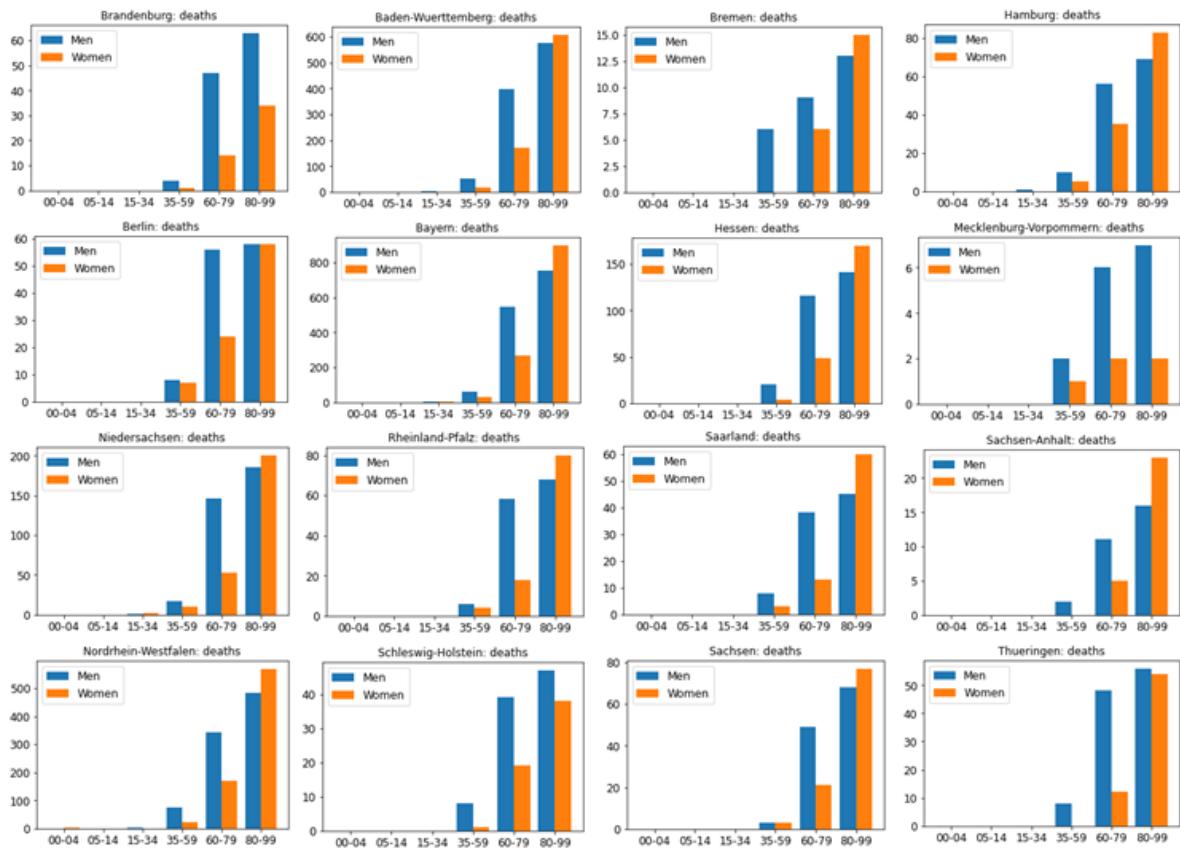


Figure 29.5: Bar charts showing the distribution of deaths for different ages and gender for all the federal states of Germany. x axis: age groups, y axis: number of cases.

As already pointed out in Section 29 the number of confirmed cases varies through the federal states and peaks in the south of Germany. Consequently we have similar results for the death and recovered statistics. To simplify the comparison between the federal states the reported number of cases are normalized by 1M residents. For example Bremen (top right corner of Germany) has less than 600k residents but the 1300 confirmed cases lead it to the federal state with the relatively highest number of conformed cases, while it has absolutely one of the lowest. The final GIFs for all three conditions can be found in our public [GitHub Repository for Week 8](#).

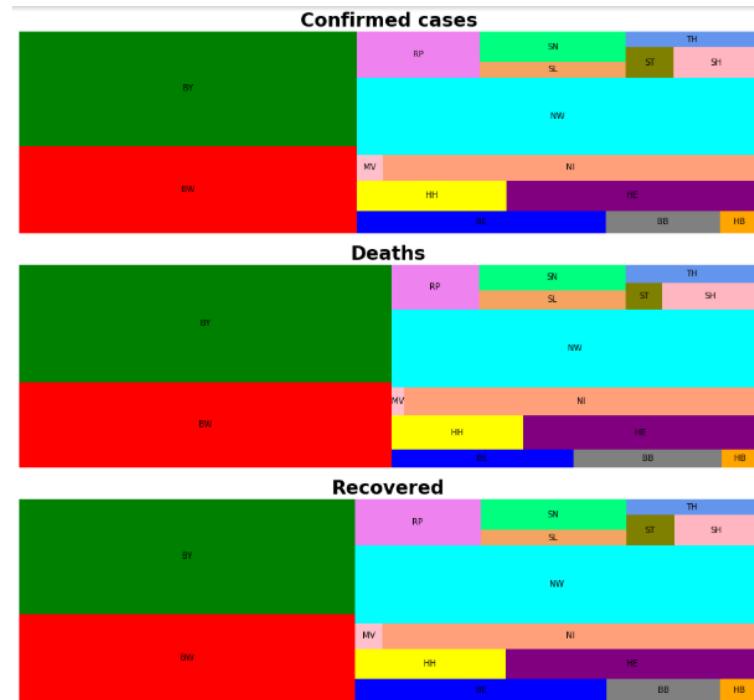


Figure 29.6: Tree map showing the distribution of confirmed death and recovered cases for all the federal states of Germany. Different color codes represents the different federal states of Germany

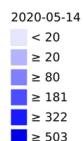


Figure 29.7: Number of confirmed cases per million residents for the federal states of Germany by the 14th of May based on the reported data in [45].



Figure 29.8: Number of cases that died to corona per million residents for the federal states of Germany by the 14th of May based on the reported data in [45].



Figure 29.9: Number of recovered cases per million residents for the federal states of Germany by the 14th of May based on the reported data in [45].

29.2 Time-Series and Prophet Prediction

Following the general outline of the curves, substantial differences in overall case numbers between different states can be seen e.g. Bavaria compared to Hessen (Figure 29.10). This might be due to the initial outbreak originating from the severely affected southern states of Europe. While the first introduction of general lockdown measures does not have an immediate effect it has to be noted that each introduction of measures should only be apparent after two weeks, as this is the time of infection till recovery of an individual is affected by corona. Even after two weeks a full exponential increases of cases can be seen. On the contrary the introduction of mask does seem to have had an apparent effect. A noticeable downslope can be seen for each state. Prophet predicts an increase in cases for each state (Figure 29.11) Especially for Berlin and Rheinland-Pfalz a noticeable increase in cases can be seen. Thus it might be advisable that the early easement of anti corona measures are leading to a recurrence of cases.

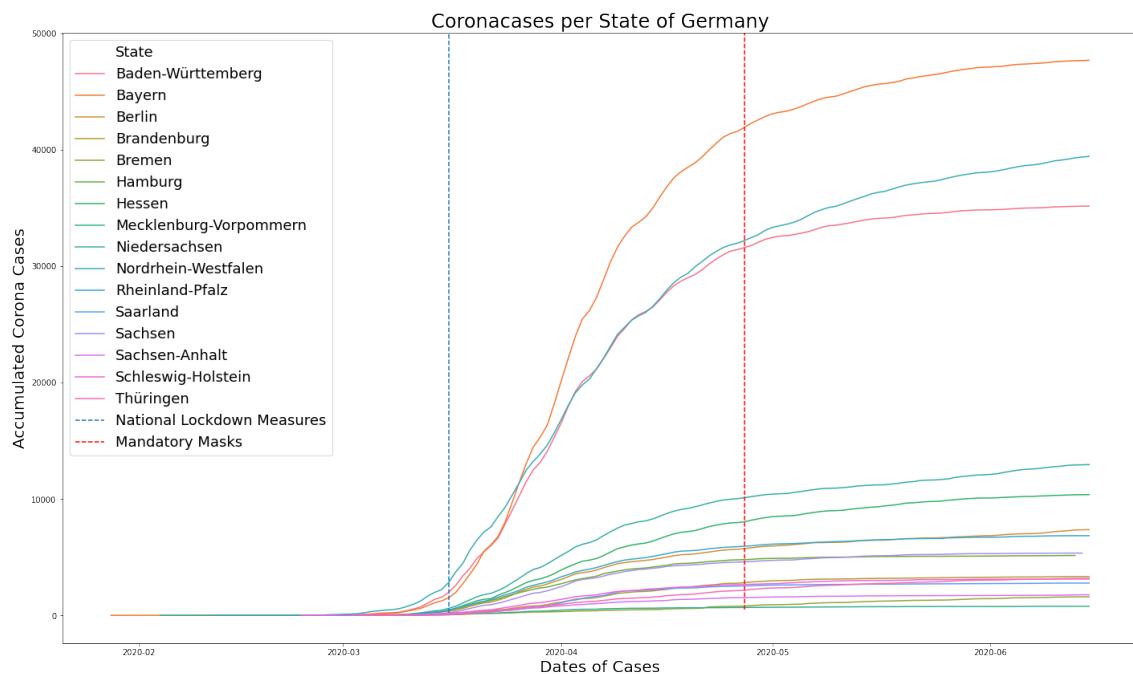


Figure 29.10: Cases trajectory plotting for each state of Germany. Two statewide coronavirus counter measures were visualized as dashed lines. The cases study follows from the 28th of January till the 15th of June. Facebook's Prophet library is used to perform the time-series prediction. All curves follow a noticeable exponential increase with a subsequent decrease in cases.

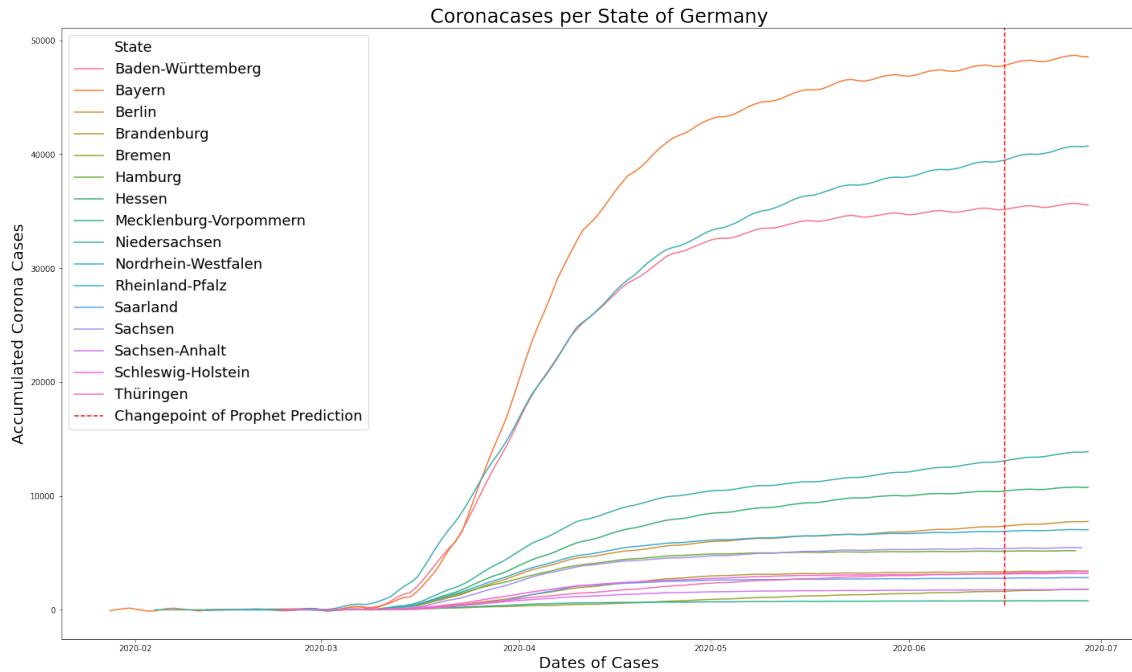


Figure 29.11: Cases trajectory plotting for each state of Germany. Predictions of the prophet library are plotted after the dashed lane, beginning with the 16th of June. Facebook's Prophet library is used to perform the time-series prediction. For each state a noticeable increase in cases is predicted.

29.3 Clustering

In order to check the data for clusters, hierarchical clustering was done first. The Figure shows the dendrogram of the hierarchical clustering of confirmed cases (see 29.12, 29.13 and 29.14).

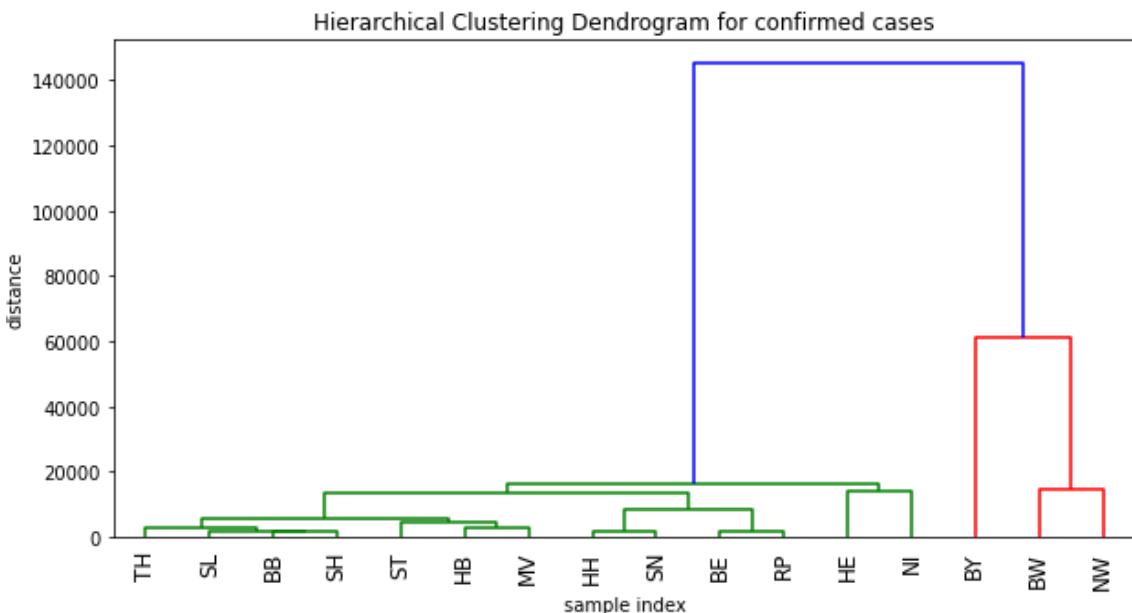


Figure 29.12: Hierarchical Clustering Dendrogram for the confirmed number of infected cases for the period March-June 2020.

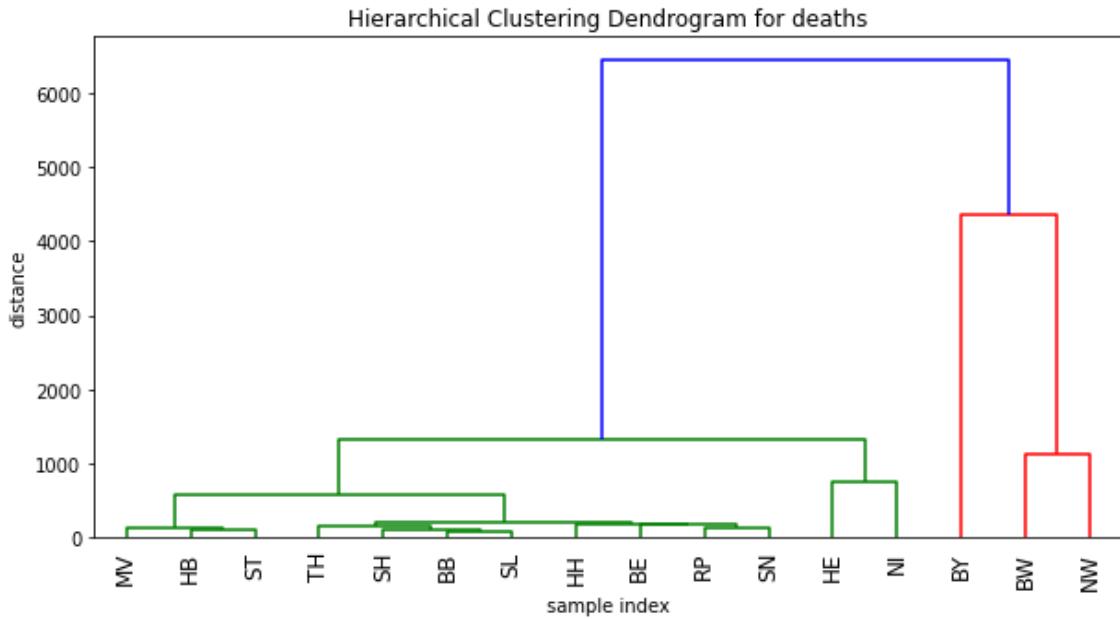


Figure 29.13: Hierarchical Clustering Dendrogram for the confirmed number patients that died to corona for the period March-June 2020.

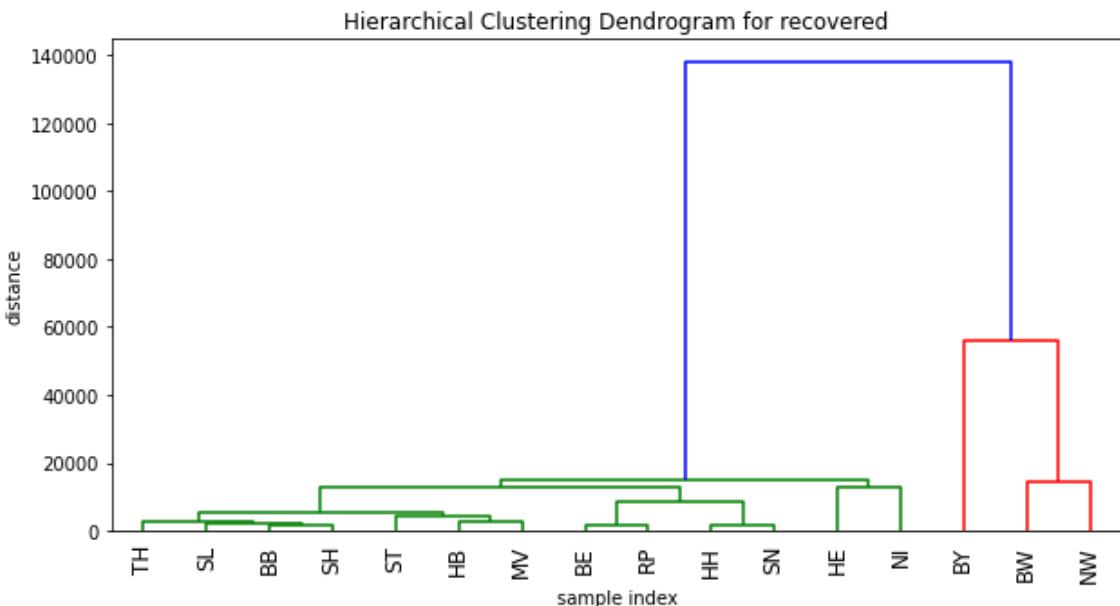


Figure 29.14: Hierarchical Clustering Dendrogram for the confirmed number patients that recovered from corona for the period March-June 2020.

The hierarchical clustering for confirmed cases, deaths and recoveries shows that there are two main clusters. One contains the states of Bayern(BY), Nordrhein-Westfalen(NW) and Baden-Württemberg(BW) and the other the remaining states, which contains four clusters (C_1 : Thüringen(TH), Saarland(SL), Brandenburg(BB), Schleswig-Holstein(SH); C_2 : Sachsen-Anhalt(ST), Bremen(HB), Mecklenburg-Vorpommern(MV); C_3 : Berlin(BE), Rheinland-Pfalz(RP); C_4 : Hessen(HE) and Niedersachsen(NI). This result corresponds to reality, because the largest number of infected, dead and recovered persons was found in the federal states BY, NW and BW, while other

federal states had lower numbers.

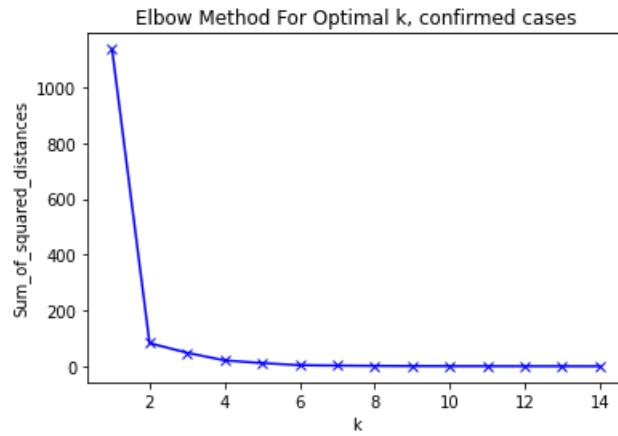


Figure 29.15: Elbow Plot to determine the optimal number of clusters for the k-means approach for confirmed cases for the period March-June 2020

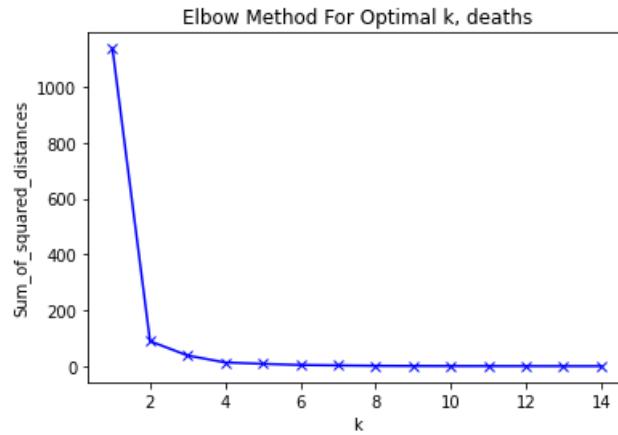


Figure 29.16: Elbow Plot to determine the optimal number of clusters for the k-means approach for the confirmed number of deaths for the period March-June 2020.

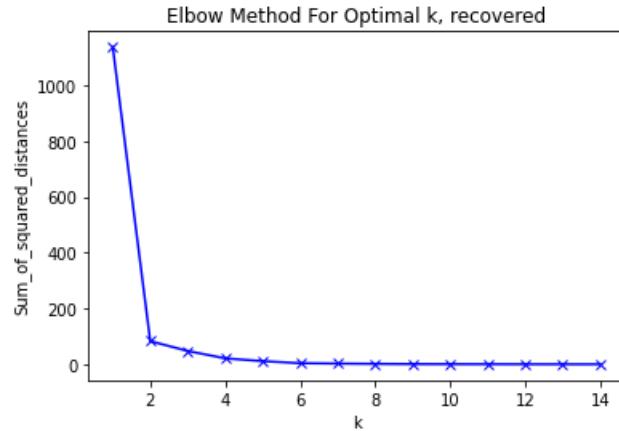


Figure 29.17: Elbow Plot to determine the optimal number of clusters for the k-means approach for the confirmed number of recovered for the period March-June 2020.

The figures above (see 29.15, 29.16 and 29.17) show the results of the elbow method for selecting the most optimal k-value, number of clusters. The best value seems to be 3 clusters.

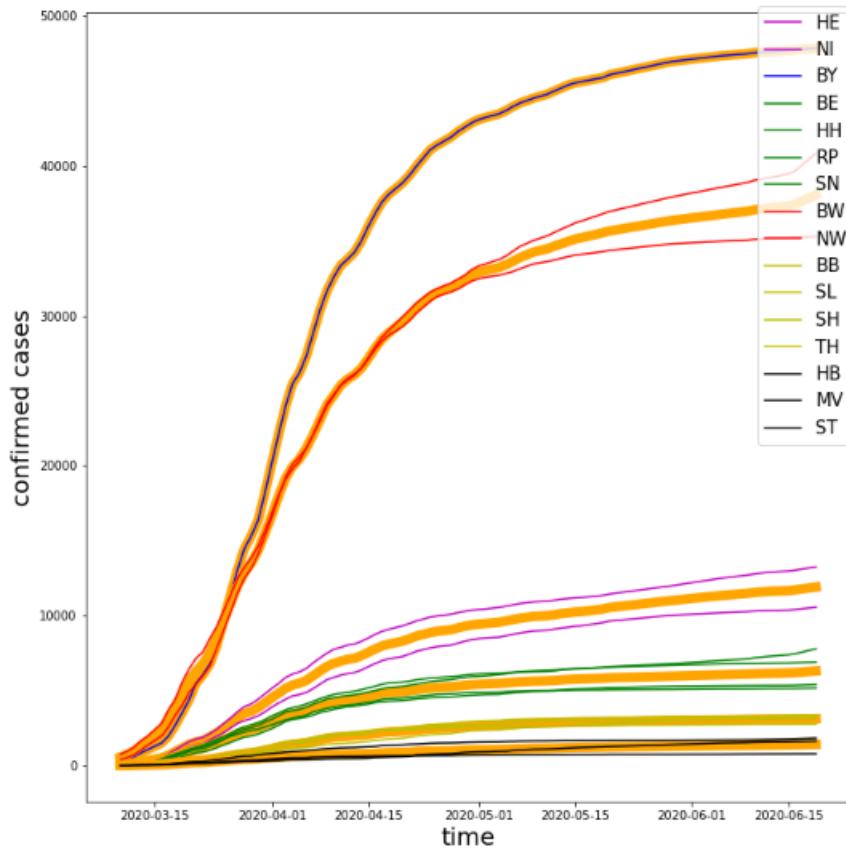


Figure 29.18: Number of confirmed and predicted cases for the period of March to June 2020 using Prophet. The Federal states are clustered using k-means-Clustering.

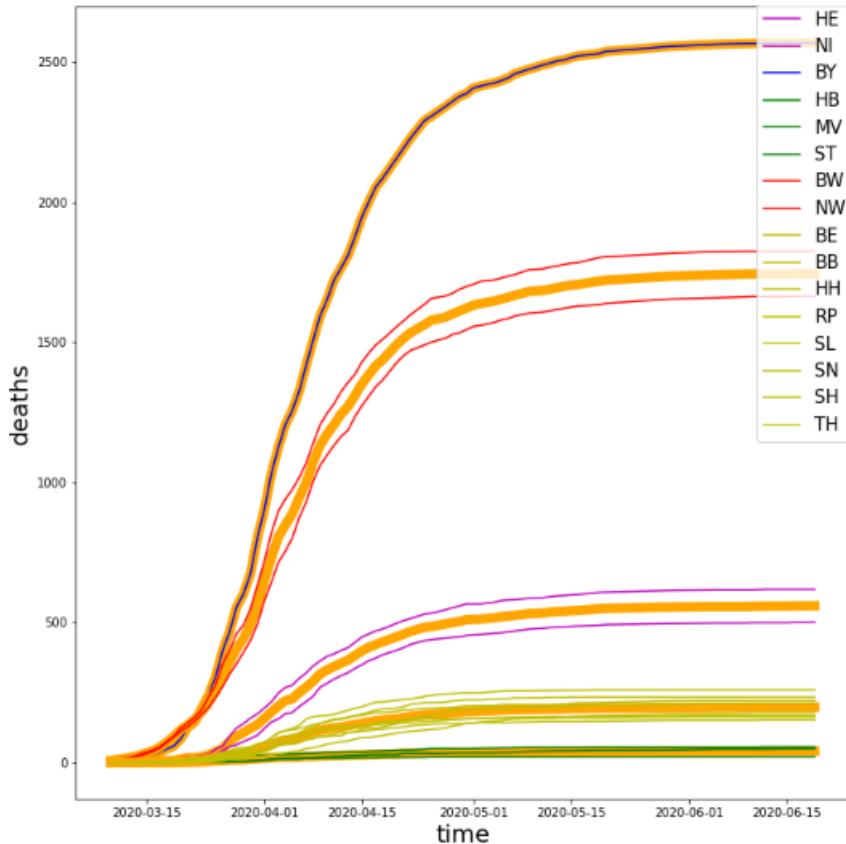


Figure 29.19: Number of confirmed and predicted deaths for the period of March to June 2020 using Prophet. The Federal states are clustered using k-means-Clustering.

In order to see also subclusters, the k value was set to 6 for confirmed cases, 5 for the deaths and 6 for recovered cases (see 29.18, 29.19 and 29.20). For confirmed cases and for recovered the same pattern of clusters could be identified (C_1 : HE, NI, C_2 : BY, C_3 : BE, HH, RP and SN; C_4 : BW and NW; C_5 : BB, SL, SH and TH; C_6 : HB, MV and ST). For the confirmed number of deaths 5 was determined to be the optimal number of clusters (see 29.19). Our results demonstrate that the k-means approach can be used to reasonable cluster the federal states of Germany according to the number of confirmed-, death-, and recovered cases within them.

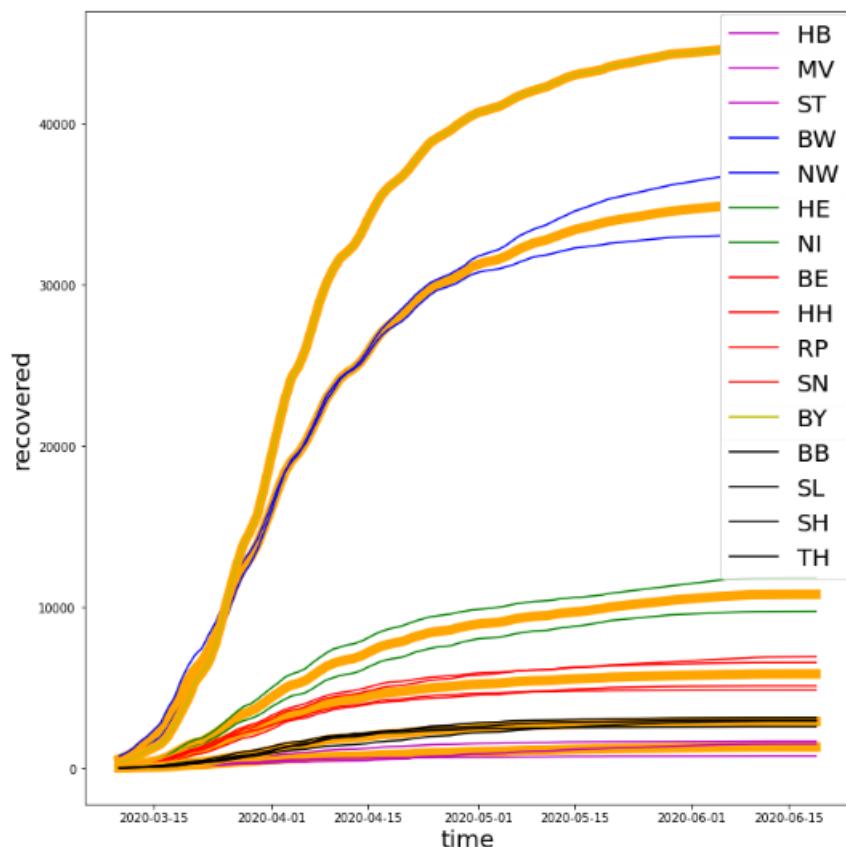


Figure 29.20: Number of confirmed and predicted recovered individuals for the period of March to June 2020 using Prophet. The Federal states are clustered using k-means-Clustering.



30. Evaluation

30.1 Project Rating

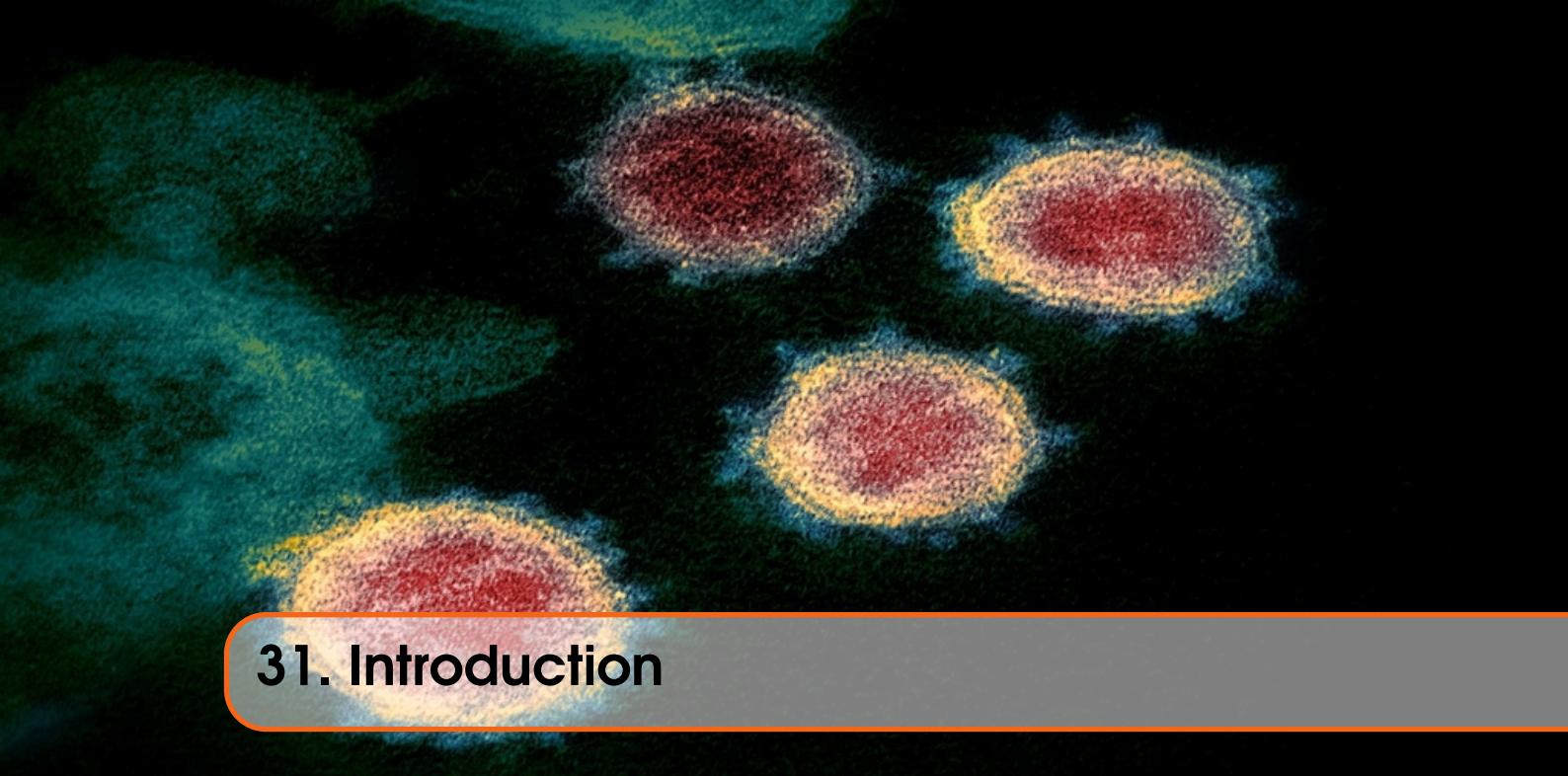
Appropriate visualization techniques for a comprehensive overview of the analysed data is an important part of being a data scientist. Finding the right graphics to emphasize certain parts of the analysis while also giving the reader a clear interpretation is not easy. Thus it is very interesting to test various methods for visualization on the recent COVID-19 cases. It is apparent that it is possible to shift the point of interest for the same data simply by using a different graphic. In this regard this project was quite interesting. On the other hand doing time series prediction again was redundant. Visualization techniques can also be explored outside of the scope of time series. Overall we rate this project favourably.

30.2 Problems

Firstly in the visualization the data, we did not face any problem as the data is directly extracted from the RKI site, and the obtaining the results with the plots was easy. Within the parts of time series prediction and curve plotting no difficulties were encountered. The preprocessing and finding of the appropriate data was easy. Furthermore, most data exploration approaches are well documented in python, such that redoing some of them for our goal of illustrating the corona cases for the federal states was not problematic either.

Project 7: Origin Analysis of COVID-19

31	Introduction	133
31.1	Background	
31.2	Goal of the project	
31.3	Outcomes	
32	Methods for Phylogenetic Analysis ...	135
32.1	Data	
32.2	Methods	
33	Analysing the Spread of SARS-CoV-2	141
33.1	Results	
33.2	Discussion	
33.3	Conclusion	



31. Introduction

31.1 Background

Phylogenetic analysis aims to reconstruct phylogenies both for a group of species and also for the individuals within those species. For pathogens, we are interested in the evolution of the diseases at the genetic level and want to derive information about their past and present diversity [56]. An important fact about the *Coronaviridae* family is that its members tend to “jump” from one species to another. When the transmission occurs from a non-human host to a human host it is called zoonosis [18]. The determination of the most recent common ancestor of the human SARS-CoV-2 and the zoonotic transmission can provide important information about biological features, key mutations and properties of the virus. A detailed understanding of how an animal virus jumped species boundaries to infect humans will help in the prevention of future zoonotic events. [2]. Analyzing the diversions of COVID-19 sequences sampled from different human hosts all over the globe can lead to information about the order of transmission chains that have taken place.

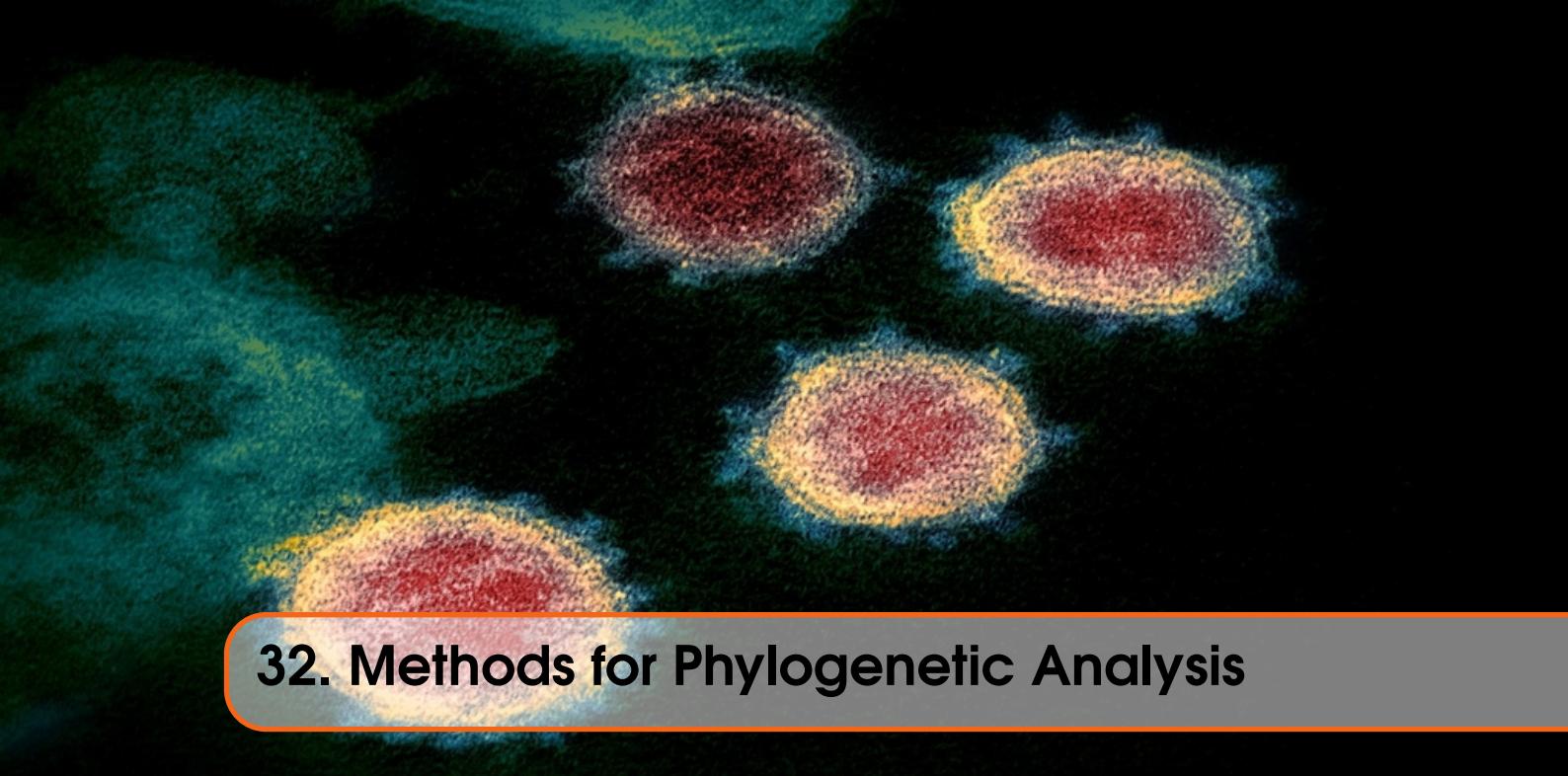
31.2 Goal of the project

We will perform two types of analysis: Origin detection and phylogeographic analysis. At first, the genetic sequence of SARS-CoV-2 is compared with six other virus sequences of the *Coronaviridae* family originated from different non-human hosts to gain information about the origin of the virus and the zoonosis. Regardless of that a second data set was created containing 14 human SARS-CoV-2 sequences from different countries to identify possible pathways of the virus spread. The phylogenetic computations are performed by the hierarchical UPGMA and the TreeTime algorithm and the results are compared and analyzed afterwards. On top of that two non-hierarchical clustering methods (k-means and k-medoids) are performed to compare its results with the outcomes of the hierarchical clustering.

31.3 Outcomes

The Rhinolophus (horseshoe bat) was identified to have the most similar genomic sequence to the human SARS-CoV-2 genome with a sequence similarity of 94% among the six *Coronaviridae*

family sequences of different hosts. For the phylogeographic analysis the UPGMA algorithm produced a very precise phylogenetic tree where clear spatial patterns can be recognized. Thus we were able to recognize possible infection pathways of the American samples. Furthermore, TreeTime offered the possibility to estimate the temporal transmission pattern of the virus spread with a higher precision in addition to the spatial patterns. The application of non-hierarchical methods on the 14 human SARS-CoV-2 sequences resulted in the construction of five different clusters. The biggest cluster contained mostly samples from China and Europe. Nevertheless, the non-hierarchical approaches performed much worse than the hierarchical approaches.



32. Methods for Phylogenetic Analysis

32.1 Data

The first data set was created to perform origin analysis (OA). The data set consists of one of the earliest sampled genetic sequence of SARS-CoV-2 from Wuhan as well as six other viruses of the *Coronaviridae* family in different hosts. The analysis is based on the Github repository from Simon Burgermeister [11] who originally downloaded the sequences from the NCBI Virus public library [22]. The Accession numbers as well as host information are given in table 32.1.

The second data set was constructed to provide information about phylogeographics (PG). Within this data set 14 sequences collected from humans with distinct geographic locations were analyzed. Each continent is represented at least once but most times with multiple sequences. To monitor the divergence of the virus at a specific time point, sequences were chosen that were all collected within March 2020 (expect one German sequence from February and another early Wuhan sequence from January). The respective metadata is listed in table 32.2.

Accession number	Host	Description
MN996528	H. Sapiens	Human SARS-CoV-2
NC_019843	H. Sapiens	Human MERS-CoV
JQ065048	Anatidae	Ducks, geese and swans
MG772934	Rhinolophus	Horseshoe bats
NC_034972	Apodemus chevrieri	Rodent
KX38909	Gallus gallus	Chicken
MT084071	Manis javanica	Pangolin

Table 32.1: Origin analysis (OA) data set containing the human SARS-CoV-2 sequences and six other viruses of the *Coronaviridae* family collected in different hosts.

Accession number	Host	Location	Collection date
MT466071	Homo Sapiens	Uruguay	2020-03-13
MT499220	Homo sapiens	Tunisia	2020-03-31
MT447176	Homo sapiens	Thailand	2020-03-20
MT531537	Homo sapiens	Italy	2020-03-01
MT470177	Homo sapiens	France	2020-03-15
MT358639	Homo sapiens	Germany	2020-02-20
MT259229	Homo sapiens	China: Hubei, Wuhan	2020-01-26
MT350282	Homo sapiens	Brazil	2020-03-18
MT407659	Homo sapiens	China: Zhejiang	2020-03-24
MT451640	Homo sapiens	Australia	2020-03-25
MT434809	Homo sapiens	USA: New York	2020-03-19
MT434808	Homo sapiens	USA: New York	2020-03-19
MT633004	Homo sapiens	USA: Washington	2020-03-23
MT632947	Homo sapiens	USA: Washington	2020-03-23

Table 32.2: Phylogeographics (PG) data set containing 14 SARS-CoV-2 sequences from different countries mostly collected in March 2020.

32.2 Methods

32.2.1 Hierarchical Approaches

Two common approaches towards constructing phylogenetic trees are the unweighted pair-group method with arithmetic mean (UPGMA) and TreeTime algorithm.

UPGMA

The UPGMA algorithm is defined by its simplicity. Based on the assumption that the rate of mutation between different lineages stays constant over time, the so called *molecular clock hypothesis*, a phylogenetic tree with equidistant leaves to the root can be constructed. UPGMA starts with a matrix of pair wise distance objects. By iterating over the matrix, a cluster with entries i and j is defined by the smallest distance between both within the matrix. These entries are then connected through a branch called the most recent common ancestor node. The distance of both to the connection node is defined as $D(i, j)/2$. Next a new cluster u is defined and its distance to each other cluster calculated as the average of the distances to all other clusters, hence the name. If no more entries are left the algorithm terminates.

TreeTime

The TreeTime approach belongs to the class of expectation maximization (EM) algorithms. EM algorithms use the divide and conquer method to divide a problem into simpler subproblems, thus decreasing computational time and increasing efficiency. The core idea of the TreeTime algorithm is a joint maximum likelihood assignment for each branch length. For each parameter e.g. leaf or node in the tree topology the assignment is calculated by finding the most likely value after summing or integrating over all unknown previous states. In practice the algorithm uses two steps: a post-order traversal and pre-order traversal. For the post-order traversal the maximum likelihood for node n to be at position t is calculated by taking the constraints of its children C and external constrains E e.g. collection data into account. Next the pre-order traversal follows where the branch length of each internal node is computed by finding the optimal value of time point t under the constraint of the parental node position.

Comparison

The assumption of a constant mutation rate for the UPGMA algorithm is also its biggest disadvantage. Because each distance from the root to the leave within the tree is the same, UPGMA frequently generates wrong tree topologies. In reality it is very unlikely that different lineages are the same length in time and thus the hypothesis is violated. Many different factors can influence the mutation rate of an organism, bacterium or virus. In contrast the TreeTime algorithm does not rely necessarily on the *molecular clock hypothesis*, thus removing one of the biggest disadvantages. Because TreeTime uses an expensive Bayesian approach the EM strategy is employed to strike a balance in computational efficiency. The heuristic nature of TreeTime can also be a disadvantage where a convergence to a wrong tree topology might be possible. UPGMA is robust in its implementation. The same distance matrix always results in the same tree topology. This might be an advantage over TreeTime where reproducibility is needed.

Implementation of UPGMA and TreeTime

The construction of the phylogenetic tree was performed by applying the preimplemented *UPGMA* function of the *Phylo* module contained in the *Biopython* package. It was executed with default parameters. The distance matrix was calculated beforehand by the *DistanceCalculator* function of the same package with the *identity* property. The multiple sequence alignment was generated by the NCBI BLAST implementation [26].

To use Nextstrain Workflow with TimeTree for phylogeographic analysis, the data sets should be created according to the Nextstrain Fauna (database tool) requirements for sequence data and sample metadata. After preprocessing of the data set according to Nextstrain Fauna the processing with Nextstrain Augur (analysis pipeline) followed, performing multiple sequence alignment with MAFFT. A phylogeny with high probabilities is derived using TreeTime, which estimates a molecular clock. Given the derived molecular clock, TreeTime then creates a time-resolved phylogeny, estimating sequence states at internal nodes and calculating the geographic migration history across the tree. The output data is exported as a JSON file that can be visualized interactively on the web with Nextstrain Auspice (visualization platform).

32.2.2 Non-Hierarchical Approaches

While there is a broad range of clustering algorithms, the field of phylogenetic analysis is dominated by hierarchical clustering approaches (e.g. UPGMA). Nevertheless other approaches like the commonly known k-means algorithm can be used to cluster sequences by similarity. Therefore we implemented two non-hierarchical clustering methods to compare its results with the ones of the hierarchical clustering.

k-means

K-means is one of the most popular clustering algorithms due to its simplicity. The data is separated by assigning each data point to cluster centers, such that the total squared error (SE) is minimized. The initial center points are set arbitrarily and adjusted at each iteration step until the SE converges or a fixed number of iterations is reached (Figure 32.1. Finding the right k (number of clusters) can be challenging. Therefore k-means is performed for several k to find the one value that strikes a balance of minimized sum of squared error and reasonable cluster number.

k-medoids

The k-medoids method is a clustering methods that is related to the k-means algorithm. Both approaches break the data sets into groups. But while k-means tries in each step to minimize the total squared error, k-medoids minimizes the sum of dissimilarities between points belonging to the same cluster and its cluster center. For k-medoids the cluster centers are set to be one of the data points. This approach is less sensitive to outliers but at the cost of runtime, since every point is set

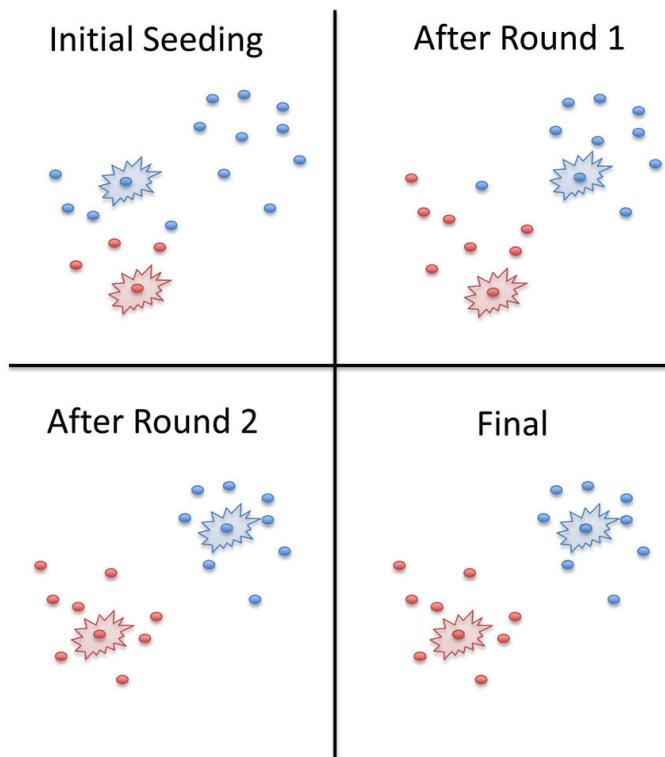


Figure 32.1: Schematic representation of the k-means algorithm for $k = 2$. The initial center points are set arbitrary and adjusted at each iteration step until the SE converges or after a fixed number of iterations. The picture is taken from [40].

to be a cluster center (medoid).

Workflow

Before we can apply both non-hierarchical clustering approaches we need to transform the sequences of the MSA such that they get comparable. Therefore we implemented the following workflow:

1) Factorize Bases:

Each sequence in the MSA contains elements of the DNA alphabet $\{A, C, T, G, N, -\}$, which – representing gaps, N unknown bases, and the remaining four characters the DNA bases. All sequences are factorized, such that to each character of the alphabet a unique number is assigned.

2) Principal Component Analysis:

The factorized bases can be seen as features. Consequently each of the 14 sequences consist of roughly 30.000 features/dimensions. Using PCA we can project our data onto two orthogonal axes, meaning that the feature space is reduced from roughly 30k to 2.

3) Clustering:

After performing PCA each data point is expressed by two principal components. The euclidean distance between the projected sequences is used to cluster them using k-means.

For k-medoids we used the euclidean, manhattan and cosine distance.

For verification of the PCA approach, the k-Medoids method was also performed on a hamming distance matrix computed between all human sequences to compare the resulting clusters.

32.2.3 Phylodynamics

A field that gives a deeper insight into how various infected diseases are evolving over time is Phylodynamics. Phylodynamics completely depends on phylogenetic inference which acts as a tool to analyse mutations patterns which are the cause for probable spillover events. The mutations alter the phenotype which in turn infect different cell types. This allows the virus to develop different possible transmission routes while being a driving factor of new epidemiological processes. Phylodynamics play an important role in filtering and acquiring the information from genetic data. The interaction of the rapidly evolving pathogens completely depends on their ecological and evolutionary dynamics which usually happens at the same time scale. As the time of sampling plays a crucial role in calibration of phylogenies it is an important information to incorporate into phylodynamics.

32.2.4 Phylogeography

Because infectious disease transmission is an inherently spatial process the geographic location of samples must be taken into account. The description of how the genetic signals are structured geographically within and among the species is called Phylogeography. Being the fastest growing field, it stands as a new technique in reconstructing the gene and genealogies [43]. Visualizing these spatial relationships over geographic locations allows us to deduce how sub species are evolving, possible transmission routes and the origin place. It completely relies on ancestral lineages thus connecting the movement through space and movement through time. A few variety of applications that are entirely dependent on Phylogeography are earth historic events, distribution models and speciation processes. By linking the patterns of divergence in the population it identifies and tests the status of the diversification in an area. It gives us insight into whether and over which spatial and temporal scales the historical and the recurrent processes have shaped.

33. Analysing the Spread of SARS-CoV-2

33.1 Results

33.1.1 Origin Analysis

The resulting phylogenetic tree (Figure 33.1) of the OA data set shows that our human SARS-CoV-2 sequence is most similar to the SARS-like coronavirus sequence of the Rhinolophus (horseshoe bat) with a similarity of 94%. The host with the next similar sequence is the *Manis javanica* (Pangolin) with a similarity of 0.89% between their genomes. The human MERS-Cov genome and the SARS-CoV-2 genome share only a sequence similarity of 0.74%.

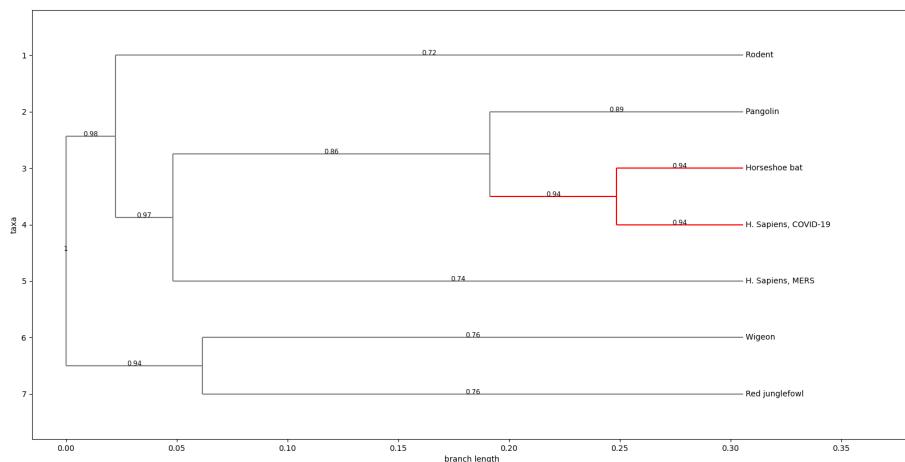


Figure 33.1: Phylogenetic tree of the origin detection analysis. The branch weights represent the sequence similarity in percent. The cluster containing the human SARS-CoV-2 sequence and the most similar sequence (horseshoe bat) is marked in red.

33.1.2 Phylodynamics and Phylogeographics

UPGMA

The UPGMA algorithm produced very plausible results (Figure 33.2). All three European samples (green) were assigned to the same cluster with a maximum sequence dissimilarity of just 0.007%. The parent branch of the European cluster connects it with the South American cluster (pink) with the Brazilian and Uruguayan genomes showing a dissimilarity of only 0.003% to each other and a dissimilarity of 0.04 to the European cluster. Observing the results of the North American samples (blue) shows a compelling outcome. The Washington samples were assigned together (dissimilarity 0.003%) as were the New York samples (dissimilarity 0.007%). However, the West Coast samples were connected to the European and South American cluster, while the East Coast samples were identified to be most similar to the Wuhan sequence with a divergence of only 0.041%. The Thai sequence was clustered together with the Zhejiang (China) sequence (red) and the Tunisian (brown) and Australian (purple) sequence forming their own single sample cluster respectively.

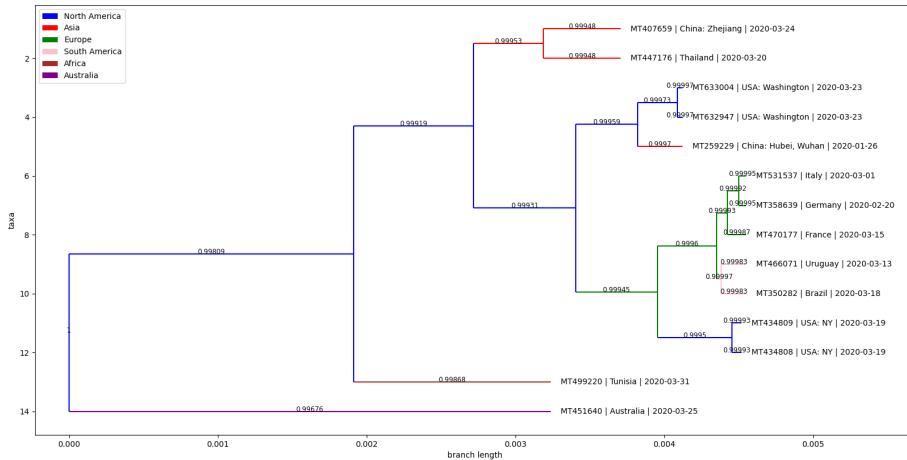


Figure 33.2: Phylogenetic tree of the phylodynamics and phylogeographics analysis. The branch weights represent the sequence similarity. The continents are marked by color as indicated in the legend while exact country and collected data is displayed in the sequence names.

33.1.3 TreeTime

A maximum likelihood phylogeography analysis of 14 SARS-CoV-2 genomes was conducted. These sequences are sampled from distinct locations around the world. Consistent with other studies it could be found out that SARS-CoV-2 moved from China, Wuhan Hubei to the other countries (see Figure 33.3). SARS-CoV-2 moved to European countries (France at first) via China Zhejiang and South America (Brazil). In parallel the North American (USA) epidemic resulted from at least two introductions, one from Europe via Tunisia and another one from South America (Uruguay). It can be estimated that SARS-CoV-2 moved from Europe (France) to Australia and from there to Germany. The introduction to Europe (Italy, Germany and France) most likely took place in late February 2020. In South America it originated once in February (Brazil, Uruguay via Thailand) and another time in North-America via Italy and Tunisia. The results are confirmed by conducted studies. By using phylogeographic analysis it was possible to estimate 5 clades based on changes in nucleotides of the virus genome, which are associated with certain countries where the virus was introduced (see Figure 33.4. The clade c1 (C3023T, C14394T, A23389G), c2(C227T), c4(G28867A, G28868A, G28869C) and c5 (C1045T, G25549T) are indicated to be distributed in European countries (Germany, France, Italy), Tunisia, Australia and North America (USA). The Clade 3 is found to be introduced in South America (Uruguay and Brazil), North America (USA) and Thailand.

Genomic epidemiology of novel coronavirus

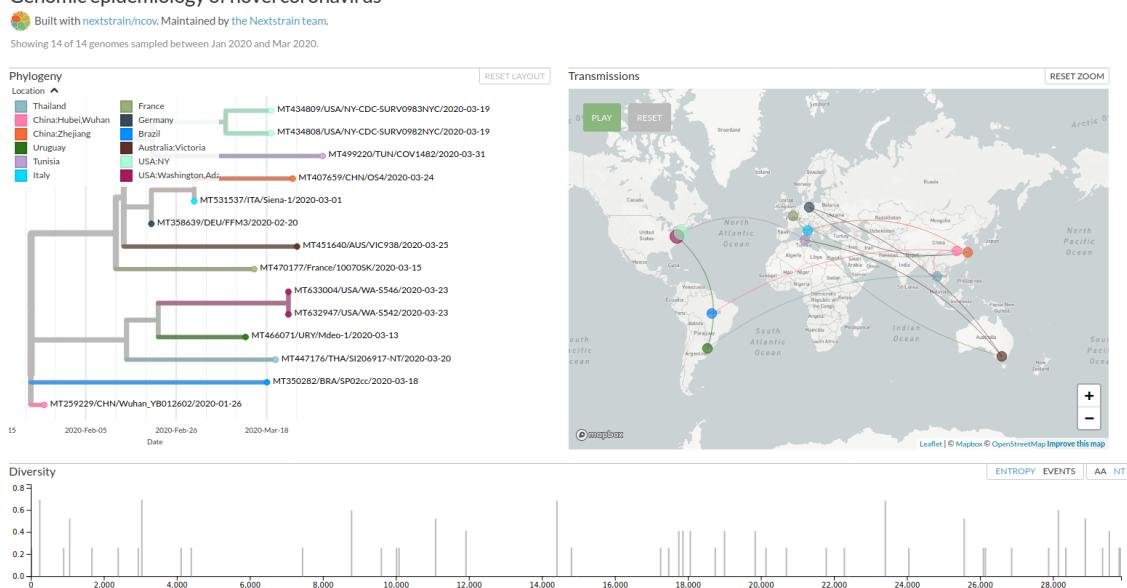


Figure 33.3: Phylogeographic analysis of 14 SARS-CoV-2 genomes. Branch colors indicate the known country of sampling and the geographic migration history.

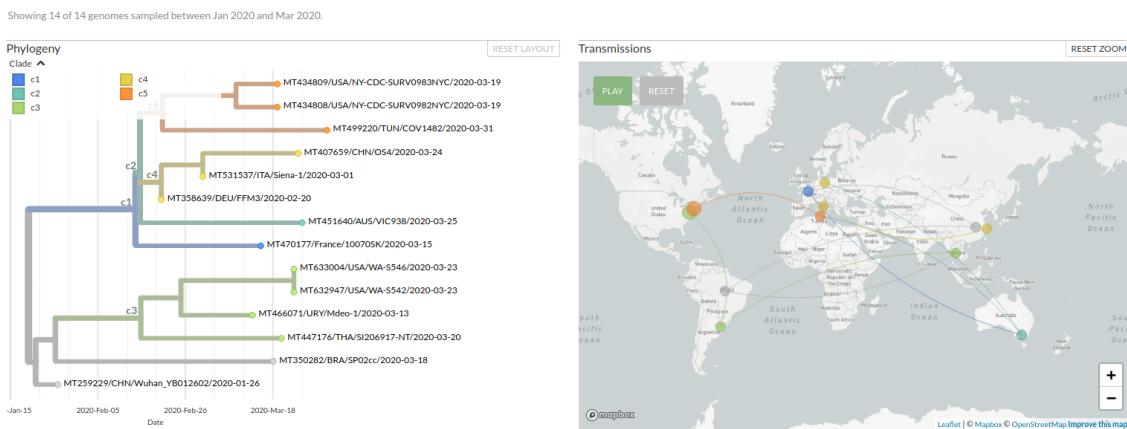


Figure 33.4: Phylogeographic analysis of 14 SARS-CoV-2 genomes. Branch colors indicate the clades. The clade c1 nucleotide changes C3023T, C14394T, A23389G, c2 C227T, c4 G28867A, G28868A, G28869C and c5 C1045T, G25549T respectively, where first the origin nucleotide, then the genomic site and at last the nucleotide to which the base mutated is given.

Non-Hierarchical Approaches

After performing PCA on human sequences the resulting principal components explain 61% of the original variance. Because the explained variance is minimal it cannot be generalized. But accounting for this would be out of the scope of this weeks project. Nevertheless in the book "Multivariate Data Analysis" Hair et al. [8] state that in many fields an explained variance of roughly 60% is sufficient. In Figure 33.5 the 2D projection of all 14 human sequences can be seen. On the left side of Figure 33.5 the biggest data point heap can be spotted that includes mostly

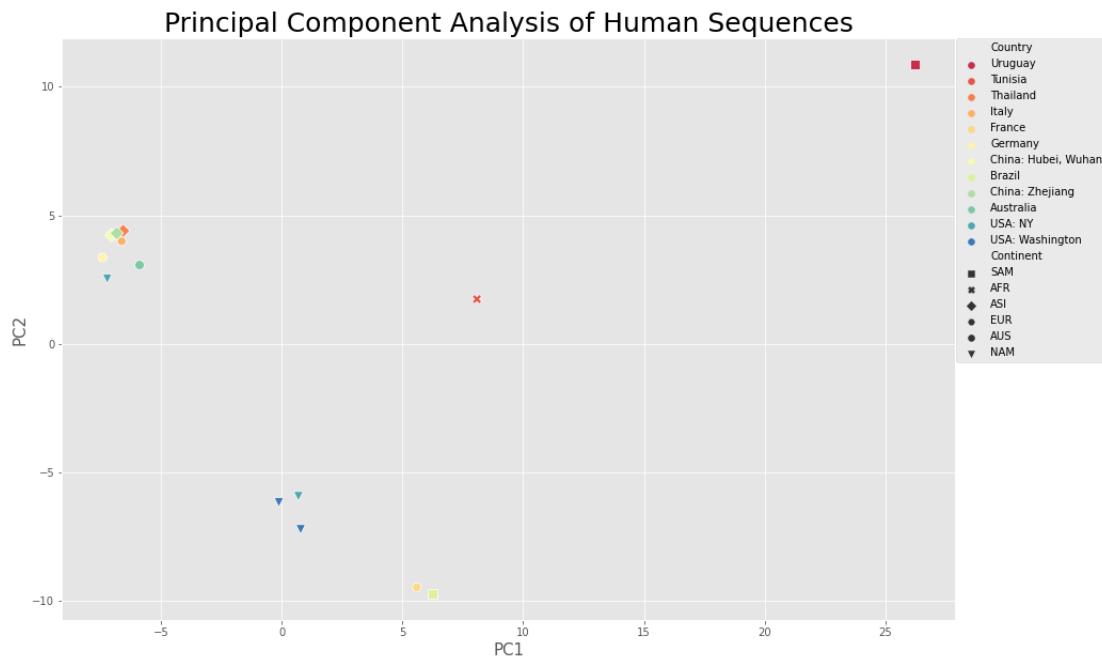


Figure 33.5: Principal Component Analysis graph of the human sequences. Each base of the roughly 30k sequences is seen as a feature and then projected on two dimensions with PCA. The shape of the markers represent the continent the samples were taken from and the color the corresponding country. The procedure is explained in Section 32.2.2.

sequences from China (Diamonds) and Europe (Circles), but also sample from the United States (Triangles). Another smaller group of three samples from the US can be spotted at the bottom left. Surprisingly another heap was formed that includes the sample from French and the one from Brazil. The African and the Uruguayan sample did not group together and are generally rather separated from the other data points.

The principal components are then used to cluster the data as described in Section 32.2.2. In Figure 33.6 we can see the resulting clusters using k-means and k-medoids with different distance metrics. Setting $k = 5$ led to minimal dissimilarities.

Comparing multiple K-Medoids metrics to K-Means and each other

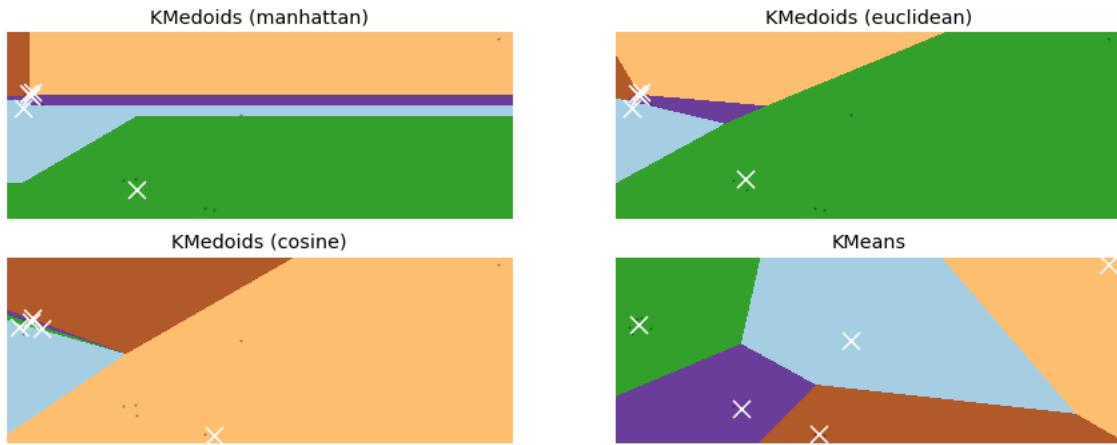


Figure 33.6: Comparison of the k-means and k-medoids clustering approaches. Each white cross represents a cluster center. The number of clusters is set to $k = 5$. The procedure is explained in Section 32.2.2.

It can be seen that the four resulting clustering approaches cluster the data quite differently. While the three k-medoids approaches split the biggest data heap containing most samples from Europe and Asia into several clusters, the k-means approach assigns the same samples to a single cluster. Applying k-medoids on the hamming distance matrix with $k = 5$ produced the same clusters as it was the case with the PCA approach.

33.2 Discussion

33.2.1 Origin Analysis

As many early cases of COVID-19 were linked to the Huanan market in Wuhan [61], it is possible that an animal source was present at this location. Given the similarity of SARS-CoV-2 to bat SARS-CoV-like coronaviruses, it is likely that bats serve as reservoir hosts for its progenitor. Although the similarity of 94% to the coronavirus sequence hosted by the *Rhinolophus*, Andersen et al. [2] identified that its spike protein diverges in the receptor binding domain (RBD), which suggests that it may not bind efficiently to the human ACE2 receptor. Furthermore, it is assumed in this and other studies [29, 33] that an intermediate host was probably involved. Nevertheless, it has to be said that more different species of the same family can still have different mutation rates and therefore the molecular clock hypothesis has most likely been violated, which was not taken into account by UPGMA.

33.2.2 Phylodynamics and Phylogeographics

UPGMA

It was previously described that the major disadvantage of UPGMA is its violation of the molecular clock hypothesis. However, in the PG data set only human samples were examined, which could only have developed apart after the zoonosis had taken place in Wuhan at the end of 2019. It can

therefore be assumed that the molecular clock hypothesis is not, or at least only slightly, violated here.

It is not surprising that the New York samples clusters together with the European samples since new research suggests the COVID-19 outbreak in New York was mainly caused by travelers from Europe, not from Asia [20]. In contrast, the outbreak on the west coast of America is assumed to be mainly triggered by Chinese travelers, which is represented by cluster of the two samples from Washington combined with the Wuhan sample. overall, the UPGMA algorithm produce a very precise phylogenetic tree in which clear spatial patterns can be recognized.

TreeTime

Using the TreeTime based approach it was possible to describe general transmission patterns and estimate the emerging of SARS-CoV-2 to distinct countries including the most probable introduction date. It was possible to find evidence that the outbreak began in China, Wuhan Hubei and was widespread with movements to countries in Europe (beginning with France, Italy, Germany), Australia, North Africa (Tunisia), North America (USA) and South America (Brazil, Uruguay). The phylogeographical analysis provides a more precise view of how the virus was spread in accordance with travel, import, export and transport patterns from country to country. In addition, the phylogeographical approach provides virus genome changing information.

Non-Hierarchical Approaches

According to the results of the PCA, one can infer that the virus first spread within China starting in Wuhan and was then carried to Europe, since those samples are clustered together closely. The most distant point of this cluster is a sample from New York, that indicates the transmission from Europe/China to the US. This makes sense as media reports about the coronavirus outbreak first started in China and then in Europe. As stated in the paper from Worobey et al. [60] the virus spread to the US in the middle of February when the pandemic already begun within Europe and China. The small cluster containing the French and Brazil samples does not make much sense for us. This cluster might be caused by the small sample size or the limits of non-hierarchical clustering methods in phylogenetic analysis. The clustering using the k-means approach clustered the samples according to our expectations, while k-medoids split the heap of Europe/Asia/US samples into several cluster, no matter which distance metric was used. This was surprising since papers using non-hierarchical approaches stated that k-medoids would be the better approach in phylogenetics than k-means, because k-medoids is less sensitive to outliers.

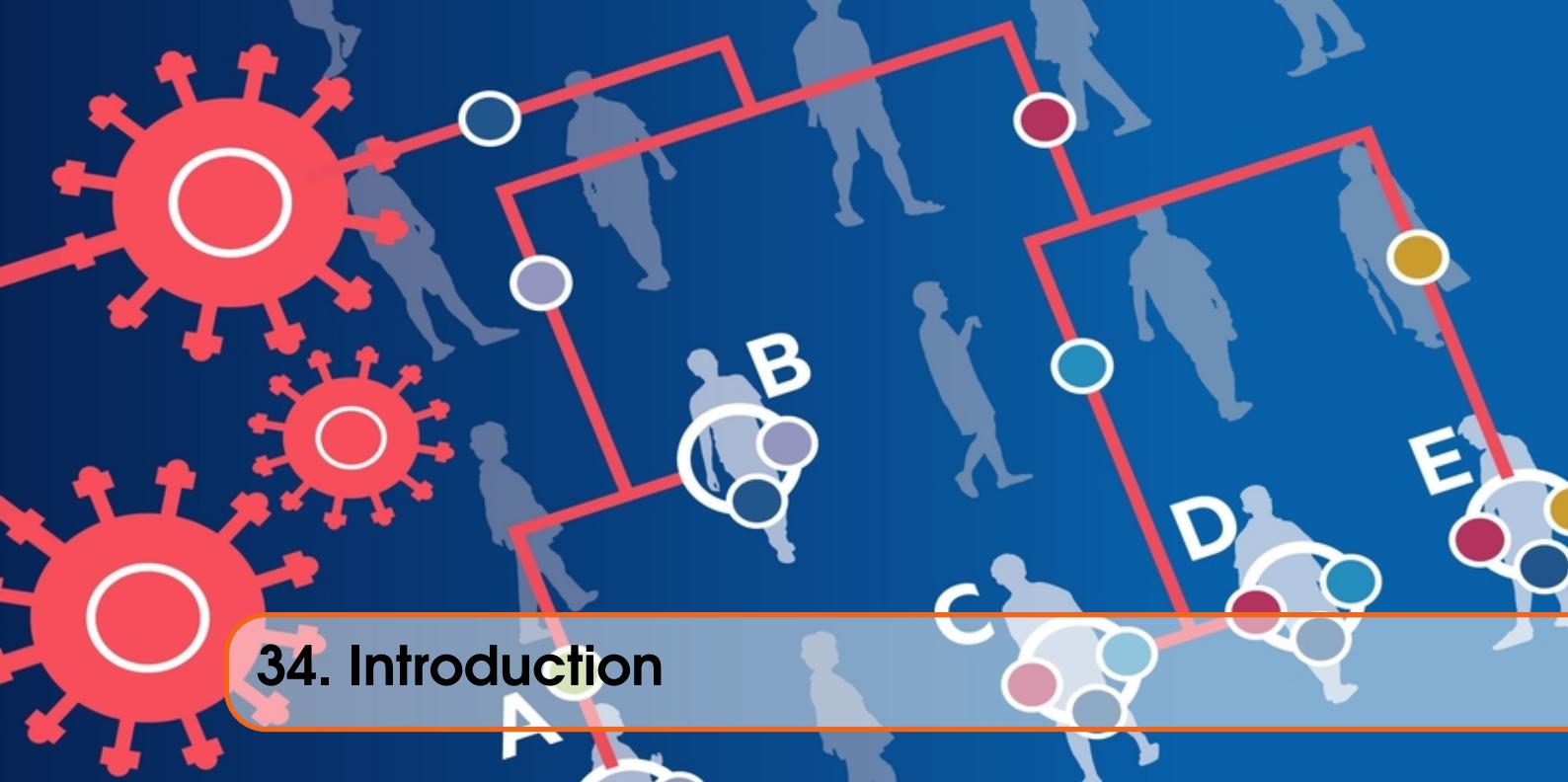
33.3 Conclusion

Even though the field of phylogenetic analysis is dominated by hierarchical clustering approaches the presented non-hierarchical methods can also provide interesting insights into the pathway of the virus. Regarding this project the hierarchical approaches produced a much more reasonable separation of the data, that allows to infer possible pathways that follow the trends described in other studies. However, non-hierarchical methods can be a better fit in the process of exploring bigger data sets for underlying structure that can afterwards be analyzed in detail since no multiple sequence alignment need to be calculated.



Project 8: Sequence-curve-based Phylogeny Analysis

34	Introduction	151
34.1	Background	
34.2	Project Description	
34.3	Outcomes	
35	Solution Approach	153
35.1	Data	
35.2	Methods	
36	Results and Discussion	157
37	Evaluation	163
37.1	Project Rating and Problems	



34. Introduction

34.1 Background

In last week's project, we could already see how classical phylogenetic models can be used to make assumptions about the spread of the virus based on their sequence identity. Recall, those models require a multiple sequence alignment, which is not only computationally expensive but also varies due to different alignment costs. Another disadvantage lies in the evolutionary assumptions that those models make. For example, UPGMA (see Section 32.2) assumes a constant rate of evolution for all branches of the tree, which is unlikely, especially for sequences that are separated by larger evolutionary distances. Those problems could be avoided by a graphical representation of the virus genome that is based on the DNA sequences themselves, rather than the multiple sequence alignment.

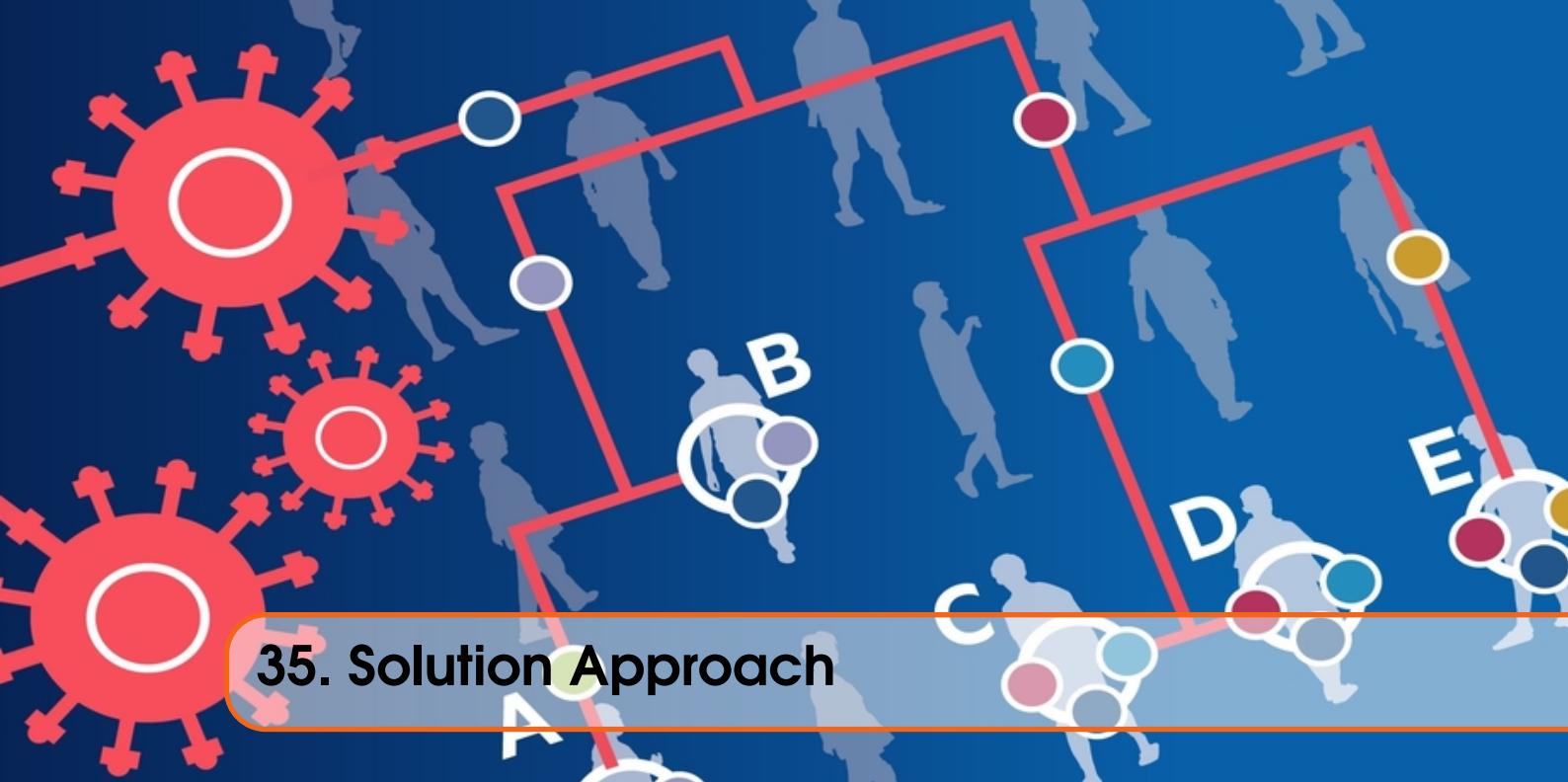
34.2 Project Description

This week's project aims to implement one such graphical method and compare its results with the one of a classical UPGMA approach. Therefore a pyrimidine–purine graph is constructed as described by Liao et al. [31]. First a classical phylogenetic tree (UPGMA) is built followed by the construction of a phylogenetic tree based on the method by Liao et al.. In the end, we compare both trees according to their visual representation and a set of metrics. This includes weighted and unweighted Robinson-Foulds distance (symmetric differences) & Euclidean distances. The trees are constructed using roughly 50 coronavirus sequences. Both models will be used to elaborate phylogenetics as well as phylogeographics on a data set containing human samples all over the globe. Additionally another data set was used containing samples from different hosts for which we perform the same analysis.

34.3 Outcomes

Upon visual inspection distinct differences between both methods were apparent. While the hypothesis of a spillover event from horseshoe bat to human is supported by the method by Liao

et al. it had difficulties performing on sequences which are very similar in alignment. All three distance metrics suggest significant deviations in tree topologies between both approaches.



35. Solution Approach

35.1 Data

For the phylogenetic analysis, we created a data set containing nine human samples from mainland China, in which we took three samples from January, February, and March, respectively. The remaining 39 human samples were all collected mid of April for different countries all over the globe. Another data set was created as described in last weeks table 7.1 to perform origin analysis (OA) with samples from different hosts.

35.2 Methods

35.2.1 2D Graphical Representation of genomic Sequences

A pyrimidine-purine graph is created to analyze the phylogenetic relationships of genomes. To construct the curve, four vectors are defined to represent the purine bases Adenine and Guanine and the pyrimidine bases Thymine and Cytosine (Figure 35.1). The vectors hereby indicate the shift from the previous to the next data point of the curve when the respective base is present at the current position. The genomic sequences are converted to a set of data point based on the formulas 35.1 and 35.2.

$$x_i = a_i \cdot m + g_i \cdot \sqrt{n} + c_i \cdot \sqrt{n} + t_i \cdot m \quad (35.1)$$

$$y_i = -a_i \cdot \sqrt{n} - g_i \cdot m + c_i \cdot m + t_i \cdot \sqrt{n} \quad (35.2)$$

The variable a_i, g_i, c_i, t_i correspond to the cumulative occurrence numbers of the bases until the current position while n and m are real numbers representing the vector length. These were set to $n = 0.5$ and $m = 0.75$ according to Yan et al. [62].

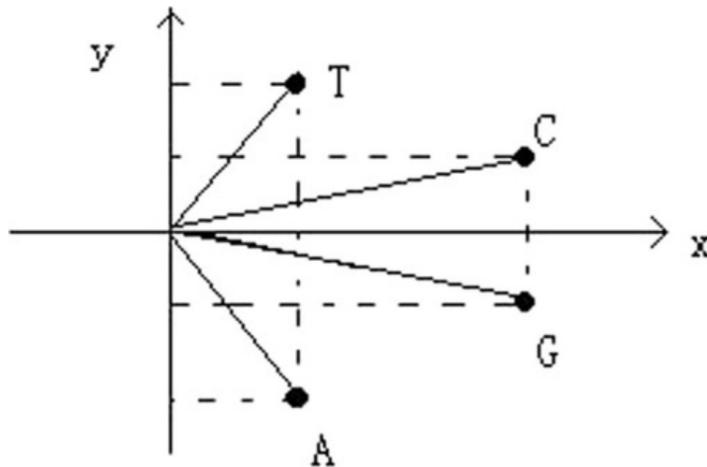


Figure 35.1: Pyrimidine-purine graph representing the four nucleotides of a DNA sequence. The vectors indicate the shift from the previous to the next data point of the curve when the respective base is present at the current position.

35.2.2 Building phylogenetic Trees based on 2D Curves

To create a tree using DNA sequences, we applied a simple 2D visualization method and created a distance matrix based on the 2D representation of sequences. The first step was to calculate the geometric center of the points using the coordinate data of DNA sequences extracted from the 2D representation:

$$x^0 = \frac{1}{N} \sum_{i=1}^N x_i, y^0 = \frac{1}{N} \sum_{i=1}^N y_i$$

The next step was to calculate the covariance matrix for the sequences:

$$\begin{cases} CM_{xx} = \frac{1}{N} \sum_{i=1}^N (x_i - x^0)(x_i - x^0) \\ CM_{xy} = \frac{1}{N} \sum_{i=1}^N (x_i - x^0)(y_i - y^0) = CM_{yx} \\ CM_{yy} = \frac{1}{N} \sum_{i=1}^N (y_i - y^0)(y_i - y^0) \end{cases}$$

Then the eigenvectors of the two eigenvalues of the covariance matrix were calculated according to the following formula:

$$EV_k^i = (EV_{k,1}^i, EV_{k,2}^i)^T, i, j = 1, 2, \dots, M, k = \lambda_1, \lambda_2$$

These are then used to calculate the arccosinus between sequences according to this formula:

$$\theta_{ij} = \arccos\left(\frac{EV_k^i EV_k^j}{|EV_k^i||EV_k^j|}\right), i, j = 1, 2, \dots, M, k = \lambda_1, \lambda_2$$

The angles are then summed for two sequences each.

$$\theta_{ij} = \theta_{ij}^{\lambda_1} + \theta_{ij}^{\lambda_2}, i, j = 1, 2, \dots, M$$

To calculate the distance between the sequences the Euclidean distance between the sequences was also needed. This was calculated with this formula.

$$d_{ij} = \sqrt{(x_i^0 - x_j^0)^2 + (y_i^0 - y_j^0)^2}, i, j = 1, 2, \dots, M$$

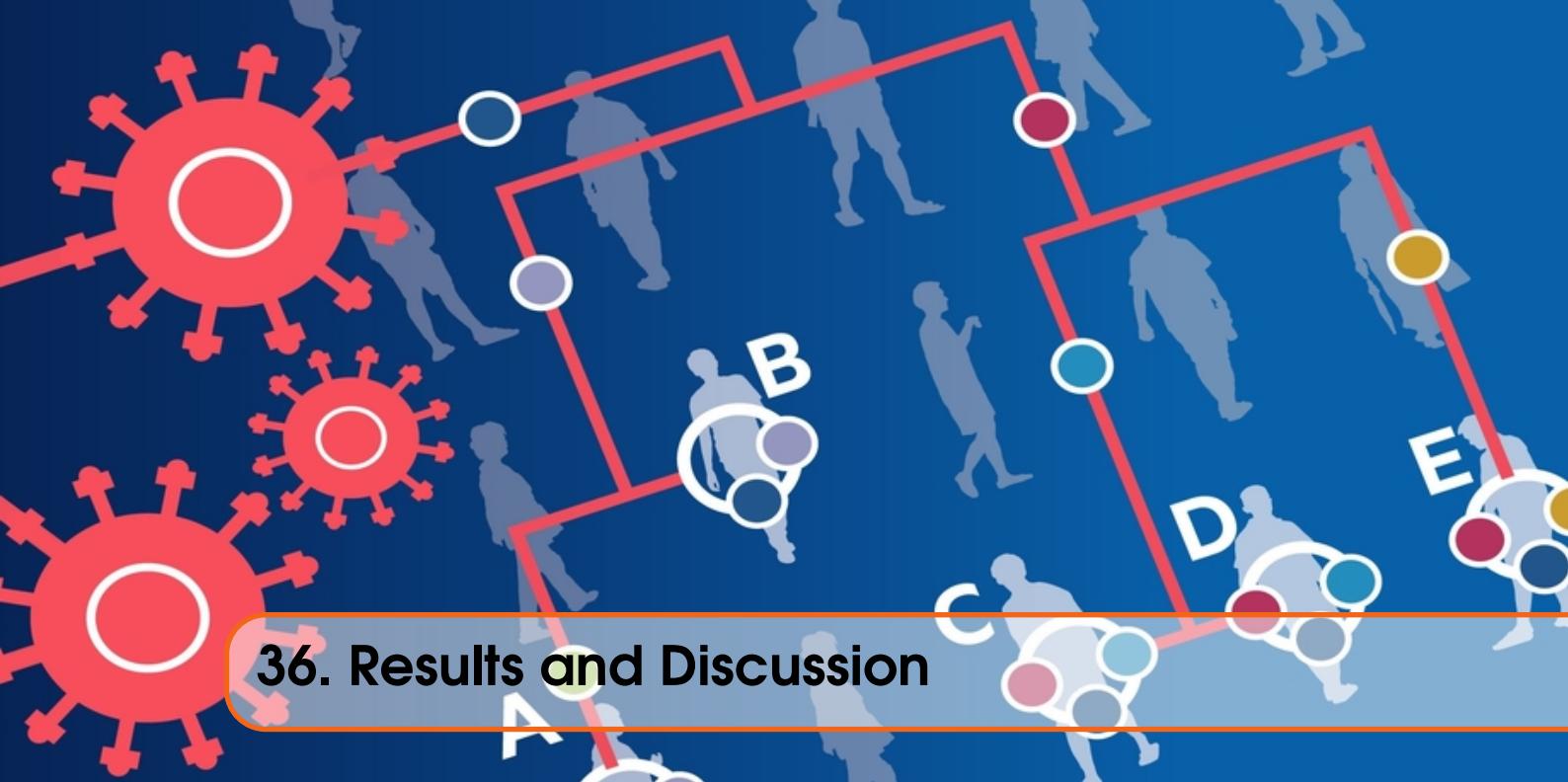
The whole distance for the sequences will be calculated as :

$$D_{ij} = d_{ij}x\theta_{ij}, i, j = 1, 2, \dots, M$$

Using these formulas the distance matrix can be estimated and the phylogenetic tree generated.

35.2.3 Metrics to compare phylogenetic Trees

Three different distance metrics were compared on both phylogenetic trees : Euclidean and Robinson foulds weighted and unweighted. Robinson foulds takes the symmetric difference between two trees into account by adding all splits that are different between tree A and tree B.



36. Results and Discussion

Host Analysis

Beginning with a first visual inspection of the produced plots for the origin analysis, both methods produce a correct reproduction of the hypothesised spillover event from the horseshoe bat host to the human host for SARS-CoV-2 (Figure 36.1 and 36.2). Visualizing the sequences as 2D curves (Figure 36.3) shows a clear separation of two groups: Wigeon, Rodent and MERS on the one hand and SARS-CoV-2 , Red junglefowl, Horeshoe bat and pangolin on the other. Because the tree topology is based on these curves it can easily be retraced. Interesting is the difference in grouping for MERS. While UPGMA puts MERS within reach of SARS-CoV-2 , the method by Liao et al. classifies it near rodents. This is unsupported by the literature as MERS originated as a possible spillover event at the interface of human and camels. Thus UPGMA classification might be more appropriate.

SARS-CoV-2 Analysis

For the larger data set of SARS-CoV-2 the visualization of sequences via the Liao et al. method fails (Figure 36.4). The sequences are too similar in alignment and thus no useful information can be extracted. Comparing the produced tree topology by UPGMA and 2D curve method distinct differences can be seen (Figure 36.5 and 36.6). UPGMA produces individual clusters of geographic similarity. Differences in time are also reproduced. On the other hand the Liao et al. method is also able to show geographic similarities and deviations in time. Both methods seem appropriate in their reproduction of spreading pathways for SARS-CoV-2.

The comparison of the two trees for different hosts and human respectively were done by different metrics which are mentioned in table 36.1. It is apparent that the low values for all three metrics for the SARS-CoV-2 tree suggest a lower deviation then for the host tree. The significance of these metric values is questionable as no reference value is present for comparison. Nevertheless all values suggest a difference in tree topologies for both approaches.

While the method by Liao et al. has clear advantages in computational time and efficiency it also has distinct disadvantages. If the sequences are very close in alignment as our testing of the SARS-CoV-2 data set suggests the method falls apart. The differences are to marginal to produce meaningful results which in turn produces banal 2D curve visualisation. On the other hand, if the

taxa are distinct enough the method performs quite well, showing an interesting way to visualise distinct differences in sequences via a 2D curve visualisation. The produced tree topology while different from UPGMA seems correct.

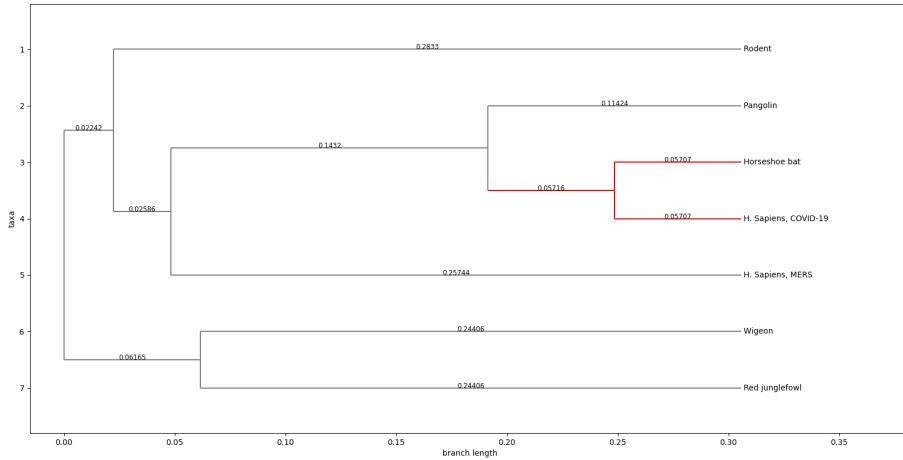


Figure 36.1: Phylogenetic tree of seven different sequences computed by the UPGMA method (see Section 32.2). Visualized as red is the closest sequence of horseshoe bat and SARS-CoV-2 in humans, suggesting a possible spillover event.

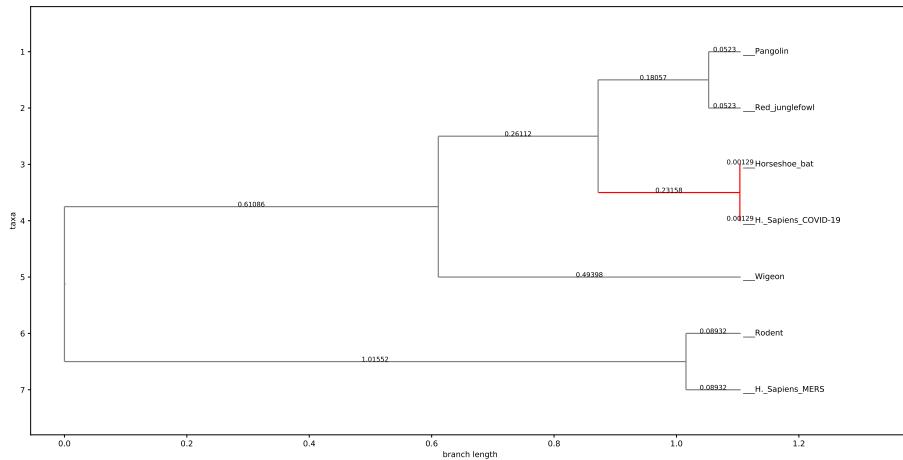


Figure 36.2: Phylogenetic tree of seven different sequences computed by the Liao et al. method (see Section 35.2). The closest classification of horseshoe bat and SARS-CoV-2 in humans is reproduced, suggesting a possible spillover event too. The phylogenetic tree was created as described in Section 33.1.2.

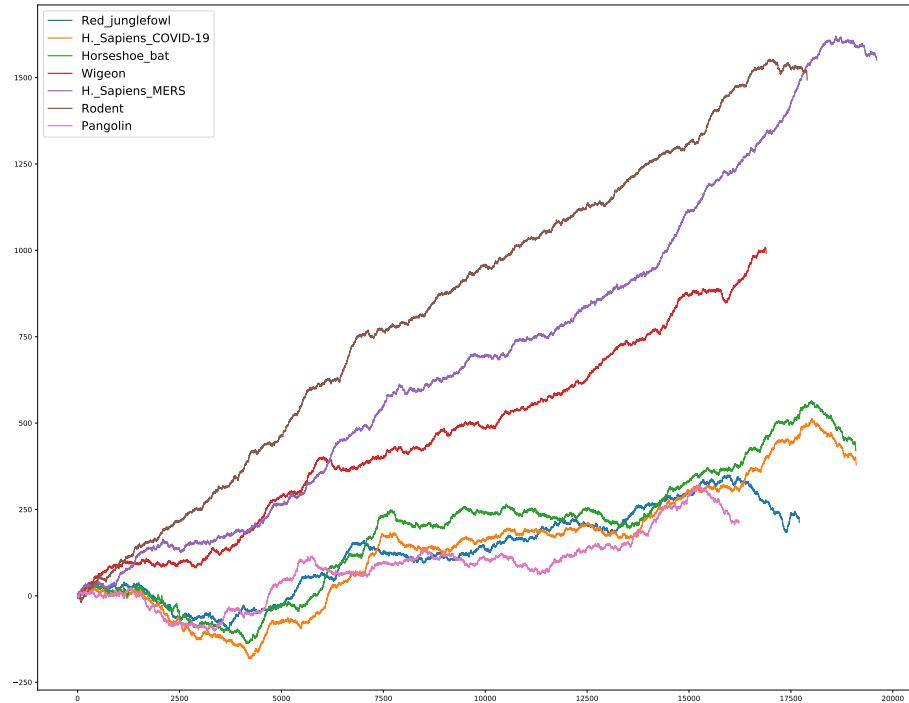


Figure 36.3: Pyrimidine-purine graph representing the shifts of individual bases for seven different sequences. While four curves are drawn together, three are deviating. The graphical representation was created as described in Section 35.2. On the x-axis u can see the position in the samples and on the y-axis the corresponding score at this position.

Metrics	phylogenetic tree different hosts	phylogenetic tree human
Robinson-Fould weighted	3.475	0.075
Robinson-Fould unweighted	6	186
Euclidean metric	1.728	0.025

Table 36.1: Comparison of two trees with three different types of metric values.

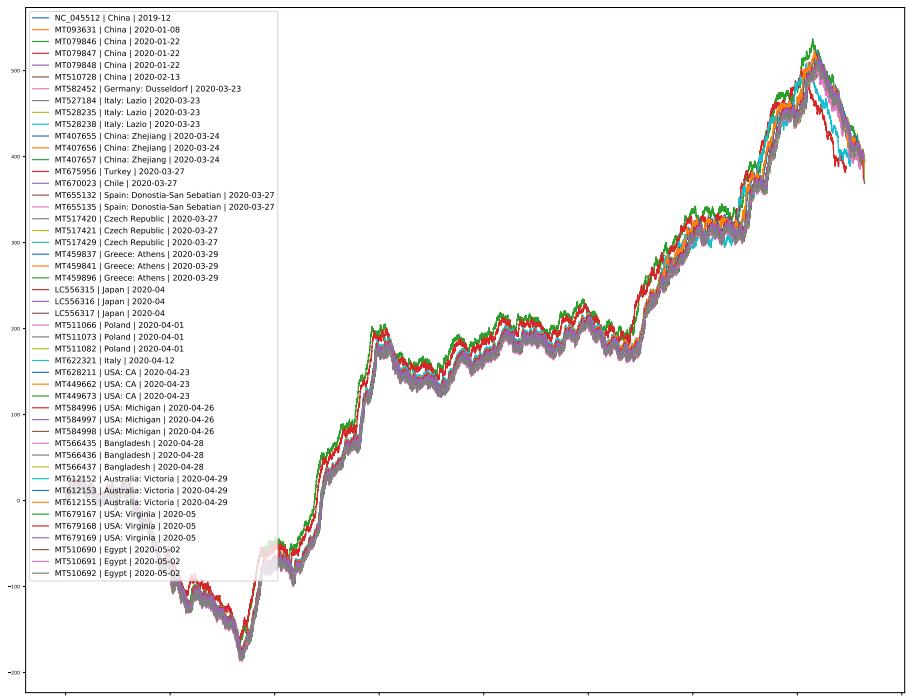


Figure 36.4: Pyrimidine-purine graph representing the shifts of individual bases for 40 different sequences of SARS-CoV-2. Because the alignment of sequences is too similar all curves are drawn together and thus it is difficult to extract useful information. The graphical representation was created as described in Section 35.2. On the x-axis u can see the position of the summed up nucleotides of the samples and on the y-axis the corresponding score at this position.

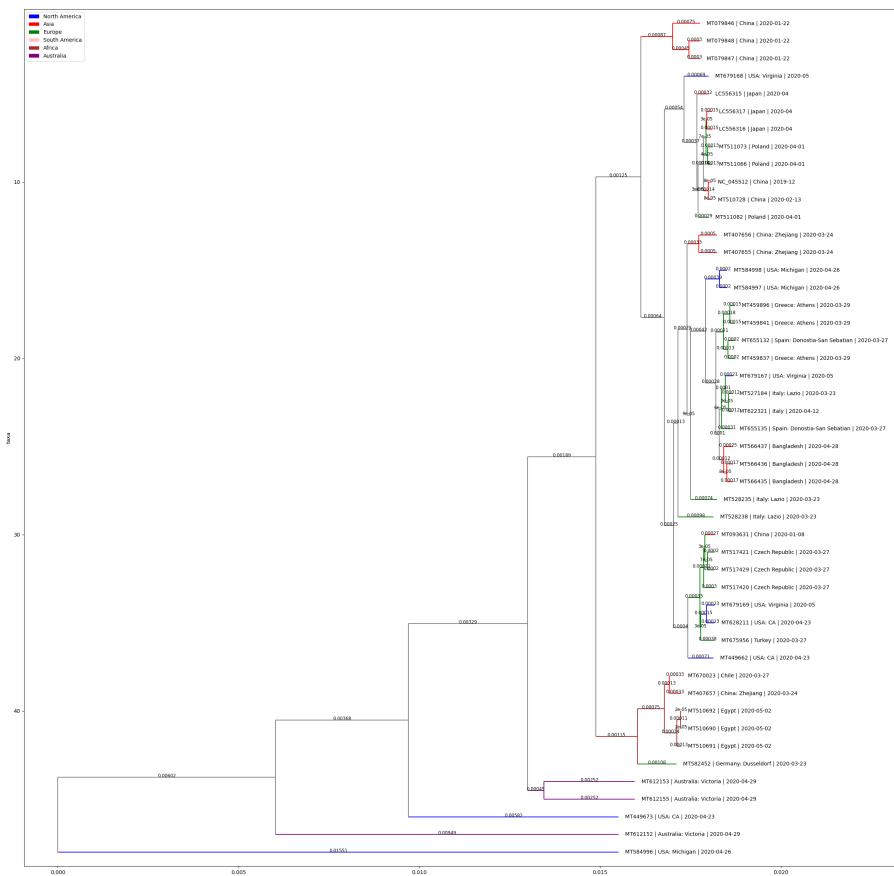


Figure 36.5: Phylogenetic tree of 40 different sequences from SARS-CoV-2 computed by the UPGMA method (see Section 32.2). Distinct geographic clusters and time-based deviations can be seen.

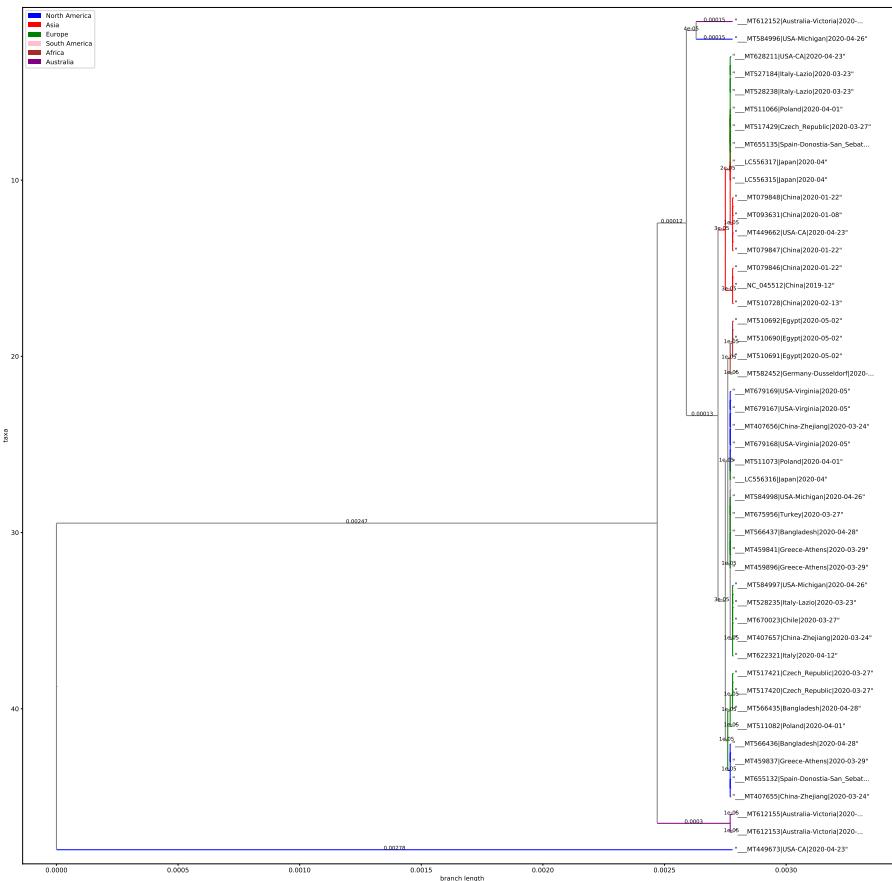
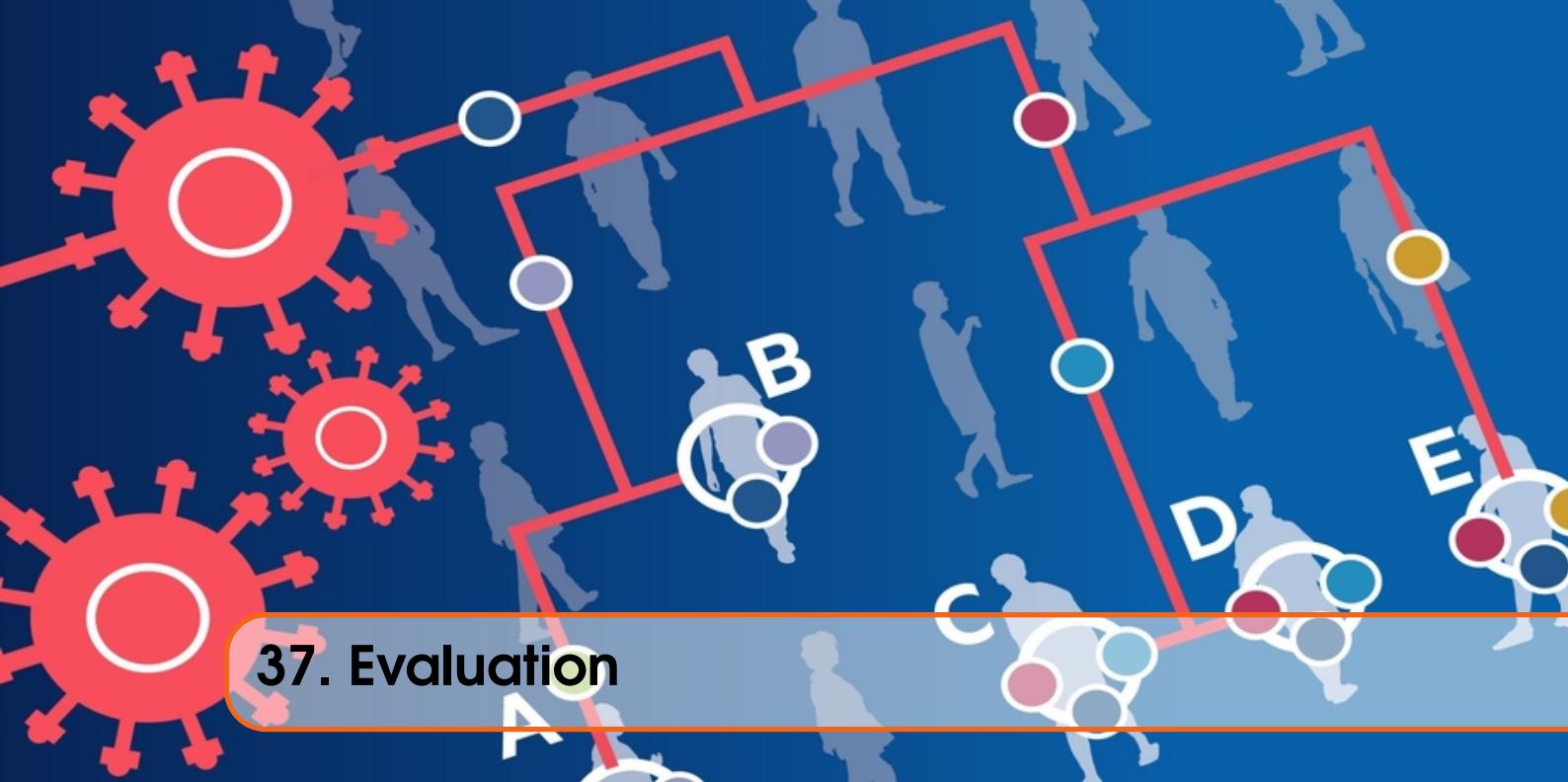


Figure 36.6: Phylogenetic tree of 40 different sequences from SARS-CoV-2 computed by the Liao et al. method (see Section 35.2). Values for the distance on each branch length had to be omitted as the computed values were to small.



37. Evaluation

37.1 Project Rating and Problems

Overall this week's project was interesting to go through because it presented an unusual approach for a common problem. On top of that, the graphical approach was easy to implement due to the paper [31] that described the implementation clearly. The biggest benefit was the short running time, which allowed us to explore the ramifications of different parameter settings or dynamics within the data set by changing the samples. Still, the graphical approach was not performing as well as the classical methods, thus we will most likely not use it in any future project, with the exception of runtime becoming a major factor, due to huge sample sizes.

Project 9: Drug Repurposing

38	Introduction to Drug Repurposing	167
38.1	Background	
38.2	Goal of the Project	
38.3	Outcomes	
39	Methods for Drug Repurposing	169
39.1	Computational Methods	
39.2	Experimental Methods	
39.3	Neural Networks that predict Drug-Target Interactions	
40	Find potential Drugs for treating COVID-19	173
40.1	Methods	
40.2	Results and Discussion	
40.3	Conclusion	



38. Introduction to Drug Repurposing

38.1 Background

New drugs are continually required by the healthcare system to either treat new emerging diseases or to improve the chances of recovery from an already existing drug. Given the high substantial costs and the slow pace of the fundamental drug discovery and development process, drug repurposing emerged as a growing research field. It represents a strategy that aims to identify new use cases for already existing and approved drugs outside the scope of the original medical indication [44]. The approach offers a better risk-versus-reward trade-off as compared with other drug development strategies, and can help to speed up the process tremendously [3].

In the context of COVID-19, where humanity is exposed to a global pandemic, there is an urgent requirement for a drug to combat the disease. The time it takes to find and develop a suitable drug is hereby a crucial factor. Since many drugs have multiple protein targets and many diseases share overlapping molecular pathways [24], reusing drugs can significantly reduce the cost, time and risks of the drug development process using the fast-growing depth of computational approaches. Based on the knowledge about other infectious diseases such as the Middle East Respiratory Syndrome (MERS) and the Severe Acute Respiratory Syndrome (SARS), several drug repurposing options are being considered and are currently under investigation to control COVID-19 [37].

38.2 Goal of the Project

In the first part of this week's project an overview of the available methodical approaches in the research field of drug repurposing is given. Afterwards, two deep neural networks are trained and performed with the aim to identify potential drugs for treating COVID-19: Deep Drug-Target Binding Affinity Prediction (deepDTA) and Molecule Transformer Drug-Target Interaction (MT-DTI).

38.3 Outcomes

The range of methods was separated into computational and experimental techniques. Unfortunately, training the deepDTA model was associated with a very long runtime. After a first run, where

the program broke due to RAM exceedance, the second run did not finish in time to analyze the results. However, MT-DTI was able to determine Atazanavir ($K_d = 94.94$ nM) to be the best fitting chemical compound in terms of inhibitory potency against the SARS-CoV-2 3C-like proteinase. Additionally, Remdesivir ($K_d = 113.13$ nM), efavirenz ($K_d = 199.17$ nM), ritonavir ($K_d = 204.05$ nM) and dolutegravir ($K_d = 336.91$ nM) also proved to be promising compounds.



39. Methods for Drug Repurposing

The goal of drug repurposing is to identify promising drug candidates for the treatment of various diseases. The advancements within the field of drug repurposing are manifold paving the way for development of systematic approaches towards compound identification [52]. Fundamentally there are two differing methods: computational and experimental. Computational methods involve the analysis of existing data such as gene expression, genotype, chemical structures and electronic health records. By using state of the art computational procedures new hypothesis for drug repurposing can be formulated. In contrast, experimental methods involve the usage of laboratory work like antibody binding, affinity chromatography and mass spectrometry to identify new chemical structures and/or binding sites of known chemical compounds.

39.1 Computational Methods

Signature matching:

Signature matching involves comparing the ‘signature’ of a drug i.e. its characteristics such as its transcriptomic, structural or adverse effect profile with that of another drug or disease phenotype. The characteristics of the drug is based on transcriptomic (RNA), proteomic or metabolomic data, chemical structures or adverse event profiles. Matching transcriptomic signatures can be used to make further drug-drug and drug-disease similarity comparisons [44].

Molecular docking:

In order to find the binding site of complementary chemical compounds one of the most common approaches is molecular docking. Molecular docking involves the prediction of the binding from a ligand of a drug and its target a protein. By finding these binding sites their binding affinity can be calculated. Two types are differentiated: 1. Conventional docking with one target and multiple ligands and 2. Inverse docking with several targets and one ligand. Using these predictions can lead to strong implications for the creation of new drug therapies and compounds cementing molecular docking as the most frequently used method in structure-based drug design.

Genetic association:

Due to advances made in genotyping technology, genome wide association studies (GWAS) found a tremendous increase over the past ten years. The main goal of GWAS is to identify the alterations in the genome associated with common diseases and give insight into its biological expression. This data in turn could be shared between the disease phenotypes and thus lead to identify novel attacking points for the drug treatment.

Pathway mapping:

Pathway mapping explores the avenue of high throughput data e.g. gene expression data to map biological pathways associated with diseases to drug networks. These drug network associated pathways for formulated drug treatments are constructed from the phenotypic drug profile. Using state of the art machine learning methods a semi-supervised relabeling of the mapping can be constructed thus finding novel drug treatments for different biological pathways.

Novel data sources:

Large-scale in vitro drug screens with paired genomic data, EHR-linked large biobanks and self-reported patient data are novel avenues to exploit for drug repurposing.

39.2 Experimental Methods

Binding assays to identify relevant target interactions:

Techniques such as affinity chromatography and mass spectrometry can be used to identify novel targets of known drugs. By using an antibody cocktail of specific binding targets in conjunction with a chemical tracer molecule, new structures can be visualized and identified.

Binding assays serve as an important experimental approach for target validation.

Phenotypic screening:

Disease relevant effects associated with the compounds can be screened with the help of phenotypic screening. Compounds need to be screened in order to get approval for the indication of repurposing. High-throughput phenotypic screening of compounds using in vitro or in vivo disease models can indicate potential for clinical evaluation.

39.3 Neural Networks that predict Drug-Target Interactions

Drug-target interactions play a significant role in the identification of suitable drugs for targeting the specific sites. Most of the computational methods mainly focus on binary classification [9] with the goal to determine the interaction of the drug-target pairs. Due to the high availability of affinity data in drug-target databases, the use of unique and advanced methods such as DeepDTA and MT-DTI are increasingly favored. These two deep-learning-based approaches paved the way for the prediction of accurate drug-target binding affinities by making use of Convolutional Neural Networks (CNN) architectures.

39.3.1 DeepDTA

Deep Drug-Target Binding Affinity Prediction (DeepDTA) is a CNN-based prediction model whose task is the prediction of the binding affinity value of drug-target pairs [17]. It automatically incorporates useful and necessary features of raw ligands and proteins into the model to predict drug-protein interactions. The model contains two separate CNN blocks, which convolute separately over the SMILES strings and protein sequence representations. CNNs make use of the kernel trick, where data is minimized in its dimensions by folding the data with a kernel matrix computation. The entire convolution is here performed by three consecutive 1D-convolutional layers, that each apply an increasing amount of filters compared to the previous layer. The convolutional layers are then followed by a max-pooling layer to calculate the most relevant features, because the

resulting SMILES representation and sequence representation are concatenated afterwards. The combined representation features are then fed into three consecutive fully-connected (FC) layers with a dropout process in between ($\text{rate}=0.1$) to avoid overfitting. This part of the pipeline is called the actual DeepDTA. It comprises of 1024 nodes in the first two FC layers, while the final layer consists of 512 nodes. The model is trained by minimizing the error between real and expected value and Rectified Linear Unit (ReLU) was used as activation function.

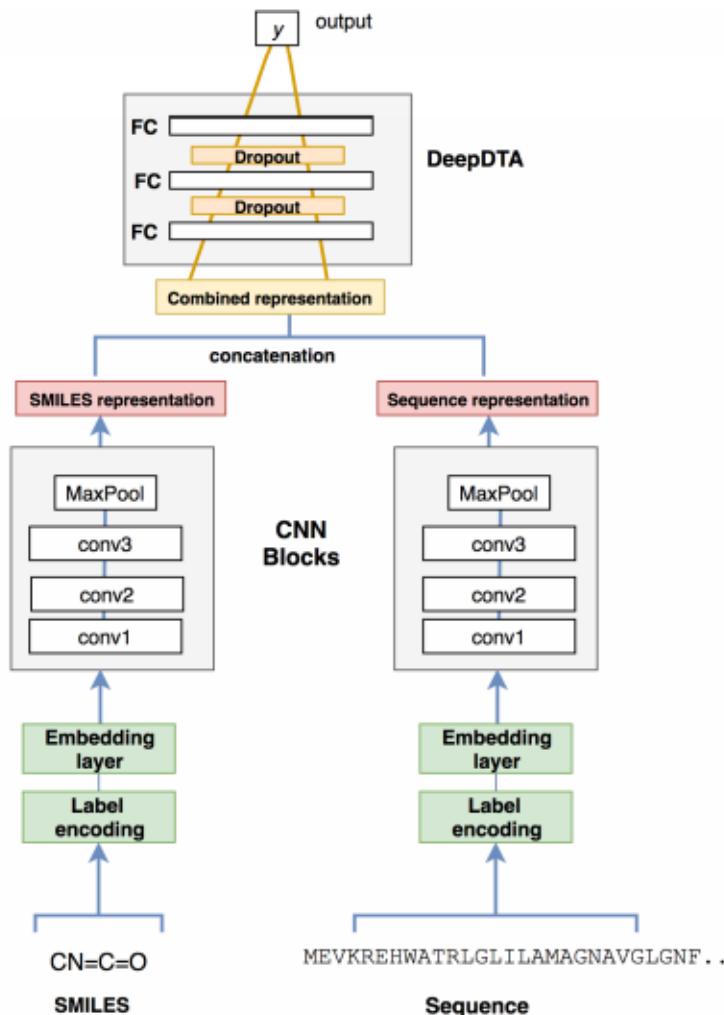


Figure 39.1: DeepDTA model with two CNN blocks to learn from both SMILES strings and protein sequences. The concatenated results are finally fed into three fully-connected layers with two dropout layers in between [17].

39.3.2 MT-DTI

Molecule Transformer-Drug Target Interaction (MT-DTI) is a method that can be used to predict binding affinity values between commercially available antiviral drugs and target proteins. The fundamental idea behind the model is, however, inspired by natural language processing (NLP): understanding a molecule sequence is in a way analogous to understanding a language. Therefore, MT-DTI was pre-trained with 'chemical language' in form of approximately 1,000,000 compounds contained in the publicly available PubChem database [9]. The goal of NLP to extract complex pattern from word sequences is here reused to identify important structures in chemical compounds with the help of a self-attention mechanism to learn the high dimensional

structure of a molecule from a given raw sequence which is an improvement to the CNN strategy of DeepDTA. The model accepts the input in SMILES and FASTA format respectively for a molecule or a protein. The molecule sequence is fed into the 'Molecule Transformer', a multi-layered bidirectional transformer encoder, that produces the molecule encoding [49]. Separately, a protein encoding is also generated. Both of the encodings are fed together into the multi-layered feed-forward network, making use of multiple interaction dense layers, followed by the last regression layer, which predicts the final binding affinity scores [49].

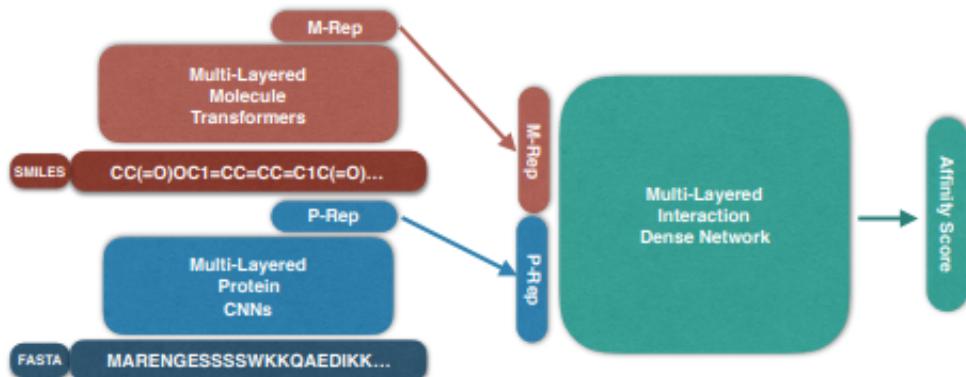


Figure 39.2: The Proposed DTI Model Architecture. Inputs are molecule (SMILES) and protein (FASTA) and the regression output is the affinity score between these two inputs [9].



40. Find potential Drugs for treating COVID-19

40.1 Methods

At first, we tried to train and perform DeepDTA [17] with the aim to find promising drugs for the treatment of COVID-19. The Kiba data set [54] was used for training which includes kinases and associated inhibitors and consists of 2111 unique drugs, 229 proteins and 118,254 interactions. The following initial parameter settings were used for the training:

Parameter	Value
<i>num_windows</i>	32
<i>seq_window_lengths</i>	8, 12
<i>smi_window_lengths</i>	4, 8
<i>batch_size</i>	256
<i>num_epoch</i>	5
<i>max_seq_len</i>	1000
<i>max_smi_len</i>	100

Table 40.1: Model parameters used for the training of DeepDTA on the Kiba data set.

Unfortunately, the process of training and finding the optimal hyper parameter values was very time-consuming and did not complete in time. After a first run, where the program broke due to RAM exceedance, the second run did not finish in time to analyze the results.

Consequently, we followed another approach. In the study of Beck et al. [5] the pre-trained MT-DTI model was used to identify commercially available drugs that could act on viral proteins of SARS-CoV-2. At this point it is important to clarify that the following results are not self-computed. The upcoming Section only describes the results that were achieved by Bo Ram Beck and the contributors.

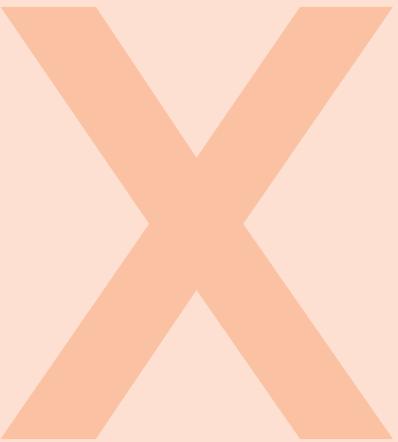
40.2 Results and Discussion

The experiments was focused on post-entry replication-associated proteins as their primary targets to suppress the viral replication. The amino acid sequences of 3C-like proteinase, RNA-dependent RNA polymerase, helicase, 3'-to-5' exonuclease, endoRNase and 2'-O-ribose methyltransferase of the SARS-CoV-2 replication complex were extracted from the SARS-CoV-2 whole genome sequence and treated as potential targets. In addition, a list of top commercially available antiviral drugs was considered that could potentially hinder the multiplication cycle of SARS-CoV-2. The subsequent aim is to get the opportunity to develop an effective drugs based on the AI-proposed drug candidates.

To underline the meaningfulness of the results, the MT-DTI model was compare to other established models (DeepDTA [17], SimBoost [23] and KronRLS [41]) and was able to achieve the best results among them. Finally, Atazanavir could be identified to be the best fitting chemical compound as it showed an inhibitory potency of $K_d = 94.94$ nM against SARS-CoV-2 3C-like proteinase as same as promising inhibitory potency values to the other four subunits of the SARS-CoV-2 replication complex. It may eventually able to inhibit all of them simultaneously. Furthermore, ganciclovir was predicted to bind to three subunits of the replication complex (RNA-dependent RNA polymerase, 3'-to-5' exonuclease, and RNA helicase), while Lopinavir and ritonavir both were predicted to have a potential affinity to the Sars-CoV-2 helicase.

40.3 Conclusion

The major advantage that already commercially available drugs identified by drug repurposing models for the treatment of COVID-19 can immediately be applied to the patient, make the technique very interesting especially in times of a global pandemic. Unfortunately, training the DeepDTA model was a highly time-consuming process and did not work out. In the end, the MT-DTI model however provided very promising results for binding affinities without domain knowledge. Nevertheless, Atazanavir, the substance that could possibly inhibit all subunits of the replication complex at the same time, need to be further tested and validated in vitro, in vivo, and in a wide range of clinical trials for efficacy and safety.



Appendix

41 Weekly Member Contribution 177

42 Bibliography 181

Bibliography 181

Articles

Books

Webpages

41. Weekly Member Contribution

Project 1

Kenrick Schulze: Average

Tobias Winterhoff: Average

Julius Tembrockhaus: Average

Natalja Amiridze: Average

Raghavendra Tikare: Average

Project 2

Kenrick Schulze: Average

Tobias Winterhoff: Above average

Julius Tembrockhaus :Above Average

Natalja Amiridze: Average

Raghavendra Tikare: Less than Average

Project 3

Kenrick Schulze: Above Average

Tobias Winterhoff: Average

Julius Tembrockhaus :Average

Natalja Amiridze: Average

Raghavendra Tikare: Average

Project 4

Kenrick Schulze: Average

Tobias Winterhoff: Above average

Julius Tembrockhaus: Average

Natalja Amiridze: Average

Raghavendra Tikare: Less than Average

Project 5

Kenrick Schulze: Average

Tobias Winterhoff: Average

Julius Tembrockhaus: Average

Natalja Amiridze: Average

Raghavendra Tikare: Average

Project 6

Kenrick Schulze: Average

Tobias Winterhoff: Average

Julius Tembrockhaus: Nothing (sickness certificate available)

Natalja Amiridze: Above Average

Raghavendra Tikare: Above Average

Project 7

Kenrick Schulze: Above Average

Tobias Winterhoff: Average

Julius Tembrockhaus: Above Average

Natalja Amiridze: Average

Raghavendra Tikare: Average

Project 8

Kenrick Schulze: Average

Tobias Winterhoff: Average

Julius Tembrockhaus: Average

Natalja Amiridze: Average

Raghavendra Tikare: Average

Project 9

Kenrick Schulze: Average

Tobias Winterhoff: Average

Julius Tembrockhaus: Average

Natalja Amiridze: Average

Raghavendra Tikare: Average



42. Bibliography

Articles

- [1] Elissa M Abrams and Stanley J Szeffler. “COVID-19 and the impact of social determinants of health”. In: *The Lancet Respiratory Medicine* (2020) (cited on page 56).
- [2] Kristian G Andersen et al. “The proximal origin of SARS-CoV-2”. In: *Nature medicine* 26.4 (2020), pages 450–452 (cited on pages 27, 28, 133, 146).
- [3] Ted T Ashburn and Karl B Thor. “Drug repositioning: identifying and developing new uses for existing drugs”. In: *Nature reviews Drug discovery* 3.8 (2004), pages 673–683 (cited on page 167).
- [4] Ali Bazghandi. “Techniques, advantages and problems of agent based modeling for traffic simulation”. In: *International Journal of Computer Science Issues (IJCSI)* 9.1 (2012), page 115 (cited on page 65).
- [5] Bo Ram Beck et al. “Predicting commercially available antiviral drugs that may act on the novel coronavirus (2019-nCoV), Wuhan, China through a drug-target interaction deep learning model”. In: *bioRxiv* (2020). DOI: 10.1101/2020.01.31.929547. eprint: <https://www.biorxiv.org/content/early/2020/02/02/2020.01.31.929547.full.pdf>. URL: <https://www.biorxiv.org/content/early/2020/02/02/2020.01.31.929547> (cited on pages 10, 173).
- [9] Keunsoo Kang Bonggun Shin Sungsoo Park and Joyce C. Ho. “Self-Attention Based Molecule Representation for Predicting Drug-Target Interaction”. In: (2019) (cited on pages 170–172).
- [10] George EP Box. “All models are wrong, but some are useful”. In: *Robustness in Statistics* 202 (1979), page 549 (cited on page 67).
- [12] Joseph Paul Cohen et al. “COVID-19 Image Data Collection: Prospective Predictions Are the Future”. In: *arXiv 2006.11988* (2020). URL: <https://github.com/ieee8023/covid-chestxray-dataset> (cited on page 23).

- [18] Kuldeep Dhamra et al. “SARS-CoV-2: Jumping the species barrier, lessons from SARS and MERS, its zoonotic spillover, transmission to humans, preventive and control measures and recent developments to counter this pandemic virus”. In: (2020) (cited on pages 27, 133).
- [20] Ana S Gonzalez-Reiche et al. “Introductions and early spread of SARS-CoV-2 in the New York City area”. In: *Science* (2020) (cited on page 147).
- [22] Eneida L Hatcher et al. “Virus Variation Resource—improved response to emergent viral outbreaks”. In: *Nucleic acids research* 45.D1 (2017), pages D482–D490 (cited on pages 27, 135).
- [23] Tong He et al. “SimBoost: a read-across approach for predicting drug–target binding affinities using gradient boosting machines”. In: *Journal of cheminformatics* 9.1 (2017), pages 1–14 (cited on page 174).
- [24] Rachel A Hodos et al. “Computational approaches to drug repurposing and pharmacology”. In: *Wiley interdisciplinary reviews. Systems biology and medicine* 8.3 (2016), page 186 (cited on page 167).
- [26] Mark Johnson et al. “NCBI BLAST: a better web interface”. In: *Nucleic acids research* 36.suppl_2 (2008), W5–W9 (cited on page 137).
- [29] Tommy Tsan-Yuk Lam et al. “Identifying SARS-CoV-2 related coronaviruses in Malayan pangolins”. In: *Nature* (2020), pages 1–6 (cited on pages 28, 146).
- [30] Thomas Lampert, Elena von der Lippe, and Stephan Müters. “Prevalence of smoking in the adult population of Germany”. In: (2013) (cited on page 56).
- [31] Bo Liao, Xuyu Xiang, and Wen Zhu. “Coronavirus phylogeny based on 2D graphical representation of DNA sequence”. In: *Journal of computational chemistry* 27.11 (2006), pages 1196–1202 (cited on pages 151, 163).
- [33] Zhixin Liu et al. “Composition and divergence of coronavirus spike proteins and host ACE2 receptors predict potential intermediate hosts of SARS-CoV-2”. In: *Journal of medical virology* 92.6 (2020), pages 595–601 (cited on pages 28, 146).
- [34] Eric Lofgren et al. “Influenza seasonality: underlying causes and modeling theories”. In: *Journal of virology* 81.11 (2007), pages 5429–5436 (cited on page 107).
- [35] Lars Lorch et al. “A spatiotemporal epidemic model to quantify the effects of contact tracing, testing, and containment”. In: *arXiv preprint arXiv:2004.07641* (2020) (cited on page 67).
- [36] Fábio Madeira et al. “The EMBL-EBI search and sequence analysis tools APIs in 2019”. In: *Nucleic acids research* 47.W1 (2019), W636–W641 (cited on page 27).
- [37] Nisha Muralidharan et al. “Computational studies of drug repurposing and synergism of lopinavir, oseltamivir and ritonavir binding with SARS-CoV-2 Protease against COVID-19”. In: *Journal of Biomolecular Structure and Dynamics* (2020), pages 1–6 (cited on page 167).
- [39] World Health Organization et al. “Coronavirus disease 2019 (COVID-19): situation report, 73”. In: (2020) (cited on pages 71, 73).
- [40] Justin Page et al. “BamBam: Genome sequence analysis tools for biologists”. In: *BMC research notes* 7 (Nov. 2014), page 829. DOI: 10 . 1186 / 1756 - 0500 - 7 - 829 (cited on page 138).
- [41] Tapio Pahikkala et al. “Toward more realistic drug–target interaction predictions”. In: *Briefings in bioinformatics* 16.2 (2015), pages 325–337 (cited on page 174).
- [42] Zsolt Vizi Péter Boldog Tamás Tekeli et al. “Risk Assessment of Novel Coronavirus COVID-19 Outbreaks Outside China”. In: *clinical medicine* (2020) (cited on pages 21, 22).

- [44] Sudeep Pushpakom et al. “Drug repurposing: progress, challenges and recommendations”. In: *Nature reviews Drug discovery* 18.1 (2019), pages 41–58 (cited on pages 167, 169).
- [49] Bonggun Shin et al. “Self-attention based molecule representation for predicting drug-target interaction”. In: *arXiv preprint arXiv:1908.06760* (2019) (cited on page 172).
- [51] Patricio Solis and Hiram Carreño. “COVID-19 Fatality and Comorbidity Risk Factors among Confirmed Patients in Mexico”. In: *medRxiv* (2020) (cited on pages 74, 86).
- [52] Munir Sudeep Pushpakom Francesco Iorio et al. “Drug repurposing: progress, challenges and recommendations”. In: (2018) (cited on page 169).
- [53] Thorsten Suess et al. “The role of facemasks and hand hygiene in the prevention of influenza transmission in households: results from a cluster randomised trial; Berlin, Germany, 2009–2011”. In: *BMC infectious diseases* 12.1 (2012), page 26 (cited on pages 59, 77).
- [54] Jing Tang et al. “Making sense of large-scale kinase inhibitor bioactivity data sets: a comparative and integrative analysis”. In: *Journal of Chemical Information and Modeling* 54.3 (2014), pages 735–743 (cited on page 173).
- [55] Sean J Taylor and Benjamin Letham. “Forecasting at scale”. In: *The American Statistician* 72.1 (2018), pages 37–45 (cited on page 93).
- [57] Samantha M Tracht, Sara Y Del Valle, and James M Hyman. “Mathematical modeling of the effectiveness of facemasks in reducing the spread of novel influenza A (H1N1)”. In: *PloS one* 5.2 (2010) (cited on pages 59, 77).
- [60] Michael Worobey et al. “The emergence of SARS-CoV-2 in Europe and the US”. In: *bioRxiv* (2020) (cited on page 147).
- [61] Fan Wu et al. “A new coronavirus associated with human respiratory disease in China”. In: *Nature* 579.7798 (2020), pages 265–269 (cited on pages 28, 146).
- [62] Stephen S-T Yau et al. “DNA sequence representation without degeneracy”. In: *Nucleic acids research* 31.12 (2003), pages 3078–3080 (cited on page 153).
- [63] Jin-jin Zhang et al. “Clinical characteristics of 140 patients infected with SARS-CoV-2 in Wuhan, China”. In: *Allergy* (2020) (cited on page 55).

Books

- [8] William C Black, Barry J Babin, Rolph E Anderson, et al. *Multivariate data analysis*. Volume 5. 3 (cited on page 145).
- [27] McKendrick Kermack. *A contribution to the mathematical theory of epidemics*. Proc. Roy. Soc. A, Band 115, 1927 (cited on page 51).
- [28] Gebhard Kirchgässner and Jürgen Wolters. *Introduction to modern time series analysis*. Springer Science & Business Media, 2007 (cited on page 89).
- [56] Michel Tibayrenc. *Genetics and evolution of infectious diseases*. Elsevier, 2017 (cited on pages 27, 133).

Webpages

- [6] *Bevölkerung - Zahl der männlichen Einwohner in Deutschland nach Altersgruppen am 31. Dezember 2018.* <https://de.statista.com/statistik/daten/studie/1112607/umfrage/maennliche-bevoelkerung-in-deutschland-nach-altersgruppen/> (cited on page 56).

- [7] *Bevölkerung - Zahl der weiblichen Einwohner in Deutschland nach Altersgruppen am 31. Dezember 2018.* <https://de.statista.com/statistik/daten/studie/1112611/umfrage/weibliche-bevoelkerung-in-deutschland-nach-altersgruppen/> (cited on page 56).
- [11] Simon Burgermeister. *covid_sequence.* https://github.com/simonjuleseric2/covid_sequence. 2020 (cited on pages 27, 135).
- [13] *Coronavirus (COVID-19) death rate in Italy as of May 20, 2020, by age group.* <https://www.statista.com/statistics/1106372/coronavirus-death-rate-by-age-group-italy/>. Accessed: 2020-05-25 (cited on page 55).
- [14] *Coronavirus Disease 2019 (COVID-19) Daily Situation Report of the Robert Koch Institute.* https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Situationsberichte/2020-04-29-en.pdf?__blob=publicationFile (cited on page 56).
- [15] *dashboard John Hopkins University Center for Systems Science and Engineering.* <https://data.humdata.org/dataset/novel-coronavirus-2019-ncov-cases/>. Accessed: 2020-08-03 (cited on page 57).
- [16] *deaths-time-series-germany.* <https://tinyurl.com/t59cgxn/>. Accessed: 2020-08-03 (cited on page 57).
- [17] *DeepDTA: Deep Drug-Target Binding Affinity Prediction.* <https://academic.oup.com/bioinformatics/article/34/17/i821/5093245/>. Accessed: 2020-07-11 (cited on pages 170, 171, 173, 174).
- [19] *Diabetes Germany.* <https://www.diabetes-news.de/nachrichten/diabetes-daten-2020-das-sind-die-zahlen>. Accessed: 2020-05-24 (cited on page 73).
- [21] *Hamburg in Zahlen.* <https://www.hamburg.de/info/3277402/hamburg-in-zahlen/> (cited on page 78).
- [25] *Hypertension RKI.* <https://edoc.rki.de/handle/176904/2663>. Accessed: 2020-05-24 (cited on page 73).
- [32] Lisphilar. *COVID-19 data with SIR model.* <https://www.kaggle.com/lisphilar/covid-19-data-with-sir-model>. 2020 (cited on page 15).
- [38] *Obesity RKI.* https://www.rki.de/DE/Content/Gesundheitsmonitoring/Gesundheitsberichterstattung/GesundAZ/Content/H/Hypertonie/Inhalt/Blutdruck_DZHK.pdf?__blob=publicationFile. Accessed: 2020-05-24 (cited on page 73).
- [43] *Phylogeography.* <https://justinbagley.org/pages/phylogeog.html>. Accessed: 2020-06-29 (cited on page 139).
- [45] Robert-Koch-Institute. *COVID-19 case count in Germany state-by-state, over time.* <https://github.com/jgehrcke/covid-19-germany-gae>. 2020 (cited on pages 57, 99, 111, 113, 119, 120).
- [46] Nabeel Sajid. *COVID-19 Detection from X Ray Images of Lungs.* <https://www.kaggle.com/nabeelsajid917/Covid-19-detection-from-x-ray-images-of-lungs>. 2020 (cited on page 23).
- [47] *SARS-CoV-2 Steckbrief zur Coronavirus-Krankheit-2019 (COVID-19).* https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Steckbrief.html (cited on pages 56, 57).

-
- [48] Leonardo Setti et al. *Airborne Transmission Route of COVID-19: Why 2 Meters/6 Feet of Inter-Personal Distance Could Not Be Enough*. 2020 (cited on page 78).
 - [50] *SIR flow diagramm*. <https://www.lewuathe.com/covid-19-dynamics-with-sir-model.html/>. Accessed: 2020-05-24 (cited on page 53).
 - [58] *Verkaufsfläche im Einzelhandel je 1.000 Einwohner in Deutschland im Jahr 2014 nach Bundesländern*. <https://www.handelsdaten.de/deutschsprachiger-einzelhandel/verkaufsflaeche-einzelhandel-je-1000-einwohner-deutschland> (cited on page 74).
 - [59] *Worldometers Kernel Description*. <https://www.worldometers.info/coronavirus/countries-where-coronavirus-has-spread/>. Accessed: 2020-05-11 (cited on page 13).

