



**FACULTY OF SCIENCE AND TECHNOLOGY**

**MASTER'S THESIS**

Study programme/specialisation: Computer Science	Spring / Autumn semester, 20.18.  Open/Confidential
Author: Sandra Moen	..... (signature of author)
Programme coordinator: Prof. Erlend Tøssebro	
Supervisor(s): Prof. Erlend Tøssebro	
Title of master's thesis:  Automated collection of multi-source spatial information for emergency management	
Credits: 30 sp	
Keywords:  Statistics, API, Data Collection, Influenza	Number of pages: .....  + supplemental material/other: .....  Stavanger, ..... date/year

# **Automated collection of multi-source spatial information for emergency management**

Tracking the influenza seasons

**Sandra Moen**

A thesis presented for the degree of  
Master of Science in Computer Science



---

**University of  
Stavanger**

Department of Electrical Engineering and  
Computer Science  
University of Stavanger  
Norway  
Spring 2018

# **Automated collection of multi-source spatial information for emergency management**

Tracking the influenza seasons

**Sandra Moen**

## **Abstract**

Influenza epidemics costs both lives and a tremendous amount of resources for any country. Citizens that become sick are less productive and the overall quality of life is drastically reduced for the amount of the individuals period of illness as well as the community during a flu season. The ability to reduce the spread of infectious diseases saves both lives and resources as well as an improvement of the quality of life.

This project aims to explore the possibilities to detect influenza outbreaks as soon as they are happening with the use of relevant datasets available. Information about different aspects of a citizens life on a grand scale reveals patterns and trends that could be linked to an epidemic outbreak, and thus prove useful for active measurements against further spread on a early début.

The results show ...

Possible solutions to ...

# Acknowledgements

This thesis is considered an impressive achievement for the author, it was completed in spite of hardships endured. Under no circumstance should this thesis be considered a Norwegian accomplishment, for the oppression suffered they are deemed unworthy.

This thesis was written for the Department of Electrical Engineering and Computer Science at the University of Stavanger. Creating a means to solve problems that limit peoples lives have always been a real motivator. Predicting the flu season and hindering it in early stages would save an enormous amount of resources and improve life quality, this would be very rewarding. A special thanks to the supervisor for this project from the University of Stavanger Professor Erlend Tøssebro for his enthusiastic guidance and involvement, and the initiator who inspired incentive to the creation of this project as well as his continuous helpful guidance and involvement Phd fellow Lars Ole Grottenberg.

# Contents

<b>1</b>	<b>Introduction</b>	<b>8</b>
1.1	Background . . . . .	8
1.2	Objectives . . . . .	9
1.3	Outline . . . . .	11
<b>2</b>	<b>Related Works</b>	<b>13</b>
2.1	Spatiotemporal information from urban systems . . . . .	13
2.2	Spatiotemporal information from VGI . . . . .	13
2.3	Twitter . . . . .	13
2.4	Data management and critical infrastructure . . . . .	15
2.5	Goompy . . . . .	15
<b>3</b>	<b>Datasets used</b>	<b>17</b>
3.1	The Norwegian Institute of Public Health . . . . .	17
3.2	The Norwegian Public Roads Administration . . . . .	17
3.3	Twitter . . . . .	18
3.4	Kolumbus . . . . .	18
3.5	Ruter . . . . .	18
<b>4</b>	<b>Implementation</b>	<b>19</b>
4.1	The Backend . . . . .	19
4.1.1	The Norwegian Institute of Public Health . . . . .	19
4.1.2	The Norwegian Public Roads Administration . . . . .	20
4.1.3	Twitter . . . . .	22
4.1.4	Kolumbus . . . . .	23
4.1.5	Ruter . . . . .	23
4.2	The Frontend . . . . .	23
4.2.1	The GUI . . . . .	24
4.2.2	The Map . . . . .	25
4.2.3	The Scrollbar . . . . .	25
4.2.4	NIPH dataframe . . . . .	25
4.2.5	NPRA dataframe . . . . .	26
<b>5</b>	<b>Results</b>	<b>36</b>
5.1	TODO . . . . .	36
5.1.1	TODO . . . . .	36
5.1.2	Twitter . . . . .	36

<b>6 Discussion</b>	<b>37</b>
6.1 TODO . . . . .	37
6.1.1 workflow . . . . .	37
<b>7 Conclusion</b>	<b>38</b>
7.1 TODO . . . . .	38
7.1.1 Future works . . . . .	38
<b>A Appendix Title</b>	<b>39</b>

# List of Figures

1.1	NIPH, 2017 . . . . .	12
2.1	Figure from Grottenberg et al. [8] . . . . .	14
2.2	A Goompy implementation of Google's static map API . . . . .	16
4.1	Influenza virus observation . . . . .	20
4.2	Influenza-like illnesses season 2016/2017 . . . . .	21
4.3	Annual traffic 2002-2015 . . . . .	22
4.4	Bergen traffic 2002-2015 . . . . .	23
4.5	Oslo traffic 2002-2015 . . . . .	24
4.6	Weekly data of the city of Bergen . . . . .	25
4.7	Weekly data of the city of Oslo . . . . .	26
4.8	Weekly data of the city of Stavanger . . . . .	27
4.9	Geospatial bounds of Bergen. The green circles show where the traffic registration stations are, and the number reveals how many there are in that general area. . . . .	28
4.10	Geospatial bounds of Oslo. The green circles show where the traffic registration stations are, and the number reveals how many there are in that general area. . . . .	29
4.11	Geospatial bounds of Stavanger. The green circles show where the traffic registration stations are, and the number reveals how many there are in that general area. . . . .	30
4.12	Geospatial hourly bounds of Bergen . . . . .	31
4.13	Geospatial hourly bounds of Oslo . . . . .	32
4.14	Geospatial hourly bounds of Stavanger . . . . .	32
4.15	Tweets concerning ILS of 2018 . . . . .	33
4.16	Monthly passenger travel with Kolumbus . . . . .	33
4.17	Daily tickets sold with Ruter, the year of 2015 does not contain Oslo's underground train service passenger data . . . . .	34
4.18	The GUI . . . . .	34
4.19	NIPH comparing buttons panel . . . . .	35
4.20	NPRA query buttons panel . . . . .	35

# List of Tables

1.1	The Norwegian surveillance system for influenza . . . . .	9
1.2	Categories of societal consumptive behaviours . . . . .	10

# Chapter 1

## Introduction

### 1.1 Background

Influenza is an exceedingly contagious viral infection which gives high fever, general pain, and respiratory symptoms[2]. An estimated five to ten percent of the population becomes infected during the yearly influenza season, which is generally in the winter. The virus is especially dangerous to the elderly and to pregnant people from the second-trimester. Annually between the months of December and April people of the northern hemisphere are struck by influenza epidemics. Since this is a seasonal occurrence mitigation or even elimination of the effects are a priority and thus observation and research are initiated. From a historical perspective, it is known that influenza can have overwhelming destructive consequences if left unreservedly to ravage the population. The last three larger pandemics were the Asian flu of 1957, the flu of 1968 which originated in Hong Kong and the H1N1 (swine flu) virus of 2009, which respectively claimed the lives of 1.1 million, 1-4 million and 284500 people [3]. The World Health Organization (WHO) estimates an annual global infection of humans to be a rate of 5-15% [1], this causes 300.000 to 650.000 deaths per year[2], and about 1700 of these are Norwegians[3]. The virus mutates often which proves immunization by a vaccine to be a seasonal effort. Infection happens via droplets in the air inhaled, and even a small exposure expands to an all-out blitz which the immune system is forced to engage.

Diseases travel with humans as they commute or travel long distances and thus spread[4][5]. The gravity and influence of an infectious disease can have is also strongly correlated to social[6] and environmental[7] circumstances. The intricate and fluctuating spread of contagious diseases within a complex and mobile human domain means that a static and a uniform approach is sub-optimal because the real grasp of the structure is a more changing operation with its own convoluted variety of variables [8][9].

One of the fundamental requirements for efficient control of urban outbreaks is to maintain situational awareness of the extent, impact, and potential of ongoing outbreaks. To accomplish this, a series of clinical indicator-based surveillance systems monitor patient-general practitioner interaction, as well as laboratory-based analysis and intensive care unit (ICU) surveillance.

The current surveillance systems are heavily based on clinical indicators, and it is of interest to establish new mechanisms that make use of other indicators. Establishing surveillance systems based on societal indicators allow for detection

System	Function
NorSySS	Indicator-based surveillance of influenza-like illness in primary health care
Hospital (all ward) surveillance	Laboratory-based surveillance of hospitalised influenza cases
ICU surveillance	ICU treated flu patients. Data collected by the Norwegian Intensive Care Registry (pilot project since 2016/17)
Virological-surveillance	(1) Submission of data and samples from Norwegian laboratories testing for influenza. (2) Sentinel system, GP-based virological surveillance.
Norwegian mortality monitoring system (NorMOMO)	Surveillance of weekly all-cause excess mortality.
Seroepidemiological analysis	Annual survey of flu immunity in the population.

Table 1.1: The Norwegian surveillance system for influenza

of non-clinical factors that indicate the presence of influenza in society. Directly monitoring behaviour at the societal level may also provide the ability to detect emerging behaviour and pattern deviations that indicate the presence of influenza at an earlier stage than what can be accomplished through patient-doctor interaction.

The power to obtain enough information to detect possible trends of influenza seasons depends on successful integration between a multitude of different participants. Automatic extraction and processing of data is paramount for efficient analysis and gives a solid basis for an autonomous pathological detection system. Scalability is important in merging new relevant datasets as they become available in an ever-growing societal infrastructure. This thesis proposes a technology that would become an influential part of a bigger foundation intertwined with a robust knowledgeable and organizational means to mobilize assets in order to respond to possible outbreaks as or even before they start. Such a system requires as many feasible input channels from different urban systems and resources as possible in order to become reliable.

## 1.2 Objectives

This thesis examines the viability of investigating, collecting and analysing relevant urban true-time data for a self-sufficient influenza seasonal recognition system. The management of seasonal influenza outbreaks is handled by public health officials and epidemiologists with the use of the national surveillance system provided by the Norwegian Institute of Public Health (NIPH)[3].

No	Indicator description
1	Public transport utilisation (Subway, trains, buses, light rail, etc.)
2	Toll road activations
3	Data traffic (internet traffic, cell phone networks)
4	Consumption of key indicator goods (Painkillers, Tamiflu, coughing medicine, etc.)
5	Utility use patterns in residential and commercial areas (Electricity, water, heating, etc.)
6	Use of key urban services (pharmacies, schools, GP offices, etc.)
7	Activity information from commercial stakeholders (stores, restaurants, etc.)

Table 1.2: Categories of societal consumptive behaviours

The Norwegian Syndromic Surveillance System (NorSySS) collects influenza-like illnesses (ILI) from general practitioners (GPs)[10], figure 1.1 shows a diagram of their process. The current NorSySS system relies upon reports of influenza-like illness from general practitioners (GPs). These subsystems compose part of the Norwegian influenza surveillance system and provide data with high reliability, but low timeliness. Typically the delay is over a week because it relies on clinical reports and laboratory endeavours, and leaves few ways to assess the societal impact of ongoing outbreaks. Measuring societal indicators based on the spatiotemporal components inherent in these data sources makes it possible to draw upon spatial epidemiological traditions to link societal behaviour to outbreaks of seasonal influenza with a significantly higher temporal resolution than found in current flu monitoring systems. The goal of this thesis is to determine whether a monitoring system of urban real-time data could do the same with less delay.

The main suggestion of this thesis is as influenza develops it reveals subtle patterns in societal behaviours that is detectable through a variety of mediums, e.g urban datasets from sewage, public transportation, medicinal purchases, recreational habits, social media and other such sources of public information, table 1.2 shows a more general view of such possible categories. With this suggestion, a tool to collect urban spatial datasets is needed and to present and visualize this information to best divulge the effect of the viral composition. This thesis focuses mainly on the Norwegian cities of Stavanger, Bergen, and Oslo. The datasets used in this thesis is explained more in chapter 3, they consist however of the NIPH ILI and virus observations, the different datasets from the NPRA showing traffic patterns, social media of Twitter reporting symptoms directly from the public of Norway and two public transportation providers of the cities Stavanger and Oslo. Unfortunately more datasets could not be obtained within the time-scope of this thesis, but nonetheless, they provide a solid basis for examination and development.

## 1.3 Outline

The thesis is structured into seven chapters.

Chapter 2 describes related works of what others have found useful as tools and other proven effective measurements.

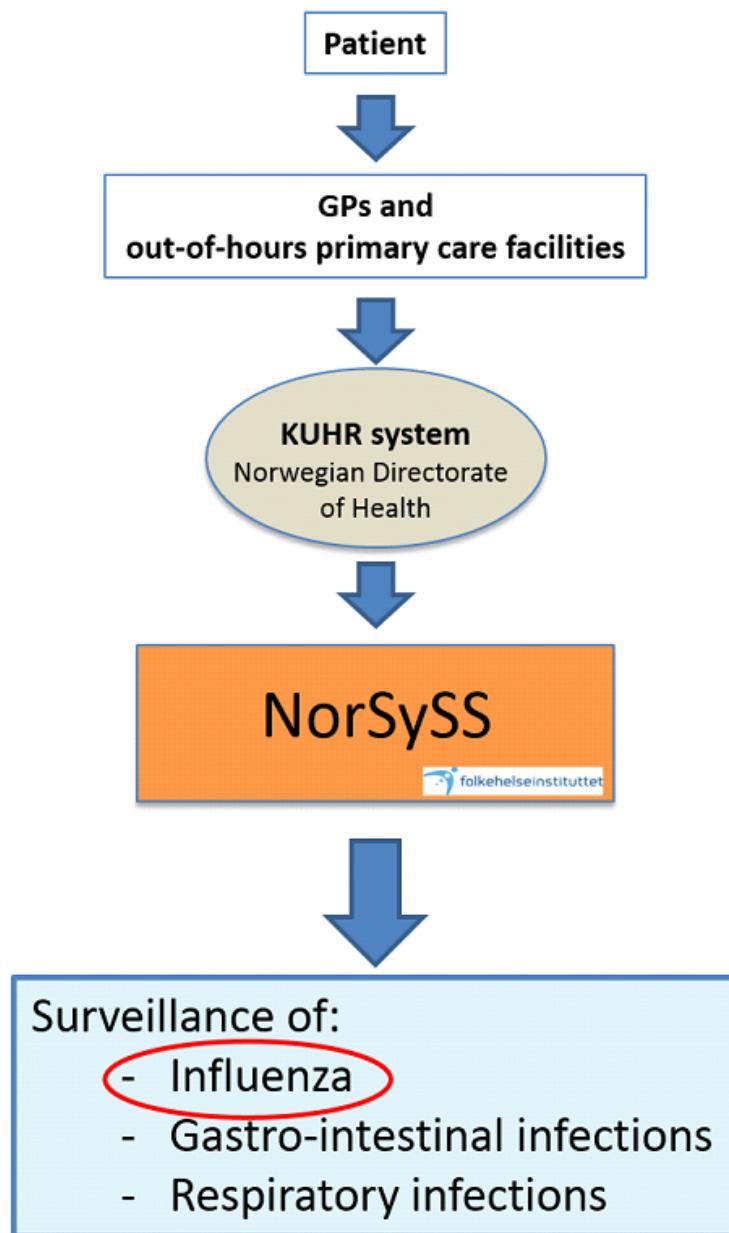
Chapter 3 marks out in detail the datasets used by this project, describes and give an explanation of relevance, challenges, limitation, and rewards.

Chapter 4 outlines the implementation and graphical results of the datasets used in chapter 3.

Chapter 5 shows the overall results.

Chapter 6 discusses the results.

Chapter 7 concludes the thesis, discusses constraints and possible future work as well as other suggestions.



(NIPH, 2017)

Figure 1.1: NIPH, 2017

# Chapter 2

## Related Works

This section looks at previous work in similar fields. It starts with presenting the paper that offer the idea that this thesis further explores, and then looks at past research on using Twitter and critical infrastructure data for similar tasks.

### 2.1 Spatiotemporal information from urban systems

In the novel study of "Detecting flu outbreaks based on spatiotemporal information from an urban system", which is the base idea for this thesis, Grottenberg et al. [8] outlines a design for a system for surveillance of flu outbreaks. Emphasis on the belief that real-time data flows could prove useful in both understanding social functions during disasters and crisis as well as give "... actionable intelligence for use in influenza management efforts.". The goal would be to extend the already implemented infrastructure with an approach to monitor human behaviour in trends throughout the influenza activity in hope for discrepancies detected through spatial analysis on important measurements. The borrowed figure 2.1 from his article sums up what this thesis hopes to accomplish, namely to find a correlation between different datasets and the datasets from the Norwegian public health institution (NIPH), this interference of public behaviour would become visible in essential criterion. This short read [8] is recommended as it gives a more in-depth understanding of the incentive for this thesis.

### 2.2 Spatiotemporal information from VGI

Volunteered Geographic Information (VGI) is peer-produced spatial data for use in crisis responses.

### 2.3 Twitter

A number of studies have been created on the information users on Twitter generate in providing valuable insights into the population by analysing millions of twitter messages (tweets). Researchers have studied tweets to reveal political opinions[11],

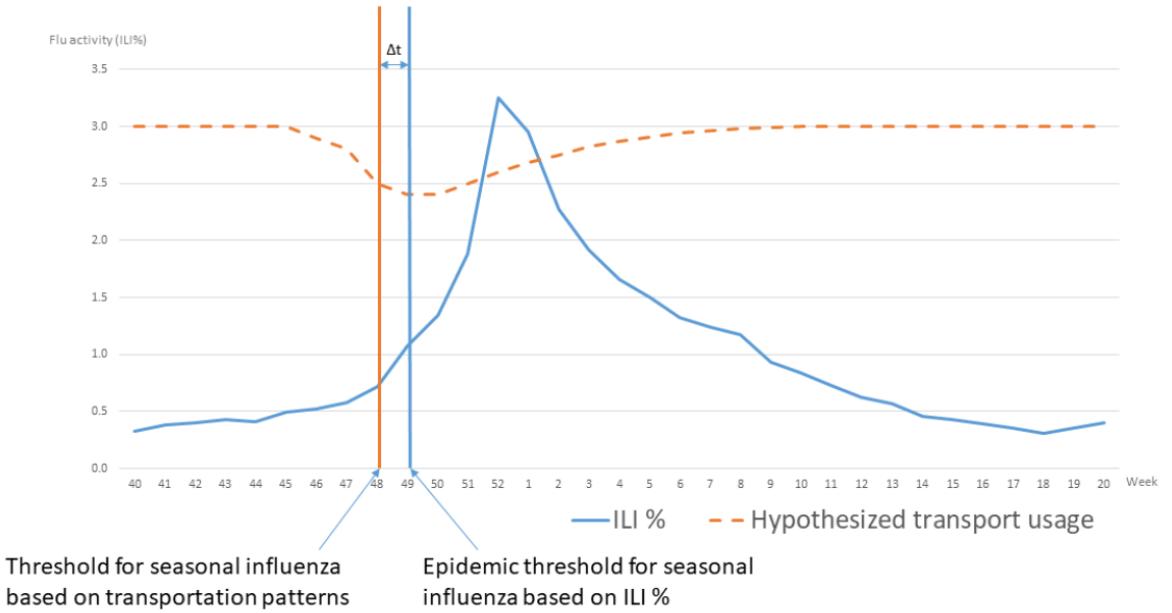


Figure 2: Theoretical correlation between weekly public transportation utilisation and flu activity (ILI %) in an urban population.

Figure 2.1: Figure from Grottenberg et al. [8]

measure public health[12], linguistic sentiments[13] and even environmental phenomena such as earthquakes[14]. Achrekar et al.[12] examines tweet flu trends and compares them with actual influenza data. The results show a high correlation between self-reported instances of flu-like illnesses (ILI) and reported ILI by public health providers. Achrekar references claims that early prevention limits the spread of infectious diseases and that twitter data is an 'untapped data source' that actually is quite reliable. This demonstrates how social media can be used to predict real-world consequences, and gives credibility to usage in this thesis.

Michal J. Paul and Mark Dredze [15] also conducted research on the usage of twitter data to measure population characteristics. In their conclusion twitter data from many users divulges reliable information about a certain topic of interest and in particular public health. They further discuss the pros and cons namely that self-reported is low cost and rapid transmission, whereas on the other side this is a 'blind authorship, lack of source citation and presentation of opinion as fact'. Certainly twitter messages may be false on an individual level, but however when taking into account thousands or even millions of messages this seems not plausible on a bigger scale. Albuquerque et al. [16] describes how they were able to extract useful information via twitter to better acquire information about a flood phenomena in German rivers, and combining this with authoritative data for disaster management. They write that social media messages gives a valuable and useful information to manage disasters, in a way this is practically the same as asking volunteers for help. For these reasons twitter data is used in this thesis as it proves an interesting and unique source of relevant information.

## 2.4 Data management and critical infrastructure

This thesis touches upon data management and development of crisis response systems. The proposed system would act as a tool in a larger system in the development of support decision making in the event of an epidemic influenza preparedness and outbreak.

Responding to extensive crisis or disasters requires coordination between a multitude of relief agencies, and this demands the right information at the right time. A system that can detect an emergence of a possible influenza outbreak would be an aiding factor to this. Gonzales et al. [17] goes into general details of how the quality of information during a crisis response is important and how to better coordinate relief agencies with the right information at the right time. They conclude that designing a computer based system for management and automation services of a work flow information conductor would better the over all quality of response and guidance. The system proposed by this thesis could be a module of such a system.

Machine-learning algorithms may also be of use in spatiotemporal analysis of social media data for disasters and damage assessment. Resch et al [18] explains how the current management of disasters have several shortcomings that can be solved by machine-learning topic models and spatiotemporal analysis. Temporal lags and limited resolution of information prevents successful and accurate resource deployment, advantages of new approaches with real-time collecting of data, like social media and other crowdsourcing networks "can significantly improve disaster management". Resch et al proposes a new approach to analyse social media with the combination of semantic machine-learning algorithms with spatio and temporal analysis. The challenge is detecting data flow continuously without prior analysis and knowledge about the event in question. Their results show remarkable improvement to accurate event tracking and other hotspots, disaster management and valuable insight to affected regions and assets.

Simulation modules could also be added to this system. This thesis is not a simulation tool but it is worth mentioning that there are several such proposed models of influenza and other disease simulation implementations. Shao et al. [19] ask the question of whether it is possible by monitoring public urban data to predict the coming outline of an overall epidemic, and simulates this. There are many more simulation tools, another is proposed by Stein et al. [20] which models an influenza outbreak in two provinces of Lao. Simulations are a way of preparing and training in order to reveal flaws and evaluation of response plans and deployment of limited health care resources.

## 2.5 Goompy

Goompy[21] is an open Github project and provides an interactive Google static map[22] for Python, it is created by Simon D. Levy. The main program, described more in chapter 4, uses this map implementation with it's own significant modifications to serve an interactive Google based map solution in order to provide visualization of information. The core Goompy file is found in the file /Frontend/goompy/\_init\_\_.py. This was heavily edited to provide the necessary functions of this thesis. The edit includes: Multithreading the fetching of Google static map images thus making Goompy about 4 times faster, dragging now changes latitude

and longitude based on x and y position of the map to better help zooming functions, having the API key fetched from a separate text file in order to hide this from misuse by other developers, support of optional map coordinates to be plotted directly in the Google static map API, using and drawing a list of coordinates as a diamond-shaped polygon with individual colors and sizes and using the mouse wheel to zoom in and out. Goompy requires a Google static map API key in order to work properly, users are asked to create the file Frontend/api\_keys.txt and paste the key there as described by the file Frontend/README.md. The original project saved the Google map images in a cache so that fetching a specific map with a familiar geolocations would be instantaneous instead of fetching them again from the Google server, this however was a violation of the terms of agreement and that function was removed from this thesis. Caching resources is a good way to quickly get often used functions, although the new implementation changes latitude and longitude often, as it allows this change, this is no longer a good strategy. For these two reasons the caching was removed. Figure 2.2 shows the Goompy map interface. In the top left corner radiobuttons change the current viewing map type. The buttons to zoom in and out are found in the bottom right corner.

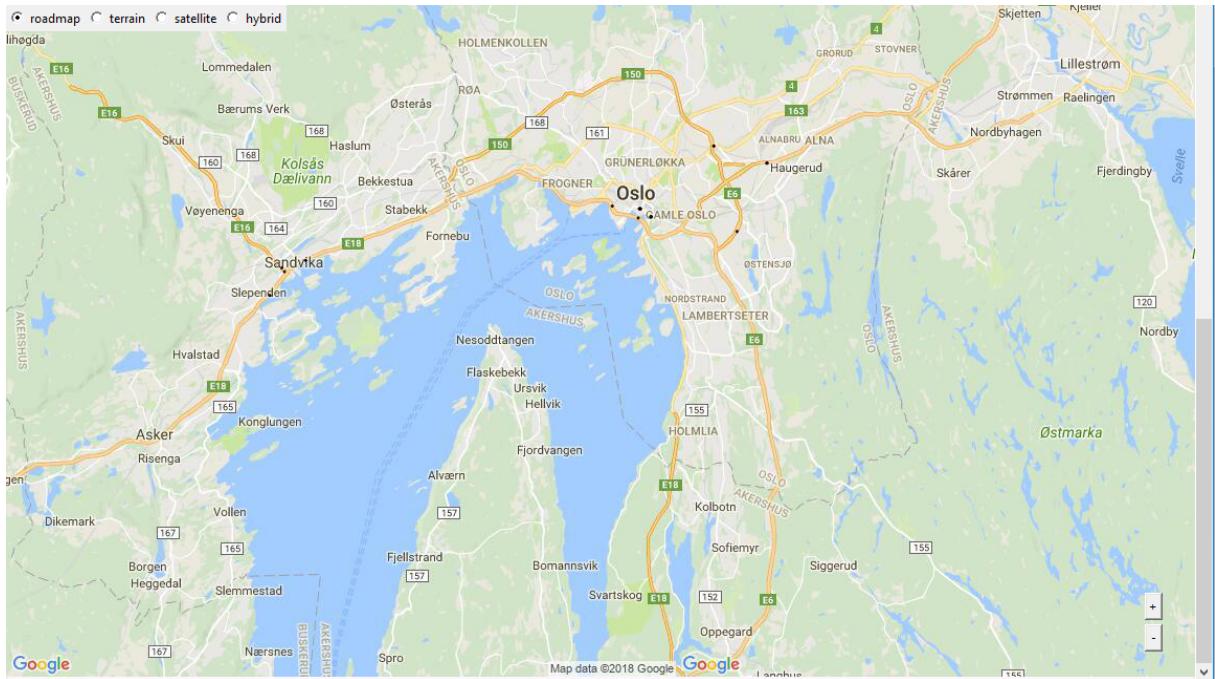


Figure 2.2: A Goompy implementation of Google’s static map API

# Chapter 3

## Datasets used

In this chapter, the different datasets used will be introduced. The goal of this thesis is to use as many datasets possible and then later evaluate them according to relevant results.

### 3.1 The Norwegian Institute of Public Health

The Norwegian Institute of Public Health (NIPH) have weekly updates[23] on the development of the current influenza season as well as previous ones. The reports include numbers of diagnoses from general practitioners (GPs) considering influenza-like illness (ILI), and hospitalized virus observations. These are the main focus and acts as a baseline for other datasets to compare against. The virus observation numbers are included in the report, ILI symptoms are not, they are however both included in graphs. Upon further request, the ILI data was provided for the season of 2016/2017, and for the cities of Oslo and Bergen of the season of 2015/2016, 2016/2017 and 2017/2018. Exact numbers of the virus observations are only included for the three last years, therefore this thesis only uses the seasons of the years 2015/2016, 2016/2017 and 2017/2018. The reports also cover what kind of influenza viruses are circulating in the country and where, vaccine status and recommendations, as well as the overall prognosis of the current season. GPs report ILI based on these characteristics: muscle pain, coughing, fever and the feeling of being sick. The ILI numbers are perhaps of more interest since they are more accessible than virus observations that only counts for hospitalization. These two datasets provide the measurement basis other datasets are held up against.

### 3.2 The Norwegian Public Roads Administration

The Norwegian Public Roads Administration (NPRA) have several different collections of data available for a number of different purposes [24]. The main motivation for traffical data in this thesis is the hypothesis that when people are ill they commute less and thus this shows when surveying statistical details. Freely on their website [24] there are a few interesting options. They have traffic information in the standard traffic management exchange data structure (DATEX), application programming interfaces (API), statistics in an extensible markup language (XML) and traffic index data relevant to the years before. It is important for this thesis

that the data collected is on a weekly basis at least in order to compare it to the influenza data. It turned out that the data on their website did not suffice for this purpose, they only had a temporal resolution of months or years while this thesis needs a temporal resolution of weeks or better. The data given contained a set of traffic registration stations throughout Norway. Data provided was on a weekly basis and also on an hourly basis for a subset of the original traffic registration stations provided. With this statistics of the traffic amount and spatial bounds can be derived showing the possible correlation influenza can have on commuting traffic. The regions of interest are the whole of Norway and the three cities of Stavanger, Bergen, and Oslo.

### 3.3 Twitter

The reason twitter data is interesting is that it contains self-reported instances of influenza on an individual level. These self-reported cases may even occur without the patient visiting a doctor, and so capture otherwise non-reported instances of ILI. The advantages are an instant notification about possible ILI and its spread, against the disadvantages of it being self-reported and thus somewhat unreliable. Twitter has several APIs available for public use, the one used in this project is the representational state transfer (REST) API or 'search API' which allows for searching against a set of keywords. The REST API is limited though, data accessible is roughly only maximum 10 days old and the search limit is on a maximum of one hundred messages called 'tweets'. The other API of interest is the stream API which continually gets the latest tweets. In order to only get Norwegian tweets, a set of geographical locations needs to be defined. The reason the stream API was not used is firstly that it requires a computer running on the internet continuously in order to get all the desired tweets. Secondly, the data collected could become large slowing down other post-processing algorithms and taking up unnecessary storage. Lastly, the stream API only provides a small set of the actual tweets tweeted, this means when searching for a specific term using the stream API some relevant tweets could go unnoticed and thus a search API is more appropriate for this task.

### 3.4 Kolumbus

Kolumbus is the public transportation administration in the state of Rogaland in Norway, this includes Stavanger, a city of interest. Unfortunately, Kolumbus provides no API, but on further request data of monthly passenger travel was provided from the years of 2015-2017.

### 3.5 Ruter

Ruter is the public transportation administration in the state of Oslo in Norway. Unfortunately, Ruter's API does not include passenger or tickets sold information, this was however provided on request for the years 2015, 2016, 2017 and up till 27 of February for the year 2018 on a daily basis.

# Chapter 4

## Implementation

This chapter describes how the use of the different datasets were implemented and presented. The program is divided into two: The backend and the frontend. The structure and functions are provided by the backend, which governs collection and manipulation of data, and the frontend presents the data in a graphical user interface (GUI) using graphs and maps.

### 4.1 The Backend

The backend is responsible for providing the frontend all the data and deeper functions it needs to visualize and administrate data to be show in graphs. The backend is partitioned into modules based on each dataset available. Each module may also be run individually for testing and easy viewing purposes. The Twitter module is unique as it requires 4 application programming interface (API) keys to work properly. The instructions for this set-up is found in the file README.md in the twitter module's directory.

#### 4.1.1 The Norwegian Institute of Public Health

There are two different sets of data, which is divided into the separate modules of NIPH.ILS.py and NIPH\_virus\_detections.py located in the same directory, and they show influenza-like illnesses (ILI) and hospitalized viral observations. They both extract data and then draw a graph using Python's matplotlib library, the graphs can be seen by running the modules individually or in the frontend main program frontend/gui.py's appropriate viewport. Figure 4.1 show the three last seasons of influenza in regards to observed viral infections. The plotting was done manually as NIPH only provides viral observational data in reports that are in pdf files on their official website[23].

Figure 4.2 shows the influenza-like illnesses (ILI) of the year 2016/2017. This was not done manually as data was provided in a simple .xlsx file which was read using Python's openpyxl module, processed and then drawn as a graph.

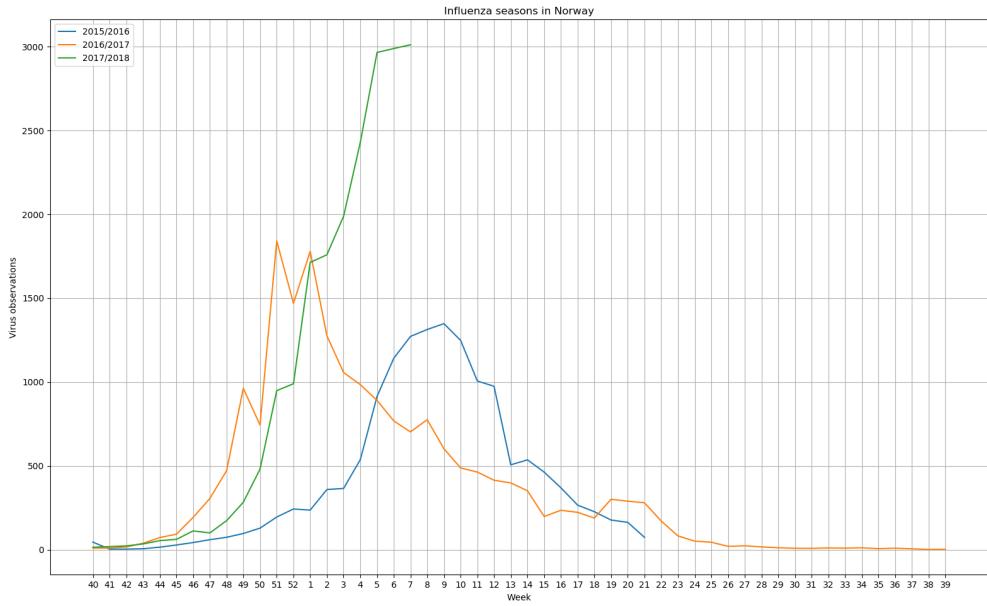


Figure 4.1: Influenza virus observation

### 4.1.2 The Norwegian Public Roads Administration

From the .xlsx files provided by the NPRA, simple graphs were created in python showing the total annual traffic on Norwegian roads from 2002 to 2015 on a monthly basis as seen in figure 4.3.

Also derived from this dataset is the annual traffic of the two cities of Bergen and Oslo, which are cities of interest.

The dataset is in an XML file structure, a module named NPRA\_monthly.py was created that reads through all rows and collects the relevant columns into an array using Python's openpyxl module and then draws a graph using Python's matplotlib module. For the annual graph, every month of every year was collected. For the towns of Bergen and Oslo the correct roads were identified and then every year of every month of those roads was collected, loaded into an array and then drawn as a graph. The separate text files 'Bergen places.txt' and 'Oslo places.txt' is to make it easy to edit should these roads change in the future. This module when run individually accepts one command argument from the user, either cities of Oslo or Bergen may be provided to specify interest, if no argument is given the annual graph will show. The problem of using these datasets is that the data is an average calculation of monthly traffic, meaning the temporal bounds are too coarse for comparison against the influenza data which in turn is on a weekly basis. For these reasons no figures of this dataset are shown in this thesis, they are however available as modules and in the frontend's main program in the programming project.

For the weekly datasets a set of traffic registration stations was needed to define the temporal bounds of each area of interest. Defined are the towns of Oslo, Stavanger, and Bergen, as well as the whole of Norway on a level 1 basis. The level 1 registrations are continuous throughout the year on an hourly basis and is exactly what this thesis requires. The module NPRA\_weekly.py captures these functions

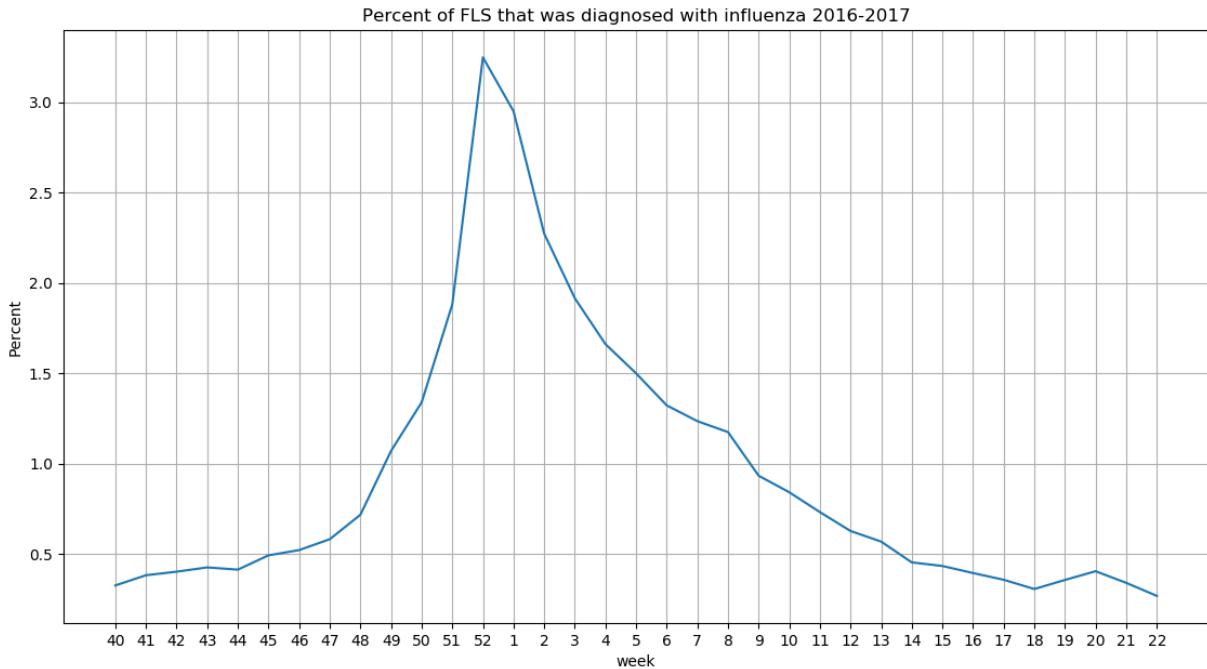


Figure 4.2: Influenza-like illnesses season 2016/2017

and also provides the user with command arguments if run individually. The commands are the cities of Bergen, Stavanger or Oslo, if no commands are given the annual graph of the whole of Norway will be drawn instead.

Figures 4.6, 4.7 and 4.8 shows the traffic on a weekly basis. This provides a better resolution for better analysis.

Figure 4.9, 4.10 and 4.11 shows the different geospatial bounds used to define the cities. The green circles with numbers inside show where and how many traffic registration stations there are.

The last NPRA dataset acquired was raw hourly data from a defined subset of all of NPRA's traffic registration stations previously used. The data contains all whole hours from all weeks over several years, number of fields available on the road (usually only two for regular roads), and how many vehicles passed by that hour and also their lengths in category. Figure 4.12, 4.13 and 4.14 shows the different geospatial hourly based bounds used. There are two modules dedicated to the hourly datasets, the NPRA\_Traffic\_Stations\_Graph.py and the NPRA\_Traffic\_Stations\_load\_data.py. The graph module is responsible for drawing a graph with specifications of hour to/from, weekday to/from, month to/from, year and field. The load data module is responsible for providing the graph with all the functions it needs to operate, like querying the dataset, the variance of the queried dataset, extracting the dataset from file and organizing it into a data structure, and reading and handling the coordinates of the traffic registration stations so that it can be shown on the map. These last hourly based datasets provide high quality information and is presented in the GUI where the user can try different queries to find different information, more explained in the frontend section of this chapter.



Figure 4.3: Annual traffic 2002-2015

#### 4.1.3 Twitter

Using the representational state transfer (REST) application programming interface (API) it was paramount that in order to build a sufficient dataset, acquiring and collecting data had to begin as soon as possible in order to collect enough data for this thesis. A simple python program was created that takes the input of the API keys provided by the file keys.txt and the keywords to be searched upon provided by the file search\_terms.txt. The program ensures that no duplicate messages are recorded, and the limit of a hundred tweets dictated by the REST API was overcome simply by searching for yet another hundred from the last date of the previous hundred until the date limit of about 10 days was reached. The output is appended to a file in this data structure on new lines: id, date, location, tweet, there is also a dotted separator for each new tweet making it more easy for humans to read. The functions described are implemented by the file twitter\_searching.py, which can be run as its own module and saves new tweets to the file twitter\_data.txt.

A straightforward analysis tool for the Twitter data in the file twitter\_data.txt was created by simply counting how many tweets there are. The idea is that during influenza seasons numbers of influenza-related tweets increases and then decrease when off the season, while the number of non-relevant tweets is constant during the whole year (or slightly increasing or decreasing based on the popularity of Twitter as a social media). A more complex tool for analyzing the tweets for relevance was elected to be too much work for this thesis. The advantage of simply counting how many possible tweets there are is that it is fast and easy to implement, the drawback is that it captures non-relevant tweets. Future work may be done to improve this quality with a better analyzing tool than this thesis chose. Figure 4.15 shows the results of the time-frame captured. The analyzing function is implemented in the file twitter\_analyzer.py, when the module is executed on its own it shows a graph

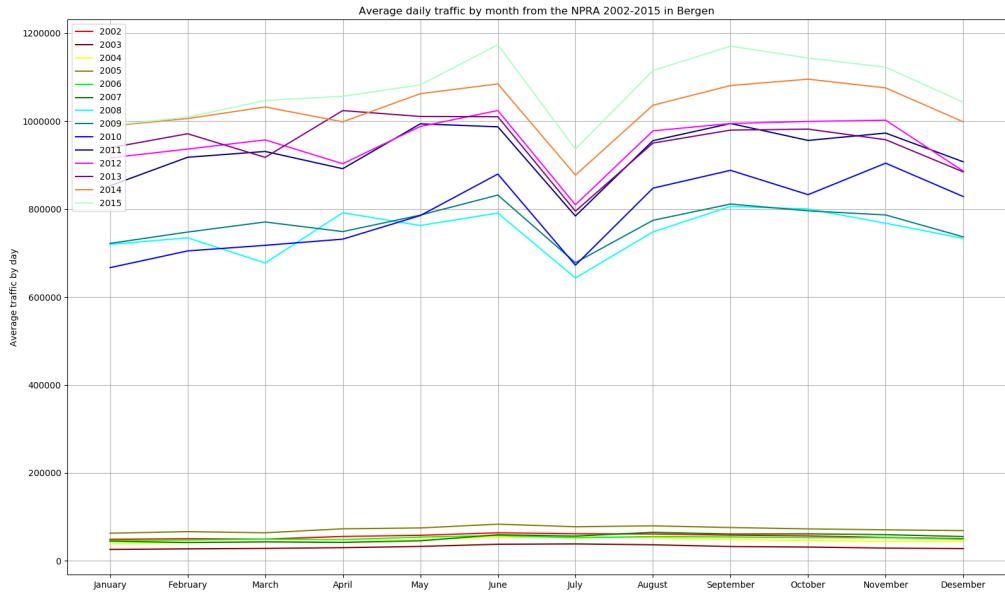


Figure 4.4: Bergen traffic 2002-2015

over the data found in the file `twitter_data.txt`. A simple batch file `twitter.bat` was created to make it easy running these programs in the desired order.

#### 4.1.4 Kolumbus

The data provided by Kolumbus was in a `.png` format needed to be converted into a more convenient (and appropriate) data structure. The chosen data structure conversion was comma separated values (CSV) stored in the file '`'15_17_månedstall_total.csv`'. From there it was a simple job to plot the data in a python script, unfortunately the data is only on a monthly basis. Figure 4.16 shows the results.

#### 4.1.5 Ruter

The data provided by Ruter was in a `.xlsx` file and could easily be read, extracted and plotted by a simple python script. Figure 4.17 shows the results. Note that with Python's `matplotlib` module a user can zoom in and out to get a more desired and uncluttered view. The data was provided by a daily basis for the years of 2015-2018. Note that the first year is lower because it does not contain Oslo's underground train service passenger data.

## 4.2 The Frontend

The thesis's program is divided into two: The backend and the frontend. The frontend is responsible for visualizing the data provided by the backend. It does so by mounting a graphical user interface (GUI) that provides everything the user needs

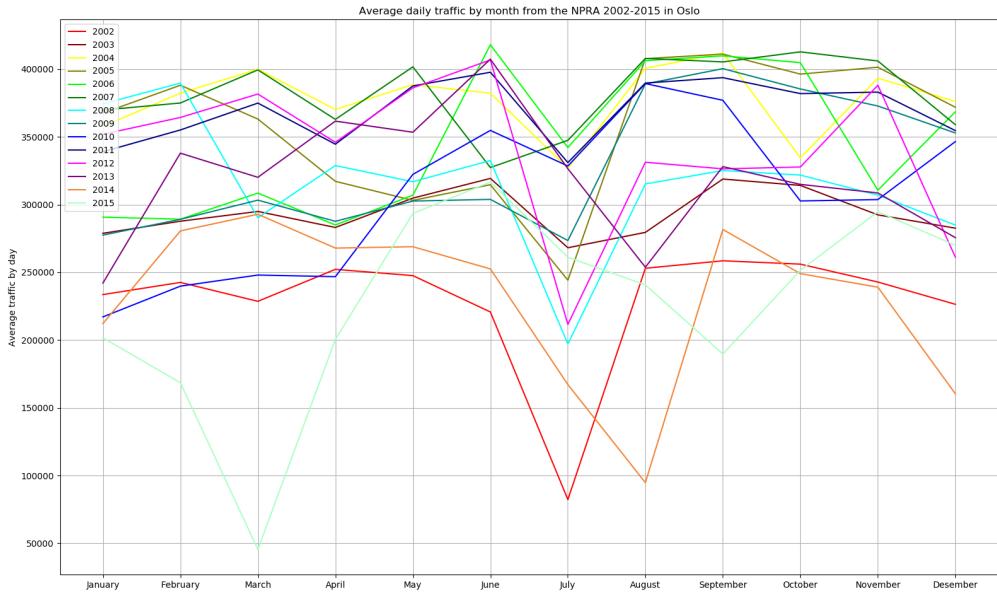


Figure 4.5: Oslo traffic 2002-2015

from this thesis. The GUI uses other frontend modules described in the following subchapters.

### 4.2.1 The GUI

The file `gui.py` is the main program. It mounts the GUI with help from backend modules and the frontend modules such as the file `map_canvas.py`, the file `scrframe.py`, the file `double_y_graphs.py`, the file `NIPH_frame.py` and the file `NPRA_frame.py`. The GUI is created using Python's standard Tkinter module, and it provides the means of a basic window creation with all the other usual GUI necessities available.

The GUI module itself is structured in two parts: The buttons frame and the data frame. The buttons frame produces a menu and simply makes available buttons to be clicked upon showing the different graphs for the respective datasets from the backend. The data frames show the graphs and if needed a map, visualizing the data from the backend. The backend takes time to load, to make this experience more user-friendly a progress bar is shown progressing relative to the loading sequence. Upon completion, the NPIH data is shown as a standard view. The user may use the mouse wheel to scroll up and down the view and click the buttons to change datasets.

In some datasets, a map is provided for further visualization. the map is interactive with its own buttons and also responds to dragging the mouse in order to move the map, double-clicking in order to zoom in and using the mouse wheel, when hovering over the map, to zoom in and out.

Figure 4.18 shows the GUI.

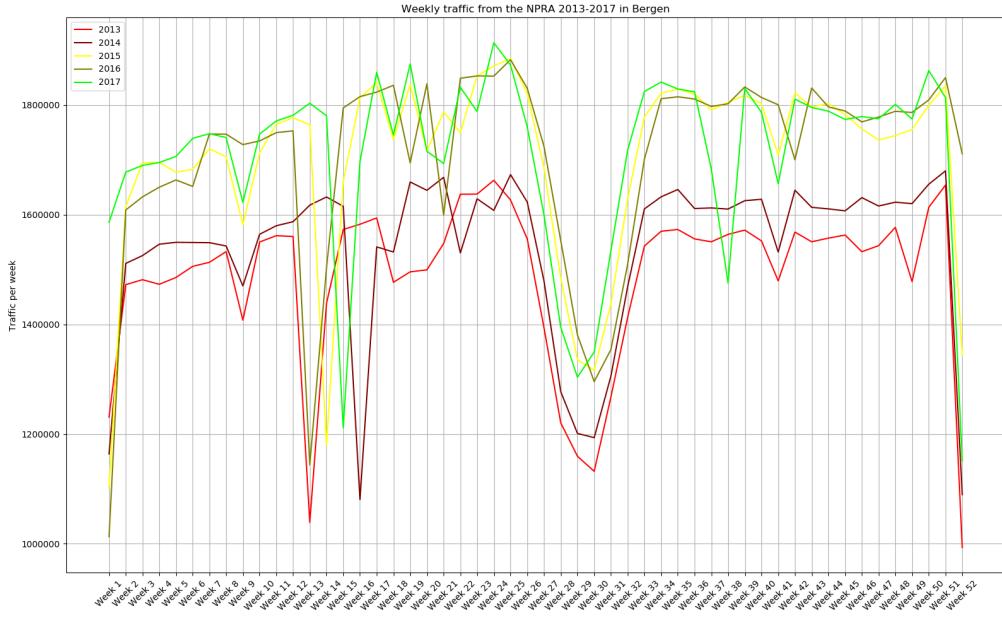


Figure 4.6: Weekly data of the city of Bergen

### 4.2.2 The Map

The file `map_canvas.py` provides the GUI a Goompy[21] map on a Tkinter canvas, as described in chapter 2. This file is also from the Goompy project, but is heavily modified to serve the purpose of this thesis. The file launches a Google static API map on a Tkinter canvas and provides basic Google map functions and user input. The functions edited for this thesis is: better zooming capabilities, coordination markers with individual colors and sizes, ability to focus on the map by will and some other minor bug fixes.

### 4.2.3 The Scrollbar

Creating a functional scrollbar that responds to mouse click and mouse wheel events in Tkinter proved difficult, which is why Eugene Bakin's Tkinter scrollable[25] frame was used. It is an open Github project. The file `Frontend/scrframe.py` contains his code with minor edits in order to be able to scroll with the mouse wheel, get the Tkinter focus, resetting scrollbar viewport and better resizing of the window. This module may also be run independently for testing purposes.

### 4.2.4 NIPH dataframe

The GUI module is structured in two parts: The buttons frame and the data frame, data frames visualize information from the backend. The NIPH data frame was further extended with the functionalities that allows for comparison of the NIPH data with all the other datasets at the different influenza seasons available. The frontend module `NIPH_frame.py` was created to be implemented by the main file `GUI.py` and the file `double_y-graphs.py` provides the necessary supportive algorithms. Both files

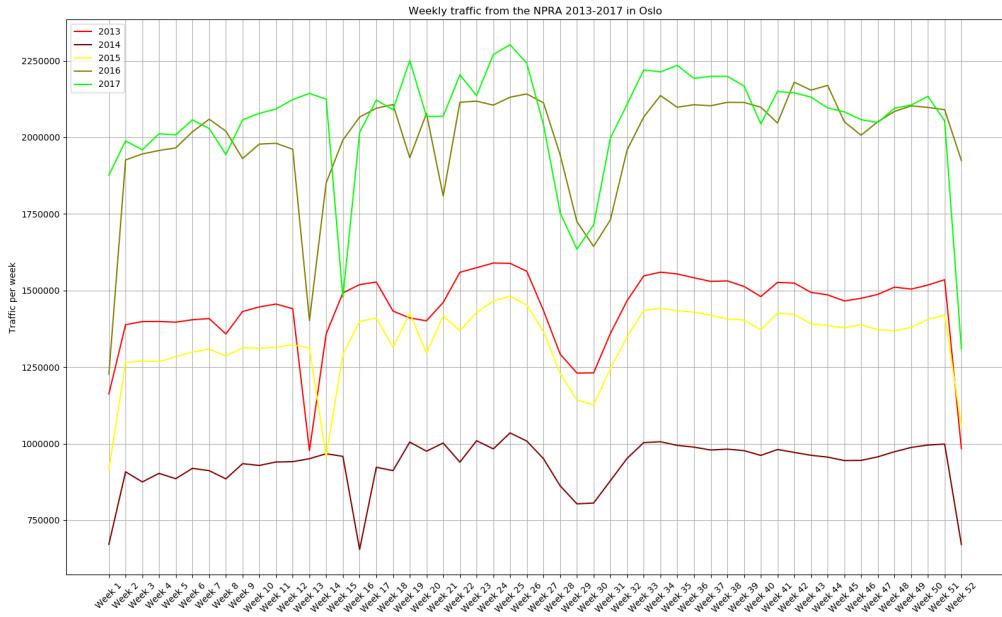


Figure 4.7: Weekly data of the city of Oslo

may be run individually for testing purposes. The comparison functions work in the way that the user selects a dataset to compare with by clicking a button in the top border. A drop-down menu will be produced giving the choices of cities and influenza seasons. Two graphs will then be drawn sharing the same x-axis but having different y-axes. This makes for easy comparison and querying the data in order to find possible correlations. Figure 4.19 shows the NIPH comparing buttons panel.

#### 4.2.5 NPRA dataframe

In addition to the monthly and weekly datasets the hourly are presented in its own GUI module implemented by the main file GUI.py. The hourly datasets contains 58 different traffic registration stations from the cities of Bergen, Stavanger and Oslo and can be queried with a buttons-panel. The dropdown buttons provide the choices of hours to/from, weekday to/from and month to/from from the year of 2017, lastly there is a show button which initiates the query. On the left border a map is shown. Figure 4.20 shows the NPRA query buttons panel.

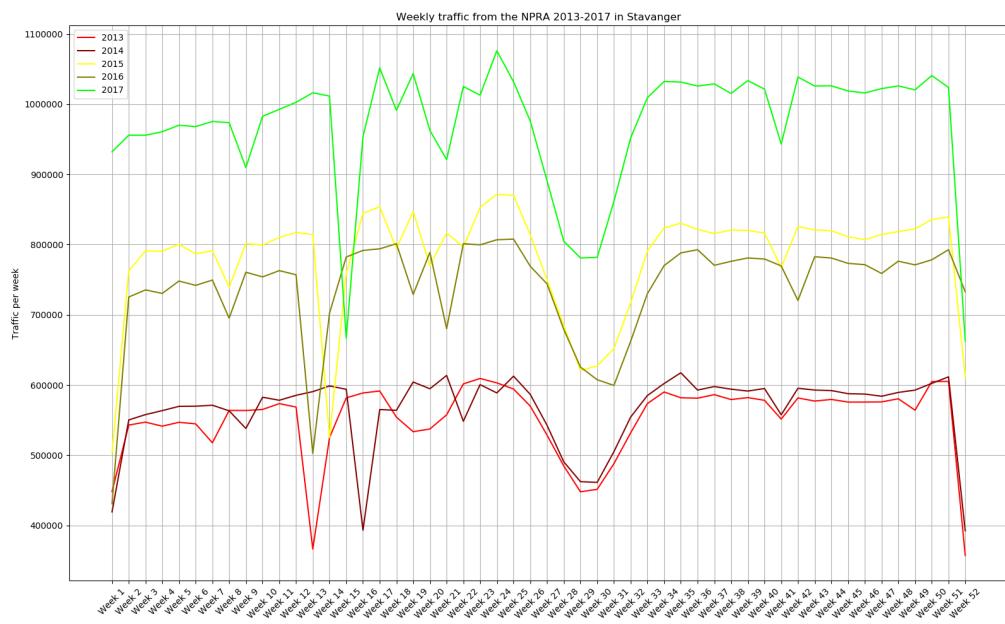


Figure 4.8: Weekly data of the city of Stavanger

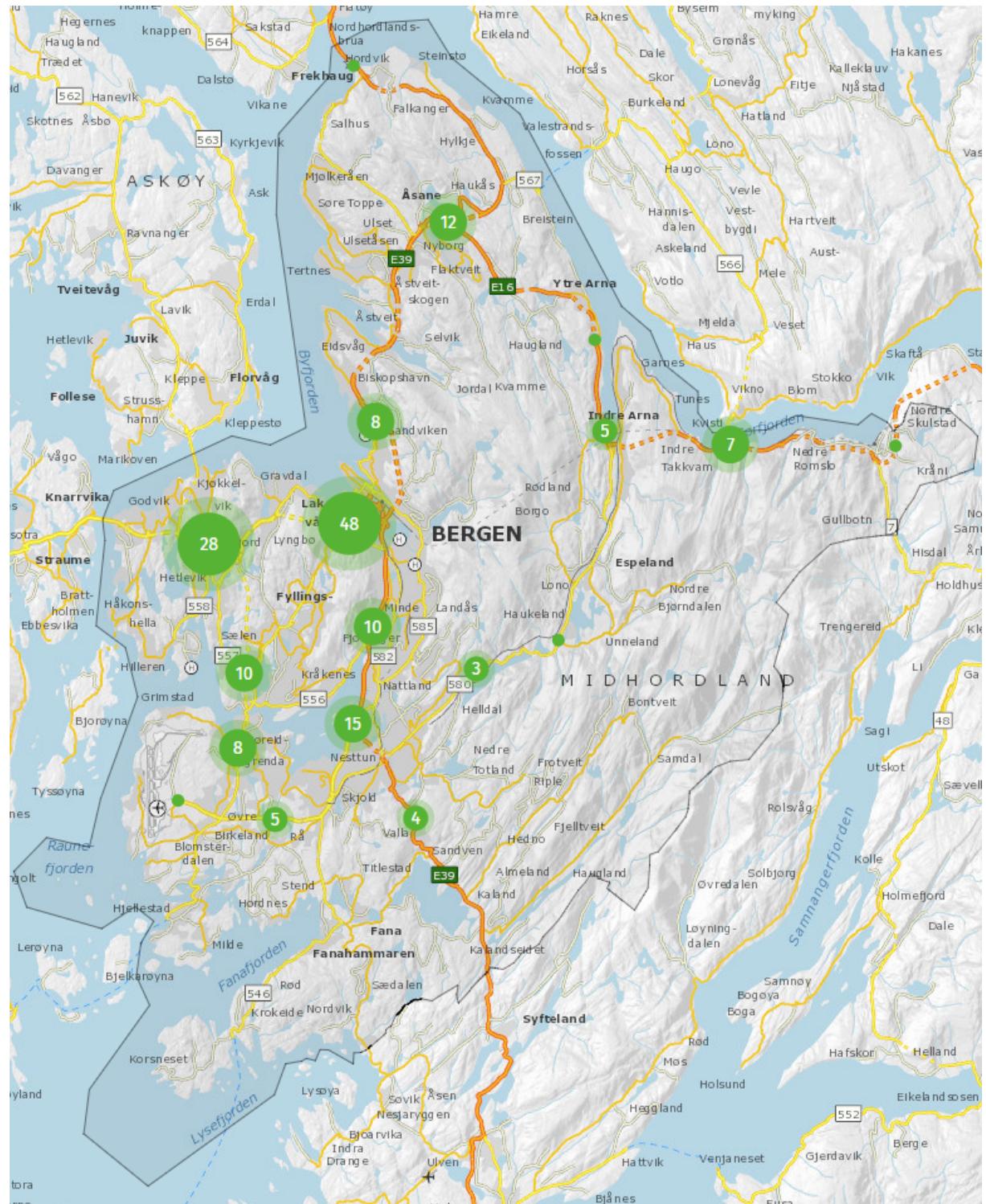


Figure 4.9: Geospatial bounds of Bergen. The green circles show where the traffic registration stations are, and the number reveals how many there are in that general area.

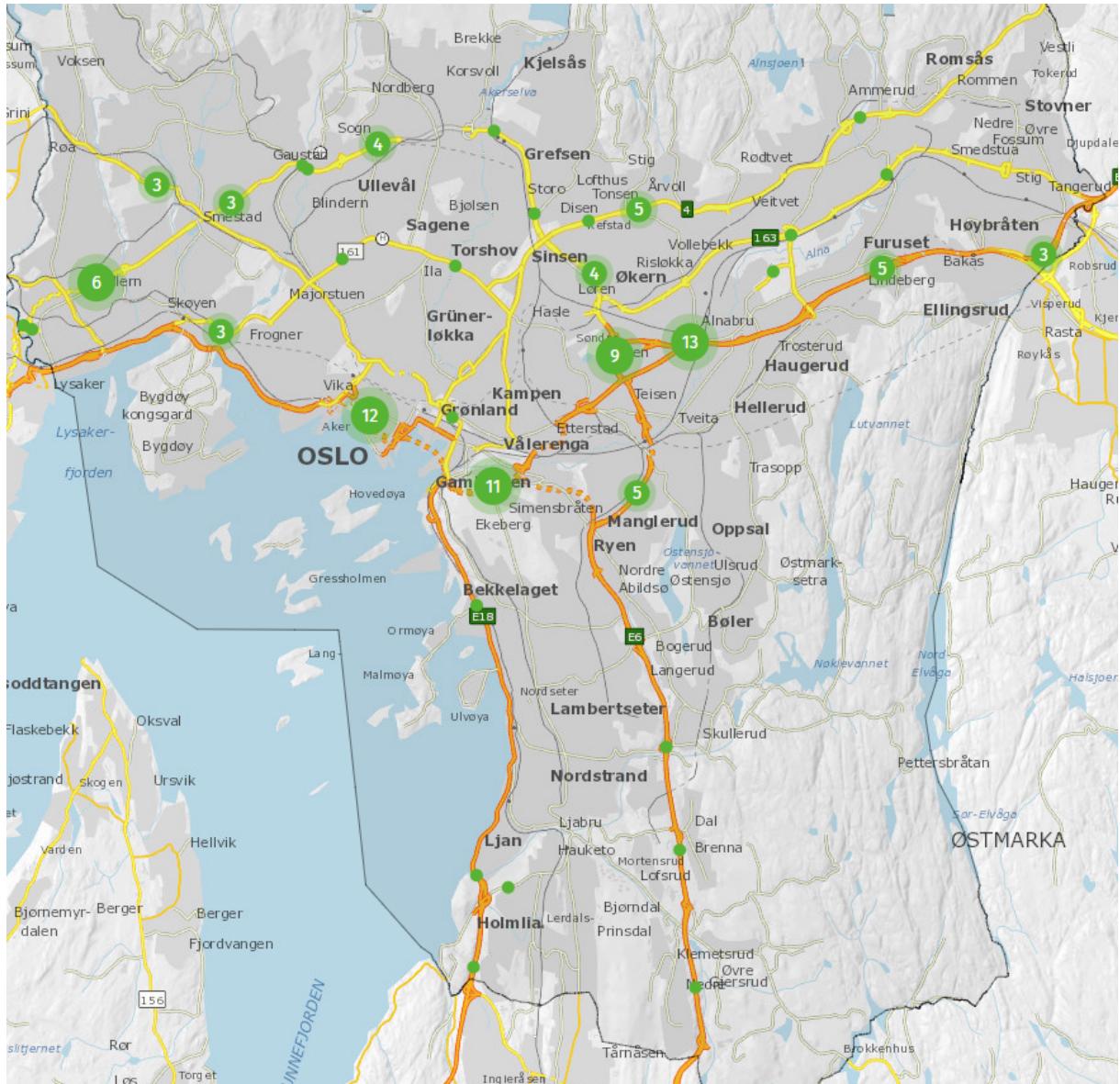


Figure 4.10: Geospatial bounds of Oslo. The green circles show where the traffic registration stations are, and the number reveals how many there are in that general area.



Figure 4.11: Geospatial bounds of Stavanger. The green circles show where the traffic registration stations are, and the number reveals how many there are in that general area.



Figure 4.12: Geospatial hourly bounds of Bergen



Figure 4.13: Geospatial hourly bounds of Oslo

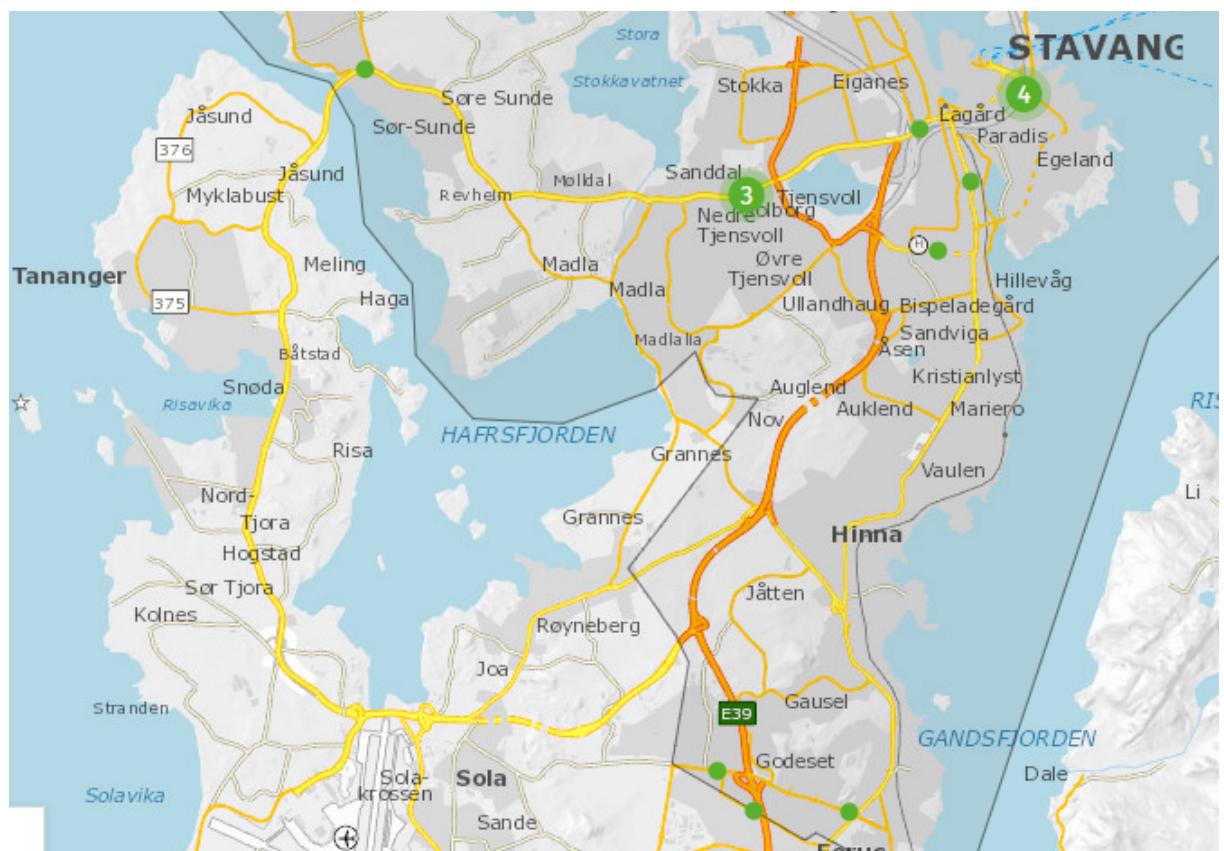


Figure 4.14: Geospatial hourly bounds of Stavanger

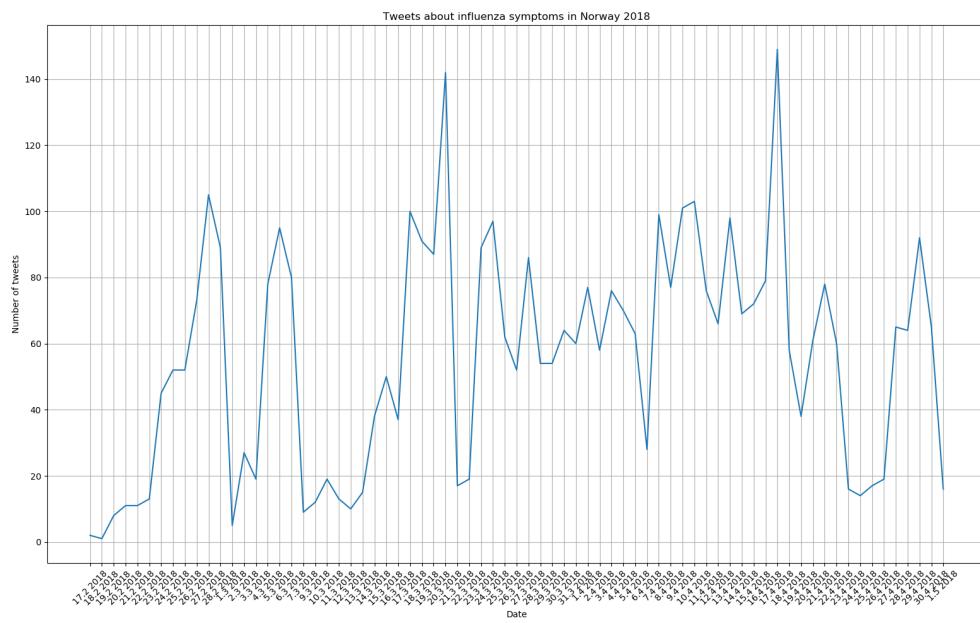


Figure 4.15: Tweets concerning ILS of 2018

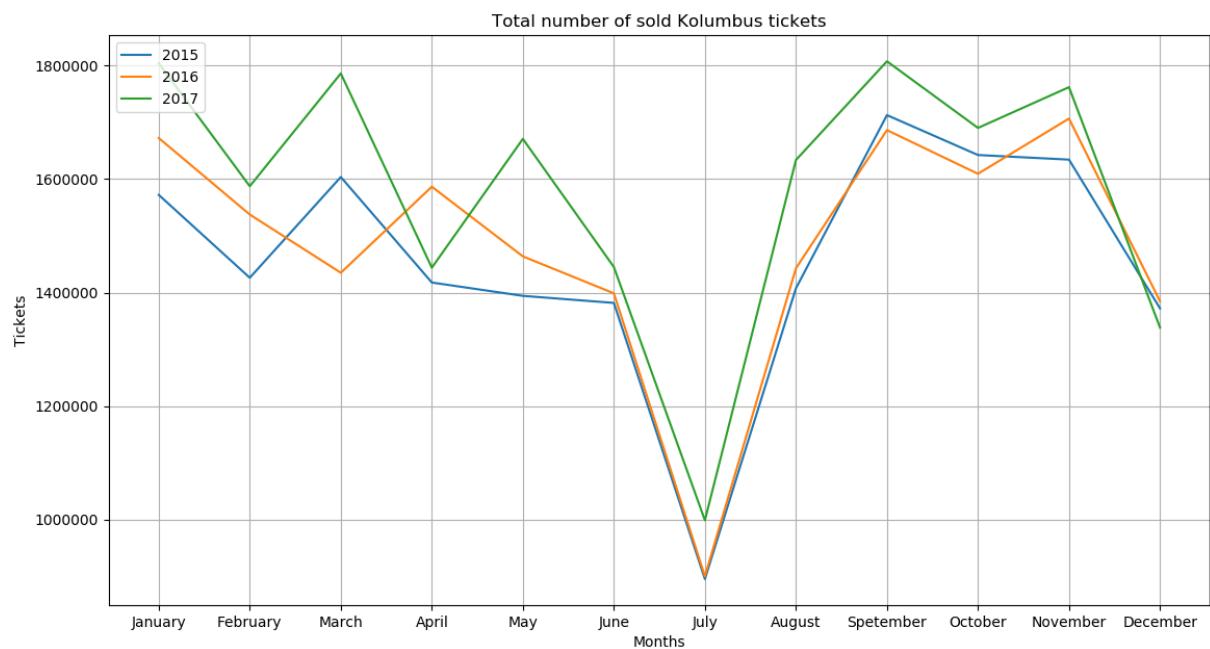


Figure 4.16: Monthly passenger travel with Kolumbus

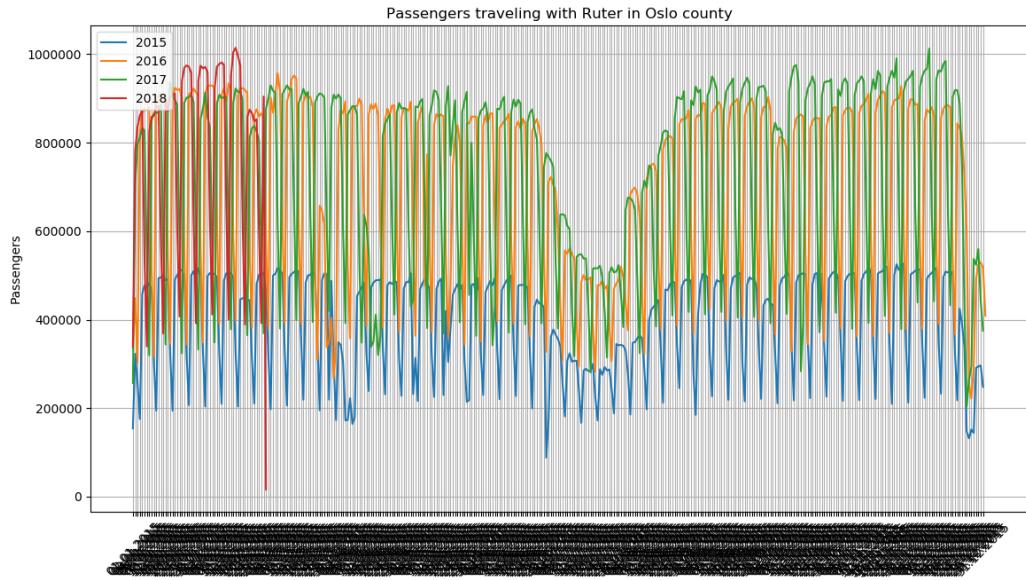


Figure 4.17: Daily tickets sold with Ruter, the year of 2015 does not contain Oslo's underground train service passenger data

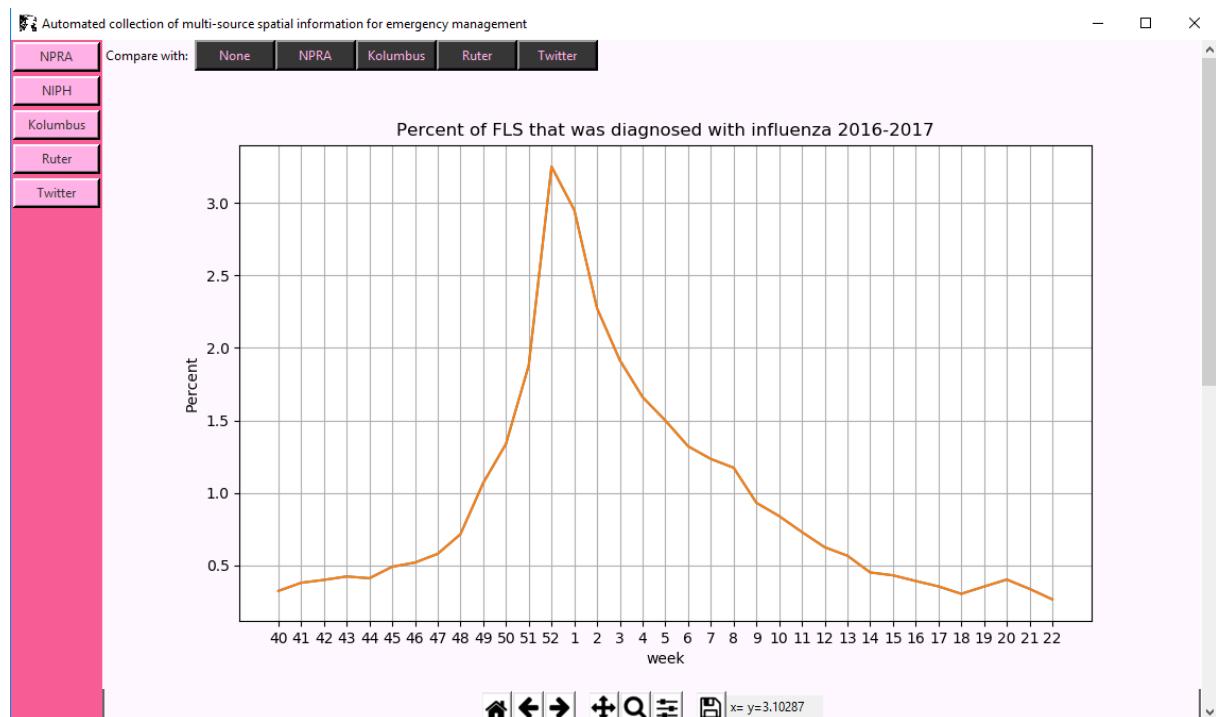


Figure 4.18: The GUI

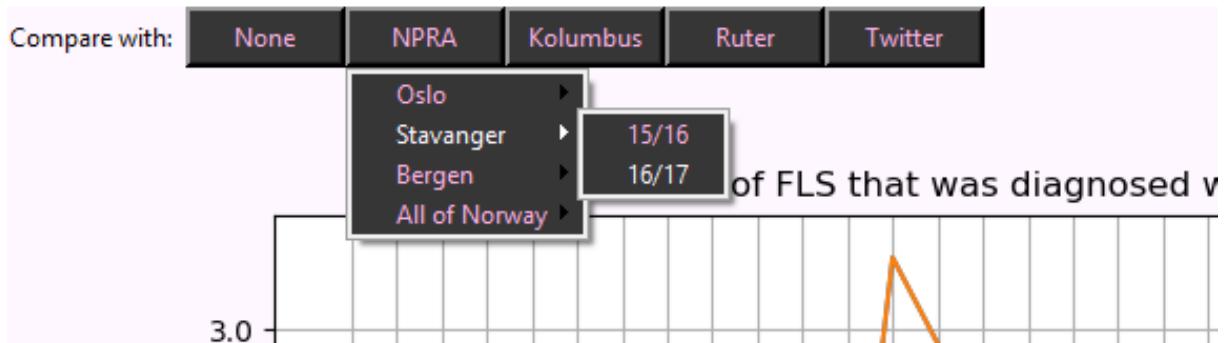


Figure 4.19: NIPH comparing buttons panel

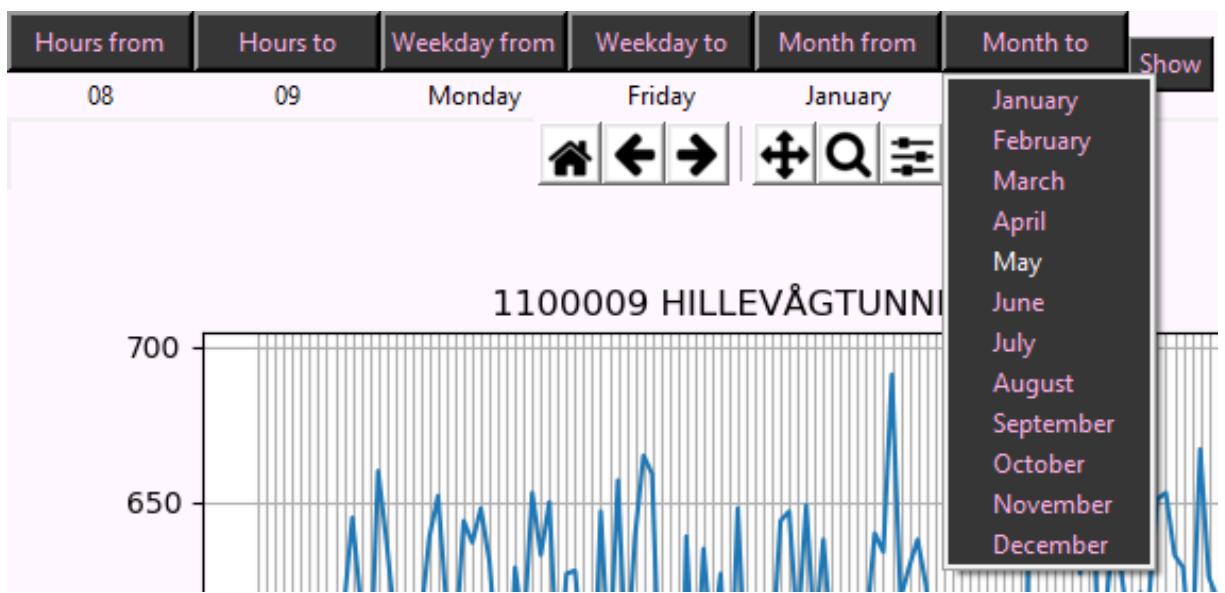


Figure 4.20: NPRA query buttons panel

# **Chapter 5**

## **Results**

### **5.1 TODO**

#### **5.1.1 TODO**

tenk over avvik i grafen, hva skyldes dette? veg.arb? ferie?

#### **5.1.2 Twitter**

duplicater, count, nevn nårtid folk vanligvis poster på sos. media og hvordan det stemmer med egne data.

# **Chapter 6**

## **Discussion**

### **6.1 TODO**

#### **6.1.1 workflow**

vi endret retning tidlig i prosjektet. ikke god nok norsk api infrastruktur. fra automasjon til manuell henting av data

# **Chapter 7**

## **Conclusion**

### **7.1 TODO**

#### **7.1.1 Future works**

svakheter, hvordan gjøre bedre? hva er mitt bidrag?

# **Appendix A**

## **Appendix Title**

# Bibliography

- [1] “Pandemic influenza risk management: A who guide to inform and harmonize national and international pandemic preparedness and response.” [http://www.who.int/influenza/preparedness/pandemic/influenza\\_risk\\_management\\_update201](http://www.who.int/influenza/preparedness/pandemic/influenza_risk_management_update201) Accessed: 2018-05-23.
- [2] A. D. Iuliano, K. M. Roguski, H. H. Chang, D. J. Muscatello, R. Palekar, S. Tempia, C. Cohen, J. M. Gran, D. Schanzer, B. J. Cowling, *et al.*, “Estimates of global seasonal influenza-associated respiratory mortality: a modelling study,” *The Lancet*, 2017.
- [3] “The norwegian institute of public health.” <https://www.fhi.no/en/>. Accessed: 2018-06-11.
- [4] C. Poletto, M. Tizzoni, and V. Colizza, “Human mobility and time spent at destination: impact on spatial epidemic spreading,” *Journal of theoretical biology*, vol. 338, pp. 41–58, 2013.
- [5] T. Wibisono, D. M. Aleman, and B. Schwartz, “A non-homogeneous approach to simulating the spread of disease in a pandemic outbreak,” in *Simulation Conference, 2008. WSC 2008. Winter*, pp. 2941–2941, IEEE, 2008.
- [6] B. Yang, H. Pei, H. Chen, J. Liu, and S. Xia, “Characterizing and discovering spatiotemporal social contact patterns for healthcare,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 8, pp. 1532–1546, 2017.
- [7] C. Robertson, “Towards a geocomputational landscape epidemiology: surveillance, modelling, and interventions,” *GeoJournal*, vol. 82, no. 2, pp. 397–414, 2017.
- [8] L. O. Grottenberg, O. Njå, E. Tøssebro, G. Braut, R. Tønnesen, and G. M. Grøneng, “Detecting flu outbreaks based on spatiotemporal information from urban systems – designing a novel study,” *Icwsm*, vol. 20, pp. 1–7, 2017.
- [9] M. K. Enduri and S. Jolad, “Dynamics of dengue disease with human and vector mobility,” *Spatial and spatio-temporal epidemiology*, vol. 25, pp. 57–66, 2018.
- [10] “The norwegian institute of public health: About the norwegian syndromic surveillance system.” <https://www.fhi.no/en/hn/statistics/NorSySS/about-the-norwegian-syndromic-surveillance-system/>. Accessed: 2018-06-11.

- [11] S. B. Elson, D. Yeung, P. Roshan, S. R. Bohandy, and A. Nader, *Using social media to gauge Iranian public opinion and mood after the 2009 election*. Rand Corporation, 2012.
- [12] H. Achrekar, A. Gandhe, R. Lazarus, S.-H. Yu, and B. Liu, “Predicting flu trends using twitter data,” in *Computer Communications Workshops (INFO-COM WKSHPS), 2011 IEEE Conference on*, pp. 702–707, IEEE, 2011.
- [13] J. Eisenstein, B. O’Connor, N. A. Smith, and E. P. Xing, “A latent variable model for geographic lexical variation,” in *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pp. 1277–1287, Association for Computational Linguistics, 2010.
- [14] T. Sakaki, M. Okazaki, and Y. Matsuo, “Earthquake shakes twitter users: real-time event detection by social sensors,” in *Proceedings of the 19th international conference on World wide web*, pp. 851–860, ACM, 2010.
- [15] M. J. Paul and M. Dredze, “You are what you tweet: Analyzing twitter for public health.,” *Icwsom*, vol. 20, pp. 265–272, 2011.
- [16] J. P. De Albuquerque, B. Herfort, A. Brenning, and A. Zipf, “A geographic approach for combining social media and authoritative data towards identifying useful information for disaster management,” *International Journal of Geographical Information Science*, vol. 29, no. 4, pp. 667–689, 2015.
- [17] R. A. Gonzalez and N. Bharosa, “A framework linking information quality dimensions and coordination challenges during interagency crisis response,” in *System Sciences, 2009. HICSS’09. 42nd Hawaii International Conference on*, pp. 1–10, IEEE, 2009.
- [18] B. Resch, F. Usländer, and C. Havas, “Combining machine-learning topic models and spatiotemporal analysis of social media data for disaster footprint and damage assessment,” *Cartography and Geographic Information Science*, vol. 45, no. 4, pp. 362–376, 2018.
- [19] H. Shao, K. Hossain, H. Wu, M. Khan, A. Vullikanti, B. A. Prakash, M. Marathe, and N. Ramakrishnan, “Forecasting the flu: designing social network sensors for epidemics,” *arXiv preprint arXiv:1602.06866*, 2016.
- [20] M. L. Stein, J. W. Rudge, R. Coker, C. van der Weijden, R. Krumkamp, P. Hanvoravongchai, I. Chavez, W. Putthasri, B. Phommasack, W. Adisasmto, *et al.*, “Development of a resource modelling tool to support decision makers in pandemic influenza preparedness: The asiaflucap simulator,” *BMC public health*, vol. 12, no. 1, p. 870, 2012.
- [21] “Interactive google maps for python, created by simon d. levy.” <https://github.com/simondlevy/GooMPy>. Accessed: 2018-06-11.
- [22] “Google static maps api.” <https://developers.google.com/maps/documentation/maps-static/intro>. Accessed: 2018-05-08.
- [23] “The norwegian institute of public health: Influenza information.” <https://fhi.no/en/id/influenta/seasonal-influenza/>. Accessed: 2018-06-11.

- [24] “The norwegian public roads administration: Open data, api for developers.” <https://www.vegvesen.no/en/the+npra/about-the-npra/open-data>. Accessed: 2018-06-11.
- [25] “Tkinter scrollable frame, created by eugene bakin.” <https://github.com/simonlevy/GooMPy>. Accessed: 2018-06-11.