

Title page

INFORMATION EXTRACTION ABOUT FLOOD EVENTS FROM TWEETS
YASER KADDOURA

This page intentionally left blank.

Abstract

This page intentionally left blank.

Contents

1	Introduction	1
2	Literature Review	3
2.1	Data Collection	3
2.2	Text Classification	4
2.3	Location Extraction	7
2.4	Text Analysis	8
2.5	Visualization	9
3	Methods	14
3.1	Data Collection	14
3.2	Text Classification	16
3.3	Location Extraction	17
3.4	Text analysis	18
3.5	Visualization	19
4	Results	25
4.1	Text Classification	25
4.2	Experiments	25
A	Diagrams	34
B	Examples	36
B.1	Nominatim output example	36
	References	38

List of Figures

2.1	RNN example [44]	6
2.2	Two Recurrent Neural Network (RNN)s making an encoder-decoder architecture [44]	6
2.3	Two RNNs making an encoder-decoder architecture with attention mechanism [44]	6
2.4	Transformer’s encoder-decoder architecture [44]	7
2.5	Global Flood Monitor application showing flood events	10
2.6	Petersen and Styve [33] application	11
2.7	Web map application with pluvial flood in Berlin by Feng and Sester [12]	12
2.8	Map with tweet markers in by Barker and Macleod [3]	12
2.9	Bubble map of tweets by Barker and Macleod [3]	13
3.1	Flow chart for the pipeline	14
3.2	Data collection and text classification steps of the pipeline	15
3.3	Flow chart for the location extraction step of the pipeline	18
3.4	Visual interface	20
3.5	Map showing clusters of tweets	21
3.6	Histogram for tweets’ creation dates	22
3.7	Table showing the tweets	22
3.8	Scatter plot for T-distributed Stochastic Neighbor Embedding (t-SNE)’s space	23
3.10	Metadata about the visual interface	23
3.9	Tables showing terms with respect their frequency and their weights	24
4.1	Map showing tweets about flood event in Gävleborg	28
4.2	Histogram showing tweets about flood event in Gävleborg	29
4.3	Tweet table, scatter plot, and Latent Dirichlet Allocation (LDA) table showing a cluster of tweets about flood event in Gävleborg	30
4.4	Tweet table showing tweets about flood event in Gothenburg	31
4.5	Map and histogram showing tweets about flood event in Swedish counties	32
4.6	Tweet table, scatter plot, and LDA table showing a cluster of tweets about flood event in Swedish counties	33
A.1	Flow chart for the pipeline	35

List of Tables

3.1	Dataset attributes	15
3.2	Tweet attributes used	16
3.3	Confusion matrix	17
4.1	Evaluation metrics	25
4.2	Miss-classified tweets	26
4.3	Miss-classified tweets for floods in Gävleborg and Dalarna	27

List of Acronyms

API Application Programming Interfaces	1
BERT Bidirectional Encoder Representations from Transformers	7
CNN Convolutional Neural Networks	5
DVC Data Version Control	16
LiU Linköping University	2
LSTM Long Short-Term Memory	7
ML Machine learning	4
NER Named-entity recognition	2
NLP Natural Language Processing	4
RNN Recurrent Neural Network	vi
SMHI Swedish Meteorological and Hydrological Institute	1
NLTK Natural Language Toolkit	18
LDA Latent Dirichlet Allocation	vi
t-SNE T-distributed Stochastic Neighbor Embedding	vi
DBSCAN Density-Based Spatial Clustering of Applications with Noise	19
SVM Support Machine Vector	5

TF-IDF Term Frequency–Inverse Document Frequency	5
ULMFit Universal Language Model Fine-tuning	7
URL Uniform Resource Locator	15
VGI Volunteered Geographic Information	7

(Might be needed to make introduction starts at odd page number)

This page intentionally left blank.

Chapter 1

Introduction

Earlier in this century, floods around Lake Vänern and Arvika have costed Sweden an estimate of 11.1 billion Swedish Krona for damages and repairs [36]. Counties of Dalarna and Gävleborg suffered from flash floods in 2021, disturbing the daily life of their citizens and damaging public and private properties [5]. Flooding is a devastating natural disaster that threatens the lively hood of people and the infrastructure of communities around the world [13].

To facilitate the process of emergency management during these hazardous events, early warning systems analyse their risk, monitor and warn the public while ensuring their readiness [4]. Traditionally, meteorologists forecast the weather by relying on tools such as gauges, satellites, and radars for data extraction. The emergence of social media platforms such as Twitter provide individuals with a public space to share their experience, effectively creating another potential data.

Researchers started harnessing this new wealth of information to aid the disaster management procedure. Twitter streaming Application Programming Interfaces ([API](#)) makes it possible to create a monitoring system for early event detection on a global [7] and local [3] scales. Another use for it would be identifying victims in real-time, locating their physical location, and communicating the information to rescue teams [40]. After the threat subsides, emergency managers can use relevant tweets to assess the impact and plan the recovery phase [3]. To prepare for future floods, authoritative entities can make informed actions by analysing historical data and determining the locations suffering from recurrent calamities. This newly acquired knowledge can augment weather warning systems' pipelines improving their accuracies such as Swedish Meteorological and Hydrological Institute ([SMHI](#)), a Swedish expert authority with a global perspective and a vital task in predicting changes in weather, water, and climate [41].

Most research addresses the problem on a global or national scale, and none addresses Sweden specifically. This thesis project covers this gap by using Swedish tweets to extract relevant knowledge about flood events in Sweden with the focus on answering the following questions:

1. What methods can be used to classify Swedish tweets related to flood events?
2. How to extract the locations mentioned in the tweets?

3. What visualizations can be used to represent the results?
4. What insights can be extracted by using text analysis on the tweets?

The project focuses on providing a proof of concept for addressing the questions above, and it won't design a solution to make a production-ready product; i.e., the engineering challenges, such as automation, scalability, and ease of deployment, are a delimitation. This thesis project is a part of a research project, AI for Climate Adaptation [27] at Linköping University ([LiU](#)) in collaboration with the [SMHI](#). It implements a pipeline that provides a visual representation of tweets related to flood events in Sweden. First, relevant tweets are pulled, processed, and classified from the Twitter API using data mining techniques. Second, physical locations are extracted from tweets mentioning flood events employing Named-entity recognition ([NER](#)) and gazetteers. Finally, the results are presented on a spatiotemporal visualization, and text analysis techniques are done on the tweets. For verification purposes, the pipeline is applied on a week's worth of tweets after past flood events.

Chapter 2

Literature Review

The massive and accessible volume of data that social media produces has attracted the attention of many researchers as a valuable data source for their research topic; however, collecting and processing data of this nature pose many challenges to extracting useful information. This section mentions what other researchers focusing on disaster management topics did to address these challenges while using Twitter; it also discusses the different approaches used for identifying relevant tweets, extracting geographical location from them, making text analysis on the text, and visualizing the results.

2.1 Data Collection

Twitter's [API](#) enables developers to retrieve historical tweets using queries that are made of operators to match a variety of tweet attributes, such as a specific keyword, having a geotag provided by the user who created the tweet, and the language classified by Twitter. Users generate around 500-700 million tweets a day [22], making it necessary to limit the number of tweets to fetch using the [API](#) to reduce computational power and downtime. Feng and Sester [12] only fetches geotagged tweets and then filters necessary them using 45 keywords in 7 languages; this approach filters out a big chunk of relevant tweets since 1% of tweets are geotagged [24]. A better approach is to fetch tweets using keywords related to the topic of interest in different languages. de Bruijn et al. [7] uses over 40 keywords associated with floods in 11 major languages in the query to fetch tweets.

In addition to using textual data, some researchers use other types of data to enhance their pipelines. Some tweets contain media attachments, such as images and videos that are potential visual information for the concerned research topic [1][37][28]; search engines are another resourceful source for images as well [12]. When it comes to flooding events, hydrological information can be a valuable source of information that can be extracted from a global precipitation dataset based on tweets' time stamps and location in the text [8]. Barker and Macleod [3] use Environment Agency flood-monitoring [API¹](#) to get river gauge levels and flood warnings to identify at-risk flooding areas.

¹<https://environment.data.gov.uk/flood-monitoring/doc/reference/#flood-warnings>

Processing text is a crucial part of any Natural Language Processing ([NLP](#)) pipeline to train an effective classifier. Research requires that the corpus is in multiple languages, so translating the text to one language (most likely English) is needed if the classifier can not handle multilingual data [40]. One of the most common text-processing tasks is removing unnecessary terms such as stopwords, Uniform Resource Locator URLs, numbers, and punctuation marks. User mentions in tweets don't provide useful information, so pipelines often remove or replace them with a generic term such as "@user" [8]. The location of the flood event is an important piece of information that is extracted from a term in the tweet, making it a potential target that includes biases in the dataset by overusing it; de Bruijn et al. [8] replaces these terms by the country name that the location is located in; on the other hand, Petersen and Styve [33] replace the terms by the word "place" if they get mentioned more than 0.5% of the size of the data set. Another way to improve the performance of the classifier is to group the terms by converting them to lower-case and transforming them to their lexeme or word stem by lemmatisation or stemming respectively.

Some tweets are noisy or redundant, making them a target for filtering out. Retweets are identical to other tweets without additional context making them unneeded. Spam bots generate similar tweets for malicious reasons, such as spreading false content to manipulate the public; other reasons could be for utility reasons, such as creating a feed for users to check updates. These tweets introduce noise to the dataset that gets reduced by removing duplicate tweets. de Bruijn et al. [7] only considers one tweet from each user in the last 14 days mentioning a specific region; they also remove tweets containing more than five consecutive words that match with those in another tweet among the previous 100 talking about a location. Singh et al. [40] approaches this problem by only extracting tweets created from mobile phones and only considers tweets from users who have followers/following < 1.

2.2 Text Classification

Identifying disaster events using social media requires a classifier to determine the relevant data. Textual data containing terms related to a disaster doesn't mean that it discusses a disastrous event since words such as "flood" can be used figuratively in sentences (e.g., a flood of joy). A binary classifier labelling the data with "on-topic" and "off-topic" labels is needed to filter out irrelevant content.

Most classifiers use supervised Machine learning ([ML](#)) algorithms requiring labelled data for training. A straightforward approach is to manually label a sample of the tweets [7][3]. Petersen and Styve [33] use Crisilext6 [29], a crowdsource labelled tweets, for training their classifiers that get evaluated on 88 million unlabelled tweets containing flood-related terms [6]. Feng and Sester [12] automatically label the tweets by checking if there is rainfall during the provided time and city location by using a weather [API](#)²; if there's a rainfall, the tweet is labelled positive, negative otherwise.

²<https://www.wunderground.com/weather/api/d/docs>

A classifier needs a numerical representation of the textual data for training. Text is often represented in a real-valued vector by encoding words and their context. There are different word embedding techniques, such as Term Frequency–Inverse Document Frequency (**TF-IDF**) [47] that reflect how important a word is to a document in a corpus. Word2vec [26] and its extension doc2Vec [20] are other word embedding techniques that capture the semantic and syntactic qualities of words via a vector space with several hundred dimensions, where each unique word in the corpus gets assigned to a vector in the space.

There are three groups of approaches for **NLP** tasks: heuristics, **ML**, and deep learning. The heuristics approach is the oldest one which builds rules manually for a specific task by using dictionaries and thesauruses. **ML** techniques, including probabilistic modelling and likelihood maximization, are used on a numerical representation of the textual data to learn a model. Neural networks are a popular choice for handling complex, and unstructured data, making them a suitable candidate for language.

There are different groups of **ML** algorithms to classify data for varying data types. Supervised algorithms are employed if the training data set is labelled; otherwise, a probabilistic approach can be used by training a naive Bayes classifier on labelled and unlabelled data [21]. Feng and Sester [12] use naive Bayes, random forest, logistic regression, Support Machine Vector (**SVM**) (RBF Kernel), and **SVM** (Linear Kernel) on labelled data transformed using **TF-IDF** with accuracies of 0.7109, 0.7582, 0.7705, 0.7712, and 0.7739, respectively. Petersen and Styve [33] results are more promising, where they train a logistic regression and random forest classifiers with 0.939 and 0.9253 accuracies, respectively. Deep learning approaches generally outperform classical algorithms; one example is Convolutional Neural Networks (**CNN**) trained on word embeddings for sentence classification. Feng and Sester [12] and Petersen and Styve [33] train a **CNN** model on word2vec embeddings with 0.7868 and 0.94611 accuracies, respectively.

RNN [16] is a common artificial neural network for **NLP** tasks, such as text classification, **NER**, and machine translation. Its memory enables it to take information from previous input to update the current input and output vector (called hidden state) as shown in Figure 2.1, making it appropriate for sequential. For common tasks such as translation an encoder-decoder architecture is needed, where the encoder encodes the input sequence into a numerical representation (called the last hidden state) that gets passed to the decoder for output sequence generation. Figure 2.2, taken from Tunstall et al.’s book[44], shows an example of translating the English statement “Transformers are great!” to the German language. **RNN** has shortcomings when it tries to capture the context for long sequences of information, where the encoder might lose the information at the start of the sequence while forming the representation. **RNN**’s weak memory can be addressed by using the attention mechanism that allows the decoder to access all the hidden states of the encoder. The main goal of attention is to enable the decoder to prioritize the states using weights it assigns at every decoding timestamp. Figure 2.3 shows an example for predicting the third token in the output sequence. Even though attention improves the accuracy of the translations, the computations are sequential and cannot be parallelized. In addition, most **NLP** tasks require trains models using a large amounts of labelled text data that might not be available. Transfer learning resolves this

problem by transferring knowledge acquired of solving one problem to other related ones.

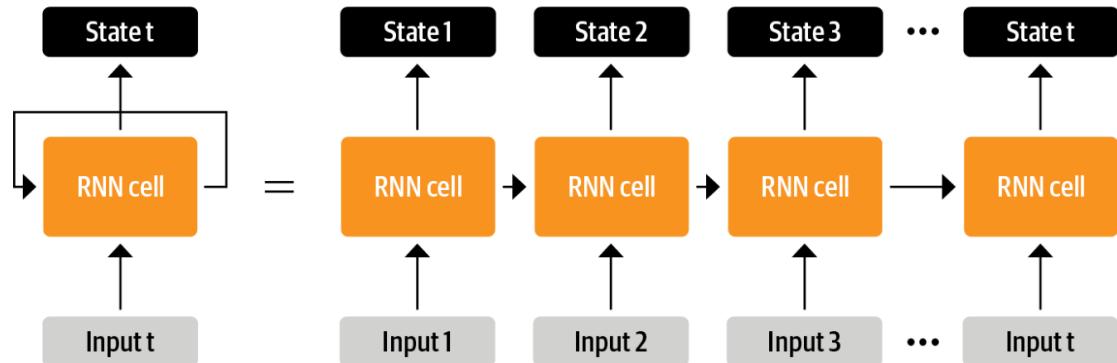


Figure 2.1: RNN example [44]

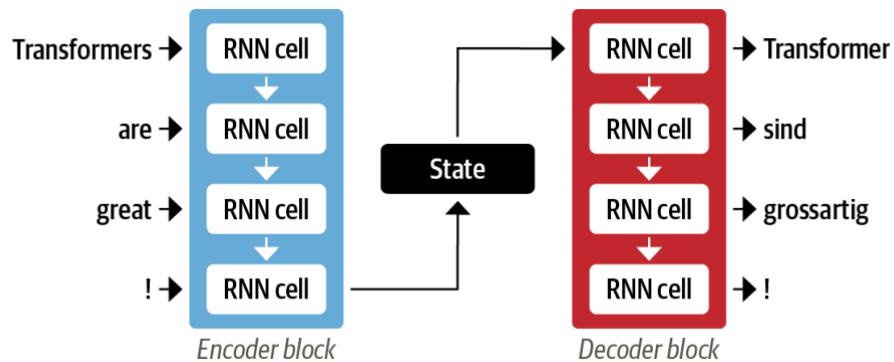


Figure 2.2: Two RNNs making an encoder-decoder architecture [44]

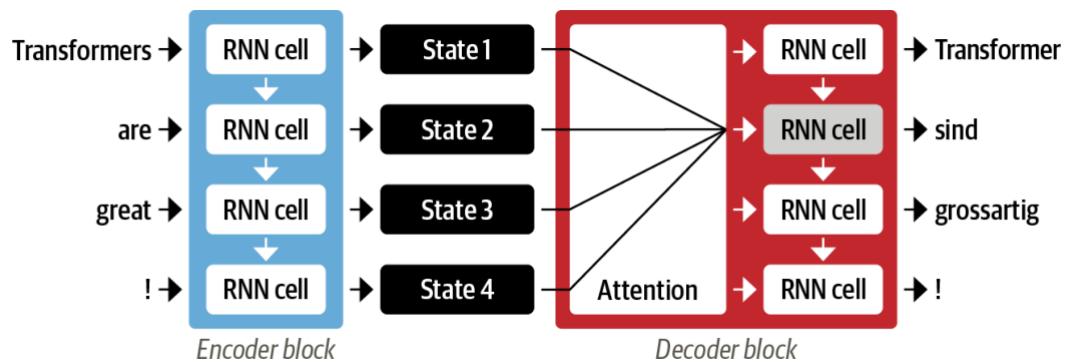


Figure 2.3: Two RNNs making an encoder-decoder architecture with attention mechanism [44]

The concept of transfer learning was used in computer vision before its introduction to [NLP](#). The models are pre-trained on large-scale datasets pre-trained on massive datasets, such as Imagenet [19] and places database, [48] to learn the basic features of images, such as edges or colours. They are fine-tuned on downstream tasks with a smaller dataset. . Feng and Sester [12] use GoogLeNet (Inception-V3 model) [42] pre-trained on ImageNet to train multilayer perceptron, random Forest, gradient boosted trees, and xg-boost with accuracies of 0.8907, 0.9133, 0.9252, and 0.9295, respectively. Ning et al. [28] uses VGGNet [39], Inception V3, ResNet [15], and DenseNet201 [18] with 0.91 accuracy.

In 2017 and 2018, several research groups proposed new approaches to use transfer learning for [NLP](#). Universal Language Model Fine-tuning ([ULMFit](#)) [17] introduced a general framework by pre-training Long Short-Term Memory ([LSTM](#)) models for various tasks. Petersen and Styve [33] fine-tunes a pre-trained [ULMFit](#) model to classify flood-relevant tweets with an accuracy of 0.9499.

Transformers with transfer learning and their self-attention architecture, proposed by google researchers [46], made the training process much faster. The idea is to use attention on all states in the same layer of the neural network. Figure 2.4 shows the self-attention mechanism on both the encoder and decoder with their outputs fed to feed-forward neural networks. Alam et al. [1] fine-tunes a pre-trained Bidirectional Encoder Representations from Transformers ([BERT](#))[10] model that works on one language with an accuracy of 0.853, and de Bruijn et al. [7] uses a multilingual model with 0.8 F1-score.

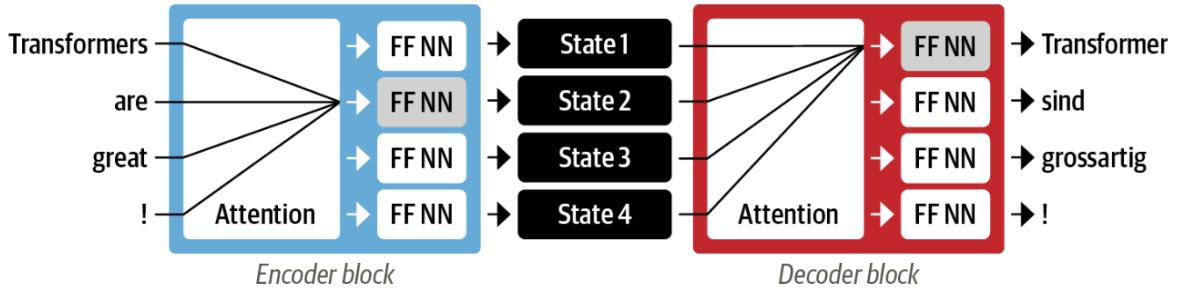


Figure 2.4: Transformer’s encoder-decoder architecture [44]

2.3 Location Extraction

Identifying the locations of disasters is helpful for the disaster management cycle. Social media enables people to generate Volunteered Geographic Information ([VGI](#)), which is more advantageous over the more expensive accuracy testing by official agencies because contributors have unique local knowledge. Detecting a disastrous event and its location as soon as possible can reduce its impact on society [7] by informing the citizens and the authority to prepare for it. During the event, the rescue teams’ task would be easier if they can locate the endangered people [40]. After the event wanes, assessing the most impacted spots can enable the authority to make informed decisions on a recovery plan.

Users can assign an accessible property to their tweets, called “geotag”, a geographical identification metadata. Adding “has:geo” to the query sent to the API will return geotagged tweets with metadata about the location, such as a display name, geo polygon, and a geo lat-log coordinate. The geotag is the most straightforward method to identify the locations [12], but unfortunately, only 1% of the tweets are geotagged [25], making it not include a massive amount of tweets mentioning locations.

Locations can be extracted using toponyms, a place’s name, in tweets’ text by using geoparsing. Geoparsing is a process of converting free-text descriptions of places (such as “twenty miles northeast of Jalalabad”) into unambiguous geographic identifiers. A toponym can have more than one location candidate, such as “Boston”, which is the name of several places, including “Boston, USA” and “Boston, UK”; this fact makes geoparsing tasks on a global scale harder than local ones. de Bruijn et al. [7] uses TAGGS de Bruijn et al. [9], a geoparsing algorithm, to extract countries, administrative areas, and settlements (i.e. cities, towns, and villages) mentioned within the tweets’ text on a global scale. The process includes toponym recognition and toponym resolution. Toponym recognition extracts the toponyms that refer to one or more locations using a gazetteer, a geographical index, or a dictionary. Toponym resolution predicts the correct location for the toponyms in several steps. A score is assigned to each possible location using metadata related to the tweet, such as the user’s timezone & hometown, the tweet’s coordinates, and mentions of nearby locations. Then, calculate the average score of grouped tweets mentioning the same toponym within a 24-hour. Finally, assign the groups of tweets with the location that has the highest score. Petersen and Styve [33] uses geotag property, geoparsing using [NER](#) on text, and user’s profile location to extract toponyms. If the text contains two toponyms, check if they are close with a distance threshold of 1500km, choose one of them randomly. They use GeoPy³ to assign geographical locations to toponyms, a Python package that is a client for several popular geocoding web services (e.g., GoogleV3 and GeoNames). Singh et al. [40] uses the fact that people visit the same locations daily to generate a Markov chain model on historical tweets created by the same user to locate them.

2.4 Text Analysis

Besides text classification and location extraction, other text analysis techniques extract valuable information from text data. In the case of disasters, disaster managers can use social media to get insights, such as how impactful an event is on society. They can visualize the results to understand the situation and act accordingly.

Gründer-Fahrer, Schlaf, and Wustmann [14] extract multiple relevant pieces of information from social media and present them to disaster managers via a searchable application. They extract the following: topics using HDP-CRF algorithm [43], locations using Openstreetmap⁴ location markers, time using the social media meta data, and names of

³<https://geopy.readthedocs.io/en/stable/>

⁴<https://www.openstreetmap.org/>

organizations using **NER**. They present the information using several interactive graphs such as pie charts, word clouds and line graphs.

Sentiment analysis is a popular text analysis technique that shows people's sentiments during an event. Lu et al. [23] extract sentiment analysis from Twitter about the Ebola virus using three different sentiment classifiers to measure the sentiment score of the tweet depending on the majority of the votes. Also, they calculate the inconsistency between the classifiers using an entropy measure [2]. The positive and negative sentiments are each presented in a density map using solid blue and red colours, respectively; the colour is blurred instead if the inconsistency score is above a certain threshold. Periñán-Pascual [32] tries to extract the sentiment by calculate three scores for the tweets: (1) the reliability of how much the tweet discusses a problem during a hazard, (2) the impact of the tweet by using the user's activity and popularity as well as how much influence the tweet is [31], (3) and the impact of the problem using the previous scores. They present the mean of the scores on a time frame basis on a line graph.

2.5 Visualization

Visualization of the results of an **NLP** pipeline is common practice for several reasons. The massive and complex data can communicate the needed knowledge for different audiences to understand the underlying situation and take action accordingly. The developers can use the visualization to validate that the pipeline is working as intended. The authority can check the Spatio-temporal data to identify places that have recurring floods and reinforce their infrastructure to prepare for future flooding. Also, the plots make event detection and monitoring much faster and more straightforward.

de Bruijn et al. [7] uses historical and real-time data to show flooding events on different levels (countries, administrative areas, and settlements). It is powered using a JavaScript library, leaflet⁵. The application, seen in figure 2.5, contains a map showing the flooding events with an adjustable timeline and a list of tweets for the selected location.

⁵<https://leafletjs.com/>

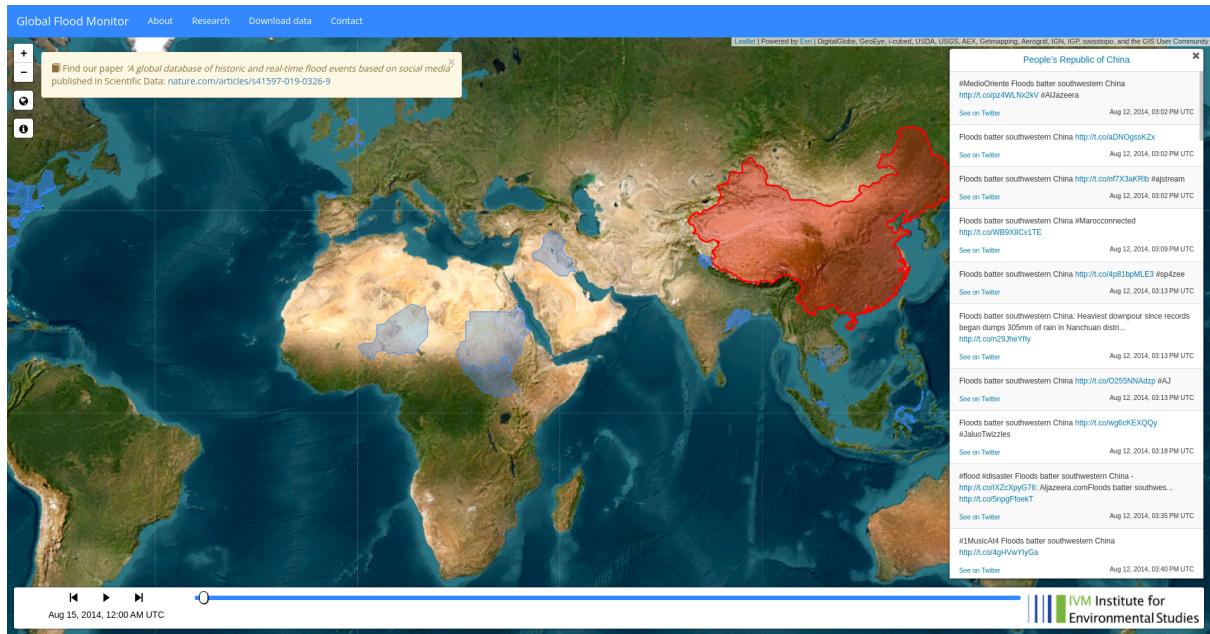


Figure 2.5: Global Flood Monitor application showing flood events

Petersen and Styve [33] provide multiple plots with sophisticated methods to configure the interface and filter the tweets. Their visualization is powered using the python libraries, Plotly⁶ and Dash⁷. The app, shown in figure 2.6, provides an interface to showcase the different aspects of the data: spatial via a map, temporal via a histogram, and textual via a list of tweets and word cloud. They use a scatter map to show the locations extracted from the tweets, where the colour of each point represents the method used to identify the location. To resolve the problem of tweets overlapping each other due to the discussion of the same location, the identical points are spread by adding Gaussian noise to their coordinates points. As for representing the timestamps, they use a histogram aggregated by each day with a time slider. Researchers can pinpoint repetitive or interesting topics by navigating the word cloud to see the most frequent keywords or manually navigating the list of tweets. The plots are interactive, where actions in one of them would influence others. The data can be filtered in different ways: keywords in the text, the method used to extract the location, tweet type (a retweet or not), a map, and a histogram. In addition, there is a drop-down to change the map graph type and the algorithm used to classify the tweets.

⁶<https://plotly.com/python/>

⁷<https://dash.plotly.com/>



Figure 2.6: Petersen and Styve [33] application

Feng and Sester [12] use leaflet to plot a map showing flooding events as observed in figure 2.7. They use Getis-Ord Gi* [30] to detect statistical hot spots and present them as a choropleth map. The light blue circles represent the Spatio-temporal clusters of events, and the circles with numbers at the centre indicate clusters of tweets in that area with their total. The markers indicate individual tweets with a pop-up containing information about it.

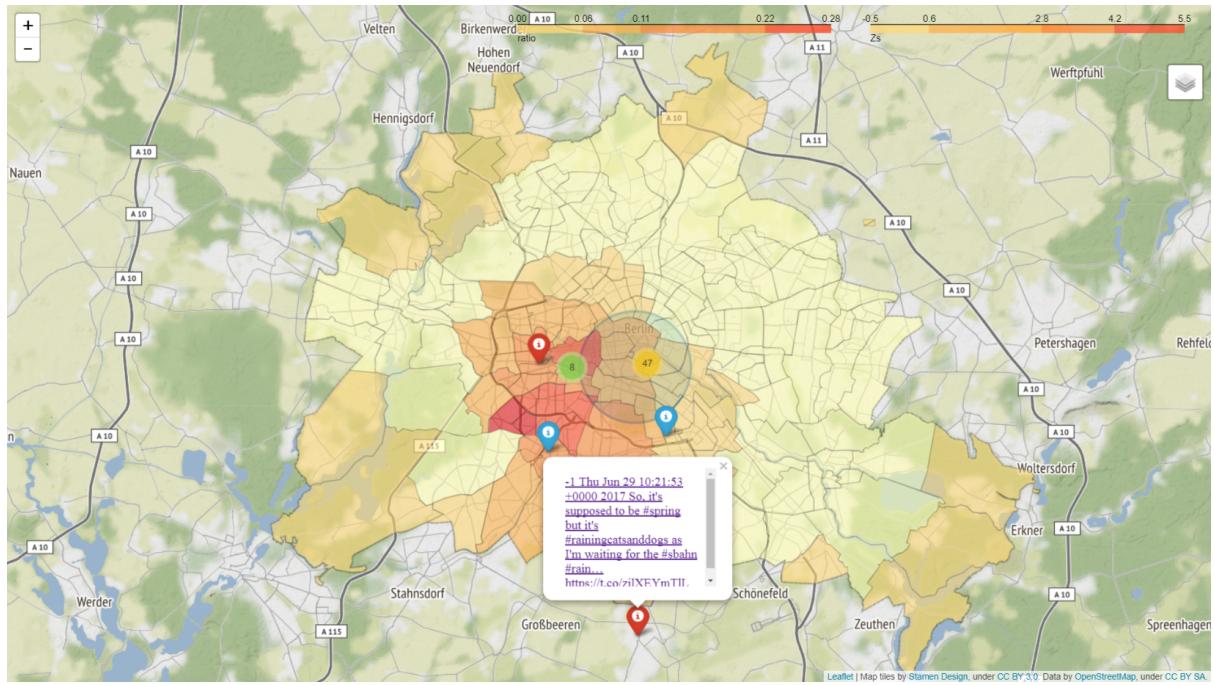


Figure 2.7: Web map application with pluvial flood in Berlin by Feng and Sester [12]

Barker and Macleod [3] visualizes the tweets using different map plots created by leaflet. The map plot in figure 2.8 consists of clickable pointers for pop-up boxes of the tweets.

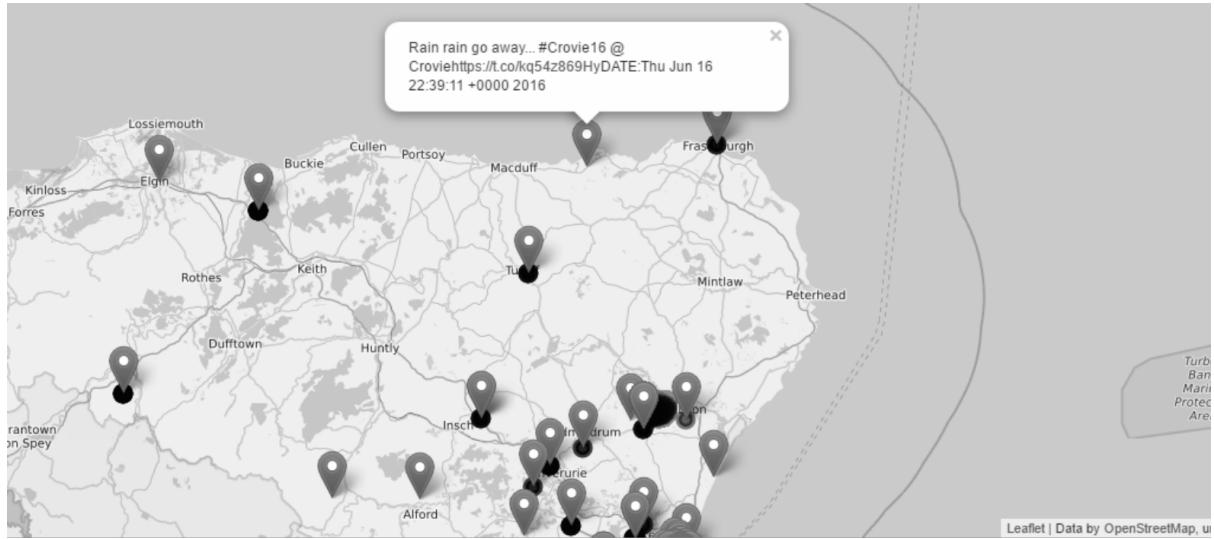


Figure 2.8: Map with tweet markers in by Barker and Macleod [3]

The bubble map in figure 2.9 displays the tweets with the size of the circles representing the area of the place and colour indicating the number of tweets talking about the

location.

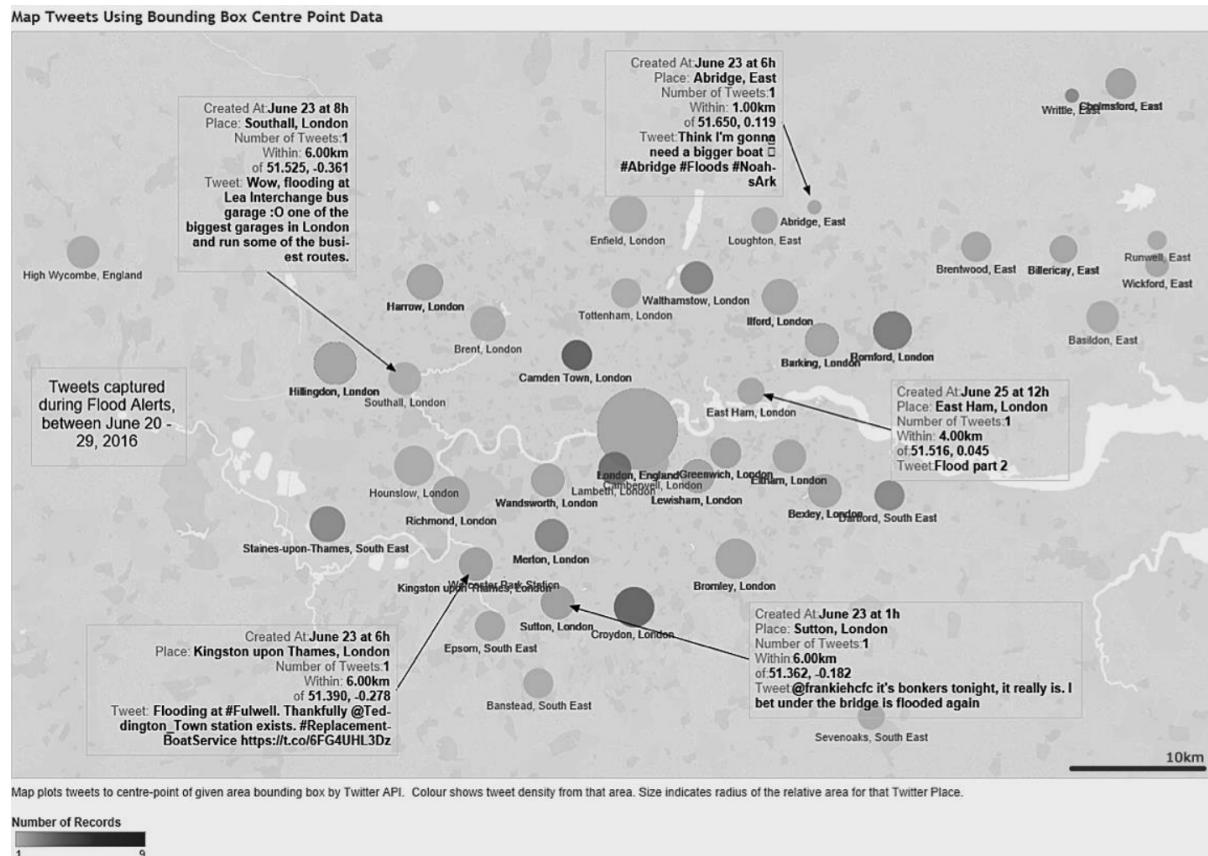


Figure 2.9: Bubble map of tweets by Barker and Macleod [3]

Chapter 3

Methods

This section discusses and motivates the methods used in the project. Figure 3.1 shows a flow chart of the steps for the pipeline (an enlarged and more detailed copy is available in Figure A.1 of appendix A). The pipeline consists of the following steps: Data collection, text classification, location extraction, and visualization. Python is the primary programming language used for the project because of the rich ecosystem surrounding it, especially when it comes to data science-related tasks. The code base is available on a GitHub repository¹ accompanied with a `README.md` containing instructions to set up the environment and run the project.

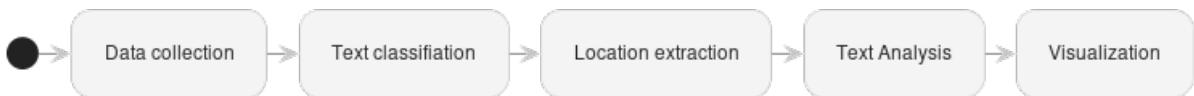


Figure 3.1: Flow chart for the pipeline

3.1 Data Collection

Finding a good quality data source is the first step to having a lean start for most research questions. The pipeline trains an ML classifier using three manually labelled datasets. Two of them are crowdsourced datasets provided by Crisisext6; the tweets are from the 2013 flood events in Alberta² and Queensland³, and there are around 10,000 records for each one with the tweet's ID, tweet's text, and a label about the relevance of the tweet regarding the event. The third dataset is about some flood events in Sweden, spanning between 2015 and 2021; it contains 4899 tweets, mostly in the Swedish language, with attributes presented in Table 3.1. The text and metadata of the tweets are extracted from Twitter's API using the IDs. The trained model performance is verified using

¹<https://github.com/YasserKa/Classification-and-visualization-of-natural-disasters-using-Twitter>

²https://en.wikipedia.org/wiki/2013_Alberta_floods

³https://en.wikipedia.org/wiki/Cyclone_Oswald

Field	Type	Description
ID	Int	ID of the tweet
On Topic	Bool	Text discusses an event
Informative sarcastic	Bool	Text contains relevant information about the event
Contains IMPACT info	Bool	Text discusses the impact of the event
Explicit location	Bool	Text mentions the location of the event

Table 3.1: Dataset attributes

tweets extracted from the API using Tweepy⁴, a python library for accessing Twitter API. Figure 3.2 shows both the source and usage of the data in the pipeline.

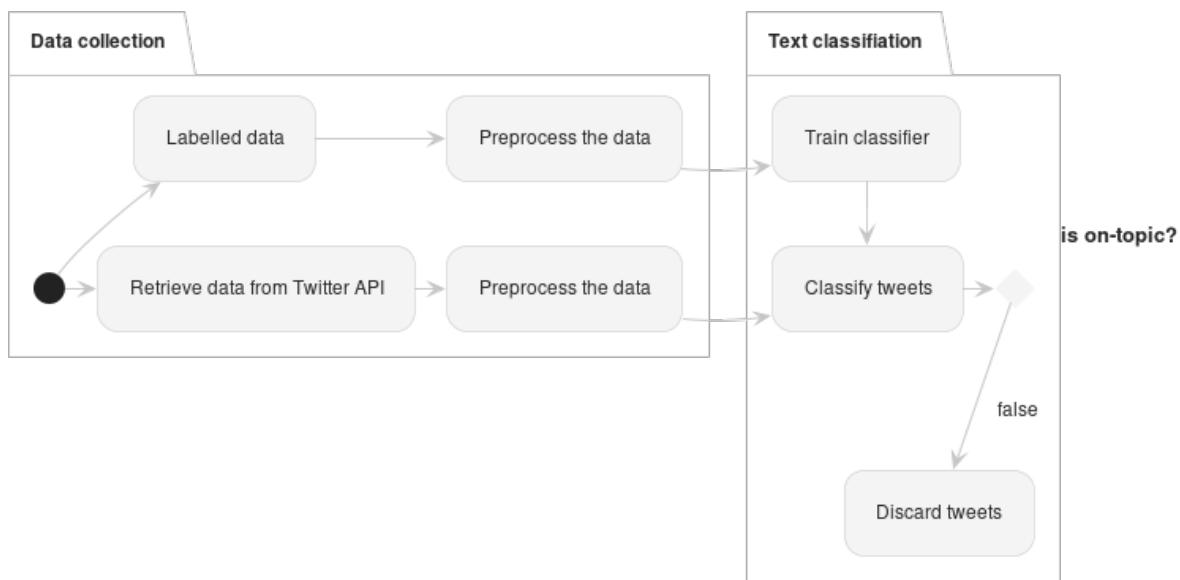


Figure 3.2: Data collection and text classification steps of the pipeline

After retrieving the data, they are pre-processed to prepare them for the upcoming tasks, such as training an ML algorithm, text analysis, and visualization. Parts of text that don't contribute to the context are removed: Uniform Resource Locator (URL)s, emojis, mentions, hashtag signs, numbers, new lines, punctuation, and stopwords (provided by spaCy⁵, an NLP python library). Afterwards, duplicate tweets, tweets containing no text, and retweets are discarded from the dataset. The trained model requires the text to be in the English language, and since Sweden is the focus of the research, most of the text is in Swedish; thus, the text is translated to English using google translate⁶ by a python library wrapper deep-translate⁷.

⁴<https://docs.tweepy.org/en/latest/index.html>

⁵<https://spacy.io/>

⁶<https://translate.google.com/>

⁷<https://deep-translator.readthedocs.io/en/latest/>

Attribute	Type	Description
id	Int	The unique identifier of the requested Tweet
text	Str	The actual UTF-8 text of the Tweet
created at	Date	Creation time of the Tweet
author id	Str	The unique identifier of the tweet creator

Table 3.2: Tweet attributes used

Data needs to be stored and managed to accommodate policies and regulations. Twitter’s developer policy⁸ has a content redistribution section stating that only the IDs of the tweets can be shared online. Thus, the tweets can’t be available publicly on such as GitHub, the service that hosts the publicly available code base. To this end, the data is stored after each step on google drive using Data Version Control ([DVC](#))’s⁹ data management capabilities.

Twitter’s API provides an extensive list of information about the tweets¹⁰. It shares the engagement metrics of the tweet, including like count, reply count, and retweet count; as well as, an [NLP](#) analysis of its own, such as the language used, and entities parsed from the text. Table 3.2 shows the tweet’s attributes used in this project for the following reasons: the id to generate the [URL](#) of the tweet, the text for [NLP](#) tasks, the created date for temporal analysis, and the author id to reduce spam.

3.2 Text Classification

This project uses the DistilBERT transformer[38], a variant of [BERT](#), for text classification. The main advantage of this model is that it achieves comparable performance to BERT while being significantly smaller than BERT while being significantly smaller and more efficient. A DistilBERT pre-trained model is provided by Hugging Face¹¹, a framework that provides a unified API for over more than 50 architectures, making it easier for users to integrate [NLP](#) models into their applications. The learning rate for the neural network is $5 \times e^{-5}$ with 100 warmup steps over four epochs using 90% of the labelled tweets as training data, 5% as test data, and 5% for validation. The text classification purpose is to identify the tweets that discuss flood events, so the “On Topic” attribute of the dataset is used as a label during training.

Training the model locally takes a long time with the available resources, so the training is done using Amazon SageMaker¹², a service that covers tools to build, train, and deploy [ML](#) models. The data is uploaded to Amazon Simple Storage Service (Amazon S3) to make it accessible for the Hugging Face training script that is executed in an instance available in the cloud. After the training is complete, the fine-tuned model and

⁸<https://developer.twitter.com/en/developer-terms/policy>

⁹<https://dvc.org/doc>

¹⁰<https://developer.twitter.com/en/docs/twitter-api/data-dictionary/object-model/tweet>

¹¹<https://huggingface.co/>

¹²<https://aws.amazon.com/sagemaker/>

		Predicted	
		Positive	Negative
Actual	Positive	TP	FP
	Negative	FN	TN

Table 3.3: Confusion matrix

the evaluation metrics are downloaded. The evaluation metrics consists of the following:

- Confusion matrix: a matrix showing the classifier’s predictions for a labelled dataset corresponding to its actual values (Table 3.3).
- Accuracy: a fraction of the number of correctly classified instances (i.e., true positives and true negatives) among all instances (i.e., whole dataset) (equation 3.1).

$$\text{Accuracy} = \frac{TN + TP}{TN + FN + TP + FP} \quad (3.1)$$

- Precision: a fraction of the number of correctly classified relevant instances (i.e., true positives) among the total number of instances classified as relevant (i.e., true positives and false positives) (equation 3.2).

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3.2)$$

- Recall: a fraction of the correctly classified relevant instances (i.e., true positives) among all relevant instances (i.e. true positives and false negatives) (equation 3.3).

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3.3)$$

- F_1 score: a harmonic mean of precision and recall (equation 3.4).

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3.4)$$

3.3 Location Extraction

The project uses a hybrid approach for geoparsing to extract locations. For toponym recognition, the tokens representing locations are extracted using the KBLab/bert-base-swedish-cased-ner model¹³. The model is based on BERT and fine-tuned for NER using The Stockholm-Umeå Corpus, a collection of Swedish texts from the 1990s that consists

¹³<https://huggingface.co/KBLab/bert-base-swedish-cased-ner>

of one million words. As for toponym resolution, the location tokens are disambiguated using Nominatim and GeoNames geocoders through Geopy¹⁴, a Python client for several popular geocoding web services. Nominatim retrieves different fields about the location from OpenStreetMap. An example of the output is available in appendix B.

The descriptions for the fields are available in the documentation¹⁵. The project uses the lat, lon, and display_name to represent the location on a map. In some cases, the text might contain two locations, the one with the smaller bounding box (area of corner coordinates) is used, which is, in most cases, a more specific place located in the bigger one (e.g. a street within a municipality). The geocoder services provide the ability to limit the search of the locations within a specific country. Since the project is limited to Sweden, the output can be limited using this option, reducing the false positives that happen when different countries have places with the same name. Tweets that don't contain location terms identifying a geographical location are discarded as shown in Figure 3.3.

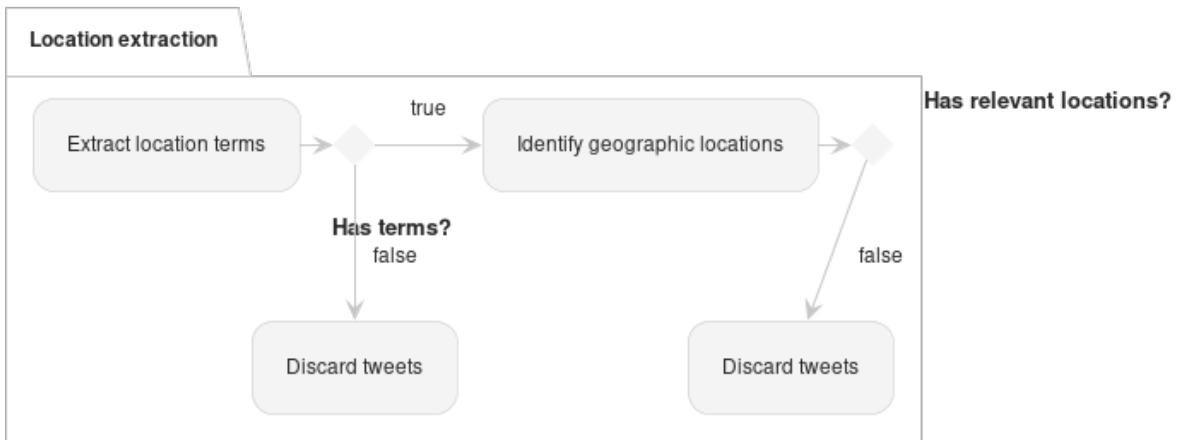


Figure 3.3: Flow chart for the location extraction step of the pipeline

3.4 Text analysis

Further pre-processing is done on the dataset to prepare for text analysis tasks. Lemmatisation is done on the text, using Natural Language Toolkit (NLTK)¹⁶, to reduce words to their lemmas. Afterwards, Tokenisation is done on the corpus. Terms occurring in less than 20 documents or 5% of the documents are removed, as well as the terms mentioned in more than 75% of the documents. Bigrams that occur more than 20 times in the corpus are included, such as traffic jams, and climate change. To reduce spam, tweets created by the same user who tweeted about the same location the past week are discarded.

The text analysis used in the project are LDA [34] [11], TF-IDF, and t-SNE[45]. LDA is a topic modelling method that generates topics (a set of terms) in a corpus and assigns

¹⁴<https://geopy.readthedocs.io/en/latest>

¹⁵<https://nominatim.org/release-docs/develop/api/Output/>

¹⁶<https://www.nltk.org/>

the relevancy of each topic in each document. The [LDA](#) model is initialized and trained using Gensim [35], where the number of discovered topics is adjustable in the visualization. The second text analysis technique used is [TF-IDF](#), using scikit-learn¹⁷, to extract interesting terms by checking their average weight and frequency in the corpus. Lastly, [t-SNE](#) is a visualization method for high-dimensional data by reducing their dimensions to two or three-dimensional maps. In this project, [t-SNE](#) reduces the dimensions of a [TF-IDF](#) matrix generated from the corpus to 2-dimensional space and then presented on a scatter plot; the points are clustered before applying [t-SNE](#) using Density-Based Spatial Clustering of Applications with Noise ([DBSCAN](#)) with adjustable eps (maximum distance between neighbours), and min_samples (number of samples in a neighbourhood for the point to be considered as a core point). [t-SNE](#), the generation of the [TF-IDF](#) matrix, and clustering are done using scikit-learn.

3.5 Visualization

Visualization is often placed at the end of the pipeline and might be the most important since it brings meaning to the results, which can be interrupted by most audiences. Also, it's a direct way to verify that the pipeline is working. The web application is made by Dash¹⁸ to create an interface for Plotly¹⁹'s visualizations. Dash Bootstrap Components²⁰ is used as well for an easier way to use Bootstrap components for Plotly Dash, such as buttons, input, and tables.

Figure 3.4 shows the visual interface containing all the graphs enabling spatial, temporal, and textual exploration of the tweets. Users can add filtering rules for the tweets in all the plots using their creation dates, location, and textual properties.

¹⁷<https://scikit-learn.org/stable/>

¹⁸<https://dash.plotly.com/>

¹⁹<https://plotly.com/python/>

²⁰<https://dash-bootstrap-components.opensource.faculty.ai/>

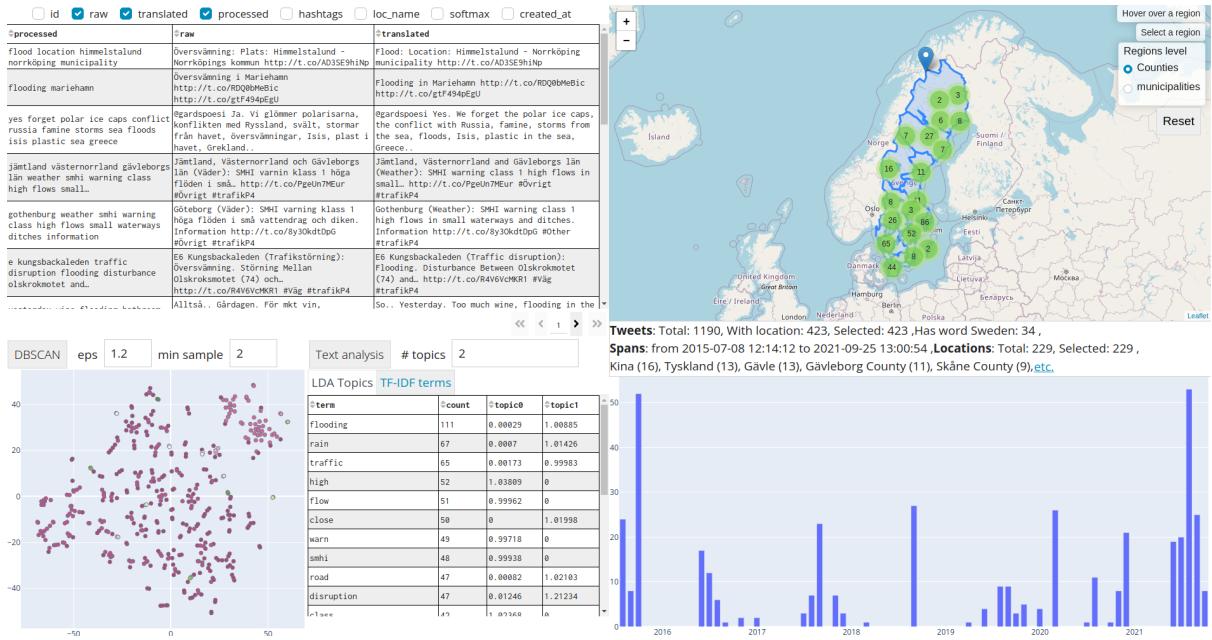


Figure 3.4: Visual interface

Figure 3.5 shows an interactive map containing clickable clustered pointers for the tweets. The clusters disperse or congregate upon zooming in or zooming out, respectively. Hovering over a pointer shows the pop-up with the location name extracted by the pipeline. Clicking on a cluster of a region select the tweets they contain; this will zoom in to cover them while filtering out the unselected pointers from the map. The top left section of the map has several elements: text elements to show the name of the hovered and selected regions, a radio element to change the level of regions that can be selected between counties and municipalities, and a reset button to remove the filter by showing all the pointers.

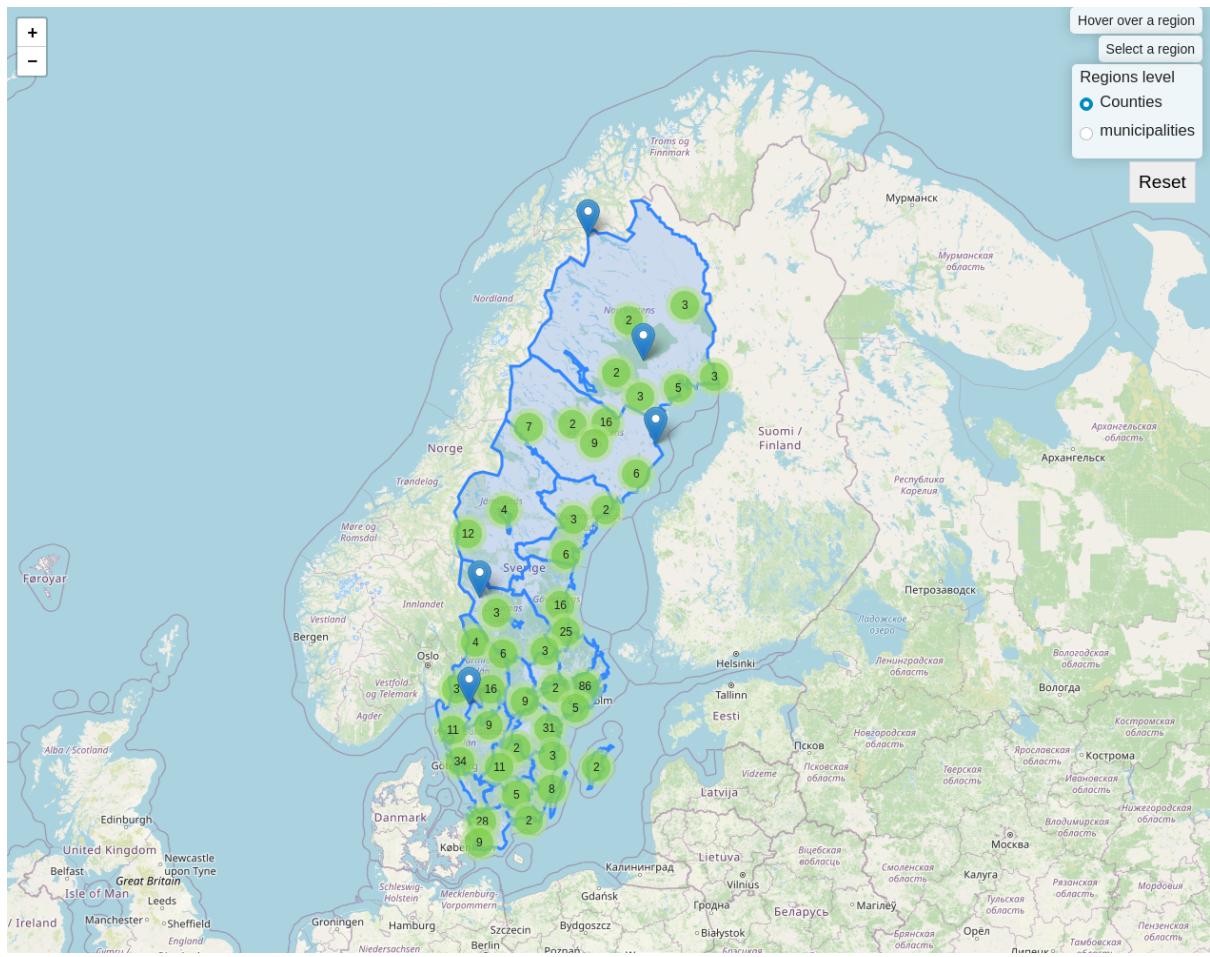


Figure 3.5: Map showing clusters of tweets

Another way to explore the tweets is by using the histogram for their creation date as seen in Figure 3.6. The dates are aggregated by day if they span a month or less; otherwise, by month. Selecting the tweets can be done by using a select box between two dates. Hovering over the bars show the date and the number of tweets, and the blue and red parts of the bars represents the selected and unselected tweets, respectively.

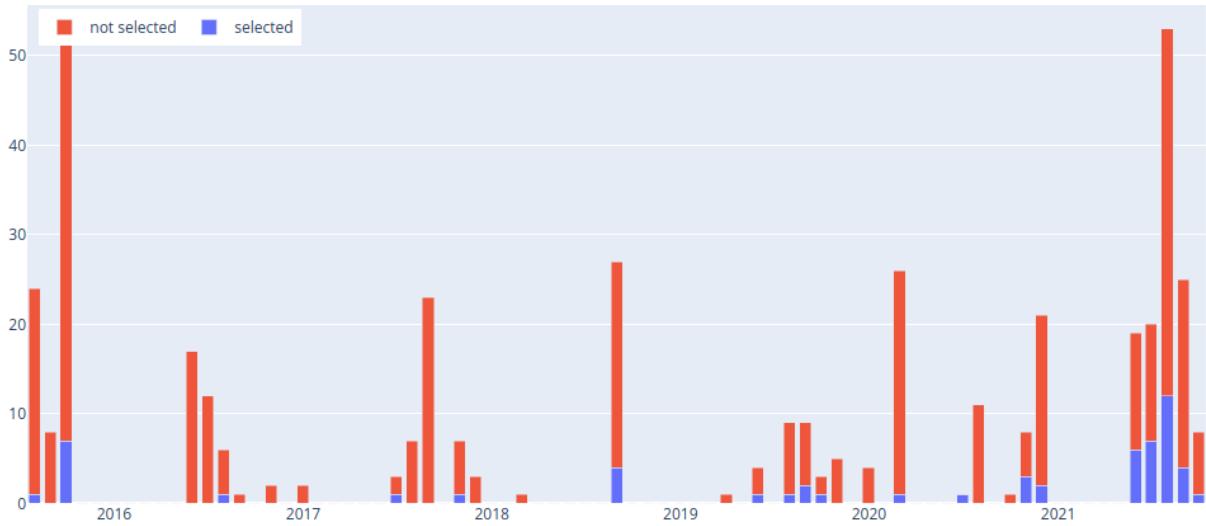


Figure 3.6: Histogram for tweets' creation dates

The table in Figure 3.7 shows the selected tweets with several of its attributes: the ID, the raw (original) text, the translated text, the processed text, the hashtags used, the location extracted, the softmax value for the prediction, and the creation date. It's sortable by column and the rows are paginated.

#created_at	#hashtags	#loc_name	#processed	raw	#softmax	#translated
2015-07-08T12:14:12+00:00		Himmelstalund	Flood location himmelstalund norrkoping municipality	översvämmning: Plats: Himmelstalund - Norrköpings kommun http://t.co/ADSE9hNp	0.945	Flood: Location: Himmelstalund - Norrköping municipality http://t.co/ADSE9hNp
2015-07-08T17:31:06+00:00		Marielhamn	flooding marielhamn	översvämmning i Marielhamn http://t.co/RDQ0BMeBic	0.95	Flooding in Marielhamn http://t.co/RDQ0BMeBic http://t.co/gtF494pEgu
2015-07-09T19:33:44+05+00		Islis	yes forgot polar ice caps conflict russia famine storms sea floods isis plastic sea greece	ljärdspoesi Ja. Vi glömmer polarisarna, konflikten med Ryssland, islänska stormen från havet, översvämningar, Isis, plast i havet, Grekland.	0.935	#ljärdspoesi Yes. We forgot the polar ice caps, the conflict with Russia, Famine, storms from the sea, floods, Isis, plastic in the sea, Greece..
2015-07-25T22:25:33+00:00	#Övrigt, #trafikP4	Gävleborg County	jämtland västerorrländ gävleborgs län weather smhi warning class high flows small..	Jämtland, Västerorrländ och Gävleborgs län (Väder): SMHI varning klass 1 höga flöden i små.. http://t.co/PgeUn7MUr #Övrigt #trafikP4	0.942	Jämtland, Västerorrländ och Gävleborgs län (Weather): SMHI warning class 1 high flows in small.. http://t.co/PgeUn7MUr #Övrigt #trafikP4
2015-07-25T22:28:34+00:00	#Övrigt, #trafikP4	Gothenburg	gothenburg weather sehi warning class high flows small waterways ditches information	Göteborg (Väder): SMHI varning klass 1 höga flöden i små vattendrag och diken. Information http://t.co/b3y0kdtp0 #Övrigt #trafikP4	0.938	Gothenburg (Weather): SMHI warning class 1 high flows in small waterways and ditches. Information http://t.co/b3y0kdtp0 #Other #trafikP4
2015-07-26T06:36:14+00:00	#Väg, #trafikP4	Olskroksmotet	e kungsbäckaleden traffic disruption flooding disturbance olskroksmotet and..	E6 Kungsbackaleden (Traffic disruption): Översvämmning. Störning Mellan Olskroksmotet (74) och... Göteborg. Väg #trafikP4	0.947	E6 Kungsbackaleden (Traffic disruption): Flooding. Disturbance Between Olskroksmotet (74) and... Göteborg. Road. Väg #trafikP4
2015-07-26T09:12:04+00:00		Turkiet	yesterday wine flooding bathroom work mom came home turkey day happy me	Alldts.. Gårdagens. För mat vin, översvämmning i badrummet och jobb på. De. Och mamma kom hem från Turkiet idag. Om hon var glad att se mig?	0.934	So.. Yesterday. Too much wine, flooding in the bathroom and work on them. And mom came home from Turkey today. If she was happy to see me?
2015-07-26T10:25:40+00:00		Dunsjöfjället	smhi warning västra götaland county sjühäradsbygden göta river warning class high flows small watercourses gothenburg	SMHI varnar - Västra Götalands län, Sjühäradsbygden och Göta älvs: Varning klass 1, höga flöden i små vattendrag i Göteborg.	0.945	SMHI warns - Västra Götaland county, Sjühäradsbygden and Göta river: Warning class 1, high flows, in small watercourses in Gothenburg.
2015-07-26T13:08:48+00:00		Gävleborg County	class warning high flows downgraded class applies counties jämtland västerorrländ gävleborg	Klass 2-varning för mycket höga flöden har gått ner till klass 1. Gäller Jämtlands, Västerorrländ och Gävleborgs län http://t.co/Oqvfd3Hep	0.946	Class 2 warning for very high flows has been downgraded to Class 1. Applies to the counties of Jämtland, Västerorrländ and Gävleborg http://t.co/Oqvfd3Hep
2015-07-26T18:24:15+00:00	#nyheter, #svetige	Molkom	floods lightning molkom	översvämmningar och blixtnedslag i Molkom - http://t.co/vs8lLT5K7g	0.941	Floods and lightning in Molkom - http://t.co/vs8lLT5K7g #nyheter #svetige
2015-07-26T20:10:05+00:00	#Övrigt, #trafikP4	Gävleborg County	jämtlan västerorrländ county amp gävleborg county weather smhi warning class high flows small..	Jämtland, Västerorrländ county amp; Gävleborg county (Weather): SMHI varning klass 1: höga flöden i små.. http://t.co/gSPMRAs0 #Övrigt #trafikP4	0.944	Jämtland, Västerorrländ county & ; Gävleborg county (Weather): SMHI warning class 1: High flows in small.. http://t.co/gSPMRAs0 #Övrigt #trafikP4
2015-07-26T22:22:41:48+00:00		Ljuskele kommun	smhi västerbotten county inland warning class high flows umesälven upstream lycksele	SMHI: Västerbottnens län inland: Varning klass 1, höga flöden, Umälvien, uppström Lycksele. http://t.co/SFnWuJyfat	0.947	SMHI: Västerbottnens county inland: Warning class 1, high flows, Umälvien, upstream Lycksele. http://t.co/SFnWuJyfat
2015-07-26T22:22:41:48+00:00		Jämtland County	smhi jämtland county mountains warning class high flows small medium..	SMHI: Jämtland län utom fjällen: Varning klass 1, höga flöden, små och medelstorlek. http://t.co/HWfvq2jsYf	0.947	SMHI: Jämtland county except the mountains: Warning class 1, high flows, small and medium.. http://t.co/HWfvq2jsYf
2015-07-26T22:22:41:48+00:00		Gävleborg County	smhi gävleborg county inland warning class high flows small medium-sized watercourses in..	SMHI: Gävleborg county inland: Varning klass 1, höga flöden, små och medelstorlek vattendrag i.. http://t.co/GwrbewBkn	0.945	SMHI: Gävleborg county inland: Warning class 1, high flows, small and medium-sized watercourses in.. http://t.co/GwrbewBkn
2015-07-26T22:22:41:48+00:00		Östergötland County	smhi östergötland county warning class high flows small medium-sized watercourses..	SMHI: Östergötlands län: Varning klass 1, höga flöden, små och medelstorlek vattendrag i.. http://t.co/9RjCneEuak	0.945	SMHI: Östergötland county: Warning class 1, high flows, small and medium-sized watercourses in.. http://t.co/9RjCneEuak
2015-07-26T22:22:41:48+00:00		Kalixälven	smhi norrbotten county inland warning class high flows kalixälven upstream..	SMHI: Norrbottens län inland: Varning klass 1, höga flöden, Kalixälven, upström.. http://t.co/fWegQW75	0.946	SMHI: Norrbottens county inland: Warning class 1, high flows, Kalixälven, upstream.. http://t.co/fWegQW75
2015-07-26T22:22:41:48+00:00		Piteälven	smhi norrbotten county inland warning class high flows piteälven upper	SMHI: Norrbottens län inland: Varning klass 1, höga flöden, Piteälven, övre delen.. http://t.co/HuL7OWNxlc	0.948	SMHI: Norrbottens county inland: Warning class 1, high flows, Piteälven, upper.. http://t.co/HuL7OWNxlc
2015-07-26T22:22:41:48+00:00		Umeälven	smhi västerbotten county southern lapland mountains warning class high flows umesälven..	SMHI: Västerbottnens län, södra Lapplandsfjällen: Varning klass 1, höga flöden, Umälvien.. http://t.co/gapJuhD00c	0.945	SMHI: Västerbottnens County, Southern Lapland Mountains: Warning class 1, high flows, Umälvien.. http://t.co/gapJuhD00c
2015-07-26T22:22:41:48+00:00		Lilla	smhi norrbotten county inland warning class high flows lilla luleälven	SMHI: Norrbottens län inland: Varning klass 1, höga flöden, Lilla Luleälven.. http://t.co/vfgh8rVMk	0.949	SMHI: Norrbottens county inland: Warning class 1, high flows, Lilla Luleälven.. http://t.co/vfgh8rVMk
2015-07-26T22:22:41:48+00:00		Jämtland County	smhi jämtland county mountains warning class high flows medium lakes..	SMHI: Jämtlands län utom fjällen: Varning klass 1, höga flöden, Medelstora sjörika.. http://t.co/2Eh30WbeQm	0.944	SMHI: Jämtland County except the mountains: Warning class 1, high flows, medium lakes.. http://t.co/2Eh30WbeQm

Figure 3.7: Table showing the tweets

Figure 3.8 shows a scatter plot for the t-SNE's 2-dimensional space with DBSCAN clustering. The text inputs adjusts the clustering properties (eps and min samples). How-

ering over the points show a pop-up of the text for the tweets, and the points can be selected using a box or lasso selection.

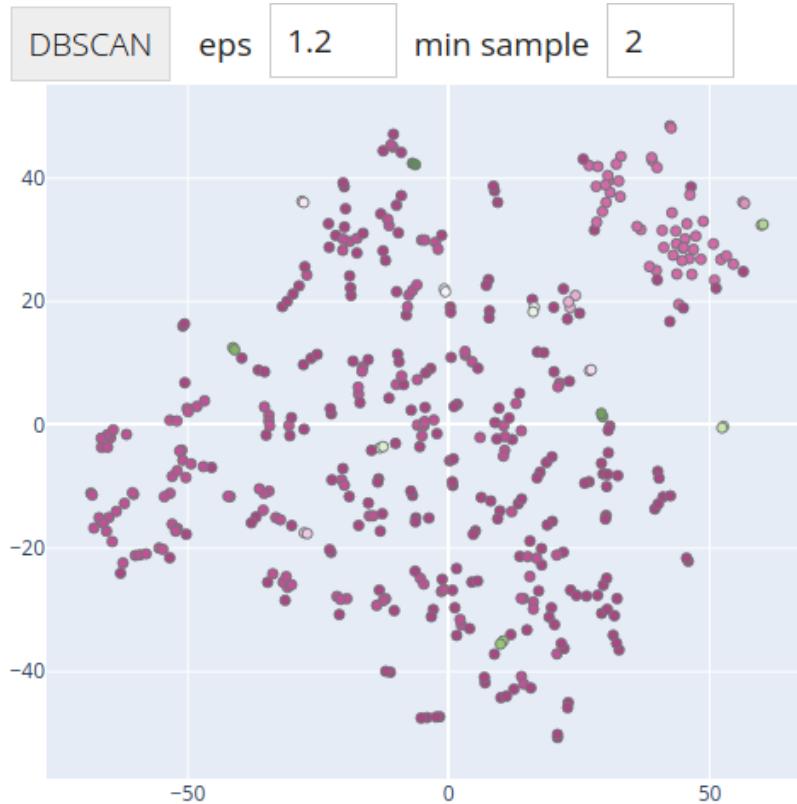


Figure 3.8: Scatter plot for t-SNE’s space

The results of LDA and TF-IDF are displayed in two tables (shown in Figure 3.9a and Figure 3.9b, respectively) showing the frequency of the terms and their mean weights. The number of topics generated by LDA can be changed using a text input, and the tables can be regenerated after changing the selected tweets by clicking the button.

Lastly, the metadata in Figure 3.10 has the following information about the interface: the total number of tweets, the number of selected tweets, the oldest and newest tweet creation dates, the total number of locations, the number of selected locations, and a list of locations’ names with the number of their occurrence. Pressing the “etc.” button shows a pop-over with the rest of the locations.

Tweets: Total: 1190, With location: 423, Selected: 423 ,Has word Sweden: 34 ,
Spans: from 2015-07-08 12:14:12 to 2021-09-25 13:00:54 ,**Locations:** Total: 229, Selected: 229 ,
 Kina (16), Tyskland (13), Gävle (13), Gävleborg County (11), Skåne County (9),[etc.](#)

Figure 3.10: Metadata about the visual interface

Text analysis # topics 2			
LDA Topics TF-IDF terms			
term	count	topic0	topic1
flooding	111	1.02611	0
rain	67	0.99787	0
traffic	65	1.03062	0
high	52	0	1.03825
flow	51	0	0.99968
close	50	1.0197	0
warn	49	0.00025	0.97907
smhi	48	0	0.99976
road	47	1.0636	0
disruption	47	1.25486	0
class	42	0	1.02373

(a) LDA topic weights

Text analysis # topics 2		
LDA Topics TF-IDF terms		
term	mean	count
flooding	0.1325	111
rain	0.16315	67
traffic	0.158	65
high	0.18865	52
flow	0.18478	51
close	0.1828	50
warn	0.19602	49
smhi	0.19863	48
disruption	0.24414	47
road	0.19971	47
class	0.20911	42

(b) TF-IDF weights

Figure 3.9: Tables showing terms with respect their frequency and their weights

Chapter 4

Results

This section presents the results of the methods used in the project to address the research questions by evaluating each step of the pipeline: (1) classifying flood-relevant tweets, (2) extracting geographical locations from tweets, (3) finding useful insights using textual analysis techniques, and (4) visualizing the results in a searchable matter.

4.1 Text Classification

Table 4.1 shows the evaluation metrics mention in Section 3.2 for the trained DistilBERT model on the dataset and a balanced version by doing undersampling using imbalanced-learn’s RandomUnderSampler method¹. Table 4.2 shows falsely classified tweets from the Swedish dataset translated to English.

4.2 Experiments

This section presents the results by applying the pipeline to three unlabelled collections and showing the most noteworthy results from the visualizations. One week’s worth of tweets are extracted from Twitter’s API starting from the date of the beginning of the

¹https://imbalanced-learn.org/stable/references/generated/imblearn.under_sampling.RandomUnderSampler.html

	Accuracy	Precision	Recall	F ₁ Score	Confusion Matrix
Original	0.9231	0.8944	0.9181	0.9061	$\begin{bmatrix} 381 & 34 \\ 568 & 45 \end{bmatrix}$
Undersampled	0.9137	0.9091	0.9138	0.9115	$\begin{bmatrix} 350 & 33 \\ 370 & 35 \end{bmatrix}$

Table 4.1: Evaluation metrics

Translated tweet	Processed tweet	Predicted label	Actual label
The road has rained away outside our driveway!! Damn storm https://t.co/wU6uuZo7El	road rained away outside driveway damn storm	0	1
Right now! Stormy weather in southern Norway. Functionality affected - all resources prioritized to save lives, correct in Vestfold.	right stormy weather southern norway functionality affected resources prioritized save lives correct vestfold	0	1
AFTER LIGHT! Basement full of water? Do you live in #Stockholm and are affected by this weekend's #flooding? Call reporter Nadya bums	light basement water live affected weekends reporter nadya bums	0	1
Impressed by efforts and people's patience. Here is the latest municipal information. #Hallsberg #Flooding #orepol #svpol http://t.co/C0sCxEDtLT	impressed efforts peoples patience latest municipal information	0	1
world Floods, war, famine, terror. Goodnight world.	floods war famine terror goodnight	1	0
Flooding in the bathtub?	flooding bathtub	1	0
A basement was flooded when a water main leaked in #Vårberga in #Borgå #borgåvatten https://t.co/zX08QDqJv9	basement flooded water main leaked	1	0
storm flood assumption years Storm flood assumption off by about 2,500 years https://t.co/v14XtEcbTC	storm flood assumption years	1	0

Table 4.2: Miss-classified tweets

Translated tweet	Processed tweet
Ovädret och det kraftiga regnandet i Gävle har tvingat Brynäs att stänga sin hemmaarena på grund av översvämning. #twittpuck #Brynäs https://t.co/hrZA9icAy7	Ovädret and the heavy rains in Gävle have forced Brynäs to close its home on the ground of overturning. #twittpuck #Brynäs https://t.co/hrZA9icAy7
Blött i Gävle sa Bull.. https://t.co/fV1ChW7ZTR	Wet in Gävle said Bull.. https://t.co/fV1ChW7ZTR
Att tänka på mycket regn bakåt i tiden o tänka på bl.a. ån i Halland som steg o ställde till det !	Thinking about a lot of rain back in time and thinking about e.g. the river in Halland that rose and made it happen!
Nån som vet om det är lite blött i Gävle?	Anyone know if it's a bit wet in Gävle?

Table 4.3: Miss-classified tweets for floods in Gävleborg and Dalarna

events using a query created by experts at a workshop in [SMHI](#) containing flood-relevant terms in Swedish:

```
"atmosfärisk flod" OR "hög vatten" OR åskskur
OR regnskur OR dagvattensystem OR dränering OR "höga vågor"
OR "höga flöden" OR dämmor
OR snösmältning OR blött OR oväder OR stormflod OR vattenstånd
OR vattennivå OR åskväder OR regnstorm"
OR "mycket regn" OR "kraftig regn" OR översvämningsskador
OR översvämnningar OR översvämning
```

Gävleborg and Dalarna counties had a flood event on the 18th of August² damaging their infrastructure, such as houses and roads. After extracting and processing the tweets, 910 are left, of which 700 are classified as flood-relevant. The classifier seems to have high precision, yet it has some low recall. Table 4.3 shows some of the false negatives.

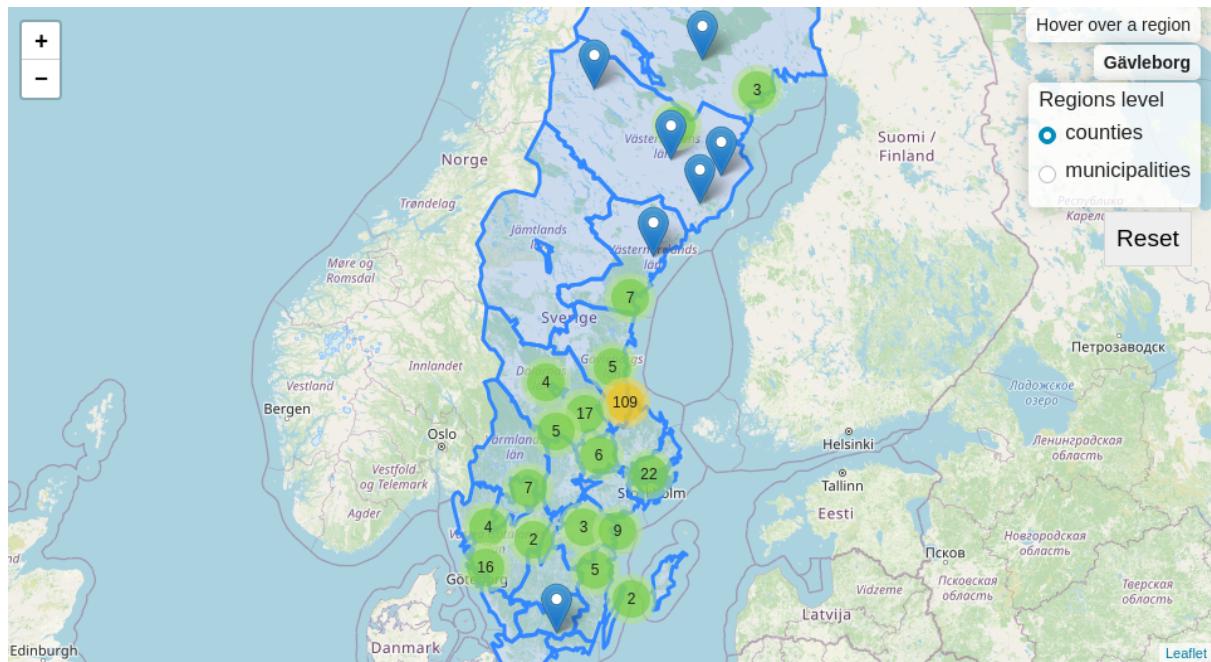
Figure 4.1 shows 114 identified locations in Gävleborg county tweets, such as Gävle (96), and according to the histogram in Figure ??, 81 out of the 114 tweets were created on the 18th of August, 22 on the 19th, and 6 on the 20th. Some locations are identified incorrectly, such as:

- **Original tweet:** Dödssiffran stiger i turkiska översvämningar #Turkiet #svpol
<https://t.co/K6kLRmxQdw>
- **Translated tweet:** Death toll rises in Turkish floods #Turkey #svpol
<https://t.co/K6kLRmxQdw>
- **Identified location:** Turkiet, a hamlet³ in Uppsala county.
- **Actual location:** Turkey, the country.

²<https://floodlist.com/europe/central-sweden-floods-august-2021>

³isolated settlement

- **Original tweet:** Information. Det kraftiga regnovädret över Gävle har orsakat översvämnningar i arenan. Detta innebär att all verksamhet i Monitor ERP Arena, vilket inkluderar bland annat aktivitet på isen samt restaurangverksamheten, tills vidare är pausad. Vi återkommer med mer information. <https://t.co/gHDfirq9VS>
- **Translated tweet:** Information. The heavy rain over Gävle has caused flooding in the arena. This means that all activities in the Monitor ERP Arena, which includes activities on the ice as well as restaurant operations, are paused until further notice. We will return with more information. <https://t.co/gHDfirq9VS>
- **Identified location:** Årena, an isolated dwelling⁴ in Kalmar county.
- **Actual location:** Gävle.



Tweets: Total: 700, With location: 247, Selected: 247 ,Has word Sweden: 31 ,

Spans: from 2021-08-17 08:15:09 to 2021-08-22 20:44:07 ,**Locations:** Total: 104, Selected: 104 ,
Gävle (96), Tyskland (7), Brynäs (5), Sundsvall (5), Gävleborg County (5),[etc.](#)

Figure 4.1: Map showing tweets about flood event in Gävleborg

⁴consist of not more than 2 households

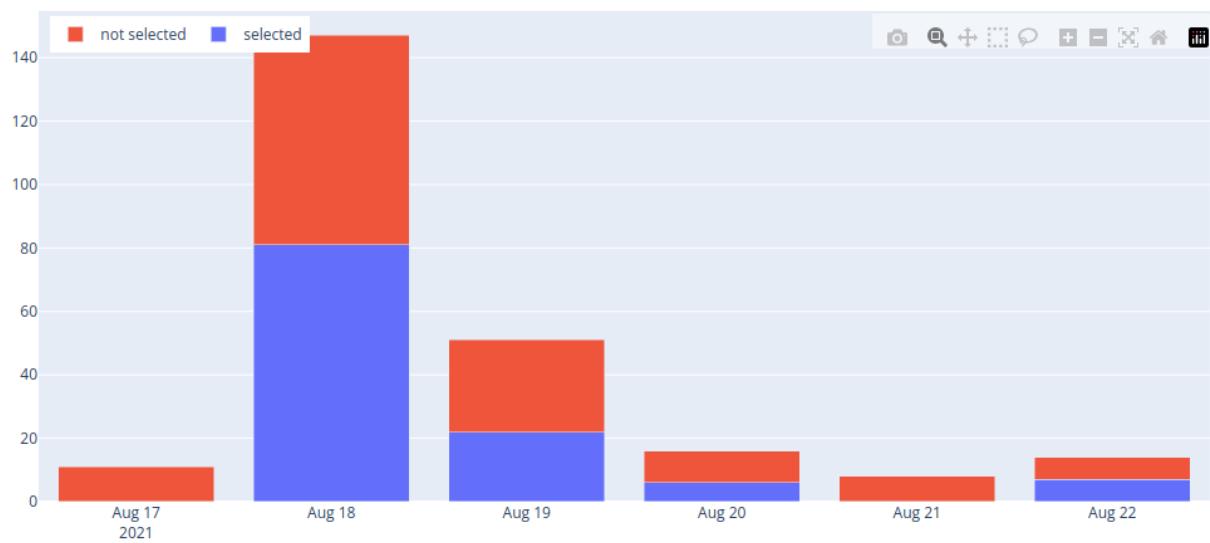


Figure 4.2: Histogram showing tweets about flood event in Gävleborg

Figure ?? shows the scatter plot, the tweet table and the [LDA](#) table after only considering a cluster of tweets in the bottom left of the scatter plot, and they are discussing traffic disruption.

id
 raw
 translated
 processed
 hashtags
 loc_name
 softmax
 created_at

processed	raw	translated
e västerås traffic disruption flooding tpl skälbymotet -tpl bäckbymotet direction enköping	E18 Västerås (Trafikstörning) Översvämnning. Tpl Skälbymotet (129)-tpl Bäckbymotet (130) i riktning mot Enköping https://t.co/AUwDv7KJyt	E18 Västerås (Traffic disruption) Flooding. Tpl Skälbymotet (129)-tpl Bäckbymotet (130) in the direction of Enköping https://t.co/AUwDv7KJyt
e västerås traffic disruption flooding tpl rocklundamotet -tpl vallbymotet direction örebro	E18 Västerås (Trafikstörning) Översvämnning. Tpl Rocklundamotet (132)-tpl Vallbymotet (131) i riktning mot Örebro https://t.co/xUAkMbyRCf	E18 Västerås (Traffic disruption) Flooding. Tpl Rocklundamotet (132)-tpl Vallbymotet (131) in the direction of Örebro https://t.co/xUAkMbyRCf
lv sundsvall traffic disturbance bad road surface storm diversion road outside sawmill	Lv 615 Sundsvall (Trafikstörning) Dålig vägbana efter oväder, på omledningsvägen. Utanför sågverket https://t.co/MAWCKqg308	Lv 615 Sundsvall (Traffic disturbance) Bad road surface after storm, on the diversion road. Outside the sawmill https://t.co/MAWCKqg308
rv borlänge-falun traffic disruption flooding height skommartjärn	Rv 50 Borlänge-Falun (Trafikstörning) Översvämnning. I höjd Skommartjärn https://t.co/wMRQQkwBLJ	Rv 50 Borlänge-Falun (Traffic disruption) Flooding. In height Skommartjärn https://t.co/wMRQQkwBLJ
rv delsbo traffic disruption flooding height staffansgården	Rv 84 Delsbo (Trafikstörning) Översvämnning I höjd med Staffansgården https://t.co/RHDhhsCG8P	Rv 84 Delsbo (Traffic disruption) Flooding At height of Staffansgården https://t.co/RHDhhsCG8P
lv säter-skenshyttan traffic disruption road closed flood	Lv 655 Säter-Skenshyttan (Trafikstörning) Vägen avstängd. Översvämnning https://t.co/KB8voETPtX	Lv 655 Säter-Skenshyttan (Traffic disruption) Road closed. Flood https://t.co/KB8voETPtX
e sandviken-gävle traffic disruption blocking road floods damaged road exit tpl	E16 Sandviken-Gävle (Trafikstörning) Sättning i vägen. Översvämnningar har skadat vägen Strax innan avfarten mot tpl Nybo i riktning	E16 Sandviken-Gävle (Traffic disruption) Blocking the road. Floods have damaged the road Just before the exit towards tpl Nybo in the

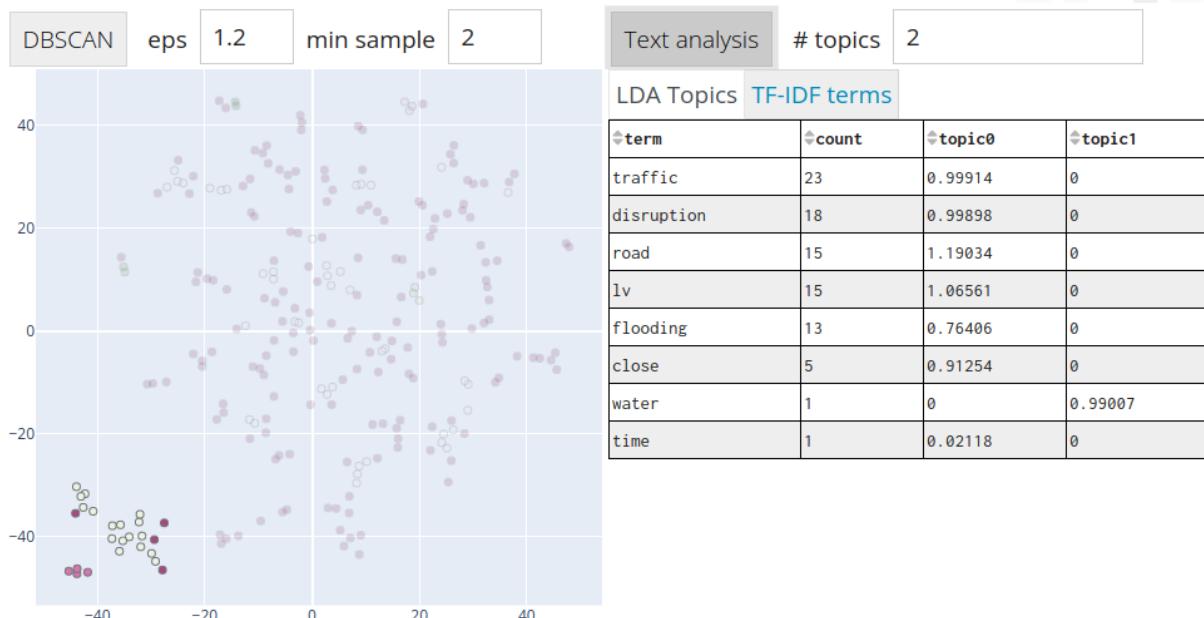


Figure 4.3: Tweet table, scatter plot, and LDA table showing a cluster of tweets about flood event in Gävleborg

Heavy rain caused flooding in Gothenburg on 11 September 2019⁵. Figure 4.4 shows that there are 18 flood-relevant tweets in Gothenburg county, of which 12 were created on the 11th and three on the 12th. There are 16 tweets containing Spanien and identifying it as an isolated dwelling in Stockholm which is incorrect; the tweets are discussing floods

⁵<https://floodlist.com/europe/sweden-flash-floods-gothenburg-september-2019>

in the country Spain⁶. There are false negatives for classifying flood-relevant tweets, such as “It was a little wet. <https://t.co/PcroA3s1A2>”, where the tweet contains a **URL** for an article mentioning the flood event.

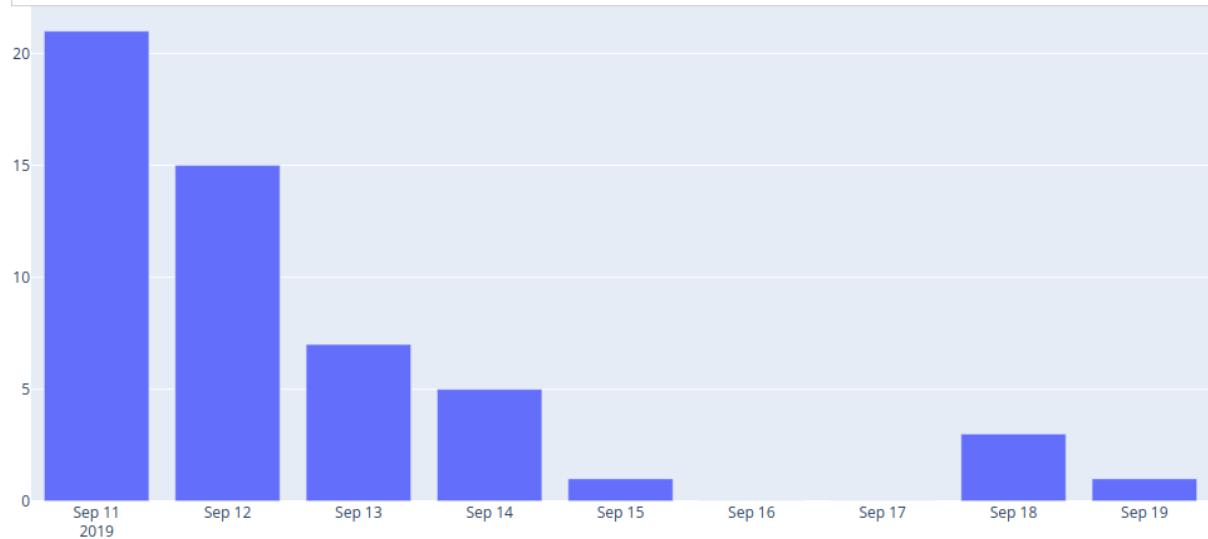
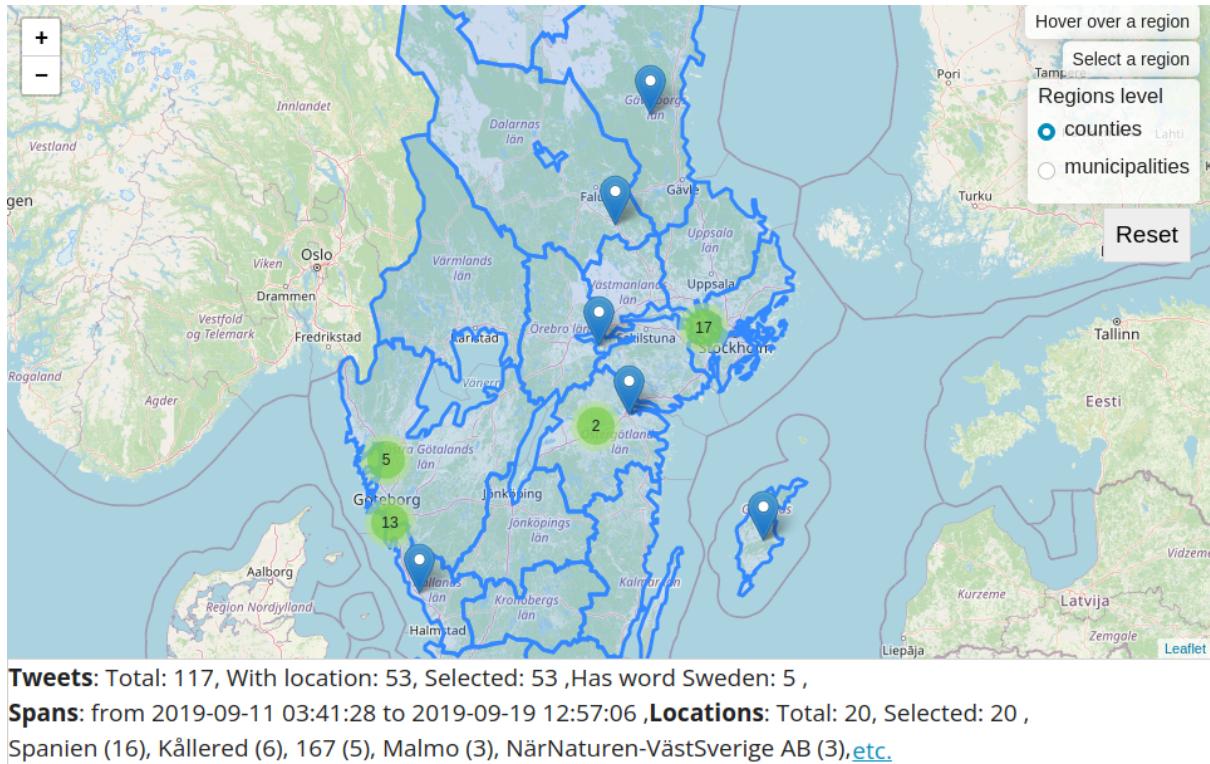


Figure 4.4: Tweet table showing tweets about flood event in Gothenburg

On 22 August 2014, Halland, Värmland and Västra Götaland counties had floods

⁶<https://www.svt.se/nyheter/utrikes/stora-oversvamningar-har-drabbat-sodra-spanien>

lasting four days caused by heavy rain. The map in Figure 4.5 shows that Halland, Värmland, and Västra Götaland counties have 11, 24, and 8 tweets, respectively. The histogram shows eight tweets created on the 22nd of the month and 12 on the 23rd.

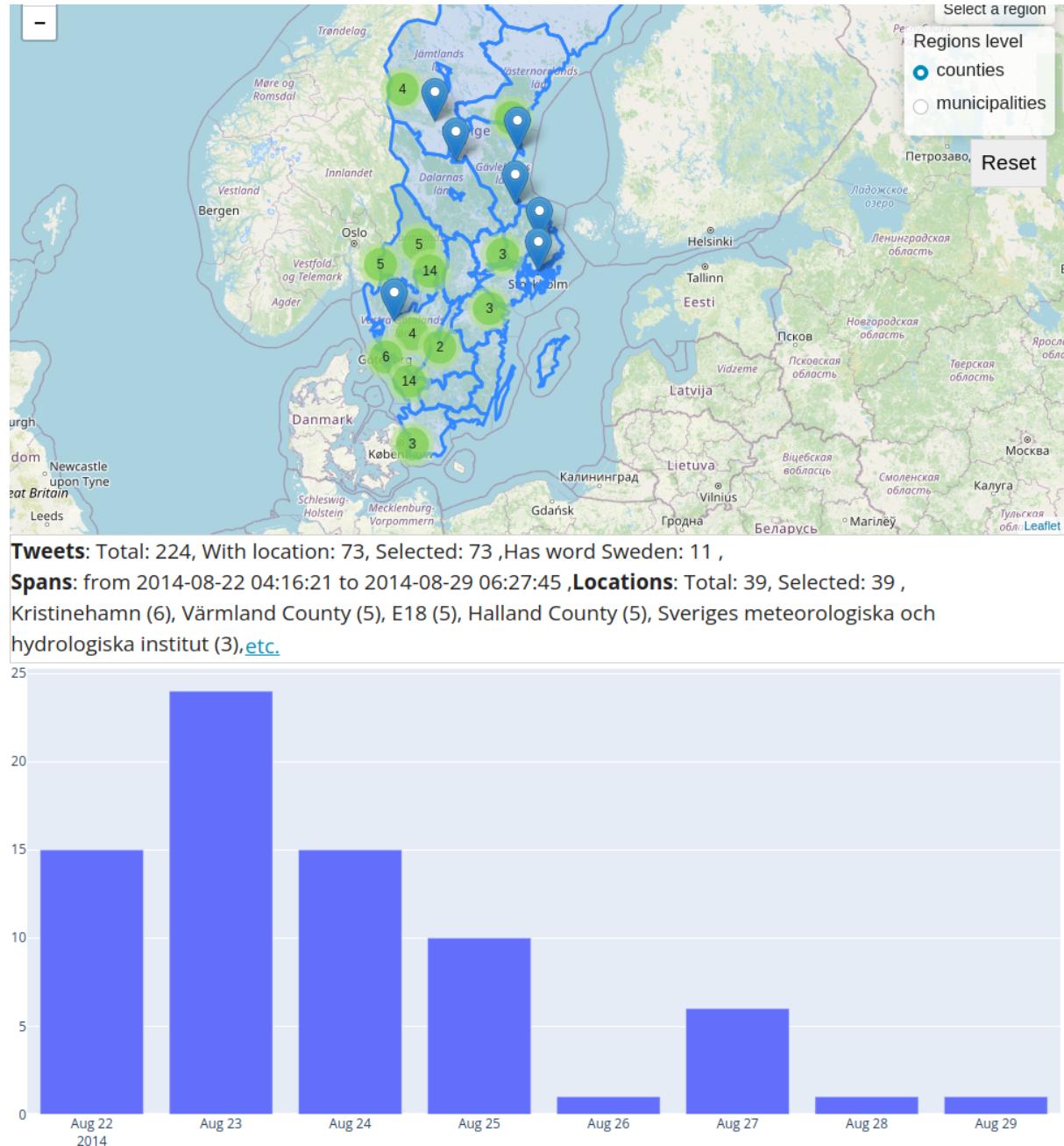


Figure 4.5: Map and histogram showing tweets about flood event in Swedish counties

The bottom left cluster in the scatter plot shown in Figure 4.6 contains tweets discussing SMHI warnings.

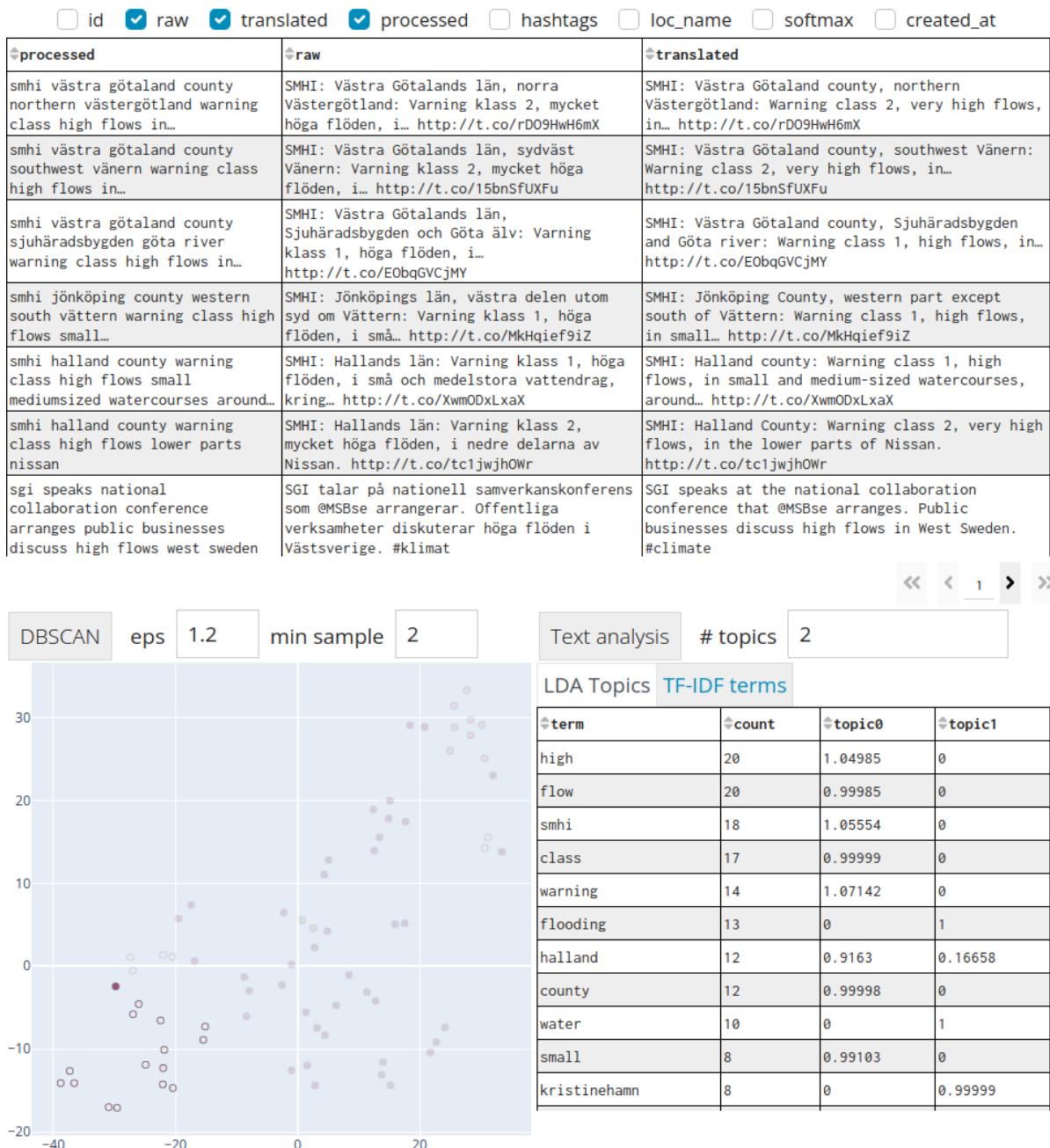


Figure 4.6: Tweet table, scatter plot, and LDA table showing a cluster of tweets about flood event in Swedish counties

Appendix A

Diagrams

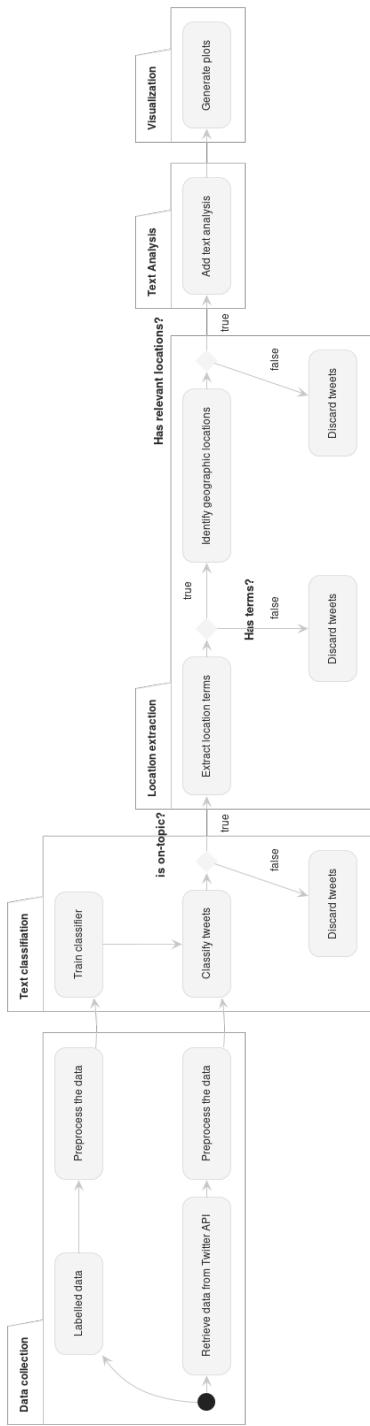


Figure A.1: Flow chart for the pipeline

Appendix B

Examples

B.1 Nominatim output example

```
{  
    "place_id": "100149",  
    "licence": "Data © OpenStreetMap contributors,  
               ODbL 1.0. https://osm.org/copyright",  
    "osm_type": "node",  
    "osm_id": "107775",  
    "boundingbox": ["51.3473219", "51.6673219",  
                   "-0.2876474", "0.0323526"],  
    "lat": "51.5073219",  
    "lon": "-0.1276474",  
    "display_name": "London, Greater London, England,  
                  SW1A 2DU, United Kingdom",  
    "class": "place",  
    "type": "city",  
    "importance": 0.9654895765402,  
    "icon": "https://nominatim.openstreetmap.org/  
            images/mapicons/poi_place_city.p.20.png",  
    "address": {  
        "city": "London",  
        "state_district": "Greater London",  
        "state": "England",  
        "ISO3166-2-lvl4": "GB-ENG",  
        "postcode": "SW1A 2DU",  
        "country": "United Kingdom",  
        "country_code": "gb"  
    },  
    "extratags": {  
        "capital": "yes",  
        "website": "http://www.london.gov.uk",  
    }  
}
```

```
    "wikidata": "Q84",
    "wikipedia": "en:London",
    "population": "8416535"
}
}
```

References

- [1] Firoj Alam et al. *Flood Detection via Twitter Streams Using Textual and Visual Features*. Version 1. Nov. 30, 2020. DOI: [10.48550/arXiv.2011.14944](https://doi.org/10.48550/arXiv.2011.14944). arXiv: [2011.14944 \[cs\]](https://arxiv.org/abs/2011.14944). URL: <http://arxiv.org/abs/2011.14944> (visited on 10/18/2022).
- [2] S. Argamon-Engelson and I. Dagan. “Committee-Based Sample Selection for Probabilistic Classifiers”. In: *Journal of Artificial Intelligence Research* 11 (Nov. 15, 1999), pp. 335–360. ISSN: 1076-9757. DOI: [10.1613/jair.612](https://doi.org/10.1613/jair.612). arXiv: [1106.0220 \[cs\]](https://arxiv.org/abs/1106.0220). URL: <http://arxiv.org/abs/1106.0220> (visited on 01/19/2023).
- [3] J.L.P. Barker and C.J.A. Macleod. “Development of a National-Scale Real-Time Twitter Data Mining Pipeline for Social Geodata on the Potential Impacts of Flooding on Communities”. In: *Environmental Modelling & Software* 115 (May 2019), pp. 213–227. ISSN: 13648152. DOI: [10.1016/j.envsoft.2018.11.013](https://doi.org/10.1016/j.envsoft.2018.11.013). URL: <https://linkinghub.elsevier.com/retrieve/pii/S136481521830094X> (visited on 09/07/2022).
- [4] Wikipedia contributors. *Early Warning System*. In: *Wikipedia*. 1119015319th ed. Wikipedia, The Free Encyclopedia, 10/30/2022, 06:41:00 AM. URL: https://en.wikipedia.org/w/index.php?title=Early_warning_system&oldid=1119015319 (visited on 11/17/2022).
- [5] Richard Davies. *Sweden – Flash Floods in Dalarna and Gävleborg After Record Rainfall*. FloodList. Aug. 19, 2021. URL: <https://floodlist.com/europe/central-sweden-floods-august-2021> (visited on 11/17/2022).
- [6] Jens de. *Flood Tweet IDs (Multilingual)*. Version V2. 2019. DOI: [10.7910/DVN/T3ZFMR](https://doi.org/10.7910/DVN/T3ZFMR). URL: <https://doi.org/10.7910/DVN/T3ZFMR>.
- [7] Jens A. de Bruijn et al. “A Global Database of Historic and Real-Time Flood Events Based on Social Media”. In: *Scientific Data* 6.1 (1 Dec. 9, 2019), p. 311. ISSN: 2052-4463. DOI: [10.1038/s41597-019-0326-9](https://doi.org/10.1038/s41597-019-0326-9). URL: <https://www.nature.com/articles/s41597-019-0326-9> (visited on 10/04/2022).
- [8] Jens A. de Bruijn et al. “Improving the Classification of Flood Tweets with Contextual Hydrological Information in a Multimodal Neural Network”. In: *Computers & Geosciences* 140 (July 2020), p. 104485. ISSN: 00983004. DOI: [10.1016/j.cageo.2020.104485](https://doi.org/10.1016/j.cageo.2020.104485). URL: <https://linkinghub.elsevier.com/retrieve/pii/S0098300419308106> (visited on 11/28/2022).

- [9] Jens A. de Bruijn et al. “TAGGS: Grouping Tweets to Improve Global Geoparsing for Disaster Response”. In: *Journal of Geovisualization and Spatial Analysis* 2.1 (Dec. 26, 2017), p. 2. ISSN: 2509-8829. DOI: [10.1007/s41651-017-0010-6](https://doi.org/10.1007/s41651-017-0010-6). URL: <https://doi.org/10.1007/s41651-017-0010-6> (visited on 10/04/2022).
- [10] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. May 24, 2019. DOI: [10.48550/arXiv.1810.04805](https://doi.org/10.48550/arXiv.1810.04805). arXiv: [1810.04805 \[cs\]](https://arxiv.org/abs/1810.04805). URL: [http://arxiv.org/abs/1810.04805](https://arxiv.org/abs/1810.04805) (visited on 11/26/2022).
- [11] Daniel Falush, Matthew Stephens, and Jonathan K Pritchard. “Inference of Population Structure Using Multilocus Genotype Data: Linked Loci and Correlated Allele Frequencies”. In: *Genetics* 164.4 (Aug. 1, 2003), pp. 1567–1587. ISSN: 1943-2631. DOI: [10.1093/genetics/164.4.1567](https://doi.org/10.1093/genetics/164.4.1567). URL: <https://academic.oup.com/genetics/article/164/4/1567/6050225> (visited on 01/26/2023).
- [12] Yu Feng and Monika Sester. “Extraction of Pluvial Flood Relevant Volunteered Geographic Information (VGI) by Deep Learning from User Generated Texts and Photos”. In: *ISPRS International Journal of Geo-Information* 7.2 (2 Feb. 2018), p. 39. ISSN: 2220-9964. DOI: [10.3390/ijgi7020039](https://doi.org/10.3390/ijgi7020039). URL: <https://www.mdpi.com/2220-9964/7/2/39> (visited on 09/07/2022).
- [13] *Floodlist*. FloodList. Aug. 19, 2021. URL: <https://floodlist.com/europe/central-sweden-floods-august-2021> (visited on 11/17/2022).
- [14] Sabine Gründer-Fahrer, Antje Schlaf, and Sebastian Wustmann. “How Social Media Text Analysis Can Inform Disaster Management”. In: *Language Technologies for the Challenges of the Digital Age*. Ed. by Georg Rehm and Thierry Declerck. Vol. 10713. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2018, pp. 199–207. ISBN: 978-3-319-73705-8 978-3-319-73706-5. DOI: [10.1007/978-3-319-73706-5_17](https://doi.org/10.1007/978-3-319-73706-5_17). URL: [http://link.springer.com/10.1007/978-3-319-73706-5_17](https://link.springer.com/10.1007/978-3-319-73706-5_17) (visited on 01/17/2023).
- [15] Kaiming He et al. *Deep Residual Learning for Image Recognition*. Dec. 10, 2015. DOI: [10.48550/arXiv.1512.03385](https://doi.org/10.48550/arXiv.1512.03385). arXiv: [1512.03385 \[cs\]](https://arxiv.org/abs/1512.03385). URL: [http://arxiv.org/abs/1512.03385](https://arxiv.org/abs/1512.03385) (visited on 01/04/2023).
- [16] J J Hopfield. “Neural Networks and Physical Systems with Emergent Collective Computational Abilities.” In: *Proceedings of the National Academy of Sciences* 79.8 (Apr. 1982), pp. 2554–2558. DOI: [10.1073/pnas.79.8.2554](https://doi.org/10.1073/pnas.79.8.2554). URL: <https://www.pnas.org/doi/10.1073/pnas.79.8.2554> (visited on 01/24/2023).
- [17] Jeremy Howard and Sebastian Ruder. *Universal Language Model Fine-tuning for Text Classification*. May 23, 2018. DOI: [10.48550/arXiv.1801.06146](https://doi.org/10.48550/arXiv.1801.06146). arXiv: [1801.06146 \[cs, stat\]](https://arxiv.org/abs/1801.06146). URL: [http://arxiv.org/abs/1801.06146](https://arxiv.org/abs/1801.06146) (visited on 01/04/2023).
- [18] Gao Huang et al. *Densely Connected Convolutional Networks*. Jan. 28, 2018. DOI: [10.48550/arXiv.1608.06993](https://doi.org/10.48550/arXiv.1608.06993). arXiv: [1608.06993 \[cs\]](https://arxiv.org/abs/1608.06993). URL: [http://arxiv.org/abs/1608.06993](https://arxiv.org/abs/1608.06993) (visited on 01/04/2023).

- [19] Alex Krizhevsky, Ilya Sutskever, and zz Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Communications of the ACM* 60.6 (May 24, 2017), pp. 84–90. ISSN: 0001-0782, 1557-7317. DOI: [10.1145/3065386](https://doi.acm.org/doi/10.1145/3065386). URL: <https://doi.acm.org/doi/10.1145/3065386> (visited on 12/15/2022).
- [20] Quoc V. Le and Tomas Mikolov. *Distributed Representations of Sentences and Documents*. May 22, 2014. DOI: [10.48550/arXiv.1405.4053](https://doi.org/10.48550/arXiv.1405.4053). arXiv: [1405.4053 \[cs\]](https://arxiv.org/abs/1405.4053). URL: [http://arxiv.org/abs/1405.4053](https://arxiv.org/abs/1405.4053) (visited on 12/30/2022).
- [21] Hongmin Li et al. “Disaster Response Aided by Tweet Classification with a Domain Adaptation Approach”. In: *Journal of Contingencies and Crisis Management* 26.1 (Mar. 2018), pp. 16–27. ISSN: 0966-0879, 1468-5973. DOI: [10.1111/1468-5973.12194](https://doi.onlinelibrary.wiley.com/doi/10.1111/1468-5973.12194). URL: <https://doi.onlinelibrary.wiley.com/doi/10.1111/1468-5973.12194> (visited on 09/11/2022).
- [22] Quanzhi Li et al. “How Much Data Do You Need? Twitter Decahose Data Analysis”. In: July 2016.
- [23] Yafeng Lu et al. “Visualizing Social Media Sentiment in Disaster Scenarios”. In: *Proceedings of the 24th International Conference on World Wide Web*. WWW ’15: 24th International World Wide Web Conference. Florence Italy: ACM, May 18, 2015, pp. 1211–1215. ISBN: 978-1-4503-3473-0. DOI: [10.1145/2740908.2741720](https://doi.acm.org/doi/10.1145/2740908.2741720). URL: <https://doi.acm.org/doi/10.1145/2740908.2741720> (visited on 12/16/2022).
- [24] Stuart E. Middleton, Lee Middleton, and Stefano Modaffer. “Real-Time Crisis Mapping of Natural Disasters Using Social Media”. In: *IEEE Intelligent Systems* 29.2 (Mar. 2014), pp. 9–17. ISSN: 1541-1672. DOI: [10.1109/MIS.2013.126](https://doi.ieeexplore.ieee.org/document/6692841/). URL: [http://ieeexplore.ieee.org/document/6692841/](https://ieeexplore.ieee.org/document/6692841/) (visited on 10/19/2022).
- [25] Stuart E. Middleton et al. “Location Extraction from Social Media: Geoparsing, Location Disambiguation, and Geotagging”. In: *ACM Transactions on Information Systems* 36.4 (June 13, 2018), 40:1–40:27. ISSN: 1046-8188. DOI: [10.1145/3202662](https://doi.org/10.1145/3202662). URL: <https://doi.org/10.1145/3202662> (visited on 10/19/2022).
- [26] Tomas Mikolov et al. *Efficient Estimation of Word Representations in Vector Space*. Sept. 6, 2013. DOI: [10.48550/arXiv.1301.3781](https://doi.org/10.48550/arXiv.1301.3781). arXiv: [1301.3781 \[cs\]](https://arxiv.org/abs/1301.3781). URL: [http://arxiv.org/abs/1301.3781](https://arxiv.org/abs/1301.3781) (visited on 12/30/2022).
- [27] Tina Nesi. *AI4ClimateAdaptation*. Linköping University. URL: <https://liu.se/en/research/ai4climateadaptation> (visited on 11/18/2022).
- [28] Huan Ning et al. “Prototyping a Social Media Flooding Photo Screening System Based on Deep Learning”. In: *ISPRS International Journal of Geo-Information* 9.2 (2 Feb. 2020), p. 104. ISSN: 2220-9964. DOI: [10.3390/ijgi9020104](https://doi.org/10.3390/ijgi9020104). URL: <https://www.mdpi.com/2220-9964/9/2/104> (visited on 09/11/2022).

- [29] Alexandra Olteanu et al. “CrisisLex: A Lexicon for Collecting and Filtering Microblogged Communications in Crises”. In: *Proceedings of the International AAAI Conference on Web and Social Media* 8.1 (May 16, 2014), pp. 376–385. ISSN: 2334-0770, 2162-3449. DOI: [10.1609/icwsm.v8i1.14538](https://doi.org/10.1609/icwsm.v8i1.14538). URL: <https://ojs.aaai.org/index.php/ICWSM/article/view/14538> (visited on 11/28/2022).
- [30] J. K. Ord and Arthur Getis. “Local Spatial Autocorrelation Statistics: Distributional Issues and an Application”. In: *Geographical Analysis* 27.4 (Sept. 3, 2010), pp. 286–306. ISSN: 00167363. DOI: [10.1111/j.1538-4632.1995.tb00912.x](https://doi.org/10.1111/j.1538-4632.1995.tb00912.x). URL: <https://onlinelibrary.wiley.com/doi/10.1111/j.1538-4632.1995.tb00912.x> (visited on 01/08/2023).
- [31] Aditya Pal and Scott Counts. “Identifying Topical Authorities in Microblogs”. In: *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining*. WSDM’11: Fourth ACM International Conference on Web Search and Data Mining. Hong Kong China: ACM, Feb. 9, 2011, pp. 45–54. ISBN: 978-1-4503-0493-1. DOI: [10.1145/1935826.1935843](https://doi.org/10.1145/1935826.1935843). URL: <https://dl.acm.org/doi/10.1145/1935826.1935843> (visited on 01/19/2023).
- [32] Carlos Periñán-Pascual. “Assessing the Impact of Tweets in Flood Events”. In: *1st International Workshop on Social Media Analysis for Intelligent Environment* (Jan. 1, 2020). URL: https://www.academia.edu/44757497/Assessing_the_Impact_of_Tweets_in_Flood_Events (visited on 12/16/2022).
- [33] Julie Maria Petersen and Lise Styve. “Identification and Exploration of Extreme Weather Events From Twitter Data”. Linköping University, 2021.
- [34] Jonathan K Pritchard, Matthew Stephens, and Peter Donnelly. “Inference of Population Structure Using Multilocus Genotype Data”. In: *Genetics* 155.2 (June 1, 2000), pp. 945–959. ISSN: 1943-2631. DOI: [10.1093/genetics/155.2.945](https://doi.org/10.1093/genetics/155.2.945). URL: <https://academic.oup.com/genetics/article/155/2/945/6048111> (visited on 01/26/2023).
- [35] Radim Řehůřek and Petr Sojka. “Software Framework for Topic Modelling with Large Corpora”. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. Valletta, Malta: ELRA, May 22, 2010, pp. 45–50.
- [36] *River Floods Sweden*. ClimateChangePost. Nov. 6, 2022. URL: <https://www.climatechangepost.com/sweden/river-floods/> (visited on 11/17/2022).
- [37] Naina Said et al. *Floods Detection in Twitter Text and Images*. Nov. 30, 2020. DOI: [10.48550/arXiv.2011.14943](https://doi.org/10.48550/arXiv.2011.14943). arXiv: [2011.14943 \[cs\]](https://arxiv.org/abs/2011.14943). URL: [http://arxiv.org/abs/2011.14943](https://arxiv.org/abs/2011.14943) (visited on 11/26/2022).
- [38] Victor Sanh et al. “DistilBERT, a Distilled Version of BERT: Smaller, Faster, Cheaper and Lighter”. In: *ArXiv* abs/1910.01108 (2019).
- [39] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image Recognition*. Apr. 10, 2015. DOI: [10.48550/arXiv.1409.1556](https://doi.org/10.48550/arXiv.1409.1556). arXiv: [1409.1556 \[cs\]](https://arxiv.org/abs/1409.1556). URL: [http://arxiv.org/abs/1409.1556](https://arxiv.org/abs/1409.1556) (visited on 12/15/2022).

- [40] Jyoti Prakash Singh et al. “Event Classification and Location Prediction from Tweets during Disasters”. In: *Annals of Operations Research* 283.1 (Dec. 1, 2019), pp. 737–757. ISSN: 1572-9338. DOI: [10.1007/s10479-017-2522-3](https://doi.org/10.1007/s10479-017-2522-3). URL: <https://doi.org/10.1007/s10479-017-2522-3> (visited on 09/07/2022).
- [41] SMHI. SMHI - Who we are. Apr. 30, 2021.
- [42] Christian Szegedy et al. “Rethinking the Inception Architecture for Computer Vision”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 2818–2826. DOI: [10.1109/CVPR.2016.308](https://doi.org/10.1109/CVPR.2016.308).
- [43] Yee Whye Teh and Michael I. Jordan. “Hierarchical Bayesian Nonparametric Models with Applications”. In: *Bayesian Nonparametrics*. Ed. by Chris Holmes et al. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press, 2010, pp. 158–207. ISBN: 978-0-521-51346-3. DOI: [10.1017/CBO9780511802478.006](https://doi.org/10.1017/CBO9780511802478.006). URL: <https://www.cambridge.org/core/books/bayesian-nonparametrics/hierarchical-bayesian-nonparametric-models-with-applications/0051298A8C5D57586096CDDF02AB1B0F> (visited on 01/19/2023).
- [44] Lewis Tunstall et al. *Natural Language Processing with Transformers: Building Language Applications with Hugging Face*. Revised edition. Sebastopol: O'Reilly, 2022. 383 pp. ISBN: 978-1-09-813679-6.
- [45] L.J.P. van der Maaten and G.E. Hinton. “Visualizing High-Dimensional Data Using t-SNE”. In: *Journal of Machine Learning Research* 9 (nov 2008), pp. 2579–2605. ISSN: 1532-4435.
- [46] Ashish Vaswani et al. *Attention Is All You Need*. Dec. 5, 2017. DOI: [10.48550/arXiv.1706.03762](https://doi.org/10.48550/arXiv.1706.03762). arXiv: [1706.03762 \[cs\]](https://arxiv.org/abs/1706.03762). URL: <http://arxiv.org/abs/1706.03762> (visited on 01/24/2023).
- [47] Wikipedia contributors. *Tf-Idf — Wikipedia, the Free Encyclopedia*. 2022. URL: <https://en.wikipedia.org/w/index.php?title=Tf%E2%80%93id&oldid=1123031029>.
- [48] Bolei Zhou et al. “Learning Deep Features for Scene Recognition Using Places Database”. In: *Advances in Neural Information Processing Systems*. Vol. 27. Curran Associates, Inc., 2014. URL: <https://papers.nips.cc/paper/2014/hash/3fe94a002317b5f9259f82690aeea4cd-Abstract.html> (visited on 12/15/2022).