# Using Latent Dirichlet Allocation and Word Embedding on Twitter Datato Gain Insight into the Mindsets of Young People.
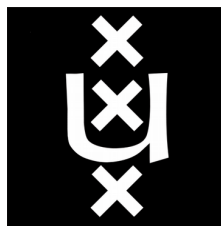
SUBMITTED IN PARTIAL FULLFILLMENT FOR THE DEGREE OF MASTER OF SCIENCE

Sije van der Veen
11422688

MASTER INFORMATION STUDIES
Data Science

FACULTY OF SCIENCE
UNIVERSITY OF AMSTERDAM

Date of defence



| | |
|---|---|
| *Academic Supervisor* | *Industry Supervisor* |
| *Dr. Maarten Marx* | *Dr. Kyle Snyder* |

# Using Latent Dirichlet Allocation and Word Embedding on Twitter Data to Gain Insight into the Mindsets of Young People

SIJE VAN DER VEEN, University of Amsterdam

Generating new topics to write about is not always easy for blog-writers, but the social media data connected to the blog can provide a solution. This paper describes a methodology for detecting word topics (combination of words that together form a topic) within Twitter messages. Latent Dirichlet allocation (LDA) is used to detect word topics in this research. Which is a novel approach because it has not been applied on 280-character-long Twitter messages. Through two sub-questions, the performance of LDA on Twitter messages is assessed. In the last experiment, the word topics are projected using a combination of Word2Vec with principal component analysis (PCA), t-distributed stochastic neighbor embedding (t-SNE) or uniform manifold approximation and projection (UMAP). These projections show, in a two-dimensional way, the links between words within a topic. These links can be used by companies and organizations that want to gain insights into their followers' thoughts.

CCS Concepts: • **Topic modeling** → **Latent Dirichlet Allocation**; • **Word embedding** → *Word2Vec*; • **Dimension reduction** → Principal component analysis ; t-distributed stochastic neighbor embedding; uniform manifold approximation and projection.

## 1 INTRODUCTION

This project was created in collaboration with Radio Netherlands Worldwide (RNW) Media. RNW Media creates communities online to contribute to social change by writing about important and sensitive topics that young people in restrictive settings care about , and having young people consume this information and discuss it [1]. Since RNW Media wants to create content that targets young people, one of the organization's challenges is the gap between the (technical or institutional) language of the organization and that of their audiences. Another problem that the organization has encountered is the costly supervised topic modeling for RNW Media makers, which also introduces bias. Unseen topics by media makers should also be provided by content on different platforms. This project aims to tackle these problems through a data-driven approach. LDA is a well-studied topic modeling method [20]. LDA learns the connections between words, topics, and documents by assuming that documents are generated by a probabilistic model. The applicability of topic modeling across languages makes the approach suitable for RNW Media since it publishes articles in Arabic, Chinese, Hindi, French, and English. In general, NGOs work with many complex languages and, therefore, this study can contribute to this field. The goal of this thesis is to gain insights into the mindset and behavior of young people from fragile countries to eventually positively influence these societies. We aim to do this by applying topic modeling to Twitter messages of the RNW Burundi Twitter platform, which contains Twitter messages from Burundi citizens. The assumption is that the topics will largely categorize and summarize what people of this platform are discussing, which can ultimately be analyzed to gain insights into the mindsets of the people who posted the messages. This thesis aims to answer the following questions:

Author's address: Sije van der Veen, University of Amsterdam.

(1) What is the optimal number of topics that can be generated on a set of Twitter messages using LDA?

(2) Are hash-tags found in the same LDA output topic? In the second sub-question, the LDA output is compared with the hash-tags that label each Twitter message. Are the messages within common hash-tags assigned to similar output groups?

(3) Can the relationship between words within a topic be visualized? In the last sub-question, a word embedding is created based on the sentences from each topic. The word embedding is used to create PCA, t-SNE, and UMAP of a Word2Vec model. In the projection, the most frequent words are shown.

The report commences by describing the related work that has been done and the ideas and models that have already been developed. Section 3 outlines the methodologies used in the study, after which Section 4 outlines the experimental setup, explaining the data and implementation of the models. The results of the performed experiments are discussed in Section 5.

## 2 RELATED WORK

### 2.1 Topic Modeling with Latent Dirichlet Allocation

Since the introduction of LDA in 2000, this topic model method has been studied extensively [20]. This section provides a brief overview of how topic modeling algorithms can be adapted to many kinds of data. Among other applications, they have been used to find patterns in genetic data, images, and social networks [2]. Different types of models have been used to estimate the continuous representation of words. LDA is a well-known topic modeling tool for exploiting the hidden thematic structures in large archives of text. LDA can become computationally expensive on larger datasets [16]. Since the introduction of LDA, it has been used for complex goals and tasks. One assumption that LDA makes is that the words are not set in a specific order, because it is based on a bag-of-words model. To make LDA robust to a non-exchangeable order of words, Wallach et al. [29] developed a topic model that relaxes the bag-of-words assumption by assuming that topics generate words conditional on the previous word. Griffiths et al. [6] developed a topic model that switches between LDA and a standard hidden Markov model. This model expands the parameter space significantly but exhibits improved language modeling performance. A second assumption of LDA is that documents are not set in a specific order. This assumption will not hold true when the documents have been published over a wide range of time, as in many situations, topics change over time. From this perspective, dynamic topic modeling has been developed [3]. In dynamic topic modeling, the order of documents is retained.

### 2.2 Hash-Tag Clustering

In the second research question, hash-tags are used as true labels to verify the output of the LDA model.A hash-tag is a type of metadata tag used on social networks such as Twitter that allows users to

apply dynamic, user-generated tagging, which makes it possible for others to easily find messages with a specific theme or content. A hash-tag archive is consequently collected into a single stream under the same hash-tag [28] [5]. In several studies, hash-tags have been clustered to analyze contextual semantics [24] [25] [17] [22], such as in a study by Stilo and Velardi in which they created a clustering algorithm for hash-tags based on temporal mining [24]. The researchers clustered hash-tags based on their temporal co-occurrence with other hash-tags. Tsur et al. created an algorithm that leverages users' practice of adding hash-tags to some messages by bootstrapping over virtual non-sparse documents [25]. These non-sparse documents were then represented as vectors in the vector space model. Rosa et al. used hash-tag clusters to achieve the topical clustering of tweets, whereby they compared these effects with each other [22].

### 2.3 Word Embeddings Projections

In the final part of this study, the application that uses Word2Vec on PCA, t-SNE, and UMAP is introduced. Word embeddings are extensively studied on different text sources from Web search engines [18] to sets of Curriculum Vitaes [7]. A combination of PCA projection using Word2Vec and Twitter messages is used in the tutorial of Leightley [14]. Another dimensional reduction method used in this tutorial is t-SNE. t-SNE is a multidimensional scaling method developed by Laurens van der Maaten and is used in dimension reduction approaches [27] [26]. Andrej Karpathy also achieved projection using t-SNE on Twitter data but in combination with JavaScript [12].

## 3 METHODS

In this section, the methodologies applied in the study are explained. LDA is explained in the first subsection, followed by the perplexity evaluation method, which is used to test the LDA model on Twitter data. An explanation of the weighted average of pairwise maximum f-score is then provided. Finally, Word2Vec, PCA, t-SNE, and UMAP, which are used for the word projections, are explained.

### 3.1 Latent Dirichlet Allocation

LDA is a generative probabilistic model for the collection of data types, such as text corpora. The model's structure can be explained as a three-level hierarchical Bayesian model. Each item or document in a collection is treated as a mixture of an underlying set of topics. Each topic is treated as a mixture of an underlying set of term probabilities. In other words, each document is treated as a mixture of different topics. A document is considered to have a set of topics that are allocated by LDA [20]. [20].

### 3.2 Evaluating the Model Using Perplexity

Perplexity indicates how well the model describes a set of documents. In this study, the documents are Twitter messages. Perplexity is an intrinsic evaluation metric and is widely used for language model evaluation. It captures how surprised a model is by new data that it has not seen before, and is measured as the normalized log-likelihood of a held-out test set [4] [19].

$$L(w) = log(w|\phi, \alpha)\Sigma_d logp(w_d|\phi, \alpha)$$

(1)

$w$ = set of Twitter messages to train and test the LDA model.
$\phi$ = the topic matrix
$\alpha$ = the hyperparameter for topic-distribution of documents
$w_d$ = a test set of Twitter messages, which are not used in the training of the model.

The measure traditionally used for topic models is the perplexity of held-out documents $w_d$, defined as shown in Formula 2. This is a decreasing function of the log-likelihood L(w) of the unseen documents $w_d$, normalized by the number of words in the set $words_N$. The lower the perplexity, the better the model [13].

$$perplexity(w_d) = exp(\frac{L(w)}{words_N})$$

(2)

$w_d$ = a test set of Twitter messages, which are not used in the training of the model.
$words_N$ = total number of words in the set.

### 3.3 Clustering Comparison using f-scores for Evaluation

In the second research question, hash-tags are used as true labels (ground truth) to verify the output of the LDA model. The weighted average of the pairwise maximum $f^a$-score is used as an evaluation method in subsection 2. The score is obtained by calculating the fa-score for each hash-tag. The $f^a$-score is based on the precision (formula 3) and recall (formula 4). Precision is defined as the ratio of correct hash-tags in each cluster compared to the total number of hash-tags in that cluster.

Precision is defined for one cluster the ratio of correct hash-tags in the cluster compared to the total number of hash-tags in the cluster.

$$precision(C_i, G_j) = \frac{G_j \uplus C_i}{C_i}$$

(3)

$G_j$ = the ground truth cluster, which is a set of messages, with the same hash-tag attached.
$C_i$ = the LDA cluster, which is a set of messages, assigned to the same LDA topic.

The recall score for a cluster is based on the number of correct hash-tags in the cluster compared to the number of correct hash-tags that should be in the cluster (i.e., the ground truth cluster). The precision of a cluster $C_i$ is calculated as follows, taking into account the ground truth cluster $G_j$:

$$recall(C_i, G_j) = \frac{G_j \uplus C_i}{G_j}$$

(4)

The $f^a$-score is based on the precision and recall and is defined as following:

$$f - score(C_i, G_j) = 2x \frac{recall(G_jC_i) * precision(G_jC_i)}{recall(G_jC_i) + precision(G_jC_i)} \quad (5)$$

Depending on the purpose of the evaluation, the final $f$-score can be calculated in two ways. In this study the overall accuracy of all clusters will be calculated. For the calculation the weighted average of pairwise maximum $f$-score are used (each pairwise maximum $f$-score is weighted by the size of the matching ground truth cluster). The $f^a$-score is formulated as following:

Depending on the purpose of the evaluation, the final $f$-score can be calculated in two ways. In this study, the overall accuracy of all clusters is calculated. For the calculation, the weighted average of pairwise maximum $f$-score is used (each pairwise maximum $f$-score by the size of the matching ground truth cluster). The $f^a$-score is formulated as follows:

$$f^a - score(C_i, G_j) = \frac{f - score(C_i, G_i^m * |G_i^m|)}{\Sigma_{i=1}^n |G_i^m|} \quad (6)$$

## 3.4 Word2Vec

Word2Vec is a two-layer neural network that is designed to process text. Its input is a text corpus and its output is a set of feature vectors for words in that corpus. Word2Vec converts text into a numerical format that can be understood by a machine.

## 3.5 PCA

PCA is mathematically defined as an orthogonal linear transformation that transforms the data into a new coordinate system such that the greatest variance of projected data comes to lie on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on [11].

## 3.6 t-SNE

The goal of t-SNE is to embed high-dimensional data in low dimensions in a way that respects similarities between data points. Nearby points in the high-dimensional space correspond to nearby embedded low-dimensional points, and distant points in high-dimensional space correspond to distant embedded low-dimensional points [14]. t-SNE does not preserve the global data structure, meaning that only distances within clusters are meaningful while similarities between clusters are not guaranteed.

## 3.7 UMAP

UMAP is a novel manifold learning technique for dimension reduction. UMAP can be used to visualize general non-linear dimension reductions. The algorithm assumes three characteristics about the data: (1) the data is uniformly distributed on a Riemannian manifold; (2) the Riemannian metric is locally constant; and (3) the manifold is locally connected. If the data has these characteristics, the manifold can be modeled with a blurry topical structure. The embedding is found by searching for a low-dimensional projection of the data [15].

## 4 EXPERIMENTAL SETUP

This section outlines the experimental setup, explaining the data and implementation of models. In the first section, the dataset is described, followed by the setups of the three sub-questions. Figure 1 presents a flowchart of the connections between the three sub-questions of this study. The blocks represent the steps taken per question.

## 4.1 RNW Media Burundi Twitter Dataset

The dataset that is used in all sub-questions is the RNW Media Burundi Twitter dataset, which comprises Twitter messages of the RNW Media Burundi Twitter platform. The messages are mainly written in French, English, and African languages. The dataset totals 464,845 Twitter messages. The Twitter message is presented in the first column and the hash-tags for that message are shown in a second column. One message can have multiple hash-tags. The Twitter messages were posted between January 2018 and May 2019. The dataset is not missing any values. In Figure 2, the distribution is shown for the number of words per message. From the figure, it can be clearly observed that the message size is small. Figure 3 presents the distribution of the number of characters per message. Most messages contain less than 280 characters. Twitter's founders created Twitter messages to have a maximum length of 160 characters, which permitted a message of 140 characters and 20 characters for a user name. The results exhibit a peak at around 160 characters, which could be related to this initial character limit. In 2018, the maximum character limit was changed to 280 characters [28]. In 2018, the maximum character limit was changed to 280 characters, which is also evident in Figure 3. In Figure 5, a word cloud of the 30 most commonly occurring hash-tags in the dataset is presented. The hash-tags are mainly related to countries, political parties and political events. In Table 1, the most frequent hash-tags in the dataset are explained. A selection of these hash-tags are used as true labels in sub-question 2.

## 4.2 RNW Media DRC Facebook Dataset

In the last subsection, the RNW Media Burundi Facebook dataset is used. This dataset consists of Facebook posts and messages from the RNW Media Facebook platform, which is active in the Democratic Republic of the Congo (DRC). The 91,579 Facebook posts and messages were posted between January 2018 and May 2019. The word cloud shown in Figure 4 displays the most common words in the dataset, among which many French words and the names of decision-makers, such as a president and a lawmaker from the DRC, are present.

## 4.3 Finding Optimal Output Topic Number

The goal of this section is to determine the settings with which LDA achieves the best results for the Twitter dataset. One setting is the number of output topics, which is determined by testing the LDA model with different output settings on the data. Testing the total dataset is costly and unnecessary. Therefore, the different settings are tested on similarly-sized subsets. A topic output size in a range between 2 and 15 is selected because the hypothesis is that the size
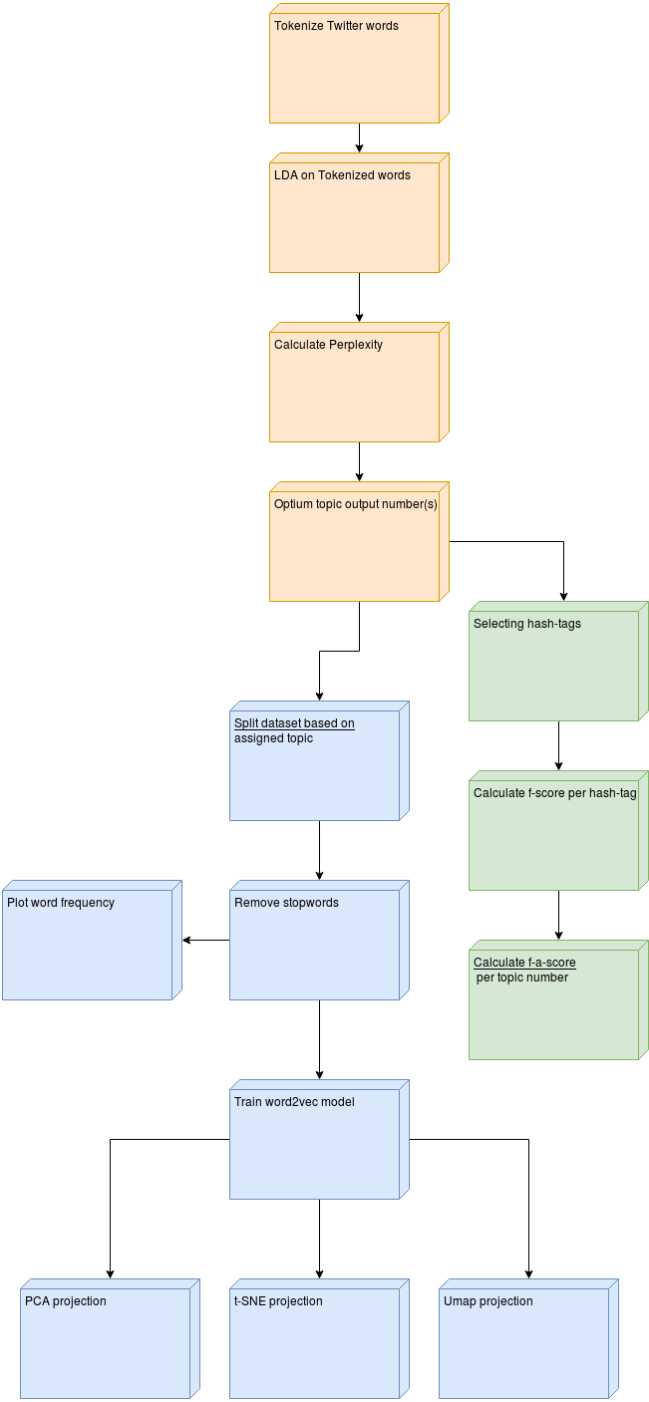
Fig. 1. Flowchart depicting the connections between the three sub-questions of this study. The yellow-colored blocks indicate the steps taken in sub-question 1, which concerns finding the optimal output topic number. The green-colored blocks indicate the steps taken for sub-question 2, which involves determining whether hash-tags are found in the same LDA output topic. The blue colored blocks indicate the steps taken in sub-question 3, which regards visualizing word relationships in Word2Vec embedding. In the third sub-question, the RNW Media Burundi Twitter dataset and the RNW Media Burundi Facebook dataset are used.
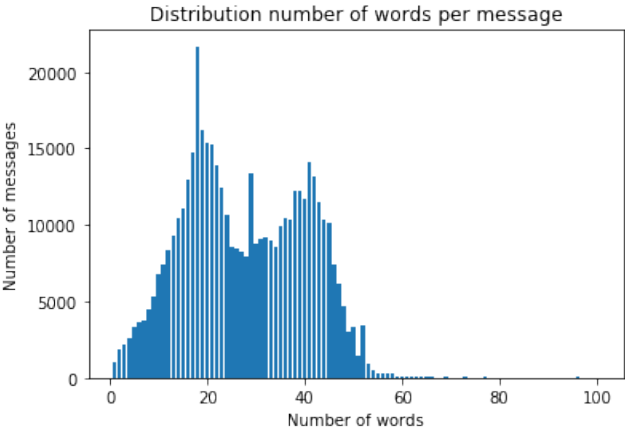
Fig. 2. The distribution of words per message. The largest group of message sizes is of 18 words. Most messages contain between 13 and 44 words. The figure clearly demonstrates that the message size is overall small, compared to documents where LDA have been applied on.



Fig. 3. The distribution of characters per message. Most messages contain less than 280 characters. Twitter's founders created Twitter messages to have a maximum length of 160 characters, which permitted a message of 140 characters and 20 characters for a user name. The results exhibit a peak at around 160 characters, which could be related to this initial character limit. In 2018, the maximum character limit was changed to 280 characters [28].

of the topic output is optimal when it is small. The perplexity score is used as a means of evaluation.

Fig. 4. Word cloud which shows the common words in RNW Media DRC Facebook dataset. The world cloud is generated before the English, French, Spanish and African stop words were removed. Which can been seen by the many French words like comme (as), aussi (also) and donc (therefore). Also decision makers are in this figure like Fayulu (businessman and lawmaker from DRC) and Kabila (President of the DRC from 2001 up to 2018).



Fig. 5. Word cloud with the 30 most common hash-tags in the dataset. The size of the word is equivalent on the number of occurrence in the dataset. The hash-tags are mainly related to country names and political parties. Explanations of the hash-tags can be found in table 1

| Hash-tag | Explanation | Frequency |
| --- | --- | --- |
| Burundi | Country | 111530 |
| Cpi | Burundi political party | 60839 |
| Burundicrisis | | 11478 |
| Burundialerte | Political tension | 11478 |
| Silentmajority | People who do not express their opinions publicly | 10460 |
| Rdc | Country | 6160 |
| Boost | Heroin-derived drug | 4716 |
| Bujumbura | Largest city and main port of Burundi | 4458 |
| Lomasleido | Political related | 3337 |
| Rwanda | Country | 2928 |
| Ndondeza | Burundi political party | 2841 |
| Nkurunziza | President of Burundi since 2005 | 2543 |
| Sindumuja | Organization | 2313 |
| Gbagbo | Ivorian politician | 1617 |
| Bjp | Indian political party | 1483 |
| Affiliate | Marketing strategy | 1456 |
| Referendum2018 | referendum about the presidents power in Burundi | 1454 |
| Jnu | University in India | 1417 |
| Imbonerakure | Political group | 1350 |
| Bsp | Indian political party | 1344 |
| Bemba | Politician in Democratic Republic of the Congo | 1310 |
| Venezuela | Country | 1305 |
| Leyestatutariajep | Law in Venezuela | 1270 |
| Cpa | Certified Public Accountants | 1250 |
| Gitega | Gitega is one of the 18 provinces of Burundi. | 1178 |
| Gabon | Country on the west coast of Central Africa | 1166 |
| Inflation | Sustained increase in the general price level of goods and services in an economy over a period of time | 931 |
| Fatoubensouda | Gambian lawyer and international criminal law prosecutor | 916 |

Table 1. Most common hash-tags in the Burundi Twitter dataset with explanations and their frequency (number of messages with this hash-tag). 65,927 messages did not have a hash-tag attached to them in the dataset. The hash-tags are mainly related to country names and political parties, but references to a university in India and the country of Venezuela and its laws also occur often in the dataset. .

### 4.4 Are Hash-tags found in the same LDA Output Topic

Can LDA distinguish between Twitter messages using hash-tags as true labels? The hypothesis is that messages with common hash-tags have a higher chance of being assigned to the same topic group. For this question, the $f^a$-score (formula 6) is used as a means of evaluation [9] [10]. This score is based on the precision (formula 3) recall and (formula 4). Since precision and recall necessarily depend on the notion of true classes, a selection of hash-tags are used as true classes. An LDA topic is assigned to a hash-tag based on the largest number of labeled hash-tags within that LDA topic. The recall and precision are used to calculate the $f$-score, which is only for one specific hash-tag (cluster). To estimate the score of all hash-tags (clusters), the $f^a$-score is used. Explanations of the selected hash-tags are provided in Table 1.

### 4.5 Visualizing word Relationships in Word2Vec Embedding

In this subsection, relationships between important words within a generated topic are visualized. This is achieved by selecting the most frequent words per topic. These words are visualized in a projection by the Word2Vec model. Three different dimension reduction methods are used, which are PCA, t-SNE, and UMAP. To obtain these words, English, French, Spanish, and African stop words are removed, and common Twitter messages are filtered out.
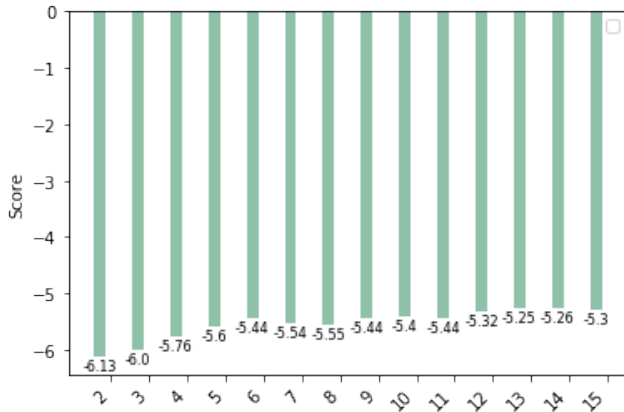
Fig. 6. Perplexity per topic number. a lower perplexity suggests a better fit. The lowest perplexity is achieved using k = 2. Which results in perplexity score of -6.13.

## 5 RESULTS AND DISCUSSION

This section contains the results of conducted experiments to answer the research questions. The study is divided into three sub-questions which will be answered separately in sub-sections.

### 5.1 Finding the Optimal Output Topic Number

What is the optimal number of topics that can be generated on a set of Twitter messages using LDA The goal of this section is to determine the settings with which LDA achieves the best results for the dataset. To evaluate the LDA model applied in this study, perplexity is used, which is the common evaluation method for LDA [4] [23] [21].

The number of topics that are generated by the LDA algorithm is written as N. To answer this research question, LDA is applied multiple times to the entire dataset. In each run, the number of generated topics is set to a specific number N. If keywords are present in multiple topics, this indicates that N is too high. In figure 6 the perplexity per topic number generated by LDA is shown. The *perplexity* of held-out documents $w_d$ is defined as shown in Formula 2. This is a decreasing function of the log-likelihood L(w) of the unseen documents $w_d$; the lower the perplexity, the better the model. Therefore a lower perplexity suggests a better fit [4]. The lowest perplexity is achieved using N = 2, which results in a perplexity score of -6.13.

### 5.2 Are Hash-Tags Found in the Same LDA Output Topic

Can LDA distinguish between Twitter messages using hash-tags as true labels? Every Twitter message contains one or more hash-tags. The hash-tags can be used to find messages with a specific theme or content. The hypothesis for this question is, therefore, that messages that contain common hash-tags have a higher chance of being assigned to the same LDA output group.

For this question, the $f^a$-score is used as evaluation metric. This score is calculated by a single cluster f-score. Since precision and recall necessarily depend on the notion of a true class, a selected

group of hash-tags is used as true classes. In Table 1, explanations can be found for the selected hash-tags.

In Table 2, hash-tags with a topic size of 2 and 3 are shown. For the topic size of 2, all $f$-scores are between 0.594 and 0.660. Which means that all $f$-scores are above random 0.5, because there is a change of 1 out of 2. Hash-tag "Jnu," used for a university in India, has the highest score. There is no direct link between this university and Burundi, which could explain this high score. For a topic output size of 3, the $f$-scores are smaller and the random score is lower. For comparison, one hash-tag is used. If a message has multiple hash-tags, the first hash-tag is selected. In Table 4, the $f^a$-scores for topic sizes 2, 3, 4, and 5 are provided and compared with a random score. The random score is the chance that a hash-tag is assigned to a cluster. The biggest difference between the $f^a$-score and the random score is for a topic size of 3, while the smallest difference between the $f^a$-score and the random score is for a topic size of 2.

| Hash-tag | Topic size 2 | Topic size 3 |
|---|---|---|
| Rwanda | 0.594 | 0.506 |
| Ndondeza | 0.607 | 0.521 |
| Nkurunziza | 0.616 | 0.531 |
| Sindumuja | 0.600 | 0.541 |
| Gbagbo | 0.637 | 0.564 |
| Bjp | 0.622 | 0.554 |
| Affiliate | 0.611 | 0.554 |
| Referendum2018 | 0.641 | 0.573 |
| Jnu | 0.660 | 0.556 |
| Imbonerakure | 0.657 | 0.600 |
| Bsp | 0.597 | 0.543 |
| Bemba | 0.596 | 0.584 |
| Venezuela | 0.658 | 0.569 |
| $f^a$-score | 0.619 | 0.547 |

Table 2. For the selected number of hash-tags, $f$-scores are given using a topic number of 2 and 3. For the topic size of 2, all scores are between 0.594 and 0.660, which means that all $f$-scores are above random (0.5). The $f$-scores for topic size 3 are lower, which is expected because the score is applied to an additional cluster.

| Hash-tag | Topic size 4 | Topic size 5 |
|---|---|---|
| Gbagbo | 0.392 | 0.280 |
| Bjp | 0.386 | 0.273 |
| Affiliate | 0.421 | 0.283 |
| Referendum2018 | 0.400 | 0.286 |
| Jnu | 0.409 | 0.286 |
| Imbonerakure | 0.418 | 0.294 |
| Bsp | 0.412 | 0.284 |
| Venezuela | 0.410 | 0.294 |
| $f^a$-score | 0.406 | 0.285 |

Table 3. For the selected number of hash-tags, the $f$-score and $f^a$-score are given using a topic number of 4 and 5. These scores are lower than the scores when using a topic number of 2 or 3, but the $f^a$-score for topic 4 (0.406) is a better score than the random score of 0.25%, because there is an increase of 62.0%.

| Topic size | $f^a$ score | random score | percent > random |
|---|---|---|---|
| 2 | 0.619 | 0.5 | +0.238 |
| 3 | 0.547 | 0.34 | +0.608 |
| 4 | 0.406 | 0.25 | +0.62 |
| 5 | 0.285 | 0.20 | +0.425 |

Table 4. The $f^a$-scores for topic sizes 2, 3, 4, and 5 when compared with the random score. The random score is the chance of a hash-tag being assigned to a cluster. Topic size 3 has the biggest difference between the $f^a$-score and the random score, while topic size 2 has the smallest difference between the $f^a$-score and the random score.

## 5.3 Visualizing Word Relationships in Word2Vec Embedding

*5.3.1 RNW Media Burundi Twitter Dataset.* To identify relationships between the words in the RNW Media Burundi Twitter dataset, the most frequent words per topic are selected. In Figure 7, the frequencies of the 25 most common words in topic 1 are shown. Many financial and politics-related words are in the top 25. For the second topic, the same steps are taken. In Figure 8, the fre quency of the most common words for topic 2 is shown. Many of these words are related to crime and humans, such as *crimes*, *militares* (soldiers), *humanitana* (humanitarian), *victimes* (victims), *police*, *regimes*, and *droits* (rights). Another observation is that topic 2 includes more Spanish and French words than topic 1. The most frequent words in both topics are shown in the PCA projection of the Word2Vec model in Figure 9. On the left of the projection is a green cluster; words such as *humanitana*, *victimes*, *commune*, *president*, and *regime* are grouped in close proximity to one another. At the bottom of the picture is a blue cluster with words such as fiscal and hourly. In Figure 10, a t-SNE projection is shown. The most green words occur on the left side of the projection, while the blue words are more spread out. The words *police*, *regime*, and *crimes* are clustered together. The UMAP projection is shown in Figure 11. Besides *chambre* (room) and *politique* (politics), the green words are clustered together. *Humanitaria*, *victimes*, *aseguro* (assure) and *bloqueen* (block) are closely grouped.



Fig. 7. Frequency of the most common words in Twitter messages that are assigned to topic 1. These words are selected from all of the Twitter messages that are assigned to topic 1. The 25 most frequently occurring words are shown in this figure on the x-axis, while the number of occurrences is shown on the y-axis. *RT*, a Russian TV channel, is the most common word in this topic. In the third position is the word *rose*. Besides that, many finance-related words occur often in this topic, such as *earnings*, *wages*, and *inflation*. Political words, such as *president*, *regime*, *cpi*, and *party*, can also be found in this topic.
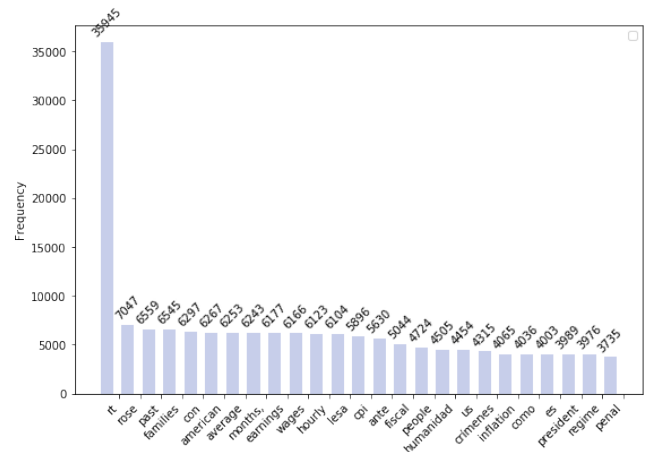


Fig. 8. Frequency of the most common words in Twitter messages that are assigned to topic 2. These words are selected from all of the Twitter messages that are assigned to topic 2. The 25 most frequently occurring words are shown on the x-axis, while the number of occurrences is shown on the y-axis. *RT*, a Russian television channel, is the most common word in this topic. This topic contains more French and Spanish words compared to topic 1. Crime and human-related words occur often in this topic, such as *crimes militares*, *humanitaria*, *victimes*, *police*, *regimes*, and *droits*.
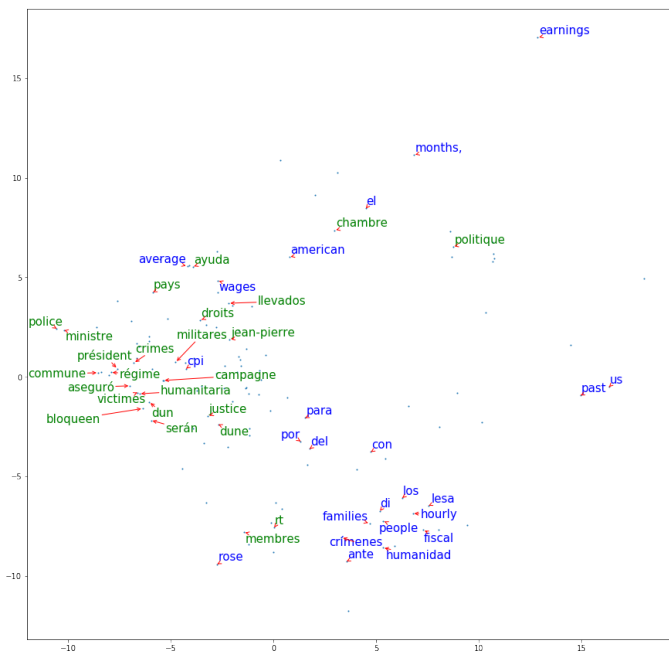
Fig. 9. PCA projection of a Word2Vec model. The 25 most frequently occurring words per topic are highlighted in green for topic 1 and blue for topic 2. The positions are shown with red arrows.
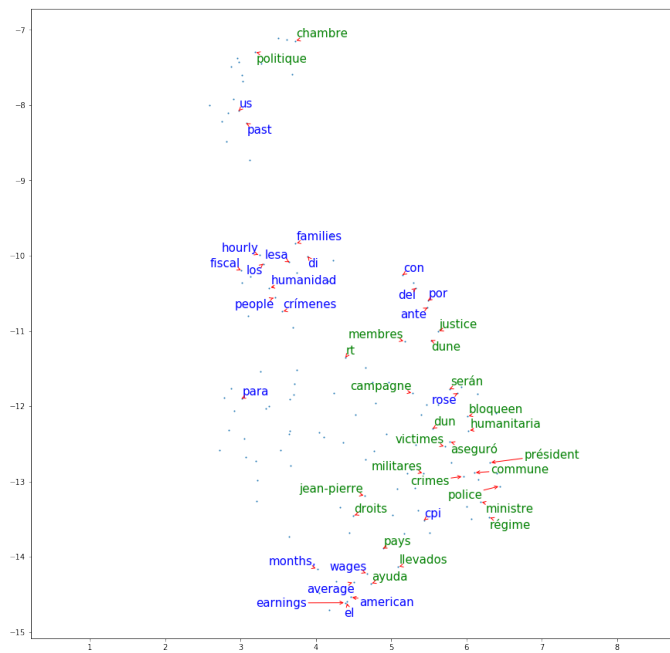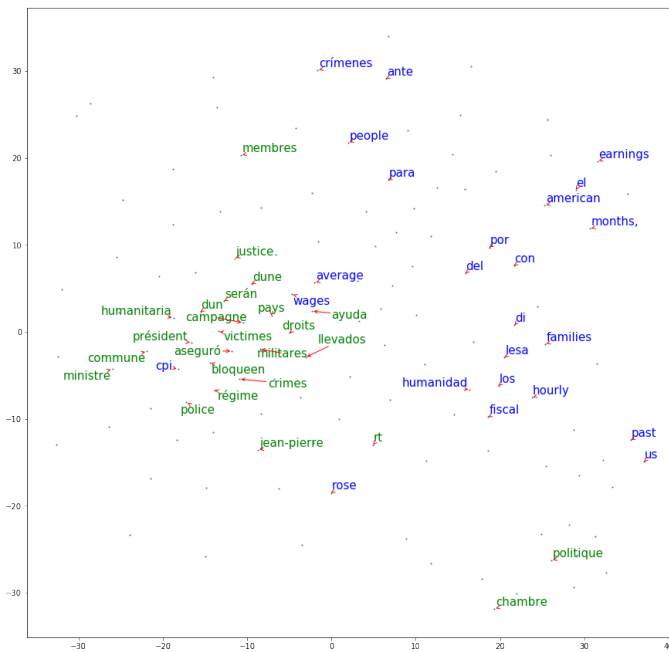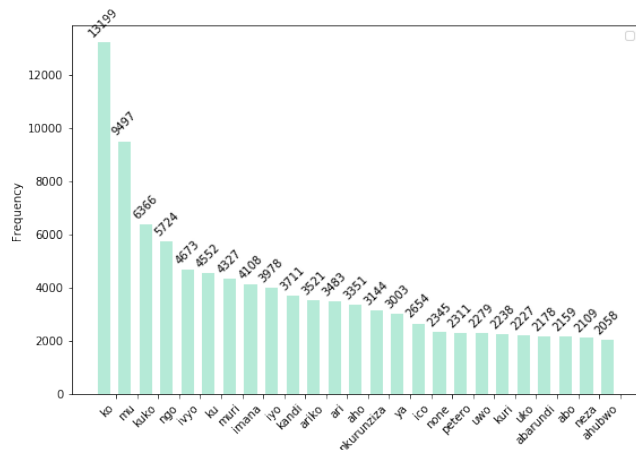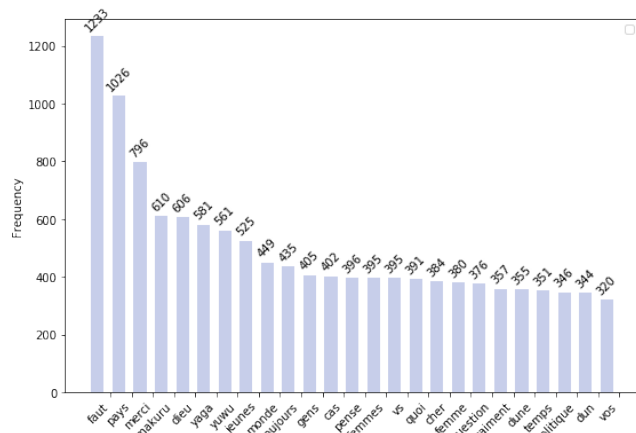


Fig. 11. UMAP projection of a Word2Vec model. The 25 most frequently occurring words per topic are highlighted in green for topic 1 and blue for topic 2. The positions are shown with red arrows.
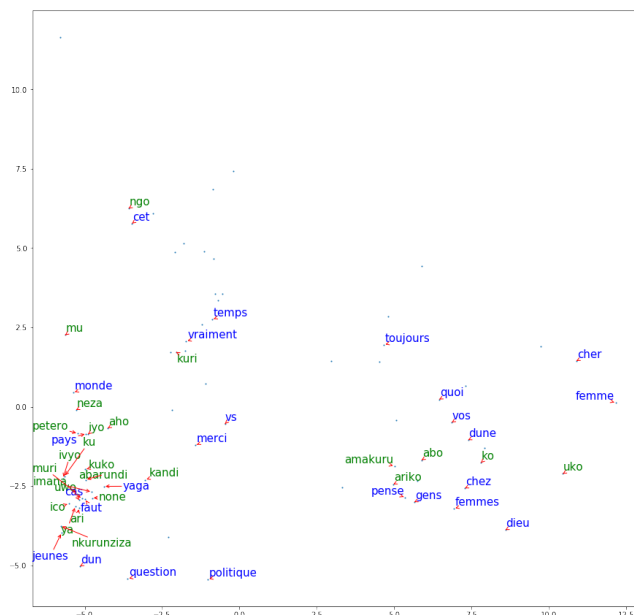


Fig. 10. t-SNE projection of a Word2Vec model. The 25 most frequently occurring words per topic are highlighted in green for topic 1 and blue for topic 2. The positions are shown with red arrows.

Fig. 12. Frequency of the most common words in Facebook messages that are assigned to topic 1. These words are selected from all of the Twitter messages that are assigned to topic 1. This topic contains mostly African words.



Table 5. requency of the most common words in Facebook messages that are assigned to topic 2. These words are selected from all of the Twitter messages that are assigned to topic 2.



Fig. 13. Scatter plot of a PCA projection of a Word2Vec model. The 25 most frequently occurring words per topic are highlighted in green for topic 1 and blue for topic 2. The positions are shown with red arrows.)
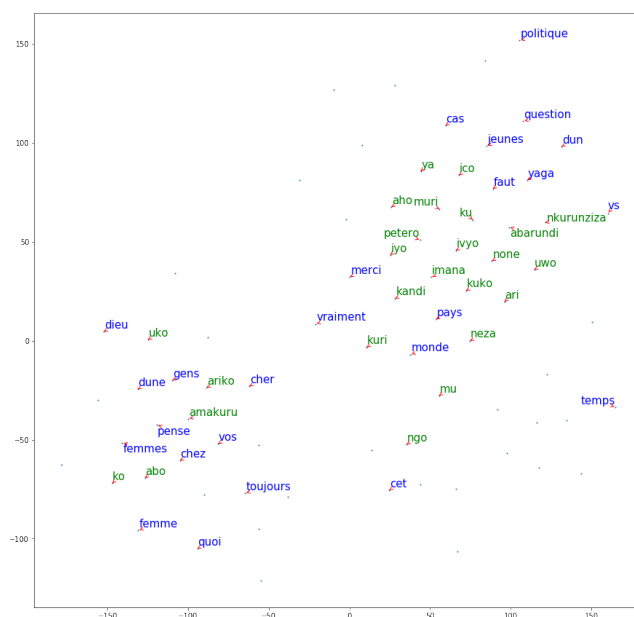


Fig. 14. Scatter plot of a t-SNE projection of a Word2Vec model. The 25 most frequently occurring words per topic are highlighted in green for topic 1 and blue for topic 2. The positions are shown with red arrows. .
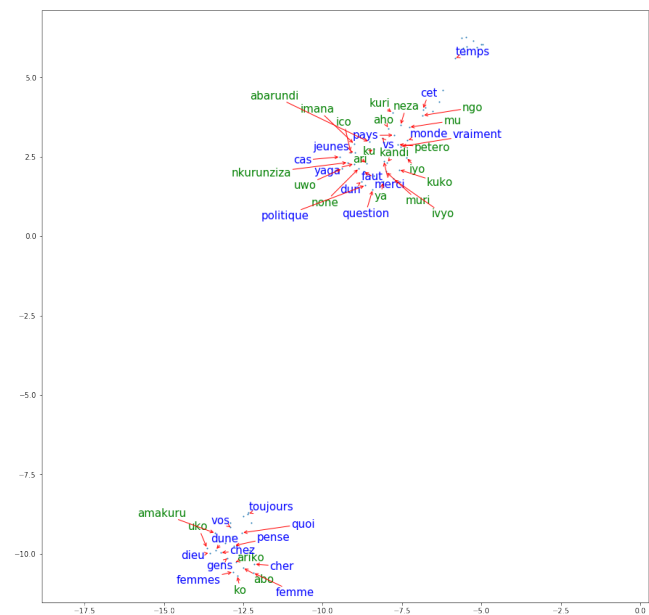
*5.3.2 RNW Media DRC Facebook Dataset.* Besides a Twitter forum, RNW Media also operates on Facebook in the DRC. Posts and messages from this Facebook platform are stored in the RNW Media DRC Facebook dataset. The steps described in section 4.5 are taken to identify word relationships. Two topics are created. The frequency of the most common words for topic 1 are shown in Figure 12 and for topic 2 in Figure 5. Topic 1 contains many African words, while topic 2 contains more French words, such as *femmes* (woman), *chère* (dear), and *pense* (thought). The most frequent words of both topics are shown in the PCA projection of a Word2Vec model in Figure 13. In this projection, words are clustered on the left side.In this figure, it can be seen that there is a relationship between the words from both topics which contain both different languages. For example,

Fig. 15. Scatter plot of a UMAP projection of a Word2Vec model. The 25 most frequently occurring words per topic are highlighted in green for topic 1 and blue for topic 2. The positions are shown with red arrows.

there is a close relationship between *amakuru* (major) and *pense* (thought).

In Figure 14, a t-SNE projection is shown. In this projection, the words are spread over the projection. The UMAP projection is shown in Figure 15, which shows two clearly separated clusters. The upper cluster contains some politics-related words, such as *politique* and *Nkurunziza* (President of Burundi since 2005).

### 5.3.3 Future improvements.

In further research, different Twitter datasets should be used for testing. For the Burundi Twitter dataset, two topics would yield solid results. It could be possible that in a larger dataset, more topics could be generated by the LDA algorithm. In the Burundi Twitter dataset, mostly English, French, and Spanish words occur. These languages are common used languages in topic modeling, for example using Python as programming language. RNW Media is also working with languages with more complex structures, such as Arabic, Hindi, and Rundi. A future study could be continued on datasets that contain these languages. In addition to a word embedding projection, the word similarity could also be visualized in a heat-map. The Word2Vec similarity function can be used to obtain a score between words. On the x-axis and y-axis of the map, the words in similar order are positioned. The color intensity would represent the similarity intensity.

### CONCLUSION

The goal of this thesis was to gain deeper insights into the Twitter and Facebook post/comments data on the social media platform of a humanitarian organization. LDA was used to create grouped word topics. This approach is unique because it has not yet been used on text with character limits, like Twitter messages. Output size was found to be optimal when it is small, which confirms the hypothesis. In addition, it was shown that messages that contain common hash-tags have a higher chance of being assigned to the same LDA output group. The word topics were subsequently projected using a combination of Word2Vec with PCA, t-SNE, and UMAP. These projections were presented in a two-dimensional way, whereby links between words from mixed languages within and between topics can be obtained. The UMAP projection resulted in the most informative visualizations.

### REFERENCES

[1] About us. (2019, 19 August). Accessed at 19 August 2019, van https://www.rnw.org/about-us/
[2] Blei, D. M. (2012). Probabilistic topic models. Communications of the ACM, 55(4), 77. https://doi.org/10.1145/2133806.2133826
[3] Blei, D. M., Lafferty, J. D. (2006). Dynamic topic models. Proceedings of the 23rd international conference on Machine learning - ICML '06, 113–120. https://doi.org/10.1145/1143844.1143859
[4] Blei, D. M., NG, A., Jordan, M. (2003). Latent Dirichlet Allocation. Journal of Machine Learning Research, 3, 993–1022. Consulted from http://jmlr.org/papers/volume3/blei03a/blei03a.pdf
[5] Doctor, V. (2012, 12 June). What Characters Can A Hashtag Include? Accessed at 4 May 2019, van https://www.hashtags.org/featured/what-characters-can-a-hashtag-include/
[6] Griffiths, T. L., Steyvers, M., Blei, D. M., Tenenbaum, J. B. (2005). Integrating Topics and Syntax. Advances in Neural Information Processing Systems, 17, 537–544. Consulted from https://cocosci.princeton.edu/tom/papers/composite.pdf
[7] Hansen, C., Tosik, M., Goossen, G., Li, C., Bayeva, L., Berbain, F., Rotaru, M. (2015). How to get the best word vectors for resume parsing. SNN Adaptive Intelligence/Symposium: Machine Learning. Consulted from https://lenabayeva.files.wordpress.com/2015/03/snn-textkernelposter-2015.pdf
[8] He, J. (2012). Exploring topic structure. ACM SIGIR Forum, 46(1), 84. https://doi.org/10.1145/2215676.2215690
[9] Javed, A., Lee, B. S. (2017). Sense-Level Semantic Clustering of Hashtags. Information Management and Big Data, 1–16. Consulted from https://arxiv.org/pdf/1802.03426v2.pdf
[10] Javed, A., Lee, B. S. (2018). Hybrid semantic clustering of hashtags. Online Social Networks and Media, 5, 23–36. https://doi.org/10.1016/j.osnem.2017.10.004
[11] Jolliffe, I. T. (2006). Principal Component Analysis. x: Springer New York.
[12] Karpathy, A. (2014, 2 juli). Visualizing Top Tweeps with t-SNE, in Javascript. Accessed at 10 June 2019, van https://karpathy.github.io/2014/07/02/visualizing-top-tweeps-with-t-sne-in-Javascript/
[13] Kim, A. (2018, 11 oktober). Perplexity Intuition (and its derivation). Geraadpleegd op 18 januari 2020, van https://towardsdatascience.com/perplexity-intuition-and-derivation-105dd481c8f3
[14] Leightley, D. (2018, 1 August). Visualizing Tweets with Word2Vec and t-SNE, in Python. Accessed at 5 May 2019, van https://leightley.com/visualizing-tweets-with-word2vec-and-t-sne-in-python/
[15] McInnes, L., Healy, J., Melville, J. (2018). UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. x. Consulted from https://arxiv.org/pdf/1802.03426.pdf
[16] Mikolov, T., Chen, K., Corrado, G. S., Dean, J. (2013). Distributed Representations of Words and Phrases and their Compositionality. Advances in neural information processing systems, 1–9. Consulted from http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf
[17] Muntean, C. I., Morar, G. A., Moldovan, D. (2012). Exploring the Meaning behind Twitter Hashtags through Clustering. Business Information Systems Workshops, 231–242.
[18] Nalisnick, E., Mitra, B., Craswell, N., Caruana, R. (2016). Improving document ranking with dual word embeddings. Proceedings of the 25th International Conference Companion on World Wide Web, 83–84. Consulted from https://www.microsoft.com/en-us/research/wp-content/uploads/2016/04/pp1291-Nalisnick.pdf
[19] Pleplé, Q. (2013 June). Perplexity To Evaluate Topic Models. Accessed at 11 May 2019, van http://qpleple.com/perplexity-to-evaluate-topic-models/
[20] Pritchard, J. K., Stephens, M., Donnelly, P. (2000). Advances in neural information processing systems. Genetics, 155, 945–959.
[21] Röder, M., Both, A., Hinneburg, A. (2015). Exploring the Space of Topic Coherence Measures. Proceedings of the Eighth ACM International Conference on Web Search and Data Mining - WSDM '15, 399–408.

https://doi.org/10.1145/2684822.2685324

[22] Rosa, K. V., Shah, R., Lin, B., Gershman, A., Frederking, R. (2011). Topical Clustering of Tweets. Proceedings of the ACM SIGIR: SWSM 63. Consulted from http://www.cs.cmu.edu/ kdelaros/sigir-swsm-2011.pdf

[23] Stevens, K., Kegelmeyer, P., Andrzejewski, D., Buttler, D. (2012). Exploring Topic Coherence over many models and many topics. Association for Computational Linguistics, 952–961. Consulted from https://www.aclweb.org/anthology/D12-1087.pdf

[24] Stilo, G., Velardi, P. (2014). Temporal Semantics: Time-Varying Hashtag Sense Clustering. Lecture Notes in Computer Science, 563–578.

[25] Tsur, O., Littman, A., Rappoport, A. (2012). Scalable multi stage clustering of tagged micro-messages. Proceedings of the 21st international conference companion on World Wide Web - WWW '12 Companion, 621–622. https://doi.org/10.1145/2187980.2188157

[26] van der Maaten, L. (2014). Accelerating t-SNE using Tree-Based Algorithms. Machine Learning, 15, 3221–3245. Consulted from http://www.jmlr.org/papers/volume15/vandermaaten14a/vandermaaten14a.pdf

[27] van der Maaten, L., Hinton, G. (2012). Visualizing non-metric similarities in multiple maps. Machine Learning, 87, 33–55. Consulted from https://link.springer.com/content/pdf/10.1007/s10994-011-5273-4.pdf

[28] Wagner, K. (2017, 7 November). Twitter's 280 character tweets are now available for everyone. Accessed at 3 June 2019, van https://www.vox.com/2017/11/7/16615914/twitter-longer-tweets-280-characters-update-available-everyone

[29] Wallach, H. M. (2006). Topic modeling. Proceedings of the 23rd international conference on Machine learning - ICML '06, 23, 977–984. https://doi.org/10.1145/1143844.1143967