

SocialSOM: Topic Detection on Twitter by Organizing Tweets on User Similarity

Bernardo Simões

Technical University of Lisbon - Taguspark Campus,
Av. Prof. Doutor Aníbal Cavaco Silva — 2744-016 Porto Salvo, Portugal
bernardo.simoes@ist.utl.pt
<http://www.ist.utl.pt/en/>

Abstract. *In this report we present a new way to find topics on Twitter, by leveraging the social network user relations as a way to find clusters of topics based on user relations. In order to achieve this we used the Neural Network algorithm Self-Organizing Maps as clustering algorithm. In the end we provide a website running our trained SOM that can categorize the tweets of any user that logs in to it with its Twitter account*

Keywords: topic detection, twitter, Self-Organizing maps, classification, clustering

1 Introduction

With the evolution of social network websites like Facebook and Twitter, the amount of pertinent content about a specific issue is increasing dramatically, which calls for new ways to make sense and catalog this data. The usage of social networks for branding quality and on-line marketing is specially compelling since 19% of all tweets [10] and 32% [19] of blog posts are about brands or products. On the other hand, finding topic sensitive information on social networks is extremely complicated due to the fact that documents have very little content, slang vocabulary and orthographically mistakes or abbreviations.

The data presented in sites like Facebook or Twitter has proven its value in papers like Asur and Huberman [1] where it is possible to predict with high precision the value of a movie box office weeks before it debuts, through real time monitoring of the velocity of usage of hashtags referencing debuting movies.

The academic and enterprise world is now starting to look at Machine Learning for new ways to achieve revenue and visualize data representing the way the world works. It is not strange to see that the Machine Learning course at Stanford is the one with more students enrolling this year ¹ with more than 760 students enrolled.

¹ <http://www.forbes.com/sites/anthonykosner/2013/12/29/why-is-machine-learning-cs-229-the-most-popular-course-at-stanford/>

With emerging new techniques like Deep Learning [4] which focuses on abstract representations that can lead to more useful representations, one example of this kind of work is [15] Building High-level Features Using Large Scale Un-supervised Learning” where a 9-layered locally connected sparse autoencoder with pooling and local contrast normalization on a large dataset of images (the model has 1 billion connections, the dataset has 10 million 200x200 pixel images downloaded from the Internet) trained using model parallelism and asynchronous SGD on a cluster with 1,000 machines (16,000 cores) during three days. Which achieved 81.7 percent accuracy in detecting human faces, 76.7 percent accuracy when identifying human body parts and 74.8 percent accuracy when identifying cats.

Even though a lot of solutions have arisen in order to automate real time searches, topic categorization and many other data intensive tasks, Twitter still uses humans in order to deliver ads to trending queries states Edwin Chen’s Data Scientist Responsible for ads quality at Twitter. On his blog post ² Edwin describes the process of Twitter to deliver real time adds to trending queries, the main problems that arise in the Twitter platform in order to identify rising topics are mainly:

- The queries people perform have never before been seen, so it’s impossible to know beforehand what they mean.
- Since the spikes in search queries are short-lived, there’s only a short window of opportunity to learn what they mean.

This means that when an event happen, people immediately come to Twitter in order to know what is happening in a determined place in real time. Twitter solves this issue by monitoring which queries are currently popular in real time, using a Storm topology ³ and after the queries are identified, they are sent to a Thrift API ⁴ that dispatches the query to Amazon’s Mechanical Turk ⁵ service where real people will be asked a variety of questions about the query.

Social Media Analytics is another raising topic which draws from Social Network Analysis, Machine Learning, Data Mining, Information Retrieval (IR), and Natural Language Processing (NLP). As stated by Melville et al. in [19] 32% of the 200 million bloggers world wide blog about opinions on products and brands, 71% of the 625 million active Internet users actually read blogs and more importantly that 78% of respondents put their trust in the opinion of other consumers where only only 57% of consumers trust advertising in traditional media and even worst only 34% of consumers put their trust in such advertising. This kind of data drives companies to Social Media Analytics in a way to know what people are saying on the web about their companies and products. This new worry

² <http://blog.echen.me/2013/01/08/improving-twitter-search-with-real-time-human-computation/>

³ <http://storm-project.net/>

⁴ <http://thrift.apache.org/>

⁵ <https://www.mturk.com/mturk/>

has brought to life a lot of new startups like Sumal⁶ or ThoughtBuzz⁷ but also solutions from the old players like IBM⁸ and SAS⁹

Its also important to notice that in the last few years Data Science/Analysis has been a trending topic, mostly due to the fact that big dot-com companies have been making lots of money through exploiting user specific information in order to deliver ads and sell products. No wonder that if you look that in the top 10 ebooks sold by O'Reilly throughout 2013, four are about data science¹⁰.

In this project we will focus on using an unsupervised learning technique based on neural networks named Self-Organizing Maps [14] in order to detect topics in Twitter posts, by using the Social Network users as base neurons for clustering. After the network is trained, it will be possible to categorize tweets in real time. This approach will be better described in subsection 1.1.

This document is organized as follows. First will be dedicated to explain some basics concepts like Document CLustering and specifically Self-Organizing Maps in Section 2. Further in, Section 3 will be dedicated to the state of the art solutions related not only to topic detection but also to twitter data analysis and Self-Organizing Maps. In Section 4 Architecture of the purposed solution and at Section 5 it will be discussed how to evaluate results achieved. Finally we will this report by referencing some possible future work and with a brief conclusion at Section 6.

1.1 Objectives

The objective of this project is clear: finding topics on Tweets by analyzing their textual specific characteristics, like number of characters in a tweet, hashtags, is a retweeted¹¹, etc., and contextualize the social network involving the person that did the tweet by retrieving user specific profile information.

After building the dataset, it will be needed to train the Self-Organizing map neural network in order to be able to have the network ready for clustering classification of each future tweet that will be added to it. After the SOM training it will be necessary to categorize the clusters in order to know which topic they belong to. When this step ha ended it will be possibly to get new tweets categorized on the moment they enter the network without further delay.

Finally it will be presented a website where a user will be able to login with his twitter account and see all of his tweets being categorization in the moment they are clustered. After all the user tweets are clustered there will a graphical presentation of the user twitter usage where it will be possible to see multiple statistical information such as the the topics a user is more interested in and his

⁶ <https://sumall.com/>

⁷ <http://www.thoughtbuzz.net/>

⁸ <http://www-01.ibm.com/software/analytics/solutions/customer-analytics/social-media-analytics/>

⁹ <http://www.sas.com/software/customer-intelligence/social-media-analytics.html>

¹⁰ http://shop.oreilly.com/category/deals/bestoforeillydotd.do?code=DEAL&cmp=tw-nabooksvideosinfoauthornote_best_of_2013

¹¹ retweet is when a user shares a tweet that is not his

own twitter Self-Organizing map network of topics with his friends, the wireframes of the website can be seen in the attachments Section A in Figure 10.

2 Basic Concepts

In this Section we will start by generally describing what Clustering is and how it works then we will outline how Self-Organizing maps [14] function, which is the Document Clustering algorithm used on this project.

2.1 Document

Document clustering is an optimal division of documents into categories without prior knowledge of the data that is being organized, based only on the similarity between them. Due to the fact that no prior knowledge of the data has to be known Document Clustering is labeled as Unsupervised Machine Learning.

Yuan-Chao Liu et Al [17] asserted that Document Clustering can be used in a variety of Computer Science fields, such as:

- Natural Language Preprocessing.
- Automatic Summarization.
- User preference mining.
- Improving text classification results.

There are two main types of Document Clustering, Hard Clustering and Soft Clustering. In Hard Clustering one document can only belong to one cluster, while in Soft Clustering one document can belong to multiple clusters.

In regard to document categorization Springorum et al. [22] performed clustering with SOMs [14] while identifying polysemous German Propositions. They used regular SOMs to create multiple clusters and used Centroid-Based or Proposition-based softening to create Soft Clusters from the Hard Clusters.

The clustering process usually works as described in 1 In the first, step a data set must be provided in order to cluster the documents. The second step is where non relevant words are removed from the documents. Kang [11] proves that improves clustering. Another way to extract keywords is to differentiate text features by analyzing the document corpora. For example if the dataset is composed from HTML or XML documents it is possible to identify more relevant features due to the characteristics of the markup. The fourth step is characterized by converting the keywords of each document into vectors, the most common model used for this task is VSM (Vector Space Model). In VSM, each vector dimension means one detected keyword and each document is represented by the vector of keywords in the feature space. This process an is described in Figure 2.

There many clustering algorithms. K-means works by randomly selecting k documents as the cluster centroids, then assigning each document to the nearest centroid, and finally recalculate the the centroid with new added documents.

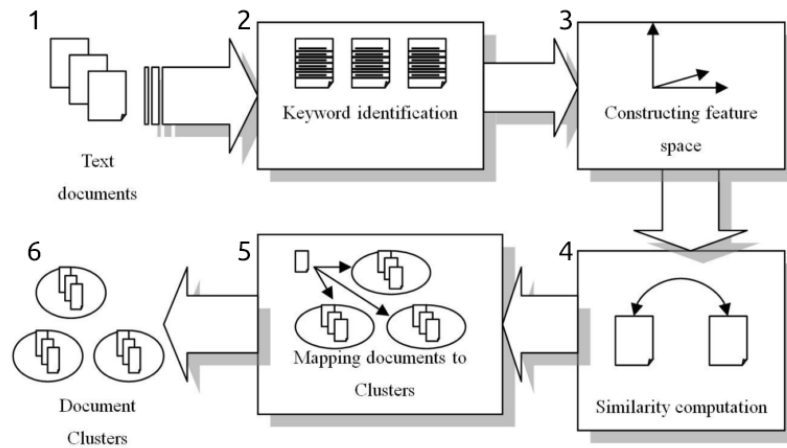


Fig. 1. Text Clustering Main Framework from Dozono [9]

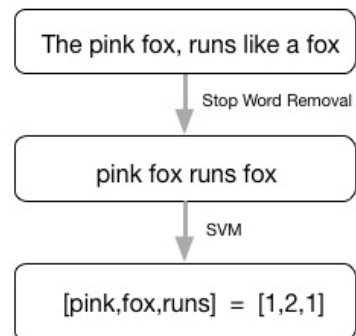


Fig. 2. Stop word removal and transformation to Vector Space Model

2.2 The Self-Organizing Map

The Self-organizing map, or SOM, is a kind of recurrent artificial neural network that has the desired property of topology preservation which mimics the way the cortex of highly developed animals brains work.

As Bação et al. [2] describes, the basic idea behind SOM is to map the data patterns into an n -dimensional grid of neurons or units. That grid is also known as the output space, as opposed to the initial space also called input space, where the input patterns are. Both spaces can be seen in Figure 4.

SOMs work similar to the way that is thought that the human brain works, by having a set of neurons that through learning experience specialize in the identification of certain types of patterns. These neurons are responsible for categorizing the input patterns for which they are responsible. Nearby neurons will be organized by similarity which will cause that similar patterns will activate similar areas of the SOM. With a topology preserving mapping, SOM organizes the information spatially where similar concepts are mapped to adjacent areas. The topology is preserved in a sense that, as far as possible, neighborhoods are preserved through the mapping process. Neurons are displayed in an N dimensional grid, generally rectangular, but other structures are possible, such as hexagonal or octagonal. The grid of neurons, also called output space, can be divided in neighborhoods, where neurons responsible for the same kind of input reside. In SOM, neurons will have the same amount of coefficients as the input patterns and can be represented as vectors through the VSM model described earlier in Section 2.1.

Before describing the algorithm it is important to define two key aspects of the SOM, the learning rate and quantization error. The learning rate is a function that will be decreased in order to converge to zero, it will be applied to winning neurons and their neighbors in order for them to move toward the corresponding input pattern. Quantization Error is the distance between a given input pattern and the associated winning neuron, it describes how well neurons represent the input pattern. The radius of the neighborhood around the winner neuron is particularly relevant to the topology of the SOM, deeply affecting the unfolding of the output space as stated by [2].

The learning phase is characterized by the training algorithm, which works the following way:

- Neurons can be initialized randomly or it is possible to select a specific initialization.
- Given an input pattern, calculate the distance between the input pattern and every neuron on the network.
- The winning neuron will be the closest neuron to the input pattern.
- The neuron will move towards the data pattern at a given learning rate, in order to improve his representation as can be seen in Figure 3.
- Neighbor neurons will also improve their representation in order to keep the network progressively organized as can be seen in Figure 4.

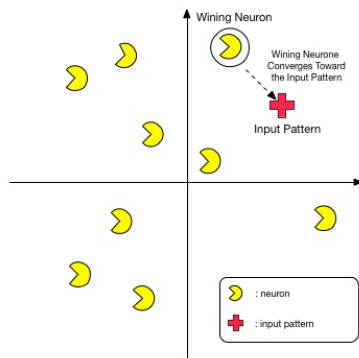


Fig. 3. Winning neuron converging at learning rate

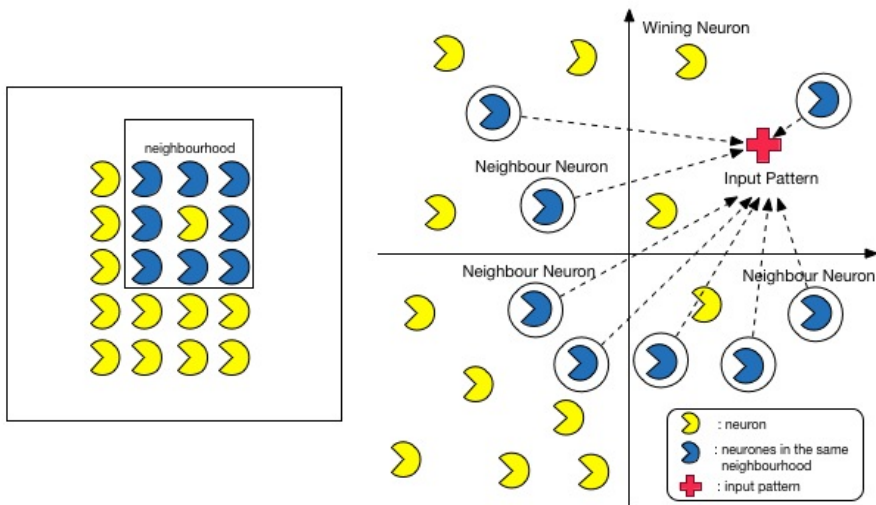


Fig. 4. On the left the output space neighbor, on the right the neighbors of the winning neuron converging

After the algorithm converges, the prediction phase starts. On the prediction phase new input patterns can be quickly assigned to the SOM, without need to apply the learning rate to the winning neuron and his neighbors. Thus it very easy and fast to classify new data now.

In order to visually interpret the result of the SOM U-matrices may be used as stated by [2]. The U-matrix is a representation of the SOM in which distances, in the input space between neurons is represented using a gray scale.

The advantages of using SOM is data noise immunity, easy to visualize the data, and parallel processing.

3 Related Work

This Section provides insight of work done in multiple research areas that are related to the project. In subsection 3.1 will be described multiple work done using Self-Organizing maps. Subsection 3.2 is dedicated to work done on topic detection on the social network Twitter ¹²

3.1 Self-Organizing Maps

Self-Organizing maps are used in a wide are of applications, from authentication systems [9] through network intrusion detection [20] and speech recognition and analysis [13].

The Geo-Som Approach The Geo-SOM Bação et al. [2] applies the first column of geography “Everything is related to everything else, but near things are more related than distant things.” to the SOM algorithm, where the winning neuron is chosen by in a radius k defined by the geo-coordinates of the data. In this way the Geo-som forces units that are close in the input space to be close in the output space. The representation of the Geo-som can be seen in Figure 5.

Detecting Hidden Patterns on Twitter Usage Cheon and Lee [8] analyzed hidden patterns created buy the natural usage of twitter by its users. In its study they started by collecting data from the twitter API of different kinds of topics like “2009 Iran Election” and “iPhone 3.0 OS launch”. They made multi level signal extraction not only from information directly present on the tweet, but also by cross referencing with other social website and with the twitter user profile information. The signals retrieved from the social network can be seen in Table 1.

¹² <http://www.twitter.com>

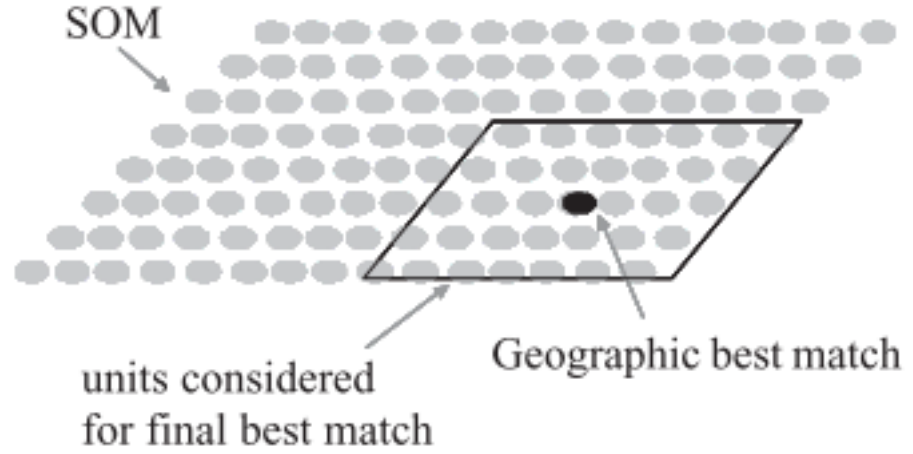


Fig. 5. Geo-SOM structure, from Bação et al. [2]

Table 1. Twitter Signals

Twitt Corpus	Twitter Profile	External Sources
Tweet Size	Gender	Other Social Network Accounts
Replies	Number of customizations	Type of website
Re-tweets	Friends to followers ratio	
Hashtags	frequency of posts	
Presence of URIs and Type of linked content	Account Age	
Type of Device	Country	
Tweet Location		

By applying a SOM, they could find 4 demographical clusters during the Iran 2009 Election. The first cluster was characterized by young web-based Iranians, with twitter accounts not older than 3 months with a high frequency of replies. The second cluster was mainly compound of web users from Iran accounts older than 3 months. The third cluster had Iranian users with mobile clients with large texts clearly trying to raise awareness. The fourth and final cluster represented the users around the world trying to raise awareness about the issue by sharing tweets with URIs. Looking at their analysis about the topic "2009 Iranian Election" it is clear to see that it was possible to describe the type of users represented in the social network and the way they interact with it.

On the iPhone 3.0 OS launch it was possible to find three main clusters. The first columnster was characterized by male users, accounts older than 90 days, coming from countries where the iPhone is marketed, with high adoption of social media clearly representing the target market of the iPhone or its customers. The second cluster had new accounts with higher rate of followers to followees, high frequency of posts per day, presence of URI linking to technology blogs or websites, no country or gender specified meaning that this cluster was clearly composed by news aggregators and technological news websites. Inside the second cluster there was a sub-cluster of Japanese users which represents the high rate of iPhone adoption in Japan. Finally the third cluster was clearly spammer accounts that where eventually deleted after a couple of months, characterized by popular social connections, posting more than 50 tweets a day with external URIs and the accounts where not older than a day or so.

In conclusion it was possible to detect Twitter usage patterns and specifically detect spammers before they where banned from the social network.

3.2 Topic Detection and Clustering

There have been many topic detection techniques. Many of them rely on the TF IDF [3] (term frequency – inverse document frequency) which is not particularly adequate for topic detection on Twitter due to the fact that tweets are very small, composed by typos or slang words and might be written in multiple languages, sometimes at the same time. In this subsection we will take a look at multiple methods of topic detection in general and specifically on the Twitter social network.

Topic and Trending Detection Due to the rapid adaptation of people to always be on-line, through the usage of cellphones on the move, desktops at work and even theTV at home, the increase of user generated content has increased tremendously in latest years. In 2006 35% of on-line adults and 57% of teenagers created content on the Internet ¹³, which in "Internet Years" was ages ago. With amount of content increasing, new real-time and scalable algorithms are needed in order to make sense of all this data. Cataldi et al. [7] propose a new technique for emerging topic detection that permits real-time retrieval of the most emergent topics expressed by a community on Twitter. Their work applies the PageRank [21] algorithm to the users follower/followee relationship in order to find the most influential user on the network, and then calculates the most trending topics by relating, social influence, word co-occurrence and time frame. In the end, an interface was created where it would be possible to navigate hot topics in a given time frame. It is important to say that topic labeling was not automatic and was implicit by the time frame of an event, if two highly social events would occur in the same time frame with word relations the results could be interpreted as the same, for example if a political candidate would win the

¹³ Data source: <http://www.pewinternet.org/Presentations/2006/UserGenerated-Content.aspx>

elections at the same of an important sports club would win a specific cup, the word *win* could be trending at the same time for two different topics and due to high temporal dependency they could be interpreted as the same topic. Weng et al. [24] also used the PageRank algorithm in order to find the most influential twitter users on a certain topic, but uses a different approach where they represent each twitter user as a bag of words comprising of all the tweets that they have posted. Afterwards it uses Latent Dirichlet Allocation [5] in order to find the topics each user is interested in. In the end it was possible to prove that follower/followee relation on twitter was not just casual, but that people actually follow other people in which they have some resemblance or common interest. This concept is called homophily and will be further explored by this project.

3.3 Data Mining in Twitter

In this subsection, we will focus on work done on the Twitter social network in order to leverage insights on how the public data available from the website can correlated within itself and with outside sources.

Enhancing the Tweet Tweet retrieval and analysis is a double edged problem. On one side the tweet is really small which makes it almost impossible to retrieve any actual sense from it. On the other hand the amount of tweets generated per day is around 140 million ¹⁴ which means that it is very hard to a deep analyses of the semantics and content of individual tweets, and that, only the more appropriate signals should be evaluated. Tao et al. [23] evaluated how the multiple signals that could be retrieved directly or indirectly from the tweet corpus could mean that a tweet is relevant for a determined topic. In his work, Tao presents premises that seem intuitively true and proves they actually are relevant through a comparison of multiple precision and recall values. Its results on feature comparison where summarized in Table 2, the first column consists of all the made hypothesis categorized by type, and the second column tells if the data used actually influenced in precision and recall results. Tau also compared result of topic characteristics, concluding that distinction between local and global events as well as temporal persistence proved to not be relevant on relevance prediction.

McCreadie and Macdonald [18] also approached the issue of having very little content on tweets in order to categorize a tweet, and tried to solve it by applying the content of linked URIs into the tweet body in order to improve precision and recall. The best fitting approach was using Field-Based weighting where for each tweet a new document is created which contains two fields; the terms in the tweet and the terms in the linked document. Afterwards a learning to rank algorithm called PL2F is used against the dataset from Microblog2011 in order to find the best weighting that should be applied to the tweet corpus and

¹⁴ <https://blog.twitter.com/2011/numbers>

Table 2. Tao et al. [23] resumed results

Hypotheses	Influence of Features
Syntactical	
Tweets that contain hashtags are more likely to be relevant than tweets that don't	Not Important
Tweets that contain an URI are more relevant than tweets that don't	Important
Tweets that are replies to other tweets are less relevant	Important
The longer the tweet is the more relevant it is	Not Important
Semantic	
The more the number of entities the more relevant a tweet is	Important
Different types of entities are of can have different amount of interest to a give topic	Important
The greater the diversity of concepts mentions in a tweet the more likely for it to be relevant	Important
The relevance of a tweet is determined buy its polarity	Important
Contextual	
The lower the temporal distance between a query and the creation of a tweet the more relevant the tweet is	Not Important
The more the number of tweets created by a user the more relevant one of his tweets will be	Not Important

the URI referenced page. With this trained model they where able to improve precision in an order of 0.9.

Rapidly Changing Trends Due to the real time nature of Twitter, using typical retrieval model that relies on term frequency models like BM25 or language modeling cannot be applied, as stated by Lin and Mishne [16]. The study of topic perdurance on the social network proved that it is presented in bursts of queries and mentions of a topic. The typical usage of twitter for search is not the same of Google. When user are searching in twitter they want to find out what is happening right now meaning that classification techniques based on past events cannot respond this kind problem. As stated by Lin and Mishne [16] this problem has not yet been solved at Twitter (or anywhere else at the time of writing this report), and issues a new kind of data analysis approach that was not taken into consideration in the past. This effect of rapidly changing topics and queries based on real time events was named "Churn", and can be clearly seen in Figure 6.

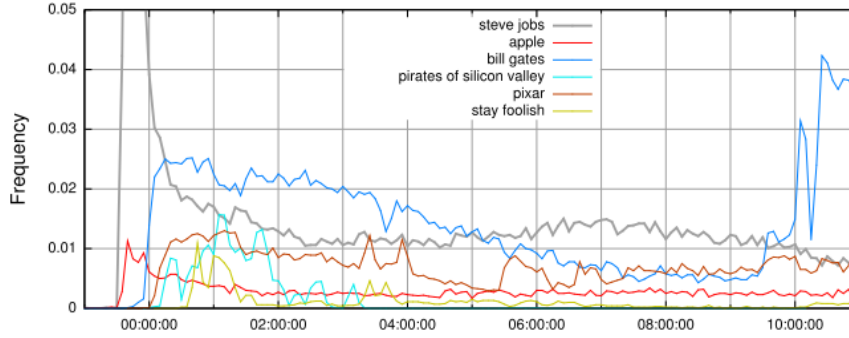


Fig. 6. The Churn effect: Frequencies of queries related to Steve Jobs death over a 12 hour period in 5-minute intervals, normalized to the total number of queries in the interval. At its peak, the query “steve jobs” reaches 0.15 (15% of the query stream); Graph taken from [16]

4 Architecture

This section is comprised of the multiple layers needed in order to successfully be able to implement the proposed project objectives. We will discuss the data gathering process, and detail the SOM algorithms that will be tried in order to achieve the best results on Topic Detection in Twitter. The description of the evaluation method will be described and finally the prototype architecture will be described.

4.1 Collecting Data

For this project we will focus on building our own dataset for SOM training, but existing datasets might be used in order to rate results. The description of existing datasets will be presented in Subsection 5.1.

Twitter API Building your own Data Set through the Twitter API has become harder with passing years with the introduction of API limits and mandatory authentication. With these new limitations, companies like *Gnip*¹⁵ or *Tweet Archivist*¹⁶ with licenses from Twitter are selling access to their archives of tweets.

In order to retrieve data from Twitter is crucial to understand how their API functions. The Twitter API right now is divided into two types, the REST API and the Streaming API. Both can be used at the same time, and have different types of limits. In a general way, the streaming API is used for subscriptions,

¹⁵ <http://gnip.com/topsy/>

¹⁶ <http://www.tweetarchivist.com/about/subscriptions>

where an application can subscribe a given hashtag or user activity on the social network, and they are automatically pulled to the subscriber app. The Streaming API has no specific limit being described in the docs as “The public streaming APIs cap the number of messages sent to your client to a small fraction of the total volume of Tweets at any given moment”¹⁷.

The REST API works by requesting resources and getting the results in a RESTful way. Here the limits are strict, an application can only get a maximum number of 3200 Tweets per user and 180 calls to the API per 15 minutes, more API limits can be found on the Twitter API Documentation¹⁸.

Crawling Twitter In order to get tweets from the Twitter API, the data will be crawled in a breath-first fashion where after selecting a first user:

- Get all Tweets from the user.
- Get user profile info.
- Get list of followers/followees.
- Select a follower/followee and repeat step one.

The algorithm will stop at a given depth level, also if API limits are exceeded the algorithm will have to stop for 15 minutes and afterwards resume. Given the API access limits, there will be no need to run the crawler asynchronously since achieving a greater level of performance will only make the algorithm achieve API limits sooner.

Storing Crawled Data While the crawler is getting data from Twitter, it will be storing it in a Redis database¹⁹. Given the amount of databases available in last few years, Redis was chosen for this project because it met the following criteria:

- Free.
- Simple to install and run.
- Can persist data to disk.
- Its really fast to write, by not granting data integrity (which is not a problem since this project is not dealing with sensitive information)
- Good documentation.
- Non relational, Key/Value store.
- Stores json.
- Client libraries for almost every programming language.
- Integrated publish/subscriber.

Given the characteristics of the Redis database, there will be no need to write a schema beforehand. With this in mind, user information and Tweets will be stored directly in json into the Database.

¹⁷ <https://dev.twitter.com/docs/faq#6861>

¹⁸ <https://dev.twitter.com/docs/rate-limiting/1.1/limits>

¹⁹ <http://redis.io/>

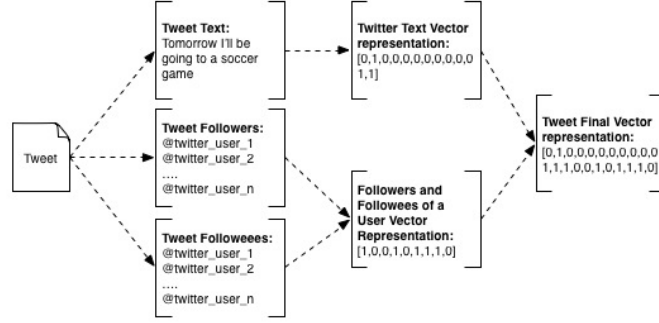


Fig. 7. Vector Space transformation of a Tweet

4.2 The SOM Algorithm

The SOM algorithm in this report will have a twofold approach. Primarily we will try to transform the tweet social characteristics and words into a vector using the vector space model, this approach will use the default SOM implementation and will be described in detail in Subsection 4.2. The second approach is inspired by Bação et al. [2] where the SOM algorithm will be altered in order to take in consideration the social network during the its training, this implementation will be described in Subsection 4.2.

Default SOM Approach In order to train the SOM first we will need to convert the tweets into the vector space model. There will be two binary vectors, the first one will represent the presence of all the words gathered in all tweets where the value 1 will represent the presence of a word and the value 0 the absence of a word. The second vector will represent the social connections between the user. On a first approach, in order to give more relevance to social connections, only followers that are followed back will be represented. This approach will give a stronger representation to the social interaction between the two users, since on the Twitter social network it possible to follow someone without the followed person accepting the request. In the end both vectors, the word representation vector and the social connections representation vector will be concatenated. Figure 7 shows the transformation from the tweet into the vector space model.

The SOM initialization will be tried in different ways and measured to see which gives the best results. The initialization characteristics will be the following:

- Random initial number of neurons with random content.
- Random initial number of neurons with random content evenly distributed.
- Each neuron will be the representation of each user that was crawled.
- Each neuron will be comprised of words relevant to a determined topic, and will be responsible to categorize that topic.

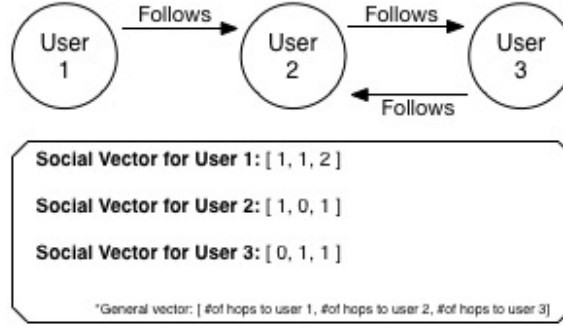


Fig. 8. Vector describing the number of hops between followers

- Each neuron will be the representation of a user that is relevant to determined topic, for example the user *Optimus Alive* would be responsible for categorizing the topic *"Music Festivals"* while *Cristiano Ronaldo* neuron representation would be responsible for detecting the topic *Soccer*.

In the initializations described above, only the last two will yield a SOM ready to classify tweets into topical categories. The first three options would build uncategorized clusters that would have to be classified *a posteriori*. Nevertheless it will be important to look at the results given by them because they could be more interesting then the results from the last two items.

Social SOM In this approach, the use of the vector that described the followers/followees of the user that created the tweet, which is described in Subsection 4.2 will be discarded. Instead there will be a new vector that describes the number of hops between twitter users, a visual representation of this vector can be seen in Figure 8. By modifying the SOM training algorithm, input patterns will only be measured against neurons that have a determined level of social affinity. This social affinity will be defined as x which will define the number of followers/followees relations will be applied. For example if $x = 1$ only neurons that belong to a follower/followee will be selected for comparison to find which is the winning neuron. If $x = 2$ not only the followers/followees of a user will be selected for comparison, but also the followers/followees of the followers/followees. Finally if x would be equal to the number of users used in the dataset, then all neurons would be used for comparison, making this solution equal to a normal SOM.

SOMS Development The development and testing of the SOM described in above will be completely independent of the prototype described in Subsection 4.3. The SOM training will be made using the datasets described in Subsection 4.1 and connection with components from the prototype will be made

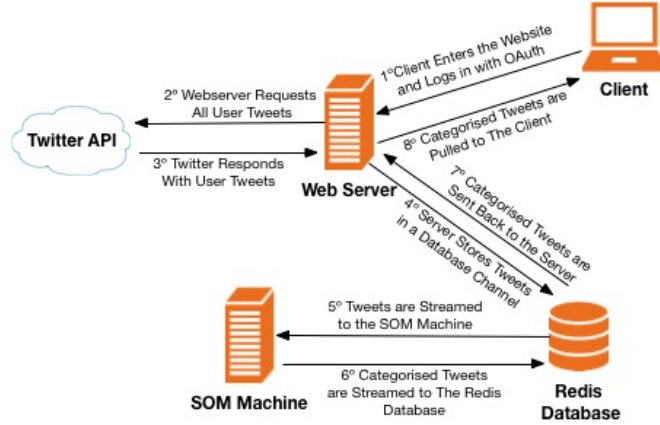


Fig. 9. Topology of the Solution

through the Publish/Subscriber Redis interface. This approach will create an highly modular solution where it will be possible to interact with the trained SOM through terminal where it will be possible to visualize and test results.

4.3 Prototype

As described in the report objectives, a prototype will be created to demonstrate the project categorizing tweets per topics. The prototype will be developed as a Web Application divided into web server and a client side application that will run on a web browser. The prototype will work in the following way:

- A user will be able to login with his twitter id.
- The browser will start displaying the user tweets, with the associated topic.
- After all the user tweets are displayed, some statistics about the user tweeting topics will be displayed.

The wire-frames for the application are provided as attachment, the architecture of the solution can be seen in Figure 9 and will be described in the following sub-subsections.

Client Side Application The client side application will be running on the browser of the user connecting to the website. This application is responsible to authenticate itself against the Twitter API, through OAuth, and after that will establish a web-socket channel that will be receiving the categorized tweets as soon as the server dispatches them. The client application will also have to display an interface to the user so he can interact with the application.

Web Server The web server will be very simple, since it will only have to get the tweets from the Twitter API and publish them in a channel through the Redis Publish/Subscriber interface, on the other side, the SOM machine will be receiving the Tweets and categorizing them. Also, the web server will have to subscribe to the channel where the SOM Machine is publishing the categorized tweets in order to be able to push them to the client through web-sockets.

Redis Database The Redis database will be used as a middle man between the Web Server and the SOM Machine as a publish/subscriber system. This is extremely useful because it will separate the login of the web server and web application from the rest of this project. In this way the SOM machine can be running the topic categorization in any programming language while the server can also be working in another language and still be able to contact each other in a simple and efficient way.

SOM Machine The SOM Machine is where the Self-Organizing Map trained with the crawled data described in Subsection 4.1 will be functioning. It will be subscribing to new tweets sent from the Web Server and will categorize them using the previously trained SOM. After the tweets are categorized they will be published back so the Web Server can deliver them to the client.

Solution Overview We now describe how everything in this solution fits together, based on the diagram in Figure 9. The steps described will be the same, but they will have more detail. For simplicity sake, the OAuth authentication from the client to the twitter API will be omitted since its only objective is to get access to the user tweets.

- **First Step** The client connects to the website through the browser and will login with his twitter account in order for the server be able to download all of his Tweets. When the login process is over it keep an open web-socket with the server in order to receive the categorized tweets.
- **Second and Third Step** The web-server will request all the user tweets through the Twitter API. Twitter will respond with all ²⁰ the tweets from the user that logged in.
- **Forth and Fifth Step** The web server will publish the tweets in the Redis database, while subscribing to the channel where the categorized tweets will come out. On the other side the SOM Machine is subscribing to the uncategorized tweets in order to categorize them.
- **Sixth and Seventh Step** After the SOM machine categorizes the tweets it will publish them to the Redis database, which is being subscribed by the web server.

²⁰ The number of tweets per user are limited to the 3200 most recent through the twitter API

- **Eighth Step** As soon as categorized tweets start hitting the server they will be immediately sent to the client through web sockets. On the client as soon as the categorized tweets hit the browser, they will be injected into the DOM and the user will start to see the tweets that he has made categorized in topics. Lastly after all the tweets have been sent to the browser, some statistics will appear with the amount of tweets per topic throughout the time.

5 Evaluation Metrics

Evaluation of the topic detection on Tweets will be made in two distinct ways. The first way will focus on binary classification using the precision and recall metrics, and will be described in Subsection 5.1. The second way will focus on statistically testing the SOM learning process and the computed trained network. This testing process will be described in Subsection 5.2.

5.1 Testing for Precision and Recall

Precision and Recall are both ways to measure the rate of right guesses made by the trained SOM network, and are defined in the following way:

- **Precision:** Fraction of retrieved instances that where relevant

$$precision = \frac{|relevant\ documents \cap retrieved\ documents|}{retrieved\ documents} \quad (1)$$

- **Recall:** Fraction of relevant instances that where retrieved

$$recall = \frac{|relevant\ documents \cap retrieved\ documents|}{relevant\ documents} \quad (2)$$

In order to calculate Precision and Recall we need to have the *relevant documents* and the *retrieved documents*. The *relevant documents* are rather hard to determine because they need to be categorized by humans, which is a tedious task. Until the writing of this document, only one free dataset categorized by hand was found²¹, but it is categorized in sentiments and not by topics which makes it irrelevant for this project. If no dataset is found until the end of this project, a small dataset might be created in order to measure precision and recall.

5.2 Statistically Testing the SOM

SOM training is always subject to some variability due to multiple causes, like the sensitivity of initial conditions, convergence to local minima and sampling variability, as stated by Bodt et al. [6]. This subsection will present statistical tools to measure the quality of the SOM, by measuring its quantization error and topology preservation.

²¹ <http://www.sananalytics.com/lab/twitter-sentiment/>

Quantization Error The SOM Quantization Error is the mean of all Euclidean distances between the observed data points and their corresponding winning neuron. This value might vary depending on the initialization neurons or the order of the input data fed into the SOM while the training is occurring. When applied to an individual input data, represents how well a neuron is representing input data. Since the SOM Quantization Error represents the mean of all quantization errors from all the input data, generally, the lower the error is the best the SOM was trained.

There is no general formula to minimize quantization error[6]. What is generally done is just to change the number and values of the starting neurons and the order of the input data in order to train multiple SOMs. In the end the SOM with the lowest quantization error is chosen. In this project since multiple approaches to the SOM algorithm and data representation will be tested, as described in Section 4.2, and the ones having the lower quantization error will be selected for the prototype.

Topology Preservation The Self-Organizing Map performs a mapping from the n -dimensional input space into the two dimensional output space and where resides one the most fascinating characteristics, which is that the output map tries to preserve the topology from the input space. This grants the SOM algorithm a way to visualize high-dimensional data that other neural networks or clustering algorithms don't have. Even though this is true, sometimes during training it is not possible to preserve the topology of the network. Thus topology preservation can be measured through the Topographic error Kiviluoto [12] which is the proportion of all data vectors for which first and second BMUs²² are not adjacent units. In this project the Topographic Error will be calculated for all SOM implementations and VSM usages in order to understand if the representation of the SOM output space is well defined.

6 Conclusions and Future Work

In this report we presented a new approach to topic detection in datasets with a short documents, by leveraging the implicit semantics behind the data. First, in Section 1, we started by introducing the concepts and problems of topic detection the Twitter social network. Afterwards at Section 2, we described basic concepts for this project, such as how document clustering and Self-Organizing Maps work. Then at Section 3 we focused on describing related work done by other researchers in the field of Self-Organizing Maps, topic detection with clustering and data mining on Twitter. At Section 4 we proposed how we where going to use and adapt the SOM algorithm to better fit for solving the problem of topic detection on Twitter. In this section we also described the architecture of our prototype which will used to demonstrate live topic categorization. Finally at Section 5 we described how we where going to evaluate our solution.

²² unit that is closest to the winning neuron. BMU Best fitting unit

Even though this project is a proof of concept on topic detection by focusing on user specific interests and relations, and on the application of Self-Organizing Maps as a tool for clustering user tweets with social features. We think that by escalating the proposed architecture in order for it to be able to process a lot more data from Twitter will render interesting results. We presume this by looking at work done in other research areas like work from Le et al. [15] that also used neural networks for unsupervised learning but leveraged giant datasets and computational power.

A Appendix

A.1 Web Site Wire-frames

A.2 Work Scheduling

Table 3. Work Scheduling

Month	Work
February	Program Twitter Crawler
February	Start Implementing the Simple SOM
March	Test Simple SOM
April	Start to implement prototype
May	Finish developing the prototype
May	Start Developing the Social SOM
June	Test the Social SOM
July	Write Dissertation



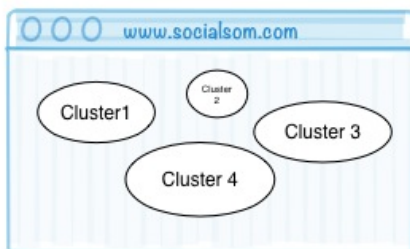
Login Screen



Display user tweets as they are categorized



Display user Statistics Based on The Topics he Tweets



Display Topic Clusters of User Followee/Follower Topics

Fig. 10. Website to display user Twitter content

Bibliography

- [1] Sitaram Asur and BA Huberman. Predicting the future with social media. ... *Agent Technology (WI-IAT), 2010 IEEE ...*, 2010. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5616710.
- [2] Fernando Baao, Victor Lobo, and Marco Painho. The self-organizing map, the Geo-SOM, and relevant variants for geosciences. *Computers & Geosciences*, 31(2):155–163, 2005. ISSN 00983004. doi: 10.1016/j.cageo.2004.06.013. URL <http://linkinghub.elsevier.com/retrieve/pii/S0098300404001918>.
- [3] Ricardo A. Baeza-Yates and Berthier Ribeiro-Neto. *Modern Information Retrieval*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1999. ISBN 020139829X.
- [4] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. (1993):1–34, 2013. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=6472238.
- [5] DM Blei, AY Ng, and MI Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003. URL <http://dl.acm.org/citation.cfm?id=944937>.
- [6] Eric De Bodt, Marie Cottrell, and Michel Verleysen. Statistical tools to assess the reliability of self-organizing maps Statistical tools to assess the reliability of self-organizing maps. pages 1–27.
- [7] Mario Cataldi, Luigi Di Caro, and Claudio Schifanella. Emerging topic detection on Twitter based on temporal and social terms evaluation. *Proceedings of the Tenth ...*, pages 1–10, 2010. URL <http://dl.acm.org/citation.cfm?id=1814245>.1814249<http://dl.acm.org/citation.cfm?id=1814249>.
- [8] Marc Cheong and Vincent Lee. A Study on Detecting Patterns in Twitter Intra-topic User and Message Clustering. *2010 20th International Conference on Pattern Recognition*, pages 3125–3128, 2010. doi: 10.1109/ICPR.2010.765. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5597282>.
- [9] Hiroshi Dozono. Application of Self Organizing Maps to Multi Modal Adaptive Authentication System Using Behavior Biometrics. 2012. URL <http://www.intechopen.com/books/applications-of-self-organizing-maps/application-of-self-organizing-maps-to-multi-modal-adaptive-authentication-system-using-behavior-bio>.
- [10] BJ Jansen and Mimi Zhang. Twitter power: Tweets as electronic word of mouth. *Journal of the American ...*, 60(11):2169–2188, 2009. doi: 10.1002/asi. URL <http://onlinelibrary.wiley.com/doi/10.1002/asi.21149/full>.
- [11] Seung-Shik Kang. Keyword-based document clustering. *Proceedings of the sixth international workshop on Information retrieval with Asian languages*

- , 11:132–137, 2003. doi: 10.3115/1118935.1118952. URL <http://portal.acm.org/citation.cfm?doid=1118935.1118952>.
- [12] K Kiviluoto. Topology Preservation in Self-Organizing Maps. *Neural Networks, 1996., IEEE International ...*, 1996. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=548907.
- [13] T. Kohonen. The 'neural' phonetic typewriter, pages 11 - 22 . 1988.
- [14] T Kohonen. The self-organizing map. *Proceedings of the IEEE*, 1990. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=58325.
- [15] QV Le, MA Ranzato, R Monga, and Matthieu Devin. Building high-level features using large scale unsupervised learning. *arXiv preprint arXiv: ...*, 2011. URL <http://arxiv.org/abs/1112.6209>.
- [16] Jimmy Lin and Gilad Mishne. A Study of" Churn" in Tweets and Real-Time Search Queries (Extended Version). *arXiv preprint arXiv:1205.6855*, 2012. URL <http://arxiv.org/abs/1205.6855>.
- [17] YC Liu, Ming Liu, and XL Wang. Application of Self-Organizing Maps in Text Clustering: A Review. 2012. URL <http://www.intechopen.com/books/applications-of-self-organizing-maps/application-of-self-organizing-maps-in-text-clustering-a-review>.
- [18] R McCreddie and Craig Macdonald. Relevance in microblogs: enhancing tweet retrieval using hyperlinked documents. ... *Research Areas in Information Retrieval*, 2013. URL <http://dl.acm.org/citation.cfm?id=2491787>.
- [19] Prem Melville, V Sindhwani, and R Lawrence. Social media analytics: Channeling the power of the blogosphere for marketing insight. *Proc. of the WIN*, pages 2–6, 2009. URL <http://people.cs.uchicago.edu/~vikass/sma-win09.pdf>.
- [20] Kulkarni P. Nikam D.M. Pachghare, V.K. Intrusion Detection System using Self Organizing Maps. 2009.
- [21] The Pagerank, Citation Ranking, and Bringing Order. 1 Introduction and Motivation 2 A Ranking for Every Page on the Web. pages 1–17, 1998.
- [22] Sylvia Springorum, SS im Walde, and Jason Utt. Detecting Polysemy in Hard and Soft Cluster Analyses of German Preposition Vector Spaces. *aclweb.org*, 1998. URL <http://www.aclweb.org/anthology/I/I13/I13-1072.pdf>.
- [23] Ke Tao, Fabian Abel, Claudia Hauff, and GJ Houben. What makes a tweet relevant for a topic? *Making Sense of Microposts (# ...*, pages 49–56, 2012. URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.232.2188&rep=rep1&type=pdf#page=54>.
- [24] Jianshu Weng, EP Lim, Jing Jiang, and Q He. Twitterrank: finding topic-sensitive influential twitterers. ... *conference on Web search and data ...*, 2010. URL <http://dl.acm.org/citation.cfm?id=1718520>.