

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

Bernardo Simões 20-25 páginas aproximadamente

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

Abstract. *70 and at most 150 words, What did I do, in a nutshell?, summarize the paper, should be written last , very short context ,what the objectives of the study were*

Keywords: topic detection, twitter, self-organizing maps, classification, clustering

1 Introduction

With the evolution of social networks websites like Facebook or Twitter throughout the last couple of years, the amount of pertinent content about a specific issue is increasing dramatically, which calls for new ways to make sense and catalog this data. In 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 value of data presented in sites like Facebook or Twitter as proven its value in papers like “Predicting the future with social media” where it is possible to predict with high precision the value of a movie box office weeks before it debuts.

This project will focus on Topic Detection on Twitter by using a new approach that will not only take in consideration the words in the corpus, but will also take in account the social network to which a tweet belongs to in order to categorize it using the concept of homophily that has been proven in past that is applicable to social networks.

Will be described the objectives of this project, at 2 we will talk about 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 finally at section 5 it will be discussed how to evaluate results achieved.

1.1 Objectives

The objective of this project is clear, finding topics on Tweets by analyzing their corpus specific characteristics, like number of characters in a tweet, hashtag, “was retweeted”, etc.. And contextualize the social network evolving the person that did the tweet.

After characterizing the tweet with information just described, we will use the unsupervised learning clustering technique Self-organizing maps in order to organize the tweets in clusters of topics. Afterwards it will be needed to categorize the clusters in order to know which topic they belong to.

Lastly the resulting topic clusters will be publicly accessible through a website to everybody that visits it.

2 Basic Concepts

In this section we will start by generally describing what Clustering is and how it works at subsection 2.1, then at subsection 3.1 it will be outlined how Self-organizing [4]maps function, which is the Document Clustering algorithm used on this project.

2.1 Document Clustering

Document clustering is an optimal division of data into categories without prior knowledge of the data that is being organized, based only on the similarity between documents. 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 [5] described that Document Clustering can be used to a variety of Computer Science fields, such as:

- Natural Language Preprocessing.
- Automatic Sumarization.
- User preference mining.
- Improve Text classification results.

In regard to document categorization 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. [6] performed hard and soft clustering with SOMs [4] while identifying polysemous German Propositions. They used regular SOMs to create multiple hard clusters and used Centroid-Based or Preposition-based softening to create Soft Clusters from the Hard Clusters.

The general mathematical description of Document Clustering can be seen in 1 In the first step a data set must be provided in order to cluster the documents. The second step “Keyword identification” is where non relevant words are

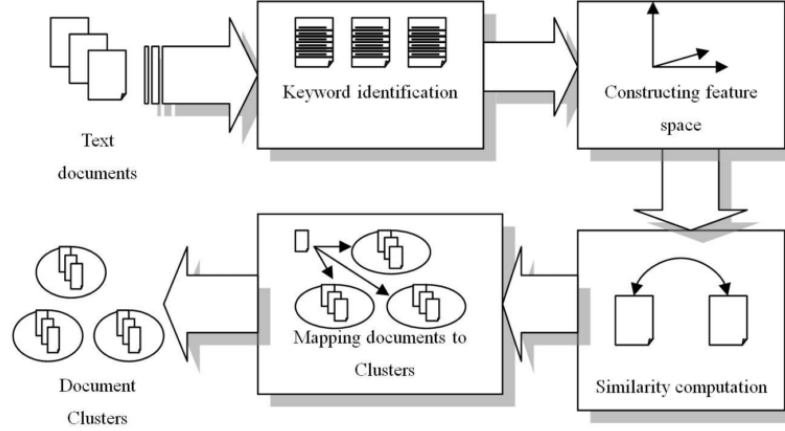


Fig. 1. Text Clustering Main Framework from Dozono [2]

removed from the documents. Kang [3] proves that keyword removal improves clustering. Another way to extract features 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 relevante features due to the characteristics of the markup. In 3.2 it will be described twitts characteristics as a document and in 4 how feature extraction will be implemented on this project. "Constructing the Feature Space" is characterized by converting the keywords of each document into vectors, the most common algorithm used for this task is SVM (Support Vector Machines). In SVM each vector dimension means one detected key word and each document is represented by the vector of keywords in the feature space. This process and keyword removal is described in Figure 2. Due to the way documents are represented in SVM it is normal that vectors become very large and full of zeros (keywords not present) and it is needed to use sparse vectors to represent the documents in a more efficient way.

Dimensional reduction is done after the construction of the vector space model, in order to reduce the size of the vector space. There are two main ways to do this PCA (Principal Component Analysis) and LSI (Latent Semantic Indexing). PCA calculates the k eigenvectors of the co-variance of the document matrix, which reduces the size of the matrix to k . LSI (Latent Semantic Indexing) works just like PCA but the eigenvectors are calculated directly from the document matrix.

There are two main strategies for Document Clustering, Complete strategie where the data set does not change and Incremental where initial number of document can increase by adding new documents. After a new document is added it can be merged into a existing cluster, or can be separated as a new cat-

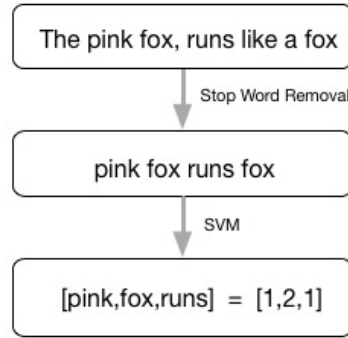


Fig. 2. Stop word removal and transformation to Vector Space Model

egory. While adding new documents it might be needed to re run the clustering algorithm.

After the algorithm converges, cluster similarity can be calculated in multiple ways:

- Shortest Distance Method: Shortest distance between two members of different clusters.
- Longest Distance Method: Longest distance between two members of different clusters.
- Group Average Method: The average distance between all elements of both clusters.
- Centric Method: The distance between the center of two clusters.

There many clustering algorithms, here we will focus on the three most popular. K-means works by randomly selecting k documents as the cluster centroids, then assign each document to the nearest centroid, and finally recalculate the the centroid with new added documents. The algorithm should be executed until convergence which reflects in the centroids stop changing. K-means has the advantage that the number of centroids must be selected before starting the algorithm. AHC or Agglomerative Hierarchical Clustering is hierarquical clustering algorithm where clusters have sub-clusters which have subclusters. Like K-means it is also a simple algorithm that starts by calculating the similarity matrix, then each document is seen as a cluster and finally merge the nearest two clusters into one and update the similarity matrix. The algorithm ends when there is only one cluster or due to clustering entropy. An AHC classic example is species taxonomy where species have subspecies which have subspecies, etc. Lastly there is Self-organizing Maps introduces by [4] which will be used in the thesis and will be detailedly described in the next subsection 3.1.

2.2 The Self-organizing Map

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

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 so called neurons are responsible for categorizing input patterns for which they are responsible. Nearby neurons will be organized by similarity which will cause that similar patterns will activate in similar areas of the SOM. With a topology preserving mapping, SOM organizes the information spatially where similar concepts are mapped to adjacent areas. Neurons are displayed in an n dimensional grid, generally rectangular, but other dimensions are possible like 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 SVM 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 learning phase is characterized by the training algorithm, which works the following way:

- Neurons can be initialized randomly or it is possible to select initialization neurons.
- 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.

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 (the learning rate converged to zero), it is very easy and fast to classify new data now. It is possible that during training the SOM gets stuck while unfolding, this kind of behavior might happen if the input patterns are very complex.

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

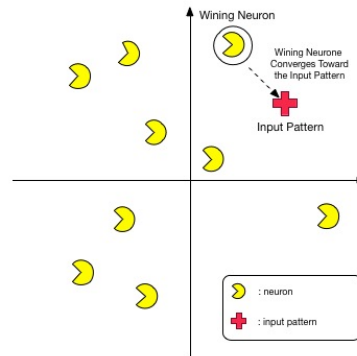


Fig. 3. Winning neuron converging at learning rate

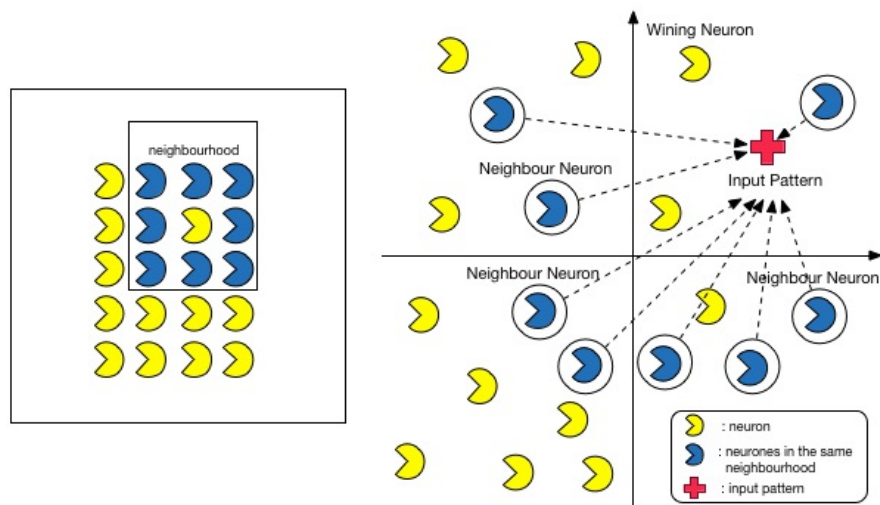


Fig. 4. On the left the output space neighbor, on the right the neighbors of the winning neuron converging

3 Related Work

3.1 Self-Organizing Maps

3.2 Topic Detection on Twitter

- What did we know about the problem before I did this study?
- What did we do different from previous works?
- Discuss the relevant primary research literature
- Works should be organized by their relevant characteristics
- Comment on why it is relevant for your work
- Comment on what your work does differently

4 Architecture

In this project we are going to apply Self Organizing Maps in order to detect clusters of Topics on Twitter.

4.1 Data Gathering

In order to retrieve data from Twitter, we will be using a ruby library called Twitter Stream [1], that enables the user to download and inspect the twitter stream of tweets. As the data is gathered it will be stored in a MongoDB database for posterior analyses. As the twitter stream is stored, another function will interact with the twitter API in order to retrieve information from a user profile and relate him with other users by analyzing his followers and who the user is following. In the end of the data-gathering process it will be possible to query the database for:

- Tweets from a user.
- Query tweets for hashtag.
- Query users followers and who he is following.
- Query for tweets that shared the same URI

Tweets will be categorized in:

- News Accounts
 - Accounts with a lot of followers
- Profile customization
- Average number of tweets a day with uri (might suggest spam)
- How am I gonna solve the problem?
- Describe the work that will be done

5 Evaluation Metrics

- How am I gonna evaluate my work?

5.1 Evaluation Criteria by Teachers

- Ability to understand the research problem
- Clear and well defined goals
- Description of the different approaches explored
- Ability to relate the state-of-the-art with the research theme Work methodology and adequate planning for the next stage Organization and quality of the written document
- Inclusion and completeness of updated and appropriate references Oral presentation and discussion

Bibliography

- [1] Steve Agalloco. TweetStream, simple Ruby access to Twitter's Streaming API. 2013. URL <https://github.com/tweetstream/tweetstream>.
- [2] 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>.
- [3] 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>.
- [4] T Kohonen. The self-organizing map. *Proceedings of the IEEE*, 1990. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=58325.
- [5] 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>.
- [6] 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>.