

# Event summarization from news and tweets correlation

Cavallari Sandro  
Giglio Marco  
Morettin Paolo

## ABSTRACT

In this report we consider the task of extracting a summary of the main event that caused a shift on the opinion of Twitter users. This paper presents several techniques which might be used to analyze these *sentiment shift* and to find the event which caused them. For each technique presented below some aspects are taken in account, which are important in the context of a data mining application, such as scalability and efficiency. A comparison of the different techniques is showed as well.

## Introduction

This work aims to find a summary of the events which caused a *sentiment shift*.

While event extraction NLP techniques are mature, their performance on tweets inevitably degrades due to the inherent sparsity in short texts. Since tweets contain heterogeneous language structures and are at most 140 characters long, to extract event from tweets is quite tricky: instead linking tweet to news allows us to extract events from news text that has a well formed structure and contains more text with respect to tweets.

Moreover NLP document compression techniques usually exploits language dependent methods, whereas we want a methodology as language independent as possible.

This paper is structured as follows: in the next chapter we present some works which are somehow related with the one we are addressing, then we present a more formal definition of our problem and three different techniques which might be used to solve it (two of this techniques are based on the *bag of words* approach - one exploiting the *SpaceSaving algorithm* presented in [7], the other exploiting *Latent Semantic Indexing* [1] - , whereas the third one exploits *n-gram graphs* [3]). After that we will show and compare the results and the perfor-

mances which characterize each of these methodologies. At the end of the paper we will present some conclusions and draw some suggestions on how this work may be extended.

## Related Work

We can distinguish two main subgoals in our project:

- compute the correlation between tweet and news
- create a summary of the main event that caused the sentiment shift

There are several works, in literature, which refer to similar tasks. In particular, due to Twitter's increasing popularity, in the last years there have been several works regarding the correlation and analysis of the data. Weiwei Guo et al. [4] proposed a framework to link tweets and news and to extract from the resulting correlation missing aspects of the tweet-event. In their research, rather than using a LSI technique, they propose a methodologies based on the Weighted Textual Matrix Factorization [5] [WTMF] model(the 2012 de facto standard). This approach revealed to return good correlation results, making WTMF a strong tool for baseline creation.

Some works which are complementary with respect to our might be found in the field of recommendation systems; in particular systems which address the problem of recommending news given some user's features. However, there is an important difference between this task and our, since tweets are much more heterogeneous and characterized by a much smaller length than news.

Recently Google presented a system for *multi-sentence compression* (MSC) [6], which aims to compress several sentences trying to maintain high readability and to not lose the semantic content of the text. Their approach exploits word gram graphs, combining them with some semantic data. The combination of the two approaches allows redundancy removal while maintaining a readable sentence as output.

## Problem Definition

We define *sentiment* the measurement of feelings performed using specific tools analyzing the stream of tweets. This kind of measurement is usually done using some supervised machine learning techniques together with NLP techniques.

Tools exist to track the *sentiment* toward a certain topic and to identify changes in it. In particular we can consider three major events in a sentiment timeline:

1. there is a change of sentiment from positive to negative (or viceversa)
2. there is no change of sentiment, but there is a peek in the graph describing it or its value crosses the average value
3. there is no change but in the volume of tweets/second

We define the first of these events as *sentiment shift* and the temporal window in which it happens is called *contradiction point* or *contradiction window*.

In addition we define *contradiction tweets* and *contradiction news* respectively the tweets and news published inside a contradiction window and *background tweets* and *background news* the set of all tweets and news regarding the addressed topic.

Given these definitions, the problem we aim to solve is to find, among the contradiction news, the ones which are the more representative of the contradiction tweets and to present to the user a summary which describes them.

## Proposed Approach

In this section we present different approaches which might be used to solve the problem we are addressing. Of course, each approach has advantages and disadvantages which we will list while showing them.

## SpaceSaving approach

### Description of the algorithm

*SpaceSaving* [7] is an algorithm which was first presented by Metwally et al. in 2005 to efficiently compute the most frequent terms in a data stream. It allows the user to find  $k$  words which are among the most frequent in the given stream. Although this is a heuristic algorithm whose result's correctness is not guaranteed, SpaceSaving is able to specify the upper bound of the error for each word presented as output.

The classic implementation is based on a fixed size list of tuple

$$(word, occurrences, error)$$

The stream is read word by word and at any time one of the following condition is verified:

1. the word is already in the list of frequent terms
2. the word is not among the frequent terms, but the list is not full
3. the word is not among the frequent terms and the list is full

If the first condition is verified, then the algorithm will proceed incrementing the number of occurrences of that term; if

condition 2 is verified, instead, we add that word to the list of frequent terms with number of occurrences 1 and number of errors 0; if condition 3 is verified, then we scan the list for the term with the lowest number of occurrences and replace it with the new word. When this replacement takes place we set the number of errors equal to the number of occurrences of the word we just replaced, then we increment the number of occurrences.

Our python implementation of this algorithm is slightly different from the classic one described above, since we test the word against a list of stop words and we added a hashmap to avoid useless replacement. In our implementation, then, the following steps are performed:

1. if the word is in the stop word lists, does nothing
2. otherwise compute the hash of the given word
3. increment the value associated to that hash
4. check whether the value stored in the hash map is above the number of occurrences of the less frequent term tracked by the SpaceSaving algorithm. If so, the replacement occurs, otherwise nothing happens.

The output of this algorithm is a couple  $(word, value)$  where the value is computed as

$$value = occurrences - error$$

In order to reduce the noise we chose to use only couples where the value is above a certain threshold.

### Description of the methodology

We have seen how we exploit the SpaceSaving algorithm to obtain the list of frequent terms in a text and their values.

To accomplish our task, we ran the program described above on the set of tweets inside a *contradiction point* [CP] (point in which a sentiment shift occur), thus obtaining the list of frequent terms inside the contradictions; then, we read all the news which were published during the contradiction (a certain error between the publishing date and the contradiction window might be taken in account) and for each of them we compute its score as

$$newsScore = \sum_{w \in W} value(w)$$

where  $W$  is the list of words in the news and the value is null if that term is not in the *frequent term list*.

The news with the highest score is selected as the one causing the shift and its first words are presented as output to the user.

This approach revealed to be very fast and from its formulation we can see that it can be applied incrementally. As a drawback, we must say that it does not take in account the background coming from the topic, hence it is not able to discern whether a given frequent term is informative or not (e.g. running this program on the Obama topic it is likely that the most frequent terms will be "President", "Barack", "Obama", "USA", hence

a news containing more repetitions of these words will have a high score, even if those words are not informative at all)

## LSI

We used a bag-of-words approach exploiting a technique called Latent Semantic Analysis (also known as Latent Semantic Indexing), formalized by Scott Deerwester et al. in 1990 [1]. LSA is a information retrieval statistical approach, that can be used to find semantical similarities between sets of documents. It takes advantage of Singular Value Decomposition on a vector space of term frequencies in order to create, given a set of documents, an *index* matrix, that can be queried with other documents to find the most similar one. In order to find correlations between the sentiment shifts and the news, our approach works as follows:

- All the tweets corresponding to a CP are merged in a single document.
- Every document is preprocessed and converted in the bag-of-word representation.
- The LSA model is trained on the whole set of news for the given topic.
- The set of candidate news to be tested is chosen accordingly to a window, and it's used to create an index.
- The news are sorted in order of similarity, optionally weighted by a factor considering the time distance of the news from the CP.
- All the candidate news over a certain threshold are considered to be correlated and a set of words calculated with TF-IDF is returned.

## Preprocessing

During the preprocessing, every document is filtered and converted from a string to a bag-of-word representation. The following steps are done in this phase:

1. The punctuation characters are substituted with a whitespace.
2. The string is tokenized.
3. Stopwords and tokens appearing less than  $t$  times in the document are removed.

## LSA and candidate news selection

Once the documents are filtered and transformed in a bag-of-words representation, a dictionary is created containing the association between tokens and a positive integer identifier. This is done in order to convert on-the-fly each document  $D$  in a vector of tuples  $(i, n(D, i))$  where  $n(D, i)$  is the number of occurrences of the token  $i$  in the document  $D$ . The LSA model is then trained with all the news available for the given topic. Given a sentiment shift over a time interval  $[T_b, T_e]$ , the news can be selected using a fixed size window  $[T_b - c, T_e]$  or a more sophisticated approach that considers the density of the tweets inside the shift. The latter seems reasonable since news that correlates with a shift are likely to be near the shift the more the *hype* increases.

## News scoring and summarization

With the candidate news a matrix index is then created and the document representing the shift is compared. This operation produces a vector of similarities of the shift with the candidate news. These scores do not take into account the hype and the temporal distance of the news w.r.t. the sentiment shift. Unfortunately, we didn't had time to investigate this aspect and find a reasonable weighting function on those parameters. The candidate news passing a threshold are considered to be correlated with the shift, so for each of them the list of the  $k$  words with higher TF-IDF values are returned.

## Implementation

The methodology described above was implemented in python 2.7 using a library called gensim [2]. It was possible to store all the preprocessed data for a single topic on main memory and perform all the computation in a single pass. The choice of the number of dimensions of the LSA space is non-trivial. For small dataset like ours, a number between 50 and 100 has proven to be optimal [?]. At the same time, for a more specific discrimination in an homogeneous topic a larger number is suggested. For those reasons we used a 200 dimensions vectorial space.

## N-Gram Graph

N-Gram Graph[NGG] is a NPL tool initially proposed by George Giannakopoulos [3] that use word or character n-gram in order to achieve documents summarization. The NGG tool basically slice the text in word or character n-gram and then represent them in a graph  $G = \{N, E, L, W\}$  according the following structure:

- $N$  is the set of nodes created for every different n-gram in the text
- $E$  represent the edge of the graph; where two node are connected by an edge if they are "near" according to a window-distance
- $L$  is a labelling function which give the label to every node and every edge (define the size of the n-gram)
- $W$  is the weight function that gives weigh to every edge according to the number of times that two n-gram appear near one to the other

Whit this implementation NGG result to be a tool with the property to be language-neutral (make no assumption about the underlying language) and allow text manipulation trough graph operations.

In particular two operation are necessary for our goal:

- the **Intersection** operator between two graphs  $G_1$  and  $G_2$ : which returns a resulting graph with only the common edges of  $G_1$  and  $G_2$  averaging the weights of the original edges assigned as the new edge weights(example in figure 1)
- the **Normalized Value Similarity**[NVS] function that for every n-gram rank, indicating how many of the edges

contained in graph  $G_i$  are also contained in graph  $G_j$ , considering also the weights of the matching edges and normalize the result respect the graph size

In particular  $NVS(G_i, G_j) = \frac{VS(G_i, G_j)}{SS(G_i, G_j)}$  where:

$$VS(G_i, G_j) = \frac{\sum_{e \in G_i} \frac{\min(w_e^i, w_e^j)}{\max(w_e^i, w_e^j)}}{\max(|G_i|, |G_j|)} \quad (1)$$

$$SS(G_i, G_j) = \frac{\min(|G_i|, |G_j|)}{\max(|G_i|, |G_j|)} \quad (2)$$

Since NGG tools allow to compute both the tweet-news correlation and summary creation, this methodology is split in different section.

### *Tweet-News Correlation*

To compute the correlation between contradiction tweet and news, this methodology use the base idea of exploiting the NVS similarity function as a correlation function: higher the similarity of the tweets n-gram graph representation and the news n-gram representation, higher will be the correlation.

More in detail this procedure perform the following step:

- Merge all the tweets corresponding to a CP are in a single document n-gram graph representation.
- Compute the news n-gram graph for the news that fall inside a time slots defined by the contradiction point windows increased on both side by a time windows duration
- for the news that are inside the time slot, the NVS is computed and represent the correlation value between news and tweet

The usage of the time windows is needed given that NVS similarity values have no time domain impact, and with out the time windows isn't possible to discharge news that are distant in time form the correlation computation.

This approach require a lot of computation, since have to compute the graph representation of all the tweet and for every news inside the time windows compare the tweet with the news text. Instead NGG allow to compute the similarity between text without the usage of grammar information and use this similarity as correlation value: even if has no time relation.

### *Summary Creation*

The second goal of this research is to create a summary of the event that cause the sentiment shift.

For this purpose we use the news with the highest correlation value obtained at the step before. Since the goal is to create a summary of 150 word, we suggest to take in consideration only the best 4 news. Increasing the summary length can require more news.

The summary algorithm perform the following action:

- sort the news according the correlation value obtained before
- save all the sentence of text best 4 news in a set; the sentence detection is perform using OpenNLP
- create the intersection graph of the best 4 news
- for every sentence saved compute the NVS between the sentence n-gram graph and the intersection graph
- create the summary using the sentence with the highest NVS value until we reach 150 word

this algorithm create a summary composed by the most significant sentence of news with the highest correlation value, but is not able to remove redundant sentence since don't use any semantic index and redundant control. Even this, a summary composed by sentences result to be more human readable than a list of key word.

## **Experimental Evaluation**

### **Experimental setup**

In order to perform the experimental evaluation we used a tool to automatically detect sentiment shift in tweets coming from year 2009 and regarding several topics. We also downloaded news from the *New York Times* and *ABC Australia* on the same topic and spanning on the same period of the tweets.

It's important to remember that tweet have different language structure that need to be normalized, so for cleaning the tweet text the following operations were performed before the computation of the different methodologies:

- URLs removal from tweets using a regular expression.
- conversion from Unicode to ASCII.

After that, we manually labelled each contradiction point with the event which caused it. The topics, contradiction points and events used for the experiments are showed in table 1

### *Experiment preparation*

## **Conclusions**

### **1. REFERENCES**

- [1] Indexing by Latent Semantic Analysis, Scott Deerwester, Susan T. Dumais\*, George W. Furnas, Thomas K. Landauer and Richard Harshman
- [2] LANDAUER, Thomas K., et al. (ed.). “*Handbook of latent semantic analysis*”. Psychology Press, 2013.
- [3] ŘEHŮŘEK, Radim, et al. “*Software framework for topic modelling with large corpora*”. 2010.
- [4] Automatic Summarization from Multiple Documents : N-Gram Graph, George Giannakopoulos and Ncsr Demokritos, 2009

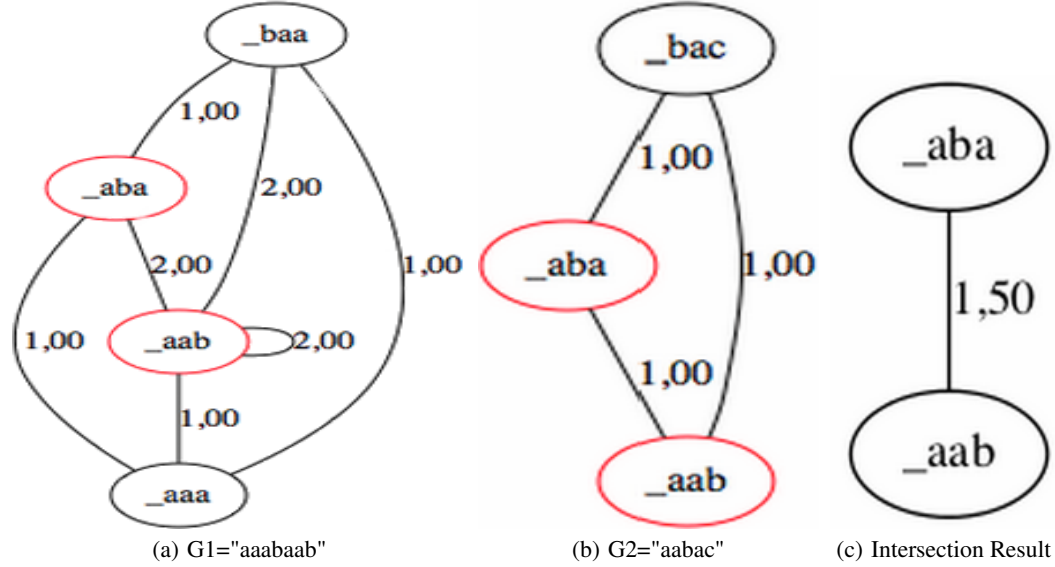


Figure 1: An example of the intersection operation

Contr. point	Contr. window	Event
Cern1	2009-10-29 - 2009-11-08	On November 3rd a bird drops a piece of bread which cause LHC overheating.
Cern2	2009-11-24 - 2009-12-04	On November 30th LHC accelerates protons to an energy of 1.18 TeV, becoming the world most powerful energy particle accelerator
FortHood	2009-11-02 - 2009-11-06	On November 5th a US marine kills 13 people
HangOver	2009-06-21 - 2009-06-27	undetermined or movie released on June 5th
Lcross	2009-10-31 - 2009-11-08	Preliminary findings from Lcross. Others announced
Jackson1	2009-06-19 - 2009-06-25	On June 25th Michael Jackson dies
Jackson2	2009-08-23 - 2009-08-29	On august 28th another popular musician, Adam Goldstein, dies. The day after, August 29th is Jackson's birthday
SwineFlu	2009-06-19 - 2009-06-23	The swine flu is recognized as a pandemic

Table 1: Contradiction points used for experimental evaluation

- [5] Linking Tweets to News: A Framework to Enrich Short Text Data in Social Media, Weiwei Guo, Hao Li, Heng Ji and Mona Diab
- [6] Modeling Sentences in the Latent Space, Weiwei Guo and Mona Diab , 2012
- [7] Multi-Sentence Compression: Finding Shortest Paths in Word Graphs, Katja Filippova
- [8] Metwally A. et al, “*Efficient Computation of Frequent and Top-k Elements in Data Streams*”, Lecture Notes in Computer Science Volume 3363, 2005, pp 398-412