

# Social Media Analytics - Italian Referendum - Temporal Analysis and Spread of Influence

VIGÈR DURAND AZIMEDEM TSAFACK, PRATUAT AMATYA

MALICK ALEXANDRE N. SARR

UNIVERSITY OF ROME "LA SAPIENZA"

July 13, 2019

## Abstract

*Social media has completely transformed the way people interact with each other, having an equally significant impact on businesses as well. There are now billions of unsolicited posts directly from consumers that brands, agencies and other organizations can use to better understand their target customer, industry landscape or brand perception. [1] In this experiment, as an illustration of how useful social media could be, we analyzed two data-sets, a network representing the relationships among twitter profiles and a massive quantity of data containing the stream of tweets across the Italian referendum day (4th of December). On those data-sets, we investigated the temporal distribution by performing a temporal analysis before Identifying potential Yes/No supporters. Additionally, we performed spread of influence analysis to evaluate the spread of the different parties over the network.*

## I. INTRODUCTION

Recent research on social media analytic has emphasized the need to adopt a business intelligence based approach to collecting, analyzing and interpreting social media data. [2] Social media presents a promising, albeit challenging, source of data for business intelligence. Customers voluntarily discuss products and companies, giving a real-time pulse of brand sentiment and adoption.[3]

Social media is one of the most important tools for marketers in the rapidly evolving media landscape. Firms have created specialized positions to handle their social media marketing. These arguments are in line with the literature on social media marketing that suggest that social media activities are interrelated and influence each other.

In this experiment, as an illustration of how

useful social media can be, we analyzed twitter data across the Italian 2016 constitutional referendum date in order to evaluate the coherence between our results and the actual outcome of that referendum.

## II. METHODS

### i. Data

The provided data-set contained a generated network consisting of over 450.193,0 twitter IDs organized in around 22.649.482,0 of edges (connections among twitter profiles). Each edge was represented on a line as follows: *source < tab > target < tab > weight*.

Additionally, We were provided a massive quantity of data containing the stream across the 4th of December divided by day, each file containing 10.000,0 tweets. The overall dimen-

sion of available data was around 10 Gigabytes.

## ii. Temporal Analysis

### ii.1 Tweet distribution over time

In order to perform analysis on the file, we needed to find a way to store, access and query the information provided within the data. For that purpose, we created 3 main indexes, related to the temporal analysis. All the text, containing relevant information have been cleaned before being saved. The built-in Lucene Italian analyser has been used along with a list of Italian stop-words taken online to omit words containing little information. All the Tweet text has been stemmed, and stripped of not needed symbols such as 'https', 'rt', etc.

The first index created was the tweet index, which contained all information about the tweet. That includes but is not limited to the tweet date, users, screen-name, text, hashtag, mentioned people. Few additional functionalities such as filtering, searching and getting tweet frequency were added to manipulate the data-set.

The next index to be created were the politician index. We first scrapped a Wikipedia page (Endorsements in the 2016 Italian constitutional referendum) containing the list of Italian politicians along with their affiliation (Yes or No) during the 2016 Italian referendum. Once that list was obtained from the web, we created a twitter lookup function that would associate the twitter username to the politician name. In order for a politician to be considered in this analysis, he was supposed to have at least 5000 followers.

Finally we separated our tweet index between politicians having voted yes and politicians having voted no. Upon doing that we then plotted the number of tweet generated by each camp over time.

### ii.2 Time series representation of tweets

In this part we performed SAX analysis on the politicians tweet for each camp (yes/no). We

can think of SAX (Symbolic Aggregate approximation) as a symbolic representation of time series. The basic idea behind SAX is the fact that it allows a time series of arbitrary length  $n$  to be reduced to a string of arbitrary length  $w$ , ( $w < n$ , typically  $w \ll n$ ). The alphabet size is also an arbitrary integer  $a$ , where  $a > 2$ . The algorithm consist of two steps: (i) it transforms the original time-series into the PAA representation and (ii) it converts the PAA data into a string. The use of PAA brings advantages of a simple and efficient dimensionality reduction while providing the important lower bounding property. The actual conversion of PAA coefficients into letters by using a lookup table is also computationally efficient and the contractive property of symbolic distance was proven by Lin et al. [?]

In order to apply the SAX analysis, for each camp, we first pulled up the earliest and latest date in our index in order to give us a proper range for the event then we extracted to top 1000 more frequent words tweeted by the politicians. Once we had the more frequent terms, we counted the occurrence of each term over time with an interval of 12 hours. Once we had a temporal distribution of the most frequent term used by each camp with a grain of 12 hours, we converted that temporal distribution to a SAX string using the SAX package for Java. The alphabet size we picked was 20. Finally we perform a unsupervised classification of the above SAX string by applying KMeans clustering on the TimeSeries.

The Kmeans algorithm is a clustering method that aims to find  $k$  centroids, one for each cluster, that minimize the sum of distance of each data point from its respective cluster centroid. It solves, for  $x_i \in X$ :  $\argmin_{C_1, \dots, C_k} \sum_{j=1}^k \sum_{x_i \in C_j} d(x_i, c_j)$  where  $(C_1, \dots, C_k)$  are  $k$  non-overlapping clusters,  $c_j$  is the representative of cluster  $C_j$ , and  $d$  is a distance function. [10]

We chose our initial centroids by picking random points withing our cluster. The distance measure used is the Jaro Winkler distance, mainly because it gave us more balanced cluster. But the clustering was tried as well

with the Euclidean distance and the Levenshtein distance. For each sax string, we calculated the distance to each centroid and we added it to the cluster belonging to the centroid with the lowest distance. Once that was done, we updated the centroids for each cluster, by averaging each char  $c$  at position  $n$  out of  $N$  of the data-points within the cluster. We repeated the same procedure until there were no changes or we reached the maximum iteration flag.

### ii.3 Time series co-occurrence graph

In order to go the co-occurrence analysis, we performed 3 main steps. We first converted the terms within the cluster to a graph, having as nodes the terms, and as edges their weight, which correspond to a normalized frequency measure obtained by extracting the frequency of occurrence of the two words. This weight will correspond to the total number of tweets that contained both term. We then applied a threshold value  $t$  which correspond to the highest frequency between the two term multiplied by a probability value  $p$ .  $t = \text{argmax}(t1, t2) * p$ . Finally an edge between two nodes/terms is added if their weight is above to the threshold value.

Once the graph has been created, we then proceeded in extracting the Largest Connected Components within the graph. A connected component of an undirected graph is a maximal set of nodes such that each pair of nodes is connected by a path. [13]. For each cluster, its largest connected component is saved, this is done with a helper graph java library called G.

Our final step consisted of extracting the k-cores of the graph. The k-core of a graph  $G$  is the maximal induced sub-graph  $H \subseteq G$  such that  $(G) \setminus H$  has degree less than  $k$ . Thus all vertices of  $H$  are adjacent to at least  $k$  other vertices in  $H$ . [12] The k-core was introduced by Steven B. Seidman in a 1983 paper entitled Network structure and minimum degree. In other words, k-core of a graph  $G$  is the largest induced sub-graph of  $G$  in which every vertex has degree at least  $k$ . The coreness

of a vertex  $v$  in  $G$  is the largest value of  $k$  such that there is a k-core of  $G$  containing  $v$ . [11]

### ii.4 Time series comparison

In order to compare the various time series, we saved the LCC and the Kcores for each cluster and for each camp(yes/no politicians). We those graph, we then went ahead and extracted the number of occurrences of each node (term) in an interval of 3 hours, then we graphically compared the behavior in order to observe if there were any similarity in the way those term occur over time.

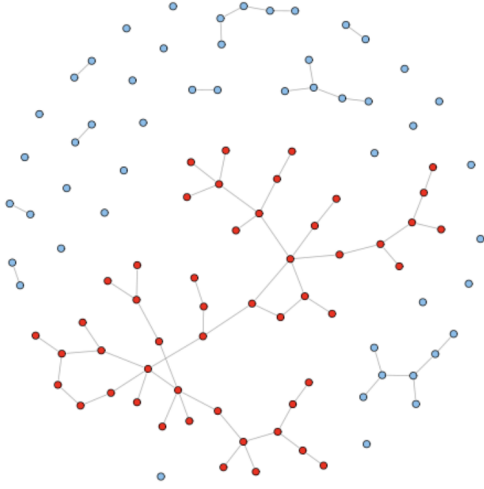
## iii. Identify mentions of candidates and Yes/No supporters

### iii.1 Largest Connected Components

Largest Connected Component (LCC) of a graph is the largest possible subset of the original graph with all the vertices connected through a path of any possible length. Real world networks often comprise of multiple disjoint component of nodes. The good majority of the nodes from original graph often forms the Largest Connected Component and hence analysis done on LCC generalises well for whole of the graph giving the advantage of reduced computational need due to smaller number of nodes and avoidance of complexity due to highly clustered or disjoint graph components.

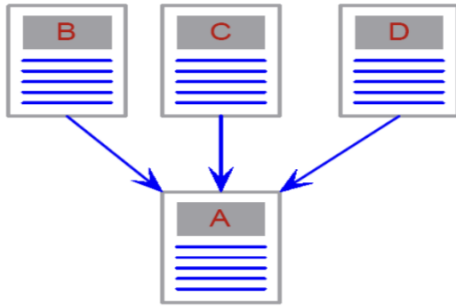
### iii.2 HITS analysis

HITS (Hyperlink Induced Topic Search) is a algorithm developed by Jon Kleinberg in IBM Almaden which is used to identify important group of nodes within the graph based upon their hubness and authority. Originally developed to represent a network of web pages where an authority is a node with high in-degree of link representing its authoritative content while a hub represents a node which work as an index with many out-links to other authority pages. The same philosophy of authority and hub is widely applicable in real



**Figure 1:** Largest connected component in a graph.

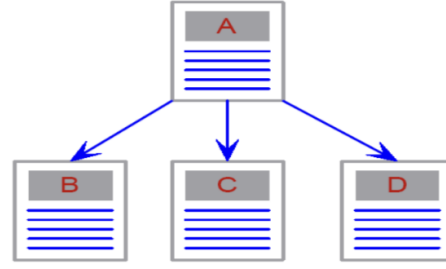
world directed graphs as in social network graph.



**Figure 2:** An authority node with high in-degree centrality.

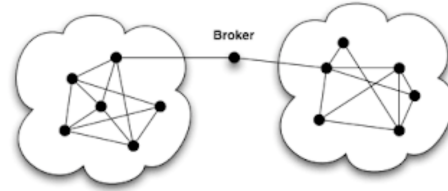
### iii.3 KPP-NEG algorithm

Key Player Problem - Negative (KPP-NEG) is a greedy algorithm used to compute a prestige measure of node being a broker between different component of a graph. A broker is a node that connects multiple segments of clusters of nodes to each other and their removal causes reduction of graph connectivity to greater extent compared to other nodes. Brokers are also



**Figure 3:** A hub node with high out-degree centrality.

called key separators in literature.



**Figure 4:** A broker node as key separator between two clusters.

### iv. Spread of influence

A social network—the graph of relationships and interactions within a group of individuals—plays a fundamental role as a medium for the spread of information, ideas, and influence among its members.[5] In this study, one of our main goals was to measure the spread of Yes/No over the provided network. To achieve this goal we focused on two main algorithms: *K – means* and *LPA* (Label Propagation Algorithm).

#### iv.1 Spread of Yes/No: K-means

K-means clustering is a type of unsupervised learning, which is used when you have unlabeled data (i.e., data without defined categories or groups). The goal of this algorithm is to find groups in the data, with the number of groups represented by the variable *K*.

The algorithm works iteratively to assign each data point to one of the  $K$  groups based on the features that are provided. Data points are clustered based on feature similarity.[4] The following is a pseudo-description of the K-means implementation we performed:

**Initialization step:** set  $K = 2$ , initialize the  $K$  clusters using the previously identified Yes/No seed nodes and randomly choose  $K$  nodes in the  $K$  initial clusters to be the initial centroids for the graph.

**Iterative steps:**

1. Associate every node to the centroid whose travel distance to the node is minimum, thus creating  $K$  clusters of nodes. The metric chosen to associate a node to a centroid can be understood by thinking that every node wants to belong with the centroid that can reach that node in the fastest way and having a similar neighborhood Yes/No labels with the node.

$$d(c, n) = \frac{1}{\cosSim(c, n)} + (1 - lSim(c, n)) \quad (1)$$

Where:

- $d(c, n)$ : is the distance between the centroid  $c$  and the node  $n$
  - $\cosSim(c, n)$ : is the cosine similarity between the centroid  $c$  and the node  $n$
  - $lSim(c, n)$ : The percentage of neighbors of the node  $n$  having the same label(Yes/No) as the centroid  $c$
2. Now that we have clusters of nodes we have to pick the best node to be the centroid of every cluster. The optimal node is the one that can reach all the other nodes with the lowest travel cost. The travel cost being the sum of the travel distances to every other node in the cluster.

**Convergence:** After a certain number of iterations, the centroids will not change anymore and the algorithm will have converged. If this does not happen, the algorithm will stop at a given maximum number of iterations.

## iv.2 Spread of Yes/No: LPA

Label Propagation is a semi-supervised machine learning algorithm that assigns labels to previously unlabeled data points. At the start of the algorithm, a (generally small) subset of the data points have labels (or classifications). These labels are propagated to the unlabeled points throughout the course of the algorithm. [6] Label propagation is an algorithm for finding communities. [7]

In comparison with other algorithms (for example with K-means) label propagation has advantages in its running time [8] and amount of a priori information needed about the network structure (no parameter is required to be known beforehand). The disadvantage is that it produces no unique solution, but an aggregate of many solutions.

In this experiment, We adopted two different strategies while using the *LabelPropagationAlgorithm*.

- The first strategy was to modify the original *LPA* in order to be able to evaluate the spread of only Yes/No labels in the network. In particular, the following are the steps we relied on:

1. **Initialization of the labels:** for each node, we set the label to be -1 if the node is already known to be a *Yes* supporter, -2 if the node is known to be a *No* supporter and a unique integer in the interval  $[0, +\infty[$  for the remaining nodes.
2. **Update of the labels:** for each node in the graph, the new label is the most popular (frequent) among all the Yes/No neighbors. if the node does not have Yes/No neighbors, we kept the previous label. This was to ensure the propagation of only Yes/No (-1/-2) labels.

- The second strategy was to run the original *LPA* 10 times and obtaining the final community for a given node by averaging the 10 different outputs. Additionally we illustrated the importance of the Normalized Mutual Information measure (NMI -

a common measure of goodness of a community detection algorithm) by computing it over the 10 *LPA* runs.

### III. RESULTS AND DISCUSSIONS

#### i. Temporal Analysis

##### i.1 Tweet distribution over time

From the 26 of November 2016 to the 6 of December 2016, the number of tweets generated in our stream at any given time in a span of 12 hour ranged in around a million. We observed that there were more tweets during the day and the we had a peak of tweets during voting day at about 1.4 million tweets. This already shows us the importance tweeter has as a social network. Figure 6 shows us the distribution of tweets between our dates of interest at any given time. The next step will be to focus our attention to politics and specifically the Italian referendum of 2016.

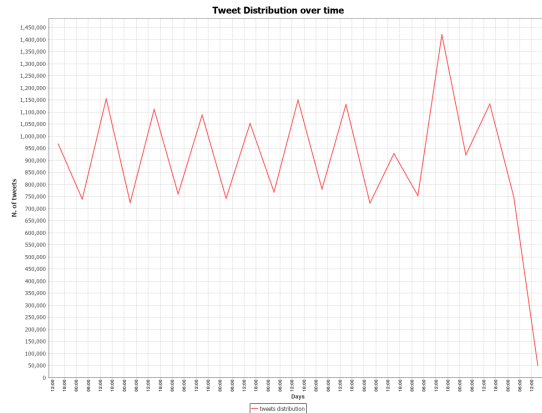


Figure 5: Tweet distribution over time for all tweets

For the temporal analysis, we worked with a total of 290 politicians. Those politicians were divided into two camps, the politicians promoting the vote of YES, referred in this paper as the Yes camp and the politicians promoting the vote of No, referred as the No camp. The Yes camp counted 157 members and the no camp counted 133 members. Figure 6 shows the total tweet of each party over time in an interval of 12 hours.

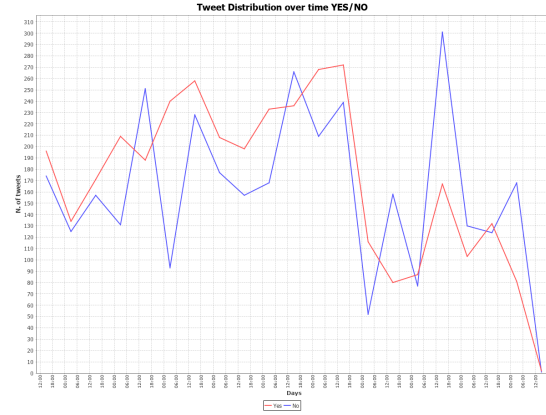


Figure 6: Tweet distribution over time for YES/NO Politicians

Some interesting fact to note here is the fact that both party tweet approximately the same amount of tweets with the people from the yes camp tweeting slightly more the days leading to the vote. During and after the voting day, we observed an increase in the amount of tweets of the no politicians.

One way we could improve was by including additional politicians basing their vote on their respective a party affiliation, this will therefor increase the amount of data we are working on. In addition to that, we could have added some popular influencers that had a say in the referendum and that worked closely with a particular party. Those include but are not limited to journalist, TV personality, athletes, etc.

##### i.2 Time series representation of tweets

The 20 most popular terms in our index are no, solo, sempr, quand, oggi, me, cosa, graz, video, via, ora, fa, te, così, ital, renzi, fare, mai, poi, giorn. We can see some popular terms used by many politician such as thank you, today, do, always, but as well the name of a politician Renzi. This suggest that he was getting in to a lot of headlines for his term to be used very frequently.

With timeserie representation of the tweet, we divided each camp into 20 clusters, each term was represented my a sax string of 20

character. The representation of 1 cluster can be seen in table 1. The rest of them are available within the files attached to this paper.

Term	SAX Representation.
risultat	ejojejejejejejejejejeje
bilanc	eeetrjoootojjejejejeje
punto	ffrrfttffrrrffflfff
social	tfotiiloosdillfddfddd
battagl	ffjfffjtjtrffrffrfff
tropp	jijdtptjtdtjdjdtdjd
staser	jijptjftfjfmffjffff
megl	pjjpjjjjdjttdddjdjdd
via	idcoilsiotllcidodocc
senator	rejjtoeoeerjteeejee
propr	jiddjdjojjdtdootdjd
smontat	gggggkrogkgktggggggg
public	ejelljptpeneejeejge
civil	ftffrrrrflfrffflfff
mille	eeeejtjteejjjeeestee
camer	jijjtrtejeejeejeeeee
comun	jijjppddtjtddddjjjd
parlamentar	tjjjjjdjptjtddddjjd

**Table 1:** Sample of SAX representation for Cluster 1 of No Politicians

One trend that we realized tend to happen is that most of the words end up clustered in 2 or 3 clusters, ergo we end up having 2-3 cluster with the bulk of the elements while the rest have very few.

A future change of improvement would be on the kmeans algorithm. Indeed we are using the very basic kmeans algorithm which has its limitation such as those unbalanced clusters. A better version was proposed by Saeid Soheily-Khah. [10] addresses those limitations.

The data and the clusters along with the sax representation of the words are available in the attached coded under resources/yesCluster and resources/noClusters.

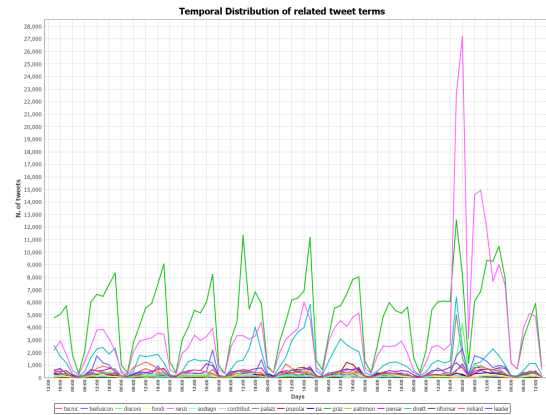
### 1.3 Time series co-occurrence graph

Since the clusters were very unbalanced, one thing that we realised is that not all clusters will have LCC or Kcores. Indeed for the No camp about 1/3 of the clusters did not have LCC and

Kcores whose weight would have been significant for analysis. And for the yes camps about 4/5 of the clusters did not have any significant co-occurrence. Another observation is that the LCC graphs are slightly bigger that the Kcores one in terms of nodes.

The co-occurrence full result can be found in the resource folder accompanying this file under resources/cooc.

Looking at some of the co-occurrence graphs and their distribution over time with a granularity of 3hours, we realize that most of the words within the kcore and LCC share a some what similar temporal behavior. For example, looking at figure 7, we we can see that there is a somewhat same temporal behavior between Renzi and Graz. We can here assuming that he's been using that word quite often and people have been thanking him as well. We can see as well the word "support" that some what behaves the same way.



**Figure 7:** Kcore temporal representation for Cluster 0 No camp

We can see as well similar type of behavior for the LCC. In this example we show the LCC of cluster 0 for the Yes camp in figure 8. In this case, we see, as most visible relation, approximately the same temporal behavior between the term analysis and contract. We can easily assume that the were taking about analysing certain contract, probably bills to be passed. Additionally, contract has somewhat the same temporal behavior as the term signature. Ergo this cluster somewhat shows a tendency for the







### ii.3 Partitioning of users in YES/NO group

The task of this section was to classify the user set  $M$  into YES and NO vote users taking into consideration the candidates they mention. For the computational convenience we initially built a Lucene index USER POLITICIAN INDEX with the documents comprising of just two fields *userScreenName* and *politicianScreenName*. The document entry was created for each single mention of a politician by a user in the USER TWEET INDEX. Then on for each unique user in set  $M$ , we counted number of mentions of both YES and NO politicians. The tally count for total number of mentions, and mentions of either group of politicians by unique users was generated and saved into *data/resources/user\_tweet\_count.csv*. Depending upon the frequency of mentions for YES and NO group, whichever is greater, we classify the user to that group, i.e,

$$\text{if } n_y > n_n \text{ then } u \in \text{YES} \quad (3)$$

$$\text{if } n_y < n_n \text{ then } u \in \text{NO} \quad (4)$$

Using above classification measure, user set  $M$  was classified into two groups and HITS analysis was performed onto respective largest connected components of the sub-graphs to obtain 1000 authority users for each group, namely  $M'$ . The sub-graph extraction and HITS analysis was performed using  $G$  library and resulting 1000 authority users for each group was saved into files *yes\_users\_top\_authorities.csv* and *no\_users\_top\_authorities.csv*.

ID	Screen Name	Score
18762875	matteorenzi	0.182002
18935802	repubblicait	0.160015
5893702	SkyTG24	0.139088
150725695	Agenzia_Anza	0.135655
963938472	Palazzo_Chigi	0.135041

**Table 4:** Top authority from YES users.

ID	Screen Name	Score
19067940	beppe_grillo	0.157597
1530798872	virginiaaggi	0.133566
289400495	Mov5Stelle	0.122574
376218450	AndreaScanzi	0.120515

**Table 5:** Top authority from NO users.

The top authority users include Matteorenzi, former Italian prime-minister who advocated for the referendum, Repubblica.it news provider and Sky\_tg24 broadcasting. Due to very high number of follower base, these user accounts certainly commands for high authority within the network. As for NO votes, top authority users are Beppe Grillo, co-founder of Five Star Movement, Virginia Raggi, current mayor of Rome and a representative of Five Star Movement and Mov5Stelle, the official twitter user for the Five Star Movement itself. The movement's stance was against the referendum and hence with high follower base, they stood out as authorities for NO vote.

### ii.4 KPP-NEG analysis

Lastly in this section we used KPP-NEG algorithm on the user sub-graph  $S(M)$  to identify key players. Because of high computational cost we applied varying threshold value for each node's out-degree centrality to reduce the graph size. The resulting users alongside their scores for each YES and NO group were saved to *yes\_users\_500KPP.csv* and *no\_users\_500KPP.csv*.

ID	Screen Name	Score
1386632492	miaudit	5358.0
437337751	Esercito	4287.0
2226802248	SquilibrateS	4286.0
16402819	andreadelogu	3214.0
18935802	repubblicait	3214.0

**Table 6:** KPP-NEG top users for YES group.

ID	Screen Name	Score
434839840	lauracesaretti1	2552.0
2651374166	ScifoGiuseppe	2552.0
16476536	ultimenotizie	1912.0
376218450	AndreaScanzi	1912.0
2455915345	tiz_mic	1912.0

**Table 7:** KPP-NEG top users for NO group.

### iii. Spread of influence

#### iii.1 Spread of Yes/No: K-means

dummy text dummy text dummy text dummy  
text dummy text dummy text dummy text  
dummy text dummy text dummy text.

#### iii.2 Spread of Yes/No: LPA

dummy text dummy text dummy text dummy  
text dummy text dummy text dummy text  
dummy text dummy text dummy text.

## REFERENCES

- [1] Crimson Hexagon The Fundamentals of Social Media Analytics: <http://www.upa.it/static/upload/the/the-fundamentals-of-social-media-analytics.pdf>
- [2] Umar Ruhi (2014), "Social Media Analytics as a Business Intelligence Practice: Current Landscape Future Prospects", Journal of Internet Social Networking Virtual Communities, Vol. 2014 (2014), Article ID 920553, DOI: 10.5171/2014.920553
- [3] Lu, Y., Wang, F., Maciejewski, R. (January 01, 2014). Business intelligence from social media: a study from the VAST Box Office Challenge. IEEE Computer Graphics and Applications, 34, 5.)
- [4] Andrea Trevino, Introduction to K-means Clustering; <https://www.datascience.com/blog/k-means-clustering>.
- [5] David Kempe, Jon Kleinberg, Éva Tardos, Published April 22, 2015 Maximizing the Spread of Influence through a Social Network
- [6] Zhu, Xiaojin, Learning From Labeled and Unlabeled Data With Label Propagation.
- [7] U.N.Raghavan – R. Albert – S. Kumara, Near linear time algorithm to detect community structures in large-scale networks.
- [8] M. E. J. Newman, Detecting community structure in networks.
- [9] Jessica Lin Eamonn Keogh Stefano Lonardi Bill Chiu A Symbolic Representation of Time Series, with Implications for Streaming Algorithms
- [10] Saeid Soheily-Khah. Generalized k-means based clustering for temporal data under time warp. Artificial Intelligence [cs.AI].
- [11] Wissam Khaouid, Marina Barsky, Venkatesh Srinivasan, Alex Thomo K-Core Decomposition of Large Networks on a Single PC
- [12] Allan Bickle The k-Cores of a Graph
- [13] Massimo Franceschet Network Connected Components

## IV. APPENDICES