# Generating a Highlight Moments Summary Video of an political even using social media speech Sentiment Analysis ontologies

Mehdi.ABID, Benayad NSIRI and Yassine SERHANE

*Abstract*— Numerous viewers choose to watch political or presidential debates highlights via TV or internet, rather than seeing the whole debate nowadays, which requires a lot of time. However, the task of making a debate summary, which can be considered neutral and does not give out a negative nor a positive image of the speaker, has never been an easy one, due to personal or political beliefs bias of the video maker.

Our study came up with a solution that generates highlights of a political event, based on twitter social network flow. We used twitter streaming API to detect an event's tweets stream using specific hashtags, and detect on a timescale the extreme changes of volume of tweets which will determine the highlight moments of our video summary at first, then we set up a process based on a group of ontologies that analyze each tweet of these moments to calculate the sentimental score of each, then classify those moments by category (positive, negative or neutral).

*Keywords*—Debate summary, API, hashtags, Twitter, Highlights Moment, Ontologies, Sentiment Analysis.

## I. INTRODUCTION

IN the 2016 Republicans primaries, CNN claimed that more than 84 million people have watched the republican candidates debating on its channel, breaking records for most events seen on CNN. FOX also cited that more than 83 million have seen the debate between the republican candidates, which made it the most watched event in the history of television.

The majority of these audiences are social media users, who respond to every controversial moment [1] on various

Platforms in real time, such as Twitter, Facebook, Snap, Instagram etc.

Mehdi ABID Laboratoire d'informatique et d'aide à la décision. Faculté des Sciences Ain Chock, Université Hassan II. Casablanca MAROCO (phone: 00212- 677-115323; e-mail: 90.abidmehdi@gmail.com).

Benayad NSIRI Laboratoire d'informatique et d'aide à la décision. Faculté des Sciences Ain Chock, Université Hassan II. Casablanca MAROCO (E-mail : nsiri2000@yahoo.fr)

Yassine SERHANE Laboratoire d'informatique et d'aide à la décision. Faculté des Sciences Ain Chock, Université Hassan II. Casablanca MAROCO (phone: 00212-677-115323; e-mail: serhane.y@gmail.com).

In our study we have used Twitter as our main audience feedback source [2], [3], because of its worldwide use (**Figure 1**), and people use it more than other social platforms to express their immediate feelings and opinions.

Several studies gave an interesting insights about the social network twitter evolution due to his dynamic nature with more than 400 million tweets posted everyday [4], using the hashtags (Significant continuation of characters without space beginning with the sign #, Which refers to a subject and inserted into a message by its author, in order to facilitate the location) we can look for trending topics and look up thousands of tweet (**Table 1**).
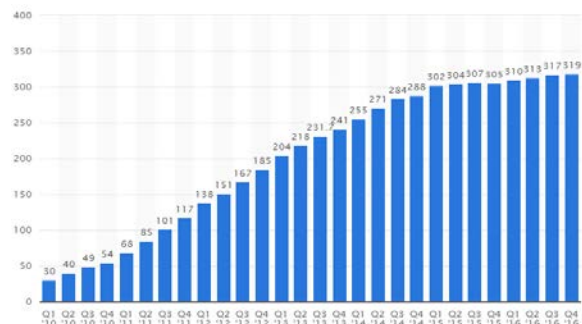


**Figure 1 : Number of quarterly active twitter users in millions**

| Topics (Hashtag) | # Tweets | Time span |
|---|---|---|
| Black Friday | 1 085 365 | 2016.11.01 -- 2016.11.01 |
| Trump | 410 854 | 2016.11.01 – 2016.11.30 |
| Iphone7 | 98 716 | 2016.06.30 – 2016.07.31 |
| FCbarcelona | 302 523 | 2017.02.01 – 2017.02.28 |

**Table 1 : Basic information of a dataset**

Twitter users usually respond to political speaker statements or point of views during a political speech, which offers a fertile ground for sentiment analysis [5], due to the outrageous tweets against the opposite political speaker or the encouraging tweets from their supporters, those tweets usually

come as a reaction to the big (Good or bad) moments of the speech, which makes their reactions a good highlights indicator for our event.

In this article, we have used the volume of these tweets in a preset amount of time as an indicator of an event highlight, gathering those highlights to wind up with a summary video generated only using random samples of tweets which grant our summary the neutrality and avoid unwanted bias.

Furthermore, we came up with an approach based on the KDD process (knowledge discovery in databases) which will be discussed in the next chapter, this approach utilizes these tweets to calculate the sentimental score of each, in order to classify these tweets into positive, negative or neutral and be able to determine the nature of the moment.

## II. Background

The construction of a summary video, which generates the highlight moments of a political, social or cultural event, is usually based on the image processing of the event, this process is basically established on objects detection [6], [7]. This detection of objects in the case of a speech, can be based on the change of the camera angle from speaker towards public, whenever the audience applauds or boo to show their disapproval or start chatting about controversial statements. However, this approach becomes delicate in cases where certain events or rallies take place where the spectators are seated behind the speaker.

Moreover, other studies rely on approaches that are based on the exploitation of audio by detecting any sudden variation of sound recording [8]–[10] caused by audience reaction, such as applauds or shouts out, this variation reflect the highlight moment during the event (**Figure2**), but unfortunately this approach can be biased and have several inconveniences such as special effects added during the event, or even the presence of noise during the whole event.



**Figure 2 : Highlight moment detection based on sound frequency**

Due to the inconveniences of objects detection and audio recording analysis approaches to detect the highlight moments of a specific event, and the evolution of social media use, we can utilize this massive quantity of data generated from these social media platforms especially twitter, to generate a summary of this event.

Similarly, several research works, such as predicting a movie success based on the reaction tweets from the trailer watching [11], and the prediction of the presidential elections established in several countries such as the USA, France and Pakistan etc.[13] – [15].

Those approaches have shown a major success in their predictions, which proves the credibility of using the social network twitter as a source of information to figure out the public tendency.

Including, the research [12] which use Twitter's data to predict the results of the Pakistan's elections in 2013, thanks to a model of classification developed in Machines Learning by using learning algorithms, in order to classify the tweets into two categories positive (Pro) or negative (Anti), through the sentiment analysis of every collected tweet, this classification is based on the contents of tweets (Hashtags and key word) eg. the use of capital letters which means a person is shouting, words, emoticons etc. and then a comparison is done by attributing every tweet to the appropriate presidential candidate.

Furthermore, other studies have also been based on the sentiment analysis process of the tweets, i.e. due to the feelings polarization of their spectators during a soccer match [15], which can be identified thanks to the use of the standardized hashtag or the the one made official by their team. This approach creates a framework that handles various reactions from numerous Twitter users during a soccer match [16], and showed -as expected- positive results, the tweets from users are positive when their team scored a goal and negative if they concede one.

In addition, some research were developed on fans swearing in tweets, while watching a soccer match and how they used it as a sentiment marker [15]. Their work concentrate heavily on the context of the tweet rather than the swearing itself, because not all swearing tweets reflect negative sentiment. They started by collecting tweets in relation with the English Premier League matches, then they linked these tweets to teams based on how many times a fan tweeted using his team hashtag the most, after that, they filtered these tweets by use of swearing, taking into consideration complication like fans using their opponent hashtags to get their attention. They conclude their work by showing that bad language is not always negative and some of the strongest sentiments expressed are self-critical.

Most of the studies described previously, have used in their approach various methods of data mining, such as KDD process, which is used widely in the research field, or using process intended for the professional area such as CRISP-DM (Cross Industry Standard Process for Data Mining) which is considered as an iterative process, and strongly used to satisfy the industrial needs (Domain of engineering, medicine, sales and marketing.) [17], [18].

In our study, we will be using the KDD process, because it is complete, precise and answers our needs, which is the search for the knowledge in big data.

Knowledge Discovery and Data Mining is a process that allows the extraction of the different information out of the massive data according to a predefined goal, in order to find oneself with a useful knowledge [19],[20] (Figure 3).
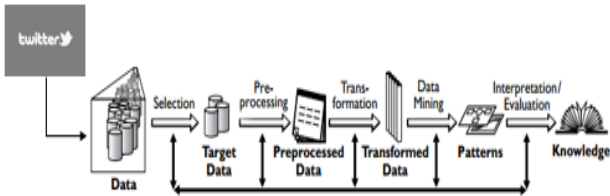
**Figure 3 : KDD processes**

This process is composed of 5 main steps (Selection, Pretreatment, Transformation, Data mining, Interpretation) [21].

- Selection: Consists of collecting and choosing the data which results in aggregating a variety of sources into a single target data.
- Pre-processing: Stage that contains the removal of noise and handles the missing values into clean target data.
- Transformation: In this phase, every data is transformed through the reduction of the database dimensions, and the transformation of the attributes, to wind up with a database that meets the requirements of our project objectives.
- Data Mining: This stage consists of choosing and adapting the algorithms of data mining based on intelligent methods in order to extract data patterns.
- Interpretation: is the final stage of this process, which includes the evaluation and the interpretation of the patterns discovered in order to determine the useful information.

### III.  METHOD

The moment there is a broadcasted political event live on television, users begin to tweet about it using related Hashtags, in order to share their opinion and symbolize them in relation with the theme of this event.

Thanks to Twitter's streaming API, we are able to recover the contents as well as the volume of tweets by their Hashtags in real time via a request sent to the twitter's servers, which allows to obtain a stream of data $\{(x_i, y_i), i = 1,..., n\}$; Taking into example two features of the data it can be represented in the form of a cloud of points of data in a (x, y) plan (**Figure 4**), where the x-axis represents speech time interval and the y-axis represents the number of tweets.
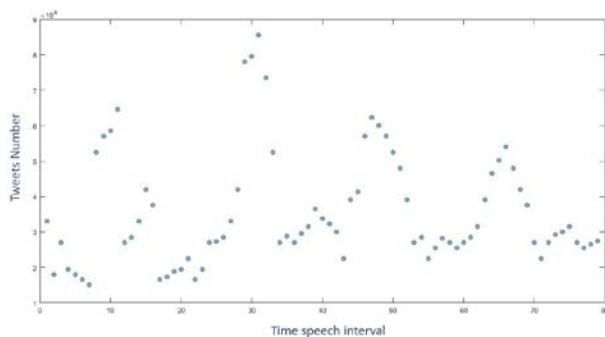


**Figure 4: Tweet Data Volume represented as a scatter graph**

Our research purpose is a summary video generation based on the highlight moment detection of an event and the analysis of the sentiment of these moments tweets, i.e. The detection of tweets volume extreme changes on a timescale at first, and then analyze the sentiment of each tweet belonging to the highlight moments in order to measure the percentage of its positivity.

To achieve that, the computing of a function   that would allow to reflect the partner of points   obtained on a graph, remains indispensable even though we did not explicitly know it. However, the mathematical approach used is the optimization of Lagrange polynomial obtained from the plot described previously, the existence of this polynomial is asserted by the following theorem [22]:

There is a unique polynomial $p_n \in \mathbb{R}_n[x]$, ($\mathbb{R}_n[x]$ being the vector space of polynomials which degree is lower or equal to n) such as:

$$p_n(x_i) = f(x_i) \; \forall i \in \{1, 2, ..., n\}$$

and  $P_n$ is given by Lagrange formula:

$$P_n(x) = \sum_{i=0}^{n} f_i \, l_i(x)$$

With

$$f_i := f(x_i) \text{ et } l_i(x) = \prod_{j \neq i} \frac{x - x_j}{x_i - x_j} \, , \, i \in \{1, 2, ..., n\}$$

$P_n$ is called the polynomial of interpolation of Lagrange in points $x_i$ for the measures $f_i$.

The theorem above allows to create a polynomial function passing by all the points obtained, e.g. as represented in **Figure 5**.
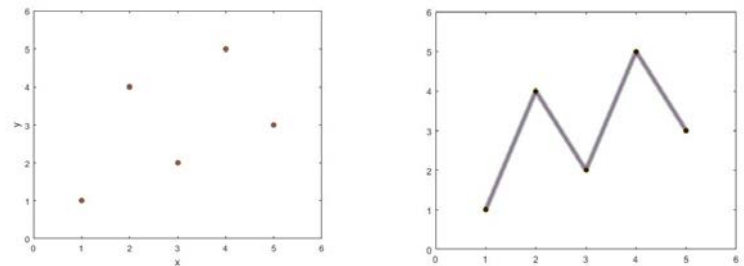


**Figure 5: Lagrange polynomial for the interpolation points**

However, the peaks of this polynomial function which varies according to the change of the tweets volume in a predefined period lead to the detection of the highlight moments. The latest can be determined through the spikes, which are the local maximums of the polynomial function $p_n$ i.e. Points that satisfy the optimality conditions:

- Condition 1 (Stationarity):

$$\frac{\partial p_n}{\partial x}(x_i) = 0$$

- Condition 2:

$$\frac{\partial^2 p_n}{\partial x^2}(x_i) < 0$$

First, we should use the method of steepest descent for the stationary points of $f$ , after that, a simple selection of the points with a positive second order derivative, will take us to our objective (Peaks Detection). (**Figure 6**)
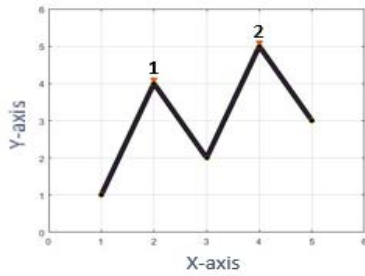


**Figure 6: Peaks detection**

Once the highlight moments is defined by generating the polynomial function and by detecting its peaks, we discover the nature of these highlight moments by applying the sentiment analysis process on each tweets that belong to every peak, in order to compute the tweeter users' sentiments scores toward the speaker**.**

The appropriate process of sentiment analysis comes down to developing a process that allows a classification of published tweet's sentiment, where data extracted from Twitter is  analyzed in a granular way, by decomposing sentences into a group of words linked to a global ontology that includes various types of terminology. The aim of our sentiment analysis process is the ability to analyze a sentence and to compute its sentiment score (Positive, negative and neutral) (Figure 7). However the use of the ontologies in our analysis will have numerous advantages, in particular with regards to the cultural, linguistic and regional expressions... [23].
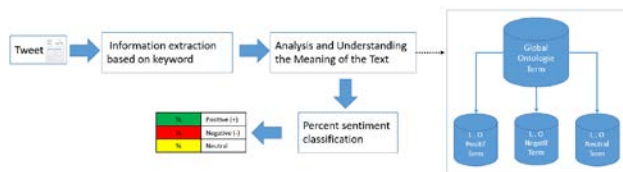


**Figure 7 : Process of classification percentage sentiment Tweets**

The global ontology used above allows to regroup different local ontologies, which describe their own local knowledge space in relation to a precise specification of each word or sentence category (positive, negative or neutral), in other words, each local ontology contains words and sentences that used to categorize the tweet components, at this level, this those local ontology becomes a class that belongs to the global

**Figure 8: Ontology life cycle**

ontology

The process of the data creation or modification, within a local ontology, is based on a specific life cycle, which starts from draft mode to the published mode. (Figure 8).

In draft mode, system users can create or edit sentences or words samples, come back to them, save them and continue to work on them until they're ready to be submitted. Once the sample is submitted, it goes into the understood approved folder, where an ontology's local manager would review it.

During the review process, the sample can be either rejected, which would then put it back into the draft, or approved, in this case the sample is published.

At same point some sentences or words samples need to be updated and transferred into another local ontology, due to their meanings or their semantic change.

When the sample is selected to be revised, it goes back into the submitted stage, where the reviewer (manager) can either, once again, reject it or approve it to be revised. In case it's rejected, it goes all the way back to draft mode and starts the process all over again.

In conclusion, our work ends with computing the sentiment score of each peak and defining the percentage of each tweet classification (Positive, Negative, Neutral), which is referred to in this research as the highlight moment, a large volume of tweets tweeted by a group of people in a specific moment. In the first stage, the calculation of each tweet sentiment score, part of a highlight moment, is done by measuring each sentiment category percentage, by using the process that allows the calculation of the sentiment classification percentage via the use of ontologies.

In the Second stage, to make the decision about the analyzed sentiment category of the highlight moment that was generated by the peak of tweets volume, we calculated the average sentiment percentage after merging the sentiment classification of each tweet into three major sentiment categories, after that, we assign the sentiment with the maximum percentage to that highlight moment (Figure 9).
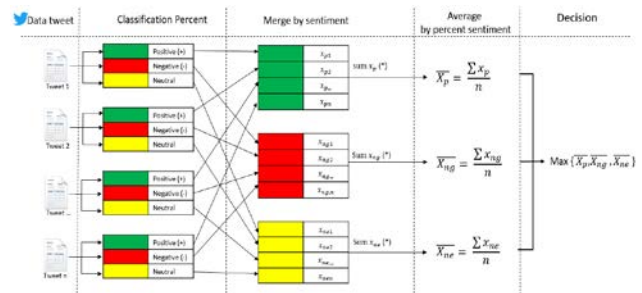


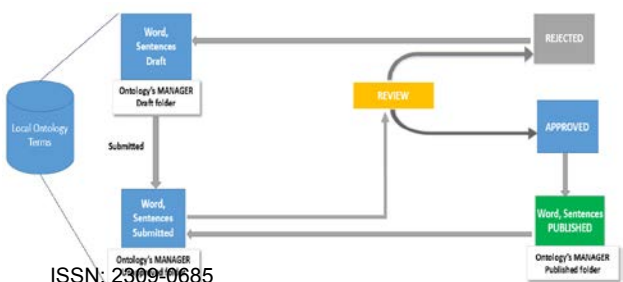**Figure 9: Highlight moment positivity decision by the calculating and merging process**

## IV. IMPLEMENTATION

To realize our objectives, which are to generate a highlight moments summary video of a live broadcasted event, and to calculate the sentiment percentage by category of each one of the highlight moments, we used the Twitter Streaming API that allows us to query Tweeter databases and get only the tweets data in regard of a specific Hashtag in real time and which
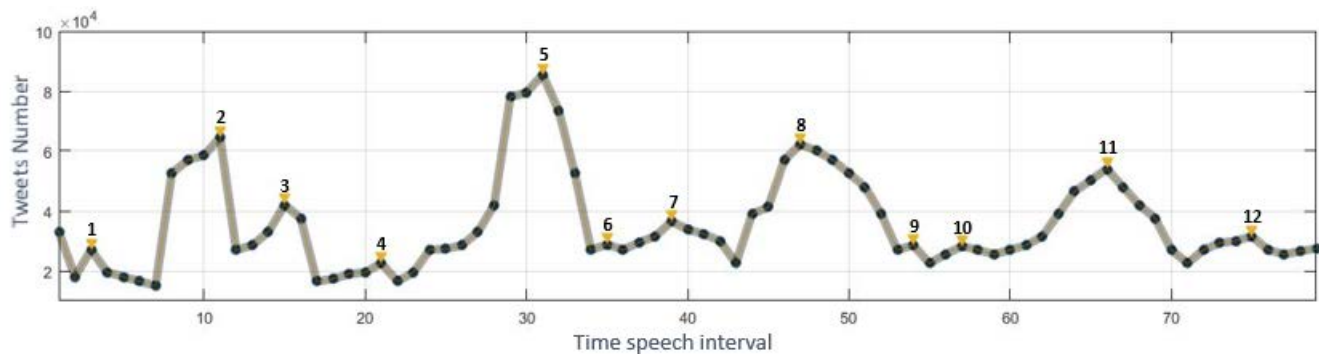
**Figure 10 : Highlights detection in political speech by detection the locals maximums of the obtain polynomial function that reflect the volume of tweets during a political event**

were generated in an exhaustive way, those hashtags have in general a relation with our speaker official account, such as #Donaldtrump, #Gop, #Maga, #Trump, #TinyTrump ...

Furthermore, thanks to LaGrange mathematical approach which has been presented in the previous chapter, we project the collected data from the twitter streaming API as a polynomial function in terms of time speech (**Figure 10**), in order to detect its local maximums (spike). The obtained peaks can be considered as the highlight moments detection key of our video summary.

Certainly, one of our work main objectives is the capacity to analyze the sentiment of every highlight moment's tweets and this by measuring computing the positivity rate forscore for each one of them by category (Positive, Negative and neutral). After using the KDD process that provided us with useful information of the big data recovered previously.

The measure of this sentiment positivity analysis can arise many challenges, due to numerous many obstacles e.g. of linguistic, cultural, regional expressions, etc.

To cope with these challenges, we came up with a reliable approach that uses ontologies, which turns out to be reliable and robust at resolving the semantic problems with regard of to the sentence or the group of word that composed the tweets.

To improve the interpretation of the sentiment analysis with regard to regarding the tweets data extracted at the semantic level, we constitute the processed tweet, as well as their rate of occurrence within every local ontology (Figure 11). Likewise, for a global sentiment classification of a single specific highlight moment. A simplified process was established, and this by merging all together each sentiment category percentage of each tweet that generates the highlight moment detected and by assigning the sentiment with the maximum percentage to the highlight moment global sentiment category (Figure12).
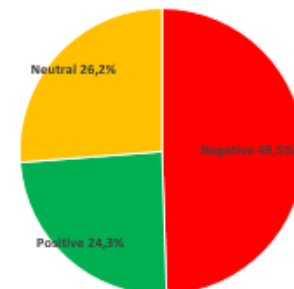


**Figure 12 : Sentiments classification percent of the fifth Highlight Moment**

## V.  CONCLUSION

In this article, a study was established on the generation of a video summary of an event, based on highlight moment detection using tweets volume changes, furthermore, we set up a process that allows measuring the sentiment expressed in the tweets of these highlight moments, then classifying those tweets by category (Positive, Negative and Neutral) to wind up by defining each highlight moment by the same categories after weighting the classification of the tweets that composed that moment.

From the results obtained, we can conclude that our proposed approach can play an important role on the detection of the citizen's sentiment in response to the speaker, which can open up a new perspective that will



**Figure 11 : Percentage measure of tweeter users' sentiment positivity**

facilitate the voters to better choose their presidential candidate during an event of a future election and not rely on media.

## REFERENCES

[1] J. E. M. de Oliveira, M. Cotacallapa, W. Seron, R. D. C. dos Santos, and M. G. Quiles, "Sentiment and Behavior Analysis of One Controversial American Individual on Twitter," Neural Inf. Process. -Springer Cham, pp. 509–518, Oct. 2016.

[2] C. Paris, H. Christensen, P. Batterham, and B. O'Dea, "Exploring Emotions in Social Media," IEEE Conf. Collab. Internet Comput., pp. 54–61, Oct. 2015.

[3] E. Lahuerta-Otero and R. Cordero-Gutiérrez, "Looking for the perfect tweet. The use of data mining techniques to find influencers on twitter," Comput. Hum. Behav., vol. 64, pp. 575– 583, Nov. 2016.

[4] F. Kalloubi, E. H. Nfaoui, and O. El Beqqali, "Graph based tweet entity linking using DBpedia," in Computer Systems and Applications (AICCSA), 2014 IEEE/ACS 11th International Conference on, 2014, pp. 501–506.

[5] M. Lai, C. Bosco, V. Patti, and D. Virone, "Debate on political reforms in Twitter: A hashtag-driven analysis of political polarization," in Data Science and Advanced Analytics (DSAA), 2015. 36678 2015. IEEE International Conference on, 2015, pp. 1–9.

[6] B. Vijayalaxmi, R. Putta, G. Shinde, and P. Lohani, "Object Detection Using Image Processing for an Industrial Robot," International Journal of Advanced Computational Engineering and Networking, 2013.

[7] J. Komala Lakshmi and M.Punithavalli, "A Survey on Performance Evaluation of Object Detection Techniques in Digital Image Processing," IJCSI Int. J. Comput. Scie Nce Issues, vol. 7, no. 6, 2010.

[8] Y. Rui, A. Gupta, and A. Acero, "Automatically Extracting Highlights for TV Baseball Programs," in Proceedings of the Eighth ACM International Conference on Multimedia, New York, NY, USA, 2000, pp. 105–115

[9] P. Suksai and P. Ratanaworabhan, "A new approach to extracting sport highlight," in Computer Science and Engineering Conference (ICSEC), 2016 International, 2016, pp. 1–6

[10] K.-S. Lin, A. Lee, Y.-H. Yang, C.-T. Lee, and H. H. Chen, "Automatic highlights extraction for drama video using music emotion and human face features," Neurocomputing, vol. 119, pp. 111–117, Nov. 2013

[11] V. Jain, "Prediction of movie success using sentiment analysis of tweets," Int. J. Soft Comput. Softw. Eng., vol. 3, no. 3, pp. 308–313, 2013.

[12] T. Mahmood, T. Iqbal, F. Amin, W. Lohanna, and A. Mustafa, "Mining Twitter big data to predict 2013 Pakistan election winner," in INMIC, 2013, pp. 49–54.

[13] K. Wegrzyn-Wolska and L. Bougueroua, "Tweets mining for french presidential election," in Computational Aspects of Social Networks (CASoN), 2012 Fourth International Conference on, 2012, pp. 138–143

[14] J. Bachhuber, C. Koppeel, J. Morina, K. Rejström, and D. Steinschulte, "US Election Prediction: A Linguistic Analysis of US Twitter Users," in Designing Networks for Innovation and Improvisation, Springer, Cham, 2016, pp. 55–63.

[15] E. Byrne and D. Corney, "Sweet FA: Sentiment, Swearing and Soccer," in 1st International Workshop on Social Multimedia and Storytelling co-located with ACM International Conference on Multimedia Retrieval (ICMR 2014), 2014.

[16] P. C. Guerra, W. Meira Jr., and C. Cardie, "Sentiment Analysis on Evolving Social Streams: How Self-report Imbalances Can Help," in Proceedings of the 7th ACM International Conference on Web Search and Data Mining, New York, NY, USA, 2014, pp. 443–452

[17] J. Pérez, E. Iturbide, V. Olivares, M. Hidalgo, A. Martínez, and N. Almanza, "A Data Preparation Methodology in Data Mining Applied to Mortality Population Databases," J. Med. Syst., vol. 39, no. 11, 2015.

[18] O. Marbán, J. Segovia, E. Menasalvas, and C. Fernández-Baizán, "Toward data mining engineering: A software engineering approach," Inf. Syst., vol. 34, no. 1, pp. 87– 107, Mar. 2009.

[19] S. K. Gupta, V. Bhatnagar, and S. K. Wasan, "Architecture for knowledge discovery and knowledge management," Knowl. Inf. Syst., vol. 7, no. 3, pp. 310–336, Mar. 2005.

[20] A. Vedder, "KDD: The challenge to individualism," Ethics Inf. Technol., vol. 1, no. 4, pp. 275–281, Dec. 1999.

[21] U. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, "The KDD Process for Extracting Useful Knowledge from Volumes of Data," Commun ACM, vol. 39, no. 11, pp. 27–34, Nov. 1996.

[22] J.-P. Berrut and L. N. Trefethen, "Barycentric Lagrange Interpolation," SIAM Rev., vol. 46, no. 3, pp. 501–517, Jan. 2004.

[23] M. Abid, B. Nsiri, and Y. Serhane, "Interoperability between different port information systems," Int. J. Math. Comput. Simul., vol. 8, pp. 156–161, 2014.

[24] S. Jai-Andaloussi, I. E. Mourabit, N. Madrane, S. B. Chaouni, and A. Sekkaki, "Soccer Events Summarization by Using Sentiment Analysis," in 2015 International Conference on Computational Science and Computational Intelligence (CSCI), 2015, pp. 398–403.

**Mehdi ABID** was born in Casablanca, Morocco in 19- Decembre 1990. He received the M.S Degree in in science: computer and internet engineering in 2013 from faculty of science university Hassan 2 Casablanca and he is currently completing a PhD in "Laboratoire d'informatique et d'aide à la décision" in Information System Integrated Logistics, he has authored and co-authored several publications in journals and international conferences (Scopus).

**Benayad NSIRI** received in 2000 his D.E.A (French equivalent of M.Sc. degree) in electronics from the Occidental Bretagne University in Brest, France; his Ph.D. degree from Telecom Bretagne in 2004. While in 2005, he received a MBI degree in computer sciences from Telecom Bretagne, and in 2010 he received HDR degree from Hassan II University, Casablanca, Morocco. Currently, he is a Professor in Ain chock faculty of sciences, Hassan II University at Casablanca, Morocco; a member in LIAD laboratory, Hassan II University and a member associate in Lab-STICC laboratory at Telecom Bretagne, Brest, France. Professor Benayad NSIRI has advised and co-advised more than 7 PhD theses, contributed to more than 60 articles in regional and international conferences and journals. His research interests include but not restricted to computer science, communication, signal and image processing, adaptive techniques, blind deconvolution, MCMC methods, seismic data and higher order statistics.

**Yassine SERHANE** was born in Essaouira, Morocco in 15th February 1990. He received his M.S Degree in science: Computer and Network Engineering in 2013 from faculty of science University of Hassan 2 Casablanca and he is currently completing a PhD in "Laboratoire d'informatique et d'aide à la décision" in Machine Learning, he has co-authored several publications in journals and international conferences (Scopus).