# USER PROFILING FOR BIG SOCIAL MEDIA DATA

S.CHOUYA[1], S.BOUREKKADI[2], S.K HOULJI[1]A.BABOUNIA[2] M.L.KERKEB[1]

1: Laboratory of  INFORMATION SYSTEMS ENGINNEERING RESEARCH GROUP
Abdelmalek Essaadi University, Faculty of Science, Tetouan, Morocco

2: LRMO national school of commerce and management, IbnTofailUniversity , Morocco

## Abstract

Online Social Networks (OSNs) have of late been the subject of various investigations that have tried to create helpful strategies for the classification and analysis of big content. A few of the key strategies of this research to current logical understanding has to do with identifying the basic subjects within content (posts and messages).

**Index Terms**: Data, Profiling, Media, University, Online Social Networks.

—————————— ◆ ——————————

## Introduction

Social media platforms started in the mid-1990s, even though the degree of their potential effect on mainstream media was not evident during that period. It was incomprehensible then that 20 years later, top networks would have a huge number of dynamic individuals or that their qualities on Wall Street would stretch out to ten-digit figures. For all intents and purposes each feature of social life and human conduct has been changed by this innovative progress, and it is presently hard to recollect a period when people couldn't express their sentiments on a scope of occasions by means of social media. Through the progress of social media, a colossal trove of information has been created that can be utilized for data mining and artificial learning. Therefore, social media has developed as the best subjects in the field of data innovation amid the most recent decade, which has been set apart by various original disclosures. Our point in this research was to build up another answer for deconstructing patterns of communication inside any online social media platform.

## User profiling

User profiling refers to all activities related with obtaining, examining, and applying information identified to the behavior of the user inside a system. The subsequent profile can encourage more exact delivery of content. The profiling procedure is made up of technical phrases involving topics that specific individuals are probably going to search. A more precise characterization deals with data on activities with common patterns, for instance, visited pages or the recurrence of logins. A few researchers trust that translating associations amongst individuals and their common communications can reveal insight into every client's veritable advantages. Such bits of knowledge have evidently improved the prescient intensity of various hypothetical models and practical systems. Even though past investigations have broadly investigated this strategy for deducing clients' interests, however we chose to expel any conceivable questions by confirming it through a test performed on extracted samples from an active network. Our selected technique was to analyze a client's attributes and exercises before utilizing this learning inside an intensive classification framework. A process manually operated was carried out to decide the main scope of attributes, and knowledge gained during this stage were implemented to mechanize the procedure during the consequent stage.

**Mining social media content**

At the point when the users of OSN share data, their common connections can be utilized for choosing content and showing calculations. It is conceivable to choose content for a specific individual depending on the impact of people who are being trailed by this individual. Selection based on relationship has its benefits, however experiences the issue of a limit in the information accessed. But, this issue isn't experienced when selection strategies concentrating on content instead of associations are utilized. A few scientists have additionally endeavored to utilize a joined technique, for instance, by building a hierarchy of classifications based on an investigation of collected information about individuals' exercises inside the system. About 80,000 new websites are made every day, with an excess of a million distributed content pieces showing up inside the same time period. It is obvious from a thought of the content made and shared through social media that the online condition contains a plenitude of data about individuals, including their preferences and inclinations on various subjects. The behavior delivered in these channels can be comprehended as indicators of areas of support or resistance connecting with an expansive range of nearby or global issues. These behaviors are frequently enunciated in extremely solid terms and can be utilized to drive different marketing activities connected to the questioned issues. The utilization of exemplary strategies can be of an incentive in carrying out task on attitude detection, however they are not suited for working with huge measures of information that are regularly assembled from social platforms. Therefore, the advancement of creative information extraction techniques is a significant task that our research endeavors to accomplish.

**Strategies and tools**

The system exhibited is special and involves a few steps that incorporate accumulating the real content, recognizing topics organized by every user, and actualizing quality control to decide the value of each tweet. Scientific calculations were done to get the estimations of the peer effect factor and also user strength coefficients. Data mining for this research depended on a few foreordained qualities. 100K clients were picked, and 2000 were selected badly from the positions of Twitter users who speak English and have an excess of 5000 followers. The first point was to gather an excess of 3200 tweets from every one of the users chosen. But, the real number of aggregated content pieces was less due to the Twitter rate limit. Special modules were used to import tweets were that we created for our framework and were saved in a prepared database to be processed. About 1600 and 3200 tweets were accumulated from every member, and the sample was concluded, with no further changes being made.

## Discussion and analysis

To affirm that the answer which uses the SOM had the capacity to deliver solid outcomes, each phase needed to be scrutinized. Due to the discoveries of this research, we noticed a few areas that required consideration for the powerful improvement of the framework. Determining the quality of user impact beneath a sensible doubt was impossible utilizing the present arrangement, and we could just expect the exactness of this variable. The strategy involved an examination of double the quantity of expected themes to represent conceivable contaminations inside the example and to stay away from words that were not really identified with the subject. Utilizing an abbreviated list, such words were allotted probabilities of relationship with the subject yet were overlooked from the last computations. The decision of the reference content for assessing dialect is additionally critically imperative. For this situation, a gathering of 10,000 words was utilized that avoided an extensive number of casual articulations.

### Restrictions and future work

### Bringing in metadata about clients

Twitter expects clients to fill forms in their profiles according to their characteristics. Thus, the profile pages of individuals' can offer significant bits of knowledge. A portion of the key fields incorporate sex, area, and age, yet other demographic details could be interesting. Given that contending OSNs additionally gather similar kinds of metadata, the arrangement could be acclimated to exploit this abundance of information sources, thus expanding the exactness of estimations of social pressure.

### Building up a more exact quality scale

A basic issue that affected the general model was the value given to each content piece. Since Twitter confines communication to 140 characters for every message, there is a restricted measure of content accessible for leading an investigation and hunting down applicable data. Past examinations have demonstrated that the greater part of all tweets convey no semantic incentive for clients past the quick hover of the author's nearest partners for whom the message was initially proposed. Automatic determination of which tweets contain noteworthy data of a more broad nature is a troublesome assignment that requires the use of different strategies. Following the aggregate message length and the normal number of characters per word and also checking for the most much of the time utilized terms as well as emoticons are a portion of the techniques that could be useful in such manner. In messages that allude to topics identified with financial aspects or science, the nearness of numerical components and scientific images could likewise demonstrate that the tweet has more than normal esteem. As this SOM-based solution rotates based on the exact estimation of each tweet's actual value, the improvement of further developed estimating frameworks for this reason would be gainful. Uniform tweet values over various topical classes represent another impediment of the present solution. Consequently, a huge update of the model empowering the assurance of discrete qualities for each topic would be especially worthwhile.

## Conclusion

The research was considered as a complicated effort that stretched beyond basic data gathering to build up an inventive arrangement that solved the problem of peer effect identifying with the opinions of Twitter users. The proposed arrangement involved following the easygoing conduct of network members and changing the information in this way acquired into efficient quantitative portrayals of specific users' online propensities. This empowered the estimation of the most conspicuous territories of individual request for each client and additionally recognition of the nearness and quality of their social effects. To achieve the goals of the research, we counseled and connected the discoveries of various examinations wherever we valued this valuable and constructive. Future work to promote this arrangement ought to be coordinated at improving the usefulness of its individual modules. Since the output of the system are dependent upon the results of every module, it is prescribed that every module ought to be streamlined with the general goal of enhancing the prescient intensity of the arrangement. The Twitter platform can be useful in this activity, as an extensive number of parameters about individuals and their exercises are followed and could be used to build up a superior comprehension of how different users are connected to each other and how they influence the reasoning of their prompt and in addition more far off contacts. These parameters would give an extra layer of information above verbal substance, in this manner constituting a multidimensional approach for tending to this issue.