

Singapore Management University

Institutional Knowledge at Singapore Management University

Research Collection School Of Computing and Information Systems

School of Computing and Information Systems

11-2017

Leveraging social analytics data for identifying customer segments for online news media

JANSEN

Bernard J

Soon-Gyo JUNG

Jisun AN

Singapore Management University, jisunan@smu.edu.sg

Haewoon KWAK

Singapore Management University, hkwak@smu.edu.sg

See next page for additional authors

Follow this and additional works at: https://ink.library.smu.edu.sg/sis_research



Part of the [Artificial Intelligence and Robotics Commons](#), and the [Databases and Information Systems Commons](#)

Citation

JANSEN; Bernard J; JUNG, Soon-Gyo; AN, Jisun; KWAK, Haewoon; and KWAK, HAEWOON. Leveraging social analytics data for identifying customer segments for online news media. (2017). *Proceedings of 2017 IEEE/ACS 14th International Conference on Computer Systems and Applications (AICCSA)*. Available at: https://ink.library.smu.edu.sg/sis_research/6654

This Conference Proceeding Article is brought to you for free and open access by the School of Computing and Information Systems at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection School Of Computing and Information Systems by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email cherylids@smu.edu.sg.

Author

JANSEN, Bernard J, Soon-Gyo JUNG, Jisun AN, Haewoon KWAK, and HAEWOON KWAK

Leveraging Social Analytics Data for Identifying Customer Segments for Online News Media

Bernard J. Jansen
Qatar Computing Research Institute
Doha, Qatar
jjansen@acm.org

Soon-gyo Jung
Qatar Computing Research Institute
Doha, Qatar
sjung@hbku.edu.qa

Joni Salminen
Qatar Computing Research Institute
Doha, Qatar
jsalminen@hbku.edu.qa

Jisun An
Qatar Computing Research
Institute
Doha, Qatar
jan@hbku.edu.qa

Haewoon Kwak
Qatar Computing Research
Institute)
Doha, Qatar
hkwak@hbku.edu.qa

Abstract— In this work, we describe a methodology for leveraging large amounts of customer interaction data with online content from major social media platforms in order to isolate meaningful customer segments. The methodology is robust in that it can rapidly identify diverse customer segments using solely online behaviors and then associate these behavioral customer segments with the related distinct demographic segments, presenting a holistic picture of the customer base of an organization. We validate our methodology via the implementation of a working system that rapidly and in near real-time processes tens of millions of online customer interactions with content posted on major social media platforms in order to identify both the distinct behavioral segments and corresponding impactful demographic segments. We illustrate the functionality of the methodology with real data from a major online content provider with millions of online interactions from more than thirty countries. We further show one possible use for such information via the automatic generation of personas for an organization, which can be used for the formulation of marketing strategy, implementation of advertising plans, or development of products. The research results offer insights into competitive marketing and product preferences for the consumers of online digital content. We conclude with a discussion of areas for future work.

Keywords— *Social analytics; social media; web analytics; online news; customer segmentation; digital marketing.*

I. INTRODUCTION

One of the uses of the web analytics data [1], and increasingly social media data also, is customer segmentation, which is an approach for separating an overall customer base into clusters or groups of individuals who are similar, as defined by some specific set of attributes. Customer segmentation is a widespread practice across industries, with the set of attributes utilized being relevant to the domain. Examples include marketing, advertising, education, and system design. E-commerce companies and other organizations employ customer segmentation to target specific customer groups with content and products that the consumers would most

likely find relevant. Additionally, customer segmentation might also lead to a deeper understanding of customers' preferences and wants by isolating what each segment finds most valuable. Based on these insights, organizations can more effectively engage with their customers.

The customer segmentation approach relies on identifying key attributes from which one can separate customers into groups [2]. Targeting customers via behavioral segmentation involves dividing the customer base based on their collective behavior. A behavior can be a single attribute (e.g., viewing online content) or a set of behaviors (e.g., viewing online content, length of video watched, etc.), but it is typically focused on the way the segment responds to, uses, or engages with a product. Targeting customers via demographic segmentation involves grouping the customers based on one or more characteristics, which often include, but are not limited to, gender, age, race, location, education, income, or career.

However, most prior work in customer segmentation has focused on using individual website data, such as that available from Google Analytics. Yet, there is an increasing effort to employ customer segmentation using social media data from the major online platforms. This data presents unique challenges, as it is typically aggregated to preserve privacy, so methods must be employed to deal with the issues this aggregation causes in inferring customer attributes.

In the ongoing research reported here, we investigate the use of aggregated social media data for isolating customer segments based on both behaviors and on demographics of those customers and then linking the two customer segments groupings for a complete representation of the customer base.

II. RELATED WORK

Customer segmentation is an ongoing research area [3], as the availability of online data has greatly increased. Website data has been used to segment customers into various revenue groupings [4], as an example of behavioral segmentation.

Search data has been used to classify the gender of searchers and then relate this demographic attribute to revenue generation [5]. Increasingly, customer segmentation processes are leveraging social media data, for both behavioral and demographic grouping [6] of customers.

However, there are many unanswered questions concerning the use of social media analytics for customer segmentation. In particular: *Can one isolate customer segments based on behavioral interaction on social media platforms? Can online data deliver rich demographic insights for customer segmentation? Can the customer segments be identified in real time? Can the customer segments be frequently updated?* These are some of the questions that motivate our research.

III. RESEARCH OBJECTIVE

In this research, our premise is that real-life user data, both behavioral and demographic, concerning consumers of a product, service, or content [6], can be rapidly collected from major social media platforms and analyzed in order to isolate customer segments that are useable for a variety of commercial purposes.

Therefore, our research objective is to develop and implement a methodology to:

- *Identify discrete customer segments based on behavioral interactions with online content posted on major social media platforms.*
- *Detect the discrete demographic customer segments associated with each of these behavioral customer segments.*

These are the research objectives that we investigate, presenting the current state of research.

IV. DATA COLLECTION AND RESEARCH DESIGN

With actual user data from AJ+ Arabic, an online media and mobile outlet, we validate our premise concerning using social media analytics data for customer segmentation. Audience leanings have been somewhat ignored by news journalists in the past, primarily due to the lack of precise readership measurements preceding the era of online news. Prior studies highlight the large differences between news production patterns and consumption patterns by using the number of articles concerning particular topics and the views of these articles posted on news websites [7]. In the competitive news and content creation field, an accurate understanding of audiences is critically important for attracting customers and increasing content consumption.

Therefore, we consider AJ+ Arabic as an excellent organization for data collection and for achieving our research objectives. Using real AJ+ Arabic user data shows the value of generating customer segments, specifically for the news sector. We believe the methodology transferable to other industry verticals that use major social media platforms for content dissemination, including e.g. e-commerce, public health, and online advertisers of various kind.

A. Data Collection Organization

AJ+ Arabic is an online news channel from Al Jazeera Media Network focusing on 34 primarily Arabic speaking countries and Arab nationals in the diaspora. AJ+ Arabic is a natively digital channel, with a presence only on social media platforms. As a media concept, AJ+ Arabic was designed from the start to provide digital content news in the medium of the viewer (i.e., YouTube, Facebook, etc.), with no links redirecting to a website. Therefore, the digital content is specifically designed to be viewed in the Facebook newsfeed, YouTube Channel, Twitter Timeline, etc. for each of the customer segments that are most active on those platforms. Digital content is created for both cross-platform consumptions (i.e., more than one social media channel) and for specific platforms (i.e., only one social media channel).

As with many organizations that rely on the creation of digital content, AJ+ Arabic has an acute need for automatic processing of large amounts of customer data for gaining customer insights in order to direct digital content generation, media planning, and strategic choices. This need is driven by the necessity for the agile and successful development of digital content in the competitive and rapidly changing online information market.

Therefore, in pursuit of our overall research objective of leveraging social media data for identifying customer segments based on both behaviors and demographics, for the research reported in this manuscript, we are specifically interested in understanding the AJ+ Arabic customers by identifying (1) whom AJ+ Arabic is reaching (i.e., market segments) and finding out (2) if there are significant differences in behaviors or demographics that are reflected in the social media analytics data.

Using these potential differences, combined with other user data [e.g., 8, 9], we validate the approach via the development of a system that uses actual customer data to aid in a variety of organizational tasks and use cases relating to content creation for social media audiences.

B. Data Collection from AJ+ Arabic YouTube Channel

We collected the data from the AJ+ Arabic YouTube channel; however, based on our experimentation with the API and data structure, the methodology is applicable to Facebook and perhaps to other social media platforms as well. The primary reason for using YouTube is that its analytics API provides detailed statistics for every video. Additionally, the data structure is comparable to other social media platforms in terms of bucketed demographic information (e.g., ‘Women, 18-24, West-Virginia’) and performance metrics (e.g., view counts, clicks). An example of an AJ+ Arabic YouTube video is shown in Fig. 1; note the likes/dislikes and views.

The YouTube analytics platform is robust and provides, for each of the AJ+ Arabic videos, aggregated data on user behavior (e.g., views, etc.) and profile (e.g., gender, age, country location, and which site the user comes from). We use these data variables to validate if social media data can provide the information needed for customer segmentation along both behavioral and demographic lines.

Fig. 1. An example of YouTube Videos from the AJ+ Arabic YouTube Channel, with Content Meta-Data, Including Views, Likes, and Dislikes.



We access the data in the AJ+ Arabic YouTube analytics platform by the YouTube APIs. Table I shows the parameters we use for calling the APIs in this current stage of research. Although there are other metrics besides viewCount (the number of views) and viewerPercentage, we deem these variables to be of the highest value. We reserve other variables, such as the number of likes, comments, etc., for future research. Also, of note, the data from a YouTube channel is private and available only to its administrator, and it is thus not accessible via public interfaces.

TABLE I. YOUTUBE API QUERY PARAMETERS AND METRICS USED FOR AJ+ ARABIC ANALYSIS.

| Demographic Dimension | Behavioral Dimension |
|-----------------------|----------------------|
| ageGroup | viewCount |
| country | viewerPercentage |
| day | |
| gender | |
| month | |

V. RESULTS

A. Exploratory Analysis of AJ+ Arabic YouTube Data

We begin by presenting some of the overall statistics from the AJ+ Arabic YouTube channel data. Due to business concerns, we do not provide the exact absolute numbers, instead providing percentages only.

During the study, there were nearly 2,000 AJ+ Arabic videos posted to the AJ+ Arabic YouTube Channel during a close to thirty-month period, February 2015 to June 2017. These videos had some twelve million views by users from 214 countries, representing 85.6 percent of the countries listed in the ISO country code listing. So, regardless of the regional focus of the AJ+ Arabic, the channel has a worldwide audience.

Overall, the AJ+ Arabic audience is concentrated in the Arabic speaking countries, although the complete reach is wider. The topmost three countries, in terms of viewership,

include Saudi Arabia, Iraq, and Morocco, with each representing between ten and twenty percent of total viewership in terms of the number of unique videos watched. Concerning the total number of views, the largest country market segment, with about twenty percent of video views, is from Saudi Arabia. It is also interesting to note that, although AJ+ Arabic was designed to target the Arabic speaking region, most viewers come from outside the traditional Arabic speaking region, making it a challenge to create a comprehensive representation of the viewers.

There are some videos that show a worldwide appeal, with about 100 videos being viewed in 100 or more countries. Equally, there were about 100 videos that were viewed by users from 5 or fewer countries. There were twenty-five countries from which viewers watch more than 75% of the AJ+ Arabic videos. Concerning the actual number of views, the view counts per individual videos follow a power law distribution, with there being a small number of videos viewed a lot and a substantial number of videos being viewed a small number of times. This finding is not surprising, as such skewed popularity of videos is one of the well-known characteristics of viewing behavior on YouTube [10].

TABLE II. PERCENTAGE OF VIDEO VIEWINGS BY AGE GROUPINGS.

| Age Grouping | Percentage |
|--------------|------------|
| 13-17 | 6.4% |
| 18-24 | 34.1% |
| 25-34 | 42.6% |
| 35-44 | 11.8% |
| 45-54 | 2.8% |
| 55-64 | 1.0% |
| 65+ | 1.2% |
| | 100.00% |

A surprising finding from the exploratory data analysis was the lack of information concerned the gender and age distribution. YouTube gathers gender information (e.g., female, male). Also, YouTube viewers are classified into one of multiple age categories (e.g., 13-17 years, 18-24 years, 25-34 years, 35-44 years, 45-54 years, 55-64 years, and 65 years and older). As AJ+ Arabic is designed to target the young generation by adopting social platforms, one would expect the audience to skew to young adult male. However, it turns that the most audience members of AJ+ Arabic login anonymously (i.e., not using a Google account); therefore, there is limited gender or age demographic information. Of the data with demographic information, the gender split was 38% females and 62% males. Table II presents the age distribution.

B. Research Objective One Analysis – Identification of Customer Behavior Segments

Moving to our specific research objectives, we first use the social media data to identify discrete customer segments based on different behaviors. The process of this is quite challenging, as the social media data from most social media

platforms is aggregated (i.e., no individual level data) to preserve the privacy of the users. Therefore, in order to isolate customer behavior patterns, the data must be disaggregated.

For segmentation, we first experimented with k-means clustering [11], but due to the aggregation of the data, it was found ineffective. Thus, we turned to matrix decomposition techniques, specifically non-negative matrix factorization (NMF), as outlined in [12], [13]. We conceptually present this matrix decomposition approach here.

To begin the decomposition, we first develop a matrix representing users' interaction with the online content products. The columns of the matrix are the online products, in this case, the AJ+ Arabic videos (e.g., c contents (C_1, C_2, \dots, C_c)). The rows of the matrix are the user groups or customer demographic segments (e.g., g user groups (G_1, G_2, \dots, G_g)). Therefore, the matrix describing the association between user groups and contents is denoted by V the $g \times c$ matrix of g user groups or customer demographic segments and c contents. The element of the matrix V , V_{ij} , is any statistic that represents the one interaction or set of interactions of the user group G_i for content C_j . In the research presented here, the user interaction element is *viewCount*.

Using this matrix approach as the basis, we can decompose (i.e., separate into simpler components) the overall matrix V into two matrices: W and H . The matrix W encodes an association between user groups and behavioral customer segments (i.e., latent content interaction patterns), and the matrix H encodes an association between behavioral customer segments and pieces of content. The resolution in finding customer segments can be adjusted by the number of columns in W or that of rows in H . To sum up, once we have the matrix H , we discover the underlying latent patterns, which describe the user interaction with content, that will become the basis of the customer demographic segments in the next step.

Relative to a basic clustering of user groups approach, decomposition has advantages in that it can find multiple behavioral patterns even from a single user group. There are several ways to decompose a given matrix. For this research, we used NMF to its high interpretability [13].

C. Research Objective Two Analysis – Identification of Related Customer Demographic Segments

Moving to our second research objective, we identify the most impactful customer demographic segments associated with the previously defined behavioral customer segments. After decomposing the matrix V , we have the matrix H (containing the customer behaviors) and another matrix, W (containing the demographic groups).

Each row in W represents how each customer demographic segment can be characterized by different consumption patterns. The columns in W show how a common consumption pattern is associated with different customer segments. A single behavioral segment can, possibly, have multiple associated customer demographic segments. Thus, for each column, the user group with the largest coefficient or weight can be interpreted as the most impactful user demographic group for that corresponding pattern. Although one can present as many behavioral

segments as the data contains, cognitive limits of the users of the system pose a limit; it is not purposeful to show them hundreds of customer segments. We present five customer behavioral segments and the associated top customer demographic segments with their weights in Table III.

TABLE III. FIVE CUSTOMER BEHAVIORAL SEGMENTS PRESENTED WITH THE TOP ASSOCIATED CUSTOMER DEMOGRAPHIC SEGMENTS. THE DEMOGRAPHIC SEGMENT WITH THE MOST WEIGHT IS BOLDDED.

| Customer Behavioral Segment | Possible Customer Demographic Segments | | | |
|-----------------------------|--|------------|---------------|---------------|
| No. | Country | Age | Gender | Weight |
| 1 | Saudi Arabia (SA) | 25 | male | 400.98 |
| | SA | 18 | male | 190.08 |
| | SA | 35 | male | 74.43 |
| | Iraq (IQ) | 25 | male | 74.07 |
| | IQ | 18 | male | 63.87 |
| 2 | Country | Age | Gender | Weight |
| | Jordan (JO) | 18 | male | 155.86 |
| | Palestine (PS) | 18 | male | 123.71 |
| | JO | 25 | male | 112.11 |
| | JO | 18 | female | 106.5 |
| | Germany (DE) | 25 | male | 91.43 |
| 3 | Country | Age | Gender | Weight |
| | SA | 25 | female | 66.07 |
| | SA | 18 | female | 38.42 |
| | SA | 25 | male | 21.63 |
| | IQ | 25 | male | 17.31 |
| | IQ | 25 | female | 17.15 |
| 4 | Country | Age | Gender | Weight |
| | Morocco (MA) | 25 | male | 57.95 |
| | MA | 18 | male | 45.14 |
| | MA | 35 | male | 25.93 |
| | MA | 18 | female | 17.01 |
| | MA | 25 | female | 14.44 |
| 5 | Country | Age | Gender | Weight |
| | IQ | 18 | male | 231.99 |
| | IQ | 25 | male | 231.63 |
| | IQ | 35 | male | 98.43 |
| | IQ | 18 | female | 91.63 |
| | IQ | 13 | male | 68.8 |

In Table III, we present five customer demographic segments identified via our matrix decomposition approach to discover

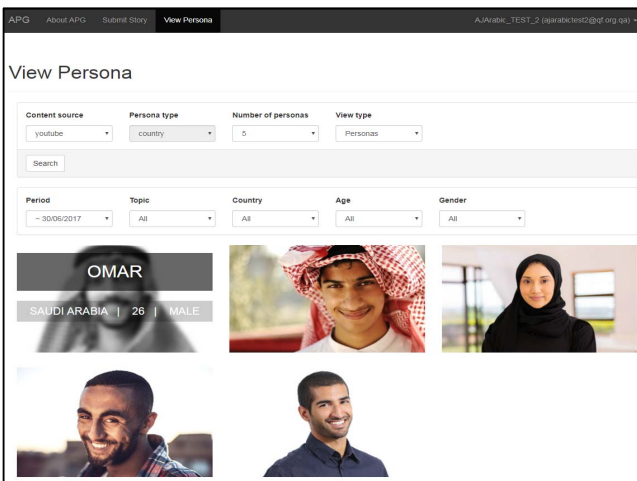
distinct customer behavioral segments based on interaction with the online content (see column 1). Then, for each behavioral pattern, Table III displays the top customer demographic segments associated with each of these five behavioral segments (see columns 2 through 5). Our decomposition approach calculates a weight for each of these demographic segments, assigning a higher weight to the most impactful (i.e., largest) demographic customer segments.

In Table III, there are at least two findings that are apparent upon analysis. First, there is a predominance of male segments. Given that the data is from YouTube, this is a reasonable and expected finding. Second, there is a clustering effect by gender, and then age, and sometimes location. For example, in the first behavioral segment, we see that the top three customer demographic segments are all young males from Saudi Arabia. This opens up an interesting research question of how granular one needs to get for each of these demographic segments, as there is apparently little behavioral difference when segregating this portion of the population into three different age brackets.

D. Use Case for Behavioral and Demographic Customer Segmentation

We believe that there are many possible use cases for the employment of behavioral and demographic segments, both separately and in conjunction by linking the behavioral segments to demographic segments. Here, we present one possible use case, using the research results for the automatic generation of personas [14] using social analytics data to first isolate customer segments, both behavioral and demographic.

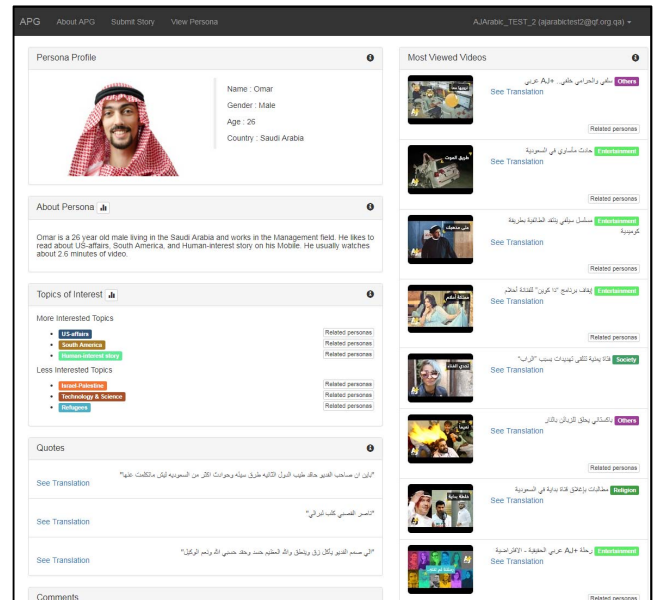
Fig. 2. Screenshot of the Automated Persona Generation System Generating Five Personas based on YouTube Data. Note the Images for Each Persona and the Demographic Information that Appears on the Cursor Rollover of One of the Images.



Using social media analytics data from a major social media platform, we have developed a system that automatically accesses aggregated data and decomposes it using the method outlined above. As shown in Fig. 2, with an example of five customer segments as personas (fictive people based on real data), the system presents demographically appropriate images, names, country, age, and gender for each persona. This demographic data is first derived from the social media

data, i.e. <gender, country, age>. Using this information, the system then accesses backend databases selecting gender, age, and country appropriate images and names. The demographic information is displayed when the cursor hovers above a persona image. When one of the persona images is clicked, the corresponding persona description is displayed, as shown in Fig. 3.

Fig. 3. Screenshot of a Persona Description that is Automatically Generated from Social Media Analytics and Contains Both Behavioral and Demographic Customer Segmentation Data.



VI. DISCUSSION AND FUTURE DIRECTION

The results of our research demonstrate that social media data from the major online platforms is quite robust in identifying customer segments for digital content creators. The implementation of the approach in a working system shows that one can use actual user behavior that is rapidly collected from social media platforms and analyze this data in order to generate customer segments for gaining insights for marketing, product development, and other critical areas of business. Yet, there are numerous research and development fronts that we are pursuing in the future to enhance the impact of this research.

A. Enhancing the Customer Attributes with Additional Behavioral and Demographic Data

The major strength of our approach is that it benefits from actual user data, reducing time and cost for generating both behavioral and demographic customer segments and providing a mechanism for linking these two types of customer groupings into a coherent cluster. We show the value of our approach by leveraging these market segments for the automatic generation of personas. Using our method, we do not need to vigilantly survey users, who may or may not be actual consumers of our content. Instead, we can leverage online social media data from millions of actual consumers to extract representative customer segments based on genuine customer data.

We consider this research is a starting point for leveraging behavioral and demographic customer segments from social media analytics data for a vast number of other applications and services with minimal manual efforts. For example, if we can leverage additional rich information concerning a consumer, such as an ethnicity, socio-economic status, and precise location, our approach and results would become even more useful. For future research, we are exploring these avenues.

As one example, it might be possible to extract demographic information using shared links on Facebook [11], via Twitter Google+ profiles. Links shared on Facebook could reveal information such as economic status, as the links reveal particular interests. There has been prior research showing that affluent customers visit more high-end luxury product websites, while budget-conscious customers visit price aggregation or discount websites. Thus, the economic status of the consumer can be distinguished by the websites they visited [15]. We are currently investigating the reliability of extracting rich demographic information for such purposes.

B. Possible Scalability Issues of Using Social Media for Automatic Customer Segmentation

In terms of possible concerns, there is the issue of scalability of a system relying on social media analytics data. However, this is a minor concern, as the API allows only a limited number of HTTP requests for each consumer, and we only need to access the ‘newer data’ once the initial dataset is collected. As such, the number of requests does not radically increase over time. The only bottleneck in this area is the daily limit of API calls, which varies by social media platform. However, we have not found this to be a major impediment at this point.

C. Benefits for News Organizations of Using Social Media for Customer Segmentation

News and other organizations that produce digital content can benefit from using social media for customer segmentation that provides both rich behavioral patterns and demographic information of those consumers who interact with the online content. In fact, the development approach presented here is derived from collaboration with journalists. Digital content creators, such as journalists, desire realistic descriptions of their actual customers, which can lead to better content selection and improved content creation. In this sense, our customer segmentation research provides a robust groundwork for future research, as the resulting customer segments, and the personas generated from them, indicate the content topics that actual audience members are interested in.

VII. CONCLUSION

In this research, we show that social media data can be effectively used to create customer market segments in near real time. As shown, we have begun the initial investigative steps to change customer segment creation from a labor-intensive analytics process to an automated procedure that takes place in near real time. Moreover, the research illustrates how individual demographic groups can be automatically

analyzed for detection of several behavioral patterns. This has the potential of leading into more granular customer insights while avoiding over-generalization of customer behavior based on demographic attributes. For future research, we are actively investigating the use of other social media data sources [16], to be used in conjunction with offline sources to provide richer demographic attributes and attitudinal characteristics. These enrichment efforts link with many subfields of computer science and marketing.

REFERENCES

- [1] B. J. Jansen, *Understanding user-web interactions via web analytics*, 1 ed. vol. 1: Morgan Claypool, 2009.
- [2] B. Coolil, L. Aksoy, and T. L. Keiningham, "Approaches to customer segmentation," *Journal of Relationship Marketing*, vol. 6, pp. 9–39, 2008.
- [3] T. V. Bonoma and B. P. Shapiro, "Evaluating market segmentation approaches," *Industrial Marketing Management*, vol. 13, pp. 257–268, 1984.
- [4] A. Ortiz-Cordova and B. J. Jansen, "Classifying Web Search Queries in Order to Identify High Revenue Generating Customers.," *Journal of the American Society for Information Science and Technology*, vol. 63, pp. 1426 – 1441, 2012.
- [5] B. J. Jansen, K. Moore, and S. Carman, "Evaluating The Performance of Demographic Targeting Using Gender in Keyword Advertising," *Information Processing & Management*, vol. 49, pp. 286-302, 2013.
- [6] B. J. Jansen, K. Sobel, and G. Cook, "Classifying Ecommerce Information Sharing Behaviour by Youths on Social Networking Sites," *Journal of Information Science*, vol. 37, pp. 120-136, 2011.
- [7] S. Abbar, J. An, H. Kwak, Y. Messaoui, and J. Borge-Holthoef, "Consumers and Suppliers: Attention asymmetries. A Case Study of Aljazeera’s News Coverage and Comments," presented at the Computation+Journalism Symposium 2015, New York, NY, 2015.
- [8] Y.-H. Hu, Y.-L. Chen, and H.-L. Chou, "Opinion mining from online hotel reviews – A text summarization approach," *Information Processing & Management*, vol. 53, pp. 436-449, 2017.
- [9] G. Vinodhini and R. M. Chandrasekaran, "A sampling based sentiment mining approach for e-commerce applications," *Information Processing & Management*, vol. 53, pp. 223-236, 2017.
- [10] M. Cha, H. Kwak, P. Rodriguez, Y.-Y. Ahn, and S. Moon, "'I tube, you tube, everybody tubes: analyzing the world’s largest user generated content video system," in *Proceedings of the 7th ACM SIGCOMM conference on Internet Measurement*, 2007, pp. 1–14.
- [11] J. An, H. Cho, H. Kwak, M. Z. Hassen, and B. J. Jansen, "Towards Automatic Persona Generation Using Social Media," in *2016 IEEE 4th International Conference on Future Internet of Things and Cloud Workshops (FiCloudW)*, 2016, pp. 206-211.
- [12] S. Jung, J. An, H. Kwak, M. Ahmad, L. Nielsen, and B. J. Jansen, "Persona Generation from Aggregated Social Media Data," in *ACM Conference on Human Factors in Computing Systems 2017 (CHI2017)*, Denver, CO, 2017.
- [13] D. D. Lee and S. H. Seung, "Learning the Parts of Objects by Non-negative Matrix Factorization," *Nature*, vol. 401, pp. 788-791, 1999.
- [14] J. An, H. Kwak, and B. J. Jansen, "Validating Social Media Data for Automatic Persona Generation," in *The Second International Workshop on Online Social Networks Technologies (OSNT-2016), 13th ACS/IEEE International Conference on Computer Systems and Applications (AICCSA2016)*, Agidar, Morocco, 2016.
- [15] A. Odlyzko, "Privacy, economics, and price discrimination on the Internet," presented at the Proceedings of the 5th International Conference on Electronic Commerce, Pittsburgh, Pennsylvania, USA, 2003.
- [16] V. Orso, T. Ruotsalo, J. Leino, L. Gamberini, and G. Jacucci, "Overlaying social information: The effects on users’ search and information-selection behavior," *Information Processing & Management*, vol. 53, pp. 1269-1286, 2017.