# Study And Analysis For The Prediction Of Human Behaviour And Comment Volume On Social Media Using Machine Learning Approaches

[1]Dr. Anuj Bhardwaj, [2]Dr. Navneet Kaur, [3]Dr. Ankur Dumka, [4]Parag Verma

[1]Professor-Computer Science & Engineering, Chandigarh University, Mohali, Punjab.
[2]Assistant Professor, Computer Science & Engineering, Bababanda Singh Bahadur Engg. College Fatehgarh Sahib
[3]Associate Professor-Computer Science & Engineering, Women Institute of Technology (Govt.), Dehradun, India
[4]Assistant Professor, Computer Science & Engineering, Uttaranchal University, Dehradun, India

Email: [1]anuj2k3@gmail.com, [2]navneet.sehgal@bbsbec.ac.in, [3]ankurdumka2@gmail.com, [4]parag_verma@yahoo.com

**Abstract**
*Utilization over Internet has been altogether expanded during most recent couple of decades. People groups saving additional time via web-based networking media locales. In this exploration proposition, we are intrigued to anticipate the character of clients by assessing their tweets. Up to this point, to accurately measure customer's characters, they predictable to get a character test. This made it unrealistic to utilize character examination in plentiful online networking spaces. In this examination proposition, we apply neural systems by which a client's character can be precisely anticipated through the freely accessible data on their Twitter profile. We will portray the sort of information gathered, our strategies for assessment, and the AI methods that permit us to effectively foresee character. This data is essential for organizations to target possible buyers or look for client suppositions in case of enhancement as a business methodology. In this way, this work examines online networking information to anticipate huge character characteristics, for example characteristics or qualities explicit to a person. The main strides towards web based life locales, raises information size and volume. The measure of information that is transferred to these person to person communication administrations is expanding step by step. Along these lines, there is gigantic prerequisite to contemplate the exceptionally unique conduct of clients towards these administrations. This is a starter work to demonstrate the client designs and to contemplate the viability of AI prescient displaying approaches on driving long range interpersonal communication administration Facebook. We demonstrated the client remark patters, over the post on FB Pages & anticipated that what number of comments a position is required to obtain in next H hours.*
*Keywords: Machine Learning, Natural Language Processing, Online Social Network, Personality Test, Profiling, Sentimental Analysis, Twitter.*

## 1. INTRODUCTION

Individuals have begun to invest their energy in sites that anybody can alter and add to. In

this way, to satisfy this require, a number of web advancements, wherever clients can intelligently team up & make a payment, are presented. These are presented for the sake of net 2.0 Technologies. These type of innovations empower clients to donate and divide contented without expecting them to have any specialized information in web programming. By the assistance of these advancements, individuals can connect with others with comparable interests.[1] During the most recent decades, some informal communication destinations are presented and have gotten profoundly mainstream in around the world. Every one of them has various destinations to convince individuals to share their encounters, thoughts or snapshots of their life kindly. Facebook gives clients a correspondence arrange comprising of their companions, families and others with whom they have associate in their genuine public activity.[2] Twitter empowers individuals to communicate their thoughts, moment pundits to others where they may most likely know each other, in actuality. LinkedIn centers on business life, and it gives a business organizing stage to representatives to impart, follow one another and help their enrollment through improved looking through offices dependent on their callings.[3] These person to person communication locales influence our reality. Numerous individuals are efficient in such stages. For example, Twitter has become a significant elective media to genuine media, it is quicker to spread news and gives more ability to speak freely. Web-based social networking is where clients present themselves to the world, detection individual subtleties and bits of acquaintance into their lives. We are preliminary to see how a segment of this data can be used to get better the clients' encounters with interfaces and with each other.

## 2. PROBLEM STATEMENT

Individuals have a characteristic need to communicate to others in the network by sharing their encounters, thoughts, exercises, and recollections. As methods, they for the most part want to utilize online life, for example, Twitter, Facebook, individual web journals, and wikis. Numerous individuals reliably add to such web based life stages by composing their own encounters, sharing photographs and status. Most of shared substance is close to home data. There are concentrates in the writing which utilize shared internet based life substance to foresee clients' Big 5 Personality Traits, for example, pleasantness, good faith, extraversion, neuroticism and receptiveness. These investigations typically use phonetic highlights, interpersonal organization data, and the recurrence of their cooperation with the stage, for example, number of posted announcements, photographs, recordings and preferences. The point of this exploration is to recognize which highlights of the mutual substance in Facebook are connected with clients' Big 5 Personality Traits and build up a model dependent on these highlights for character forecast Second, we show that the consideration of data viewing clients' companions, for example, their Big 5 Personality data improves the exactness contrasted with different strategies in the writing. Character forecasts customarily rely on client's profile data, notices, messages posted by them, and so on. This work presents a character forecast framework that works with a gathering of tweets as opposed to a solitary line of text. This framework involves three modules. The primary module extricates meta-information from Tweets and makes a 'meta-base'. It doesn't think about client's profile data. The subsequent module changes multi-name issue into five twofold grouping issues. At long last, a multilayer neural system arrangement calculation is applied to decide character characteristics. Since the internet based life information is being produced at quicker rate, this framework uses various named information to characterize the unlabeled information. To exhibit the working of this framework, set of Tweets specific to a client are acquired from Twitter API and applied to the framework. The Multilayer Perception Neural Network AI calculation is utilized as the classifier. The framework produces constructive outcomes and precisely predicts the character of 'Tweeters' without thinking

about their profile data, which is an alternate methodology from those looked into in the writing. The main patterns towards interpersonal interaction administrations had drawn monstrous open consideration from one and half' decade.[4] The converging up of processing with the physical things had empowered the transformation of regular articles into data machines. These administrations are going about as a multi-device with routine applications e.g.: news, notices, correspondence, remarking, banking, showcasing and so forth. These administrations are altering step by step and significantly more are in transit. These all administrations share day by day enormous substance age practically speaking, that is bound to be put away on Hadoop bunches. As in Facebook, 500+ terabytes of new information ingested into the databases consistently, 100+ petabytes of plate space in one of FBs biggest Hadoop (HDFS) groups what's more, their is 2.5 billion substance things shared every day (announcements + divider posts + photographs + recordings + remarks). Flickr highlights 5.5 billion pictures as that of January 31,2011 and around 3k-5k pictures are including every moment. This exploration proposition, additionally uncovered upon the main informal communication administration FB, in especially 'FB Pages' , for programmed investigation of patterns and examples of clients. Along these lines, for this work, we built up a product model that comprises of crawler, Information extractor, data processor and information revelation module. Our exploration is arranged towards the remark volume expectation (CVP) that a record is required to get in next H hours.

## 3.    BEHAVIOUR MODEL

We chose to utilize the 5 Model in this examination, since it is at present the most across the board and by and large acknowledged model of character and its capacity to foresee human conduct has been all around contemplated.[5][6] This model has-been appeared to subsume the most known character qualities and gives a terminology and a theoretical system that brings together a significant part of the exploration discoveries in brain science of individual contrasts and character.  The Five Factor Model partitions character into five dimensional characteristics: Openness to Experience, Conscientiousness, Extraversion, Agreeable, and Neuroticism (OCEAN). Each measurement has its delegate qualities. Receptiveness to encounter quantifies an individual's creative mind, awareness, looking for of novel encounter and eagerness for culture, thoughts, and style. Reliability mirrors how much an individual is sorted out, determined and trustworthy. Extraversion quantifies an individual's propensity to look for incitement in the outer world, organization of others, and express positive feelings. Appropriateness quantifies the degree to which a personality is centered on maintenance up optimistic social relations, mirroring a propensity to be trustful, thoughtful and helpful. The five characteristics have been seen to be hereditarily heritable, stable after some time and predictable across sexual orientations, societies, and races.[7]

Table 1 summarize the huge 5 personality traits along with their delegate expressive terms for both low down and high scorers.

| Portrayal | Transparency Openness is identified with creative mind, innovativeness, interest, resilience, political progressivism, and gratefulness for culture. Individuals scoring high on Openness like change, acknowledge new and irregular thoughts, and have a decent feeling of style |
|---|---|
| Conscientiousness | Conscientiousness actions leaning for a collected way to treaty with life as different to an unrestrained one. Individuals scoring high on Conscientiousness are clear to be competent, solid, and consistent. |
| Extroversion | Extroversion gauges a propensity to look for incitement in the external world, the association of others, and to converse positive feelings. |

| | Individuals scoring high on sociability will in general be all the more cordial, pleasant, and communally dynamic. |
|---|---|
| Agreeableness | Agreeableness recognize with an concentration on observance up positive social relations, being benevolent, empathetic, and agreeable. Individuals scoring high on Agreeableness will in general trust others and adjust to their provisions. |
| Neuroticism | Emotional constancy, contrarily alluded to as Neuroticism, quantifies the propensity to encounter temperament swings and feelings, for example, blame, outrage, tension, and despondency. |

## 4. BACKGROUND AND RELATED WORK

Earlier examination has indicated that character can proficiently clarify a significant measure of fluctuation in human inclinations and conduct across various spaces, for instance media and social inclinations [8][9], and long range interpersonal communication sites utilization.[10] According to data preparing hypothesis, the fulfillment individuals get from outside incitement, relies upon their ideal or favored excitement levels. One's inclination more than one thing is believed to be influenced by the relating data preparing limit and full of feeling directions.[11] Character is along these lines seen as significant for understanding people's energy about human expressions, for instance, artistic creations and music.[8][11] Ongoing exploration proposed that character qualities might be measured as significant middle people of media content inclinations. They found that receptiveness plainly empowers an enthusiasm for mind boggling and energizing recreational practices.[12] Good faith and benevolence (agreeableness) tend to effectsly affect exercises that are either troublesome or unusual, while enthusiastic security adversely impacts more unsurprising ways to get out from regular day to day existence. The work in [13] demonstrated that site inclinations are impacted by character attributes, similar to those for substance in genuine planet. The creators establish that site crowds regularly have unmistakable character profile, and the connection among character and inclinations identified with site and site classes is mentally significant. As of late, online networking sites (e.g., FB, Twitter) have risen as a significant media individuals speak with one another and express their sincere beliefs. Specialist shave become inspired by how character impacts client cooperation's on those web based life sites. The work in [14] demonstrated that Extroverts will in general discover internet based life website simple to utilize and valuable. Clients are probably going to choose contacts with comparative character attributes, and they for the most part will in general incline toward individuals high in Agreeableness.[15] Current investigation premiums have been more centered around the relations among character and clients' use practices (e.g., the quantity of posts, likes) and profiles (e.g., the quantity of companions/followings/adherents, age, sex) in social websites.[10][16] Additionally, expanding consideration has been paid on the expectation of character attributes scores dependent on those publically accessible conduct and profile information.[16][17] Golbeck et al. [17] demonstrated that clients with various character will in general utilize different words in their posts and portrayals. Quercia et al. [16] contemplated Twitter clients and found that both mainstream clients and powerful are outgoing individuals and genuinely steady. They further found that famous clients are 'creative' (high in Openness), while persuasive will in general be 'composed' (high in Conscientiousness). In [10], Querciaet al. inspected the connection between sociometric notoriety (number of Face book contacts) and character attributes on an alternate long range informal communication stage, Facebook. They reasoned that mainstream Facebook clients will in general have a similar character as individuals well known in reality. Also, [18] showed a huge association between character characteristics and different highlights of Facebook profiles. As far as we could possibly know, scarcely anSy investigations have been

done on the impacts of character on clients' conduct in client inclination demonstrating. In this examination proposition, we are attempting to respond to this focal exploration question: how much does a character factor influence rating practices? Client produced content on Twitter (e.g., tweets) additionally gives a significant wellspring of data for gathering clients' character qualities. One of the Twitter datasets frequently utilized in the writing is gathered through them Personality venture. Among a large number of members engaged with the my Personality venture, just a couple many clients presented joins on their Twitter accounts, which frames the substance of this dataset. This informational index has been utilized for the errand of consequently anticipating the characters of the clients, just as for client conduct examinations.[16][17][19] For example, Quercia et al. [16] found that social butterflies and sincerely stable individuals are famous just as powerful clients on Twitter. It was additionally seen that mainstream clients are creative, while powerful individuals on Twitter are more composed.

## 5. METHODOLOGY TO PREDICT    PERSONALITY    ON    SOCIAL    MEDIA        PLATFORMS

This work makes obtainable a performance prediction system for social network data analysis. The following diagram flow of the scheme depicted in Figure 1. Here the scheme consists of 4 modules: such as Data collection, pre-processing, transformation and classification



Figure 1: Flow chart of Personality trait prediction system

These modules are clarified underneath.

*1) Information Collection:*

For showing the framework, we'd like tweets announce by Associate in Nursing individual(s). For this, tweets are non heritable utilizing Twitter API. Twitter API[1] provides Associate in Nursing entrance to twitter info together with information concerning the purchasers, tweets distributed by a consumer, list things on twitter so on. Tweet object is in .json style.

*2) Pre-handling module:*

The tweets are 1st gotten from the tweet object. At that time the framework separates meta-properties from the tweets. the information removed are often isolated into social conduct and grammar data. The linguistic information incorporates: traditional length of text, traditional range of positive and negative words, traditional range of uncommon characters like comma, punctuation so forth. The social conduct information incorporates: traditional range of connections, traditional range of hash labels, traditional range of notices so forth. the conventional is decided by separating the all out of a linguistic and social conduct classification attribute by absolute length of the tweets. within the wake of removing the meta-traits, they're sent to the modification module.

*3) Transformation module:*

This module changes the "multi-name issue into twofold grouping issues". This module gets the meta properties separated from the past module. Utilizing this information, it builds a part

---

[1] Twitter API available on https://apps.Twitter.com/

vector. every position within the vector compares to a meta-quality. This vector is then sent to the grouping module whose yield are going to be either '1' (yes) or '0' (no). this transformation is important on the grounds that multilayer recognition neural system that's used as classifier will simply distinguish numbers.

*4) Classification module:*

A Multilayer Perception (MLP) Neural Network is employed for order. There are 5 neural systems (classifiers),one for each character attribute. The arrangement module gets a preparation set antecedently modified additionally to check set. every neural system contrasts the check embody vector and its preparation highlight vectors. The yield of each classifier is either '1'(yes) or '0' (no) contingent upon whether or not the vectors matches or not, more surmising that the individual has the character characteristic or no.

## 6. METHODOLOGY TO PREDICT COMMENT VOLUME ON SOCIAL MEDIA PLATFORMS

We focused on fine grained discerning demonstrating procedures. For fine grained expectation, we tend to address this issue as a relapse issue. Given some posts that showed up in past, whose target esteems (remarks got) are currently legendary, we tend to reproduced things. The assignment is to anticipate that what range of remarks that a post is relied upon to induce in next H hrs. For this, we tend to crept the Face book pages for crude data, pre-handled it, and created a transient split of the data to line up the preparation and testing set. At that time, this preparation set is used to organize the agent and execution of repressor's at that time assessed utilizing testing data(whose target esteem is roofed up) utilizing some assessment measurements. This complete procedure is exhibited in Figure 2 and purpose by point during this space.
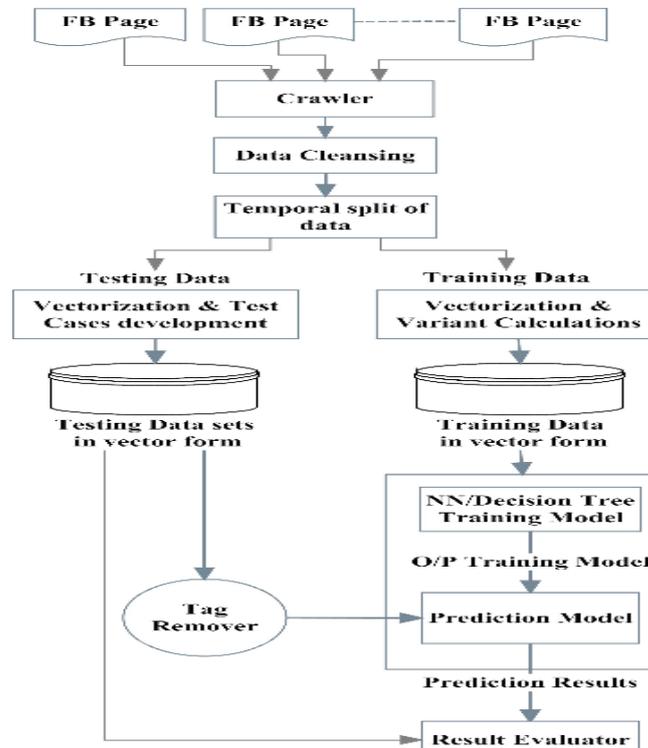


Figure 2: Comment Volume Prediction Process Demonstration

*(I) Feature set utilized for this work*

We had distinguished sixty nine highlights Associate in nursing one as target an incentive for every post and ordered these highlights as:

**Page Features:** we tend to distinguished four highlights of this classification that comes with

highlights that characterize the notoriety/Likes, classification, checking's and discussing of wellspring of report. Page likes: it's an element that characterizes purchasers support for express remarks, pictures, divider posts, statuses, or pages.

**Page Category:** This characterized the classification of wellspring of report e.g.: native business or spot, whole or item, organization or institution, craftsman, band, amusement, network so forth. Page Checking's: it's an indication of indicating closeness at specific spot and beneath the category of spot, organization pages because it were.

Page talking About: this can be the real tally of purchasers who are 'locked in' and collaborating therewith Face book Page. The purchasers who extremely come to the page, within the wake of preferring the page. This incorporates exercises, for instance, remarks, likes to a post, and offers by guests to the page.

**2) Essential Features:** This remembers the instance of remark for the post in several time stretches w.r.t to the haphazardly selected base date/time exhibited in Figure 3, named as C1 to C5.
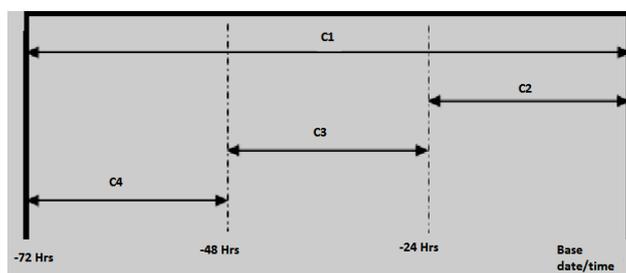


Figure 3: Demonstrating the essential feature details.

C1: Total remark tally before chosen base date/time.

C2: Comment embody in last twenty four hrs w.r.t to selected base date/time.

C3: Comment tally is last 48hrs to last twenty four hrs w.r.t to base date/time.

C4: Comment embody in initial twenty four hrs behind distributing the report, but before the selected base date/time.

C5: the excellence somewhere within the vary of C2 and C3. Besides, we have a tendency to accumulated these highlights by supply and engineered up some determined highlights by ascertaining min, max, normal, middle and variance of five antecedently mentioned highlights. On these lines, together with the five basic highlights and twenty five determined basic highlights, we have a tendency to got thirty highlights of this category.

3) Weekday Features: Binary pointers (0, 1) are utilized to talk to the day on that the post was distributed and also the day on selected base date/time. Twenty five highlights of this type are distinguished.

4) alternative essential Features: This incorporate some report connected highlights like length of record, delay between chosen base date/time and archive distributed date/time ranges from (0, 71), archive advancement standing esteems (0, 1) and post share tally. Five highlights of this classification are recognized.

(ii) Travel The information begins from Facebook pages. The crude data is crept utilizing crawler that's meant for this examination work. This crawler is planned utilizing JAVA and Face book search language (FQL). The crude data is slithered by crawler and cleansed on premise of following rules:

• we have a tendency to thought-about, simply those remarks that was distributed in most up-to-date 3 days w.r.t to 1base date/time because it is traditional that the more experienced posts for the foremost half don't get to any extent further thought.

• we have a tendency to excluded posts whose remarks or another essential subtleties are absent. on these lines we have a tendency to delivered the cleansed data for examination.

(iii) Pre-preparing

The slithered data can't be utilised foursquare for investigation. during this manner, it's

brought out through various procedures like split and vectorization. we have a tendency to created short split on this corpus to accumulate making ready and testing informational assortment as we will utilize the past data(Training information) to organize the model to form expectations for the long run data(Testing data) [20][21] this can be finished by selecting a position time and partitions the whole corpus in 2 sections. At that time this data is exposed to vectorization. To utilize the knowledge for calculations it's needed to alter that information in to vector structure. For this alteration, we have a tendency to had recognized a number of highlights as of currently talked concerning during this space, on that remark volume depends and adjusted the accessible data to vector structure for calculations. The procedure of vectorization is distinctive in making ready and testing set:

**1) Coaching set vectorization:** beneath the preparation set, the vectorization procedure goes in corresponding with the variation age method. Variation is characterized as, what variety of occurrences of conclusive making ready set is gotten from single example/post of preparing set. This can be finished by selecting distinctive base date/time for same post indiscriminately and forms them solely as pictured in Figure 4 a pair of. Variation - X, characterizes that, X occurrences are determined structure single making ready case as represented just in case of facebook official page id: 107684304856555with post id: 219515841718793, denote on TueFeb 19 08:14:19 IST 2018, post crept on SatFeb2414:51:48 IST 2018. It got complete of 515 remark Sat time of locomotion as appeared in Figure 4.
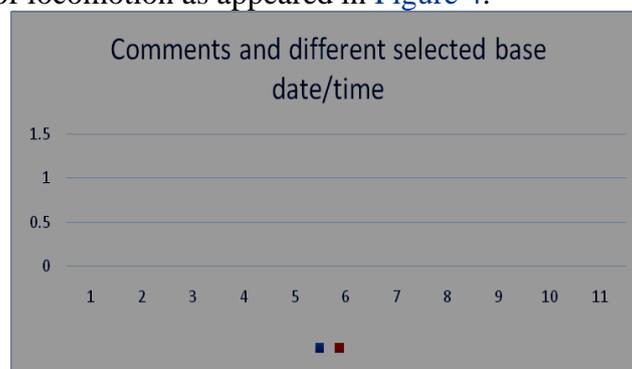


Figure 4: Collective Comments and different selected base date/time

**REFERENCES**

[1]  A. Kamilaris and A. Pitsillides, "Social networking of the smart home," in *21st Annual IEEE International Symposium on Personal, Indoor and Mobile Radio Communications*, 2010, pp. 2632–2637.

[2]  K. Shvachko, H. Kuang, S. Radia, and R. Chansler, "The hadoop distributed file system," in *2010 IEEE 26th symposium on mass storage systems and technologies (MSST)*, 2010, pp. 1–10.

[3]  I. Polato, R. Ré, A. Goldman, and F. Kon, "A comprehensive view of Hadoop research—A systematic literature review," *J. Netw. Comput. Appl.*, vol. 46, pp. 1–25, 2014.

[4]  T. Reuter, P. Cimiano, L. Drumond, K. Buza, and L. Schmidt-Thieme, "Scalable event-based clustering of social media via record linkage techniques," 2011.

[5]  T. P. Costa Jr, "The NEO-PI-R professional manual: Revised NEO Five-Factor Inventory.(NEO-FFI)," *Psychol. Assess. Resour.*, 1992.

[6]  O. P. John and S. Srivastava, "The Big Five trait taxonomy: History, measurement, and theoretical perspectives," *Handb. Personal. Theory Res.*, vol. 2, no. 1999, pp. 102–138, 1999.

[7]  O. P. John, R. W. Robins, and L. A. Pervin, *Handbook of personality: Theory and research*. Guilford Press, 2010.

[8]  P. J. Rentfrow and S. D. Gosling, "The do re mi's of everyday life: the structure and personality correlates of music preferences.," *J. Pers. Soc. Psychol.*, vol. 84, no. 6, p. 1236, 2003.

[9]  G. Kraaykamp and K. Van Eijck, "Personality, media preferences, and cultural participation," *Pers. Individ. Dif.*, vol. 38, no. 7, pp. 1675–1688, 2005.

[10] D. Quercia, R. Lambiotte, D. Stillwell, M. Kosinski, and J. Crowcroft, "The personality of popular facebook users," in *Proceedings of the ACM 2012 conference on computer supported cooperative work*, 2012, pp. 955–964.

[11] H. Ganzeboom, "Explaining differential participation in high-cultural activities: a confrontation of information-processing and status-seeking theories," *Theor. Model. Empir. Anal. Contrib. to Explan. Individ. actions Collect. Phenom.*, pp. 186–205, 1982.

[12] M. Zuckerman, R. S. Ulrich, and J. McLaughlin, "Sensation seeking and reactions to nature paintings," *Pers. Individ. Dif.*, vol. 15, no. 5, pp. 563–576, 1993.

[13] M. Kosinski, Y. Bachrach, P. Kohli, D. Stillwell, and T. Graepel, "Manifestations of user personality in website choice and behaviour on online social networks," *Mach. Learn.*, vol. 95, no. 3, pp. 357–380, 2014.

[14] P. A. Rosen and D. H. Kluemper, "The impact of the big five personality traits on the acceptance of social networking website," *AMCIS 2008 Proc.*, p. 274, 2008.

[15] J. Schrammel, C. Köffel, and M. Tscheligi, "Personality traits, usage patterns and information disclosure in online communities," *People Comput. XXIII Celebr. People Technol.*, pp. 169–174, 2009.

[16] D. Quercia, M. Kosinski, D. Stillwell, and J. Crowcroft, "Our twitter profiles, our selves: Predicting personality with twitter," in *2011 IEEE third international conference on privacy, security, risk and trust and 2011 IEEE third international conference on social computing*, 2011, pp. 180–185.

[17] J. Golbeck, C. Robles, and K. Turner, "Predicting personality with social media," in *CHI'11 extended abstracts on human factors in computing systems*, 2011, pp. 253–262.

[18] Y. Bachrach, M. Kosinski, T. Graepel, P. Kohli, and D. Stillwell, "Personality and patterns of Facebook usage," in *Proceedings of the 4th annual ACM web science conference*, 2012, pp. 24–32.

[19] D. J. Hughes, M. Rowe, M. Batey, and A. Lee, "A tale of two sites: Twitter vs. Facebook and the personality predictors of social media usage," *Comput. Human Behav.*, vol. 28, no. 2, pp. 561–569, 2012.

[20] E. Cambria, B. Schuller, B. Liu, H. Wang, and C. Havasi, "Knowledge-based approaches to concept-level sentiment analysis," *IEEE Intell. Syst.*, vol. 28, no. 2, pp. 12–14, 2013.

[21] E. C. Tupes and R. E. Christal, "Recurrent personality factors based on trait ratings," *J. Pers.*, vol. 60, no. 2, pp. 225–251, 1992.