# A User Modeling Oriented Analysis of Cultural Backgrounds in Microblogging

Elena Ilina

Postbus 3017, 2601DA, Delft, The Netherlands

Email:Elena@Ilina.nl

## ABSTRACT

Adaptive applications rely on the knowledge of their users, with their needs and differences. For instance, training processes can be adapted to user origins using information on their cultural background. Our goal is to gather culture-specific information from publicly available microblogging content. For this, we analyze in this paper culture-specific microblogging behavior patterns. We monitor the usage of Twitter-specific elements including hashtags, web links and user mentions. We analyze how users from different cultural groups employ these elements when they tweet. On the analyzed user groups from different regions we identify distinctive microblogging patterns. Our findings reveal a culture-specific user behavior on Twitter which we explain in terms of previous sociological research. Since our results enable us to distinguish between different cultural origins of user groups, we propose a culture-oriented user modeling approach which enables us to capture user microblogging behavior patterns. User microblogging behavior provides an outlook into user preferences towards sharing content with others, time preferences and social networking attitudes.

## I  INTRODUCTION

Adaptive applications such as e-learning environments benefit from the knowledge of the cultural backgrounds of users. For instance, e-learning applications aiming to work with students from different cultural backgrounds benefit from a representation of culture-related aspects of the users. One of the case-studies of the ImReal[1] project involves learning how to effectively communicate with people from other cultural backgrounds. In this case, culture-oriented user modeling could take place by considering cultural aspects of users and using them in adapting the application behavior according to the user needs. However, cultural-oriented user modeling is not a trivial task, since it requires an in-depth understanding of user characteristics in relation to the concept of culture including nationality, religious and political views, education level, country of living and other residence locations which influence the real-life user experience [1].

As result of user modeling, user profiles representing user characteristics are created and used to adapt applications to user needs. In case user-related information cannot be retrieved directly from the user, or is not available, adaptive applications might exploit user data derived from external sources like social networks and microblogs. For instance, Twitter content can be used to create user profiles describing user interests [2]. Twitter profile data can provide information on a user's geographic locations and use of languages. Related data may also be extracted from microblogging content published by the user.

User preferences according to user's location can be extracted from microblogging content and this information stored in the user profiles. However, would it also be possible to derive culture-specific behavioral traits based on user microblogging activities? Can we ascertain culture-oriented behavioral patterns of user behavior on microblogs? Does the information derived from Twitter allow us to differentiate users belonging to different cultural groups? These questions motivated us to investigate how to mine cultural patterns of user behavior on Twitter.

In this work, we analyze microblogging behavioral patterns and relate our findings with sociology research on intercultural communication. We adopt the well-known Lewis model [3] of Cultures, used for describing differences in communication of people belonging to different cultural groups defined by nationalities. We assume that communication differences could be reflected in the way people blog. For this, we create stereotypical cultural background models reflecting their behavioral patterns on Twitter, based on a set of users with defined geographic locations. These stereotypical models allow us to get insights on user microblogging behavior and its differences among cultural groups. Our main contributions include the following:

- An analysis of user behavior on Twitter for five user groups of different cultural origin.

---

[1] http://www.imreal-project.eu/

- Culture-specific microblogging patterns as explained by culture-related characteristics from sociology research by Lewis [3].

- A Culture-oriented User Modeling (CUL-UM) approach based on user behavior in microblogs and its experimental assessment.

In the next section we briefly outline the scope of the project and related work. In section 3 we outline the background of the Lewis model of Cultures, describe our research methodology and the experimental setup for our Culture-oriented User Modeling (CUL-UM) experiments. In section 4 we provide an analysis of Twitter features usage for the selected countries and user groups, and report on the quality of the created CUL-UM, which is based on experiments predicting users cultural origins. We conclude with our main findings on user behavior for the five cultural groups and provide insights on further user modeling research directions considering cultural behavioral patterns.

## II  RELATED WORK

Previous research works on personalization and adaptive systems exploit information published in social network platforms in order to collect information on user traits and interests. For instance, [2] uses Twitter for creating content-based user profiles, which are further aligned with news articles in their news recommendation experiments. [4] demonstrated that information from several social networks, including Twitter, Facebook and LinkedIn can be used to provide improved recommendations. As a possible application, [5] proposed a generic adaptive system based on Twitter data.

Even though existing approaches can be used for gathering information on particular user traits, only a few selected works investigate how the user microblogging behavior differs between distinctive cultural user groups. [6] studied microblogging behavior by analyzing amongst other language usage and network-related features for ten countries on Twitter. Japanese and English blogs analysis was performed by [7] identifying opinion differences between both cultural groups on selected topics using Wikipedia as reference. Trending topics analysis was performed by [8] revealing international differences of users interests to news and topic popularity across cultures of six countries based on a large tweets dataset. Tweets as a source of information on music genres popularity was studied by [9] finding that user preferences differ amongst countries and cities. On Facebook social networking web site, functionality employed and usage time differ across cultures [10]. The recent work by [11] compares user behavior on Twitter and Weibo, linking identified behavioral patterns with the culture model by Hofstede. Hofstede studies cultural differences and their impact on social interactions by relating people from different cultural origins to personality traits including power distance, individualism or collectivism, uncertainty avoidance [12]. The Hofstede model was also adopted in the study by [13] investigating correlation of social network sites' functionality usage and cultural user backgrounds among countries including the USA, Korea and China.

The Hofstede model is widely applied in studies comparing social networking with the help of cultural dimensions such as individualism and collectivism or uncertainty avoidance [14]. Such cultural dimensions can be analyzed to design components of social networking sites customized to related user traits. Considering cultural differences is important for businesses operating on the Global market and when localization to certain countries/cultures is preferred over standardization. When implementing web sites targeting certain cultures, the functionality and design adaptation is paramount for improving users experience as previous studies such as [15] indicate [14]. [10] also emphasizes needs for localization based on the findings revealing cultural differences of using Facebook features.

Another social model developed by Lewis [3] is based on creating stereotypical communication profiles in relation to different nationalities. Lewis defines three cultural dimensions in respect of how people from different cultures communicate, whether their focus is on people or factual information. Since our main goal is to model user behavior based on microblogging activities, we consider the Lewis model due to its focus on communication patterns. We extend our research on mining cultural patterns from Twitter described in [16] and provide insights into culture-oriented user modeling and adaptation.

## III  RESEARCH METHODOLOGY

In this section, first we outline our conceptual framework for CUL-UM, which is based on the Lewis Model of Cultures. Second, we formulate research questions in relation to findings from previous research. Third, we describe the experimental setup and dataset used.
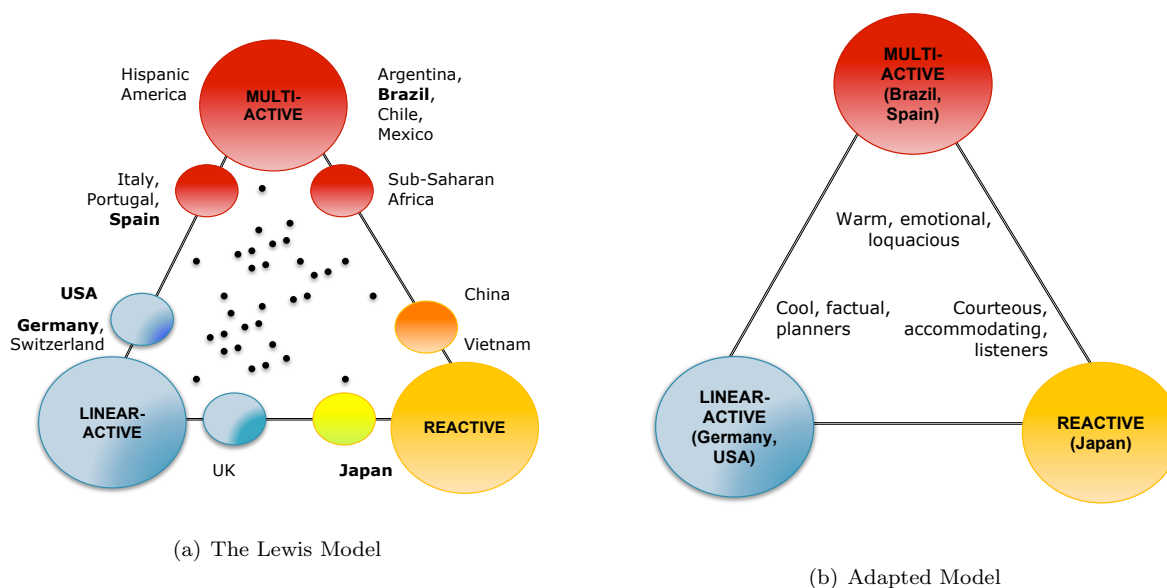
(a) The Lewis Model

(b) Adapted Model

fig.1 The Lewis Model of Culture (simplified and adopted from [3])

# 1 CONCEPTUAL FRAMEWORK AND APPROACH

In [3] Lewis analyses the cultural personality dimensions in relation to the country of origin or nationality. For instance, the multi-active Hispanic users group includes citizens from Argentina, Spain and Brazil, the linear-active group from Germany and the USA. The reactive dimension reflects the group of citizens from countries such as Vietnam and, to a lesser extent, Japan. Persons from the linear-active group share some similarities amongst each other such as a focus on planning activities, factual information and respect towards authorities. The reactive group can be associated with polite behavior and conflict avoidance. Citizens from other countries are placed between these extreme groups, and each person can be described in terms of reactivity, linear-activity and multi-activity traits [3].

The Lewis model of Culture is represented in the form of a triangle with corner points depicting the cultural dimensions mentioned above, as shown in figure 1 (a). These cultural dimensions reflect differences in the way people with different cultural backgrounds communicate [3]. In our experiments, as shown in figure 1 (b), we selected users from Germany and Brazil located in the apexes of the Lewis model of Cultures and representing linear-active and multi-active user

groups respectively [16]. The USA and Spain were added to the respective user groups even though these countries are not located directly at the apexes of the triangle. This enabled to analyze the behavior of the aforementioned user groups and how their behavior differs in respect of the Lewis research findings. We selected Japan for representing a reactive user group even though it is not depicted in an apex of the Lewis model, since it is listed as one of the top most active countries on Twitter[2]. This is why we selected Twitter users from Germany and the USA for representing the linear-active group, users from Japan for representing the reactive group, and users from Brazil and Spain for representing the multi-active group, as shown in figure 1 [16]. The inclusion of five countries enabled a comparison the user groups originating from these countries. User groups from countries belonging to the same cultural profile, corresponding to the Lewis model, are further aggregated for the purpose of comparison. We believe that this approach can be used for further modeling user profiles of users from different countries into the three cultural profiles according to the Lewis model. This could be advantageous for applications targeting cultural differences.

In order to acquire knowledge on user traits related to the cultural background of a user, we propose to mine them from microblogging activities of the user. Based on the analysis of microblogging behavioral

---

[2]http://semiocast.com/publications/2012_01_31_Brazil_becomes_2nd_country_on_Twitter_superseds_Japan

patterns, culture-oriented user modeling can be performed. In result, user profiles with information on culture-specific user traits and preferences can be created and used in the adaptation process as shown in figure 2.
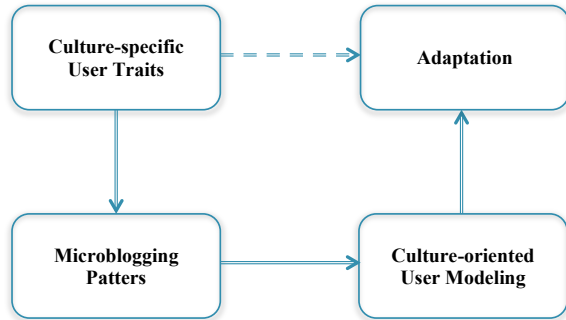


fig.2 Inferring Cultural Characteristics for CUL-UM

## 2 RESEARCH QUESTIONS

We base our investigation on the idea that personality traits as defined by Lewis [3] are also reflected in the way how people blog on Twitter. The previous works by [9], [10] inform us about cultural differences in music listening patterns derived from the Twitter posts, usage of Facebook features and time spent differ respectively across countries.

### 2.1 SHARING CONTENT

[17] stated that hashtags can be used not only for conversational purpose, but also for organizing content, in which case the hashtag standard deviation time is higher than the standard deviation time of hashtags used in conversations. Also, in [11] we read, that users from less individualistic societies might refrain from using hashtags since they do not like their tweets to appear in trending topics. [6] stated that the USA is the first country in their list, leading in Uniform resource locator (URL) sharing. We assume that this can be explained by the linear-activity characteristics of these users. In his work [3] Lewis wrote that linear-active persons are "data-oriented" and prefer to work with factual information, while reactive users are more "dialogue-oriented". This is why we consider Hashtags and URLs for comparing linear-active users with others.

### 2.2 FOREIGN LANGUAGES USAGE

[6] found that the English language is the most popular language on their dataset, created from Twitter content and accounts for more than half of tweets published by users in the ten countries analyzed.

Taking into account the widespread usage of the English language on Twitter and the challenges regarding automatic language detection as stated in [6], we investigate how the number of detected languages differs between the analyzed user groups. We are interested in comparing foreign languages usage of multi-active and reactive persons, which in accord to Lewis model [3] are both people-oriented, while multi-active persons are more extraverted and loquacious.

### 2.3 TWEETING MOBILITY

Since tweeting mobility is also interesting to investigate as mentioned in [6], we analyze how users from different countries use Twitter while travelling. Since linear-activity persons can be very conscious about effectively allocating their time [3], we assume that linear-active persons might use Twitter on their travels. We therefore compare their tweeting behavior with reactive persons.

### 2.4 POSTING TIME

Tweeting time during weekends or weekdays was important to relate with the reactive user, who generally has a different perception of time, being very punctual and polite, as outlined in [3]. This is why we assume that reactive persons might employ Twitter more on weekends compared to other persons, particularly linear-active persons, since multi-active persons tend to do things "at the same time", as described in [3].

### 2.5 REFERRING TO OTHER USERS

Since findings by [6] show that Japanese users mention other users the least, we investigate user mentions employed by the analyzed user groups. Also, [3] states that persons from Japan generally employ less names than persons from Western countries.

## 2.6 SOCIAL NETWORK SIZE AND CONVERSATIONS

We also consider conversation and social network features for analyzing user communication on Twitter. [18] stated that retweets can be used to give credit to other bloggers or even self-promotion. A study by [6] showed that Japanese users retweet the least. We investigate the retweet frequency of the analyzed user groups and compare the results with the previous research by [6]. Besides, since multi-active and reactive persons are people-oriented as stated in [3], we compare these two user groups. As stated in [3], since reactive cultures value silence and more in-depth content, we assume that reactive persons might refrain from retweets in opposite to multi-active users. This is why we hypothesise that reactive users might reply more, since they are indicated as being very good listeners in [3].

In a nutshell, our main research goal consists in analyzing Twitter microblogging behavior for users from different cultural groups as defined in the Lewis model. We study how users from linear-active countries (Germany and the USA), reactive countries (Japan) and multi-active countries (Spain and Brazil) use Twitter, and investigate how Twitter behavior differs between these different cultural groups. We focus on content-based, activity-based and social network-based features. We explore the following research questions referring to the usage of aforementioned Twitter features:

### Content-based characteristics

- RQ1 (Hashtags usage): How does the usage of hashtags differ between cultural groups?

- RQ2 (URLs): How often do users from different cultural groups share URLs?

- RQ3 (Languages): How do users from different cultural groups make use of foreign languages in their posts?

### Activity-based characteristics

- RQ4 (Mobility): To what extent do the different groups of users publish tweets from different geographic locations?

- RQ5 (Weekends): How frequently do users post on weekends compared to weekdays? Do these temporal Twitter traits differ between the different cultural groups?

### Social Network-based characteristics

- RQ6 (Friends): How does the number of friends differ between the cultural groups?

- RQ7 (Followers): Is there a relation between the number of followers that a user has on Twitter and the cultural background of the user?

### Conversation characteristics

- RQ8 (Mentions): How often do users from different cultural groups refer to other users?

- RQ9 (Replies): How often do users from different cultural groups reply to other users?

- RQ10 (Retweets): To what extent do users from different cultural groups retweet?

The above research questions refer to different features which describe certain aspects of the users' behavior on Twitter. In our analysis, we compare for which features the cultural groups exhibit the most respectively least profound differences. Following the feature analysis, we model stereotype user profiles and perform a series of experiments for predicting a user belonging to a specific stereotype profile. This helps us to investigate how well our model works for different cultural user stereotypes and how can we describe user activities on Twitter in relation to the Lewis' model.

## 3 EXPERIMENTAL SETUP

In order to perform culture-oriented user modeling on Twitter, first of all we identified differences in microblogging behavior of people from different countries. For this, we selected Twitter users who indicated their location in Twitter profiles. It is important to mention however, that we do not profile users into gender and age groups. We analyse instead all users having countries defines in their profile. In [19], we read that cultural statistics on personality traits for 26 countries showed similar personality levels for men and women. Additionally, age was profiled in a similar way across countries. Therefore, we assume that users of different age and gender groups can be combined to profile aggregate cultural groups.

For our experiments, we selected five countries, including Japan, Germany, the USA, Brazil and Spain,

which match with the cultural dimensions of Lewis' model as shown in figure 1. In order to find Twitter users belonging to the selected countries, we employed Twitter Streaming API[3] providing samples from public data streams. We selected users having more than ten friends and tweets, and having less than 5000 followers. For all selected users the location field mentioned the corresponding country. We also define geographic locations to include large cities such as Berlin, Hamburg and Munich for Germany, Tokyo, Yokohama and Osaka for Japan, New York, San Francisco and Washington D.C. for the USA, São Paulo and Rio de Janeiro for Brazil. For the tweets collected for the defined user groups, we analyze the use of Twitter-specific features. Our aim was to find behavioral patterns for these cultural user groups and explain them in relation to the Lewis' model of cultures. Overall, we performed the following steps:

- STEP 1: For the defined culture groups Germany, Japan, Spain, Brazil and the USA we selected a set of users tweeting from their respective geographic locations[4].

- STEP 2: Using User Twitter Streams, during two months we collected tweets published by the users selected in STEP 1. In order to get a solid understanding of users' behavioral characteristics on Twitter, we limited our dataset to users who published at least 100 tweets as shown in figure 3. Our threshold of 100 tweets enabled us to aggregate user microblogging behavior for 11998 users. This allowed us to analyze the user behavior for more than 1000 of users for each country. In addition, during the two months crawling period, we analysed users mobility defined as tweeting from different geographic locations. We identified the country name using the Geonames API[5] and Google Geocoding API[6].

- STEP 3: After completing the crawling process, we summarized, based on 100 randomly selected tweets published by each user, the tweeting behavior of each user in a user profile including Twitter-specific characteristics such as the use of hashtags, user mentions and link sharing. On the user profiles created, we analyzed with descriptive statistics how the Twitter-specific behavior differs between culture groups.

We employed t-tests for identifying which user groups behave differently and the level of significance.

- STEP 4: Next, we created classification tree models based on the features set analysed on the previous step. For this, we used a set of user profiles created in STEP 3. The classification experiments allowed us to assess the predictive value of the analyzed features. We used our set of features as a set of numeric variables for building the decision tree classifiers.

- STEP 5: Finally, we evaluated the classification tree using a resubstitution method and ten-fold cross-validation. This allowed us to assess the classification accuracy and quality of generated CUL-UM user profiles.

| Country | Number of Users | Users Posted at least 100 Tweets |
|---------|-----------------|----------------------------------|
| Japan | 4885 | 2984 |
| Spain | 4906 | 3119 |
| Brazil | 4910 | 2935 |
| USA | 1714 | 1316 |
| Germany | 2823 | 1644 |

fig.3 Users Dataset
(crawled from 2012-03-26 to 2012-06-01)

## IV ANALYSIS OF BEHAVIORAL TWITTER FEATURES

Next, in order to investigate how Twitter users from different culture groups behave on Twitter, we provide an analysis of Twitter-specific features:

- *Content-based* features including Twitter characteristics such as the usage of URLs, hashtags and the number of automatically detected languages[7] in the user content;

- *Activity-based* features describing Twitter activities such as the balance of tweeting during weekends against weekdays and the number of tweets from different geographic locations;

---

[3]https://dev.twitter.com/docs/streaming-apis/streams/public
[4]https://dev.twitter.com/docs/streaming-apis/parameters#locations
[5]http://www.geonames.org/export/web-services.html
[6]https://developers.google.com/maps/documentation/geocoding/
[7]http://code.google.com/p/language-detection

| Hypothesis and related Research Questions | $t$ | $df$ | $\mu_1$ | $\mu_2$ | Result |
|---|---|---|---|---|---|
| Content-based features: Hashtags, URLs, Number of Languages Detected (RQ1 to RQ3) | | | | | |
| $H_{1(a)}$ Linear-active users use Hashtags the most | 21.8 | 4188.3 | 31.9 | 17.4 | $\mu_1 > \mu_2$ |
| $H_{1(b)}$ Reactive users share Hashtags the least | -41.6 | 10379 | 7.6 | 25.6 | $\mu_1 < \mu_2$ |
| $H_{2(a)}$ Linear-active users include URLs the most | 14.4 | 5109 | 39.6 | 31.6 | $\mu_1 > \mu_2$ |
| $H_{2(b)}$ Reactive users share URLs the least | -3.7 | 6471.4 | 32.1 | 34.0 | $\mu_1 < \mu_2$ |
| $H_{3(a)}$ Multi-active users employ the most foreign languages | 51.4 | 11145 | 1.1 | 0.4 | $\mu_1 > \mu_2$ |
| $H_{3(b)}$ Reactive users employ the least of foreign languages | -9.8 | 6044.7 | 0.16 | 0.8 | $\mu_1 < \mu_2$ |
| Activity-based features: Mobility and Weekends (RQ4 to RQ5) | | | | | |
| $H_{4(a)}$ Reactive users tweet the least from different locations | -30.3 | 5791.3 | 0.6 | 0.9 | $\mu_1 < \mu_2$ |
| $H_{4(b)}$ Linear-active users tweet the most from diff. loc. | 15.4 | 4703.3 | 0.9 | 0.8 | $\mu_1 > \mu_2$ |
| $H_{5(a)}$ Reactive users tweet the most on weekends | 22.2 | 6109.4 | 28.6 | 24.3 | $\mu_1 > \mu_2$ |
| $H_{5(b)}$ Linear-active users tweet mostly during weekdays | -6.1 | 5395.3 | 24.5 | 25.7 | $\mu_1 < \mu_2$ |
| Social Network-based features: Friends and Followers (RQ6 to RQ7) | | | | | |
| $H_{6(a)}$ Multi-active users have the larger number of friends | -6.1 | 12315 | 310.2 | 355.2 | $\mu_1 < \mu_2$ |
| $H_{6(b)}$ Linear-active users have the smaller number of friends | 6.2 | 4836.1 | 375.1 | 319.6 | $\mu_1 > \mu_2$ |
| $H_{7(a)}$ Multi-active users have the larger number of followers | -6.4 | 12853 | 315.1 | 376.8 | $\mu_1 < \mu_2$ |
| $H_{7(b)}$ Linear-active users have the least number of followers | 9.7 | 4234.6 | 442.7 | 316.2 | $\mu_1 > \mu_2$ |
| Conversation-based features: User Mentions, Replies and Retweets (RQ8 to RQ10) | | | | | |
| $H_{8(a)}$ Reactive users employ user mentions the least | -40.3 | 8052.8 | 46.5 | 71.0 | $\mu_1 < \mu_2$ |
| $H_{8(b)}$ Multi-active users mention other users the most | 22.6 | 13037 | 71.6 | 57.5 | $\mu_1 > \mu_2$ |
| $H_{9(a)}$ Reactive users have the most replies | 3.6 | 5456.8 | 27.2 | 25.8 | $\mu_1 > \mu_2$ |
| $H_{9(b)}$ Multi-active users have the least replies | -7.5 | 12837 | 24.9 | 27.4 | $\mu_1 < \mu_2$ |
| $H_{10(a)}$ Reactive users have less retweets | -37.8 | 7889.5 | 8.2 | 17.7 | $\mu_1 < \mu_2$ |
| $H_{10(b)}$ Multi-active users have the most retweets | 30.2 | 12802 | 18.9 | 11.4 | $\mu_1 > \mu_2$ |

fig.4 Research Questions and Hypothesis Test Results for Comparing Cultural User Groups (with $p < 0.001$)

- *Social Network-based* features including Twitter user mentions, number of friends and followers in the user network;

- *Conversational* features reflect the number of retweets and replies posted by a user;

Figure 4 above outlines the analyzed features and their relation with the research questions and hypotheses. These user features were derived from the Twitter profile of the users from the chosen country groups. Assuming that users from Japan belong to the reactive user group, the USA and Germany belong to the linear-active user group, and Brazil and Spain belong to the multi-active user group, we created user profiles based on the data collected from the user content. For establishing our hypothesis we assumed that user behavior on Twitter reflects the user's cultural background. For instance, tweeting time during weekends or weekdays was important to relate with the reactive user, generally having a different perception of time as explained in [3]. In addition, we also considered conversation and social network features for analyzing user communication on Twitter.

# 1 RESULTS OF FEATURES ANALYSIS

In order to facilitate our comparison of the defined features and user groups, we performed 2-sample t-test statistics assuming non-equal variances. Figure 4 shows Hypothesis and t-test results for the feature categories comprising the Content-, Activity-, Social Network- and Conversation-based categories. The results provide $t$ statistic values, $df$ values for associated degrees-of-freedom, values $\mu_1$ and $\mu_2$ representing mean values for the compared user groups on the culture-level. On the country-level, user groups were compared in [16] and further t-test results [8] shown in figure 5, where countries are denoted by their two-letter ISO 3166-1 country codes.

## 1.1 RQ1: HASHTAGS USAGE

T-test results show that mean values for linear-active user groups are significantly higher than means of users from other groups. This supports our hypothesis $H_{1(a)}$ that, in average, linear-active users use hashtags the most compared to other user groups.

---

[8]Supplementary material is available at `http://ilina.nl/projects/cultures`

| $G_1$ | $G_2$ | $t$ | $df$ | $p$ | $\mu_1$ | $\mu_2$ | $G_1$ | $G_2$ | $t$ | $df$ | $p$ | $\mu_1$ | $\mu_2$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| \multicolumn{14}{c}{RQ1: Hashtags} |
| ES | BR | 22.82 | 6500.3 | < 0.001 | 29.56 | 14.76 | **ES** | **US** | 0.88 | 2192.8 | > 0.05 | 29.56 | 28.67 |
| ES | JP | 38.15 | 5755.4 | < 0.001 | 29.56 | 7.63 | ES | DE | -4.89 | 2942.4 | < 0.001 | 29.56 | 34.41 |
| BR | US | -14.36 | 1973.6 | < 0.001 | 14.76 | 28.67 | BR | JP | 13.64 | 5713.1 | < 0.001 | 14.76 | 7.63 |
| BR | DE | -20.49 | 2655.8 | < 0.001 | 14.76 | 34.40 | US | JP | 22.84 | 1642.9 | < 0.001 | 28.67 | 7.63 |
| US | DE | -4.68 | 3005.1 | < 0.001 | 28.67 | 34.41 | JP | DE | -29.40 | 2217.7 | < 0.001 | 7.63 | 34.41 |
| \multicolumn{14}{c}{RQ2: URLs} |
| **ES** | **BR** | -1.93 | 6069 | > 0.05 | 30.78 | 32.09 | ES | US | -13.96 | 2279.9 | < 0.001 | 30.78 | 42.46 |
| ES | JP | -2.15 | 6917.4 | < 0.05 | 30.78 | 32.09 | ES | DE | -8.64 | 3287 | < 0.001 | 30.78 | 37.49 |
| BR | US | -11.55 | 2798.2 | < 0.001 | 32.09 | 42.46 | **BR** | **JP** | 0.001 | 6193.2 | > 0.05 | 32.09 | 32.09 |
| BR | DE | -6.40 | 3934.9 | < 0.001 | 32.09 | 37.49 | US | JP | 12.27 | 2356.8 | < 0.001 | 42.46 | 32.09 |
| US | DE | 5.12 | 2900.6 | < 0.001 | 42.46 | 37.49 | JP | DE | -6.86 | 3398 | < 0.001 | 32.09 | 37.49 |
| \multicolumn{14}{c}{RQ3: Languages} |
| ES | BR | -4.45 | 5930 | < 0.001 | 1.067 | 1.16 | ES | US | 39.78 | 4235.4 | < 0.001 | 1.066 | 0.20 |
| ES | JP | 48.54 | 4399.1 | < 0.001 | 1.06 | 0.16 | ES | DE | 3.01 | 4346.7 | < 0.01 | 1.06 | 0.993 |
| BR | US | 48.89 | 3732.2 | < 0.001 | 1.16 | 0.20 | BR | JP | 61.85 | 4589.7 | < 0.001 | 1.16 | 0.16 |
| BR | DE | 7.76 | 3781.7 | < 0.001 | 1.16 | 0.99 | US | JP | 2.55 | JP3 | < 0.05 | 0.20 | 0.16 |
| US | DE | -36.85 | 2937.6 | < 0.001 | 0.20 | 0.99 | JP | DE | -45.30 | 2394.4 | < 0.001 | 0.16 | 0.99 |
| \multicolumn{14}{c}{RQ4: Mobility} |
| ES | BR | -2.96 | 6464.5 | < 0.05 | 0.87 | 0.90 | ES | US | -3.26 | 2610.1 | < 0.05 | 0.87 | 0.91 |
| ES | JP | 22.08 | 6907.9 | < 0.001 | 0.87 | 0.601 | ES | DE | -6.89 | 2942.5 | < 0.001 | 0.87 | 0.98 |
| **BR** | **US** | -1.13 | 2234.8 | > 0.05 | 0.90 | 0.91 | BR | JP | 26.59 | 6412.7 | < 0.001 | 0.90 | 0.60 |
| BR | DE | -5.16 | 2606.1 | < 0.001 | 0.90 | 0.98 | US | JP | 20.99 | 2755.1 | < 0.001 | 0.91 | 0.60 |
| US | DE | -3.56 | 3080.7 | < 0.001 | 0.91 | 0.98 | JP | DE | -22.47 | 3073.7 | < 0.001 | 0.60 | 0.98 |
| \multicolumn{14}{c}{RQ5: Weekends} |
| **ES** | **BR** | -1.28 | 6453.4 | > 0.05 | 24.04 | 24.35 | **ES** | **US** | 1.81 | 2576.8 | > 0.05 | 24.04 | 23.50 |
| ES | JP | -19.34 | 6924.7 | < 0.001 | 24.04 | 28.57 | ES | DE | -4.3 | 3577.3 | < 0.001 | 24.04 | 25.26 |
| BR | US | 2.79 | 2719.3 | < 0.05 | 24.35 | 23.53 | BR | JP | -17.41 | 6473.5 | < 0.001 | 24.35 | 28.57 |
| BR | DE | -3.15 | 3727.2 | < 0.01 | 24.35 | 25.26 | US | JP | -17.09 | 2573.1 | < 0.001 | 23.50 | 28.57 |
| US | DE | -5.23 | 2949.2 | < 0.001 | 23.50 | 25.26 | JP | DE | 11.63 | 3574.2 | < 0.001 | 28.57 | 25.26 |
| \multicolumn{14}{c}{RQ6: Friends} |
| ES | BR | 6.01 | 6545 | < 0.001 | 335.37 | 282.14 | ES | US | -4.69 | 2074.3 | < 0.001 | 335.37 | 400.5 |
| **ES** | **JP** | -0.20 | 6524.7 | > 0.05 | 335.37 | 337.5 | **ES** | **DE** | -1.70 | 3138.6 | > 0.05 | 335.37 | 356.2 |
| BR | US | -8.63 | 1987.4 | < 0.001 | 282.14 | 400.5 | BR | JP | -5.36 | 6200.2 | < 0.001 | 282.14 | 337.5 |
| BR | DE | -6.13 | 2981.4 | < 0.001 | 282.14 | 356.2 | US | JP | 4.25 | 2591 | < 0.001 | 400.55 | 337.5 |
| US | DE | 2.75 | 2807.3 | < 0.01 | 400.55 | 356.2 | **JP** | **DE** | -1.41 | 3930.4 | > 0.05 | 337.47 | 356.2 |
| \multicolumn{14}{c}{RQ7: Followers} |
| ES | BR | -3.13 | 6398.3 | < 0.01 | 296.45 | 335.8 | ES | US | -9.19 | 1771.3 | < 0.001 | 296.45 | 501.6 |
| **ES** | **JP** | -1.82 | 6925 | > 0.05 | 296.45 | 318.3 | ES | DE | -6.27 | 3078 | < 0.001 | 296.45 | 398.9 |
| BR | US | -7.31 | 1870.1 | < 0.001 | 335.81 | 501.6 | **BR** | **JP** | 1.39 | 6425.6 | > 0.05 | 335.81 | 318.3 |
| BR | DE | -3.75 | 3308.2 | < 0.001 | 335.81 | 398.9 | US | JP | 8.21 | 1774.2 | < 0.001 | 501.6 | 318.3 |
| US | DE | 4.11 | 2421 | < 0.001 | 501.6 | 398.9 | JP | DE | -4.93 | 3087.1 | < 0.001 | 318.27 | 398.9 |
| \multicolumn{14}{c}{RQ8: User Mentions} |
| ES | BR | 30.23 | 6520.2 | < 0.001 | 83.91 | 57.94 | ES | US | 7.18 | 2215.8 | < 0.001 | 83.912 | 75.13 |
| ES | JP | 48.54 | 6555.9 | < 0.001 | 83.91 | 46.51 | ES | DE | 17.57 | 3613.8 | < 0.001 | 83.91 | 65.83 |
| BR | US | -14.02 | 2221 | < 0.001 | 57.94 | 75.13 | BR | JP | 14.76 | 6036.8 | < 0.001 | 57.94 | 46.51 |
| BR | DE | -7.65 | 3588 | < 0.001 | 57.94 | 65.83 | US | JP | 24.56 | 1869.7 | < 0.001 | 75.13 | 46.516 |
| US | DE | 6.88 | 2674.4 | < 0.001 | 75.13 | 65.83 | JP | DE | -20.13 | 2960.3 | < 0.001 | 46.516 | 65.83 |
| \multicolumn{14}{c}{RQ9: Replies} |
| ES | BR | 12.63 | 6447.3 | < 0.001 | 27.47 | 22.04 | ES | US | 2.20 | 2242 | < 0.05 | 27.474 | 26.18 |
| **ES** | **JP** | 0.50 | 6673.2 | > 0.05 | 27.47 | 27.24 | ES | DE | -2.11 | 3111.2 | < 0.05 | 27.47 | 28.66 |
| BR | US | -6.92 | 2367.5 | < 0.001 | 22.04 | 26.18 | BR | JP | -10.89 | 6540.9 | < 0.001 | 22.04 | 27.24 |
| BR | DE | -11.53 | 3266.8 | < 0.001 | 22.04 | 28.66 | **US** | **JP** | -1.70 | 2698.7 | > 0.05 | 26.18 | 27.24 |
| US | DE | -3.55 | 2959.1 | < 0.001 | 26.18 | 28.66 | JP | DE | -2.37 | 3705.3 | < 0.05 | 27.24 | 28.66 |
| \multicolumn{14}{c}{RQ10: Retweets} |
| ES | BR | 24.73 | 6475.3 | < 0.001 | 23.16 | 14.26 | ES | US | 17.77 | 2832 | < 0.001 | 23.16 | 15.00 |
| ES | JP | 44.29 | 6278.1 | < 0.001 | 23.16 | 8.22 | ES | DE | 18.89 | 3973.3 | < 0.001 | 23.16 | 14.95 |
| **BR** | **US** | -1.68 | 2415.4 | > 0.05 | 14.26 | 15.0 | BR | JP | 19.76 | 6256.1 | < 0.001 | 14.26 | 8.22 |
| **BR** | **DE** | -1.66 | 3418.3 | > 0.05 | 14.26 | 14.95 | US | JP | 16.19 | 2107.9 | < 0.001 | 15.00 | 8.22 |
| **US** | **DE** | 0.11 | 2919 | > 0.05 | 15.0 | 14.95 | JP | DE | -17.18 | 3001.8 | < 0.001 | 8.22 | 14.95 |

fig.5 Comparison of User Groups by Countries (rows shown in bold font indicate the cases with no significant differences between group means and $p > 0.05$)

User Group Germany has higher mean of hashtags usage compared with the USA user group ($\mu_1 = 34.4$, $\mu_2 = 28.7$, $p < 0.001$). It is important to mention that the means of hashtag usage are close for the user groups of the USA and Spain ($\mu_1 = 28.7$, $\mu_2 = 29.6$, $p > 0.05$), sharing a bit more of hashtags compared with users from the USA. Our experiments support the hypothesis $H_{1(b)}$ stating that reactive users use the least of hashtags compared to other user groups. The results of t-tests show the acceptance of null hypothesis, that users from Japan employ hashtags the least, in average at the very high significance level ($\mu_1 = 7.6$, $\mu_2 = 25.6$, $p < 0.001$). The country-level tests reveal that the user group from Japan shares less hashtags compared to the other four countries, in average.

## 1.2  RQ2: URLS USAGE

The results of the tests support the hypothesis $H_{2(a)}$. Linear-active users use URLs the most compared to other user groups. Country-level statistics reveal that users from the USA ($\mu_1 = 42.5$) employ more URLs compared to users from Germany ($\mu_2 = 37.5$, $p < 0.001$), in average.

Our tests support hypothesis $H_{2(b)}$ stating that reactive users from Japan share the least of URLs ($\mu_1 = 32.1$), in average. However, country-level statistics for users from Spain (multi-active) indicate that they share less URLs compared to users from Japan (reactive) ($\mu_1 = 30.8$, $\mu_2 = 32.1$, $p < 0.05$). Tests show a similar hashtag usage for Japan and Brazil users ($\mu_1 = 32.1$, $\mu_2 = 32.1$, $p > 0.05$).

## 1.3  RQ3: FOREIGN LANGUAGES

The hypothesis $H_{3(a)}$ is supported by our experiments, indicating that multi-active users employ the most of foreign languages automatically detected from the user-generated content compared to other user groups. The hypothesis $H_{3(b)}$ is also supported since our experiments show that reactive users from Japan employ the least of foreign languages in their tweets ($\mu = 0.16$) compared to other users. On the country-level, Japanese users employ less foreign languages followed by the USA, Germany, Spain and Brazil. Users from Brazil employ the most of foreign languages.

## 1.4  RQ4: MOBILITY

The hypothesis $H_{4(a)}$ can be accepted, since our statistic shows reactive users in average tweet less from different geographic locations compared to other user groups ($\mu_1 = 0.6$).

The hypotheses $H_{4(b)}$ can also be supported, since linear-active users (USA: $\mu_1 = 0.9$, Germany: $\mu_2 = 1$) tweet the most from different geographic locations, in average. It is important to note that all other country-level user groups except Brazil have smaller mean values for the mobility feature compared to the linear-active group. Brazil and USA user group means do not differ significantly in our tests ($\mu_1 = 0.9$, $\mu_2 = 0.9$, $p > 0.05$).

## 1.5  RQ5: WEEKENDS

Statistic of the tweets fraction published on weekends demonstrates that the hypothesis $H_{5(a)}$ is supported at the very high significance level. Reactive users from Japan have a larger amount of tweets on weekends ($\mu_1 = 28.6$) compared to other user groups, in average.

Our tests indicate that the hypothesis $H_{5(b)}$ can also be accepted, since linear-active users have a smaller fraction of tweets on weekends compared to other users in average. The same trend holds on country-level statistic indicating that users from Germany ($\mu_1 = 25.3$) and the USA ($\mu_2 = 23.5$) tweet less on weekends than others in average. Interestingly, mean values for Spain and Brazil ($\mu_1 = 24.0$, $\mu_2 = 24.3$, $p > 0.05$), and mean values for Spain and the USA ($\mu_1 = 24.0$, $\mu_2 = 23.5$, $p > 0.05$) do not differ significantly, which indicates a similar attitude of tweeting on weekends. Brazil and the USA users tweet less than German users on weekends.

## 1.6  RQ6: FRIENDS

The hypothesis $H_{6(a)}$ could not be supported, since the multi-active users have a smaller number of friends compared to other user groups, in average. Moreover, the tests also do not support the hypothesis $H_{6(b)}$, since linear-active users (USA: $\mu_1 = 400.5$, Germany: $\mu_2 = 356.2$) mostly have greater means of the number of friends compared with other user groups. On the country-level, means for the groups of Spain and Japan ($\mu_1 = 335.4$, $\mu_2 = 337.5$, $p > 0.05$), Spain and Germany ($\mu_1 = 335.4$, $\mu_2 = 356.2$, $p > 0.05$), Japan and Germany ($\mu_1 = 337.5$, $\mu_2 = 356.2$,

$p > 0.05$) do not differ significantly.

## 1.7 RQ7: FOLLOWERS

Similarly, the hypothesis $H_{7(a)}$ and $H_{7(b)}$ cannot be supported, since the multi-active users have smaller number of followers compared to other user groups, while linear-active users have greater number followers compared with other user groups, in average. On the country-level, users from Spain and Japan ($\mu_1 = 296.4$, $\mu_2 = 318.3$, $p > 0.05$), users from Brazil and Japan ($\mu_1 = 335.8$, $\mu_2 = 318.3$, $p > 0.05$) do not differ significantly in the number of followers they have in average.

## 1.8 RQ8: USER MENTIONS

The hypothesis $H_{8(a)}$ can be supported, since reactive users from Japan indeed mention other users the least, in average, compared to other user groups on the cultural-group and country-group levels. The hypothesis $H_{8(b)}$ can be supported, since multi-active users have greater means for user mentions compared to other users, in average. On the country-level, however, German users mention other users more than users from Brazil ($\mu_1 = 65.8$, $\mu_2 = 57.9$, $p < 0.001$).

## 1.9 RQ9: REPLIES

The hypothesis $H_{9(a)}$, stating that reactive users from Japan should have more replies in average compared to other cultural user groups, can be accepted at the very high significance level. On the country-level, users from Japan ($\mu = 27.2$) behave similarly to users from Spain ($\mu = 27.5$) and the USA ($\mu = 26.2$). The hypothesis $H_{9(b)}$ can also be accepted, since the average number of replies from the multi-active users is lower compared to other users. On the country-level, however, users from Spain replied more in average compared to users from the USA ($\mu_1 = 27.5$, $\mu_2 = 26.2$, $p < 0.05$). Statistics show, that users from the USA reply less actively compared to other users, for the exception of Brazil ($\mu = 22.0$).

## 1.10 RQ10: RETWEETS

The hypothesis $H_{10(a)}$ can be accepted at the very high significance level. Reactive users from Japan have a smaller number of retweets in average compared with other user groups on the culture-group and country-group levels. The hypothesis $H_{10(b)}$

is also supported in our experiments, showing that multi-active users retweet the most compared to other users in our dataset. On the country-level, however, Brazilian users ($\mu = 14.3$), belonging to the multi-active culture group, retweet less (not significantly) compared to users from the USA and Germany ($\mu_1 = 15.0$, $\mu_2 = 14.9$), which are linear-active. Overall, users from Brazil, Germany and the USA exploited the retweeting functionality in a similar way. Spanish users retweeted the most ($\mu = 23.2$) compared to other users groups.

Overall, for all our tests on the culture-level shown in figure 4, the calculated $p$ statistic was less than 0.001, indicating very highly significant differences between user groups on the culture-level. This corresponds to the chance of 99.9% that mean values significantly differ. Country-level tests indicate that mean values of features for country groups differ significantly in the majority of cases. Country groups of Spain and the USA, Brazil and Japan have comparable means of hashtags and URLs usage. Spain and Brazil, Spain and the USA user groups have similar mean values of number of tweets published during weekends. Spain and Japan user groups have comparable values of mean values for number of friends and followers. Spain and Germany, Japan and Germany have comparable means of number of friends, while Brazil and Japan have comparable means for number of followers. Spain and Japan, the USA and Japan have comparable means for number of replies. Brazil and the USA, Brazil and Japan, Japan and the USA employ retweets similarly.

## 2 USER GROUP MEAN DISTANCES

Based on the Multivariate Analysis of Variance, we draw scatter plots showing clusters of user groups by countries and cultural user groups in figure 6 (a) and figure 6 (b) respectively. The scatter plots help to visualize the differences between the user groups. Two canonical variables are used to distinguish between user groups. They are calculated from the means of the feature values analyzed.

The first canonical c1 variable helps to separate clusters for the country-level user groups of Spain, Japan, the USA and Brazil. As can be seen from the figure 6 (a), the clusters for the user groups Spain and Japan are separated vertically, while user groups from the USA and Brazil are located on about the same level.
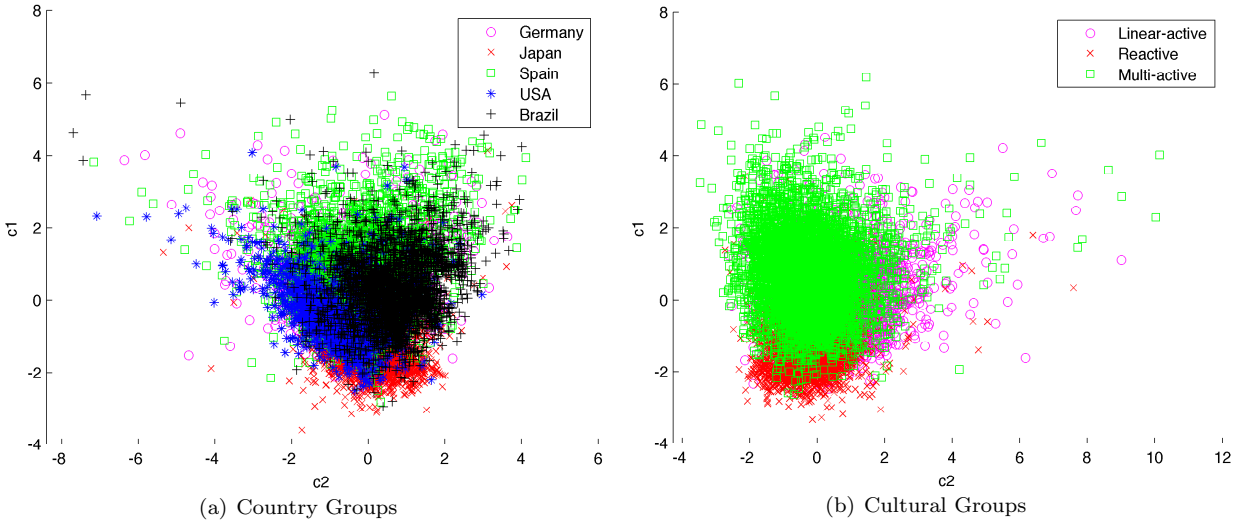
(a) Country Groups      (b) Cultural Groups

fig.6 Clusters Separation

On the culture-level, c1 helps to separate reactive users group depicted in red cluster below from other two clusters, multi-active users and linear-active users. This indicates that reactive users from Japan behave differently on Twitter, when considering the features set analyzed.

Similarly, the canonical variable c2 helps to separate user group clusters on the horizontal axis. On the country-level, c2 variable helps to distinguish clusters for users from USA and Brazil on the horizontal axis. On the culture-group level, c2 assists in separating multi-active users from the linear-active users. Figure 6 (b). demonstrated that the feature set enables a relatively good separation between reactive and two other cultural user groups. It is noted however, that the multi-active and reactive user group clusters overlap considerably.

Next, we calculate mean distances between user group means shown in the figure 7 and figure 8. As seen from figure 8. showing distances between each pair of group means for the mix of the aforementioned features, the distance between linear-active groups and multi-active group means (1.09) is much smaller than the distance between multi-active and reactive groups (4.06). For instance, the distance between German and Spain means is about 0.9, while distance between the Spain and Japan is about 4.65 as seen from figure 7.

|         | Japan | Spain | USA  | Brazil |
|---------|-------|-------|------|--------|
| Germany | **3.51** | 0.90  | 1.7  | 1.20   |
| Japan   |       | **4.65** | 2.19 | **3.17** |
| Spain   |       |       | **2.23** | 1.12   |
| USA     |       |       |      | 2.74   |

fig.7 Distances between Country Group Means

|              | Reactive | Multi-active |
|--------------|----------|--------------|
| Linear-active | **2.54** | 1.09         |
| Reactive     |          | **4.06**     |

fig.8 Distances between Cultural Group Means

Interestingly, distances between both linear-active groups (distance of 1.20 between Germany and Brazil, and 2.74 between the USA and Brazil) and Brazil are larger than between the linear-active groups and Spain (distance of 0.9 between Germany and Spain, and 2.23 between the USA and Spain). This coincides with the Lewis model in that Spain is more close to the linear-active triangle corner than Brazil, considered the "extreme" multi-active country. Therefore, we can conclude, that the linear-active and multi-active user groups are more similar in their behavior, while reactive users behave differently on Twitter in respect of the analyzed features.

## 3  PREDICTION QUALITY

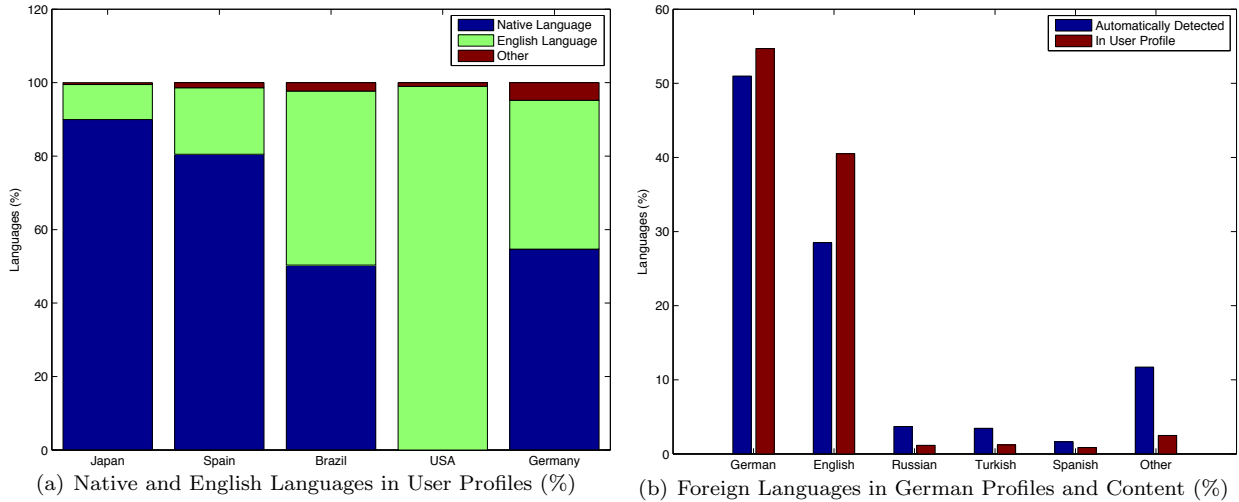In order to assess the quality of user profiling based on the analyzed feature set, we created six decision

(a) Native and English Languages in User Profiles (%)     (b) Foreign Languages in German Profiles and Content (%)

fig.9 Fraction of Languages Automatically Detected in the Tweets and Twitter user Profiles

| Country-level | | | | | Culture-level | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Test | Features | R.Err. | Nodes | C.Err. | Test | Features | R.Err. | Nodes | C.Err. |
| 1 | $LANG$ | 0.22 | 6 | 0.22 | 2 | $LANG$ | 0.17 | 4 | 0.17 |
| 3 | $DEF$ | 0.17 | 51 | 0.42 | 4 | $DEF$ | 0.10 | 47 | 0.29 |
| 5 | $DEF + LANG$ | 0.02 | 680 | 0.06 | 6 | $DEF + LANG$ | 0.01 | 511 | 0.04 |

fig.10 Resubstitution (R.Err.) and Cross-validation (C.Err.) Error Rates for Predicting User Groups with Decision Tree Classification (feature sets include the $DEF$ - features analysed in the section IV.A, $LANG$ - language in the user profile, $DEF + LANG$ - combination of previous two.)

tree classification models. The first two models (1 and 2) were created based on the language defined in the user profile. However, languages specified in the Twitter user profile could be misleading. For instance, in our dataset a large fraction of users from Germany specified their preferred language as "English", as show in the figure 9 . This is why we also created classification models (3 and 4) based on the selected features set while excluding languages defined in the user profile. Models 5 and 6 were created by combining features set and languages specified in the user profile.

The classification models enabled us to predict users belonging to a user group on country-level or culture-level. The classification models were assessed by calculating resubstitution error rate and testing error rate. For cross-validation, we split our sample into ten almost equally sized parts used for finding out the testing error rate. Figure 10 shows resubstitution errors, number of terminal nodes for pruned trees, and cross-validation errors for aforementioned tests and feature sets defined. As it can be seen from the table, when the profile information on languages is not available, the $DEF$ features set can be used to predict a user belonging to cultural dimensions or

one of the five countries, analyzed with a relatively high cross-validation error rate of 0.29 and 0.42, respectively. This indicates that the $DEF$ features set might be further extended with languages, other features when available in the profile or tweets content of the user. The combination of the $DEF$ and $LANG$ feature sets enables the lowest cross-validation error for culture-level and country-level classifications. The cross-validation error for feature set $DEF + LANG$ decreased in half compared to the cross-validation error when employing only languages defined in the user profile.

## 4 INTERPRETATION OF RESULTS AND DISCUSSION

### 4.1 CULTURAL DIFFERENCES AND COUNTRY-LEVEL SIMILARITIES

Based on descriptive statistics and comparison of mean values of features for different cultural groups, we found distinct differences between the reactive user group and other user groups. Japanese belong to the reactive users group, they share the least of hashtags and user mentions. Japanese reply however more

than other user groups, for the exception of Germans. Japanese reply quite a lot, and retweet less compared with other user groups. This can be explained by their good listening skills and "high-context" orientation as explained in [3]. Japanese users also tweet the least from different geographic locations. Moreover, we detected the least of foreign languages in the content published by reactive users compared to others. Japanese also tweet more on weekends compared with other user groups.

Interestingly, even though we initially hypothesized that multi-active people as more people-oriented persons might have larger social networks of friends and followers, tests showed that linear-active users from the USA and Germany have, in average, more followers [16]. They also have more friends compared to other user groups, except for users from Japan, for which they show a comparable mean value. Linear-active users also generally share more URLs compared with other user groups. Interestingly, German users belonging to the linear-active group have the greatest mean for replies compared to other users. Overall, linear-active users share also more hashtags compared with other groups but Spain. The means of hashtag usage are similar for Spanish users and users from the USA.

Moreover, multi-active users have similarities with linear-active user groups and are therefore difficult to separate. Considering the multi-active users group, Spanish refer the most to other users (mentions usage) and are quite similar in their behavior with the USA group, while Brazilians share fewer links, and only refer more to other users than Japanese. For multi-active users, we detected more foreign languages in average compared with linear-active users. They also have a smaller number of followers and friends compared to others.

Our findings agree with the study of [11] indicating that persons from Eastern countries are less individualistic, refraining from the usage of hashtags. In [6], users from South Korea and Japan have a smaller fraction of hashtags in their tweets. Our experiments also correspond with findings of [6] in that Japanese persons employ less user mentions than persons from Western countries. Our findings reveal that Japanese users retweet the least, which corresponds with [6], while they reply the most. This corresponds with [3] stating that reactive persons are generally good listeners and prefer in-depth content.

Our findings also correspond with the study of Lewis [3] in that linear-active Western persons are "data-oriented". We found a similar pattern of URLs usage as in [6] where users from the USA share the most URLs compared to others. In opposite, as explained in [3], multi-active and reactive persons are "people-oriented". Our experiments support this idea, since persons from Spain, Japan and Brazil share less URLs compared to the "data-oriented" persons from the USA and Germany. Multi-active persons are described as locatious in [3], in our experiments, Brazil and Spanish users employ also the most of foreign languages. To summarise, some of the findings correspond with the previous studies by [6] and [3]. This indicates that we found similar microblogging culture-specific behavior patterns even though working with different data-sets of Twitter users. It therefore appears that human communication in social networks could be influenced by cultural differences, which could be further explored in future studies to facilitate better user experience in social or virtual environments.

## 4.2 CLUSTER ANALYSIS

The cluster analysis showed that the distance between Spain and Germany is smaller than the distance between clusters of Germany and the USA. Also, linear-active users behave similarly to multi-active users when analyzing clusters formed from the multivariate analysis of their variances for the analyzed features. The user group from Germany is difficult to separate from the user groups of the USA, Spain and Brazil. We explain it by possible cultural similarities between these user groups and how they behave on Twitter. It is also reasonable to assume that this could be explained by the peculiarity of our dataset or in relation with Lewis model. This is why we cannot confirm the strict relation with the Lewis' model.

Moreover, [3] stated that Spanish people coming from different regions might behave very similarly to linear-active people in the sense of productivity. The geographic proximity have also a strong impact on personality across cultures [20]. This implies that there are more variables and relationships which might be considered for creating cultural user models based on microblogs. For instance, the features set can be further extended with topics derived from the tweets content and user opinions mined in a process similarly to works such as [7] and [8]. Cross-cultural topics analysis in tweets can be considered as a direction of future research.

In addition, the study [11] informs us about different fractions of positive posts for users from Sina Weibo and Twitter. This is why more features, as for instance emoticons could be added to the classification model to reflect differences in expressing feelings and moods. Real-life communication differences between people of different cultures as explained by sociological models thus can be further analyzed in the context of microblogging behavior and self-expression. A possible research direction could be to investigate how could we mine affective states from microblogs and how they reflect real-life communication patterns.

## 4.3 LOCALISATION AND ADAPTATION ASSUMPTIONS

Nevertheless, microblogging patterns on the country-level still can reveal users' attitudes on how they use the Twitter functionality. The insight that linear-active users from Germany and the USA tend to share more URLs and hashtags, have a larger contact network might suggest that the related microblogging functionality can be further enhanced for these users. For instance, a reply button functionality could be more visible for reactive users willing to participate in a more substantial dialogue, instead of providing a button for retweeting, which might be preferred by users from Brazil, Spain and the USA.

Furthermore, the distance between clusters for linear-active and multi-active users is about 1, while the distances between reactive and multi-active, between reactive and linear-active user clusters is about 4 and 2.5 respectively. It seems that the features analysis shows us that reactive users stay apart from the other two groups. As it was suggested by Lewis in [3], marketing efforts should not neglect reactive and multi-active persons, which worldwide are more than linear-active persons. The design and functionality of social networking websites and other applications can be tailored to the particular cultural user groups to reflect their preferences. In this sense, our findings agree with [10] and [14] on localization benefits for social network services targeting users from different cultural origins.

## 4.4 DATA COLLECTION AND EXPERIMENTAL SETUP

It is important to note that our study was based on the users having indicated their geographic location in their user profile. We have restricted our crawling process to the big cities in five selected countries. As

it is advised by [14], more in-depth research is needed to analyze more countries and social networking services. We agree with this and in future work we aim to extend our framework with more countries/users to allow analysis on a larger scale.

Our original dataset included in average more than 600 tweets per user. For building individual user profiles we considered however only 100 tweets, since otherwise we would only be able to model less than 300 users from the USA user group. Therefore, following our assumption that classification performance increases given more users, we selected 100 tweets as a starting point in our experiments based on more than thousand users per country group. In further experiments we plan to extend our users dataset and investigate the number of users/tweets required to build representative user profiles for modeling cultural origins. This would allow to better understand how classification performance scales with number of users and tweets included into the user profile. We believe that increasing the number of users would enable better prediction outcomes in the classification experiments we performed. Besides, we do not distinguish between age and gender of our users, which could be of interest for a cross-cultural analysis of user behavior on Twitter or other microblogging services.

Moreover, our experiments on classifying user profiles showed that we can employ classification methods such as decision trees to classify users into particular user groups on the culture and country levels. The analyzed feature set extended with the language defined in the twitter profile of the user enables a low cross-validation error rate. However, in case when language information is not available, the language can be inferred automatically from user content. Alternatively, more features derived from the user profile/content can be further analyzed to improve quality of users classification, which can be performed using other methods such as logistic regression or ensemble classifiers. In further work, we aim to implement and analyze other classifiers in order to facilitate separation of users from linear-active and multi-active countries. Such classification can be further exploited by adaptive applications when knowledge on user cultural background is needed as mentioned in [16].

Nevertheless, whilst it is challenging to assess adaptation outcomes next to a statistical observation [1], user modeling efforts could be beneficial for improving user experience. Previous studies have shown that

users of adaptive applications can benefit from adaptive functionality features. An empirical study by [21] has shown that simple user modeling introduced into a commercial application influenced positively user perception over software capabilities. [22] found a positive correlation between learning outcomes and adaptability to a learner state of uncertainty in a dialog-based tutoring system providing adapted feedback to learner answers.

## V   CONCLUSIONS

In the foregoing, we analyzed microblogging behavior on Twitter for user groups from Germany, USA, Spain, Brazil and Japan. We found, that Japanese users behave very differently from the rest of the user groups. In comparison, they tweet more on weekends, reply more and share the least of hashtags and user mentions. In contrast, users from the USA and Germany generally share more URLs and have more friends compared with the other user groups. Users from Spain and Brazil stay apart in a way that they have some similarities with the rest of groups, but are difficult to differentiate when using the analyzed set of features. Multi-active users however appear to employ more foreign languages than others.

We reflected on the results with the help of the sociological model by Lewis. Whilst it was not possible to explicitly map cultural-related communication patterns to microblogging behavior on Twitter, some of the derived microblogging patterns enabled us to distinguish between different cultural groups on Twitter. Based on the found microblogging patterns, we proposed an approach of culture-oriented user modeling that considers cultural/country differences of the users. The information on user microblogging activities, preferences to information sharing and/or dialogs can be further exploited for designing adaptable applications which suit to user needs based on her cultural/country origins. In further work, we will perform user analysis on a larger user dataset and investigate other Twitter-specific features in order to find further insights on cultural differences of microblogging behavior. We aim to create more accurate cultural user models that might be exploited in adaptive applications such as micro-blogging services or recommender systems.

## ACKNOWLEDGMENT

## References

[1] K. Reinecke, G. Reif, and A. Bernstein, "Cultural user modeling with cumo: An approach to overcome the personalization bootstrapping problem", in *Workshop on Cultural Heritage on the Semantic Web, International Semantic Web Conference*, 2007.

[2] F. Abel, Q. Gao, G.J. Houben, and K. Tao, "Analyzing user modeling on twitter for personalized news recommendations", *User Modeling, Adaption and Personalization*, pp. 1–12, 2011.

[3] R.D. Lewis, *When cultures collide: Managing successfully across cultures*, Nicholas Brealey Publishing, 2000.

[4] F. Abel, E. Herder, G.J. Houben, N. Henze, and D. Krause, "Cross-system user modeling and personalization on the social web", *User Modeling and User-Adapted Interaction (UMUAI), Special Issue on Personalization in Social Web Systems*, vol. 22, no. 3, pp. 1–42, 2011.

[5] J. Hannon, E. Knutov, P. De Bra, M. Pechenizkiy, K. McCarthy, and B. Smyth, "Bridging recommendation and adaptation: Generic adaptation framework-twittomender compliance case-study.", in *Second international Workshop on Dynamic and Adaptive Hypertext*, 2011, p. 1.

[6] B. Poblete, R. Garcia, M. Mendoza, and A. Jaimes, "Do all birds tweet the same?: characterizing twitter around the world", in *Proceedings of the 20th ACM international conference on Information and knowledge management*. ACM, 2011, pp. 1025–1030.

[7] H. Nakasaki, M. Kawaba, T. Utsuro, and T. Fukuhara, "Mining cross-lingual/cross-cultural differences in concerns and opinions in blogs", *Computer Processing of Oriental Languages. Language Technology for the Knowledge-based Economy*, pp. 213–224, 2009.

[8] D. Wilkinson and M. Thelwall, "Trending twitter topics in english: An international comparison", *Journal of the American Society for Information Science and Technology*, 2012.

[9] M. Schedl and D. Hauger, "Mining microblogs to infer music artist similarity and cultural listening patterns", in *Proceedings of the 21st international conference companion on World Wide Web*. ACM, 2012, pp. 877–886.

[10] A. Vasalou, A.N. Joinson, and D. Courvoisier, "Cultural differences, experience with social networks and the nature of "true commitment" in facebook", *International journal of human-computer studies*, vol. 68, no. 10, pp. 719–728, 2010.

[11] Q. Gao, F. Abel, G.J. Houben, and Y. Yu, "A comparative study of users' microblogging behavior on sina weibo and twitter", *User Modeling, Adaptation, and Personalization*, pp. 88–101, 2012.

[12] G. Hofstede, "A european in asia ", *Asian Journal of Social Psychology*, vol. 10, no. 1, pp. 16–21, 2007.

[13] Y.G. Ji, H. Hwangbo, J.S. Yi, P.L.P. Rau, X. Fang, and C. Ling, "The influence of cultural differences on the use of social network services and the formation of social capital", *Intl. Journal of Human–Computer Interaction*, vol. 26, no. 11-12, pp. 1100–1121, 2010.

[14] E. Vitkauskaitė, "Cultural adaptation issues in social networking sites", *Economics and Management*, vol. 16, pp. 1348–1355, 2011.

[15] N. Singh, H. Zhao, and X. Hu, "Analyzing the cultural content of web sites: A cross-national comparision of china, india, japan, and us", *International Marketing Review*, vol. 22, no. 2, pp. 129–146, 2005.

[16] E. Ilina, F. Abel, and G.J. Houben, "Mining twitter for cultural patterns", in *ABIS 2012 Workshop on Personalization and Recommendation on the Web and Beyond*, 2012.

[17] J. Huang, K.M. Thornton, and E.N. Efthimiadis, "Conversational tagging in twitter", in *Proceedings of the 21st ACM conference on Hypertext and hypermedia*. ACM, 2010, pp. 173–178.

[18] D. Boyd, S. Golder, and G. Lotan, "Tweet, tweet, retweet: Conversational aspects of retweeting on twitter", in *System Sciences (HICSS), 2010 43rd Hawaii International Conference on*. IEEE, 2010, pp. 1–10.

[19] R.R. McCrae, "Trait psychology and culture: Exploring intercultural comparisons", *Journal of personality*, vol. 69, no. 6, pp. 819–846, 2001.

[20] J. Allik and R.R. McCrae, "Toward a geography of personality traits", *Journal of Cross-Cultural Psychology*, vol. 35, no. 1, pp. 13–28, 2004.

[21] L. Strachan, J. Anderson, M. Sneesby, and M. Evans, "Pragmatic user modelling in a commercial software system", *COURSES AND LECTURES-INTERNATIONAL CENTRE FOR MECHANICAL SCIENCES*, pp. 189–200, 1997.

[22] K. Forbes-Riley and D. Litman, "Adapting to student uncertainty improves tutoring dialogues", in *Proc. Intl. Conf. on Artificial Intelligence in Education*, 2009.