# TWITTER/X DATASET – ENTITY RELATIONSHIP DIAGRAM

**Display Name Dimensions**
- Characters
- Hashtags (#)
- Emojis

**Date Dimensions**
- Year
- Month
- Day

**Time Dimensions**
- Hour
- Minute
- Second

**Text Dimensions**
- Characters
- Usernames (@)
- Hashtags (#)
- Emojis
- URL Links

**Referenced Tweets Dimensions**
- Referenced Tweet ID
- Type
- Referenced User ID
- Referenced Tweet Text

**Engagement Options**
- Retweet Count
- Like Count
- Quote Count
- Reply Count

**Created At Dimensions**
- Date
- Time

**Public Metric Options**
- Followers Count
- Following Count
- Tweet Count
- Listed Count

**Profile**
- Display Name
- Username
- Created At
- Description
- Location
- Pinned Tweet ID
- Protected
- Public Metrics
- Verification

**Main Table**

| PK | Respondent ID |
|----|---------------|
| | Profile |
| | Tweets |
| | Likes |
| | Follows |

**Tweets**
- Text
- Created At
- Engagement Metrics (Public)
- Referenced Tweets
- Attachments
- Conversation ID
- Edit History Tweet IDs
- Reply Settings
- Language
- Location

**Type Options**
- Quoted
- Replied To
- Retweeted

**Location Dimensiions**
- Geo Location ID
- Full Place Name

**Poll Keys Dimensions**
- Poll Key ID
- Text
- Options

**Attachments Dimensions**
- Poll Keys
- Media Keys

**Reply Settings Options**
- Everyone
- Mentioned Users
- Following

**Likes**
- Liked Tweet ID
- Liked Tweet Author ID

**Poll Options Dimensions**
- Position
- Label
- Votes

**Media Type Options**
- Video
- Image
- GIF

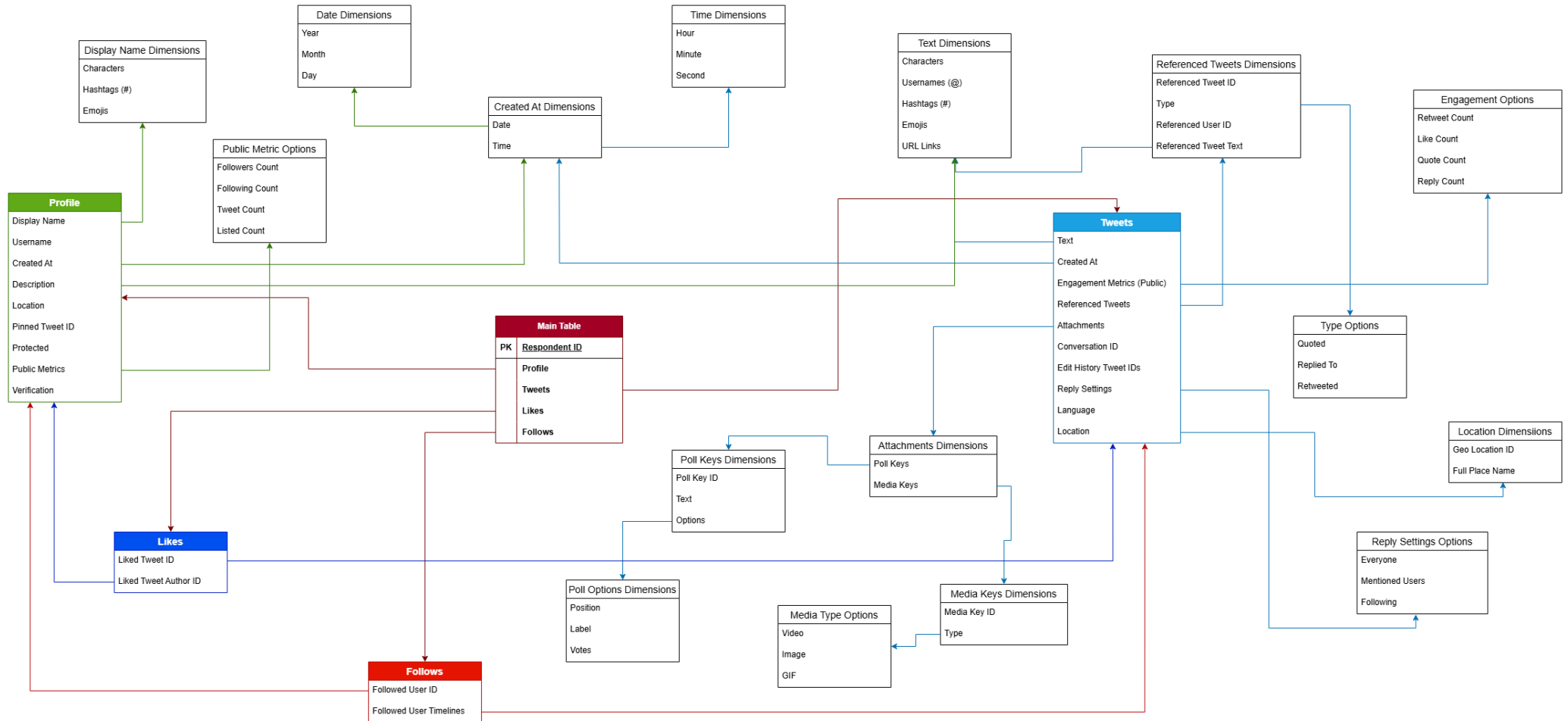**Media Keys Dimensions**
- Media Key ID
- Type

**Follows**
- Followed User ID
- Followed User Timelines

# DATASET 1a/1b – DiCED VARIABLES

This document is a (currently non-exhaustive) list of potential digital variables we could derive from **Dataset 1** (DiCED).  There are two tables. The first contains a list of **structural variables** which are designed to describe the general structure and content of a user's digital data. The second table contains a list of **substantive variables** which are designed to estimate new variables of interest about the user based on their digital data such as ideological position or degree of populism. Some of these variables also contain **validation variables** from the survey data where we can use the digital variable to validate what a user has said in the survey. The specific metrics we can use in the digital (Twitter/X) data to derive each variable are listed in the table which corresponds to the entity relationship diagram above.

**STRUCTURAL VARIABLES**: 22

**SUBSTANTIVE VARIABLES:** 17

**VARIABLES WITH VALIDATION:** 17

**TOTAL:** 56

**VALUE/COMPLEXITY/RELIABILITY:**

The final three columns in the table all contain a score between 1 and 5 based on an estimation of the overall **value** of the variable to researchers, the degree of **complexity** in building the variable (this could be quantified in time, resources and/or difficulty), and the **reliability** of the method to accurately capture the concept of interest. Higher scores for value and reliability indicate greater value/reliability, higher scores for complexity indicate a greater degree of complexity.

*Note:* Whenever referring to "distribution" in the measures column, this means general descriptive statistics such as: min, max, mean, median, mode, lower Q, upper Q, skewness, kurtosis, etc.

# STRUCTURAL VARIABLES

| No. | Variable Name | Main Table | Field Table | Field Dimensions | Description | Value | Complex. | Rel. |
|-----|---------------|------------|-------------|------------------|-------------|-------|----------|------|
| 1 | **Length of Membership** | Profile | Created At | Date, Time | Measures:<br><br>1. Age of account (in days/weeks/months)<br><br>I don't know how granular we can get with this date. We could in theory just post the created at date/time as is, or the date at the very least, but does this potentially risk a breach of anonymity(?).<br><br>This is also useful variable for quantifying tweet frequencies. (Tweet Count / Length of Mem.)<br><br>Possible Research Questions:<br><br>1. How long has a user had their acccount?<br><br>**SURVEY VALIDATION VARIABLE:**<br><br>**[uom_uspre_23]** And how long roughly have you held your account with this network? | **5** | **1** | **5** |
| 2 | **Profile Metrics** | Profile | Public Metrics | Followers Count, Following Count, Tweet | I imagine we can just release these numeric values *as is* seeing as they are not identifiable information? These variables are important descriptive statistics about each user's profile. | **5** | **1** | **5** |

| | | | | Count, Listed Count, Like Count | | | | |
|---|---|---|---|---|---|---|---|---|
| 3 | **Verification Status** | Profile | Verification | - | We can also release whether or not they are verified *as is*. TRUE/FALSE. I don't know how this changes with the new verification system on X.<br><br>Possible Research Questions:<br>Is a user verified or not? Presuming the data was taking before the change in verification status meaning after Twitter's rebranding to X, this can tell if if a user's account was *"active, authentic and notable"* – the previous verification criteria: https://help.x.com/en/managing-your-account/about-x-verified-accounts<br><br>If this is after changes to verification, it can still tell us a lot about whether or not a user has paid for a premium account or not and how this affects their usage and engagement statistics. What types of people are more likely to pay for a Twitter/X premium account? | **5** | **1** | **5** |
| 4 | **Location** | Profile | Location | - | Where a location can be matched using the manual location field (fuzzy matching), we could provide a high-level location that they have stated in their bio. Maybe at the city or regional level? Match this to the geographical level provided in the survey data.<br><br>**SURVEY VALIDATION VARIABLE:** | **5** | **1** | **3** |

| | | | | | Stated Geographical Location (if this is asked)<br><br>Do users report their correct location on their digital profiles? And to what level? | | | |
|---|---|---|---|---|---|---|---|---|
| 5 | **Named Entities** | Profile | Display Name, Username, Description | **Description:** Characters, Usernames, Emojis, Hashtags, URL Links | Similar to the location field, we could expand the use of a Named Entity Recogniser (NER) to identity any named entities (names/locations/ organisations/occupations/dates/times etc.)  in a user's profile metadata. As suggested by **RIZA:**<br><br>A named entity recogniser (NER) that can handle location/place names might be useful. An example is https://www.tandfonline.com/doi/full/ 10.1080/13658816.2022.2133125#d1e1297<br><br>We might not be able to make some of these publicly available if they are identifiable (Keep it high-level, only match what is reported in the survey data?).<br><br>**SURVEY VALIDATION VARIABLE:**<br><br>Stated Demographic Variables (if these are asked).<br><br>How much does the information a user gives away in their profile match what is stated in the survey? Where does it differ, and *how*? | **5** | **2** | **4** |

| 6 | Profile Formality | Profile | Display Name, Description, Location | **Description:** Characters, Hashtags, Emojis | 1. Character length distribution<br>2. Word length distribution<br>3. Score of grammaticality (could use a simple spellchecker for this?)<br>4. % of abbreviations/acronyms<br>5. % of capitalisation<br>6. % of hashtags<br>7. % of emojis<br><br>W could combine all of these measures together and create an equation that produces a single "formality" score for their profile description from 0 (extremely informal) to 1 (extremely formal). How "professionalised" is their bio?<br><br>**RIZA:**<br><br>An alternative measure might be the prediction of a classifier for formal/informal language, considering that even BERT-based models can obtain F-scores of 80+. See: https://aclanthology.org/2023.ranlp-1.31.pdf<br><br>Possible Research Questions:<br><br>How professionalised is a user's account?<br>Does appear to be more of an official 'work'-style account or something more casual? | 5 | 2 | 4 |

| 7 | **Post Frequency** | Tweets | Created At | Date, Time | 1. Overall number of tweets posted<br>2. Daily posting distribution<br>3. Time of day posting distribution<br><br>This could also be expanded to weekly, monthly, yearly tweeting patterns too if interested.<br><br>Possible Research Questions:<br><br>1. How frequently does a user post?<br>2. Are they a "power-user" or a "lurker"? (We could calculate this using the ratio of authored tweets to rewets, for example. See: https://www.pewresearch.org/short-reads/2022/03/16/5-facts-about-twitter-lurkers/)<br>3. Are they more of a morning tweeter of evening? How does this affect the way they post?<br>4. Is there any particular pattern to their tweeting? Does it increase in response to certain offline events? | 5 | 1 | 5 |
| 8 | **Post Categorisation** | Tweets | Referenced Tweets | Type | We can use the referenced tweet type field to categorise whether a tweet is a retweet, reply, quote or standalone tweet. We can then split standalone tweets into @mentions for any tweets that contain another user's handle. We can then provide high-level stats about user tweet types: | 5 | 1 | 5 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | 1. % Standalone Tweets<br>2. % Retweets<br>3. % Replies (direct/indirect)<br>4. % Quote Tweets<br>5. % Mentions<br><br>Possible Research Questions:<br><br>1. What type of general user are they? Are they primarily a "solo author" (mostly original tweets), an "amplifier" (more RTs), or an engager (Replies/QTs/Mentions)?<br>2. *How* do they engage with other users? We could combine this with the following variables about toxicity/sentimentality to assess if they negatively engage with users, if they regularly attack them, insult them, etc. | | | |

| 9 | **Post Engagement** | Tweets | Engagement Metrics | Retweet Count, Like Count, Quote Count, Reply Count | 1. Retweet Count Distribution<br>2. Like Count Distribution<br>3. Quote Count Distribution<br>4. Reply Count Distribution<br><br>These could also be combined into a single "general engagement" statistic. We could also consider looking at ratios to generate a positive/negative response variable. E.g: a tweet that receives more replies than likes/retweets is considered to be negatively received.<br><br>Helps to answer many questions around user impact and notability.<br><br>Possible Research Questions:<br><br>1. How much engagement do users get on average?<br>2. Do their posts get a lot of engagement and thus are of a higher quality, or do they get little interest?<br>3. How does engagement vary by the types of posts are user publishes? Do negative posts get more engagement that positive etc. | **5** | **1** | **5** |
| 10 | **Post Topics** | Tweets | Text, Attachments | **Text:** Characters, Hashtags, Usernames, Emojis | Could use an out-of-the-box tool for topic modelling on the textual data such as *BERTopic:* https://maartengr.github.io/BERTopic/index.html<br><br>This could allow us to broadly categorise the general things that users tweet about. Is it primarily political, | | | |

| | | | | | **Attachments:** Media Keys | sport, music, film, work, general life? etc. A similar method could also be applied to media content such as videos and images using a multimodal model (*BERTopic also offers this).*<br><br>We could even categorise sub-topics within politics if we wanted to focus in on the political aspect of online behaiour.<br><br>This variable can be combined with many other variables to get an understanding of how digital behaviour vary depending on what users talk about.<br><br>Possible Research Questions:<br><br>1. What topics do users tweet about the most? Are they a political user or something more casual like sport or film etc.?<br>2. How does topic moderate other variables such as toxicity, engagement, notability? For instance, do political posts get more engagement? Are politcal posters more active than general users etc?<br><br>**SURVEY VALIDATION VARIABLE:**<br><br>**[uom_uspre_26]** How often, if ever, do you post messages about politics on Twitter?<br><br>**Most Important Issue:** | **5** | **3** | | **4** |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | [uom_uspre_10] And of the issues you have said are important for you right now, which one would you say is the **most important**? | | | |
| 11 | **Post Attachments** | Tweets | Text, Attachments | **Text:** URLs<br><br>**Attachments:** Poll Keys, Media Keys | 1. **% of tweets containing URLs**<br>2. **% of tweets containing a poll**<br>3. **% of tweets containing media (Image, Video, GIF) –** these can be separated<br><br>By %, this could be quantified as either the overall % of tweets that contain each of these types, or could be combined into one general "attachment" variable. Higher use of URLs could also suggest a user type: "information propagator". Depends on what the URLs lead to.<br><br>Possible Research Questions:<br><br>1. How much does a user rely on non-textual content when posting?<br>2. How much does the use of non-textual content correlate with other types of digital behaviour? I.e: is it correlated with higher levels of conspiratorial content. | **5** | **1** | **5** |
| 12 | **Post Formality** | Tweets | Text | Characters, Hashtags, Emojis | 1. **Character length distribution**<br>2. **Word length distribution**<br>3. **Score of grammaticality** (could use a simple spellchecker for this?)<br>4. **% of abbreviations/acronyms** | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | 5. **% of capitalisation**<br>6. **% of hashtags**<br>7. **% of emojis**<br><br>By %, this could be quantified as either the overall % of tweets that contain each of these types, or could be median number per tweet. We could also combine all of these measures together and write an equation that produces a single "formality" score for each tweet from 0 (extremely informal) to 1 (extremely formal).<br><br>**RIZA:**<br><br>An alternative measure might be the prediction of a classifier for formal/informal language, considering that even BERT-based models can obtain F-scores of 80+. See: https://aclanthology.org/2023.ranlp-1.31.pdf | **5** | **2** | **4** |
| 13 | **Post Targets** | Tweets | Referenced Tweets | Referenced User ID | Along with the referenced tweet type, we can also get additional information about who they have referenced specifically, and also potentially the tweet they have quoted/retweeted/replied to. We could extract additional information from this if we wanted. (E.g: what is the type of content they typically reply to/retweet?)<br><br>Possible Research Questions: | **5** | **2** | 4 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | 1. What type of users do they engage with most? Are they mostly ordinary users or public figures?<br>2. Are the targets of their tweets mostly political types or something else?<br>3. How does their posting change depending on who their targeting? | | | |
| 14 | **Post Toxicity** | Tweets | Text | Characters, Emojis (?) | Use *Google's* Perspective API: https://perspectiveapi.com/<br><br>Provides a flagship **"Toxicity"** score, but can also provide additional attributes:<br><br>   1. **Severe Toxicity**<br>   2. **Insult**<br>   3. **Profanity**<br>   4. **Identity Attack**<br>   5. **Threat**<br>   6. **Sexually Explicit**<br><br>We could provide an overall mean/median toxicity score for each respondent's set of tweets, as well as for each individual attribute if of interest. This could give us a high-level idea of how "toxic" a respondent is.<br><br>Possible Research Questions:<br><br>   1. How 'toxic' are a user's posts overall? | **5** | **3** | **3** |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | 2. How does this breakdown by particular types of toxicity? <br><br> We could combine this metric with other variables such as tweet topics to see how toxicity levels change depending on what users primarily tweet about. Are political users more toxic than sports users etc.? | | | |
| 15 | **Post Sentimentality** | Tweets | Text | Characters, (Emojis?) | Can use an out-of-the-box tool for sentiment classification such as *Vader:* https://github.com/cjhutto/vaderSentiment <br><br> Provides a high-level measure of tweet sentiment which can be used to categorise tweets as **Positive, Negative, Neutral.** We could use this to provide an aggregated measure of general posting sentimentality. Some models also provide further categories such as anger, joy, sadness, sarcasm etc. <br><br> **RIZA:** <br><br> It would be good to try VADER; however, note that some benchmark studies have shown machine learning (ML)-based approaches to significantly outperform it, e.g., https://pmc.ncbi.nlm.nih.gov/articles/PMC10987730/. <br><br> Having said that, it is also possible to use VADER scores as features for ML; see https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0245909 | 5 | 2 | 4 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | Detecting these [anger, joy, sadness, sarcasm etc] would require emotion detection/classification models, the selection of which would depend on the particular types of emotions we might be interested in. For example, this study chose 8 emotion types which the authors then justified in the paper: https://www.frontiersin.org/journals/ psychology/articles/10.3389/fpsyg.2022.931921/full<br><br>Possible Research Questions:<br><br>Similar to toxicity, we can use this variable to quantify the type of user someone is: are they a mostly positive or negative user? Or are there more neutral in their language?<br><br>How does sentimentality change depending on the topics they are discussing?<br><br>Sentimentality could also be an important variable for potentially measuring ideological position. Sentiment towards particualr political topics can be informative of their ideological position or stance on a specific issue. | | | |
| 16 | **Post Curation** | Tweets | Edit History Tweet IDs | - | A tweet's edit history can be used as an indicator for how often a user edits their tweets, and how many times. This can give an indication of how much a user curates their timelines, if at all: | | | |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | 1. **% of Edited Tweets**<br>2. **Median no of times a tweet is edited**<br><br>This could also be combined with tweet formality to give an indication of how "professionalised" a user's timeline is. (Although this wouldn't be able to account for how many tweets they have deleted altogether). Possible Research Questions:<br><br>This variable can contribute to our understanding of how carefully a user curate their profile.<br><br>Are more contentious tweets edited more often?<br><br>Is there a correlation between the formality of a user's tweets and how often they edit their tweets? | **3** | **1** | **5** |
| 17 | **Post Openness** | Tweets | Reply Settings | - | A user's reply settings can provide a high-level view of their general openness to response from other users. There are three possible settings:<br><br>1. **% Open to Everyone**<br>2. **% Open to Mentioned Accounts Only**<br>3. **% Open to Followed Accounts Only**<br><br>Possible Research Questions:<br><br>How "closed off" is a user to replies/engagement from other users? | **2** | **1** | **5** |

| 18 | Types of Accounts Followed | Follows | Followed User ID | Follower Profile Metadata, Follower Tweet Timelines | We could train an ML model to categorise accounts into N types based on information in their profiles and possibly also their tweets. We could use NER to categorise the account types via the named entities in their profile bios (organisations/job titles/usernames) E.g: Are they political accounts, media, academic/scientist, sport etc? We could also apply a classifier to their tweets to gather information about exactly what type of account they are. Possible Research Questions:<br><br>What sorts of accounts does a user primarily follow?<br><br>This could be combined or compared with the topics they mostly tweet about/engage with to assess whether there is a discrepancy between what users like to consume vs. what they like to engage with vs. what they like to propogate. | **5** | **2** | **4** |
|---|---|---|---|---|---|---|---|---|
| 19 | Network Formality | Follows | Followed User ID | Followed Account Profile Metadata | Similar to the above, but looks to quantify the degree of "formality" to a user's followed accounts. What % of their followed accounts are verified? What are their average number of followers?<br><br>We could use these simple metrics to quantify how "formal" their following network is. I.e: do they mostly follow official accounts like media organisations and influencers or is it more casual friends and family? | **5** | **1** | **5** |

| # | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | |
| 20 | **Network Credibility (Verified Media)** | Follows | Followed User ID | Followed Account Profile Metadata | We can verify the credibility of their follow network by assessing the quality of the "official" accounts that they follow. We can ascertain official accounts using either whether they are verified or have over N number of followers. We can then assess the credibility of the verified media/public figures they follow by using an offline media bias/quality checker such as https://mediabiasfactcheck.com/<br><br>**SURVEY VALIDATION VARIABLES:**<br><br>**[uom_uspre_20]** Looking through the following list of media sources, could you tell us how often you get news and information about politics from each one?<br>+<br>-[uom_uspre_20_7]    My social media feeds | **5** | **2** | **4** |
| 21 | **Following Political Influencers** | Follows | Followed User ID | - | We could calculate the top *N* number of accounts followed by respondents in the dataset to identify prominent political influencers in the digital data and then measure how many of those are followed by a particular respondent. We could also employ a spatial mapping technique like PCA or correspondence analysis to quantify the relative distances between each user based on the similarity/dissimilarity of the accounts they follow. This would allow us to identify the users with the most varied information diets (relative to the other users). | **5** | **1** | **5** |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | How similar or disimilar is a user's followed accounts relative to other users in the dataset? | | | |
| 22 | **Influencer Status** | Profile, Tweets | **Profile:** Verification, Public Metrics<br><br>**Tweets:** Engagement Metrics | Verification, Followers Count, Listed Count<br><br>Liked Count, Retweet Count, Reply Count, Quote Count | We can use the combination of whether they are verified or not with their number of followers and the number of lists they appear in to devise a measure of influencer status (or "notability"). This could also include overall tweet count and tweet engagement statistics to get an overall measure of how they are received on the platform ("impact score").<br><br>**Possible Research Questions:**<br><br>How 'influential' is a user generally? This could be quantified on a sliding scale and potentially standardised for all users in the dataset so we could ascertain who are the most influential users and whether there are particualr characteristics in either their Twitter usage or in their demographic information that correlate with this. | **5** | **1** | **5** |

# SUBSTANTIVE VARIABLES

| No. | Variable Name | Main Table | Field Table | Field Dimensions | Description | Val. | Compl. | Rel. |
|-----|---------------|-----------|-------------|------------------|-------------|------|--------|------|
| 23 | **Ideological Position (User)** | Tweets, Follows, Likes | **Tweets:** Text, Attachments, Referenced Tweets<br><br>**Follows:** Followed User ID, Followed User Timelines<br><br>**Likes:** Liked Tweet ID, Liked Tweet Author ID | **Tweets:** Characters, Usernames, Hashtags, Emojis, URL Links, Media Keys<br><br>**Follows:** Profile Metadata, Tweet Text<br><br>**Likes:** Profile Metadata, Tweet Text | There are a number of different ways we can estimate the ideological position of a user based on their digital Twitter data. We can derive left-right positions via their **tweet timelines** e.g:<br><br>https://aclanthology.org/P17-1068/<br><br>https://collaborate.princeton.edu/en/publications/quantifying-political-leaning-from-tweets-retweets-and-retweeters<br><br>https://www.cambridge.org/core/journals/political-analysis/article/ideological-scaling-of-social-media-users-a-dynamic-lexicon-approach/0173A3145A67CB89ACFC8DE09B30C482<br><br>Or from their follow networks:<br><br>https://journals.sagepub.com/doi/abs/10.1177/0956797615594620<br><br>https://www.cambridge.org/core/journals/british-journal-of-political-science/article/estimating-ideal-points-of-british-mps-through-their-social-media-followership/1627B42FE1A547458DB1ED860CE502F1<br><br>https://github.com/pablobarbera/twitter_ideology *(Tweetscores)*<br><br>https://web.cs.ucla.edu/~yzsun/papers/2017_SBP_Ideology.pdf | **5** | **3** | **3** |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | Or from their likes/retweets, where we could combine a mixture of tweet content and network content together for estimation

**SURVEY VALIDATION VARIABLE:**

**Ideological Congruence:**

**[uom_uspre_2]** In politics people sometimes talk of 'liberals' and 'conservatives'. Using the scale from 0 to 10 below where would you place yourself on this scale, where 0 means 'very liberal' and 10 means 'very conservative'? | | | |
| **24**

**3** | **Ideological Exposure (Network)** | Follows | Followed User Timelines | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | We can estimate the ideological leaning of the content that users are exposed to based on the tweet timelines of the users that they follow. (This makes an assumption that users are evenly exposed to all the content that every one of the accounts they follow has posted.) We can estimate this via the same method as listed in the variable above (21). Additionally, we could take a more formal approach and only estimate ideology for the official accounts that they follow and their ideological leaning can be quantified via an offline media bias checker such as https://mediabiasfactcheck.com/

**SURVEY VALIDATION VARIABLE:**

**Filter Bubble Perception:**

**[uom_uspre_28]** Thinking about the accounts that you follow on Twitter we would like you to rate the similarity of their political viewpoints using a scale of 0 to 100. | **5** | **3** | **3** |

| 25 | Ideological Diversity | Follows | Followed User Timelines | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | Builds on the variable above (23) to not only consider the average ideological leaning of their informational network, but also the degree of diversity (spread). For instance, two users can have the same average ideological leaning of their follow networks but one can have a very narrow spread around the average, where others may show much higher degrees of diversity. This can be measured via a number of straightforward descriptive statistics. (This may perhaps be the more effective way to validate **filter bubble perception** than the above variable). | 5 | 3 | 3 |
|----|----|----|----|----|----|----|----|----|
| 3 | | | | | | | | |

| 26 | Populist Propagation | Tweets | Text, Attachments | Characters, Usernames, Hashtags, Emojis, URL Links, Media Keys | We can estimate the degree of populism a user propagates via the text and attachments in the tweets they post on their timelines. There a number of ways in which we can train an ML classifier to identify populist rhetoric (which I believe Riza has already begun developing). Some examples found online are: https://doi.org/10.1111/weng.12496 https://journals.sagepub.com/doi/full/10.1177/00491241221122317 https://journals.sagepub.com/doi/10.1177/1461444820976970 https://www.tandfonline.com/doi/full/10.1080/10584609.2022.2025505 | 5 | 4 | 3 |
|----|----|----|----|----|----|----|----|----|
| 3 | | | | | | | | |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | **SURVEY VALIDATION VARIABLE:** <br><br>**Populism and Trust:** <br><br>**[uom_uspre_8]** On the whole, how satisfied or dissatisfied are you with the way that democracy works in the U.S.? <br><br>**[uom_uspre_6a]** Thinking about some of the **main political institutions and actors in the U.S.**, using a scale of 0 to 10, how much trust would you say you have in each of the following to do the right thing? <br><br>**[uom_uspre_7a]** Turning now to the **current campaign and election process** and using the same scale of 0 to 10 | | | |
| 27 | **Populism Exposure** <br><br>3 | Follows | Followed User Timelines | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | Following the same methodology as the above variable (24), we can measure the degree in populism in the tweet timelines of the accounts that users follow. We can then develop a measure of populism exposure. | **5** | **4** | **3** |
| 28 | **Misinformation Propagation** <br><br>2 | Tweets | Text, Attachments | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | https://link.springer.com/article/10.1007/s00521-022-07797-y This paper explores the state of the art on misinformation classification, especially on social media. (specifically, COVID-19). I think this is something that would be incredibly useful to build but will likely be very complex. I don't know if there's a way we could be our own classification ("fact checker") model that can check the validity of tweets based on real-world information. This will require a bit more reading I think. | **5** | **5** | **2** |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | **SURVEY VALIDATION VARIABLE:**<br><br>**[uom_usddc_16]** And when you have seen content online that you think is **misleading** or not fully accurate online how have you responded to it? (RTs, Likes, QTs, Replies) | | | |
| 29 | **Misinformation Exposure**<br><br>2 | Follows | Text, Attachments | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | Again, this builds directly on the above variable () but assesses the degree of misinformation that is present in the timelines of the accounts that a user follows.<br><br>**SURVEY VALIDATION VARIABLE:**<br><br>**[uom_usddc_15]** Have you ever personally seen or come across content **online** that you thought was misleading or not fully accurate?<br><br>**[uom_usddc_14_1]** There is a lot of discussion nowadays about how **misleading and inaccurate information may be circulating online**. On a scale of 0 to 100 where 0 is not accurate at all and 100 is fully accurate how would you rate the content put out by each of the following organizations or groups of people?<br><br>**+**<br><br>**[uom_usddc_14_5]** **Members of my online social networks ** | **5** | **5** | **2** |

| 30 | Credibility of Link Sharing | Tweets | Text | URL Links | This variable is similar to the misinformation variable, but differs in that we could directly assess the credibility of a link via the domain level URL. We could then check this against a media bias checker such as https://mediabiasfactcheck.com/ to get a straightforward score of link credibility based on the site. (This may then limit this variable to only media websites, but I think this is okay) | 5 | 3 | 4 |
|----|-----------------------------|--------|------|-----------|---|---|---|---|
| 31 | Bias of Link Sharing | Tweets | Text | URL Links | Same as the above, but instead of assessing credibility, we look at the ideological slant of the site URL instead. This site contains a score of bias as well as credibility which we can use. https://mediabiasfactcheck.com/.  Perhaps the only difficulty with this is figuring out a suitable way to automate it. | 5 | 3 | 4 |
| 32 | Conspiracy Theory Propagation | Tweets | Text | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | There are a few papers that have looked at building a conspiracy theory classifier. For instance:<br><br>https://link.springer.com/article/10.1007/s42001-020-00086-5<br><br>https://www.jmir.org/2020/5/e19458<br><br>https://ceur-ws.org/Vol-3181/paper67.pdf<br><br>Most of which is focused on COVID-19. This is different from the misinformation variable in that conspiracy theories, while a form of misinformation, are much more elaborate and connected ideas. It would be interesting to see if we can | 5 | 5 | 2 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | classify how much a user propagates certain conspiracies (and which ones), and how this correlates with other important variables. We could potentially even build a simple dictionary of keywords relating to different conspiracies and classify that way (may be too simplistic).<br><br>**SURVEY VALIDATION VARIABLE:**<br><br>**[uom_uspre_16]** The following are some statements about the current COVID-19 or Coronavirus outbreak. Please indicate if you think they are true or false. | | | |
| 33 | **Conspiracy Theory Exposure** | Follows | Text, Attachments | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | Same as the above variable, how often do the accounts that users follow post conspiracy theory information? | 5 | 5 | 2 |
| 34 | **Online Political Attention** | Follows | Profile Metadata Text, Attachments | **Profile:** Username, Display Name, Description<br><br>**Text:** Characters, Usernames, Hashtags, | We can construct a variable that approximates political attention based on a combination of content that users engage with online along with the accounts they follow. If we employ classifier models to classify types of posts and types of accounts as "political", we can derive a proxy measure of how much overall attention a user pays to politics, which we can then use to validate their response to the survey question relating to political attention. | 5 | 2 | 4 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | Emojis, URL Links, Attachments | (Although I am unsure whether this would be a valid measure of political attention. Maybe this would be political "exposure" instead? Unless we just use the content that they explicitly engage with.)<br><br>**SURVEY VALIDATION VARIABLES:**<br><br>**[uom_uspre_5]** Based on a scale of 0 to 10 (where 0 equals no attention and 10 is pay a great deal of attention), how much attention do you generally pay to politics? | | | |
| **35**<br><br>**4** | **Online Political Discussion** | Tweets | Text | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | We can ascertain whether or not a respondent has engaged in political discussion with someone online by first identifying posts which directly reference another user (reply/quote tweet/mention) and then categorise whether any of these tweets are in any way political. Using this data, we can then validate whether or not a respondent has in fact discussed politics online or not.<br><br>**SURVEY VALIDATION VARIABLES:**<br><br>**[uom_uspre_3b]** And again, during the past 12 months, have you done any of the following?<br><br>**+**<br><br>**[uom_uspre_4_5]**    Gotten into a political debate or discussed politics with someone<br><br>**+** | **5** | **2** | **4** |

| | | | | | **<1>** Online<br><br>**[uom_uspre_27]** And how often, if ever, do you respond to messages or content about politics that you see on Twitter. This could include retweeting, replying, messaging or liking it. | | | |
|---|---|---|---|---|---|---|---|---|
| **36**<br><br>**3** | **Political Campaign Exposure**<br><br>**4** | Follows | Text, Attachments | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | We could quantify how much campaign content they have been exposed to. We can derive this from the timelines of the accounts that users follow and establish a way to classify campaign specific content (?).<br><br>**SURVEY VALIDATION VARIABLE:**<br><br>**[uom_usddc_11]** For each of the following types of content could you tell us how often you recall seeing it when you have been online in the _past month_?<br><br>This may not be perfectly accurate as it talks about all online exposure. We may want to only use the answers to –<br><br>**[uom_usddc_11_4]**    Non-sponsored content about the election or political issues posted by people or organizations I don't know personally, but that I follow on like social media | **5** | **3** | **4** |
| **37** | **Political Campaign Engagement** | Tweets, Likes | Referenced Tweet | Type, Text | Similar to the above, if we can classify tweets as campaign focused tweets, we can then measure the degree to which a | **5** | **3** | **4** |

| | | | | | user has engaged with these tweets (RT, QT, Reply, Like). This will allow us to validate the following variable:<br><br>**SURVEY VALIDATION VARIABLE:**<br><br>**[uom_usddc_12]** Thinking back specifically to the campaign adverts from candidates, parties and other groups that you have seen online in the past month can you tell us a bit more about how you have responded or engaged with them, if at all? | | | |
| 38<br><br>2<br><br>4 | **Level of Anonymity** | Profile, Tweets, Follows | **Profile:** Username, Display Name, Description, Location, Profile Image, Cover Image<br><br>**Tweets:** Text, Attachments, Location<br><br>**Follows:** Followed User ID | **Profile:** Characters, Usernames, Hashtags, Emojis, URL Links,<br><br>**Tweets:** Characters, Usernames, Hashtags, Emojis, URL Links, Attachments<br><br>**Follows:** Profile Metadata | Is their display name an actual name or a pseudonym? Do they use an actual location in the location field or something else? We could potentially also incorporate analysis of their profile picture and cover image. Are these human faces or anonymous pictures? Do they specify where they work or anything else in their descriptions that could identify them? Do they give away anything in their tweets that could potentially identify them? We could also consider analysing their followed accounts to see if they follow anyone who could be used to identify them. Level of anonymity could be measured on a sliding scale: 0 (Identifiable) to 1 (completely anonymous).<br>**Possible Research Questions:**<br><br>Which tyes of users are more likely to anonymise their accounts?<br><br>An interesting validation variable might be how this lines up with their attitudes towards digital privacy. Do people protect their data/anonymity as much as they claim to? | **5** | **2** | **4** |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | Are users who tweet about more contentious subjects or use more toxic/negative lanuage more likely to anonymise their profiles?<br><br>**SURVEY VALIDATION VARIABLES:**<br><br>**Privacy Paradox:**<br><br>[uom_uspre_25x] - Some people prefer their personal details are never collected or shared when they are online while other people do not mind if these details are used to help personalize the content they see. Using the following scale could you please indicate your preference between these two views:<br><br>[uom_usddc_23] - There are a number of things people can do to help protect their privacy online. Before taking this survey, were you aware of, or had you done any of the following? | | | |
| **39**<br><br>**2** | **Content/Person Endorsement**<br><br>**4** | Tweets | Text, Attachments | Characters, Usernames, Hashtags, Emojis, URL Links, Attachments | This variable is broad. We can validate content or person endorsement based on the sentimentality of tweets they have posted about the particular issue or person. The most obvious ones from the survey are candidate endorsement, so we could validate their sentiment towards Biden/Trump against the sentimentality of the tweets they have posted about them.<br><br>**SURVEY VALIDATION VARIABLES:**<br><br>**[uom_uspre_11_1 and 2]** When it comes to the candidates in the U.S. Presidential race we would like to know how positively or negatively you feel toward them on a scale of zero to 100. | **5** | **2** | **4** |

| | | | | | Could potentially also validate the following:<br><br>**[uom_uspre_19]** How would you rate President Trump's overall handling of the Covid-19 / Coronavirus outbreak in the U.S. to date?<br><br>**[uom_uspre_8]** On the whole, how satisfied or dissatisfied are you with the way that democracy works in the U.S.? | | | |