

The IMPED Model of Information Quality

Marco Bastos @toledobastos – *University College Dublin & City, University of London*

Shawn Walker @walkeroh – *School of Social and Behavioral Sciences, Arizona State University*

Michael Simeone @michael_simeone – *ASU University Library, Arizona State University*

Abstract

This paper introduces a model for detecting low-quality information we refer to as the Index of Measured-diversity, Partisan-certainty, Ephemerality, and Domain (IMPED). The model purports that low-quality information is characterized by ephemerality, as opposed to quality content that is designed for permanence. The IMPED model leverages linguistic and temporal patterns in the content of social media messages and linked webpages to estimate a parametric survival model and the likelihood the content will be removed from the Internet. We review the limitations of current approaches for the detection of problematic content, including misinformation and false news, which are largely based on fact-checking and machine learning, and detail the requirements for a successful implementation of the IMPED model. The paper concludes with a review of examples taken from the 2018 election cycle and the performance of the model in identifying low-quality information as a proxy for problematic content.

Keywords: content moderation; diversity index; partisanship; misinformation; web archive

Introduction

The weaponization of social media platforms for misinformation campaigns pose a fundamental challenge to collective discussion based on a consensual and fact-based reality. Successful misinformation campaigns take advantage of the biases (Comor, 2001; Innis, 2008) intrinsic to social media platforms, particularly the underlying attention economy and its supply chain of viral content (Jenkins et al., 2012). These tactics are leveraged to polarize and alienate voters from the electoral process (Benkler et al., 2018), a strategy that may take the form of misinformation—unintentional behavior that is inadvertently misleading or inaccurate (Karlova & Fisher, 2013)—or disinformation, the intentional distribution of fabricated stories to advance political goals (Bennett & Livingston, 2018).

Social media platforms have implemented a range of measures to identify false amplification (Weedon et al., 2017) and remove “fake accounts” seeding problematic content (Twitter, 2018b), an operational term employed to describe accounts, posts, and links to content selected for removal. Similarly, we rely on the purposely broad term “low-quality information” as an umbrella concept that refers to misinformation (and subcategories such as disinformation), influence operations that mimic the appearance of news outlets, false or fabricated news items, and user-generated hyperpartisan news—i.e., polarized narratives reinforcing partisan identity. The style, messaging, lifespan, and provenance of low-quality information and hyperpartisan content can vary depending on the strategies employed to sow division. It is against this challenging backdrop that social platforms have implemented community standards that define “problematic content” and the enforcement of such policies (Facebook, 2018a, 2018b; Weedon et al., 2017).

We define low-quality information as textual or audio-visual content with the potential to interfere with online deliberation by deceiving, confusing, or misinforming. This working definition is agnostic to the veracity or falsehood of statements, drawing instead from the field of risk innovation to identify threats to values within social and organizational contexts (Maynard, 2015). While social platforms define problematic content as information selected for removal, our framework deems information to be of low-quality whenever a repertoire of tactics is employed to hinder online deliberation, including fast access to information, opportunities for self-publishing, and symmetrical conversations among users (Halpern & Gibbs, 2013).

In summary, low-quality information refers to high-risk, divisive content that blurs the lines among misinformation, disinformation, and propaganda, but also user-generated and tabloid-like content marked by sensationalized language and speculative reasoning (Bastos, 2016). This content is often deleted or blocked from social platforms, only to resurface again through sockpuppet or surrogate accounts. We seek to contribute to this multidimensional problem by proposing the Index of Measured-diversity, Partisan-certainty, Ephemerality, and Domain (IMPED), a model that identifies low-quality, partisan information on social media at scale and in near real- time. The model employs automated text analysis and web archiving to gauge whether online information fulfills the parameters of quality content.

Previous work

Control mechanisms

Historically, the challenges posed by modern propaganda were partially addressed with the consolidation of highly regulated networks that prevented mechanisms for sharing information. Individuals who wanted to share content faced high costs of production and distribution, as with

print, or contended for access to scarce resources, such as the electromagnetic spectrum required for broadcasting television and radio (Bastos et al., 2013). Gatekeeping encapsulated this control-communication mechanism based upon sender–receiver roles, and the information that passed from sender to receiver had a source–destination direction predefined by the distribution system (Barzilai-Nahon, 2008, 2009).

The concept of gatekeeping as a progressive control system provided the perfect metaphor to an information ecosystem supported by mass media. The gatekeeper became a powerful image describing a newspaper editor, who determined which inputs were “newsworthy” and controlled the flow of reported events. The emergence of digital, decentralized networks invited scholars to review the concept of gatekeeping, but it continued to refer to a process of selection with control points between the channel and its external environment (Shoemaker & Vos, 2009; Welbers & Opgenhaffen, 2018). In other words, it remained rooted in a centralized network topology characterized by a bottleneck of interconnections determining the flow of information.

While gatekeeping is historically tied to broadcasting and highly centralized networks that prevent mechanisms for sharing information horizontally, social media platforms allow peripheral users to seed viral messages. Indeed, the relative centrality of users in the network is not a condition for triggering information cascades (González-Bailón et al., 2013). Committed minorities can be resilient to influence and may reverse a prevailing majority opinion in the population by consistently proselytizing the opposing opinion (Xie et al., 2011). Social media platforms compound these effects by utilizing a social infrastructure where users receive information from various sources that crowd out the limited content from traditional news outlets (Wu et al., 2011). By forging an information infrastructure that bypasses editorial gatekeepers,

social platforms simultaneously provide exposure to opinions that have been traditionally suppressed and allow for large-scale, inexpensive, and horizontal distribution of misinformation.

The decentralized network topology of social media platforms is not the only cause for the prevalent circulation of low-quality information. Low levels of trust in news organizations, particularly in the United States and parts of Europe (Newman et al., 2016), along with broader societal shifts associated with low-trust societies (Fukuyama, 1995), are largely at odds with the consensus enforced by broadcasting networks. Though social media platforms are not sufficient cause for democratic destabilization or the fragmentation of public discourse (Benkler et al., 2018), they have accelerated a process in which traditional institutions are increasingly seen as less trustworthy (Zuckerman, 2017).

Fact-checking

It is against this backdrop of polarization and hyperpartisanship that low-quality information flourishes, with politically skewed misinformation resurfacing in the current debate about “fake news” (Lazer et al., 2018). Literature dedicated to identifying and removing mis/disinformation from social media platforms has relied largely on fact-checking (Vosoughi et al., 2018). Fact-checking is often posited as the diametrical opposite of misinformation, providing evidence to rebut the inaccuracies advanced to mislead individuals (Jiang & Wilson, 2018). While we accept that linguistic and rhetorical devices (including emotional and topical markers) may be used to detect the incidence of misinformation, we contend that fact-checking inevitably fails to address the problem, both operationally and epistemologically.

Fact-checking fails operationally because: a) it is time-consuming; b) it is available with significant delay after misinformation cascades are triggered; and c) it focuses on popular

content, ignoring a wealth of deleted content that constitutes the long tail of misinformation. This approach is particularly unsuitable to the problem given the relative short shelf-life of misinformation campaigns on social media platforms, where considerable damage may occur before the fact-checking process can even begin. Other serious shortcomings of fact-checking and rumor correction are that they are shared at a much slower rate and their penetration is considerably lower than the original story (Arif et al., 2017; Starbird et al., 2014).

A framework relying on fact-checking is also epistemologically objectionable because the tacit presupposition is that there cannot be genuine political debate about facts, which are assumed to be unambiguous and not subject to interpretation, selection and confirmation bias, or ideological coherence (Uscinski & Butler, 2013). Indeed, Marietta and Barker (2019) found that fact-checking fails to ameliorate polarization. The partisan divide in trusting fact checks is particularly pronounced in those with strong commitments to their values, who are more certain than others that their perceptions are correct and rarely read fact-checks. In other words, while the literature on detecting false news draws from a framework in which facts are opposed to misinformation, we understand that both can coexist and that influence operations often layers true information with false (Starbird, 2019), thereby exploiting different interpretations, superinterpretations, and oftentimes manipulation of otherwise objective facts.

Ephemeral hyperpartisanship

In contrast, we have found that low-quality information circulating on social media is marked by a short shelf life, which makes it difficult to retrospectively rebuild the public conversation (Bastos & Mercea, 2019). This problem is compounded by the extensive use of images on social platforms, which are associated with information cascades (Cheng et al., 2014; Dow et al., 2013),

and the Terms of Service requiring content deleted by a user be removed from social platforms (Twitter, 2018a). These policies enable the disappearance of messages and URLs from the public view and prevent research on misinformation campaigns. Public web archives, such as the Internet Archive, rarely contain records of specific tweets, their attached images, and linked content. As a result, it is often impossible to determine what the original social media post and associated image conveyed at the time of posting, or how far afield it cascaded.

User-generated partisan content, often drafted in support of contentious political issues, is remarkably ephemeral, disappearing or significantly changing shortly after being posted.

Previous research on the Brexit referendum campaign has found that a significant share of the URLs posted on Twitter disappeared after the ballot (Bastos & Mercea, 2019). The URLs, which could not be resolved after the referendum, either linked to a Twitter account that had been removed, blocked, or deleted or to a webpage that no longer existed. Nearly one third (29%) of the URLs monitored in the study linked to multimedia content, such as Twitter statuses and pictures, that was no longer available and whose original posting account had also been deleted.

We contribute to this literature by proposing a scalable model that identifies low-quality information in near real-time. The model leverages the temporal patterns of posting activity and embedded webpage content, particularly content modification and deletion, along with rhetorical and linguistic devices associated with low-quality information. The model requires real-time data collection and archiving and employs methods from information science, social network analysis, and text analysis. The model is benchmarked against quality information, thereby flagging false news, a subset of content often referred to with the contested and ideologically inflected notion of “fake news,” but also misinformation, or content that is inaccurate, but not necessarily spread with a political agenda. To a lesser extent, the model may also flag

disinformation, or the deliberate spread of inaccurate content. Disinformation campaigns may however be carefully crafted to meet the requirements of quality content, thereby limiting the predictive power of the model.

IMPED: a probabilistic model for detection of low-quality information

Rationale

The IMPED model assumes that low-quality information shares a repertoire of features that allow for classifying such content at scale. No single similarity measure or training dataset can account for the variety of approaches available to misinformation outlets, but we expect misinformation campaigns to source content of quantifiably lower quality compared with mainstream media articles, which are labor intensive, likely reviewed by an editorial gatekeeper, and therefore more expensive. Accordingly, instead of quantifying misinformation with a predefined set of metrics, we address this challenge by parametrizing common traits of high-quality content. We rely on these metrics to reverse engineer misinformation content based on stylistic and temporal qualities resulting from reduced gatekeeping structures and practices.

The IMPED model is distinct from techniques that train machine learning classifiers to detect false news based on ground truth “real” news and ground truth “false” news. While the latter helps to make distinctions based on stylistic similarities of documents based on terms, term combinations, sentence length, punctuations, or parts of speech, the IMPED model is based on parametrizing information quality rooted in a range of communication and information theories. In other words, our approach deviates from machine learning algorithms based on predictive analytics and is etiologically similar to statistical models based on probability distributions. In

summary, the IMPED model is entrenched in theory, relies on finite samples, is based on deduction, and is assessed in terms of confidence intervals (Breiman, 2001).

We contend that political misinformation has four identifiers. Firstly, ephemerality is a trait broadly observed across misinformation sources. While quality content requires extensive editorial work at various levels of the news industry and is designed for durability, low-quality content is produced in bulk for quick turnover and is characterized by a short shelf life. Secondly, political misinformation outlets exploit the pre-existing audience biases through partisan content, with a lesser focus on building trust through conventional mechanisms. Thirdly, we contend that the fast-paced production of low-quality information presents linguistic signatures that sits in opposition to quality content. Fourth, current misinformation tactics are designed to hijack information literacy training by impersonating sources and professionally designed websites to legitimize low-quality information. While this content is recurrently taken down by web hosting services and social platforms, such campaigns explicitly play a numbers game by repeatedly reposting the content. In the next section we parameterize this repertoire.

Key metrics

A key challenge in detecting and studying hyperpartisan news is the relatively high level of ephemerality and short shelf life of these articles. With current research focusing on the impact, reach, and spread of misinformation (Jiang & Wilson, 2018; Lazer et al., 2018; Vosoughi et al., 2018), little attention has been given to their remarkably short lifespan. As a result, we lack a detailed catalogue of metrics and high-fidelity real-time datasets needed to understand how problematic content is operationalized and shared before disappearing from social media platforms. Indeed, despite multiple efforts to understand the diffusion of problematic content,

there is scant literature on how user-generated partisan content is drafted to support political issues, subsequently disappearing shortly after being posted.

The incidence of ephemeral and partisan content is not restricted to contentious issues such as the 2016 U.S. presidential election or the 2016 U.K. E.U membership referendum. After monitoring national elections in 2018, we found that over 6% of user-generated content disappeared after the ballot. Content that tends to disappear from Twitter is largely hosted by social media platforms and content curation services. Other than twitter.com, the most common domains include bit.ly,youtu.be, goo.gl, instagram.com, facebook.com, Breitbart.com, and paper.li, with nearly half of the content published using paper.li disappearing shortly after it was made public. Paper.li is only one of many services that generate a “professional looking newspaper” from user-generated content, another marker of low-quality content as misinformation sources often model their website layouts after established news outlets.

Given the above, we posit that false and hyperpartisan content can be modeled and identified based on a catalogue of metrics. We expect content sourced from mainstream media to be less likely to include misinformation compared with user-generated content. We also assume that low-quality information is considerably cheaper to produce and more likely to be recycled or removed altogether from the internet. We therefore expect a key predictor of the model to be the time elapsed from the date the article was published to the date the article is significantly modified or ultimately removed. The ephemerality score of content can also be aggregated up to the domain level, thereby providing another temporal metric for identifying unstable content across an entire website. Partisan certainty is naturally an element of hyperpartisan news, and we expect it to be another predictor of low-quality information. Lastly, the content of the social

media post and that of the webpage(s) linked in the post can be analyzed for measurements of diversity and rhetorical devices.

The IMPED model

The metrics catalogued above allow us to parametrize a misinformation threshold as a combination of five scores based on linguistic diversity, partisan certainty, user or mainstream-generated content, ephemerality index, and domain stability. The misinformation threshold indicates the optimal point after which we expect the seeding of such content to impede the flow of quality information. These metrics are formalized as: a) measured-diversity score; b) partisan-certainty score; c) source score; d) domain stability score; and e) ephemerality score, that is, the difference in the URL content between the time it was referenced in a social media post and the time the content was significantly changed or became inaccessible. These metrics are leveraged to identify user-generated information that is likely to meet the classification of low-quality content. Further versions of the model could also incorporate properties based on user and post metadata as well as network metrics associated with the message diffusion.

We refer to this model as IMPED (Index of Measured-diversity, Partisan-certainty, Ephemerality, and Domain), in which the elements of linguistic variance, partisan certainty, media source, ephemerality, and domain stability are used to estimate a parametric survival model where the likelihood of content disappearing is an approximation to low-quality information. The unit of analysis of the IMPED model is the combination of the social media post (e.g., tweet or Facebook post) and the content of URLs embedded in the post (if available). We refer to this unit of analysis as the *u-content*. Implementation of the model requires real-time archiving of URLs posted on social media platforms at the time of posting and the archiving of

the social media post for analysis of the *u-content*. The IMPED model relies on five key scores to estimate the likelihood the *u-content* is problematic, detailed in Table 1 below.

TABLE 1 HERE

First, the model relies on the Pielou's and Shannon diversity indices (*s-score*) to estimate the linguistic diversity of *u-content* as a function of both vocabulary richness and evenness (Pielou, 1966; Shannon & Weaver, 1962). Pielou's species evenness and the Shannon-Wiener diversity index are commonly employed in the natural sciences to measure the level of complexity of a community structure, particularly diversity and patterns of distribution for species in ecosystems. Shannon's measure of information content also identifies the "surprisal" of rare messages, which are expected to be more informative compared with ordinary messages. In other words, the unit of information $I(w_n)$ associated with outcome w_n with probability $P(w_n)$ is formalized as $I(w_n) = -\log(P(w_n)) = \log(\frac{1}{P(w_n)})$. In our field of inquiry, Brugnoli et al. (2019) explored the association between lexical entrainment, or the convergence of linguistic properties, and group polarization. Correspondingly, the *s-score* is an index that characterizes linguistic diversity in the ecosystem of *u-content*, with the expectation that low-quality information will be less diverse compared with quality content which is marked by accuracy, grammar and spelling consistency, and a richer vocabulary (Lijffijt et al., 2016).

Second, we estimate partisan certainty (*p-score*) based on linguistic signals employed in the compositions, such as the absence of subjunctive mood to express uncertainty (Jiang & Wilson, 2018), specific rhetorical indicators of certainty (Kuklinski et al., 2000), and affective validation (Rucker et al., 2014). We posit that the predominance of nouns and the paucity of adjectives foreground less-nuanced content expressing partisan certainty, whereas quality information contrasts established facts and conditional possibilities. This metric draws from

scholarship identifying common expressions in hate speech, false news, and texts associated with sex, death, and anxiety, as opposed to words related to work, business, and the economy (Pérez-Rosas et al., 2017). While partisanship is typically associated with bias or ideological alignment, the *p-score* is indifferent to political leanings and focuses instead on rhetorical features that express certainty, speculation, and stylistic qualities of text data, such as the mood of verbs and use of conditional phrases or future tense. The *p-score* calculation sums the ratio of conditional words to all tokens with the ratio of nouns to adjectives. Because the use of nouns or adjectives is not assumed to be interdependent with subjunctive mood verbs, we sum the ratios to measure the overall certainty of texts.

Third, we classify the *u-content* as user-generated or produced by mainstream news outlets (*u-score*) based on the comScore classification of media sites under the category “information/news” with a minimum percentage reach of .01% (comScore, 2013a, 2013b). The *u-score* is therefore the simplest classifier in the model and is limited to scoring the source as mainstream or user-generated, a necessary control mechanism that regulates the results of the *e-score* for news outlets publishing dynamic content. As such, the *u-score* is instrumental in normalizing the *e-score* for news outlets constantly updating their stories, providing live coverage to events, or publishing additional information and corrections that would otherwise trigger the *e-score*.

Fourth, we estimate the ephemerality score of the content (*e-score*) as a continuum of ephemerality level (Walker, 2015) in which a zero score is content that has not been changed (i.e., it is stable) since posting, a score of .01–.99 indicates content that has been altered, and 1 (the maximum *e-score* possible) indicates the content is no longer accessible. Naturally, the extent to which social platforms succeed at identifying and taking down problematic content is a

potential confounding factor of the *e-score*, so it is important to calculate the *e-score* on the *u-content* instead of the social media post alone. Fifth, and in close connection to the above, we calculate the domain stability score (*d-score*) based on the average ephemerality score of URLs hosted by the domain sourcing the weblink embedded in the *u-content*. The *d-score* thus updates a stable list of blacklisted and domain profile information and represents the time-weighted, aggregate ephemerality score for a given domain name.

These five scores constitute the value space where the IMPED model is parametrized, with the *s-score* and *p-score* calculated by processing the combined corpus of social media post (e.g., tweet) and webpage content (when available), and the *e-score* and *d-score* calculated based on whether the post was edited or removed and the aggregate score for the domain sourcing the webpage. Figure 1 shows the five scores used to parametrize the IMPED model in relation to *u-content* and the misinformation threshold, which is the optimal operating point of this linear classifier.

FIGURE 1 HERE

Model Implementation

General considerations

Implementation of the IMPED model requires social media posts and webpage(s) linked in the post to be archived and parsed in real-time, thus creating baseline copies of content and images embedded in the message. In addition to archiving web pages on a rolling basis, it is also necessary to perform regular checks on the URLs to ascertain whether the resource is no longer accessible (URL decay). Without these steps, it is not possible to calculate the *e-score* and the *d-score*. Once the content has been archived, the Pielou-Shannon diversity index and partisan-certainty scores (*s-score* and *p-score*, respectively) can be calculated on the *u-content*. The *u-*

score, based on the URL domain, is the only metric that can be calculated prior to archiving, but it may be necessary to resolve shortened URLs to identify user-generated and mainstream media sources.

Our preliminary implementation of the model explored weblinks that circulated on Twitter during the U.S. gubernatorial, House of Representatives, and Senate elections; the Irish presidential election and abortion referendum; the Italian, Brazilian, Mexican, Egyptian, and Russian general elections; and the United Kingdom, German, and Swedish local elections. We archived content from Twitter Streaming API and took snapshots of images and URLs embedded in the tweet. While the IMPED model can also be implemented with previously collected data, the archiving of URLs needs to commence as soon as a social media post is created. This is because the ephemerality of social media posts will prevent dynamic forensic analysis on items that disappeared or significantly changed, a subset to which no baseline is available, and thus no *s-score* and *d-score* can be calculated. In summary, while implementation of the IMPED model is relatively simple, it requires the real-time capture of social media posts and their embedded images and webpages.

Examples

The five examples detailed in Table 2 were archived in the run-up to the 2018 Elections in the U.S. Two of the five original tweets are no longer available (accounts suspended), and all five original URLs embedded in the tweets are no longer retrievable. As such, the *e-score* calculated from the *u-content* for these examples varies from .75 to 1, which is the maximum ephemerality score in the model. Sample 01 is sourced from msn.com, which is a relatively authoritative source of content, but the remainder of the sampled tweets were sourced from user-generated

sources, including michaelsnyderforidaho.com, conspiracyoutpost.com, and news-info.net. The *u-score* of these items is necessarily high and offers another indication that the content is likely of low quality.

The aggregate score calculated for these websites shows that content hosted by these domains is significantly unstable and likely to be altered or deleted, therefore prompting a high *d-score*. The content of the webpages, as shown in Table 2, is marked by stylistic devices and vague statements common to daily communication, in contrast to curated content sourced from mainstream news outlets. Similarly, the content featured in the text body of the tweets often includes several hashtags, a marker of emphatic partisan loyalty. These textual features and rhetorical devices would trigger a high *s-score* and, similarly, a positive *p-score*, therefore placing the five examples as likely sources of low-quality information, which we posit is a proxy to problematic content.

Consistent with the central assumption driving the IMPED model, the content of these posts often remains available on Twitter even after the original seeding account is blocked or the webpage sourcing the content is removed. This is because deleted content resurfaces via other accounts that repost the original webpage on other similarly hyperpartisan websites. For example, the original tweet ID 991023408816250880 in Table 2 featured the headline “Outrageously Unreasonable Arizona Teachers Strike Is Illegal.” The tweet is no longer available on Twitter, nor is the post from conspiracyoutpost.com identified with the shortened URL t.co/jQWtQmYcPW. However, at the time of this April 2018 post, five tweets featured the same headline and directing to various hyperpartisan websites, including “Grumpy Opinions” at grumpyelder.com and “Moonbattery” at moonbattery.com. As shown in Figure 2, and despite the original *u-content* being no longer available, the original hyperpartisan content continues to live

on Twitter through a process of continuous reposting, recycling, and resourcing to other sister accounts and URL domains.

FIGURE 2 HERE

These examples were drawn from a database of 13,770,019 unique webpages tweeted in the period leading up to 2018 elections, in which 6.3% (869,053) of the webpages disappeared after the election cycle. As detailed above, we expect the deleted webpages to have generated several sister URL webpages, thereby placing the ecosystem of low-quality information at several million webpages that continue to live on social media platforms through reposting and resharing.

Validation

The *u-score* and *d-score* are simple measures based on lists of domains. The former is based on the comScore ranking and the latter on the cumulative *e-score* of domain names that updates a list of blacklisted domains. The remainder components of the model (*s-score*, *p-score*, and *e-score*), however, require greater justification for those seeking to implement the model. We validated these scores using a test dataset of the 2018 elections comprising 37,052 tweets posted by 27,213 unique users, a subset of content for which we have a snapshot of the webpages embedded to the tweet. The test dataset includes 37,052 unique URLs, along with a high-fidelity snapshot of the original content, posted between January-April 2018 across eight national and regional elections in Brazil, USA, Ireland, UK, Italy, Russia, Mexico, and Egypt, though the majority of the content is associated with the United States elections which account for 94% of the collected data in the test dataset. Seven percent of the posted weblinks in this test dataset are

no longer available, a deletion rate similar to that observed across the complete database and therefore suitable for testing the assumptions of the IMPED model.

We proceed by fitting a series of models in which *s-score*, *p-score*, and *e-score* are variables predicting whether the embedded URL will be removed and, ultimately, whether the user account posting the content will be blocked, deleted, or suspended from social platforms. While the former is a key metric of our model, the latter meets the working definition of “problematic content” employed by social media platforms and our working definition of “low-quality information:” content that blurs the lines among misinformation, disinformation, and propaganda and that is likely to be taken down by social platforms and web hosting services.

We relied on two datasets to parametrize the *s-score*. Firstly, we calculate the linguistic diversity of web pages embedded to posts by benchmarking the observed Pielou-Shannon diversity index against Google’s Trillion Word Corpus (Fletcher, 2012), sampled to the 250K most common words in English. Secondly, we relied on a dataset of the 250K most common words tweeted by 24.7M Twitter Verified Accounts (Baumgartner, 2019) to benchmark individual Twitter accounts in the test dataset. Over 56% of the terms in Google’s Trillion Word Corpus also appear in the Twitter Verified Account Corpus, and 70% of the latter appear in the former. We further benchmarked the results against a dataset including two weeks of headlines and blurbs published by The Guardian and The New York Times (Bastos, 2015) and influential Twitter accounts.

The *s-score* algorithm first calculates the number of characters per word and the number of words employed in a sample of text. Words are not stemmed, but punctuation is removed. The algorithm returns the mean and variance results along with the “EvennessJ” and Gini scores (Pielou, 1966), with the latter increasing as the number of words in the text corpus rises. As Gini

is significantly more sensitive to the size of the text corpora compared with EvennessJ ($r=.578$ and $r=.003$, $p<.001$ and $p=.985$, respectively), the central metric of the *s-score* is the evenness of words we refer to as EvennessJ, a parameter that is less sensitive to the size of the corpus, but which requires a minimum of 100 words to perform the calculations reliably. We parametrized the test dataset against the benchmark sets and established a .90 point-threshold for the *s-score*, after which we expect the content to be indicative of low-quality information.

Validation of the *s-score* shows that low-J words' decay curves stay lower longer and that high-J words decrease quickly. This is consistent with the fundamentals of the *s-score*: more specialized web domains have more evenness of words. The content is less centralized, as news outlets strive to produce diverse content instead of republishing the same story repeatedly. Similarly, high-quality content, measured by established sources in our benchmark dataset, often focus on a predefined set of issues that define their coverage and topical focus. This is reflected in Pielou's evenness index ($J, \frac{H}{H_{max}}$) and the larger number of unique words in the samples. Inversely, tweets in our sample dataset, particularly those from accounts that disappeared and that sourced hyperpartisan websites, are more centralized: posts tend to repeat the talking points of a given partisan alliance that then echoes that same content. The results are therefore consistent with our hypothesis, with user-generated websites being more likely to score above the .90 threshold established by the EvennessJ of our benchmark datasets.

We completed the validation of the *s-score* by performing a user-lookup on Twitter API to identify user accounts that have been removed, blocked, or suspended since the tweet was posted. We assume this cohort to be probable sources of "problematic content," as Twitter has removed this set of accounts for violations to their Terms of Service. Using the *s-score* as independent variable accounting for tweet deletion, but also controlling for deleted content and

short tweets where the *s-score* cannot be reliably calculated, a significant model was observed [$F(2033)=37048$; $p<2.2e-16$; $R^2_{adj}=.14$]. While only one-seventh of the variance in account deletion is explained by the model, it is substantial for an endogenous construct such as the *s-score*. These results are more clearly observed when analyzing the *s-score* of larger population of users. Indeed, the highest *s-score* measured from tweets of a random sample of real-world users ($N=2000$) is .78, which is the lowest *s-score* for users in the test dataset. The score of .78 is also considerably higher than the benchmark established by measuring term diversity in verified accounts and the other datasets employed to benchmark the model.

The *s-score* performance is considerably improved on *u-content* that includes webpage articles due to the larger text corpora available. In our tests, the misinformation threshold remained at .90, with The New York Times and The Guardian blurbs (standfirst) scoring a relatively low .83 and .87, respectively. In comparison, a sample of 48 full articles retrieved from Breitbart scored on average .95. Even when running the algorithm on the entire Breitbart corpus of 48 articles, comprising 3809 words amassed into one single document, Breitbart content continues to score above .90 ($x=.9027231$), which is a substantial departure from content sourced from The New York Times and The Guardian. Figure 3a unpacks the parameters of the *s-score* and shows the results for the benchmark datasets and the Twitter account of Donald Trump. It also includes the five most prolific accounts in our test dataset. This cohort of accounts scored .90 or higher for *s-score* and only @nuuzfeed has not been blocked by Twitter at the time of this writing.

FIGURE 3 HERE

Validation of the *p-score* proved challenging, as the classifier requires larger sets of text to perform reliably. The sample dataset includes several partisan tweets that nonetheless scored low

on the *p-score* scale. These tweets include almost exclusively hashtags, often many hashtags, thereby lacking sentence structure with conditional verbs that could be processed by the POS tagger, an indication that the algorithm may have to be tailored to the particularities of social platforms. We addressed this problem by testing the *p-score* against Andrew Thompson’s All the News dataset, which contains 143,000 news headlines from 15 U.S. news outlets (Thompson, 2017) and the “Fake News in the 2016 Election” dataset (Allcott & Gentzkow, 2017), which includes 150 labeled “fake news” headlines according to the codebook published alongside (data access was provided through ICPSR, the Inter-university Consortium for Political and Social Research). The classifier assigns a high *p-score* to titles containing many nouns but few adjectives as well as a great deal of conditional language.

Consequently, headlines and text corpora with more descriptors and a lower incidence of conditional verbs present a lower *p-score*, which we found to be associated with partisan sources. Figure 3b shows the *p-score* of news headlines taken from the “Fake News in the 2016 Election” dataset of 1800 headlines and the 150 headlines included in the ICPSR data. While the ranking shown in Figure 3b may not represent a consensus understanding of the relative partisanship of each publications, it does show the *p-score* potential in identifying partisan-driven certainty as a communication style, with parameters that are easy to calculate when evaluating the model. This approach is in no way the only available method for measuring (un)certainty in language (Isenegger et al., 2019; Rubin et al., 2006), but the current implementation shows promise while also being relatively straightforward. The datasets the *p-score* was tested against indicate a normalized .40 point-threshold, after which we expect the content to include distinctive markers of partisan certainty. The threshold can be observed as a simple cut-off point; no conventional

news outlets surpassed this mark. Subsequent versions of the model should aim for a more nuanced heuristic for rating and weighting partisanship.

Lastly, we sought to validate the *e-score* by looking into temporal anomalies in the rate to which webpages embedded to tweets became inaccessible and the association between web content deletion rate and Twitter account termination or suspension. Unfortunately we do not have temporal information about user account suspension, but there are substantial anomalies in the deletion of web content and a significant correlation between tweet deletion and URL decay was observed [$r=.29$; $n=37050$, $p<2.2e-16$]. Indeed, the universe of 2733 webpages in our test dataset that we are no longer available were disproportionately tweeted by accounts that since then have been blocked, deleted, or suspended by Twitter. In other words, while this cohort of blocked and deleted accounts is relatively small, they posted 42% of the links to webpages and websites that are no longer available.

Conclusion

The assumption underlying the IMPED model is that the editorial and curation processes required to produce quality information can be parametrized. The model quantifies the extent to which content has been edited, curated, or reviewed for quality, accuracy, and persistence. The absence of such indicators is used to benchmark low-quality information as a proxy for problematic content. These parameters allow us to apply the model to a set of content that is broader and more diverse than news articles, despite the model foregrounding news values and linguistic features. The reliance on the news production framework and syntax is deliberate since misinformation content is intentionally modeled after established news outlets and the semantics

of news production, often relying on headlines, captions, and the use of quotes to frame an unfolding story.

This leads to a key shortcoming of the IMPED model: it cannot identify sophisticated disinformation campaigns that appear well-reasoned and supported by evidence. One such example is Andrew Wakefield's article linking MMR vaccine and autism, which was published in a prestigious academic journal and fashioned in scientific language. Another fundamental shortcoming of the IMPED model is that it is more likely to classify alternative sources as low-quality information compared with mainstream media sources. This is due to the model drawing from the assumption that quality information can be undermined by networked communication and declining public trust in the press. In other words, the parameters of the model rely on linguistic markers and post metadata to identify less-well-edited text, highly partisan discourse, and user-generated content that is quickly modified or erased; measures that are likely to identify not only misinformation and the broader universe of low-quality information, but also user-generated posts from marginalized communities, activist groups, and grassroots organizations in the counterpublics where younger, less educated, second-language speakers are more likely to participate and voice potentially radical political ideas.

We are nonetheless cognizant of the positive roles that user-generated communication networks play in online deliberation (Bennett & Pfetsch, 2018). These networks support gatewatching (Bruns, 2005) and practices in citizen journalism essential to a diverse media ecosystem (Hermida, 2010), with citizens auditing the gatekeeping power of mainstream media and holding elite interests to account (Tufekci & Wilson, 2012). Much of the information generated within the scope of citizen journalism is ephemeral, uses informal language, and displays partisan traits that differ from the standards of mainstream news coverage. These

sources of information may not only deviate from any definition of misinformation; they also play a central role in increasing the transparency of the democratic process. Conversely, news outlets enforcing selective gatekeeping may inadvertently trigger misinformation, as specific perspectives are systematically prevented from reaching larger audiences.

These shortcomings can only be properly addressed by parametrizing the model with real-world data that includes instances of citizen journalism and unintentional behavior that is misleading. As the model draws from probability distributions, it should be possible to address these shortcomings while also avoiding the limitations of current approaches for the detection of problematic content based on fact-checking and predictive analytics. The underlying assumption is that high-quality content is tailored to remain available over time as opposed to low-quality content, largely optimized for fast turnover. This leads to one of the central parameters of the model: the ephemerality score that can help identifying low-quality information circulating on social media at scale. The relatively low computational requirements of the proposed heuristic and score-based model has the added benefit of allowing low-quality content to be identified as it reappears on social media platforms through multiple iterations, even after the original content has been blocked or the webpage removed.

We nonetheless expect challenges to arise when implementing the model at scale. Implementations of the model will require the real-time archiving of the posts and URLs posted on social media platforms, but the output of the model can only be offered in near real-time. This is because the calculations of the ephemerality scores for domain and *u-content* require continuous processing of archived data to identify when posts were modified. These posts need to be checked regularly to identify the point in time when the resource was modified or became inaccessible (URL decay). More sophisticated implementations may incorporate metadata from

users and posts and sharing metrics that can be used to prioritize the workflow of the model, but a time-lag between social media posting and the IMPED calculation is likely to remain.

Acknowledgements

The authors are thankful to Michael Levy for his help in developing the s-score algorithm and for the support offered to this project by the Research Computing team at Arizona State University.

Funding

This research was supported by Twitter, Inc. research grant 50069SS “The Brexit Value Space and the Geography of Online Echo Chambers.”

References

- Allcott, H., & Gentzkow, M. (2017). Social Media and Fake News in the 2016 Election. *Journal of Economic Perspectives*, 31(2), 211-236. <https://doi.org/10.1257/jep.31.2.211>
- Arif, A., Robinson, J. J., Stanek, S. A., Fichet, E. S., Townsend, P., Worku, Z., & Starbird, K. (2017). A closer look at the self-correcting crowd: Examining corrections in online rumors. 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing, Portland, Oregon, USA.
- Barzilai-Nahon, K. (2008, Jul). Toward a theory of network gatekeeping: A framework for exploring information control. *Journal of the American Society for Information Science and Technology*, 59(9), 1493-1512. <https://doi.org/10.1002/asi.20857>
- Barzilai-Nahon, K. (2009). Gatekeeping: A Critical Review. *Annual Review of Information Science and Technology*, 43, 433-478. <Go to ISI>://000260638900011
- Bastos, M. T. (2015, 2015/05/04). Shares, Pins, and Tweets: News readership from daily papers to social media. *Journalism Studies*, 16(3), 305-325. <https://doi.org/10.1080/1461670x.2014.891857>
- Bastos, M. T. (2016). Digital Journalism and Tabloid Journalism. In B. Franklin & S. Eldridge (Eds.), *Routledge Companion to Digital Journalism Studies* (pp. 217-225). Routledge.

Bastos, M. T., & Mercea, D. (2019). The Brexit Botnet and User-Generated Hyperpartisan News. *Social Science Computer Review*, 37(1), 38-54.

<https://doi.org/10.1177/0894439317734157>

Bastos, M. T., Raimundo, R. L. G., & Travitzki, R. (2013, March 1, 2013). Gatekeeping Twitter: message diffusion in political hashtags. *Media, Culture & Society*, 35(2), 260-270.

<https://doi.org/10.1177/0163443712467594>

Baumgartner, J. (2019). *24.7 Million Original Tweets from All Verified Twitter Accounts*.

http://pushshift.io/twitter/twitter_word_frequency.csv

<https://twitter.com/jasonbaumgartne/status/1008959965162299394>

http://pushshift.io/twitter/TU_verified.ndjson.xz

Benkler, Y., Faris, R., & Roberts, H. (2018). *Network Propaganda: Manipulation, Disinformation, and Radicalization in American Politics*. Oxford University Press.

Bennett, W. L., & Livingston, S. (2018). The disinformation order: Disruptive communication and the decline of democratic institutions. *European Journal of Communication*, 33(2), 122-139.

Bennett, W. L., & Pfetsch, B. (2018). Rethinking Political Communication in a Time of Disrupted Public Spheres. *Journal of Communication*, 68(2), 243-253.

<https://doi.org/10.1093/joc/jqx017>

Breiman, L. (2001). Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical Science*, 16(3), 199-231.

Brugnoli, E., Cinelli, M., Zollo, F., Quattrocioni, W., & Scala, A. (2019). *Lexical convergence inside and across echo chambers*.

Bruns, A. (2005). *Gatewatching: Collaborative online news production* (Vol. 26). Peter Lang Pub Inc.

Cheng, J., Adamic, L. A., Dow, P. A., Kleinberg, J., & Leskovec, J. (2014). Can cascades be predicted? 23rd International Conference on World Wide Web (WWW'14), Seoul, Korea.

Comor, E. (2001). Harold Innis and 'The bias of communication'. *Information, Communication & Society*, 4(2), 274-294.

comScore. (2013a). *2013 UK Digital Future in Focus* (Future in Focus, Issue.

comScore. (2013b). *2013 US Digital Future in Focus* (Future in Focus, Issue.

Dow, P. A., Adamic, L. A., & Friggeri, A. (2013, 8-11 July 2013). The Anatomy of Large Facebook Cascades. 7th International AAAI Conference on Weblogs and Social Media (ICWSM13), Boston.

Community Standards, (2018a). <https://www.facebook.com/help/975828035803295>

Understanding the Facebook: Community Standards Enforcement Report, (2018b).

https://fbnewsroomus.files.wordpress.com/2018/05/understanding_the_community_standards_enforcement_report.pdf

Fletcher, W. H. (2012). Corpus analysis of the world wide web. *The encyclopedia of applied linguistics*.

Fukuyama, F. (1995). *Trust: The social virtues and the creation of prosperity*. Free Press Paperbacks.

González-Bailón, S., Borge-Holthoefer, J., & Moreno, Y. (2013, March 8, 2013). Broadcasters and Hidden Influentials in Online Protest Diffusion. *American Behavioral Scientist*, 57(7), 943-965. <https://doi.org/10.1177/0002764213479371>

Halpern, D., & Gibbs, J. (2013). Social media as a catalyst for online deliberation? Exploring the affordances of Facebook and YouTube for political expression. *Computers in Human Behavior*, 29(3), 1159-1168.

Hermida, A. (2010, 2010/08/01). Twittering The News: The emergence of ambient journalism. *Journalism Practice*, 4(3), 297-308. <https://doi.org/10.1080/17512781003640703>

Innis, H. A. (2008). *The bias of communication*. University of Toronto Press.

Isenegger, K., Dong, Y., Shang, M., Furst, J., & Stan-Raicu, D. (2019). Characterizing and Quantifying Diagnostic (Un) Certainty in Medical Reports through Natural Language Processing. 2019 International Conference on Computational Science and Computational Intelligence (CSCI),

Jenkins, H., Ford, S., & Green, J. (2012). *Spreadable media: Creating value and meaning in a networked culture*. NYU Press.

Jiang, S., & Wilson, C. (2018). *Linguistic Signals under Misinformation and Fact-Checking: Evidence from User Comments on Social Media* Proceedings of the ACM on Human-Computer Interaction,

Karlova, N. A., & Fisher, K. E. (2013). A social diffusion model of misinformation and disinformation for understanding human information behaviour Information Research, 18 (1) paper 573. *Information Research*, 18(1). <http://informationr.net/ir/18-1/paper573.html>

Kuklinski, J. H., Quirk, P. J., Jerit, J., Schwieder, D., & Rich, R. F. (2000). Misinformation and the currency of democratic citizenship. *Journal of Politics*, 62(3), 790-816.

- Lazer, D. M. J., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., Metzger, M. J., Nyhan, B., Pennycook, G., Rothschild, D., Schudson, M., Sloman, S. A., Sunstein, C. R., Thorson, E. A., Watts, D. J., & Zittrain, J. L. (2018). The science of fake news. *Science*, 359(6380), 1094-1096. <https://doi.org/10.1126/science.aao2998>
- Lijffijt, J., Nevalainen, T., Säily, T., Papapetrou, P., Puolamäki, K., & Mannila, H. (2016). Significance testing of word frequencies in corpora. *Literary and Linguistic Computing*, 31(2), 374-397.
- Marietta, M., & Barker, D. C. (2019). *One Nation, Two Realities: Dueling Facts in American Democracy*. Oxford University Press.
- Maynard, A. D. (2015). Why we need risk innovation. *Nature nanotechnology*, 10(9), 730-731.
- Newman, N., Fletcher, R., Levy, D. A. L., & Nielsen, R. K. (2016). *Reuters Institute Digital News Report 2016*.
- Pérez-Rosas, V., Kleinberg, B., Lefevre, A., & Mihalcea, R. (2017). *Automatic detection of fake news*.
- Pielou, E. C. (1966). The measurement of diversity in different types of biological collections. *Journal of theoretical biology*, 13, 131-144.

Rubin, V. L., Liddy, E. D., & Kando, N. (2006). Certainty identification in texts: Categorization model and manual tagging results. In *Computing attitude and affect in text: Theory and applications* (pp. 61-76). Springer.

Rucker, D. D., Tormala, Z. L., Petty, R. E., & Briñol, P. (2014). Consumer conviction and commitment: An appraisal-based framework for attitude certainty. *Journal of Consumer Psychology*, 24(1), 119-136.

Shannon, C. E., & Weaver, W. (1962). *The mathematical theory of communication*. University of Illinois Press.

Shoemaker, P., & Vos, T. (2009). *Gatekeeping theory*. Routledge.

Starbird, K. (2019). Disinformation's spread: bots, trolls and all of us. *Nature*, 571(7766), 449.

Starbird, K., Maddock, J., Orand, M., Achterman, P., & Mason, R. M. (2014). Rumors, false flags, and digital vigilantes: Misinformation on twitter after the 2013 boston marathon bombing. IConference 2014, Berlin, Germany.

Thompson, A. (2017). *All the news: 143,000 articles from 15 American publications* Version 4) Kaggle. <https://www.kaggle.com/snapcrack/all-the-news/data>

Tufekci, Z., & Wilson, C. (2012). Social media and the decision to participate in political protest: Observations from Tahrir Square. *Journal of Communication*, 62(2), 363-379.

Twitter. (2018a). *Retweet FAQs* <https://help.twitter.com/en/using-twitter/retweet-faqs>

Twitter. (2018b). *Update on Twitter's Review of the 2016 U.S. Election* (Global Public Policy, Issue. https://blog.twitter.com/official/en_us/topics/company/2018/2016-election-update.html

Uscinski, J. E., & Butler, R. W. (2013, 2013/06/01). The Epistemology of Fact Checking. *Critical Review*, 25(2), 162-180. <https://doi.org/10.1080/08913811.2013.843872>

Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online. *Science*, 359(6380), 1146-1151. <https://doi.org/10.1126/science.aap9559>

Walker, S. (2015). *The Complexity of Collecting Social Media Data in Ephemeral Contexts* Internet Research 16, Phoenix, AZ.

Weedon, J., Nuland, W., & Stamos, A. (2017). *Information Operations and Facebook*.

Welbers, K., & Opgenhaffen, M. (2018). Social media gatekeeping: An analysis of the gatekeeping influence of newspapers' public Facebook pages. *New Media & Society*, 20(12), 4728-4747. <https://doi.org/10.1177/1461444818784302>

Wu, S., Hofman, J. M., Mason, W. A., & Watts, D. J. (2011). Who Says What to Whom on Twitter. 20th international conference on World Wide Web, New York.

Xie, J., Sreenivasan, S., Korniss, G., Zhang, W., Lim, C., & Szymanski, B. K. (2011). Social consensus through the influence of committed minorities. *Physical Review E*, 84(1), 011130. <https://doi.org/citeulike-article-id:9578014>

Zuckerman, E. (2017). Mistrust, efficacy and the new civics: Understanding the deep roots of the crisis of faith in journalism.

Table 1: Catalogue of metrics of IMPED model

| | | |
|---|----------------|------------------------------------------------------------|
| 1 | <i>s-score</i> | Shannon diversity index of the linguistic diversity |
| 2 | <i>p-score</i> | Lexicon-based classifier of partisan-certainty |
| 3 | <i>u-score</i> | User-generated vs. mainstream media classifier |
| 4 | <i>e-score</i> | Ephemerality level from stable to modified to inaccessible |
| 5 | <i>d-score</i> | Mean deletion rate of webpages from the source domain |

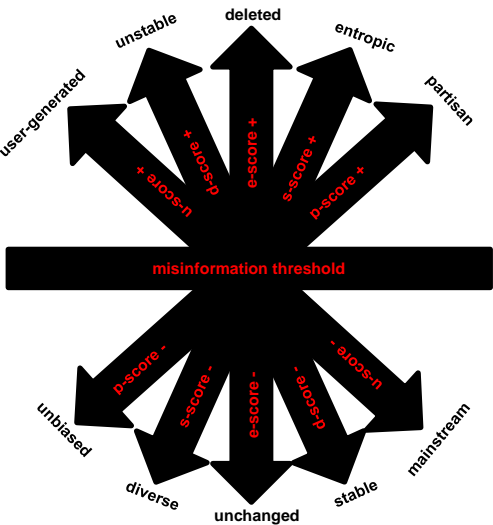


Figure 1: IMPED theoretical model

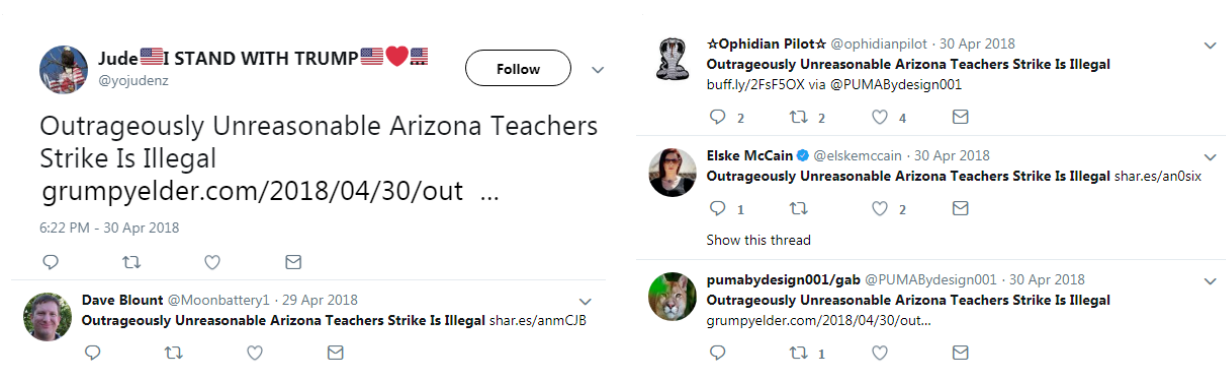


Figure 2: Reposting of "Outrageously Unreasonable Arizona Teachers Strike Is Illegal" hosted by other domains

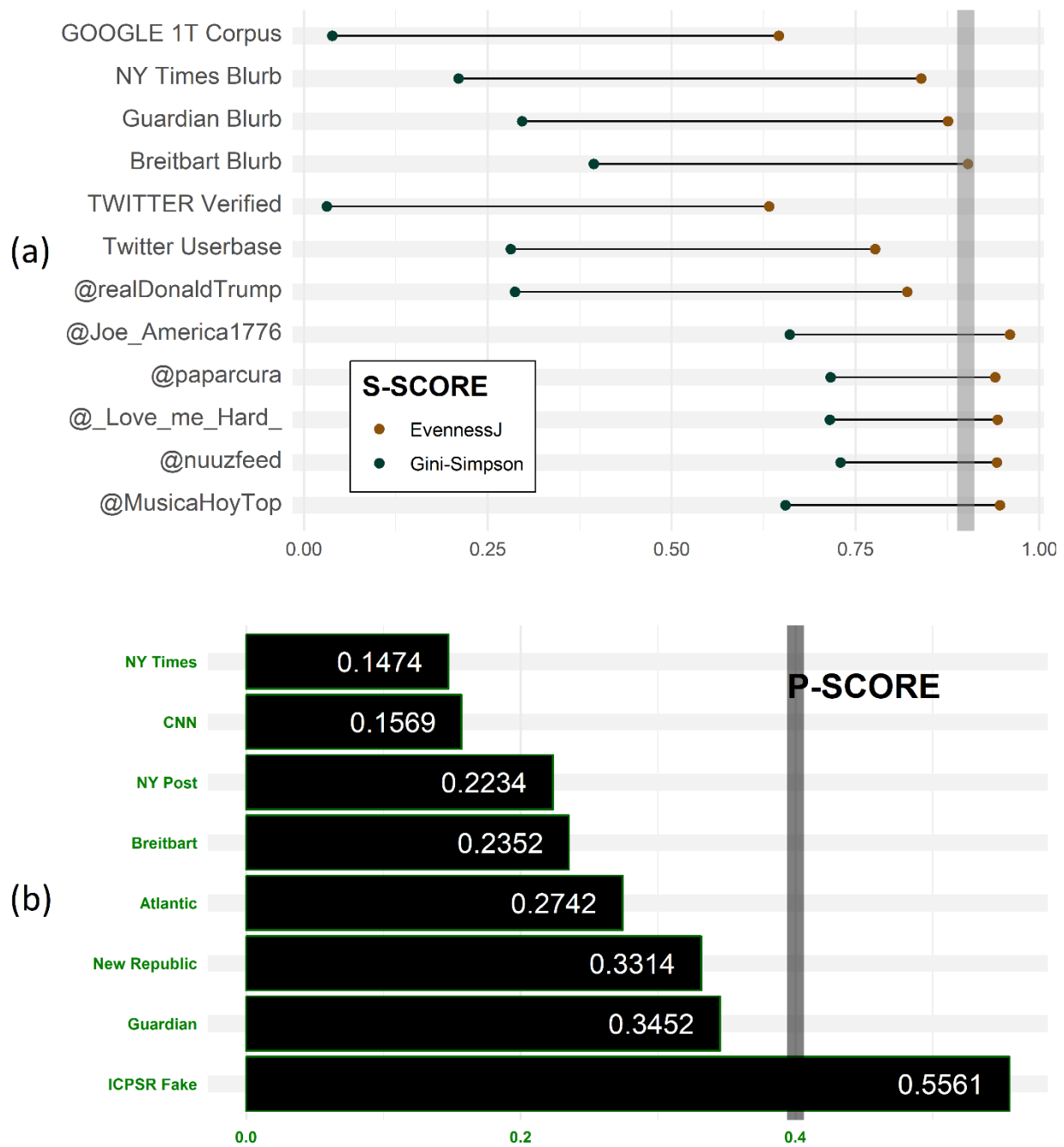


Figure 3: Validation of a) *s-score* and b) *p-score*