



# Style over substance: A psychologically informed approach to feature selection and generalisability for author classification

Isabel Holmes<sup>a,\*</sup>, Timothy Cribbin<sup>a</sup>, Nelli Ferenczi<sup>b</sup>

<sup>a</sup> Department of Computer Science, Brunel University London, UK

<sup>b</sup> Department of Psychology, Brunel University London, UK

## ARTICLE INFO

### Keywords:

Author profiling  
Political affiliation classification  
Stylistic feature sets  
Model generalisability  
Political psychology  
Feature development  
Interdisciplinarity  
Domain-specific knowledge

## ABSTRACT

Author profiling, or classifying user generated content based on demographic or other personal attributes, is a key task in social media-based research. Whilst high-accuracy has been achieved on many attributes, most studies tend to train and test models on a single domain only, ignoring cross-domain performance and research shows that models often transfer poorly into new domains as they tend to depend heavily on topic-specific (i.e., lexical) features. Knowledge specific to the field (e.g., Psychology, Political Science) is often ignored, with a reliance on data driven algorithms for feature development and selection.

Focusing on political affiliation, we evaluate an approach that selects stylistic features according to known psychological correlates (personality traits) of this attribute. Training data was collected from Reddit posts made by regular users of the political subreddits r/republican and r/democrat. A second, non-political dataset, was created by collecting posts by the same users but in different subreddits.

Our results show that introducing domain specific knowledge in the form of psychologically informed stylistic features resulted in better out of training domain performance than lexical or more commonly used stylistic features.

## 1. Introduction

Researchers are increasingly interested in what we can discover about a person from their writing. What can a person's posts on social media, for example, reveal about their social group, attitudes or personality? For instance, can we group individuals by gender purely based on their blog posts? These questions fall under the heading of author profiling, a sub task of author analysis, which involves inferring demographic and other personal attributes about the authors of a text. This area has become increasingly diverse in terms of target attributes, with studies now covering a range of domains including political science and psychology (Hinds & Joinson, 2018; Oberlander & Gill, 2004; Yu et al., 2008). In particular, the topic of political affiliation classification has been addressed many times. Here the task is to label an author (or speaker if speeches are used) by their political affiliation or outlook. In the United States, for example, it might be a binary task - Republican or Democrat - although in some circumstances potential political affiliations or outlooks may sometimes involve a higher cardinality i.e., a multiclass task (Gu & Jiang, 2021; Yu & Diermeier, 2010). The task of inferring political affiliation typically involves the use of machine

learning algorithms which must be trained using features extracted from text.

Developing a good feature-set is key for ensuring a model performs well, as summarised by the axiom "garbage in, garbage out". Approaches in political classification have been varied. For example, researchers have made wide use of 'bag-of words' methods such as TFIDF, a way of weighting the importance of words by their frequency, when vectorising text and selecting which words to use as inputs (Yu & Diermeier, 2010; Yu et al., 2008). Vector representations comprised of literal word or unigram counts can in some cases make up the entirety of the feature-set. More recent work has utilised more sophisticated word-embedding approaches such as GloVe (Pennington et al., 2014) and also stylistic features and non-textual features such as retweeting or mentioning (Das et al., 2021). Typically, it is model performance rather than a priori hypotheses (i.e., a data-driven approach) that is used to determine which features are likely to be most effective at discriminating the classes. However, our research shows that little to no attention has been paid to what attribute domain-specific knowledge could do to benefit data science work in this area.

Several traits are known to correlate with conservative or liberal

\* Corresponding author.

E-mail address: [isabel.holmes@brunel.ac.uk](mailto:isabel.holmes@brunel.ac.uk) (I. Holmes).

<https://doi.org/10.1016/j.chbr.2022.100267>

Received 25 October 2021; Received in revised form 14 December 2022; Accepted 20 December 2022

Available online 13 January 2023

2451-9588/© 2023 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

beliefs, and previous work has found that traits such as grandiose narcissism manifest in writing (Cutler et al., 2021; Jost et al., 2003; Kruglanski, 1996; Zavala et al., 2010). We therefore argue here that more valid models might result if we use this knowledge of predictive traits and their likely expressive manifestations to inform the specification of features.

A second, related problem in political inference modelling, we argue, is that model performance is usually assessed only on text from a very similar topical domain to that used training. This means it is often difficult to know how well a model will generalise to new observations or, indeed, if is actually measuring political affiliation rather than some other trait or attitude. In this paper, our experimental results show how model performance within a training domain of political discourse is not necessarily predictive of performance on the same authors writing in a different context.

To summarise, in the present study we introduce a feature-set developed by examining measures for three traits that have relationships with political beliefs; social dominance orientation, need for cognitive closure and need for cognition (Cacioppo & Petty, 1982; Pratto et al., 1994; Webster & Kruglanski, 1994). We compare a support vector machines model trained using these features against a vectorised text-only approach and an approach that uses stylistic features chosen without reference to psychological, political or sociological literature. We then test each model on a non-political data set, to try and capture which model is truly classifying based on political stance as opposed to context specific clues. Our aim is to highlight the possibilities that even a light-touch approach to domain-specific knowledge, in this case from the field of psychology, can offer researchers, whilst also offering a potential avenue of research that address model generalisability issues.

## 2. Related work

In this section we begin by reviewing work on political affiliation classification, before discussing the importance of considering model generalisability as part of the model testing process. We then introduce our psychologically informed approach to feature selection, citing relevant empirical evidence of traits associated with political affiliation. Finally, we define our experimental aims and hypotheses.

### 2.1. Political affiliation classification

Political affiliation classification can be defined as the task of determining an author's political stance from their written (or oral) communications. Much of the early work in this area focused as classifying authors, often politicians, by membership of a political party (Dahlhof, 2012; Diermeier et al., 2012; Yu & Diermeier, 2010; Yu et al., 2008). Historically, researchers used classic machine learning techniques such as Support Vector Machines with 'bag of words' (BoW) feature sets. Feature selection was performed using formulas such as term-frequency inverse document frequency (TFIDF).

These early efforts yielded promising results. For example, researchers were able to classify congressional speeches correctly as Republican or Democrat in 80% of instances (Yu et al., 2008). Work classifying social media users, particularly Twitter users, also appeared to have good results (Makazhanov & Rafiei, 2013; Pennacchiotti & Popescu, 2011). For example, Joshi et al. (2016) gave an accuracy figure of 68% when classifying twitter uses as Republican or Democrat. More recently, researchers have reported accuracies over 90% when classifying tweets (Ullah et al., 2021). Researchers have also been able to classify celebrities by their political affiliation using tweets (Das et al., 2021). Here, features specific to Twitter have been utilised by researchers, including hashtag usage alongside stylistic measures and text vectors. In addition, work has been carried out in various languages, such as Chinese (Gu & Jiang, 2021).

Despite the good results shown in many of these studies, it is rare to see any form of theoretical justification for the features used. Instead,

feature sets mostly appear to be decided in a data-driven way, that is based on experimental results or received wisdom in natural language processing practice, rather than based on any empirical evidence of psychological traits known to be associated with political affiliation or relevant theoretical frameworks. In Section 2.3, we discuss the extensive psychological research in this area, which forms the foundation of our methodological approach, detailed in Section 3.

### 2.2. Generalisability

Whilst the literature provides us with many examples of models exhibiting high accuracy results, improving the generalisability of models across time or topic has not been prioritised. In Psychology, generalisability refers to the ability to extrapolate from findings of a study to the target population at large. However, in this case, generalisability refers to the performance of the model on a different dataset, perhaps collected at a different time, or containing text that covers different topics or is written by a different author, and still achieve good results. A replication of studies that classified Twitter users found that accuracy dropped by as much as 30% when classifiers were used on everyday users, rather than political figures (Cohen & Ruths, 2013). In this case the model failed to generalise across author-type, as well as topics. This could suggest that the model is using features inherent to political speech to assign class labels, as opposed to some inherent writing style linked to political affiliation. The issue of generalisation has also been addressed in the PAN (a long-running series of tasks and events focusing on text exploration and classification) 2020 task, where fandoms of fan-fiction were varied in an cross-domain authorship verification task (Bevendorff et al., 2020).

In a further example that highlights the need to pay attention to generalisability, a model ostensibly trained to classify orators in the Canadian parliament by political party instead appeared to have labelled them by party political status: in or out of power (Hirst et al., 2010). This confound was discovered when the authors applied their model to data collected in a different period of time, to test its generalisability. Perhaps this is why other models of this kind have failed to maintain accuracy across time (Yu et al., 2008). Despite this risk, most models are not tested on datasets that feature data covering different topics or timeframes, meaning it is difficult to discern, from published results, how well these models are likely to generalise.

### 2.3. Domain-specific knowledge

We posit that many of the issues that lead to poor generalisation could be addressed by introducing domain-specific knowledge into the feature selection process. Over-fitting of a model to the training domain occurs because the optimal models tend to be biased towards surface features, such as topical key words, rather than features that are inherent or typical to the domain or attribute of interest. Whilst stylistic and other non-lexical features have been widely used, to our knowledge, previous approaches tend to be data-driven, selecting features based on algorithmic evidence, rather than domain knowledge and theory. It is our objective to explore how introducing foundational domain-specific knowledge can improve model performance. In addition, models often rely on text which contains topics that may only be relevant to a certain type of user or point in time. As an example, post the 2020 United States Presidential election voter fraud became a popular topic amongst Republicans, with as many as 77% of voters for the Republican candidate, Donald Trump, believing this type of fraud was commonplace (Pennycook & Rand, 2021). Words related to voter fraud word be highly useful features therefore for a model training on text written post 2020. However, the same model might struggle to categorise text written in 2015, when the topic was far less popular. Therefore, there is a need to introduce features free from the influence of topic, that draw upon relevant theory. We suggest a stylistic feature set created with reference to psychological traits correlated with conservatism and liberalism.

Stylistic features are commonly defined as features that represent distinctive patterns or trends in an author's writing, rather than the content or topic of the text. Much like authorship identification, stylometry has a long history and has often been used to aid author classification (Holmes, 1998). Examples can include counts or ratios of parts of speech or punctuation usage, with the idea that these features tap into authorship style over content, and are therefore able to tell us something about the 'who' of the author, as opposed to the 'what' of the text content (Kavuri & Kavitha, 2020; Lagutina et al., 2019). Stylistic features focus on pervasive and often unconscious forms of expression and may vary less than content-based features with topic or subject matter. These features should tap into the traits underpinning belief, and therefore allow a model to remain relevant across topic and time. In the case of the present study, we theorised that using stylistic features would improve model generalisability across topic where the authors remained consistent.

Below, a sentence is broken down into parts of speech, a common stylistic feature.

The	quick	brown	fox	jumped	over	the	lazy	dog
<i>determiner</i>	<i>adjective</i>	<i>adjective</i>	<i>noun</i>	<i>verb</i>	<i>preposition</i>	<i>determiner</i>	<i>adjective</i>	<i>noun</i>

To test this approach we selected three psychological traits of interest due to their relationships to political belief evidenced in the literature. These are Social Dominance Orientation (SDO), Need for Cognitive Closure (NFCC) and Need For Cognition (NFC) (Cacioppo & Petty, 1982; Pratto et al., 1994; Webster & Kruglanski, 1994). To our knowledge, with the exception of one paper limited to the use of nouns and NFCC, there has been no work that has examined how traits linked to political belief might manifest in writing (Cichočka et al., 2016). Therefore, we decided to draw upon traits shown to have relationships with political affiliation, and extrapolate from them. We seek to demonstrate that minimal reference to relevant domain knowledge can improve model performance, even without the costs associated with recruiting participants to create a primary data-set. For this reason, we used measures of these three traits as references to justify feature selection.

### 2.3.1. Social dominance orientation

Social dominance orientation (SDO) reflects a person's preference for hierarchy. A person scoring high in this trait would prefer for society to be organised in such a way that some groups are higher than others, and they believe that there is a natural order to society (Pratto et al., 1994). SDO has been shown to predict conservatism (Harnish et al., 2018; Pratto et al., 1994; Wilson & Sibley, 2013). Given that conservatism has been defined in the literature as a reluctance to change, a desire to maintain existing order, and an acceptance that society will always be to some extent unequal, the parallels with SDO are clear and it is not surprising that the two are linked (Huntington, 1957; Jost et al., 2003). More recent research suggests that SDO can be seen not only as a preference for hierarchy, but as a strategy for gaining power and maintaining ingroup dominance (Sinn & Hayes, 2018). An example of how this trait might manifest in Republican policy is encapsulated particularly in the anti-immigration policies of the party, such as the so called 'Muslim Ban', where then President Donald Trump prevented residents of several predominantly Muslim countries from entering the United States of America (ACLU, 2017). This policy fits neatly with research which found SDO to be strongly associated with low warmth towards immigrants, as well as anti-immigration attitudes (Satherley & Sibley, 2016). In our approach, we relied on a measure of Social Dominance Orientation developed by Ho et al. (2015) for the present study (appendix 1).

### 2.3.2. Need for cognitive closure

Need for cognitive closure (NFCC) (Webster & Kruglanski, 1994) reflects an individual's preferences and motivations for making judgments and interpreting information. Those high in the trait seek quick answers to questions and dislike ruminating on an issue. They feel uncomfortable when faced with ambiguity, and conversely comfort when given certainty. Once they have found an answer, they are resistant to change even if their view is proven to be factually inaccurate (Kruglanski, 1996). Need for cognitive closure has been shown to be higher in those with conservative views; indeed, a meta-analysis conducted by Jost et al. (2003) found that need for cognitive closure correlated significantly with self-reported conservatism. We used another short-form measure (Roets & Van Hiel, 2011) for inspiration, and again relied on sub-facets as well as individual questions (see appendix 2 For scale).

### 2.3.3. Need for Cognition

Need for cognition (NFC) can be summarised as a drive to think deeply about and fully comprehend a subject or problem (Cacioppo &

Petty, 1982). Those high in this trait enjoy exploring the facets of an argument, in almost direct contrast to those high in need for cognitive closure. For example, a person high in NFC would report putting more effort into thinking about a task, and also recall multiple argument messages post-task (Cacioppo et al., 1983). Need for cognition has been found to be positively correlated with liberal views and attitudes, and negatively correlated with conservatism; however, it is important to note that the correlation, whilst significant, is small (Ksiazkiewicz et al., 2016). As with the scales used above, we use a short form version of an original scale, namely the six-item need for cognition scale developed by Lins de Holanda Coelho et al. (2018) (see appendix 3 for scale).

In summary, we posit that we can map from the kinds of traits identified above to specific stylistic features and that the features inspired by these traits should be similarly present, and discriminatory, in both political and non-political writing, as the traits themselves remain consistent across time and setting when topics do not. We therefore expected models containing such features would generalise better than models that did not.

Following the above, we developed the following formal hypotheses:

H1 The text-only model will be the weakest performer on the test set.

H2 The model trained using theory informed features will outperform the non-theory informed feature-set.

## 2.4. Experimental design

To determine the effectiveness of this approach a modified testing approach is required. Work in the field of data science often relies on results of k-fold cross-validation as a measure of performance or hold-out test set performance, and in particular cross validation is favoured when datasets are relatively small as in the present study (Yadav & Shukla, 2016). Fig. 1 describes the process of k-fold cross validation, where results are given as the mean of performance across the various folds. Fig. 2, in contrast, shows a hold-out test set approach, where a model is trained on the training set alone and then performance is measured on the unseen test set. However a hold-out test set is typically drawn from the same domain as the training data and therefore is likely to contain the same topical characteristics. Both approaches, we argue, run the risk of the model being over-fitted to the idiosyncratic properties of the training data, rather than the attribute domain itself, which can result in poor generalisability of the model.

To address this problem, we applied a dataset that features non-

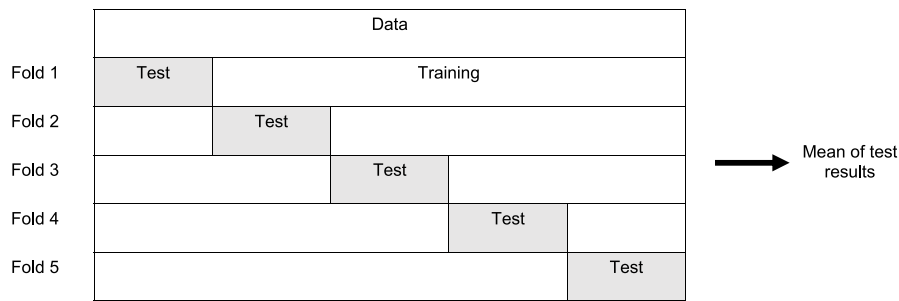


Fig. 1. A representation of k-fold cross validation. Results are calculated by finding the mean performance on each fold.

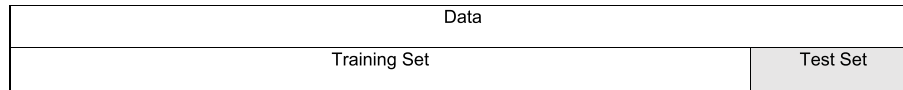


Fig. 2. Holdout test set. Here performance on the test set, which is unseen by the model during training, is reported.

political speech, posted by the same authors, as a test set. The test set unseen by the model during training and does not contain the same topics or themes as the training data. In this way we hope to provide a better assessment of generalisability, in a similar fashion to the approach taken at PAN 2020 (Bevendorff et al., 2020; PAN, 2020). We compare models trained using a text-only approach, a text and standard stylistic feature approach, and a model trained using our domain specific stylistic features and text. The aim is to explore how synthesising knowledge from different fields (in this case author profiling and psychology), can be of benefit to data scientists. In all other aspects, we try to use standard practices in the field to investigate the impact of just the addition of the psychologically informed features.

### 3. Methodology

#### 3.1. Dataset creation

Two datasets were created by collecting posts from the Reddit API and a python script using the PRAW wrapper (Boe, 2016). The first

dataset was made up of posts made in the r/democrats and r/republican subreddit. All authors with more than four posts in the dataset were retained, however in the case of a user having made fewer posts, if they had more downvoted than upvoted posts they were dropped. This was to filter out potential troll posters who might post infrequently to elicit a negative reaction. Comments were then concatenated by author, and only those who had written more than 100 words were kept. A flow chart of these steps can be seen in Fig. 3

To create a second dataset made up of non-political posts, we collected posts made by the authors in the original dataset in other subreddits: r/amitheasshole, where users ask Redditors to provide feedback on morally ambiguous situations, and r/todayilearned, where users share interesting knowledge they have learned. These subreddits were chosen as there was a large number of users in the political dataset who had made posts in them. Again, documents were concatenated by author and dropped if they were less than 100 words in length. This gave us 811 Democrat authors and 424 Republican authors in this non-political dataset. To balance the classes, we dropped half of the Democrat authors at random, giving us a new total of 406. The sample method included with the Pandas library was used.

Table 1 below shows key information for the training dataset. The mean comment length was 60 words (0dp) and the mean document (concatenated comments) length was 452 (0dp). Some authors were super-contributors; 14 users made over 200 comments, and two made over 1000. The max number of posts was 2887.

Across both datasets, Democrat authors were coded as 0, and Republican authors were coded as 1.

#### 3.2. Feature-set development

The three measures chosen were the Social Dominance Orientation Short Scale (Ho et al., 2015) (appendix 1), the short form Need for Cognitive Closure Scale (Roets & Van Hiel, 2011) (appendix 2), and the Need for Cognition Scale (Lins de Holanda Coelho et al., 2018) (appendix 3). A detailed list of all features, relevant trait and extraction method can be found in appendix 4.

Table 2 shows all features intended to tap into Need for Cognition. In keeping with the statements in the measure shown in appendix 3, we tried to select features that would convey a sense of openness and complex thought. For example, we scored posts using several measures of readability, as we hypothesised that a higher level of writing might indicate more complex thinking and argumentation.

Table 3 shows features intended to map to Need for Cognitive Closure. Here we tried to capture the sense of certainty craved by individuals high in this trait. For example, we selected modal verbs of

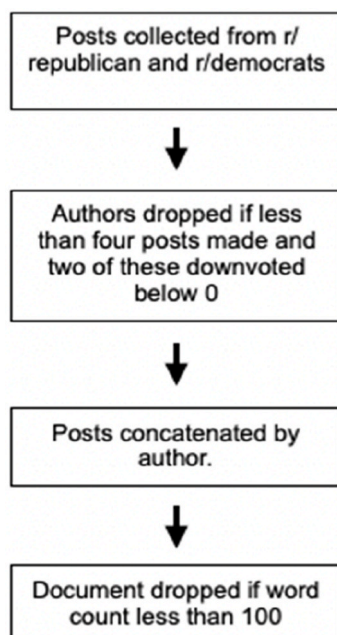


Fig. 3. Diagram to show dataset creation.

**Table 1**  
Training dataset basics.

Subreddit	Number of Authors/Documents	Number of Posts	Number of Documents in Non-Political Dataset	Mean Number of Posts per Author (2dp)	Mode Number of Posts per Author
r/democrats	4366	50,751	811	11.62	3
r/republican	4242	44,157	424	10.41	3
Total	8608	94,908			

**Table 2**  
Features inspired by Need for Cognition.

Features			
Mean comment length	Dash	Conjunctions	Dale-Chall
Mentions of subreddits	Flesch-Kincaid Grade Level	Question mark	Mean syllables per sentence
Mentions of users	Gunning-Fog	Colons	Number of hapax legomena
Pronouns	Automated Readability Index	Semicolons	Average sentence length
Urls and Emails	Coleman-Liau Index	Commas	Average characters per sentence

**Table 3**  
Features inspired by Need for Cognitive Closure (\* indicates hypothesised negative correlation with trait).

Features			
Nouns	First person plural pronouns	Proper Nouns	Adverbs of certainty high*
Possessive Nouns	Third person plural pronouns	Exclamation mark	Adverbs of certainty low
Determiners	Modal verbs of obligation	Modal verbs of possibility*	Adverbs of frequency high*

**Table 4**  
Features inspired by Social Dominance Orientation (\* indicates hypothesised negative correlation with trait).

Features	
Money	Comparative adjectives
Possessive pronouns - first person singular	Superlative adjectives
Possessive pronouns - first person plural	Emojis*
Possessive pronouns - third person singular	Smileys*
Possessive pronouns - third person plural	Possessive pronouns - second person*

obligation (must, should etc.), as these have a definite feel to them.

Table 4 sets out the features linked to Social Dominance Orientation. As an example here, third person plural pronouns such as “they” were intended to map onto the desire for group divisions.

Further examples of features linked to each construct are given in the next section to illustrate extraction techniques.

### 3.2.1. Feature extraction and data preparation

A variety of techniques were used for feature extraction. Where feasible, Python code was used to calculate word counts. For slightly more complex extraction, such as counts of types of punctuation, regular expressions were used in a script written by the authors. At the level above this, we made use of prebuilt Python libraries. For example, we used the Readability library (van Cranenburgh, 2019) for Gunning-Fog scores, Automated Readability Index, and Flesch-Kincaid grade-level measure.

In order to obtain parts of speech counts we used two separate parts of speech taggers: TweetNLP (Owoputi et al., 2013) and the NLTK parts of speech taggers (Bird et al., 2009). TweetNLP deals well with slang and the short posts made on social media, however the NLTK tagger provided extra tags such as determiners. For a full list of all features and the extraction techniques used, see appendix 4. All non-text features were normalised using the SciKit Learn (Pedregosa et al., 2011) libraries normalise function which applies L2 normalisation (values are scaled so that the sum of squares is 1). Normalisation improves performance when using a distance-based method such as SVM (Ali & Smith-Miles, 2006).

Below is a brief description of a selection of features, to illustrate our rationale as well as the extraction process.

### 3.2.2. Nouns

Work by Cichocka et al. (2016) found that conservatives prefer nouns over adjectives. It is hypothesised that individuals scoring higher on NFCC prefer to use nouns over adjectives as a way of defining or stereotyping people or other entities.

“Small fact<sub>1</sub>, though the Navajo<sub>2</sub> wind<sub>3</sub> talkers<sub>4</sub> was the biggest group<sub>5</sub>, they weren’t the only group<sub>6</sub> of natives<sub>7</sub> who used their language<sub>8</sub> to create an unbreakable code<sub>9</sub>.”

A sentence from the non-political data set with nouns numbered.

In the above example, nouns are numbered. The taggers used take into account the context of the word to estimate the correct part of speech. The final figure for this piece of text would be  $\frac{n}{w}$ , where n is number of nouns and w is total number of words in the document.

### 3.2.3. Mentions of subreddits

As those high in need for cognition prefer more complex thinking, we searched for mentions of subreddits in posts. The idea here is that referencing other sources is a more complex form of argumentation. We used a regular expression to extract these features.

Tables 5 and 6 set out how this process works. Again, the final figure is found by dividing the number of subreddit mentions by the total number of terms in the document.

### 3.2.4. Comparative adjectives

We hypothesised that an individual high in SDO and, more specifically, the dominance sub-facet, may tend to make more comparisons and seek to define things and other groups as better or worse than, because comparisons allow them to define one groups as dominant and another as subservient. We therefore added comparative adjectives to the feature set. Here again, a part-of-speech tagger was used. Below is an example of a post from the non-political test set with the comparative adjectives numbered.

“I’m sure what you said is very true, and it’s made complicated by the fact that some American products do have better<sub>1</sub> quality. I work for a manufacturer with operations in both the U.S. and overseas. The products sold overseas are sold under a different brand, worse<sub>2</sub> quality, and are cheaper<sub>3</sub> because that’s what the people there want. Americans expect higher<sub>4</sub> quality, so that’s what they get (along with a higher<sub>5</sub> price). Knowing which products are better<sub>6</sub> (and by how much) ... that’s a tricky question.”

Again, a final figure is found by dividing the number of comparative adjectives in a document by the total number of terms.

## 3.3. Text preparation

In order to make the word terms useful as features, they must be represented numerically. We did this using the following standard pre-processing steps.

1. Stopword removal
2. Lemmatization

**Table 5**  
Breakdown of regex phrase.

Regular Expression	Matches
<code>\s / .+</code>	Any whitespace
<code>\s</code>	“r/”
<code>r/</code>	Any single character
<code>.</code>	One or more of the preceding item
<code>+</code>	

**Table 6**  
Matches to the Regex phrase.

Test Phrase	Match
r/test	Yes
r/1test	Yes
r/!test	Yes
/r/test	No
rtest	No

## 3. TF-IDF vectorization

### 3.3.1. Stop word removal

This is the process of removing words that do not carry meaning and are very common and therefore unlikely to be useful for modelling. We used the list of stop words that comes as a part of the NLTK python library. There are 179 words altogether, and examples include “it”, “am” and “is”.

### 3.3.2. Lemmatization

This refers to reducing words with the same basic root meaning to one form. An example of this is shown in Table 7.

### 3.3.3. TF-IDF vectorization

This is a method of numerically representing every word in a corpus (collection of documents). The below formula is used to give each term a score that represents how important it is.

$$TF(t(\text{term of interest}), d(\text{document})) = \frac{\text{number of times } t \text{ appears in } d}{\text{total number of terms in } d}$$

$$IDF(t) = \log \frac{\text{total number of documents in corpus}}{1 + \text{number of documents containing } t}$$

$$TF - IDF = TF * IDF$$

## 3.4. Modelling

We trained our models using a support vector machine with a linear kernel as it is relatively simple to understand and a common approach in the field. Fig. 4 shows a basic depiction of an SVM model, where the aim is to find the optimal hyperplane, where the distance between the hyperplane and the closest data points, or support vectors, is maximised.

As the present study did not specifically seek to maximise performance but instead demonstrate the impact of feature-set, we tuned for C and performed no feature selection. C is an optimization parameter that effects the size of the margin in the model. A larger C will give a smaller margin, and a smaller C a larger margin, as shown in Fig. 5. We used the gridsearch feature in SciKit Learn, which inputs multiple values of given parameters and uses cross validation to determine the best performer, to find C for each model type. A C of 1 was selected for the text only and random stylistic feature model, whereas 0.1 was selected for the theory driven dataset.

The output of the model is a label of Republican or Democrat for each author in the dataset, based on the inputs or feature-set.

**Table 7**  
Three words and the lemmatized output.

Word	Reduced Form
Loudly	Loud
Louder	
Loudest	

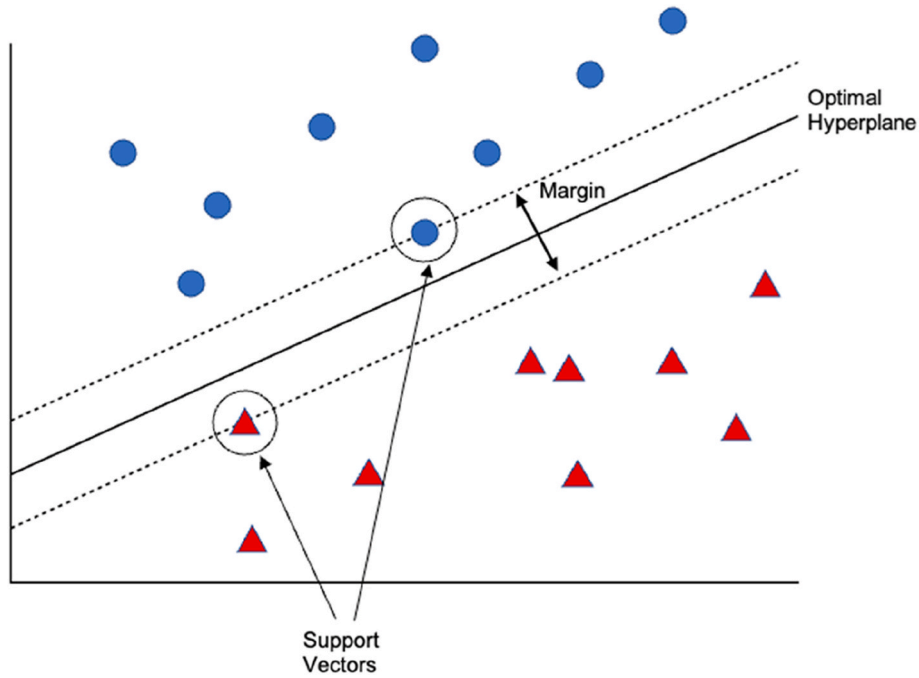


Fig. 4. Graphical depiction of the basic principles of the support vector machine algorithm.

#### 4. Results

This study sought to examine the usefulness of stylistic features, informed by psychological theory as a guard against poor generalisability in text classification. Having created a series of models using multiple feature sets, we present our findings below. We use performance on the non-political data set as an indicator of a model’s generalisability across topics.

We use accuracy, which is the percentage of authors assigned the

correct label, as our main performance metric as the test set was balanced. However, we also report F1, as this is commonly used in classification tasks. This is the harmonic mean of Recall and Precision and is preferred when a dataset is unbalanced. The lower the score, the poorer the performance. An F1 of 1 would be considered perfect performance. F1 is calculated for each class. Here we report the mean F1 score for both classes.

$$F1 = \frac{2}{\frac{1}{Recall} + \frac{1}{Precision}}$$

$$Precision = \frac{true\ positives}{true\ positives + false\ positives}$$

$$Recall = \frac{true\ positives}{true\ positives + false\ negatives}$$

Table 8 shows the results during training on the political posts (5-fold cross validation average) as well as performance on the non-political test set.

During training, the model that used random stylistic features and text is the better performer with an accuracy of 81.6%. This is a similar result to previous work in the area and is not surprising. This feature set contained additional stylistic features that may map onto political affiliation or speech in a way we did not explore.

However, as shown in Figs. 6 and 7, the performance of all the models drops when tested on the non-political dataset. The model that includes our psychologically informed features suffers from the smallest

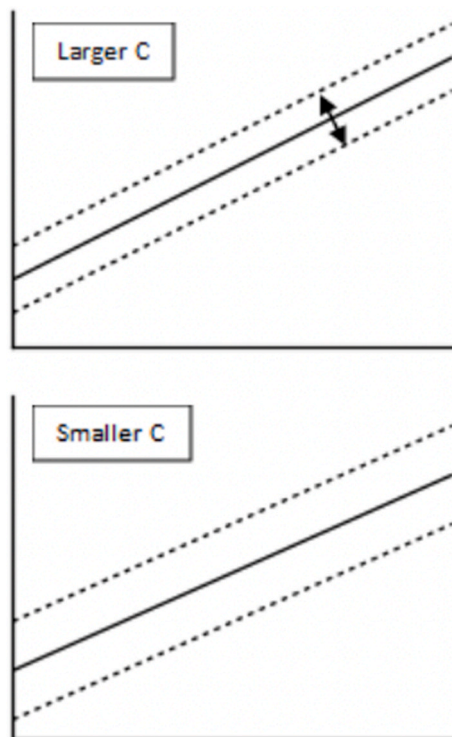


Fig. 5. Hyperplanes and margins with differing values of C.

Table 8  
-Table of results for models on cross validation and non-political test.

Features Used	Accuracy (%)		F1	
	Cross Validation	Non-Political Test Set	Cross Validation	Non-Political Test Set
Theory Driven Features and Text	80.89%	53.86%	0.801	0.538
Text Only	81%	52.77%	0.81	0.528
Random Stylistic and Text	81.60%	51.08%	0.816	0.51

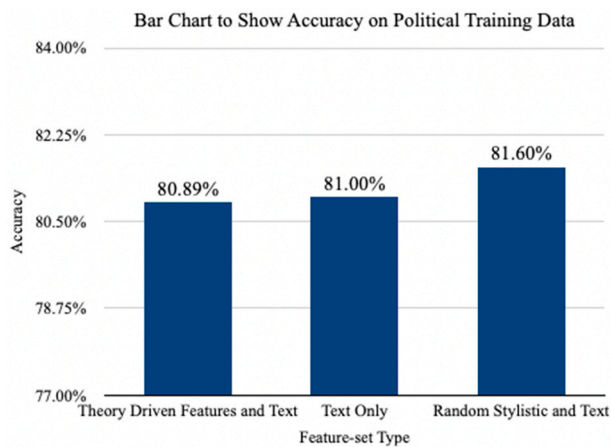


Fig. 6. A bar chart of accuracy scores (%) on training data.

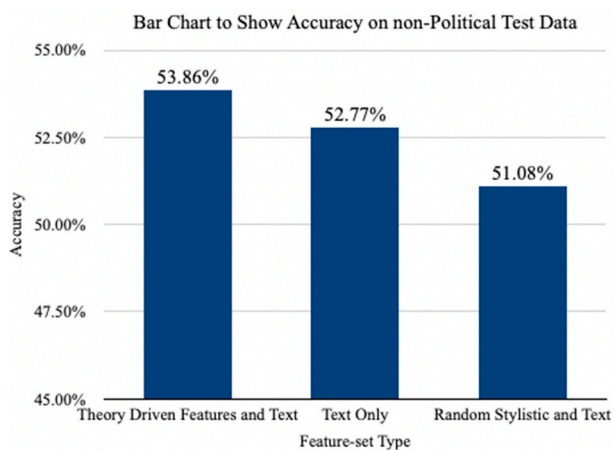


Fig. 7. A bar chart of accuracy scores (%) on test data.

drop in performance, outperforming both of the other models, albeit by a small margin. The model that was previously the best performer is now the worst.

In addition to the modelling, we carried out t-tests to examine differences between the two groups for the stylistic markers of each trait. Variables were reverse scored where they were predicted to be negatively associated with conservatism (see appendix 1 for details). We then z-scored the variables and summed all the variables associated with each trait for every author in the dataset. This gave us an overall score for NFCC markers, NFC markers and SDO markers. Republican authors ( $M = -0.026$ ,  $SD = 3.591$ ) and Democrat authors ( $M = 0.025$ ,  $SD = 3.469$ ) were not significantly different on markers of NFCC,  $t(8,606) = -0.678$ ,  $p = 0.498$ . Similarly, Republicans ( $M = 0.02$ ,  $SD = 7.808$ ) and Democrats ( $M = -0.02$ ,  $SD = 8.099$ ) did not differ significantly on markers of NFC,  $t(8,606) = -0.234$ ,  $p = 0.815$ . However, Republican authors ( $M = -0.122$ ,  $SD = 3.277$ ) and Democrat authors ( $M = 0.119$ ,  $SD = 3.015$ ) did differ significantly on markers of SDO,  $t(8500.152) = -3.551$ ,  $p = 0.000$ . In the case of the SDO variables, Levene's test was significant ( $p = 0.03$ ) which is why Welch's  $t$ -test was used.

## 5. Discussion

The results of the study show that a feature set created with domain-specific knowledge, in this case psychological traits linked to political affiliation, resulted in small but measurable gains in model generalisability. The model trained using the psychologically informed stylistic set outperformed both the text only model as well as the model that had the added benefit of non-informed stylistic features. We argue that this

suggests that the performance gain is not merely an artifact of the use of stylistic features but is in fact linked to the knowledge behind the features. By using a minimal approach that did not involve collecting primary data, we show that while this type of data may be preferable in many ways, it is not necessary for performance gains, lowering the bar in terms of accessibility to researchers from the field of data science. However, future work should seek to optimise the selection of such features through a combination of theory and experimental feedback.

In addition, the performance fall observed on the non-political test set calls into question claims made by authors such as [Diermeier et al. \(2012\)](#) that models were sorting authors by underlying ideology; if that were the case, the model should continue to detect ideology in the non-political text. Our findings support the idea that models may be categorising authors based on unknown confounding variables or may be overfitted to overtly political speech ([Cohen & Ruths, 2013](#); [Hirst et al., 2010](#)). Furthermore, the results raise questions about the usefulness of stylistic features chosen without reference to theory as a means to improve generalisability. Whilst it is true that these types of features do improve results, it is possible they could be further enhanced, and with relatively low cost to researchers.

Some may argue that models should be trained to achieve the highest accuracy possible, with generalisability less of a concern. We would argue that both objectives are equally important, and that we must think carefully about how models are to be used. If the goal is to create a model that performs as well as possible on one dataset, then the traditional approach is appropriate. On the other hand, if we want to create a model that will generalise across time and topic, we believe it would be sensible for researchers to introduce domain specific knowledge and to also use an alternative test-set, as has been done in other fields ([Yin & Zubiaga, 2021](#)). Whilst machine learning can deliver impressive results, there is value in understanding relevant theory, as shown in our results. In addition, whilst our feature-set was not the best performer on the training data, it cannot be said to be a poor performer. With greater tuning or the use of other modelling techniques any penalisation could be minimised further.

We found no significant differences between  $r$ /Republican authors and  $r$ /Democrat on the features we tied to NFC and NFCC, whilst the difference between scores for the SDO features was significant. This finding is in contrast to the several studies in psychological literature, including recent work that found that cognitive style was a better predictor of ideological preference than demographic predictors, ([Chirumbolo et al., 2004](#); [Jost et al., 2003](#); [Ksiazkiewicz et al., 2016](#); [Zavala et al., 2010](#); [Zmigrod et al., 2021](#)). However, this result is congruent with work in the field of political science suggesting that conservatives and liberals are fairly close cognitively. For example, there is little to separate conservatives and liberals when it comes to physiological response to threats, disgust sensitivity, susceptibility to fake news or, perhaps most intriguingly, the cognitive precedents of populist attitudes ([Bakker et al., 2020](#); [Clifford et al., 2022](#); [Erisen et al., 2021](#); [Strandberg et al., 2020](#)). In addition, those higher in political knowledge have been found to be more likely to engage in cognitively complex thinking when evaluating statements incongruent with their political beliefs ([Erisen et al., 2018](#)). Given that the members of a political subreddit would almost certainly have more political knowledge than most, this perhaps also explains why the  $r$ /Republican members would have little to distinguish them from the  $r$ /Democrat users in terms of our two measures of cognitive complexity.

Again, we would suggest that a more in-depth exploration of the stylistic features linked to these traits would be useful here, to rule out the possibility that these findings were merely the result of poorly chosen features. For example, examples of text with corresponding scores on the relevant measures. Although are approach is lightweight and low cost, without this kind of analysis it is too much like guesswork. Having features selected in this way may also improve the model performance overall, as our gains were very minimal.

Furthermore, the mean score on our measure of SDO were higher for



the r/Democrat authors, which is in direct contrast to well-established precedents in the literature (Harnish et al., 2018; Wilson & Sibley, 2013). Here there are two possibilities. It could be that the features selected were not tapping SDO but rather some other unknown variable. Or perhaps party opposition status played a role here as it has done in past analyses (Hirst et al., 2010). Data was collected in mid to late 2020, which meant that redditors in r/Republican had their chosen leader in the White House and also controlled the Senate, as well as having a majority of conservative judges in the Supreme court. In contrast, Democrats only controlled the House, and detested Donald Trump (in fact, the top reason Biden supporters gave for voting for him was that he was not Trump (Atske, 2020). We posit that this context may have meant r/Democrat members felt a sense of continuous threat to their ingroup as they discussed and evaluated Republican policies, which led to traits of SDO being expressed in their writing. They had a powerful outgroup to rail against, whilst the r/Republicans users were in a position of power. Indeed, intergroup threat perception and racism were both found to be higher following manipulation of threat perception, with SDO acting as a moderator (Uenal et al., 2021). Perhaps here a similar effect is occurring, with the threat of Republican dominance increasing the expression of certain stylistic features moderated by SDO. This theory could be tested by collecting posts made in both subreddits since the election of President Biden, and observing the differences in SDO for any change.

### 5.1. Future directions

In terms of future work, the introduction of more complex modelling techniques would be a logical extension of this work. In this study, we took a very simplistic high-level approach as we were concerned with showing the usefulness of our approach, rather than developing a state-of-the-art model. We chose an SVM model as that was the approach used by early researchers in the field (Diermeier et al., 2012; Yu & Diermeier, 2010; Yu et al., 2008). We also did not tune any of the model parameters apart from C, again to keep the methodology as simple as possible. However, random forests, logistic regression, naïve bayes, and KNN are all commonly used algorithms for text classification (Pranckevičius & Marcinkevičius, 2017; Shah et al., 2020). Therefore, it would be prudent to explore how a feature set such as the one developed here would affect performance in these cases.

In addition, we did not explore the impact of feature selection to our results. We would suggest a feature ablation study, perhaps using SVM recursive feature elimination (RFE) (Sanz et al., 2018). Here, whilst the number of input variables remains greater than two, a model is trained and features are ranked by the weight of their coefficients squared. The feature with the lowest value is dropped and the model is retrained. Once the process is complete, a ranked list of variables is created. Not only would this aid performance as it would allow unhelpful variables to be removed, it could also reveal interesting results relevant to psychologists. For example, if exclamation marks were found to be a highly useful variable, this would raise interesting questions as to why, opening avenues for future experimental work.

To strengthen the interdisciplinary nature of our approach, we would also suggest using primary data to improve the feature development process. This could involve recruiting participants to complete writing tasks and measures of traits of interest. The text they produce could then be explored for any differences that are linked to trait score. Indeed, previous work has been carried out exploring differences in writing style for those high and low in personality traits such as the Big Five and narcissism (Chung & Pennebaker, 2013; Cutler et al., 2021; Stillwell & Kosinski, 2012). This approach could be extended into other domains and traits, with results collated and shared across disciplines for use by researchers of different backgrounds.

Finally, this approach could be explored across languages to demonstrate that its usefulness is not limited to an American context. There is already work that explores classifying authors by political

affiliation in multiple languages and we would hope that here too reference to domain specific knowledge would be of use (Abd et al., 2020; Kapočiūtė-Dzikiene et al., 2014; Laponi et al., 2018).

### 5.2. Practical applications

There are also potential practical applications of the present study. In security research, there may be a desire to flag forum users as extremists so they can be tracked online (Ellen & Parameswaran, 2011). Here we can imagine that it would be vitally important that a model tap into an underlying trait and be generalisable across context. In this way, a user could be identified as dangerous regardless of the topics of their posts. This is especially important given that social media plays a role in the recruitment process for almost 90% of extremists in 2016 (START, 2021). Tools that can provide an early warning of such activity to the appropriate security services are of great value (Gaikwad et al., 2021).

However, the practical applications of our methodology also raise important ethical concerns. In this case, by posting in the r/Democrat and r/Republicans subreddits, users are outing their own political affiliation. However, when we work to develop a model that can classify users who post in non-political spaces, are we violating their privacy? The sanctity of the voting booth is enshrined in the universal declaration of human rights (United Nations, 1948), and if a political party, government, or organization were able to determine a person's political affiliation without their permission, there could be dangerous ramifications. For example, imagine an autocratic regime that imprisons supporters of rival political group: how could an individual stay safe when the regime could determine their political position, just from posts made in non-political spaces? Further to this, is it appropriate to label a person as extremist, with all the associated connotations, if they have not broken the law? Widescale implementation of this kind of methodology could have a chilling effect on free speech. However, given how underprepared the U.K. government, for example, is in terms of tackling issues such as far-right extremism, perhaps there is an argument to made here about the greater benefit for society at large unprepared (Ozduzen et al., 2021).

### 5.3. Limitations

In terms of the limitations of our methodology, as previously discussed, we used a very simplistic approach that does not make use of the plethora of state-of-the-art techniques available. Again, this was a deliberate choice made to allow the impact of the feature-set to be more clearly understood. It should also be noted that we looked for correlates and predictors of conservatism and liberalism, whilst our dataset is labelled as Republican or Democrat, respectively. We feel confident that these party affiliations are appropriate proxies for the relevant ideologies given that the definitions given by the literature and the policies of the parties are well-matched (Caplan, 2016; Graham et al., 2009; Jost et al., 2003; Saad et al., 2019). However, there are Conservatives and Liberals who do not identify as Republican or Democrat and vice versa. Indeed, a recent Gallup poll (Saad, 2022) found that 12% of Democrats identified as Conservative, and 4% of Republicans identified as Liberal. The solution here would be to create a dataset of posts for authors alongside measures of their political ideology, however there would be heavy financial costs associated with this approach.

### 5.4. Conclusion

Author profiling remains a popular and enduring task for data scientists. The field of political affiliation classification in particular has a long history, stemming from the classification of politicians to more recent work looking at users of social media (Gu & Jiang, 2021; Yu et al., 2008). In the present study, we have attempted to show how considering field-specific knowledge, in this case psychological theory relating to personality traits, can be helpful to political affiliation inference

research. This approach could be especially helpful with reference to the increasingly relevant issues of model generalisability, as highlighted by the recent PAN authorship attribution tasks (Bevendorff et al., 2020). In addition, our results suggest that past work may have been tapping into confounding variables, as previously suggested by other authors (Hirst et al., 2010). The psychologically informed feature-set we developed showed superior performance to the two approaches that did not involve domain-specific knowledge on the task of determining author political affiliation using non-political text. Future work should seek to extend this approach into other topics and using more sophisticated and nuanced methods.

**Declaration of competing interest**

None.

**Data availability**

Data will be made available on request.

**Acknowledgments**

We thank the reviewers for their helpful suggestions: their insight was invaluable. We would also like to thank the Editors.

**Appendix A. Supplementary data**

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.chbr.2022.100267>.

**Appendix 1**

**Table 9**  
Items on the short Social Dominance Orientation scale(\*\* indicates reverse scoring)

Statements found on the Social Dominance Orientation short scale (Ho et al., 2015)	
1.	An ideal society requires some groups to be on top and others to be on the bottom.
2.	Some groups of people are simply inferior to other groups.
3.	Groups at the bottom are just as deserving as groups at the top. **
4.	No one group should dominate in society. **
5.	Group equality should not be our primary goal.
6.	It is unjust to try to make groups equal.
7.	We should do what we can to equalize conditions for different groups. **
8.	We should work to give all groups an equal chance to succeed. **

Key: Dominance – Yellow, Antiegalitarianism - Blue.

**Appendix 2**

**Table 10**  
Items on the short form Need for Cognitive Closure scale

1.	I don't like situations that are uncertain.
2.	I dislike questions which could be answered in many different ways.
3.	I find that a well-ordered life with regular hours suits my temperament.
4.	I feel uncomfortable when I don't understand the reason why an event occurred in my life.
5.	I feel irritated when one person disagrees with what everyone else in a group believes.
6.	I don't like to go into a situation without knowing what I can expect from it.
7.	When I have made a decision, I feel relieved.
8.	When I am confronted with a problem, I'm dying to reach a solution very quickly.
9.	I would quickly become impatient and irritated if I would not find a solution to a problem immediately.
10.	I don't like to be with people who are capable of unexpected actions.
11.	I dislike it when a person's statement could mean many different things.
12.	I find that establishing a consistent routine enables me to enjoy life more.
13.	I enjoy having a clear and structured mode of life.
14.	I do not usually consult many different opinions before forming my own view.
15.	I dislike unpredictable situations.

Facet	Colour
Order	Grey
Predictability	Light Grey
Decisiveness	Yellow
Ambiguity	Light Grey
Closed-mindedness	Blue

Appendix 3

**Table 11**  
Items on the short 6 item need for cognition scale (\*\* indicates reverse scoring)

Statements found on the six item Need for Cognition Scale (Lins de Holanda Coelho et al., 2018)	
1.	I prefer complex to simple problems.
2.	I like to have the responsibility of handling a situation that requires a lot of thinking.
3.	Thinking is not my idea of fun.**
4.	I would rather do something that requires little thought than something that is sure to challenge my thinking abilities.**
5.	I really enjoy a task that involves coming up with new solutions to problems.
6.	I would prefer a task that is intellectual, difficult, and important to one that is somewhat important but does not require much thought.

Appendix 4

**Table 12**  
Feature-set breakdown

Feature	Related To	Normalised by Word Count	Relationship to Republicanism	Extraction Technique
Mean comment length	Need for Cognition	No	Negative	Author created code
Mentions of subreddits	Need for Cognition	Yes	Negative	Regular Expression
Mentions of users	Need for Cognition	Yes	Negative	Regular Expression
Pronouns	Need for Cognition	Yes	Negative	PoS Tagger <sup>a</sup>
Urls and Emails	Need for Cognition	Yes	Negative	Regular Expression
Conjunctions	Need for Cognition	Yes	Negative	PoS Tagger
Question mark	Need for Cognition	Yes	Negative	Regular Expression
Colons	Need for Cognition	Yes	Negative	Regular Expression
Semicolons	Need for Cognition	Yes	Negative	Regular Expression
Commas	Need for Cognition	Yes	Negative	Regular Expression
Dash	Need for Cognition	Yes	Negative	Regular Expression
Flesch-Kincaid Grade Level	Need for Cognition	No	Negative	Readability library
Gunning-Fog	Need for Cognition	No	Negative	Readability library
Automated Readability Index	Need for Cognition	No	Negative	Readability library
Coleman-Liau Index	Need for Cognition	No	Negative	Readability library
Dale-Chall	Need for Cognition	No	Negative	Readability library
Mean syllables per sentence	Need for Cognition	No	Negative	Syllapy library
Number of hapax legomena	Need for Cognition	Yes	Negative	NLTK library
Average sentence length	Need for Cognition	No	Negative	Author created code
Average characters per sentence	Need for Cognition	No	Negative	Author created code
Nouns	Need for Cognitive Closure	Yes	Positive	PoS tagger
Possessive Nouns	Need for Cognitive Closure	Yes	Positive	PoS tagger
Determiners	Need for Cognitive Closure	Yes	Positive	PoS tagger
Proper Nouns	Need for Cognitive Closure	Yes	Positive	PoS tagger
Exclamation mark	Need for Cognitive Closure	Yes	Positive	Regular Expression
First person plural pronouns	Need for Cognitive Closure	Yes	Positive	Regular Expression
Third person plural pronouns	Need for Cognitive Closure	Yes	Positive	Regular Expression
Modal verbs of obligation	Need for Cognitive Closure	Yes	Positive	Regular Expression
Modal verbs of possibility	Need for Cognitive Closure	Yes	Negative	Regular Expression
Adverbs of certainty high	Need for Cognitive Closure	Yes	Positive	Regular Expression
Adverbs of certainty low	Need for Cognitive Closure	Yes	Negative	Regular Expression
Adverbs of frequency high	Need for Cognitive Closure	Yes	Positive	Regular Expression
Money	SDO	Yes	Positive	Regular Expression
Possessive pronouns	SDO	Yes	Positive	Regular Expression
- first person singular				
Possessive pronouns	SDO	Yes	Positive	Regular Expression
- first person plural				
Possessive pronouns	SDO	Yes	Positive	Regular Expression
- third person singular				
Possessive pronouns	SDO	Yes	Positive	Regular Expression
- third person plural				
Comparative adjectives	SDO	Yes	Positive	PoS tagger
Superlative adjectives	SDO	Dominance	Positive	PoS tagger
Emojis	SDO	Dominance	Negative	Emojis library
Smileys	SDO	Dominance	Negative	Regular Expression
Possessive pronouns	SDO	Yes	Negative	Regular Expression
- second person				

<sup>a</sup> TweetNLP and NLTK parts of speech taggers used.

References

Abd, D. H., Sadiq, A. T., & Abbas, A. R. (2020). Classifying political Arabic articles using support vector machine with different feature extraction. In M. I. Khalaf, D. Al-Jumeli, & A. Lisitsa (Eds.), *Applied computing to support industry: Innovation and technology* (pp. 79–94). Communications in Computer and Information Science.

Springer International Publishing. [https://doi.org/10.1007/978-3-030-38752-5\\_7](https://doi.org/10.1007/978-3-030-38752-5_7). Cham, 2020.

ACLU. (2017). Timeline of the Muslim ban. Available at: <https://www.aclu-wa.org/pages/timeline-muslim-ban>. (Accessed 11 June 2021).

Ali, S., & Smith-Miles, K. A. (2006). Improved support vector machine generalization using normalized input space. In *Proceedings of the 19th Australian joint conference on artificial intelligence: Advances in artificial intelligence, Berlin, Heidelberg, 4 december 2006* (pp. 362–371). Springer-Verlag. [https://doi.org/10.1007/11941439\\_40](https://doi.org/10.1007/11941439_40). Al'06.

- Atske, S. (2020). Perceptions of Trump and Biden. In *Pew research center - U.S. Politics & policy*. Available at: <https://www.pewresearch.org/politics/2020/08/13/perceptions-of-trump-and-biden/>. (Accessed 2 December 2022).
- Bakker, B. N., Schumacher, G., Gothreau, C., et al. (2020). Conservatives and liberals have similar physiological responses to threats. *Nature Human Behaviour*, 4(6), 613–621. <https://doi.org/10.1038/s41562-020-0823-z>, 6. Nature Publishing Group.
- Bevendoff, J., Ghanem, B., Giachanou, A., et al. (2020). Shared tasks on authorship analysis at PAN 2020. In J. M. Jose, E. Yilmaz, J. Magalhães, et al. (Eds.), *Advances in information retrieval* (pp. 508–516). Lecture Notes in Computer Science. Springer International Publishing. [https://doi.org/10.1007/978-3-030-45442-5\\_66](https://doi.org/10.1007/978-3-030-45442-5_66). Cham, 2020.
- Bird, S., Loper, E., & Klein, E. (2009). *Natural Language processing with Python*. Newton, Massachusetts, USA: O'Reilly Media Inc.
- Boe, B. (2016). *PRAW: The Python Reddit API Wrapper* (7.6.1) [Computer Software]. <https://github.com/praw-dev/praw>.
- Cacioppo, J. T., & Petty, R. E. (1982). The need for cognition. *Journal of Personality and Social Psychology*, 42(1). <https://doi.org/10.1037/0022-3514.42.1.116>. US: American Psychological Association: 116–131.
- Cacioppo, J., Petty, R., & Morris, K. (1983). Effects of need for cognition on message evaluation, recall, and persuasion. <https://doi.org/10.1037/0022-3514.45.4.805>
- Caplan, D. (2016). Log cabin republicans: GOP platform the 'most anti-LGBT' in party's history. Available at: <https://abcnews.go.com/Politics/log-cabin-republicans-gop-party-platform-anti-lgbt/story?id=40564850>. (Accessed 4 June 2020).
- Chirumbolo, A., Areni, A., & Sensales, G. (2004). Need for cognitive closure and politics: Voting, political attitudes and attributional style. *International Journal of Psychology*, 39(4). <https://doi.org/10.1080/00207590444000005>. Routledge: 245–253.
- Chung, C. K., & Pennebaker, J. W. (2013). Linguistic inquiry and word count (LIWC). In *Applied natural language processing* (pp. 206–229). <https://doi.org/10.4018/978-1-60960-741-8.ch012>
- Cichočka, A., Bilewicz, M., Jost, J. T., et al. (2016). On the grammar of politics— or why conservatives prefer nouns. *Political Psychology*, 37(6), 799–815. Wiley Online Library.
- Clifford, S., Erisen, C., Wendell, D., et al. (2022). Disgust sensitivity and support for immigration across five nations. In *Politics and the life sciences* (pp. 1–16). Cambridge University Press. <https://doi.org/10.1017/pls.2022.6>.
- Cohen, R., & Rutch, D. (2013). Classifying political orientation on twitter: It's not easy. *Proceedings of the International AAAI Conference on Web and Social Media*, 7(1), 91–99, 1.
- Cutler, A. D., Carden, S. W., Dorrough, H. L., et al. (2021). Inferring grandiose narcissism from text: LIWC versus machine learning. *Journal of Language and Social Psychology*, 40(2), 260–276. <https://doi.org/10.1177/0261927X20936309>. SAGE Publications Inc.
- Dahllof, M. (2012). Automatic prediction of gender, political affiliation, and age in Swedish politicians from the wording of their speeches—a comparative study of classifiability. *Literary and Linguistic Computing*, 27(2), 139–153.
- Das, K. G., Patra, B. G., & Naskar, S. K. (2021). Profiling celebrity profession from twitter data. In *2021 international conference on asian language processing (IALP), december 2021* (pp. 207–212). <https://doi.org/10.1109/IALP54817.2021.9675260>
- Diermeier, D., Godbout, J.-F., Yu, B., et al. (2012). Language and ideology in congress. *British Journal of Political Science*, 42(1), 31–55. Cambridge University Press.
- Ellen, J., & Parameswaran, S. (2011). Machine learning for author affiliation within web forums – using statistical techniques on NLP features for online group identification. In *2011 10th international conference on machine learning and applications and workshops, december 2011* (pp. 100–105). <https://doi.org/10.1109/ICMLA.2011.90>
- Erisen, C., Guidi, M., Martini, S., et al. (2021). Psychological correlates of populist attitudes. *Political Psychology*, 42(S1), 149–171. <https://doi.org/10.1111/pops.12768>
- Erisen, C., Redlawsk, D. P., & Erisen, E. (2018). Complex thinking as a result of incongruent information exposure. *American Politics Research*, 46(2), 217–245. <https://doi.org/10.1177/1532673X17725864>. SAGE Publications Inc.
- Gaikwad, M., Ahirrao, S., Phansalkar, S., et al. (2021). Online extremism detection: A systematic literature review with emphasis on datasets, classification techniques, validation methods, and tools. *IEEE Access*, 9, 48364–48404. <https://doi.org/10.1109/ACCESS.2021.3068313>
- Graham, J., Haidt, J., & Nosek, B. A. (2009). Liberals and conservatives rely on different sets of moral foundations. *Journal of Personality and Social Psychology*, 96(5), 1029–1046. <https://doi.org/10.1037/a0015141>
- Gu, F., & Jiang, D. (2021). Prediction of political leanings of Chinese speaking twitter users. arXiv:2110.05723. arXiv. Available at: <http://arxiv.org/abs/2110.05723>. (Accessed 24 September 2022).
- Harnish, R., Bridges, K., & Gump, J. (2018). Predicting economic, social, and foreign policy conservatism: The role of right-wing authoritarianism, social dominance orientation, moral foundations orientation, and religious fundamentalism. *Current Psychology*, 37. <https://doi.org/10.1007/s12144-016-9552-x>
- Hinds, J., & Joinson, A. N. (2018). What demographic attributes do our digital footprints reveal? A systematic review. *PLoS One*, 13(11), 1–40. <https://doi.org/10.1371/journal.pone.0207112>. Public Library of Science.
- Hirst, G., Riabinin, Y., & Graham, J. (2010). Party status as a confound in the automatic classification of political speech by ideology. In *Proceedings of the 10th international conference on statistical analysis of textual data (JADT 2010)*, 2010 (pp. 731–742).
- Holmes, D. I. (1998). The evolution of stylometry in humanities scholarship. *Literary and Linguistic Computing*, 13(3), 111–117. <https://doi.org/10.1093/lc/13.3.111>
- Ho, A. K., Sidanius, J., Kteily, N., et al. (2015). The nature of social dominance orientation: Theorizing and measuring preferences for intergroup inequality using the new SDO<sub>r</sub> scale. *Journal of Personality and Social Psychology*, 109(6), 1003. American Psychological Association.
- Huntington, S. P. (1957). Conservatism as an ideology. *American Political Science Review*, 51(2), 454–473. <https://doi.org/10.2307/1952202> [American Political Science Association, Cambridge University Press].
- Joshi, A., Bhattacharyya, P., & Carman, M. (2016). Political issue extraction model: A novel hierarchical topic model that uses tweets by political and non-political authors. In *Proceedings of the 7th workshop on computational approaches to subjectivity, sentiment and social media analysis, san Diego, California, june 2016* (pp. 82–90). Association for Computational Linguistics. <https://doi.org/10.18653/v1/W16-0415>.
- Jost, J. T., Glaser, J., Kruglanski, A. W., et al. (2003). Political conservatism as motivated social cognition. *Psychological Bulletin*, 129(3), 339. American Psychological Association.
- Kapočiūtė-Dzikiene, J., Utka, A., & Šarkutė, L. (2014). Feature exploration for authorship attribution of Lithuanian parliamentary speeches. In P. Sojka, A. Horák, I. Kopeček, et al. (Eds.), *Text, Speech and dialogue. Lecture notes in computer science* (pp. 93–100). Cham: Springer International Publishing. [https://doi.org/10.1007/978-3-319-10816-2\\_12](https://doi.org/10.1007/978-3-319-10816-2_12).
- Kavuri, K., & Kavitha, M. (2020). A stylistic features based approach for author profiling. In H. Sharma, A. K. S. Pundir, N. Yadav, et al. (Eds.), *Recent trends in communication and intelligent systems* (pp. 185–193). Algorithms for Intelligent Systems. Springer. [https://doi.org/10.1007/978-981-15-0426-6\\_20](https://doi.org/10.1007/978-981-15-0426-6_20). Singapore, 2020.
- Kruglanski, A. W. (1996). Motivated social cognition: Principles of the interface. In *Social psychology: Handbook of basic principles* (pp. 493–520). New York, NY, US: The Guilford Press.
- Ksiazkiewicz, A., Ludeke, S., & Krueger, R. (2016). The role of cognitive style in the link between genes and political ideology. *Political Psychology*, 37(6), 761–776. <https://doi.org/10.1111/pops.12318>
- Lagutina, K., Lagutina, N., Boychuk, E., et al. (2019). A survey on stylometric text features. In *2019 25th conference of open innovations association (FRUCT), november 2019* (pp. 184–195). <https://doi.org/10.23919/FRUCT48121.2019.8981504>
- Lapponi, E., Søyland, M. G., Veldal, E., et al. (2018). The talk of Norway: A richly annotated corpus of the Norwegian parliament, 1998–2016. *Language Resources and Evaluation*, 52(3), 873–893. <https://doi.org/10.1007/s10579-018-9411-5>
- Lins de Holanda Coelho, G. H. P., Hanel, P., & Wolf L, J. (2018). The very efficient assessment of need for cognition: Developing a six-item version. *Assessment*. SAGE Publications Inc, Article 1073191118793208. <https://doi.org/10.1177/1073191118793208>
- Makazhanov, A., & Rafiei, D. (2013). Predicting political preference of Twitter users. In *Social network analysis and mining, Niagara falls, 2013* (p. 193). IEEE. <https://doi.org/10.1007/s13278-014-0193-5>.
- Oberlander, J., & Gill, A. J. (2004). Individual differences and implicit language: Personality, parts-of-speech and pervasiveness. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 26(26). Available at: <https://escholarship.org/uc/item/94c490mq>. (Accessed 16 February 2020).
- Owoputi, O., O'Connor, B., Dyer, C., et al. (2013). Improved part-of-speech tagging for online conversational text with word clusters. In *Proceedings of NAACL 2013, atlanta, GA, USA, 2013* (p. 11).
- Ozduzen, O., Ferenczi, N., Holmes, I., Rosun, N., Liu, K., & Alsayednoor, S. (2021). Stakeholders of (De)-Radicalisation in the UK. (D3.1). Horizon 2020. <https://dradproject.com/wp-content/uploads/2021/06/D.Rad-D3.1-UK.pdf>.
- PAN (2020). PAN shared tasks Available at: <https://pan.webis.de/shared-tasks.html>. (Accessed 7 May 2020).
- Pedregosa, F., Varoquaux, G., Gramfort, A., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12(85), 2825–2830.
- Pennacchiotti, M., & Popescu, A.-M. (2011). Democrats, republicans and starbucks aficionados: User classification in twitter. In *Proceedings of the ACM SIGKDD international conference on knowledge discovery and data mining, 21 august 2011* (pp. 430–438). <https://doi.org/10.1145/2020408.2020477>
- Pennington, J., Socher, R., & Manning, C. (2014). GloVe: Global vectors for word representation. Available at: <https://nlp.stanford.edu/projects/glove/>. (Accessed 13 October 2022).
- Pennycook, G., & Rand, D. G. (2021). Research note: Examining false beliefs about voter fraud in the wake of the 2020 Presidential Election. Harvard Kennedy School Misinformation Review. <https://doi.org/10.37016/mr-2020-51>
- Pranckevičius, T., & Marcinkevičius, V. (2017). Comparison of naive bayes, random forest, decision tree, support vector machines, and logistic regression classifiers for text reviews classification. *Baltic Journal of Modern Computing*, 5(2). <https://doi.org/10.22364/bjmc.2017.5.2.05>
- Pratto, F., Sidanius, J., Stallworth, L. M., et al. (1994). Social dominance orientation: A personality variable predicting social and political attitudes. *Journal of Personality and Social Psychology*, 67(4), 741–763. <https://doi.org/10.1037/0022-3514.67.4.741>
- Roets, A., & Van Hiel, A. (2011). Item selection and validation of a brief, 15-item version of the Need for Closure Scale. *Personality and Individual Differences*, 50(1), 90–94. Elsevier.
- Saad, L. (2022). U.S. Political ideology steady; conservatives, moderates tie. Available at: <https://news.gallup.com/poll/388988/political-ideology-steady-conservatives-moderates-tie.aspx>. (Accessed 25 September 2022).
- Saad, L., Jones, J., & Brennan, M. (2019). Understanding shifts in democratic party ideology. Available at: <https://news.gallup.com/poll/246806/understanding-shifts-democratic-party-ideology.aspx>. (Accessed 4 June 2020).
- Sanz, H., Valim, C., Vegas, E., et al. (2018). SVM-RFE: Selection and visualization of the most relevant features through non-linear kernels. *BMC Bioinformatics*, 19(1), 432. <https://doi.org/10.1186/s12859-018-2451-4>
- Satherley, N., & Sibley, C. G. (2016). A Dual Process Model of attitudes toward immigration: Predicting intergroup and international relations with China.

- International Journal of Intercultural Relations*, 53, 72–82. <https://doi.org/10.1016/j.ijintrel.2016.05.008>
- Shah, K., Patel, H., Sanghvi, D., et al. (2020). A comparative analysis of logistic regression, random forest and KNN models for the text classification. *Augmented Human Research*, 5(1), 12. <https://doi.org/10.1007/s41133-020-00032-0>
- Sinn, J. S., & Hayes, M. W. (2018). Is political conservatism adaptive? Reinterpreting right-wing authoritarianism and social dominance orientation as evolved, sociofunctional strategies. *Political Psychology*, 39(5), 1123–1139. <https://doi.org/10.1111/pops.12475>
- START. (2021). Profiles of individual radicalization in the United States (PIRUS). Available at: <https://www.start.umd.edu/data-tools/profiles-individual-radicalization-united-states-pirus>. (Accessed 12 July 2021).
- Stillwell, D. J., & Kosinski, M. (2012). myPersonality project: Example of successful utilization of online social networks for large-scale social research. *American Psychologist*, 59(2), 93–104.
- Strandberg, T., Olson, J. A., Hall, L., et al. (2020). Depolarizing American voters: Democrats and Republicans are equally susceptible to false attitude feedback. *PLoS One*, 15(2), Article e0226799. <https://doi.org/10.1371/journal.pone.0226799>. Public Library of Science.
- Uenal, F., Sidanius, J., Roozenbeek, J., et al. (2021). Climate change threats increase modern racism as a function of social dominance orientation and ingroup identification. *Journal of Experimental Social Psychology*, 97. <https://doi.org/10.1016/j.jesp.2021.104228>
- Ullah, H., Ahmad, B., Sana, I., et al. (2021). Comparative study for machine learning classifier recommendation to predict political affiliation based on online reviews. *CAA Transactions on Intelligence Technology*, 6(3), 251–264. <https://doi.org/10.1049/cit.12046>
- United Nations. Universal Declaration of Human Rights. United Nations. <https://www.un.org/sites/un2.un.org/files/2021/03/udhr.pdf>.
- van Cranenburgh, A. (2019). *readability: Measure the readability of a given text using surface characteristics (Version 1) [Computer Software]*. <https://github.com/andreascv/readability/>.
- Webster, D. M., & Kruglanski, A. W. (1994). Individual differences in need for cognitive closure. *Journal of Personality and Social Psychology*, 67(6), 1049. American Psychological Association.
- Wilson, M. S., & Sibley, C. G. (2013). Social dominance orientation and right-wing authoritarianism: Additive and interactive effects on political conservatism. *Political Psychology*, 34(2), 277–284. <https://doi.org/10.1111/j.1467-9221.2012.00929.x>
- Yadav, S., & Shukla, S. (2016). Analysis of k-fold cross-validation over hold-out validation on colossal datasets for quality classification. In *2016 IEEE 6th international conference on advanced computing (IACC), February 2016* (pp. 78–83). <https://doi.org/10.1109/IACC.2016.25>
- Yin, W., & Zubiaga, A. (2021). *Towards generalisable hate speech detection: A review on obstacles and solutions*. *arXiv:2102.08886 [cs]*. Available at: <http://arxiv.org/abs/2102.08886>. (Accessed 25 March 2021).
- Yu, B., & Diermeier, D. (2010). A longitudinal study of language and ideology in congress. In *The 68th national conference of Midwest political science association, Chicago, IL, 2010*.
- Yu, B., Kaufmann, S., & Diermeier, D. (2008). Classifying party affiliation from political speech. *Journal of Information Technology & Politics*, 5(1). <https://doi.org/10.1080/19331680802149608>. Routledge: 33–48.
- Zavala, A. G. D., Cislak, A., & Wesolowska, E. (2010). Political conservatism, need for cognitive closure, and intergroup hostility. *Political Psychology*, 31(4), 521–541. <https://doi.org/10.1111/j.1467-9221.2010.00767.x>
- Zmigrod, L., Eisenberg, I. W., Bissett, P. G., et al. (2021). The cognitive and perceptual correlates of ideological attitudes: A data-driven approach. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 376(1822). <https://doi.org/10.1098/rstb.2020.0424>. Royal Society: 20200424.