

# Beyond Binary Labels: Political Ideology Prediction of Twitter Users

**Daniel Preoțiu-Pietro**

Positive Psychology Center  
University of Pennsylvania  
danielpr@sas.upenn.edu

**Daniel J. Hopkins**

Political Science Department  
University of Pennsylvania  
danhop@sas.upenn.edu

**Ye Liu\***

School of Computing  
National University of Singapore  
liuye@comp.nus.edu.sg

**Lyle Ungar**

Computing & Information Science  
University of Pennsylvania  
ungar@cis.upenn.edu

## Abstract

Automatic political preference prediction from social media posts has to date proven successful only in distinguishing between publicly declared liberals and conservatives in the US. This study examines users' political ideology using a seven-point scale which enables us to identify politically moderate and neutral users – groups which are of particular interest to political scientists and pollsters. Using a novel data set with political ideology labels self-reported through surveys, our goal is two-fold: a) to characterize the political groups of users through language use on Twitter; b) to build a fine-grained model that predicts political ideology of unseen users. Our results identify differences in both political leaning and engagement and the extent to which each group tweets using political keywords. Finally, we demonstrate how to improve ideology prediction accuracy by exploiting the relationships between the user groups.

## 1 Introduction

Social media is used by people to share their opinions and views. Unsurprisingly, an important part of the population shares opinions and news related to politics or causes they support, thus offering strong cues about their political preferences and ideologies. In addition, political membership is also predictable purely from one's interests or demographics — it is much more likely for a religious person to be conservative or for a younger person to lean liberal (Ellis and Stimson, 2012).

User trait prediction from text is based on the assumption that language use reflects a user's demographics, psychological states or preferences. Applications include prediction of age (Rao et al., 2010; Flekova et al., 2016b), gender (Burger et al., 2011; Sap et al., 2014), personality (Schwartz et al., 2013; Preoțiu-Pietro et al., 2016), socio-economic status (Preoțiu-Pietro et al., 2015a,b; Liu et al., 2016c), popularity (Lampos et al., 2014) or location (Cheng et al., 2010).

Research on predicting political orientation has focused on methodological improvements (Pennacchiotti and Popescu, 2011) and used data sets with publicly stated dichotomous political orientation labels due to their easy accessibility (Sylwester and Purver, 2015). However, these data sets are not representative samples of the entire population (Cohen and Ruths, 2013) and do not accurately reflect the variety of political attitudes and engagement (Kam et al., 2007).

For example, we expect users who state their political affiliation in their profile description, tweet with partisan hashtags or appear in public party lists to use social media as a means of popularizing and supporting their political beliefs (BarberASA, 2015). Many users may choose not to publicly post about their political preference for various social goals or perhaps this preference may not be strong or representative enough to be disclosed online. Dichotomous political preference also ignores users who do not have a political ideology. All of these types of users are very important for researchers aiming to understand group preferences, traits or moral values (Lewis and Reiley, 2014; Hersh, 2015).

The most common political ideology spectrum in the US is the conservative – liberal (Ellis and Stimson, 2012). We collect a novel data set of Twitter users mapped to this seven-point spectrum which allows us to:

---

\* Work carried out during a research visit at the University of Pennsylvania

1. Uncover the differences in language use between ideological groups;
2. Develop a user-level political ideology prediction algorithm that classifies all levels of engagement and leverages the structure in the political ideology spectrum.

First, using a broad range of language features including unigrams, word clusters and emotions, we study the linguistic differences between the two ideologically extreme groups, the two ideologically moderate groups and between both extremes and moderates in order to provide insight into the content they post on Twitter. In addition, we examine the extent to which the ideological groups in our data set post about politics and compare it to a data set obtained similarly to previous work.

In prediction experiments, we show how accurately we can distinguish between opposing ideological groups in various scenarios and that previous binary political orientation prediction has been oversimplified. Then, we measure the extent to which we can predict the two dimensions of political leaning and engagement. Finally, we build an ideology classifier in a multi-task learning setup that leverages the relationships between groups.<sup>1</sup>

## 2 Related Work

Automatically inferring user traits from their online footprints is a prolific topic of research, enabled by the increasing availability of user generated data and advances in machine learning. Beyond its research oriented goals, user profiling has important industry applications in online marketing, personalization or large-scale audience profiling. To this end, researchers have used a wide range of types of online footprints, including video (Subramanian et al., 2013), audio (Alam and Riccardi, 2014), text (Preoțiuc-Pietro et al., 2015a), profile images (Liu et al., 2016a), social data (Van Der Heide et al., 2012; Hall et al., 2014), social networks (Perozzi and Skiena, 2015; Rout et al., 2013), payment data (Wang et al., 2016) and endorsements (Kosinski et al., 2013).

Political orientation prediction has been studied in two related, albeit crucially different scenarios, as also identified in (Zafar et al., 2016). First, researchers aimed to identify and quantify orientation of words (Monroe et al., 2008), hashtags (Weber et al., 2013) or documents (Iyyer et al., 2014),

or to detect bias (Yano et al., 2010) or impartiality (Zafar et al., 2016) at a document level.

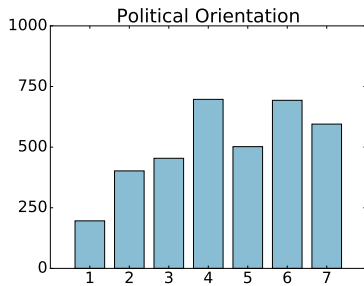
Our study belongs to the second category, where political orientation is inferred at a user-level. All previous studies study labeling US conservatives vs. liberals using either text (Rao et al., 2010), social network connections (Zamal et al., 2012), platform-specific features (Conover et al., 2011) or a combination of these (Pennacchiotti and Popescu, 2011; Volkova et al., 2014), with very high reported accuracies of up to 94.9% (Conover et al., 2011).

However, all previous work on predicting user-level political preferences are limited to a binary prediction between liberal/democrat and conservative/republican, disregarding any nuances in political ideology. In addition, as the focus of the studies is more on the methodological or interpretation aspects of the problem, another downside is that the user labels were obtained in simple, albeit biased ways. These include users who explicitly state their political orientation on user lists of party supporters (Zamal et al., 2012; Pennacchiotti and Popescu, 2011), supporting partisan causes (Rao et al., 2010), by following political figures (Volkova et al., 2014) or party accounts (Sylwester and Purver, 2015) or that retweet partisan hashtags (Conover et al., 2011). As also identified in (Cohen and Ruths, 2013) and further confirmed later in this study, these data sets are biased: most people do not clearly state their political preference online – fewer than 5% according to Priante et al. (2016) – and those that state their preference are very likely to be political activists. Cohen and Ruths (2013) demonstrated that predictive accuracy of classifiers is significantly lower when confronted with users that do not explicitly mention their political orientation. Despite this, their study is limited because in their hardest classification task, they use crowdsourced political orientation labels, which may not correspond to reality and suffer from biases (Flekova et al., 2016a; Carpenter et al., 2016). Further, they still only look at predicting binary political orientation. To date, no other research on this topic has taken into account these findings.

## 3 Data Set

The main data set used in this study consists of 3,938 users recruited through the Qualtrics platform ( $\mathcal{D}_1$ ). Each participant was compensated

<sup>1</sup>Data is available at <http://www.preotiuc.ro>



**Figure 1:** Distribution of political ideology in our data set, from 1 – Very Conservative through 7 – Very Liberal.

with 3 USD for 15 minutes of their time. All participants first answered the same demographic questions (including political ideology), then were directed to one of four sets of psychological questionnaires unrelated to the political ideology question. They were asked to self-report their political ideology on a seven point scale: *Very conservative* (1), *Conservative* (2), *Moderately conservative* (3), *Moderate* (4), *Moderately liberal* (5), *Liberal* (6), *Very liberal* (7). In addition, participants had the option of choosing *Apathetic* and *Other*, which have ambiguous fits on the conservative – liberal spectrum and were removed from our analysis (399 users). We also asked participants to self-report their gender (2322 female, 1205 male, 12 other) and age. Participants were all from the US in order to limit the impact of cultural and political factors. The political ideology distribution in our sample is presented in Figure 1.

We asked users their Twitter handle and downloaded their most recent 3,200 tweets, leading to a total of 4,833,133 tweets. Before adding users to our 3,938 user data set, we performed the following checks to ensure that the Twitter handle was the user’s own: 1) *after* compensation, users were if they were truthful in reporting their handle and if not, we removed their data from analysis; 2) we manually examined all handles marked as verified by Twitter or that had over 2000 followers and eliminated them if they were celebrities or corporate/news accounts, as these were unlikely the users who participated in the survey. This study received approval from the Institutional Review Board (IRB) of the University of Pennsylvania.

In addition, to facilitate comparison to previous work, we also use a data set of 13,651 users with overt political orientation ( $\mathcal{D}_2$ ). We selected popular political figures unambiguously associated with US liberal politics (@SenSanders,

@JoeBiden, @CoryBooker, @JohnKerry) or US conservative politics (@marcorubio, @tedcruz, @RandPaul, @RealBenCarson). Liberals in our set ( $N_l = 7417$ ) had to follow on Twitter all of the liberal political figures and none of the conservative figures. Likewise, conservative users ( $N_c = 6234$ ) had to follow all of the conservative figures and no liberal figures. We downloaded up to 3,200 of each user’s most recent tweets, leading to a total of 25,493,407 tweets. All tweets were downloaded around 10 August 2016.

## 4 Features

In our analysis, we use a broad range of linguistic features described below.

**Unigrams** We use the bag-of-words representation to reduce each user’s posting history to a normalised frequency distribution over the vocabulary consisting of all words used by at least 10% of the users (6,060 words).

**LIWC** Traditional psychological studies use a dictionary-based approach to representing text. The most popular method is based on Linguistic Inquiry and Word Count (LIWC) (Pennebaker et al., 2001), and automatically counts word frequencies for 64 different categories manually constructed based on psychological theory. These include different parts-of-speech, topical categories and emotions. Each user is thereby represented as a frequency distribution over these categories.

**Word2Vec Topics** An alternative to LIWC is to use automatically generated word clusters i.e., groups of words that are semantically and/or syntactically similar. The clusters help reducing the feature space and provides additional interpretability.

To create these groups of words, we use an automatic method that leverages word co-occurrence patterns in large corpora by making use of the distributional hypothesis: similar words tend to co-occur in similar contexts (Harris, 1954). Based on co-occurrence statistics, each word is represented as a low dimensional vector of numbers with words closer in this space being more similar (Deerwester et al., 1990). We use the method from (Preoțiuc-Pietro et al., 2015a) to compute topics using word2vec similarity (Mikolov et al., 2013a,b) and spectral clustering (Shi and Malik, 2000; von Luxburg, 2007) of different sizes (from 30 to 2000). We have tried other alternatives to building clusters: using other word similarities to

generate clusters – such as NPMI (Lampos et al., 2014) or GloVe (Pennington et al., 2014) as proposed in (Preoțiuc-Pietro et al., 2015a) – or using standard topic modelling approached to create soft clusters of words e.g., Latent Dirichlet Allocation (Blei et al., 2003). For brevity, we present experiments with the best performing feature set containing 500 Word2Vec clusters. We aggregate all the words posted in a users’ tweets and represent each user as a distribution of the fraction of words belonging to each cluster.

**Sentiment & Emotions** We hypothesise that different political ideologies differ in the type and amount of emotions the users express through their posts. The most studied model of discrete emotions is the Ekman model (Ekman, 1992; Strapparava and Mihalcea, 2008; Strapparava et al., 2004) which posits the existence of six basic emotions: anger, disgust, fear, joy, sadness and surprise. We automatically quantify these emotions from our Twitter data set using a publicly available crowd-sourcing derived lexicon of words associated with any of the six emotions, as well as general positive and negative sentiment (Mohammad and Turney, 2010, 2013). Using these lexicons, we assign a predicted emotion to each message and then average across all users’ posts to obtain user level emotion expression scores.

**Political Terms** In order to select unigrams pertaining to politics, we assigned the most frequent 12,000 unigrams in our data set to three categories:

- **Political words:** mentions of political terms (234);
- **Political NEs:** mentions of politician proper names out of the political terms (39);
- **Media NEs:** mentions of political media sources and pundits out of the political terms (20).

This coding was initially performed by a research assistant studying political science with good knowledge of US politics and were further filtered and checked by one of the authors.

## 5 Analysis

First, we explore the relationships between language use and political ideological groups within each feature set and pairs of opposing user groups. To illustrate differences between ideological groups we compare the two political extremes (Very Conservative – Very Liberal) and the political moderates (Moderate Conservative – Moderate

Liberal). We further compare outright moderates with a group combining the two political extremes to study if we can uncover differences in political engagement and extremity, regardless of the conservative–liberal leaning.

We use univariate partial linear correlations with age and gender as co-variates to factor out the influence of basic demographics. For example, in  $\mathcal{D}_1$ , users who reported themselves as very conservative are older and more likely males ( $\mu_{age} = 35.1, \text{pct}_{male} = 44\%$ ) than the data average ( $\mu_{age} = 31.2, \text{pct}_{male} = 35\%$ ). Additionally, prior to combining the two ideologically extreme groups, we sub-sampled the larger class (Very Liberal) to match the smaller class (Very Conservative) in age and gender. In the later prediction experiments, we do not perform matching, as this represents useful signal for classification (Ellis and Stimson, 2012). Results with unigrams are presented in Figure 2 and with the other features in Table 1. These are selected using standard statistical significance tests.

### 5.1 Very Conservatives vs. Very Liberals

The comparison between the extreme categories reveals the largest number of significant differences. The unigrams and Word2Vec clusters specific to conservatives are dominated by religion specific terms (‘praying’, ‘god’, W2V-485, W2V-018, W2V-099, L-RELIG), confirming a well-documented relationship (Gelman, 2009) and words describing family relationships (‘uncle’, ‘son’, L-FAMILY), another conservative value (Lakoff, 1997). The emphasis on religious terms among conservatives is consistent with the claim that many Americans associate ‘conservative’ with ‘religious’ (Ellis and Stimson, 2012). Extreme liberals show a tendency to use more adjectives (W2V-075, W2V-110), adverbs (L-ADVERB), conjunctions (L-CONJ) and comparisons (L-COMPARE) which indicate more nuanced and complex posts. Extreme conservatives post tweets higher in all positive emotions than liberals (L-POSEMO, Emot-Joy, Emot-Positive), confirming a previously hypothesised relationship (Napier and Jost, 2008). However, extreme liberals are not associated with posting negative emotions either, only using words that reflect more anxiety (L-ANX), which is related to neuroticism in which the liberals are higher (Gerber et al., 2010).

Political term analysis reveals the partisan terms



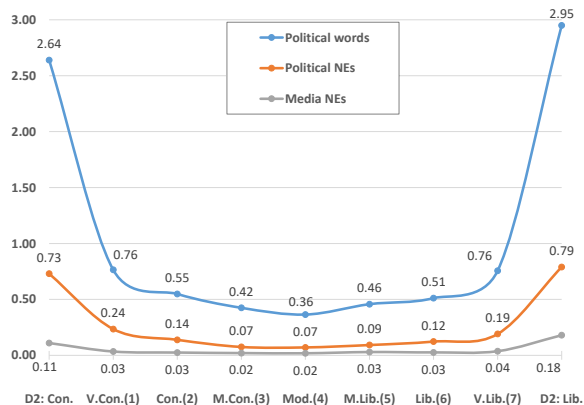
#tcot), while extreme liberals focus on issues ('gay', 'racism', 'feminism', 'transgender'). This perhaps reflects the desire for conservatives on Twitter to identify like-minded individuals, as extreme conservatives are a minority on the platform. Liberals, by contrast, use the platform to discuss and popularize their causes.

## 5.2 Moderate Conservatives vs. Moderate Liberals

Comparing the two sides of moderate users reveals a slightly more nuanced view of the two ideologies. While moderate conservatives still make heavy use of religious terms and express positive emotions (Emot-Joy, L-DRIVES), they also use affiliative language (L-AFFILIATION) and plural pronouns (L-WE). Moderate liberals are identified by very different features compared to their more extreme counterparts. Most striking is the use of swear and sex words (L-SEXUAL, L-ANGER, W2V-316), also highlighted by [Sylwester and Purver \(2015\)](#). Two word clusters relating to British culture (W2V-458) and art (W2V-373) reflect that liberals are more inclined towards arts ([Dollinger, 2007](#)). Statistically significant political terms are very few compared to the previous comparison, probably due to their lower overall usage, which we further investigate later.

## 5.3 Moderates vs. Extremists

Our final comparison looks at outright moderates compared to the two extreme groups combined, as we hypothesise the existence of a difference in overall political engagement. Moderates are not characterized by many features besides a topic of casual words (W2V-098), indicating the heterogeneity of this group of users. However, regardless of their orientation, the ideological extremists stand out from moderates. They use words and word clusters related to political actors (W2V-309), issues (W2V-237) and laws (W2V-296, W2V-288). LIWC analysis uncovers differences in article use (L-ARTICLE) or power words (L-POWER) specific of political tweets. The overall sentiment of these users is negative (Emot-Fear, Emot-Disgust, Emot-Sadness, L-DEATH) compared to moderates. This reveals – combined with the finding from the first comparison – that while extreme conservatives are overall more positive than liberals, both groups share negative expression. Political terms are almost all significantly correlated with the extreme ideological groups,



**Figure 3:** Distribution of political word and entity usage across political categories in % from the total words used. Users from data set  $\mathcal{D}_2$  who are following the accounts of the four political figures are prefixed with D2. The rest of the categories are from data set  $\mathcal{D}_1$ .

confirming the existence of a difference in political engagement which we study in detail next.

## 5.4 Political Terms

Figure 3 presents the use of the three types of political terms across the 7 ideological groups in  $\mathcal{D}_1$  and the two political groups from  $\mathcal{D}_2$ . We notice the following:

- $\mathcal{D}_2$  has a huge skew towards political words, with an average of more than three times more political terms across all three categories than our extreme classes from  $\mathcal{D}_1$ ;
- Within the groups in  $\mathcal{D}_1$ , we observe an almost perfectly symmetrical U-shape across all three types of political terms, confirming our hypothesis about political engagement;
- The difference between 1–2/6–7 is larger than 2–3/5–6. The extreme liberals and conservatives are disproportionately political, and have the potential to give Twitter’s political discussions an unrepresentative, extremist hue ([Fiorenza, 1999](#)). It is also possible, however, that characterizing one as an extreme liberal or conservative indicates as much about her level of political engagement as it does about her placement on a left-right scale ([Converse, 1964](#); [Broockman, 2016](#)).

## 6 Prediction

In this section we build predictive models of political ideology and compare them to data sets obtained using previous work.

## 6.1 Cross-Group Prediction

First, we experiment with classifying between conservatives and liberals across various levels of political engagement in  $\mathcal{D}_1$  and between the two polarized groups in  $\mathcal{D}_2$ . We use logistic regression classification to compare three setups in Table 2 with results measured with ROC AUC as the classes are slightly imbalanced:

- 10-fold cross-validation where training is performed on the same task as the testing (principal diagonal);
- A train–test setup where training is performed on one task (presented in rows) and testing is performed on another (presented in columns);
- A domain adaptation setup (results in brackets) where on each of the 10 folds, the 9 training folds (presented in rows) are supplemented with all the data from a different task (presented in columns) using the EasyAdapt algorithm (Daumé III, 2007) as a proof on concept on the effects of using additional distantly supervised data. Data pooling lead to worse results than EasyAdapt.

Each of the three tasks from  $\mathcal{D}_1$  have a similar number of training samples, hence we do not expect that data set size has any effects in comparing the results across tasks.

The results with both sets of features show that:

- Prediction performance is much higher for  $\mathcal{D}_2$  than for  $\mathcal{D}_1$ , with the more extreme groups in  $\mathcal{D}_1$  being easier to predict than the moderate groups. This confirms that the very high accuracies reported by previous research are an artifact of user label collection and that on regular users, the expected accuracy is much lower (Cohen and Ruths, 2013). We further show that, as the level of political engagement decreases, the classification problem becomes even harder;
- The model trained on  $\mathcal{D}_2$  and Word2Vec word clusters performs significantly worse on  $\mathcal{D}_1$  tasks even if the training data is over 10 times larger. When using political words, the  $\mathcal{D}_2$  trained classifier performs relatively well on all tasks from  $\mathcal{D}_1$ ;
- Overall, using political words as features performs better than Word2Vec clusters in the binary classification tasks;
- Domain adaptation helps in the majority of cases, leading to improvements of up to .03 in AUC (predicting 2v6 supplemented with 3v5 data).

Train	Test			
	1v7	2v6	3v5	D2
1v7	<b>.785</b>	.639 (.681)	.575 (.598)	.705 (.887)
2v6	.729 (.789)	<b>.662</b>	.574 (.586)	.663 (.889)
3v5	.618 (.778)	.617 (.690)	<b>.581</b>	.684 (.887)
D2	.708 (.764)	.627 (.644)	.571 (.574)	<b>.891</b>

(a) Word2Vec 500

Train	Test			
	1v7	2v6	3v5	D2
1v7	<b>.785</b>	.657 (.679)	.589 (.616)	.928 (.976)
2v6	.739 (.773)	<b>.679</b>	.593 (.612)	.920 (.976)
3v5	.727 (.766)	.636 (.670)	.590	.891 (.976)
D2	.766 (.789)	.677 (.683)	<b>.625</b> (.613)	<b>.972</b>

(b) Political Terms

**Table 2:** Prediction results of the logistic regression classification in ROC AUC when discriminating between two political groups across different levels of engagement and both data sets. The binary classifier from data set  $\mathcal{D}_2$  is represented by *D2*, the rest of the categories are from data set  $\mathcal{D}_1$ . Results on the principal diagonal represent 10-fold cross-validation results (training in-domain). Results off-diagonal represent training the classifier from the column and testing on the problem indicated in the row (training out-of-domain). Numbers in brackets indicate performance when the training data was added in the 10-fold cross-validation setup using the EasyAdapt algorithm (domain adaptation). Best results without domain adaptation are in bold, while the best results with domain adaptation are in italics.

## 6.2 Political Leaning and Engagement Prediction

Political leaning (Conservative – Liberal, excluding the Moderate group) can be considered an ordinal variable and the prediction problem framed as one of regression. In addition to the political leaning prediction, based on analysis and previous prediction results, we hypothesize the existence of a separate dimension of political engagement regardless of the partisan side. Thus, we merge users from classes 3–5, 2–6, 1–7 and create a variable with four values, where the lowest value is represented by moderate users (4) and the highest value is represented by either very conservative (1) or very liberal (7) users.

We use a linear regression algorithm with an Elastic Net regularizer (Zou and Hastie, 2005) as implemented in ScikitLearn (Pedregosa et al., 2011). To evaluate our results, we split our data into 10 stratified folds and performed cross-validation on one held-out fold at a time. For all our methods we tune the parameters of our models on a separate validation fold. The overall performance is assessed using Pearson correlation between the set of predicted values and the user-reported score. Results are presented in Table 3.

The same patterns hold when evaluating the results with Root Mean Squared Error (RMSE).

Features	# Feat.	Political Leaning	Political Engagement
Unigrams	6060	.294	.165
LIWC	73	.286	.149
Word2Vec Clusters	500	.300	.169
Emotions	8	.145	.079
Political Terms	234	.256	.169
All (Ensemble)	5	<b>.369</b>	<b>.196</b>

**Table 3:** Pearson correlations between the predictions and self-reported ideologies using linear regression with each feature category and a linear combination of their predictions in a 10-fold cross-validation setup. Political leaning is represented on the 1–7 scale removing the moderates (4). Political engagement is a scale ranging from 4 through 3–5 and 2–6 to 1–7.

The results show that both dimensions can be predicted well above chance, with political leaning being easier to predict than engagement. Word2Vec clusters obtain the highest predictive accuracy for political leaning, even though they did not perform as well in the previous classification tasks. For political engagement, political terms and Word2Vec clusters obtain similar predictive accuracy. This result is expected based on the results from Figure 3, which showed how political term usage varies across groups, and how it is especially dependent on political engagement. While political terms are very effective at distinguishing between two opposing political groups, they can not discriminate as well between levels of engagement within the same ideological orientation. Combining all classifiers’ predictions in a linear ensemble obtains best results when compared to each individual category.

### 6.3 Encoding Class Structure

In our previous experiments, we uncovered that certain relationships exist between the seven groups. For example, extreme conservatives and liberals both demonstrate strong political engagement. Therefore, this class structure can be exploited to improve classification performance. To this end, we deploy the sparse graph regularized approach (Argyriou et al., 2007; Zhou et al., 2011) to encode the structure of the seven classes as a graph regularizer in a logistic regression framework.

In particular, we employed a multi-task learning paradigm, where each task is a one-vs-all classification. Multi-task learning (MTL) is a learning paradigm that jointly learns multiple related

Method	Accuracy
Baseline	19.6%
LR	22.2%
GR–Engagement	24.2%
GR–Leaning	26.2%
GR–Learnt	<b>27.6%</b>

**Table 4:** Experimental results for seven-way classification using multi-task learning (GR–Engagement, GR–Leaning, GR–Learnt) and 500 Word2Vec clusters as features.

tasks and can achieve better generalization performance than learning each task individually, especially when presented with insufficient training samples (Liu et al., 2015, 2016b,d). The group structure is encoded into a matrix  $\mathbf{R}$  which codes the groups which are considered similar. The objective of the sparse graph regularized multi-task learning problem is:

$$\min_{\mathbf{W}, \mathbf{c}} \sum_{t=1}^{\tau} \sum_{i=1}^N \log(1 + \exp(-\mathbf{Y}_{t,i}(\mathbf{W}_{i,t}^T \mathbf{X}_{t,i} + c_t))) + \gamma \|\mathbf{WR}\|_F^2 + \lambda \|\mathbf{W}\|_1,$$

where  $\tau$  is the number of tasks,  $|N|$  the number of samples,  $\mathbf{X}$  the feature matrix,  $\mathbf{Y}$  the outcome matrix,  $\mathbf{W}_{i,t}$  and  $c_t$  is the model for task  $t$  and  $\mathbf{R}$  is the structure matrix.

We define three  $R$  matrices: (1) codes that groups with similar political engagement are similar (i.e. 1–7, 2–6, 3–5); (2) codes that groups from each ideological side are similar (i.e. 1–2, 1–3, 2–3, 5–6, 5–7, 6–7); (3) learnt from the data. Results are presented in Table 4. Regular logistic regression performs slightly better than the majority class baseline, which demonstrates that the 7-class classification is a very hard problem although most miss-classifications are within one ideology point. The graph regularization (GR) improves the classification performance over logistic regression (LR) in all cases, with political leaning based matrix (GR–Leaning) obtaining 2% in accuracy higher than the political engagement one (GR–Engagement) and the learnt matrix (GR–Learnt) obtaining best results.

## 7 Conclusions

This study analyzed user-level political ideology through Twitter posts. In contrast to previous work, we made use of a novel data set where fine-grained user political ideology labels are obtained through surveys as opposed to binary self-reports. We showed that users in our data set are far less



likely to post about politics and real-world fine-grained political ideology prediction is harder and more nuanced than previously reported. We analyzed language differences between the ideological groups and uncovered a dimension of political engagement separate from political leaning.

Our work has implications for pollsters or marketers, who are most interested to identify and persuade moderate users. With respect to political conclusions, researchers commonly conceptualize ideology as a single, left-right dimension similar to what we observe in the U.S. Congress (Ansolabehere et al., 2008; Bafumi and Herron, 2010). Our results suggest a different direction: self-reported political extremity is more an indication of political engagement than of ideological self-placement (Abramowitz, 2010). In fact, only self-reported extremists appear to devote much of their Twitter activity to politics at all.

While our study focused solely on text posted by the user, follow-up work can use other modalities such as images or social network analysis to improve prediction performance. In addition, our work on user-level modeling can be integrated with work on message-level political bias to study how this is revealed across users with various levels of engagement. Another direction of future study will look at political ideology prediction in other countries and cultures, where ideology has different or multiple dimensions.

## Acknowledgments

The authors acknowledge the support of the Templeton Religion Trust, grant TRT-0048. We wish to thank Prof. David S. Rosenblum for supporting the research visit of Ye Liu.

## References

Alan I Abramowitz. 2010. *The Disappearing Center: Engaged Citizens, Polarization, and American Democracy*. Yale University Press.

Firoj Alam and Giuseppe Riccardi. 2014. Predicting Personality Traits using Multimodal Information. In *Workshop on Computational Personality Recognition (WCPR)*. MM, pages 15–18.

Stephen Ansolabehere, Jonathan Rodden, and James M Snyder. 2008. The strength of issues: Using multiple measures to gauge preference stability, ideological constraint, and issue voting. *American Political Science Review* 102(02):215–232.

Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. 2007. Multi-task Feature Learning. In *Advances in Neural Information Processing Systems*. NIPS, pages 41–49.

Joseph Bafumi and Michael C Herron. 2010. Leapfrog Representation and Extremism: A Study of American Voters and their Members in Congress. *American Political Science Review* 104(03):519–542.

Pablo BarberASa. 2015. Birds of the Same Feather Tweet Together: Bayesian Ideal Point Estimation using Twitter Data. *Political Analysis* 23(1):76–91.

David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet Allocation. *Journal of Machine Learning Research* 3:993–1022.

David E Broockman. 2016. Approaches to Studying Policy Representation. *Legislative Studies Quarterly* 41(1):181–215.

D. John Burger, John Henderson, George Kim, and Guido Zarrella. 2011. Discriminating Gender on Twitter. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*. EMNLP, pages 1301–1309.

Jordan Carpenter, Daniel Preoțiuc-Pietro, Lucie Flekova, Salvatore Giorgi, Courtney Hagan, Margaret Kern, Anneke Buffone, Lyle Ungar, and Martin Seligman. 2016. Real Men don’t say ‘Cute’: Using Automatic Language Analysis to Isolate Inaccurate Aspects of Stereotypes. *Social Psychological and Personality Science*.

Zhiyuan Cheng, James Caverlee, and Kyumin Lee. 2010. You are where you Tweet: A Content-Based Approach to Geo-Locating Twitter Users. In *Proceedings of the 19th ACM Conference on Information and Knowledge Management*. CIKM, pages 759–768.

Raviv Cohen and Derek Ruths. 2013. Classifying Political Orientation on Twitter: It’s Not Easy! In *Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media*. ICWSM, pages 91–99.

Michael D Conover, Bruno Gonçalves, Jacob Ratkiewicz, Alessandro Flammini, and Filippo Menczer. 2011. Predicting the Political Alignment of Twitter Users. In *IEEE Third International Conference on Privacy, Security, Risk and Trust (PASSAT) and the IEEE Third International Conference on Social Computing (SocialCom)*. pages 192–199.

Philip E Converse. 1964. The Nature of Belief Systems in Mass Publics. In David Apter, editor, *Ideology and Discontent*, Free Press, New York.

Hal Daumé III. 2007. Frustratingly Easy Domain Adaptation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*. ACL, pages 256–263.

- Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. 1990. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science* 41(6):391–407.
- Stephen J Dollinger. 2007. Creativity and Conservatism. *Personality and Individual Differences* 43(5):1025–1035.
- Paul Ekman. 1992. An Argument for Basic Emotions. *Cognition & Emotion* 6(3-4):169–200.
- Christopher Ellis and James A Stimson. 2012. *Ideology in America*. Cambridge University Press.
- Morris P Fiorina. 1999. Extreme Voices: A Dark Side of Civic Engagement. In Morris P. Fiorina and Theda Skocpol, editors, *Civic engagement in American democracy*, Washington, DC: Brookings Institution Press, pages 405–413.
- Lucie Flekova, Jordan Carpenter, Salvatore Giorgi, Lyle Ungar, and Daniel Preoțiu-Pietro. 2016a. Analyzing Biases in Human Perception of User Age and Gender from Text. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*. ACL, pages 843–854.
- Lucie Flekova, Lyle Ungar, and Daniel Preoțiu-Pietro. 2016b. Exploring Stylistic Variation with Age and Income on Twitter. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*. ACL, pages 313–319.
- Andrew Gelman. 2009. *Red State, Blue State, Rich State, Poor State: Why Americans Vote the Way they Do*. Princeton University Press.
- Alan S Gerber, Gregory A Huber, David Doherty, Conor M Dowling, and Shang E Ha. 2010. Personality and Political Attitudes: Relationships across Issue Domains and Political Contexts. *American Political Science Review* 104(01):111–133.
- Jeffrey A Hall, Natalie Pennington, and Allyn Lueders. 2014. Impression Management and Formation on Facebook: A Lens Model Approach. *New Media & Society* 16(6):958–982.
- Z. Harris. 1954. Distributional Structure. *Word* 10(23):146 – 162.
- Eitan D Hersh. 2015. *Hacking the Electorate: How Campaigns Perceive Voters*. Cambridge University Press.
- Mohit Iyyer, Peter Enns, Jordan Boyd-Graber, and Philip Resnik. 2014. Political Ideology Detection using Recursive Neural Networks. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*. ACL, pages 1113–1122.
- Cindy D Kam, Jennifer R Wilking, and Elizabeth J Zechmeister. 2007. Beyond the Narrow Data base: Another Convenience Sample for Experimental Research. *Political Behavior* 29(4):415–440.
- Michal Kosinski, David Stillwell, and Thore Graepel. 2013. Private Traits and Attributes are Predictable from Digital Records of Human Behavior. *PNAS* 110(15):5802–5805.
- George Lakoff. 1997. *Moral Politics: What Conservatives Know that Liberals Don't*. University of Chicago Press.
- Vasileios Lamos, Nikolaos Aletras, Daniel Preoțiu-Pietro, and Trevor Cohn. 2014. Predicting and Characterising User Impact on Twitter. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*. EACL, pages 405–413.
- Randall A Lewis and David H Reiley. 2014. Online Ads and Offline Sales: Measuring the Effect of Retail Advertising via a Controlled Experiment on Yahoo! *Quantitative Marketing and Economics* 12(3):235–266.
- Leqi Liu, Daniel Preoțiu-Pietro, Zahra Riahi Samani, Mohsen E. Moghaddam, and Lyle Ungar. 2016a. Analyzing Personality through Social Media Profile Picture Choice. In *Proceedings of the Tenth International AAI Conference on Weblogs and Social Media*. ICWSM, pages 211–220.
- Ye Liu, Liqiang Nie, Lei Han, Luming Zhang, and David S Rosenblum. 2015. Action2Activity: Recognizing Complex Activities from Sensor Data. In *Proceedings of the International Joint Conference on Artificial Intelligence*. IJCAI, pages 1617–1623.
- Ye Liu, Liqiang Nie, Li Liu, and David S Rosenblum. 2016b. From Action to Activity: Sensor-based Activity Recognition. *Neurocomputing* 181:108–115.
- Ye Liu, Luming Zhang, Liqiang Nie, Yan Yan, and David S Rosenblum. 2016c. Fortune Teller: Predicting your Career Path. In *Proceedings of the AAI Conference on Artificial Intelligence*. AAI, pages 201–207.
- Ye Liu, Yu Zheng, Yuxuan Liang, Shuming Liu, and David S. Rosenblum. 2016d. Urban Water Quality Prediction Based on Multi-task Multi-view Learning. In *Proceedings of the International Joint Conference on Artificial Intelligence*. IJCAI, pages 2576–2582.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems*. NIPS, pages 3111–3119.
- Tomas Mikolov, Wen tau Yih, and Geoffrey Zweig. 2013b. Linguistic Regularities in Continuous Space Word Representations. In *Proceedings of the 2010 annual Conference of the North American Chapter of the Association for Computational Linguistics*. NAACL, pages 746–751.

- Saif M. Mohammad and Peter D. Turney. 2010. Emotions Evoked by Common Words and Phrases: Using Mechanical Turk to Create an Emotion Lexicon. In *Proceedings of the Workshop on Computational Approaches to Analysis and Generation of Emotion in Text*. NAACL, pages 26–34.
- Saif M. Mohammad and Peter D. Turney. 2013. Crowdsourcing a Word-Emotion Association Lexicon. *Computational Intelligence* 29(3):436–465.
- Burt L Monroe, Michael P Colaresi, and Kevin M Quinn. 2008. Fightin’ Words: Lexical Feature Selection and Evaluation for Identifying the Content of Political Conflict. *Political Analysis* 16(4):372–403.
- Jaime L Napier and John T Jost. 2008. Why are Conservatives Happier than Liberals? *Psychological Science* 19(6):565–572.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine Learning in Python. *JMLR* 12.
- Marco Pennacchiotti and Ana-Maria Popescu. 2011. A Machine Learning Approach to Twitter User Classification. In *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*. ICWSM, pages 281–288.
- James W. Pennebaker, Martha E. Francis, and Roger J. Booth. 2001. *Linguistic Inquiry and Word Count*. Mahway: Lawrence Erlbaum Associates.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. EMNLP, pages 1532–1543.
- Bryan Perozzi and Steven Skiena. 2015. Exact Age Prediction in Social Networks. In *Proceedings of the 24th International Conference on World Wide Web*. WWW, pages 91–92.
- Daniel Preoȃuc-Pietro, Jordan Carpenter, Salvatore Giorgi, and Lyle Ungar. 2016. Studying the Dark Triad of Personality using Twitter Behavior. In *Proceedings of the 25th ACM Conference on Information and Knowledge Management*. CIKM, pages 761–770.
- Daniel Preoȃuc-Pietro, Vasileios Lampos, and Nikolaos Aletras. 2015a. An Analysis of the User Occupational Class through Twitter Content. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*. ACL, pages 1754–1764.
- Daniel Preoȃuc-Pietro, Svitlana Volkova, Vasileios Lampos, Yoram Bachrach, and Nikolaos Aletras. 2015b. Studying User Income through Language, Behaviour and Affect in Social Media. *PLoS ONE* .
- Anna Priante, Djoerd Hiemstra, Tijs van den Broek, Aaqib Saeed, Michel Ehrenhard, and Ariana Need. 2016. #WhoAmI in 160 Characters? Classifying Social Identities Based on Twitter. In *Proceedings of the Workshop on Natural Language Processing and Computational Social Science*. EMNLP, pages 55–65.
- Delip Rao, David Yarowsky, Abhishek Shreevats, and Manaswi Gupta. 2010. Classifying Latent User Attributes in Twitter. In *Proceedings of the 2nd International Workshop on Search and Mining User-generated Contents*. SMUC, pages 37–44.
- Dominic Rout, Daniel Preoȃuc-Pietro, Bontcheva Kalina, and Trevor Cohn. 2013. Where’s @wally: A Classification Approach to Geolocating Users based on their Social Ties. In *Proceedings of the 24th ACM Conference on Hypertext and Social Media*. HT, pages 11–20.
- Maarten Sap, Gregory Park, Johannes C. Eichstaedt, Margaret L. Kern, David J. Stillwell, Michal Kosinski, Lyle H. Ungar, and Hansen Andrew Schwartz. 2014. Developing Age and Gender Predictive Lexica over Social Media. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. EMNLP, pages 1146–1151.
- H Andrew Schwartz, Johannes C Eichstaedt, Margaret L Kern, Lukasz Dziurzynski, Stephanie M Ramones, Megha Agrawal, Achal Shah, Michal Kosinski, David Stillwell, and Martin EP Seligman. 2013. Personality, Gender, and Age in the Language of Social Media: The Open-vocabulary Approach. *PLoS ONE* 8(9).
- Jianbo Shi and Jitendra Malik. 2000. Normalized Cuts and Image Segmentation. *Transactions on Pattern Analysis and Machine Intelligence* 22(8):888–905.
- Carlo Strapparava and Rada Mihalcea. 2008. Learning to Identify Emotions in Text. In *Proceedings of the 2008 ACM Symposium on Applied Computing*. pages 1556–1560.
- Carlo Strapparava, Alessandro Valitutti, et al. 2004. WordNet Affect: an Affective Extension of WordNet. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation*. volume 4 of *LREC*, pages 1083–1086.
- Ramanathan Subramanian, Yan Yan, Jacopo Staiano, Oswald Lanz, and Nicu Sebe. 2013. On the Relationship between Head Pose, Social Attention and Personality Prediction for Unstructured and Dynamic Group Interactions. In *Proceedings of the 15th ACM on International Conference on Multimodal Interaction*. ICMI, pages 3–10.
- Karolina Sylwester and Matthew Purver. 2015. Twitter Language Use Reflects Psychological Differences between Democrats and Republicans. *PLoS ONE* 10(9).

- Brandon Van Der Heide, Jonathan D D'Angelo, and Erin M Schumaker. 2012. The Effects of Verbal versus Photographic Self-presentation on Impression Formation in Facebook. *Journal of Communication* 62(1):98–116.
- Svitlana Volkova, Glen Coppersmith, and Benjamin Van Durme. 2014. Inferring User Political Preferences from Streaming Communications. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*. ACL, pages 186–196.
- Ulrike von Luxburg. 2007. A Tutorial on Spectral Clustering. *Statistics and Computing* 17(4):395–416.
- Pengfei Wang, Jiafeng Guo, Yanyan Lan, Jun Xu, and Xueqi Cheng. 2016. Your Cart tells You: Inferring Demographic Attributes from Purchase Data. In *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*. WSDM, pages 173–182.
- Ingmar Weber, Venkata Rama Kiran Garimella, and Asmelash Teka. 2013. Political Hashtag Trends. In *European Conference on Information Retrieval*. ECIR, pages 857–860.
- Tae Yano, Philip Resnik, and Noah A Smith. 2010. Shedding (a Thousand Points of) Light on Biased Language. In *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk*. NAACL, pages 152–158.
- Muhammad Bilal Zafar, Krishna P Gummadi, and Cristian Danescu-Niculescu-Mizil. 2016. Message Impartiality in Social Media Discussions. In *Proceedings of the Tenth International AAAI Conference on Weblogs and Social Media*. ICWSM, pages 466–475.
- Faiyaz Al Zamal, Wendy Liu, and Derek Ruths. 2012. Homophily and Latent Attribute Inference: Inferring Latent Attributes of Twitter Users from Neighbors. In *Proceedings of the Sixth International AAAI Conference on Weblogs and Social Media*. ICWSM, pages 387–390.
- Jiayu Zhou, Jianhui Chen, and Jieping Ye. 2011. MAL-SAR: Multi-Task Learning via Structural Regularization. *Arizona State University* .
- Hui Zou and Trevor Hastie. 2005. Regularization and Variable Selection via the Elastic Net. *Journal of the Royal Statistical Society, Series B* .