

# Detecting Latent Ideology in Expert Text: Evidence From Academic Papers in Economics

Zubin Jelveh<sup>1</sup>, Bruce Kogut<sup>2</sup>, and Suresh Naidu<sup>3</sup>

<sup>1</sup>Dept. of Computer Science & Engineering, New York University

<sup>2</sup>Columbia Business School and Dept. of Sociology, Columbia University

<sup>3</sup>Dept. of Economics and SIPA, Columbia University

`zj292@nyu.edu`, `bruce.kogut@columbia.edu`, `sn2430@columbia.edu`

## Abstract

Previous work on extracting ideology from text has focused on domains where expression of political views is expected, but it's unclear if current technology can work in domains where displays of ideology are considered inappropriate. We present a supervised ensemble  $n$ -gram model for ideology extraction with topic adjustments and apply it to one such domain: research papers written by academic economists. We show economists' political leanings can be correctly predicted, that our predictions generalize to new domains, and that they correlate with public policy-relevant research findings. We also present evidence that unsupervised models can under-perform in domains where ideological expression is discouraged.

## 1 Introduction

Recent advances in text mining demonstrate that political ideology can be predicted from text – often with great accuracy. Standard experimental settings in this literature are ones where ideology is explicit, such as speeches by American politicians or editorials by Israeli and Palestinian authors. An open question is whether ideology can be detected in arenas where it is strongly discouraged. A further consideration for applied researchers is whether these tools can offer insight into questions of import for policymakers. To address both of these issues, we examine one such domain that is both policy-relevant and where ideology is not overtly expressed: research papers written by academic economists.

Why economics? Economic ideas are important for shaping policy by influencing the public debate and setting the range of expert opinion on various policy options (Rodrik, 2014). Economics also

views itself as a science (Chetty, 2013) carefully applying rigorous methodologies and using institutionalized safe-guards such as peer review. The field's most prominent research organization explicitly prohibits researchers from making policy recommendations in papers that it releases (National Bureau of Economic Research, 2010). Despite these measures, economics' close proximity to public policy decisions have led many to see it as being driven by ideology (A.S., 2010). Does this view of partisan economics have any empirical basis?

To answer the question of whether economics is politicized or neutral, we present a supervised ensemble  $n$ -gram model of ideology extraction with topic adjustments.<sup>1</sup> Our methodology is most closely related to Taddy (2013) and Gentzkow and Shapiro (2010), the latter of which used  $\chi^2$  tests to find phrases most associated with ideology as proxied by the language of U.S. Congresspersons. We improve on this methodology by accounting for ideological word choice within topics and incorporating an ensemble approach that increases predictive accuracy. We also motivate the need to adjust for topics even if doing so does not improve accuracy (although it does in this case). We further provide evidence that fully unsupervised methods (Mei et al., 2007; Lin et al., 2008; Ahmed and Xing, 2010; Paul and Girju, 2010; Eisenstein et al., 2011; Wang et al., 2012) may encounter difficulties learning latent ideological aspects when those aspects are not first order in the data.

Our algorithm is able to correctly predict the ideology of 69.2% of economists in our data purely from their academic output. We also show that our predictions generalize and are predictors of responses by a panel of top economists on issues of economic importance. In a companion paper (Jelveh et al., 2014), we further show that

<sup>1</sup>Grimmer and Stewart (2013) provide an overview of models used for ideology detection.

predicted ideologies are significantly correlated to economists' research findings. The latter result shows the relevance and applicability of these tools beyond the task of ideology extraction.

## 2 Data

### Linking Economists to Their Political Activity:

We obtain the member directory of the American Economics Association (AEA) and link it to two datasets: economists' political campaign contributions and petition signing activities. We obtain campaign contribution data from the Federal Election Commission's website and petition signing data from Hedengren et al. (2010). From this data, we construct a binary variable to indicate the ground-truth ideologies of economists. See our companion paper (Jelveh et al., 2014) for further details on the construction of this dataset. Revealed ideology through contributions and petitions is largely consistent. Of 441 economists appearing in both datasets, 83.4% showed agreement between contributions and petitions. For the final dataset of ground-truth authors we include all economists with campaign contributions and/or petition signatures, however, we drop those economists whose ideologies were different across the contribution and petition datasets. Overall, 60% of 2,204 economists with imputed ideologies in this final dataset are left-leaning while 40% lean rightwards.

**Economic Papers Corpus:** To create our corpus of academic writings by economists, we collect 17,503 working papers from NBER's website covering June 1973 to October 2011. We also obtained from JSTOR the fulltext of 62,888 research articles published in 93 journals in economics for the years 1991 to 2008. Combining the set of economists and papers leaves us with 2,171 authors with ground truth ideology and 17,870 papers they wrote. From the text of these papers we create  $n$ -grams of length two through eight. While  $n$ -grams greater than three words in length are uncommon, Margolin et al. (2013) demonstrate that ideological word choice can be detected by longer phrases. To capture other expressions of ideology not revealed in adjacent terms, we also include skipgrams of length two by combining non-adjacent terms that are three to five words apart. We remove phrases used by fewer than five authors.

**Topic Adjustments:** Table 1 presents the top

20 most conservative and liberal bigrams ranked by  $\chi^2$  scores from a Pearson's test of independence between phrase usage by left- and right-leaning economists. It appears that top ideological phrases are related to specific research subfields. For example, right-leaning terms 'free\_bank', 'stock\_return', and 'feder\_reserv' are related to finance and left-leaning terms 'mental\_health', 'child\_care', and 'birth\_weight' are related to health care. This observation leads us to ask: Are apparently ideological phrases merely a by-product of an economist's research interest rather than reflective of true ideology?

To see why this is a critical question, consider that ideology has both direct and indirect effects on word choice, the former of which is what we wish to capture. The indirect pathway is through topic: ideology may influence the research area an economist enters into, but not the word choice within that area. In that case, if more conservative economists choose macroeconomics, the observed correlation between macro-related phrases and right-leaning ideology would be spurious. The implication is that accounting for topics may not necessarily improve performance but provide evidence to support an underlying model of how ideology affects word choice. Therefore, to better capture the direct effect of ideology on phrase usage we adjust our predictions by topic by creating mappings from papers to topics. For a **topic mapping**, we predict economists' ideologies from their word choice *within* each topic and combine these results to form an overall prediction. We compare different supervised and unsupervised topic mappings and assess their predictive ability.

To create supervised topic mappings, we take advantage of the fact that economics papers are manually categorized by the Journal of Economic Literature (JEL). These codes are hierarchical indicators of an article's subject area. For example, the code C51 can be read, in increasing order of specificity, as Mathematical and Quantitative Methods (C), Econometric Modeling (C5), Model Construction and Estimation (C51). We construct two sets of topic mappings: **JEL1** derived from the 1st-level codes (e.g. C) and **JEL2** derived from the 2nd-level codes (e.g. C5). The former covers broad areas (e.g. macroeconomics, microeconomics, etc.) while the latter contains more refined ones (e.g. monetary policy, firm behavior, etc.).

For unsupervised mappings, we run Latent

Left-Leaning Bigrams	Right-Leaning Bigrams
mental_health	public_choic
post_keynesian	stock_return
child_care	feder_reserv
labor_market	yes_yes
health_care	market_valu
work_time	journal_financi
keynesian_econom	bank_note
high_school	money_suppli
polici_analys	free_bank
analys_politiqu	liquid_effect
politiqu_vol	journal_financ
birth_weight	median_voter
labor_forc	law_econom
journal_post	vote_share
latin_america	war_spend
mental_ill	journal_law
medic_care	money_demand
labour_market	gold_reserv
social_capit	anna_j
singl_mother	switch_cost

Table 1: Top 20 bigrams and trigrams.

Dirichlet Allocation (Blei et al., 2003) on our corpus. We use 30, 50, and 100 topics to create **LDA30**, **LDA50**, and **LDA100** topic mappings. We use the topic distributions estimated by LDA to assign articles to topics. A paper  $p$  is assigned to a topic  $t$  if the probability that  $t$  appears in  $p$  is greater than 5%. While 5% might seem to be a lower threshold, the topic distributions estimated by LDA tend to be sparse. For example, even with 50 topics to ‘choose’ from in **LDA50** and a threshold of 5%, 99.5% of the papers would be assigned to five or fewer topics. This compares favorably with JEL2 codings where 98.8% of papers have five or fewer topics.

### 3 Algorithm

There are two components to our topic-adjusted algorithm for ideology prediction. First, we focus on  $n$ -grams and skipgrams that are most correlated with ideology in the training data. For each topic within a topic mapping, we count the total number of times each phrase is used by all left- and all right-leaning economists. Then, we compute Pearson’s  $\chi^2$  statistic and associated p-values and keep phrases with  $p \leq 0.05$ . As an additional filter, we split the data into ten folds and perform the

$\chi^2$  test within each fold. For each topic, we keep phrases that are consistently ideological across all folds. This greatly reduces the number of ideological phrases. For **LDA50**, the mean number of ideological phrases per topic before the cross validation filter is 12,932 but falls to 963 afterwards.

With the list of ideological phrases in hand, the second step is to iterate over each topic and predict the ideologies of economists in our test set. To compute the predictions we perform partial least squares (PLS): With our training data, we construct the standardized frequency matrix  $\mathbf{F}_{t,train}$  where the  $(e, p)$ -th entry is the number of times economist  $e$  used partisan phrase  $p$  across all of  $e$ ’s papers in  $t$ . This number is divided by the total number of phrases used by  $e$  in topic  $t$ . For papers with multiple authors, each author gets same count of phrases. About 5% of the papers in our dataset are written by authors with differing ideologies. We do not treat these differently. Columns of  $\mathbf{F}_{t,train}$  are standardized to have unit variance. Let  $\mathbf{y}$  be the vector of ground-truth ideologies, test set ideologies are predicted as follows:

- 1) Compute  $\mathbf{w} = Corr(\mathbf{F}_{t,train}, \mathbf{y})$ , the correlations between each phrase and ideology
- 2) Project to one dimension:  $\mathbf{z} = \mathbf{F}_{t,train} \mathbf{w}$
- 3) Regress ideology,  $\mathbf{y}$ , on the constructed variable  $\mathbf{z}$ :  $\mathbf{y} = b_1 \mathbf{z}$
- 4) Predict ideology  $\hat{y}_e$  of new economist by  $\hat{y}_e = b_1 \tilde{\mathbf{f}}_e' \mathbf{w}$ , ( $\tilde{\mathbf{f}}_e$  is scaled frequency vector)

To avoid over-fitting we introduce an ensemble element: For each  $t$ , we sample from the list of significant  $n$ -grams in  $t$  and sample with replacement from the authors who have written in  $t$ .<sup>2</sup> PLS is performed on this sample data 125 times. Each PLS iteration can be viewed as a vote on whether an author is left- or right-leaning. We calculate the vote as follows. For each iteration, we predict the ideologies of economists in the *training data*. We find the threshold  $f$  that minimizes the distance between the true and false positive rates for the current iteration and the same rates for the perfect classifier: 1.0 and 0.0, respectively. Then, an author in the test set is voted left-leaning if  $y_{t,test} \leq f$  and right-leaning otherwise.

For a given topic mapping, our algorithm returns a three-dimensional array with the  $(e, t, c)$ -th entry representing the number of votes economist  $e$  received in topic  $t$  for ideology  $c$  (left- or right-

<sup>2</sup>The number of phrases sampled each iteration is the square root of the number of ideological phrases in the topic.

leaning). To produce a final prediction, we sum across the second dimension and compute ideology as the percentage of right-leaning votes received across all topics within a topic-mapping. Therefore, ideology values closer to zero are associated with a left-leaning ideology and values closer to one are associated with a rightward lean.

To recap, we start with a topic mapping and then for each topic run an ensemble algorithm with PLS at its core.<sup>3</sup> The output for each topic is a set of votes. We sum across topics to compute a final prediction for ideology.

#### 4 Validation and Results

We split our ground-truth set of 2,171 authors into training (80%) and test sets (20%) and compute predictions as above. As our data exhibits skew with 1.5 left-leaning for every right-leaning economist, we report the area under the curve (AUC) which is robust to class skew (Fawcett, 2006). It's worth noting that a classifier that randomly predicts a liberal economist 60% of the time would have an AUC of 0.5. To compare our model with fully unsupervised methods, we also include results from running the Topic-Aspect Model (TAM) (Paul and Girju, 2010) on our data. TAM decomposes documents into two components: one affecting topics and one affecting a latent aspect that influences all topics in a similar manner. We run TAM with 30 topics and 2 aspects (**TAM2/30**). We follow Paul and Girju and use the learned topic and aspect distributions as training data for a SVM.<sup>4</sup>

Columns 2 to 4 from Table 2 show that our models' predictions have a clear association with ground-truth ideology. The LDA topic mappings outperform the supervised mappings as well as a model that does not adjust for topics (**NoTopic**). Perhaps not surprisingly, TAM does not perform well in our domain. A drawback of unsupervised methods is that the learned aspects may not be related to ideology but some other hidden factor.

For further insight into how well our model generalizes, we use data from Gordon and Dahl (2013) to compare our predictions to potentially ideological responses of economists on a survey

<sup>3</sup>Other predictions algorithms could be dropped in for PLS. Logistic regression and SVM produced similar results.

<sup>4</sup>Authors are treated as documents. TAM is run for 1,000 iterations with the following priors:  $\alpha = 1.0$ ,  $\beta = 1.0$ ,  $\gamma_0 = 1$ ,  $\gamma_1 = 1$ ,  $\delta_0 = 20$ ,  $\delta_1 = 80$ .

(1)	(2)	(3)	(4)
Topic	Accu-	Corr. w/	AUC
Map	racy(%)	Truth	
LDA50	<b>69.2</b>	<b>0.381</b>	<b>0.719</b>
LDA100	66.3	0.364	0.707
LDA30	65.0	0.313	0.674
NoTopic	63.9	0.290	0.672
JEL1	61.0	0.263	0.647
JEL2	61.8	0.240	0.646
TAM2/30	61.5	0.228	0.580

Table 2: Model comparisons

	(1)	(2)	(3)
LDA50	1.814***	2.457***	2.243***
Log-Lik.	-1075.0	-758.7	-740.6
JEL1	1.450***	2.128***	1.799***
Log-Lik.	-1075.3	-757.4	-740.5
No Topic	0.524***	0.659***	0.824***
Log-Lik.	-1075.3	-760.5	-741.0
Question	No	Yes	Yes
Demog./Prof.	No	No	Yes
Observations	715	715	715
Individuals	39	39	39

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 3: IGM correlations. Column (1) shows results of regression of response on predicted ideology. Column (2) adds question dummies. Column (3) adds demographic and professional variables.

conducted by the University of Chicago.<sup>5</sup> Each survey question asks for an economists opinion on an issue of political relevance such as minimum wages or tax rates. For further details on the data see Gordon and Dahl. Of importance here is that Gordon and Dahl categorize 22 questions where agreement (disagreement) with the statement implies belief a conservative (liberal) viewpoint.

To see if our predicted ideologies are correlated with survey responses, we run an ordered-logistic regression (McCullagh, 1980). Survey responses are coded with the following order: Strongly Disagree, Disagree, Uncertain, Agree, Strongly Agree. We regress survey responses onto predicted ideologies. We also include question-level dummies and explanatory variables for a re-

<sup>5</sup><http://igmchicago.org>

spondent's gender, year of Ph.D., Ph.D. university, NBER membership, and experience in federal government. Table 3 shows the results of these regressions for three topic mappings. The correlation between our predictions and survey respondents are all strongly significant.

One way to interpret these results is to compare the change in predicted probability of providing an Agree or Strongly Agree answer (agreeing with the conservative view point) if we change predicted ideology from most liberal to most conservative. For **NoTopic**, this predicted probability is 35% when ideology is set to most liberal and jumps to 73.7% when set to most conservative. This difference increases for topic-adjusted models. For **LDA50**, the probability of a conservative answer when ideology is set to most liberal is 14.5% and 93.8% for most conservative.

Figure 1 compares the predicted probabilities of choosing different answers when ideology is set to most liberal and most conservative. Our topic-adjusted models suggest that the most conservative economists are much more likely to strongly agree with a conservative response than for the most liberal economists to strongly agree with a liberal response. It is worthwhile to note from the small increase in log-likelihood in Table 3 when controls are added, suggesting that our ideology scores are much better predictors of IGM responses than demographic and professional controls.

## 5 Conclusions and Future Work

We've presented a supervised methodology for extracting political sentiment in a domain where it's discouraged and shown how it even predicts the partisanship calculated from completely unrelated IGM survey data. In a companion paper (Jelveh et al., 2014) we further demonstrate how this tool can be used to aid policymakers in de-biasing research findings. When compared to domains where ideological language is expected, our predictive ability is reduced. Future work should disentangle how much this difference is due to modeling decisions and limitations versus actual absence of ideology. Future works should also investigate how fully unsupervised methods can be extended to match our performance.

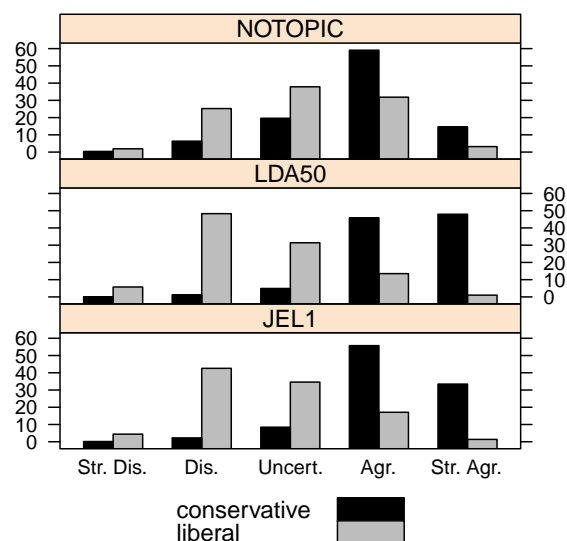


Figure 1: The predicted probability of agreeing with a conservative response when ideology is set to most liberal (gray) and most conservative (black).

## Acknowledgement

This work was supported in part by the NSF (under grant 0966187). The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of any of the sponsors.

## References

- Amr Ahmed and Eric P. Xing. 2010. Staying informed: supervised and semi-supervised multi-view topical analysis of ideological perspective. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 1140–1150. Association for Computational Linguistics.
- A.S. 2010. Is economics a right-wing conspiracy? *The Economist*, August.
- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022.
- Raj Chetty. 2013. Yes, economics is a science. *The New York Times*, October.
- Jacob Eisenstein, Amr Ahmed, and Eric P. Xing. 2011. Sparse additive generative models of text. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 1041–1048.
- T. Fawcett. 2006. An introduction to ROC analysis. *Pattern recognition letters*, 27(8):861–874.
- Matthew Gentzkow and Jesse M. Shapiro. 2010. What drives media slant? evidence from U.S. daily newspapers. *Econometrica*, 78(1):35–71.
- Roger Gordon and Gordon B Dahl. 2013. Views among economists: Professional consensus or point-counterpoint? *American Economic Review*, 103(3):629–635, May.
- J. Grimmer and B. M. Stewart. 2013. Text as data: The promise and pitfalls of automatic content analysis methods for political texts. *Political Analysis*, 21(3):267–297, January.
- David Hedengren, Daniel B. Klein, and Carrie Milton. 2010. Economist petitions: Ideology revealed. *Econ Journal Watch*, 7(3):288–319.
- Zubin Jelveh, Bruce Kogut, and Suresh Naidu. 2014. Political language in economics.
- Wei-Hao Lin, Eric Xing, and Alexander Hauptmann. 2008. A joint topic and perspective model for ideological discourse. In *Machine Learning and Knowledge Discovery in Databases*, pages 17–32. Springer.
- Drew Margolin, Yu-Ru Lin, and David Lazer. 2013. Why so similar?: Identifying semantic organizing processes in large textual corpora.
- Peter McCullagh. 1980. Regression models for ordinal data. *Journal of the royal statistical society. Series B (Methodological)*, pages 109–142.
- Qiaozhu Mei, Xu Ling, Matthew Wondra, Hang Su, and ChengXiang Zhai. 2007. Topic sentiment mixture: modeling facets and opinions in weblogs. In *Proceedings of the 16th international conference on World Wide Web*, pages 171–180. ACM.
- National Bureau of Economic Research. 2010. Amended and restated by-laws of national bureau of economic research, inc.
- Michael Paul and Roxana Girju. 2010. A two-dimensional topic-aspect model for discovering multi-faceted topics. *Urbana*, 51.
- Dani Rodrik. 2014. When ideas trump interests: Preferences, worldviews, and policy innovations. *Journal of Economic Perspectives*, 28(1):189–208, February.
- Matt Taddy. 2013. Multinomial inverse regression for text analysis. *Journal of the American Statistical Association*, 108.
- William Yang Wang, Elijah Mayfield, Suresh Naidu, and Jeremiah Dittmar. 2012. Historical analysis of legal opinions with a sparse mixed-effects latent variable model. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 740–749. Association for Computational Linguistics.