



Predicting Political Ideology Using Campaign Finance Data

Keri A. McKiernan and Joe A. Napoli

Department of Chemistry, Stanford University

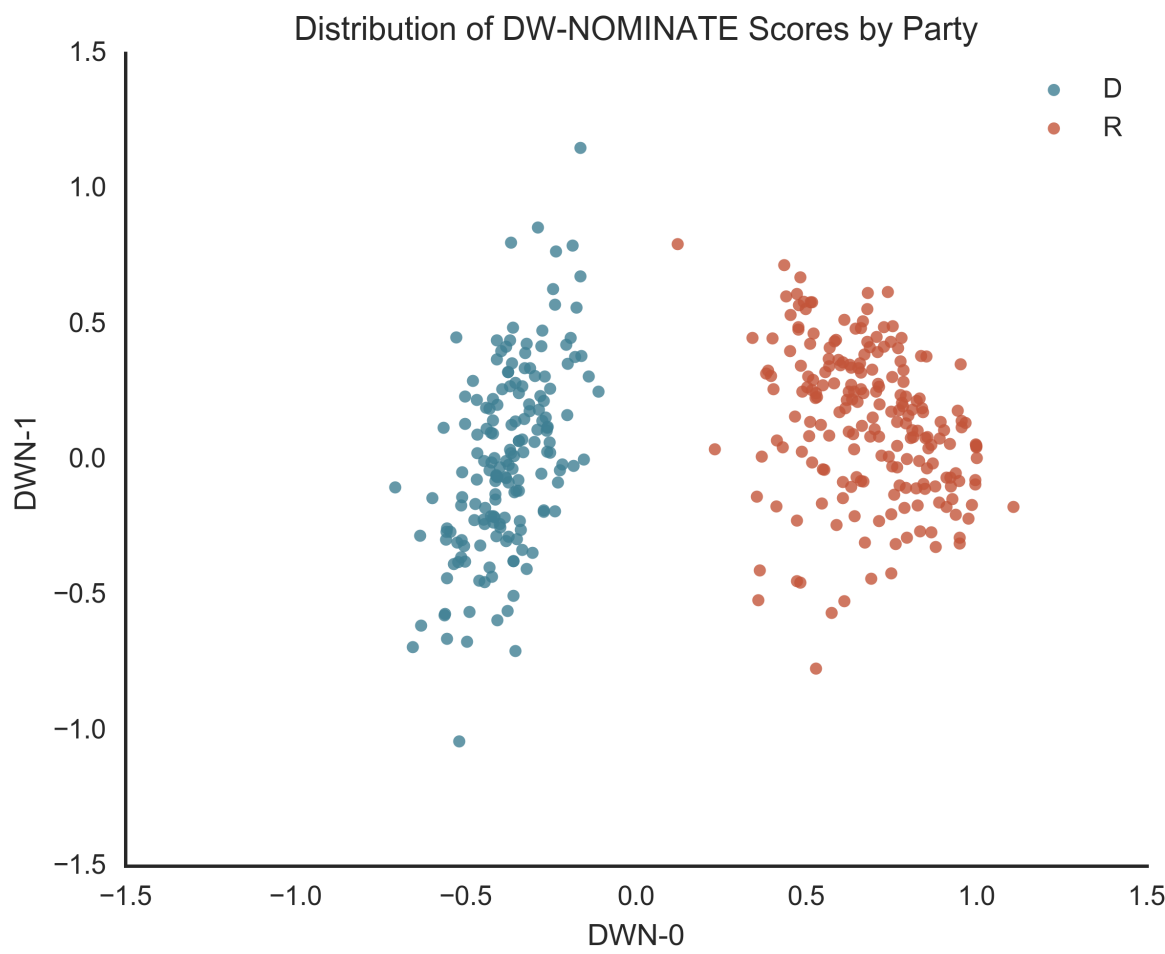


Background

The research of Poole and Rosenthal has focused on quantifying political ideology [2] via the DW-NOMINATE method. For this method, ‘ideal point’ coordinates were obtained for legislators by maximizing the log likelihood function

$$\mathcal{L} = \sum_{t=1}^T \sum_{i=1}^{p_t} \sum_{j=1}^{q_t} \sum_{\tau=1}^2 C_{ij\tau t} \ln P_{ij\tau t}$$

where $P_{ij\tau t}$ is the probability of voting for choice τ and $C_{ij\tau t} = 1$ if that probability accurately predicts the vote [1]. Indices j , i , and t sum over roll call votes, legislators, and legislative sessions, respectively. A common interpretation of the first coordinate is that it reflects the divide between the Republican and Democratic parties, whereas the second coordinate is more highly correlated with intra-party division. The distribution of these scores is illustrated below.



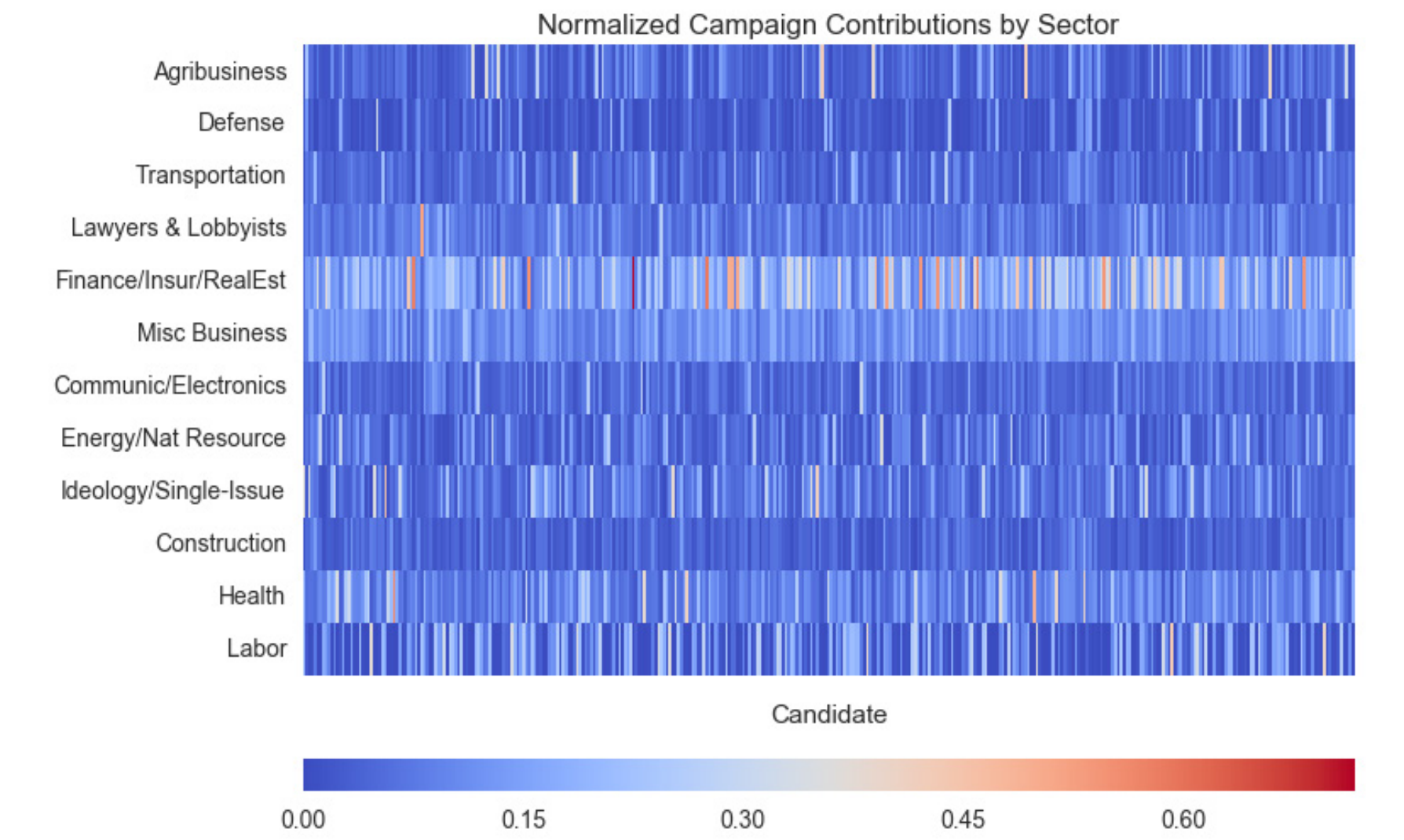
Here we aim to use machine learning methods to predict DW-NOMINATE scores using a candidate’s campaign finance information. It would be useful to be able to predict the ideal point of a candidate even before they have established a congressional voting record. Furthermore, the ability to do so would help elucidate a relationship between monetary contributions to candidates and the voting patterns those contributions may effect.

References

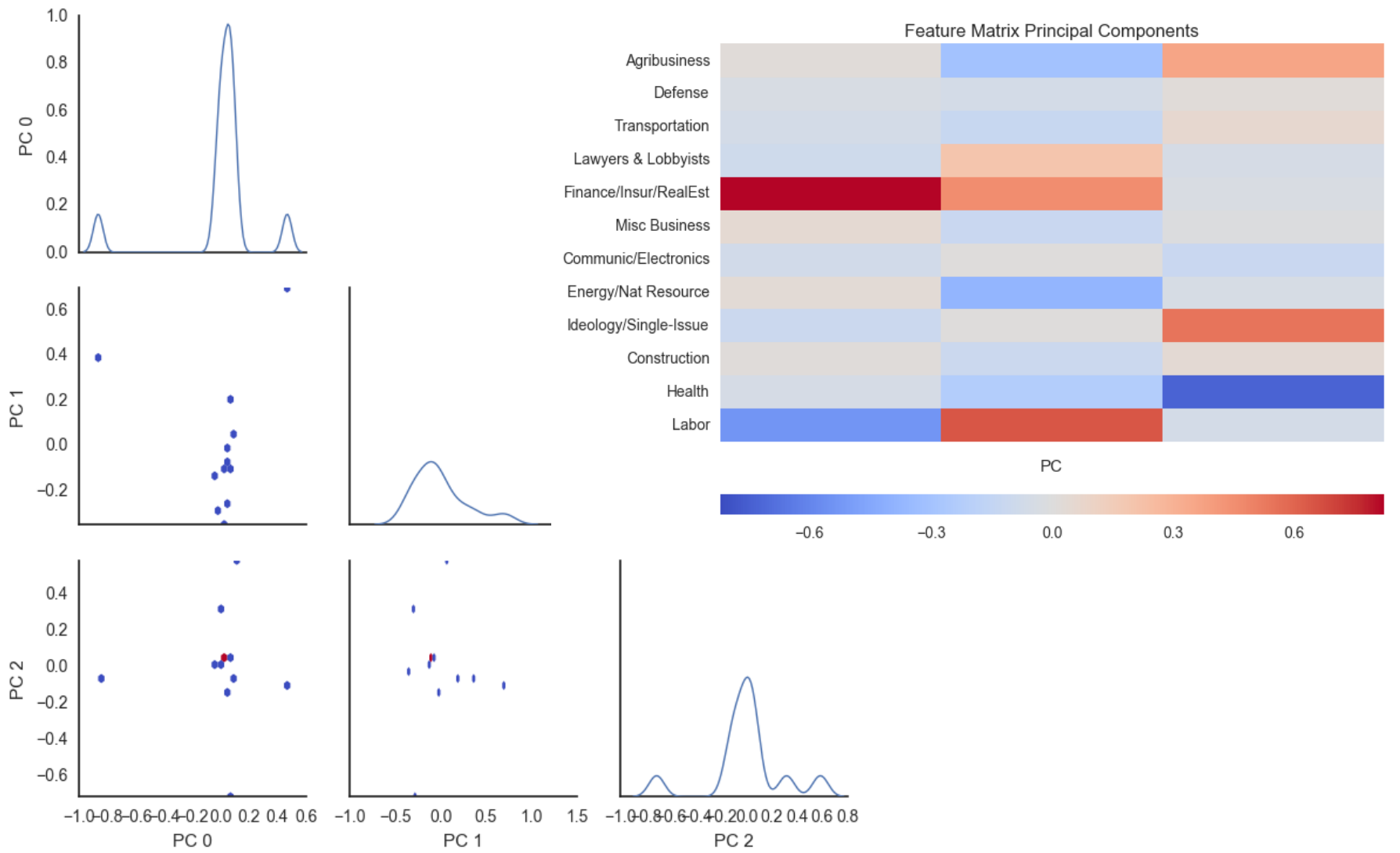
[1] Royce Carroll et al. “Measuring bias and uncertainty in DW-NOMINATE ideal point estimates via the parametric bootstrap”. In: Political Analysis 17.3 (2009), pp. 261–275.
[2] K.T. Poole and H. Rosenthal. Congress: A Political-Economic History of Roll Call Voting. Oxford University Press, 2000.
[3] Scikit-learn: Machine Learning in Python, Pedregosa et al., JMLR 12, pp. 2825-2830, 2011.

Training

Campaign finance data for a set of candidates (as determined by sampling from the DW-NOMINATE scores) was obtained using the Open Secrets python API [4]. Of the available data, we chose to partition finance streams by sector. For each sector, individual and PAC contributions were combined and normalized on a per candidate basis. A heatmap of the total feature set is illustrated below.



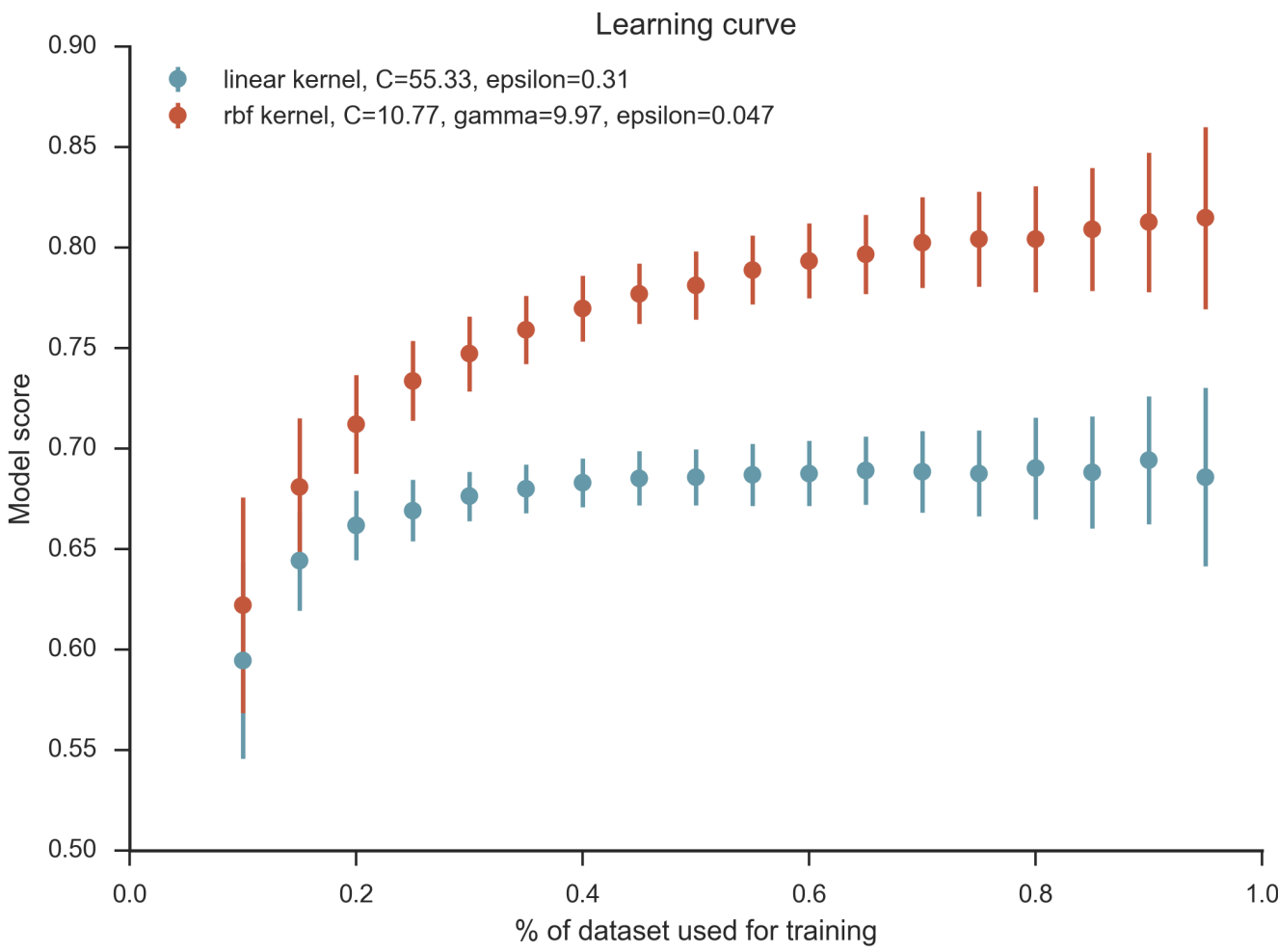
In order to evaluate the structure of the feature data set, PCA was performed on the above matrix. It appears the financial, labor, and health sectors specify the dimensions of largest variance.



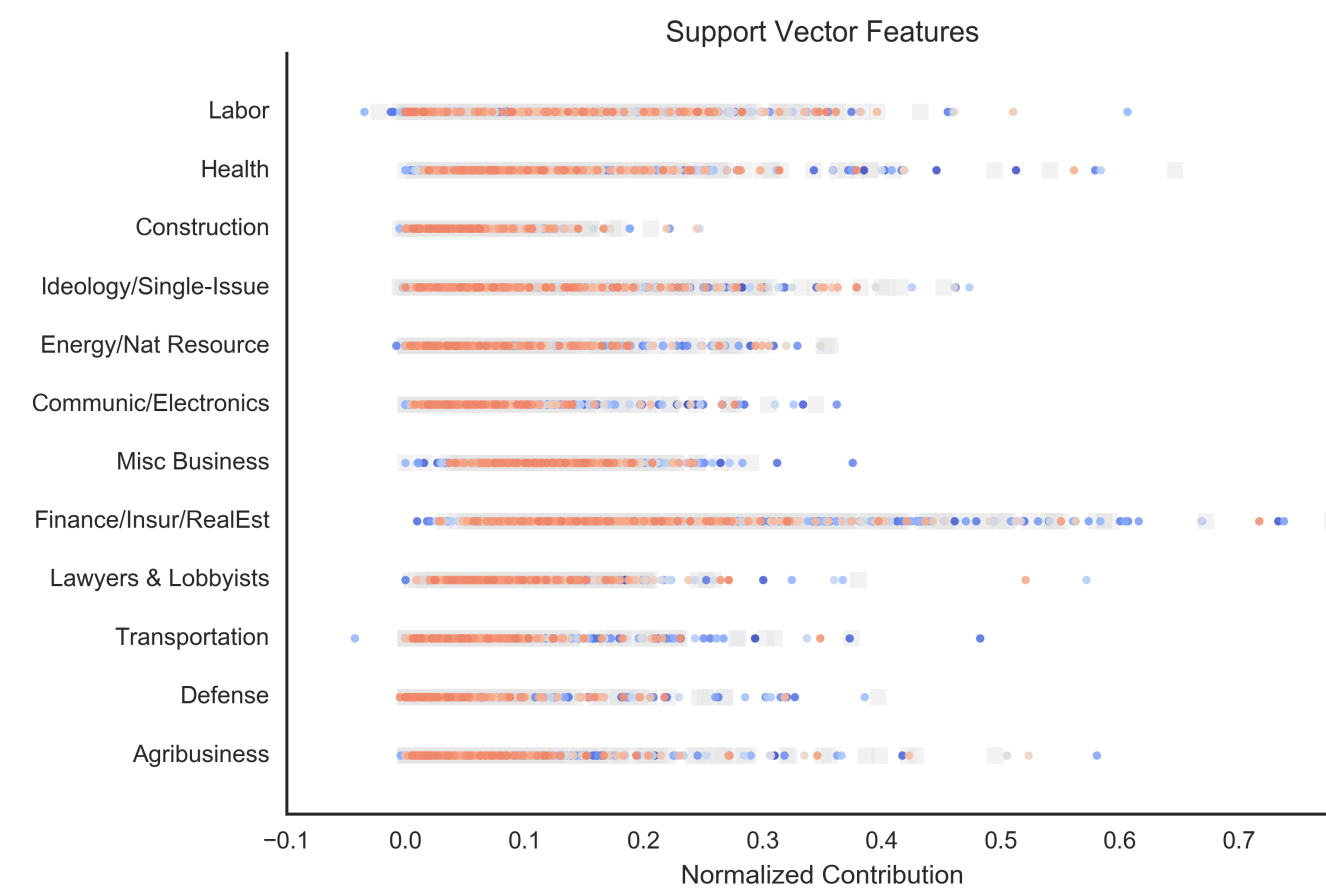
This data was used to train a range of support vector regression (SVR) models with variable kernel types, over a range of hyperparameters using the SciKitLearn cross_validation and GridSearch methods [3].

Validation

We performed a randomized search over the SVR hyperparameters in order to arrive at the final model, with each iteration of the search performing a 5-fold cross validation in order to estimate the training score. The final model had a training score of 82% and yielded a test score of 83% on held-out data. Learning curves, as a function of training set size, were calculated for each model type, and are shown below. It was found that the gaussian kernel had the best performance with respect to model score and time required to train the model. Analysis showed polynomial kernels performed comparably, but required significantly training time.



The support vectors for the RBF SVR model are depicted below. We see that the SVs are diffuse over the entire training set.



It was found that we were able to build a reasonably accurate model given the available data. Our learning curve suggests that these results would improve if we were not data limited.