

# Legislator Arithmetic \*

Daniel Argyle    *FiscalNote*

---

See intro...

*Keywords:* ideal point estimation

---

## Introduction

We propose a neural network implementation of ideal-point estimation that scales well to large datasets and allows incorporation of additional metadata. Neural networks are well-suited for these models, and the performance benefit, along with distributed computing capabilities, allows application of ideal point estimation to pooled datasets where computation was previously infeasible due to scale. We demonstrate the algorithm on two different datasets, the complete history of US Congressional roll call votes and modern cosponsorship networks, and compare the results against standard ideal point estimation techniques.

To evaluate algorithmic performance, we test the resulting estimates on both training and test data by holding out a subset of legislators' votes. This allows us to compare the quality of different model parameterizations and choice of dimensions while still guarding against overfitting. Specifically, we directly compare the performance of different ideal point parameterizations such as DW-NOMINATE and the conventional Bayesian parameterization.

We demonstrate the algorithms in two ways. First, we jointly estimate ideal points over the pooled set of US Congressional roll call votes from 1789-2018. Unidimensional ideal points from the neural network implementation are similar to the conventional DW-NOMINATE results. However, cross validation scores indicate that the data are better explained with more than one dimension. Clustering the multidimensional ideal points yields intuitive temporal and ideological groupings and provides a more nuanced picture of ideological polarization.

Second, we take advantage of the fact that many more bills are sponsored than actually come to a vote and estimate an ideal point distribution over a large set of sponsorship and cosponsorship decisions in the 93rd-114th Congresses. Cosponsorship provides a different perspective on legislators' beliefs, independent of strategic voting or administrative votes of little ideological salience. We treat cosponsorship as a clear endorsement of a bill's content and assume that a choice not to cosponsor a bill can be interpreted as something less than full support. When compared to traditional ideal points, cosponsorship ideal points show somewhat different trends in polarization and result in a higher number of optimal dimensions.

## Implementing ideal points as a neural network

This work is inspired by two similar lines of research in political science and computer science. Speaking broadly<sup>1</sup>, political scientists—from the days of (really old cite), through Poole and Rosenthal [Poole and Rosenthal \(1985\)](#), and up to and including modern contributions—have focused primarily on ideal point estimation as a means to study the ideology space implied by the ideal

---

\*The code for this method is available at the author's github repository.

<sup>1</sup>A complete review of the literature in either field is beyond the scope of this work, and there are, unsurprisingly, some exceptions to this assertion.

points themselves. That these methods also predict votes is somewhat of an afterthought. On the other hand, computer science implementations largely focus on using an ideal point framework to predict legislator votes, without concern regarding the ideal points themselves.

Combining the insights of these two fields, we wish to have the most predictive power without overfitting. We can then interpret the most predictive ideal points for insights.

We suggest that 1. the model that predicts the best *on a held out sample of votes* provides the most insight into ideal points and 2. that there is a clear trade-off between explanatory power and ease of interpretation of ideal point models that should be exploit.<sup>2</sup>

We implement ideal point models in a neural network framework because it's easily extensible and transparent. Our results suggest that the two dimensional model relied on sacrifices explanatory power in voting decisions.

### Model frameworks

There are two broad frameworks that have been commonly used for ideal point estimation. Both tie back to spatial voting theory and the idea that a legislator will prefer something that is “near” to them in the ideology space. This ideology space can be parameterized with an arbitrary number of dimensions (denoted with  $K$ ), but in most applications  $K$  is set to be 1 or 2.

The first framework, the NOMINATE family, posits that the probability that a legislator votes yes on a given bill is given by:

$$prob(vote = yea) = \Phi(u_{ijy} - u_{ijn}) = \Phi\left(\beta \left[ \exp\left(-\frac{1}{2} \sum_{k=1}^s w_k^2 d_{ijyk}^2\right) - \exp\left(-\frac{1}{2} \sum_{k=1}^s w_k^2 d_{ijnk}^2\right) \right]\right) \quad (1)$$

where  $u_{ijy}$  is the utility legislator  $i$  receives from voting yes ( $y$ ) on proposal  $j$ .<sup>3</sup> This utility is expressed in terms of the squared distance  $d_{ijyk}^2$  between a legislator's ideal point and the “yes” outcome point (and similarly  $d_{ijnk}^2$  is the distance between the legislator's ideal point and the “no” outcome point). There are also a set of salience weights  $w_k^2$  which allow different dimensions to have more impact on the final voting decision. The function  $\Phi$  is commonly a normal or logistic cdf, which makes this model, once the weights have been obtained, essentially a probit or logit regression with a single parameter  $\beta$ .

The item-response framework, commonly associated with [Clinton, Jackman and Rivers \(2004\)](#), predicts a vote as follows:

$$prob(vote = yea) = \Phi(\beta_j \cdot \mathbf{x}_i + \alpha_j) \quad (2)$$

where  $\mathbf{x}_i$  is a  $K$  dimensional vector of ideal points and again  $\Phi$  can correspond to either a logistic or probit model. The bill level parameter  $\beta_j$ , sometimes referred to as polarity, is multiplied with the ideal point vector and determines the lean of the bill. The parameter  $\alpha_j$ , sometimes referred to as popularity, determine a baseline approval rate for a specific proposal.

<sup>2</sup>We will provide evidence that the optimal number of dimensions for ideal point estimation is larger than the most common one or two. It is perfectly reasonable to choose to rely on two dimensions, but the trade-offs of that choice should be clear.

<sup>3</sup>For simplicity, we are suppressing the time dimension that is often present in these models. However, NN-NOMINATE generalizes to include dynamic ideal points as is discussed later in the paper. See also [Carroll et al. \(2009\)](#) for more detail about the background of this model.

## The neural network

While neural networks have become most strongly associated with artificial intelligence, they are, in essence, flexible optimization systems fitted using various forms of gradient decent algorithms.<sup>4</sup> All that is required to estimate an ideal point model using a neural network is to set the network structure to match the underlying framework of the ideal point model. One example can be seen in Figure 1, which represents the architecture of a WNOMINATE model. We will walk through each layer of this example in some detail as a baseline for further exploration.

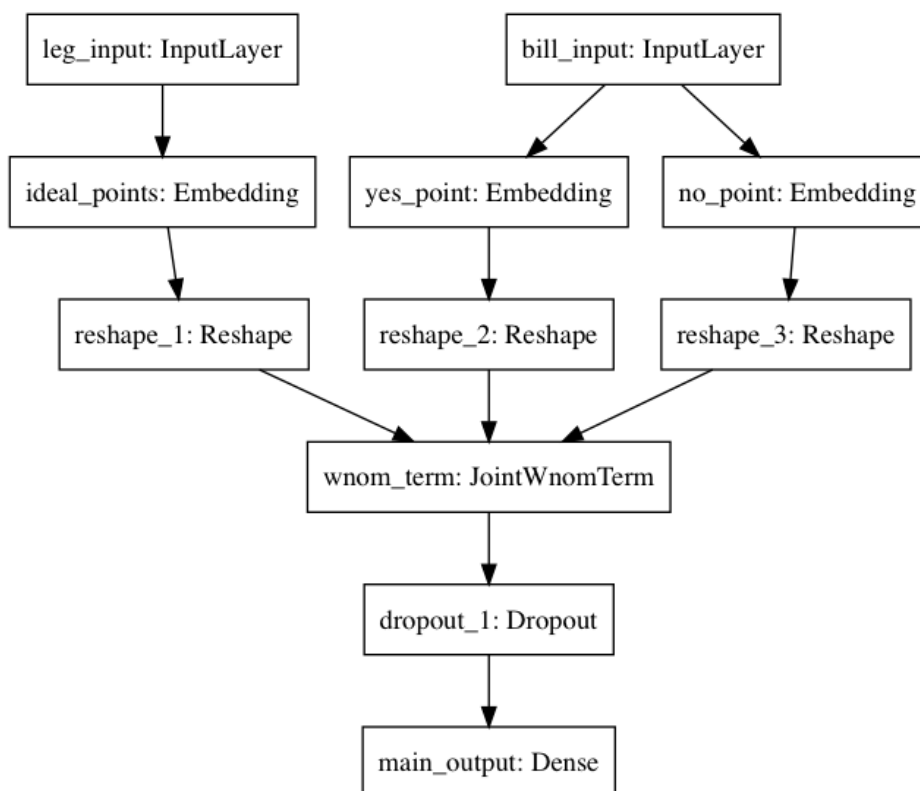


Figure 1: Base Model

- **Input layers:** The inputs to this model consist of numeric ids for both legislators and bills. The data is given as a batch of triples,  $(legislator\_id, bill\_id, vote)$ .
- **Embedding layers:** These layers take the numeric ids and convert them into continuous measures, for example mapping a specific legislator id to her ideal point. In this model there is one ideal point embedding for legislators, and two embeddings for bills, the yes-point and the no point, that correspond to numeric values associated with WNOMINATE style estimation. The ideal point layer is restricted in two ways. First, the model includes orthogonality regularization on the ideal points (Brock et al., 2016; Vorontsov et al., 2017). This penalizes correlation between ideal point dimensions and ensures that the resulting estimates are orthogonal. Ideal point embeddings also have a maximum norm constraint,

<sup>4</sup>Indeed, while not directly applicable to the models contained here, there are several papers proving that feed forward neural networks can approximate any continuous function (Cybenko 1989, approximation, Hornik 1991)

which ensures that the norm of the ideal point vector for any legislator lies within the unit hypersphere.<sup>5</sup>

- Reshape layers: Because they are often used for text analysis, by default embedding layers assume that there is a sequence of ids to transform into the embedding space. Since we have only a single id to embed, these layers drop the dimension associated with the sequence to form a two dimensional tensor.
- JointWnomTerm layer: The embedded and reshaped data is fed into a custom neural network layer that implements the conventional WNOMINATE estimation (cite Poole and Rosenthal and the WNOMINATE package). Specifically, the layer implements the inner portion of model in 1, sepcifically calculating

$$\exp\left(-\frac{1}{2}\sum_{k=1}^s w_k^2 d_{iyk}^2\right) - \exp\left(-\frac{1}{2}\sum_{k=1}^s w_k^2 d_{ink}^2\right) \quad (3)$$

where  $K$  is the number of dimensions in the model,  $w_k$  is the salience weight for dimension  $k$ , and  $d_{iyk}^2$  is the distance between a legislator's ideal point and the yes point of the bill and  $d_{ink}^2$  is the distance between the legislator's ideal point and the no point of the bill.[<sup>10</sup>]

- Dropout layer: Dropout regularization, as proposed by [Srivastava et al. \(2014\)](#), has become an extremely common way of limiting model overfitting, often with the additional benefit of improved behavior during model training. Dropout layers set model weights to 0 at random during iterations of the training process. This prevents any single weight dominating the prediction of the algorithm. In this specific instance, the dropout layer sets the salience weight for any given ideology dimension to 0 for a specific iteration of the optimization algorithm relying upon the other dimensions to make a vote prediction. In practice, this has resulted in better model performance in terms of observable metrics and by limiting overfitting during the training process..
- Dense layer: The quantity obtained in the JointWnomLayer is then used as an input to a Dense layer, parameterised as a logistic regression (i.e. sigmoid activation and binary\_-crossentropy loss). The single model weight estimated here is denoted  $\beta$  in most NOMINATE models and represents the signal to noise ratio. A high value corresponds to votes being largely deterministic, a lower value suggest more randomness.

## The neural network implementation works

We demonstrate that the neural network implementation works at least as well as existing methods through two simple tests. The first is to test the model on synthetic data, generated with known parameters, and determine how well the method recovers these known values. The second is to check that the results match the commonly used implementations of WNOMINATE ([Poole et al., 2011](#)) and item response ([Jackman, 2017](#)) ideal points.

<sup>5</sup>Note that the maximum norm constraint does not require that any legislator actually attain a unit norm. This differs somewhat from existing implementations which seem to ensure all dimensions fill the range from [-1, 1]. If this behavior is desired, it can be easily fixed in postprocessing.

### Recovering known ideal points

The synthetic data is generated to have 3 known dimensions, where the random predictions are generated in the NOMINATE parameterization. To guard against overfitting to specific votes a legislator took, we train the model on a subset of the data and test it on 20% of votes as a held out set. This means that a legislator’s ideal point is determined only by the votes in the training sample and that a bill’s parameters are determined only by the subset of legislators included who voted on this bill.

There are two primary benefits to using a held out sample. First, relying on an out of sample evaluation set is often used to determine when a neural network has converged while estimating a model. There is little point in continuing to optimize a model where out of sample performance has plateaued and doing so usually results in overfitting. Second, performance on a held out sample can be used to make additional decisions about the model structure, including parameterization choices or for evaluating distributional assumptions. For example, how does making the ideal points dynamic affect the model performance? What is the optimal number of ideal point dimensions? An additional dimension may offer an improvement in out of sample prediction, but it almost always increases in sample performance. The decision to add a dimension has always been somewhat arbitrary, an out of sample performance test provides an objective metric to evaluate.<sup>6</sup>

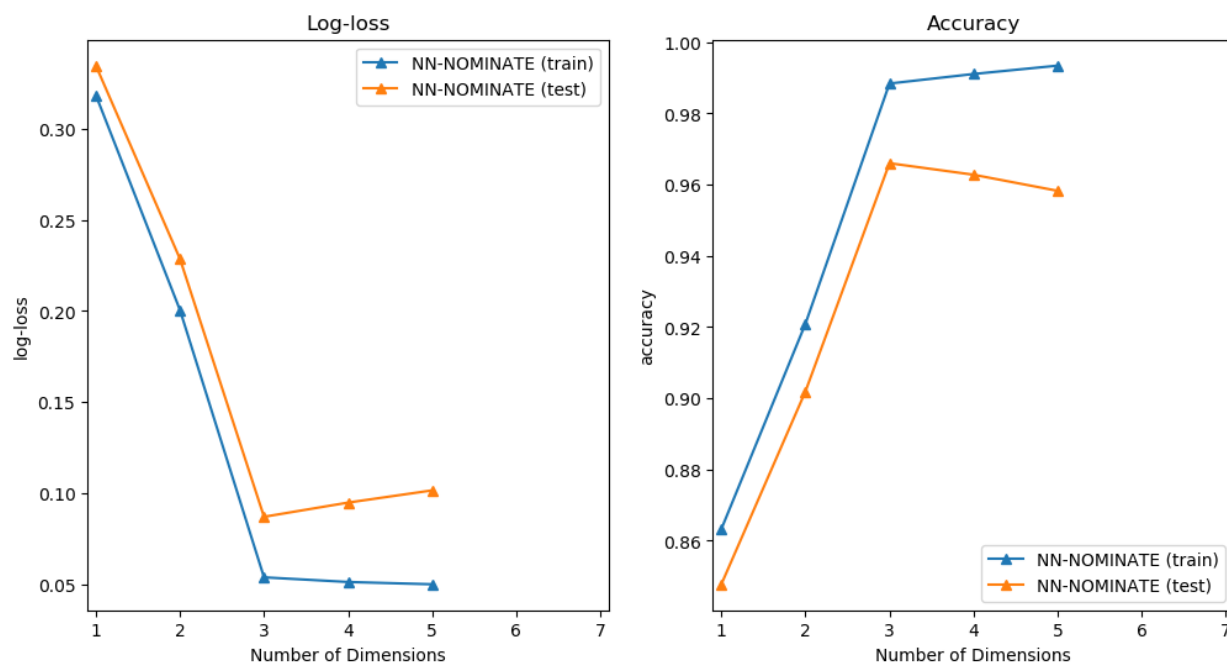


Figure 2: Performance metrics on Synthetic Data

Figure 2 shows the model performance of NN-NOMINATE on synthetic data, varying the number of dimensions in the model. We evaluate two metrics, accuracy (whether or not the vote was predicted correctly) and log-loss (a common metric used in classification problems). Log-loss is useful because it provides an idea of how good the probability estimates of a model are, and penalizes over-confident predictions. We see that the log-loss decreases rapidly is minimized at the optimal number of dimensions on the test data, and the test log-loss begins to increase even

<sup>6</sup>It is not as simple as saying that “the model that fits best on test data” is the right answer. In some cases it is preferable to say, this is the first model that attains a near minimum.

as in sample . Training accuracy shows similar patterns, attaining best performance on the known optimal number of dimensions.

### Comparison to existing packages

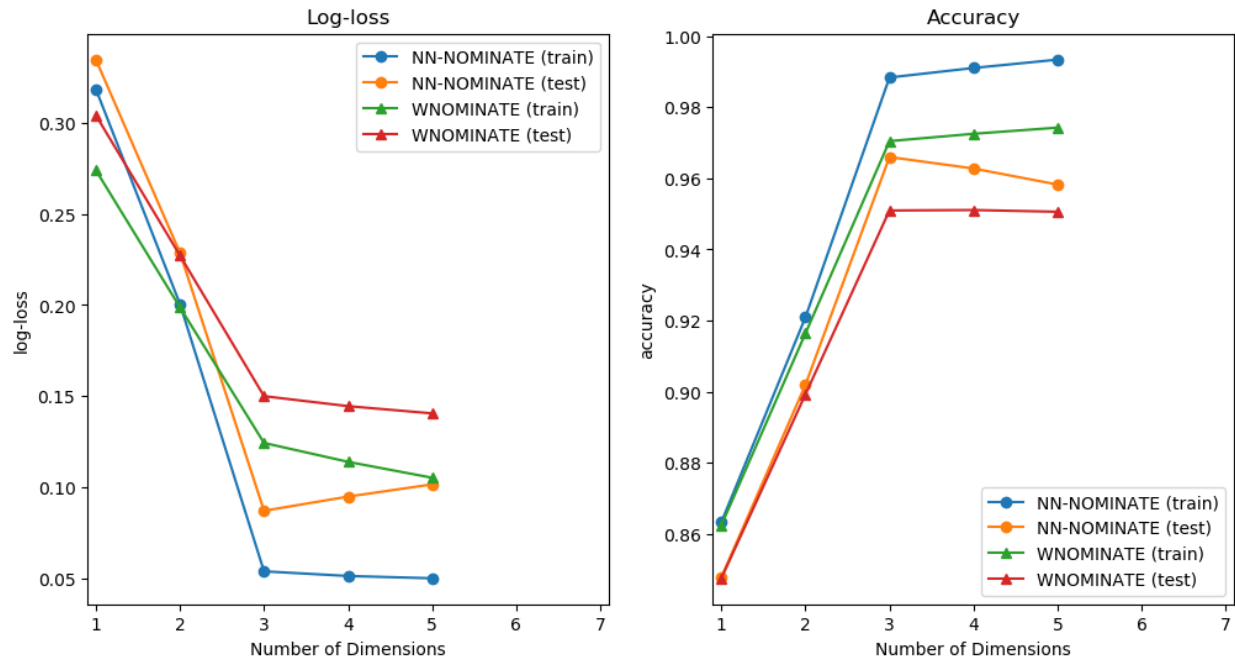


Figure 3: Performance metrics on Synthetic Data

Since it has not yet been common to estimate ideal point models on a training and test set we verify that the model works by comparing the results of NN-NOMINATE to the WNOMINATE package available in R. The results for the synthetic data are shown in Figure 3, where WNOMINATE shows similar patterns over training and test data. It appears that for the synthetic dataset at least, NN-NOMINATE (slightly) outperforms the commonly used WNOMINATE package attaining lower log loss and higher accuracy on test data.

In Table XX, we show the same metrics on real world data from the 110th-115th US Senate. In addition to similar performance metrics, the results of the two methods are strongly correlated.<sup>7</sup>

### Results and applications

While it is comforting to know that a new optimization of an existing model returns similar results, NN-NOMINATE is capable of more than the existing implementations. The first is large scale ideal point estimation, for example fitting models on the entire history of the US Congress locally on a laptop. The second is to demonstrate the extensibility of the framework, where the model can be extended very simply to include covariates as well as more complicated representations of a bill than existing representations.

<sup>7</sup>Since ideal points are not independently identified, including up to dimension switching and other problems, we use simple correlations between the metrics as our outcome of choice. It is often the case that NN-NOMINATE results in ideal points that are smaller in magnitude on average than the WNOMINATE package.

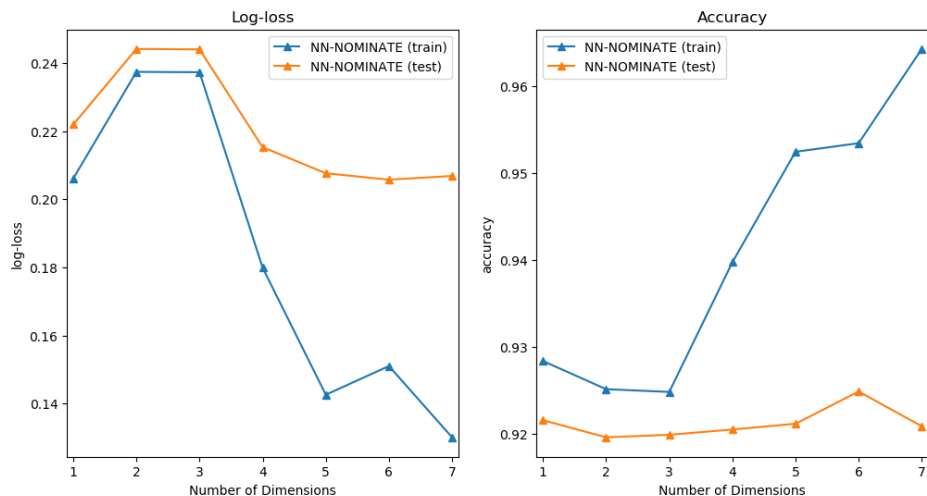


Figure 4: Log-loss on 110th-115th Senate

We examine two large datasets, all of which can be fitted locally on a 3 year old mac laptop in less than an hour.<sup>8</sup> The first is simply the entire history of US rollcalls. The dimension tradeoff calculation from above is shown in Figure XX, which shows that the optimal number of dimensions is at least 8 over this entire dataset.

The second dataset is the set of all cosponsorship decisions from the 93rd-114th US Congresses, a task that was also attempted in (cite that earlier paper). While that work hypothesized that cosponsorship would show fewer dimensions than voting, we believe that it is also plausible that there should be more. Since so many voting decisions include strategic concerns, cosponsorship, a potentially less impactful choice, could allow legislators to espouse a wider variety of policy opinions (and thus require additional dimensions in the modeling process).

In this case, a “yes” vote is considered to be sponsoring or cosponsoring a bill, whereas everyone who did not is assumed (at least by opportunity cost) to be a “no”. Note that this data is much noisier than votes. Since so many bills are introduced that do not proceed in a given session, there are likely many cases where a legislator would have preferred that a bill pass even if they did not have the opportunity to copponsor it. Because of this, we add additional dropout layers on the bill parameters themselves. This limits the extent to which the model can overfit to noise in the dataset, although the results in Table XX show that this is only somewhat successful as the training accuracy somewhat rapidly reaches 1.

While the above results could be construed as evidence, as was the assertion made (in that other paper), that there are fewer dimensions in cosponsoring decisions than there are in voting it seems equally plausible that the model is more easily fitting to noise in the cosponsorship decisions so rapidly that the true policy preference dimensions.

<sup>8</sup>Formal performance comparisons to existing methods is forthcoming. Note that in the case of the WNOMINATE and pscl packages, performance comparisons are impossible because they require computation of a roll call matrix that is  $n \times m$  which is more than can easily fit in memory (assuming integer values for all votes this still takes up more than XXGB).]



We also extend the model to include a simple covariate. It is plausible that a legislator might behave differently while in the party that controls the chamber. One way this might be observed is that a legislator in the party in power is more likely to vote yes, on average, than someone in the minority party. This could even be the case when the match with their ideal point is very close, but due to party pressure a legislator could deviate from this to some extent. As such we estimate the models from the previous section including a variable indicating that a legislator is in the party in power so that the logistic layer now becomes:

$$prob(vote = yea) = \Phi \left( \beta \left[ \exp \left( -\frac{1}{2} \sum_{k=1}^s w_k^2 d_{iyk}^2 \right) - \exp \left( -\frac{1}{2} \sum_{k=1}^s w_k^2 d_{ink}^2 \right) \right] + \gamma * pip_{ik} \right) \quad (4)$$

where  $\gamma$  represents the effect of being in the party in power. The results show that this factor is positive, as expected, but that the magnitude represents only a very small difference in vote probability. Note that the NN-Nominate framework does not easily lend itself to hypothesis testing. One way to achieve this would be a parametric bootstrap (ala Poole and Rosenthal) but this has not been attempted.<sup>9</sup>

Additionally, while a discussion of using a text representation of a bill, rather than a simple indicator value, is beyond the scope of this work this is a prime area of current research (see e.g. Kornilova). These models work by building embedding representations of the bills themselves and then determining legislators ideal points relative to this more robust representation of a bill. While the accuracy of these models does not yet match that of the bill indicator approach, these models have the strong advantage of allowing out of sample prediction on new bill text. This means that we can determine how a legislator would vote on new bills, as they are introduced (for more we refer you to the Kornilova paper).

## Conclusion

1. I implemented DW-NOMINATE as a neural network
  - Why? Because I could! But also because it's a nice platform for this kind of optimization.
  - It scales much better than existing implementations
  - It's extensible in very interesting ways
2. All ideal point models are (a bit) overfit.
  - At some point the algorithm starts to make marginal improvements to the parameters that don't improve out of sample performance
  - Out of sample performance matters much more than in sample (e.g. if we only cared about in sample we'd just add dimensions until we can predict it perfectly)
  - Out of sample performance is a useful metric for evaluating modeling choices (adding another dimension, adding a time component, adding an external variable)

this is inline 42 and

---

<sup>9</sup>Blei paper about neural networks as Bayesian inference.



```
py$answer
```

```
## [1] "42"
```

1. How to implement in a neural network
  - a. Discuss both parameterizations
  - b. Show convergence plots
  - c. Discuss out of sample fit as the metrics of importance
2. A Neural Network implementation works
  - a. It estimates correctly on synthetic data (accuracy tables for both of these)
  - b. It matches results from standard libraries (wnominate and pscl) on both synthetic and real data
    - i. Show both in sample and out of sample accuracy (will require writing a python wrapper to get predictions from the R results)
3. Results and applications of the neural network implementation
  - a. What are the optimal number of dimensions? Plot these and discuss tradeoffs
  - b. Applying the model to cosponsorship
  - c. Polarization differences
    - i. Votes models over time (how to do in multiple dimensions?)
    - ii. Compare between cosponsorship and votes

Dropout in two places: 1. Each dimension to ensure fitting even if first dimensions dominates.  
2. On bill embeddings (to not overfit to a specific bill). Why not on legislators? The structure of the dynamic part of the model doesn't allow this easily.

## References

- Brock, Andrew, Theodore Lim, James M. Ritchie and Nick Weston. 2016. "Neural Photo Editing with Introspective Adversarial Networks." *CoRR* abs/1609.07093.  
**URL:** <http://arxiv.org/abs/1609.07093>
- Carroll, Royce, Jeffrey B Lewis, James Lo, Keith T Poole and Howard Rosenthal. 2009. "Measuring bias and uncertainty in DW-NOMINATE ideal point estimates via the parametric bootstrap." *Political Analysis* 17(3):261–275.
- Clinton, Joshua, Simon Jackman and Douglas Rivers. 2004. "The statistical analysis of roll call data." *American Political Science Review* 98(2):355–370.
- Jackman, Simon. 2017. *pscl: Classes and Methods for R Developed in the Political Science Computational Laboratory*. Sydney, New South Wales, Australia: United States Studies Centre, University of Sydney. R package version 1.5.2.  
**URL:** <https://github.com/atahk/pscl/>
- Poole, Keith, Jeffrey Lewis, James Lo and Royce Carroll. 2011. "Scaling Roll Call Votes with wnominate in R." *Journal of Statistical Software* 42(14):1–21.  
**URL:** <http://www.jstatsoft.org/v42/i14/>
- Poole, Keith T and Howard Rosenthal. 1985. "A spatial model for legislative roll call analysis." *American Journal of Political Science* pp. 357–384.
- Srivastava, Nitish, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever and Ruslan Salakhutdinov. 2014. "Dropout: A Simple Way to Prevent Neural Networks from Overfitting." *Journal of Machine Learning Research* 15:1929–1958.  
**URL:** <http://jmlr.org/papers/v15/srivastava14a.html>
- Vorontsov, Eugene, Chiheb Trabelsi, Samuel Kadoury and Chris Pal. 2017. "On orthogonality and learning recurrent networks with long term dependencies." *CoRR* abs/1702.00071.  
**URL:** <http://arxiv.org/abs/1702.00071>