

# Congressional Survival

SpringBoard Capstone 2  
Neil Horning

# Contents

Client / Stakeholder  
Case

Exploratory Data  
Analysis

Machine Learning

Findings,  
Recommendations  
and Next Steps

# Problem:

Based on voting history, what type of congressional member is most likely to survive? One that meets partisan preferences of a district or the expectations of their own partisan base?

# Client / Stakeholder:



**AUDREY  
DENNEY**

**(D, CA-1)**

Audrey Denney is a Democratic congresswoman from California District 1, she was recently elected after a very close contest with Doug LaMalfa in 2020, riding the historic blue wave that swept Trump out of office. CA-1 is a rural, traditionally Republican district with a partisan lean of R+22.53. She wants to know what type of voting strategy will help her survive next election. Another stakeholder could be the DNC

# Exploratory Data Analysis

	bill	date	democratic	description	document_number
(116, 'Senate', 1, 1)	{'bill_id': 's1-116', 'number': 'S.1', 'sponsor_id': 'R000595', 'api_uri': 'https://api.propublica.org/congress/v1/116/bills/s1.json', 'title': 'A bill to make improvements to certain defense and security assistance provisions and to authorize the appropriation of funds to Israel, to reauthorize the United States-Jordan Defense Cooperation Act of 2015, and to halt the wholesale slaughter of the Syrian people, and for other purposes.', 'latest_action': 'Held at the desk.'}	2019-01-08	{'yes': 4, 'no': 41, 'present': 0, 'not_voting': 0, 'majority_position': 'No'}	A bill to make improvements to certain defense and security assistance provisions and to authorize the appropriation of funds to Israel, to reauthorize the United States-Jordan Defense Cooperation Act of 2015, and to halt the wholesale slaughter of the Syrian people, and for other purposes.	1

Meta Data

			party	D			ID	D		
			state	MA	CA	NJ	VT	WI	MA	HI
			dw_nominate	-0.774	-0.710	-0.611	-0.526	-0.512	-0.506	-0.498
			member_id	W000817	H001075	B001288	S000033	B001230	M000133	H001042
			name	Elizabeth Warren	Kamala Harris	Cory Booker	Bernard Sanders	Tammy Baldwin	Edward J. Markey	Mazie Hirono
congress	chamber	session	roll_call							
116	Senate	1	1	No	No	No	No	No	No	No
			2	No	No	No	No	No	No	No
			3	No	No	No	No	No	No	No
			4	No	No	No	No	No	No	No
			5	Yes	Yes	Yes	Yes	Yes	Yes	Yes

5 rows x 100 columns

Member Votes

## Data Sets:

- congressional voting records via the [ProPublica Congress API](<https://projects.propublica.org/api-docs/congress-api/>).
- Partisan lean of [districts](#) and [states](#) available from [FiveThirtyEight](#):

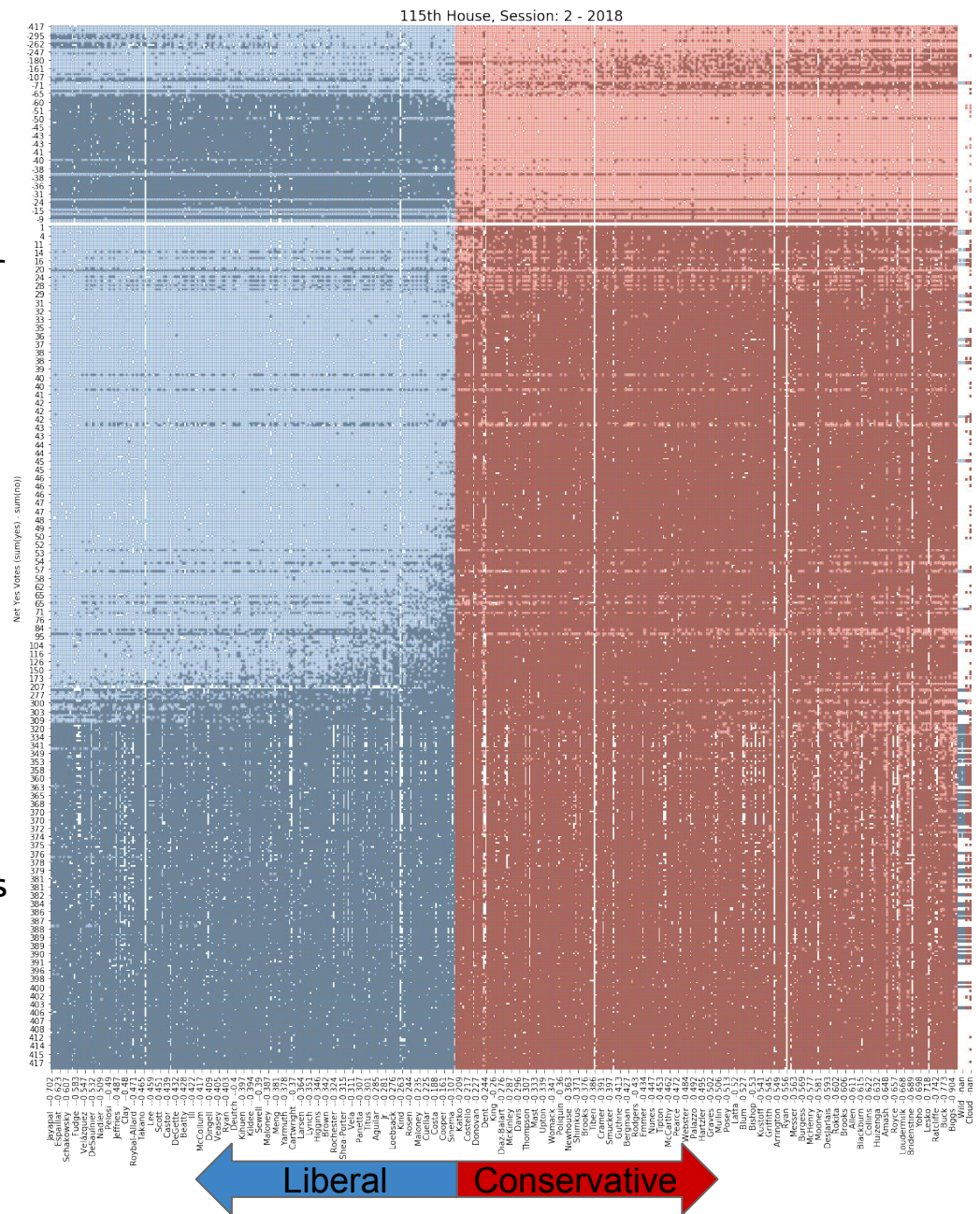
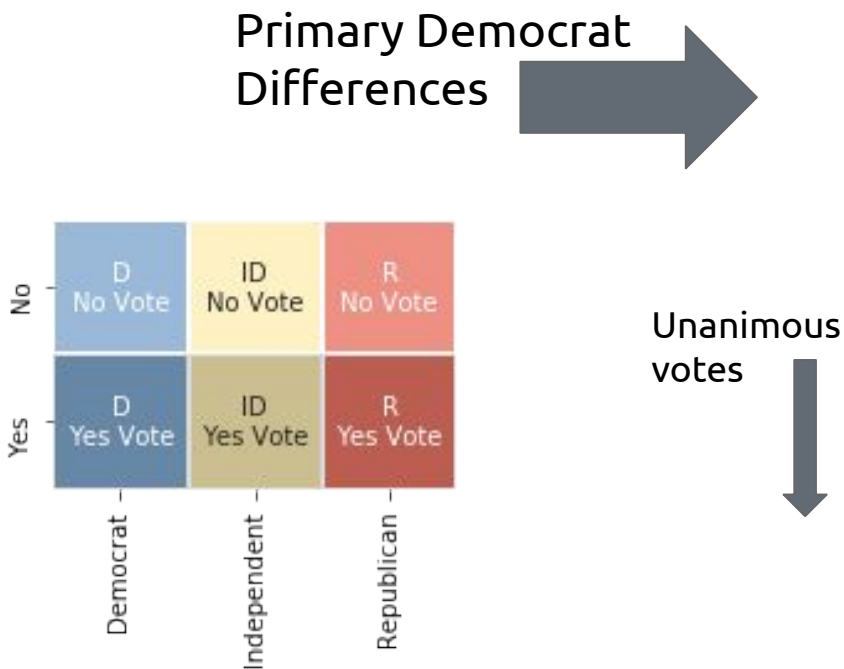
## Cleaning/Wrangling:

Created Python functions to:

- Request the ProPublica API for voting metadata for all months for a given congress chamber and year, and return a [tidy](#) Pandas DataFrame.
- Take the index of the metadata for a given year and use it to request each member's position for each vote in that year.
- Handle missing vote position data by filling rows with 'NaN' values
- Take a given year (or congress session) and chamber, return corresponding DataFrame of all the metadata and vote positions, cache this data into csv format for automatic quick loading of additional requests.
- Detect and update out of date csv cache, by fetching only the missing data.
- Convert district lean from FiveThirtyEight from positive (R+x|D+x) values into continuous -1 to +1 scale compatible with ideological score.



# What patterns are apparent in the data?



The 115th congress, session 2 of 2018, has been displayed here in a heat map with each column representing a member of congress, and each row representing a given roll call vote. Because the Republicans controlled this chamber, they fully supported nearly all of the votes that passed, and appear to be a solid red rectangular block of yes votes. However, Democrats' voting habits can be differentiated by the narrow curve of increasing support from right to left, as the bills become more popular.

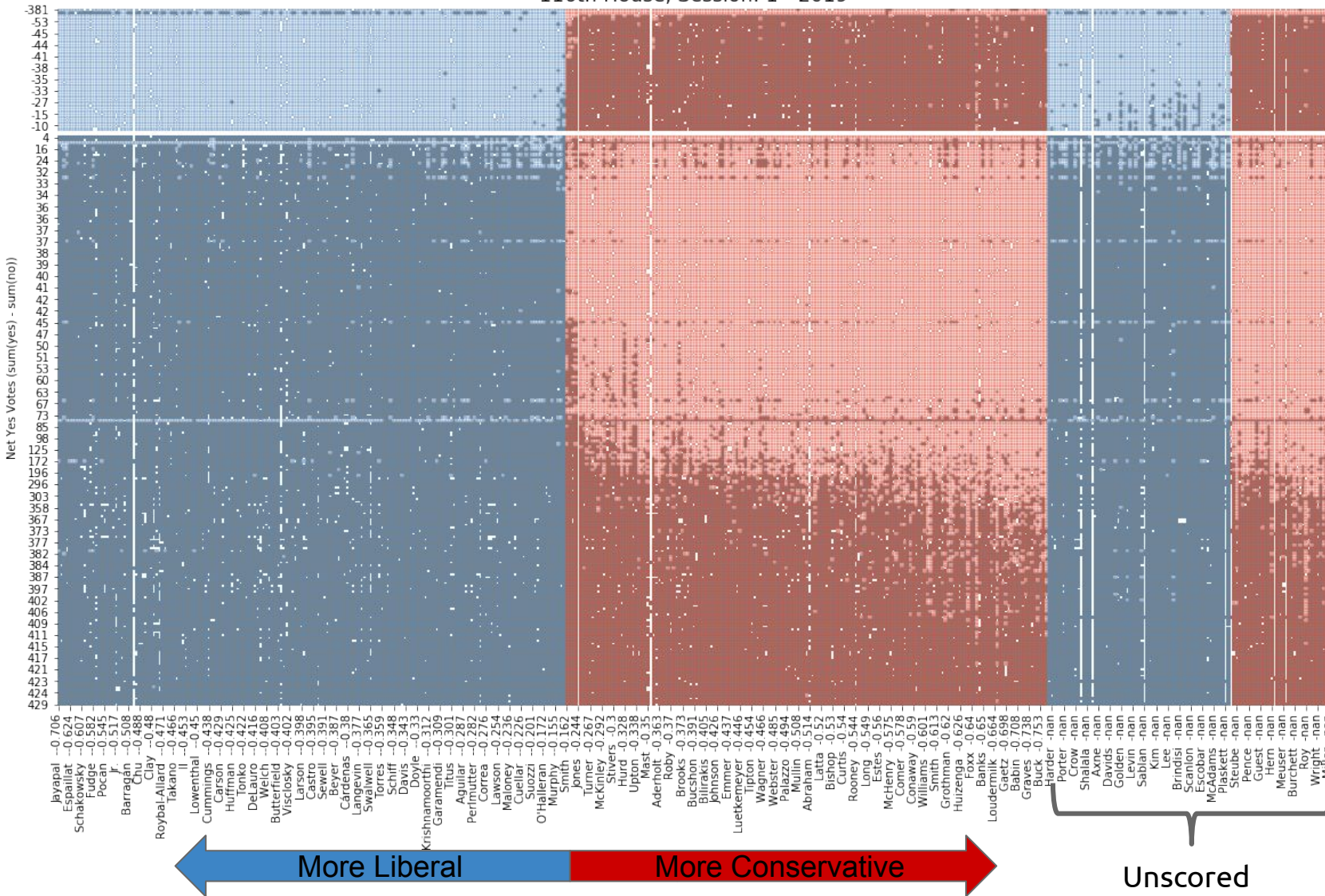


# What patterns are apparent in the data?

116th House, Session: 1 - 2019

No	D No Vote	ID No Vote	R No Vote
Yes	D Yes Vote	ID Yes Vote	R Yes Vote
	Democrat	Independent	Republican

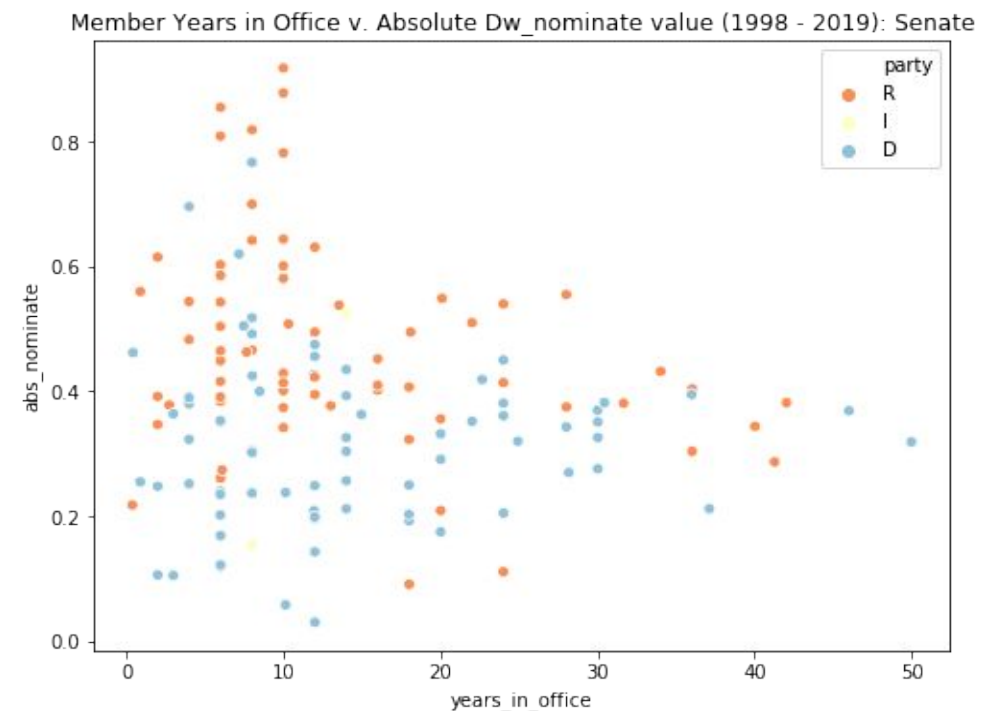
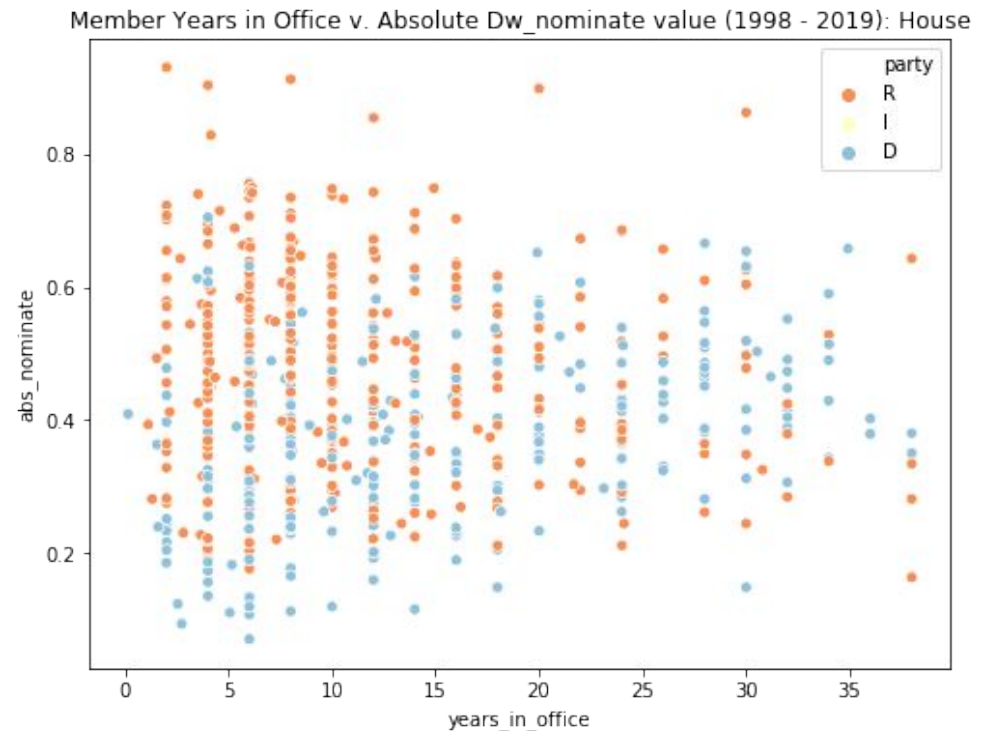
Primary  
Republican  
Differences



The same pattern can be seen in the current house (2019), except that the democrats control the chamber. Liberal / conservative scores are determined by pre-calculated "dw\_nominate" values included in the ProPublica API, some of which haven't been calculated yet.

# Which members tend to survive?

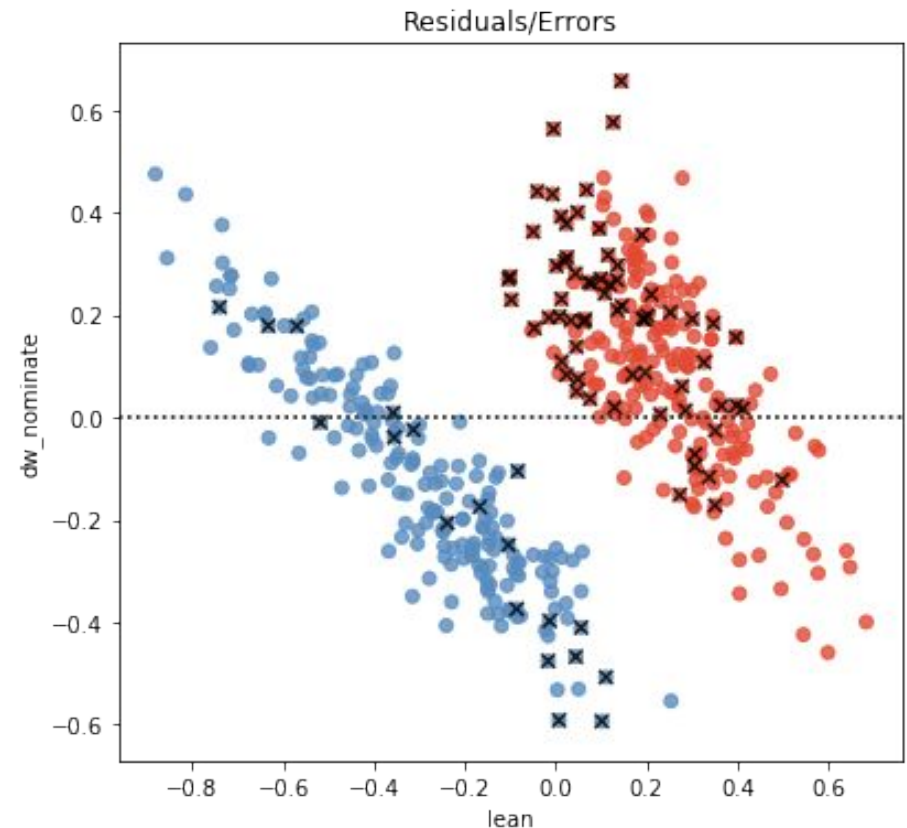
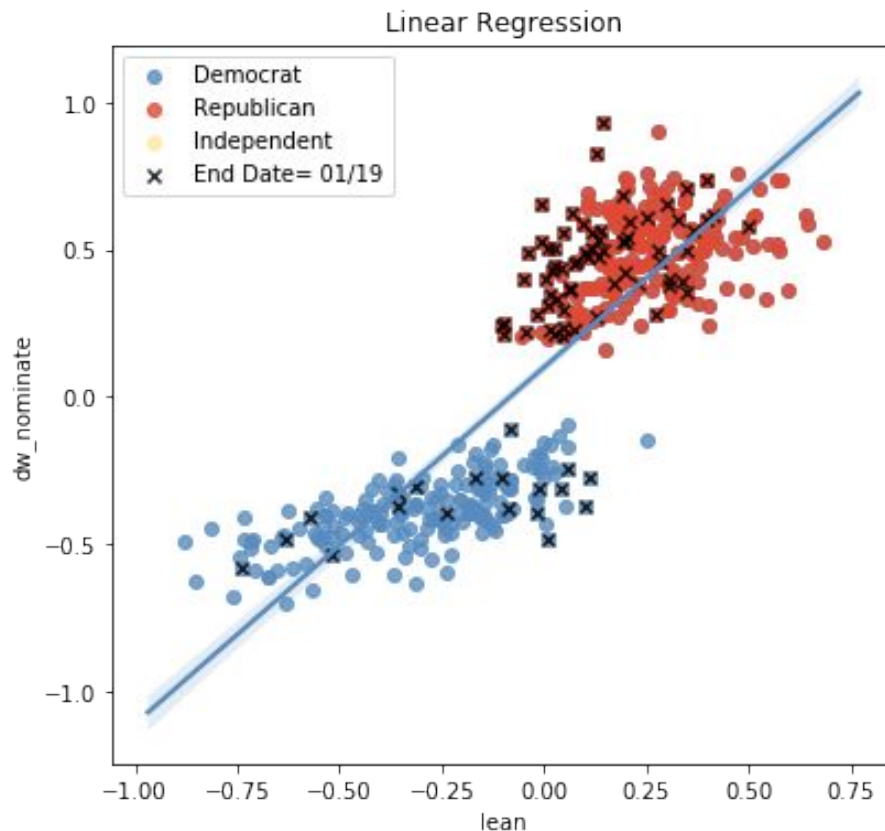
Longer lasting members tend to be more toward the center of their respective side of the political spectrum. Senate dw\_nominate scores appear to approach .37 distance from 0 as years in office increase. House, scores appear to converge toward .5 distance from 0. Outlier dw\_nominate scores are all held by house Republicans.





# Hypothesis: Members who vote more in line with the party preferences of their constituents survive longer

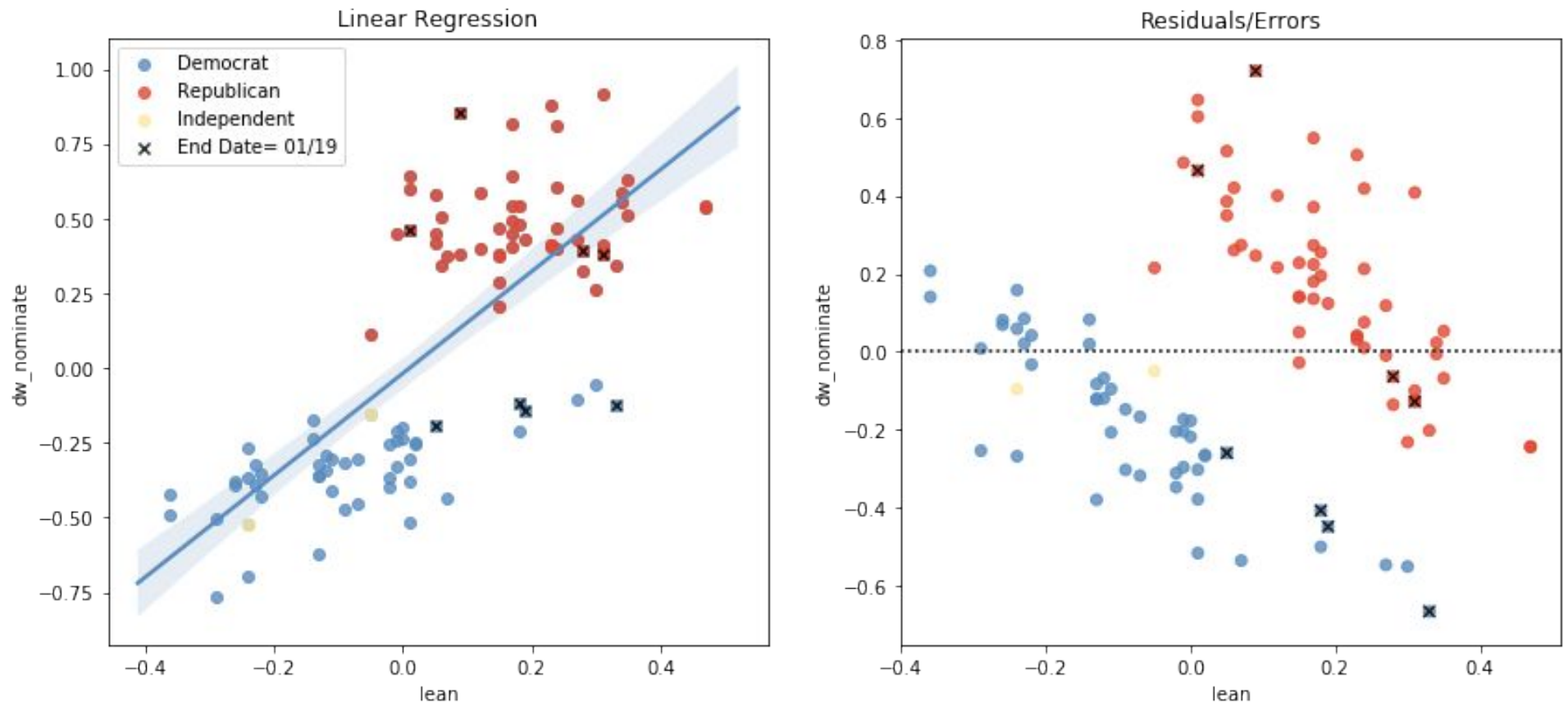
District Lean v. Ideological Score, 2018 House



Both parties are less moderate than district lean would predict by linear regression in swing districts/states, and less extreme than district lean would predict in highly partisan areas. In the residual plots on the right, members who are the most out of line with what is predicted by partisan lean are toward the upper and lower ends of the plot. It appears that those members most out of line with their district lean have been eliminated more frequently. This is particularly stark among House Democrats in districts near a partisan lean of 0

# Hypothesis: Members who vote more in line with the party preferences of their constituents survive longer

District Lean v. Ideological Score, 2018 Senate



Those more extreme than predicted, seem to be eliminated with more frequency than those more moderate than predicted. This apparent trend can be investigated further with data from more years. Note that in the Senate, this plot does not take into account which members stood for re-election.

# Machine Learning





# Dimension Reduction (PCA)

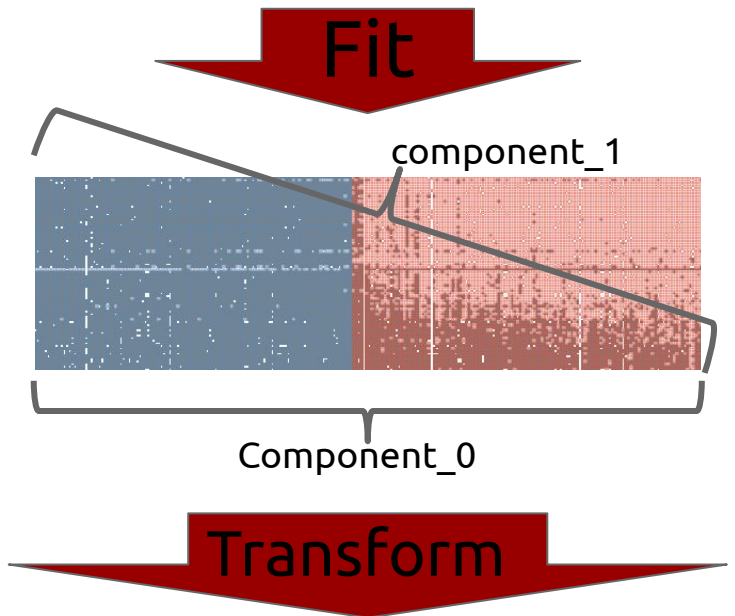
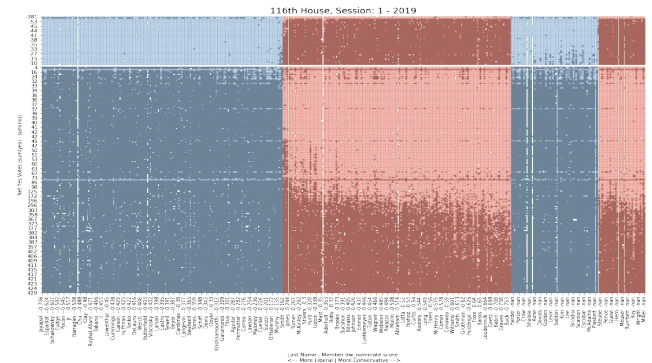
"Dw\_nominate" is a multi-dimensional scaling method, developed in the 80's where congressional voting records are interpreted along two dimensions. Let's see if we can "roll our own" dw\_nominate scores with Principal Component Analysis.

Most of the variation in the data is explained by component '0' with component '1' providing much less, but probably enough to be useful, and component '2' providing even less than that. This means it is likely possible to interpret results in 2 dimensions quite easily

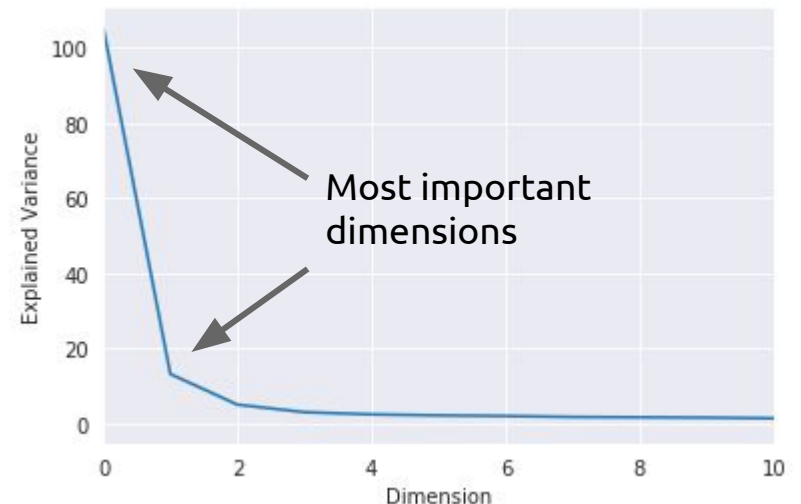
>100d:

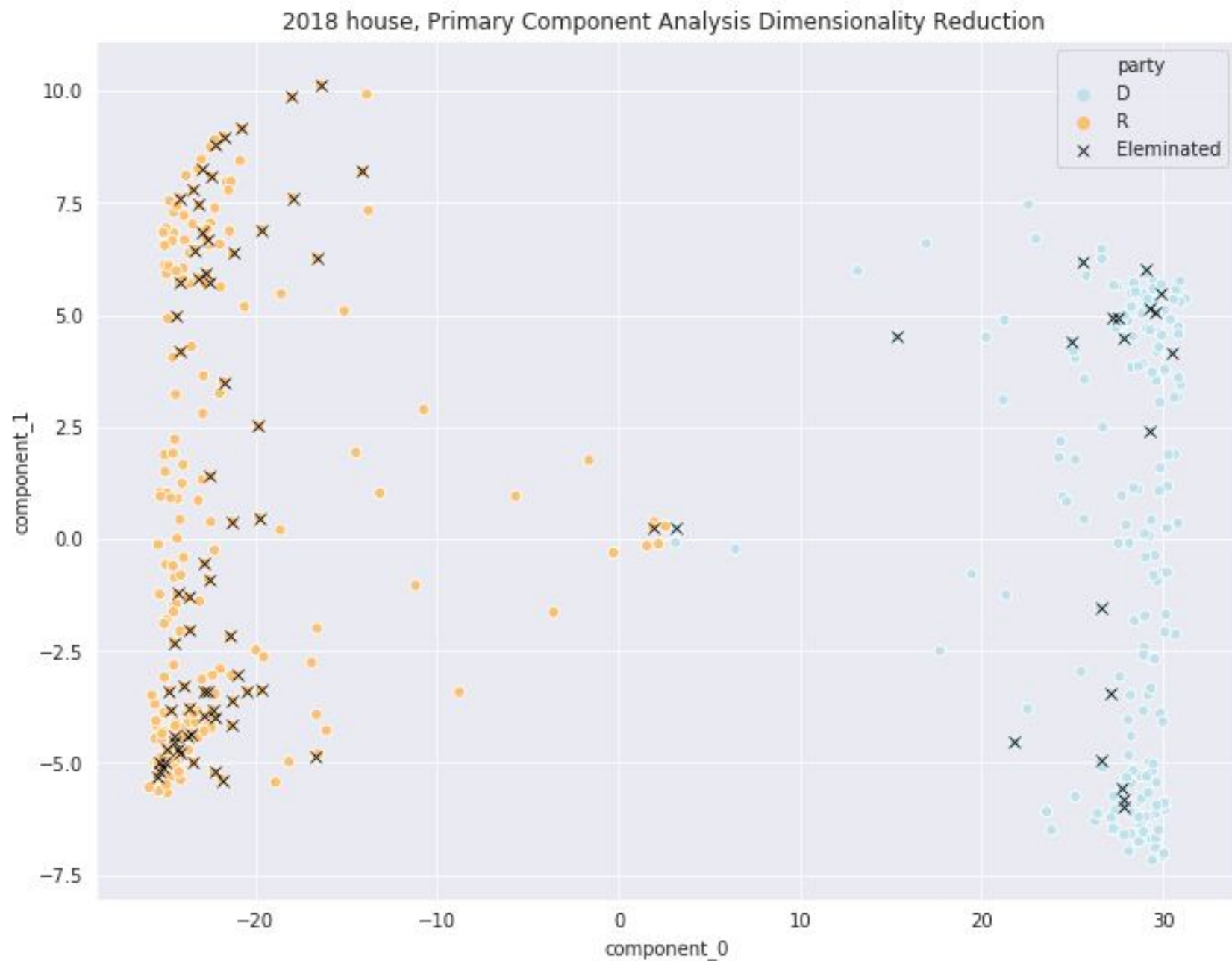
PCA reduces the hundreds of dimensions represented by individual votes by repeatedly finding an axis that represents the largest variance in the data, and projecting it into its own dimension.

Crude artistic interpretation.

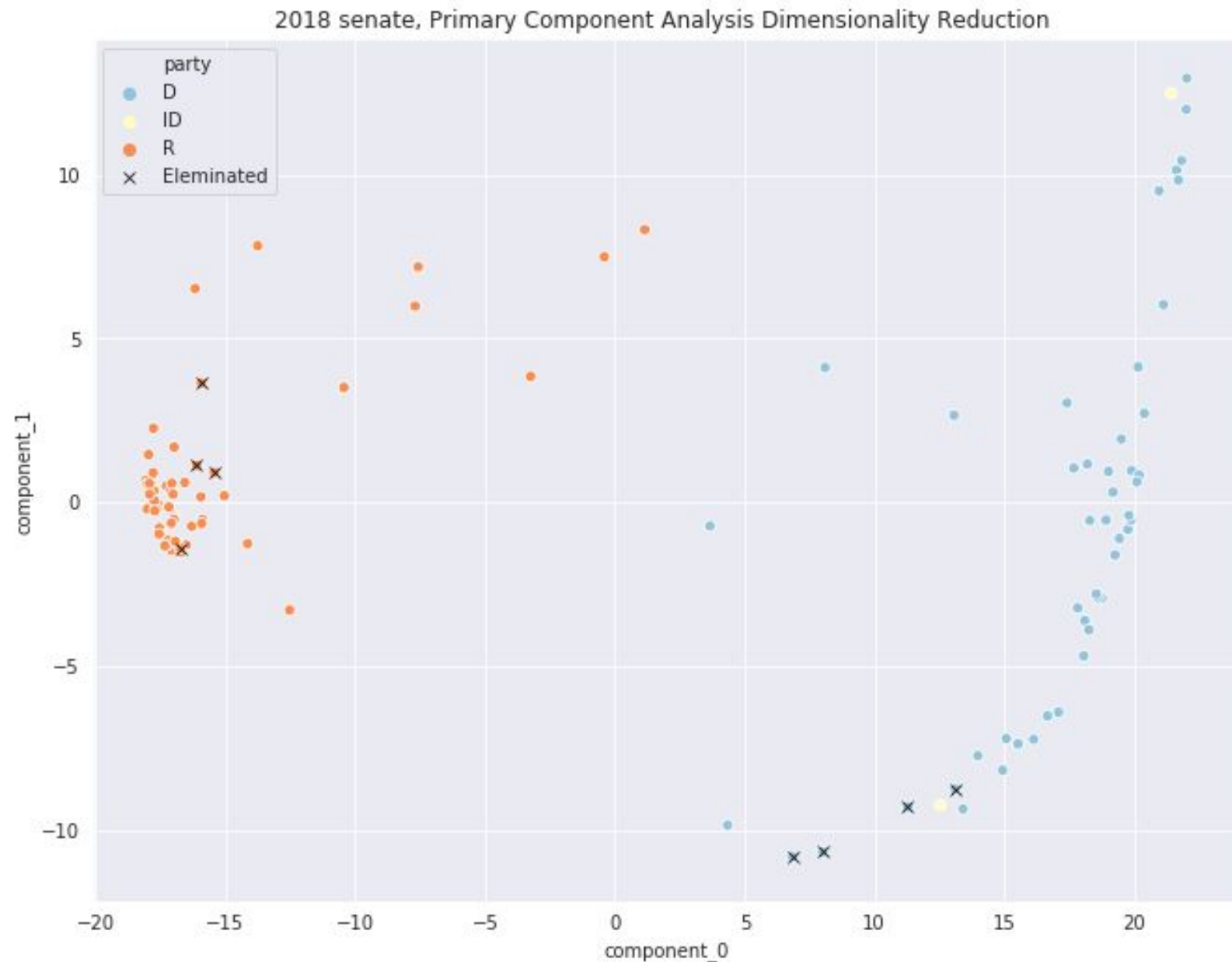


2d:





As shown, PCA can neatly separate the parties from one another based entirely on voting records, for both the house ...

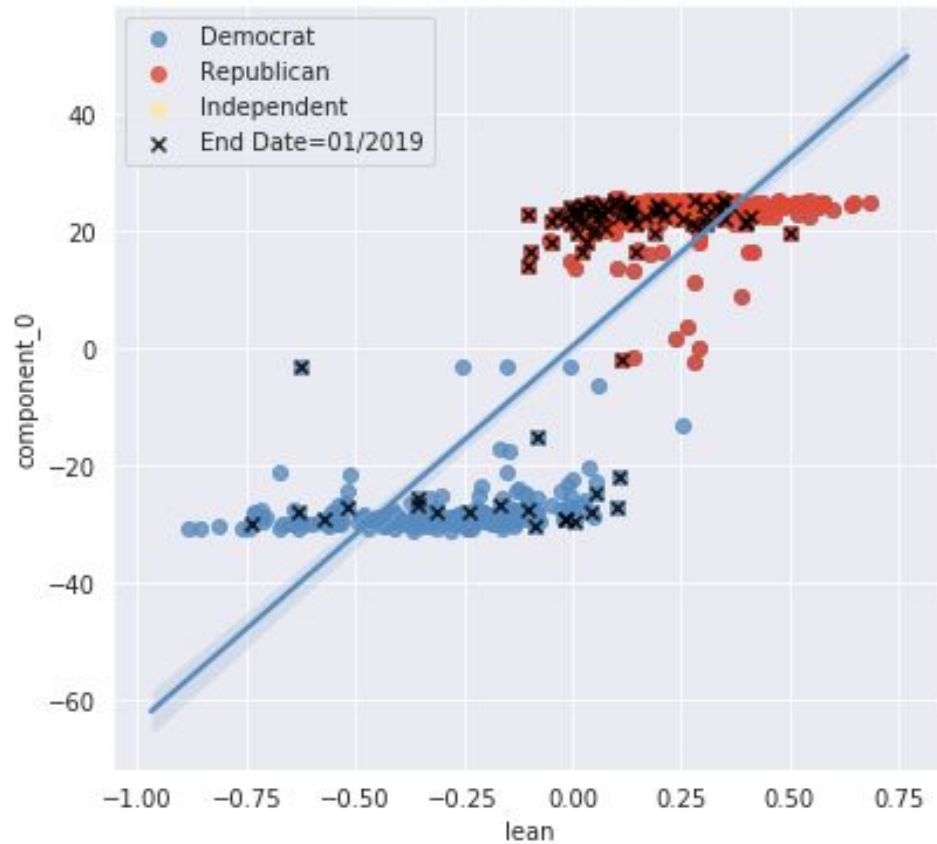


... and the Senate. It can probably be assumed that "component\_0" roughly corresponds with liberal vs. conservative voting records. "Component\_1" might align with social values, but this won't be clear without further analysis

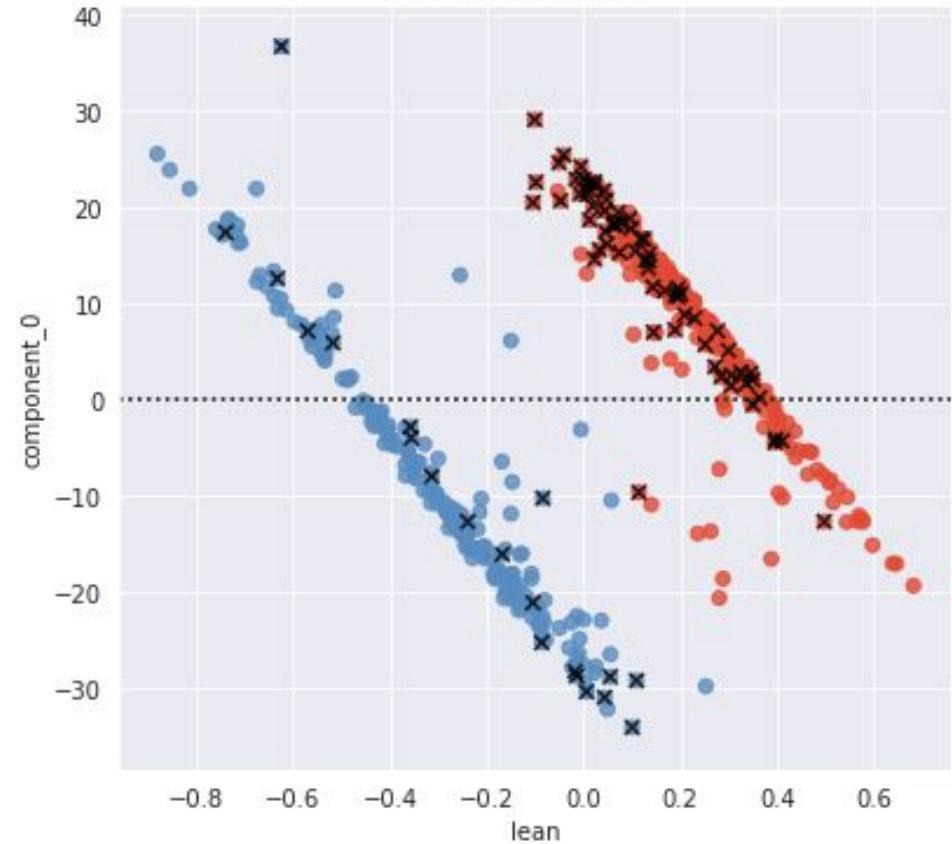


Partisan Lean v. Ideological Score, 2018 House

Linear Regression



Residuals/Errors

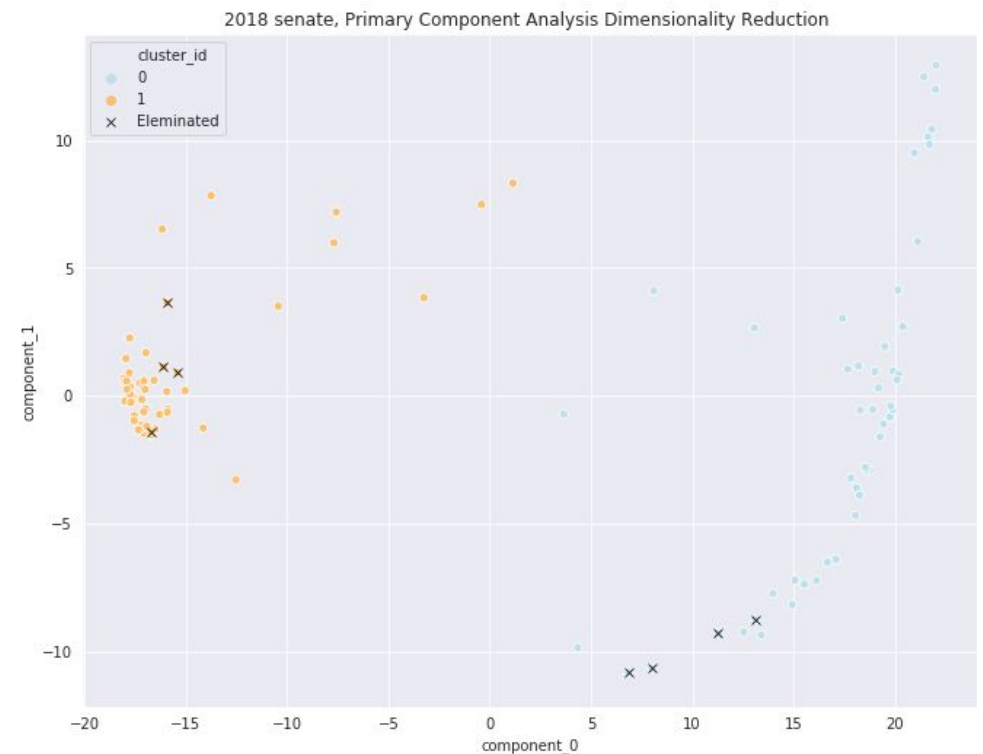
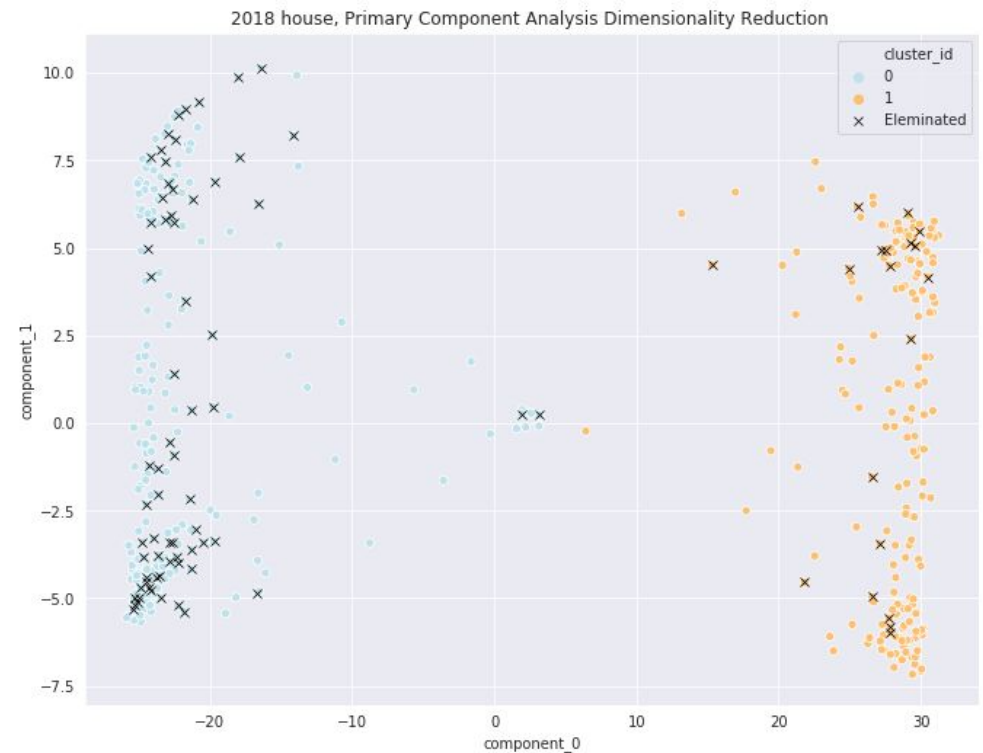


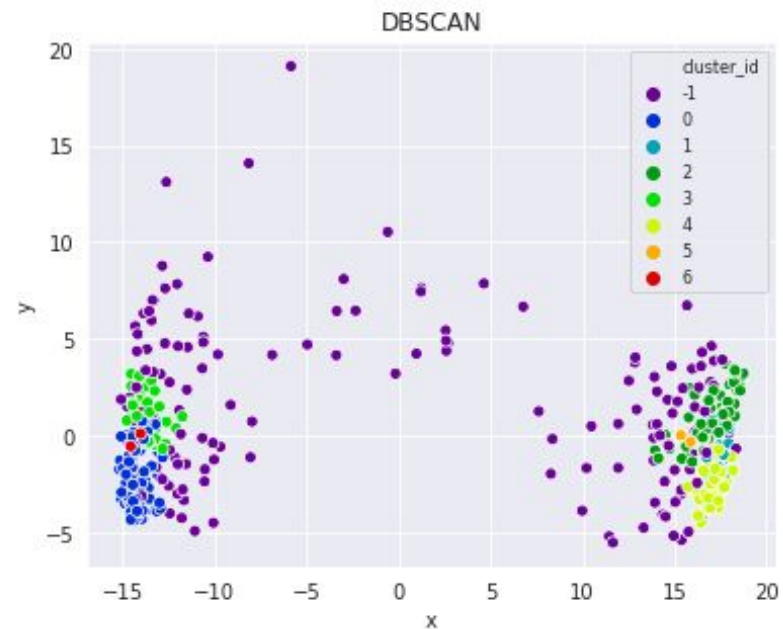
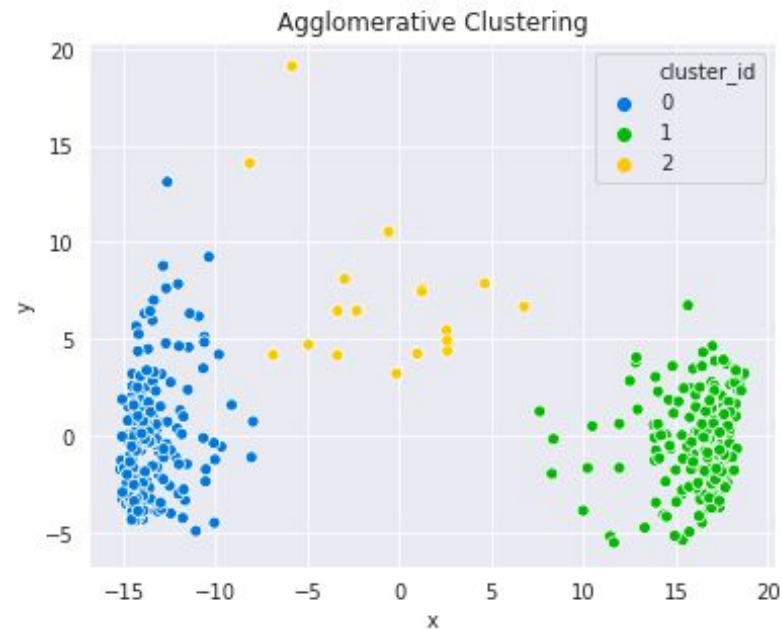
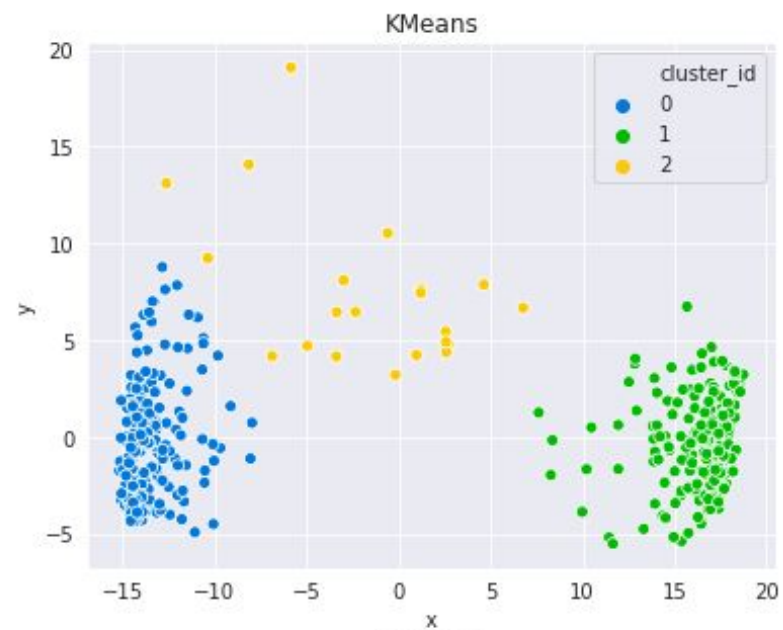
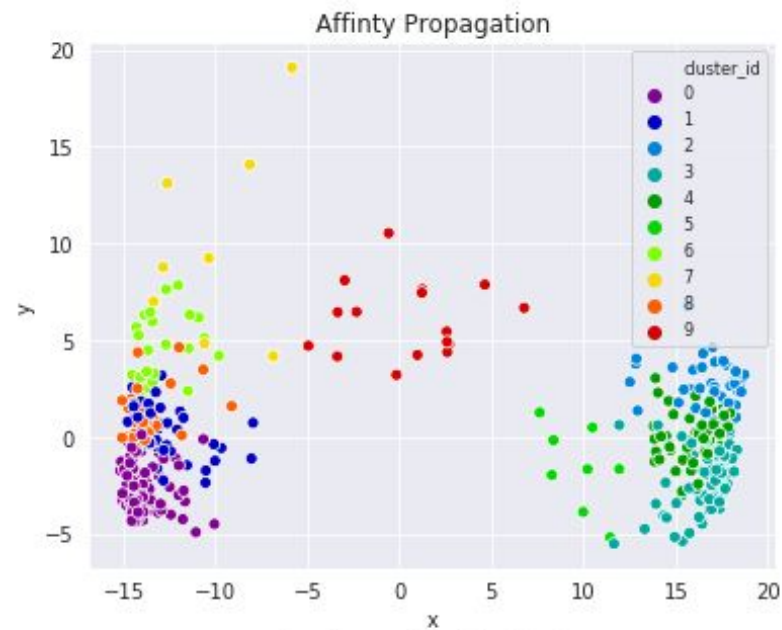
While primary "component\_0" seems to leave the parties a bit flat compared to `dw_nominate` scores, it appears that it may be an even better predictor of a given member's chance of being eliminated, when combined with partisan lean.

# Unsupervised Learning (Clustering)

Unsupervised machine learning distinguishes patterns in data without pre-labing it for training.

In this case, KMeans clustering is being used to easily distinguish the two parties from one another, based only on their dimension reduced voting behavior. Kmeans is only being told to find two clusters without labeling the political parties for it ahead of time, and has done so with only two misclassifications (both in the house).





Here, multiple unsupervised clustering methods are being used to label dimension reduced voting from the 2018 house. With more analysis, this method could possibly be used to tease out the voting differences between different distinct groups, and perhaps determine their differential survival.



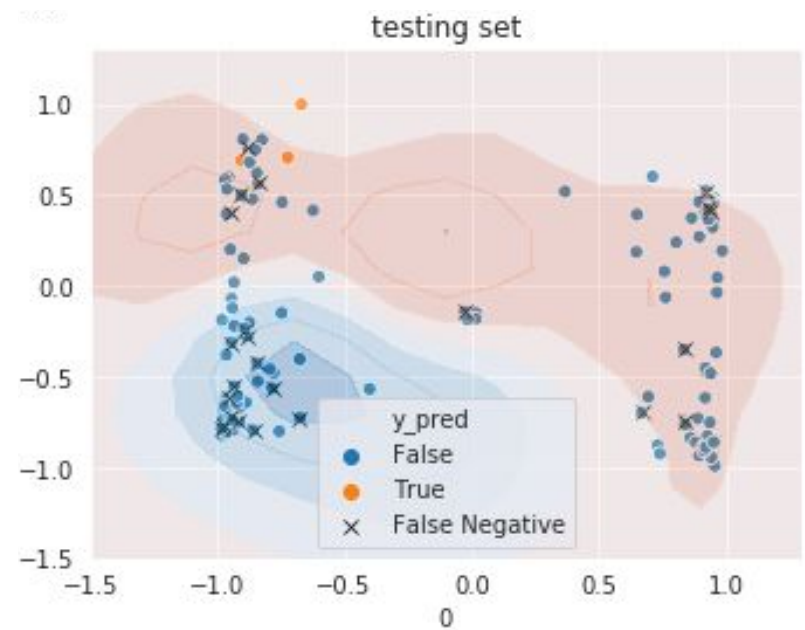
# Supervised learning (Classification)

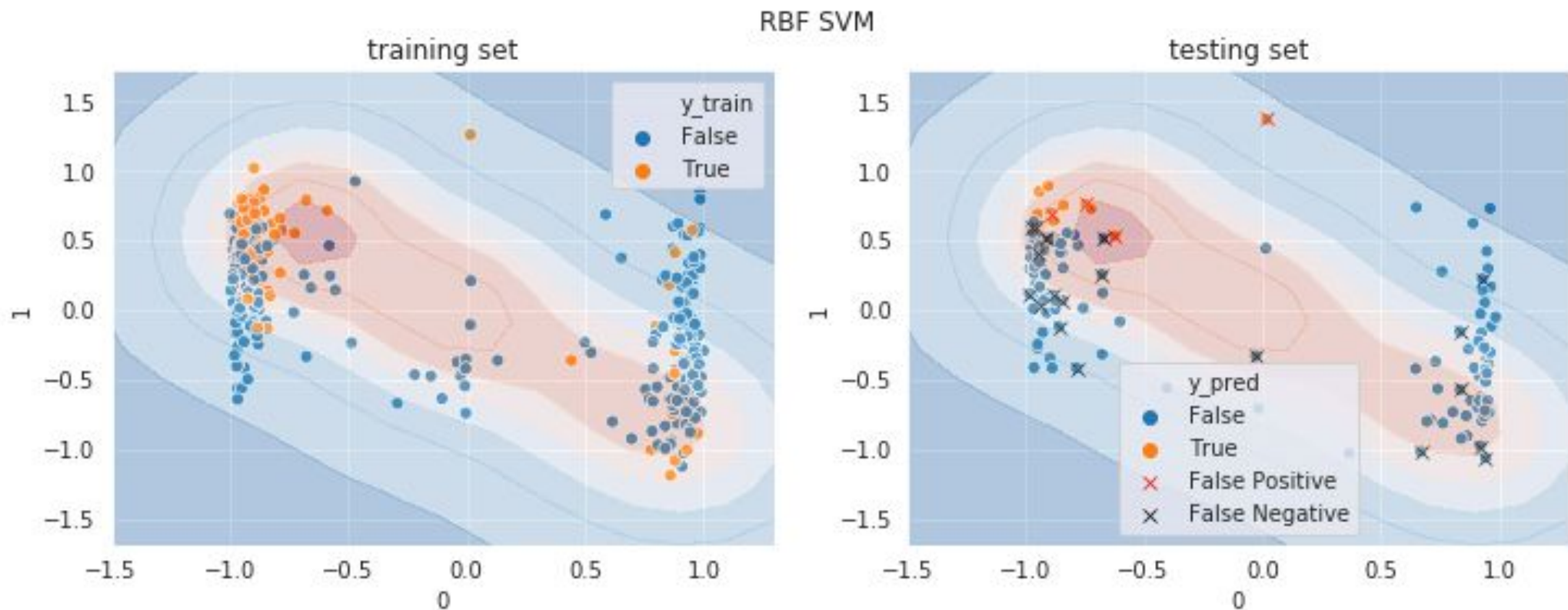
Here the RBF SVM classifier is being fed the same dimension reduced voting data from the 115th house, but is being trained to predict which members will be eliminated. In the top training set plot (top) members labeled “True”, were eliminated, and “False” were not.

In the testing plot (bottom), the classifier is giving it's best guess with the remaining 3rd of the data, with incorrect guesses being crossed out. The surrounding contours represent the decision function of the model on the two dimensional plane we can see.



Predict





In the first set of plots, the classifier is using 'component\_0' and 'component\_1' combined with a 3rd partisan lean dimension to classify who was eliminated. However, in this second set of plots, 'component\_1' has been replaced with a dimension corresponding to the residuals of a linear regression of component\_0 and district lean. This is to additionally represent how far out of line a congressperson's voting record is with the partisan lean of their district.

# Improvement against baseline - Precision

The training and testing sequence is done with each 3rd of the data, using the other two thirds for training, to create a 3-fold crossvalidation score average, which has been done here across several different classifiers

Replacing 'component\_1' with lean vs. component\_0 residuals results in an improvement in the metrics of most of the models, even when partisan lean is already included in the model.

"Precision" is the proportion of positive classifications that are correct.



Mean cross validation scores - without residual Dimension



Mean cross validation scores - with residual Dimension

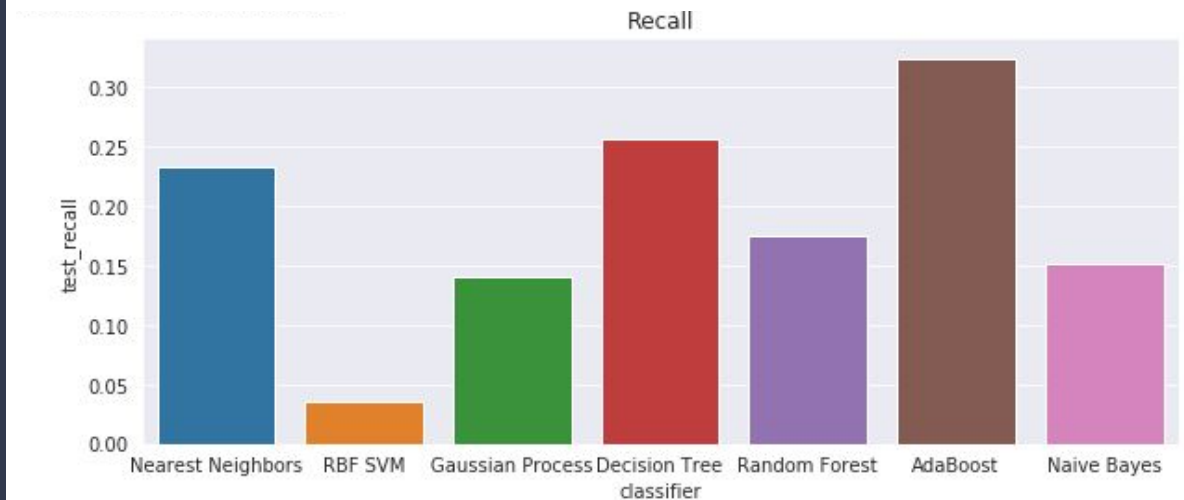


# Improvement against baseline - Recall

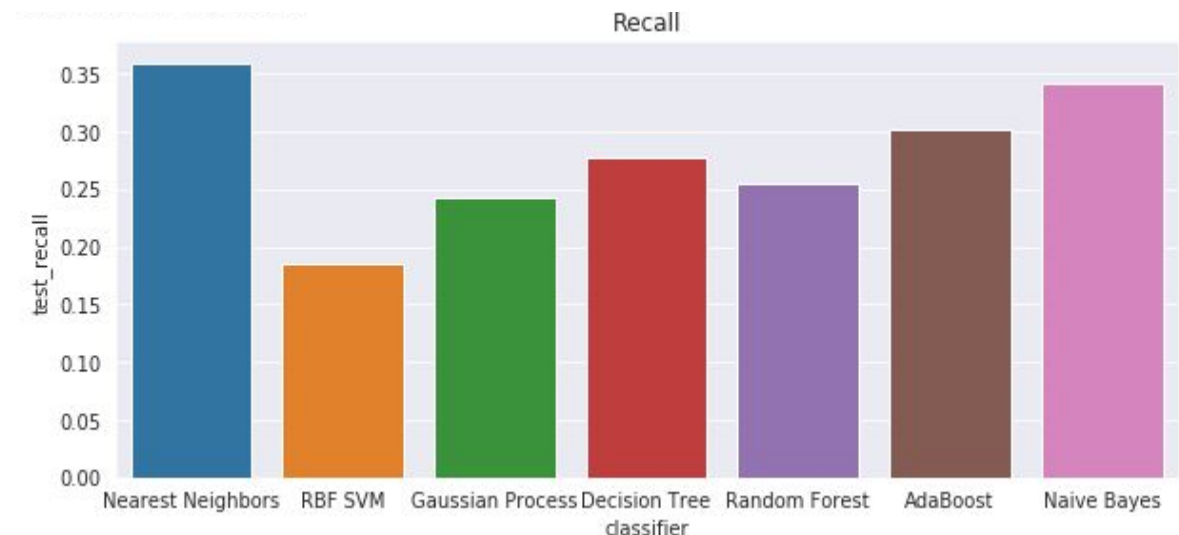
Here we can see the same improvement occurs with cross validated recall scores.

"Recall" is the proportion of total positive values that have been correctly classified.

Nearest Neighbors precision of .6 and recall of .35 means when it predicts a given house member is eliminated it is correct 60% of the time, and is able to identify 35% of those which were eliminated.



Mean cross validation scores - without residual dimension



Mean cross validation scores - with residual dimension

# Findings:

1. Voting behavior in line with the partisan lean of constituents is an important predictive indicator of congressional survival.
2. Members with voting records more “extreme” than the partisan lean of their constituents are less likely to survive than those with a more moderate one
3. There is a correlation between longer surviving members and voting records closer to the center of the ideological spectrum of their party.

# Recommendations:

1. To survive in an R+22 district, pursue one of the most moderate voting records in the house. Otherwise survival is highly unlikely.
2. House members take on less risk by voting more “moderate” than the partisan lean in their district than more “radical.” So it is better for the survival of the delegation for members in safe districts to vote less radical than it is for members in swing districts to be less moderate.

# Possible Next Steps

1. Create a fully fledged predictive model using 20 years of data.

As historical partisan lean from FiveThirtyEight and the Cook Political Report is not readily available, it may need to be constructed from federal elections returns. Select and tune the best performing classifier. Identify precise percentage chances of elimination for a given voting score and partisan lean

2. Investigate the delta (difference) between different unsupervised

**clusters** and different ends of the political spectrum to find precisely which bills are different, and therefore what dimension reduction is measuring. Explore the advantages/disadvantages of custom dimension reduction over pre-calculated `dw_nominate` scores.

3. Map which individual votes effected chances of elimination the most. This can probably done with Naive Bayes classifier without using dimension reduction.

# Appendix

1. Datasets
  - a. [ProPublica Congress API](#)
  - b. Partisan Lean - FiveThirtyEight
    - i. [District](#)
    - ii. [State](#)
2. Project on git-hub
  - a. [https://github.com/Nhorning/Congressional\\_Survival](https://github.com/Nhorning/Congressional_Survival)