

# State Legislature Election Modeling

In this project I aim to create models able to predict the outcome of political campaigns for state legislative seats in all 50 states. Through the work below I created models that predict when a seat is likely to be a tossup. This is especially useful to political organizers who can use this information in order to prioritize resource allocation in favor of close elections.

But first, a definition of terms. 49 of 50 states have bicameral legislatures, i.e. a higher house with fewer members and a lower house with more. Nebraska stands alone as the only state with a unicameral legislature. Its legislature is treated as an upper house in the data I used in this project. Every other state calls their upper house the Senate.

However, for lower houses things get a bit more creative. Most states simply call them Houses of Representatives, mirroring the federal legislature. A number of states, mostly older east coast states, have more unique names such as Assemblies or Houses of Delegates. The differences are largely vestigial, harkening back to colonial governments

There are other similar political bodies included in the census data used, including the District of Columbia Council and Puerto Rico Legislative Assembly. The purpose of this project is simply to model state legislatures in US states and these other bodies have their own separate dynamics that aren't necessarily relevant."

In [1]:

Requirement already satisfied: censusdata in /usr/local/lib/python3.7/dist-packages (1.15)

Requirement already satisfied: category\_encoders in /usr/local/lib/python3.7/dist-packages (2.3.0)

Requirement already satisfied: geopandas in /usr/local/lib/python3.7/dist-packages (0.10.2)

Requirement already satisfied: boruta in /usr/local/lib/python3.7/dist-packages (0.3)

Requirement already satisfied: catboost in /usr/local/lib/python3.7/dist-packages (1.0.3)

Requirement already satisfied: shap in /usr/local/lib/python3.7/dist-packages (0.40.0)

Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from censusdata) (2.23.0)

Requirement already satisfied: pandas in /usr/local/lib/python3.7/dist-packages (from censusdata) (1.1.5)

Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.7/dist-packages (from category\_encoders) (1.4.1)

Requirement already satisfied: patsy>=0.5.1 in /usr/local/lib/python3.7/dist-packages (from category\_encoders) (0.5.2)

Requirement already satisfied: numpy>=1.14.0 in /usr/local/lib/python3.7/dist-packages (from category\_encoders) (1.19.5)

Requirement already satisfied: scikit-learn>=0.20.0 in /usr/local/lib/python3.7/dist-packages (from category\_encoders) (1.0.1)

Requirement already satisfied: statsmodels>=0.9.0 in /usr/local/lib/python3.7/dist-packages (from category\_encoders) (0.10.2)

Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.7/dist-packages (from pandas->censusdata) (2018.9)

Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-packages (from pandas->censusdata) (2.8.2)

Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from patsy>=0.5.1->category\_encoders) (1.15.0)

Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.20.0->category\_encoders) (3.0.0)

Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.20.0->category\_encoders) (1.1.0)

Requirement already satisfied: pyproj>=2.2.0 in /usr/local/lib/python3.7/dist-packages (from geopandas) (3.2.1)

Requirement already satisfied: shapely>=1.6 in /usr/local/lib/python3.7/dist-packages (from geopandas) (1.8.0)

Requirement already satisfied: fiona>=1.8 in /usr/local/lib/python3.7/dist-packages (from geopandas) (1.8.20)

Requirement already satisfied: click-plugins>=1.0 in /usr/local/lib/python3.7/dist-packages (from fiona>=1.8->geopandas) (1.1.1)

Requirement already satisfied: munch in /usr/local/lib/python3.7/dist-packages (from fiona>=1.8->geopandas) (2.5.0)

Requirement already satisfied: setuptools in /usr/local/lib/python3.7/dist-packages (from fiona>=1.8->geopandas) (57.4.0)

Requirement already satisfied: cligj>=0.5 in /usr/local/lib/python3.7/dist-packages (from fiona>=1.8->geopandas) (0.7.2)

Requirement already satisfied: attrs>=17 in /usr/local/lib/python3.7/dist-packages (from fiona>=1.8->geopandas) (21.2.0)

Requirement already satisfied: click>=4.0 in /usr/local/lib/python3.7/dist-packages (from fiona>=1.8->geopandas) (7.1.2)

Requirement already satisfied: certifi in /usr/local/lib/python3.7/dist-packages

```
t-packages (from fiona>=1.8->geopandas) (2021.10.8)
Requirement already satisfied: plotly in /usr/local/lib/python3.7/dist-
-packages (from catboost) (4.4.1)
Requirement already satisfied: graphviz in /usr/local/lib/python3.7/di
st-packages (from catboost) (0.10.1)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.7/
dist-packages (from catboost) (3.2.2)
Requirement already satisfied: packaging>20.9 in /usr/local/lib/python
```

```
In [2]: 1 # enable text wrapping on Google Colab
2 from IPython.display import HTML, display
3
4 def set_css():
5     display(HTML('''
6     <style>
7         pre {
8             white-space: pre-wrap;
9         }
10    </style>
11    '''))
12 get_ipython().events.register('pre_run_cell', set_css)
13
14 from warnings import simplefilter
15
16 # ignore all warnings
17 simplefilter(action='ignore')
18
19 import pandas as pd
20 import censusdata
21 import openpyxl
22 import re
23 import os
24 import numpy as np
25
26 import matplotlib.pyplot as plt
27
28 from sklearn.preprocessing import StandardScaler
29 from sklearn.ensemble import *
30 from sklearn.tree import DecisionTreeClassifier
31 from sklearn.linear_model import LogisticRegression
32 from sklearn.model_selection import RandomizedSearchCV,
33    train_test_split, validation_curve
34 from sklearn.metrics import classification_report,
35    plot_confusion_matrix, recall_score
36
37 from datetime import datetime
38
39 import category_encoders as ce
40
41 import geopandas as gpd
42
43 from boruta import BorutaPy
44
45 from imblearn.pipeline import Pipeline
46
47 import xgboost as xgb
```

```
46 import lightgbm as lgb
47
```

## Data Exploration

I started with election result data gathered by Carl Klarner and available on online from [Harvard University \(https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/3WZFK9&widget=dataverse@harvard\)](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/3WZFK9&widget=dataverse@harvard). This dataset aims to capture every state-level legislative election from 1967 to 2016. I decided to limit the scope of this project to within the last 10 years for both technical and theoretical reasons. On the technical side, the way in which the Census Bureau presents their American Community Survey data products has changed since their earliest offering on their API in 2005, making it difficult to compare past and present data. In terms of theory, I subscribe to the notion that the ways in which demographics and politics interact are changing. [One example of this is the sharp turn towards the Democratic Party in American suburbs which used to be reliably Republican. \(https://fivethirtyeight.com/features/why-the-suburbs-have-shifted-blue/\)](https://fivethirtyeight.com/features/why-the-suburbs-have-shifted-blue/)

As a result of these two factors I fear including too much historical data could compromise my models' ability to understand the data in such a way as to usefully predict future election results.

```
In [3]: 1 harvard_data = pd.read_csv('/content/drive/MyDrive/capstone
        /harvard_data_2000+.csv', index_col=0)
```

```
Out[3]:
```

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
2825	2002	al	1	0	g	starkey, nelson jr.	democrat	inc	8540.0	w
2826	2002	al	1	0	g	franklin, joey	libertarian	none	304.0	l
2827	2002	al	1	0	g	mcnatt, william dale	modernrepublican	none	4772.0	l
2828	2002	al	1	0	g	writein	99999	none	28.0	l
2829	2002	al	2	0	g	pettus, mary	democrat	none	5718.0	l
...	...	...	...	...	...	...	...	...	...	...
378340	2016	wy	28	1	g	anderson, james lee	republican	inc	5216.0	w
378341	2016	wy	28	1	g	scattering	writein	none	22.0	l
378342	2016	wy	30	1	g	ford, robert	democrat	none	1521.0	l
378343	2016	wy	30	1	g	scott, charles k.	republican	inc	5831.0	w
378344	2016	wy	30	1	g	scattering	writein	none	45.0	l

142905 rows × 10 columns



```
In [4]:
```

```
Out[4]:
```

```
Index(['year', 'sab', 'ddez', 'sen', 'etype', 'cand', 'party', 'exper',  
      'vote',  
      'outcome'],
```

In [5]:

```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 142905 entries, 2825 to 378344  
Data columns (total 10 columns):  
#   Column      Non-Null Count  Dtype  
---  -  
0   year        142905 non-null  int64  
1   sab         142905 non-null  object  
2   ddez        142905 non-null  object  
3   sen         142905 non-null  int64  
4   etype       142905 non-null  object  
5   cand        142905 non-null  object  
6   party       142748 non-null  object  
7   exper       142905 non-null  object  
8   vote        137078 non-null  float64  
9   outcome     142905 non-null  object  
dtypes: float64(1), int64(2), object(7)  
memory usage: 12.0+ MB
```

In [6]:

Out[6]:

	year	sen	vote
count	142905.000000	142905.000000	137078.000000
mean	2009.537287	0.225730	10098.985629
std	5.239091	0.418064	15806.774287
min	2000.000000	0.000000	0.000000
25%	2005.000000	0.000000	1583.250000
50%	2010.000000	0.000000	4885.000000
75%	2014.000000	0.000000	12522.000000
max	2017.000000	1.000000	326755.000000

In [7]:

```
1 harvard_data[(harvard_data['year'] == 2012) & (harvard_data['sab'] ==  
2 'ny') & (harvard_data['sen'] == 1)  
3 & (harvard_data['ddez'] == '1')]
```

Out[7]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
246627	2012	ny	1	1	g	fleming, bridget m.	democrat	none	46783.0	l
246628	2012	ny	1	1	g	fleming, bridget m.	workingfamilies	none	4240.0	l
246629	2012	ny	1	1	g	lavalle, kenneth p.	independent	inc	4191.0	w

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>246630</b>	2012	ny	1	1	g	lavalle, kenneth p.	conservative	inc	10462.0	w
<b>246631</b>	2012	ny	1	1	g	lavalle, kenneth p.	modernrepublican	inc	61130.0	w

```
In [8]: 1 harvard_data[(harvard_data['sab'] == 'al') & (harvard_data['sen'] ==  
2         0)  
        & (harvard_data['ddez'] == '1')  
        ]
```

Out[8]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>2825</b>	2002	al	1	0	g	starkey, nelson jr.	democrat	inc	8540.0	w
<b>2826</b>	2002	al	1	0	g	franklin, joey	libertarian	none	304.0	l
<b>2827</b>	2002	al	1	0	g	mcnatt, william dale	modernrepublican	none	4772.0	l
<b>2828</b>	2002	al	1	0	g	writein	99999	none	28.0	l
<b>3105</b>	2006	al	1	0	g	irons, tammy	democrat	none	8410.0	w
<b>3106</b>	2006	al	1	0	g	smith, william	modernrepublican	none	4511.0	l
<b>3257</b>	2010	al	1	0	g	burdine, greg	democrat	none	7083.0	w
<b>3258</b>	2010	al	1	0	g	hanson, quinton	modernrepublican	none	6877.0	l
<b>3259</b>	2010	al	1	0	g	writein	independentwritein	none	14.0	l
<b>3517</b>	2014	al	1	0	dprfsettled	burdine, greg	democrat	inc	NaN	w
<b>3518</b>	2014	al	1	0	rprfsettled	pettus, phillip	modernrepublican	none	1172.0	w
<b>3519</b>	2014	al	1	0	rprfsettled	statom, sterling (josh)	modernrepublican	none	909.0	l
<b>3727</b>	2014	al	1	0	g	burdine, greg	democrat	inc	4652.0	l
<b>3728</b>	2014	al	1	0	g	pettus, phillip	modernrepublican	none	4933.0	w
<b>3729</b>	2014	al	1	0	g	scattering	9999	none	10.0	l

```
In [9]: 1 dupes = harvard_data[(harvard_data['cand'] != 'scattering') &  
        (harvard_data['etype'] == 'g') & (harvard_data.drop(['vote', 'party'],  
axis=1).duplicated(keep=False))] .dropna(subset=['vote'])
```

Out[9]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>3730</b>	2014	al	2	0	g	betterton, andew (andy)	democrat	none	549.0	l

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>3731</b>	2014	al	2	0	g	betterton, andrew (andy)	democrat	none	4675.0	l
<b>3732</b>	2014	al	2	0	g	greer, lynn	modernrepublican	inc	1428.0	w
<b>3733</b>	2014	al	2	0	g	greer, lynn	modernrepublican	inc	7133.0	w
<b>3736</b>	2014	al	3	0	g	black, marcel	democrat	inc	262.0	w
...	...	...	...	...	...	...	...	...	...	...
<b>378322</b>	2016	wy	20	1	g	agar, wyatt	republican	none	1225.0	w
<b>378323</b>	2016	wy	20	1	g	agar, wyatt	republican	none	1817.0	w
<b>378324</b>	2016	wy	20	1	g	agar, wyatt	republican	none	3157.0	w
<b>378330</b>	2016	wy	22	1	g	kinskey, dave	republican	none	3746.0	w
<b>378331</b>	2016	wy	22	1	g	kinskey, dave	republican	none	3857.0	w

In [10]:

```
Out[10]: ny      7839
         ct      1590
         va      1046
         me      1021
         in       929
         il       914
         wv       884
         nc       747
         nj       735
         ky       715
         Name: sab, dtype: int64
```

In [11]:

Out[11]: 47

The data as collected presents a number of issue's we have to overcome. First, individual rows appear to reflect individual candidates but split between their vote totals in different counties or in the case of fusion voting in New York, accross different parties. Rows will have to be combined in order to reflect individual elections so we can project elections outcomes.

In [12]:

```
1 nj_multimember = harvard_data[(harvard_data['outcome'] == 'w') &
    (harvard_data['sab'] == 'nj')][harvard_data.duplicated(subset=['year',
    'sab', 'ddez', 'sen', 'etype', 'outcome'], keep=False)]
```

Out[12]:

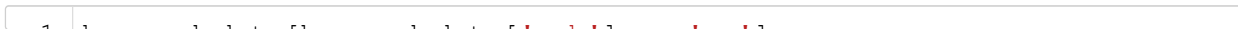
	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>217767</b>	2001	nj	1	0	g	vandrew, jeff	democrat	none	32271.0	w
<b>217768</b>	2001	nj	1	0	g	asselta, nicholas	modernrepublican	inc	36392.0	w
<b>217773</b>	2001	nj	2	0	g	blee, frank	modernrepublican	inc	29010.0	w

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>217774</b>	2001	nj	2	0	g	damato, paul r.	modernrepublican	none	29427.0	w
<b>217775</b>	2001	nj	3	0	g	burzichelli, john j.	democrat	none	30213.0	w
...	...	...	...	...	...	...	...	...	...	...
<b>221114</b>	2013	nj	40	1	g	otoole, kevin j.	modernrepublican	inc	3751.0	w
<b>221115</b>	2013	nj	40	1	g	otoole, kevin j.	modernrepublican	inc	14674.0	w
<b>221116</b>	2013	nj	40	1	g	otoole, kevin j.	modernrepublican	inc	16752.0	w
<b>221118</b>	2015	nj	5	1	gs	cruzperez, nilsa	democrat	pastother	7979.0	w
<b>221119</b>	2015	nj	5	1	gs	cruzperez, nilsa	democrat	pastother	11171.0	w

1037 rows × 10 columns

Second, some states use mutli-member districts, i.e. more than one winner per election. The original dataset accounted for this through a dtype feature which aimed to keep track of individual seats by assigning arbitrary numbers to them. However, given that each seat in multimember districts has the exact same electorate, any artificial consistency on individual seats seems to me to be unnecessary information as far as the type of demographic analysis I aim to do goes.

In [13]:



Out[13]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>140623</b>	2000	ma	First Barnstable	0	g	george, thomas n.	modernrepublican	inc	19302.0	w
<b>140624</b>	2000	ma	Second Barnstable	0	g	atsalis, demetrius	democrat	inc	17044.0	w
<b>140625</b>	2000	ma	Second Barnstable	0	g	zalis, lawrence a.	modernrepublican	none	6327.0	l
<b>140626</b>	2000	ma	Third Barnstable	0	g	patrick, matthew	democrat	none	13285.0	w
<b>140627</b>	2000	ma	Third Barnstable	0	g	vieira, david joseph	modernrepublican	none	10508.0	l
...	...	...	...	...	...	...	...	...	...	...
<b>145423</b>	2016	ma	Worcester & Norfolk	1	g	fattman, ryan c.	republican	inc	64665.0	w
<b>145424</b>	2016	ma	Worcester & Norfolk	1	g	scattering	scattering	none	1021.0	l
<b>145425</b>	2016	ma	Plymouth & Norfolk	1	s	meschino, joan	democratic	none	8108.0	l

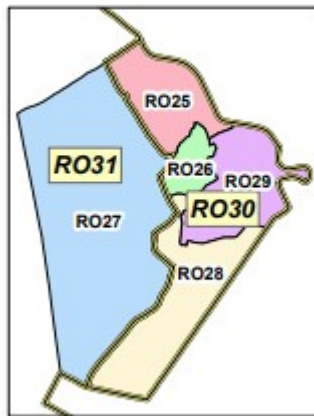


	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>145426</b>	2016	ma	Plymouth & Norfolk	1	s	oconnor, patrick m.	republican	none	9047.0	w
<b>145427</b>	2016	ma	Plymouth & Norfolk	1	s	scattering	scattering	none	30.0	l

2467 rows x 11 columns

Third, certain New England states don't number their districts, preferring county or locality designations. This is further complicated by New Hampshire's floterial districts. Put simply, some districts in New Hampshire's lower house consist of numerous other lower house districts put together. For example, see this map of Portsmouth's house districts.

### Portsmouth House Districts



```
In [14]: 1 harvard_data[harvard_data['outcome'] ==
```

```
Out[14]: modernrepublican      35132
democrat      32710
republican      5121
workingfamilies  1341
conservative      1104
independent      998
democratfarmerlabor    892
democratic      709
independence      342
democraticparty    337
indep      247
4002      225
reform      157
working      110
republicananddemocrat    95
liberal      93
unspecifiednonmajor    90
democraticfarmerlabor    90
libertarian      72
democratandrepublican    67
Name: party, dtype: int64
```

Yet another issue is the way candidates' parties are tracked. As previously mentioned, New York

uses a fusion voting system allowing for a single candidate to run with multiple parties in the same race. The results record the number of votes each candidate gets with each respective party but as far as determining a winner the only thing that matters is the individual running, not their disparate party-lines.

However, other states do seem to allow for cross endorsement. One especially common feature in this data is candidates running as both democrats and republicans in Pennsylvania. They might be listed as "republicananddemocrat", "democratandrepublican" or other less common combinations.

Furthermore even pulling a list of people who only ran as Democrats isn't straight forward. Democrats in this data were recorded as "democrat," "democratfarmlabor" (Minnesota), "democratic", "democraticparty", and "democraticfarmlabor" (also Minnesota). And that's just among the 20 most common winning party labels.

To deal with this I had to do some research on various candidates that had uncommon party labels that by definition seemed to imply one of the two major parties but was written in some different way. I assembled lists of ways I would accept party listings as either with the Democratic or Republican parties.

Candidates running as a member of one of the major parties as well as some minor parties were counted as running with that major party. However, candidates running with endorsements from both major parties were excluded from the data.

```
In [15]: 1 republican = ['republican',
2                 'republican/independent',
3                 'republicanwritein',
4                 'writeinrepublican',
5                 'modernrepublican',
6                 'independentrepublicanparty',
7                 'independentrepublican',
8                 'gop/independentparty',
9                 'independentandrepublican'
10                ]
11
12 democratic = ['progressive/democratic/workingfamily',
13               'progressive/democratic',
14               'progressivedemocrat',
15               'writeindemocrat',
16               'independentanddemocrat',
17               'independentdemocrat',
18               'independentdemparty', 'democrat',
19               'democrat/progressive',
20               'democrat/workingfamilies',
21               'democratfarmerlabor',
22               'democratic',
23               'democratic/progressive',
24               'democraticfarmerlabor',
25               'democraticparty',
26               'democraticwritein',
27               'democratparty'
28              ]
```

First we'll deal with the issue of results being split across counties and parties as a result of different reporting methods. For sanity's sake, I've highlighted a candidate whose rows need merging to make sure the function works as intended.

In [16]:

Out[16]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>237384</b>	2010	ny	1	0	g	losquadro, daniel p.	unspecifiednonmajor	none	353.0	w
<b>237385</b>	2010	ny	1	0	g	losquadro, daniel p.	conservative	none	4752.0	w
<b>237386</b>	2010	ny	1	0	g	losquadro, daniel p.	modernrepublican	none	18755.0	w
<b>238124</b>	2012	ny	2	0	g	losquadro, daniel p.	independent	inc	1941.0	w
<b>238125</b>	2012	ny	2	0	g	losquadro, daniel p.	conservative	inc	5206.0	w
<b>238126</b>	2012	ny	2	0	g	losquadro, daniel p.	modernrepublican	inc	27158.0	w

In [17]:

```

1 def fusion_county_voting(row):
2     '''Takes in groupby-ed rows, returns single row with single
3     standard major party ID.'''
4     if row['party'].isin(democratic).any() and not
5     row['party'].isin(republican).any():
6         row['party'] = 'democratic'
7     elif row['party'].isin(republican).any() and not
8     row['party'].isin(democratic).any():
9         row['party'] = 'republican'
10    else:
11        row['party'] = np.nan
12
13    row['vote'] = row['vote'].sum()
14    return row
15
16 merged_df = dupes.groupby(['year', 'cand', 'sab', 'ddez', 'sen',
17                             'outcome'],
18                             as_index=False).apply(fusion_county_voting).drop_duplicates()

```

Out[17]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
<b>3730</b>	2014	al	2	0	g	betterton, andew (andy)	democratic	none	5224.0	l
<b>3732</b>	2014	al	2	0	g	greer, lynn	republican	inc	8561.0	w
<b>3736</b>	2014	al	3	0	g	black, marcel	democratic	inc	7993.0	w
<b>3739</b>	2014	al	3	0	g	joly, fred	republican	none	5357.0	l
<b>3745</b>	2014	al	4	0	g	hammon, micky	republican	inc	8473.0	w

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
...	...	...	...	...	...	...	...	...	...	...
378303	2016	wy	16	1	g	kusaba, richard	democratic	none	1989.0	l
378306	2016	wy	16	1	g	dockstader, dan	republican	inc	7208.0	w
378315	2016	wy	20	1	g	norskog, mary jane	democratic	none	1546.0	l
378320	2016	wy	20	1	g	agar, wyatt	republican	none	6893.0	w
378330	2016	wy	22	1	g	kinskey, dave	republican	none	7603.0	w

And after...

In [18]:

Out[18]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
237384	2010	ny	1	0	g	losquadro, daniel p.	republican	none	23860.0	w
238124	2012	ny	2	0	g	losquadro, daniel p.	republican	inc	34305.0	w

In [19]:

Out[19]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
3730	2014	al	2	0	g	betterton, andew (andy)	democratic	none	5224.0	l
3732	2014	al	2	0	g	greer, lynn	republican	inc	8561.0	w
3736	2014	al	3	0	g	black, marcel	democratic	inc	7993.0	w
3739	2014	al	3	0	g	joly, fred	republican	none	5357.0	l
3745	2014	al	4	0	g	hammon, micky	republican	inc	8473.0	w
...	...	...	...	...	...	...	...	...	...	...
5771	2014	al	29	1	g	mcclendon, melinda	republican	none	16145.0	w
5777	2014	al	30	1	g	morgan	NaN	none	5653.0	l
5782	2014	al	30	1	g	chambliss, clyde	republican	none	22916.0	w
5792	2014	al	31	1	g	greenwood, larry	democratic	none	8627.0	l
5796	2014	al	31	1	g	holley, jimmy w.	republican	inc	23067.0	w

125 rows × 10 columns

Looks good to me! Let's merge these results with the results that didn't have this duplicate problem.

In [20]:

1

```
election_data_merged = pd.concat([merged_df,
harvard_data.drop(dupes.index)])
```

Out[20]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome
3730	2014	al	2	0	g	betterton, andrew (andy)	democratic	none	5224.0	l
3732	2014	al	2	0	g	greer, lynn	republican	inc	8561.0	w
3736	2014	al	3	0	g	black, marcel	democratic	inc	7993.0	w
3739	2014	al	3	0	g	joly, fred	republican	none	5357.0	l
3745	2014	al	4	0	g	hammon, micky	republican	inc	8473.0	w
...	...	...	...	...	...	...	...	...	...	...
378340	2016	wy	28	1	g	anderson, james lee	republican	inc	5216.0	w
378341	2016	wy	28	1	g	scattering	writein	none	22.0	l
378342	2016	wy	30	1	g	ford, robert	democrat	none	1521.0	l
378343	2016	wy	30	1	g	scott, charles k.	republican	inc	5831.0	w
378344	2016	wy	30	1	g	scattering	writein	none	45.0	l

123816 rows × 10 columns

```
In [21]: 1 election_data_merged[(election_data_merged['cand'] != 'scattering') &
(election_data_merged['etype'] ==
'g') & (election_data_merged.drop(['vote', 'party'],
axis=1)
```

Out[21]:

```
year  sab  ddez  sen  etype  cand  party  exper  vote  outcome
```

Now I'll combine rows so each one represents one election. At this point I'll be dropping exclusively third party candidates. While they do occasionally win elections and in some states are even viable rivals, they remain rare exceptions. Where they do exist they tend to be regional and canvass with one of the two major parties. Including them would do little to inform about the dynamics of demographics as they influence partisan American elections.

```
In [22]: 1 election_data_merged['district'] = election_data_merged.apply(lambda
row: ('u' if row['sen'] == 1 else 'l') + '{:0>3}'.format(row['ddez']))
axis=1)
```

Out[22]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
3730	2014	al	2	0	g	betterton, andrew (andy)	democratic	none	5224.0	l	l002
3732	2014	al	2	0	g	greer, lynn	republican	inc	8561.0	w	l002
3736	2014	al	3	0	g	black, marcel	democratic	inc	7993.0	w	l003

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
<b>3739</b>	2014	al	3	0	g	joly, fred	republican	none	5357.0	l	l003
<b>3745</b>	2014	al	4	0	g	hammon, micky	republican	inc	8473.0	w	l004
...	...	...	...	...	...	...	...	...	...	...	...
<b>378340</b>	2016	wy	28	1	g	anderson, james lee	republican	inc	5216.0	w	u028
<b>378341</b>	2016	wy	28	1	g	scattering	writein	none	22.0	l	u028
<b>378342</b>	2016	wy	30	1	g	ford, robert	democrat	none	1521.0	l	u030
<b>378343</b>	2016	wy	30	1	g	scott, charles k.	republican	inc	5831.0	w	u030
<b>378344</b>	2016	wy	30	1	g	scattering	writein	none	45.0	l	u030

For rows not included in the county/fusion voting cleanup, they need their party names standardized.

```
In [23]: 1 party_map = {}
2
3 for each in democratic:
4     party_map[each] = 'democratic'
5 for each in republican:
6     party_map[each] = 'republican'
7
8 election_data_merged['party'] =
```

```
In [24]: 1 # drop non-general elections
2 election_data_merged =
  election_data_merged[election_data_merged['etype'] ==
```

```
In [25]:
```

Out[25]:

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
<b>0</b>	2014	al	2	0	g	betterton, andew (andy)	democratic	none	5224.0	l	l002
<b>1</b>	2014	al	2	0	g	greer, lynn	republican	inc	8561.0	w	l002
<b>2</b>	2014	al	3	0	g	black, marcel	democratic	inc	7993.0	w	l003
<b>3</b>	2014	al	3	0	g	joly, fred	republican	none	5357.0	l	l003
<b>4</b>	2014	al	4	0	g	hammon, micky	republican	inc	8473.0	w	l004
...	...	...	...	...	...	...	...	...	...	...	...
<b>106945</b>	2016	wy	28	1	g	anderson, james lee	republican	inc	5216.0	w	u028

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
106946	2016	wy	28	1	g	scattering	NaN	none	22.0	l	u028
106947	2016	wy	30	1	g	ford, robert	democratic	none	1521.0	l	u030
106948	2016	wy	30	1	g	scott, charles k.	republican	inc	5831.0	w	u030
106949	2016	wy	30	1	g	scattering	NaN	none	45.0	l	u030

106950 rows x 11 columns

## Feature Generation

### Restricting Data by Years

At this point I decided to create a date cutoff for my data. District lines are generally redrawn every 10 years and in the year following the decennial census. As a result years ending in 2 tend to be the first ones featuring updated maps that last for roughly the next 10 years. I know I want this project to reflect current election trends first and foremost so it was never an option to use the completed 2010s decade cycle.

For the end date unfortunately no readily available database appears to exist that lists results for state-level legislative elections past 2016. I opted to collect some further data about New York (my home state) to use as holdout data but it would take a significant effort to do the same for all remaining data.

```
In [26]: 1 years = [2012, 2013, 2014, 2015, 2016]
2 states = election_data_merged['sab'].unique()
3 final_election_data = pd.DataFrame()
4
5 for year in years:
6     for state in states:
7         # districts vary by state
8         districts = election_data_merged[(election_data_merged['sab']
== state) & (election_data_merged['year'] == year)]
9         ['district'].unique()
10        for district in districts:
11            # excludes third party exclusive candidates. multiparty
candidates should already be combined into their largest representing
party
12            election =
election_data_merged[(election_data_merged['year'] == year) &
(election_data_merged['sab'] == state) &
(election_data_merged['district'] == district) &
(election_data_merged['party'].isin(democratic+republican))]
13            # skip if election dataframe doesn't have at least 2
parties and at least 1 winner
14            if election.empty or 'w' not in
election['outcome'].unique() or len(election['party'].unique()) == 1:
15                pass
16            else:
```

```

        # number of seats is defined by the number of winning
major party candidates per district per year
17         num_of_seats = election['outcome'].value_counts()['w']
18         if len(election) < num_of_seats * 2 or (num_of_seats
!= election['party'].isin(republican).value_counts()[True] or
num_of_seats != election['party'].isin(democratic).value_counts()
[True]):
19             pass
20         else:
21             for seat in range(num_of_seats):
22                 dem_vote =
election[election['party'].isin(democratic)].iloc[seat]['vote'] /
(election[election['party'].isin(democratic)].iloc[seat]['vote'] +
election[election['party'].isin(republican)].iloc[seat]['vote'])
23                 dem_incumbent =
election[election['party'].isin(democratic)]['exper'].iloc[0]
24                 if dem_incumbent == 'inc':
25                     dem_incumbent = True
26                 else:
27                     dem_incumbent = False
28                 rep_incumbent =
election[election['party'].isin(republican)]['exper'].iloc[0]
29                 if rep_incumbent == 'inc':
30                     rep_incumbent = True
31                 else:
32                     rep_incumbent = False
33
34                 if dem_incumbent == True:
35                     previous_winner_dem = True
36                 elif rep_incumbent == True:
37                     previous_winner_dem = False
38                 else:
39                     # if no incumbent - check who won 2 years
ago unless that would be a different district map
40                     if year - 2 > 2011:
41                         prev_election =
election_data_merged[(election_data_merged['year'] == year - 2) &
(election_data_merged['sab'] == state) &
(election_data_merged['district'] == district)]
42                     # west virginia has multimember
staggared state senate elections
43                     if prev_election.empty or state ==
'wv' and district.startswith('u'):
44                         # or in some cases, 4 years ago
45                         if year - 4 > 2011:
46                             prev_election =
election_data_merged[(election_data_merged['year'] == year - 4) &
(election_data_merged['sab'] == state) &
(election_data_merged['district'] == district)]
47                             if prev_election.empty:
48                                 # or in the case of
certain Montana Senate seats, 6 years
49                                 if state == 'mt' and year
- 6 > 2011:
50

```



```

prev_election =
election_data_merged[(election_data_merged['year'] == year - 6) &
(election_data_merged['sab'] == state) &
(election_data_merged['district'] == district)]
51         if
prev_election.empty:
52     previous_winner_dem = np.nan
53         else:
54     previous_winner_dem =
np.nan
55         # we only care about
democratic vs republican
56         elif
prev_election[(prev_election['party'].isin(democratic+republican))].e
pty:
57     previous_winner_dem =
np.nan
58         else:
59         if
len(prev_election[prev_election['outcome'] == 'w'])-1 < seat:
60     previous_winner_dem =
np.nan
61         else:
62     party_of_prev_winner
prev_election[prev_election['outcome'] == 'w'].iloc[seat]['party']
63     if party_of_prev_winner i
democratic:
64     previous_winner_dem =
True
65         elif party_of_prev_winner
in republican:
66     previous_winner_dem =
False
67         else:
68     previous_winner_dem = np.nan
69         elif
prev_election[(prev_election['party'].isin(democratic+republican))].e
pty:
70     previous_winner_dem = np.nan
71         else:
72         if
len(prev_election[prev_election['outcome'] == 'w'])-1 < seat:
73     previous_winner_dem = np.nan
74         else:
75     party_of_prev_winner =
prev_election[prev_election['outcome'] == 'w'].iloc[seat]['party']
76     if party_of_prev_winner in
democratic:
77     previous_winner_dem = Tru
78     elif party_of_prev_winner in
republican:
79     previous_winner_dem =
False
80         else:
81

```

```

            if
len(prev_election[prev_election['outcome'] == '1'])-1 < seat:
82         previous_winner_dem =
np.nan
83         else:
84         if
prev_election[prev_election['outcome'] == '1'].iloc[seat]['party'] in
democratic:
85         previous_winner_dem = False
86         else:
87         previous_winner_dem = np.nan
88         else:
89         previous_winner_dem = np.nan
90         # keep both incumbency info and previous winner
info as they're two related but distinct pieces of info
91         final_election_data =
final_election_data.append({'year': year,
92                             'state': state,
93                             'temp_district': district
94                             'district': district + '-'
+ str(seat),
95                             'dem_vote': dem_vote,
96                             'previous_winner_dem':
previous_winner_dem,
97                             'dem_incumbent':
dem_incumbent,
98                             'rep_incumbent':

```

Out[26]:

```

dem_incumbent dem_vote district previous_winner_dem rep_incumbent state temp_dist
99 0 0.0 0.459653 I006-1 0.0 1.0 az True
100 1 0.0 0.476306 I008-1 0.0 1.0 az
101 2 0.0 0.532791 I010-1 NaN 0.0 az
3 0.0 0.389172 I014-1 0.0 1.0 az
4 0.0 0.005714 I016-1 NaN 0.0 az
...
9686 0.0 0.528484 u009-0 0.0 1.0 hi u
9687 1.0 0.813933 u011-0 1.0 0.0 hi u
9688 0.0 0.738679 u013-0 1.0 0.0 hi u
9689 1.0 0.616010 u019-0 1.0 0.0 hi u
9690 1.0 0.673618 u025-0 1.0 0.0 hi u

```

9691 rows × 8 columns



In [27]:

```
1 # check to make sure no rows contain no vote data
```

Out[27]:

```
dem_incumbent dem_vote district previous_winner_dem rep_incumbent state temp_district
```

```
In [28]: 1 # cases where there was no major party representative winner in the
          prior run for the seat
```

```
Out[28]:
```

	dem_incumbent	dem_vote	district	previous_winner_dem	rep_incumbent	state	temp_dist
2	0.0	0.532791	l010-1	NaN	0.0	az	l
4	0.0	0.005714	l016-1	NaN	0.0	az	l
6	0.0	0.480042	l020-1	NaN	0.0	az	l
8	0.0	0.583422	l026-1	NaN	0.0	az	l
10	0.0	0.287712	u005-0	NaN	0.0	az	u
...	...	...	...	...	...	...	...
9483	0.0	0.428948	u6-1-0	NaN	0.0	wv	u
9484	0.0	0.588194	u7-1-0	NaN	0.0	wv	u
9486	0.0	0.481926	u9-1-0	NaN	0.0	wv	u
9487	0.0	0.452720	u10-1-0	NaN	0.0	wv	u
9488	0.0	0.488941	u11-1-0	NaN	0.0	wv	u

1300 rows × 8 columns

It looks like some previous\_winner\_dem values are NaN. My intention is to exclude these when the prior winner isn't labeled as either a Democrat or Republican or the seat has changed due to redistricting. Let's spot check these to make sure that the only reason these are Nans.

```
In [29]: 1 election_data_merged[(election_data_merged['year'] < 2011) &
          (election_data_merged['sab'] == 'ms') &
```

```
Out[29]:
```

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
52460	2003	ms	40	0	g	barnett, ron	democratic	none	1164.0	l	l040
52461	2003	ms	40	0	g	mayhall, w. t. (ted)	republican	none	3173.0	w	l040
52640	2007	ms	40	0	g	mayhall, w. t. (ted)	NaN	inc	1847.0	w	l040

```
In [30]: 1 election_data_merged[(election_data_merged['year'] < 2011) &
          (election_data_merged['sab'] == 'va') &
```

```
Out[30]:
```

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
97888	2001	va	59	0	g	hagenau, h. p.	democratic	none	6829.0	l	l059

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
97889	2001	va	59	0	g	abbitt, watkins m. jr.	NaN	inc	11782.0	w	l059
98045	2003	va	59	0	g	hale, a. m.	democratic	none	5786.0	l	l059
98046	2003	va	59	0	g	abbitt, watkins m. jr.	NaN	inc	11834.0	w	l059
98196	2005	va	59	0	g	abbitt, watkins m. jr.	NaN	inc	16398.0	w	l059
98348	2007	va	59	0	g	brennan, connie	democratic	none	9136.0	l	l059
98349	2007	va	59	0	g	abbitt, watkins m. jr.	NaN	inc	13874.0	w	l059
98575	2009	va	59	0	g	abbitt, watkins m.	NaN	inc	16896.0	w	l059

```
In [31]: 1 # drop rows where prior winner couldn't be determined or the election
          2 was before 2012
          3 final_election_data = final_election_data.dropna(subset=
            ['previous_winner_dem'])
          4 final_election_data = final_election_data[final_election_data['year']
            >= 2012]
```

## Combining with census data

In order to match election data with census data by legislative district some examples need to be converted in specific ways. To accomplish this I'll be making a census matching column

```
In [32]: 1 final_election_data['census_matching_temp_district'] =
          2 final_election_data.apply(lambda row: 'l' +
            '{:0>3}'.format(row['temp_district'][1:-2]) if row['temp_district'][0]
            == 'l' and row['state'] == 'id' else
          3 (row['temp_district']
            [0]+'{:0>2}'.format(row['temp_district'][1:2]) +
            row['temp_district'][-1] if row['temp_district'][0] == 'l' and
            row['state'] == 'mn' else
          4 (row['temp_district'][:-2] +
            row['temp_district'][-1] if row['temp_district'][-1] in ['A', 'B'] and
            row['state'] == 'sd' else
          5 (row['temp_district'][0] +
            '{:0>3}'.format(row['temp_district'][1]) if row['temp_district'][0] ==
            'l' and row['state'] == 'wa' else
          6 (row['temp_district']
            [0]+'{:0>3}'.format(row['temp_district'][1:-2]) if
            row['temp_district'][-2] == '-' and row['state'] == 'wv' else
          7 ('l' + '{:0>2}'.format(row['temp_district']
            [1:-2]) + row['temp_district'][-1] if row['temp_district'][-1] in
            ['A', 'B', 'C'] and row['state'] == 'md' else
```

```

        (row['temp_district'].replace('lGrandisle',
        'lGrand-Isle') if row['temp_district'][:10] == 'lGrandisle' and
        row['state'] == 'vt' else
8         (row['temp_district'].replace(' ', '-') if
        row['temp_district'][0] == 'l' and row['state'] == 'vt' else
        row['temp_district'].replace(' ', '-')))

```

In [33]:

```

1  # American Community Survey coded data to grab
2  values_to_grab = {'B01001001': "Total Pop",
3                    'B01001002': "Male",
4                    'B01001026': "Female",
5                    'B02001002': "White alone",
6                    'B02001003': "Black or African American alone",
7                    'B02001004': "American Indian and Alaska Native
8                    alone",
9                    'B02001005': "Asian alone",
10                   'B02001006': "Native Hawaiian and Other Pacific
11                   Islander alone",
12                   'B02001007': "Some other race alone",
13                   'B02001008': "Two or more races:",
14                   'B02001009': "Two races including Some other race",
15                   'B02001010': "Two races excluding Some other race,
16                   and three or more races",
17                   'B05001002': "U.S. citizen, born in the United
18                   States",
19                   'B05001003': "U.S. citizen, born in Puerto Rico or
20                   U.S. Island Areas",
21                   'B06007002': "Speak only English",
22                   'B06007003': "Speak Spanish",
23                   'B06007006': "Speak other languages",
24                   'B06008002': "Never married",
25                   'B06008003': "Now married, except separated",
26                   'B06008004': "Divorced",
27                   'B06008005': "Separated",
28                   'B06008006': "Widowed",
29                   'B06009002': "Less than high school graduate",
30                   'B06009003': "High school graduate (includes
31                   equivalency)",
32                   'B06009004': "Some college or associate's degree",
33                   'B06009005': "Bachelor's degree",
34                   'B06009006': "Graduate or professional degree",
35                   'B06010002': "No income",
36                   'B06010004': "$1 to $9,999 or loss",
37                   'B06010005': "$10,000 to $14,999",
38                   'B06010006': "$15,000 to $24,999",
39                   'B06010007': "$25,000 to $34,999",
40                   'B06010008': "$35,000 to $49,999",
41                   'B06010009': "$50,000 to $64,999",
42                   'B06010010': "$65,000 to $74,999",
43                   'B06010011': "$75,000 or more",
44                   'B16004002': "5 to 17 years",
45                   'B16004024': "18 to 64 years",
46                   'B16004046': "65 years and over"

```

In [34]:

```

1  # for converting Massachusetts lower house names
2  ordinal_dict = {

```

```
3      '1st': 'First',
4      '2nd': 'Second',
5      '3rd': 'Third',
6      '4th': 'Fourth',
7      '5th': 'Fifth',
8      '6th': 'Sixth',
9      '7th': 'Seventh',
10     '8th': 'Eighth',
11     '9th': 'Ninth',
12     '10th': 'Tenth',
13     '11th': 'Eleventh',
14     '12th': 'Twelfth',
15     '13th': 'Thirteenth',
16     '14th': 'Fourteenth',
17     '15th': 'Fifteenth',
18     '16th': 'Sixteenth',
19     '17th': 'Seventeenth',
20     '18th': 'Eighteenth',
21     '19th': 'Nineteenth',
22     '20th': 'Twentieth',
23     '21st': 'Twenty-First',
24     '22nd': 'Twenty-Second',
25     '23rd': 'Twenty-Third',
26     '24th': 'Twenty-Fourth',
27     '25th': 'Twenty-Fifth',
28     '26th': 'Twenty-Sixth',
29     '27th': 'Twenty-Seventh',
30     '28th': 'Twenty-Eighth',
31     '29th': 'Twenty-Ninth',
32     '30th': 'Thirtieth',
33     '31st': 'Thirty-First',
34     '32nd': 'Thirty-Second',
35     '33rd': 'Thirty-Third',
36     '34th': 'Thirty-Fourth',
37     '35th': 'Thirty-Fifth',
38     '36th': 'Thirty-Sixth',
39     '37th': 'Thirty-Seventh'
```

In [35]:

```
1  # states and their census bureau numbers
2
3  states = {
4  '28': ["Mississippi", "ms"],
5  '29': ["Missouri", "mo"],
6  '30': ["Montana", "mt"],
7  '31': ["Nebraska", "ne"],
8  '32': ["Nevada", "nv"],
9  '33': ["New Hampshire", "nh"],
10 '34': ["New Jersey", "nj"],
11 '35': ["New Mexico", "nm"],
12 '36': ["New York", "ny"],
13 '37': ["North Carolina", "nc"],
14 '38': ["North Dakota", "nd"],
15 '39': ["Ohio", "oh"],
16 '40': ["Oklahoma", "ok"],
17 '41': ["Oregon", "or"],
18 '42': ["Pennsylvania", "pa"],
```

```

19 '44': ["Rhode Island", "ri"],
20 '45': ["South Carolina", "sc"],
21 '46': ["South Dakota", "sd"],
22 '47': ["Tennessee", "tn"],
23 '48': ["Texas", "tx"],
24 '50': ["Vermont", "vt"],
25 '49': ["Utah", "ut"],
26 '51': ["Virginia", "va"],
27 '53': ["Washington", "wa"],
28 '54': ["West Virginia", "wv"],
29 '55': ["Wisconsin", "wi"],
30 '56': ["Wyoming", "wy"],
31 '01': ["Alabama", "al"],
32 '02': ["Alaska", "ak"],
33 '04': ["Arizona", "az"],
34 '05': ["Arkansas", "ar"],
35 '06': ["California", "ca"],
36 '08': ["Colorado", "co"],
37 '10': ["Delaware", "de"],
38 '13': ["Georgia", "ga"],
39 '09': ["Connecticut", "ct"],
40 '12': ["Florida", "fl"],
41 '16': ["Idaho", "id"],
42 '15': ["Hawaii", "hi"],
43 '17': ["Illinois", "il"],
44 '18': ["Indiana", "in"],
45 '19': ["Iowa", "ia"],
46 '20': ["Kansas", "ks"],
47 '21': ["Kentucky", "ky"],
48 '22': ["Louisiana", "la"],
49 '23': ["Maine", "me"],
50 '24': ["Maryland", "md"],
51 '25': ["Massachusetts", "ma"],
52 '26': ["Michigan", "mi"],
53 '27': ["Minnesota", "mn"]

```

In [36]:

```

1 from drive.MyDrive.capstone.api_keys import census
2
3 census_data = pd.DataFrame()
4 for year in [int(str(x)[:4]) for x in
5 final_election_data['year'].value_counts().index]:
6     for state in states.keys():
7         demo_data_senate = censusdata.download('acs5', year,
8 censusdata.censusgeo([('state', state),
9 ('state%20legislative%20district%20(upper%20chamber)', '*')]),
10 [each[:6] + '_' + each[6:] + 'E' for each in values_to_grab.keys()],
11 key=census)
12         # convert columns from acs codes
13         demo_data_senate =
14 demo_data_senate.rename(columns=dict(zip([each[:6] + '_' + each[6:] +
15 'E' for each in values_to_grab.keys()], values_to_grab.values())))
16         demo_data_senate = demo_data_senate.reset_index(drop=False)
17         # Massachusetts specific rules
18         if state == '25':

```

```
demo_data_senate['index'] =
demo_data_senate['index'].apply(lambda x: 'u' +
x.name[:x.name.find('District')-1].replace(',',''))
13     # Vermont specific Rules
14     elif state == '50':
15         demo_data_senate['index'] =
demo_data_senate['index'].apply(lambda x: 'u' + x.name.split()[0])
16     else:
17         demo_data_senate['index'] =
demo_data_senate['index'].apply(lambda x: 'u'+x.geo[1][1].replace(',',''))
18     demo_data_senate = demo_data_senate.rename(columns={'index':
'census_matching_temp_district'})
19     demo_data_senate['year'] = year
20     demo_data_senate['state'] = states[state][1]
21     # Nebraska has no lower house
22     if state != '31':
23         demo_data_assembly = censusdata.download('acs5', year,
censusdata.censusgeo([('state', state),
('state%20legislative%20district%20(lower%20chamber)', '*')]),
[each[:6] + '_' + each[6:] + 'E' for each in values_to_grab.keys()],
key=census)
24         demo_data_assembly =
demo_data_assembly.rename(columns=dict(zip([each[:6] + '_' + each[6:]
+ 'E' for each in values_to_grab.keys()], values_to_grab.values())))
25         demo_data_assembly =
demo_data_assembly.reset_index(drop=False)
26
27         # Massachusetts specific rules
28         if state == '25':
29             demo_data_assembly['index'] =
demo_data_assembly['index'].apply(lambda x: '1' + x.name.split('(')
[0].strip()[:-9].replace(x.name.split()[0],
ordinal_dict[x.name.split()[0]]) if x.name[0].isnumeric() else '1' +
x.name[:x.name.find('District')-1].replace(',',''))
30         # New Hampshire specific rules
31         elif state == '33':
32             demo_data_assembly['index'] =
demo_data_assembly['index'].apply(lambda x: '1' + '
'.join(x.name.split()[3:7:3]).replace(',',''))
33         # Vermont specific rules
34         elif state == '50':
35             demo_data_assembly['index'] =
demo_data_assembly['index'].apply(lambda x: '1' + x.name.split()
[0].replace(',',''))
36         else:
37             demo_data_assembly['index'] =
demo_data_assembly['index'].apply(lambda x: '1'+x.geo[1]
[1].replace(',',''))
38         demo_data_assembly = demo_data_assembly.rename(columns=
{'index': 'census_matching_temp_district'})
39         demo_data_assembly['year'] = year
40         demo_data_assembly['state'] = states[state][1]
41
42
43
```



```
        census_data = pd.concat([census_data, demo_data_assembly,
demo_data_senate])
44     else:
45         census_data = pd.concat([census_data, demo_data_senate])
46
47
```

Out[36]:

	census_matching_temp_district	Total Pop	Male	Female	White alone	Black or African American alone	American Indian and Alaska Native alone	Asian alone	H
0	l090	24519	12017	12502	16048	8229	3	29	
1	l091	24156	11596	12560	9804	14118	6	6	
2	l092	24381	12026	12355	19174	4806	9	237	
3	l093	25147	12755	12392	20440	3903	115	86	
4	l094	24227	12325	11902	7044	16098	33	55	
...	...	...	...	...	...	...	...	...	
62	u024	79732	40138	39594	72959	2718	312	703	
63	u025	79685	38793	40892	69584	3299	105	4235	
64	u026	80323	39463	40860	70061	3850	162	3735	
65	u027	79187	39239	39948	74425	1323	206	936	
66	u028	79064	38863	40201	76235	710	199	875	

33615 rows × 42 columns



```
In [37]: 1 # mapping New Hampshire's floterial districts to their component part
          2 in order to assemble entire districts
          3 nh_floterial_map = {
          4     'lBelknap' : [['8', '5', '7'],
          5                    ['9', '3', '6']],
          6
          7     'lCarroll' : [['7', '1', '2', '3'],
          8                    ['8', '4', '5']],
          9
          10     'lCheshire': [['14', '11', '9'],
          11                    ['15', '10', '12', '13'],
          12                    ['16', '4', '5', '6', '7', '8']],
          13
          14     'lCoos' : [['7', '2', '4', '5']],
          15
          16     'lGrafton' : [['14', '1', '2'],
          17                    ['15', '3', '4'],
          18                    ['16', '11', '6'],
          19                    ['17', '10', '9']],
          20
          21     'lHillsborough' : [['38', '1', '3', '4'],
          22                    ['39', '2', '6'],
          23                    ['40', '23', '27', '5'],
```

```

23         ['41', '22', '7'],
24         ['42', '10', '8', '9'],
25         ['43', '11', '12', '13', '14'],
26         ['44', '15', '16', '20'],
27         ['45', '17', '18', '19']],
28     'lMerrimack' : [['25', '1', '7'],
29                    ['26', '3', '8', '9'],
30                    ['27', '11', '12', '13', '14', '15', '16'],
31                    ['28', '17', '18', '19'],
32                    ['29', '21', '22']],
33     'lRockingham': [['30', '25', '26', '28', '29'],
34                    ['31', '22', '23', '27'],
35                    ['32', '1', '2'],
36                    ['33', '10', '11', '12'],
37                    ['34', '13', '14'],
38                    ['35', '15', '16'],
39                    ['36', '17', '18', '19'],
40                    ['37', '20', '21']],
41     'lStrafford': [['19', '13', '14'],
42                   ['20', '15', '16'],
43                   ['21', '17', '18'],
44                   ['22', '7', '8'],
45                   ['23', '10', '9'],
46                   ['24', '11', '12'],
47                   ['25', '4', '5']],
48     'lSullivan': [['9', '1', '2', '6'],
49                  ['10', '3', '4', '5'],
50                  ['11', '7', '8']]

```

```

In [38]: 1 def nh_floterial_district_census_combiner(nh_floterial_map,
          census_data):
          2     '''Takes census_data and the floterial district map, returns
          census_data rows for the floterial districts'''
          3     flot_dists = pd.DataFrame()
          4     for county, numbers in nh_floterial_map.items():
          5         for floterial in range(len(numbers)):
          6             new_district =
          census_data[(census_data['census_matching_temp_district'].isin((count
          + ' ' + c for c in numbers[floterial][1:])))].groupby('year').sum()
          7             new_district['census_matching_temp_district'] = county +
          ' + numbers[floterial][0]
          8             new_district['state'] = 'nh'
          9             new_district['made_from'] = str(numbers[floterial][1:])
          10            new_district = new_district.reset_index(drop=False)
          11            flot_dists = flot_dists.append(new_district,
          ignore_index=True)
          12
          13            return(flot_dists)
          14 flots = nh_floterial_district_census_combiner(nh_floterial_map,

```

```

In [39]: 1 # add floterial districts to census_data

```

In [40]:

Out[40]:

	dem_incumbent	dem_vote	district	previous_winner_dem	rep_incumbent	state	temp_dist
0	0.0	0.459653	I006-1	0.0	1.0	az	I
1	0.0	0.476306	I008-1	0.0	1.0	az	I
3	0.0	0.389172	I014-1	0.0	1.0	az	I
5	0.0	0.471326	I018-1	0.0	1.0	az	I
7	0.0	0.395057	I021-1	0.0	1.0	az	I
...	...	...	...	...	...	...	...
9686	0.0	0.528484	u009-0	0.0	1.0	hi	u
9687	1.0	0.813933	u011-0	1.0	0.0	hi	u
9688	0.0	0.738679	u013-0	1.0	0.0	hi	u
9689	1.0	0.616010	u019-0	1.0	0.0	hi	u
9690	1.0	0.673618	u025-0	1.0	0.0	hi	u

8391 rows × 9 columns

In [41]:

```
1 sample = pd.merge(left = final_election_data,  
2                     right = census_data_plus,  
3                     how = 'left',  
4                     on = ['census_matching_temp_district', 'year', 'state'])  
5
```

Out[41]:

	dem_incumbent	dem_vote	district	previous_winner_dem	rep_incumbent	state	temp_dist
0	0.0	0.459653	I006-1	0.0	1.0	az	I
1	0.0	0.476306	I008-1	0.0	1.0	az	I
2	0.0	0.389172	I014-1	0.0	1.0	az	I
3	0.0	0.471326	I018-1	0.0	1.0	az	I
4	0.0	0.395057	I021-1	0.0	1.0	az	I
...	...	...	...	...	...	...	...
8386	0.0	0.528484	u009-0	0.0	1.0	hi	u
8387	1.0	0.813933	u011-0	1.0	0.0	hi	u
8388	0.0	0.738679	u013-0	1.0	0.0	hi	u
8389	1.0	0.616010	u019-0	1.0	0.0	hi	u

```
dem incumbent dem vote district previous winner dem rep incumbent state temp dist
```

Now that we've combined our election data with census data, let's double-check to make sure every district we're interested in was properly mapped to census data.

In [42]:

Out [42]:

```
dem_incumbent dem_vote district previous_winner_dem rep_incumbent state temp_district
```



The generic ballot is an indicator that is commonly used by political observers to gauge the general partisan sway of the country as a whole. It's gauged in political polls by asking voters which party they'd like to see control congress. It's generally believed that a very Republican-friendly environment can help Republican candidates down the ballot and vice versa in a Democratic-friendly environment. For this reason I'm including the national generic ballot averaged over the handful of polls preceding each year's November election as reflected on [RealClearPolitics' website \(https://www.realclearpolitics.com/epolls/other/2022-generic-congressional-vote-7361.html\)](https://www.realclearpolitics.com/epolls/other/2022-generic-congressional-vote-7361.html).

```
In [43]: 1 generic_ballot = {2011: 1.8,
2             2012: -0.2,
3             2013: -0.2,
4             2014: -2.4,
5             2015: 0,
6             2016: 0.6,
7             2018: 7.3,
8             2020: 6.8}
9
10 sample['generic_ballot'] = sample.apply(lambda row:
```

Another factor thought to be a strong indicator of partisanship is population density. To understand this I've grabbed the shapefiles for every legislative district from 2018 - well after each map should be settled.

```
In [44]: 1 leg_dist_areas = pd.DataFrame()
2
3 for state in states.keys():
4     if state not in ['72', '11']:
5         #         senates
6
```

```

        shapefile = gpd.read_file("drive/MyDrive/capstone
/SLDU/tl_2018_" + state + "_sldu/tl_2018_" + state + "_sldu.shp")
7         # must change projection to cea for accurate area
8         shapefile = shapefile.to_crs({'proj': 'cea'})
9         shapefile["area"] = shapefile['geometry'].area/ 10**6
10        # Massachusetts naming
11        if state == '25':
12            shapefile['census_matching_temp_district'] =
shapefile['NAMELSAD'].apply(lambda x: 'u' + x[:x.find('District')-
1].replace(',',' '))
13            # Vermont naming
14            elif state == '50':
15                shapefile['census_matching_temp_district'] =
shapefile['NAMELSAD'].apply(lambda x: 'u' + x.split()[0])
16            else:
17                shapefile['census_matching_temp_district'] =
shapefile['SLDUST'].apply(lambda x: 'u'+x.replace(',',' '))
18
19
20
21
22        shapefile['state'] = shapefile['STATEFP'].apply(lambda x:
states[state][1])
23        leg_dist_areas = leg_dist_areas.append(shapefile[['area',
'census_matching_temp_district', 'state']], ignore_index=True)
24        # lower houses
25        # excluding Nebraska which has no lower house
26        if state != '31':
27            shapefile = gpd.read_file("drive/MyDrive/capstone
/SLDL/tl_2018_" + state + "_sldl/tl_2018_" + state + "_sldl.shp")
28            shapefile = shapefile.to_crs({'proj': 'cea'})
29            shapefile["area"] = shapefile['geometry'].area/ 10**6
30            # Massachusetts naming
31            if state == '25':
32                shapefile['census_matching_temp_district'] =
shapefile['NAMELSAD'].apply(lambda x: 'l' + x.split('(')[0].strip()
[:-9].replace(x.split()[0], ordinal_dict[x.split()[0]]) if
x[0].isnumeric() else 'l' + x[:x.find('District')-1].replace(',',' '))
33            # New Hampshire naming
34            elif state == '33':
35                shapefile['census_matching_temp_district'] =
shapefile['NAMELSAD'].apply(lambda x: 'l' + ' '.join(x.split()
[3:7:3]).replace(',',' '))
36            # Vermont naming
37            elif state == '50':
38                shapefile['census_matching_temp_district'] =
shapefile['NAMELSAD'].apply(lambda x: 'l' + x.split()[0].replace(',','
'))
39            else:
40                shapefile['census_matching_temp_district'] =
shapefile['SLDLST'].apply(lambda x: 'l'+x.replace(',',' '))
41                shapefile['state'] = shapefile['STATEFP'].apply(lambda x:
states[state][1])
42                leg_dist_areas = leg_dist_areas.append(shapefile[['area',
'census_matching_temp_district', 'state']], ignore_index=True)
43

```

	area	census_matching_temp_district	state
0	495.375785	u049	ms
1	1930.382029	u017	ms
2	947.814366	u006	ms
3	912.550467	u020	ms
4	3303.149472	u031	ms
...	...	...	...
6731	258.593410	l54B	mn
6732	46.139066	l53A	mn
6733	4587.942110	l06B	mn
6734	1047.503175	l03B	mn
6735	19.052332	l64A	mn

6736 rows × 3 columns

```
In [45]: 1 # similar to the census function, I needed to create a function to
2         combine districts to account for NH's floterial districts
3         def nh_floterial_district_area_combiner(nh_floterial_map,
4         leg_dist_areas):
5             '''Takes district areas and the floterial district map, returns
6             census_data rows for the floterial districts'''
7
8             flot_dists = pd.DataFrame()
9             for county, numbers in nh_floterial_map.items():
10                 for floterial in range(len(numbers)):
11                     new_district =
12                     leg_dist_areas[(leg_dist_areas['census_matching_temp_district'].isin(
13                     county + ' ' + c for c in numbers[floterial][1:]))]
14                     new_district['census_matching_temp_district'] = county +
15                     ' ' + numbers[floterial][0]
16                     new_district['state'] = 'nh'
17                     new_district['made_from'] = str(numbers[floterial][1:])
18                     new_district = new_district.reset_index(drop=True)
19                     flot_dists = flot_dists.append(new_district,
20                     ignore_index=True)
21
22             return flot_dists
23         flots = nh_floterial_district_area_combiner(nh_floterial_map,
```

In [46]:

Out[46]:

	area	census_matching_temp_district	state	made_from
0	369.821472	lBelknap 8	nh	['5', '7']
1	116.391494	lBelknap 8	nh	['5', '7']
2	82.575713	lBelknap 9	nh	['3', '6']
3	67.631465	lBelknap 9	nh	['3', '6']

	area	census_matching_temp_district	state	made_from
4	406.823180	ICarroll 7	nh	['1', '2', '3']
...	...	...	...	...
107	19.273284	ISullivan 10	nh	['3', '4', '5']
108	36.040656	ISullivan 10	nh	['3', '4', '5']
109	58.775621	ISullivan 10	nh	['3', '4', '5']
110	410.128183	ISullivan 11	nh	['7', '8']
111	98.522632	ISullivan 11	nh	['7', '8']

112 rows × 4 columns

```
In [47]: 1 leg_dist_areas_plus =  
        leg_dist_areas.append(flots.drop(columns='made_from'))
```

```
Out[47]:
```

	area	census_matching_temp_district	state
0	495.375785	u049	ms
1	1930.382029	u017	ms
2	947.814366	u006	ms
3	912.550467	u020	ms
4	3303.149472	u031	ms
...	...	...	...
107	19.273284	ISullivan 10	nh
108	36.040656	ISullivan 10	nh
109	58.775621	ISullivan 10	nh
110	410.128183	ISullivan 11	nh
111	98.522632	ISullivan 11	nh

6848 rows × 3 columns



```
In [48]: 1 sample = pd.merge(left=sample,  
2         right=leg_dist_areas_plus,  
3         how='left',
```

```
In [49]: 1 # confirm that every row has an area
```

```
Out[49]: array([], dtype=object)
```

```
In [50]: 1 # derive population density in population per kilometer
          2
          3 sample['pop_density_km'] = sample['Total Pop']/sample['area']
```

Out[50]:

	dem_incumbent	dem_vote	district	previous_winner_dem	rep_incumbent	state	temp_dist
0	0.0	0.459653	l006-1	0.0	1.0	az	l
1	0.0	0.476306	l008-1	0.0	1.0	az	l
2	0.0	0.389172	l014-1	0.0	1.0	az	l
3	0.0	0.471326	l018-1	0.0	1.0	az	l
4	0.0	0.395057	l021-1	0.0	1.0	az	l
...	...	...	...	...	...	...	...
8539	0.0	0.528484	u009-0	0.0	1.0	hi	u
8540	1.0	0.813933	u011-0	1.0	0.0	hi	u
8541	0.0	0.738679	u013-0	1.0	0.0	hi	u
8542	1.0	0.616010	u019-0	1.0	0.0	hi	u
8543	1.0	0.673618	u025-0	1.0	0.0	hi	u

8544 rows × 51 columns



I'd like to convert my population counting rows into proportions of the total population. Otherwise my results will likely be skewed by the fact that senate seats will always be bigger than house seats within a given state.

```
In [51]: 1 for each in list(values_to_grab.values())[1:]:
          2     sample[each] = sample[each]/sample['Total Pop']
          3
```

Out[51]:

	dem_incumbent	dem_vote	district	previous_winner_dem	rep_incumbent	state	temp_dist
0	0.0	0.459653	l006-1	0.0	1.0	az	l
1	0.0	0.476306	l008-1	0.0	1.0	az	l



	dem_incumbent	dem_vote	district	previous_winner_dem	rep_incumbent	state	temp_dist
2	0.0	0.389172	I014-1	0.0	1.0	az	I
3	0.0	0.471326	I018-1	0.0	1.0	az	I
4	0.0	0.395057	I021-1	0.0	1.0	az	I
...	...	...	...	...	...	...	
8539	0.0	0.528484	u009-0	0.0	1.0	hi	u

In [52]: 1 *# as Total Pop and area are completely implied by other features, I'm dropping them*

In [53]: 1

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 8544 entries, 0 to 8543
Data columns (total 49 columns):
 #   Column                                     Non
-Null Count  Dtype
---  -
0    dem_incumbent                             854
4 non-null   float64
1    dem_vote                                   854
4 non-null   float64
2    district                                  854
4 non-null   object
3    previous_winner_dem                       854
4 non-null   float64
4    rep_incumbent                             854
4 non-null   float64
5    state                                     854
4 non-null   object
6    temp_district                             854
4 non-null   object
7    year                                      854
4 non-null   float64
8    census_matching_temp_district             854
4 non-null   object
9    Male                                      854
4 non-null   float64
10   Female                                   854
4 non-null   float64
11   White alone                             854
4 non-null   float64
12   Black or African American alone          854
4 non-null   float64
13   American Indian and Alaska Native alone  854
4 non-null   float64
14   Asian alone                             854
4 non-null   float64
15   Native Hawaiian and Other Pacific Islander alone  854
4 non-null   float64
16   Some other race alone                    854
4 non-null   float64
17   Two or more races:                       854
4 non-null   float64
18   Two races including Some other race      854
4 non-null   float64
19   Two races excluding Some other race, and three or more races  854
4 non-null   float64
20   U.S. citizen, born in the United States  854
4 non-null   float64
21   U.S. citizen, born in Puerto Rico or U.S. Island Areas  854
4 non-null   float64
22   Speak only English                       854
4 non-null   float64
23   Speak Spanish                            854
4 non-null   float64
24   Speak other languages                     854

```

```

4 non-null    float64
25  Never married                                854
4 non-null    float64
26  Now married, except separated                854
4 non-null    float64
27  Divorced                                    854
4 non-null    float64
28  Separated                                  854
4 non-null    float64
29  Widowed                                    854
4 non-null    float64
30  Less than high school graduate              854
4 non-null    float64
31  High school graduate (includes equivalency) 854
4 non-null    float64
32  Some college or associate's degree          854
4 non-null    float64
33  Bachelor's degree                          854

```

In [54]:

Out[54]:

	dem_incumbent	dem_vote	previous_winner_dem	rep_incumbent	year	I
<b>count</b>	8544.000000	8544.000000	8544.000000	8544.000000	8544.000000	8544.000000
<b>mean</b>	0.372074	0.485887	0.445342	0.454588	2014.216175	0.490000
<b>std</b>	0.483386	0.157199	0.497033	0.497963	1.580333	0.010000
<b>min</b>	0.000000	0.000000	0.000000	0.000000	2012.000000	0.410000
<b>25%</b>	0.000000	0.371840	0.000000	0.000000	2012.000000	0.480000
<b>50%</b>	0.000000	0.466854	0.000000	0.000000	2014.000000	0.490000
<b>75%</b>	1.000000	0.586498	1.000000	1.000000	2016.000000	0.500000
<b>max</b>	1.000000	0.999818	1.000000	1.000000	2016.000000	0.590000

Ultimately what I hope to be able to predict is which seats are likely to be tossups, i.e. could go either Republican or Democratic. I've defined that here as within 5 percentage points of 50%.

In [55]:

```
1 sample['dem_vote'] = sample['dem_vote'].apply(lambda x: -1 if x<= .45
```

In [56]:

```
1
```

## Holdout Data

## Generating some holdout data from after 2016

```
In [57]: 1 def process_ny_boe_data(document_name):
2         output = pd.DataFrame()
3
4         for district in pd.read_excel('drive/MyDrive/capstone/ny_data/' +
document_name, sheet_name=None).keys():
5             sheet = pd.read_excel('drive/MyDrive/capstone/ny_data/' +
document_name, sheet_name=district, header=1)
6             if 'Candidate Name (Party)' in sheet.columns:
7                 sheet = sheet.dropna()
8                 sheet['party'] = sheet['Candidate Name
(Party)'].apply(lambda x: 'democratic' if x[-4:-1] == 'DEM' else
('republican' if x[-4:-1] == 'REP' else None))
9                 sheet = sheet[['Candidate Name (Party)', 'Total Votes by
Candidate', 'party']]
10                dnumber, dtype = district.split()
11                dnumber = re.findall(r'\d+', dnumber)[0]
12                if dtype == 'AD':
13                    district = 'l' + '{:0>3}'.format(dnumber)
14                elif dtype == 'SD':
15                    district = 'u' + '{:0>3}'.format(dnumber)
16                sheet['district'] = pd.Series([district, district])
17                year = int(document_name[:4])
18                sheet['year'] = pd.Series([year, year])
19                output = pd.concat([output, sheet])
20            else:
21                pass
22
23
24        output = output.rename(columns={'Candidate Name (Party)': 'cand',
'Total Votes by Candidate': "vote"})
```

```
In [58]: 1 # importing election data downloaded from NY's Board of Elections
2 holdout_ny = pd.DataFrame()
3
4 # because of covid-imposed limitations on census data, 2020 had to be
excluded
5 for file in ['2018Assembly.xlsx', '2018NYSenate.xlsx']:
6     holdout_ny = holdout_ny.append(process_ny_boe_data(file),
ignore_index=True)
7
8 holdout_ny['sab'] = 'ny'
```

```
Out[58]:
```

	cand	vote	party	district	year	sab
0	Fred W. Thiele, Jr. (DEM)	31961.0	democratic	l001	2018.0	ny
1	Patrick M. O'Connor (REP)	19953.0	republican	l001	2018.0	ny
2	Rona Smith (DEM)	21533.0	democratic	l002	2018.0	ny
3	Anthony H. Palumbo (REP)	31242.0	republican	l002	2018.0	ny
4	Clyde E. Parker (DEM)	17822.0	democratic	l003	2018.0	ny

	cand	vote	party	district	year	sab
...	...	...	...	...	...	...
396	Joan Elizabeth Seamans (DEM)	51471.0	democratic	u061	2018.0	ny
397	Michael H. Ranzenhofer (REP)	60780.0	republican	u061	2018.0	ny
398	Robert G. Ort (REP)	69118.0	republican	u062	2018.0	ny
399	Peter A. Diachun (GRE)	10539.0	None	NaN	NaN	ny
400	Timothy M. Kennedy (DEM)	70221.0	democratic	u063	2018.0	ny

401 rows x 6 columns

```
In [59]: 1 # grabbing data for the 2016 election in order to generate
          2 previous_win_dem column
          3 ny_2016 = election_data_merged[(election_data_merged['sab'] == 'ny')
          4 (election_data_merged['year'] == 2016)]
```

```
Out[59]:
```

	year	sab	ddez	sen	etype	cand	party	exper	vote	outcome	district
5777	2016	ny	1	0	g	thiele, fred w. jr.	democratic	inc	35246.0	w	l001
5778	2016	ny	2	0	g	palumbo, anthony h.	republican	inc	39795.0	w	l002
5779	2016	ny	3	0	g	murray, dean	republican	inc	29087.0	w	l003
5780	2016	ny	4	0	g	englebright, steven	democratic	inc	31941.0	w	l004
5781	2016	ny	4	0	g	weissbard, steven	republican	none	21994.0	l	l004
...	...	...	...	...	...	...	...	...	...	...	...
71336	2016	ny	61	1	g	scattering	NaN	none	15.0	l	u061
71337	2016	ny	62	1	g	scattering	NaN	none	18.0	l	u062
71338	2016	ny	62	1	g	scattering	NaN	none	56.0	l	u062
71339	2016	ny	62	1	g	scattering	NaN	none	295.0	l	u062
71340	2016	ny	63	1	g	scattering	NaN	none	0.0	l	u063

750 rows x 11 columns

```
In [60]: 1 # largely a duplicate of the code that generated final_election_data
          2 but specific to the needs of this dataset
          3 years = [2018]
          4 states = ['ny']
          5 holdout_data = pd.DataFrame()
          6
          7 for year in years:
          8     for state in states:
          9         districts = holdout_data[(holdout_data['sab'] == state) &
          10 (holdout_data['year'] == year)][['district']].unique()
```

```

9         for district in districts:
10             election = holdout_ny[(holdout_ny['year'] == year) &
(holdout_ny['sab'] == state) & (holdout_ny['district'] == district) &
(holdout_ny['party'].isin(democratic+republican))]
11             if len(election[(election['party'].isin(republican)) |
(election['party'].isin(democratic))]) == 1:
12                 pass
13             else:
14                 dem_net_vote =
election[election['party'].isin(democratic)].iloc[0]['vote'] /
(election[election['party'].isin(democratic)].iloc[0]['vote'] +
election[election['party'].isin(republican)].iloc[0]['vote'])
15                 election = election.sort_values('vote',
ascending=False).reset_index(drop=True)
16                 election['outcome'] = 'l'
17                 election.at[0, 'outcome'] = 'w'
18                 prev_election = ny_2016[(ny_2016['year'] == year - 2)
& (ny_2016['sab'] == state) & (ny_2016['district'] == district)]
19
20                 party_of_prev_winner =
prev_election[prev_election['outcome'] == 'w'].iloc[0]['party']
21                 if party_of_prev_winner in democratic:
22                     previous_winner_dem = True
23                 elif party_of_prev_winner in republican:
24                     previous_winner_dem = False
25                 else:
26                     if len(prev_election[prev_election['outcome'] ==
'1'])-1 < 0:
27                         previous_winner_dem = np.nan
28                     else:
29                         if prev_election[prev_election['outcome'] ==
'1'].iloc[0]['party'] in democratic:
30                             previous_winner_dem = False
31                         else:
32                             previous_winner_dem = np.nan
33
34                 holdout_data = holdout_data.append({'year': year,
35                                                     'state': state,
36                                                     'temp_district': district,
37                                                     'district': district + '-' +
str(0),
38                                                     'dem_net_vote': dem_net_vote,
39                                                     'previous_winner_dem':
previous_winner_dem,
40                                                     'dem_incumbent':
dem_incumbent,
41                                                     'rep_incumbent': rep_incumbent
42                                                     }, ignore_index=True)
43

```

Out[60]:

	dem_incumbent	dem_net_vote	district	previous_winner_dem	rep_incumbent	state	temp_c
0	1.0	0.615653	l001-0	1.0	0.0	ny	
1	1.0	0.408015	l002-0	0.0	0.0	ny	
2	1.0	0.460065	l003-0	0.0	0.0	ny	

	dem_incumbent	dem_net_vote	district	previous_winner_dem	rep_incumbent	state	temp_d
3	1.0	0.606159	l004-0	1.0	0.0	ny	
4	1.0	0.413354	l005-0	0.0	0.0	ny	
...	...	...	...	...	...	...	
125	1.0	0.480922	u055-0	0.0	0.0	ny	
126	1.0	0.444585	u056-0	0.0	0.0	ny	
127	1.0	0.405021	u058-0	0.0	0.0	ny	
128	1.0	0.442403	u060-0	0.0	0.0	ny	
129	1.0	0.458535	u061-0	0.0	0.0	ny	

```
In [61]: 1 census_data = pd.DataFrame()
2 # acs data doesn't exist 2020 due to COVID
3 for year in [2018]:
4     # NY is state 36 in the Census codebook
5     for state in ['36']:
6         demo_data_senate = censusdata.download('acs5', year,
7 censusdata.censusgeo([('state', state),
8 ('state%20legislative%20district%20(upper%20chamber)', '*')]),
9 [each[:6] + '_' + each[6:] + 'E' for each in values_to_grab.keys()],
10 key=census)
11 demo_data_senate =
12 demo_data_senate.rename(columns=dict(zip([each[:6] + '_' + each[6:] +
13 'E' for each in values_to_grab.keys()], values_to_grab.values()))
14 demo_data_senate = demo_data_senate.reset_index(drop=False)
15 demo_data_senate['index'] =
16 demo_data_senate['index'].apply(lambda x: 'u'+x.geo[1][1].replace(', '
17 ''))
18 demo_data_senate = demo_data_senate.rename(columns={'index':
19 'temp_district'})
20 demo_data_senate['year'] = year
21 demo_data_senate['state'] = 'ny'
22
23 demo_data_assembly = censusdata.download('acs5', year,
24 censusdata.censusgeo([('state', state),
25 ('state%20legislative%20district%20(lower%20chamber)', '*')]),
26 [each[:6] + '_' + each[6:] + 'E' for each in values_to_grab.keys()],
27 key=census)
28 demo_data_assembly =
29 demo_data_assembly.rename(columns=dict(zip([each[:6] + '_' + each[6:]
30 + 'E' for each in values_to_grab.keys()], values_to_grab.values()))
31 demo_data_assembly =
32 demo_data_assembly.reset_index(drop=False)
33 demo_data_assembly['index'] =
34 demo_data_assembly['index'].apply(lambda x: 'l'+x.geo[1]
35 [1].replace(', ', ''))
36 demo_data_assembly = demo_data_assembly.rename(columns=
37 {'index': 'temp_district'})
38 demo_data_assembly['year'] = year
39 demo_data_assembly['state'] = 'ny'
```

```
        census_data = pd.concat([census_data, demo_data_assembly,
demo_data_senate])
23
24
25
```

Out[61]:

	temp_district	Total Pop	Male	Female	White alone	Black or African American alone	American Indian and Alaska Native alone	Asian alone	Native Hawaiian and Other Pacific Islander alone	Some other race alone
0	l032	136376	63997	72379	6779	82696	443	15281	86	26486
1	l061	120894	59189	61705	64496	32340	371	9255	12	9198
2	l140	128611	61121	67490	108673	7542	463	3629	18	3740
3	l018	131400	63047	68353	30913	60487	1223	2330	0	26569
4	l029	129461	61024	68437	6302	87277	674	17719	9	13623
...	...	...	...	...	...	...	...	...	...	...
58	u042	293926	148790	145136	229571	26549	1184	6375	79	17948
59	u043	296988	146472	150516	275265	7422	457	5732	25	2168
60	u044	293263	140757	152506	210092	44241	532	19138	133	5027
61	u045	298350	154072	144278	272849	10505	4315	2591	160	3585
62	u046	291782	145530	146252	257200	13296	532	6521	17	6433

213 rows × 42 columns



```
In [62]: 1 holdout_data = pd.merge(left = holdout_data,
2         right = census_data,
3         how = 'left',
4         on = ['temp_district', 'year', 'state'])
5
```

Out[62]:

	dem_incumbent	dem_net_vote	district	previous_winner_dem	rep_incumbent	state	temp_c
0	1.0	0.615653	l001-0	1.0	0.0	ny	
1	1.0	0.408015	l002-0	0.0	0.0	ny	
2	1.0	0.460065	l003-0	0.0	0.0	ny	
3	1.0	0.606159	l004-0	1.0	0.0	ny	



	dem_incumbent	dem_net_vote	district	previous_winner_dem	rep_incumbent	state	temp_c
4	1.0	0.413354	l005-0	0.0	0.0	ny	
...	...	...	...	...	...	...	...
125	1.0	0.480922	u055-0	0.0	0.0	ny	

In [63]:

```
1 generic_ballot = {2011: 1.8,  
2                   2012: -0.2,  
3                   2013: -0.2,  
4                   2014: -2.4,  
5                   2015: 0,  
6                   2016: 0.6,  
7                   2018: 7.3,  
8                   2020: 6.8}  
9 holdout_data['generic_ballot'] = holdout_data.apply(lambda row:
```

```

In [64]: 1 leg_dist_areas = pd.DataFrame()
          2
          3 for state in ['36']:
          4     # NY senate
          5     shapefile = gpd.read_file("drive/MyDrive/capstone/SLDU/tl_2018_"
state + "_sldu/tl_2018_" + state + "_sldu.shp")
          6     # change projection to cea for accurate area
          7     shapefile = shapefile.to_crs({'proj': 'cea'})
          8     shapefile["area"] = shapefile['geometry'].area / 10**6
          9     shapefile['temp_district'] = shapefile['SLDUST'].apply(lambda x:
'u'+x.replace(', ', ''))
         10
         11
         12
         13
         14     shapefile['state'] = "ny"
         15     leg_dist_areas = leg_dist_areas.append(shapefile[['area',
'temp_district', 'state']], ignore_index=True)
         16
         17     # NY assembly
         18     shapefile = gpd.read_file("drive/MyDrive/capstone/SLDL/tl_2018_"
state + "_sldl/tl_2018_" + state + "_sldl.shp")
         19     shapefile = shapefile.to_crs({'proj': 'cea'})
         20     shapefile["area"] = shapefile['geometry'].area / 10**6
         21     shapefile['temp_district'] = shapefile['SLDLST'].apply(lambda x:
'l'+x.replace(', ', ''))
         22     shapefile['state'] = "ny"
         23     leg_dist_areas = leg_dist_areas.append(shapefile[['area',
'temp_district', 'state']], ignore_index=True)
         24
         25

```

```

Out[64]:

```

	area	temp_district	state
0	10253.021636	u049	ny
1	16.347054	u017	ny
2	131.168312	u006	ny
3	13.213151	u020	ny
4	1622.232442	u056	ny
...	...	...	...
208	7.380144	l053	ny
209	996.896394	l105	ny
210	174.316577	l129	ny
211	5.627293	l042	ny
212	162.113485	l014	ny

213 rows × 3 columns



```

In [65]: 1 holdout_data = pd.merge(left=holdout_data,

```

```
2         right=leg_dist_areas,
3         how='left',
4         on=['temp_district', 'state'])
5 holdout_data['pop_density_km'] = holdout_data['Total
  Pop']/holdout_data['area']
```

Out[65]:

	dem_incumbent	dem_net_vote	district	previous_winner_dem	rep_incumbent	state	temp_c
0	1.0	0.615653	l001-0	1.0	0.0	ny	
1	1.0	0.408015	l002-0	0.0	0.0	ny	
2	1.0	0.460065	l003-0	0.0	0.0	ny	
3	1.0	0.606159	l004-0	1.0	0.0	ny	
4	1.0	0.413354	l005-0	0.0	0.0	ny	
...	...	...	...	...	...	...	
125	1.0	0.480922	u055-0	0.0	0.0	ny	
126	1.0	0.444585	u056-0	0.0	0.0	ny	
127	1.0	0.405021	u058-0	0.0	0.0	ny	
128	1.0	0.442403	u060-0	0.0	0.0	ny	
129	1.0	0.458535	u061-0	0.0	0.0	ny	

130 rows × 50 columns



```
In [66]: 1 # converting census data to percentage of total pop
2 for each in list(values_to_grab.values())[1:]:
3     holdout_data[each] = holdout_data[each]/holdout_data['Total Pop']
4
```

Out[66]:

	dem_incumbent	dem_net_vote	district	previous_winner_dem	rep_incumbent	state	temp_c
0	1.0	0.615653	l001-0	1.0	0.0	ny	
1	1.0	0.408015	l002-0	0.0	0.0	ny	
2	1.0	0.460065	l003-0	0.0	0.0	ny	
3	1.0	0.606159	l004-0	1.0	0.0	ny	

	dem_incumbent	dem_net_vote	district	previous_winner_dem	rep_incumbent	state	temp_c
4	1.0	0.413354	l005-0	0.0	0.0	ny	
...	...	...	...	...	...	...	
125	1.0	0.480922	u055-0	0.0	0.0	ny	

```
In [67]: 1 holdout_data = holdout_data.drop(columns=['Total Pop', 'area',
```

```
In [68]: 1 holdout_data['dem_net_vote'] =
holdout_data['dem_net_vote'].apply(lambda x: -1 if x<= .45 else (1 if
>= .55 else 0))
2 holdout_data = holdout_data.dropna()
```

Out[68]:

	dem_incumbent	dem_net_vote	previous_winner_dem	rep_incumbent	state	temp_district
0	1.0	1	1.0	0.0	ny	l001 2
1	1.0	-1	0.0	0.0	ny	l002 2
2	1.0	0	0.0	0.0	ny	l003 2
3	1.0	1	1.0	0.0	ny	l004 2
4	1.0	-1	0.0	0.0	ny	l005 2
...	...	...	...	...	...	...
125	1.0	0	0.0	0.0	ny	u055 2
126	1.0	-1	0.0	0.0	ny	u056 2
127	1.0	-1	0.0	0.0	ny	u058 2
128	1.0	-1	0.0	0.0	ny	u060 2
129	1.0	0	0.0	0.0	ny	u061 2

129 rows × 47 columns



## Model and Feature Selection

```
In [138]: 1 import pickle
2
3 class BorutaFeatureSelection():
4     def __init__(self):
5         pass
6
7     def fit(self, X, y):
8         self.target_encoder =
ce.target_encoder.TargetEncoder().fit(X[['state', 'temp_district']],
y)
9         Xtr = self.target_encoder.transform(X[['state',
'temp_district']])
10        # print(Xtr)
11        X['state'] = Xtr['state']
12        X['temp_district'] = Xtr['temp_district']
13        # print(X['state'])
14        perc = 100
15        while True:
16            boruta = BorutaPy(
17                estimator = RandomForestClassifier(),
18                n_estimators = 'auto',
19                perc=perc
20            ).fit(np.array(X), pd.Series(y))
21
22            self.green_area = X.columns[boruta.support_].to_list()
23            if perc == 10:
24                break
25            if len(self.green_area) >= 5:
26                break
27            perc -= 10
28
29        return(self)
30
31
32    def transform(self, X):
33        if X['state'].dtype == 'O':
34            Xtr = self.target_encoder.transform(X[['state',
'temp_district']])
35            X['state'] = Xtr['state']
36            X['temp_district'] = Xtr['temp_district']
37            return(X[self.green_area])
38
39    def fit_transform(self, X, y):
40        self.fit(X, y)
41        return(self.transform(X))
42
43 X = sample.drop(["census_matching_temp_district", 'district',
'dem_vote'], axis=1)
44 y = pd.Series(sample['dem_vote'])
45
46 b = BorutaFeatureSelection()
47
48 X = b.fit_transform(X, y)
49
50
```

```

X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y,
random_state=12)
51
52 X_holdout = holdout_data.drop(['dem_net_vote'], axis=1)
53 y_holdout = holdout_data['dem_net_vote']
54
55 target_encoder = ce.target_encoder.TargetEncoder().fit(X[['state',
'temp_district']], y)
56 Xtr = b.target_encoder.transform(X_holdout[['state',
'temp_district']])
57 # print(Xtr)
58 X_holdout['state'] = Xtr['state']
59 X_holdout['temp_district'] = Xtr['temp_district']
60

```

### Baseline model - Guessing the mode

```

In [70]: 1 from scipy import stats
2
3 # "Learn" the mode from the training data
4 mode_train = stats.mode(y_train)
5 # Get predictions on the test set
6 baseline_predictions = np.ones(y_test.shape) * mode_train.mode
7 # Compute MAE
8 print(classification_report(y_test, baseline_predictions))
9 # Recall_baseline = recall_score(y_test, baseline_predictions,
10 #                               average='micro')

```

	precision	recall	f1-score	support
-1	0.46	1.00	0.63	980
0	0.00	0.00	0.00	473
1	0.00	0.00	0.00	683
accuracy			0.46	2136
macro avg	0.15	0.33	0.21	2136
weighted avg	0.21	0.46	0.29	2136

```

In [71]: 1 national_models = []
2
3 models = [
4     RandomForestClassifier(),
5     DecisionTreeClassifier(),
6     LogisticRegression(),
7     xgb.XGBClassifier(),
8     lgb.LGBMClassifier(),
9     CatBoostClassifier(verbose=False)
10 ]
11
12 for model in models:
13
14     pipe = Pipeline([
15         ('scaler', StandardScaler()),
16         ("classifier", model)

```

```

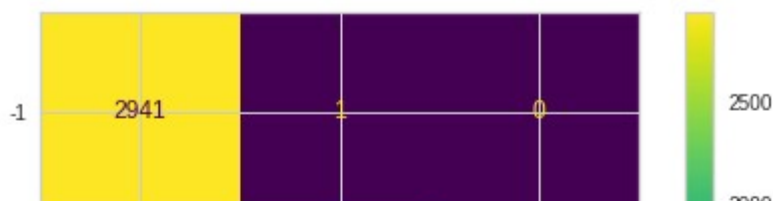
17     ])
18
19     start = datetime.now()
20     pipe.fit(X_train, y_train)
21     print('Fit time: ', datetime.now() - start)
22
23     train_preds = pipe.predict(X_train)
24     test_preds = pipe.predict(X_test)
25     holdout_preds = pipe.predict(X_holdout)
26
27     print('Train Results')
28
29     print(classification_report(y_train, train_preds))
30     plot_confusion_matrix(pipe, X_train, y_train)
31     plt.show()
32
33     print('Test Results')
34
35     print(classification_report(y_test, test_preds))
36     plot_confusion_matrix(pipe, X_test, y_test)
37     plt.show()
38
39     print('Holdout Results')
40
41     print(classification_report(y_holdout, holdout_preds))
42     plot_confusion_matrix(pipe, X_holdout, y_holdout)
43     plt.show()
44
45

```

Fit time: 0:00:03.474945

Train Results

	precision	recall	f1-score	support
-1	1.00	1.00	1.00	2942
0	1.00	0.99	1.00	1419
1	1.00	1.00	1.00	2047
accuracy			1.00	6408
macro avg	1.00	1.00	1.00	6408
weighted avg	1.00	1.00	1.00	6408



In [72]:

Out[72]:

```
{'auto_class_weights': 'None',
 'bagging_temperature': 1,
 'bayesian_matrix_reg': 0.10000000149011612,
 'best_model_min_trees': 1,
 'boost_from_average': False,
 'boosting_type': 'Plain',
 'bootstrap_type': 'Bayesian',
 'border_count': 254,
 'class_names': [-1, 0, 1],
 'classes_count': 0,
 'depth': 6,
 'eval_metric': 'MultiClass',
 'feature_border_type': 'GreedyLogSum',
 'force_unit_auto_pair_weights': False,
 'grow_policy': 'SymmetricTree',
 'iterations': 1000,
 'l2_leaf_reg': 3,
 'leaf_estimation_backtracking': 'AnyImprovement',
 'leaf_estimation_iterations': 1,
 'leaf_estimation_method': 'Newton',
 'learning_rate': 0.08698900043964386,
 'loss_function': 'MultiClass',
 'max_leaves': 64,
 'min_data_in_leaf': 1,
 'model_shrink_mode': 'Constant',
 'model_shrink_rate': 0,
 'model_size_reg': 0.5,
 'nan_mode': 'Min',
 'penalties_coefficient': 1,
 'pool_metainfo_options': {'tags': {}},
 'posterior_sampling': False,
 'random_seed': 0,
 'random_strength': 1,
 'rsm': 1,
```

All these classifiers are performing pretty well but most are showing serious overfitting, as indicated by the fact that the train scores in all but the Logistic Regression model are significantly better than either the test or holdout scores. Given these results I feel that both the XGB and LGBM models show the strongest likelihood of improving given hyperparameter tuning.

## Model Iteration

In [78]:

```
1 # ...
```

Out[78]:



```
{'auto_class_weights': 'None',
 'bagging_temperature': 1,
 'bayesian_matrix_reg': 0.10000000149011612,
 'best_model_min_trees': 1,
 'boost_from_average': False,
 'boosting_type': 'Plain',
 'bootstrap_type': 'Bayesian',
 'border_count': 254,
 'class_names': [-1, 0, 1],
 'classes_count': 0,
 'depth': 6,
 'eval_metric': 'MultiClass',
 'feature_border_type': 'GreedyLogSum',
 'force_unit_auto_pair_weights': False,
 'grow_policy': 'SymmetricTree',
 'iterations': 1000,
 'l2_leaf_reg': 3,
 'leaf_estimation_backtracking': 'AnyImprovement',
 'leaf_estimation_iterations': 1,
 'leaf_estimation_method': 'Newton',
 'learning_rate': 0.08698900043964386,
 'loss_function': 'MultiClass',
 'max_leaves': 64,
 'min_data_in_leaf': 1,
 'model_shrink_mode': 'Constant',
 'model_shrink_rate': 0,
 'model_size_reg': 0.5,
 'nan_mode': 'Min',
 'penalties_coefficient': 1,
 'pool_metainfo_options': {'tags': {}},
```

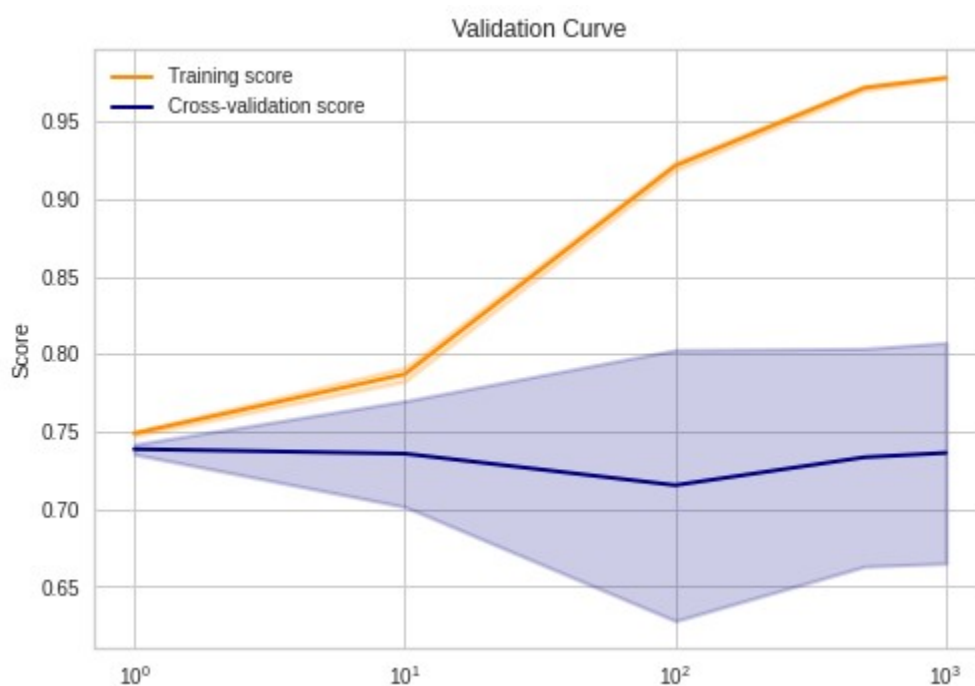
Seeing as the CatBoost model default is overfit, let's try reducing the number of trees through the iterations parameter.

```
In [79]: 1
2 pipe = Pipeline([
3     ('scaler', StandardScaler())])
4
5
6 param_range = [1, 10, 100, 500, 1000]
7
8 cb = national_models[-1]
9
10 train_scores, test_scores = validation_curve(estimator = cb,
11                                             X = pipe.fit_transform(X, y),
12                                             y = y,
13                                             param_name='iterations',
14                                             param_range=param_range)
15
16 train_scores_mean = np.mean(train_scores, axis=1)
17 train_scores_std = np.std(train_scores, axis=1)
18 test_scores_mean = np.mean(test_scores, axis=1)
19 test_scores_std = np.std(test_scores, axis=1)
20
21 plt.title("Validation Curve")
```

```

22
23 plt.ylabel("Score")
24
25 lw = 2
26 plt.semilogx(param_range, train_scores_mean, label="Training score",
27               color="darkorange", lw=lw)
28 plt.fill_between(param_range, train_scores_mean - train_scores_std,
29                  train_scores_mean + train_scores_std, alpha=0.2,
30                  color="darkorange", lw=lw)
31 plt.semilogx(param_range, test_scores_mean, label="Cross-validation
32               score",
33               color="navy", lw=lw)
34 plt.fill_between(param_range, test_scores_mean - test_scores_std,
35                  test_scores_mean + test_scores_std, alpha=0.2,
36                  color="navy", lw=lw)
37 plt.legend(loc="best")

```



The cross-validation score looks to peak around 700 iterations. Let's see how that performs.

```

In [98]: 1 cb = CatBoostClassifier(verbose = False, iterations=700)
2
3 pipe = Pipeline([
4         ('scaler', StandardScaler()),
5         ('classifier', cb)
6     ])
7
8 start = datetime.now()
9 pipe.fit(X_train, y_train)
10 print('Fit time: ', datetime.now() - start)
11
12 train_preds = pipe.predict(X_train)
13 test_preds = pipe.predict(X_test)
14 holdout_preds = pipe.predict(X_holdout)

```

```

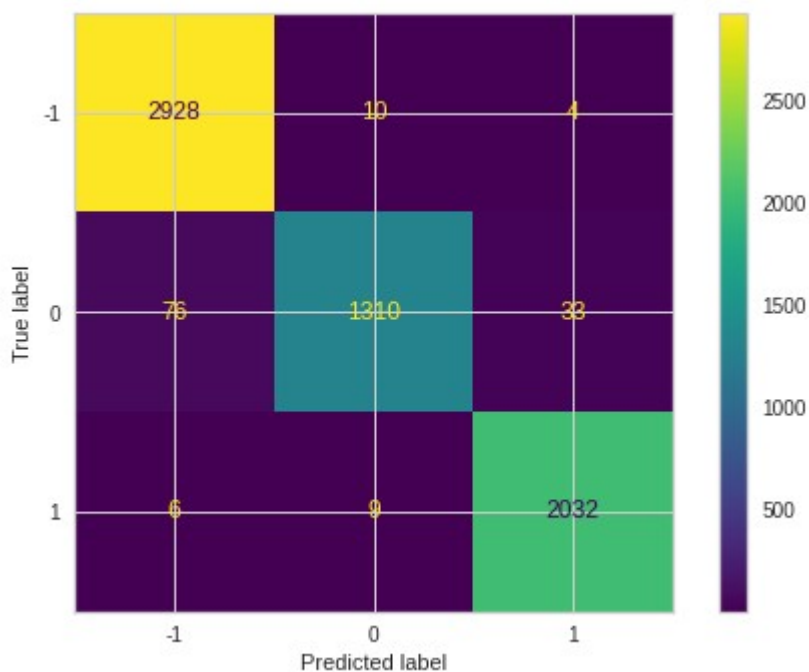
15
16 print('Train Results')
17
18 print(classification_report(y_train, train_preds))
19 plot_confusion_matrix(pipe, X_train, y_train)
20 plt.show()
21
22 print('Test Results')
23
24 print(classification_report(y_test, test_preds))
25 plot_confusion_matrix(pipe, X_test, y_test)
26 plt.show()
27
28 print('Holdout Results')
29
30 print(classification_report(y_holdout, holdout_preds))
31 plot_confusion_matrix(pipe, X_holdout, y_holdout)
32 plt.show()

```

Fit time: 0:00:17.490548

Train Results

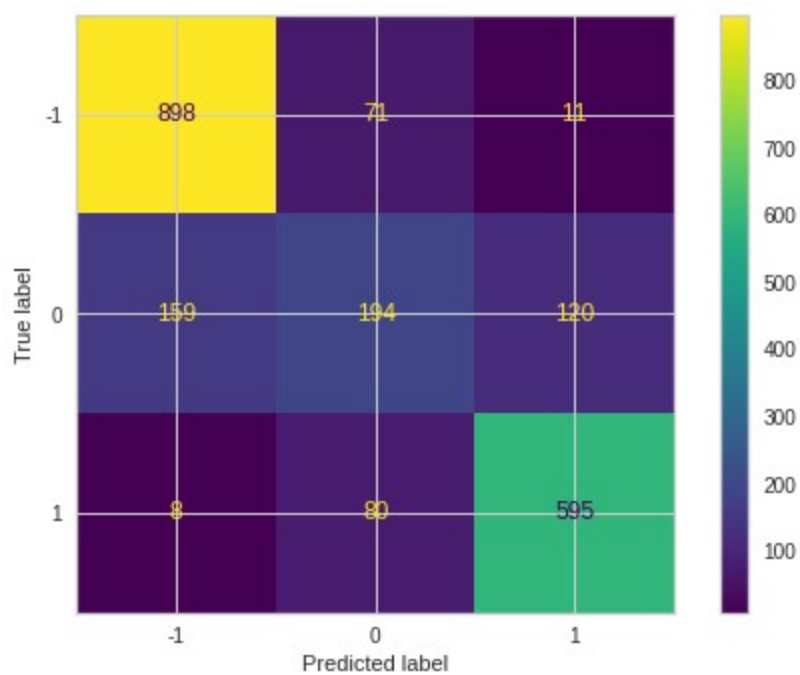
	precision	recall	f1-score	support
-1	0.97	1.00	0.98	2942
0	0.99	0.92	0.95	1419
1	0.98	0.99	0.99	2047
accuracy			0.98	6408
macro avg	0.98	0.97	0.97	6408
weighted avg	0.98	0.98	0.98	6408



Test Results

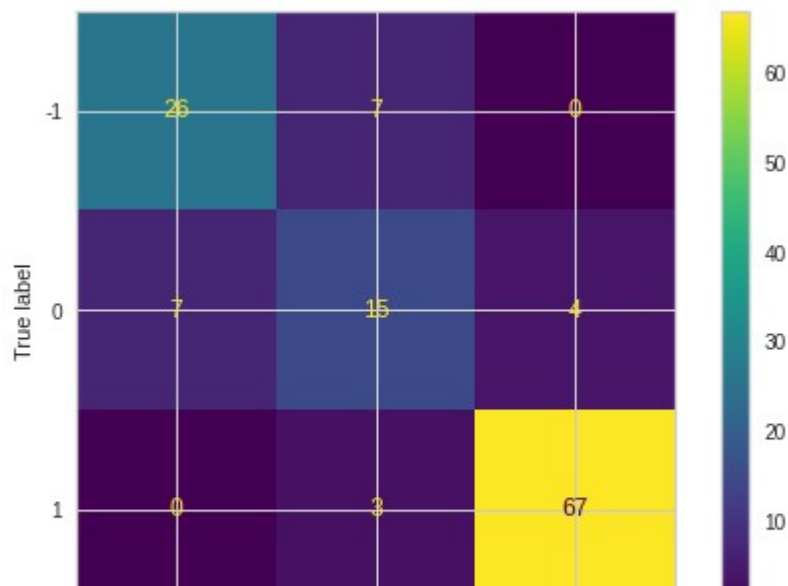
	precision	recall	f1-score	support
-1	0.84	0.92	0.88	980
0	0.56	0.41	0.47	473

	1	0.82	0.87	0.84	683
accuracy				0.79	2136
macro avg		0.74	0.73	0.73	2136
weighted avg		0.77	0.79	0.78	2136



#### Holdout Results

	precision	recall	f1-score	support
-1	0.79	0.79	0.79	33
0	0.60	0.58	0.59	26
1	0.94	0.96	0.95	70
accuracy			0.84	129
macro avg	0.78	0.77	0.78	129
weighted avg	0.83	0.84	0.84	129



The test score stayed the same but the holdout score went down a little - from 88 to 84. There's not a lot to read into here other than the fact that overfitting is still present. Per CatBoost's documentation, let's try decreasing the learning rate.

```
In [100]: 1 # pre-existing parameters
```

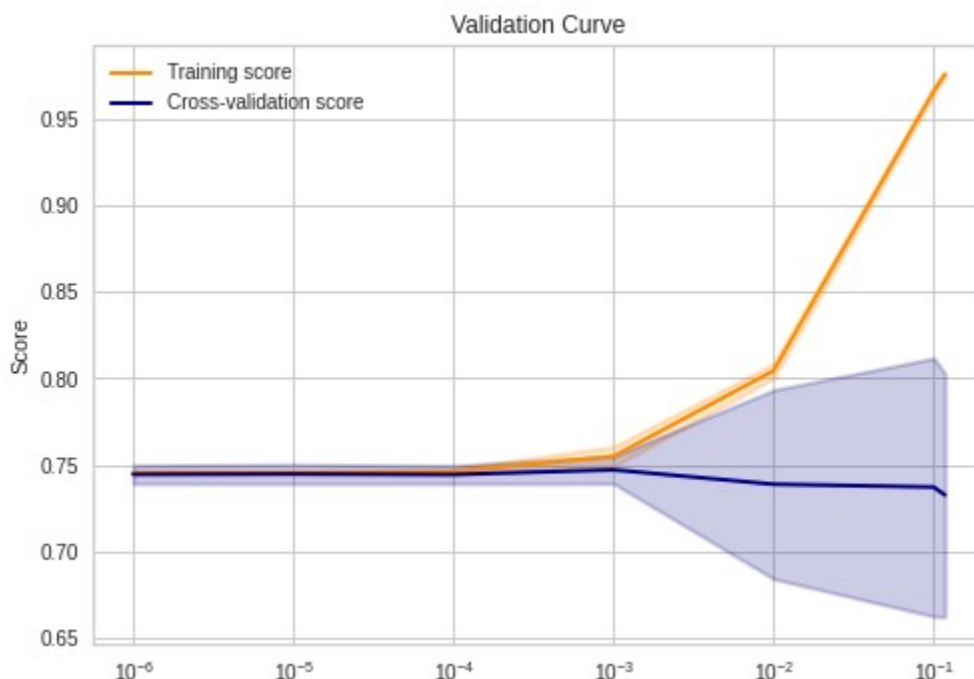
```
Out[100]: {'auto_class_weights': 'None',
            'bagging_temperature': 1,
            'bayesian_matrix_reg': 0.10000000149011612,
            'best_model_min_trees': 1,
            'boost_from_average': False,
            'boosting_type': 'Plain',
            'bootstrap_type': 'Bayesian',
            'border_count': 254,
            'class_names': [-1, 0, 1],
            'classes_count': 0,
            'depth': 6,
            'eval_metric': 'MultiClass',
            'feature_border_type': 'GreedyLogSum',
            'force_unit_auto_pair_weights': False,
            'grow_policy': 'SymmetricTree',
            'iterations': 700,
            'l2_leaf_reg': 3,
            'leaf_estimation_backtracking': 'AnyImprovement',
            'leaf_estimation_iterations': 1,
            'leaf_estimation_method': 'Newton',
            'learning_rate': 0.11758700013160706,
            'loss_function': 'MultiClass',
            'max_leaves': 64,
            'min_data_in_leaf': 1,
            'model_shrink_mode': 'Constant',
            'model_shrink_rate': 0,
            'model_size_reg': 0.5,
            'nan_mode': 'Min',
            'penalties_coefficient': 1,
            'pool_metainfo_options': {'tags': {}},
            'posterior_sampling': False,
            'random_seed': 0,
            'random_strength': 1,
            'rsm': 1,
            'sampling_frequency': 'PerTree',
            'score_function': 'Cosine',
            'sparse_features_conflict_fraction': 0,
            'task_type': 'CPU',
            'use_best_model': False}
```

```
In [101]: 1 pipe = Pipeline([
            2 ('scaler', StandardScaler())])
            3
            4
            5 param_range = [0.11758700013160706, .1, .01, .001, .0001, .00001,
            6 .000001]
```

```

7 cb = CatBoostClassifier(verbose=False, iterations=700)
8
9 train_scores, test_scores = validation_curve(estimator = cb,
10                                             X = pipe.fit_transform(X, y),
11                                             y = y,
12                                             param_name='learning_rate',
13                                             param_range=param_range)
14
15 train_scores_mean = np.mean(train_scores, axis=1)
16 train_scores_std = np.std(train_scores, axis=1)
17 test_scores_mean = np.mean(test_scores, axis=1)
18 test_scores_std = np.std(test_scores, axis=1)
19
20 plt.title("Validation Curve")
21
22 plt.ylabel("Score")
23
24 lw = 2
25 plt.semilogx(param_range, train_scores_mean, label="Training score",
26              color="darkorange", lw=lw)
27 plt.fill_between(param_range, train_scores_mean - train_scores_std,
28                 train_scores_mean + train_scores_std, alpha=0.2,
29                 color="darkorange", lw=lw)
30 plt.semilogx(param_range, test_scores_mean, label="Cross-validation
31              score",
32              color="navy", lw=lw)
33 plt.fill_between(param_range, test_scores_mean - test_scores_std,
34                 test_scores_mean + test_scores_std, alpha=0.2,
35                 color="navy", lw=lw)
36 plt.legend(loc="best")

```



It looks like the cross-validation score is highest around .0001. Let's see where that gets us.

```

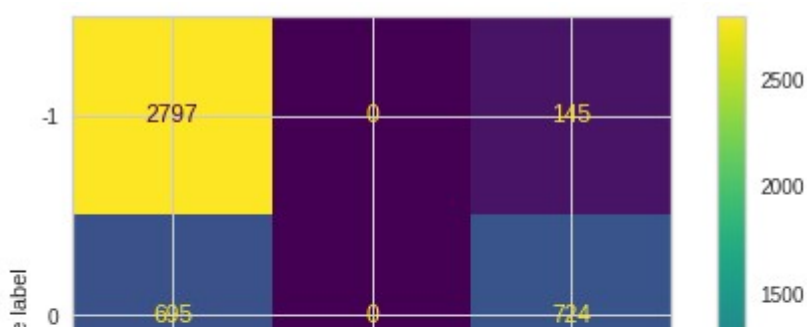
In [139]: 1 cb = CatBoostClassifier(verbose = False, iterations=700,
           2 learning_rate=.0001)
           3
           4 pipe = Pipeline([
           5     ('scaler', StandardScaler()),
           6     ("classifier", cb)
           7 ])
           8
           9 start = datetime.now()
          10 pipe.fit(X_train, y_train)
          11 print('Fit time: ', datetime.now() - start)
          12
          13 train_preds = pipe.predict(X_train)
          14 test_preds = pipe.predict(X_test)
          15 holdout_preds = pipe.predict(X_holdout)
          16
          17 print('Train Results')
          18
          19 print(classification_report(y_train, train_preds))
          20 plot_confusion_matrix(pipe, X_train, y_train)
          21 plt.show()
          22
          23 print('Test Results')
          24
          25 print(classification_report(y_test, test_preds))
          26 plot_confusion_matrix(pipe, X_test, y_test)
          27 plt.show()
          28
          29 print('Holdout Results')
          30
          31 print(classification_report(y_holdout, holdout_preds))
          32 plot_confusion_matrix(pipe, X_holdout, y_holdout)

```

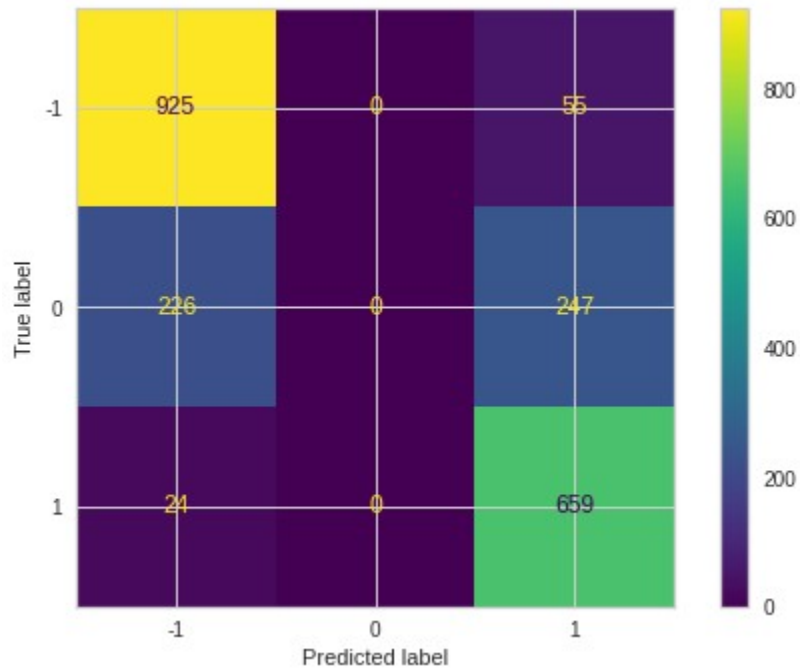
Fit time: 0:00:16.956833

Train Results

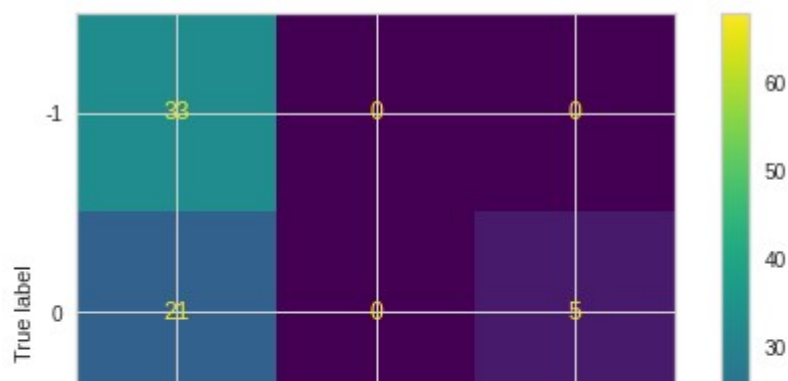
	precision	recall	f1-score	support
-1	0.78	0.95	0.86	2942
0	0.00	0.00	0.00	1419
1	0.69	0.96	0.81	2047
accuracy			0.74	6408
macro avg	0.49	0.64	0.56	6408
weighted avg	0.58	0.74	0.65	6408



Test Results					
		precision	recall	f1-score	support
	Predicted label				
	-1	0	1		
-1	0.79	0.94	0.86	980	
0	0.00	0.00	0.00	473	
1	0.69	0.96	0.80	683	
accuracy				0.74	2136
macro avg	0.49	0.64	0.55		2136
weighted avg	0.58	0.74	0.65		2136



Holdout Results					
		precision	recall	f1-score	support
	Predicted label				
	-1	0	1		
-1	0.59	1.00	0.74	33	
0	0.00	0.00	0.00	26	
1	0.93	0.97	0.95	70	
accuracy				0.78	129
macro avg	0.51	0.66	0.56	129	
weighted avg	0.66	0.78	0.71	129	



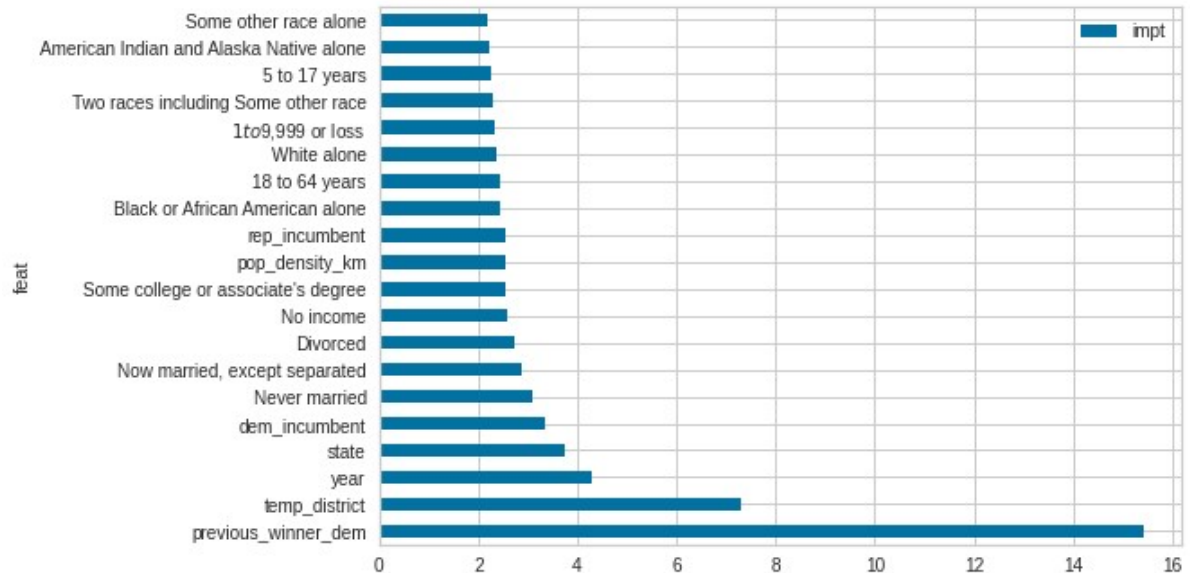


In [99]:

```

1 # feature importances
2
3 importances = sorted(list(zip(cb.get_feature_importance(),
4                               X.columns)))
5
6 impts = pd.DataFrame(importances, columns=['impt', 'feat'])
7 impts = impts.set_index('feat')
8 impts = impts.sort_values(['impt'], ascending=False)
9 # graph 20 most impactful words to given model
10 impts.iloc[0:20].plot(kind='barh')

```



In [140]:

```

1 from sklearn.preprocessing import scale
2
3 def class_feature_importance(X, Y, feature_importances):
4     N, M = X_test.shape
5     X = scale(X_test)
6
7     out = {}
8     for c in set(test_preds):
9         out[c] = dict(
10             zip(range(N), np.mean(X.test[test_preds==c, :],
11                                   axis=0)*feature_importances)
12         )
13

```

The model's now more generalizable but also undoubtedly worse. It's completely avoiding the toss-up option and is therefor useless.

Let's try working with Random Forest, another model that did pretty well when we first ran it.

In [103]:

```

1 params = [{'criterion': ['gini', 'entropy'],
2                 'bootstrap': [True, False],
3                 'max_depth': [10, 50, 100],
4                 'max_features': ['auto', 'sqrt'],

```

```

5         'min_samples_split': [2, 20, 40],
6         'n_estimators': [10, 100, 1000]
7     }]
8
9     grid_search = RandomizedSearchCV(estimator=RandomForestClassifier(),
10                                     param_distributions=params,
11                                     scoring='recall_micro',
12                                     n_jobs=-1,
13                                     verbose =4
14                                     )
15
16     pipe = Pipeline([
17         ('scaler', StandardScaler()),
18         ("classifier", grid_search)
19     ])
20
21     start = datetime.now()
22     pipe.fit(X, y)
23     print('Fit time: ', datetime.now() - start)
24
25     train_preds = pipe.predict(X_train)
26     test_preds = pipe.predict(X_test)
27     holdout_preds = pipe.predict(X_holdout)
28
29     print('Train Results')
30
31     print(classification_report(y_train, train_preds))
32     plot_confusion_matrix(pipe, X_train, y_train)
33     plt.show()
34
35     print('Test Results')
36
37     print(classification_report(y_test, test_preds))
38     plot_confusion_matrix(pipe, X_test, y_test)
39     plt.show()
40
41     print('Holdout Results')
42
43     print(classification_report(y_holdout, holdout_preds))
44     plot_confusion_matrix(pipe, X_holdout, y_holdout)
45

```

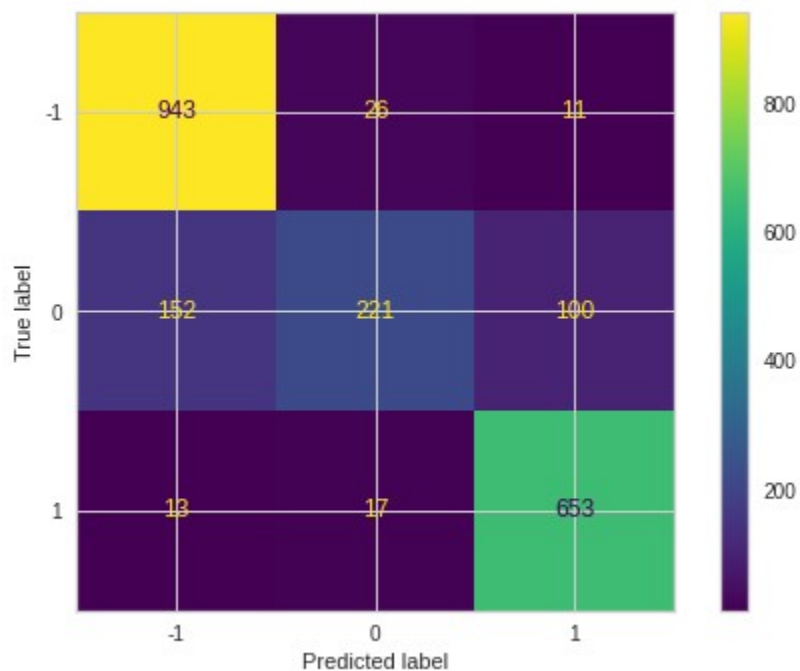
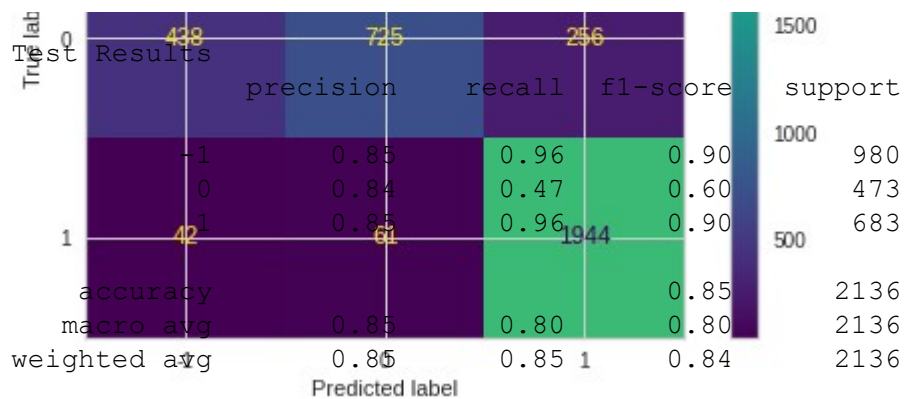
Fitting 5 folds for each of 10 candidates, totalling 50 fits

Fit time: 0:16:35.883930

Train Results

	precision	recall	f1-score	support
-1	0.86	0.96	0.91	2942
0	0.86	0.51	0.64	1419
1	0.87	0.95	0.91	2047
accuracy			0.86	6408
macro avg	0.86	0.81	0.82	6408
weighted avg	0.86	0.86	0.85	6408





### Holdout Results

	precision	recall	f1-score	support
-1	0.63	0.88	0.73	33
0	0.67	0.08	0.14	26
1	0.86	0.99	0.92	70
accuracy			0.78	129
macro avg	0.72	0.65	0.60	129
weighted avg	0.76	0.78	0.71	129



-1                      0                      1  
Predicted label

The improvements here, if they can be called that, are extremely marginal. Let's see if manual adjustments can do better.

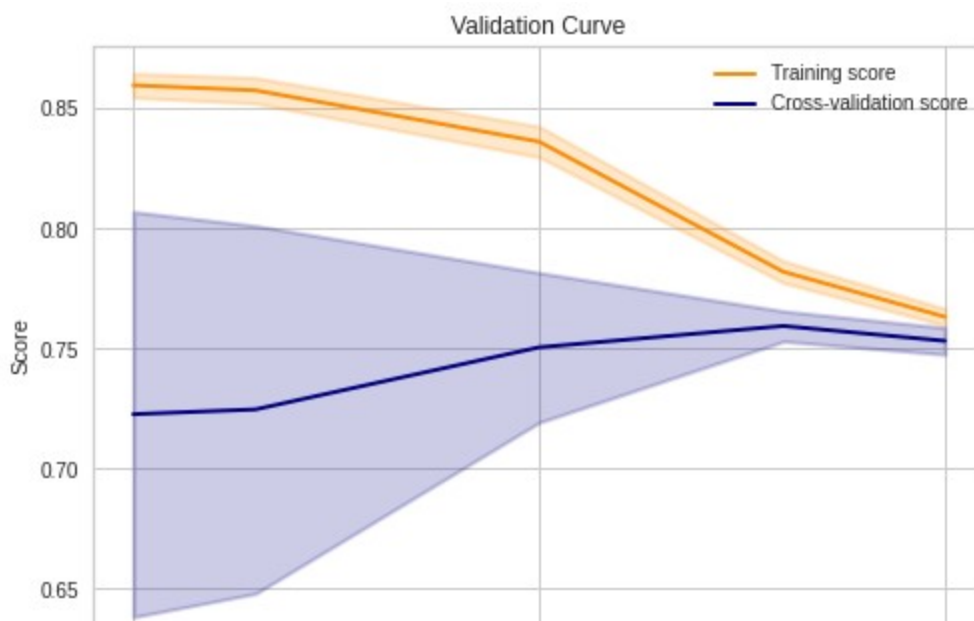
In [104]:

```
Out[104]: {'bootstrap': True,  
           'criterion': 'entropy',  
           'max_depth': 100,  
           'max_features': 'sqrt',  
           'min_samples_split': 40,  
           'n_estimators': 1000}
```

```

In [105]: 1 pipe = Pipeline([
2           ('scaler', StandardScaler())])
3
4
5 param_range = [1, 2, 10, 40, 100]
6
7 rf = RandomForestClassifier().set_params(**grid_search.best_params_)
8
9 train_scores, test_scores = validation_curve(estimator = rf,
10                                              X = pipe.fit_transform(X, y),
11                                              y = y,
12                                              param_name='min_samples_leaf',
13                                              param_range=param_range)
14
15 train_scores_mean = np.mean(train_scores, axis=1)
16 train_scores_std = np.std(train_scores, axis=1)
17 test_scores_mean = np.mean(test_scores, axis=1)
18 test_scores_std = np.std(test_scores, axis=1)
19
20 plt.title("Validation Curve")
21
22 plt.ylabel("Score")
23
24 lw = 2
25 plt.semilogx(param_range, train_scores_mean, label="Training score",
26              color="darkorange", lw=lw)
27 plt.fill_between(param_range, train_scores_mean - train_scores_std,
28                 train_scores_mean + train_scores_std, alpha=0.2,
29                 color="darkorange", lw=lw)
30 plt.semilogx(param_range, test_scores_mean, label="Cross-validation
31              score",
32              color="navy", lw=lw)
33 plt.fill_between(param_range, test_scores_mean - test_scores_std,
34                 test_scores_mean + test_scores_std, alpha=0.2,
35                 color="navy", lw=lw)
36 plt.legend(loc="best")

```



It looks like the cross-validation score peaks around the 40 parameter. Let's run it!

```

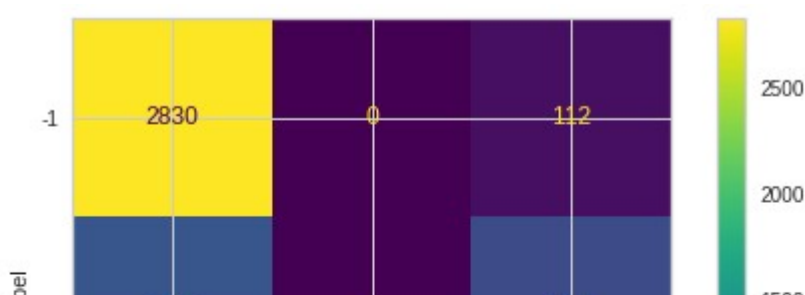
In [124]: 1 rf =
           RandomForestClassifier(min_samples_leaf=400).set_params(**grid_search
           best_params_)
           2
           3 pipe = Pipeline([
           4             ('scaler', StandardScaler()),
           5             ("classifier", rf)
           6         ])
           7
           8 start = datetime.now()
           9 pipe.fit(X_train, y_train)
          10 print('Fit time: ', datetime.now() - start)
          11
          12 train_preds = pipe.predict(X_train)
          13 test_preds = pipe.predict(X_test)
          14 holdout_preds = pipe.predict(X_holdout)
          15
          16 print('Train Results')
          17
          18 print(classification_report(y_train, train_preds))
          19 plot_confusion_matrix(pipe, X_train, y_train)
          20 plt.show()
          21
          22 print('Test Results')
          23
          24 print(classification_report(y_test, test_preds))
          25 plot_confusion_matrix(pipe, X_test, y_test)
          26 plt.show()
          27
          28 print('Holdout Results')
          29
          30 print(classification_report(y_holdout, holdout_preds))
          31 plot_confusion_matrix(pipe, X_holdout, y_holdout)
          32 plt.show()

```

Fit time: 0:00:17.348502

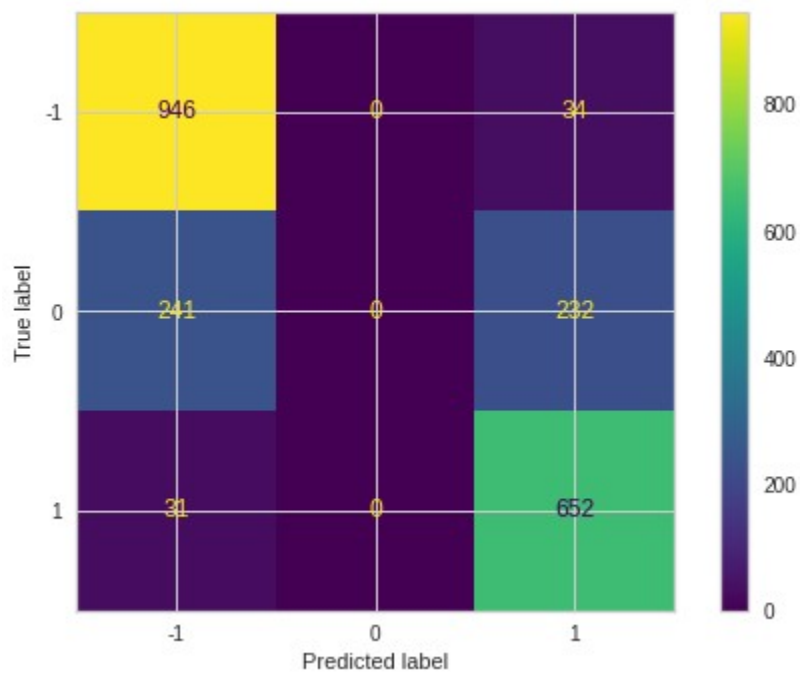
Train Results

	precision	recall	f1-score	support
-1	0.77	0.96	0.85	2942
0	0.00	0.00	0.00	1419
1	0.72	0.95	0.82	2047
accuracy			0.75	6408
macro avg	0.49	0.64	0.56	6408
weighted avg	0.58	0.75	0.65	6408



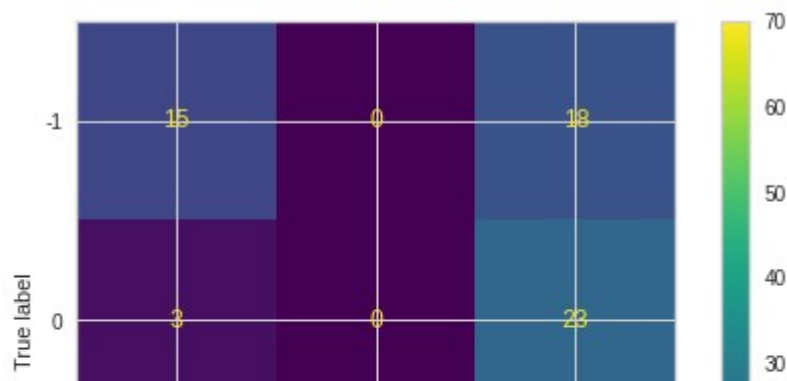
Test Results

	precision	recall	f1-score	support
-1	0.78	0.97	0.86	980
0	0.00	0.00	0.00	473
1	0.71	0.95	0.81	683
accuracy			0.75	2136
macro avg	0.50	0.64	0.56	2136
weighted avg	0.58	0.75	0.66	2136



Holdout Results

	precision	recall	f1-score	support
-1	0.83	0.45	0.59	33
0	0.00	0.00	0.00	26
1	0.63	1.00	0.77	70
accuracy			0.66	129
macro avg	0.49	0.48	0.45	129
weighted avg	0.56	0.66	0.57	129







The scores are across the board worse. While we have fixed our overfitting problem we're once again seeing a model chose to ignore the smallest class, the toss-ups.

One final thing I want to try is simply modeling each state individually to see if those work better than one national model.

```
In [126]: 1 state_models = pd.DataFrame()
2
3 models = [
4     DecisionTreeClassifier(),
5     RandomForestClassifier(),
6     LogisticRegression(),
7     xgb.XGBClassifier(),
8     lgb.LGBMClassifier(),
9     CatBoostClassifier(verbose=False)
10 ]
11
12
13 for state in sample['state'].unique():
14     print(state)
15     X = sample[sample['state'] ==
16 state].drop(["census_matching_temp_district", 'district', 'dem_vote']
17 axis=1)
18     y = pd.Series(sample[sample['state'] == state]['dem_vote'])
19     if state == 'ny':
20         X_holdout = holdout_data.drop(['dem_net_vote'], axis=1)
21         y_holdout = holdout_data['dem_net_vote']
22
23     X_train, X_test, y_train, y_test = train_test_split(X, y,
24 stratify=y, random_state=12)
25
26
27     for model in models:
28
29         pipe = Pipeline([
30             ('feature_selection', BorutaFeatureSelection()),
31             ('scaler', StandardScaler()),
32             ("classifier", model)
33         ])
34
35         start = datetime.now()
36         pipe.fit(X_train, y_train)
37         print('Fit time: ', datetime.now() - start)
38
39         train_preds = pipe.predict(X_train)
40         test_preds = pipe.predict(X_test)
41
42
43         print('Train Results')
44
45         print(classification_report(y_train, train_preds))
46         plot_confusion_matrix(pipe, X_train, y_train)
```

```

46         plt.show()
47
48         print('Test Results')
49
50         print(classification_report(y_test, test_preds))
51         plot_confusion_matrix(pipe, X_test, y_test)
52         plt.show()
53
54         if state == 'ny':
55             target_encoder =
ce.target_encoder.TargetEncoder().fit(X[['state', 'temp_district']],
y)
56             Xtr = b.target_encoder.transform(X_holdout[['state',
'temp_district']])
57             X_holdout['state'] = Xtr['state']
58             X_holdout['temp_district'] = Xtr['temp_district']
59
60             holdout_preds = pipe.predict(X_holdout)
61             print('Holdout Results')
62
63             print(classification_report(y_holdout, holdout_preds))
64             plot_confusion_matrix(pipe, X_holdout, y_holdout)
65             plt.show()
66             state_models = state_models.append(pd.Series({'estimator':
pipe['classifier'], 'test_recall': recall_score(y_test, test_preds,
average='micro'), 'state': state
67                                                         })),

```

az

Fit time: 0:00:19.393298

Train Results

	precision	recall	f1-score	support
-1	1.00	1.00	1.00	25
0	0.92	1.00	0.96	11
1	1.00	0.91	0.95	11
accuracy			0.98	47
macro avg	0.97	0.97	0.97	47
weighted avg	0.98	0.98	0.98	47



In [127]:

Out[127]:

	estimator	state	test_recall
0	DecisionTreeClassifier()	az	0.812500
1	(DecisionTreeClassifier(max_features='auto', r...	az	0.875000

	estimator	state	test_recall
2	LogisticRegression()	az	0.812500
3	XGBClassifier(objective='multi:softprob')	az	0.875000
4	LGBMClassifier()	az	0.812500
...	...	...	...
283	(DecisionTreeClassifier(max_features='auto', r...	ms	0.909091
284	LogisticRegression()	ms	0.818182
285	XGBClassifier(objective='multi:softprob')	ms	0.818182
286	LGBMClassifier()	ms	0.636364
287	<catboost.core.CatBoostClassifier object at 0x...	ms	0.818182

```
In [133]: 1 best_models = pd.DataFrame()
2
3 for state in sample['state'].unique():
4     row = state_models[state_models['state'] ==
5     state].sort_values(['test_recall'], ascending=False).iloc[0]
```

```
In [136]:
```

```
Out[136]:
```

	estimator	state	test_recall
199	(DecisionTreeClassifier(max_features='auto', r...	tn	0.939394
16	LGBMClassifier()	ca	0.934783
163	(DecisionTreeClassifier(max_features='auto', r...	ok	0.923077
282	DecisionTreeClassifier()	ms	0.909091
215	<catboost.core.CatBoostClassifier object at 0x...	ut	0.909091
44	LogisticRegression()	ga	0.909091
109	(DecisionTreeClassifier(max_features='auto', r...	mo	0.905660
49	(DecisionTreeClassifier(max_features='auto', r...	id	0.900000
207	XGBClassifier(objective='multi:softprob')	tx	0.897436
259	(DecisionTreeClassifier(max_features='auto', r...	nj	0.892857
138	DecisionTreeClassifier()	ny	0.886364
62	LogisticRegression()	in	0.880000
97	(DecisionTreeClassifier(max_features='auto', r...	mi	0.879518
1	(DecisionTreeClassifier(max_features='auto', r...	az	0.875000
155	<catboost.core.CatBoostClassifier object at 0x...	nd	0.869565
190	LGBMClassifier()	sc	0.863636
253	(DecisionTreeClassifier(max_features='auto', r...	hi	0.857143
225	XGBClassifier(objective='multi:softprob')	wa	0.857143

	estimator	state	test_recall
115	(DecisionTreeClassifier(max_features='auto', r...	mt	0.854839
177	XGBClassifier(objective='multi:softprob')	pa	0.844156
265	(DecisionTreeClassifier(max_features='auto', r...	va	0.842105
31	(DecisionTreeClassifier(max_features='auto', r...	de	0.842105
248	LogisticRegression()	ak	0.833333
157	(DecisionTreeClassifier(max_features='auto', r...	oh	0.833333
58	LGBMClassifier()	il	0.829268
235	(DecisionTreeClassifier(max_features='auto', r...	wi	0.818182
93	XGBClassifier(objective='multi:softprob')	ma	0.818182
145	(DecisionTreeClassifier(max_features='auto', r...	nc	0.813559
20	LogisticRegression()	co	0.804878
277	(DecisionTreeClassifier(max_features='auto', r...	md	0.800000
270	DecisionTreeClassifier()	al	0.800000
29	<catboost.core.CatBoostClassifier object at 0x...	ct	0.800000
73	(DecisionTreeClassifier(max_features='auto', r...	ks	0.790323
105	XGBClassifier(objective='multi:softprob')	mn	0.780702
169	(DecisionTreeClassifier(max_features='auto', r...	or	0.771429
70	LGBMClassifier()	ia	0.769231
131	<catboost.core.CatBoostClassifier object at 0x...	nh	0.763158
7	(DecisionTreeClassifier(max_features='auto', r...	ar	0.760000
122	LogisticRegression()	nv	0.750000
133	(DecisionTreeClassifier(max_features='auto', r...	nm	0.740741
240	DecisionTreeClassifier()	wy	0.739130
37	(DecisionTreeClassifier(max_features='auto', r...	fl	0.714286
221	<catboost.core.CatBoostClassifier object at 0x...	vt	0.689655
85	(DecisionTreeClassifier(max_features='auto', r...	me	0.679245
83	<catboost.core.CatBoostClassifier object at 0x...	ky	0.666667
192	DecisionTreeClassifier()	sd	0.652174
228	DecisionTreeClassifier()	wv	0.550000

In [137]:

Out[137]:

	test_recall
count	48.000000
mean	0.812041
std	0.091557

	test_recall
min	0.538462
25%	0.767713
50%	0.831301
75%	0.879639
max	0.939304

The results here are very much mixed. The mean micro recall (and, due to the multiclass nature of the problem, the accuracy) is lower than we had for our best RandomForestClassifier. Though certain states did outperform this, we can't conclude that individual state models would outperform one national model.

## Conclusion

Election modeling is a challenging field with near endless directions to explore. What's clear from my analysis is that demographics and certain types of political context are useful in understanding how an election might go, it's far from the whole story. With the vast and ever growing ways people's behavior is tracked and analyzed the amount of features that could be considered useful is hard to imagine. While my models definitely passed their benchmark of 46% accuracy, their inability to reliably identify toss-ups indicates there's still much more work to be done.

Further steps include:

- more exhaustive data manipulation including Feature Unions
- expanding the number of years analyzed
- adding demographic data from more sources, especially commercial ones
- analysis of the impact that other simultaneous races may have