

Data Processing: Ballot Processing

Group 42: Elise Penn, Manish Vuyyuru, Victor Sheng, Yajaira Gonzalez

Processes general congressional ballot data for input into our model. Needs access to the files available on the team GitHub folder, which are too large and numerous to upload here.

```
In [ ]: 1 import pandas as pd
        2 import datetime
        3 import numpy as np
        4 import pickle
```

```
In [ ]: 1 def drop_rows(df, column, value):
        2     ''' Drop a row where 'value' is in 'column'. Only grabs first item.
        3     arguments:
        4         df -- dataframe to modify
        5         column -- the column you want to search
        6         value -- if you find this value in the column, drop the row
        7     returns:
        8         dataframe without specified rows
        9     '''
        10    try:
        11        new_df = df.drop(df[column][df[column]==value].index[0])
        12    except IndexError:
        13        new_df = df
        14    return new_df
```

```
In [ ]: 1 def get_election_day(year):
        2     ''' Get the date of election day in a given year
        3     arguments:
        4         year -- year as int
        5     returns:
        6         datetime object of election day. '''
        7     if year%2 == 1:
        8         raise ValueError('No election in even years.')
        9     # possible days = Nov. 2 - Nov. 8
        10    possible_days = [datetime.datetime(year, 11, d) for d in range(2,9)]
        11    for day in possible_days:
        12        if day.weekday() == 1: # return if it's a Tuesday
        13        return day
```

```

In [ ]: 1 def get_mean_spread(year,n_days_before_election):
2         ''' Cleans data and converts each file into mean spread of all the polls
3         N days before the election
4         arguments:
5             year -- (int) the year you want to pull poll data from
6                   there must be a file with the name
7                   'Datasets/YYYY_generic_congressional_vote.csv'
8             n_days_before_election -- (int) maximum number of days before an
9                                       election a poll should end to be included
10                                      in your estimate
11         returns:
12             mean of the spread N days before the election in that year (float)
13         '''
14         # read in data file
15         ballot_df = pd.read_csv('Datasets/'+str(year)+'_generic_congressional_vote.csv')
16
17         # Data Cleaning
18         ballot_df = drop_rows(ballot_df, 'Poll', 'Final Results')
19         ballot_df = drop_rows(ballot_df, 'Poll', 'RCP Average')
20
21         election_day = get_election_day(year)
22
23         # make spread standardized around 0
24         if year >= 2014: # they changed their column names after 2014
25             ballot_df['Spread'] = ballot_df['Democrats (D)'] - ballot_df['Republicans (R)']
26         else:
27             ballot_df['Spread'] = ballot_df['Democrats'] - ballot_df['Republicans']
28
29         spread = []
30         for index, row in ballot_df.iterrows():
31             # clean up the date format
32             dates = row['Date'].split('-')
33             start = datetime.datetime.strptime(str(year)+'/'+dates[0].strip(), '%Y/%m/%d')
34             #end = datetime.datetime.strptime(str(year)+'/'+dates[1].strip(), '%Y/%m/%d')
35
36             # take all the polls which started less than 4 weeks ago
37             if (start - election_day).days <= n_days_before_election:
38                 spread.append(row['Spread'])
39
40         # find the mean of the spread over the last 4 weeks
41         return np.mean(spread)

```

```
In [ ]: 1 def format_national_polls(years, n_days_before_election=28):
2         ''' Format generic congressional vote into a dataframe with indices
3         of AA_00_0000 (state abbr., district, year). For example, WI_04_2016.
4         There must be a file named 'Datasets/YYYY_generic_congressional_vote.csv'
5         for each year you want to process data.
6         inputs:
7         years -- (list) list of years you want to put in the dataframe
8         n_days_before_election -- (int) maximum number of days before an
9         election a poll should end to be included
10        in your estimate
11        returns:
12        None.
13        For each year, the mean of the spread N days before the election in that
14        calculated. Then we throw it into all the districts for that year.
15        Dumps a dataframe with the proper indexing into 'Datasets/national_poll.p'
16
17        formatted_poll_df = pickle.load(open('Datasets/master_index.p', 'rb'))
18        formatted_poll_df['national_poll'] = np.nan # add new column to the df
19        for year in years:
20            spread = get_mean_spread(year, n_days_before_election)
21            formatted_poll_df.loc[formatted_poll_df['year']==year, 'national_poll'] =
22            spread
23        pickle.dump(formatted_poll_df, open('Datasets/national_poll.p', 'wb'))
```

```
In [ ]: 1 # make the clean data file
2 years = [2002, 2004, 2006, 2008, 2010, 2012, 2014, 2016, 2018]
3 format_national_polls(years, n_days_before_election=28)
```

```
In [ ]: 1 # test the clean data file
2 pickle.load(open('Datasets/national_poll.p', 'rb'))
```

```
In [ ]: 1
```