

Final Model

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

```
In [1]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.image as mpimg
4 import numpy as np
5 import pandas as pd
6 import pickle
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.linear_model import LogisticRegressionCV
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import accuracy_score
11 from sklearn.ensemble import RandomForestClassifier
12 from sklearn.metrics import confusion_matrix
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.ensemble import AdaBoostClassifier
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.linear_model import LassoCV
17 from sklearn.base import clone
18 from sklearn.metrics import r2_score
19 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
20 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis as QDA
21 from sklearn.metrics import log_loss
22 import seaborn as sns
23 import warnings
24
25 warnings.filterwarnings('ignore')
26 %matplotlib inline
```

Functions

```

In [2]: 1 def upsample(test, monsterframe=False):
2         '''
3         input: test (pandas.dataframe)
4         output: pandas.dataframe
5
6         input dataframe MUST have the columns dem_win_prev and dem_win
7         function will ensure that there are an equal number of flips and no flips in
8         '''
9
10        #check if party flip as per usual
11        test['party_flip'] = (test.dem_win_prev != test.dem_win)*1
12        #drop 2018 rows as we want to test on these, don't mess with them
13        test = test[test['year'] != 2018]
14        #count number of flips and no flips
15        print(np.sum((test['party_flip']==1)), np.sum((test['party_flip']==0)))
16        #we want to match these counts, so we need target number of bootstrapped samples
17        target = (np.sum((test['party_flip']==0))) - (np.sum((test['party_flip']==1)))
18        #collect target bootstrap samples where we see a flip
19        samples = test[test['party_flip'] == 1].sample(n=target, replace=True)
20        #add them back to the dataframe
21        test = pd.concat([test, samples])
22        print(np.sum((test['party_flip']==1)), np.sum((test['party_flip']==0)))
23        print(len(set(test[test['party_flip'] == 1].index)))
24        test = test.drop('party_flip', axis=1)
25
26        if monsterframe:
27            target = 5*(len(test))
28            samples = test.sample(n=target, replace=True)
29            test = pd.concat([test, samples])
30        print(len(test))
31
32
33
34        return test
35
36    #full_dataset = bootstrap(full_dataset)

```

```

In [3]: 1 # format data items for input into the model
2
3 def format_model_input(filename, x_cols, y_col, state_hot_encoder=False, upsample
4
5     # load full dataset
6     full_dataset = pickle.load(open(filename, 'rb'))
7
8     # train data is everything except 2018
9     if upsample:
10         pre_2018_dataset = full_dataset.loc[full_dataset['year']!=2018]
11         pre_2018_dataset = upsample(pre_2018_dataset, monsterframe=False)
12     else:
13         pre_2018_dataset = full_dataset.loc[full_dataset['year']!=2018]
14
15     X_train = pre_2018_dataset[cols_to_use]
16     y_train = pre_2018_dataset[y_col]
17     flip_train = np.abs(pre_2018_dataset['dem_win']-pre_2018_dataset['dem_win_pre
18
19     # test data is 2018
20     the_2018_dataset = full_dataset.loc[full_dataset['year']==2018]
21     X_test = the_2018_dataset[cols_to_use]
22     y_test = the_2018_dataset[y_col]
23     flip_test = np.abs(the_2018_dataset['dem_win']-the_2018_dataset['dem_win_prev
24
25     if state_hot_encoder:
26         #add one hot encoder for states (with or without dropping first)
27         X_train.loc[:, 'state'] = pre_2018_dataset['state']
28         X_train = pd.get_dummies(X_train, prefix='state', columns=['state'], drop_f
29         X_test.loc[:, 'state'] = the_2018_dataset['state']
30         X_test = pd.get_dummies(X_test, prefix='state', columns=['state'], drop_fix
31
32     # Make sure train and test have all the same states
33     # the problem is that PA is no longer in the test set...
34     states_missing_in_test = np.array([x for x in X_train.columns.values if x
35     if states_missing_in_test.size!=0:
36         for state_missing in states_missing_in_test:
37             X_test[state_missing] = 0
38
39     return X_train, y_train, X_test, y_test, flip_train, flip_test

```

```

In [4]: 1 # plot all metric specified on the model and dataset specificed
2 # outputs subplot
3
4 def plot_metrics(ax,model_dict, X_train_input, y_train_input, X_test_input, y_test_input,
5                 flip_train, flip_test, score='accuracy_score'):
6
7     df_for_plotting = pd.DataFrame()
8     for model_name, model in model_dict.items():
9         dict_for_plotting = {}
10        # copy arrays so you don't accidentally change them
11        X_train = X_train_input.copy()
12        y_train = y_train_input.copy()
13        X_test = X_test_input.copy()
14        y_test = y_test_input.copy()
15
16        # get the flip and noflip data
17        #flip_mask_train = np.array(flip_train)==1
18        flip_mask_test = np.array(flip_test)==1
19
20        # calculate metrics
21        if score=='accuracy_score':
22            # predict test and train data
23            y_pred_test = model.predict(X_test)
24            y_pred_train = model.predict(X_train)
25
26            # calculate accuracy score
27            y_pred_train = y_pred_train.round()
28            y_pred_test = y_pred_test.round()
29            train_score = accuracy_score(y_train, y_pred_train)
30            test_score = accuracy_score(y_test, y_pred_test)
31            percent_noflip_correct = accuracy_score(y_test[~flip_mask_test], y_pred_test[~flip_mask_test])
32            percent_flip_correct = accuracy_score(y_test[flip_mask_test], y_pred_test[flip_mask_test])
33
34            # colors and names for the plot
35            metric_name = 'accuracy'
36            palette = {'train_'+metric_name : sns.color_palette("Paired")[0],
37                      'test_'+metric_name : sns.color_palette("Paired")[1],
38                      'noflip_'+metric_name : sns.color_palette("Paired")[2],
39                      'flip_'+metric_name : sns.color_palette("Paired")[3] }
40
41        elif score=='log_loss':
42            # predict_proba test and train data
43            y_prob_test = model.predict_proba(X_test)
44            y_prob_train = model.predict_proba(X_train)
45
46            # calculate logloss
47            train_score = log_loss(y_train, y_prob_train)
48            test_score = log_loss(y_test, y_prob_test)
49            percent_noflip_correct = log_loss(y_test[~flip_mask_test], y_prob_test[~flip_mask_test])
50            percent_flip_correct = log_loss(y_test[flip_mask_test], y_prob_test[flip_mask_test])
51
52            # colors and names for the plot
53            metric_name = 'logloss'
54            palette = {'train_'+metric_name : sns.color_palette("Paired")[4],
55                      'test_'+metric_name : sns.color_palette("Paired")[5],
56                      'noflip_'+metric_name : sns.color_palette("Paired")[6],
57                      'flip_'+metric_name : sns.color_palette("Paired")[7] }
58            ax.set_ylim(0,35)
59
60        elif score=='r2_score':
61            # predict_proba test and train data
62            y_prob_test = model.predict_proba(X_test)[:,:1]
63            y_prob_train = model.predict_proba(X_train)[:,:1]
64
65            # calculate r2 score

```

```

In [5]: 1 # plot flips predicted by a model
2
3 def plot_flips(X_test_input, y_test_input, y_pred_test_input, flip_train,
4               flip_test, model_name):
5
6     # copy arrays so you don't accidentally change them
7     X_test = X_test_input.copy()
8     y_test = y_test_input.copy()
9     y_pred_test = y_pred_test_input.copy()
10
11     # add flip data to the array
12     X_test['party_flip'] = flip_test
13
14     #looking at the missclassifications with logisticRegressionCV
15     miss_class_df = X_test[y_test != y_pred_test]
16     good_class_df = X_test[y_test == y_pred_test]
17
18     #plot the flips and non flips for each data
19     fig, ax = plt.subplots(1,2, figsize=(14,4))
20     ax[0].set_title(model_name+' Proportion of party flips on wellclassified test
21     ax[1].set_title(model_name+' Proportion of party flips missclassified test da
22
23     if not good_class_df.empty:
24         good_class_df.groupby('party_flip')['party_flip'].count().plot.bar(ax=ax[
25         for name, group in (good_class_df.groupby('party_flip')['party_flip']):
26             ax[0].text(int(name)+.25, group.count() + .25, \
27                 '{0:.2f}'.format(group.count()/good_class_df.shape[0]*100)+'%', color
28
29     if not miss_class_df.empty:
30         miss_class_df.groupby('party_flip')['party_flip'].count().plot.bar(ax=ax[
31         for name, group in (miss_class_df.groupby('party_flip')['party_flip']):
32             ax[1].text(int(name)+.15, group.count() + .25, \
33                 '{0:.2f}'.format(group.count()/miss_class_df.shape[0]*100)+'%', color
34
35     normnoflip = len(X_test[X_test['party_flip'] == 0])
36     normflip = len(X_test[X_test['party_flip'] == 1])
37     print('predicted correctly. does not flip.: {}'.format(len(good_class_df[good_c
38     print('predicted wrongly. does not flip.: {}'.format(len(miss_class_df[miss_c
39
40     print('predicted correctly. does flips.: {}'.format(len(good_class_df[good_cl
41     print('predicted wrongly. does flips.: {}'.format(len(miss_class_df[miss_clas
42
43     plt.show()

```

```
In [6]: 1 # plot a scatter plot of the probabilities predicted by the model.
2
3 def scatter_results(ax, fitted_model_dict, X_test, y_test, flip, model_name='LogF
4     model = fitted_model_dict[model_name]
5     y_prob = model.predict_proba(X_test)[:,-1]
6
7     flip = flip.values.squeeze()
8     y_test = y_test.values.squeeze()
9     y_dem_flip = y_prob[np.logical_and(y_test==1, flip==1)]
10    y_dem_noflip = y_prob[np.logical_and(y_test==1, flip==0)]
11    y_rep_flip = y_prob[np.logical_and(y_test==0, flip==1)]
12    y_rep_noflip = y_prob[np.logical_and(y_test==0, flip==0)]
13
14    ax.plot(y_dem_noflip, 'b.', label='D, No Flip')
15    ax.plot(y_dem_flip, 'b*', label = 'D, Flip')
16    ax.plot(y_rep_noflip, 'r.', label = 'R, No Flip')
17    ax.plot(y_rep_flip, 'r*', label = 'R, Flip')
18
19    plt.xlabel('index')
20    plt.ylabel('rep-dem')
21    plt.legend()
22    ax.set_title(model_name+' Results')
23    ax.axhline(0.5)
24
25    print('We predicted democrats got {} spots.'.format(np.sum(y_prob.round()==1)
26    print('We predicted republicans got {} spots.'.format(np.sum(y_prob.round()==
```

```

In [7]: 1 # print out R2, accuracy, and plot # flipped seats correctly predicted for a give
2
3 def report_model_stats(filename, x_cols, y_col, model_dict, state_hot_encoder=False):
4
5     # read in the desired data
6     X_train,y_train, X_test,y_test, flip_train,flip_test = \
7         format_model_input(filename, cols_to_use, y_col, state_hot_encoder)
8
9     fitted_model_dict = {}
10
11     for model_name, model in model_dict.items():
12         model_copy = clone(model) # deep copy model to prevent fitting it twice
13         # print(model_name)
14
15         fitted_model = model_copy.fit(X_train,y_train)
16         y_pred_test = fitted_model.predict(X_test)
17
18         # print("Confusion Matrix: \n",confusion_matrix(y_test,y_pred_test.round(
19         # print("\n{} R2: ".format(model_name), fitted_model.score(X_train,y_train)
20         # print("{} Test Accuracy".format(model_name), accuracy_score(y_test,y_pred_test)
21         # plot_flips(X_test, y_test, y_pred_test.round(), flip_train, flip_test,
22
23         fitted_model_dict[model_name] = fitted_model
24
25
26     # try:
27     #     coef = fitted_model.coef_.squeeze()
28     #     print(len(X_test.columns), len(coef))
29     #     print('intercept = {}'.format(fitted_model.intercept_))
30     #     print([i for i in zip(X_test.columns,coef)])
31     # except AttributeError:
32     #     print('doesnt have any coefficients')
33
34     # plot all the metrics
35     fig, ax = plt.subplots(1,3,figsize=(16,4))
36     plot_metrics(ax[0],fitted_model_dict, X_train, y_train, X_test, y_test, flip_train, flip_test)
37     plot_metrics(ax[1],fitted_model_dict, X_train, y_train, X_test, y_test, flip_train, flip_test)
38     scatter_results(ax[2], fitted_model_dict, X_test, y_test, flip_test)
39
40     # rotate tickmarks 45 degrees
41     for ax in fig.axes:
42         plt.sca(ax)
43         plt.xticks(rotation=45)
44
45     fig.suptitle(title)
46     plt.show()
47
48     return fitted_model_dict

```

```

In [8]: 1 # makes a plot to visualize collinearity between predictors
2
3 def collinear(filename, x_cols, state_hot_encoder=False,y_col='dem_win'):
4     X_train,y_train, X_test,y_test, flip_train,flip_test = \
5         format_model_input(filename, cols_to_use, y_col, state_hot_encoder)
6     sns.heatmap(np.abs(X_train[x_cols].corr()),xticklabels=True,yticklabels=True,
7                 vmin = 0, vmax=1)
8     plt.rcParams['figure.figsize'] = (12,8)
9     plt.title('collinearity of predictors considered')

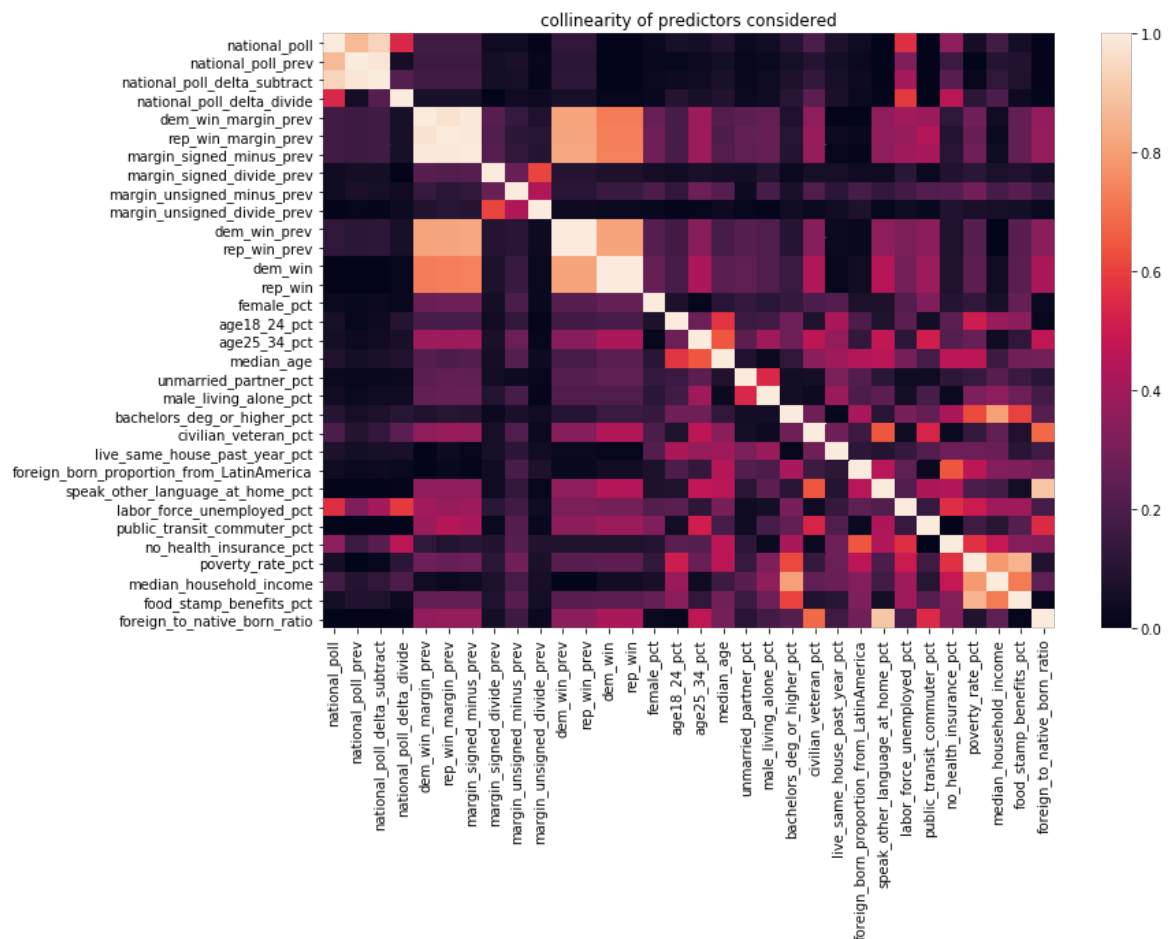
```

Visualization of predictors

```

In [11]: 1 filename = 'Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICS_2010_2018.p'
2 cols_to_use = [ 'national_poll', 'national_poll_prev',
3                 'national_poll_delta_subtract',
4                 'national_poll_delta_divide', 'previous_party',
5                 'dem_win_margin_prev', 'rep_win_margin_prev',
6                 'margin_signed_minus_prev', 'margin_signed_divide_prev',
7                 'margin_unsigned_minus_prev', 'margin_unsigned_divide_prev',
8                 'dem_win_prev', 'rep_win_prev', 'dem_win', 'rep_win',
9                 'female_pct', 'age18_24_pct',
10                'age25_34_pct', 'median_age', 'unmarried_partner_pct',
11                'male_living_alone_pct', 'bachelors_deg_or_higher_pct',
12                'civilian_veteran_pct', #'past_year_births_to_unmarried_women_pct'
13                'live_same_house_past_year_pct',
14                'foreign_born_proportion_from_LatinAmerica',
15                'speak_other_language_at_home_pct', 'labor_force_unemployed_pct',
16                'public_transit_commuter_pct', 'no_health_insurance_pct',
17                'poverty_rate_pct', 'median_household_income', # 'median_housing'
18                'food_stamp_benefits_pct', 'foreign_to_native_born_ratio'
19            ] # everything
20
21 collinear(filename,cols_to_use)
22
23
24 # filename = 'Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICSIMPUTED_2004_2018.p'
25 # collinear(filename,cols_to_use)

```



Try out different models

Change the variables below to try out different models.

These are potential predictors you could include in the model:

```
['district', 'state', 'year', 'party', 'candidatevotes', 'totalvotes', 'candidate', 'national_poll', 'national_poll_prev',  
'national_poll_delta_subtract', 'national_poll_delta_divide', 'previous_party', 'dem_win_margin_prev',  
'rep_win_margin_prev', 'margin_signed_minus_prev', 'margin_signed_divide_prev',  
'margin_unsigned_minus_prev', 'margin_unsigned_divide_prev', 'dem_win_prev', 'rep_win_prev', 'dem_win',  
'rep_win', 'female_pct', 'age18_24_pct', 'age25_34_pct', 'median_age', 'unmarried_partner_pct',  
'male_living_alone_pct', 'bachelors_deg_or_higher_pct', 'past_year_births_to_unmarried_women_pct',  
'civilian_veteran_pct', 'live_same_house_past_year_pct', 'foreign_born_proportion_from_LatinAmerica',  
'speak_other_language_at_home_pct', 'labor_force_unemployed_pct', 'public_transit_commuter_pct',  
'no_health_insurance_pct', 'poverty_rate_pct', 'median_housing_value', 'median_household_income',  
'food_stamp_benefits_pct', 'foreign_to_native_born_ratio']
```

These are potential datasets you could use for the model:

data_FEC_NATIONALPOLL_DEMOGRAPHICSIMPUTED_DISTRICTMIXED_2004_2018.p

data_FEC_NATIONALPOLL_DEMOGRAPHICSIMPUTED_2004_2018_REDISTRICTDROP.p

data_FEC_NATIONALPOLL_DEMOGRAPHICSIMPUTED_2004_2018.p

data_FEC_NATIONALPOLL_DEMOGRAPHICS_2010_2018_REDISTRICTDROP.p

data_FEC_NATIONALPOLL_DEMOGRAPHICS_2010_2018.p

data_FEC_NATIONALPOLL_2004_2018_REDISTRICTDROP.p

data_FEC_NATIONALPOLL_2004_2018.p

If you are running this from the Github repo, remember that all of these datasets are in the Datasets/ folder.

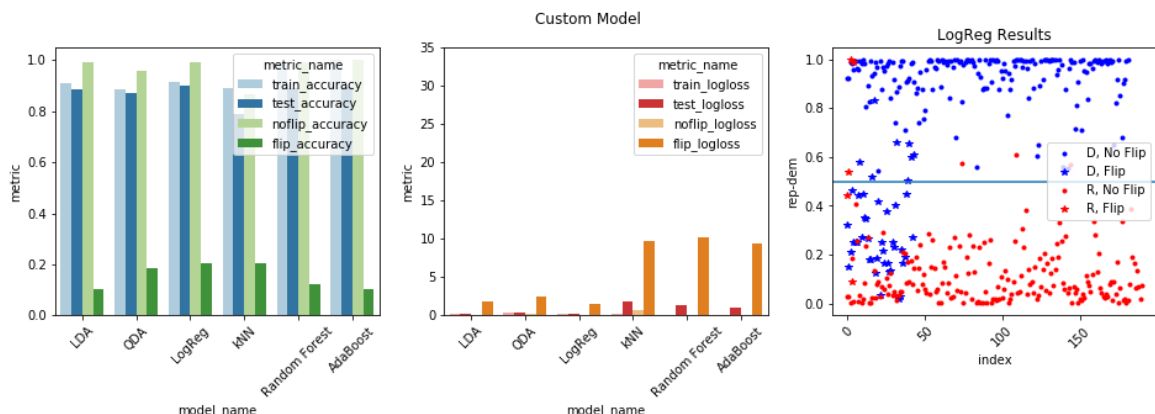
```

In [10]: 1 *****
2 # COMMON VARIABLES TO CHANGE
3 *****
4 # specify the dataset you want to use
5 filename = 'Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICSIMPUTED_DISTRICTMIXED_2004
6 # specify which predictors to use here
7 cols_to_use = ['national_poll', 'margin_signed_minus_prev',
8               'female_pct', 'foreign_to_native_born_ratio',
9               'age18_24_pct', 'age25_34_pct']
10
11 # do you want to try upsampling flipped districts?
12 upsample = False
13 # do you want to one-hot-encode states?
14 state_hot_encoder = False
15
16
17 *****
18 # LESS COMMON VARIABLES TO CHANGE
19 *****
20 # specify which column you are predicting
21 y_col = ['dem_win']
22 # specify your models to use here:
23 model_dict = {
24     'LDA' : LDA(),
25     'QDA' : QDA(),
26     'LogReg' : LogisticRegression(n_jobs=4),
27     'kNN' : KNeighborsClassifier(n_neighbors=4,n_jobs=4),
28     'Random Forest' : RandomForestClassifier(n_jobs=4,random_state=209),
29     'AdaBoost' : AdaBoostClassifier(DecisionTreeClassifier(max_depth=10)),
30 }
31
32
33 *****
34 # The model & visualizations run here
35 *****
36 fitted_model_dict = report_model_stats(
37     filename, cols_to_use, y_col, model_dict, state_hot_encoder=False, upsample=
38     title='Custom Model')

```

We predicted democrats got 196 spots.

We predicted republicans got 226 spots.



In []:

1