

Data Processing: Shapefile Processing

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

Processes shapefiles to get population overlap for use in imputing data.

Needs to be in the folder /shapefiles in our github repo, where it has access to the shapefile data (too big to upload here.)

Takes ~5 hours to run all the way through.

```
In [ ]: 1 import numpy as np
        2 import pickle
        3 import cartopy.io.shapereader as shpreader
        4 from cartopy.feature import ShapelyFeature
        5 from shapely.prepared import prep
        6 import matplotlib.pyplot as plt
        7 import cartopy.crs as ccrs
        8 import pandas as pd
```

```
In [ ]: 1 # Projection is lat/lon (unprojected/cylindrical equidistant)
        2 # The proj.4 string:
        3 # +proj=longlat +ellps=GRS80 +towgs84=0,0,0,0,0,0 +no_defs
```

Functions

```
In [ ]: 1 # some useful arrays to translate between naming conventions
        2
        3 election_year_list = np.array([1992, 1994, 1996, 1998, 2000, 2002, 2004, 2006, 20
        4                                2014, 2016, 2018])
        5 congress_ID_list = np.array([103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 11
        6
        7 state_names = np.array(['ALABAMA', 'ALASKA', 'ARIZONA', 'ARKANSAS', 'CALIFORNIA',
        8                        'COLORADO', 'CONNECTICUT', 'DELAWARE', 'FLORIDA', 'GEORGIA',
        9                        'HAWAII', 'IDAHO', 'ILLINOIS', 'INDIANA', 'IOWA', 'KANSAS',
        10                       'KENTUCKY', 'LOUISIANA', 'MAINE', 'MARYLAND', 'MASSACHUSETTS',
        11                       'MICHIGAN', 'MINNESOTA', 'MISSISSIPPI', 'MISSOURI', 'MONTANA',
        12                       'NEBRASKA', 'NEVADA', 'NEW HAMPSHIRE', 'NEW JERSEY', 'NEW MEXICO',
        13                       'NEW YORK', 'NORTH CAROLINA', 'NORTH DAKOTA', 'OHIO', 'OKLAHOMA',
        14                       'OREGON', 'PENNSYLVANIA', 'RHODE ISLAND', 'SOUTH CAROLINA',
        15                       'SOUTH DAKOTA', 'TENNESSEE', 'TEXAS', 'UTAH', 'VERMONT',
        16                       'VIRGINIA', 'WASHINGTON', 'WEST VIRGINIA', 'WISCONSIN', 'WYOMING'])
        17
        18 state_abbrs = np.array(['AL', 'AK', 'AZ', 'AR', 'CA', 'CO', 'CT', 'DE', 'FL', 'GA', 'HI', 'I
        19                        'IN', 'IA', 'KS', 'KY', 'LA', 'ME', 'MD', 'MA', 'MI', 'MN', 'MS', 'MO', 'MT',
        20                        'NE', 'NV', 'NH', 'NJ', 'NM', 'NY', 'NC', 'ND', 'OH', 'OK', 'OR', 'PA', 'RI',
        21                        'SC', 'SD', 'TN', 'TX', 'UT', 'VT', 'VA', 'WA', 'WV', 'WI', 'WY'])
        22
        23 state_fips = np.array([1, 2, 4, 5, 6, 8, 9, 10, 12, 13, 15, 16, 17, 18, 19, 20, 2
        24                        24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38
        25                        41, 42, 44, 45, 46, 47, 48, 49, 50, 51, 53, 54, 55, 56])
```

```

In [ ]: 1 def read_shapefiles(election_years,verbose=True):
2         """
3         Reads in shapefiles from UCLA database (pre-2016) US Census TigerLine files
4         data sources: http://cdmaps.polisci.ucla.edu/
5                       https://www.census.gov/geo/maps-data/data/cbf/cbf_cds.html
6         Note that file names must be in the same folder as the code with the followi
7         -pre-2016: districtShapes-NNN/districtsNNN.shp
8         -2016 on: must be named tl_YYYY_us_cdNNN/tl_YYYY_us_cdNNN.shp
9         ...where NNN is congress ID and YYYY is election year
10
11         input:
12             election_years -- (list) Elections years of the shape files you want to
13             verbose -- (bool) set to True if you want it to print every file as it r
14         output:
15             district_df -- pandas data frame with index in format ST_00_YYYY (e.g. A
16                           shapefiles stored in a column named 'shape'.
17         """
18         # read in the standard dictionary
19         district_df = pickle.load(open('../Datasets/master_index.p','rb'))
20         district_df['shape'] = [np.nan]*district_df.shape[0] # make a blank column
21         district_df['shape'] = district_df['shape'].astype(object) # reassign to obj
22
23         for election_year in election_years:
24             # convert election year to "Nth Congress"
25             congress_ID = congress_ID_list[election_year_list==election_year][0]
26
27             # pre-2016 files come from http://cdmaps.polisci.ucla.edu/
28             # they are a bit cleaner
29             if election_year < 2016:
30                 # read in the shapefile (must be named 'districtsN.shp' in a folder
31                 shpfilename = 'districtShapes-{0}/districts{0}.shp'.format(congress_
32                 reader = shpreader.Reader(shpfilename)
33                 districts = reader.records() # get full records
34                 geometries = reader.geometries() # get just the shape
35
36                 # put the shapefiles into the standard dictionary
37                 for record in reader.records(): # loop over districts
38                     attr = record.attributes # dictionary of information about the d
39                     poly = record.geometry # coordinates of the district as a shapel
40
41                     # 1) get the state abbr of the district
42                     if any(state_names==attr['STATENAME'].upper()): # filter out dis
43                         ST = state_abbrs[state_names==attr['STATENAME'].upper()][0]
44                     else: # pretty much just Washington, DC
45                         print('{0} is not a state.'.format(attr['STATENAME'].upper()))
46                         continue
47                     # 2) get the id of the district
48                     id_int = int(attr['DISTRICT'])
49                     if id_int == 0: # change at-large district ID from 0 to 1 to pla
50                         id_int = 1 # todo: do we want to change this back? Looks lik
51                     ID = '{0:02d}'.format(id_int)
52                     # 3) reformat the polygon if it intersects with itself
53                     if not poly.is_valid:
54                         # if it does, use buffer to correct this
55                         if verbose:
56                             print('The following polygons intersected with themselve
57                             print(ind)
58                             poly = poly.buffer(0)
59                     # 4) put it all together into an index
60                     ind = '{0}_{1}_{2}'.format(ST, ID, election_year)
61                     if verbose:
62                         print('{0} was read in.'.format(ind))
63                     # put the polygon in the dictionary
64                     district_df.at[ind,'shape'] = poly

```

```

In [ ]: 1 # compute overlap percent between this district and last year's districts
2 def district_overlap(this_year, district_df, threshold_for_change=0.1):
3     """
4     Finds the fractional overlap between this year's district and the previous y
5     i.e., if a Florida's 1st district has changed its borders from 2014 to 2016,
6         20% of its area may come from district 2 in 2014
7         40% of its area may come from district 3 in 2014
8         10% of its area may come from district 4 in 2014
9
10    input:
11        this_year -- (int) Elections years of the districts you want to check. E
12        district_df -- (pd dataframe) dataframe with default indices and shapefi
13                       a column named 'shape'
14
15    output:
16        district_df -- pandas data frame with index in format ST_00_YYYY (e.g. A
17
18        fractional overlap stored in a column called 'fractional_
19        Each district has a dictionary where the keys are the ind
20        districts which overlap with our district, and the values
21
22        change from last year stored in a column called 'border_c
23        Statuses:
24            0 - this district has not changed at all since the pr
25            1 - this district has either changed area > threshold
26                the previous year, or it is new.
27
28    """
29    # add a column
30    if 'overlap_frac' not in district_df.columns:
31        district_df['overlap_frac'] = [np.nan]*district_df.shape[0] # make a bla
32        district_df['overlap_frac'] = district_df['overlap_frac'].astype(object)
33
34    if 'border_change' not in district_df.columns:
35        district_df['border_change'] = [np.nan]*district_df.shape[0] # make a bl
36
37    # loop over states so you only have to compare districts in-state
38    # otherwise, comparing each district to 434 other districts would be super s
39    for ST in state_abbrs:
40        prev_year = this_year-2
41        # get the relevant districts
42        districts = district_df.loc[np.logical_and(district_df['state']==ST,
43                                                  district_df['year']==this_year)
44        districts_prev = district_df.loc[np.logical_and(district_df['state']==ST
45                                                  district_df['year']==prev_year)
46
47        for ind,district in districts.iterrows(): # loop over districts in your
48            overlap_dict = {}
49            shape = district['shape']
50            area = shape.area # area of this district
51
52            # check if shapes intersect with themselves
53            if not shape.is_valid:
54                # if they do, use buffer to correct this
55                print('The following polygons intersected with themselves. Attem
56                print(ind)
57                shape = shape.buffer(0)
58
59            for ind_prev,district_prev in districts_prev.iterrows(): # loop over
60                shape_prev = district_prev['shape']
61
62            # check if shapes intersect with themselves
63            if not shape_prev.is_valid:
64                # if they do, use buffer to correct this
65                print('The following polygons intersected with themselves. A
66                print(ind_prev)
67                shape_prev = shape_prev.buffer(0)

```

```
In [ ]: 1 def get_centroid(district_df):
2         """
3         Finds the centroid of a district in lon,lat.
4
5         input:
6         district_df -- (pd dataframe) dataframe with default indices and shapefile
7                        a column named 'shape'
8
9         output:
10        district_df -- pandas data frame with index in format ST_00_YYYY (e.g. AL
11                      with centroid stored in a column called 'centroid'.
12                      The centroid is calculated as a lon,lat on a Cartesian plane.
13                      It ignores spherical geometry.
14                      The centroid is stored as a tuple in the form (lon,lat)
15
16        """
17        # add column
18        district_df['centroid'] = [np.nan]*district_df.shape[0] # make a blank column
19        district_df['centroid'] = district_df['centroid'].astype(object) # reassign type
20
21        for ind, district in district_df.iterrows():
22            shape = district['shape']
23            if pd.isnull(district['shape']): # if there's no shape, fill it with a name
24                district_df.at[ind, 'centroid'] = np.nan
25            else:
26                centroid = shape.centroid.coords
27                district_df.at[ind, 'centroid'] = centroid # in units of lat/lon
28
29        return district_df
```

```

In [ ]: 1 # compute overlap percent between this district and last year's districts
2 def population_overlap(this_year, district_df, threshold_for_change=0.1):
3     """
4     Estimates percent of population coming from previous district boundaries using
5     the equation:
6
7     
$$\text{population\_overlap} = (\text{overlap\_area} / \text{area}) * (1 / \text{area\_prev})$$

8
9     Where overlap_area is the area of overlap between the new and old districts,
10    area is the area of the new district, and area_prev is the area of the previous
11    district.
12
13    This equation assumes that the voting population in each district is constant
14    but area changes. Although each district has roughly 711,000 people, variation in
15    voter turnout between districts means that this assumption is false.
16
17    density_weighted_overlap is calculated for each previous district which intersects
18    the new district. density_weighted_overlap is then scaled so that the sum of all
19    previous districts which intersect with the new district is 1.
20
21    input:
22        this_year -- (int) Elections years of the districts you want to check. Example: 2018
23        district_df -- (pd dataframe) dataframe with default indices and shapefile
24                      a column named 'shape'
25    output:
26        district_df -- pandas data frame with index in format ST_00_YYYY (e.g. AL002018)
27
28        population_overlap is stored in a column called 'population_overlap'
29        Each district has a dictionary where the keys are the indices of the
30        districts which overlap with our district, and the values are the
31        population_overlap values.
32    """
33
34    # add a column
35    if 'population_overlap' not in district_df.columns:
36        district_df['population_overlap'] = [np.nan]*district_df.shape[0] # make
37        district_df['population_overlap'] = district_df['population_overlap'].astype(float)
38
39    # loop over states so you only have to compare districts in-state
40    # otherwise, comparing each district to 434 other districts would be super slow
41    for ST in state_abbrs:
42        prev_year = this_year-2
43        # get the relevant districts
44        districts = district_df.loc[np.logical_and(district_df['state']==ST,
45                                                  district_df['year']==this_year)]
46        districts_prev = district_df.loc[np.logical_and(district_df['state']==ST,
47                                                      district_df['year']==prev_year)]
48
49        for ind,district in districts.iterrows(): # loop over districts in your current state
50            overlap_dict = {}
51            shape = district['shape']
52            area = shape.area # area of this district
53
54            # check if shapes intersect with themselves
55            if not shape.is_valid:
56                # if they do, use buffer to correct this
57                print('The following polygons intersected with themselves. Attempting to correct...')
58                print(ind)
59                shape = shape.buffer(0)
60
61            for ind_prev,district_prev in districts_prev.iterrows(): # loop over
62                shape_prev = district_prev['shape']
63
64                # check if shapes intersect with themselves
65                if not shape_prev.is_valid:

```

```

In [ ]: 1 def inverse_population_overlap(this_year, district_df):
2         """
3         Exactly the same as population_overlap, except it gives each district_prev a
4         where the keys are current year's districts, and the values are the estimated
5         population overlap of district_prev and district.
6
7         This allows you to impute values from (year) into (year-2) in your dataset.
8
9         """
10
11        # add a column
12        if 'inverse_population_overlap' not in district_df.columns:
13            district_df['inverse_population_overlap'] = [np.nan]*district_df.shape[0]
14            district_df['inverse_population_overlap'] = district_df['inverse_population_overlap']
15
16        # loop over states so you only have to compare districts in-state
17        # otherwise, comparing each district to 434 other districts would be super slow
18        for ST in state_abbrs:
19            prev_year = this_year-2
20            # get the relevant districts
21            districts = district_df.loc[np.logical_and(district_df['state']==ST,
22                                                       district_df['year']==this_year)]
23            districts_prev = district_df.loc[np.logical_and(district_df['state']==ST,
24                                                           district_df['year']==prev_year)]
25
26            for ind_prev,district_prev in districts_prev.iterrows(): # loop over districts
27                overlap_dict = {}
28                shape_prev = district_prev['shape']
29                area_prev = shape_prev.area # area of this district
30
31                # check if shapes intersect with themselves
32                if not shape_prev.is_valid:
33                    # if they do, use buffer to correct this
34                    print('The following polygons intersected with themselves. Attempting to fix...')
35                    print(ind_prev)
36                    shape_prev = shape_prev.buffer(0)
37
38                for ind,district in districts.iterrows(): # loop over districts in year
39                    shape = district['shape']
40
41                    # check if shapes intersect with themselves
42                    if not shape.is_valid:
43                        # if they do, use buffer to correct this
44                        print('The following polygons intersected with themselves. Attempting to fix...')
45                        print(ind)
46                        shape = shape.buffer(0)
47
48                    # calculate frac overlap
49                    area = shape.area # area of current district
50                    density = 1./area # assume population roughly the same, but area
51                    overlap_area = shape.intersection(shape_prev).area # area of overlap
52                    frac_overlap = overlap_area/area_prev # fractional overlap between
53                                                                # round to 3 decimal places
54
55                    if frac_overlap > 10**-3: # use threshold of 0.1% to avoid trivial
56                        population_overlap = frac_overlap*density
57                        overlap_dict[ind] = population_overlap
58
59                # make sure areas add up to 1.
60                dict_sum = np.sum(list(overlap_dict.values()))
61                # rescale so they add up to 1.
62                overlap_dict = {key : np.around(val/dict_sum, decimals=3)
63                                for key,val in overlap_dict.items() if val/dict_sum > 10**-3}
64
65        return overlap_dict

```

```

In [ ]: 1 # find if a district has changed between last year and this year
2 def check_if_districts_changed(this_year, district_df, threshold_for_change=0.1):
3     """
4     *** DEPRECIATED ***
5     might have some bugs.
6     use district_overlap to perform this function instead.
7
8     Checks if each district changed more than the set threshold since the last el
9     Change is a fraction of the district which overlaps with the previous year of
10    If the total area of the district increases, the overlap area is divided by t
11    (this prevents counting a district as unchanged if its area has increased
12
13    input:
14    this_year -- (int) Elections years of the districts you want to check. Ev
15    district_df -- (pd dataframe) dataframe with default indices and shapefil
16                  a column named 'shape'
17
18    output:
19    district_df -- pandas data frame with index in format ST_00_YYYY (e.g. AL
20                  with change from last year stored in a column called 'border
21                  Statuses:
22                  'same' - this district has not changed at all since th
23                  'new' - indicates a district with this ID was not in t
24                  'changed' - indicates the borders have changed from la
25                          this district was present in its state las
26    """
27    # loop over states so you only have to compare districts in-state
28    # otherwise, comparing each district to 434 other districts would be super si
29    for ST in state_abbrs:
30        prev_year = this_year-2
31        # get the relevant districts
32        districts = district_df.loc[np.logical_and(district_df['state']==ST,
33                                                    district_df['year']==this_year)]
34        districts_prev = district_df.loc[np.logical_and(district_df['state']==ST,
35                                                        district_df['year']==prev_year)]
36        # loop over districts in your current year
37        for ind,district in districts.iterrows():
38            # find previous year's district
39            district_prev = districts_prev.loc[districts_prev['district']==district]
40            ind_prev = '{}{}'.format(ind[:6],prev_year)
41
42            # determine whether district is new or borders have changed
43            if ind_prev in districts_prev.index: # if the district didn't exist i
44                # then this district is new this year
45                district_df.loc[ind,'border_change'] = 'new'
46            else:
47                # check if the borders have changed
48                shape_prev = district_prev['shape'].values[0]
49                shape = district['shape']
50
51                # check if shapes intersect with themselves
52                if not (shape.is_valid and shape_prev.is_valid):
53                    # if they do, use buffer to correct this
54                    print('The following polygons intersected with themselves. At
55                    if not shape.is_valid:
56                        print(ind)
57                        shape = shape.buffer(0)
58                    if not shape_prev.is_valid:
59                        print(ind_prev)
60                        shape_prev = shape_prev.buffer(0)
61
62                # calculate overlap percent
63                area = shape.area # area of this district
64                area_prev = shape_prev.area
65                overlap_area = shape.intersection(shape_prev).area # area of over

```


Run code here

```
In [ ]: 1 # get centroid coords for everything
2 years = [2002, 2004, 2006, 2008, 2010, 2012, 2014, 2016, 2018]
3
4 centroid_df = read_shapefiles(years, verbose=False)
5 for year in years:
6     centroid_df = get_centroid(centroid_df)
7 pickle.dump(centroid_df.drop('shape', axis=1), open('centroid.p', 'wb'))
```

```
In [ ]: 1 # find out the area overlap and border change between this year's and last year's
2 # takes ~2hrs to run for all 8 years.
3
4 # you can only check a district if you have last year's data, so you read one mo
5 years_to_read = [2002, 2004, 2006, 2008, 2010, 2012, 2014, 2016, 2018]
6 years_to_check = [2004, 2006, 2008, 2010, 2012, 2014, 2016, 2018]
7
8 # read in the data (make a fresh df)
9 #overlap_df = read_shapefiles(years_to_read, verbose=False)
10 for year in years_to_check:
11     overlap_df = district_overlap(year, overlap_df)
12     pickle.dump(overlap_df.drop('shape', axis=1), open('overlap_frac_{}.p'.format
13 pickle.dump(overlap_df.drop('shape', axis=1), open('overlap_frac.p', 'wb'))
```

```
In [ ]: 1 # estimate % population which derives from each previous district
2
3 # you can only check a district if you have last year's data, so you read one mo
4 years_to_read = [2002, 2004, 2006, 2008, 2010, 2012, 2014, 2016, 2018]
5 years_to_check = [2004, 2006, 2008, 2010, 2012, 2014, 2016, 2018]
6
7 # read in the data (make a fresh df)
8 pop_df = read_shapefiles(years_to_read, verbose=False)
9 for year in years_to_check:
10     pop_df = population_overlap(year, pop_df)
11     pickle.dump(pop_df.drop('shape', axis=1), open('pop_frac_{}.p'.format(year),
12 pickle.dump(pop_df.drop('shape', axis=1), open('pop_frac.p', 'wb'))
```

```
In [ ]: 1 # Read in the files written here and add them to a master overlap set
2 all_overlap_data_df = pickle.load(open('overlap_frac.p', 'rb'))
3 all_overlap_data_df['centroid'] = pickle.load(open('centroid.p', 'rb'))['centroid']
4 all_overlap_data_df['population_overlap'] = \
5     pickle.load(open('pop_frac.p', 'rb'))['population_overlap']
6 all_overlap_data_df['inverse_population_overlap'] = \
7     pickle.load(open('inv_pop_frac.p', 'rb'))['inverse_population_overlap']
8 pickle.dump(all_overlap_data_df, open('all_overlap_data.p', 'wb'))
```

```
In [ ]: 1
```