

Data Processing: ACS Demographics

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

Processes the ACS demographics. Needs access to the files available on the team GitHub folder, which are too large and numerous to upload here.

```
In [1]: 1 #imports
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib.image as mpimg
5 import seaborn as sns
6 import numpy as np
7 import pandas as pd
8 import pickle
9 %matplotlib inline

In [2]: 1 # dictionary to translate numeric codes to state abbrev. names
2 state_codes = {
3 1: "AL", 2: "AK", 4: "AZ", 5: "AR", 6: "CA", 8: "CO", 9: "CT", 10: "DE",
4 11: "DC", 12: "FL", 13: "GA", 15: "HI", 16: "ID", 17: "IL", 18: "IN", 19: "IA",
5 20: "KS", 21: "KY", 22: "LA", 23: "ME", 24: "MD", 25: "MA", 26: "MI", 27: "MN",
6 28: "MS", 29: "MO", 30: "MT", 31: "NE", 32: "NV", 33: "NH", 34: "NJ", 35: "NM",
7 36: "NY", 37: "NC", 38: "ND", 39: "OH", 40: "OK", 41: "OR", 42: "PA", 44: "RI",
8 45: "SC", 46: "SD", 47: "TN", 48: "TX", 49: "UT", 50: "VT", 51: "VA", 53: "WA",
9 54: "WV", 55: "WI", 56: "WY"}

In [3]: 1 # raw ACS data file paths
2 path10 = 'Datasets/demographics/ACS_10_1YR_S0201_with_ann.csv'
3 path12 = 'Datasets/demographics/ACS_12_1YR_S0201_with_ann.csv'
4 path14 = 'Datasets/demographics/ACS_14_1YR_S0201_with_ann.csv'
5 path16 = 'Datasets/demographics/ACS_16_1YR_S0201_with_ann.csv'
6 path17 = 'Datasets/demographics/ACS_17_1YR_S0201_with_ann.csv'
7
8 # read in the raw files
9 df10 = pd.read_csv(path10, header = 1)
10 df12 = pd.read_csv(path12, header = 1)
11 df14 = pd.read_csv(path14, header = 1)
12 df16 = pd.read_csv(path16, header = 1)
13 df17 = pd.read_csv(path17, header = 1)
```

```

In [4]: 1 def clean(df, year):
        2
        3     # insert year column
        4     df.insert(0, 'year', year)
        5
        6     # insert state column (as 2-char code, like: MA)
        7     df.insert(1, 'state', df.Id2.apply(lambda x: int(x / 100)))
        8     df.state = df.state.apply(lambda x: state_codes[x])
        9
        10    # insert district column, as an integer
        11    # districts should start with district 1, not 0
        12    df.insert(2, 'district', df.Id2.apply(lambda x: int(str(x)[-2:])))
        13    df.district = df.district.replace(0, 1)
        14
        15    # filter out Margin of Error columns
        16    columns_to_keep = [col for col in df.columns if not 'Margin of Error;' in col]
        17    df = df[columns_to_keep]
        18
        19    # create index as: state_district_year
        20    df.index = ['{0}_{1:02d}_{2}'.format(row['state'], row['district'], row['year'])
        21                for row in df.iterrows()]
        22    return df

```

```

In [5]: 1 acs10 = clean(df10, 2010)
        2 acs12 = clean(df12, 2012)
        3 acs14 = clean(df14, 2014)
        4 acs16 = clean(df16, 2016)
        5 acs17 = clean(df17, 2017)

```

```

In [6]: 1 # check order of unemployment columns
        2 [col for col in acs10.columns if 'EMPLOYMENT STATUS - In labor force - Civilian l

```

```

Out[6]: ['Estimate; EMPLOYMENT STATUS - In labor force - Civilian labor force - Unemployed
',
'Estimate; EMPLOYMENT STATUS - In labor force - Civilian labor force - Unemployed
- Percent of civilian labor force',
'Estimate; EMPLOYMENT STATUS - In labor force - Civilian labor force - Unemployed
.1',
'Estimate; EMPLOYMENT STATUS - In labor force - Civilian labor force - Unemployed
- Percent of civilian labor force.1']

```

```

In [7]: 1 len(acs17.columns)

```

```

Out[7]: 316

```

```

In [8]: 1 # find the common set of columns between two dataframes
        2 len(set(acs16.columns).intersection(set(acs17.columns)))

```

```

Out[8]: 287

```

```
In [9]: 1 set(acs10.columns).intersection(set(acs12.columns))
```

```
Out[9]: {'Estimate; CLASS OF WORKER - Civilian employed population 16 years and over',
'Estimate; CLASS OF WORKER - Government workers',
'Estimate; CLASS OF WORKER - Private wage and salary workers',
'Estimate; CLASS OF WORKER - Self-employed workers in own not incorporated busine
ss',
'Estimate; CLASS OF WORKER - Unpaid family workers',
'Estimate; COMMUTING TO WORK - Car, truck, or van - carpooled',
'Estimate; COMMUTING TO WORK - Car, truck, or van - drove alone',
'Estimate; COMMUTING TO WORK - Mean travel time to work (minutes)',
'Estimate; COMMUTING TO WORK - Other means',
'Estimate; COMMUTING TO WORK - Public transportation (excluding taxicab)',
'Estimate; COMMUTING TO WORK - Walked',
'Estimate; COMMUTING TO WORK - Worked at home',
'Estimate; COMMUTING TO WORK - Workers 16 years and over',
'Estimate; DISABILITY STATUS - Civilian noninstitutionalized population 18 to 64
years',
'Estimate; DISABILITY STATUS - Civilian noninstitutionalized population 65 years
and older',
'Estimate; DISABILITY STATUS - Civilian noninstitutionalized population under 18
-----'}
```

```
In [10]: 1 #[col for col in df.columns if not 'Margin of Error;' in col]
2
```

```

In [11]: 1 # 2010 and 2012
2 raw_predictor_name_subset_10_12 = [
3 "Estimate; SEX AND AGE - Female",
4 "Estimate; SEX AND AGE - 18 to 24 years",
5 "Estimate; SEX AND AGE - 25 to 34 years",
6 "Estimate; SEX AND AGE - Median age (years)",
7 "Estimate; RELATIONSHIP - Nonrelatives - Unmarried partner",
8 "Estimate; HOUSEHOLDS BY TYPE - Nonfamily households - Male householder - Living
9 "Estimate; EDUCATIONAL ATTAINMENT - Bachelor's degree or higher",
10 "Estimate; FERTILITY - Women 15 to 50 years who had a birth in the past 12 months
11 "Estimate; VETERAN STATUS - Civilian veteran",
12 "Estimate; RESIDENCE 1 YEAR AGO - Same house",
13 "Estimate; PLACE OF BIRTH, CITIZENSHIP STATUS AND YEAR OF ENTRY - Native",
14 "Estimate; PLACE OF BIRTH, CITIZENSHIP STATUS AND YEAR OF ENTRY - Foreign born",
15 "Estimate; WORLD REGION OF BIRTH OF FOREIGN BORN - Latin America",
16 "Estimate; LANGUAGE SPOKEN AT HOME AND ABILITY TO SPEAK ENGLISH - Language other
17 "Estimate; EMPLOYMENT STATUS - In labor force - Civilian labor force - Unemployed
18 "Estimate; COMMUTING TO WORK - Public transportation (excluding taxicab)",
19 "Estimate; HEALTH INSURANCE COVERAGE - No health insurance coverage",
20 "Estimate; POVERTY RATES FOR FAMILIES AND PEOPLE FOR WHOM POVERTY STATUS IS DETER
21 "Estimate; OWNER CHARACTERISTICS - Median value (dollars)"]
22
23 raw_predictor_name_supplement_10 = [
24 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2010 INFLATION-ADJUSTED DOLLARS) - Me
25 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2010 INFLATION-ADJUSTED DOLLARS) - Wi
26
27 raw_predictor_name_supplement_12 = [
28 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2012 INFLATION-ADJUSTED DOLLARS) - Me
29 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2012 INFLATION-ADJUSTED DOLLARS) - Wi
30
31 raw_predictor_names_10 = raw_predictor_name_subset_10_12 + raw_predictor_name_sup
32 raw_predictor_names_12 = raw_predictor_name_subset_10_12 + raw_predictor_name_sup
33
34
35 # 2014 and beyond
36 raw_predictor_name_subset_14_beyond = [
37 "Estimate; SEX AND AGE - Total population - Female",
38 "Estimate; SEX AND AGE - 18 to 24 years",
39 "Estimate; SEX AND AGE - 25 to 34 years",
40 "Estimate; SEX AND AGE - Median age (years)",
41 "Estimate; RELATIONSHIP - Population in households - Nonrelatives - Unmarried par
42 "Estimate; HOUSEHOLDS BY TYPE - Households - Nonfamily households - Male househol
43 "Estimate; EDUCATIONAL ATTAINMENT - Bachelor's degree or higher",
44 "Estimate; FERTILITY - Women 15 to 50 years - Women 15 to 50 years who had a birt
45 "Estimate; VETERAN STATUS - Civilian population 18 years and over - Civilian vete
46 "Estimate; RESIDENCE 1 YEAR AGO - Population 1 year and over - Same house",
47 "Estimate; PLACE OF BIRTH, CITIZENSHIP STATUS AND YEAR OF ENTRY - Native",
48 "Estimate; PLACE OF BIRTH, CITIZENSHIP STATUS AND YEAR OF ENTRY - Foreign born",
49 "Estimate; WORLD REGION OF BIRTH OF FOREIGN BORN - Foreign-born population excluc
50 "Estimate; LANGUAGE SPOKEN AT HOME AND ABILITY TO SPEAK ENGLISH - Population 5 ye
51 "Estimate; EMPLOYMENT STATUS - Population 16 years and over - In labor force - Ci
52 "Estimate; COMMUTING TO WORK - Workers 16 years and over - Public transportation
53 "Estimate; HEALTH INSURANCE COVERAGE - Civilian noninstitutionalized population -
54 "Estimate; POVERTY RATES FOR FAMILIES AND PEOPLE FOR WHOM POVERTY STATUS IS DETER
55 "Estimate; OWNER CHARACTERISTICS - Owner-occupied housing units - Median value (c
56
57 raw_predictor_name_supplement_14 = [
58 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2014 INFLATION-ADJUSTED DOLLARS) - Ho
59 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2014 INFLATION-ADJUSTED DOLLARS) - Wi
60
61 raw_predictor_name_supplement_16 = [
62 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2016 INFLATION-ADJUSTED DOLLARS) - Ho
63 "Estimate; INCOME IN THE PAST 12 MONTHS (IN 2016 INFLATION-ADJUSTED DOLLARS) - Wi
64
65

```

```
In [12]: 1 renamed_predictors = [
2         "female_pct",
3         "age18_24_pct",
4         "age25_34_pct",
5         "median_age",
6         "unmarried_partner_pct",
7         "male_living_alone_pct",
8         "bachelors_deg_or_higher_pct",
9         "past_year_births_to_unmarried_women_pct",
10        "civilian_veteran_pct",
11        "live_same_house_past_year_pct",
12        "native_born_population",
13        "foreign_born_population",
14        "foreign_born_proportion_from_LatinAmerica",
15        "speak_other_language_at_home_pct",
16        "labor_force_unemployed_pct",
17        "public_transit_commuter_pct",
18        "no_health_insurance_pct",
19        "poverty_rate_pct",
20        "median_housing_value",
21        "median_household_income",
22        "food_stamp_benefits_pct"]
```

```
In [13]: 1 # check number of predictors
2         len(acs17[raw_predictor_names_17].columns)
```

Out[13]: 21

```
In [14]: 1 identifier_cols = ['state', 'district', 'year']
2
3         acs10 = acs10[['state', 'district', 'year'] + raw_predictor_names_10]
4         acs12 = acs12[['state', 'district', 'year'] + raw_predictor_names_12]
5         acs14 = acs14[['state', 'district', 'year'] + raw_predictor_names_14]
6         acs16 = acs16[['state', 'district', 'year'] + raw_predictor_names_16]
7         acs17 = acs17[['state', 'district', 'year'] + raw_predictor_names_17]
```

```
In [15]: 1 acs10.columns = identifier_cols + renamed_predictors
2         acs10.head()
```

Out[15]:

	state	district	year	female_pct	age18_24_pct	age25_34_pct	median_age	unmarried_partner_pct	n
AL_01_2010	AL	1	2010	51.0	9.0	12.3	38.4	1.6	
AL_02_2010	AL	2	2010	52.1	10.1	12.6	37.7	1.7	
AL_03_2010	AL	3	2010	51.7	12.5	12.1	36.7	1.6	
AL_04_2010	AL	4	2010	50.6	8.4	11.3	40.2	1.2	
AL_05_2010	AL	5	2010	51.0	9.8	12.4	38.3	1.2	

5 rows × 24 columns

```
In [16]: 1 # convert foreign / native born population to ratio
2         for df in [acs10, acs12, acs14, acs16, acs17]:
3             df.columns = identifier_cols + renamed_predictors
4             df['foreign_to_native_born_ratio'] = df['foreign_born_population'] / df['native_born_population']
5             df.drop(['native_born_population', 'foreign_born_population'], axis = 1, inplace = True)
```

```
In [17]: 1 # use 2017 demographics as 2018 data
2 acs18 = acs17
3 acs18.year = 2018
4 acs18.index = ['{0}_{1:02d}_{2}'.format(row['state'], row['district'], row['year']
```

```
In [18]: 1 acs18
```

```
Out[18]:
```

	state	district	year	female_pct	age18_24_pct	age25_34_pct	median_age	unmarried_partner_pct	ma
AL_01_2018	AL	1	2018	51.8	8.2	12.9	40.0	1.4	
AL_02_2018	AL	2	2018	51.2	9.7	13.0	38.5	1.7	
AL_03_2018	AL	3	2018	51.2	10.7	13.0	38.1	2.1	
AL_04_2018	AL	4	2018	51.4	8.8	11.7	40.7	1.3	
AL_05_2018	AL	5	2018	50.8	9.3	13.2	39.5	1.6	
AL_06_2018	AL	6	2018	52.3	8.6	12.3	39.1	1.3	
AL_07_2018	AL	7	2018	52.3	12.5	14.2	36.3	1.6	
AK_01_2018	AK	1	2018	47.9	9.9	16.0	34.5	2.7	
AZ_01_2018	AZ	1	2018	49.5	10.7	12.2	37.6	2.3	
AZ_02_2018	AZ	2	2018	50.5	9.5	12.8	41.4	3.0	
AZ_03_2018	AZ	3	2018	50.0	13.1	14.1	31.8	2.6	

```
In [19]: 1 # stack 5 dataframes from each year on top of each other
2 demographics_data = acs18.append([acs16, acs14, acs12, acs10])
```

```
In [20]: 1 demographics_data
```

```
Out[20]:
```

	state	district	year	female_pct	age18_24_pct	age25_34_pct	median_age	unmarried_partner_pct	ma
AL_01_2018	AL	1	2018	51.8	8.2	12.9	40.0	1.4	
AL_02_2018	AL	2	2018	51.2	9.7	13.0	38.5	1.7	
AL_03_2018	AL	3	2018	51.2	10.7	13.0	38.1	2.1	
AL_04_2018	AL	4	2018	51.4	8.8	11.7	40.7	1.3	
AL_05_2018	AL	5	2018	50.8	9.3	13.2	39.5	1.6	
AL_06_2018	AL	6	2018	52.3	8.6	12.3	39.1	1.3	
AL_07_2018	AL	7	2018	52.3	12.5	14.2	36.3	1.6	
AK_01_2018	AK	1	2018	47.9	9.9	16.0	34.5	2.7	
AZ_01_2018	AZ	1	2018	49.5	10.7	12.2	37.6	2.3	
AZ_02_2018	AZ	2	2018	50.5	9.5	12.8	41.4	3.0	
AZ_03_2018	AZ	3	2018	50.0	13.1	14.1	31.8	2.6	

```
In [21]: 1 demographics_data.to_csv('Datasets/demographics_data_2010_to_2018.csv')
```

```
In [22]: 1 demographics_data.columns
```

```
Out[22]: Index(['state', 'district', 'year', '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'],  
              dtype='object')
```

```
In [23]: 1 fec = pickle.load(open('Datasets/data_FEC_NATIONALPOLL_2004_2018.p', 'rb'))
```

```
In [24]: 1 fec.index
```

```
Out[24]: Index(['AK_01_2004', 'AL_01_2004', 'AL_02_2004', 'AL_03_2004', 'AL_04_2004',  
              'AL_05_2004', 'AL_06_2004', 'AL_07_2004', 'AR_01_2004', 'AR_02_2004',  
              ...  
              'WV_03_2018', 'WI_01_2018', 'WI_02_2018', 'WI_03_2018', 'WI_04_2018',  
              'WI_05_2018', 'WI_06_2018', 'WI_07_2018', 'WI_08_2018', 'WY_01_2018'],  
              dtype='object', length=3431)
```

```
In [25]: 1 demographics_data = demographics_data.drop(['district', 'state', 'year'], axis=1)
```

```
In [26]: 1 # try joining FEC election data with ACS demographic data  
2  
3  
4 fec_demographics_data = fec.join(demographics_data, how='inner')
```

```
In [27]: 1 set(fec_demographics_data['year'])
```

```
Out[27]: {2010, 2012, 2014, 2016, 2018}
```

```
In [28]: 1 # pickle.dump(fec_demographics_data, open('Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICS_2010_2018.p', 'wb'))  
2 # fec_demographics_data.to_csv('Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICS_2010_2018.csv')  
3 # test = pickle.load(open('Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICS_2010_2018.p', 'rb'))  
4 # fec_demographics_data.shape, test.shape
```

```
In [29]: 1 fec.shape, demographics_data.shape
```

```
Out[29]: ((3431, 22), (2180, 20))
```

```
In [30]: 1 fec_demographics_imputed_data = fec.join(demographics_data, how='outer')  
2  
3 for index, row in fec_demographics_imputed_data.iterrows():  
4     if index not in fec.index:  
5         fec_demographics_imputed_data = fec_demographics_imputed_data[fec_demographics_imputed_data.index == index]  
6  
7 fec_demographics_imputed_data.shape
```

```
Out[30]: (3431, 42)
```

```
In [31]: 1 assert(sorted(fec_demographics_imputed_data.index) == sorted(fec.index))
```

```
In [32]: 1 demographics_data.columns
```

```
Out[32]: Index(['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'],
              dtype='object')
```

```
In [33]: 1 for index, row in fec_demographics_imputed_data.iterrows():
          2     if row['year'] < 2010:
          3         impute_index = index[:-5] + '_2010'
          4         fec_demographics_imputed_data.loc[fec_demographics_imputed_data.index ==
```

```
In [34]: 1 fec_demographics_imputed_data[fec_demographics_imputed_data.index == 'AK_01_2010']
```

```
Out[34]:
```

	district	state	year	party	candidatevotes	totalvotes	candidate	national_poll	national_poll_p
AK_01_2010	1.0	AK	2010.0	republican	175384	254335	Don Young	-2.622642	9.824

1 rows x 42 columns

```
In [35]: 1 fec_demographics_imputed_data[fec_demographics_imputed_data.index == 'AK_01_2004']
```

```
Out[35]:
```

	district	state	year	party	candidatevotes	totalvotes	candidate	national_poll	national_poll_p
AK_01_2004	1.0	AK	2004.0	republican	213216	299996	Don Young	3.680556	-0.989

1 rows x 42 columns

```
In [36]: 1 for col in demographics_data.columns:
          2     if fec_demographics_imputed_data[col].isnull().values.any():
          3         raise ValueError
```

```
In [37]: 1 set(fec_demographics_imputed_data['year'].values)
```

```
Out[37]: {2004.0, 2006.0, 2008.0, 2010.0, 2012.0, 2014.0, 2016.0, 2018.0}
```

```
In [38]: 1 pickle.dump(fec_demographics_imputed_data, open('Datasets/data_FEC_NATIONALPOLL_I
          2 fec_demographics_imputed_data.to_csv('Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICS
          3 test = pickle.load(open('Datasets/data_FEC_NATIONALPOLL_DEMOGRAPHICSIMPUTED_2004_
          4 fec_demographics_imputed_data.shape, test.shape)
```

```
Out[38]: ((3431, 42), (3431, 42))
```

```
In [39]: 1 set(test['year'].values)
```

```
Out[39]: {2004.0, 2006.0, 2008.0, 2010.0, 2012.0, 2014.0, 2016.0, 2018.0}
```

```
In [40]: 1 for col in demographics_data.columns:
          2     if test[col].isnull().values.any():
          3         raise ValueError
```



```
In [ ]: 1 fec[fec.index == 'WI_06_2018']
```

```
In [ ]: 1 demographics_data[demographics_data.index == 'WI_06_2018']
```

```
In [ ]: 1 fec_demographics_data[fec_demographics_data.index == 'WI_06_2018']
```

```
In [ ]: 1 test.columns
```

```
In [ ]: 1 fec.shape
```

```
In [ ]: 1 demographics_data.shape
```

```
In [ ]: 1 demographics_data.columns
```

```
In [ ]: 1
```