

# Capstone-AdaReg-Report

## Prediction of United States' Counties Poverty Rates¶¶

Mark Vervuurt  
19-01-2018

### Executive Summary¶¶

This document presents the results of the regression analysis to predict poverty rates of United States' Counties. The result of this regression analysis is the creation of an AdaBoostRegressor. This AdaBoostRegressor is able to predict United States' Counties poverty rates with an RMSE of 2.7853.

In the data understanding phase was discovered that the following features play a significant role in predicting United States' counties poverty rates. They have moderate positive and negative Pearson correlation coefficients. Furthermore the Boxplots show enough variation and separation of the data with respect to the target variable poverty\_rate. However recursive function elimination was ultimately used to select the optimum number of features for the lowest RMSE score.

Significant Features	Short Description
area_rucc	Rural urban continuum code of county
econ_economic_typology	economic dependence type of county
au_i_pct65y_cat	created categorical feature combining area_urban_influence and categorical percentage of 65 years old per county
demo_pct_adults_less_than_a_high_school_diploma	percentage of adults with less than high school diploma per county
health_homicides_per_100k	homicides per 100k inhabitants per county
econ_pct_unemployment	percentage of unemployment per county
health_pct_low_birthweight	percentage of low birth weight per county
econ_pct_uninsured_adults	percentage of uninsured adults per county
health_pct_diabetes	percentage of diabetes per county
demo_pct_non_hispanic_african_american	percentage of African Americans per

econ\_\_pct\_civilian\_labor

county

percentage of civilian labor per  
county

The CRISP-DM Methodology was used in order to create an accurate regression model:

- **Business Understanding:** read through the '[Rural Poverty & Well-being](#)' report to better understand the circumstances of poverty.
- **Data Understanding:** explore the quantitative and categorical variables that play a key role in predicting poverty rates. Create new, better and informative features.
- **Data Preparation:** drop redundant and uninformative features, fill missing values, etc.
- **Modeling:** create and select the best regression model.
- **Evaluation:** evaluate the regression models using nested cross validation.
- **Deployment:** the deployment of the regression model is not strictly applicable here. However presenting the results of the regression analysis with this report can be considered as the deployment step.

## Business Understanding¶

As described in the online report the '[Rural Poverty & Well-being](#)': "Concentrated poverty contributes to poor housing and health conditions, higher crime and school dropout rates, as well as employment dislocations". With this information the data will be explored to see how health, crime, education and employment related factors contribute to poverty.

Another important feature of poverty is time. An area that doesn't have a high level of poverty in two following years is likely better off than an area that has a high level of poverty in both years. It will not be possible to construct a feature with this information because we cannot compare the state's poverty rate over year 'a' and 'b' within this data set. We don't have a unique key to identify counties.

Counties are generally compared by their Non-Metro and Metro status. There is more poverty in Non-Metro areas than Metro areas. Poverty is also higher under certain ages and ethnicities. Here also the data will be explored on the basis of this information.

In [72]:

```
import re
import bs4
import time
import plyfile
import html5lib
import multiprocessing
import itertools

import numpy as np
import pandas as pd
```

```
import seaborn as sns
from scipy import misc
import scipy.io.wavfile as wavfile
```

```
import scipy
from math import sqrt
from scipy import stats
from pprint import pprint
from sklearn import tree
from sklearn.svm import SVC
from sklearn import manifold
from tempfile import mkdtemp
from textwrap import wrap
from matplotlib import cm as cm
```

```
import sklearn.metrics as metrics
from pandas.plotting import scatter_matrix
from scipy.stats import randint as sp_randint
from sklearn.pipeline import TransformerMixin
from sklearn.metrics.scorer import make_scorer
from sklearn.pipeline import Pipeline
from sklearn.decomposition import PCA
from sklearn.datasets import load_iris
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import Binarizer
from sklearn.naive_bayes import MultinomialNB
from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.feature_selection import RFECV, SelectFromModel, f_regression,
SelectKBest
from sklearn.ensemble import RandomForestClassifier, AdaBoostRegressor,
AdaBoostClassifier
from sklearn.dummy import DummyClassifier, DummyRegressor
from sklearn.cluster import AgglomerativeClustering, KMeans
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.linear_model import LinearRegression, LassoCV, RidgeCV, Lasso,
Ridge
from sklearn.preprocessing import MaxAbsScaler, MinMaxScaler, Normalizer,
RobustScaler, StandardScaler
from sklearn.metrics import recall_score, accuracy_score, confusion_matrix,
roc_curve, roc_auc_score, mean_squared_error, accuracy_score
from sklearn.model_selection import train_test_split, GridSearchCV,
RandomizedSearchCV, cross_val_score, cross_validate, cross_val_predict,
KFold, ShuffleSplit, StratifiedShuffleSplit
```

```
import matplotlib
import matplotlib.pyplot as plt
from matplotlib import cm as cm
from mpl_toolkits.mplot3d import Axes3D
from pandas.plotting import parallel_coordinates, andrews_curves
```

```
%matplotlib inline
matplotlib.style.use('ggplot')

pd.set_option('display.max_columns', None)
```

## Data Understanding¶

In order to build this regression model and determine its most significant features a thorough data exploration was done to understand the relationship between poverty rates and other features.

### Initial Data Exploration¶

The dataset consists of 3198 records about United States' counties. Each record contains socioeconomic indicators about a United States' county for a given year. Besides the 'row\_id', 'yr' and the target value 'poverty\_rate', the dataset contains 32 features about socioeconomic indicators.

```
In [2]:
poverty_train = pd.read_csv('./Microsoft_
_DAT102x_Predicting_Poverty_in_the_United_States_-_Training_values.csv')
```

```
In [3]:
train_shape_tmp = poverty_train.shape
```

```
In [4]:
train_dtypes_tmp = poverty_train.dtypes
```

### Individual Feature Statistics¶

Here are the summary statistics for all the socioeconomic features:

- summary statistics of categorical variables: the total count (count), number of unique elements (unique), most frequent element (top) and the frequency of the most frequent element (frequent)
- summary statistics of quantitative variables: the mean, the standard deviation (std), the minimum value (min), 25% percentile, 50% percentile (median), 75% percentile and the maximum value (max).

```
In [73]:
poverty_train.drop(columns=['row_id'], axis=1).describe(include='all').T
```

Out[73]:

	co u nt	un iq ue	fr e q	mea n	std	min	25%	50 %	75 %	ma x
area_rucc	31	9	Nonm etro -	60	NaN	NaN	NaN	NaN	NaN	NaN

	9 8		Urban popul ation of 2,500 to 19,99 9...	8							
area_urban_influence	3 1 9 8	12	Small- in a metro area with fewer than 1 millio. ..	6 9 2	NaN	NaN	NaN	NaN	NaN	NaN	NaN
econ_economic_typology	3 1 9 8	6	Nonsp ecializ ed	1 2 6 6	NaN	NaN	NaN	NaN	NaN	NaN	NaN
econ_pct_civilian_labor	3 1 9 8	NaN	NaN	NaN	0.46 707 1	0.07 454 1	0.21 7	0.42	0.4 67	0.5 14	1
econ_pct_unemployment	3 1 9 8	NaN	NaN	NaN	0.05 961 04	0.02 284 97	0.00 8	0.04 4	0.0 57	0.0 71	0.2 4
econ_pct_uninsured_adults	3 1 9 6	NaN	NaN	NaN	0.21 753 4	0.06 737 18	0.04 6	0.16 6	0.2 16	0.2 62	0.4 95
econ_pct_uninsured_children	3 1 9 6	NaN	NaN	NaN	0.08 592 02	0.04 000 46	0.00 9	0.05 7	0.0 77	0.1 05	0.2 85
demo_pct_female	3 1 9 6	NaN	NaN	NaN	0.49 878 1	0.02 425 08	0.29 4	0.49 3	0.5 03	0.5 12	0.5 76
demo_pct_below_18_years_of_age	3 1 9	NaN	NaN	NaN	0.22 776	0.03 429	0.09 8	0.20 7	0.2 26	0.2 452	0.4 17

	6			N	3	09				5	
demo_pct_aged_65_years_and_older	3	Na	NaN	N	0.17	0.04	0.04	0.14	0.1	0.1	0.3
	1	N		a	013	359	3	2	67	94	55
	9			N	7	37					
	6										
demo_pct_hispanic	3	Na	NaN	N	0.09	0.14	0	0.01	0.0	0.0	0.9
	1	N		a	023	270		9	35	88	45
	9			N	34	7					
	6										
demo_pct_non_hispanic_african_american	3	Na	NaN	N	0.09	0.14	0	0.00	0.0	0.0	0.8
	1	N		a	111	710		6	22	962	55
	9			N	67	4				5	
	6										
demo_pct_non_hispanic_white	3	Na	NaN	N	0.77	0.20	0.06	0.64	0.8	0.9	0.9
	1	N		a	020	790		8	54	36	98
	9			N	7	3					
	6										
demo_pct_american_indian_or_alaskan_native	3	Na	NaN	N	0.02	0.08	0	0.00	0.0	0.0	0.8
	1	N		a	465	463		2	07	14	52
	9			N	86	41					
	6										
demo_pct_asian	3	Na	NaN	N	0.01	0.02	0	0.00	0.0	0.0	0.3
	1	N		a	330	536		3	07	13	46
	9			N	35	56					
	6										
demo_pct_adults_less_than_a_high_school_diploma	3	Na	NaN	N	0.14	0.06	0.01	0.09	0.1	0.1	0.4
	1	N		a	879	825	612	746	335	951	668
	9			N	4	47	9	83	01	71	67
	8										
demo_pct_adults_with_high_school_diploma	3	Na	NaN	N	0.35	0.07	0.07	0.30	0.3	0.3	0.5
	1	N		a	03	053	282	591	557	991	516
	9			N		42	05	5	01	97	89
	8										
demo_pct_adults_with_some_college	3	Na	NaN	N	0.30	0.05	0.11	0.26	0.3	0.3	0.4
	1	N		a	136	249	282	536	015	359	742
	9			N	6	76	1	2	95	72	16
	8										
demo_pct_adults_bachelors_or_higher	3	Na	NaN	N	0.19	0.08	0.01	0.13	0.1	0.2	0.7
	1	N		a	954	915	398	884	772	332	948
	9			N		77	6		47	58	72
	8										
demo_birth_rate_per_1k	3	Na	NaN	N	11.6	2.73	4	10	11	13	29
	1	N		a	77	952					

	9			N								
	8											
demo_death_rate_per_1k	3198	Na	NaN	NaN	10.3011	2.78614	0	8	10	12	27	
health_pct_adult_obesity	3196	Na	NaN	NaN	0.30759	0.0434	0.14	0.284	0.309	0.334	0.484	
health_pct_adult_smoking	2734	Na	NaN	NaN	0.21351	0.06309	0.05	0.171	0.211	0.2497	0.526	
health_pct_diabetes	3196	Na	NaN	NaN	0.10928	0.02319	0.033	0.094	0.109	0.124	0.197	
health_pct_low_birthweight	316	Na	NaN	NaN	0.08353	0.02238	0.025	0.068	0.08	0.095	0.232	
health_pct_excessive_drinking	220	Na	NaN	NaN	0.16483	0.05023	0.038	0.129	0.164	0.196	0.358	
health_pct_physical_inactivity	3196	Na	NaN	NaN	0.27730	0.05294	0.097	0.243	0.28	0.313	0.443	
health_air_pollution_particulate_matter	3170	Na	NaN	NaN	11.6265	1.54493	7	10	12	13	15	
health_homicides_per_100k	1231	Na	NaN	NaN	5.95075	5.06337	-0.39	2.66	4.84	7.825	51.49	
health_motor_vehicle_crash_deaths_per_100k	2781	Na	NaN	NaN	21.161	10.517	3.09	13.46	19.63	26.47	110.45	
health_pop_per_denti	2	Na	NaN	NaN	343	256	339	181	269	408	281	

st	9	N		a	1.44	9.44		2.25	0	9.7	29
	5			N						5	
	4										
health_pop_per_prim	2	Na	NaN	N	255	210	189	141	199	285	234
ary_care_physician	9	N		a	1.35	0.48		9	9	9	00
	6			N							
	8										
yr	3	2	b	1	NaN	NaN	NaN	NaN	Na	Na	Na
	1			5					N	N	N
	9			9							
	8			9							

```
In [6]:
poverty_labels = pd.read_csv('Microsoft_
_DAT102x_Predicting_Poverty_in_the_United_States_-_Training_labels.csv')
```

```
In [7]:
lbl_shape_tmp = poverty_labels.shape
```

```
In [8]:
lbl_dtype_tmp = poverty_labels.dtypes
```

Here are the summary statistics for the target variable which is quantitative:

```
In [74]:
poverty_labels.drop(columns=['row_id'], axis=1).describe().T
```

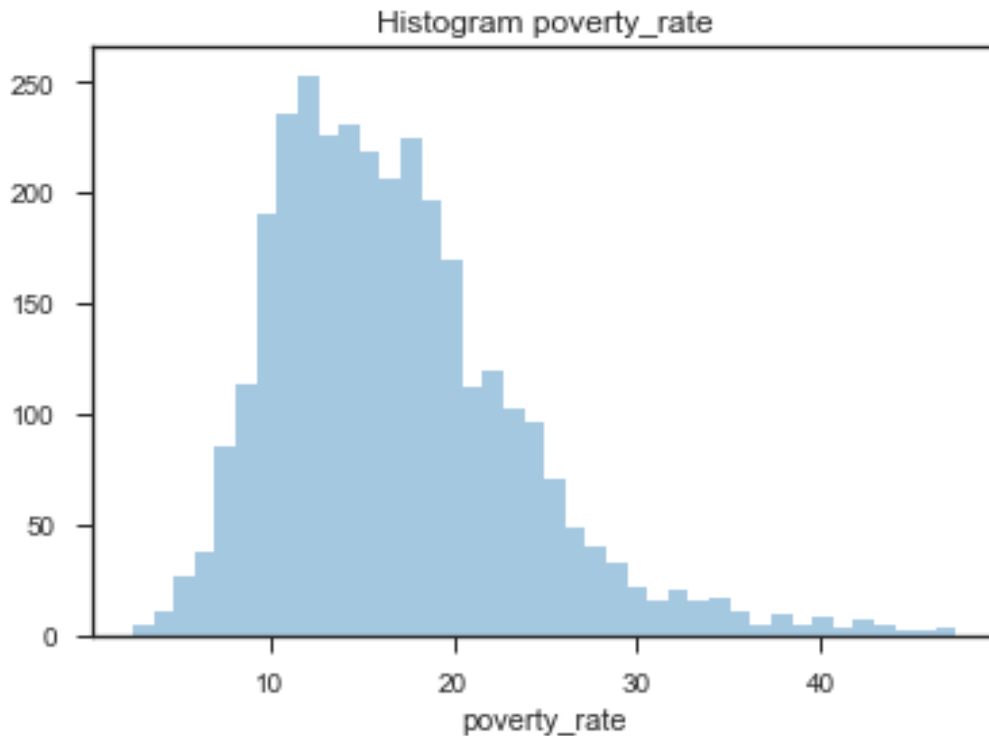
```
Out[74]:
```

	count	mean	std	min	25%	50%	75%	max
poverty_rate	3198.0	16.817136	6.697969	2.5	12.0	15.8	20.3	47.4

Poverty rates are right or positively skewed with a skew value of 1.048357. We can recognize a slight bell curve in the data. The mean and median are relatively close to each other and the standard deviation is relatively low which indicates low variability in the poverty rates. Most United States' counties have a poverty\_rate between 10% and 20% poverty.

```
In [150]:
ht_pov = sns.distplot(poverty_labels.drop(columns=['row_id'], axis=1),
kde=False, norm_hist=False, color='#1f77b4',
axlabel='poverty_rate').set_title('Histogram poverty_rate')
```





In [11]:

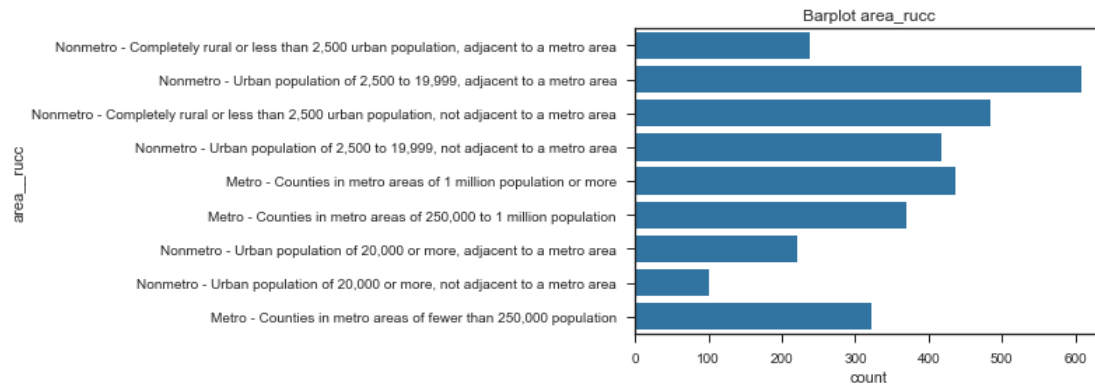
```
poverty = pd.merge(poverty_train, poverty_labels, on='row_id')
pov_shape_tmp = poverty.shape
```

From the summary statistics above, should be clear that there are three categorical variables included in the dataset:

- area\_rucc with 9 values:
  - 'Nonmetro - Urban population of 2,500 to 19,999, adjacent to a metro area' counties are most frequent with 608 counties.
  - 'Nonmetro - Urban population of 20,000 or more, not adjacent to a metro area' counties are most infrequent with 100 counties.
- area\_urban\_influence with 12 values:
  - 'Small-in a metro area with fewer than 1 million residents' counties are most frequent with 692 counties.
  - 'Noncore not adjacent to a metro/micro area and contains a town of 2,500 or more residents' counties are most infrequent with 122.
- econ\_economic\_typology with 6 values:
  - 'Non specialized' economic typology counties are most frequent with 1266 counties.
  - 'Mining-dependent' economic typology counties are most infrequent with 254 counties.

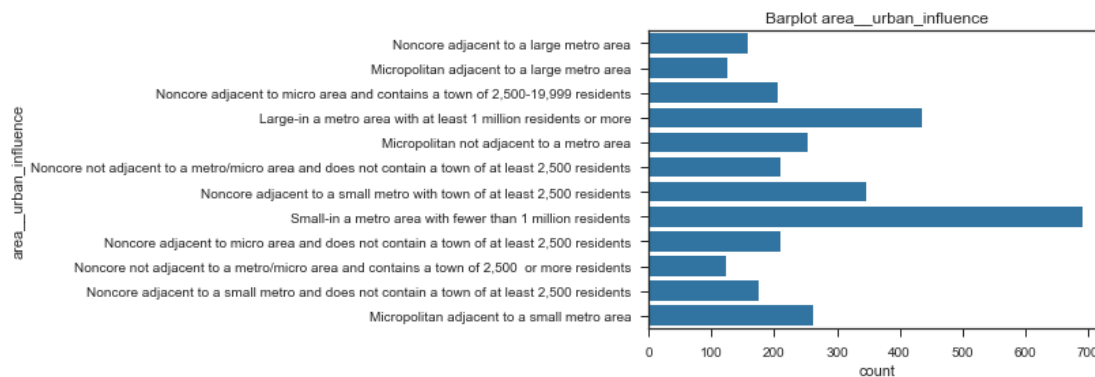
In [130]:

```
bh_ar = sns.countplot(y='area_rucc', data=poverty,
color='#1f77b4').set_title("Barplot area_rucc")
```



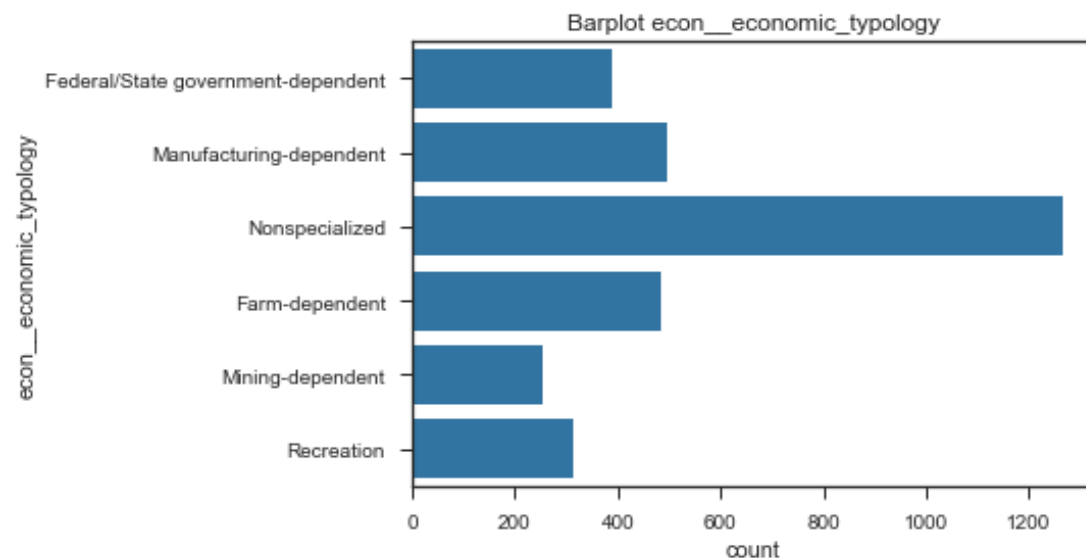
In [129]:

```
bh_aui = sns.countplot(y='area__urban_influence', data=poverty,
color='#1f77b4').set_title("Barplot area__urban_influence")
```



In [153]:

```
bh_eet = sns.countplot(y='econ__economic_typology', data=poverty,
color='#1f77b4').set_title("Barplot econ__economic_typology")
```



## Data Exploration and Visualization of Categorical Variables¶

Here the predictive value of the categorical variables 'econ\_\_economic\_typology', 'area\_\_urban\_influence', 'area\_\_rucc' and 'yr' is explored. Box plots are used to explore these categorical variables.

The boxplots of the categorical variables "econeconomic\_typology", "areaurban\_influence" and "area\_\_rucc" show interesting variation:

- 'Farm-dependent' counties have the lowest poverty rates and 'Federal/State government-dependent' counties have the highest poverty rates.
- 'Large-in a metro area with at least 1 million residents or more' counties have the lowest poverty rates.
- 'Metro - Counties in metro areas with 1 million population or more' counties have the lowest poverty rates.

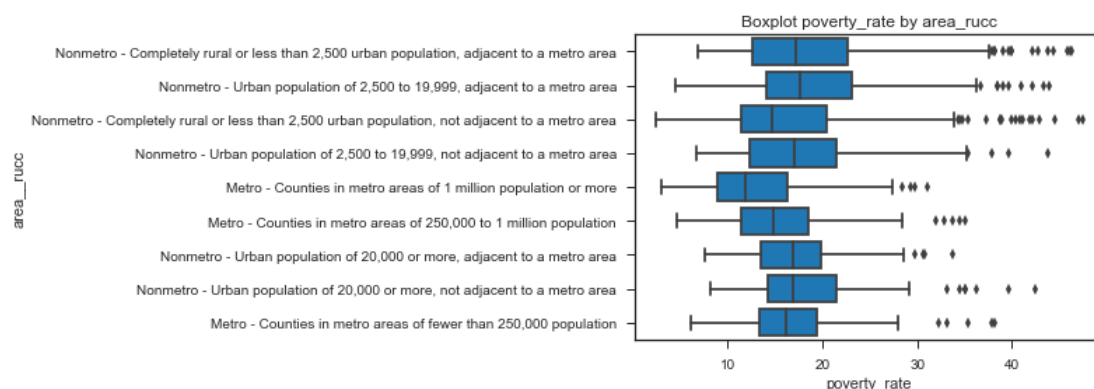
Furthermore by combining features more interesting categorical variables can be created explaining much more of the variance in poverty rates.

- "demopct\_aged\_65\_years\_and\_older" and "areaurban\_influence". The general trend is that counties with a low percentage population of "65 years or older" have a higher poverty rate.

The difference in poverty over year 'a' and 'b' is really minimal. Furthermore it doesn't make sense to use this feature to predict poverty rates. This feature will be dropped at the cleaning stage.

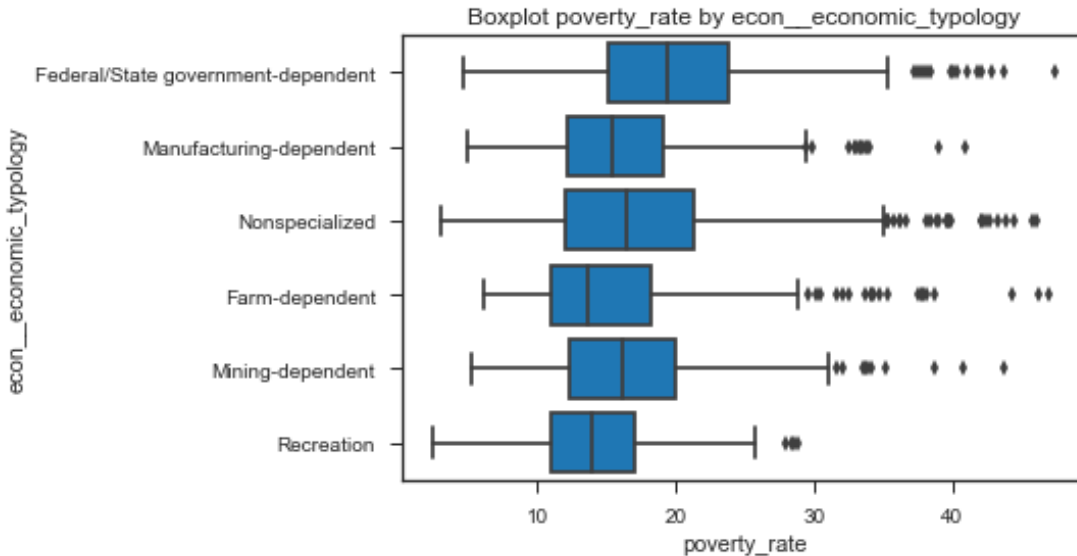
In [144]:

```
bpd_ar = sns.boxplot(orient="h", x='poverty_rate', y='area__rucc',
data=poverty, color='#1f77b4', saturation=1.0).set_title('Boxplot
poverty_rate by area__rucc')
```



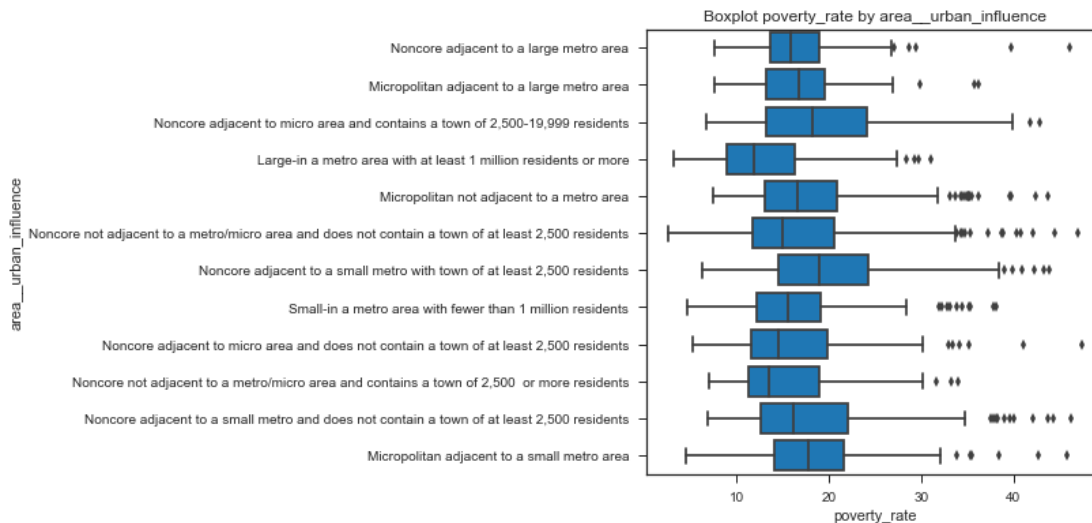
In [143]:

```
bpd_eet = sns.boxplot(orient="h", x='poverty_rate',
y='econ__economic_typology', data=poverty, color='#1f77b4',
saturation=1.0).set_title('Boxplot poverty_rate by econ__economic_typology')
```



In [142]:

```
plt.figure(figsize=(6,6))
bpd_aui = sns.boxplot(orient="h", x='poverty_rate',
y='area__urban_influence', data=poverty, color='#1f77b4',
saturation=1.0).set_title('Boxplot poverty_rate by area__urban_influence')
```



In [17]:

```
def create_old_age_cat(input_df):
    low_pct_olds = poverty.demo__pct_aged_65_years_and_older < 0.167000
    high_pct_olds = poverty.demo__pct_aged_65_years_and_older >= 0.167000
    input_df.loc[low_pct_olds,'pct_65years_cat'] = 'low_pct_65years'
    input_df.loc[high_pct_olds,'pct_65years_cat'] = 'high_pct_65years'

    age_old_cats = ['low_pct_65years','high_pct_65years']
    input_df.loc[:, 'pct_65years_cat'] =
input_df.pct_65years_cat.astype('category')
    input_df.loc[:, 'pct_65years_cat'] =
```

```
input_df.pct_65years_cat.cat.set_categories(age_old_cats, ordered=True)
return input_df
```

In [18]:

```
poverty = create_old_age_cat(poverty)
```

In [19]:

```
def create_aui_pct65y_cat(input_df):
    aui_cats = input_df.area__urban_influence.unique()
    pct65y_cats = input_df.pct_65years_cat.cat.categories

    aui_pct65y_masks = [ ((input_df.area__urban_influence == aui) &
(input_df.pct_65years_cat == pct65y)
                        , aui + ', ' + pct65y)
                        for (aui, pct65y) in list(itertools.product(aui_cats,
pct65y_cats)))]

    aui_pct65y_lbls = [aui + ', ' + pct65y for (aui, pct65y)
                        in list(itertools.product(aui_cats, pct65y_cats))]

    for mask, aui_pct65y_lb in aui_pct65y_masks:
        input_df.loc[mask, 'aui_pct65y_cat'] = aui_pct65y_lb

    input_df.loc[:, 'aui_pct65y_cat'] =
input_df.aui_pct65y_cat.astype('category')
    input_df.loc[:, 'aui_pct65y_cat'] =
input_df.aui_pct65y_cat.cat.set_categories(aui_pct65y_lbls)
    return input_df
```

In [79]:

```
poverty = create_aui_pct65y_cat(poverty)
```

## Data Exploration and Visualization of Quantitative Variables¶

For the quantitative variables the correlation matrix is computed first followed by the visual display of the scatter plot matrices.

### Correlation Matrix¶

The strongest correlations observed are moderate positive and negative for the following features:

- demo\_pct\_adults\_less\_than\_a\_high\_school\_diploma
- health\_homicides\_per\_100k
- econ\_pct\_unemployment
- health\_pct\_low\_birthweight
- econ\_pct\_uninsured\_adults
- health\_pct\_diabetes

- demo\_\_pct\_non\_hispanic\_african\_american
- econ\_\_pct\_civilian\_labor
- demo\_\_pct\_non\_hispanic\_white
- demo\_\_pct\_adults\_bachelors\_or\_higher

The whole correlation matrix of interest is shown here under.

Features	Pearson Correlation Coefficient
econ__pct_civilian_labor	-0.670417
demo__pct_non_hispanic_white	-0.499974
demo__pct_adults_bachelors_or_higher	-0.467134
demo__pct_adults_with_some_college	-0.363875
health__pct_excessive_drinking	-0.353254
demo__pct_asian	-0.163033
demo__pct_aged_65_years_and_older	-0.088123
demo__pct_female	-0.068065
demo__pct_below_18_years_of_age	0.039237
health__air_pollution_particulate_matter	0.058582
econ__pct_uninsured_children	0.098882
demo__pct_hispanic	0.105574
demo__birth_rate_per_1k	0.127506
health__pop_per_primary_care_physician	0.156942
demo__pct_adults_with_high_school_diploma	0.202928
demo__pct_american_indian_or_alaskan_native	0.236508
demo__death_rate_per_1k	0.244093
health__pop_per_dentist	0.268996
health__pct_adult_smoking	0.395457
health__motor_vehicle_crash_deaths_per_100k	0.420348
health__pct_physical_inactivity	0.437680
health__pct_adult_obesity	0.444293
demo__pct_non_hispanic_african_american	0.507048
health__pct_diabetes	0.537038
econ__pct_uninsured_adults	0.541712
health__pct_low_birthweight	0.565456
econ__pct_unemployment	0.592022
health__homicides_per_100k	0.621399
demo__pct_adults_less_than_a_high_school_diploma	0.680360

## Scatter Plot Matrices¶

After reading the '[Rural Poverty & Well-being](#)' report it is clear that education, ethnicity and health related issues play an important role in predicting poverty. In this dataset are also added economic indicators of United States' counties. The scatter plot matrices of these four groups of socioeconomic indicators are shown here under. The scatter plot matrices visually confirms the finding of the correlation matrix.

NB: linear statistical transformations (sqrt, square, exponential, etc) were also applied to the target variable 'poverty\_rate' but they did not improve substantially the correlation coefficients and the scatter plot matrices.

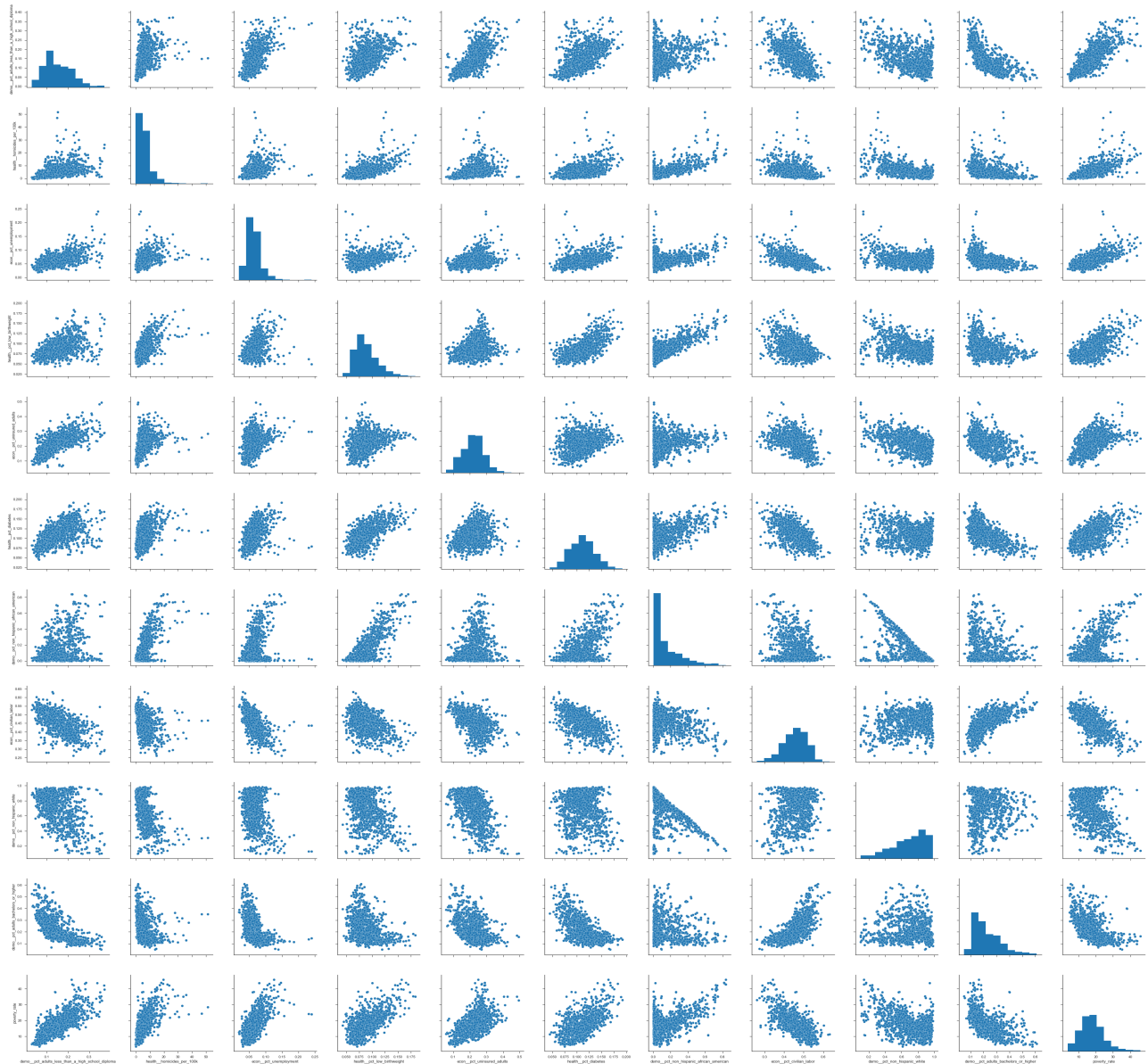
The correlation matrices and scatter plot matrices visually confirm that the variables correlate moderately strong with the target variable 'poverty\_rate' seem to have a linear relationship.

--> limit to most important or special scatterplot matrices

### *Moderately Strong Correlating Features Scatter Plot Matrices¶*

In [156]:

```
sns.set(style="ticks")
scatter_top = sns.pairplot(poverty.loc[:,
['demo__pct_adults_less_than_a_high_school_diploma'
, 'health__homicides_per_100k', 'econ__pct_unemployment'
, 'health__pct_low_birthweight', 'econ__pct_uninsured_adults'
, 'health__pct_diabetes',
'demo__pct_non_hispanic_african_american'
, 'econ__pct_civilian_labor', 'demo__pct_non_hispanic_white'
, 'demo__pct_adults_bachelors_or_higher'
, 'poverty_rate']]).dropna(),
size=4
, plot_kws={'color': '#1f77b4'},
diag_kws={'color': '#1f77b4'})
```



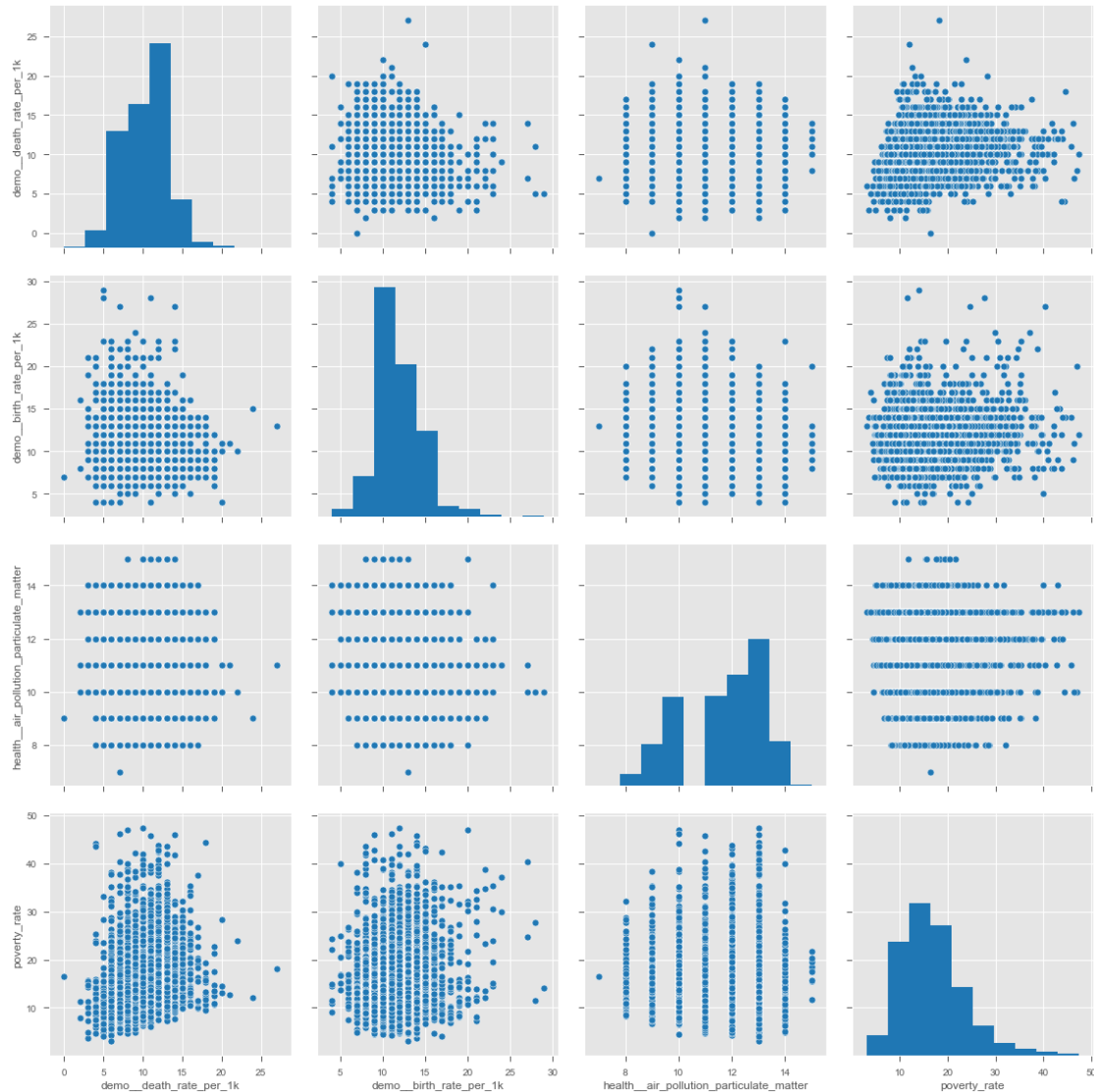
### Categorical Features Scatter Plot Matrices ¶

These quantitative features behave like categorical features.

In [124]:

```
scatter_non_ln = sns.pairplot(poverty.loc[:, ['demo_death_rate_per_1k',
                                              'demo_birth_rate_per_1k',
                                              'health_air_pollution_particulate_matter',
                                              'poverty_rate']]).dropna(),
size=4
, plot_kws={'color': '#1f77b4'},
diag_kws={'color': '#1f77b4'})
```





## Creation and Visualization of New Categorical Variables¶

From observing the scatter plot matrices of the features 'demo\_death\_rate\_per\_1k', 'demo\_birth\_rate\_per\_1k' and 'health\_air\_pollution\_particulate\_matter' is clear that these quantitative variables behave like categorical variables. They will be transformed into categorical features by binning them.

In [27]:

```
def create_birtherate_cat(input_df):
    cats = [ "birtherate {0} - {1}".format(i, i + 5) for i in range(0, 40, 5) ]

    input_df.loc[:, 'birth_rate_cat'] =
pd.cut(input_df.demo__birth_rate_per_1k, range(0, 45, 5), right=False,
include_lowest=True, labels=cats)
    input_df.loc[:, 'birth_rate_cat'] =
input_df.birth_rate_cat.astype('category')
```

```

input_df.loc[:, 'birth_rate_cat'] =
input_df.birth_rate_cat.cat.set_categories(cats, ordered=True)
return input_df

```

In [28]:

```

def create_deathrate_cat(input_df):
    cats = [ "deathrate {0} - {1}".format(i, i + 5) for i in range(0, 40, 5)
]
    input_df.loc[:, 'death_rate_cat'] =
pd.cut(input_df.demo__death_rate_per_1k, range(0, 45, 5), right=False,
include_lowest=True, labels=cats)
    input_df.loc[:, 'death_rate_cat'] =
input_df.death_rate_cat.astype('category')
    input_df.loc[:, 'death_rate_cat'] =
input_df.death_rate_cat.cat.set_categories(cats, ordered=True)
return input_df

```

In [29]:

```

def create_air_poll_cat(input_df):
    cats = [ "airpoll {0} - {1}".format(i, i + 5) for i in range(0, 35, 5) ]
    input_df.loc[:, 'air_poll_cat'] =
pd.cut(input_df.health__air_pollution_particulate_matter, range(0, 40, 5),
right=False, include_lowest=True, labels=cats)
    input_df.loc[:, 'air_poll_cat'] = input_df.air_poll_cat.astype('category')
    input_df.loc[:, 'air_poll_cat'] =
input_df.air_poll_cat.cat.set_categories(cats, ordered=True)
return input_df

```

In [30]:

```

def create_features(input_df):
    input_df = create_birthrate_cat(input_df)
    input_df = create_deathrate_cat(input_df)
    input_df = create_air_poll_cat(input_df)
    return input_df

```

In [31]:

```

poverty = create_features(poverty)

```

## Data Preparation¶

This phase involves mostly the cleaning, scaling and one hot encoding of features:

- dropping redundant features
- converting features to the right type
- missing values are replaced by the respective median value of the feature. The median is preferred over the mean because it is less sensible to skewed data and gives a better measure of centrality.
- features are scaled to have the same scale. The MinMaxScaler is applied to the features "healthhomicides\_per\_100k" and "healthmotor\_vehicle\_crash\_deaths\_per\_100k" to

scale them the same way as other quantitative variables that are in percentages between 0 and 1.

- One hot encoding of the categorical variables is performed

In [32]:

```
def drop_features(input_df):
    result_df =
input_df.drop(columns=['health__air_pollution_particulate_matter'
                        , 'demo__death_rate_per_1k',
                        'demo__birth_rate_per_1k'
                        , 'pct_65years_cat', 'au_i_pct65y_cat', 'yr'], axis=1)
    return result_df
```

In [33]:

```
poverty_clean = drop_features(poverty)
```

In [34]:

```
def convert_to_cat(input_df):
    input_df.loc[:, 'area__rucc'] = input_df.area__rucc.astype("category")
    input_df.loc[:, 'area__urban_influence'] =
input_df.area__urban_influence.astype("category")
    input_df.loc[:, 'econ__economic_typology'] =
input_df.econ__economic_typology.astype("category")
    return input_df
```

In [35]:

```
poverty_clean = convert_to_cat(poverty_clean)
```

In [36]:

```
dtypes_tmp = poverty_clean.dtypes
```

In [37]:

```
def clean_nans(input_df):
    result_df = input_df.fillna(poverty_clean.median())
    return result_df
```

In [38]:

```
poverty_clean = clean_nans(poverty_clean)
```

In [39]:

```
clean_tmp = poverty_clean.isnull().sum()
```

In [40]:

```
def scale_features(input_df):
    input_scale = input_df.loc[:, ['health__homicides_per_100k'
                                   , 'health__motor_vehicle_crash_deaths_per_100k'
                                   , 'health__pop_per_dentist'
                                   , 'health__pop_per_primary_care_physician']]
```

```
input_scaled = pd.DataFrame(MinMaxScaler().fit_transform(input_scale),
columns=input_scale.columns)
```

```
input_df.loc[:, 'health__homicides_per_100k'] =
input_scaled.loc[:, 'health__homicides_per_100k']
input_df.loc[:, 'health__motor_vehicle_crash_deaths_per_100k'] =
input_scaled.loc[:, 'health__motor_vehicle_crash_deaths_per_100k']
input_df.loc[:, 'health__pop_per_dentist'] =
input_scaled.loc[:, 'health__pop_per_dentist']
input_df.loc[:, 'health__pop_per_primary_care_physician'] =
input_scaled.loc[:, 'health__pop_per_primary_care_physician']
return input_df
```

```
In [41]:
poverty_clean = scale_features(poverty_clean)
```

```
In [42]:
def cat_to_dummies(input_df):
    result_df = pd.get_dummies(input_df, dummy_na=True,
columns=['area__rucc', 'econ__economic_typology'

, 'birth_rate_cat', 'death_rate_cat', 'air_poll_cat'

, 'area__urban_influence'])
    return result_df
```

```
In [43]:
poverty_clean = cat_to_dummies(poverty_clean)
```

```
In [44]:
shape_tmp = poverty_clean.shape
```

## Modeling and Evaluation¶

In this phase two models are compared with each other using the RMSE evaluation metric:

- Least Square Linear Model after applying recursive feature selection to create a linear model with the most important features
- An AdaBoostRegressor which is an ensemble learning model of decision trees.

The best RMSE scores obtained by these two models are:

Model	RMSE
AdaBoostRegressor	2.7847
Linear Regression	2.9297

The AdaBoostRegressor happens to be more precise than the Least Squares Linear Model because it can handle non linear relationships. The AdaBoostRegressor is chosen as the regression model to predict poverty rates for United States Counties.

In [45]:

```
rng = np.random.RandomState(0)
```

In [46]:

```
poverty_X = poverty_clean.drop(columns=['row_id', 'poverty_rate'], axis=1)
poverty_y = poverty_clean.poverty_rate
```

In [47]:

```
# use r2 adjusted in the future
scoring = {'r2': 'r2', 'mse': make_scorer(mean_squared_error,
greater_is_better=False)}
```

In [48]:

```
inner_cv = ShuffleSplit(n_splits=5, test_size=0.3, random_state=rng)
outer_cv = ShuffleSplit(n_splits=5, test_size=0.3, random_state=rng)
```

## Recursive Feature Selection Linear Regression¶

In [49]:

```
cachedir = mkdtemp()
```

```
mse = make_scorer(mean_squared_error, greater_is_better=False)
```

```
rfecv = RFECV(estimator=LinearRegression(), step=1, cv=inner_cv, scoring=mse)
```

```
rfecv.fit(poverty_X, poverty_y)
```

```
print("Optimal number of features : %d" % rfecv.n_features_)
```

```
Optimal number of features : 73
```

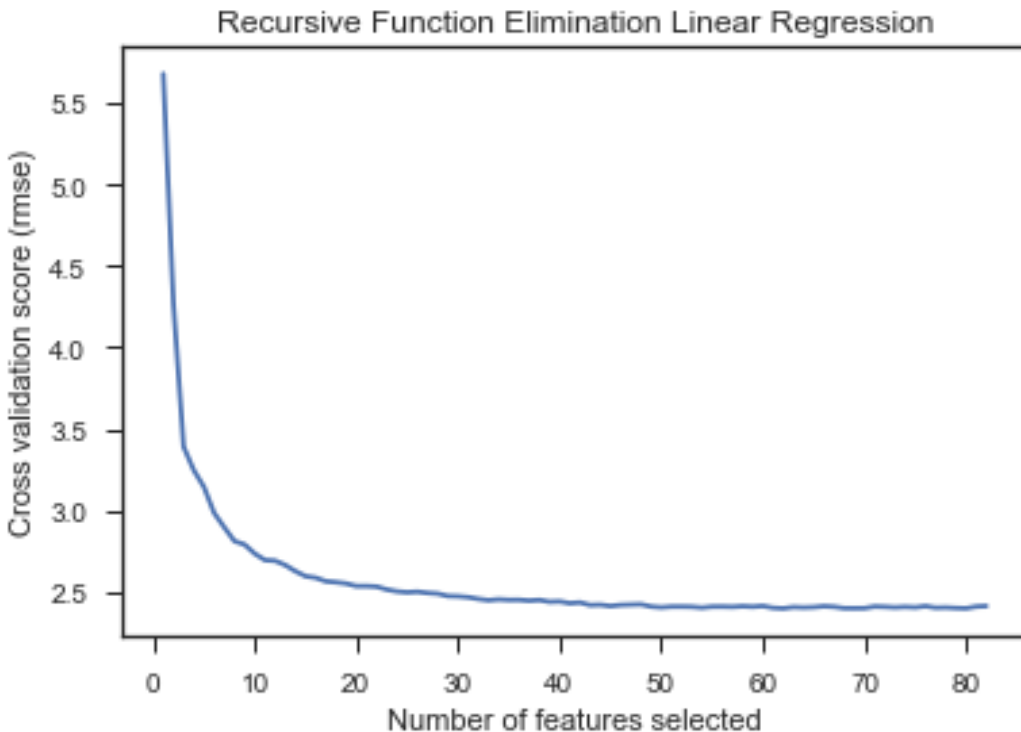
In [50]:

```
pprint('RMSE score: %f' %
np.sqrt(np.abs(rfecv.grid_scores_[rfecv.n_features_])))
```

```
'RMSE score: 3.121806'
```

In [158]:

```
plt.figure()
plt.title('Recursive Function Elimination Linear Regression')
plt.xlabel("Number of features selected")
plt.ylabel("Cross validation score (rmse)")
plt.plot(range(1, len(rfecv.grid_scores_) + 1),
np.sqrt(np.abs(rfecv.grid_scores_)))
plt.show()
```



### Nested Cross Validation AdaBoostRegressor Vs Linear Regression¶

Using nested cross validation the AdaBoostRegressor is compared with the linear regression model. The AdaBoostRegressor wins the lowest RMSE score.

In [52]:

```
cachedir = mkdtemp()
estimators = [('reg_model', LinearRegression())]
regr_pipe = Pipeline(estimators, memory=cachedir)
```

In [53]:

```
adaReg = AdaBoostRegressor(base_estimator=DecisionTreeRegressor(max_depth=13,
                                                                , random_state=rng)
                           , n_estimators=600, loss='linear',
                           learning_rate=1, random_state=rng)
```

In [54]:

```
param_grid = dict(reg_model=[rfecv, adaReg])
```

In [55]:

```
reg_grid = GridSearchCV(estimator=regr_pipe, param_grid=param_grid,
                        scoring=scoring, cv=inner_cv
                        , error_score=0, refit='mse')
reg_pred = reg_grid.fit(poverty_X, poverty_y)
```

Out[55]:

```
Pipeline(memory='/var/folders/_j/vyb4dyfx2wq850vj9vh25wy40000gn/T/tmpe8qvy037',
        steps=[('reg_model',
AdaBoostRegressor(base_estimator=DecisionTreeRegressor(criterion='mse',
max_depth=13, max_features=None,
                    max_leaf_nodes=None, min_impurity_decrease=0.0,
                    min_impurity_split=None, min_samples_leaf=1,
                    min_samples_split=2, min_weight_fraction_leaf=0....oss='linear',
n_estimators=600,
                    random_state=<mtrand.RandomState object at 0x1a2f38ecf0>))])])
```

```
In [56]:
train_scores = cross_validate(reg_grid, poverty_X, poverty_y, cv=outer_cv,
scoring=scoring, return_train_score=True)
pprint('RMSE score: %f' %
np.average(np.sqrt(np.abs(train_scores['test_mse']))))

'RMSE score: 2.428183'
```

## Recursive Feature Selection AdaBoostRegressor¶

To improve the AdaBoostRegressor even more, the best features are selected using recursive feature elimination or backwards elimination.

```
In [57]:
mse = make_scorer(mean_squared_error, greater_is_better=False)

rfecv = RFECV(estimator=adaReg, step=1, cv=inner_cv, scoring=mse)

rfecv.fit(poverty_X, poverty_y)

print("Optimal number of features : %d" % rfecv.n_features_)

Optimal number of features : 62
```

```
In [58]:
pprint('RMSE score: %f' %
np.sqrt(np.abs(rfecv.grid_scores_[rfecv.n_features_])))

'RMSE score: 2.406639'
```

```
In [59]:
train_scores = cross_validate(rfecv, poverty_X, poverty_y, cv=outer_cv,
scoring=scoring, return_train_score=True)
pprint('RMSE score: %f' %
np.average(np.sqrt(np.abs(train_scores['test_mse']))))

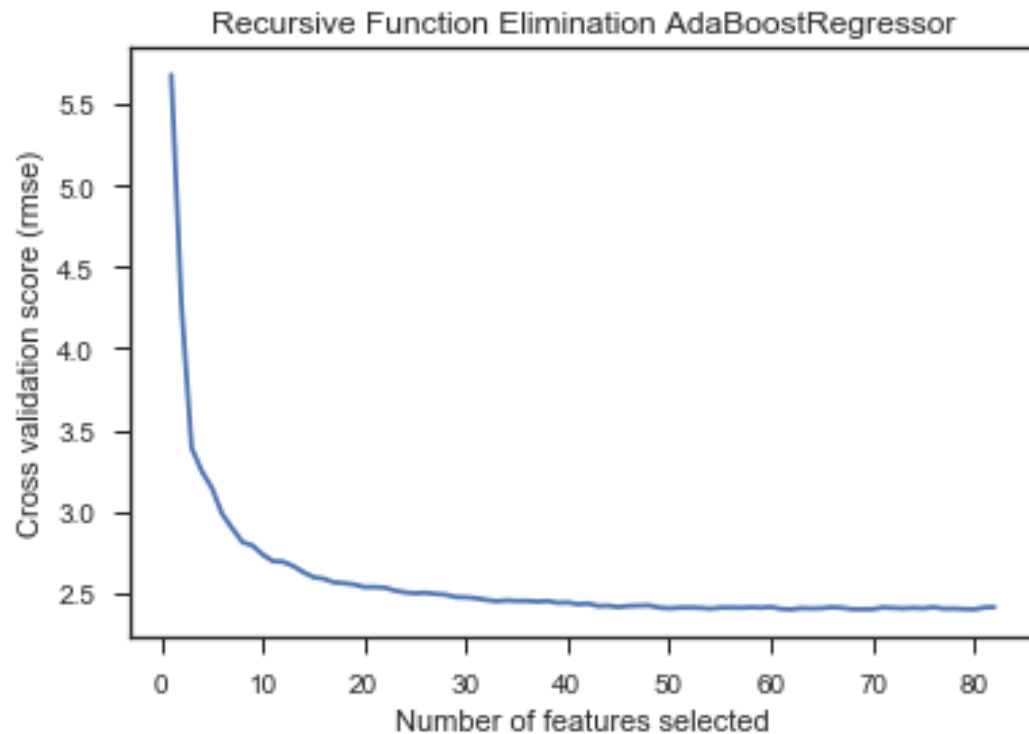
'RMSE score: 2.366428'
```

```
In [60]:
```

```

plt.figure()
plt.title('Recursive Function Elimination AdaBoostRegressor')
plt.xlabel("Number of features selected")
plt.ylabel("Cross validation score (rmse)")
plt.plot(range(1, len(rfecv.grid_scores_) + 1),
np.sqrt(np.abs(rfecv.grid_scores_)))
plt.show()

```



## Analysis of Predictions and Residuals¶

The Analysis of the quality of the predictions and residuals shows that the accuracy of the AdaBoostRegressor is high. However the AdaBoostRegressor probably slightly overfits the data.

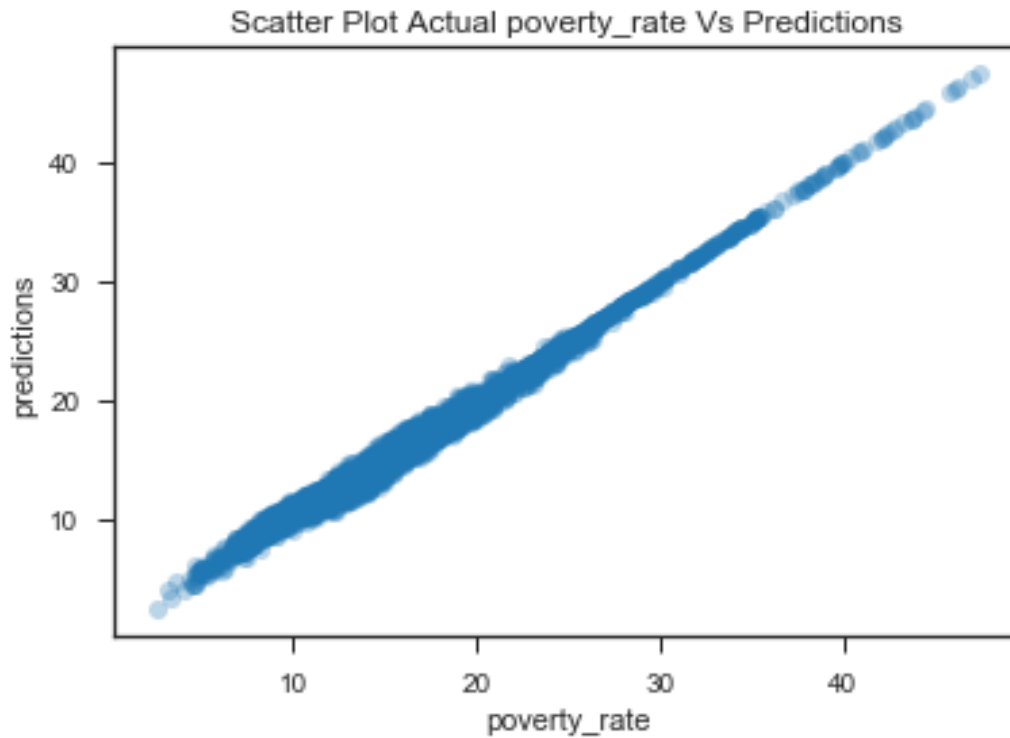
In [159]:

```

sc1_tmp = sns.regplot(x=poverty_y, y=rfecv.predict(poverty_X), fit_reg=False,
color='#1f77b4', scatter_kws={'alpha':0.3})
tmp = sc1_tmp.set_ylabel('predictions')
tmp = sc1_tmp.set_title('Scatter Plot Actual poverty_rate Vs Predictions')

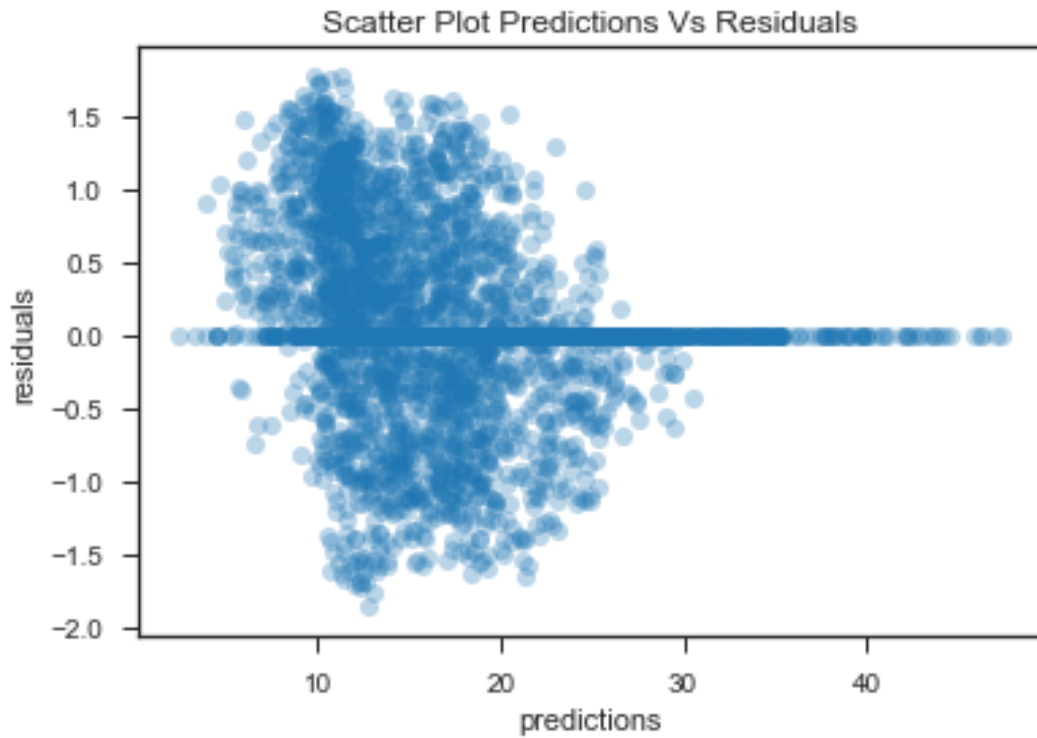
```



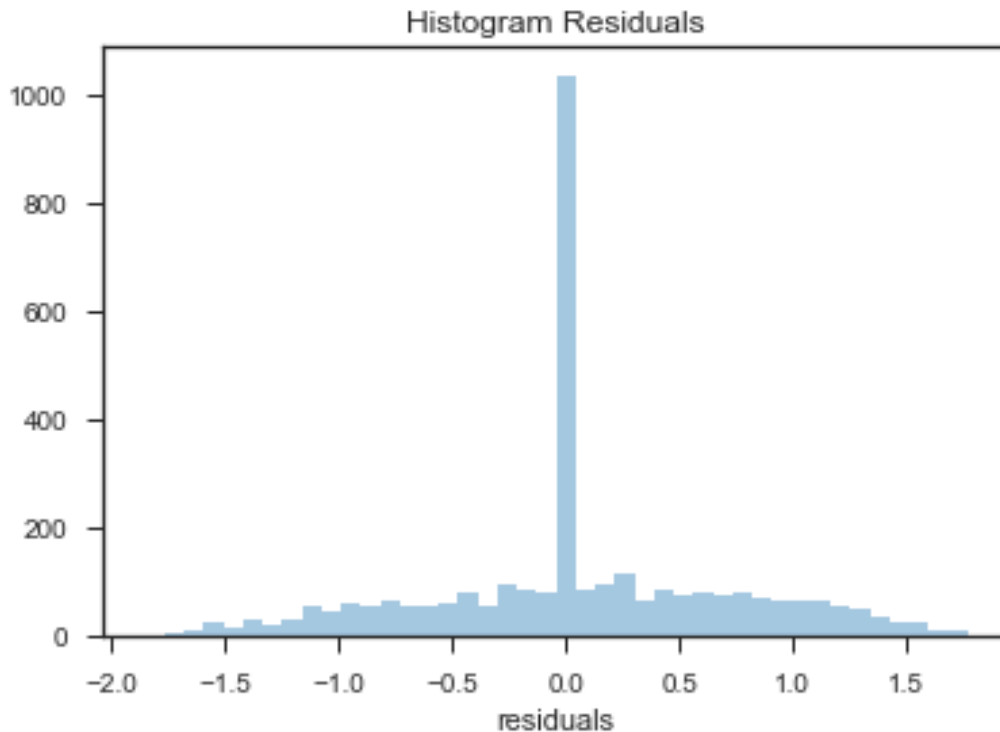


In [160]:

```
sc2_tmp = sns.regplot(x=rfev.predict(poverty_X), y=rfev.predict(poverty_X)
- poverty_y
                        , fit_reg=False, color='#1f77b4',
scatter_kws={'alpha':0.3})
tmp = sc2_tmp.set_xlabel('predictions')
tmp = sc2_tmp.set_ylabel('residuals')
tmp = sc2_tmp.set_title('Scatter Plot Predictions Vs Residuals')
```



```
In [161]:  
residuals = rfecv.predict(poverty_X) - poverty_y  
residuals = residuals.rename('residuals')  
ht_pov = sns.distplot(residuals, color='#1f77b4',  
kde=False).set_title('Histogram Residuals')
```



In [67]:

```
poverty_test = pd.read_csv('./Microsoft_-  
_DAT102x_Predicting_Poverty_in_the_United_States_-_Test_values.csv')
```

In [68]:

```
#Create Features  
poverty_test = create_old_age_cat(poverty_test)  
poverty_test = create_aui_pct65y_cat(poverty_test)  
poverty_test = create_features(poverty_test)  
  
#Convert to correct type  
poverty_test = convert_to_cat(poverty_test)  
  
#Drop Features  
poverty_row_id = poverty_test.row_id  
poverty_test = drop_features(poverty_test)  
poverty_test = poverty_test.drop(columns=['row_id'], axis=1)  
  
#Replace NANS  
poverty_test_clean = poverty_test.fillna(poverty_test.median())  
poverty_test_clean = cat_to_dummies(poverty_test_clean)  
  
#Scale Features  
poverty_test_clean = scale_features(poverty_test_clean)
```

In [69]:

```
#Create Prediction
submission = pd.DataFrame(rfecv.predict(poverty_test_clean))
submission = np.clip(submission,2.50, 48.00)
```

In [70]:

```
poverty_submission = pd.concat([poverty_row_id, submission], axis=1)
poverty_submission = poverty_submission.rename(index=str, columns={0:
'poverty_rate'})
poverty_submission = poverty_submission.round({'poverty_rate':2})
```

In [71]:

```
poverty_submission.to_csv(path_or_buf='./MV_Poverty_Submission_AdaReg_tmp.csv
', index=False)
```

## Conclusion¶

The Regression analysis shows that is possible to build an accurate regression model to predict poverty rates of United States' counties using an AdaBoostRegressor. From the data exploration phase it is clear that economical, educational, ethnical and health related factors play an important role in predicting poverty. However is recursive function elimination is used to determine the optimal number of features that leads to the best prediction.