

```
In [3]: import pandas as pd
import numpy as np
import statsmodels.api as sm
from scipy import stats
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import KFold

from sklearn.metrics import r2_score

import seaborn as sns
import matplotlib.pyplot as plt

import re

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

In [4]: df = pd.read_csv('Existing_Buildings_Energy_Performance_Ordinance_Report.csv')
```

```
In [65]: df.columns

Out[65]: Index(['Parcel(s)', 'Building Name', 'Building Address', 'Postal Code',
        'Full Address', 'Floor Area', 'Property Type',
        'Property Type - Self Selected', 'PIM Link', 'Year Built',
        ...,
        '2012 Weather Normalized Source EUI (kBtu/ft2)',
        '2011 ENERGY STAR Score', '2011 Site EUI (kBtu/ft2)',
        '2011 Source EUI (kBtu/ft2)',
        '2011 Percent Better than National Median Site EUI',
        '2011 Percent Better than National Median Source EUI',
        '2011 Total GHG Emissions (Metric Tons CO2e)',
        '2011 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2011 Weather Normalized Site EUI (kBtu/ft2)',
        '2011 Weather Normalized Source EUI (kBtu/ft2)',
        '2011 Weather Normalized Source EUI (kBtu/ft2)'],
      dtype='object', length=113)
```

Explore the raw data

```
In [5]: [x for x in df.columns]

Out[5]: ['Parcel(s)',
        'Building Name',
        'Building Address',
        'Postal Code',
        'Full Address',
        'Floor Area',
        'Property Type',
        'Property Type - Self Selected',
        'PIM Link',
        'Year Built',
        'Energy Audit Due Date',
        'Energy Audit Status',
        'Previous Year ENERGY STAR Score',
        'Two Years Ago ENERGY STAR Score',
        'Current Year YoY Change',
        'Age In Years',
        'Recalculated GHG Emissions Intensity (kgCO2e/ft2)',
        'Is DownTown',
        'ENERGY STAR Score',
        '2012 Weather Normalized Source EUI (kBtu/ft2)',
        '2012 Source EUI (kBtu/ft2)',
        '2012 Percent Better than National Median Site EUI',
        '2012 Percent Better than National Median Source EUI',
        '2012 Total GHG Emissions (Metric Tons CO2e)',
        '2012 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2012 Weather Normalized Site EUI (kBtu/ft2)',
        '2012 Weather Normalized Source EUI (kBtu/ft2)',
        '2012 Weather Normalized Source EUI (kBtu/ft2)',
        '2017 ENERGY STAR Score',
        '2017 Site EUI (kBtu/ft2)',
        '2017 Percent Better than National Median Site EUI',
        '2017 Percent Better than National Median Source EUI',
        '2017 Total GHG Emissions (Metric Tons CO2e)',
        '2017 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2017 Weather Normalized Site EUI (kBtu/ft2)',
        '2017 Weather Normalized Source EUI (kBtu/ft2)',
        '2017 Weather Normalized Source EUI (kBtu/ft2)',
        '2016 ENERGY STAR Score',
        '2016 Site EUI (kBtu/ft2)',
        '2016 Percent Better than National Median Site EUI',
        '2016 Percent Better than National Median Source EUI',
        '2016 Total GHG Emissions (Metric Tons CO2e)',
        '2016 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2016 Weather Normalized Site EUI (kBtu/ft2)',
        '2016 Weather Normalized Source EUI (kBtu/ft2)',
        '2016 Weather Normalized Source EUI (kBtu/ft2)',
        '2015 ENERGY STAR Score',
        '2015 Site EUI (kBtu/ft2)',
        '2015 Percent Better than National Median Site EUI',
        '2015 Percent Better than National Median Source EUI',
        '2015 Total GHG Emissions (Metric Tons CO2e)',
        '2015 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2015 Weather Normalized Site EUI (kBtu/ft2)',
        '2015 Weather Normalized Source EUI (kBtu/ft2)',
        '2015 Weather Normalized Source EUI (kBtu/ft2)',
        '2014 ENERGY STAR Score',
        '2014 Site EUI (kBtu/ft2)',
        '2014 Percent Better than National Median Site EUI',
        '2014 Percent Better than National Median Source EUI',
        '2014 Total GHG Emissions (Metric Tons CO2e)',
        '2014 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2014 Weather Normalized Site EUI (kBtu/ft2)',
        '2014 Weather Normalized Source EUI (kBtu/ft2)',
        '2014 Weather Normalized Source EUI (kBtu/ft2)',
        '2013 ENERGY STAR Score',
        '2013 Site EUI (kBtu/ft2)',
        '2013 Percent Better than National Median Site EUI',
        '2013 Percent Better than National Median Source EUI',
        '2013 Total GHG Emissions (Metric Tons CO2e)',
        '2013 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2013 Weather Normalized Site EUI (kBtu/ft2)',
        '2013 Weather Normalized Source EUI (kBtu/ft2)',
        '2013 Weather Normalized Source EUI (kBtu/ft2)',
        '2012 ENERGY STAR Score',
        '2012 Site EUI (kBtu/ft2)',
        '2012 Source EUI (kBtu/ft2)',
        '2012 Percent Better than National Median Site EUI',
        '2012 Percent Better than National Median Source EUI',
        '2012 Total GHG Emissions (Metric Tons CO2e)',
        '2012 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2012 Weather Normalized Site EUI (kBtu/ft2)',
        '2012 Weather Normalized Source EUI (kBtu/ft2)',
        '2012 Weather Normalized Source EUI (kBtu/ft2)',
        '2011 ENERGY STAR Score',
        '2011 Source EUI (kBtu/ft2)',
        '2011 Percent Better than National Median Site EUI',
        '2011 Percent Better than National Median Source EUI',
        '2011 Total GHG Emissions (Metric Tons CO2e)',
        '2011 Total GHG Emissions Intensity (kgCO2e/ft2)',
        '2011 Weather Normalized Site EUI (kBtu/ft2)',
        '2011 Weather Normalized Source EUI (kBtu/ft2)',
        '2011 Weather Normalized Source EUI (kBtu/ft2)']
```

```
In [6]: df.sample()

Out[6]:
```

| Parcel(s) | Building Name | Building Address | Postal Code | Full Address | Floor Area | Property Type | Property Type - Self Selected | PIM Link | Year Built | Energy Audit Due Date |
|-----------|---------------|------------------|--------------|-------------------------------|-------------------------------------------------|---------------|-------------------------------|----------|------------|-----------------------|
| 1913 | 4750031 | 900 EARI ST | 9000 EURL ST | 900 EARL ST 94124 94124 | 900 EARL ST STYUSAN ST FRANCISCO 94124 | 697C1 | Mixed Residential | NaN | 1900.0 | NaN |

```
In [7]: df.shape

Out[7]: (2629, 113)
```

Retabulate data to building-year level

```
In [3]: # Unstack the data just the year-based data columns
df = df.groupby('Parcel(s)').unstack()

In [15]: # Extract the years and make them their own column
year_starts = [x.find('201') for x in temp2.level_0.values]

years = []
for year_start, value in zip(year_starts, temp2.level_0.values):
    years.append(value/year_start:year_start+1)

temp2['year'] = years

In [60]: # Take out the mention of years to normalize the values of the columns
years_list = ['2010', '2011', '2012', '2013', '2014', '2015', '2016', '2017', '2018', '2019']
col_names = temp2.level_0.values

res = []
for col_name in col_names:
    for sub in years_list:
        if sub in col_name:
            col_name = col_name.replace(sub, '')
            res.append(" ".join(col_name.split(' ')))
        else:
            continue

temp2['metric_year_stripped'] = res

In [65]: temp3 = temp2.pivot(index='Parcel(s)', 'year'), columns='metric_year_stripped', values='0')

In [72]: # Rejoin the non-year-level data
cols_no_year = [x for x in df.columns if ('201' not in x)]
temp4 = temp3.reset_index().merge(df[cols_no_year], left_on='Parcel(s)', right_on='Parcel(s)', how='left')

In [76]: temp4.to_csv('energy_building_data_retabulated.csv', index=False)
```

Spot check a few

```
In [91]: df.sample()

Out[91]:
```

| Parcel(s) | Building Name | Building Address | Postal Code | Full Address | Floor Area | Property Type | Property Type - Self Selected | PIM Link | Year Built |
|-----------|---------------|------------------|-------------------|-------------------------------------------------------|-------------------------------------------------------|---------------|-------------------------------|----------|------------|
| 1567 | 0259026 | 655 California | 655 CALIFORNIA ST | 555 CALIFORNIA ST STYUSAN ST FRANCISCO 94104 | 555 CALIFORNIA ST STYUSAN ST FRANCISCO 94104 | 197B104 | Commercial | Office | 1969.0 |

```
In [92]: temp4[temp4['Parcel(s)']=='0259026']

Out[92]:
```

| Parcel(s) | year | Benchmark Status | ENERGY STAR Score | Percent Better than National Median Site EUI | Percent Better than National Median Source EUI | Percentage Better than National Median Source EUI | Reason for Exemption | Site EUI (kBtu/ft2) | Source EUI (kBtu/ft2) | Total GHG Emissions (Metric Tons CO2e) | Total GHG Emissions Intensity (kgCO2e/ft2) |
|-----------|---------|------------------|-------------------------------|----------------------------------------------|------------------------------------------------|---------------------------------------------------|----------------------|---------------------|-----------------------|----------------------------------------|--------------------------------------------|
| 3670 | 0259026 | 2010 | Completed | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 3671 | 0259026 | 2011 | Completed | 70 | -19.1 | NaN | NaN | 96.7 | NaN | 15049.7 | NaN |
| 3672 | 0259026 | 2012 | Completed | 71 | -20.6 | NaN | NaN | 91.7 | 234 | 14336.9 | NaN |
| 3673 | 0259026 | 2013 | Violation - Insufficient Data | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| 3674 | 0259026 | 2014 | Completed | 68 | -19.4 | -19.4 | NaN | 83.4 | 204.7 | 11499.5 | 6.4 |
| 3675 | 0259026 | 2015 | Completed | 69 | -19.7 | -19.7 | NaN | 87.7 | 205.7 | 11956.8 | 6.6 |
| 3676 | 0259026 | 2016 | Completed | 77 | -28 | -28 | NaN | 75.9 | 179.7 | 10798.1 | 6 |
| 3677 | 0259026 | 2017 | Completed | 78 | -29.5 | -29.5 | NaN | 75.2 | 177.7 | 9804 | 5.4 |
| 3678 | 0259026 | 2018 | Completed | 58 | -10.6 | -10.6 | NaN | 75.4 | 161.7 | 9385.5 | 5.2 |
| 3679 | 0259026 | 2019 | Completed | 62 | -15.3 | NaN | -15.3 | 80 | 166 | 9921.7 | 5.5 |

```
In [3]: df = pd.read_csv('energy_building_data_final.csv')
```

Modeling (retabulated data)

```
In [13]: df[df['ENERGY STAR Score'].notnull()].isnull().mean()

Out[13]: Unnamed: 0 0.000000
Parcel(s) 0.000000
Benchmark Status 0.000000
ENERGY STAR Score 0.000000
Percent Better than National Median Site EUI 0.000000
Percent Better than National Median Source EUI 0.253865
Percentage Better than National Median Source EUI 0.863443
Reason for Exemption 0.968788
Site EUI (kBtu/ft2) 0.017228
Source EUI (kBtu/ft2) 0.090455
Total GHG Emissions (Metric Tons CO2e) 0.081819
Total GHG Emissions Intensity (kgCO2e/ft2) 0.151849
Weather Normalized Site EUI (kBtu/ft2) 0.012580
Weather Normalized Source EUI (kBtu/ft2) 0.068351
Building Name 0.000000
Building Address 0.000000
Postal Code 0.000000
Full Address 0.000000
Floor Area 0.000000
Property Type 0.000000
Property Type - Self Selected 0.000000
PIM Link 0.000000
Year Built 0.009548
Energy Audit Due Date 0.843346
Energy Audit Status 0.044711
Previous Year ENERGY STAR Score 0.210654
Two Years Ago ENERGY STAR Score 0.366475
Current Year YoY Change 0.210654
Previous Year YoY Change 0.366475
Age In Years 0.009548
Recalculated GHG Emissions Intensity (kgCO2e/ft2) 0.068181
Is DownTown 0.000000
dtype: float64

In [14]: df.dtypes

Out[14]: Unnamed: 0 object
Parcel(s) int64
Benchmark Status object
ENERGY STAR Score float64
Percent Better than National Median Site EUI float64
Percent Better than National Median Source EUI float64
Reason for Exemption object
Site EUI (kBtu/ft2) float64
Source EUI (kBtu/ft2) float64
Total GHG Emissions (Metric Tons CO2e) float64
Total GHG Emissions Intensity (kgCO2e/ft2) float64
Weather Normalized Site EUI (kBtu/ft2) float64
Weather Normalized Source EUI (kBtu/ft2) float64
Building Name object
Building Address object
Postal Code int64
Full Address object
Floor Area int64
Property Type object
Property Type - Self Selected object
PIM Link object
Year Built float64
Energy Audit Due Date object
Energy Audit Status object
Previous Year ENERGY STAR Score float64
Two Years Ago ENERGY STAR Score float64
Current Year YoY Change float64
Previous Year YoY Change float64
Age In Years float64
Recalculated GHG Emissions Intensity (kgCO2e/ft2) float64
Is DownTown bool
dtype: object
```

Modeling

```
In [15]: df['Year Built'] = df['Year Built'].replace(0, np.NaN)
```

Linear Regression

Using all the features

```
In [6]: dummy_features = ['Property Type', 'Postal Code']

In [7]: features = [col for col in df_train.numeric.columns if col not in ['Reason for Exemption', 'Percent Better than National Median Source EUI', 'Percent Better than National Median Site EUI', 'ENERGY STAR Score']]

# drop nans
df_train_dropped_nas = df_train.numeric.dropna(subset=features + ['ENERGY STAR Score'])

X_train_no_constant = df_train_dropped_nas[features]
y_train = sm.add_constant(df_train_dropped_nas['ENERGY STAR Score'])

X_train = sm.add_constant(X_train_no_constant)

print(X_train.shape)

est = sm.OLS(y_train, X_train).fit()

print(est.summary())

OLS Regression Results
=====
Dep. Variable: ENERGY STAR Score R-squared: 0.347
Model: OLS Adj. R-squared: 0.345
Method: Least Squares F-statistic: 185.2
Date: Sun, 29 Nov 2020 Prob. F-statistic: 0.000
Time: 16:05:37 Log-Likelihood: -17246.
No. Observations: 3849 AIC: 3.452e+04
Of Residuals: 3847 BIC: 0.068181
R-squared: 0.000000
Covariance Type: nonrobust
=====
[0.925 0.975] coef std err t P>|t [0.025 0.975]
-----
const 2.587e+04 4018.487 6.439 0.000 1.8e+04
year 0.1585 0.160 0.992 0.321 -0.155
Percent Better than National Median Site EUI 0.0024 0.000 -1.618 0.000 -0.924
Site EUI (kBtu/ft2) -0.4172 0.258 -1.637 0.106 -0.922
Source EUI (kBtu/ft2) -0.7315 0.166 -4.414 0.000 -1.056
Total GHG Emissions (Metric Tons CO2e) 0.0002 0.001 0.296 0.767 -0.001
Weather Normalized Site EUI (kBtu/ft2) -0.1947 0.123 -1.567 0.000 2.811
Weather Normalized Source EUI (kBtu/ft2) -0.1419 0.254 -0.561 0.575 -0.638
Postal Code 0.7212 0.166 4.351 0.000 0.396
Floor Area -0.2787 0.042 -6.582 0.000 -0.362
Year Built 1.849e-05 3.18e-06 5.813 0.000 1.22e-05
2.47e-05 0.0588 0.000 0.732 0.464 -0.099
years_since_built 0.0993 0.000 1.243 0.214 -0.057
=====
Omnibus: 639.774 Durbin-Watson: 0.750
Prob(Omnibus): 0.000 Jarque-Bera (JB): 9142.267
Skew: -0.321 Kurtosis: 10.521
Cond. No.: 8.14e+16
=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The smallest eigenvalue is 4.26e-20. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.
```

```
In [14]: df.columns

Out[14]: Index(['Parcel(s)', 'year', 'Benchmark Status', 'ENERGY STAR Score',
        'Percent Better than National Median Site EUI',
        'Percent Better than National Median Source EUI',
        'Percentage Better than National Median Source EUI',
        'Reason for Exemption', 'Site EUI (kBtu/ft2)', 'Source EUI (kBtu/ft2)',
        'Total GHG Emissions (Metric Tons CO2e)',
        'Total GHG Emissions Intensity (kgCO2e/ft2)',
        'Weather Normalized Site EUI (kBtu/ft2)',
        'Weather Normalized Source EUI (kBtu/ft2)',
        'Building Name',
        'Building Address', 'Full Address', 'Floor Area',
        'Property Type', 'Property Type - Self Selected',
        'Year Built', 'Energy Audit Due Date', 'Energy Audit Status'],
      dtype='object')
```

```
In [3]: df.head()

Out[3]:
```

| Unnamed: 0 | Parcel(s) | year | Benchmark Status | ENERGY STAR Score | Percent Better than National Median Site EUI | Percent Better than National Median Source EUI | Percentage Better than National Median Source EUI | Reason for Exemption | Site EUI (kBtu/ft2) | Source EUI (kBtu/ft2) | Total GHG Emissions (Metric Tons CO2e) | Total GHG Emissions Intensity (kgCO2e/ft2) |
|------------|-----------|---------|------------------|-------------------|----------------------------------------------|------------------------------------------------|---------------------------------------------------|----------------------|---------------------|-----------------------|----------------------------------------|--------------------------------------------|
| 0 | 2 | 0010001 | 2012 | Completed | 81.0 | -34.5 | NaN | NaN | NaN | 71.0 | 149.9 | 696.4 |
| 1 | 3 | 0010001 | 2013 | Completed | 74.0 | -26.2 | -26.2 | NaN | NaN | 81.3 | 166.9 | 786.5 |
| 2 | 4 | 0010001 | 2014 | Completed | 56.0 | -6.5 | -6.5 | NaN | NaN | 73.4 | 158.5 | 668.9 |
| 3 | 5 | 0010001 | 2015 | Completed | 72.0 | -23.0 | -23.0 | NaN | NaN | 72.0 | 153.9 | 653.2 |
| 4 | 7 | 0010001 | 2016 | Completed | 75.0 | -26.3 | -26.3 | NaN | NaN | 68.2 | 149.7 | 652.4 |

```
In [5]: # make full list of all relevant features
features_relevant = ['Benchmark Status',
                    'Percent Better than National Median Site EUI',
                    'Percent Better than National Median Source EUI',
                    'Site EUI (kBtu/ft2)',
                    'Source EUI (kBtu/ft2)',
                    'Total GHG Emissions (Metric Tons CO2e)',
                    'Total GHG Emissions Intensity (kgCO2e/ft2)', # recalculated below
                    'Weather Normalized Site EUI (kBtu/ft2)',
                    'Weather Normalized Source EUI (kBtu/ft2)',
                    'Postal Code', # dummyify or just use downtown
                    'Floor Area',
                    'Property Type', # dummyify
                    'Property Type - Self Selected',
                    'Previous Year ENERGY STAR Score',
                    'Two Years Ago ENERGY STAR Score',
                    'Current Year YoY Change',
                    'Age In Years',
                    'Recalculated GHG Emissions Intensity (kgCO2e/ft2)',
                    'Is DownTown']

In [5]: df[features_relevant].head()

Out[5]:
```

| Unnamed: 0 | Parcel(s) | year | Benchmark Status | Percent Better than National Median Site EUI | Percent Better than National Median Source EUI | Total GHG Emissions (Metric Tons CO2e) | Weather Normalized Site EUI (kBtu/ft2) | Weather Normalized Source EUI (kBtu/ft2) | Postal Code | Floor Area | Property Type | Property Type - Self Selected | Previous Year ENERGY STAR Score | Two Years Ago ENERGY STAR Score |
|------------|-----------|-------|------------------|----------------------------------------------|------------------------------------------------|----------------------------------------|----------------------------------------|------------------------------------------|-------------|------------|---------------|-------------------------------|---------------------------------|---------------------------------|
| 0 | Completed | -34.5 | -34.5 | 71.0 | 149.9 | 696.4 | 70.7 | 149.7 | 94109 | 133675 | Commercial | Office | 8 | N |
| 1 | Completed | -26.2 | -26.2 | 81.3 | 166.9 | 786.5 | 81.6 | 167.2 | 94109 | 133675 | Commercial | Office | 8 | N |
| 2 | Completed | -6.5 | -6.5 | 73.4 | 158.5 | 668.9 | 73.7 | 162.4 | 94109 | 133675 | Commercial | Office | 5 | 7 |
| 3 | Completed | -23.0 | -23.0 | 72.0 | 153.9 | 653.2 | 71.7 | 155.6 | 94109 | 133675 | Commercial | Office | 5 | 7 |
| 4 | Completed | -26.3 | -26.3 | 68.2 | 149.7 | 652.4 | 71.5 | 153.2 | 94109 | 133675 | Commercial | Office | 7 | 7 |

```
In [6]: df[features_relevant].corr()

Out[6]:
```

```
lm_data1 = lm_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           ], axis=1)

lm_data1 = lm_data1.dropna()

X_train = lm_data1[lm_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = lm_data1[lm_data1['year'].isin(['2018', '2019'])]['ENERGY STAR Score']

X_test = lm_data1[lm_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = lm_data1[lm_data1['year'].isin(['2018', '2019'])]['ENERGY STAR Score']

X_train = sm.add_constant(X_train)

est = sm.OLS(y_train, X_train).fit()
print(est.summary())

y_pred = est.predict(X_test)
```

```
In [33]: sns.heatmap(df[features_relevant + ['ENERGY STAR Score']].corr(), [[['ENERGY STAR Score']], cmap='YlGnB
```

```
Out[33]: <matplotlib.axes._subplots.AxesSubplot at 0x7f024b69a340>
```

```
In [34]: df[features_relevant].dtypes

Out[34]: Benchmark Status object
Percent Better than National Median Site EUI float64
Percent Better than National Median Source EUI float64
Site EUI (kBtu/ft2) float64
Source EUI (kBtu/ft2) float64
Total GHG Emissions (Metric Tons CO2e) float64
Weather Normalized Site EUI (kBtu/ft2) float64
Weather Normalized Source EUI (kBtu/ft2) float64
Postal Code int64
Floor Area int64
Property Type object
Property Type - Self Selected object
Previous Year ENERGY STAR Score float64
Two Years Ago ENERGY STAR Score float64
Current Year YoY Change float64
Previous Year YoY Change float64
Age In Years float64
Recalculated GHG Emissions Intensity (kgCO2e/ft2) float64
Is DownTown bool
dtype: object
```

```
In [3]: # dummyify features
property_dummies = pd.get_dummies(df['Property Type'])
benchmark_dummies = pd.get_dummies(df['Benchmark Status'])
df['Is DownTown'] = df['Is DownTown'].astype(int)

dummies_all = pd.concat([property_dummies, benchmark_dummies], axis=1)
df_dummies = pd.concat((df, dummies_all), axis=1)
```

```
In [35]: lm_data = df_dummies[features_relevant + dummies_all.columns.to_list() + ['ENERGY STAR Score']]
lm_data = lm_data.drop(['Benchmark Status',
                        'Postal Code',
                        'Property Type - Self Selected',
                        ], axis=1)
lm_data = lm_data.dropna()
```

```
X_train = lm_data[~lm_data['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = lm_data[~lm_data['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = lm_data[~lm_data['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = lm_data[~lm_data['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_train = sm.add_constant(X_train)

est = sm.OLS(y_train, X_train).fit()

y_pred = est.predict(X_test)

print(r2_score(y_test, y_pred))
```

```
OLS Regression Results
=====
Dep. Variable: ENERGY STAR Score R-squared: 0.000
Model: OLS Adj. R-squared: 1.000
Method: Least Squares F-statistic: 1.000
Date: Mon, 14 Dec 2020 Prob. F-statistic: 0.000
Time: 16:35:04 Log-Likelihood: 63991.
No. Observations: 2765 AIC: -1.261e+05
Of Residuals: 2748 BIC: -1.260e+05
Covariance Type: nonrobust
=====
[0.025 0.975] coef
```



```
In [32]: # Take out multicollinear variables
lm_data1 = df.dummies[features_relevant+ dummies_all.columns.to_list()] + ['ENERGY STAR Score']

lm_data1 = lm_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           'Two Years Ago ENERGY STAR Score',
                           'Weather Normalized Source EUI (kBtu/ft2)',
                           'Percent Better than National Median Source EUI',
                           'Current Year YoY Change',
                           'Total GHG Emissions (Metric Tons CO2e)',
                           'Source EUI (kBtu/ft2)',
                           ], axis=1)

lm_data1 = lm_data1.dropna()

X_train = lm_data1[~lm_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = lm_data1[~lm_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = lm_data1[lm_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = lm_data1[lm_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_train = sm.add_constant(X_train)

est = sm.OLS(y_train, X_train).fit()
print(est.summary())

y_pred = est.predict(X_test)
print(r2_score(y_test, y_pred))

=====
OLS Regression Results
=====
Dep. Variable:          ENERGY STAR Score    R-squared:                0.876
Model:                  OLS                  Adj. R-squared:           0.875
Method:                 Least Squares        F-statistic:             1762.
Date:                   Mon, 14 Dec 2020      Prob (F-statistic):       0.00
Time:                   16:34:02              Log-Likelihood:          -9973.2
No. Observations:       2765                AIC:                     1.997e+04
Df Residuals:           2765                BIC:                     2.004e+04
Df Model:                11
Covariance Type:        nonrobust

=====
coef      std err      t      P>|t|
-----
Percent Better than National Median Site EUI -0.1597    0.006   -26.205    0.000 -
0.172    -0.148
Site EUI (kBtu/ft2) -0.1482    0.104   -1.422    0.155 -
0.352    0.056
Weather Normalized Site EUI (kBtu/ft2) 0.1115    0.104    1.075    0.282 -
0.092    0.315
Floor Area 2.516e-06  6.88e-07  3.657    0.000  1.17
e-06    3.87e-06
Previous Year ENERGY STAR Score 0.634    0.675    0.939    0.345
0.634    0.675
Previous Year YoY Change -0.0378    0.016   -2.420    0.016 -
0.068    -0.007
Age In Years 0.0002    0.000    0.469    0.639 -
0.001    0.001
Recalculated GHG Emissions Intensity (kgCO2e/ft2) 0.2946    0.139    2.115    0.035
0.021    0.568
Is Downtown 1.0287    0.353    2.914    0.004
0.337    1.721
Commercial 18.5215    2.566    7.219    0.000  1
3.491    23.552
Mixed Residential -3.074e-16  1.76e-15 -0.175    0.861 -
e-15    3.15e-15
Multifamily 1.925e-16  4.38e-16  0.440    0.660 -6.65
e-16    1.05e-15
Residential 0 0 nan nan
0 0
Complied 3.2800    2.509    1.267    0.191 -
1.640    8.296
Data Not Verified -1.4084    6.001   -0.207    0.836 -1
4.744    11.927
Exempt 0 0 nan nan
0 0
Violation - Insufficient Data 16.6500    4.836    4.126    0.000
0.573    24.563
Omnibus: 1520.982 Durbin-Watson: 1.654
Prob(Omnibus): 0.000 Jarque-Bera (JB): 124010.017
Skew: 1.735 Prob(JB): 0.00
Kurtosis: 35.624 Cond. No. inf

=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The smallest eigenvalue is 0. This might indicate that there are
strong multicollinearity problems or that the design matrix is singular.
0.810461640486299
```

```
Out[59]: Dep. Variable:          ENERGY STAR Score    R-squared:                0.876
Model:                  OLS                  Adj. R-squared:           0.875
Method:                 Least Squares        F-statistic:             1939.
Date:                   Tue, 15 Dec 2020      Prob (F-statistic):       0.00
Time:                   06:55:08              Log-Likelihood:          -9973.3
No. Observations:       2765                AIC:                     1.997e+04
Df Residuals:           2765                BIC:                     2.003e+04
Df Model:                10
Covariance Type:        nonrobust

=====
coef      std err      t      P>|t|
-----
Percent Better than National Median Site EUI -0.1598    0.006   -26.276    0.000 -
0.172    -0.148
Site EUI (kBtu/ft2) -0.1500    0.104   -1.440    0.150 -
0.352    0.056
Weather Normalized Site EUI (kBtu/ft2) 0.1131    0.104    1.092    0.275 -
0.090    0.316
Floor Area 2.518e-06  6.88e-07  3.661    0.000  1.17
e-06    3.87e-06
Previous Year ENERGY STAR Score 0.634    0.675    0.939    0.345
0.634    0.675
Previous Year YoY Change -0.0376    0.016   -2.407    0.016 -
0.068    -0.007
Recalculated GHG Emissions Intensity (kgCO2e/ft2) 0.2999    0.139    2.160    0.031
0.021    0.572
Is Downtown 1.0326    0.353    2.927    0.003
0.341    1.724
Complied 21.8190    2.582    7.903    0.000  2
0.286    23.352
Data Not Verified 17.1038    0.875    1.906    0.057 -
0.494    34.762
Violation - Insufficient Data 35.2652    4.530    7.785    0.000  2
6.383    44.147
Omnibus: 1520.500 Durbin-Watson: 1.654
Prob(Omnibus): 0.000 Jarque-Bera (JB): 123926.442
Skew: 1.734 Prob(JB): 0.00
Kurtosis: 35.614 Cond. No. 1.63e+07

=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.63e+07. This might indicate that there are
strong multicollinearity or other numerical problems.
0.8109496927539122
```

```
In [59]: X_train.columns
Out[59]: Index(['Percent Better than National Median Site EUI', 'Site EUI (kBtu/ft2)',
              'Weather Normalized Site EUI (kBtu/ft2)', 'Floor Area',
              'Previous Year ENERGY STAR Score', 'Previous Year YoY Change',
              'Recalculated GHG Emissions Intensity (kgCO2e/ft2)', 'Is Downtown',
              'Complied', 'Data Not Verified', 'Violation - Insufficient Data'],
              dtype='object')

In [60]: X_train.columns
Out[60]: Index(['const', 'Percent Better than National Median Site EUI',
              'Site EUI (kBtu/ft2)', 'Weather Normalized Site EUI (kBtu/ft2)',
              'Floor Area', 'Previous Year ENERGY STAR Score',
              'Previous Year YoY Change',
              'Recalculated GHG Emissions Intensity (kgCO2e/ft2)', 'Is Downtown',
              'Complied', 'Data Not Verified', 'Violation - Insufficient Data'],
              dtype='object')
```

Random Forest

With all relevant features, default hyperparameters

```
In [16]: rf_data1 = df.dummies[features_relevant+ dummies_all.columns.to_list()] + ['ENERGY STAR Score']

rf_data1 = rf_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           ], axis=1)

rf_data1 = rf_data1.dropna()

X_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

regr = RandomForestRegressor()

regr.fit(X_train, y_train.values.ravel())
print(regr.score(X_train, y_train))
print(regr.score(X_test, y_test))

for (x, y) in zip(X_train.columns, regr.feature_importances_):
    print(x, round(y, 3))

0.99265977805959
0.979522382758568
Percent Better than National Median Site EUI 0.464
Percent Better than National Median Source EUI 0.526
Site EUI (kBtu/ft2) 0.001
Source EUI (kBtu/ft2) 0.001
Total GHG Emissions (Metric Tons CO2e) 0.0
Weather Normalized Site EUI (kBtu/ft2) 0.001
Weather Normalized Source EUI (kBtu/ft2) 0.001
Floor Area 0.0
Previous Year ENERGY STAR Score 0.002
Two Years Ago ENERGY STAR Score 0.0
Current Year YoY Change 0.001
Previous Year YoY Change 0.0
Age In Years 0.0
Recalculated GHG Emissions Intensity (kgCO2e/ft2) 0.001
Is Downtown 0.0
Commercial 0.0
Mixed Residential 0.0
Residential 0.0
Multifamily 0.0
Complied 0.0
Data Not Verified 0.0
Exempt 0.0
Violation - Insufficient Data 0.0
```

With cross validation of the complexity parameter

Note that we are using only X_train, so data from prior to 2018

Looks like they are all pretty good, with accuracy 0.99

```
In [50]: # With everything
rf_data1 = df.dummies[features_relevant+ dummies_all.columns.to_list()] + ['ENERGY STAR Score']

rf_data1 = rf_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           ], axis=1)

rf_data1 = rf_data1.dropna()

X_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

CPs = [0.0, 0.005, 0.01, 0.015, 0.02, 0.025]

kf = KFold(n_splits=5)

for cp in CPs:
    regr = RandomForestRegressor(ccp_alpha=cp)
    scores = []

    for k, (train_index, test_index) in enumerate(kf.split(X_train.values)):
        X_train_kfold, y_test_kfold = X_train.values[train_index], X_train.values[test_index]
        y_train_kfold, y_test_kfold = y_train.values[train_index], y_train.values[test_index]

        regr.fit(X_train_kfold, y_train_kfold.ravel())

        scores.append(regr.score(X_test_kfold, y_test_kfold))

    print('Mean score for cp {}: {}'.format(cp, np.round(np.mean(scores), 2)))
    print()

Mean score for cp 0.0: 0.99

Mean score for cp 0.005: 0.99

Mean score for cp 0.01: 0.99

Mean score for cp 0.015: 0.99

Mean score for cp 0.02: 0.99

Mean score for cp 0.025: 0.99
```

With cross validation of the max number of features to sample

```
In [52]: # Starting at all the features, moving towards just 1
# With everything
rf_data1 = df.dummies[features_relevant+ dummies_all.columns.to_list()] + ['ENERGY STAR Score']

rf_data1 = rf_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           ], axis=1)

rf_data1 = rf_data1.dropna()

X_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

features_to_sample = np.arange(X_train.shape[1], 0, -1)

for f in features_to_sample:
    clf = RandomForestRegressor(max_features=f)
    scores = []

    for k, (train_index, test_index) in enumerate(kf.split(X_train.values)):
        X_train_kfold, y_test_kfold = X_train.values[train_index], X_train.values[test_index]
        y_train_kfold, y_test_kfold = y_train.values[train_index], y_train.values[test_index]

        regr.fit(X_train_kfold, y_train_kfold.ravel())

        scores.append(regr.score(X_test_kfold, y_test_kfold))

    print('Mean score for max_features {}: {}'.format(f, np.round(np.mean(scores), 2)))
    print()

Mean score for max_features 23: 0.99

Mean score for max_features 22: 0.99

Mean score for max_features 21: 0.99

Mean score for max_features 20: 0.99

Mean score for max_features 19: 0.99

Mean score for max_features 18: 0.99

Mean score for max_features 17: 0.99

Mean score for max_features 16: 0.99

Mean score for max_features 15: 0.99

Mean score for max_features 14: 0.99

Mean score for max_features 13: 0.99

Mean score for max_features 12: 0.99

Mean score for max_features 11: 0.99

Mean score for max_features 10: 0.99

Mean score for max_features 9: 0.99

Mean score for max_features 8: 0.99

Mean score for max_features 7: 0.99

Mean score for max_features 6: 0.99

Mean score for max_features 5: 0.99

Mean score for max_features 4: 0.99

Mean score for max_features 3: 0.99

Mean score for max_features 2: 0.99

Mean score for max_features 1: 0.99
```

```
In [53]: # Do the same thing on max_depth (max depth of the trees)
# Starting at all the features, moving towards just 1
rf_data1 = df.dummies[features_relevant+ dummies_all.columns.to_list()] + ['ENERGY STAR Score']

rf_data1 = rf_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           ], axis=1)

rf_data1 = rf_data1.dropna()

X_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

max_depths = np.arange(1, 10)

for d in max_depths:
    clf = RandomForestRegressor(max_depth=d)
    scores = []

    for k, (train_index, test_index) in enumerate(kf.split(X_train.values)):
        X_train_kfold, y_test_kfold = X_train.values[train_index], X_train.values[test_index]
        y_train_kfold, y_test_kfold = y_train.values[train_index], y_train.values[test_index]

        regr.fit(X_train_kfold, y_train_kfold.ravel())

        scores.append(regr.score(X_test_kfold, y_test_kfold))

    print('Mean score for max_depth {}: {}'.format(d, np.round(np.mean(scores), 2)))
    print()

['Percent Better than National Median Site EUI', 'Percent Better than National Median Source EUI', 'Site EUI (kBtu/ft2)', 'Source EUI (kBtu/ft2)', 'Total GHG Emissions (Metric Tons CO2e)', 'Weather Normalized Site EUI (kBtu/ft2)', 'Weather Normalized Source EUI (kBtu/ft2)', 'Floor Area', 'Previous Year ENERGY STAR Score', 'Two Years Ago ENERGY STAR Score', 'Current Year YoY Change', 'Previous Year YoY Change', 'Age In Years', 'Recalculated GHG Emissions Intensity (kgCO2e/ft2)', 'Is Downtown', 'Commerce', 'Mixed Residential', 'Multifamily', 'Residential', 'Complied', 'Data Not Verified', 'Exempt', 'Violation - Insufficient Data']
Mean score for max_depth 1: 0.99

Mean score for max_depth 2: 0.99

Mean score for max_depth 3: 0.99

Mean score for max_depth 4: 0.99

Mean score for max_depth 5: 0.99

Mean score for max_depth 6: 0.99

Mean score for max_depth 7: 0.99

Mean score for max_depth 8: 0.99

Mean score for max_depth 9: 0.99
```

Suppose the building is new, and we must take out the previous 1 and 2 years' scores and YoY change, how do they models perform in this case?

Linear Regression

```
In [38]: # Take out multicollinear variables
lm_data1 = df.dummies[features_relevant+ dummies_all.columns.to_list()] + ['ENERGY STAR Score']

lm_data1 = lm_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           'Two Years Ago ENERGY STAR Score',
                           'Weather Normalized Source EUI (kBtu/ft2)',
                           'Percent Better than National Median Source EUI',
                           'Current Year YoY Change',
                           'Total GHG Emissions (Metric Tons CO2e)',
                           'Source EUI (kBtu/ft2)',
                           'Previous Year YoY Change',
                           ], axis=1)

lm_data1 = lm_data1.dropna()

X_train = lm_data1[~lm_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = lm_data1[~lm_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = lm_data1[lm_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = lm_data1[lm_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_train = sm.add_constant(X_train)

est = sm.OLS(y_train, X_train).fit()
print(est.summary())

y_pred = est.predict(X_test)
print(r2_score(y_test, y_pred))

=====
OLS Regression Results
=====
Dep. Variable:          ENERGY STAR Score    R-squared:                0.208
Model:                  OLS                  Adj. R-squared:           0.200
Method:                 Least Squares        F-statistic:             122.5
Date:                   Tue, 15 Dec 2020      Prob (F-statistic):       1.20e-227
Time:                   06:06:57              Log-Likelihood:          -21348.
No. Observations:       4604                AIC:                     4.272e+04
Df Residuals:           4603                BIC:                     4.279e+04
Df Model:                10
Covariance Type:        nonrobust

=====
coef      std err      t      P>|t|
-----
Percent Better than National Median Site EUI -0.0012    0.000   -3.113    0.002 -
0.002    -0.000
Site EUI (kBtu/ft2) -0.5443    0.170   -3.202    0.001 -
0.878    0.211
Weather Normalized Site EUI (kBtu/ft2) 0.2369    0.169    1.400    0.161 -
0.095    0.569
Floor Area 1.595e-05  1.41e-06  11.323    0.000  1.32
e-05    1.87e-05
Age In Years 0.0001    0.000    0.814    0.667 -
0.001    0.002
Recalculated GHG Emissions Intensity (kgCO2e/ft2) 3.4486    0.122    28.319    0.000
3.210    3.687
Is Downtown 6.0706    0.692    8.757    0.000
0.644    6.327
Commercial 56.1895    5.291    10.621    0.000  4
5.817    66.562
Mixed Residential -6.73e-15  1.86e-15 -3.628    0.000 -1.04
e-14    -3.09e-15
Multifamily -3.673e-15  3.11e-15 -1.181    0.238 -0.97
e-15    2.42e-15
Residential 1.431e-15  2.32e-15  0.617    0.537 -3.11
e-15    25.837
Complied 15.4958    5.275    2.938    0.003
3.187    49.955
Data Not Verified 13.3845    18.654    0.718    0.473 -2
5.155    29.571
Exempt 14.4052    7.736    1.862    0.063 -
0.161    29.575
Violation - Insufficient Data 12.9041    9.585    1.346    0.178 -
5.887    31.695
Omnibus: 697.912 Durbin-Watson: 0.705
Prob(Omnibus): 0.000 Jarque-Bera (JB): 1185.587
Skew: -0.084 Prob(JB): 3.57e-258
Kurtosis: 4.463 Cond. No. 7.17e+26

=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The smallest eigenvalue is 7.65e-40. This might indicate that there are
strong multicollinearity problems or that the design matrix is singular.
0.679792718143806
```

```
In [7]: # Previous Year ENERGY STAR Score
# Two Years Ago ENERGY STAR Score
# Current Year YoY Change
# Previous Year YoY Change
```

```
In [7]: # No previous years' data
# Default hyperparameters
rf_data1 = df.dummies[features_relevant+ dummies_all.columns.to_list()] + ['ENERGY STAR Score']

rf_data1 = rf_data1.drop(['Benchmark Status',
                           'Postal Code',
                           'Property Type',
                           'Property Type - Self Selected',
                           'Previous Year ENERGY STAR Score',
                           'Two Years Ago ENERGY STAR Score',
                           'Current Year YoY Change',
                           'Previous Year YoY Change',
                           ], axis=1)

rf_data1 = rf_data1.dropna()

X_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_train = rf_data1[~rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

X_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])].drop(['year', 'ENERGY STAR Score'], axis=1)
y_test = rf_data1[rf_data1['year'].isin(['2018', '2019'])][['ENERGY STAR Score']]

regr = RandomForestRegressor()

regr.fit(X_train, y_train.values.ravel())
print(regr.score(X_train, y_train))
print(regr.score(X_test, y_test))

for (x, y) in zip(X_train.columns, regr.feature_importances_):
    print(x, round(y, 3))

0.9999187520005361
0.9745513961848253
```

```
In [15]: for (x, y) in zip(X_train.columns, regr.feature_importances_):
    print(x, round(y, 3))

Percent Better than National Median Site EUI 0.444
Percent Better than National Median Source EUI 0.545
Site EUI (kBtu/ft2) 0.002
Source EUI (kBtu/ft2) 0.002
Total GHG Emissions (Metric Tons CO2e) 0.001
Weather Normalized Site EUI (kBtu/ft2) 0.002
Weather Normalized Source EUI (kBtu/ft2) 0.002
Floor Area 0.001
Age In Years 0.001
Recalculated GHG Emissions Intensity (kgCO2e/ft2) 0.001
Is Downtown 0.0
Commercial 0.0
Mixed Residential 0.0
Multifamily 0.0
Residential 0.0
Complied 0.0
Data Not Verified 0.0
Exempt 0.0
Violation - Insufficient Data 0.0
```

```
In [17]: frame = { 'Feature': X_train.columns, 'Feature importance': regr.feature_importances_ }

In [22]: pd.DataFrame(frame)

Out[22]:
```

| | Feature | Feature importance |
|----|---------------------------------------------------|--------------------|
| 0 | Percent Better than National Median Site EUI | 0.463618 |
| 1 | Percent Better than National Median Source EUI | 0.525894 |
| 2 | Site EUI (kBtu/ft2) | 0.001293 |
| 3 | Source EUI (kBtu/ft2) | 0.000923 |
| 4 | Total GHG Emissions (Metric Tons CO2e) | 0.000491 |
| 5 | Weather Normalized Site EUI (kBtu/ft2) | 0.001181 |
| 6 | Weather Normalized Source EUI (kBtu/ft2) | 0.001027 |
| 7 | Floor Area | 0.000403 |
| 8 | Previous Year ENERGY STAR Score | 0.001688 |
| 9 | Two Years Ago ENERGY STAR Score | 0.000443 |
| 10 | Current Year YoY Change | 0.001369 |
| 11 | Previous Year YoY Change | 0.000225 |
| 12 | Age In Years | 0.000402 |
| 13 | Recalculated GHG Emissions Intensity (kgCO2e/ft2) | 0.000961 |
| 14 | Is Downtown | 0.000055 |
| 15 | Commercial | 0.000000 |
| 16 | Mixed Residential | 0.000000 |
| 17 | Multifamily | 0.000000 |
| 18 | Residential | 0.000000 |
| 19 | Complied | 0.000003 |
| 20 | Data Not Verified | 0.000002 |
| 21 | Exempt | 0.000000 |
| 22 | Violation - Insufficient Data | 0.000002 |

```
In [7]:
```