- 3.3.0 EHOUGE CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

1 Business Problem:

▼ 1.1 Background

Every year, the federal government releases large amounts of data on U.S. schools, school multiple datasets that are often difficult to access, and changes in data structure complicated (https://educationdata.urban.org/documentation/index.html (<a href="https://educationdata.urban.u

Using the Urban institutes consolidated data platform we will be combining datasets describing characteristics at the school district level.

We will use this data to build a classification model which classifies high schools as either low and high graduation rates is based on the federal government's standard that those s graduation rate schools.

▼ 1.2 Limitations

- Year: For our examination we only used data from 2015 as this year was the most dat
- · Our data is limited to only those schools which reported their graduation rates

1.3 Problem Statement

Predict which schools have high and which schools have low high school graduation rates and low graduation rates so that school districts know where to focus resources when att

2 Import Libraries

- 3.3.0 ELICOUE CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
import pandas as pd
In [1]:
        import matplotlib.pyplot as plt
        import matplotlib.ticker as mtic
        import seaborn as sns
        import numpy as np
        import scipy.stats as stats
        import statsmodels.api as sm
        import catboost
        import time
        import warnings
        warnings.filterwarnings('ignore')
        from sklearn.utils import class weight
        from sklearn.metrics import accuracy score, confusion matrix, class
        from catboost import Pool, sum models
        from catboost import CatBoostClassifier
        from statsmodels.formula.api import ols
        from sklearn.feature selection import RFE
        from sklearn.linear model import LinearRegression
        from sklearn.linear model import LogisticRegression
        from sklearn.metrics import mean squared error, r2 score, mean abso
        from sklearn.model selection import KFold, cross val score, Stratif
        from sklearn.model selection import train test split
        from sklearn.model selection import GridSearchCV
        from sklearn.model selection import RandomizedSearchCV
        from sklearn.preprocessing import LabelEncoder, OneHotEncoder, Mir
        from sklearn.tree import DecisionTreeRegressor
        from sklearn.ensemble import RandomForestRegressor, RandomForestCla
        from sklearn.ensemble import GradientBoostingClassifier
        from sklearn import metrics
        from sklearn.preprocessing import OneHotEncoder
        from sklearn.preprocessing import LabelEncoder
        from sklearn.impute import SimpleImputer, MissingIndicator
        from sklearn.compose import ColumnTransformer
        from sklearn.pipeline import FeatureUnion
        from sklearn.pipeline import Pipeline
        from scipy.stats import uniform, truncnorm, randint
        executed in 2.99s, finished 14:34:07 2021-06-17
```

3 Data Exploration and Prep

Contents **₽** ♥

3.3.0 Encode Charter

3.3.7 Reading and Math Tests

3.3.8 Numeric Columns - nulls

3.3.9 Visualizations

▼ 3.4 df_district

▼ 3.4.1 Data Fields

3.4.1.1 Data Field Cleanup

3.4.2 Columns to drop

3.4.3 english_language_learners

▼ 4 Join Datasets

4.1 Check for NaNs

▼ 5 Numeric Columns - cleanup

5.1 Cleanup

5.2 Engineer % Columns

5.3 Reset Num List

▼ 6 Category Columns - cleanup

6.1 Final Check for NaNs

7 Train Test Split

▼ 8 Encode Features

8.1 X train Encode

8.2 X test Encode

▼ 9 Model Development

▼ 9.1 Logistic Regression

9.1.1 Check for Overfit

9.1.2 Model reiteration - parameter tuning

▼ 9.2 Random Forest Classifier

9.2.1 Check for Overfit

9.2.2 Feature Importance

9.2.3 Model reiteration - parameter tunin

▼ 9.3 Gradient Boosting Classifier

9.3.1 Check for Overfit

9.3.2 Model reiteration - parameter tunin

10 Conclusions

Our school level data contains information on the school location, degree of urbanization, assessments, number of allegations for harassment/bullying, number of students enrolled number of students participation in ACT/SAT tests.

Our district level data contains financial information for each school district. This includes of the fiscal year, district expenditures, and revenue. District data also includes number of instruction.

3.1 Data Load

In [2]: #Northeastern United States school and district data
 df_school_northeast = pd.read_csv('data/EducationDataPortal_schools
 df_district_northeast = pd.read_csv('data/EducationDataPortal_distr
 executed in 125ms, finished 14:34:07 2021-06-17

In [4]: #Midwest United States school and district data
 df_school_midwest = pd.read_csv('data/EducationDataPortal_schools_n
 df_district_midwest = pd.read_csv('data/EducationDataPortal_distric
 executed in 183ms, finished 14:34:08 2021-06-17

In [5]: #Southern Atlantic United States school and district data
df_school_south_atl = pd.read_csv('data/EducationDataPortal_schools
df_district_south_atl = pd.read_csv('data/EducationDataPortal_distr
executed in 61ms, finished 14:34:08 2021-06-17

In [6]: #Southern Central United States school and district data
df_school_south_central = pd.read_csv('data/EducationDataPortal_school_district_south_central = pd.read_csv('data/EducationDataPortal_central_educationDataPortal_central_educationDataPortal_central_educationDataPortal_central_educationDataPortal_central_educationDataPortal_central_educationDataPortal_central_educationDataPortal_educat

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

3.2 Functions

3.3 df_school

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
df school['state name'].value counts(normalize=True, dropna=False)
In [11]:
          executed in 19ms, finished 14:34:08 2021-06-17
                                    0.012206
          Oregon
          Louisiana
                                    0.011815
          Maryland
                                    0.010741
          New Mexico
                                    0.009423
          Idaho
                                    0.008935
          South Dakota
                                    0.008886
          Utah
                                    0.008691
          Montana
                                    0.008349
          North Dakota
                                    0.008154
          West Virginia
                                    0.005761
          Nevada
                                    0.005517
          Maine
                                    0.005371
          New Hampshire
                                    0.004589
          Wyoming
                                    0.004101
          Vermont
                                    0.002929
          Rhode Island
                                    0.002734
          Delaware
                                    0.001855
          District of Columbia
                                    0.001660
          Name: state name, dtype: float64
In [12]:
         df school['grad rate midpt'].value counts(normalize=True, dropna=Fa
          executed in 8ms, finished 14:34:08 2021-06-17
Out[12]: 97
                  0.106288
          NaN
                  0.104726
          92
                  0.091593
          95
                  0.082219
          90
                  0.061664
          35
                  0.000098
          23
                  0.000049
          26
                  0.000049
          16
                  0.000049
          20
                  0.000049
          Name: grad rate midpt, Length: 99, dtype: float64
```

In [13]: df school

executed in 38ms, finished 14:34:08 2021-06-17

Out[13]:

urban_cei	zip_location	lea_name	state_name	school_name	ncessch	year	
(6106	Connecticut Technical High Sc	Connecticut	A. I. Prince Technical High School	90000201136	2015	0
(6610	Connecticut Technical High Sc	Connecticut	Bullard- Havens Technical High School	90000201137	2015	1
S	6053	Connecticut Technical High Sc	Connecticut	E. C. Goodwin Technical High School	90000201138	2015	2
	6340	Connecticut Technical High Sc	Connecticut	Ella T. Grasso Southeastern Technical High School	90000201139	2015	3
S	6514	Connecticut Technical	Connecticut	Eli Whitney Technical	90000201140	2015	4

Litala Oa

Litaria Ondanal

3.3.1 Data Fields

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
- ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
- ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
- ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
- 10 Conclusions

- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

141:	df	_school	head(,
] -			, (. ,

In [

						_	1 -
		COMMONICAL	17	4:34:08 2021-06-	20ms, finished 1	cuted in	exe
City, r	6106	Technical High Sc	Connecticut	Technical High School	90000201136	2015	0
City, r	6610	Connecticut Technical High Sc	Connecticut	Bullard- Havens Technical High School	90000201137	2015	1
Suburt	6053	Connecticut Technical High Sc	Connecticut	E. C. Goodwin Technical High School	90000201138	2015	2
Rural	6340	Connecticut Technical High Sc	Connecticut	Ella T. Grasso Southeastern Technical High School	90000201139	2015	3
Suburt	6514	Connecticut Technical High Sc	Connecticut	Eli Whitney Technical High School	90000201140	2015	4

5 rows × 64 columns

executed in 31ms, finished 14:34:08 2021-06-17

Contents & &
3.3.7 Reading and Math Tests
3.3.8 Numeric Columns - nulls
3.3.9 Visualizations
▼ 3.4 df_district
▼ 3.4.1 Data Fields
3.4.1.1 Data Field Cleanup
3.4.2 Columns to drop
3.4.3 english_language_learners
▼ 4 Join Datasets
4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
5.1 Cleanup
5.2 Engineer % Columns
5.3 Reset Num List
▼ 6 Category Columns - cleanup
6.1 Final Check for NaNs
7 Train Test Split
▼ 8 Encode Features
8.1 X_train Encode
8.2 X_test Encode
▼ 9 Model Development
▼ 9.1 Logistic Regression
9.1.1 Check for Overfit
9.1.2 Model reiteration - parameter tuning
▼ 9.2 Random Forest Classifier
9.2.1 Check for Overfit
9.2.2 Feature Importance
9.2.3 Model reiteration - parameter tunin
▼ 9.3 Gradient Boosting Classifier
9.3.1 Check for Overfit
9.3.2 Model reiteration - parameter tunin10 Conclusions

```
In [15]: df school.info()
         executed in 52ms, finished 14:34:08 2021-06-17
          28 transfers_alt_sch_disc
                                                18961 non-null float64
              days suspended
                                                18549 non-null float64
          30
              suspensions instances preschool 18549 non-null float64
          31 suspensions instances
                                                18549 non-null float64
          32 corpinstances preschool
                                                18549 non-null float64
          33 corpinstances
                                                18549 non-null float64
          34 salaries teachers
                                                18549 non-null object
                                                18337 non-null float64
          35 cohort num
              grad rate high
                                                18337 non-null object
          36
          37
              grad rate low
                                                18337 non-null object
          38
              grad rate midpt
                                                18337 non-null object
              allegations harass sex
                                                18961 non-null object
              allegations_harass_race
          40
                                                18961 non-null object
          41 allegations_harass_disability
                                                18961 non-null object
          42 allegations harass orientation
                                                18961 non-null object
              allegations harass religion
                                                18961 non-null object
          44 students disc harass dis
                                                18961 non-null float64
              students disc harass race
                                                18961 non-null float64
              students_disc_harass_sex
                                                                 float64
                                                18961 non-null
          47 students report harass dis
                                                18961 non-null float64
In [16]: df school.shape
         executed in 3ms, finished 14:34:08 2021-06-17
Out[16]: (20482, 64)
         3.3.1.1 Data Field Cleanup
In [17]: | column names = list(df school.columns)
         executed in 1ms, finished 14:34:08 2021-06-17
In [18]: for c in column names:
             df school[c].replace('Suppressed data', np.NaN,inplace =True)
```

Contents *⊋* ❖

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

▼ 3.3.2 Columns to drop

In [21]: |df_school.info()

executed in 31ms, finished 14:34:08 2021-06-17

<class 'pandas.core.frame.DataFrame'> Int64Index: 20482 entries, 0 to 3966

# Column Non-Null Count Dtype	
0 year 20482 non-null int64 1 ncessch 20482 non-null int64 2 school_name 20482 non-null object 3 state_name 20482 non-null object 4 lea_name 20482 non-null object 5 zip_location 19473 non-null object 6 urban_centric_locale 20457 non-null object 7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
1 ncessch 20482 non-null int64 2 school_name 20482 non-null object 3 state_name 20482 non-null object 4 lea_name 20482 non-null object 5 zip_location 19473 non-null object 6 urban_centric_locale 20457 non-null object 7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
2 school_name 20482 non-null object 3 state_name 20482 non-null object 4 lea_name 20482 non-null object 5 zip_location 19473 non-null object 6 urban_centric_locale 20457 non-null object 7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
3 state_name 20482 non-null object 4 lea_name 20482 non-null object 5 zip_location 19473 non-null object 6 urban_centric_locale 20457 non-null object 7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
4 lea_name 20482 non-null object 5 zip_location 19473 non-null object 6 urban_centric_locale 20457 non-null object 7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
5 zip_location 19473 non-null object 6 urban_centric_locale 20457 non-null object 7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
6 urban_centric_locale 20457 non-null object 7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
7 school_level 20482 non-null object 8 school_type 20482 non-null object 9 charter 19368 non-null object	
8 school_type 20482 non-null object 9 charter 19368 non-null object	
9 charter 19368 non-null object	
•	
10 enrollment 20001 non-null object	
11 read_test_num_valid 18073 non-null float6	
12 read_test_pct_prof_low 17434 non-null object	
13 read_test_pct_prof_high 17434 non-null object	
14 read_test_pct_prof_midpt 17434 non-null object	
15 math_test_num_valid 18132 non-null float6	4
16 math_test_pct_prof_low 17432 non-null object	,
17 math_test_pct_prof_high 17432 non-null object	
18 math_test_pct_prof_midpt 17432 non-null object	
19 students_susp_in_sch 18732 non-null object	
20 students_susp_out_sch_single 18898 non-null object	
21 students_susp_out_sch_multiple 18898 non-null object	
22 expulsions_no_ed_serv 18898 non-null object	
23 expulsions_with_ed_serv 18732 non-null object	
24 expulsions_zero_tolerance 18898 non-null object	
25 students corporal punish 954 non-null object	
26 students_arrested 18305 non-null object	
27 students referred law enforce 18728 non-null object	
28 transfers alt sch disc 18961 non-null float6	4
29 days_suspended 18549 non-null float6	4
30 suspensions_instances_preschool 18549 non-null float6	
31 suspensions_instances 18549 non-null float6	
32 corpinstances preschool 18549 non-null float6	
33 corpinstances 18549 non-null float6	
34 salaries teachers 18049 non-null object	

- 3.3.0 Elicode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

Contents 2 🌣 3.3.0 Encode Charter 3.3.7 Reading and Math Tests 3.3.8 Numeric Columns - nulls 3.3.9 Visualizations ▼ 3.4 df_district ▼ 3.4.1 Data Fields 3.4.1.1 Data Field Cleanup 3.4.2 Columns to drop 3.4.3 english_language_learners ▼ 4 Join Datasets 4.1 Check for NaNs ▼ 5 Numeric Columns - cleanup 5.1 Cleanup 5.2 Engineer % Columns 5.3 Reset Num List ▼ 6 Category Columns - cleanup 6.1 Final Check for NaNs 7 Train Test Split ▼ 8 Encode Features 8.1 X train Encode 8.2 X test Encode ▼ 9 Model Development ▼ 9.1 Logistic Regression 9.1.1 Check for Overfit 9.1.2 Model reiteration - parameter tuning ▼ 9.2 Random Forest Classifier 9.2.1 Check for Overfit 9.2.2 Feature Importance 9.2.3 Model reiteration - parameter tunin ▼ 9.3 Gradient Boosting Classifier 9.3.1 Check for Overfit 9.3.2 Model reiteration - parameter tunin 10 Conclusions

```
18337 non-null
                                                      float64
    cohort num
 36
    grad rate high
                                      17812 non-null
                                                      object
 37
    grad rate low
                                      17812 non-null
                                                      object
 38
    grad rate midpt
                                      17812 non-null
                                                      object
 39
     allegations harass sex
                                      18836 non-null
                                                      object
    allegations harass race
                                      18836 non-null
                                                      object
     allegations harass disability
                                      18836 non-null
 41
                                                      object
     allegations harass orientation
                                      18817 non-null
                                                      object
     allegations harass religion
                                      18817 non-null
                                                      object
    students disc harass dis
                                      18961 non-null
                                                      float64
    students disc_harass_race
 45
                                      18961 non-null
                                                      float64
    students_disc_harass_sex
 46
                                      18961 non-null
                                                      float64
    students report harass dis
                                      18961 non-null
                                                      float64
    students report harass race
                                      18961 non-null
                                                      float64
 49
    students report harass sex
                                      18961 non-null
                                                      float64
 50
    enrl biology
                                      17067 non-null
                                                      float64
    enrl chemistry
 51
                                      15532 non-null
                                                      float64
 52
    enrl advanced math
                                      14161 non-null float64
 53
    enrl calculus
                                      11670 non-null
                                                      float64
 54
    enrl algebra2
                                      16321 non-null float64
    enrl physics
                                      13303 non-null
                                                      float64
    enrl geometry
                                                      float64
 56
                                      16722 non-null
    instances mech restraint
                                      18424 non-null float64
 57
    instances phys restraint
                                      18827 non-null float64
    instances_seclusion
                                      18827 non-null
                                                      float64
    students mech restraint
                                      18429 non-null
                                                      float64
 61 students phys_restraint
                                      18813 non-null float64
 62 students_seclusion
                                      18828 non-null
                                                      float64
 63 students SAT ACT
                                      18961 non-null
                                                     float64
dtypes: float64(29), int64(2), object(33)
memory usage: 10.2+ MB
```

```
Contents 2 🌣
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
```

```
capstone-graduation-analysis - Jupyter Notebook
In [23]: df school['students corporal punish'].value counts(normalize=True,
           executed in 7ms, finished 14:34:08 2021-06-17
Out[23]: NaN
                     0.953423
                     0.007714
           2
                     0.005908
                     0.003955
                     0.002246
           78
                     0.000049
           171
                     0.000049
           1653
                     0.000049
           111
                     0.000049
           120
                     0.000049
           Name: students_corporal_punish, Length: 116, dtype: float64

    Drop year as this is the same for all rows

    Drop corpinstances_preschool this column is mostly -1 which is not interpretable

In [24]: df_school = df_school.drop(['year','corpinstances preschool','stude
           executed in 11ms, finished 14:34:08 2021-06-17
```

3.3.3 school_type

9.3.2 Model reiteration - parameter tunin

10 Conclusions

Contents *⊋* ❖

- 3.3.0 Elicone Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

In [26]:	<pre>df_school_regular = df_school.loc[df_school['school_type'] == 'Regu</pre>
	df_school_regular

execute	executed in 39ms, finished 14:34:08 2021-06-17								
38	90001800834	Explorations	Connecticut	DISTRICT	Town, fringe				
	•••								
3960	484668005347	YSLETA H S	Texas	YSLETA ISD	City, large				
3962	484668006496	DEL VALLE H S	Texas	YSLETA ISD	City, large				
3964	484668011982	VALLE VERDE EARLY COLLEGE H S	Texas	YSLETA ISD	City, large				
3965	484671005351	ZAPATA H S	Texas	ZAPATA COUNTY ISD	Rural, fringe				
3966	489913005627	BENAVIDES SECONDARY	Texas	BENAVIDES ISD	Rural, remote				

16365 rows × 60 columns

Out[27]:

Out[27]:							
		ncessch	school_name	state_name	lea_name	urban_centric_locale	school_l
	16	90000300343	Bridgeport Correctional Center	Connecticut	UNIFIED SCHOOL DISTRICT #1	City, midsize	
	17	90000300344	Brooklyn Correctional Institution	Connecticut	UNIFIED SCHOOL DISTRICT #1	Rural, fringe	
	18	90000300347	Cheshire Correctional Institution	Connecticut	UNIFIED SCHOOL DISTRICT #1	Suburb, large	
	19	90000300361	Hartford Correctional Center	Connecticut	UNIFIED SCHOOL DISTRICT #1	City, midsize	
	20	90000300374	New Haven Correctional Center	Connecticut	UNIFIED SCHOOL DISTRICT #1	City, midsize	
	3932	484578000736	HARRELL ACCELERATED LEARNING CENTER	Texas	WICHITA FALLS ISD	City, midsize	
	3936	484578009422	WICHITA COUNTY JUVENILE JUSTICE AEP	Texas	WICHITA FALLS ISD	City, midsize	
	3954	484668003897	CESAR CHAVEZ ACADEMY	Texas	YSLETA ISD	City, large	
	3961	484668005938	TEJAS SCHOOL OF CHOICE	Texas	YSLETA ISD	City, large	

Contents <i>⊋</i> ❖
3.3.0 Encode Charter
3.3.7 Reading and Math Tests
3.3.8 Numeric Columns - nulls
3.3.9 Visualizations
▼ 3.4 df_district
▼ 3.4.1 Data Fields
3.4.1.1 Data Field Cleanup
3.4.2 Columns to drop
3.4.3 english_language_learners
▼ 4 Join Datasets
4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
5.1 Cleanup
5.2 Engineer % Columns
5.3 Reset Num List
▼ 6 Category Columns - cleanup
6.1 Final Check for NaNs
7 Train Test Split
▼ 8 Encode Features
8.1 X_train Encode
8.2 X_test Encode
▼ 9 Model Development
▼ 9.1 Logistic Regression
9.1.1 Check for Overfit
9.1.2 Model reiteration - parameter tuning

▼ 9.2 Random Forest Classifier
 9.2.1 Check for Overfit
 9.2.2 Feature Importance

▼ 9.3 Gradient Boosting Classifier 9.3.1 Check for Overfit

10 Conclusions

9.2.3 Model reiteration - parameter tunin

9.3.2 Model reiteration - parameter tunin

	ncessch	school_name	state_name	lea_name	urban_centric_locale	school_l
3963	484668008548	PLATO ACADEMY	Texas	YSLETA ISD	City, large	

2797 rows × 60 columns

- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
- ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
- ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
- ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
- 10 Conclusions

In [28]: df_school_vocational = df_school.loc[df_school['school_type'] == '\
 df_school_vocational.head(10)
 executed in 76ms, finished 14:34:08 2021-06-17

Out[28]:		noossah	sobool name	state name	loa namo	urban contrio localo	school level
		ncessch	school_name	state_name	lea_name	urban_centric_locale	SCHOOL_level
	0	90000201136	A. I. Prince Technical High School	Connecticut	Connecticut Technical High Sc	City, midsize	High
	1	90000201137	Bullard- Havens Technical High School	Connecticut	Connecticut Technical High Sc	City, midsize	High
	2	90000201138	E. C. Goodwin Technical High School	Connecticut	Connecticut Technical High Sc	Suburb, large	High
	3	90000201139	Ella T. Grasso Southeastern Technical High School	Connecticut	Connecticut Technical High Sc	Rural, fringe	High
	4	90000201140	Eli Whitney Technical High School	Connecticut	Connecticut Technical High Sc	Suburb, large	High
	5	90000201141	Emmett OBrien Technical High School	Connecticut	Connecticut Technical High Sc	Suburb, large	High
	6	90000201142	H. C. Wilcox Technical High School	Connecticut	Connecticut Technical High Sc	Suburb, large	High
	7	90000201143	H. H. Ellis Technical High School	Connecticut	Connecticut Technical High Sc	Suburb, large	High
	8	90000201144	Henry Abbott Technical High School	Connecticut	Connecticut Technical High Sc	City, small	High
	9	90000201145	Howell Cheney Technical High School	Connecticut	Connecticut Technical High Sc	Suburb, large	High

Contents	\mathcal{C}	*
----------	---------------	---

- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
- ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
- ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
- ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
- 10 Conclusions

10 rows × 60 columns

- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

Out[29]:

out[29]:						
		ncessch	school_name	state_name	lea_name	urban_centric_locale
	46	90007001588	Center for Autism Spectrum and Development Dis	Connecticut	AREA COOPERATIVE EDUCATIONAL	Suburb, large
	64	90051001515	Community/Vocational Program	Connecticut	BRISTOL SCHOOL DISTRICT	Suburb, large
	71	90070001567	Lincoln Academy Regional Interdistrict Collabo	Connecticut	CAPITOL REGION EDUCATION COUN	City, small
	73	90070001693	Soundbridge @ Wethersfield High School	Connecticut	CAPITOL REGION EDUCATION COUN	Suburb, large
	74	90070001802	STRIVE (Southern Transition Real-World and Ind	Connecticut	CAPITOL REGION EDUCATION COUN	Suburb, large
	77	90075001766	Cheshire Quinnipiac University Transition Coll	Connecticut	CHESHIRE SCHOOL DISTRICT	Suburb, large
	92	90123001514	The Learning Center at East Hampton	Connecticut	EAST HAMPTON SCHOOL DISTRICT	Town, fringe
	98	90132001730	Post High School Transition Program	Connecticut	EAST LYME SCHOOL DISTRICT	Suburb, midsize
	100	90132001843	Medically Fragile Program (WAVES)	Connecticut	EAST LYME SCHOOL DISTRICT	NaN
	108	90147001523	Enfield Transitional Learning Academy	Connecticut	ENFIELD SCHOOL DISTRICT	Suburb, large

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- **▼** 9 Model Development
- ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
- ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
- ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
- 10 Conclusions

10 rows × 60 columns

```
Contents 2 🌣
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [30]: #Proportion of nulls regular school type
          percent null df(df school regular, 'df school regular')
          executed in 39ms, finished 14:34:08 2021-06-17
Out[30]: 'percent of nulls in df school regular is 4%'
In [31]: #Proportion of nulls vocational school type
          percent null df(df school vocational, 'df school vocational')
          executed in 8ms, finished 14:34:08 2021-06-17
Out[31]: 'percent of nulls in df school vocational is 61%'
In [32]: #Proportion of nulls alternative school type
          percent null df(df school alternative, 'df school alternative')
          executed in 12ms, finished 14:34:08 2021-06-17
Out[32]: 'percent of nulls in df school alternative is 24%'
In [33]: #Proportion of nulls special education school type
          percent_null_df(df_school_special')
          executed in 6ms, finished 14:34:08 2021-06-17
Out[33]: 'percent of nulls in df school special is 38%'
```

Our Other/alternative school in the school_type field includes jails and detention centers. In other data fields. We will remove all rows with alternative schools since there seems to be students in these types of schools are experiencing circumstances and educational exper

We will also remove vocational and special education school types as these school types nulls

```
In [34]: df_school.shape[0]
executed in 3ms, finished 14:34:08 2021-06-17

Out[34]: 20482
```

```
Contents ₽ ❖
```

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- **▼** 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

Out[35]: 16365

3.3.4 Graduation Rate - format the target

- The graduation rates are numerical values separated into three columns.
 - grad rate high
 - grad_rate_midpt
 - grad_rate_low
- Looking at the data documentation from EDFacts https://www2.ed.gov/about/inits/ed/edfacts/index.htmlThe) we see that the reason for graduation rates in a range so that student privacy is protected. The range creates a limit middle of the low and high values.
- We will use grad_rate_midpt to measure graduation rates as this is the most balanced
- The federal government defines graduation rates as low when less than 2/3 of a coho
- We will classify high schools with midpoint graduation rates 66 and below as low and high.

In [37]: df_school[grad_rate_cols].head(10)

executed in 9ms, finished 14:34:08 2021-06-17

Out[37]:

	school_name	cohort_num	grad_rate_high	grad_rate_midpt	grad_rate_l
32	Walter G. Cady School	NaN	NaN	NaN	N
33	E. O. Smith High School	255.0	95	95	
36	Common Ground High School	34.0	100	95	
37	The Bridge Academy	34.0	79	74	
38	Explorations	27.0	79	69	
39	Connecticut Valley Hospital	NaN	NaN	NaN	N
40	Stamford Academy	53.0	29	24	
42	Ansonia High School	125.0	89	87	
47	Avon High School	239.0	97	97	
48	Berlin High School	248.0	95	95	

In [38]: df_school = df_school.drop(['grad_rate_high','grad_rate_low'], axis executed in 7ms, finished 14:34:08 2021-06-17

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [39]: | df school.info()
          executed in 25ms, finished 14:34:08 2021-06-17
          <class 'pandas.core.frame.DataFrame'>
          Int64Index: 16365 entries, 32 to 3966
          Data columns (total 58 columns):
               Column
                                                   Non-Null Count
                                                                    Dtype
           0
                                                   16365 non-null
                                                                    int64
               ncessch
           1
                                                   16365 non-null
                                                                    object
               school name
                                                   16365 non-null
                                                                    object
               state_name
           3
               lea name
                                                   16365 non-null
                                                                    object
               urban centric locale
                                                   16364 non-null
                                                                    object
               school level
                                                   16365 non-null
                                                                    object
               school type
                                                   16365 non-null
                                                                    object
               charter
                                                   15404 non-null
                                                                    object
           8
               enrollment
                                                   16328 non-null
                                                                    object
           9
               read test num valid
                                                   15747 non-null
                                                                    float64
               read test pct prof low
                                                   15500 non-null
                                                                    object
           11
               read test pct prof high
                                                   15500 non-null
                                                                    object
           12
               read test pct prof midpt
                                                   15500 non-null
                                                                    object
               math_test_num_valid
                                                   15819 non-null
                                                                    float64
In [40]:
         df school.update(df school[['grad rate midpt']].fillna(0))
          executed in 17ms, finished 14:34:08 2021-06-17
         df school['grad rate midpt'].value counts(normalize=True,dropna=Fal
In [41]:
          executed in 7ms, finished 14:34:09 2021-06-17
Out[41]: 97
                0.127345
          92
                0.112679
          95
                0.100214
          90
                0.073816
          87
                0.071983
          25
                0.000061
          38
                0.000061
          0
                0.000061
          18
                0.000061
          23
                0.000061
          Name: grad rate midpt, Length: 92, dtype: float64
```

```
Contents 2 🌣
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [42]: df school.shape
           executed in 3ms, finished 14:34:09 2021-06-17
Out[42]: (16365, 58)
In [43]: 19198*(1-.43)
           executed in 4ms, finished 14:34:09 2021-06-17
Out[43]: 10942.86
In [44]: df_school = df_school[df_school.grad_rate midpt != 0]
           executed in 12ms, finished 14:34:09 2021-06-17
In [45]: df_school.shape[0]
           executed in 3ms, finished 14:34:09 2021-06-17
Out[45]: 15830
In [46]: # convert grad rate columns to numeric
           df school[['grad rate midpt']] = df school[['grad rate midpt']].ast
           executed in 9ms, finished 14:34:09 2021-06-17
In [47]: # Let's take a look at the summary statistics for the grad rate mic
          print(df school['grad rate midpt'].describe())
           executed in 6ms, finished 14:34:09 2021-06-17
                     15830.000000
           count
                        85.590587
          mean
           std
                        15.379323
          min
                          0.000000
           25%
                        82.000000
           50%
                        90.000000
           75%
                        95.000000
                         99.000000
          max
          Name: grad rate midpt, dtype: float64
```

```
Contents ⊋ ‡
```

3.3.0 Encode Charter

3.3.7 Reading and Math Tests

3.3.8 Numeric Columns - nulls

3.3.9 Visualizations

▼ 3.4 df_district

▼ 3.4.1 Data Fields

3.4.1.1 Data Field Cleanup

3.4.2 Columns to drop

3.4.3 english_language_learners

▼ 4 Join Datasets

4.1 Check for NaNs

▼ 5 Numeric Columns - cleanup

5.1 Cleanup

5.2 Engineer % Columns

5.3 Reset Num List

▼ 6 Category Columns - cleanup

6.1 Final Check for NaNs

7 Train Test Split

▼ 8 Encode Features

8.1 X train Encode

8.2 X test Encode

▼ 9 Model Development

▼ 9.1 Logistic Regression

9.1.1 Check for Overfit

9.1.2 Model reiteration - parameter tuning

▼ 9.2 Random Forest Classifier

9.2.1 Check for Overfit

9.2.2 Feature Importance

9.2.3 Model reiteration - parameter tunin

▼ 9.3 Gradient Boosting Classifier

9.3.1 Check for Overfit

9.3.2 Model reiteration - parameter tunin

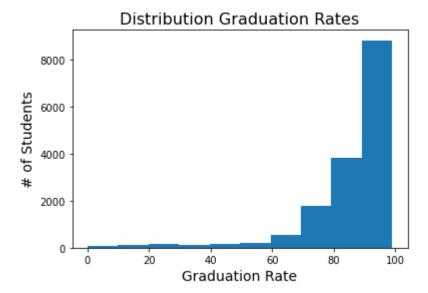
10 Conclusions

```
grad_rate_midpt
  97    0.131649
92    0.116488
95    0.103601
90    0.076311
87    0.074416
Name: grad_rate_midpt, dtype: float64
```

executed in 142ms, finished 14:34:09 2021-06-17

```
In [49]: # Plot of the target status group

plt.hist(df_school['grad_rate_midpt'])
plt.xlabel('Graduation Rate', fontsize=14)
plt.ylabel('# of Students', fontsize=14)
plt.title("Distribution Graduation Rates", fontsize=16)
plt.show()
```



It looks like there is a strong skew towards higher performing high schools. Lets go ahead again.

▼ 3.3.4.1 Bin the target

As stated above based on how the federal government defines low graduation rates we w

- Low: 66 and belowHigh: 67 and above
- riigiii or ana abor
- In [50]: #Based on the federal government's definition we will define the bu
 grad_rate_bins = [0,66,100]
 executed in 3ms, finished 14:34:09 2021-06-17
- In [52]: # and we will also take a look at the distribution across bins
 print(df_school['grad_rate_midpt'].value_counts(normalize=True, dro
 executed in 7ms, finished 14:34:09 2021-06-17

 (66.0, 100.0] 0.931522

```
(0.0, 66.0] 0.068414

NaN 0.000063

Name: grad_rate_midpt, dtype: float64
```

- In [53]: df_school = df_school.dropna(subset=['grad_rate_midpt'])
 df_school.shape
 executed in 11ms, finished 14:34:09 2021-06-17
- Out[53]: (15829, 58)

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- **▼** 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
Contents 2 &
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

10 Conclusions

```
In [54]: print(df school['grad rate midpt'].value counts(normalize=True, dro
          executed in 5ms, finished 14:34:09 2021-06-17
           (66, 1001
                         0.931581
           (0, 661)
                         0.068419
          Name: grad_rate_midpt, dtype: float64
In [55]:
          # We will need to address this class imbalance
           executed in 2ms, finished 14:34:09 2021-06-17
          3.3.4.2 Encode the target
In [56]: df school['grad rate midpt'].value counts(normalize=True, dropna=Fa
           executed in 6ms, finished 14:34:09 2021-06-17
Out[56]: (66, 100]
                         0.931581
           (0, 661
                         0.068419
          Name: grad_rate_midpt, dtype: float64
In [57]:
           df school['grad rate midpt'] = df school['grad rate midpt'].astype
           executed in 51ms, finished 14:34:09 2021-06-17
In [58]: df school['grad rate midpt'].value counts(normalize=True, dropna=Fa
          executed in 7ms, finished 14:34:09 2021-06-17
Out[58]: (66, 100]
                         0.931581
           (0, 661
                         0.068419
          Name: grad_rate_midpt, dtype: float64
In [59]:
           df school['grad rate midpt'] = df school['grad rate midpt'].replac
           executed in 4ms, finished 14:34:09 2021-06-17
           df school['grad rate midpt'] = df school['grad rate midpt'].replac
In [60]:
           executed in 6ms, finished 14:34:09 2021-06-17
```

3.3.5 subject enrollment - drop all enri columns

```
In [61]: #Let's look at how many nulls we are dealing with now
percent_null_df(df_school,'df_school')
executed in 36ms, finished 14:34:09 2021-06-17
```

Out[61]: 'percent of nulls in df school is 3%'

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

In [62]: df school.info()

executed in 25ms, finished 14:34:09 2021-06-17

<class 'pandas.core.frame.DataFrame'> Int64Index: 15829 entries, 33 to 3966 Data columne (total 58 columne).

columns (total 58 columns):			
Column	Non-Nu	ıll Count	Dtype
ncessch	15829	non-null	 int64
			object
_			object
_			object
_			object
			object
_			object
_			object
			object
			float64
			object
			object
			object
			float64
	15343	non-null	object
	15343	non-null	object
math_test_pct_prof_midpt	15343	non-null	object
students_susp_in_sch	15583	non-null	object
students_susp_out_sch_single	15598	non-null	object
students_susp_out_sch_multiple	15598	non-null	object
expulsions_no_ed_serv	15598	non-null	object
expulsions_with_ed_serv	15583	non-null	object
expulsions_zero_tolerance	15598	non-null	object
students_arrested	15218	non-null	object
students_referred_law_enforce	15579	non-null	object
transfers_alt_sch_disc	15617	non-null	float64
days_suspended	15265	non-null	float64
suspensions_instances_preschool	15265	non-null	float64
suspensions_instances	15265	non-null	float64
corpinstances	15265	non-null	float64
salaries_teachers	15107	non-null	object
cohort_num	15829	non-null	float64
<pre>grad_rate_midpt</pre>	15829	non-null	int64
allegations_harass_sex	15547	non-null	object
allegations_harass_race	15547	non-null	object
	column ncessch school_name state_name lea_name urban_centric_locale school_level school_type charter enrollment read_test_num_valid read_test_pct_prof_low read_test_pct_prof_high read_test_pct_prof_low math_test_pct_prof_low math_test_pct_prof_high math_test_pct_prof_midpt students_susp_in_sch students_susp_in_sch students_susp_out_sch_multiple expulsions_no_ed_serv expulsions_with_ed_serv expulsions_zero_tolerance students_arrested students_referred_law_enforce transfers_alt_sch_disc days_suspended suspensions_instances_preschool suspensions_instances corpinstances salaries_teachers cohort_num grad_rate_midpt allegations_harass_sex	Column Non-Non-Non-Non-Non-Non-Non-Non-Non-Non-	Column Non-Null Count

- 3.3.0 Elicode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- **▼** 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
allegations harass disability
                                      15547 non-null
                                                     object
    allegations harass orientation
                                      15533 non-null
                                                     object
    allegations harass religion
                                                     object
 37
                                      15533 non-null
 38 students disc harass dis
                                                     float64
                                      15617 non-null
    students disc harass race
                                      15617 non-null
                                                     float64
    students disc harass sex
                                                     float64
                                      15617 non-null
    students report harass dis
                                      15617 non-null
                                                     float64
    students_report_harass_race
                                                     float64
                                      15617 non-null
    students report harass sex
                                      15617 non-null float64
 44
    enrl biology
                                      14892 non-null float64
    enrl chemistry
 45
                                      14262 non-null float64
    enrl advanced math
                                      13362 non-null float64
 46
    enrl calculus
                                      11327 non-null float64
    enrl algebra2
                                      14580 non-null float64
 48
 49
    enrl physics
                                      12400 non-null float64
 50
    enrl geometry
                                      14621 non-null float64
    instances mech restraint
                                      15196 non-null float64
    instances phys restraint
 52
                                      15538 non-null float64
    instances seclusion
 53
                                      15538 non-null float64
 54 students mech restraint
                                      15197 non-null float64
    students phys restraint
                                      15526 non-null
                                                     float64
    students seclusion
                                      15539 non-null float64
    students SAT ACT
                                      15617 non-null float64
dtypes: float64(28), int64(2), object(28)
memory usage: 7.1+ MB
```

The class subject enrollment columns stand out as having the most null values let's drop

226 Engada Chartar

```
In [65]: df_school['charter'].value_counts(normalize=True, dropna=False)
           executed in 6ms, finished 14:34:09 2021-06-17
Out[65]: No
                   0.872386
           Yes
                   0.069556
                   0.058058
           NaN
           Name: charter, dtype: float64
            df_school['charter'] = df_school['charter'].replace('No', 0)
In [66]:
           executed in 3ms, finished 14:34:09 2021-06-17
           df_school['charter'] = df_school['charter'].replace('Yes', 1)
In [67]:
           executed in 8ms, finished 14:34:09 2021-06-17
In [68]: df_school['charter'].value_counts(normalize=True, dropna=False)
           executed in 6ms, finished 14:34:09 2021-06-17
Out[68]: 0.0
                   0.872386
           1.0
                   0.069556
           NaN
                   0.058058
           Name: charter, dtype: float64
           We will group the small proportion of NaN values with the majority class which is zero or a
In [69]: | df_school.update(df_school[['charter']].fillna(0))
           executed in 13ms, finished 14:34:09 2021-06-17
```

3.3.7 Reading and Math Tests

```
Contents 2 🌣
```

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
- ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
- ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
- ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
- 10 Conclusions

In [70]: df school.info()

executed in 22ms, finished 14:34:09 2021-06-17

<class 'pandas.core.frame.DataFrame'> Int64Index: 15829 entries, 33 to 3966 Data columns (total 51 columns):

Data #	columns (total 51 columns): Column	Non-Ni	ıll Count	Dtype
		1011-110		ьсуре
0	ncessch	15829	non-null	int64
1	school name		non-null	object
2	state name		non-null	object
3	lea name		non-null	object
4	urban_centric_locale		non-null	object
5	school_level	15829	non-null	object
6	school_type	15829	non-null	object
7	charter	15829	non-null	float64
8	enrollment	15822	non-null	object
9	read_test_num_valid	15478	non-null	float64
10	read_test_pct_prof_low	15305	non-null	object
11	read_test_pct_prof_high	15305	non-null	object
12	read_test_pct_prof_midpt	15305	non-null	object
13	math_test_num_valid	15552	non-null	float64
14	math_test_pct_prof_low	15343	non-null	object
15	math_test_pct_prof_high	15343	non-null	object
16	math_test_pct_prof_midpt	15343	non-null	object
17	students_susp_in_sch	15583	non-null	object
18	students_susp_out_sch_single	15598	non-null	object
19	students_susp_out_sch_multiple		non-null	object
20	expulsions_no_ed_serv		non-null	object
21	expulsions_with_ed_serv		non-null	object
22	expulsions_zero_tolerance		non-null	object
23	students_arrested		non-null	object
24	students_referred_law_enforce		non-null	object
25	transfers_alt_sch_disc		non-null	float64
26	days_suspended		non-null	float64
27	suspensions_instances_preschool		non-null	float64
28	suspensions_instances		non-null	float64
29	corpinstances		non-null	float64
30	salaries_teachers		non-null	object
31	cohort_num		non-null	float64
32	grad_rate_midpt		non-null	int64
33	allegations_harass_sex		non-null	object
34	allegations_harass_race	15547	non-null	object

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
allegations harass disability
                                      15547 non-null
                                                      object
    allegations harass orientation
                                      15533 non-null
                                                      object
    allegations harass religion
 37
                                      15533 non-null
                                                      object
    students disc harass dis
                                      15617 non-null
                                                      float64
    students disc harass race
                                      15617 non-null
                                                      float64
    students disc harass sex
                                                      float64
                                      15617 non-null
    students report harass dis
                                      15617 non-null
 41
                                                      float64
    students report harass race
                                                      float64
                                      15617 non-null
    students_report_harass_sex
                                      15617 non-null
                                                      float64
    instances mech restraint
                                      15196 non-null
                                                      float64
 45
    instances phys restraint
                                      15538 non-null
                                                      float64
 46
    instances_seclusion
                                      15538 non-null
                                                      float64
    students mech restraint
                                      15197 non-null
                                                      float64
    students phys restraint
                                      15526 non-null
                                                      float64
    students seclusion
                                      15539 non-null
                                                      float64
    students SAT ACT
                                      15617 non-null
                                                     float64
dtypes: float64(22), int64(2), object(27)
memory usage: 6.3+ MB
```

In [72]: df_school[read_test].head(20)
executed in 10ms, finished 14:34:09 2021-06-17

Out[72]:

	school_name	read_test_num_valid	read_test_pct_prof
33	E. O. Smith High School	260.0	
36	Common Ground High School	39.0	
37	The Bridge Academy	139.0	
38	Explorations	21.0	
40	Stamford Academy	25.0	
42	Ansonia High School	137.0	
47	Avon High School	244.0	
48	Berlin High School	231.0	
49	Path Academy	6.0	
50	Bethel High School	199.0	
51	Bloomfield High School	143.0	

executed in 2ms, finished 14:34:09 2021-06-17

Contents 2	₩.
------------	----

- 3.3.0 ELICOUE CHARLET
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

In [74]: df school[math test].head(20)

executed in 11ms, finished 14:34:09 2021-06-17

Out[74]:

	school_name	math_test_num_valid	math_test_pct_prc
33	E. O. Smith High School	260.0	
36	Common Ground High School	39.0	
37	The Bridge Academy	139.0	
38	Explorations	21.0	
40	Stamford Academy	25.0	
42	Ansonia High School	137.0	
47	Avon High School	246.0	
48	Berlin High School	231.0	
49	Path Academy	6.0	
50	Bethel High School	199.0	
51	Bloomfield High School	143.0	
52	Learning Academy at Bloomfield	1.0	
53	Bolton High School	67.0	
54	Branford High School	233.0	
56	Bassick High School	119.0	
57	Central High School	249.0	
58	Harding High School	176.0	
59	Information Technology and Software Engineerin	132.0	
60	Biotechnology Research and Zoological Studies	134.0	
61	Aerospace/Hydrospace Engineering and Physical \dots	130.0	

The read_test_num_valid column and the math_test_num_valid column desribe the Numb and for whom a proficiency level was assigned

3.3.0 Encode Charter

3.3.7 Reading and Math Tests

3.3.8 Numeric Columns - nulls

3.3.9 Visualizations

▼ 3.4 df_district

▼ 3.4.1 Data Fields

3.4.1.1 Data Field Cleanup

3.4.2 Columns to drop

3.4.3 english_language_learners

▼ 4 Join Datasets

4.1 Check for NaNs

▼ 5 Numeric Columns - cleanup

5.1 Cleanup

5.2 Engineer % Columns

5.3 Reset Num List

▼ 6 Category Columns - cleanup

6.1 Final Check for NaNs

7 Train Test Split

▼ 8 Encode Features

8.1 X_train Encode

8.2 X test Encode

▼ 9 Model Development

▼ 9.1 Logistic Regression

9.1.1 Check for Overfit

9.1.2 Model reiteration - parameter tuning

▼ 9.2 Random Forest Classifier

9.2.1 Check for Overfit

9.2.2 Feature Importance

9.2.3 Model reiteration - parameter tunin

▼ 9.3 Gradient Boosting Classifier

9.3.1 Check for Overfit

9.3.2 Model reiteration - parameter tunin

10 Conclusions

The low, high, and midpt columns describe the low, high, and midpoint of the range used language arts assessment (0–100 scale)

3.3.8 Numeric Columns - nulls

Let's take a look at our columns and check back in on the proportion of nulls. If our proportion of nulls. If our proportion of nulls.

```
In [76]: | df school.info()
         executed in 23ms, finished 14:34:09 2021-06-17
                                                IJUZJ HUH-HUIL IIUUU
          8
                                                15822 non-null
                                                                object
              enrollment
              read_test_num_valid
                                                15478 non-null float64
          10 read_test_pct_prof_midpt
                                                                object
                                                15305 non-null
          11 math test num valid
                                                15552 non-null
                                                                float64
              math test pct prof midpt
                                                                object
          12
                                                15343 non-null
          13
              students susp in sch
                                                15583 non-null
                                                                object
              students susp out sch single
                                                15598 non-null
                                                                object
              students susp out sch multiple
                                                15598 non-null
                                                                object
              expulsions no ed serv
          16
                                                15598 non-null
                                                                object
              expulsions with ed serv
                                                                object
          17
                                                15583 non-null
              expulsions zero tolerance
                                                15598 non-null
                                                                object
          18
          19
              students arrested
                                                15218 non-null
                                                                object
              students referred law enforce
                                                                object
                                                15579 non-null
              transfers alt sch disc
          21
                                                15617 non-null float64
              days suspended
          22
                                                15265 non-null
                                                                float64
              suspensions instances preschool 15265 non-null float64
              suspensions instances
                                                15265 non-null
                                                               float64
              corpinstances
                                                15265 non-null float64
                                                15107 non-null object
              calaries teachers
```

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- **▼** 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
In [77]: #How do our nulls look?
    percent_null_df(df_school,'df_school')
    executed in 30ms, finished 14:34:09 2021-06-17

Out[77]: 'percent of nulls in df_school is 1%'

In [78]: column_names = list(df_school.columns)
    executed in 2ms, finished 14:34:09 2021-06-17

In [79]: categorical = ['year', 'ncessch', 'school_name', 'state_name', 'lea_nam 'urban_centric_locale', 'school_level', 'school_type', 'corpinstances_preschool']

executed in 2ms, finished 14:34:09 2021-06-17
```

In [80]:

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
numerical = set(column names) - set(categorical)
         numerical = list(numerical)
         numerical
          executed in 4ms, finished 14:34:09 2021-06-17
Out[80]: ['days_suspended',
           'students disc harass sex',
           'expulsions with ed serv',
           'students arrested',
           'salaries teachers',
           'students seclusion',
           'instances phys restraint',
           'math test num valid',
           'students disc harass dis',
           'instances mech restraint',
           'math test pct prof midpt',
           'read test pct prof midpt',
           'corpinstances',
           'read_test_num_valid',
           'allegations harass race',
           'students susp out sch single',
           'students report harass race',
           'allegations harass disability',
           'students susp out sch multiple',
           'enrollment',
           'suspensions instances',
           'suspensions instances preschool',
           'expulsions_zero_tolerance',
           'students mech restraint',
           'transfers_alt_sch_disc',
           'cohort num',
           'allegations harass orientation',
           'students report harass sex',
           'allegations harass religion',
           'students phys restraint',
           'students SAT ACT',
           'students disc harass race',
           'instances seclusion',
           'grad rate midpt',
           'students report harass dis',
           'students susp in sch',
           'expulsions no ed serv',
```

We've dropped some of our numerical columns let's update our nume

```
'allegations_harass_sex',
'students_referred_law_enforce']
```

3.3.0 Encode Charter 3.3.7 Reading and Math Tests 3.3.8 Numeric Columns - nulls 3.3.9 Visualizations ▼ 3.4 df_district ▼ 3.4.1 Data Fields 3.4.1.1 Data Field Cleanup 3.4.2 Columns to drop 3.4.3 english_language_learners ▼ 4 Join Datasets 4.1 Check for NaNs ▼ 5 Numeric Columns - cleanup 5.1 Cleanup 5.2 Engineer % Columns 5.3 Reset Num List ▼ 6 Category Columns - cleanup 6.1 Final Check for NaNs 7 Train Test Split ▼ 8 Encode Features 8.1 X train Encode 8.2 X test Encode ▼ 9 Model Development ▼ 9.1 Logistic Regression 9.1.1 Check for Overfit 9.1.2 Model reiteration - parameter tuning ▼ 9.2 Random Forest Classifier 9.2.1 Check for Overfit 9.2.2 Feature Importance 9.2.3 Model reiteration - parameter tunin ▼ 9.3 Gradient Boosting Classifier 9.3.1 Check for Overfit 9.3.2 Model reiteration - parameter tunin 10 Conclusions

```
In [81]: numerical.remove('grad rate midpt')
          executed in 2ms, finished 14:34:09 2021-06-17
In [82]: nan values = df school.isna()
          nan columns = nan values.any()
          columns with nan = df school.columns[nan columns].tolist()
          print(columns_with_nan)
          executed in 16ms, finished 14:34:09 2021-06-17
          ['enrollment', 'read test num valid', 'read test pct prof midpt', '
          t', 'students_susp_in_sch', 'students_susp_out_sch single', 'studer
          rv', 'expulsions with ed serv', 'expulsions zero tolerance', 'stude
          'transfers alt sch disc', 'days suspended', 'suspensions instances
          nces', 'salaries teachers', 'allegations harass sex', 'allegations
          'allegations_harass_orientation', 'allegations_harass_religion', '&
          race', 'students disc harass sex', 'students report harass dis', '
          harass sex', 'instances mech restraint', 'instances phys restraint'
          nt', 'students phys restraint', 'students seclusion', 'students SAT
In [83]: #Replace NaNs with median for each column
          df_school[numerical] = df_school[numerical].fillna(df_school[numeri
          executed in 70ms, finished 14:34:09 2021-06-17
In [84]: nan values = df school.isna()
          nan columns = nan values.any()
          columns with nan = df school.columns[nan columns].tolist()
          print(columns with nan)
          executed in 16ms, finished 14:34:09 2021-06-17
          []
In [85]: df school[numerical] = df school[numerical].astype(str).astype(floating)
          executed in 242ms, finished 14:34:09 2021-06-17
```

```
In [86]: #How do our nulls look?
    percent_null_df(df_school,'df_school')
    executed in 26ms, finished 14:34:10 2021-06-17
Out[86]: 'percent of nulls in df_school is 0%'
```

3.3.9 Visualizations

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

In [87]: df school.info()

executed in 13ms, finished 14:34:10 2021-06-17

<class 'pandas.core.frame.DataFrame'> Int64Index: 15829 entries, 33 to 3966 Data columns (total 47 columns):

Data #	Columns (total 4/ columns):	Non-Ni	ull Count	Dtype
0	ncessch	15829	non-null	int64
1	school_name	15829	non-null	object
2	state_name	15829	non-null	object
3	lea_name	15829	non-null	object
4	urban_centric_locale	15829	non-null	object
5	school_level	15829	non-null	object
6	school_type	15829	non-null	object
7	charter	15829	non-null	float64
8	enrollment	15829	non-null	int64
9	read_test_num_valid	15829	non-null	int64
10	read_test_pct_prof_midpt	15829	non-null	int64
11	math_test_num_valid	15829	non-null	int64
12	math_test_pct_prof_midpt	15829	non-null	int64
13	students_susp_in_sch	15829	non-null	int64
14	students_susp_out_sch_single	15829	non-null	int64
15	students_susp_out_sch_multiple	15829	non-null	int64
16	expulsions_no_ed_serv	15829	non-null	int64
17	expulsions_with_ed_serv	15829	non-null	int64
18	expulsions_zero_tolerance	15829	non-null	int64
19	students_arrested	15829	non-null	int64
20	students_referred_law_enforce	15829	non-null	int64
21	transfers_alt_sch_disc	15829	non-null	int64
22	days_suspended	15829	non-null	int64
23	suspensions_instances_preschool	15829	non-null	int64
24	suspensions_instances	15829	non-null	int64
25	corpinstances	15829	non-null	int64
26	salaries_teachers	15829	non-null	int64
27	cohort_num	15829	non-null	int64
28	<pre>grad_rate_midpt</pre>	15829	non-null	int64
29	allegations_harass_sex	15829	non-null	int64
30	allegations_harass_race	15829	non-null	int64
31	allegations_harass_disability	15829	non-null	int64
32	allegations_harass_orientation		non-null	int64
33	allegations_harass_religion		non-null	int64
34	students_disc_harass_dis	15829	non-null	int64

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

Contents *⊋* ❖

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
students disc harass race
                                     15829 non-null
                                                     int64
   students disc harass sex
                                     15829 non-null
                                                     int64
   students report harass dis
37
                                     15829 non-null
                                                     int64
   students_report_harass_race
                                     15829 non-null
                                                     int64
   students report harass sex
                                     15829 non-null
                                                     int64
   instances mech restraint
                                     15829 non-null
                                                     int64
   instances phys restraint
                                     15829 non-null
41
                                                     int64
   instances seclusion
                                     15829 non-null
                                                     int64
   students mech restraint
                                     15829 non-null
                                                     int64
   students phys restraint
                                     15829 non-null
                                                     int64
   students_seclusion
                                     15829 non-null int64
46 students SAT_ACT
                                     15829 non-null int64
```

dtypes: float64(1), int64(40), object(6)

memory usage: 5.8+ MB

executed in 2ms, finished 14:34:10 2021-06-17

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district
    ▼ 3.4.1 Data Fields
        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X_train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

10 Conclusions

```
In [89]: for variable in graph feats:
               ax, figure = plt.subplots(1,1,figsize=(12,12))
               plt.ylim(-100,2500)
               sns.boxplot(x='grad_rate_midpt', y=variable, data=df_school, st
               plt.title("{} vs. Graduation Rate".format(variable))
           executed in 799ms, finished 14:34:10 2021-06-17
              1500
            days_suspended
              1000
               500
In [90]: numerical_teacher = ['salaries_teachers']
           executed in 3ms, finished 14:34:10 2021-06-17
```

```
Contents ⊋ ❖
```

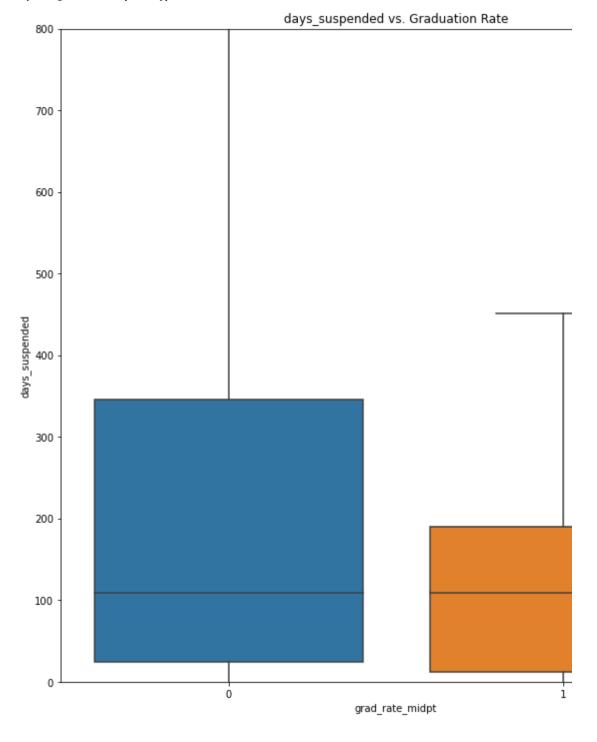
- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
In [91]: for variable in numerical_teacher:
    ax, figure = plt.subplots(1,1,figsize=(12,12))
    plt.ylim(-100,10000000)
    sns.boxplot(x='grad_rate_midpt', y=variable, data=df_school, sh
    plt.title("{} vs. Graduation Rate".format(variable))
executed in 118ms, finished 14:34:10 2021-06-17
```

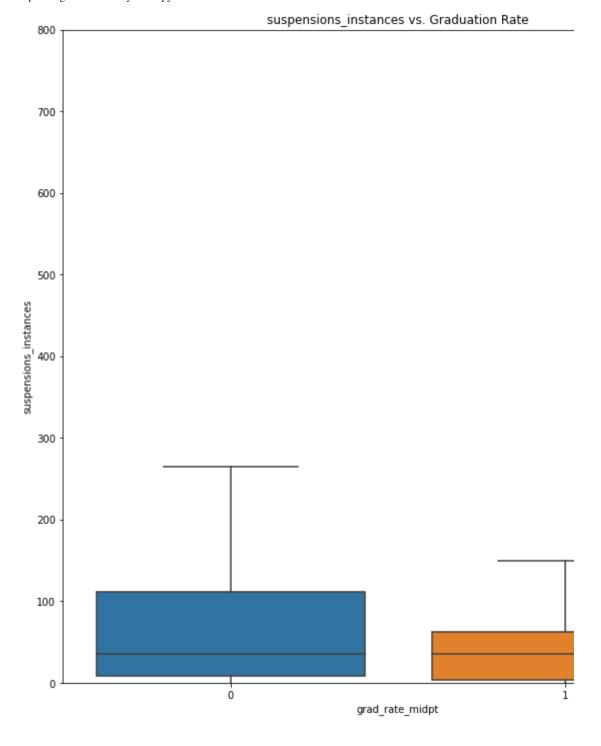
- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
In [93]: for variable in graph_feats:
    ax, figure = plt.subplots(1,1,figsize=(12,12))
    plt.ylim(0,800)
    sns.boxplot(x='grad_rate_midpt', y=variable, data=df_school, sh
    plt.title("{} vs. Graduation Rate".format(variable))
executed in 956ms, finished 14:34:11 2021-06-17
```

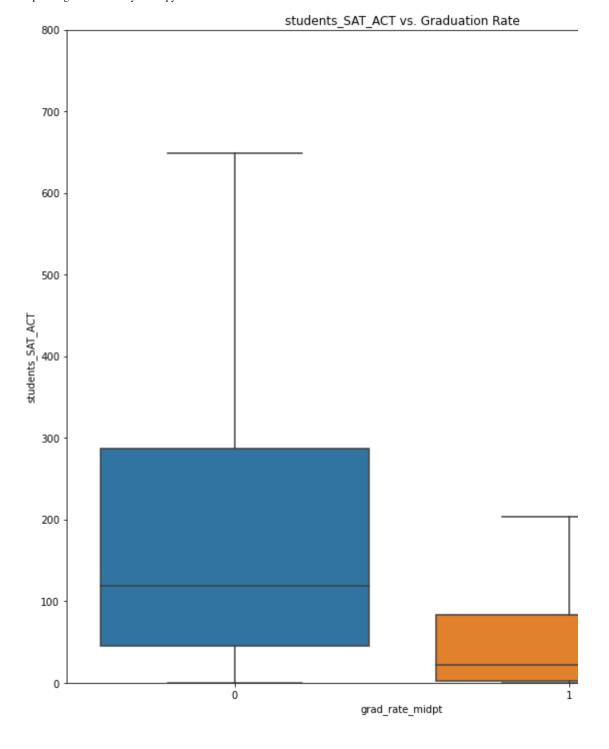
- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions



- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

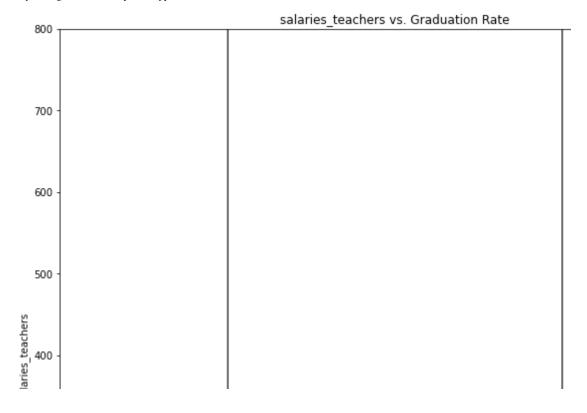


- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

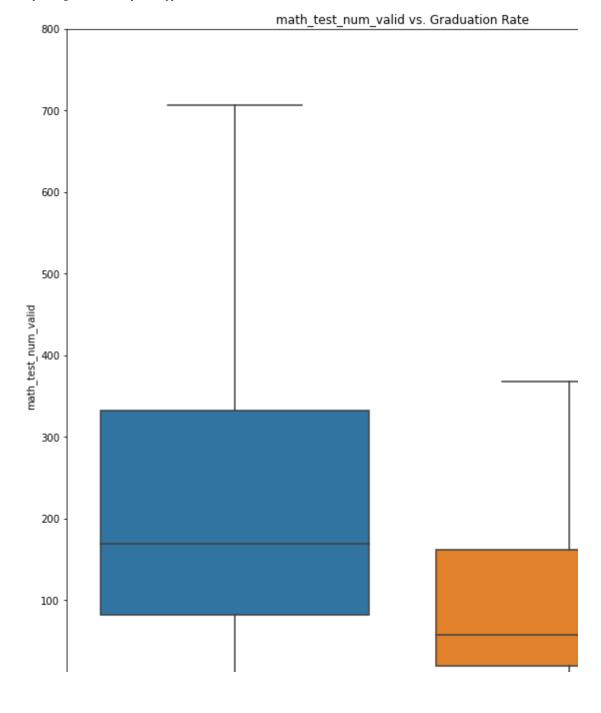


- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X_test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

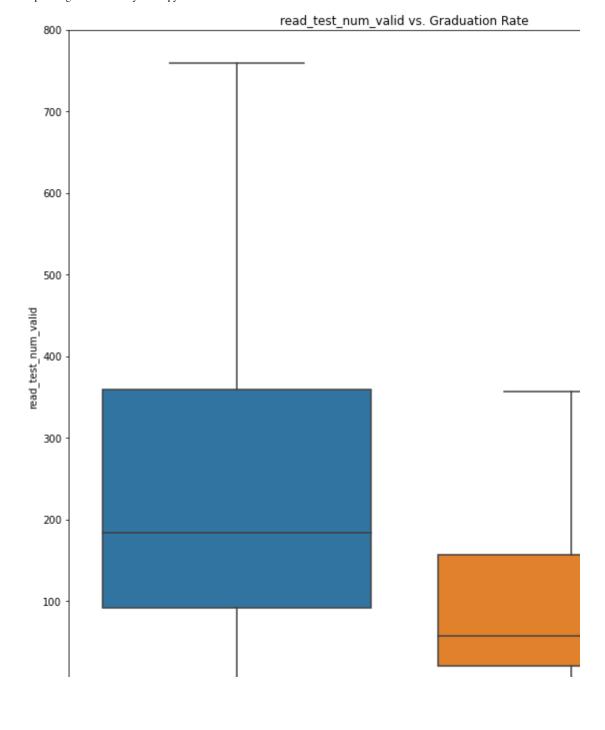
- 3.3.0 ELICOUE CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions



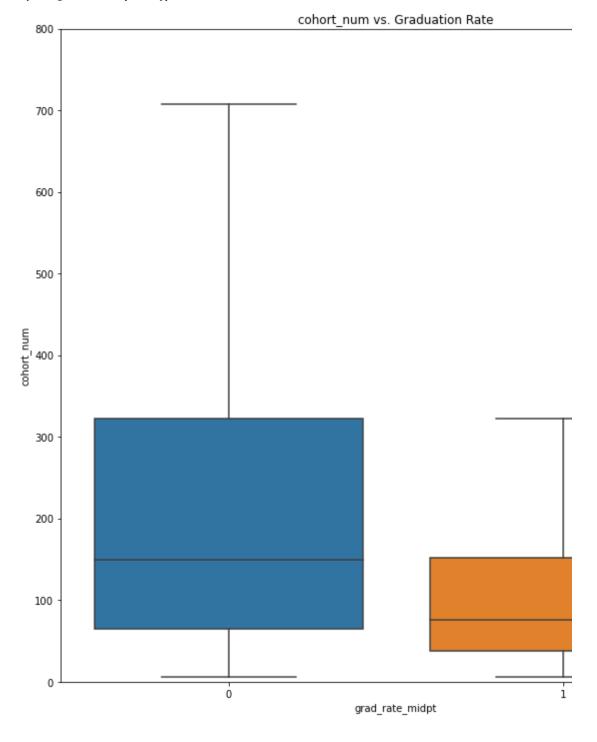
- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions



- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions



- 3.3.0 ELICOUE CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions



- 3.3.0 ELICOUR CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X_test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

3.4 df district

In [94]: #Concatenate district info for all three states into one dataframe df district = pd.concat([df district northeast, df district south & df_district_midwest, df_district_west, df_

executed in 24ms, finished 14:34:11 2021-06-17

Contents & *
3.3.7 Reading and Math Tests
3.3.8 Numeric Columns - nulls
3.3.9 Visualizations
▼ 3.4 df_district
▼ 3.4.1 Data Fields
3.4.1.1 Data Field Cleanup
3.4.2 Columns to drop
3.4.3 english_language_learners
▼ 4 Join Datasets
4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
5.1 Cleanup
5.2 Engineer % Columns
5.3 Reset Num List
▼ 6 Category Columns - cleanup
6.1 Final Check for NaNs
7 Train Test Split
▼ 8 Encode Features
8.1 X_train Encode
8.2 X_test Encode
▼ 9 Model Development
▼ 9.1 Logistic Regression
9.1.1 Check for Overfit
9.1.2 Model reiteration - parameter tuning
▼ 9.2 Random Forest Classifier
9.2.1 Check for Overfit
9.2.2 Feature Importance
9.2.3 Model reiteration - parameter tunin
▼ 9.3 Gradient Boosting Classifier
9.3.1 Check for Overfit
9.3.2 Model reiteration - parameter tunin
10 Conclusions

In [95]:	df_district['state_name'].value_counts(normalize=True,	dropna =Fals €
	executed in 6ms, finished 14:34:11 202		
Out[95]:	Texas	0.067443	
	California	0.062879	
	Ohio	0.060570	
	Illinois	0.057134	
	New York	0.055308	
	Michigan	0.049992	
	Pennsylvania	0.042958	
	Arizona	0.038125	
	New Jersey	0.037427	
	Oklahoma	0.032594	
	Minnesota	0.031305	
	Missouri	0.030768	
	Montana	0.026687	
	Wisconsin	0.025345	
	Indiana	0.023090	
	Massachusetts	0.021962	
	Vermont	0.019492	
	Iowa	0.018633	
	Washington	0.017720	
	Kansas	0.017022	
	North Carolina	0.016968	
	New Hampshire	0.016109	
	Arkansas	0.015626	
	Nebraska	0.015250	
	Maine	0.014391	
	Colorado	0.014230	
	North Dakota	0.012189	
	Georgia	0.011974	
	Virginia	0.011974	
	Louisiana	0.011974	
	Oregon	0.011867	
	Connecticut	0.011115	
	Kentucky	0.009988	
	Alabama South Dakota	0.009665	
	Utah	0.009182 0.008914	
	Mississippi	0.008914	
	Idaho	0.008592	
	New Mexico	0.008538	
	Tennessee	0.007840	
	Tennessee	0.00/040	

3.3.0 ELICOUR CHARLE

0.0.0
3.3.7 Reading and Math Tests
3.3.8 Numeric Columns - nulls
3.3.9 Visualizations
▼ 3.4 df_district
▼ 3.4.1 Data Fields
3.4.1.1 Data Field Cleanup
3.4.2 Columns to drop
3.4.3 english_language_learners
▼ 4 Join Datasets
4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
5.1 Cleanup
5.2 Engineer % Columns
5.3 Reset Num List
▼ 6 Category Columns - cleanup
6.1 Final Check for NaNs
7 Train Test Split
▼ 8 Encode Features
8.1 X_train Encode
8.2 X_test Encode
▼ 9 Model Development
▼ 9.1 Logistic Regression
9.1.1 Check for Overfit

9.1.2 Model reiteration - parameter tuning

9.2.3 Model reiteration - parameter tunin

9.3.2 Model reiteration - parameter tunin

▼ 9.2 Random Forest Classifier
 9.2.1 Check for Overfit
 9.2.2 Feature Importance

▼ 9.3 Gradient Boosting Classifier 9.3.1 Check for Overfit

10 Conclusions

South Carolina	0.005477
Florida	0.004081
District of Columbia	0.003974
Rhode Island	0.003437
Wyoming	0.003276
West Virginia	0.003061
Delaware	0.002846
Maryland	0.001342
Nevada	0.001020
Name: state_name, dtype	: float64

3.4.1 Data Fields

In [96]: df_district.head()

executed in 19ms, finished 14:34:11 2021-06-17

Out[96]:

	year	leaid	lea_name	state_name	state_leaid	city_location	urban_centric_loca
0	2015	900001	UNIFIED SCHOOL DISTRICT #3	Connecticut	349	HARTFORD	City, midsi:
1	2015	900002	Connecticut Technical High Sc	Connecticut	900	MIDDLETOWN	Suburb, lar
2	2015	900003	UNIFIED SCHOOL DISTRICT #1	Connecticut	336	WETHERSFIELD	Suburb, ları
3	2015	900004	UNIFIED SCHOOL DISTRICT #2	Connecticut	347	HARTFORD	City, sm
4	2015	900005	REGIONAL SCHOOL DISTRICT 19	Connecticut	219	STORRS	Suburb, lar

5 rows × 55 columns

```
Contents 2 &
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [97]: df district.info()
           executed in 56ms, finished 14:34:12 2021-06-17
               exp nonersec other
                                                    10023 HOH-HULL
                                                                     UD TECL
                exp textbooks
                                                    18623 non-null
                                                                     object
                exp utilities energy
                                                    18623 non-null
                                                                     object
                exp tech supplies services
                                                    18623 non-null
                                                                     object
            41
                exp tech equipment
                                                    18623 non-null
                                                                     object
                outlay capital total
                                                    18623 non-null
                                                                     object
                outlay capital construction
                                                    18623 non-null
                                                                     object
                outlay capital land structures
                                                    18623 non-null
                                                                     object
                outlay capital instruc equip
                                                    18623 non-null
                                                                     object
            46 outlay capital other equip
                                                    18623 non-null
                                                                     object
                outlay capital nonspec equip
                                                    18623 non-null
                                                                     object
            48
                salaries total
                                                    18623 non-null
                                                                     object
            49
                salaries instruction
                                                    18623 non-null
                                                                     object
                benefits employee total
                                                    18623 non-null
                                                                     object
                                                    12235 non-null
                                                                     float.64
            51
                cohort num
            52
                grad rate high
                                                    12235 non-null
                                                                     object
            53
                grad rate low
                                                    12235 non-null
                                                                     object
                grad rate midpt
                                                    12235 non-null
                                                                     object
           dtypes: float64(4), int64(2), object(49)
           memory usage: 8.0+ MB
 In [98]:
           #What proportion of our data frame is nulls?
           percent null df(df district, 'df district')
           executed in 63ms, finished 14:34:12 2021-06-17
 Out[98]: 'percent of nulls in df_district is 4%'
           3.4.1.1 Data Field Cleanup
 In [99]:
           column names = list(df district.columns)
           executed in 2ms, finished 14:34:12 2021-06-17
In [100]: categorical = ['year', 'leaid', 'leaname', 'state name', 'state leaid',
                           'urban centric locale', 'agency type' l
           executed in 2ms, finished 14:34:12 2021-06-17
```

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

10 Conclusions

```
In [101]: for ele in categorical:
                try:
                     column names.remove(ele)
                except ValueError:
                     pass
            executed in 2ms, finished 14:34:12 2021-06-17
In [102]: numerical = column_names
            executed in 2ms, finished 14:34:12 2021-06-17
In [103]: column_names = list(df_district.columns)
            executed in 2ms, finished 14:34:12 2021-06-17
In [104]: for c in column names:
                df district[c].replace('Suppressed data', np.NaN,inplace =True)
            executed in 35ms, finished 14:34:12 2021-06-17
In [105]: for c in column names:
                df district[c].replace('Not applicable', np.NaN,inplace =True)
            executed in 37ms, finished 14:34:12 2021-06-17
In [106]: for c in column names:
                df district[c].replace('Missing/not reported', np.NaN,inplace =
            executed in 35ms, finished 14:34:12 2021-06-17
```

3.4.2 Columns to drop

In [107]: df district.info()

executed in 37ms, finished 14:34:12 2021-06-17

<class 'pandas.core.frame.DataFrame'> Int64Index: 18623 entries, 0 to 3049 Data columns (total 55 columns).

Data	columns (total 55 columns):		
#	Column	Non-Null Count	Dtype
0	year	18623 non-null	 int64
1	leaid	18623 non-null	int64
2	lea name	18623 non-null	object
3	state name	18623 non-null	object
4	state leaid	18623 non-null	object
5	city location	18623 non-null	object
6	urban centric locale	18592 non-null	object
7	agency type	18623 non-null	object
8	enrollment	17047 non-null	object
9	<pre>english_language_learners</pre>	12381 non-null	object
10	est_population_total	13111 non-null	float64
11	est_population_5_17_poverty	13111 non-null	float64
12	<pre>est_population_5_17_poverty_pct</pre>	13107 non-null	float64
13	rev_total	16883 non-null	object
14	rev_fed_total	16883 non-null	object
15	rev_state_total	16883 non-null	object
16	rev_local_total	16883 non-null	object
17	rev_local_prop_tax	12577 non-null	object
18	exp_total	16883 non-null	object
19	<pre>exp_current_elsec_total</pre>	16883 non-null	object
20	<pre>exp_current_instruction_total</pre>	16883 non-null	object
21	<pre>exp_current_supp_serve_total</pre>	16883 non-null	object
22	exp_current_pupils	16883 non-null	object
23	exp_current_instruc_staff	16883 non-null	object
24	<pre>exp_current_general_admin</pre>	16883 non-null	object
25	exp_current_sch_admin	16883 non-null	object
26	<pre>exp_current_operation_plant</pre>	16883 non-null	object
27	<pre>exp_current_student_transport</pre>	16883 non-null	object
28	exp_current_bco	16883 non-null	object
29	exp_current_supp_serv_nonspec	16883 non-null	object
30	exp_current_other	16883 non-null	object
31	exp_current_food_serv	16883 non-null	object
32	exp_current_enterprise	16883 non-null	object
33	exp_current_other_elsec	16883 non-null	object
34	exp_nonelsec	16883 non-null	object

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
exp nonelsec community serv
                                     16883 non-null
                                                     object
36
    exp nonelsec adult education
                                     16883 non-null
                                                     object
    exp nonelsec other
37
                                     16883 non-null
                                                     object
    exp textbooks
38
                                     16883 non-null
                                                     object
39
    exp utilities energy
                                     16883 non-null
                                                     object
    exp tech supplies services
                                     16883 non-null object
    exp tech equipment
                                     16883 non-null
                                                     object
41
    outlay capital total
42
                                     16883 non-null
                                                     object
    outlay capital construction
                                     16883 non-null
                                                     object
    outlay capital land structures
                                     16883 non-null
                                                     object
    outlay capital instruc equip
                                     16883 non-null
                                                     object
46 outlay capital other equip
                                                     object
                                     16883 non-null
    outlay capital nonspec equip
                                     16883 non-null
                                                     object
    salaries total
                                     16883 non-null object
48
    salaries instruction
                                     16883 non-null object
    benefits_employee_total
50
                                     16883 non-null
                                                     object
                                                     float64
    cohort num
                                     12235 non-null
52
    grad rate high
                                     11972 non-null object
53
    grad rate low
                                     11972 non-null
                                                     object
   grad rate midpt
                                     11972 non-null object
dtypes: float64(4), int64(2), object(49)
```

executed in 11ms, finished 14:34:12 2021-06-17

memory usage: 8.0+ MB

executed in 10ms, finished 14:34:12 2021-06-17

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [110]: df district['enrollment'].value counts(normalize=True, dropna=False
           executed in 11ms, finished 14:34:12 2021-06-17
Out[110]: NaN
                     0.084627
                     0.020459
           177
                     0.001557
           202
                     0.001450
           147
                     0.001289
                        . . .
           83648
                     0.000054
           4334
                     0.000054
           8824
                     0.000054
           23529
                     0.000054
           4308
                     0.000054
           Name: enrollment, Length: 5443, dtype: float64
In [111]: # Enrollment is mostly zeros or NaNs
           df_district = df_district.drop(['enrollment'], axis=1)
           executed in 10ms, finished 14:34:12 2021-06-17
In [112]: # year is the same for all rows, state name same in school data, st
           df_district = df_district.drop(['year','state_name','state_leaid'],
           executed in 10ms, finished 14:34:12 2021-06-17
In [113]: #What proportion of our data frame is nulls?
           percent null df(df district, 'df district')
           executed in 54ms, finished 14:34:12 2021-06-17
Out[113]: 'percent of nulls in df district is 9%'
```

3.4.3 english_language_learners

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [114]: df district['english language learners'].value counts(normalize=Tru
            executed in 8ms, finished 14:34:12 2021-06-17
Out[114]: NaN
                      0.335177
            3
                      0.102669
                      0.021694
                      0.020459
                      0.018955
                        . . .
            5423
                      0.000054
           1553
                      0.000054
           1996
                      0.000054
                      0.000054
           13796
            675
                      0.000054
           Name: english language learners, Length: 1581, dtype: float64
In [115]: ell = ['english language learners']
           executed in 2ms, finished 14:34:12 2021-06-17
In [116]: | for c in ell:
                df district[c].replace('Missing/not reported', np.NaN,inplace =
            executed in 4ms, finished 14:34:12 2021-06-17
           df district['english language learners'].value counts(normalize=Tru
In [117]:
            executed in 7ms, finished 14:34:12 2021-06-17
Out[117]: NaN
                      0.335177
                      0.102669
            3
                      0.021694
                      0.020459
                      0.018955
                      0.000054
           5423
           1553
                      0.00054
           1996
                      0.000054
           13796
                      0.000054
           675
                      0.000054
           Name: english language learners, Length: 1581, dtype: float64
```

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district
    ▼ 3.4.1 Data Fields
        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

10 Conclusions

```
In [118]: # fill na with median
           df district['english language learners'] = df district['english lar
           executed in 6ms, finished 14:34:12 2021-06-17
In [119]: df_district['english language learners'].value_counts(normalize=Tru
           executed in 8ms, finished 14:34:12 2021-06-17
Out[119]: 26.0
                     0.335177
           3
                     0.102669
           4
                     0.021694
                     0.020459
                     0.018955
                     0.00054
           735
           62575
                     0.000054
           2085
                     0.000054
           991
                     0.000054
           1812
                     0.000054
           Name: english language learners, Length: 1581, dtype: float64
```

```
Contents 2 &
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [120]: df district.info()
           executed in 33ms, finished 14:34:12 2021-06-17
           <class 'pandas.core.frame.DataFrame'>
           Int64Index: 18623 entries, 0 to 3049
           Data columns (total 42 columns):
                Column
                                                  Non-Null Count Dtype
            0
                leaid
                                                  18623 non-null
                                                                   int64
            1
                                                  18623 non-null
                                                                   object
                lea name
                                                  18623 non-null object
                agency type
            3
                english language learners
                                                  18623 non-null object
                rev total
                                                  16883 non-null object
                rev fed total
                                                  16883 non-null
                                                                   object
                rev state total
                                                  16883 non-null
                                                                   object
                rev local total
                                                  16883 non-null
                                                                   object
                rev local prop tax
                                                  12577 non-null
                                                                   object
            9
                exp total
                                                  16883 non-null
                                                                   object
            10
               exp current elsec total
                                                  16883 non-null
                                                                 object
            11
               exp current instruction total
                                                  16883 non-null object
            12
                exp current supp serve total
                                                  16883 non-null
                                                                   object
                exp current pupils
                                                  16883 non-null
                                                                   object
In [121]:
          df district = df district.drop duplicates(subset=['lea name'])
           executed in 14ms, finished 14:34:12 2021-06-17
```

4 Join Datasets

```
Contents ⊋ ❖
```

- 3.3.0 ELICOUE CHARLE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
- ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
- ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
- ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
- 10 Conclusions

In [125]: df_train = df_school.merge(df_district, on = 'lea_name')
 df_train.head()
 executed in 48ms, finished 14:34:12 2021-06-17

Out[125]:

	state_name	lea_name	urban_centric_locale	charter	enrollment	read_test_num_
(O Connecticut	REGIONAL SCHOOL DISTRICT 19	Suburb, large	0.0	1153	
	1 Connecticut	COMMON GROUND HIGH SCHOOL DIS	City, midsize	1.0	186	
:	2 Connecticut	THE BRIDGE ACADEMY DISTRICT	City, midsize	1.0	279	
;	3 Connecticut	EXPLORATIONS DISTRICT	Town, fringe	1.0	92	
	4 Connecticut	Stamford Academy	City, midsize	1.0	150	

5 rows × 83 columns

```
Contents 2 &
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

10 Conclusions

```
In [126]: df train.info()
           executed in 48ms, finished 14:34:12 2021-06-17
                                                  15/85 non-null
                                                                   int64
               students report harass sex
               instances mech restraint
                                                  15785 non-null
                                                                   int64
               instances phys restraint
                                                  15785 non-null
                                                                   int64
            36
            37
                instances seclusion
                                                  15785 non-null
                                                                   int64
               students mech restraint
                                                  15785 non-null
                                                                   int64
               students phys restraint
                                                  15785 non-null
                                                                   int64
                students seclusion
                                                  15785 non-null
                                                                   int64
                students SAT ACT
                                                  15785 non-null
                                                                   int64
            41
            42
                leaid
                                                  15785 non-null
                                                                   int64
            43
                agency_type
                                                  15785 non-null
                                                                   object
                english language learners
                                                  15785 non-null
                                                                   object
            45
               rev total
                                                  15331 non-null
                                                                   object
               rev fed total
                                                  15331 non-null
                                                                   object
               rev state total
                                                  15331 non-null
                                                                   object
            48
               rev local total
                                                  15331 non-null
                                                                   object
               rev local prop tax
                                                  13039 non-null
                                                                   object
            50
                                                                   object
                exp_total
                                                  15331 non-null
            51
               exp current elsec total
                                                  15331 non-null
                                                                   object
            52
               exp_current_instruction_total
                                                  15331 non-null
                                                                   object
                                                  15331 non-null
               exp current supp serve total
                                                                   object
In [127]: # Columns to drop - we don't need these anymore
           df train = df train.drop(['leaid','lea name'], axis=1)
           executed in 11ms, finished 14:34:12 2021-06-17
```

4.1 Check for NaNs

```
In [128]: nan_values = df_train.isna()
nan_columns = nan_values.any()
executed in 25ms, finished 14:34:12 2021-06-17
```

Contents 2 ❖

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
In [129]: columns_with_nan = df_train.columns[nan_columns].tolist()
    print(columns_with_nan)
    executed in 3ms, finished 14:34:12 2021-06-17
```

['rev_total', 'rev_fed_total', 'rev_state_total', 'rev_local_total' nt_elsec_total', 'exp_current_instruction_total', 'exp_current_supp nt_instruc_staff', 'exp_current_general_admin', 'exp_current_sch_ac nt_student_transport', 'exp_current_bco', 'exp_current_supp_serv_nc serv', 'exp_current_enterprise', 'exp_current_other_elsec', 'exp_nc nelsec_adult_education', 'exp_nonelsec_other', 'exp_textbooks', 'exs', 'exp_tech_equipment', 'outlay_capital_total', 'outlay_capital_c'outlay_capital_instruc_equip', 'outlay_capital_other_equip', 'outlay_capital_other_equip', 'outlay_capital_instruction', 'benefits_employee_total']

We still have NaN values in our numerical fields

5 Numeric Columns - cleanup

▼ 5.1 Cleanup

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

▼ 5.2 Engineer % Columns

Out[136]:

	state_name	urban_centric_locale	charter	enrollment	read_test_num_valid	read_test_p
0	Connecticut	Suburb, large	0	1153	260	_
1	Connecticut	City, midsize	1	186	39	
2	Connecticut	City, midsize	1	279	139	
3	Connecticut	Town, fringe	1	92	21	
4	Connecticut	City, midsize	1	150	25	

5 rows × 83 columns

```
In [137]: pct_graph = ['SAT_ACT_%_enrollment','teacher_%_total_salary','read_
executed in 2ms, finished 14:34:13 2021-06-17
```

3.3.0 Encode Charter 3.3.7 Reading and Math Tests 3.3.8 Numeric Columns - nulls 3.3.9 Visualizations ▼ 3.4 df_district ▼ 3.4.1 Data Fields 3.4.1.1 Data Field Cleanup 3.4.2 Columns to drop 3.4.3 english_language_learners ▼ 4 Join Datasets 4.1 Check for NaNs ▼ 5 Numeric Columns - cleanup 5.1 Cleanup 5.2 Engineer % Columns 5.3 Reset Num List ▼ 6 Category Columns - cleanup 6.1 Final Check for NaNs 7 Train Test Split ▼ 8 Encode Features 8.1 X train Encode 8.2 X test Encode ▼ 9 Model Development ▼ 9.1 Logistic Regression 9.1.1 Check for Overfit 9.1.2 Model reiteration - parameter tuning ▼ 9.2 Random Forest Classifier 9.2.1 Check for Overfit 9.2.2 Feature Importance

9.2.3 Model reiteration - parameter tunin

9.3.2 Model reiteration - parameter tunin

▼ 9.3 Gradient Boosting Classifier

9.3.1 Check for Overfit

10 Conclusions

```
In [138]: for variable in pct graph:
                ax, figure = plt.subplots(1,1,figsize=(12,12))
                plt.ylim(0,100)
                sns.boxplot( x='grad rate midpt', y=variable, data=df train, sf
                plt.title("{} vs. Graduation Rate".format(variable))
            executed in 475ms, finished 14:34:13 2021-06-17
                80
                60
             SAT_ACT_%_enrollment
In [139]: df_train = df_train.drop(['salaries_instruction','salaries_total',
            executed in 18ms, finished 14:34:14 2021-06-17
```

5.3 Reset Num List

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- **▼** 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

```
In [142]: num_col_names = set(col_names) - set(col_names_category)
    num_col_names = list(num_col_names)
    executed in 3ms, finished 14:34:14 2021-06-17
In [143]: df_train[num_col_names] = df_train[num_col_names].fillna(df_train[resecuted in 42ms, finished 14:34:14 2021-06-17)
```

▼ 6 Category Columns - cleanup

3.3.0 Encode Charter 3.3.7 Reading and Math Tests 3.3.8 Numeric Columns - nulls 3.3.9 Visualizations ▼ 3.4 df_district ▼ 3.4.1 Data Fields 3.4.1.1 Data Field Cleanup 3.4.2 Columns to drop 3.4.3 english_language_learners ▼ 4 Join Datasets 4.1 Check for NaNs ▼ 5 Numeric Columns - cleanup 5.1 Cleanup 5.2 Engineer % Columns 5.3 Reset Num List ▼ 6 Category Columns - cleanup 6.1 Final Check for NaNs 7 Train Test Split ▼ 8 Encode Features 8.1 X train Encode 8.2 X test Encode ▼ 9 Model Development ▼ 9.1 Logistic Regression 9.1.1 Check for Overfit 9.1.2 Model reiteration - parameter tuning ▼ 9.2 Random Forest Classifier 9.2.1 Check for Overfit 9.2.2 Feature Importance 9.2.3 Model reiteration - parameter tunin ▼ 9.3 Gradient Boosting Classifier 9.3.1 Check for Overfit 9.3.2 Model reiteration - parameter tunin

10 Conclusions

```
In [148]: | df train.info()
           executed in 20ms, finished 14:34:15 2021-06-17
                                                  15785 non-null
                                                                  int64
               students_disc_harass_race
               students_disc_harass_sex
                                                  15785 non-null
                                                                  int64
               students report harass dis
                                                  15785 non-null
                                                                  int64
               students report harass race
                                                  15785 non-null
                                                                  int64
               students_report_harass_sex
                                                  15785 non-null
                                                                  int64
               instances_mech_restraint
                                                  15785 non-null
                                                                  int64
               instances phys restraint
                                                  15785 non-null
                                                                  int64
               instances seclusion
                                                  15785 non-null
                                                                  int64
               students mech restraint
                                                 15785 non-null
                                                                  int64
                                                                  int64
               students phys restraint
                                                  15785 non-null
               students seclusion
                                                  15785 non-null
                                                                  int64
                                                 15785 non-null
                agency type
                                                                  object
            41
               english language learners
                                                  15785 non-null
                                                                  int64
               rev total
                                                  15785 non-null
                                                                  int64
            43
               rev fed total
                                                  15785 non-null
                                                                  int64
               rev state total
                                                  15785 non-null
                                                                  int64
               rev local total
                                                  15785 non-null
                                                                  int64
               rev local prop tax
                                                  15785 non-null
                                                                  int64
            47
               exp total
                                                  15785 non-null
                                                                  int64
                exp current elsec total
                                                 15785 non-null
                                                                  int64
```

▼ 6.1 Final Check for NaNs

```
In [149]: nan_values = df_train.isna()
    nan_columns = nan_values.any()

    executed in 6ms, finished 14:34:15 2021-06-17

In [150]: columns_with_nan = df_train.columns[nan_columns].tolist()
    print(columns_with_nan)
    executed in 4ms, finished 14:34:15 2021-06-17
```

7 Train Test Split

[]

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations

▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

10 Conclusions

```
In [151]:
             X=df train.drop(columns = ['grad rate midpt']) # Features
             y=df train['grad rate midpt'] #Target
             executed in 7ms, finished 14:34:15 2021-06-17
In [152]: X train, X test, y train, y test = train_test_split(X, y, test_size
             executed in 13ms, finished 14:34:15 2021-06-17
In [153]: X_train
             executed in 22ms, finished 14:34:15 2021-06-17
Out[153]:
                      state_name urban_centric_locale charter enrollment read_test_num_valid read_te
                                          Rural, distant
                                                            0
                                                                      183
                                                                                           96
               7105
                       Minnesota
                                         Town, distant
                                                            0
                                                                      241
                                                                                           50
              13034
                        Alabama
               6979
                        Michigan
                                         Suburb, small
                                                            0
                                                                     1265
                                                                                          300
                                                                                            7
                                                                      121
               9785
                         Arizona
                                          Town, fringe
                                                            0
                                                                      269
               6744
                        Michigan
                                            City, large
                                                                                           44
                                                                                            ...
                                         Town, distant
                                                            0
                                                                      992
                                                                                          243
               7404
                       Minnesota
              10737
                        California
                                         Suburb, large
                                                            0
                                                                     1377
                                                                                          312
                        Colorado
                                            City, large
                                                            1
                                                                      501
                                                                                          177
              11490
                      Washington
                                          Town, distant
                                                            0
                                                                      713
                                                                                          163
              12604
```

11049 rows × 79 columns

Illinois

In [154]: X train.shape executed in 4ms, finished 14:34:15 2021-06-17

City, midsize

0

1240

Out[154]: (11049, 79)

5545

```
In [155]: | X_test.shape
             executed in 3ms, finished 14:34:15 2021-06-17
Out[155]: (4736, 79)
```


- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

8 Encode Features

```
In [156]: # define transformers
           ss = StandardScaler()
          ohe = OneHotEncoder(handle_unknown='ignore', sparse=False,)
           # set up pipelines for each column group
           numeric pipe = Pipeline([('ss', ss)])
          categorical_pipe = Pipeline([('ohe', ohe)])
           # set up columnTransformer
          col_transformer = ColumnTransformer(
                                transformers=[
                                    ('nums', numeric pipe, num col names),
                                    ('cats', categorical_pipe, col_names_ohe),
                                ],
                                remainder='drop',
                                n_{jobs=-1}
           executed in 4ms, finished 14:34:15 2021-06-17
```

8.1 X train Encode

```
In [157]: X train = col_transformer.fit_transform(X_train)
            executed in 1.93s, finished 14:34:17 2021-06-17
```

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
```

9.2.3 Model reiteration - parameter tunin

9.3.2 Model reiteration - parameter tunin

▼ 9.3 Gradient Boosting Classifier

9.3.1 Check for Overfit

10 Conclusions

```
In [158]: col transformer.named transformers ['cats'].named steps['ohe']\
                .get feature names()
          executed in 5ms, finished 14:34:17 2021-06-17
Out[158]: array(['x0_Alabama', 'x0_Arizona', 'x0_Arkansas', 'x0_California',
                  'x0 Colorado', 'x0 Connecticut', 'x0 Delaware',
                  'x0 District of Columbia', 'x0 Florida', 'x0 Georgia', 'x0 ]
                  'x0 Illinois', 'x0_Indiana', 'x0_Iowa', 'x0_Kansas', 'x0_Ker
                  'x0_Louisiana', 'x0_Maine', 'x0_Maryland', 'x0_Massachusetts
                  'x0 Michigan', 'x0 Minnesota', 'x0 Mississippi', 'x0 Missour
                  'x0 Montana', 'x0 Nebraska', 'x0 Nevada', 'x0 New Hampshire'
                  'x0 New Jersey', 'x0 New Mexico', 'x0 New York',
                  'x0 North Carolina', 'x0 North Dakota', 'x0 Ohio', 'x0 Oklah
                  'x0_Oregon', 'x0_Pennsylvania', 'x0_Rhode Island',
                  'x0 South Carolina', 'x0 South Dakota', 'x0 Tennessee', 'x0
                  'x0_Utah', 'x0_Vermont', 'x0_Virginia', 'x0_Washington',
                  'x0_West Virginia', 'x0_Wisconsin', 'x0_Wyoming', 'x1_City,
                  'x1_City, midsize', 'x1_City, small', 'x1_Rural, distant',
                  'x1 Rural, fringe', 'x1_Rural, remote', 'x1_Suburb, large',
                  'x1_Suburb, midsize', 'x1_Suburb, small', 'x1_Town, distant'
                  'x1_Town, fringe', 'x1_Town, remote', 'x2_Charter agency',
                  'x2 Local school district that is a component of a superviso
                  'x2 Other education agency',
                  'x2 Regional education service agency',
                  'x2 Regular local school district', 'x2 State-operated agenc
                  'x2 Supervisory union'], dtype=object)
In [159]: X train.shape
           executed in 3ms, finished 14:34:17 2021-06-17
Out[159]: (11049, 144)
```

8.2 X_test Encode

```
In [160]: X_test = col_transformer.transform(X_test)
executed in 843ms, finished 14:34:17 2021-06-17
```

```
Contents 2 &
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
col transformer.named transformers ['cats'].named steps['ohe']\
In [161]:
                .get feature names()
           executed in 5ms, finished 14:34:17 2021-06-17
Out[161]: array(['x0 Alabama', 'x0 Arizona', 'x0 Arkansas', 'x0 California',
                  'x0 Colorado', 'x0 Connecticut', 'x0 Delaware',
                  'x0 District of Columbia', 'x0 Florida', 'x0 Georgia', 'x0 ]
                  'x0 Illinois', 'x0 Indiana', 'x0 Iowa', 'x0 Kansas', 'x0 Ker
                  'x0 Louisiana', 'x0 Maine', 'x0 Maryland', 'x0 Massachusetts
                  'x0 Michigan', 'x0 Minnesota', 'x0 Mississippi', 'x0 Missour
                  'x0 Montana', 'x0 Nebraska', 'x0 Nevada', 'x0 New Hampshire'
                  'x0 New Jersey', 'x0 New Mexico', 'x0 New York',
                  'x0 North Carolina', 'x0 North Dakota', 'x0 Ohio', 'x0 Oklah
                  'x0 Oregon', 'x0 Pennsylvania', 'x0 Rhode Island',
                  'x0 South Carolina', 'x0 South Dakota', 'x0 Tennessee', 'x0
                  'x0_Utah', 'x0_Vermont', 'x0_Virginia', 'x0_Washington',
                  'x0 West Virginia', 'x0 Wisconsin', 'x0 Wyoming', 'x1 City,
                  'x1 City, midsize', 'x1 City, small', 'x1 Rural, distant',
                  'x1 Rural, fringe', 'x1 Rural, remote', 'x1 Suburb, large',
                  'x1 Suburb, midsize', 'x1 Suburb, small', 'x1 Town, distant'
                  'x1 Town, fringe', 'x1 Town, remote', 'x2 Charter agency',
                  'x2 Local school district that is a component of a superviso
                  'x2 Other education agency',
                  'x2 Regional education service agency',
                  'x2 Regular local school district', 'x2 State-operated agenc
                  'x2 Supervisory union'], dtype=object)
In [162]: X test.shape
          executed in 4ms, finished 14:34:17 2021-06-17
Out[162]: (4736, 144)
In [163]: X test.shape
          executed in 4ms, finished 14:34:17 2021-06-17
Out[163]: (4736, 144)
```

9 Model Development

9.1 Logistic Regression

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [164]: # Logistic model
           log clf = LogisticRegression(class weight='balanced')
           log model = log clf.fit(X train, y train)
           y pred train = log model.predict(X train)
           y pred = log model.predict(X test)
           executed in 237ms, finished 14:34:18 2021-06-17
In [165]: #Confusion matrix for Logistic Regression
           log matrix = confusion_matrix(y test, y pred)
           print('Confusion Matrix:\n', log matrix)
           executed in 10ms, finished 14:34:18 2021-06-17
           Confusion Matrix:
            [[3926 466]
            [ 54 290]]
In [166]: print('The f1 score for the training model is:', f1 score(y train, y
           print('The f1 score for the model is:',f1 score(y test, y pred))
           executed in 13ms, finished 14:34:18 2021-06-17
           The f1 score for the training model is: 0.5368209255533198
           The f1 score for the model is: 0.5272727272727273
In [167]: print(metrics.classification report(y test, y pred, labels=[0,1],
                                                  target names=['normal grad rate
           executed in 13ms, finished 14:34:18 2021-06-17
                              precision
                                             recall f1-score
                                                                  support
           normal grad rate
                                               0.89
                                                          0.94
                                                                     4392
                                    0.99
              low grad rate
                                    0.38
                                               0.84
                                                          0.53
                                                                      344
                                                          0.89
                                                                     4736
                    accuracy
                                                          0.73
                                                                     4736
                  macro avq
                                    0.69
                                               0.87
               weighted avg
                                    0.94
                                               0.89
                                                          0.91
                                                                     4736
```

▼ 9.1.1 Check for Overfit

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

executed in 19ms, finished 14:34:18 2021-06-17

	precision	recall	f1-score	support	
normal grad rate	0.99	0.90	0.94	10313	
low grad rate	0.38	0.91	0.54	736	
accuracy			0.90	11049	
macro avg	0.69	0.90	0.74	11049	
weighted avg	0.95	0.90	0.91	11049	

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

```
target names=['normal grad rate
executed in 12ms, finished 14:34:18 2021-06-17
                    precision
                                   recall f1-score
                                                         support
normal grad rate
                          0.99
                                      0.89
                                                 0.94
                                                             4392
   low grad rate
                          0.38
                                      0.84
                                                 0.53
                                                              344
                                                 0.89
                                                             4736
         accuracy
       macro avq
                          0.69
                                      0.87
                                                 0.73
                                                             4736
    weighted avg
                          0.94
                                      0.89
                                                 0.91
                                                             4736
```

In [171]: print(metrics.classification report(y test, y pred, labels=[0,1],

9.1.2 Model reiteration - parameter tuning

```
Contents 2 *
      3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

```
In [174]: # Instantiate and fit the LogisticReg Model
           log clf = LogisticRegression(class weight='balanced')
           log clf = GridSearchCV(log clf, model params, scoring = 'recall')
           log clf.fit(X train,y train)
           executed in 25.8s, finished 14:34:44 2021-06-17
Out[174]: GridSearchCV(estimator=LogisticRegression(class weight='balanced'),
                        param grid={'max iter': [10, 50, 100], 'random state':
                                      'solver': ['newton-cg', 'sag', 'saga', 'lk
                         scoring='recall')
In [175]: LogisticRegression(class weight='balanced').get params().keys()
           executed in 4ms, finished 14:34:44 2021-06-17
Out[175]: dict_keys(['C', 'class_weight', 'dual', 'fit_intercept', 'intercept
           s', 'n jobs', 'penalty', 'random state', 'solver', 'tol', 'verbose'
In [176]: print(log_clf.best_estimator_.get_params())
           executed in 17ms, finished 14:34:44 2021-06-17
           {'C': 1.0, 'class weight': 'balanced', 'dual': False, 'fit intercer
           one, 'max_iter': 50, 'multi_class': 'auto', 'n_jobs': None, 'penalt
           'tol': 0.0001, 'verbose': 0, 'warm start': False}
In [177]: print('The f1 score for the training model is:', f1 score(y train, y
           print('The f1 score for the model is:',f1 score(y test, y pred))
           executed in 11ms, finished 14:34:44 2021-06-17
           The fl score for the training model is: 0.5368209255533198
           The f1 score for the model is: 0.5272727272727273
```

Contents € 🌣	
3.3.0 Elicode Cliaitei	
3.3.7 Reading and Math Tests	
3.3.8 Numeric Columns - nulls	
3.3.9 Visualizations	
▼ 3.4 df_district	
▼ 3.4.1 Data Fields	
3.4.1.1 Data Field Cleanup	
3.4.2 Columns to drop	

3.4.3 english_language_learners

- ▼ 4 Join Datasets4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X_test Encode
- **▼** 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

In [178]:	print(metrics.cla	ssification_		train, y_pr rget_names	_	
	executed in 17ms, finished 1	14:34:44 2021-06-17				
		precision	recall	f1-score	support	
	normal grad rate	0.99	0.90	0.94	10313	
	low grad rate	0.38	0.91	0.54	736	
	accuracy			0.90	11049	
	macro avg	0.69	0.90	0.74	11049	
	weighted avg	0.95	0.90	0.91	11049	
In [179]:	print(metrics.cla			test, y_pre		
	executed in 12ms, linished	precision	recall	f1-score	support	
	normal grad rate	0.99	0.89	0.94	4392	
	low grad rate	0.38	0.84	0.53	344	
	accuracy			0.89	4736	
	macro avg	0.69	0.87	0.73	4736	
	weighted avg		0.89		4736	

9.2 Random Forest Classifier

 <pre># Instantiate and fit the RandomForestClassifier rfc=RandomForestClassifier(n_estimators=10,class_weight='balanced') rfc.fit(X_train,y_train)</pre>
executed in 247ms, finished 14:34:44 2021-06-17

Out[180]: RandomForestClassifier(class_weight='balanced', n_estimators=10)

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [181]: training start = time.perf counter()
           rfc.fit(X train, y train)
           training end = time.perf counter()
           prediction start = time.perf counter()
          y pred = rfc.predict(X test)
          y pred train = rfc.predict(X train)
           prediction end = time.perf counter()
          acc rfc = (y pred == y test).sum().astype(float) / len(y pred)*100
          rfc train time = training end-training start
          rfc prediction time = prediction end-prediction start
          print("Scikit-Learn's Random Forest Classifier's prediction accurac
          print("Time consumed for training: %4.3f seconds" % (rfc train time
          print("Time consumed for prediction: %6.5f seconds" % (rfc predicti
           executed in 265ms, finished 14:34:44 2021-06-17
           Scikit-Learn's Random Forest Classifier's prediction accuracy is: 9
           Time consumed for training: 0.239 seconds
           Time consumed for prediction: 0.02154 seconds
In [182]: print('The f1 score for the training model is:', f1 score(y train, y
          print('The f1 score for the model is:',f1 score(y test, y pred))
           executed in 11ms, finished 14:34:44 2021-06-17
          The f1 score for the training model is: 0.9560906515580737
           The fl score for the model is: 0.5279383429672447
In [183]: print(metrics.classification_report(y_test, y_pred, labels=[0,1],
                                                 target names=['normal grad rate
           executed in 10ms, finished 14:34:44 2021-06-17
                             precision
                                           recall f1-score
                                                                support
           normal grad rate
                                   0.95
                                              0.99
                                                        0.97
                                                                   4392
              low grad rate
                                   0.78
                                              0.40
                                                        0.53
                                                                    344
                                                        0.95
                                                                   4736
                   accuracy
                                   0.87
                                              0.69
                                                        0.75
                                                                   4736
                  macro avq
               weighted avg
                                   0.94
                                              0.95
                                                        0.94
                                                                   4736
```

9.2.1 Check for Overfit

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [184]: # View confusion matrix for train data and predictions
           confusion_matrix(y_train, y pred_train)
           executed in 11ms, finished 14:34:44 2021-06-17
Out[184]: array([[10312,
                                1],
                       61,
                             67511)
In [185]: # View confusion matrix for test data and predictions
           confusion_matrix(y_test, y_pred)
           executed in 8ms, finished 14:34:44 2021-06-17
Out[185]: array([[4354, 38],
                  [ 207, 137]])
In [186]: rfc training preds = rfc.predict(X train)
           rfc training f1 = f1 score(y train, rfc training preds)
           rfc val preds = rfc.predict(X test) # y hat
           rfc val f1 = f1 score(y test, rfc val preds)
           print(rfc training f1)
           print(rfc val f1)
           executed in 31ms, finished 14:34:44 2021-06-17
           0.9560906515580737
           0.5279383429672447
In [187]: print(metrics.classification_report(y_train, y_pred_train, labels=[
                                                  target names=['normal grad rate
           executed in 16ms, finished 14:34:44 2021-06-17
                               precision
                                             recall f1-score
                                                                  support
           normal grad rate
                                    0.99
                                               1.00
                                                          1.00
                                                                    10313
              low grad rate
                                    1.00
                                               0.92
                                                          0.96
                                                                      736
                                                          0.99
                                                                    11049
                    accuracy
                  macro avg
                                    1.00
                                               0.96
                                                          0.98
                                                                    11049
               weighted avg
                                    0.99
                                               0.99
                                                          0.99
                                                                    11049
```

In [188]:

Contents *⊋* ❖

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X_train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
- ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
- ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
- ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
- 10 Conclusions

<pre>print(metrics.classification_report(y_test, y_pred, labels=[0,1],</pre>
target_names=['normal grad rate'
executed in 12ms, finished 14:34:45 2021-06-17

	precision	recall	f1-score	support	
normal grad rate low grad rate	0.95 0.78	0.99 0.40	0.97 0.53	4392 344	
accuracy macro avg weighted avg	0.87 0.94	0.69 0.95	0.95 0.75 0.94	4736 4736 4736	

Definitely overfit we see this especially in the low grad rate class

▼ 9.2.2 Feature Importance

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

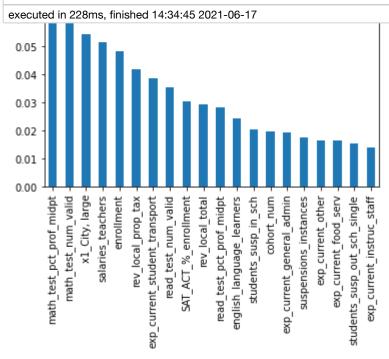
        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
```

```
In [189]: eed to make a list of feature names for feat importance graph
          eed function since data is in numpy arrays
           get column names from ColumnTransformer(column transformer):
           col name = []
           for transformer in columns in column transformer.transformers [:-1
                raw col name = transformer in columns[2]
                if isinstance(transformer in columns[1],Pipeline):
                    transformer = transformer in columns[1].steps[-1][1]
                else:
                    transformer = transformer in columns[1]
                try:
                    names = transformer.get feature names()
               except AttributeError: # if no 'get feature names' function, u
                    names = raw col name
                if isinstance(names,np.ndarray): # eq.
                    col_name += names.tolist()
                elif isinstance(names, list):
                    col name += names
                elif isinstance(names,str):
                    col name.append(names)
           return col name
           executed in 6ms, finished 14:34:45 2021-06-17
In [190]: col name num = get column names from ColumnTransformer(col transfor
           executed in 3ms, finished 14:34:45 2021-06-17
In [191]: col name cat = list(col transformer.named transformers ['cats'].nam
                .get_feature_names())
           executed in 3ms, finished 14:34:45 2021-06-17
In [192]: rfc_columns = col_name_num + col_name_cat
           executed in 2ms, finished 14:34:45 2021-06-17
```

Contents *⊋* ❖

- 3.3.0 EHOUGE CHAILE
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

In [193]: feat_importances = pd.Series(rfc.feature_importances_, index = rfc_feat_importances.nlargest(20).plot(kind='bar')



9.2.3 Model reiteration - parameter tuning

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [194]: model params = {
               # number of
               'n estimators': [5,10,20,50,100,200],
               # number of max features
               'max features': [10,15,20,50,100],
               # max number of levels in each decision tree
               'max depth': [5,10,15],
               # minimum amount of samples in a node
               'min samples leaf': [10],
               # minimum amount of samples to split
               'min samples split' : [100,1000],
               # random state
               'random state': [1]
           executed in 3ms, finished 14:34:45 2021-06-17
In [195]: # Instantiate and fit the RandomForestClassifier
          rfc=RandomForestClassifier(class weight='balanced')
          rfc = GridSearchCV(rfc, model params, scoring = 'recall')
          rfc.fit(X train,y train)
           executed in 35m 3s, finished 15:09:48 2021-06-17
Out[195]: GridSearchCV(estimator=RandomForestClassifier(class weight='balance
                        param grid={'max depth': [5, 10, 15],
                                      'max features': [10, 15, 20, 50, 100],
                                     'min samples leaf': [10],
                                     'min samples split': [100, 1000],
                                     'n_estimators': [5, 10, 20, 50, 100, 200],
                                     'random state': [1]},
                        scoring='recall')
In [196]: print(rfc.best estimator .get params())
           executed in 6ms, finished 15:09:48 2021-06-17
           {'bootstrap': True, 'ccp alpha': 0.0, 'class weight': 'balanced', '
           res': 100, 'max leaf nodes': None, 'max samples': None, 'min impuri
           'min samples leaf': 10, 'min samples split': 1000, 'min weight frac
          None, 'oob score': False, 'random state': 1, 'verbose': 0, 'warm st
```

Contents <i>⊋</i> ♦
3.3.0 Elicode Chartel
3.3.7 Reading and Math Tests
3.3.8 Numeric Columns - nulls
3.3.9 Visualizations
▼ 3.4 df_district
▼ 3.4.1 Data Fields
3.4.1.1 Data Field Cleanup
3.4.2 Columns to drop
3.4.3 english_language_learners
▼ 4 Join Datasets
4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
5.1 Cleanup
5.2 Engineer % Columns
5.3 Reset Num List
▼ 6 Category Columns - cleanup
6.1 Final Check for NaNs
7 Train Test Split
▼ 8 Encode Features
8.1 X_train Encode
8.2 X_test Encode
▼ 9 Model Development
▼ 9.1 Logistic Regression
9.1.1 Check for Overfit
9.1.2 Model reiteration - parameter tuning
▼ 9.2 Random Forest Classifier
9.2.1 Check for Overfit
9.2.2 Feature Importance
9.2.3 Model reiteration - parameter tunin
▼ 9.3 Gradient Boosting Classifier
9.3.1 Check for Overfit
9.3.2 Model reiteration - parameter tunin
10 Conclusions

```
print('The f1 score for the model is:',f1 score(y test, y pred))
           executed in 13ms, finished 15:09:48 2021-06-17
           The fl score for the training model is: 0.9560906515580737
           The f1 score for the model is: 0.5279383429672447
In [198]: print(metrics.classification_report(y_train, y_pred_train, labels=[
                                                  target names=['normal grad rate
           executed in 16ms, finished 15:09:48 2021-06-17
                              precision
                                             recall f1-score
                                                                  support
                                                                    10313
           normal grad rate
                                    0.99
                                               1.00
                                                          1.00
              low grad rate
                                    1.00
                                               0.92
                                                          0.96
                                                                      736
                    accuracy
                                                          0.99
                                                                    11049
                  macro avg
                                    1.00
                                               0.96
                                                          0.98
                                                                    11049
               weighted avg
                                    0.99
                                               0.99
                                                          0.99
                                                                    11049
In [199]: print(metrics.classification_report(y_test, y_pred, labels=[0,1],
                                                  target_names=['normal grad rate
           executed in 37ms, finished 15:09:48 2021-06-17
                              precision
                                             recall f1-score
                                                                  support
           normal grad rate
                                    0.95
                                               0.99
                                                          0.97
                                                                     4392
              low grad rate
                                    0.78
                                               0.40
                                                          0.53
                                                                      344
                                                          0.95
                                                                     4736
                    accuracy
                  macro avg
                                    0.87
                                               0.69
                                                          0.75
                                                                     4736
               weighted avg
                                    0.94
                                               0.95
                                                          0.94
                                                                     4736
```

In [197]: print('The f1 score for the training model is:', f1 score(y train, y

9.3 Gradient Boosting Classifier

In [200]: gb clf = GradientBoostingClassifier()

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
```

```
gb_clf.fit(X_train, y train)
           y pred = gb clf.predict(X test)
           y pred train = gb clf.predict(X train)
           # print("Accuracy score (training): {0:.3f}".format(qb clf.score(X)
           # print("Accuracy score (validation): {0:.3f}".format(qb clf.score)
           executed in 10.6s, finished 15:09:59 2021-06-17
In [201]: print('The fl score for the training model is:', fl score(y train, y
           print('The f1 score for the model is:',f1 score(y test, y pred))
           executed in 23ms, finished 15:09:59 2021-06-17
           The f1 score for the training model is: 0.7556596409055425
           The f1 score for the model is: 0.6079447322970639
In [202]: print(metrics.classification report(y test, y pred, labels=[0,1],
                                                  target names=['normal grad rate
           executed in 10ms, finished 15:09:59 2021-06-17
                              precision
                                             recall f1-score
                                                                  support
           normal grad rate
                                    0.96
                                               0.99
                                                          0.97
                                                                     4392
              low grad rate
                                    0.75
                                               0.51
                                                          0.61
                                                                      344
                                                          0.95
                                                                     4736
                    accuracy
                  macro avg
                                    0.86
                                               0.75
                                                          0.79
                                                                     4736
               weighted avg
                                    0.95
                                               0.95
                                                          0.95
                                                                     4736
```

9.3.1 Check for Overfit

9.3.2 Model reiteration - parameter tunin

Contents & A
3.3.7 Reading and Math Tests
3.3.8 Numeric Columns - nulls
3.3.9 Visualizations
▼ 3.4 df_district
▼ 3.4.1 Data Fields
3.4.1.1 Data Field Cleanup
3.4.2 Columns to drop
3.4.3 english_language_learners
▼ 4 Join Datasets
4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
5.1 Cleanup
5.2 Engineer % Columns
5.3 Reset Num List
▼ 6 Category Columns - cleanup
6.1 Final Check for NaNs
7 Train Test Split
▼ 8 Encode Features
8.1 X_train Encode
8.2 X_test Encode
▼ 9 Model Development
▼ 9.1 Logistic Regression
9.1.1 Check for Overfit
9.1.2 Model reiteration - parameter tuning
▼ 9.2 Random Forest Classifier
9.2.1 Check for Overfit
9.2.2 Feature Importance
9.2.3 Model reiteration - parameter tunin
▼ 9.3 Gradient Boosting Classifier
9.3.1 Check for Overfit
9.3.2 Model reiteration - parameter tunin
10 Conclusions

```
In [204]: # View confusion matrix for test data and predictions
           confusion matrix(y test, y pred)
           executed in 8ms, finished 15:09:59 2021-06-17
Out[204]: array([[4333, 59],
                  [ 168, 176]])
In [205]: print(metrics.classification report(y train, y pred train, labels=[
                                                  target names=['normal grad rate
           executed in 16ms, finished 15:09:59 2021-06-17
                              precision
                                            recall f1-score
                                                                 support
           normal grad rate
                                    0.98
                                               0.99
                                                          0.98
                                                                   10313
              low grad rate
                                    0.89
                                                          0.76
                                                                      736
                                               0.66
                                                          0.97
                                                                   11049
                    accuracy
                                                          0.87
                  macro avg
                                    0.93
                                               0.83
                                                                   11049
               weighted avg
                                                          0.97
                                    0.97
                                               0.97
                                                                   11049
In [206]: print(metrics.classification_report(y_test, y_pred, labels=[0,1],
                                                  target_names=['normal grad rate
           executed in 11ms, finished 15:09:59 2021-06-17
                              precision
                                            recall f1-score
                                                                 support
           normal grad rate
                                    0.96
                                               0.99
                                                          0.97
                                                                    4392
              low grad rate
                                    0.75
                                               0.51
                                                          0.61
                                                                     344
                    accuracy
                                                          0.95
                                                                    4736
                                                          0.79
                  macro avg
                                    0.86
                                               0.75
                                                                    4736
               weighted avg
                                    0.95
                                               0.95
                                                          0.95
                                                                    4736
```

9.3.2 Model reiteration - parameter tuning

```
3.3.0 Encode Charter
      3.3.7 Reading and Math Tests
      3.3.8 Numeric Columns - nulls
      3.3.9 Visualizations
  ▼ 3.4 df_district

▼ 3.4.1 Data Fields

        3.4.1.1 Data Field Cleanup
      3.4.2 Columns to drop
      3.4.3 english_language_learners
▼ 4 Join Datasets
    4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
    5.1 Cleanup
    5.2 Engineer % Columns
    5.3 Reset Num List
▼ 6 Category Columns - cleanup
    6.1 Final Check for NaNs
  7 Train Test Split
▼ 8 Encode Features
    8.1 X train Encode
    8.2 X test Encode
▼ 9 Model Development
  ▼ 9.1 Logistic Regression
      9.1.1 Check for Overfit
      9.1.2 Model reiteration - parameter tuning
  ▼ 9.2 Random Forest Classifier
      9.2.1 Check for Overfit
      9.2.2 Feature Importance
      9.2.3 Model reiteration - parameter tunin
  ▼ 9.3 Gradient Boosting Classifier
      9.3.1 Check for Overfit
      9.3.2 Model reiteration - parameter tunin
  10 Conclusions
```

```
In [207]: model params = {
               # number of boosting stages
               'n_estimators': [5,10,20,50,100,200],
               # number of max features
               'max features': [10,15,20,50,100],
                # minimum amount of samples in a node
               'min samples leaf': [10],
               # Learning rate
               'learning rate':[.25,.5,.75,1],
               #The minimum number of samples required to split an internal no
               'min samples split' : [100,1000],
                # random state
               'random state':[1]
           executed in 3ms, finished 15:09:59 2021-06-17
In [208]: | gb clf = GradientBoostingClassifier()
          gb clf = GridSearchCV(gb clf,model params, scoring = 'recall')
           gb clf.fit(X train,y train)
           executed in 30m 25s, finished 15:40:24 2021-06-17
Out[208]: GridSearchCV(estimator=GradientBoostingClassifier(),
                        param_grid={'learning_rate': [0.25, 0.5, 0.75, 1],
                                      'max features': [10, 15, 20, 50, 100],
                                      'min samples_leaf': [10],
                                      'min_samples_split': [100, 1000],
                                      'n_estimators': [5, 10, 20, 50, 100, 200],
                                      'random state': [1]},
                        scoring='recall')
```

Contents & \$
3.3.7 Reading and Math Tests
3.3.8 Numeric Columns - nulls
3.3.9 Visualizations
▼ 3.4 df_district
▼ 3.4.1 Data Fields
3.4.1.1 Data Field Cleanup
3.4.2 Columns to drop
3.4.3 english_language_learners
▼ 4 Join Datasets
4.1 Check for NaNs
▼ 5 Numeric Columns - cleanup
5.1 Cleanup
5.2 Engineer % Columns
5.3 Reset Num List
▼ 6 Category Columns - cleanup
6.1 Final Check for NaNs
7 Train Test Split
▼ 8 Encode Features
8.1 X_train Encode
8.2 X_test Encode
▼ 9 Model Development
▼ 9.1 Logistic Regression
9.1.1 Check for Overfit
9.1.2 Model reiteration - parameter tuning
▼ 9.2 Random Forest Classifier
9.2.1 Check for Overfit
9.2.2 Feature Importance
9.2.3 Model reiteration - parameter tunin
▼ 9.3 Gradient Boosting Classifier
9.3.1 Check for Overfit
9.3.2 Model reiteration - parameter tunin
J.J. Hoder followers paramoter terms

In [209]:	print(metrics.cla	ssification_			red_train, l =['normal gr	
	executed in 27ms, finished	15:40:24 2021-06-17				
		precision	recall	f1-score	support	
	normal grad rate	0.98	0.99	0.98	10313	
	low grad rate	0.89	0.66	0.76	736	
	accuracy			0.97	11049	
	macro avg	0.93	0.83	0.87	11049	
	weighted avg	0.97	0.97	0.97	11049	
In [210]:	print(metrics.cla				ed, labels=[=['normal gr	
	executed in Tims, illistical	precision	recall	f1-score	support	
	normal grad rate	0.96	0.99	0.97	4392	
	low grad rate	0.75	0.51	0.61	344	
	accuracy			0.95	4736	
	macro avg	0.86	0.75	0.79	4736	
	macio avq					

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- **▼** 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

executed in 16.4s, finished 15:42:51 2021-06-17

```
Learning rate: 0.05
Accuracy score (training): 0.942
Accuracy score (validation): 0.938
The f1 score for the training model is: 0.7556596409055425
The f1 score for the model is: 0.6079447322970639
Learning rate: 0.075
Accuracy score (training): 0.949
Accuracy score (validation): 0.942
The fl score for the training model is: 0.7556596409055425
The f1 score for the model is: 0.6079447322970639
Learning rate: 0.1
Accuracy score (training): 0.953
Accuracy score (validation): 0.945
The f1 score for the training model is: 0.7556596409055425
The f1 score for the model is: 0.6079447322970639
Learning rate: 0.25
Accuracy score (training): 0.962
Accuracy score (validation): 0.948
The fl score for the training model is: 0.7556596409055425
```

- Our modeling struggles with overfitting and is better at classifying the majority class (high enough performance that we feel confident in making recommendations based our modeling.
- · The number of students who completed math and reading tests and received a profic

- The number of students participating in ACT/SAT tests was a strong indicator for our normal graduation rate group.
- Higher % of total salary devoted to teacher salaries was an indicator of normal high s

- 3.3.0 Encode Charter
- 3.3.7 Reading and Math Tests
- 3.3.8 Numeric Columns nulls
- 3.3.9 Visualizations
- ▼ 3.4 df_district
 - ▼ 3.4.1 Data Fields
 - 3.4.1.1 Data Field Cleanup
 - 3.4.2 Columns to drop
 - 3.4.3 english_language_learners
- ▼ 4 Join Datasets
 - 4.1 Check for NaNs
- ▼ 5 Numeric Columns cleanup
 - 5.1 Cleanup
 - 5.2 Engineer % Columns
 - 5.3 Reset Num List
- ▼ 6 Category Columns cleanup
 - 6.1 Final Check for NaNs
 - 7 Train Test Split
- ▼ 8 Encode Features
 - 8.1 X train Encode
 - 8.2 X test Encode
- ▼ 9 Model Development
 - ▼ 9.1 Logistic Regression
 - 9.1.1 Check for Overfit
 - 9.1.2 Model reiteration parameter tuning
 - ▼ 9.2 Random Forest Classifier
 - 9.2.1 Check for Overfit
 - 9.2.2 Feature Importance
 - 9.2.3 Model reiteration parameter tunin
 - ▼ 9.3 Gradient Boosting Classifier
 - 9.3.1 Check for Overfit
 - 9.3.2 Model reiteration parameter tunin
 - 10 Conclusions

11 Recommendations

- Mathematics and Reading: Focus resources on providing extra support for mathemat to improve these scores whether it be through offering remedial opportunities additionareas.
- ACT/SAT: As with many education metrics ACT/SAT participation likely also has a relationary regardless of other factors ACT/SAT prep can be a motivating factor for students in high school is to prepare themselves so they will be accepted into
- · Teacher Salaries:
 - This feature may be confounded by other factors. For example a school with less workers if they are unable to receive funding for newer facilities.
 - We do not recommend making teacher salaries a higher proportion in budgeting factors.

12 Future Work

- Perform feature selection steps to improve overfit of the modeling.
- Investigate adding new features. Now that we have a better handle on which types of features from our reporting sources and also engineering new features with the data \(\)
- Investigate which types of schools did not report graduation rates. We need to know graduation rates. Is the type of schools not reporting graduation rates similar to those