



# New York City High School Success



Mark Brennan



# Intro

---

- In this project, I attempt to classify the success of NYC public high schools.
- Success is relative term, but given three years of NYC high school graduation outcomes (2015, 2016, 2017), I try to classify success as a graduation rate above the mean (73% of a given cohort graduates).
- This is as much an exploration of the rich feature set as it is exercise in prediction.

# Data

---

- NYC Open Data

A rich trove of education reports in .xlsx and/or .csv formats

- Two main sets:

- High School Outcomes - denotes percent of graduating cohorts for each year
- High School Quality Reports - yearly reports capturing rich qualitative and quantitative metrics - 40+ columns.

# NYC Open Data - NYC DOE

---

“The Quality Review is a process that evaluates how well schools are organized to support student learning and teacher practice. It was developed to assist New York City Department of Education (NYCDOE) schools in raising student achievement by looking behind a school’s performance statistics to ensure that the school is engaged in effective methods of accelerating student learning.”

# Data Prep

---

- Select observations and “labels”:
  - There are only about 480 NYC public high schools, so I selected three years of graduation outcomes (2015, 2016, and 2017) for 400+ high schools, yielding **1,225 observations**.
- Join targets and features
- Clean data
- Create label encodings for categorical data

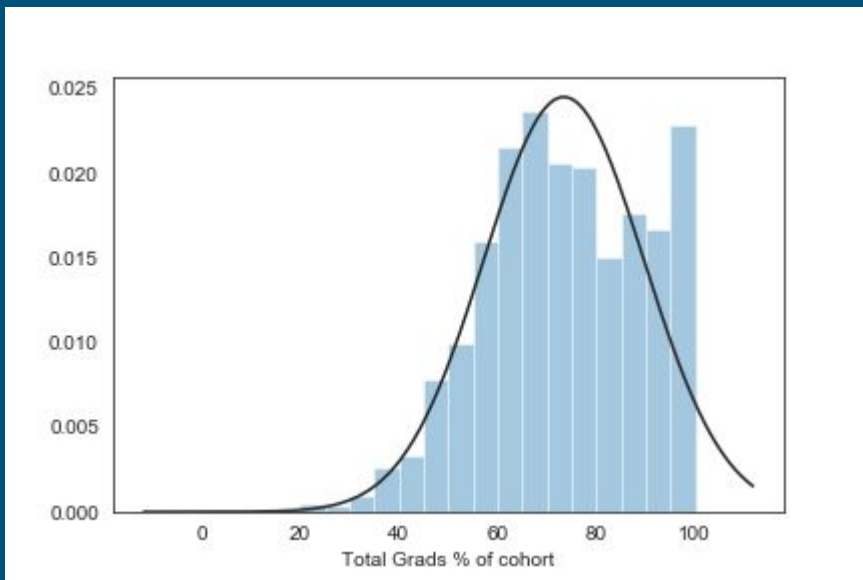
# Data Prep

---

- Select observations and “labels”:
  - There are only about 480 NYC public high schools, so I selected three years of graduation outcomes (2015, 2016, and 2017) for 400+ high schools, yielding **1,225 observations**.
- Join targets and features (approx 37!)
- Clean data
- Create label encodings for categorical data

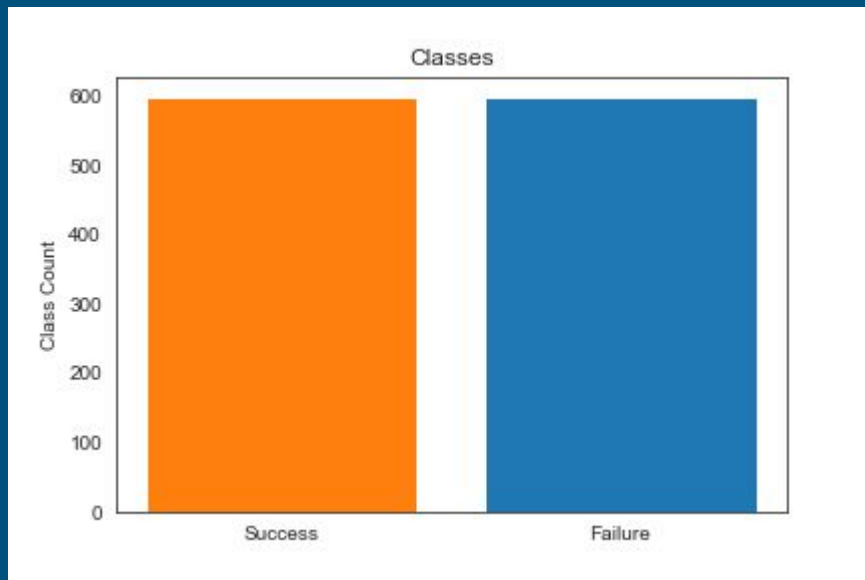
# Target distribution

- The target is the total graduates as a percentage of their cohort:



# Classification

- We create two classes - success or failure, based on graduation rates. Given the cutoff at the mean (73% graduation rate), there is an even split of the two classes, “Success” and “Failure”:





# Baseline, “dummy” model

=====

STRATIFIED DUMMY CLASSIFIER

Accuracy is: 51.02040816326531

AUC is: 0.52

-----

Confusion Matrix

col_0	0	1	All
-------	---	---	-----

target
--------

0	65	47	112
---	----	----	-----

1	73	60	133
---	----	----	-----

All	138	107	245
-----	-----	-----	-----

-----

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.47	0.58	0.52	112
---	------	------	------	-----

1	0.56	0.45	0.50	133
---	------	------	------	-----

accuracy	0.51	245
----------	------	-----

macro avg	0.52	0.52	0.51	245
-----------	------	------	------	-----

weighted avg	0.52	0.51	0.51	245
--------------	------	------	------	-----

# KNN

KNN

Accuracy is: 59.183673469387756

AUC is: 0.59

Confusion Matrix

col_0	0	1	All
-------	---	---	-----

target			
--------	--	--	--

0	64	48	112
---	----	----	-----

1	52	81	133
---	----	----	-----

All	116	129	245
-----	-----	-----	-----

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.55	0.57	0.56	112
---	------	------	------	-----

1	0.63	0.61	0.62	133
---	------	------	------	-----

accuracy		0.59		245
----------	--	------	--	-----

macro avg	0.59	0.59	0.59	245
-----------	------	------	------	-----

weighted avg	0.59	0.59	0.59	245
--------------	------	------	------	-----

# Decision Tree

```
=====
DECISION TREE
```

```
Accuracy is: 83.26530612244898
```

```
AUC is: 0.83
```

```
-----
Confusion Matrix
```

```
col_0  0   1 All
```

```
target
```

```
0      90  22 112
```

```
1      19 114 133
```

```
All    109 136 245
```

```
-----
precision recall f1-score support
```

```
0      0.83   0.80   0.81    112
```

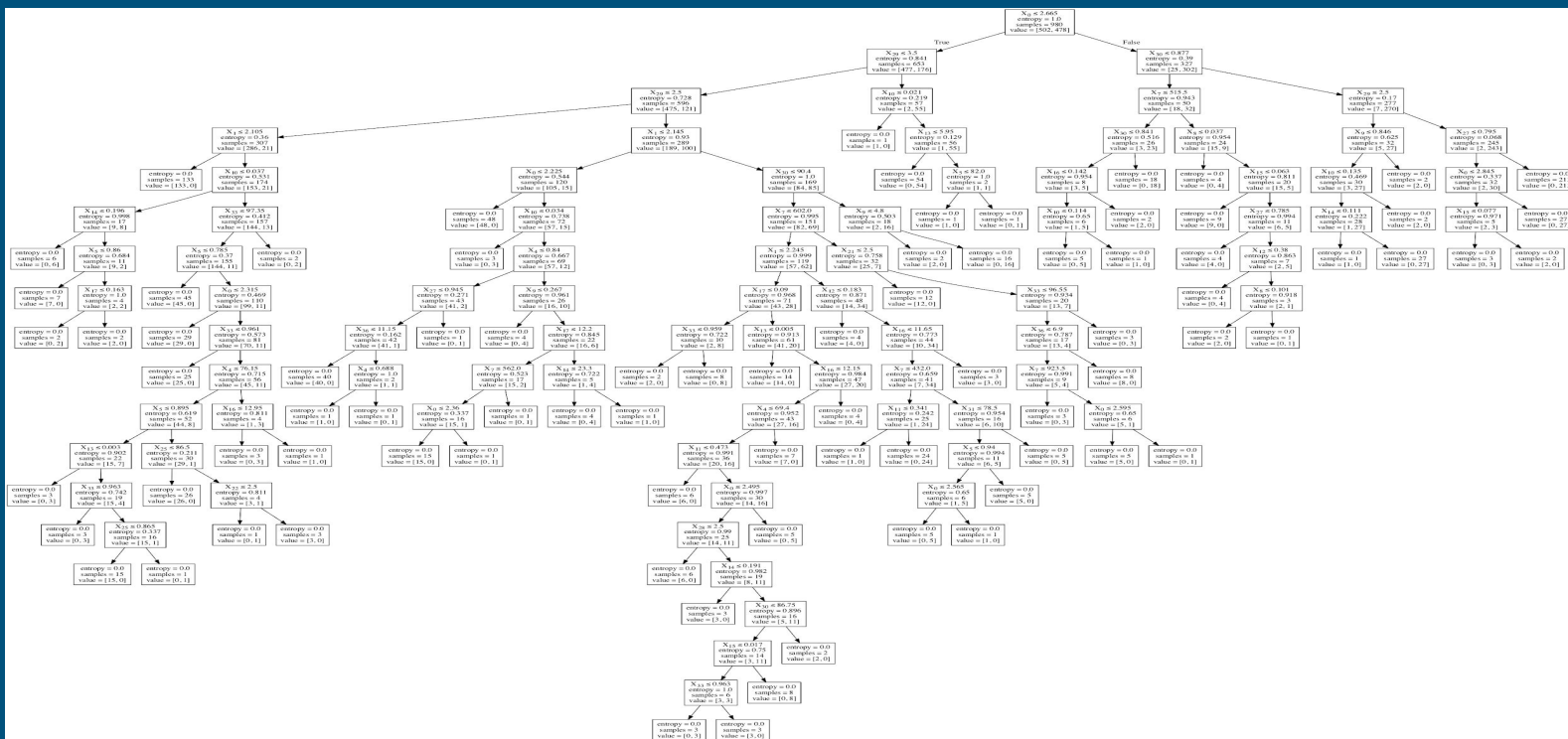
```
1      0.84   0.86   0.85    133
```

```
accuracy          0.83    245
```

```
macro avg         0.83   0.83   0.83    245
```

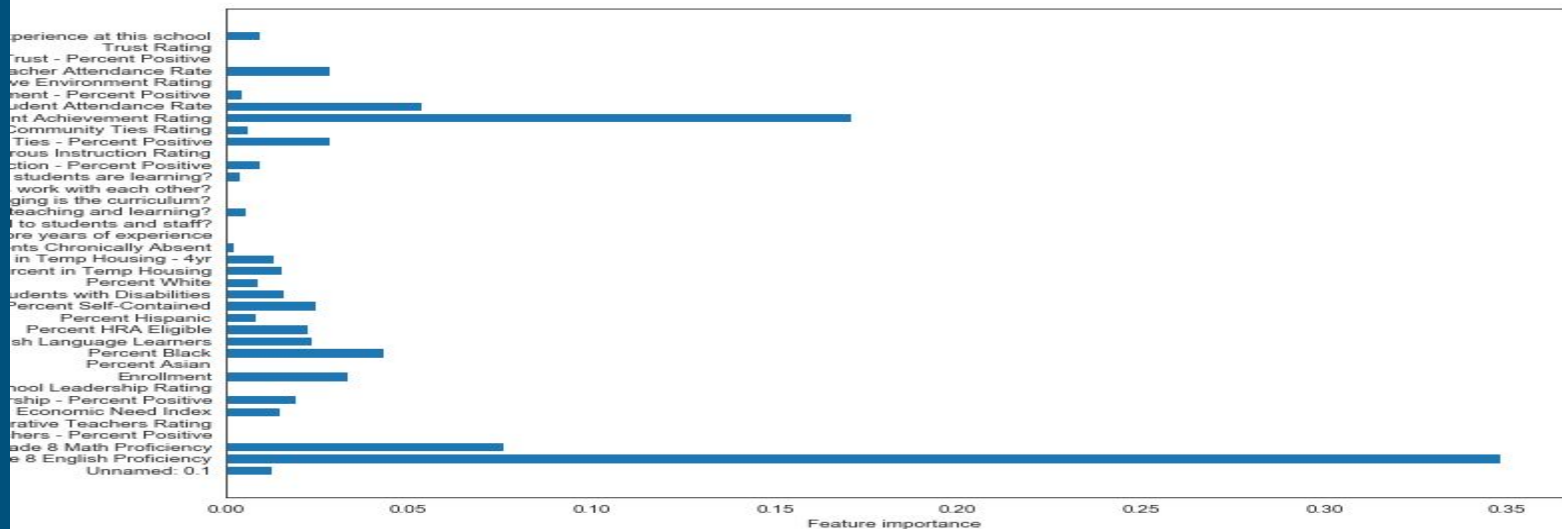
```
weighted avg      0.83   0.83   0.83    245
```

# Decision Tree



# Decision Tree Feature Selection

- A handful of features stand out as driving the classification:



# Random Forest

```
=====
DEFAULT RANDOM FOREST
```

```
Accuracy is: 84.89795918367346
```

```
AUC is: 0.85
```

```
-----
Confusion Matrix
```

```
col_0  0  1 All
```

```
target
```

```
0      97  15 112
```

```
1      22 111 133
```

```
All   119 126 245
```

```
-----
precision  recall f1-score  support
```

```
0      0.82    0.87    0.84    112
```

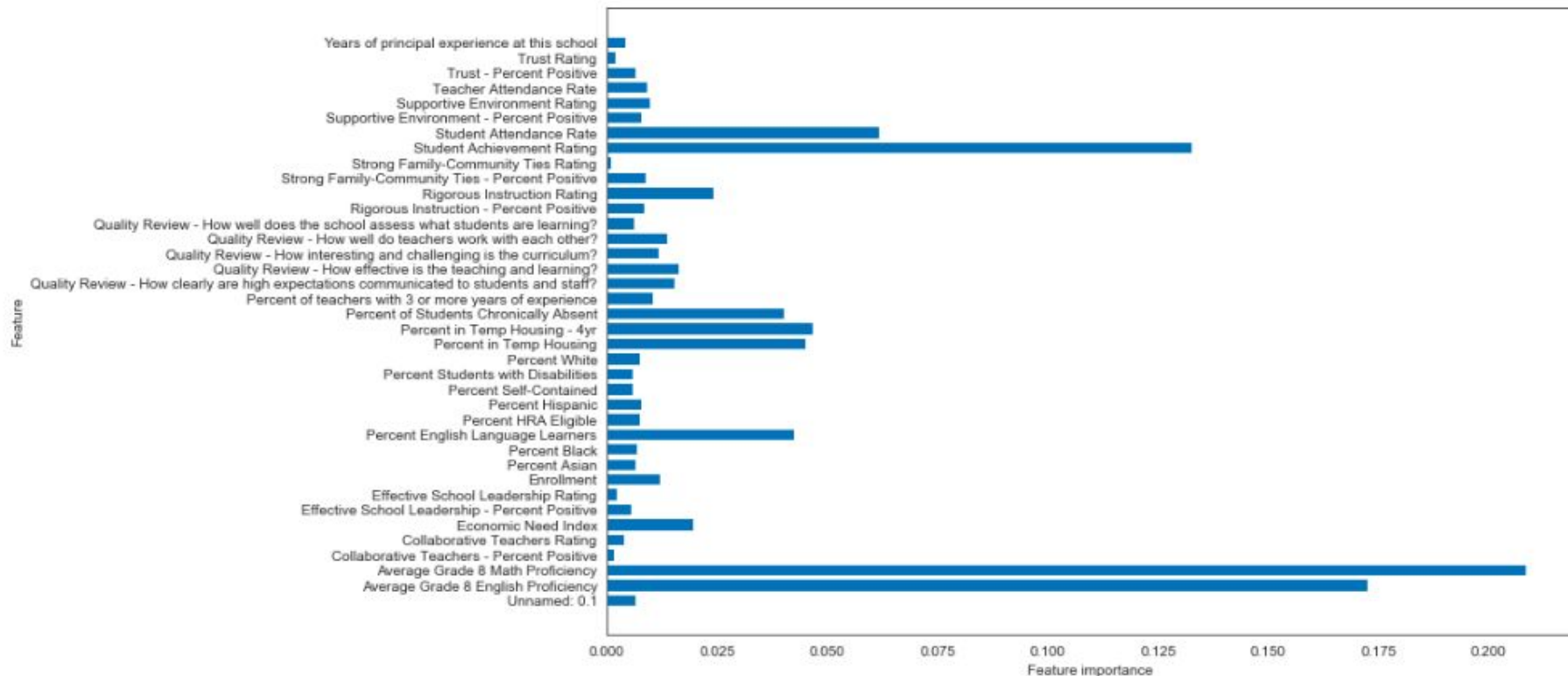
```
1      0.88    0.83    0.86    133
```

```
accuracy                0.85    245
```

```
macro avg      0.85    0.85    0.85    245
```

```
weighted avg    0.85    0.85    0.85    245
```

# Random Forest Feature Selection



# XGBoost

=====

XGBOOST

Accuracy is: 86.12244897959184

AUC is: 0.86

-----

Confusion Matrix

col_0	0	1	All
-------	---	---	-----

target
--------

0	98	14	112
---	----	----	-----

1	20	113	133
---	----	-----	-----

All	118	127	245
-----	-----	-----	-----

-----

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.83	0.88	0.85	112
---	------	------	------	-----

1	0.89	0.85	0.87	133
---	------	------	------	-----

accuracy	0.86	245
----------	------	-----

macro avg	0.86	0.86	0.86	245
-----------	------	------	------	-----

weighted avg	0.86	0.86	0.86	245
--------------	------	------	------	-----



# Conclusion And Next Steps

---

Not surprisingly, core metrics of educational performance, including “Average Grade 8 English Proficiency”, and “Average Grade 8 Math Proficiency” were among the main drivers of a school’s success!

- Next Steps
  - The feature engineering and modeling exercises were fairly coarse;
  - Engineering new features, or tuning the models differently, may yield less obvious insights.