# 2013-2019 School Attendance

By Jeremy and Frank

# Background Information

Attendance data includes students in district 1-32, 75 (Special Education), district 79 (Alternative Schools & Programs), charter schools, home schooling. Home and hospital instruction are excluded. Pre-K data does not include NYC Early Education Centers or District Pre-K Centers therefore data is limited to those who attend K-12 schools that offer Pre-K. Transfer school counts are not in school level file. Attendance is registered to school student is attending at the time. If a student attend multiple schools in a school year the data will be reflected in multiple schools. Chronical absence is defined if a student has an attendance rate of less than 90 percent ( students who are absent 10 percent or more of the total days).
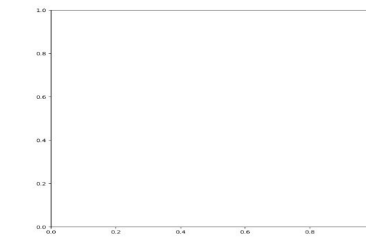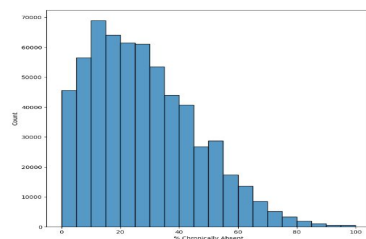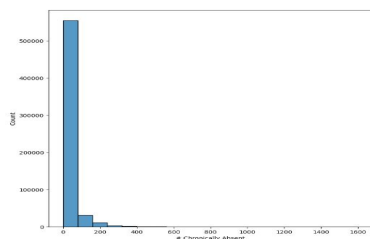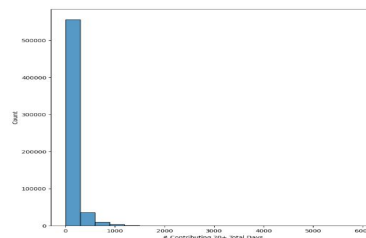
Link to Dataset: https://data.cityofnewyork.us/Education/2013-2019-Attendance-Results-School/vww9-qguh

# Features

- #total day: The combination of #days present and # days absent. It is also the total number of days of school for all students on that row.
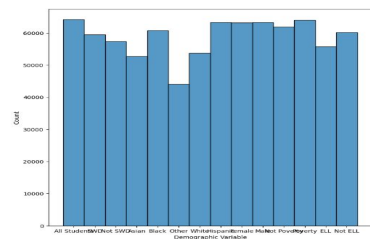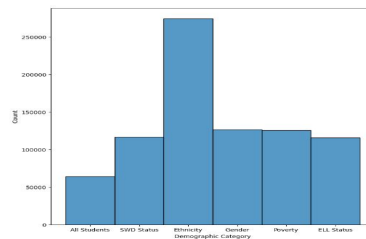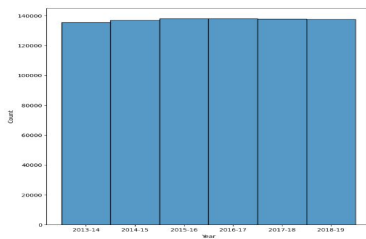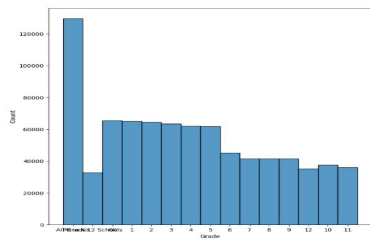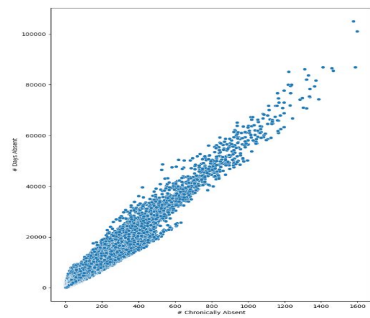- #days absent: The total number of days all students on that row are absent.
- #days present: The total number of days that all students on that row are present.(unsure how data is collected for this, needs more research)
- % Attendance: Percentage of total days that all students on that row are present.
- # Contributing 20+ Total Days: The total number of students in the specific row that attended school for at least 20 days.
- # Chronically Absent: Total number of students who were chronically absent.
- % Chronically Absent: The percent of students who are chronically absent out of the total number of students.
- Demographics Category: Contains the types of demographics used when collecting the data. There is All students, SWD status, Ethnicity, Gender, Poverty and ELL status.

# Demographic Categories

- All Students: This demographic consists of all students. It does not separate the student data into separate rows: Female and Male
- SWD status: Refers to whether a student has a disability or not.
- Ethnicity: Refers to the students' ethnicity: white, asian, black, hispanic, and other.
- Gender: Separates student data based on gender, male and female.
- Poverty: Separates students into two different categories: poverty and not poverty
- ELL status: Separates students into whether English is their first language or not.
- DBN: The distribute borough number of the students at that row
- Year: The year of the students at that row
- Grade: The grade of the students at that row
- School: The school at which those students attend.

# Objective

The goal is to find whether a student is classified as SWD(Student with Disability) or not SWD. The 3 independent variables are the number of days absent, the total number of students in the specific row that attended school for at least 20 days and the grade that group of students are in.

The dependent variable is a classification of a student who is classified as SWD or not SWD

0 = Not SWD
1 = SWD

Z is an ordinal encoding of the Grade column

Z is an ordinal encoding of the Grade column

Z is an ordinal
encoding of the
Grade column

Y is a Binary
Encoding of the
SWD category in
the Demographic
Variable Column

Y is a Binary Encoding of the SWD category in the Demographic Variable Column

Y is a Binary Encoding of the SWD category in the Demographic Variable Column

# Objective Model 1

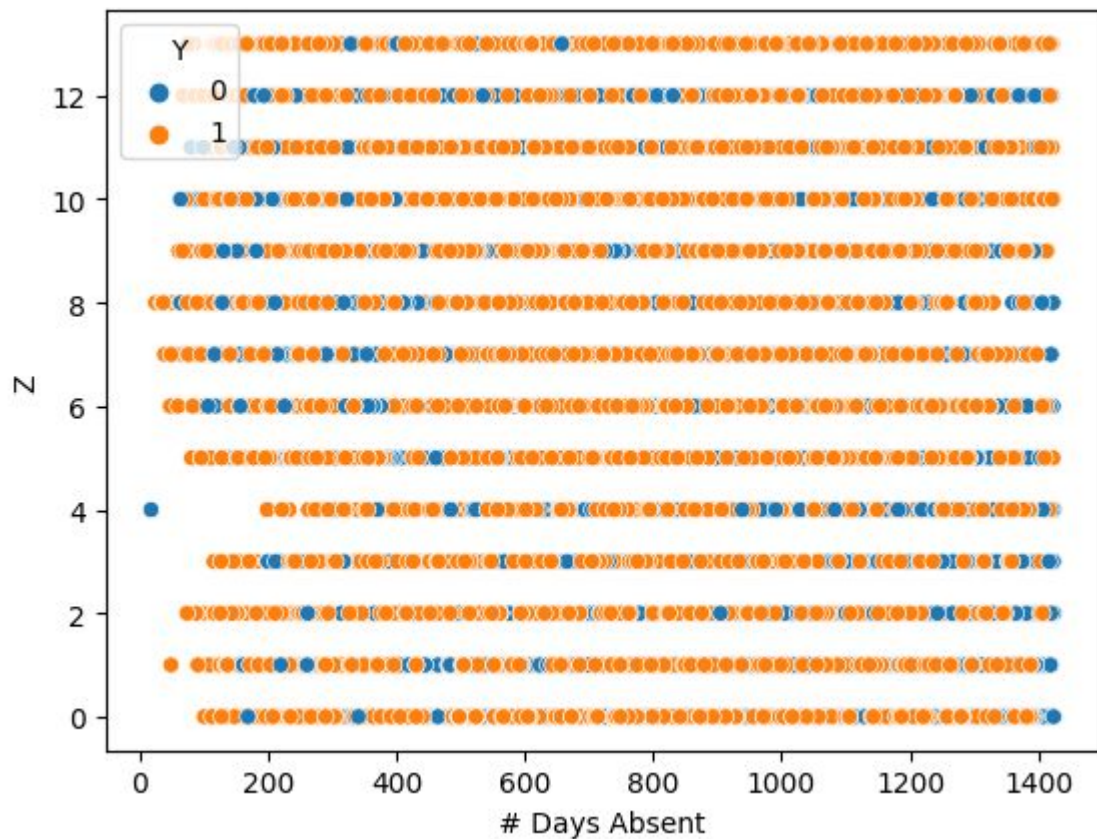We will use Linear Regression to portray the relationship between the three independent variables, the number of days absent, the total number of students in the specific row that attended school for at least 20 days and the grade that group of students are in. The dependent variable is the SWD status of that group of students surveyed.

# Linear Regression Model

R-squared Score: 0.3547830570047862

Coefficients:

    # Days Absent: 0.00033351

    # Contributing 20+ Total Days: -0.01525098

    Z: 0.02973281

Y Intercept: 0.8137906450124841

[ ]

```
X_train_LR, X_test_LR, y_train_LR,
y_test_LR = train_test_split(X, y,
test_size=0.3,)

from sklearn.linear_model import
LinearRegression
lr = LinearRegression()
model_LR=lr.fit(X_train_LR,
y_train_LR)

from sklearn.metrics import
r2_score

y_pred_LR =
model_LR.predict(X_test_LR)
r2_score_LR=r2_score(y_test_LR,
y_pred_LR)
slopes_LR=model_LR.coef_
intercept_LR=model_LR.intercept_
print(r2_score_LR)
print(slopes_LR)
print(intercept_LR)
```

```
sns.regplot(x="# Days
Absent", y=y_train_LR,
data=X_train_LR,
line_kws={"color":
"red"})
```

```
sns.regplot(x="#
Contributing 20+ Total
Days", y=y_train_LR,
data=X_train_LR,
line_kws={"color":
"red"})
```

```
sns.regplot(x="Z",
y=y_train_LR,
data=X_train_LR,
line_kws={"color":
"red"})
```

# Objective Model 2

We will use Logistic Regression to portray the relationship between the three independent variables, the number of days absent, 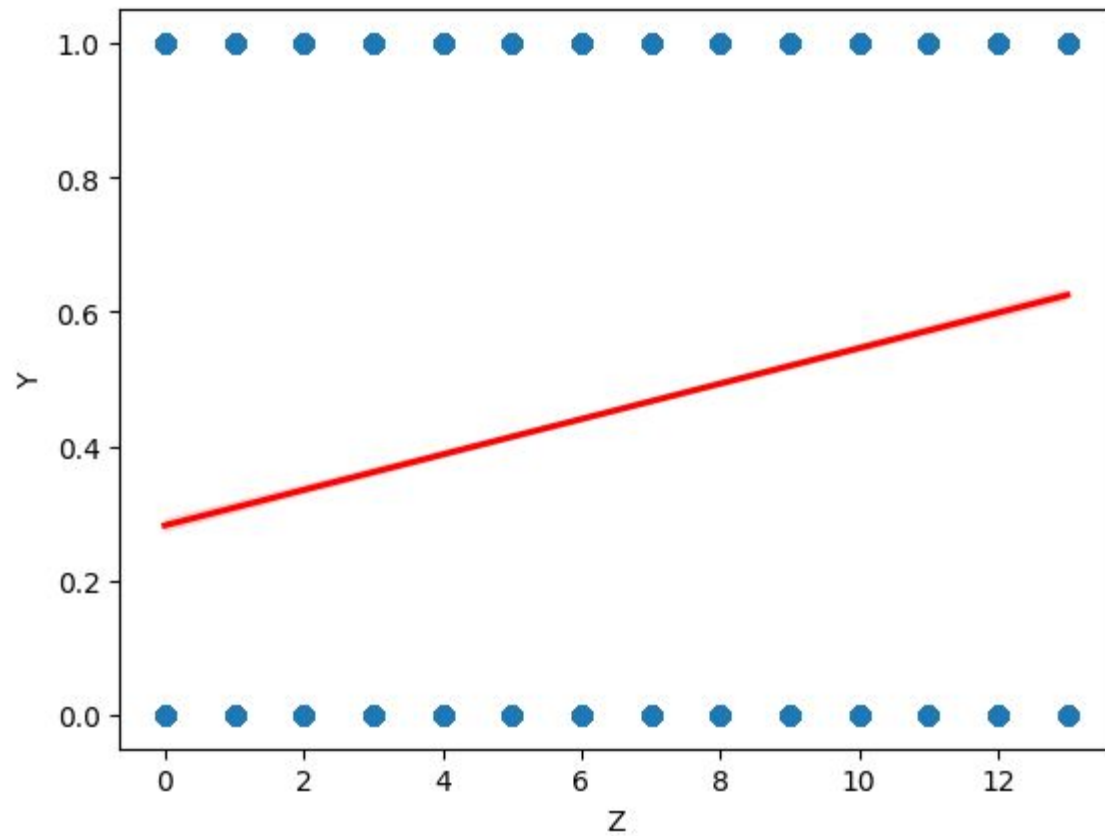the total number of students in the specific row that attended school for at least 20 days and the grade that group of students are in. The dependent variable is the SWD status of that group of students surveyed.

# Logistic Regression Model

Accuracy Score:   0.7836212801583073

Coefficients:

   Z: 0.20169294

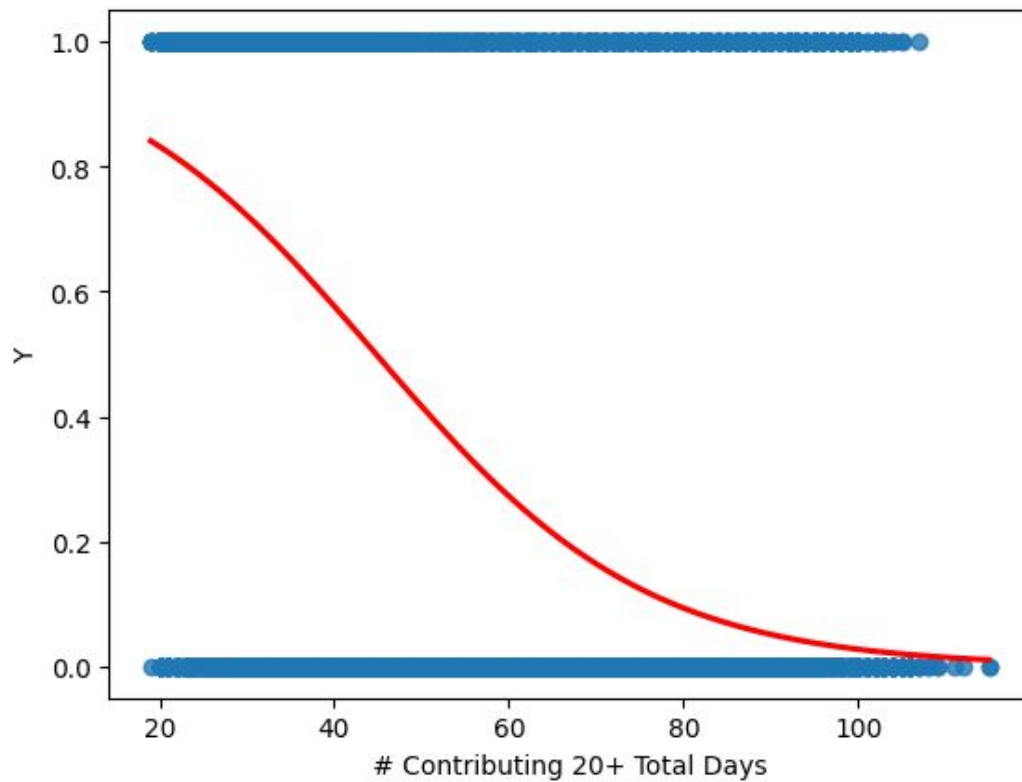   # Days Absent: 0.00254206

   # Contributing 20+ Total Days: -0.1006931

Y Intercept:   1.67466691

```python
from sklearn.model_selection import
train_test_split
X_train_LG, X_test_LG, y_train_LG, y_test_LG =
train_test_split(X, y, test_size=0.3,)

from sklearn.linear_model import
LogisticRegression
lgr = LogisticRegression()
model_LG=lgr.fit(X_train_LG, y_train_LG)

from sklearn.metrics import accuracy_score
y_pred_LG = model_LG.predict(X_test_LG)
accuracy_LG=accuracy_score(y_test_LG,y_pred_LG
,normalize=True)
slopes_LG=model_LG.coef_
intercept_LG=model_LG.intercept_
```
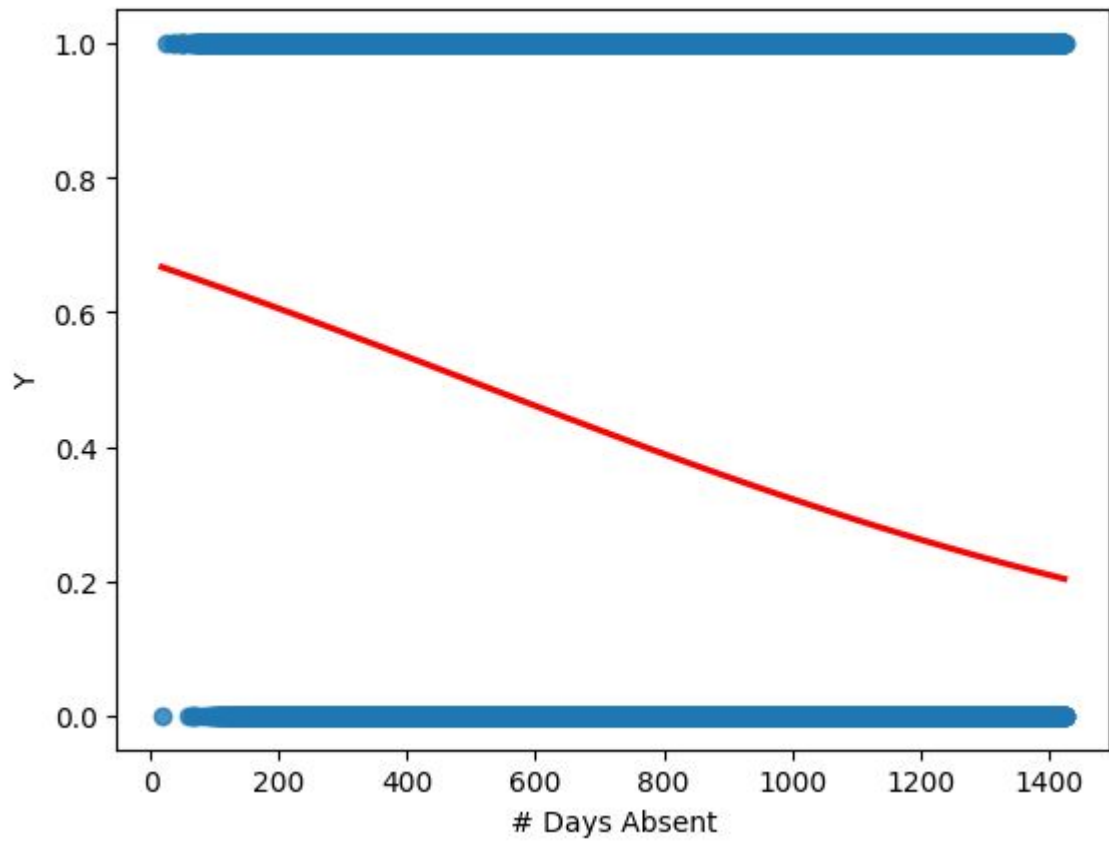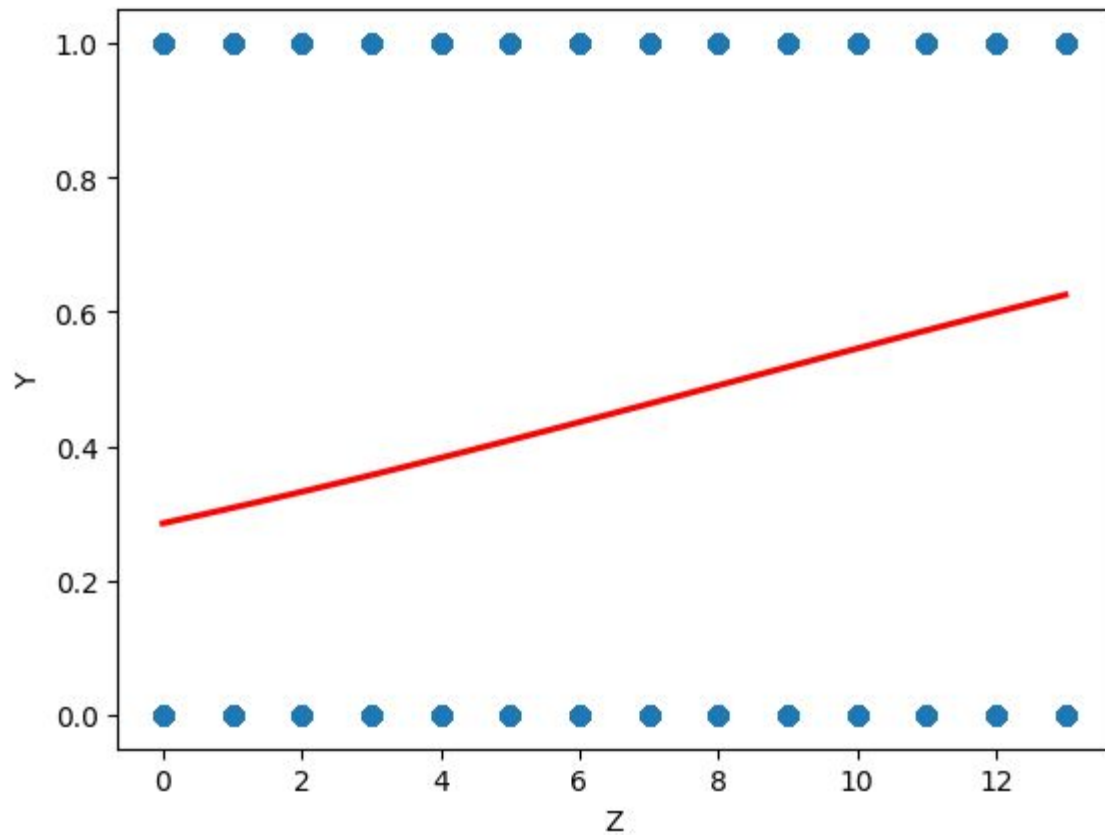
```
sns.regplot(x=log_data.loc[:,"#
Contributing 20+ Total Days"], y=y,
data=log_data, logistic=True,
ci=None, line_kws={"color": "red"})
```

```
sns.regplot(x=log_data.loc[:,"#
Days Absent"], y=y, data=log_data,
logistic=True, ci=None,
line_kws={"color": "red"})
```

```
sns.regplot(x=log_data.
loc[:,"Z"], y=y,
data=log_data,
logistic=True, ci=None,
line_kws={"color":
"red"})
```

# Conclusion

The R2 score of Linear Regression is lower than 0.5 while the accuracy score of Logistic Regression is greater than 0.5. Not only is the accuracy score higher for Logistic Regression, the graph for Linear Regression shows that the model will eventually be useless as it goes outside of the upper and lower bounds (0, 1). Therefore, the Logistic Regression model is the better model to use to predict the probability that a group of students that were surveyed were SWD.