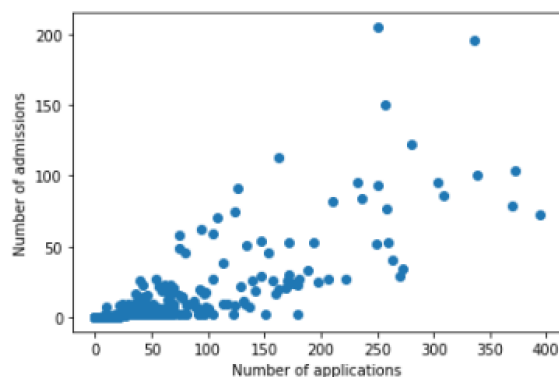
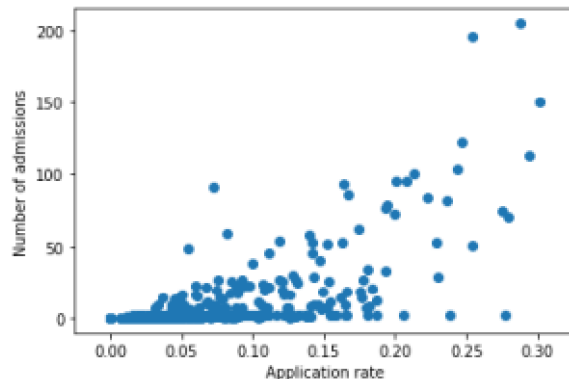


IDS Final Project

- 1) In order to find the correlation between the number of applications and admissions I turned each appropriate column into its own series and used `corr()` to find the correlations between the two variables. The correlation between the number of applications and admissions to HSPHS is $r = 0.8017265370719315$. This means that there is a strong correlation between the number of applications and the number of admissions. As more applications are sent from a school, more applications are accepted at HSPHS from that school.



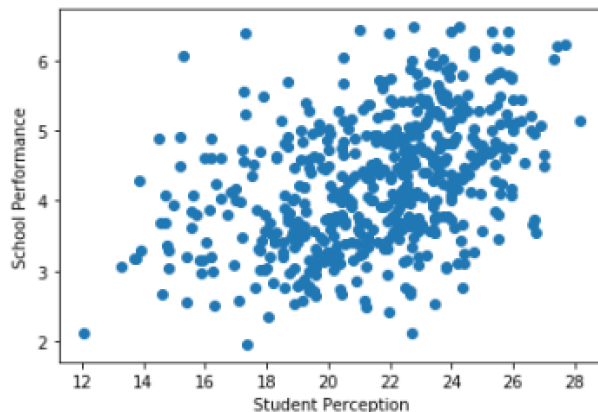
- 2) In order to compare applications and application rates' correlation with acceptances, I did the same steps as Question 1, but with the application rate, which was calculated by dividing applications by the school size. The correlation between the application rate and admissions is $r = 0.658750752900268$. Being lower than 0.80171, the application rate has a lower correlation with admissions than the raw number of applications. Therefore, the raw number of applications from a school is a better predictor for admissions than the application rate.



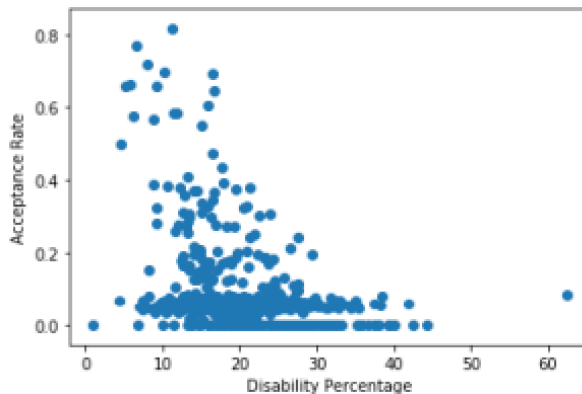
- 3) In order to find the best *per student* odds of sending someone to HSPHS, I first created a separate column with the per-student probability of getting into HSPHS by taking the acceptance column value and dividing it by the school size column. Then, using that column, I used the `max()` function to find that the highest per-students probability of getting accepted from a single school is 0.23482245131729668. Using the highest per-student odds, I located the index where the per

student odds equaled the max_value, and used that index to locate the school name using loc(). The result is The Christa McCauliffe School.

- 4) Finding the relationship between student perceptions about their schools and school performance based on achievements required that there were separate columns created in the dataframe combining the different factors that make up each category. For student perception, the columns that represent rigorous instruction, collaborative teachers, supportive environment, effective school leadership, and strong family community were added together to create singular student perceptions values. For school performance, the columns representing students achievement, reading scores exceed, and math scores exceed were added together to create singular school performance values. Then the correlation between the two new values was calculated, which turned out to be $r = 0.41874194039783613$.

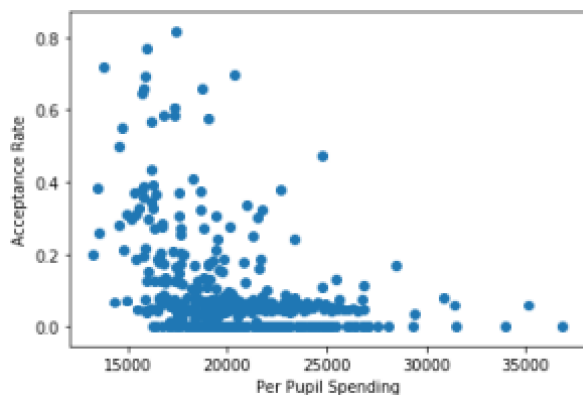


- 5) While looking through the variables available to compare, the disability percentage stood out to me. Disability can include learning disabilities as well, which would logically result in lower test scores. I wanted to test whether a higher disability percentage would lead to a lower acceptance rate, which may not be as strong as schools with lower disability percentages. The hypothesis that I decided to test was the higher the percentage rate of students with disabilities (independent variable), the lower the acceptance rate to HSPHS would be (dependent variable). In order to do this, I simply calculated the acceptance rate by dividing the acceptances by the number of applications and found the correlation between that value and disability percentage, which turned out to be $r = -0.37334474181581684$, which shows a moderate negative correlation between the two variables, meaning that my hypothesis that the higher the disability percentage, the lower the acceptance rate to HSPHS is possibly correct.



- 6) In order to see whether the availability of material resources impacts objective measures of achievement or admission to HSPHS, I chose to find the correlation between acceptance rate

and per pupil spending. This yields a correlation coefficient of $r = -0.39581588119318084$, which is a moderate correlation, suggesting that in some way, per pupil spending affects acceptance rate to HSPHS.



- 7) To find the proportion of schools that make up 90% of admissions to HSPHS, I first sorted the data from highest to lowest number of admissions to HSPHS. From there, I found the total number of admitted students to HSPHS, and found the threshold number that would be 90% of admitted students. I then iterated through the data, adding each admitted value onto a sum and adding 1 to the number of schools with a flag at the threshold. Once it hits the threshold the loop breaks and produces the number of schools it took to reach that threshold amount. With that number, I found the proportion by dividing it by the total number of schools. 0.765993265993266 or 76.5993265993266% of NYC middle schools make up 90% of the admissions to HSPHS.
- 8) In order to do a multiple regression, I decided to use a numpy array and matrices. Although my code did not yield any valid results, the steps I took where to find the linear regression of each possible dependent variable to each of the two independent variables, acceptances and school performance (a combination of student achievement, reading scores, and math scores), then did a regression model with all the factors combined. The code I attempted is below in the appendix.
- 9) The school characteristics that appear most relevant in acceptances to HSPHS are applications and application rate and student perception of their own school. This is due to the higher correlation between these variables and the acceptances to HSPHS. I expected to find a positive correlation between acceptance rate and per pupil spending, but instead there is a moderate negative correlation, suggesting that as per pupil spending goes up, acceptance rate goes down. Perhaps, more spending means a richer and better performing school district, so even the high-achievers at these schools would choose to continue their education with their own highly-funded public high school instead of HSPHS.
- 10) I would recommend three things to the New York City Department of Education:
 - a) Spread awareness about HSPHS to get application numbers to be higher. The largest correlation is between applications and acceptances, so if there are simply more applications from a school, there should be higher acceptances.
 - b) Since the correlation between student perception of their school and school performance is pretty moderate, I would suggest focusing on school pride for each of these schools. Increasing community feel especially would help students feel more comfortable in their learning environment and would result in better performance.
 - c) Finally, in order to boost both acceptances and school performance, I would suggest hiring more teachers in schools with large class sizes to lower that number. That way, students would receive more 1-on-1 attention, and school performance overall should increase.

11)

Appendix (Code)

Question 1:

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Mon Dec 21 19:19:22 2020
5
6  @author: oviya
7  """
8
9  import numpy as np
10 import pandas as pd
11
12 import matplotlib.pyplot as plt
13
14
15 df = pd.read_csv("middleSchoolData.csv")
16 data = np.genfromtxt(df, delimiter = ',')
17
18
19 # Q1: Correlation b/w applications and admissions
20
21 applications = df["applications"]
22 admissions = df["acceptances"]
23 app_ad_corr = applications.corr(admissions)
24 corr = str(app_ad_corr)
25
26 print("Question 1: The correlation between the number of applications and admissions is " + corr)
27
28 # Plot
29 f1 = plt.figure(1)
30 x = df.applications
31 y = df.acceptances
32 plt.scatter(x,y)
33 plt.xlabel('Number of applications')
34 plt.ylabel('Number of admissions')
35
36
```

Question 2:

```
37
38 # Q2: Which is better predictor for admission to HSPS: raw number of applications or application rate
39
40 df['application_rate'] = df['applications']/df['school_size']
41 application_rate = df['application_rate']
42 apprate_ad_corr = application_rate.corr(admissions)
43 corr2 = str(apprate_ad_corr)
44 print("\nQuestion 2: The correlation between the application rate and admissions is " + corr2)
45
46 ## Compare correlations
47 if app_ad_corr > apprate_ad_corr:
48     print("Raw number of applications is a better predictor for admissions.")
49 else:
50     print ("Application rate is a better predictor for admissions.")
51
52
53 # Plot
54 f2 = plt.figure(2)
55 a = df.application_rate
56 b = df.acceptances
57 plt.scatter(a,b)
58 plt.xlabel('Application rate')
59 plt.ylabel('Number of admissions')
60
61
```

Question 3:

```
60
61
62 # Q3: Best per-student odds of sending someone to HSPHS
63
64 df['per_student_odds'] = df['acceptances']/df['school_size']
65 per_student_odds = df['per_student_odds']
66 max_value = per_student_odds.max()
67 max_val_str = str(max_value)
68 print("\nQuestion 3: The highest per-student odds of sending someone to HSPHS is " + max_val_str)
69
70
71 name_odds = df[['school_name', 'per_student_odds']]
72 index = df[df['per_student_odds'] == max_value].index.values
73 name = df.loc[index, 'school_name']
74 print(name)
75 print("The school with the highest per-student odds of sending a student to HSPHS is shown above.")
76
77
```

Question 4:

```
77
78
79 # Q4: Relationship b/w students' perception of school and school performance
80
81 df['student_perception'] = df['rigorous_instruction'] + df['collaborative_teachers'] + \
82 df['supportive_environment'] + df['effective_school_leadership'] + df['strong_family_community_ties'] + df['trust']
83 student_perception = df['student_perception']
84 df['school_performance'] = df['student_achievement'] + df['reading_scores_exceed'] + df['math_scores_exceed']
85 school_performance = df['school_performance']
86 perc_perf_corr = student_perception.corr(school_performance)
87 corr3 = str(perc_perf_corr)
88 print("\nQuestion 4: The correlation between student perception and school performance is " + corr3)
89
90 abs_corr = abs(perc_perf_corr)
91
92 if abs_corr == 1:
93     print("There is a perfect correlation between student perception and school performance.")
94 elif abs_corr >= 0.5:
95     print("There is a high/strong correlation between student perception and school performance.")
96 elif abs_corr >= 0.3:
97     print("There is a moderate correlation between student perception and school performance.")
98 elif abs_corr >= 0:
99     print("There is a low/weak correlation between student perception and school performance.")
100 else:
101     print("Error: Correlation is out of range.")
102
103
104 # Plot
105 f3 = plt.figure(3)
106 j = df.student_perception
107 k = df.school_performance
108 plt.scatter(j,k)
109 plt.xlabel('Student Perception')
110 plt.ylabel('School Performance')
111
```

Question 5:

```
113
114 # Q5: Testing a custom hypothesis: The higher the disability percent, the lower the acceptance rate
115
116 df['acceptance_rate'] = df['acceptances']/df['applications']
117 acceptance_rate = df['acceptance_rate']
118 disability = df['disability_percent']
119 accrate_dis_corr = acceptance_rate.corr(disability)
120 corr4 = str(accrate_dis_corr)
121 print("\nQuestion 5: The correlation between acceptance rate and disability percentage is " + corr4)
122
123 abs_corr4 = abs(accrate_dis_corr)
124
125 if abs_corr4 == 1:
126     print("There is a perfect correlation between acceptance rate and disability percentage.")
127 elif abs_corr4 >= 0.5:
128     print("There is a high/strong correlation between acceptance rate and disability percentage.")
129 elif abs_corr4 >= 0.3:
130     print("There is a moderate correlation between acceptance rate and disability percentage.")
131 elif abs_corr4 >= 0:
132     print("There is a low/weak correlation between acceptance rate and disability percentage.")
133 else:
134     print("Error: Correlation is out of range.")
135
136 # Plot
137 f4 = plt.figure(4)
138 q = df.disability_percent
139 r = df.acceptance_rate
140 plt.scatter(q,r)
141 plt.xlabel('Disability Percentage')
142 plt.ylabel('Acceptance Rate')
143
144
```

Question 6:

```

145
146 # Q6: Finding evidence that the availability of resources (per student spending) impacts achievement (admissions in this case)
147
148 df['acceptance_rate'] = df['acceptances']/df['applications']
149 acceptance_rate = df['acceptance_rate']
150 spending = df['per_pupil_spending']
151 spending_accrate_corr = acceptance_rate.corr(spending)
152 corr5 = str(spending_accrate_corr)
153 print("\nQuestion 6: The correlation between acceptance rate and per-student spending is " + corr5)
154
155 abs_corr5 = abs(spending_accrate_corr)
156
157 if abs_corr5 == 1:
158     print("There is a perfect correlation between acceptance rate and disability percentage.")
159 elif abs_corr5 >= 0.5:
160     print("There is a high/strong correlation between acceptance rate and disability percentage.")
161 elif abs_corr5 >= 0.3:
162     print("There is a moderate correlation between acceptance rate and disability percentage.")
163 elif abs_corr5 >= 0:
164     print("There is a low/weak correlation between acceptance rate and disability percentage.")
165 else:
166     print("Error: Correlation is out of range.")
167
168 # Plot
169 f5 = plt.figure(5)
170 a = df.per_pupil_spending
171 b = df.acceptance_rate
172 plt.scatter(a,b)
173 plt.xlabel('Per Pupil Spending')
174 plt.ylabel('Acceptance Rate')
175

```

Question 7:

```

178 # Q7: What proportion of schools accounts for 90% of all students accepted to HSPHS?
179
180 # Sorted acceptances list in descending order
181 accept = df.sort_values('acceptances', ascending=False)['acceptances']
182 # Get number of acceptances that make up 90%
183 total_acceptances = df['acceptances']
184 threshold = total_acceptances.sum() * 0.9
185 total_schools = len(accept)
186 num_schools = 0
187 num_acceptances = 0
188
189 for i in range(len(accept)):
190     if num_acceptances < threshold:
191         num_acceptances += accept[i]
192         num_schools += 1
193     else:
194         break
195 prop = num_schools/ total_schools
196 perc = prop * 100
197
198 prop_str = str(prop)
199 perc_str = str(perc)
200
201 print("\nQuestion 7: " + prop_str + " or " + perc_str + "% of all NYC schools make up 90% of admissions to HSPHS.")
202
203

```

Question 8:

```

204
205 ### Q8: Multiple regression model
206
207 from sklearn import linear_model
208
209 X_df = df[['applications', 'acceptances', 'per_pupil_spending', 'avg_class_size', 'asian_percent', 'black_percent', \
210           'hispanic_percent', 'multiple_percent', 'white_percent', 'rigorous_instruction', 'collaborative_teachers', \
211           'supportive_environment', 'effective_school_leadership', 'strong_family_community_ties', 'trust', \
212           'disability_percent', 'poverty_percent', 'ESL_percent', 'school_size', 'student_achievement', \
213           'reading_scores_exceed', 'math_scores_exceed']]
214 Y1_df = df[['acceptances']]
215 Y2_df = df[['school_performance']]
216
217
218 # Descriptives
219 d1 = np.mean(data,axis=0)
220 d2 = np.median(data,axis=0)
221 d3 = np.std(data,axis=0)
222 d4 = np.corrcoef(data[:,0],data[:,1])
223
224
225 # For acceptances regression
226 X1 = np.transpose([data[:,2],data[:,4],data[:,5], data[:,6], data[:,6], data[:,7], data[:,8], data[:,9], \
227                  data[:,10], data[:,11], data[:,12], data[:,13], data[:,14], data[:,15], data[:,16], \
228                  data[:,17], data[:,18], data[:,19], data[:,20], data[:,21], data[:,22], data[:,23]])
229 Y1 = data[:,3]
230 regr = linear_model.LinearRegression()
231 regr.fit(X1,Y1)
232 r_sqr = regr.score(X1,Y1)
233 betas = regr.coef_
234 y_int = regr.intercept_
235 # Plot
236 y_hat = betas[0]*data[:,0] + betas[1]*data[:,1] + betas[2]*data[:,2] + y_int
237 plt.plot(y_hat,data[:,3],'o',markersize=.75)
238 plt.xlabel('Prediction')
239 plt.ylabel('Actual acceptances')
240 plt.title('R^2: {:.3f}'.format(r_sqr))
241
242
243 # For achievement regression
244 X2 = np.transpose([data[:,2], data[:,3], data[:,4],data[:,5], data[:,6], data[:,6], data[:,7], data[:,8], \
245                  data[:,9], data[:,10], data[:,11], data[:,12], data[:,13], data[:,14], data[:,15], data[:,16], \
246                  data[:,17], data[:,18], data[:,19], data[:,20]])
247 Y2 = np.transpose([data[:,21], data[:,22], data[:,23]])
248 regr = linear_model.LinearRegression()
249 regr.fit(X2,Y2)
250 r_sqr = regr.score(X2,Y2)
251 betas = regr.coef_
252 y_int = regr.intercept_
253 # Plot
254 y_hat = betas[0]*data[:,0] + betas[1]*data[:,1] + betas[2]*data[:,2] + y_int
255 plt.plot(y_hat,data[:,21],'o',markersize=.75)
256 plt.xlabel('Prediction')
257 plt.ylabel('Actual school performance')
258 plt.title('R^2: {:.3f}'.format(r_sqr))
259

```

Output:

```

In [155]: runfile('/Users/oviya/Desktop/Data Science/OviyaAdhan-FinalProject.py', wdir='/Users/oviya/Desktop/Data
Science')
Question 1: The correlation between the number of applications and admissions is 0.8017265370719315

Question 2: The correlation between the application rate and admissions is 0.658750752900268
Raw number of applications is a better predictor for admissions.

Question 3: The highest per-student odds of sending someone to HSPHS is 0.23482245131729668
304 THE CHRISTA MCAULIFFE SCHOOL\I.S. 187
Name: school_name, dtype: object
The school with the highest per-student odds of sending a student to HSPHS is shown above.

Question 4: The correlation between student perception and school performance is 0.41874194039783613
There is a moderate correlation between student perception and school performance.

Question 5: The correlation between acceptance rate and disability percentage is -0.37334474181581684
There is a moderate correlation between acceptance rate and disability percentage.

Question 6: The correlation between acceptance rate and per-student spending is -0.39581588119318084
There is a moderate correlation between acceptance rate and disability percentage.

Question 7: 0.765993265993266 or 76.5993265993266% of all NYC schools make up 90% of admissions to HSPHS.

```