

NCAA APR Database

Exploratory Analysis

Jack Huff, jhuff@bellarmine.edu
Brett Tully, btully@bellarmine.edu

I. INTRODUCTION

Our data set includes a list of all NCAA schools and their sports teams. It includes data on what conference and division they are in, as well as the sports teams respective APR score, retention rate, eligibility, and the amount of 4 year players they have. Our data set was found on Kaggle. We chose this dataset because we were interested in academic performance and its relationship to specific conferences, sports, and teams.

II. DATA SET DESCRIPTION

Our data set has a total of 6511 rows and 57 columns. The 57 columns have various data types including int64, float64, and object.

Table 1: Data Types and Missing Data

<i>Variable Name</i>	<i>Data Type</i>	<i>Missing Data (%)</i>
V1 SCHOOL_ID	nominal/int64	0%
V2 SCHOOL_NAME	nominal/object	0%
V3 SCHOOL_TYPE	nominal/int64	0%
V4 ACADEMIC_YEAR	ordinal/int64	0%
V5 SPORT_CODE	nominal/int64	0%
V6 SPORT_NAME	nominal./object	0%
V7 NCAA_DIVISION	nominal/int64	0%
V8 NCAA_SUBDIVISION	nominal/int64	0%
V9 NCAA_CONFERENCE	nominal/object	0%
V10 FOURYEAR_ATHLETES	interval/int64	0%
V11 FOURYEAR_SCORE	interval/int64	9.6%
V12 FOURYEAR_ELIGIBILITY	interval/float64	9.6%

V13 FOURYEAR_RETENTION	interval/float64	9.6%
V14 2014_ATHLETES	interval/int64	9.6%
V15 2014_SCORE	interval/int64	9.6%
V16 2014_ELIGIBILITY	interval/float64	9.6%
V17 2014_RETENTION	interval/float64	9.6%
V18 2013_ATHLETES	interval/int64	9.6%
V19 2013_SCORE	interval/int64	9.6%
V20 2013_ELIGIBILITY	interval/float64	9.6%
V21 2013_RETENTION	interval/float64	9.6%
V22 2012_ATHLETES	interval/int64	9.6%
V23 2012_SCORE	interval/int64	9.6%
V24 2012_ELIGIBILITY	interval/float64	9.6%
V25 2012_RETENTION	interval/float64	9.6%
V26 2011_ATHLETES	interval/int64	9.6%
V27 2011_SCORE	interval/int64	9.6%
V28 2011_ELIGIBILITY	interval/float64	9.6%
V29 2011_RETENTION	interval/float64	9.6%
V30 2010_ATHLETES	interval/int64	9.6%
V31 2010_SCORE	interval/int64	9.6%
V32 2010_ELIGIBILITY	interval/float64	9.6%
V33 2010_RETENTION	interval/float64	9.6%
V34 2009_ATHLETES	interval/int64	9.6%
V35 2009_SCORE	interval/int64	9.6%
V36 2009_ELIGIBILITY	interval/float64	9.6%
V37 2009_RETENTION	interval/float64	9.6%

V38 2008_ATHLETES	interval/int64	9.6%
V39 2008_SCORE	interval/int64	9.6%
V40 2008_ELIGIBILITY	interval/float64	9.6%
V41 2008_RETENTION	interval/float64	9.6%
V42 2007_ATHLETES	interval/int64	9.6%
V43 2007_SCORE	interval/int64	9.6%
V44 2007_ELIGIBILITY	interval/float64	9.6%
V45 2007_RETENTION	interval/float64	9.6%
V46 2006_ATHLETES	interval/int64	9.6%
V47 2006_SCORE	interval/int64	9.6%
V48 2006_ELIGIBILITY	interval/float64	9.6%
V49 2006_RETENTION	interval/float64	9.6%
V50 2005_ATHLETES	interval/int64	9.6%
V51 2005_SCORE	interval/int64	9.6%
V52 2005_ELIGIBILITY	interval/float64	9.6%
V53 2005_RETENTION	interval/float64	9.6%
V54 2004_ATHLETES	interval/int64	9.6%
V55 2004_SCORE	interval/int64	9.6%
V56 2004_ELIGIBILITY	interval/float64	9.6%
V57 2004_RETENTION	interval/float64	9.6%

III. Data Set Summary Statistics

Here, we took a look at the major college sports conferences and their APR scores. There are 25 teams in each conference, and we did a simple statistical analysis on each conference and their scores.

Table 2: Summary Statistics for Pivot (name of dataset)

<i>Variable Name</i>	<i>Count</i>	<i>Mean</i>	<i>Standard Deviation</i>	<i>Min</i>	<i>25th</i>	<i>50th</i>	<i>75th</i>	<i>Max</i>
Atlantic Coast Conference	25	980.269	8.320	960.94 5	975.74 5	981.15 1	985.16 2	992.948
Atlantic Sun Conference	25	966.136	12.370	932.63 6	966.13 6	966.13 6	975.57 5	983.909
Big 12 Conference	25	967.764	10.716	945.41 8	961.68 1	967.76 3	976.18 1	985.747
Big East Conference	25	983.851	9.106	960.07 2	978.36 3	986.58 4	989.89 3	996.090
Big Ten Conference	25	978.441	8.462	958.44 1	974.14 1	979.74 5	985.21 2	991.654
Pac-12 Conference	25	973.336	9.937	949.93 9	968.80 3	974.12 8	981.13 2	991.136
Southeastern Conference	25	971.081	11.650	947.34 4	965.83 6	971.08 1	978.37 0	997.363
The Ivy League	25	993.721	2.779	986.59 0	992.17 0	993.95 4	995.78 4	998.077

We wanted to narrow down the Conferences to the main conferences that people generally pay more attention to.

There should be a table for **EACH** categorical variable.

Table 3: Proportions for Sports (n=yyy)

<i>Category: Sports Teams</i>	<i>Frequency</i>	<i>Proportion (%)</i>
<i>Men's Basketball</i>	<i>322</i>	<i>5.91%</i>
<i>Women's Basketball</i>	<i>321</i>	<i>5.89%</i>
<i>Women's Cross Country</i>	<i>273</i>	<i>5.01%</i>
<i>Women's Volleyball</i>	<i>303</i>	<i>5.56%</i>
<i>Women's Track, Outdoor</i>	<i>282</i>	<i>5.17%</i>
<i>Women's Soccer</i>	<i>290</i>	<i>5.32%</i>

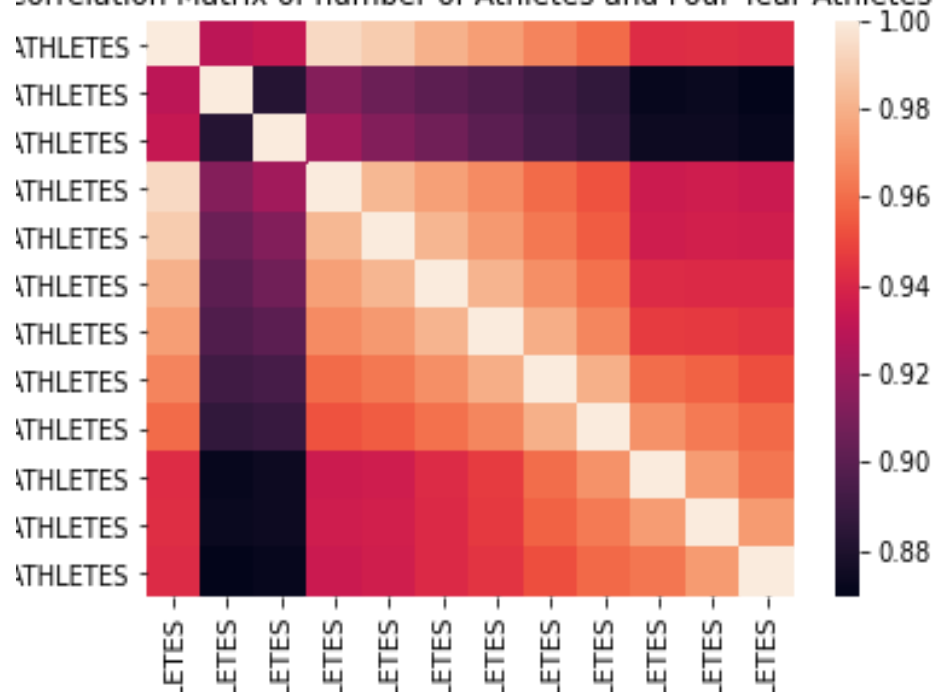
<i>Women's Tennis</i>	260	4.77%
<i>Women's Track, Indoor</i>	274	5.03%
<i>Men's Cross Country</i>	209	3.83%
<i>Men's Golf</i>	189	3.47%
<i>Baseball</i>	270	4.96%
<i>Women's Softball</i>	257	4.72%
<i>Men's Track, Outdoor</i>	232	4.26%
<i>Women's Golf</i>	189	3.47%
<i>Men's Tennis</i>	211	3.87%
<i>Men's Track, Indoor</i>	212	3.89%
<i>Football</i>	228	4.19%
<i>Men's Soccer</i>	183	3.36%
<i>Women's Swimming</i>	165	3.03%
<i>Men's Swimming</i>	114	2.09%
<i>Women's Lacrosse</i>	71	1.30%
<i>Women's Rowing</i>	65	1.19%
<i>Women's Field Hockey</i>	71	1.30%
<i>Men's Wrestling</i>	67	1.23%

<i>Men's Lacrosse</i>	<i>49</i>	<i>.89%</i>
<i>Women's Gymnastics</i>	<i>60</i>	<i>1.10%</i>
<i>Men's Ice Hockey</i>	<i>56</i>	<i>1.02%</i>
<i>Women's Ice Hockey</i>	<i>25</i>	<i>.46%</i>
<i>Women's Water Polo</i>	<i>25</i>	<i>.46%</i>
<i>Women's Bowling</i>	<i>16</i>	<i>.29%</i>
<i>Women's Fencing</i>	<i>11</i>	<i>.20%</i>
<i>Men's Water Polo</i>	<i>17</i>	<i>.31%</i>
<i>Men's Volleyball</i>	<i>17</i>	<i>.31%</i>
<i>Mixed Rifle</i>	<i>15</i>	<i>.27%</i>
<i>Men's Fencing</i>	<i>12</i>	<i>.22%</i>
<i>Men's Gymnastics</i>	<i>14</i>	<i>.26%</i>
<i>Women's Skiing</i>	<i>10</i>	<i>.18%</i>
<i>Men's Skiing</i>	<i>9</i>	<i>.16%</i>

Table 4: Correlation Table/Tables

	FOURYE R_ATHLE TES	2014_A THLET ES	2013_A THLET ES	2012_A THLET ES	2011_A THLET ES	2010_A THLET ES	2009_A THLET ES	2008_A THLET ES	2007_A THLET ES	2006_A THLET ES	2005_A THLET ES	2004_A THLET ES
FOURYE R_ATHLE TES	1.000000	0.92985 2	0.93280 2	0.99343 5	0.98918 6	0.98060 9	0.97437 1	0.96669 1	0.95954 8	0.94218 7	0.94290 4	0.94172 4
2014_ATH LETES	0.929852	1.00000 0	0.88250 5	0.91326 8	0.90568 3	0.90102 2	0.89674 3	0.89150 5	0.88635 0	0.87159 3	0.87307 7	0.86952 0
2013_ATH LETES	0.932802	0.88250 5	1.00000 0	0.92174 3	0.91262 3	0.90701 7	0.90051 5	0.89395 7	0.88849 6	0.87384 1	0.87361 0	0.87203 2
2012_ATH LETES	0.993435	0.91326 8	0.92174 3	1.00000 0	0.98226 5	0.97453 7	0.96843 3	0.95950 6	0.95269 9	0.93517 6	0.93599 6	0.93449 7
2011_ATH LETES	0.989186	0.90568 3	0.91262 3	0.98226 5	1.00000 0	0.98168 5	0.97290 8	0.96322 6	0.95530 5	0.93593 8	0.93726 5	0.93633 3
2010_ATH LETES	0.980609	0.90102 2	0.90701 7	0.97453 7	0.98168 5	1.00000 0	0.98140 5	0.96944 0	0.96076 2	0.94188 1	0.94114 0	0.94131 0
2009_ATH LETES	0.974371	0.89674 3	0.90051 5	0.96843 3	0.97290 8	0.98140 5	1.00000 0	0.97912 8	0.96723 9	0.94660 6	0.94608 4	0.94449 1
2008_ATH LETES	0.966691	0.89150 5	0.89395 7	0.95950 6	0.96322 6	0.96944 0	0.97912 8	1.00000 0	0.98006 7	0.95976 4	0.95682 2	0.95177 0
2007_ATH LETES	0.959548	0.88635 0	0.88849 6	0.95269 9	0.95530 5	0.96076 2	0.96723 9	0.98006 7	1.00000 0	0.97059 3	0.96367 2	0.95914 9
2006_ATH LETES	0.942187	0.87159 3	0.87384 1	0.93517 6	0.93593 8	0.94188 1	0.94660 6	0.95976 4	0.97059 3	1.00000 0	0.97395 6	0.96270 7
2005_ATH LETES	0.942904	0.87307 7	0.87361 0	0.93599 6	0.93726 5	0.94114 0	0.94608 4	0.95682 2	0.96367 2	0.97395 6	1.00000 0	0.97315 9
2004_ATH LETES	0.941724	0.86952 0	0.87203 2	0.93449 7	0.93633 3	0.94131 0	0.94449 1	0.95177 0	0.95914 9	0.96270 7	0.97315 9	1.00000 0

Correlation Matrix of number of Athletes and Four Year Athletes



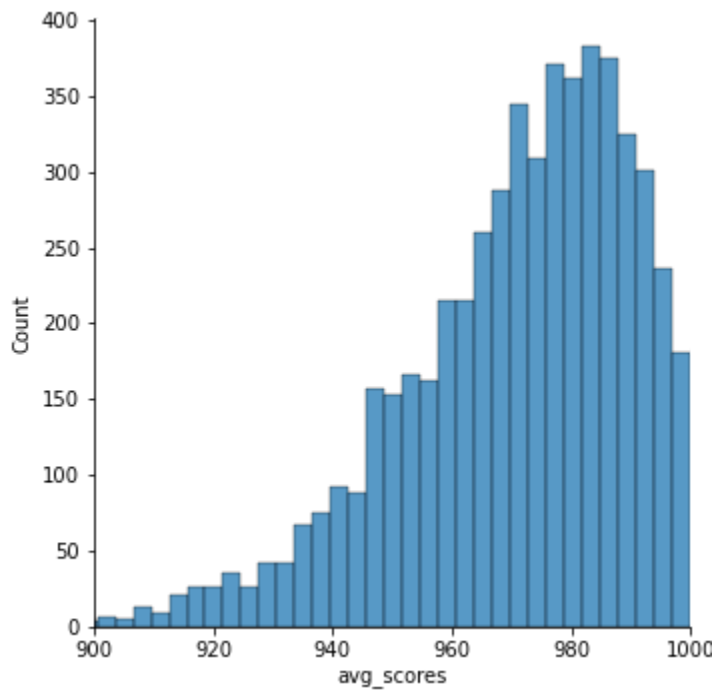
IV. DATA SET GRAPHICAL EXPLORATION

After we narrowed down our data set for what we were interested in looking at, we decided to visualize the data. We graphed distributions, scatterplots, barcharts, and heatmaps. The visualization of the data helped us see some interesting patterns.

Titles of some of the charts are cut-off/not showing in the document

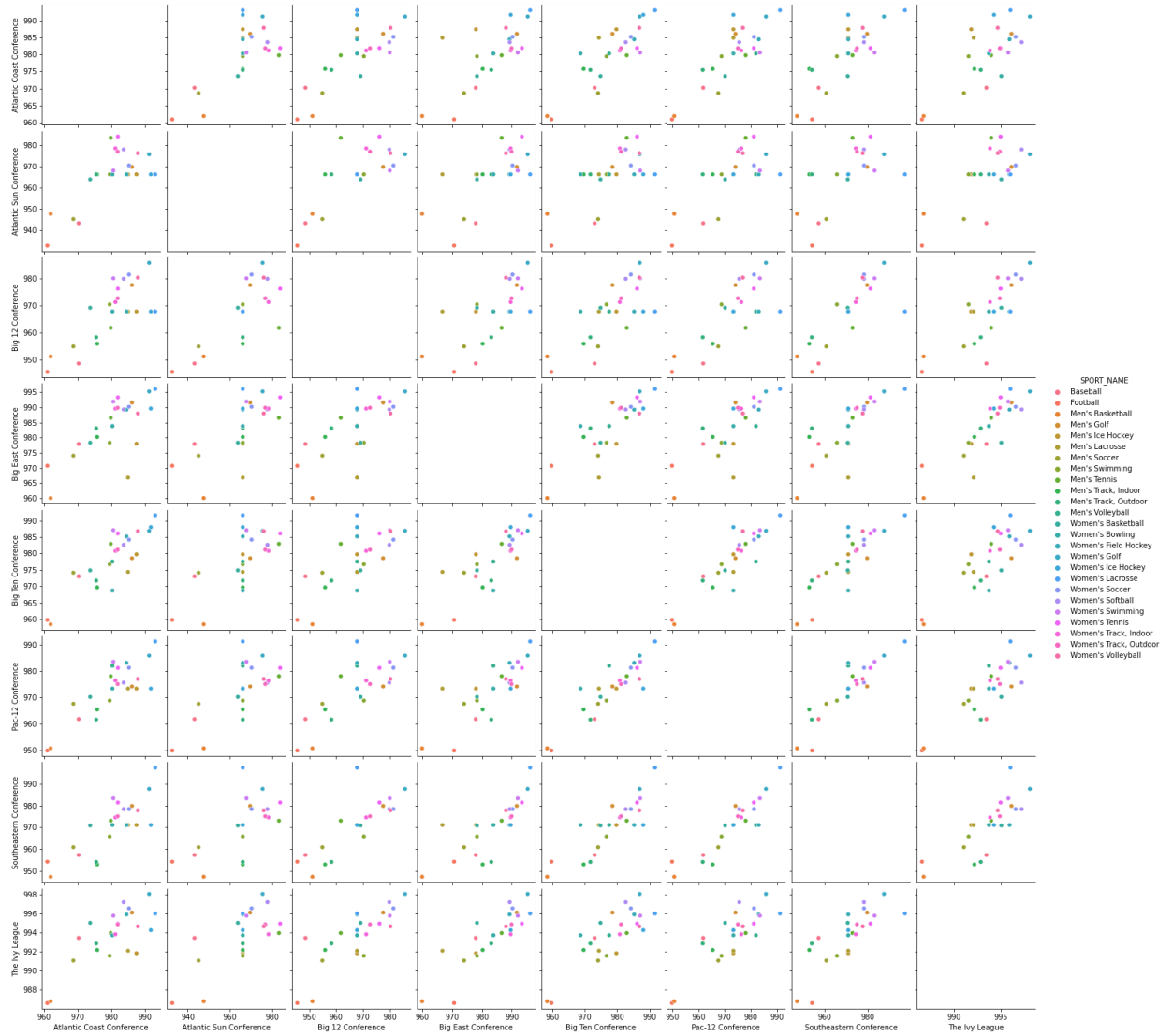
A. Distributions

Figure 1: Histogram of Distribution from Dataset



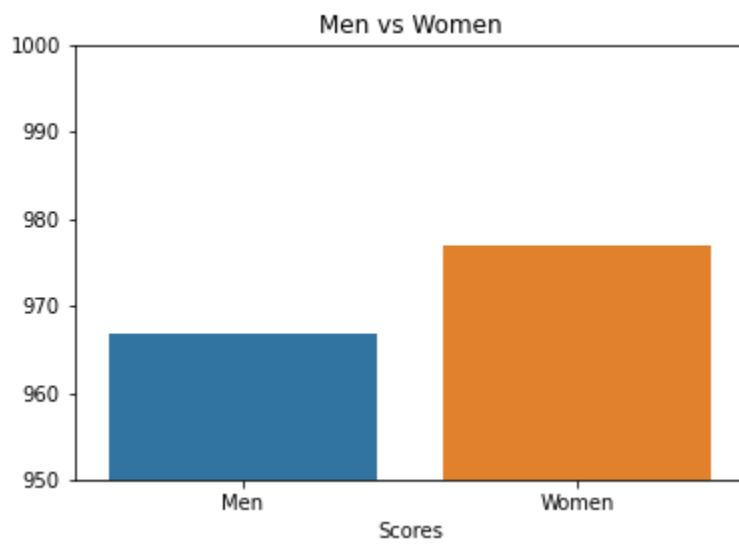
B. ScatterPlots / Pairwise Plots (continuous variables)

Figure 2: Pairwise Plots for Sports and comparing each Conference to each other



C. *Barcharts (categorical variables)*

Figure 3: Comparison between Men and Women scores



D. *Other Plots*

Figure 4: Scatter Plot of 7 specific Conferences

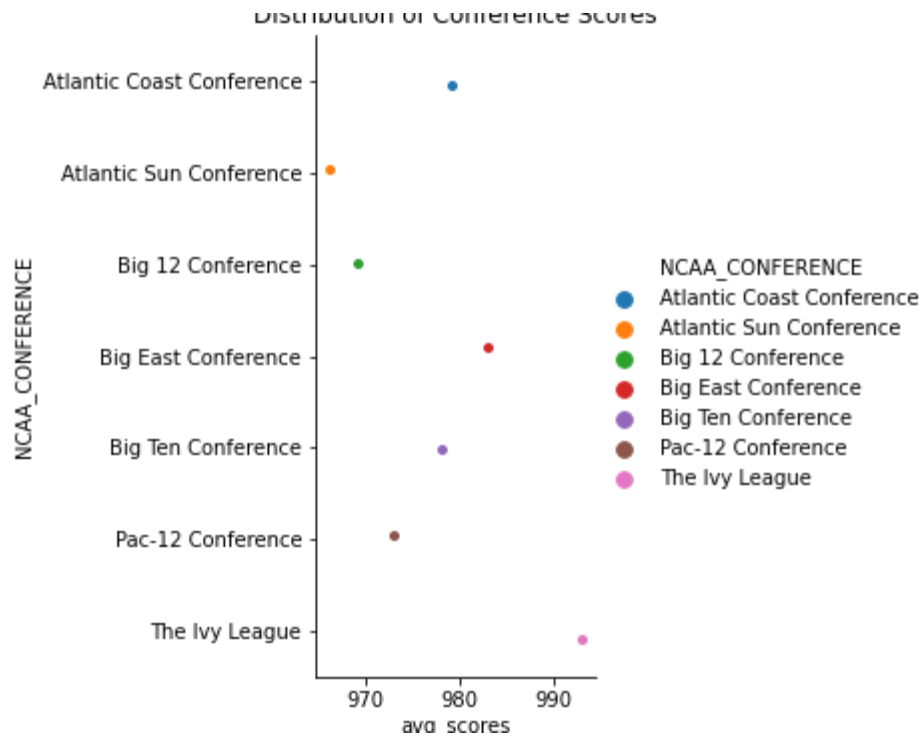


Figure 5: Heatmap of Men's sports from every Conference in NCAA

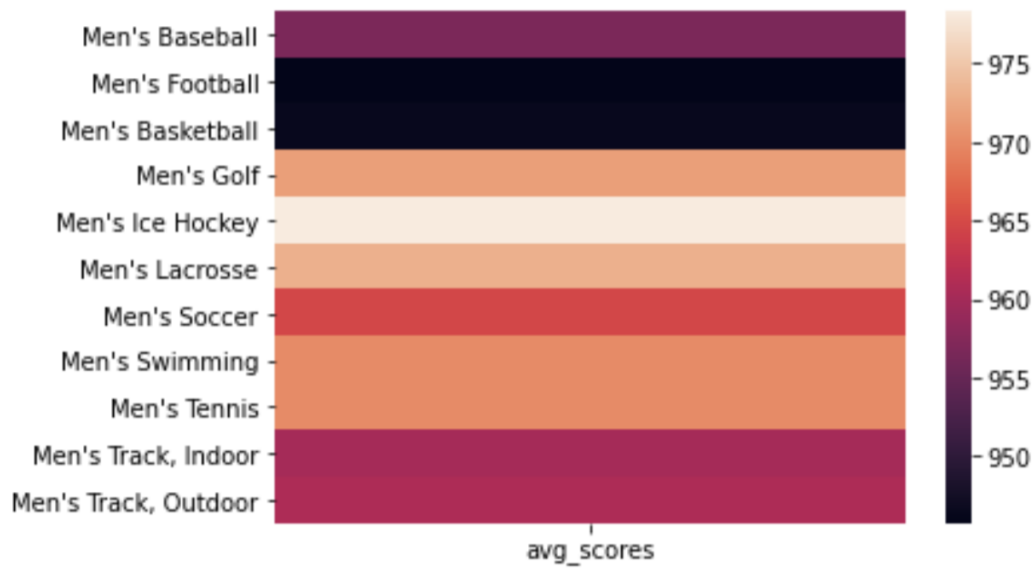
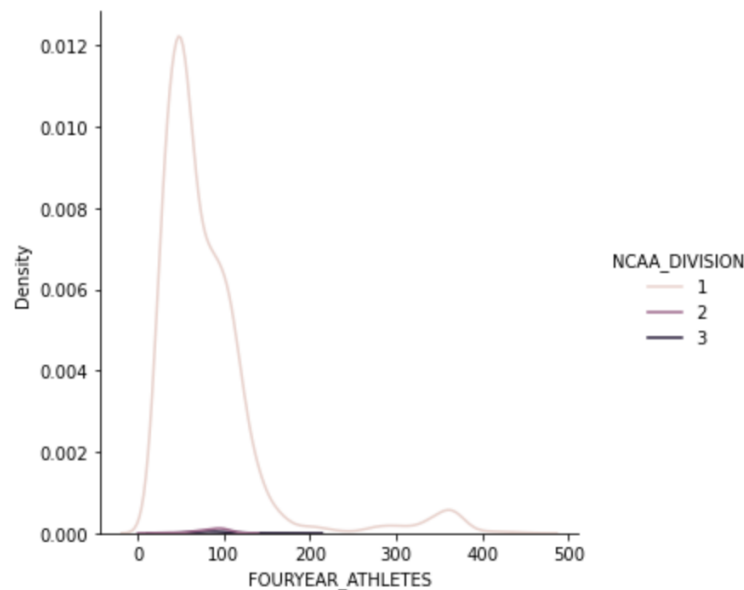


Figure 7: KDE of the amount of Four Year Athletes in all of NCAA



V. SUMMARY OF FINDINGS

In **Figure 1**, we graphed the distribution of the average APR score for each sports team. It seems to have a pretty normal distribution. A perfect APR score is 1000. **Figure 2**, the pairwise plot, shows how the APR score in each sport in each Conference pairs up to the scores of the same sports in other Conferences. **Figure 3**, the bar chart shows the average APR scores for men's versus women's sports. It appears that on average, women's sports teams tend to outperform men's sports teams when it comes to academic performance. **Figure 4** shows the average APR score for the most talked about Conferences, with Ivy League added in to view its academic superiority and Atlantic Sun Conference for obvious reasons. Unsurprisingly, the Ivy League has the highest average APR score, and unfortunately, the A-Sun conference has the lowest. **Figure 5**; This heat map shows the average APR scores for the main mens sports throughout the entire NCAA. Basketball and Football seem to have the lowest average, and Golf, Hockey, and Lacrosse appear to have the highest scores. This is a good way to quickly compare sports averages against each other. **Figure 6**; This heat map is similar to the last one, but this one shows each sports team's average APR score within its respective conference. This would be useful if one wanted to compare how a particular sport in a particular conference compares to another sport in a different conference. **Figure 7**; This kernel density chart shows the density of four year athletes in division 1, 2, and 3. It is clear that division one sports have a much higher density of four year athletes when compared to division 2 and 3 sports. This most likely is caused by scholarship options that Division 1 players receive.

Working with this data set was very interesting. Coming into it, we both had a few ideas in mind about what sports would have good/bad academic performance, and it was cool to see if we were correct. We found that men's basketball tended to have the worst academic performance across all conferences, while men's ice hockey had the best academic performance. We also found that, generally speaking, women's sports teams outperform men's sports teams academically. To no surprise, the Ivy league sports teams substantially outperformed all other conferences, and the Atlantic sun conference had the worst academic performance. When it comes to 4 year athletes, Division 1 schools had a much larger amount than Division 2 and 3 schools. We assume this has to do with the

scholarship opportunities Division 1 players have that Division 2 and 3 players don't. Overall, working with this data set was interesting because it did reveal that there are some relationships between specific sports, conferences, and gender when it comes to academic performance within the NCAA.