# Advanced methods in applied statistics: Problem set 1

Tomás Fernández Bouvier

February 17, 2021



Total: 9.0 points

Figure 1: Histogram representation of AdjD for each conference

2 +3 pts



Figure 2: Scatter plot of the imporvement of each team from 2009 to 2014

| Conference | means | medians |
|------------|-------|---------|
| ACC | -0.625 | -0.05 |
| SEC | 2.283 | 1.65 |
| B10 | 2.673 | 4.30 |
| BSky | 2.322 | 3.20 |
| A10 | 4.336 | 4.90 |
| other | 2.595 | 1.90 |

Table 1: Table of medians of the difference between 2014 and 2009 in AdjO

Figure 3: Histogram representation of AdjD for each conference including BE
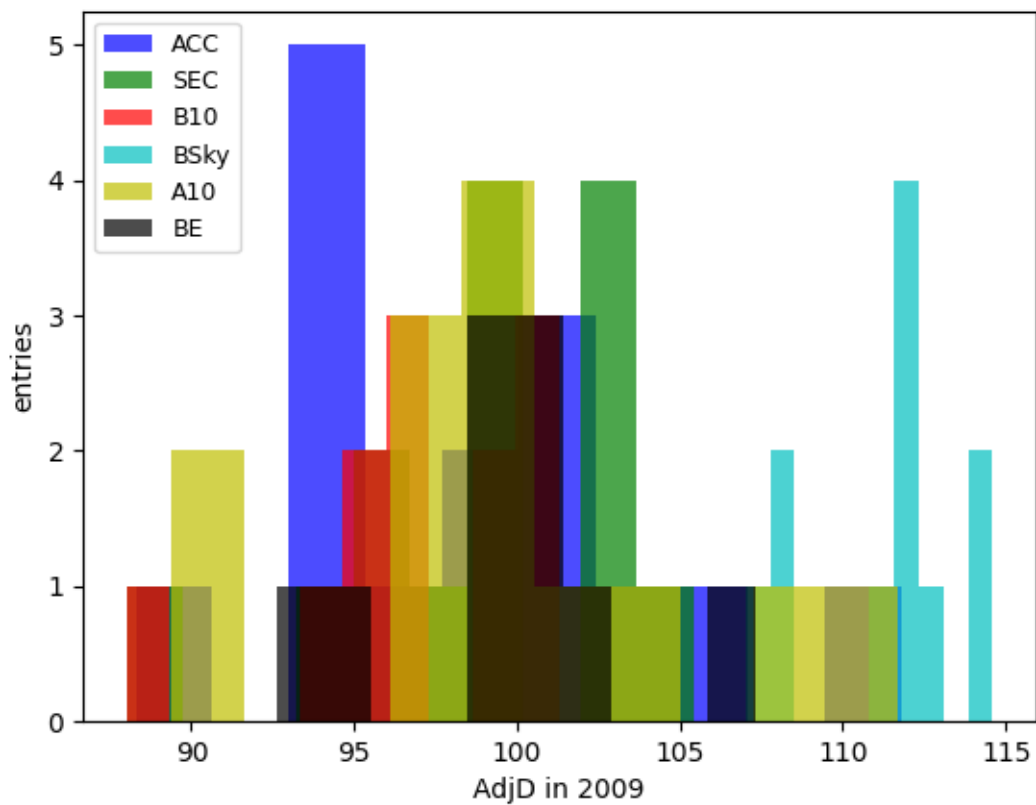
Figure 4: Scatter plot of the imporvement of each team from 2009 to 2014 including BE

| Conference | means | medians |
|------------|-------|---------|
| ACC | -0.625 | -0.05 |
| SEC | 2.283 | 1.65 |
| B10 | 2.673 | 4.30 |
| BSky | 2.322 | 3.20 |
| A10 | 4.336 | 4.90 |
| BE | 1.143 | 1.90 |
| other | 2.624 | 1.85 |

Table 2: Table of medians of the difference between 2014 and 2009 in AdjO

# 4   (EXTRA) +0.5 pts

After parsing and depuring the data in the pdf document we were able to obtain a list of names of authors. The number of unique authors in this list is 3522. After ordering the list aphabetically we also obtained that the mid point is situated around the author "K. DE"

I looked through the code — nice solution! Good thinking with choosing the "set" for uniqueness, and with removing the collaboration names using re < > expressions — please report this in text next time as well. Based on your result, I am guessing that you sorted the list by the first name+last name combination rather than just the last name. It would be great if you included more details about your solutions in the report! (e.g. how you ensured uniqueness, treated the sorting of special characters, whether you sorted in ascending/descending order [I see you did reverse = True; why?], etc so it's easier to benchmark. Good job otherwise!)