

Research into factors affecting the performance of tennis players in ATP

You Peng

2022-02-23

Introduction

Tennis is a popular sport around the world, and those professional tennis players who maintained a high ranking in ATP seem to have something in common. There are already research studying the Association between body height and serve speed in elite tennis players. For example, body height of the men explained 27% of the variance of fastest serve in a match (Vaverka & Cernosek, 2013). Serving speed may not directly contribute to a tennis player's ranking in tournaments, and there could be more potential factors affecting the performance of tennis players. Therefore, we raise our research question as how could some physical factors and experience of a tennis player be used to explain the variation observed in rankings of 639 players participated in 2017 ATP tournaments? By exploring the relationship between tennis players' physical characteristics, experience, and their performance, we could help coaches to better recognize if someone could be a good tennis player and inspire tennis players what they can improve on in order to making progress. The goal of this research is to find a model that is not overly complicated, but also having reasonable properties required to make good predictions.

Methods

The dataset we used is the ATP World Tour tennis data. "This dataset contains tennis data from the ATP World Tour website. The data contains ATP tournaments, rankings and player's overview. The latest available data is for 2017." (ATP World Tour, 2018) The dataset can be downloaded from the url: 'https://datahub.io/sports-data/atp-world-tour-tennis-data/datapackage.json' or website 'https://datahub.io/sports-data/atp-world-tour-tennis-data/#r', we only require the 9th and 10th csv file from this json file. We provided two methods to read in the dataset. The first one gives the url path provided by the json file to read.csv(), which will takes a long time to download from the website. The second one requires the downloaded csv file stored in the same path as the rmd file, and then apply read.csv() directly on the path 'rankings_1973-2017_csv.csv' and 'player_overviews_unindexed_csv.csv'.

The outcome variable and predictors were contained in two different datasets we loaded. The first dataset contains the ranking of tennis players every week from 1973 to 2017, which is a lot of data. In order to perform the most up-to-date analysis as well as to obtain a continuous outcome variable, I averaged the ranking of tennis players over the whole year 2017, so that each tennis player has one corresponding average ranking. Similarly, I averaged the age and number of tournaments played in 2017 as well. This process is done by using filter() from dplyr to leave only those observations about 2017 in the dataset, and aggregate() to calculate the mean of rank, age, and tournaments played to be new variables in the dataset. Then we used merge() from dplyr by player_id to merge two dataset into one, and only kept the following variables after applying select() on the joint dataset:

"weight_kg": body weight in kg of the tennis player in 2017. "height_cm": body height in cm of the tennis player in 2017. "handedness": whether the tennis player is left-handed or right-handed. "average_age": the average age of the tennis player during the 2017 tournament period. For example, some players were

24 years old during the first few tournaments held in 2017 and grew to 25 years old later in the 2017. In those cases, their average age would become decimals. “average_tourneys_played”: the average number of tournaments played during the 2017 period. “backhand”: whether the player is using one-hand backhand or two-hand backhand. “average_rank”: the average rank a tennis player achieved during the whole 2017 tournament period. This is a continuous variable since the rank is averaged over the whole period. This is our response variable.

The joint dataset now only contains predictors and outcome variables that I need, and I dropped all rows that contained null values or contained 0 height or weight since these observations were missing the predictor values we need. Note that I also dropped one observation because that tennis player had a height of 3cm. Another tennis player had a weight of 675kg, which was modified as 67.5 after the cleaning.

Then we conducted exploratory data analysis on this cleaned dataset. The exploratory data analysis provides a general idea about the distribution of the data we have as well as a note of anything odd such as skewness or outliers. In the exploratory data analysis. We first took a look at the dimension of our dataset, as well as a summary of mean and standard deviation for variables in the dataset. Moreover, we generated boxplots and histograms for each of the numeric variables using `boxplot()` and `hist()`. For categorical variables, we provided bar plots for them by using `geom_bar()` in the `ggplot2`. In order to visualize the correlation between predictors and response variable, scatter plots between each of the predictors and response were produced by `plot()` functions. Boxplots for average ranks after grouped by categorical variables were also generated by `plot()` functions in order to visualize the effect of these categorical variables on the response. We also provided a table presenting the means and standard deviations of our response variables after grouping by each of the categorical variable. Moreover, a scatterplot for average ranking vs. height after grouping by handedness and backhand, as well as a smoothing linear line, is generated by `geom_point()` and `geom_smooth()` from `ggplot2`. Lastly, we provided a max-min plot on handedness and average ranking and a max-min plot on backhand and average ranking, using `stat_summary()` in `ggplot2`.

We continued our research by checking the violations of 4 assumptions needed for fitting a multiple linear regression model. We checked these four assumptions by generating the QQ plot using `qqnorm()` and `qqline()`, and generating residual plots by using `resid()`. If any assumption is violated, we will apply `powerTransform()` from `car` package, which is an automated power transformation, to transform predictors and response simultaneously, then check if the transformed model satisfies the assumptions. We have to be careful not making the transformation too complicated.

Then we move to the stage of model comparison. We built another GAM model where we put a cubic regression spline on height using `gam()` from `mgcv` package. Then we compared the adjusted R square of these two models to see which one is better. We would prefer a model that has fairly large adjusted R² with appropriate number of predictors. A plot for linear and non-linear associations between rankings and heights of athletes is also generated by using `ggplot2`. At this point, it would just be a preliminary comparison.

Preliminary Results

After cleaning and wrangling, there are 639 observations and 8 variables in the dataset, and the response variable is the average rank a tennis player achieved during the whole 2017 tournament period. The summary statistics for five numerical variables are given below, we can see that the mean height for all players is 185.3, which is much higher than the average height of men between 20 to 39 in US, which is 176.1 in that case (Fryer, 2018). This may suggest that a higher height could take advantage in tennis.

Table 1: Summary statistics for numerical variables in the dataset

| Variable | mean (s.d.) in dataset |
|-------------------------|------------------------|
| average_rank | 555.877 (490.213) |
| weight_kg | 79.154 (6.783) |
| height_cm | 185.252 (6.631) |
| average_tourneys_played | 17.624 (8.455) |
| average_age | 25.991 (4.522) |

For numerical variables in dataset, their histograms are given in Figure 1. We found that the response variable is heavily right skewed, and the distribution of “average tournaments played” has a heavy tail on the left. These observations suggest that we may need to transform some variables later.

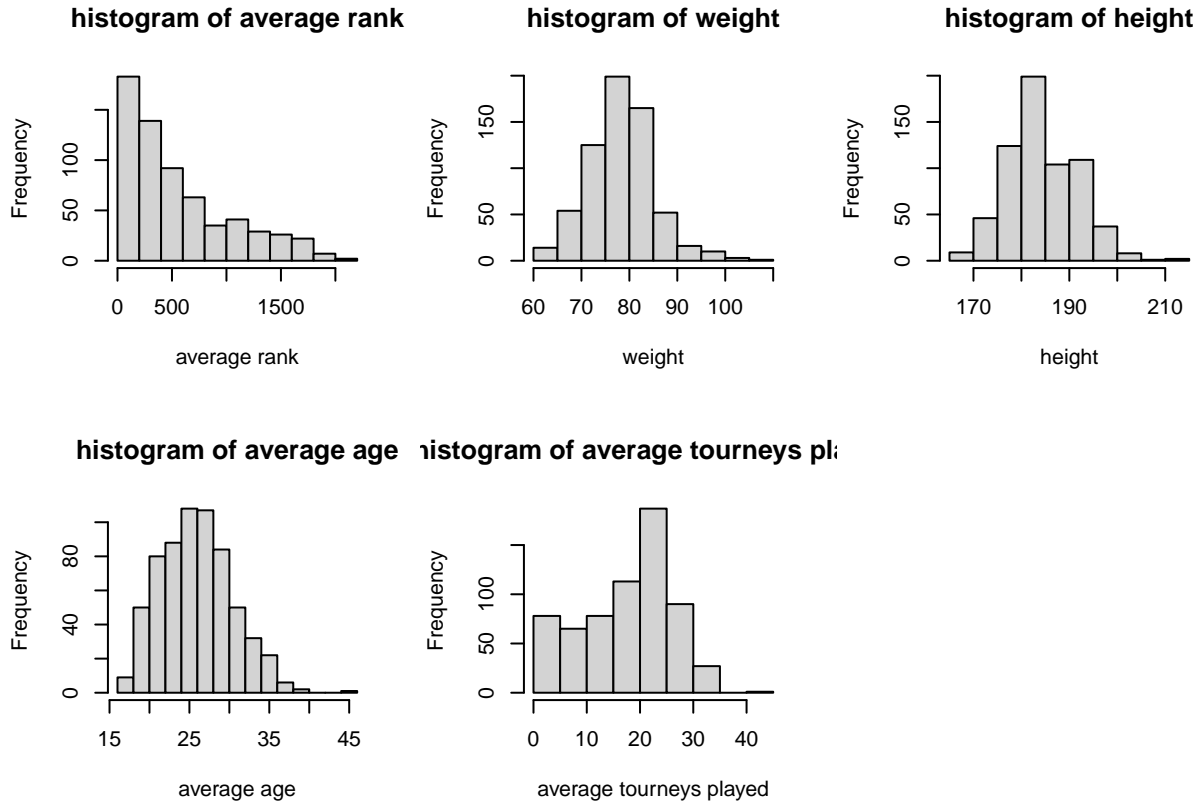
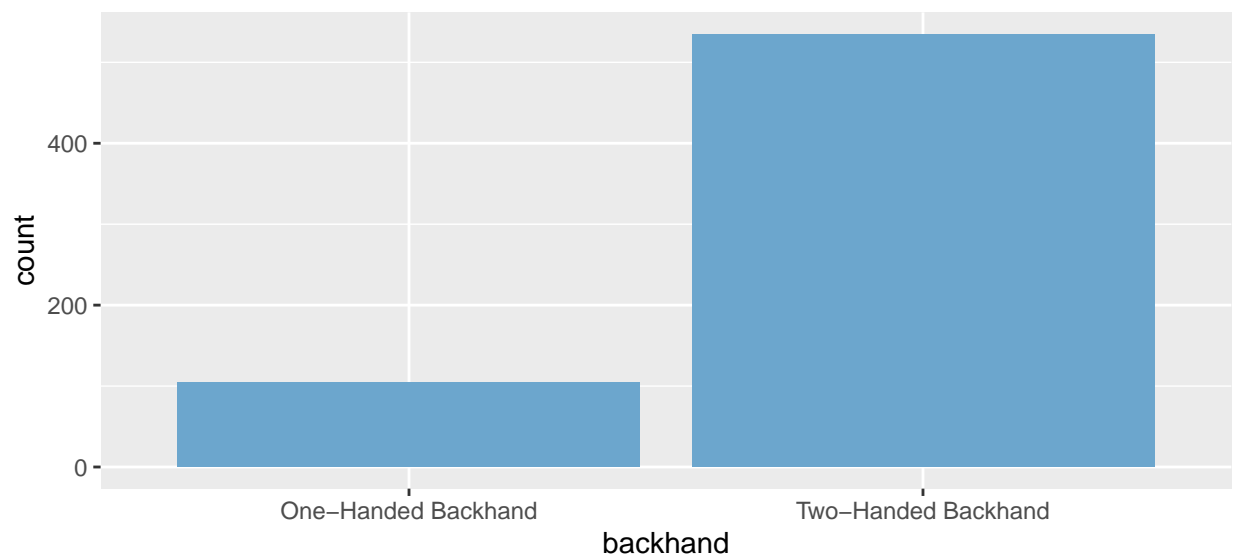
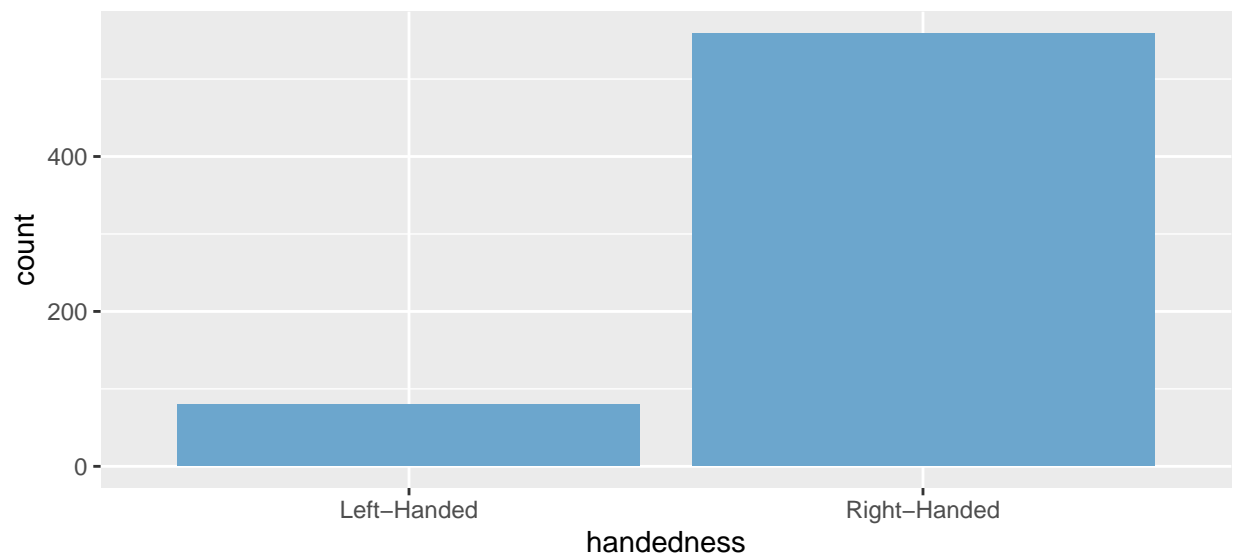


Figure 1: histogram for numerical variables

For binary variables in the dataset, their bar plots are given below. There are only 80 left-hand players where majority being right-hand players, and 104 one-handed backhand players where majority being two-handed backhand players. This indicates that data are very imbalanced in both variables.



From plots presenting below, we can see a clear linear pattern in the scatterplot between average rank and average tournaments played. Since the slope is negative, this suggest that more tournaments a player has played, the better rank he is expected to have. The correlations in other plots are not very obvious to see from the plot, while the right-handed players may perform a little bit better than left-handed players since the median of them is a little bit larger.

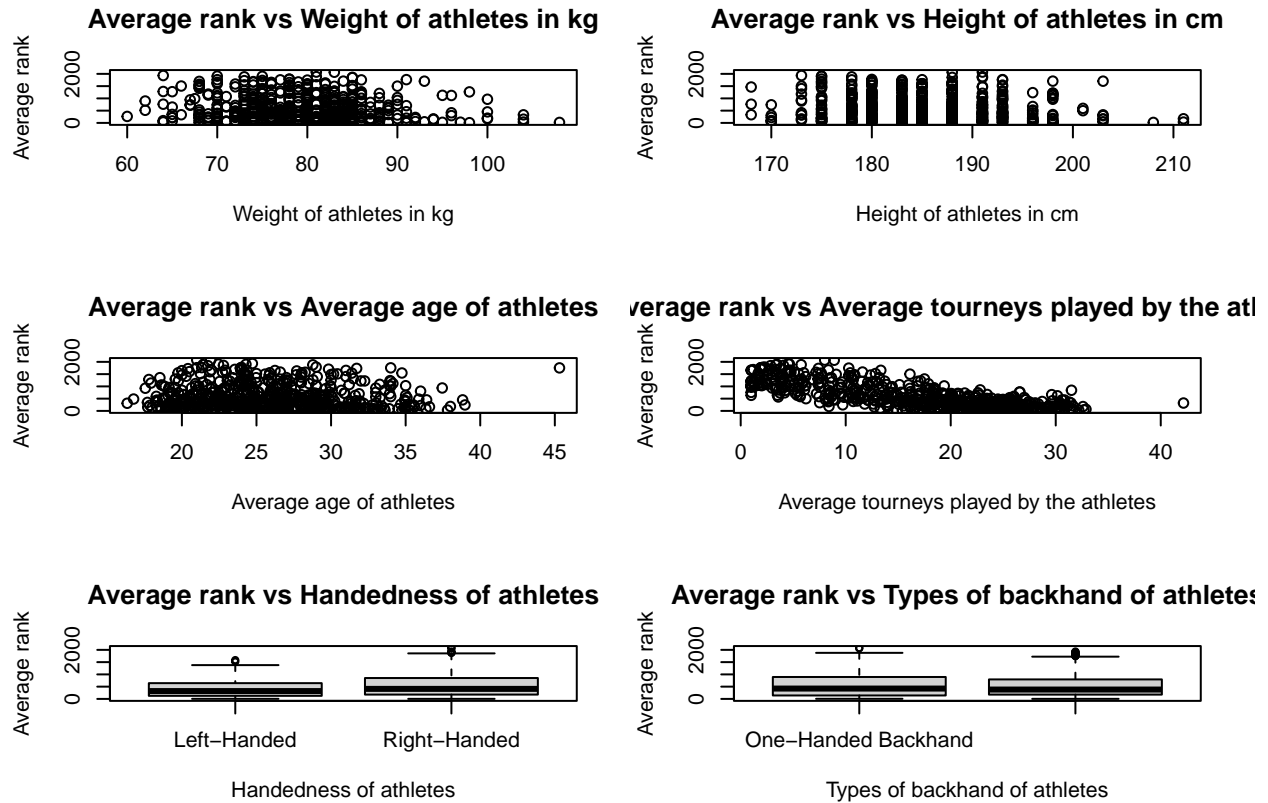


Figure 2: plots showing correlations between each predictor and the response variable

Boxplots for other numerical variables are also shown below. It's clear that there are some outliers in weight and height variables. Moreover, height and average tournaments played variables are left skewed since the first quantile deviates from median more. And the average rank is right skewed, which is consistent with the conclusion from histogram.

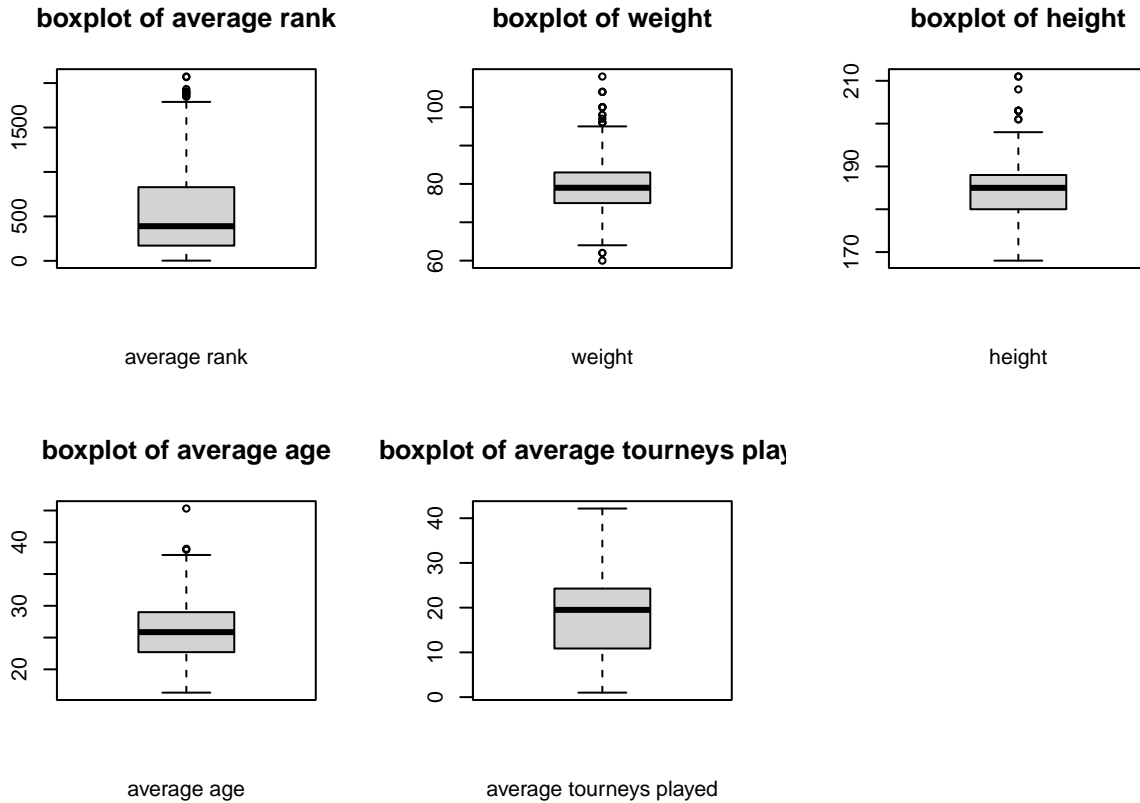


Figure 3: box plots for numerical variables

The following plot shows the scatterplot of average height and average rank in each category of handedness. We can see that negative linear relations in both category, while the slope seems to be more negative in the right-handed group. This trend suggest that higher players are expected to have better rank, and a right handed player is more affected by height. This conclusion lines up with the hypothesis we had from the summary statistics.

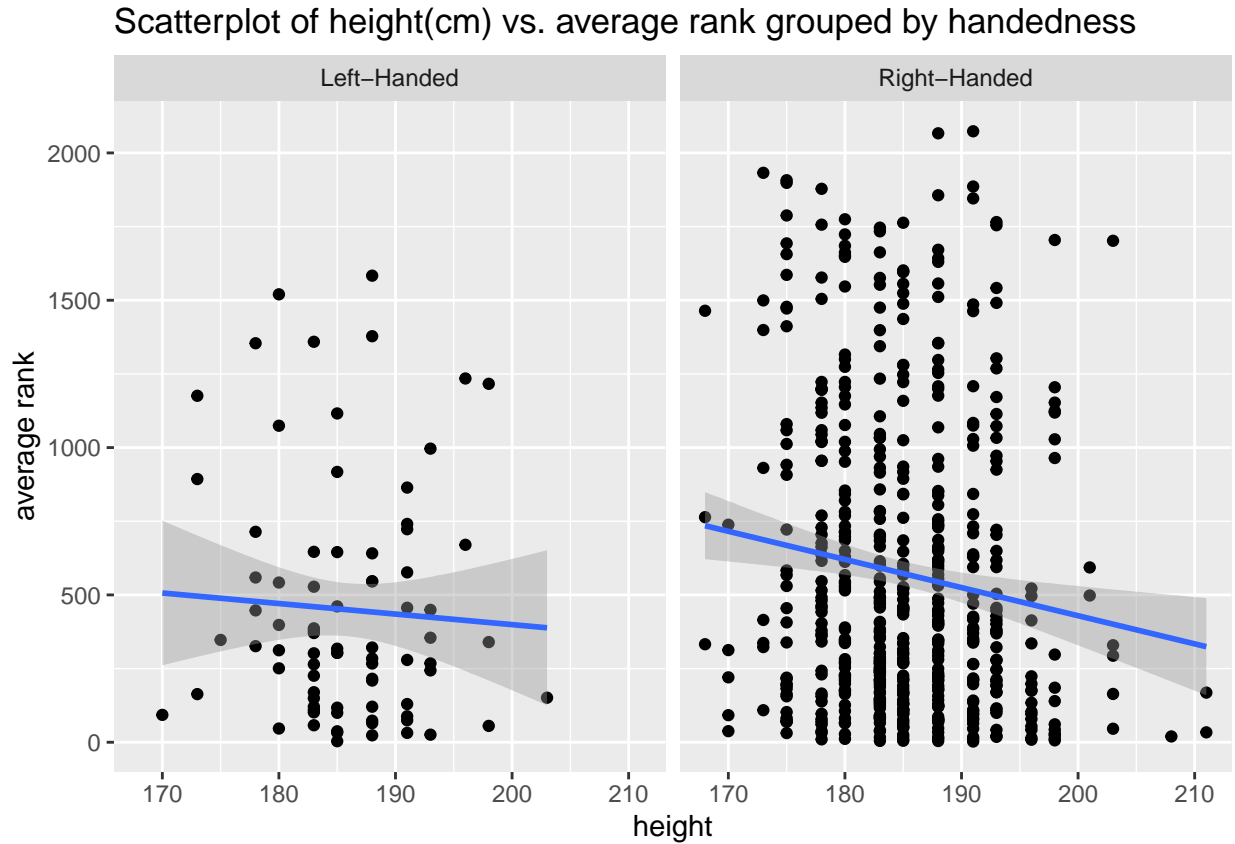


Figure 4: Scatterplot of height(cm) vs. average rank grouped by handedness

The following plot shows the scatterplot of average height and average rank in each category of backhand. This pair of plots present a interesting pattern, where we can see a negative linear relations in two-handed backhand category, while the slope become a slight positive one in the one-handed backhand group. This trend may suggest that the conclusion we had before saying that higher players are expected to have better rank is not suitable on one-handed backhand players. However, it's also possible that the slight positive trend is due to a lack of observations, since there are only 104 players in this category.

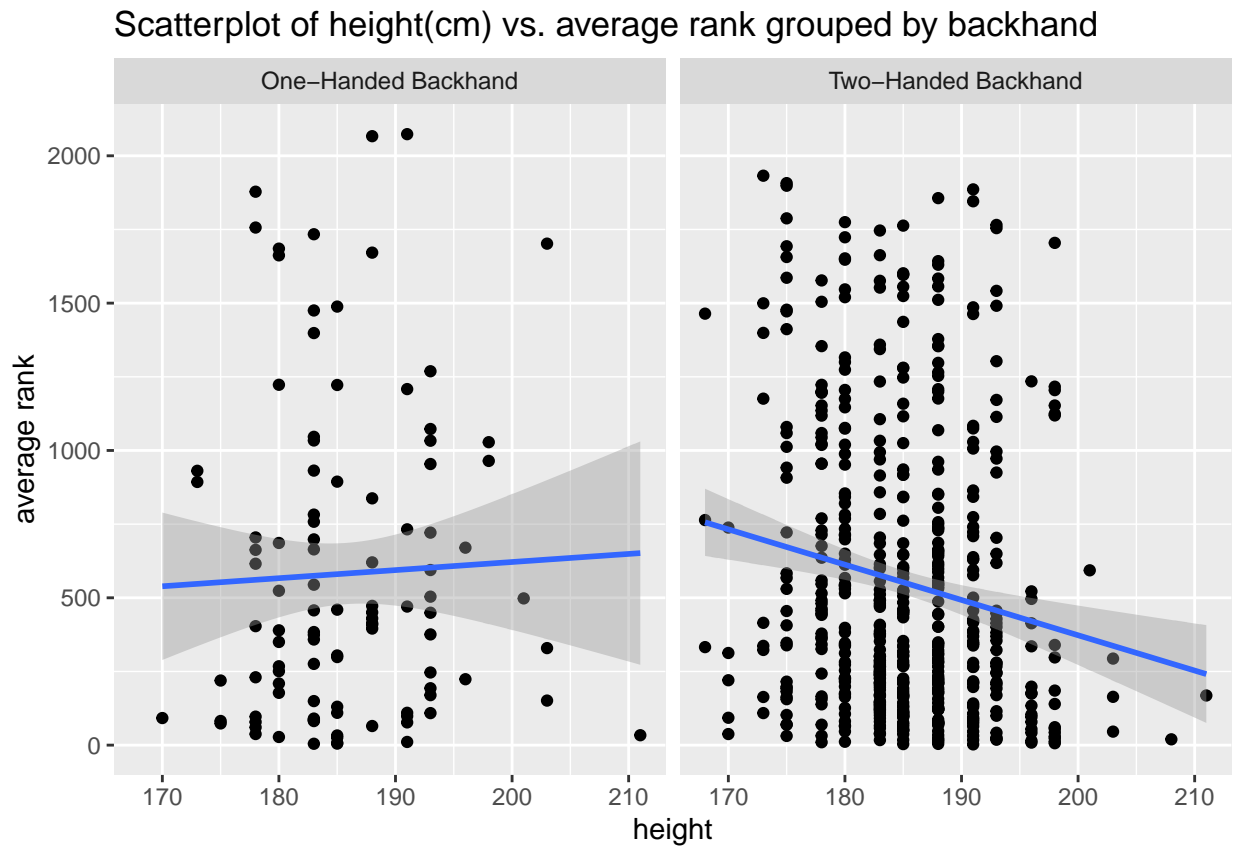


Figure 5: Scatterplot of height(cm) vs. average rank grouped by backhand

From the following max-min plot, we can see that the mean rank of left-handed players is better, while the difference between max and min is larger in right-handed players group. The means in two backhand categories are almost the same, while the difference between max and min is larger than in one-handed players group. This is interesting since one-handed players are the minority in the dataset, however, they had a larger max-min difference.

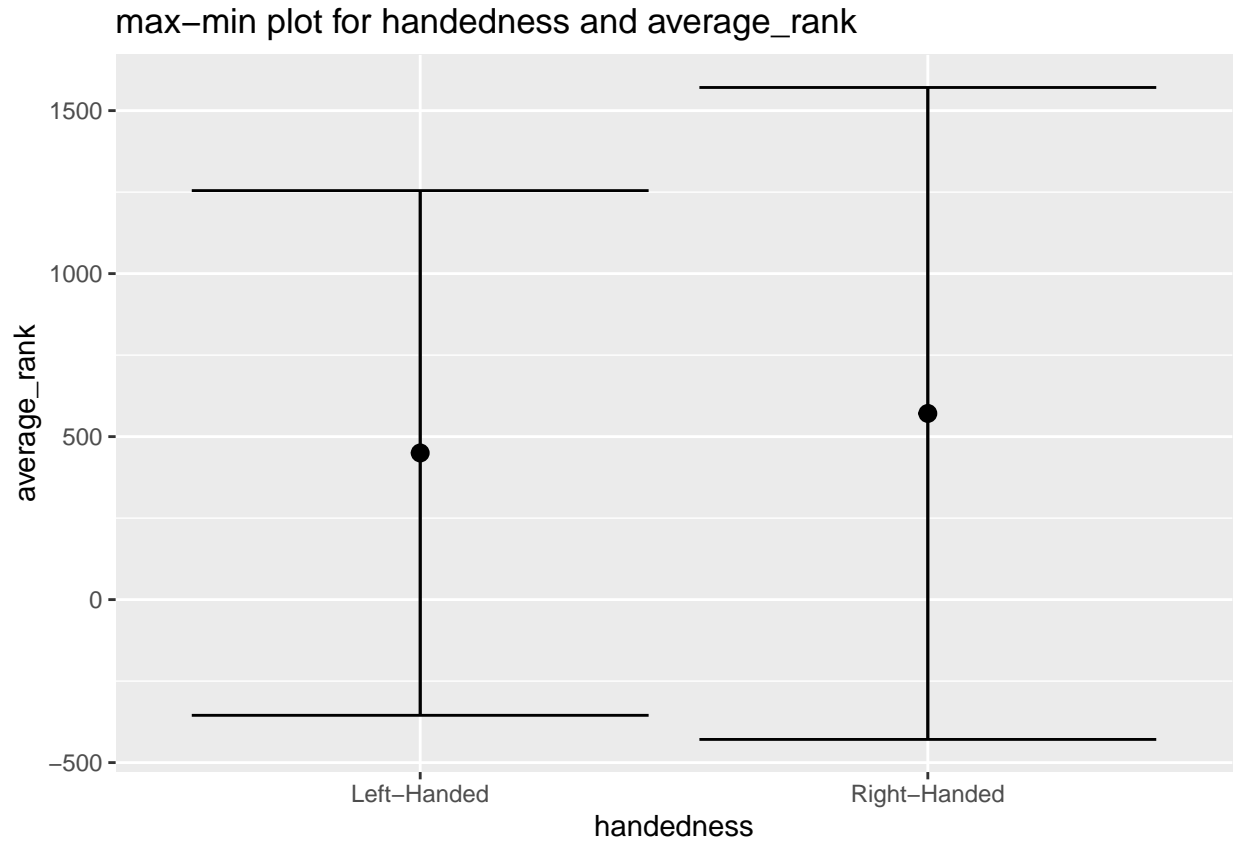


Figure 6: Max-min plots for binary variables

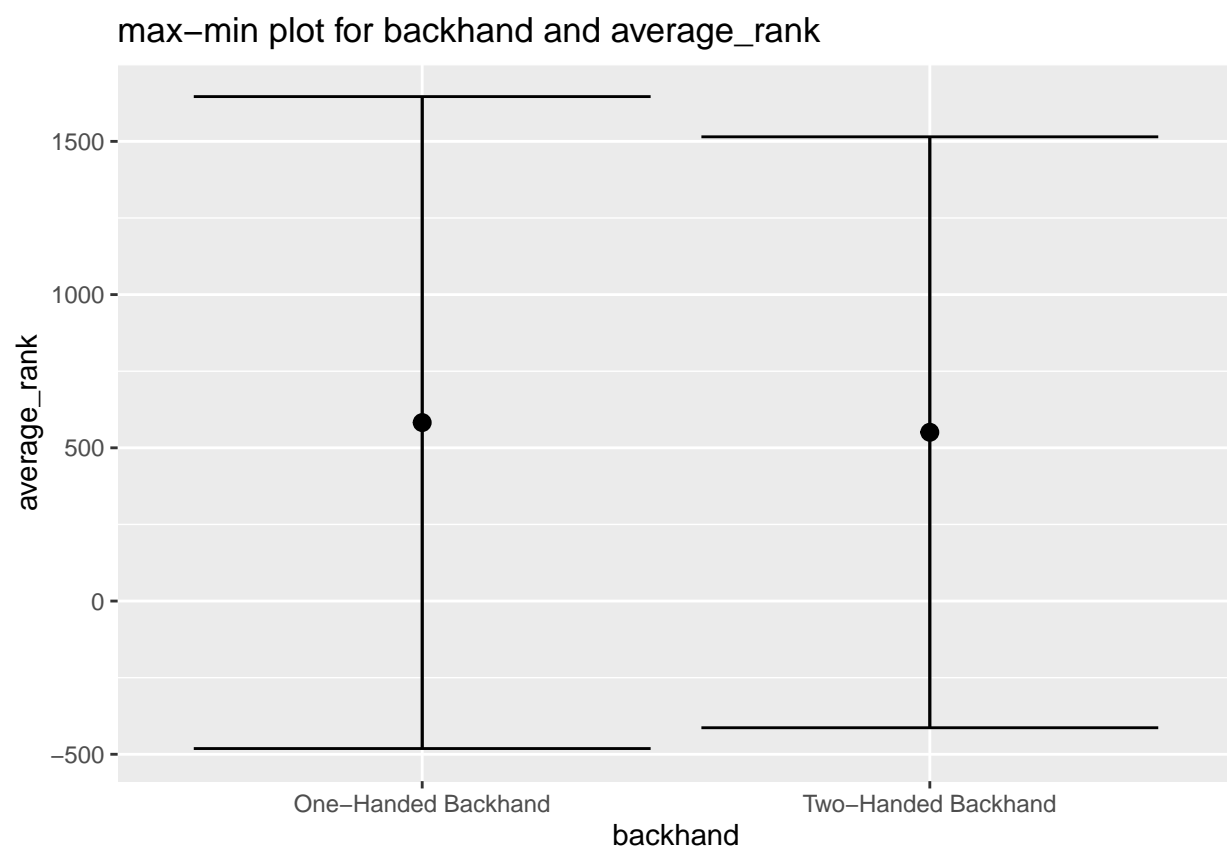
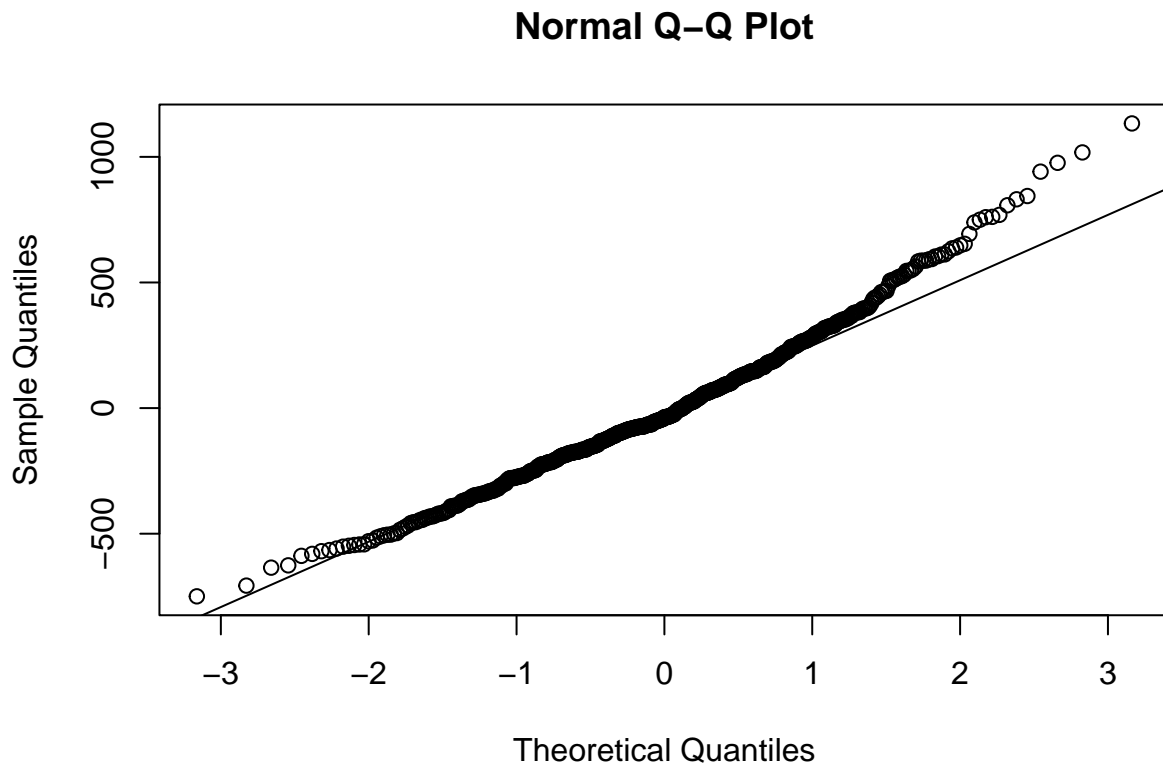


Figure 7: Max-min plots for binary variables

We first built a linear regression model taking account into all predictors. From the following QQ plot, we can see large deviations at the right end, which means that the normality assumption is not satisfied currently and we need to apply some transformations to our variables.



The automated transformation result in applying a power of -0.5 to height, and a power of 0.5 to average rank. Then we refit the model and here is the summary of the model. From this summary, we can see that the square root of rank is expected to be -4.98 if given all other predictors as 0. This interpretation will not be so useful since it's impossible for all other predictors to be 0 and having a negative square root. We can also see that the square root of rank is expected to decrease 0.96 given a one-unit increase in the average tournaments played while holding other variables in the model constant. Moreover, the square root of rank is expected to decrease 0.51 given a one-unit increase in the average age while holding other variables in the model constant, and the square root of rank is expected to increase 821.80 given a one-unit increase in the negative square root of height while holding other variables in the model constant. Note that since the height is taking a negative power, the relation between height and rank should be negative instead. From the p values of variables, we can see that only average tournaments played, age, and height are significant linear predictors in the model. We also obtained the R-squared to be 0.6434, which means that the linear model explains 64.34% of variance in FEV. The adjusted R^2 is 0.64, which is not bad.

| | Estimate | Pr(> t) |
|--------------------------------|-------------|-----------|
| (Intercept) | -4.9822069 | 0.8362737 |
| weight_kg | -0.0434122 | 0.4320916 |
| average_tournaments_played | -0.9632150 | 0.0000000 |
| average_age | -0.5106662 | 0.0000000 |
| backhandTwo-Handed Backhand | -0.7591655 | 0.2666604 |
| handednessRight-Handed | 0.0645296 | 0.9316181 |
| negative_square_root_of_height | 821.8036938 | 0.0038217 |

The QQ plot and residual plots for the model after transformation is given below. We can see that the QQ plot looks better and the residual plots looks randomly scattered around 0.

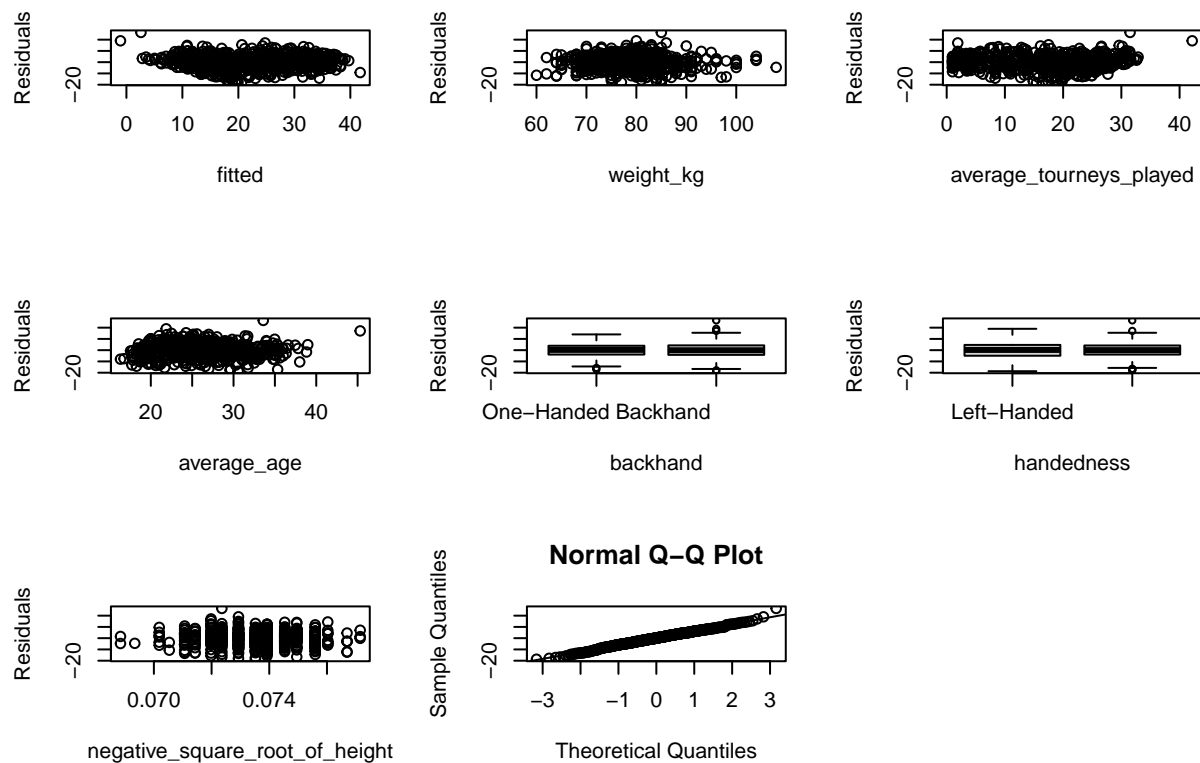


Figure 8: Plots checking assumptions for linear model

The following plot presents a linear regression line and gam regression line for negative square root of height and square root of rank. These two lines are very close which means that the relationship between negative square root of height and square root of rank could be explained by linear regression very well, and the gam regression don't need a non-linear line to explain this relationship. Moreover, we built the gam model where we put a cubic regression spline on negative square root of height. The adjusted R square for this gam model is 0.64, which is the same to the linear model.

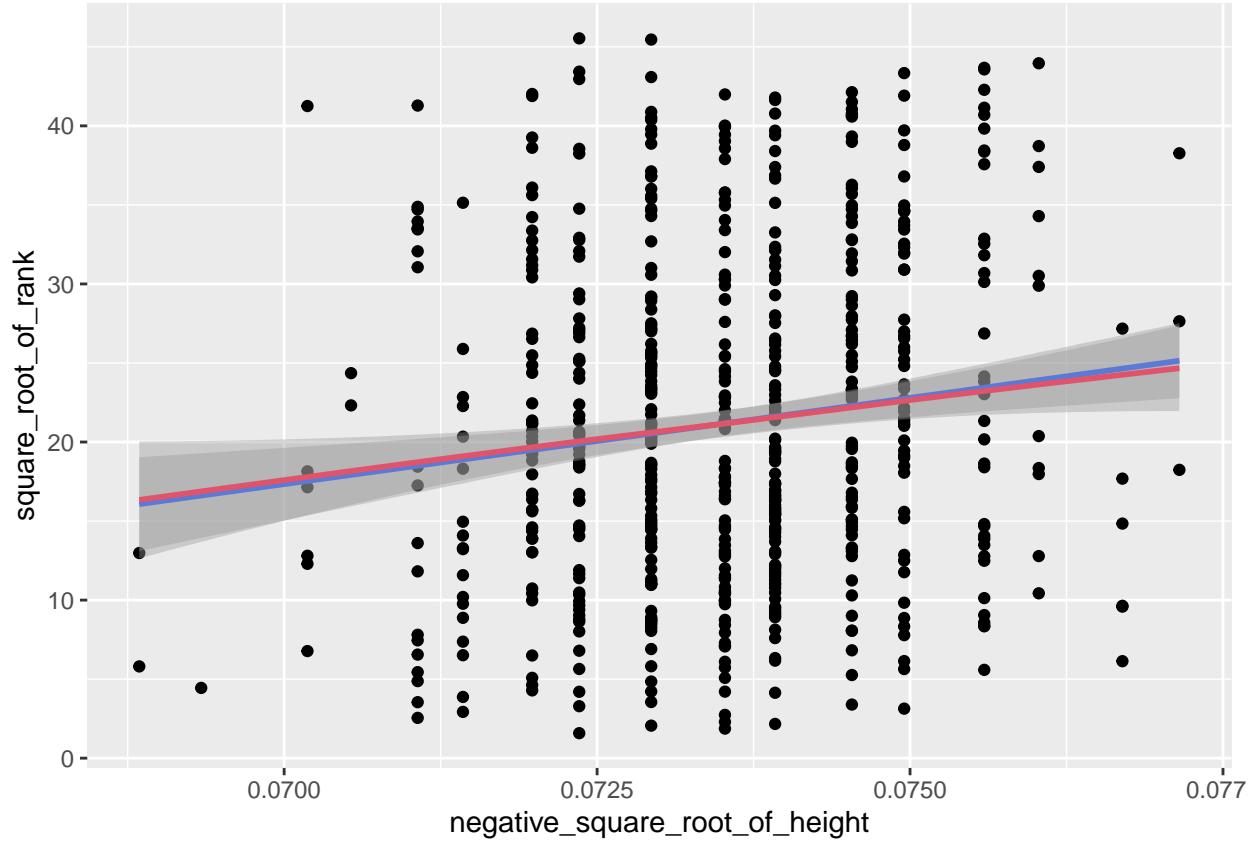


Figure 9: Scatterplot of $\text{height}(\text{cm})^{(-0.5)}$ vs. $\text{average rank}^{(0.5)}$

Conclusion

The linear model is telling that athletes are expected to have better rank if he played more tournaments while holding other variables the same, or he is older holding other variables the same, or he is higher holding other variables the same. Note that this model aligns with the study mentioned in introduction since height of the athletes could explain the variation of the serve speed which further explains the variation of ranking. The gam model has the same adjusted R square, which means that the linear model could fit the data well enough. And the gam model also suggest that a higher tennis athlete is expected to have a better ranking. In the future research, this linear model could be further improved to obtain a larger adjusted R square by performing model selection process on it. Similarly, the gam model could be further improved by trying cubic regression spline on different predictors.

The linear model also has some limitations such as the QQ plot has some deviations at right end and the residual plot for average age has a cluster on the left. Moreover, the comparison between two models are preliminary since we are only looking at the adjusted R square. We could compare their AICs and BICs in

the future. Lastly, the research is conducted on the 2017 dataset since this is the latest dataset we found, the result concluded from this dataset may not be the case 10 years later.

References

- ATP World Tour. (2018, October). ATP World Tour tennis data. Retrieved from DataHub: <https://datahub.io/sports-data/atp-world-tour-tennis-data#r>
- Fryer CD, et al. (2018). Mean body weight, height, waist circumference, and body mass index among adults: United States, 1999–2000 through 2015–2016.
- Peng, Y. (2020). Research into the performance of tennis players [Unpublished paper]. University of Toronto.
- R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Vaverka, F., & Cernosek, M. (2013, March 12). Association between body height and serve speed in elite tennis players. Retrieved from Scholar Portal Journals: https://journals-scholarsportal-info.myaccess.library.utoronto.ca/details/14763141/v12i0001/30_abbhassietp.xml