

Springboard

Data Science Career Track

Capstone Project 1

**What are the factors that contribute to a home team
winning an NCAA basketball game?**

Milestone Report (May 2020)

by

Adeyemi Adejuwon

Abstract

From 1985-2019, teams in the regular season of NCAA Division 1 Men's basketball games won 92,732 of their home games but only won 47,547 of their away games, and 15,805 of matches played in neutral locations. This home court advantage is further confirmed with a Tukey HSD test, where the statistically significant difference between the pairs home and away games, the pairs home and neutral games, versus the pair away and neutral games was calculated.

Many articles have explained that the reason for a team's home-court success is the presence of fans and the arena, as opposed to the fact that they are simply the better team in a particular matchup. Assuming fans are the reasons for this home court advantage, the odds are already stacked against a visiting team, when playing these basketball matches. This report tries to analyse some of the differences between playing on the road and playing at home.

In this report, evaluation of the in game statistics was done from matches played during the year 2003-2019. The in-game statistics showed that the average points posted for the winning team when playing at home was about three points higher than when playing away or at a neutral location; the average three-pointers scored by the winning team was similar when playing at home versus when playing away even though there was a general trend in increasing three-pointers scored from season season to season, and the turnover conceded by the winning team was 0.5 points higher when playing at away matches than when playing at home or at a neutral location.

In-game analysis of the data also showed strong positive correlations between the field goals made and assists for the winning teams, while there were weak correlations between turnovers conceded to the winning team to the score made by the winning team.

A comparison of the teams with the highest winning percentages in the seasons 1985 -2019 and from 2003-2019 were made. The reason for this is that the in-game statistics was only available for seasons 2003-2019. In these comparisons, Kansas and Duke had the highest winning percentages for the seasons 1985- 2019, while Gonzaga and Kansas had the highest winning percentages for the years 2003-2019.

These results showed that Gonzaga must have improved a lot from the year 2003-2019, Duke might have dropped in form during this period, while Kansas maintained their consistency during both time periods of comparison.

Based on this information a comparison was made between the average winning scores in matches for Kansas and Gonzaga. The results showed that the average winning scores when these teams play their opponents could not be proven to be statistically significant. The average winning score is 75 points in matches won by both teams.

Abstract	2
Introduction	5
Objective	5
Dataset	5
Data Cleaning & Wrangling	7
Data Type	7
Handling Missing Data	7
Exploratory Data Analysis	8
Box Plot - points scored by winning team at playing locations	9
Box Plot - points scored by losing team at playing locations	10
Bar Chart - Proportion of games played at different locations in all seasons	11
Line Chart - Points scored per winning team per season	13
Line Chart - Three-Pointers scored per winning team per season	14
Line Chart - Turnovers conceded per winning team per season	15
Bar Chart - Ranking top 15 teams by win percentage for 2003 -2019	16
Bar Chart - Ranking top 15 teams by win percentage for 1985 -2019	16
Applications of Inferential Statistics	17
ANOVA test & Tukey HSD	18
Heat map with correlation matrix and p value	19
Pair plot matrix	21
Bootstrap Replicates	22
Summary and Next Steps	24

Introduction

The phenomenon of home field advantage is nothing new in the world of sports¹. Home teams generally have the advantage of having the support of larger, more enthusiastic crowds, and it has been suggested in some studies, that it is the home court atmosphere that enhances the home teams opportunities for winning matches².

Objective

The goal of my project is to study game factors that affect a home team winning an NCAA Division 1 Men's basketball game, and likewise affect an away team losing a match. The information derived from these analysis can help a basketball coach tune his game play tactics and thereby increase his odds of winning a game when playing at an away location.

Dataset

The datasets for this analysis will be obtained from Kaggle³. These data is associated with an annual competition sponsored by Google. The datasets explored in our analysis were from the tables

Regular SeasonDetailed Results.csv

The regular season detailed results file identifies the game by game match play data and results for regular season matches for the years 2003 - 2019. The dataset for the regular detailed season games consists of 87,366 data points for 350 college basketball games playing in the NCAA since 2002. This dataset includes 34 variables such as number of assists, three-point percentages per game, win and loss records per season, location of matches played.

¹ "Home-Field Advantage (SOCIAL PSYCHOLOGY"
<http://psychology.iresearchnet.com/social-psychology/control/home-field-advantage/>. Accessed 14 May. 2020.

² "The Home Court Advantage: Evidence from Men's College" 9 Mar. 2017,
<http://thesportjournal.org/article/the-home-court-advantage-evidence-from-mens-college-basketball/>. Accessed 14 May. 2020.

³ <https://www.kaggle.com/ncaa/ncaa-basketball>

Regular Season Compact Results.csv

The regular season compact results identify just the team losses and wins from 1985-2016. The dataset for regular season compact games consists of 156,089 entries and 8 variables. These variables do not include in-game data.

Teamspellings.csv

The team spellings file is used to correlate TeamID numbers with their associated names.

- WTeamID - this identifies the id number of the team that won the game.
- WScore - this identifies the number of points scored by the winning team.
- LTeamID - this identifies the id number of the team that lost the game.
- LScore - this identifies the number of points scored by the losing team.
- WLoc - this identifies the "location" of the winning team. The home team is given the value "H", while the visiting team is given the value "A", and the value "N" is given to a match played on a neutral location.
- NumOT - this indicates the number of overtime periods in the game, an integer 0 or higher.
- WFGA - field goals attempted (by the winning team)
- WFGM3 - three pointers made (by the winning team)
- WFGA3 - three pointers attempted (by the winning team)
- WFTM - free throws made (by the winning team)
- WFTA - free throws attempted (by the winning team)
- WOR - offensive rebounds (pulled by the winning team)
- WDR - defensive rebounds (pulled by the winning team)
- WAst - assists (by the winning team)
- WTO - turnovers committed (by the winning team)
- WStl - steals (accomplished by the winning team)
- WBlk - blocks (accomplished by the winning team)
- WPF - personal fouls committed (by the winning team)

Data Cleaning & Wrangling

Before beginning analysis of the data, it is essential to explore there are no missing values in our dataset. It is also essential for the data in our table to be of the correct data type. This upfront work will give more confidence in interrogating the data and would allow better conclusions to be made as regards the dataset. The libraries used for the data cleaning and wrangling of the data sets are:

- numpy for scientific computing of the numerical arrays
- pandas for data analysis and manipulation ,
- matplotlib for visualization

Data Type

The next step in the data wrangling stage was to determine the type of data type in the dataset. As can be observed in the table below, there are 87,366 entries, with no missing values in any of the 34 columns. Additionally, all but one column takes integer values , whereas the lone column (WLoc) takes a string entry.

In fact from the data set description we know that the WLoc column will take only three values each one representing the location of games played. To confirm the entry of the WLoc column we call the unique () function on that column.

- “H” stands for “Home game
- “A” stands for away (visiting to opponent’s site)
- “N” is the location of games played at a neutral location

Handling Missing Data

The next step in the data wrangling stage is to check for any gaps in the dataset. This is confirmed with the function “is null” and “value_counts”. The isnull function finds the null value in the data set, while the value_counts function displays the amount of the categorical variables in WLoc.

Based on these results, we can observe that there are **no missing** values in the dataset.

It should be noted that if there were missing values in the column we would either drop them or fill them in. This is because some of the techniques in the data exploratory will not allow for missing data.

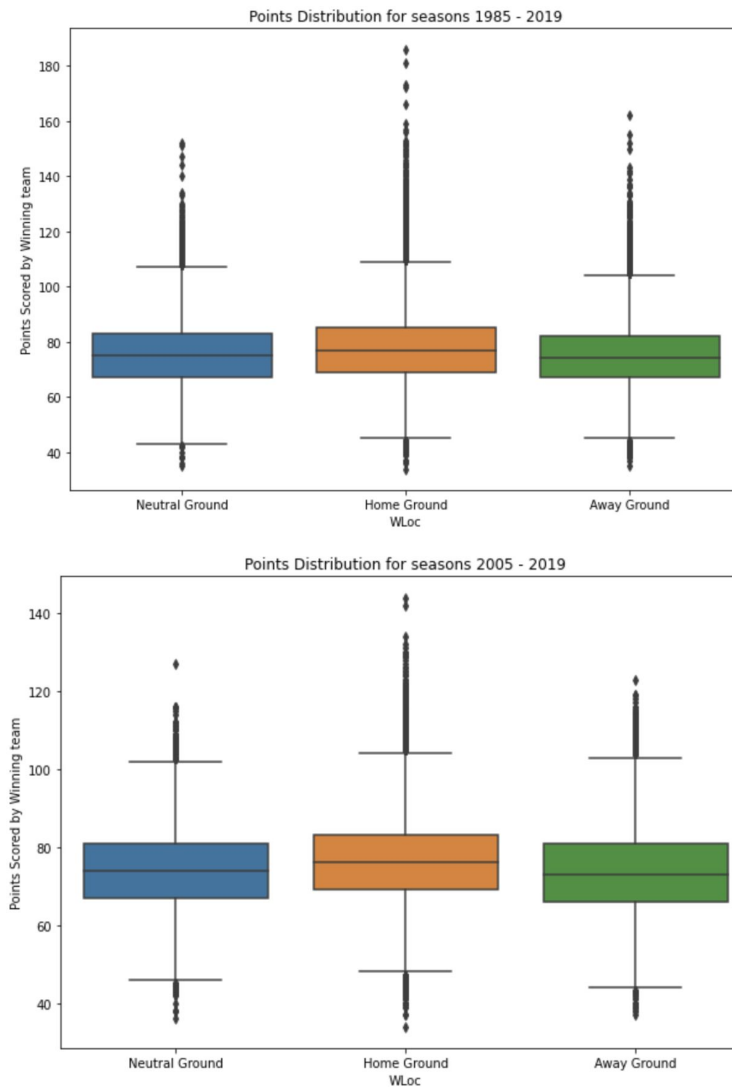
Exploratory Data Analysis

Following the data wrangling of the data, the next step is to interrogate the dataset and ask a series of questions of the dataset. These questions will help identify the contributing factors affecting home games winning matches. The questions being asked of the data are :

1. Does a winning team score more points when playing at home, than when playing at either a neutral ground or an away ground?
2. Does a losing team score more points when playing at home, than when playing at either a neutral ground or an away ground?
3. Is there a difference in the amount of matches a team wins at home, away or a neutral location?
4. What is the average variation in points scored for the winning team per season?
5. What is the average variation in three-points scored by the winning team when playing at home, away or a neutral location per season?
6. What is the average turnover by the winning team when playing at home, away or a neutral location per season?
7. What is the ranking of the top 15 teams based on the winning percentage per season for seasons 1985-2019?
8. What is the ranking of the top 15 teams based on the winning percentage per season for seasons 2003-2019?

1. Does a winning team score more points when playing at home, than when playing at either a neutral ground or an away ground?

Box Plot - points scored by winning team at playing locations

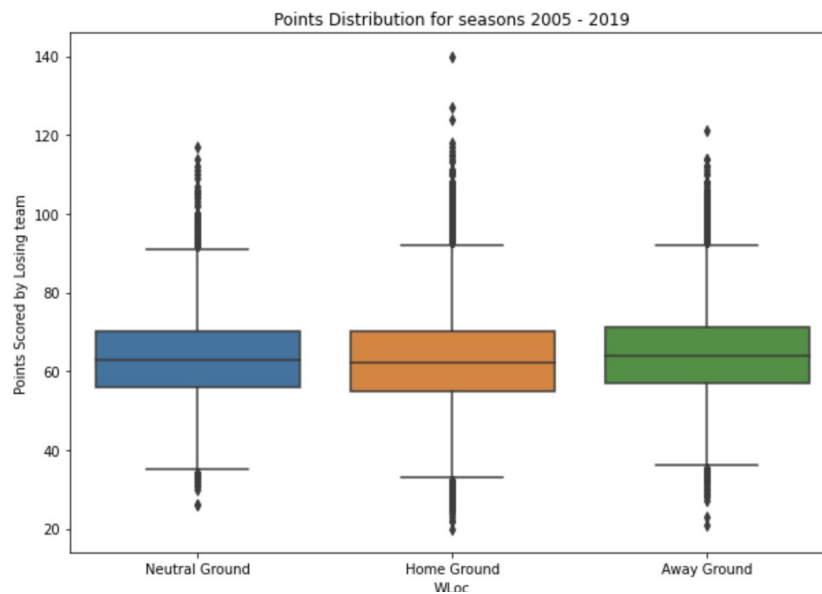
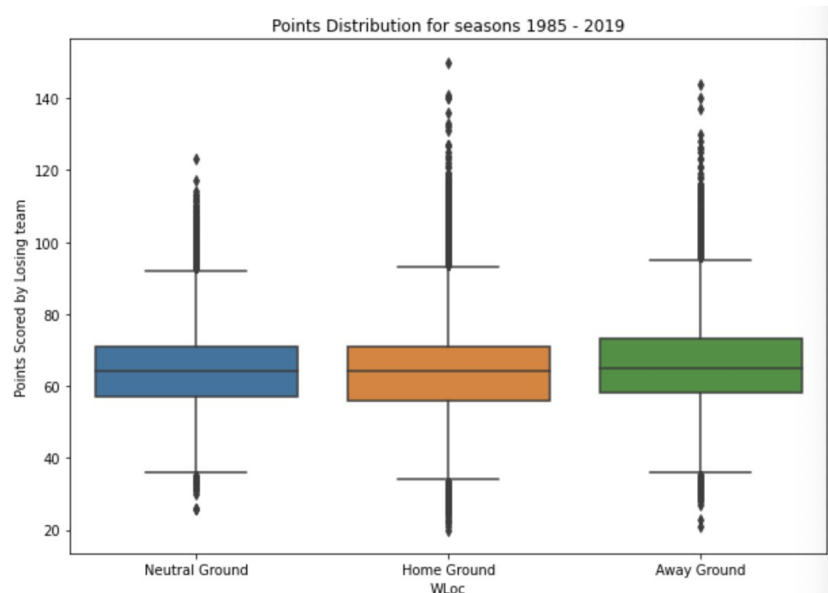


Conclusion

The box plots show that the average number of points for the winning team at their home ground is higher than the points scored when playing at either away or neutral grounds. This result will be further investigated in the applications inferential statistics section of the report. This box plot confirms the phenomenon of home advantage being present for home teams.

2. Does a losing team score more points when playing at home than when playing at either a neutral ground or an away ground?

Box Plot - points scored by losing team at playing locations



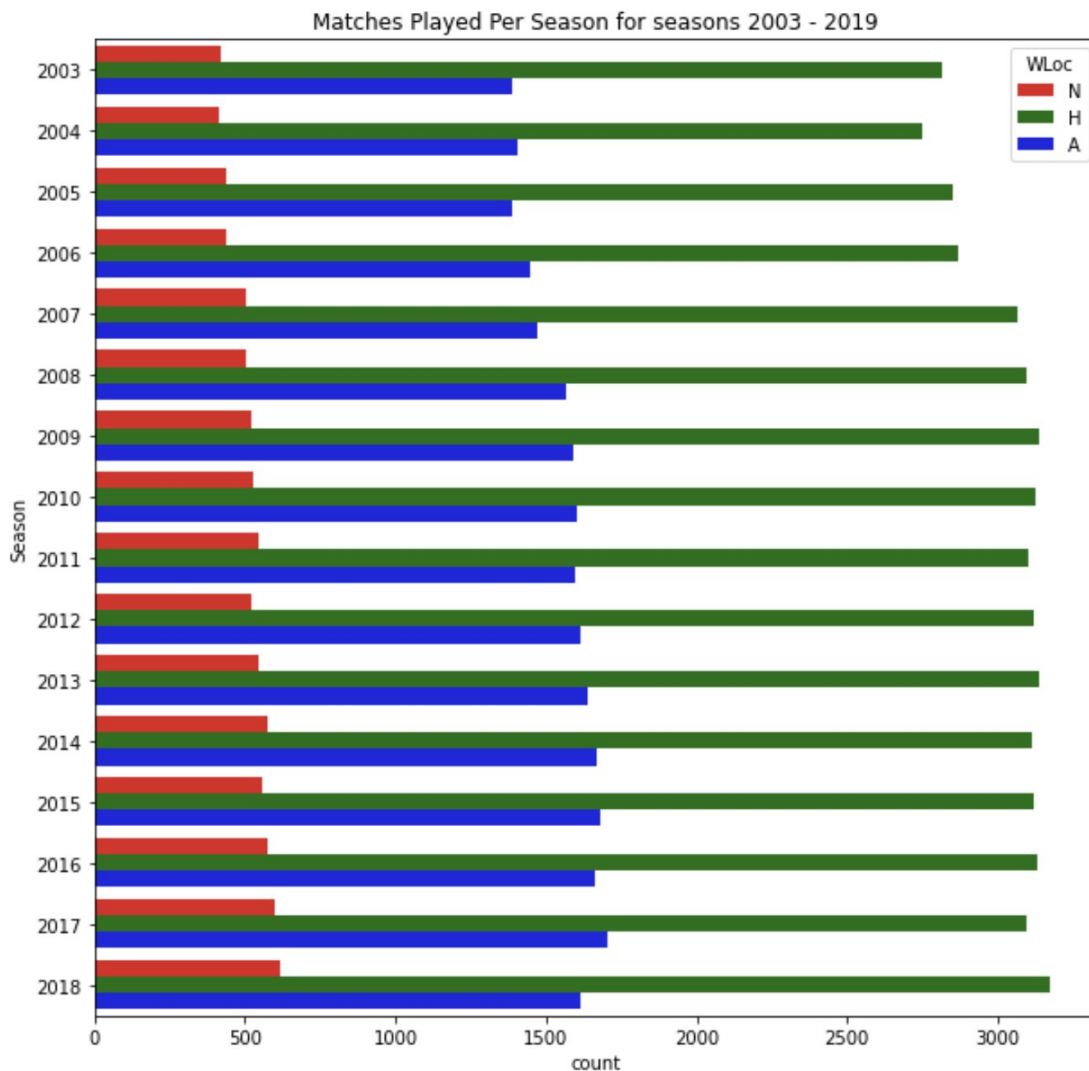
Conclusion

The plots seem to show that the losing teams score more points when playing at an away ground versus when playing at their neutral or home ground. However, there seems to be no difference to the average number of points scored by the losing team when playing in either of the neutral or home locations.

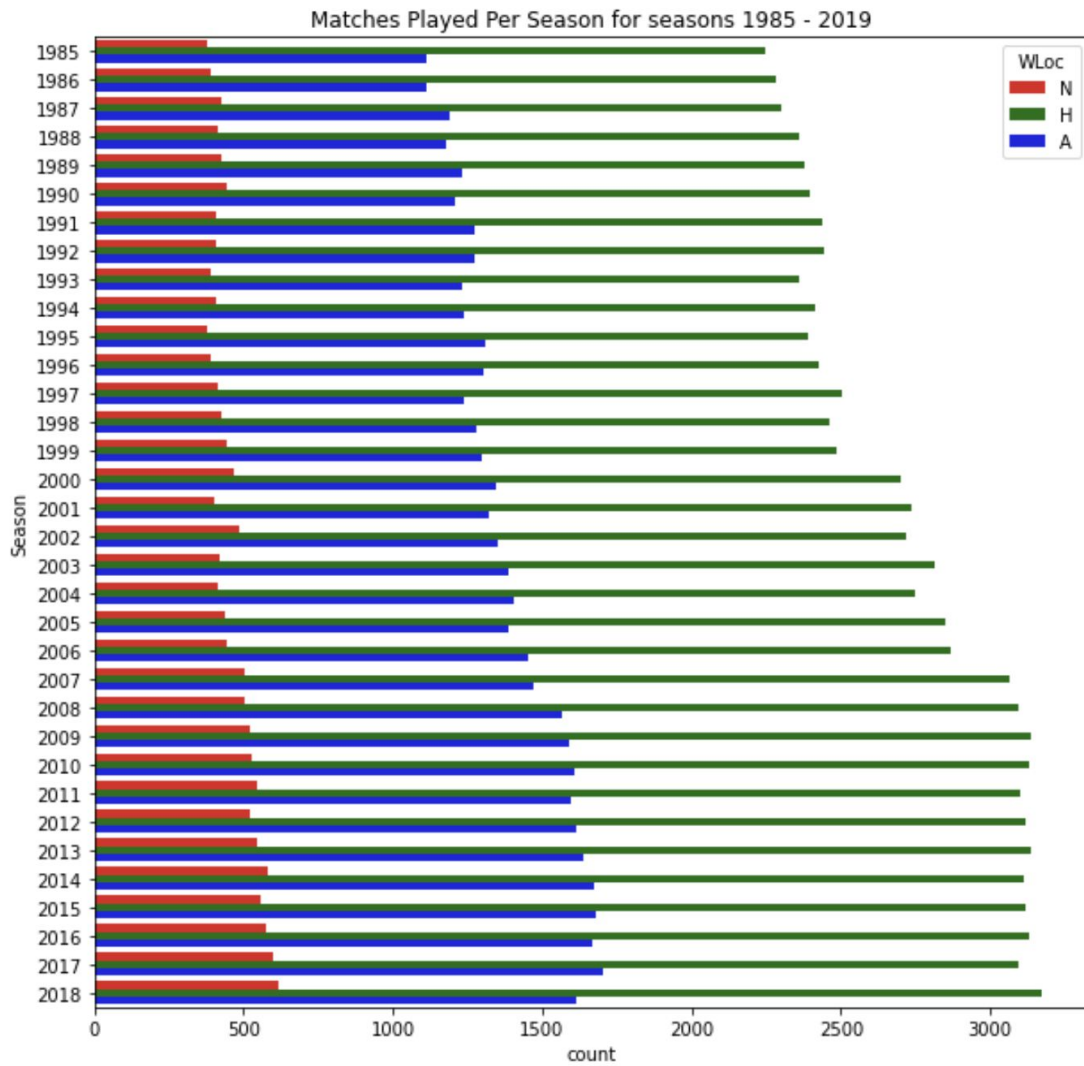
3. *What is the difference in the amount of games a team wins either home, away or at a neutral location?*

Bar Chart - Proportion of games played at different locations in all seasons

[19]



[20]



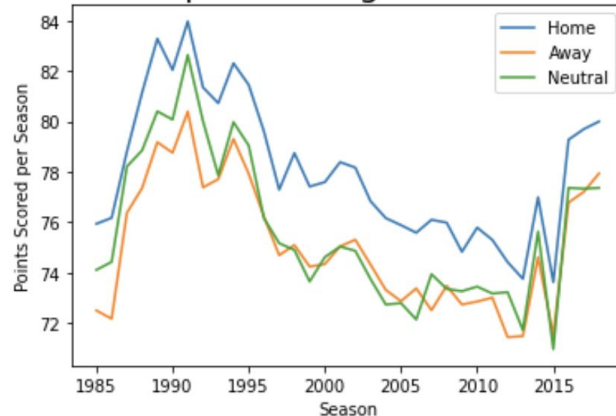
Conclusion

The proportion of all games won at home, away and in a neutral ground were all similar across all seasons.

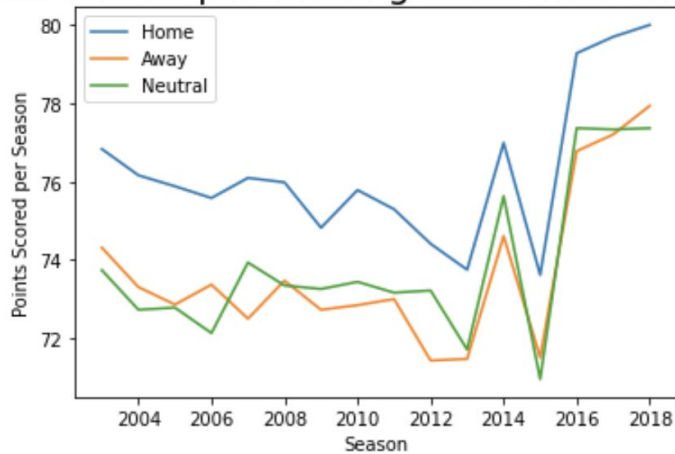
4. What is the average variation in points scored for the winning team per season?

Line Chart - Points scored per winning team per season

☞ Points Scored per winning team from 1985-2019



☞ Points Scored per winning team from 2003-2019



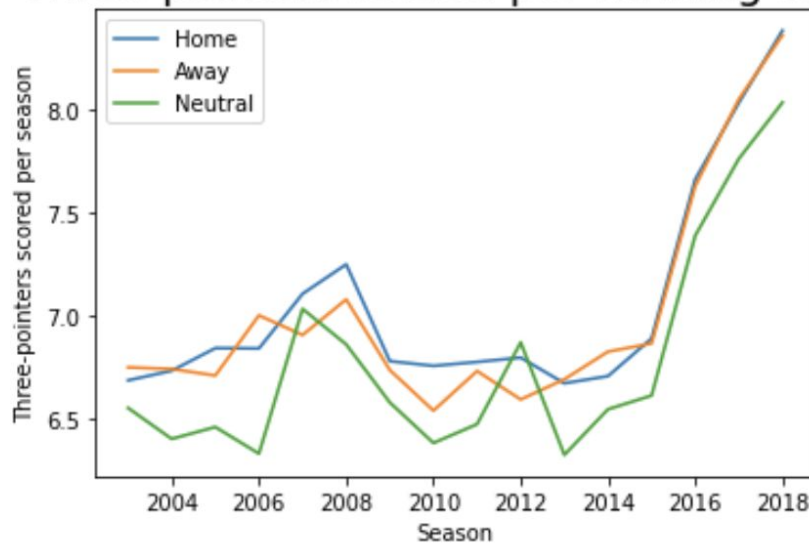
Conclusion

For all seasons the number of points scored by the winning team while playing at home is consistently greater than points scored when playing at an away or a neutral location. After an initial upward trend in points scored from 1985 -1990, the total number of points scored per winning team started decreasing, and later started following a general upward trend in the points scored per season, from year 2016 onwards. Some suggestions for the uptick in scoring are that maybe teams are becoming more efficient in scoring, or it could also be that the rules of basketball have changed to favor the scoring team. Another observation from this analysis is that in 2015 there was a dip in scoring in all locations. This observation would require further investigation as to why there was a specific drop in scoring this particular year.

5. What is the average variation in three-pointers scored by the winning team when playing at home, away or a neutral location per season?

Line Chart - Three-Pointers scored per winning team per season

☞ Three-pointers scored per winning team

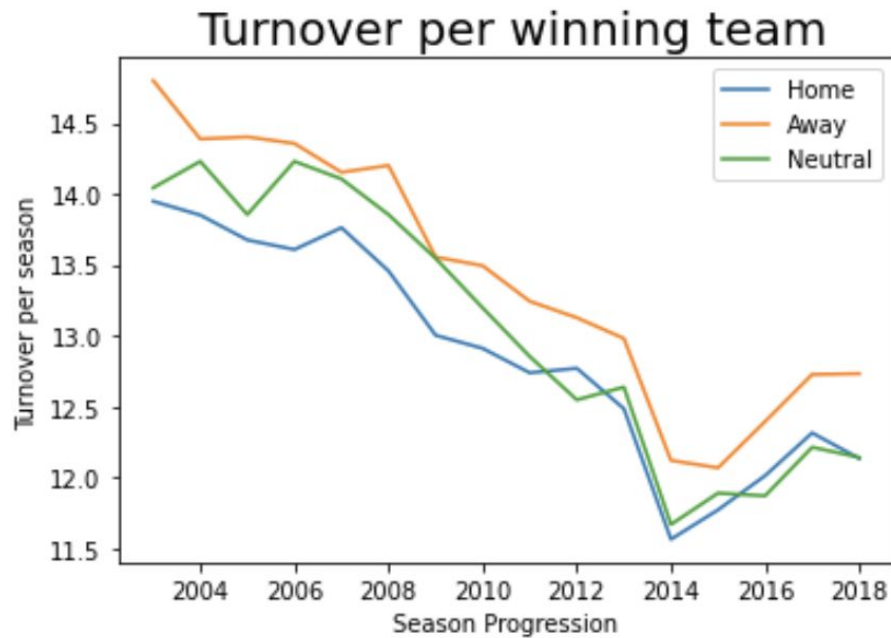


Conclusion

There seems to be a general upward trend in three pointers scored as the years go by. It can also be observed for most seasons the number of three-pointers scored by the winning team while playing at home is greater than three-pointers scored when playing at an away, or a neutral location. Some of the reasons for this general increase in three-pointers could be the game of basketball is evolving to more teams scoring more three pointers.

6. What is the average amount of turnovers per season?

Line Chart - Turnovers conceded per winning team per season

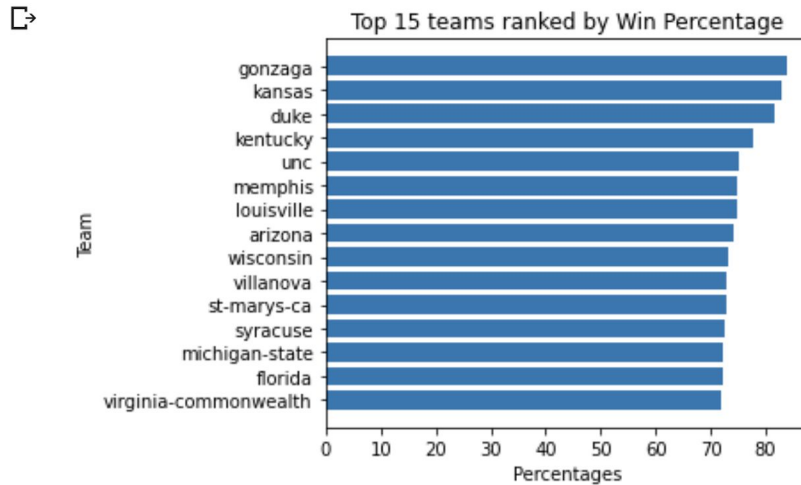


Conclusion

For all seasons the number of turnovers conceded while playing at an away location for the winning team is greater than turnovers conceded when playing at an home or a neutral location. There also seems to be a decrease in number of turnovers through the years.

7. What is the ranking of the top 15 teams based on the winning percentage per season for seasons 2003 -2019? Winning Percentage is defined as (matches won/total matches played *100)

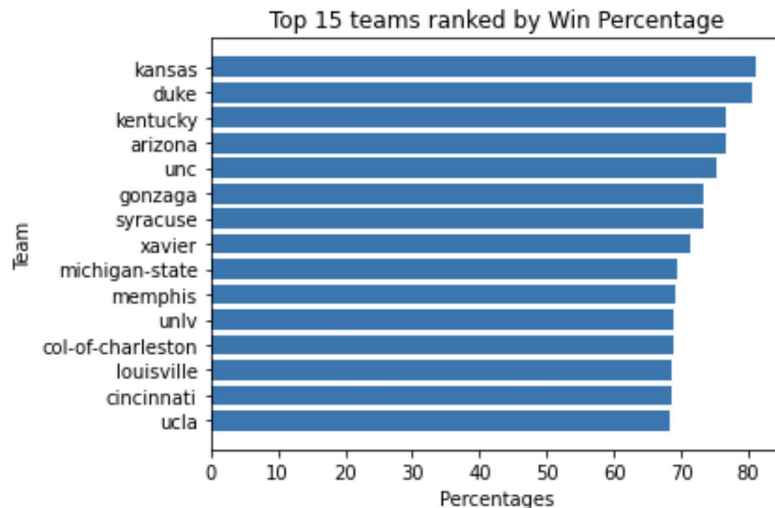
Bar Chart - Ranking top 15 teams by win percentage for 2003 -2019



The team with the highest winning percentage is Gonzaga followed by Kansas

8. What is the ranking of the top 15 teams based on the winning percentage per season for seasons 1985 -2019? Winning Percentage is defined as (matches won/total matches played *100)

Bar Chart - Ranking top 15 teams by win percentage for 1985 -2019



The team with the highest winning percentage is Duke followed by Kansas

Applications of Inferential Statistics

Upon exploration of the data and making initial observations about the data, the next step will be to find out whether there are any correlations within the in game statistics, and also confirm whether some of our initial observations are statistically significant. The questions that was asked of the data was:

1. Is the difference between a team winning college basketball matches at home locations versus playing at away or neutral locations statistically significant ? This question will be answered with the aid of an ANOVA (Analysis of Variance) test, followed by a pairwise Tukey HSD (Honest Significant Difference) correlation.
2. What is the correlation between game-by-game data for a winning team? A heat map will be used to show the correlations between the game by game data. The correlation matrix from the heat data will show the correlations that exist between in game statistics during the season.
3. What is the 95% confidence interval for the difference between the standard deviations of the winning scores for the two top performing teams? This analysis will be done with a

bootstrap sampling analysis, calculating these differences over 10000 replicates. The teams being considered will be Gonzaga and Kansas.

1. Is the difference between a team winning college basketball matches at home locations statistically significant from when playing at away or a neutral location? This question will be answered with the aid of an ANOVA (Analysis of Variance) test, followed by a pairwise Tukey HSD correlation.

When comparing the three numerical datasets, ANOVA (Analysis of Variance) was used to test the null hypothesis that all of the datasets have the same mean. If we reject the null hypothesis with ANOVA, we're saying that at least one of the sets has a different mean; however, it does not tell us which datasets are different.

We can use the SciPy function `f_oneway` to perform ANOVA on multiple datasets of winning scores of teams playing in home, away and neutral locations. The `f_oneway` function takes in each dataset as a different input and returns the **t-statistic and the p-value**.

ANOVA test & Tukey HSD

```
[>] Multiple Comparison of Means - Tukey HSD, FWER=0.05
=====
group1 group2 meandiff p-adj  lower  upper  reject
-----
      a      b  -2.5265  0.001 -2.7284 -2.3247   True
      a      c  -2.253   0.001 -2.561  -1.9449   True
      b      c   0.2736  0.1247 -0.055   0.6022  False
-----
```

H_0 = all three locations have the same mean score in the basketball game

H_a = at least one of the locations have different means in the basketball game.

We will reject this null hypothesis for the pairs a & b and a & c (since we are getting a p-value less than 0.05). We are reasonably confident that a pair of datasets is statistically significantly different. After using only ANOVA, we can't make any conclusions on which two populations between the home, away and neutral locations have a significant difference.

There is a significant difference between the pairs home and away games, and the pairs home and neutral games, but there is not a significant difference between the pair away and neutral games.

2. What is the correlation between game-by-game data for a winning team?

A heat map analysis was conducted on the data to find correlations between the game by game data. The correlation matrix from the heat map data displayed the correlations that existed between in game statistics for the **winning** team during the season.

This heat map displays the correlation coefficient, and also excludes correlations that have a p-value less than 0.05

Heat map with correlation matrix and p value

In the heat map displayed below orange means positive, and blue means negative. The stronger the color, the larger the correlation magnitude.

WScore	1														
NumOT	0.13	1													
WFGM	0.82	0.068	1												
WFGA	0.5	0.21	0.64	1											
WFGM3	0.44	0.016	0.33	0.2	1										
WFGA3	0.27	0.084	0.19	0.39	0.75	1									
WFTM	0.33	0.12	-0.21	-0.17	-0.2	-0.18	1								
WFTA	0.29	0.13	-0.19	-0.15	-0.22	-0.2	0.92	1							
WOR	0.11	0.086	0.13	0.58	-0.1	0.099	0.057	0.12	1						
WDR	0.16	0.097	0.14	0.21	-0.02	0.042	0.09	0.12	0.054	1					
WAsT	0.56	-0.016	0.63	0.33	0.44	0.29	-0.17	-0.16	0.0074	0.095	1				
WTO	-0.018	0.083	-0.07	-0.21	-0.084	-0.16	0.12	0.13	0.071	0.22	-0.0085	1			
WStI	0.12	0.011	0.15	0.22	-0.032	0.033	0.018	0.049	0.12	-0.2	0.12	0.14	1		
WBIK	0.066	0.015	0.083	0.11	-0.05	-0.021	0.017	0.043	0.099	0.22	0.085	0.076	0.038	1	
WPF	0.23	0.15	0.066	0.076	0.012		0.3	0.33	0.049	0.064	-0.022	0.22	0.041	-0.025	1
	WScore	NumOT	WFGM	WFGA	WFGM3	WFGA3	WFTM	WFTA	WOR	WDR	WAsT	WTO	WStI	WBIK	WPF

The heatmap can be used to investigate possible collinearity among the multiple variables in the dataset. The `.corr()` function was used in the code to show the correlation between the values. This is where we want to set our independent or target variable. Our target variable is “*WScore*”. This is the number of points scored by the winning team. We want to find out how all of the other variables affect the points scored by the winning team. In the heatmap, the dark red areas represent a positive correlation, while light blue represents a negative correlation. It is also normal that the darkest areas are a 1:1 ratio since $WScore=WScore$, $NumOT=NumOT$, etc.

While *WScore* is still our independent variable, we can see in the map below that there is little to no correlation between the WTO (-0.018), though a high correlation between WFGM and WAsT (0.82, 0.56). these relationships are obvious (assists in a basketball game and field goals made positively correlates with the scores in a match)

Features with high correlation are more linearly dependent and hence have almost the same effect on the dependent variable. So, when two features have high correlation, we can drop one of the two features. In order to decide to find which features to drop when we decide to create our machine learning model, we compare the relationship between WFGM and WAsT, and we get a correlation number of (0.63). Dropping either of these variables when tuning our eventual machine learning model might produce a more accurate prediction model.

In addition to plotting the correlation coefficients in the heat map, only significant p-value correlations ($\alpha = .05$) were plotted. This was achieved with the `def corr_sig` function. It can

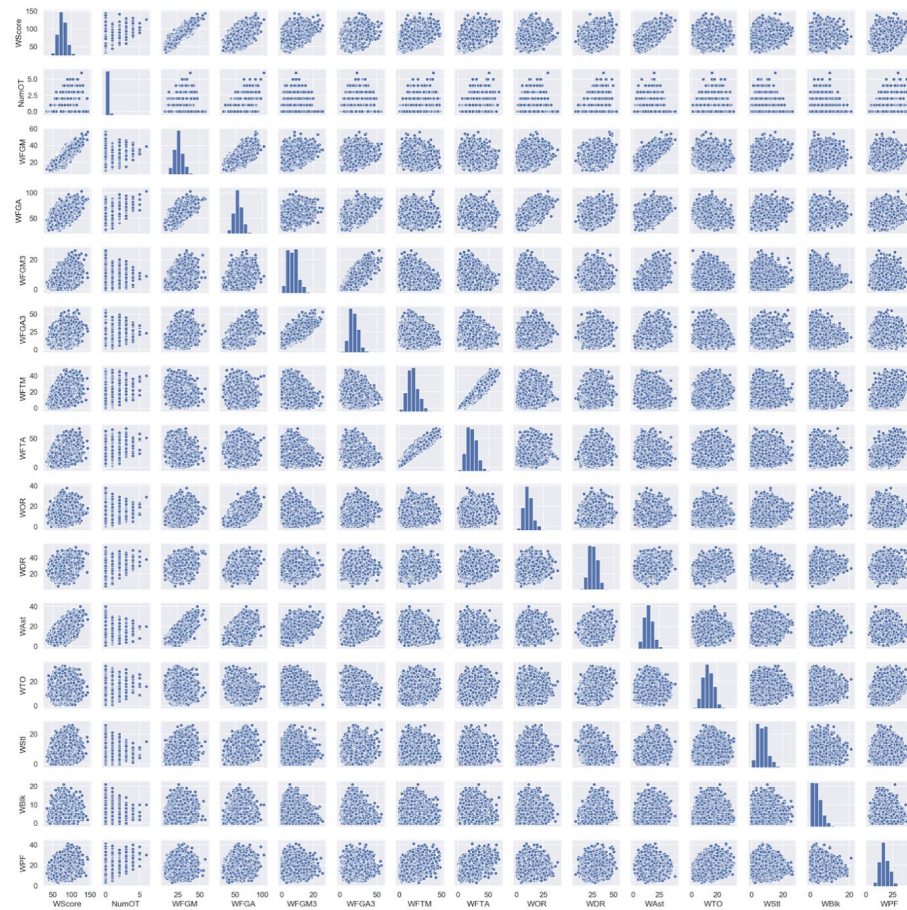
be seen from the plot that the relationship between WFGA3 and WPF was not statistically significant.

Pair plot matrix

A pairs plot allows us to see both the distribution of single variables and relationships between pairs of variables. A visual scan through of the pairs plot shows variables that seem highly correlated. The relationship between (WScore and WFGM), and (WFTA and WFTM) shows what seems to be a strong linear correlation between these variables. This positive correlation was also confirmed by the heatmap.

```
In [27]: sns.pairplot(capstone_drop,height=1.5)
```

```
Out[27]: <seaborn.axisgrid.PairGrid at 0x1a28349fd0>
```



3. What is the 95% confidence interval for the difference between the standard deviations of the winning scores for the two top performing teams? This analysis will be done with a bootstrap sampling analysis, calculating these differences over 10000 replicates. The teams being considered will be Gonzaga and Kansas.

A bootstrap method will be used to compare the means of winning scores of two teams over multiple seasons. The two chosen are Gonzaga and Kansas. These two teams were chosen

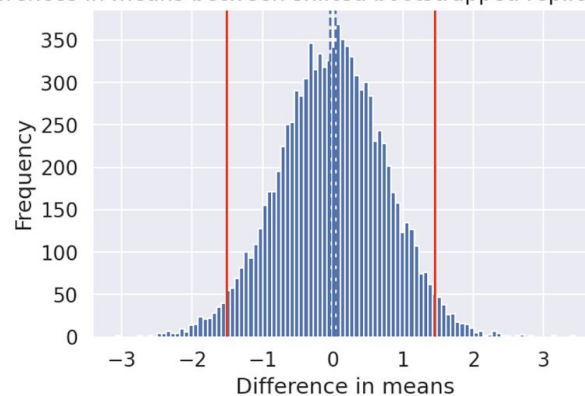
because they have the highest winning percentages i.e. winning Percentage is defined as (matches won/total matches played *100) . We will proceed by defining both the null and the alternative hypothesis for our calculations

H_0 : there is no difference in standard deviations in the winning scores between Kansas and Gonzaga

H_a : there is a difference in standard deviations in the winning scores between Kansas and Gonzaga

Bootstrap Replicates

↳ Distribution of differences in means between shifted bootstrapped replicates of Gonzaga and Kansas



The solid red vertical lines correspond to the 95% lower and upper confidence intervals of expected random differences in means of bootstrap replicates of Gonzaga and Kansas samples. The dashed blue vertical lines correspond to the observed difference in means between Gonzaga and Kansas samples.

Our Null and Alternative Hypotheses were as follows:

H_0 : there is no statistically significant difference in the means of the winning scores between Kansas and Gonzaga

H_a : there is a statistically significant difference in the means of the winning scores between Kansas and Gonzaga

The calculated p value of 0.16 is greater than the significance value of 0.05 i.e. $0.16 > 0.05$. Therefore we fail to reject the null hypothesis, and we cannot say that there is a statistically significant difference in the means of the winning scores between Kansas and Gonzaga.

Our bootstrap replicates with a 95% confidence interval indicate that the difference in means between the two groups have a 95% chance of lying within $[-1.4236463156146701, 1.4045422356365809]$. Our calculated difference in means is 0.71. Since the value is within the 95% confidence range, we therefore fail to reject the null hypothesis, and say there is no statistically significant difference in means between the Kansas and Gonzaga winning scores.

Summary and Next Steps

In summary, we found that:

- Home court does give an advantage in basketball games played in the NCAA
- Gonzaga's and Kansas have the highest winning percentage in the seasons where we have the ingame statistics. This is seasons 2005 - 2019, while Kansas and Duke have the highest winning percentage from seasons 1985-2019.
- The winning scores by the teams with the highest winning percentages , Kansas and Gonzaga, were not statistically significant.

There was also a strong correlation between in game statistics of field goal made by the winning team and the assist performed by the winning team, while there was little correlation between the scores made by the winning team and turnover conceded by the winning team

The next steps will cover the development of classification machine learning models to estimate the probability of a home team winning a game, which will be used to understand how game metrics affect the teams winning or losing games.