

Final Report

Tim Vigers

03 December 2019

Introduction

Basketball has come a long way since James Naismith

The Data

Team passing data were manually downloaded from <https://stats.nba.com/teams/passing/> and concatenated into a “long” dataset. These data were relatively well-organized to begin with and required minimal cleaning, but unfortunately only go back as far as 2013. Traditional statistics going back to the beginning of the NBA and ABA were downloaded using an HTML scraping tool developed for this project (see Appendix for code). These data were also relatively clean, but teams which moved or changed names were assigned a unique three letter code corresponding to their current location (e.g. observations from the New Orleans Jazz were given the code “UTA” in order to group them with the rest of the Utah Jazz data). Also, seasons were designated using the numeric year of the first game of the season, (e.g. 2018 for the 2018-2019 season) in order to treat time as a continuous variable. There were no missing or excluded observations in these data, and counting statistics such as points, turnovers, etc. were converted to per-game measures in order to account for shortened seasons in 1998 and 2011. For these analyses I considered only data from after the ABA and NBA merger in 1976.

Passing

Mixed Model Selection

Prior to modeling the number of passes over time, I created a spaghetti plot of passes over time with a line for each team (see Figure A1). There did not appear to be much of an overall trend. The total number of passes in a season appears to follow a normal distribution (Figure A2), so this outcome was modeled using a linear mixed model.

In order to test for a fixed effect of season on total number of passes made, I compared four linear mixed models. In the following models, i indexes team and j indexes season.

Model 1: Random Intercept Only

$$Y_{ij} = \beta_0 + \beta_1 x_j + b_{0i} + \epsilon_{ij}$$

$$b_{0i} \sim N(0, \sigma_{Team}^2) \text{ and } \epsilon_{ij} \sim N(0, \sigma_\epsilon^2)$$

Model 2: Random Intercept and AR(1) Structure for Repeated Measures

$$Y_{ij} = \beta_0 + \beta_1 x_j + b_{0i} + \epsilon_{ij}$$

$$b_{0i} \sim N(0, \sigma_{Team}^2) \text{ and } \epsilon_{ij} \sim N(0, R_i)$$

$$R_i = \sigma_\epsilon^2 \begin{bmatrix} 1 & \phi & \phi^2 & \phi^3 & \dots \\ \phi & 1 & \phi & \phi^2 & \\ \phi^2 & \phi & 1 & \phi & \\ \phi^3 & \phi^2 & \phi & 1 & \\ \vdots & & & & \ddots \end{bmatrix}$$

Models 3 & 4: Random Slope for Season

The last two models are the same as models 1 and 2, but with the addition of a random slope for season, so the random effects were

$$b_{0i} + b_{1j}x_j$$

with

$$b_{0i} \sim N(0, \sigma_{T_{eam}}^2) \text{ and } b_{1j} \sim N(0, \sigma_{Season}^2)$$

The model with random intercept and random slope did not converge without the AR(1) structure for repeated measures, and the model with random intercept and AR(1) structure was the best by the Akaike information criterion (AIC) (Table A1).

Using loess smoothing to plot total number of passes made suggested a potential cubic trend in the data. So once the final model was selected, I also tested the polynomial effects of season, up to a quadratic term:

$$Y_{ij} = \beta_0 + \beta_1x_j + \beta_2x_j^2 + \beta_3x_j^3 + \beta_4x_j^4 + b_{0i} + \epsilon_{ij}$$

$$b_{0i} \sim N(0, \sigma_{T_{eam}}^2) \text{ and } \epsilon_{ij} \sim N(0, R_i)$$

Piecewise Model

In addition to a linear mixed model, I also tried a linear spline model with a knot at 2015, including random intercept and AR(1) structure for repeated measures:

$$Y_{ij} = \beta_0 + \beta_1x_j + \beta_2\max(x_j - 2015, 0) + b_{0i} + \epsilon_{ij}$$

$$b_{0i} \sim N(0, \sigma_{T_{eam}}^2) \text{ and } \epsilon_{ij} \sim N(0, R_i)$$

Results

Table 1: The Effect of Time on Total Passes Made

	Value	Std.Error	DF	t-value	p-value
(Intercept)	24349.989	235.835	146	103.250	<1e-04
Season	-40.362	2004.510	146	-0.020	0.984
Season^2	-1941.549	1404.499	146	-1.382	0.169
Season^3	360.020	1088.741	146	0.331	0.741
Season^4	-465.829	925.730	146	-0.503	0.616

According to the linear mixed model, passing has not changed significantly since 2013.

Table 2: Change in Total Passes Made After the 2015 Season

	Value	Std.Error	DF	t-value	p-value
(Intercept)	164836.019	213588.290	148	0.772	0.441
Season	-69.854	106.029	148	-0.659	0.511
Change in Slope	0.181	0.154	148	1.178	0.241

Passing appears to increase slightly after 2015 according to the linear spline model, but the change in slope is not statistically significant ($p = 0.24$).

Assists

Model Selection

Winning percentage appears to be reasonably normally distributed (Figure A3), so I used normal theory linear mixed models to determine whether increasing assists results in more wins. Model selection for this question followed a similar process to the passing question. I compared models with random intercept for team and random intercept for team and random slope for season, both with and without an AR(1) structure for repeated measures. However, in these models the outcome is regular season win percentage and the fixed effects are average team age (“Age”); average team height (“Ht.”); average team weight (“Wt.”); team field goal percentage (“FG%”); and assists (“APG”), steals (“SPG”), blocks (“BPG”), and turnovers (“TPG”) per game. Once again, the model with random intercept for team and AR(1) structure for repeated measures was the best by AIC (Table A2).

However, during model selection, I realized that there was a significant positive association between assists per game and winning percentage, but that this effect goes away when adjusting for field goal percentage (Table A3). So, I conducted a mediation analysis (see Appendix for code) to try and determine whether field goal percentage mediates the effect of assists on winning [1]. The “mediation” package in R requires models without the AR(1) structure for repeated measures, so the mediation analysis was conducted using only a random intercept for team.

Results

Table 3: The Effect of Assists per Game on Winning Percentage

Table 3: Fixed Effects

	Value	Std.Error	DF	t-value	p-value
(Intercept)	-210.826	46.224	1111	-4.561	<1e-04
APG	1.924	0.156	1111	12.316	<1e-04
Age	3.553	0.222	1111	15.971	<1e-04
Ht.	0.643	0.594	1111	1.083	0.279
Wt.	0.320	0.066	1111	4.865	<1e-04
SPG	2.574	0.362	1111	7.106	<1e-04
BPG	3.668	0.367	1111	9.984	<1e-04
TPG	-2.393	0.229	1111	-10.467	<1e-04

Without adjusting for FG%, on average increasing assists by 5 per game can lead to a statistically significant ($p = <1e-04$) increase in winning percentage of 9.62 points on the season (or about 7.9 games). After adjustment for FG%, this effect is no longer significant (Table A3).

Mediation Analysis

There is a significant positive association between APG and FG%, FG% and winning, and between APG and winning (Figure A4). So, FG% is mediating a significant proportion of the effect of assists on winning ($p < 0.0001$, see code in Appendix for detailed results).

Discussion

Passing hasn't changed, but assists help through increasing FG%.

References

1. Tingley D, Yamamoto T, Hirose K, Keele L, Imai K. mediation: R package for causal mediation analysis. UCLA Stat Stat Assoc. August 2014. <https://dspace.mit.edu/handle/1721.1/91154>. Accessed December 2, 2019.

Appendix

Figure A1: Total Passes by Season

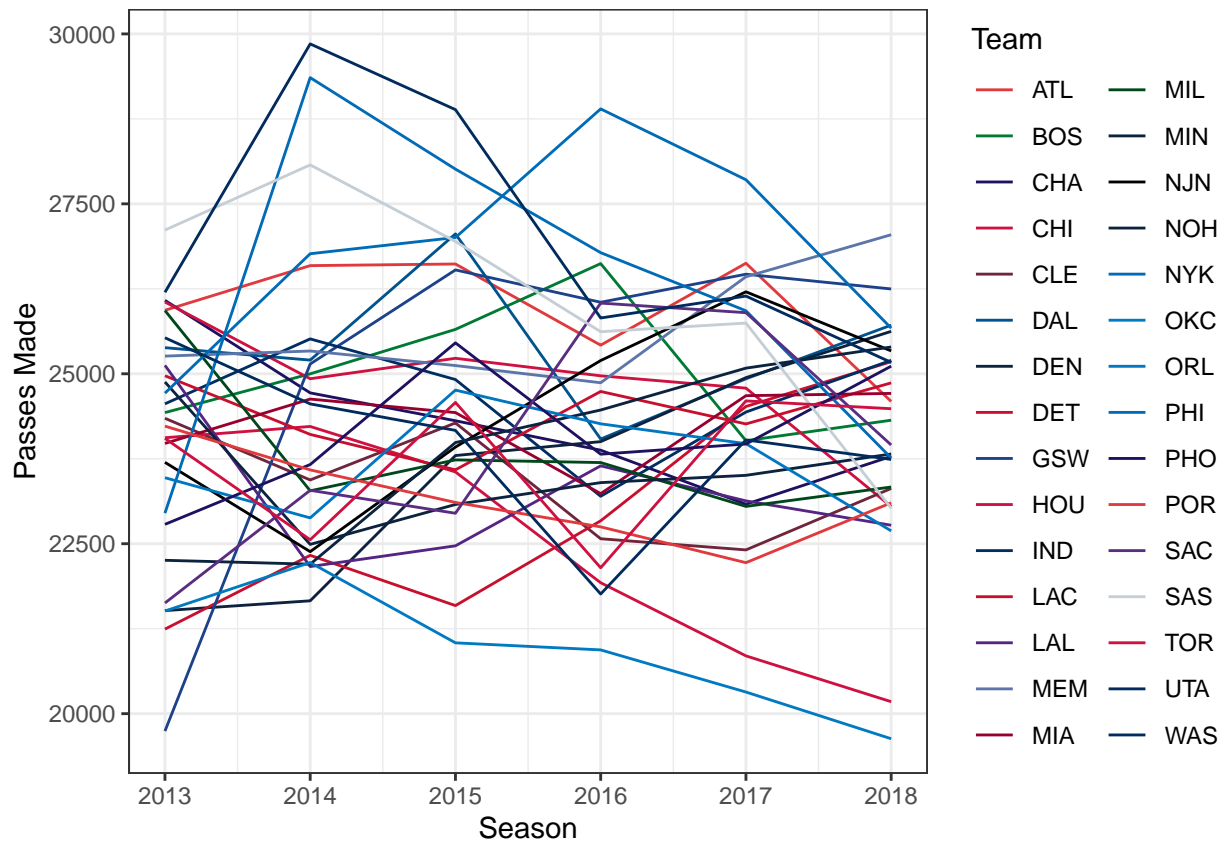


Figure A2: Distribution of Total Passes

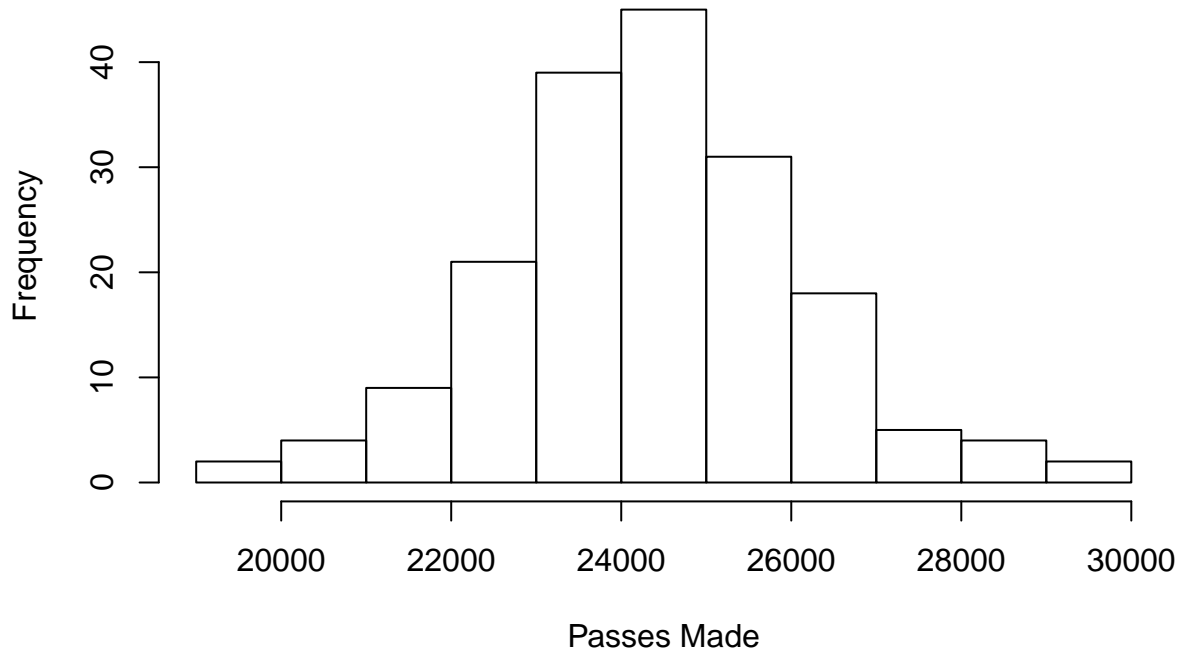


Table A1: AIC of Passes Made Models

All models fit using ML estimation.

	df	AIC
RI Only	4	3159.781
RI and AR(1)	5	3120.225
RI, RS, and AR(1)	7	3124.225

Figure A3: Distribution of Win Percentage

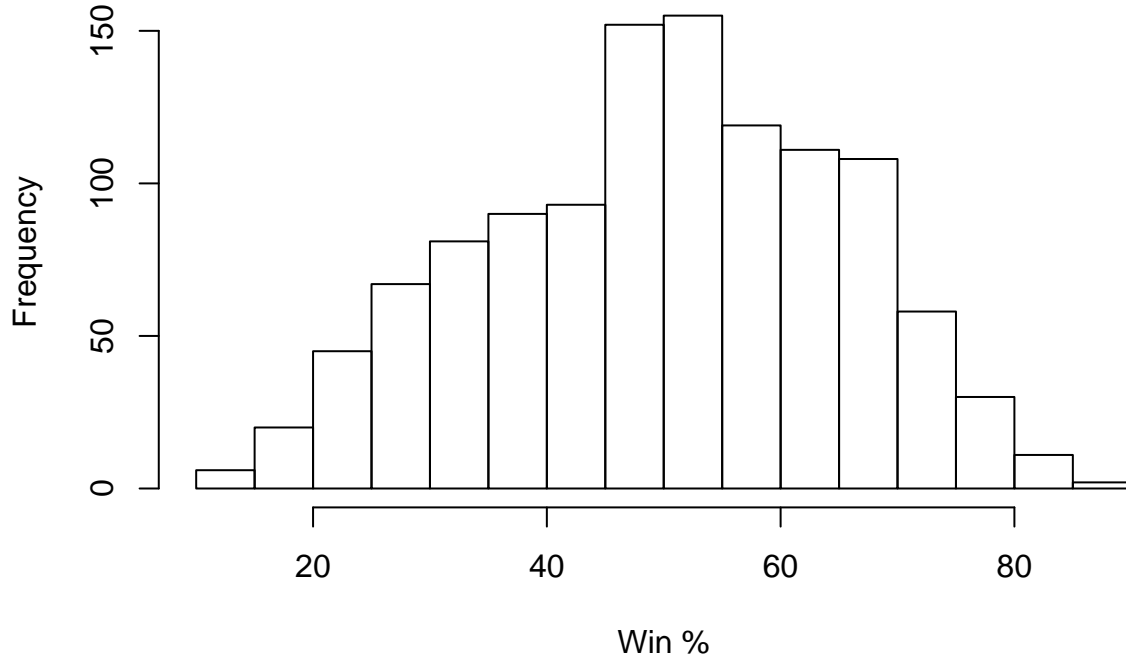


Table A2: AIC of Win Percentage Models

All models fit using ML estimation.

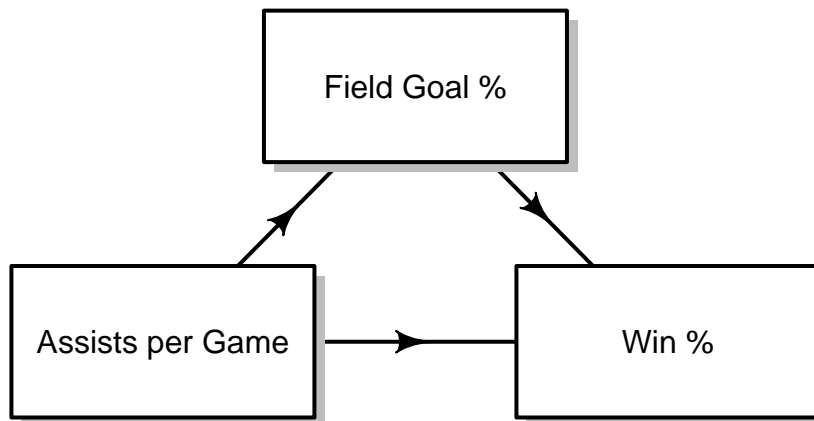
	df	AIC
RI Only	4	9336.655
RI and RS	6	9340.606
RI and AR(1)	5	8813.803
RI, RS, and AR(1)	7	8817.803

Table A3: Effect of Assists on Win Percentage, Adjusted for FG%

Table 6: Fixed Effects

	Value	Std.Error	DF	t-value	p-value
(Intercept)	-291.625	40.631	1110	-7.177	<1e-04
APG	-0.247	0.179	1110	-1.379	0.168
Age	3.257	0.195	1110	16.672	<1e-04
Ht.	-0.147	0.521	1110	-0.282	0.778
Wt.	0.424	0.058	1110	7.333	<1e-04
SPG	3.178	0.319	1110	9.970	<1e-04
BPG	3.478	0.322	1110	10.816	<1e-04
TPG	-2.994	0.203	1110	-14.782	<1e-04
FG%	398.631	21.277	1110	18.735	<1e-04

Figure A4: Mediation Diagram



Code

HTML Scraping Tool

```

library(rvest)
library(tidyverse)
teams <- c("ATL","BOS","NJN","CHA","CHI","CLE","DAL","DEN","DET","GSW","HOU",
          "IND","LAC","LAL","MEM","MIA","MIL","MIN","NOH","NYK","OKC","ORL",
          "PHI","PHO","POR","SAC","SAS","TOR","UTA","WAS")

# Scrape each team page
all_seasons <- data.frame()
for (team in teams) {
  url <- paste0("https://www.basketball-reference.com/teams/",team,"/stats_basic_totals.html")
  table <- url %>%
    read_html() %>%
    html_nodes("table") %>%
    html_table()
  df <- as.data.frame(table[[1]])
  df <- df[colnames(df) != ""] %>%
    filter(Season != "Season", Season != "2019-20")
  df[df == ""] <- NA
  df <- as.data.frame(lapply(df, as.character))
  colnames(df) <- c("Season","Lg","Tm","W","L","Finish","Age","Ht.","Wt.",
                    "G","MP","FG","FGA","FG%", "3P","3PA",
                    "3P%","2P","2PA","2P%","FT","FTA","FT%","ORB","DRB","TRB",
                    "AST","STL","BLK","TOV","PF","PTS")
  df$Team <- team
  all_seasons <- rbind.data.frame(all_seasons,df)
}

```

Mediation

```

# Mediation with FGP as mediator
mod.y <- lmer(w_perc ~ AST_game + Age + Ht. + Wt. + STL_game + BLK_game + TOV_game +
             FGP + (1|Team), data = post_merger)
mod.m <- lmer(FGP ~ AST_game + Age + Ht. + Wt. + STL_game + BLK_game + TOV_game +
             (1|Team), data = post_merger)
med_fgp <- mediate(mod.m, mod.y, treat = "AST_game", mediator = "FGP")

```



```
# Mediation with AST as mediator
mod.m <- lmer(AST_game ~ FGP + Age + Ht. + Wt. + STL_game + BLK_game + TOV_game +
              (1|Team),data = post_merger)
med_ast <- mediate(mod.m,mod.y,treat = "FGP",mediator = "AST_game")
# Mediation summary
summary(med_fgp)
```

```
##
## Causal Mediation Analysis
##
## Quasi-Bayesian Confidence Intervals
##
## Mediator Groups: Team
##
## Outcome Groups: Team
##
## Output Based on Overall Averages Across Groups
##
##           Estimate 95% CI Lower 95% CI Upper p-value
## ACME             2.168      1.921      2.46 <2e-16 ***
## ADE             -0.246     -0.594      0.07   0.15
## Total Effect      1.922      1.608      2.23 <2e-16 ***
## Prop. Mediated    1.119      0.964      1.35 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Sample Size Used: 1148
##
##
## Simulations: 1000
```

```
summary(med_ast)
```

```
##
## Causal Mediation Analysis
##
## Quasi-Bayesian Confidence Intervals
##
## Mediator Groups: Team
##
## Outcome Groups: Team
##
## Output Based on Overall Averages Across Groups
##
##           Estimate 95% CI Lower 95% CI Upper p-value
## ACME          -19.2346    -46.3442      6.50   0.16
## ADE           399.2345    358.0665    440.92 <2e-16 ***
## Total Effect  379.9999    349.3348    412.06 <2e-16 ***
## Prop. Mediated -0.0487     -0.1237     0.02   0.16
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Sample Size Used: 1148
##
```

```
##  
## Simulations: 1000
```