

Preparation of a NBA dataset with SAS

Chenpeng Guan, William Schlemmel, Jonah Somers, Krti Tallam

Introduction

The data files we used are CSV files from basketball-reference.com, a renowned basketball data website. They are relatively clean, with some missing values and repetitions which gave the opportunity to test data cleaning and organizing skills. There are 30 variables and approximately 3000 observations. Our goal was to combine per-game, per-player statistics from multiple seasons (2012-2016) and upon doing so, utilize the data to identify some interesting trends and notable players throughout the past half-decade. Preparation of these data sets also holds great business value for NBA franchises.

The important variables in this data set include identifying variables (name, position, team) as well as the averages and percentages of offensive metrics (field goals, three pointers, two pointers, free throws, points per game) and various other per-game statistics (rebounds, assists, blocks). We use these to identify notable players, like the point leaders (total and three's only), the defensive rebound leaders, and players who averaged a 20-5-5 combination (points, rebounds, and assists).

Datasets:

https://www.basketball-reference.com/leagues/NBA_2016_per_game.html

https://www.basketball-reference.com/leagues/NBA_2015_per_game.html

https://www.basketball-reference.com/leagues/NBA_2014_per_game.html

https://www.basketball-reference.com/leagues/NBA_2013_per_game.html

https://www.basketball-reference.com/leagues/NBA_2012_per_game.html

Note: Must download as CSV files, and there may or may not be extra information at the end of the files. This must be deleted before running our code if it appears.

Methods

Each CSV file has 30 columns and between 500 and 600 rows. This adds up to a combined 3,008 observations when concatenated. The files contain per-game data on each NBA player, as well as general information on the player such as name, position, and team. All variables are numeric except those three. The numeric variables read in include: age, games, games started, minutes played, a count, attempts, and percentage variable for each of field goals, three pointers, two pointers, and free throws, effective field goal percentage, rebounds (offensive and defensive), assists, steals, blocks, turnovers, personal fouls, and points per game. Multiple

entries for a single player are possible, due to trades and a player's presence in more than one of the seasons included.

Each file was read in with using several options on the infile statement. While combining the five data files, a "Year" variable was added to specify observational seasons. The extra text following each player's name (e.g. Quincy Acy\acyqu01) was also deleted using a simple scan function. These additions as well as variables for total points and total three pointers were created with PROC SQL. We combined the five datasets within the same procedure with union operator..The table was modified with appropriate formats and labels.

Some data exploration techniques then revealed several issues to address. A PROC FREQ step with nlevels allowed us to notice any missing values and any incorrect entries. Three problems arose: all percentage variables had missing values, the number of teams was 33 (there are 30 in the league), and some players had more than 82 games played in a season (normal length of regular season).

The first problem arose due to the fact that if a player never attempted a free-throw, 3-point shot, or 2-point shot, their corresponding percentage variable would be missing due to a divide-by-zero error when calculating the percentage. We fixed this by replacing any missing values with 0, and recalculating percentages based on their equations in the data set glossary (<https://www.basketball-reference.com/about/glossary.html>).

There were 33 levels for team, even though 30 teams exist in the league. This is due to the fact that the New Orleans Hornets were renamed the Pelicans, the Charlotte Bobcats became the Charlotte Hornets, and the TOT level represents total statistics for those players who were traded and played for multiple teams during one season. Consequently, no fixes were needed. A similar situation arose for position; 13 levels were found yet there are only 5 positions. This was due to the fact that the data set included hyphenated combinations of positions.

For the games issue, we found two players (Josh Smith and Ramon Sessions) that had 83 games played in a year. These were rare and caused by complications from trading. Therefore, they did indeed play 83 games.

A data step addressed these issues and included wholesale data fixes, which were needed because of the problems in the calculations of percentages. The new and cleaned data set was named nba_clean. Another PROC FREQ confirmed our fixes.

We used SQL to count all the unique teams a player has been on. The first 10 observations were printed in descending order.

Next we accumulated number of points and points from three-point shots in a data step. The leading scorer on each team was subsequently then also found using PROC MEANS, SORT, and a data step. Finally we displayed with a SQL several tables, including, the top 10 three points leaders, the top 10 points leaders, the top 10 field goal shooting percentages, the top 10 number of defensive rebounds, and the players who averaged the 20-5-5 combination.

Results

After reading in the data, a PROC CONTENTS (seen below) allowed us to view the dimensions of our new table as well as confirm our assignments of labels and formats.

The CONTENTS Procedure

Data Set Name	NBA.NBA	Observations	3008
Member Type	DATA	Variables	32
Engine	V9	Indexes	0
Created	12/20/2017 19:18:19	Observation Length	296
Last Modified	12/20/2017 19:18:19	Deleted Observations	0
Protection		Compressed	NO
Data Set Type		Sorted	NO
Label			
Data Representation	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
Encoding	utf-8 Unicode (UTF-8)		

Alphabetic List of Variables and Attributes

#	Variable	Type	Len	Format	Label
23	AST	Num	8		Assists per Game
2	Age	Num	8		
25	BLK	Num	8		Blocks per Game
21	DRB	Num	8		Defensive Rebounds per Game
7	FG	Num	8		Field Goals Made
8	FGA	Num	8		Field Goal Attempts per Game
9	FG_Pct	Num	8	PERCENT8.3	Field Goal Percentage
17	FT	Num	8		Free Throws Made
18	FTA	Num	8		Free Throw Attempts
19	FT_Pct	Num	8	PERCENT8.3	Free Throw Percentage
4	G	Num	8		Games Played
5	GS	Num	8		Games Started
6	MP	Num	8		Average Minutes Played per Game
20	ORB	Num	8		Offensive Rebounds per Game
27	PF	Num	8		Personal Fouls per Game
30	Player_Name	Char	40	\$40.	Player
1	Position	Char	8		
28	Pts_G	Num	8		Points per Game
32	Pts_Tot	Num	8		Total Points
24	STL	Num	8		Steals per Game
26	TOV	Num	8		Turnovers per Game
22	TRB	Num	8		Total Rebounds per Game
3	Team	Char	10		
10	ThreePoints	Num	8		3-Pointers Made
11	ThreePoints_A	Num	8		3-Pointers Attempts
12	ThreePoints_Pct	Num	8	PERCENT8.3	3-Pointers Percentage
31	ThreePoints_Tot	Num	8		Total 3-Pointers
13	TwoPoints	Num	8		2-Pointers Made
14	TwoPoints_A	Num	8		2-Pointers Attempted
15	TwoPoints_Pct	Num	8	PERCENT8.3	2-Pointers Percentage
29	Year	Num	8	8.	
16	eFG_Pct	Num	8	PERCENT8.3	Effective Field Goal Percentage

Running a PROC FREQ step with nlevels allowed us to notice any missing values and any incorrect entries, all of which belonged to percent variables. A frequency procedure was used rather than a PROC MEANS since it also revealed an unexpected number of levels for position and team.

Overall Missing Levels				
The FREQ Procedure				
Number of Variable Levels				
Variable	Label	Levels	Missing Levels	Nonmissing Levels
Position		13	0	13
Age		22	0	22
Team		33	0	33
G	Games Played	83	0	83
GS	Games Started	83	0	83
MP	Average Minutes Played per Game	371	0	371
FG	Field Goals Made	104	0	104
FGA	Field Goal Attempts per Game	206	0	206
FG_Pct	Field Goal Percentage	383	1	382
ThreePoints	3-Pointers Made	39	0	39
ThreePoints_A	3-Pointers Attempts	88	0	88
ThreePoints_Pct	3-Pointers Percentage	314	1	313
TwoPoints	2-Pointers Made	92	0	92
TwoPoints_A	2-Pointers Attempted	174	0	174
TwoPoints_Pct	2-Pointers Percentage	364	1	363
eFG_Pct	Effective Field Goal Percentage	390	1	389
FT	Free Throws Made	82	0	82
FTA	Free Throw Attempts	95	0	95
FT_Pct	Free Throw Percentage	483	1	482
ORB	Offensive Rebounds per Game	51	0	51
DRB	Defensive Rebounds per Game	99	0	99
TRB	Total Rebounds per Game	134	0	134
AST	Assists per Game	99	0	99
STL	Steals per Game	26	0	26
BLK	Blocks per Game	33	0	33
TOV	Turnovers per Game	48	0	48
PF	Personal Fouls per Game	43	0	43
Pts_G	Points per Game	264	0	264
Year		5	0	5
Player_Name	Player	816	0	816
ThreePoints_Tot	Total 3-Pointers	208	0	208
Pts_Tot	Total Points	1065	0	1065

The 3 extra levels in the team variable were identified in the following table as TOT and the two renamed franchises (NOH/NOP and CHA/CHO). These 4 teams were combined into 2 teams respectively.

Unusual Levels in Teams

The FREQ Procedure

Team	Frequency	Percent	Cumulative Frequency	Cumulative Percent
ATL	89	2.96	89	2.96
BOS	92	3.06	181	6.02
BRK	93	3.09	274	9.11
CHA	37	1.23	311	10.34
CHI	83	2.76	394	13.10
CHO	53	1.76	447	14.86
CLE	99	3.29	546	18.15
DAL	97	3.22	643	21.38
DEN	89	2.96	732	24.34
DET	83	2.76	815	27.09
GSW	82	2.73	897	29.82
HOU	94	3.13	991	32.95
IND	81	2.69	1072	35.64
LAC	88	2.93	1160	38.56
LAL	86	2.86	1246	41.42
MEM	102	3.39	1348	44.81
MIA	91	3.03	1439	47.84
MIL	93	3.09	1532	50.93
MIN	93	3.09	1625	54.02
NOH	19	0.63	1644	54.65
NOP	89	2.96	1733	57.61
NYK	89	2.96	1822	60.57
OKC	91	3.03	1913	63.60
ORL	85	2.83	1998	66.42
PHI	104	3.46	2102	69.88
PHO	95	3.16	2197	73.04
POR	79	2.63	2276	75.66
SAC	93	3.09	2369	78.76
SAS	85	2.83	2454	81.58
TOR	86	2.86	2540	84.44
TOT	292	9.71	2832	94.15
UTA	85	2.83	2917	96.97
WAS	91	3.03	3008	100.00

Our suspicions on the number of levels for position were relieved once we saw the data included hyphenated combinations of positions (accounting for the 13 possible positions), as seen with the below PROC FREQ output.

Unusual Levels in Position				
The FREQ Procedure				
Position	Frequency	Percent	Cumulative Frequency	Cumulative Percent
C	541	17.99	541	17.99
C-PF	1	0.03	542	18.02
PF	618	20.55	1160	38.56
PF-C	3	0.10	1163	38.66
PF-SF	3	0.10	1166	38.76
PG	603	20.05	1769	58.81
PG-SG	5	0.17	1774	58.98
SF	586	19.48	2360	78.46
SF-PF	3	0.10	2363	78.56
SF-SG	1	0.03	2364	78.59
SG	631	20.98	2995	99.57
SG-PG	8	0.27	3003	99.83
SG-SF	5	0.17	3008	100.00

The two players with unusual game counts are outputted below, via PROC MEANS.

Unusual Total Games						
The MEANS Procedure						
Analysis Variable : G Games Played						
Player	N Obs	N	Mean	Std Dev	Minimum	Maximum
Josh Smith	1	1	83.0000000	.	83.0000000	83.0000000
Ramon Sessions	1	1	83.0000000	.	83.0000000	83.0000000

The following two tables were used to confirm our fixes.

Confirms Miscellaneous Fixes

The FREQ Procedure

Number of Variable Levels		
Variable	Label	Levels
Position		13
Age		22
Team		31
G	Games Played	83
GS	Games Started	83
MP	Average Minutes Played per Game	371
FG	Field Goals Made	104
FGA	Field Goal Attempts per Game	206
FG_Pct	Field Goal Percentage	915
ThreePoints	3-Pointers Made	39
ThreePoints_A	3-Pointers Attempts	88
ThreePoints_Pct	3-Pointers Percentage	274
TwoPoints	2-Pointers Made	92
TwoPoints_A	2-Pointers Attempted	174
TwoPoints_Pct	2-Pointers Percentage	696
eFG_Pct	Effective Field Goal Percentage	389
FT	Free Throws Made	82
FTA	Free Throw Attempts	95
FT_Pct	Free Throw Percentage	345
ORB	Offensive Rebounds per Game	51
DRB	Defensive Rebounds per Game	99
TRB	Total Rebounds per Game	134
AST	Assists per Game	99
STL	Steals per Game	26
BLK	Blocks per Game	33
TOV	Turnovers per Game	48
PF	Personal Fouls per Game	43
Pts_G	Points per Game	264
Year		5
Player_Name	Player	816
ThreePoints_Tot	Total 3-Pointers	208
Pts_Tot	Total Points	1065

Confirms Fixes for Renamed Franchises

The FREQ Procedure

Team	Frequency	Percent	Cumulative Frequency	Cumulative Percent
ATL	89	2.96	89	2.96
BOS	92	3.06	181	6.02
BRK	93	3.09	274	9.11
CHI	83	2.76	357	11.87
CHO/CHA	90	2.99	447	14.86
CLE	99	3.29	546	18.15
DAL	97	3.22	643	21.38
DEN	89	2.96	732	24.34
DET	83	2.76	815	27.09
GSW	82	2.73	897	29.82
HOU	94	3.13	991	32.95
IND	81	2.69	1072	35.64
LAC	88	2.93	1160	38.56
LAL	86	2.86	1246	41.42
MEM	102	3.39	1348	44.81
MIA	91	3.03	1439	47.84
MIL	93	3.09	1532	50.93
MIN	93	3.09	1625	54.02
NOH/NOP	108	3.59	1733	57.61
NYK	89	2.96	1822	60.57
OKC	91	3.03	1913	63.60
ORL	85	2.83	1998	66.42
PHI	104	3.46	2102	69.88
PHO	95	3.16	2197	73.04
POR	79	2.63	2276	75.66
SAC	93	3.09	2369	78.76
SAS	85	2.83	2454	81.58
TOR	86	2.86	2540	84.44
TOT	292	9.71	2832	94.15
UTA	85	2.83	2917	96.97
WAS	91	3.03	3008	100.00

We also used SQL to show us how many distinct teams each player has been on. We sorted the table by descending order and only showed the first 10 observations below. These are interesting to know as they can give insight into why a player is traded (demand for talent, bad deals, etc.)

Unique Teams per Player (Top 10 Obs.)

Player	Distinct_Teams
D.J. Augustin	7
Ish Smith	7
Lance Stephenson	6
Wayne Ellington	6
Gary Neal	6
Ersan Ilyasova	6
Jose Calderon	6
Beno Udrih	6
Jared Cunningham	6
Toney Douglas	6

Next, we found the points leader per team and printed the resulting data set below.

**NBA Data
Points Leader Per Team**

Team	Player	Number of Years	Total Points
HOU	James Harden	5	10,828
GSW	Stephen Curry	5	9,939
OKC	Russell Westbrook	5	9,228
POR	Damian Lillard	5	8,885
TOR	DeMar DeRozan	5	8,336
NYK	Carmelo Anthony	5	8,229
SAC	DeMarcus Cousins	5	7,595
NOH/NOP	Anthony Davis	5	7,499
CHO/CHA	Kemba Walker	5	7,342
CLE	Kyrie Irving	5	7,286
WAS	John Wall	5	7,214
UTA	Gordon Hayward	5	6,904
LAC	Blake Griffin	5	6,902
IND	Paul George	5	6,813
DAL	Dirk Nowitzki	5	6,125
CHI	Jimmy Butler	5	6,099
BRK	Brook Lopez	5	6,068
SAS	Kawhi Leonard	5	6,004
MEM	Mike Conley	5	5,802
ORL	Nikola Vucevic	5	5,524
MIA	Dwyane Wade	4	5,228
ATL	Paul Millsap	4	5,178
DET	Andre Drummond	5	5,114
MIN	Andrew Wiggins	3	4,998
MIL	Giannis Antetokounmpo	4	4,737
BOS	Avery Bradley	5	4,476
DEN	Kenneth Faried	5	4,385
PHO	Eric Bledsoe	4	4,163
LAL	Kobe Bryant	4	4,155
PHI	Robert Covington	3	2,667

Most of these players in the table are all-star players or are renown in the NBA community, but this table can explain a lot about a team's history over the last five years. Teams like the Philadelphia 76ers or the New Orleans Hornets have scoring leaders that have scored very few points (probably wouldn't even be a top 3 scorer on some of the top teams), and this makes sense since these are teams that are currently rebuilding their rosters or are without a superstar player. Top team-leaders for scoring like James Harden (Houston Rockets) and Stephen Curry (Golden State Warriors) are superstars on teams that have been contending for the championship for each of the past five years.

In a following PROC SQL step, tables for top 10 three points leaders, the top 10 points leaders, the top 10 field goal shooting percentages, the top 10 number of defensive rebounds, and the players who averaged the 20-5-5 combination. They are below.

Top 10 Three Points Leaders

Player	Three Point Total
Stephen Curry	1,545
Klay Thompson	1,224
Kyle Korver	1,078
James Harden	1,062
Damian Lillard	1,058
J.J. Redick	1,010
J.R. Smith	959
Isaiah Thomas	914
Wesley Matthews	907
Kyle Lowry	833

Top 10 Points Leaders

Player	Total Points
James Harden	10,828
DeMarcus Cousins	9,954
Stephen Curry	9,939
LeBron James	9,747
Russell Westbrook	9,228
Kevin Durant	9,140
Damian Lillard	8,885
Rudy Gay	8,846
Isaiah Thomas	8,773
DeMar DeRozan	8,336

These two tables show an interesting difference. Of the top ten 3-point leaders, only Stephen Curry, James Harden, Damian Lillard, and Isaiah Thomas are top ten in overall points scored in the last five years. With a lot of relevant articles and media claiming the NBA has become very focused on 3-point scoring, a many players are still scoring lots of points without relying on the 3-point shot. DeMarcus Cousins, LeBron James, and Russell Westbrook are top five in points scored, but aren't even in the top ten for 3-pointers made. While this shows that 3-point shots aren't everything, the fact that Stephen Curry and Klay Thompson are both on the Golden State Warriors, who have won two championships in the past five years, and are the best three-point scorers in the NBA is a compelling counter-argument.

The following table displays the best FG percentages.

Top 10 Field Goal Shooting Percentages

Player	Position	Year	Effective Field Goal Percentage	Field Goal Attempts per Game
Edy Tavares	C	2016	75.00%	4.0
Brandan Wright	PF	2014	74.80%	5.0
DeAndre Jordan	C	2016	71.40%	7.1
DeAndre Jordan	C	2014	71.10%	6.5
DeAndre Jordan	C	2015	70.30%	6.6
Axel Toupane	SF	2016	68.80%	4.0
Brandan Wright	C	2013	67.70%	5.7
Troy Daniels	SG	2013	67.70%	6.2
DeAndre Jordan	C	2013	67.60%	6.3
Brandan Wright	PF	2015	67.30%	4.3

These results show that a majority of the most efficient scorers are centers and power forwards (only one guard makes this list, with the rest being relatively taller and larger players), which would make sense because these big men play close to the basket and often score from close range. DeAndre Jordan seems to be a stand out here as one of the most consistent efficient scorers.

Top 10 Number of Defensive Rebounds

Player	Year	Defensive Rebounds per Game
Earl Barron	2012	12.0
Kevin Love	2012	10.4
DeAndre Jordan	2015	10.3
Hassan Whiteside	2016	10.3
DeMarcus Cousins	2016	10.2
DeAndre Jordan	2014	10.1
DeAndre Jordan	2016	10.1
Andre Drummond	2015	9.9
Kevin Love	2013	9.6
DeAndre Jordan	2013	9.5

In this table, we pulled out the most defensive rebounds (when defense obtains possession of the ball after a missed shot by the offense) per game per player for the top 10

players in our list, for any year including and between 2012-2016. We sorted from greatest to least number of defensive rebounds, demonstrating values ranging between 12.0 and 9.5.

The table for the 20-5-5 combination averages is displayed next.

Player	Position	Year	Total Rebounds per Game	Points per Game	Assists per Game
Giannis Antetokounmpo	SF	2016	8.8	22.9	5.4
Kobe Bryant	SG	2012	5.6	27.3	6.0
Jimmy Butler	SF	2016	6.2	23.9	5.5
Kobe Bryant	SG	2014	5.7	22.3	5.6
Stephen Curry	PG	2015	5.4	30.1	6.7
Kevin Durant	SF	2013	7.4	32.0	5.5
James Harden	PG	2016	8.1	29.1	11.2
James Harden	SG	2015	6.1	29.0	7.5
Blake Griffin	PF	2014	7.6	21.9	5.3
James Harden	SG	2014	5.7	27.4	7.0
LeBron James	PF	2012	8.0	26.8	7.3
LeBron James	SF	2016	8.6	26.4	8.7
LeBron James	SF	2015	7.4	25.3	6.8
LeBron James	PF	2013	6.9	27.1	6.3
LeBron James	SF	2014	6.0	25.3	7.4
Russell Westbrook	PG	2012	5.2	23.2	7.4
Russell Westbrook	PG	2015	7.8	23.5	10.4
Russell Westbrook	PG	2016	10.7	31.6	10.4
Russell Westbrook	PG	2013	5.7	21.8	6.9
Russell Westbrook	PG	2014	7.3	28.1	8.6

These players are considered superstar-caliber players who are very versatile. Most of these players listed here are candidates for the Most Valuable Player (MVP) award or are former MVPs, which shows that 20-5-5 can be a good metric for identifying potential MVP candidates. LeBron James and Russell Westbrook are increasingly considered to be among the greatest players to ever play in the NBA, and it shows with their consistent performances throughout the years with both having averaged a 20-5-5 in all of the years the data set covers. Note that Russell Westbrook average over 10 rebounds, 10 assists, and 30 points per game during the 2016-2017 season, which from a pure statistical perspective is one of the most impressive performances for a player in NBA history.