# Steph Curry and the dried up 3s

Using data analysis for predictive modelling

# So what is data science, anyway?

When you hear the term data scientist, what do you think of? If you're like most people, you might think of something really complex, with statistical terms and programming languages no one can understand. You might think that you might need to go to university to learn how to do data science.
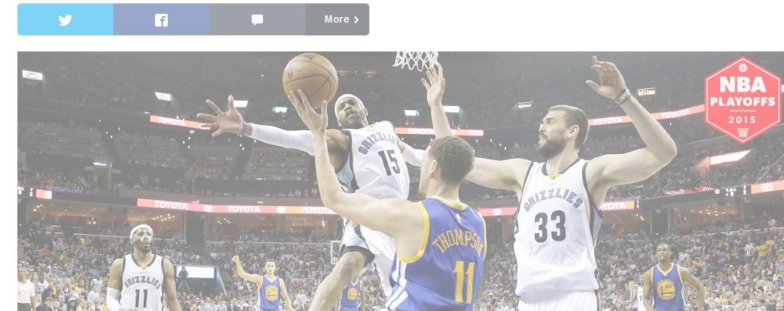
Data science is just a 21$^{st}$ century spin on maths that people have been doing for centuries. Big data, data science and analytics are just fancy ways of saying we use information available to gain insight and improve a business.

Whether it's a small Excel spreadsheet or a 100 million records in a database, the goal is always the same: to find value.

# Step 1: Have a question or something you're curious about

In the 2015 NBA Playoffs, Steph Curry and the Golden State Warriors were down 2 games to 1 against the Memphis Grizzlies, and Curry's 3-point shooting was down in the previous two games which the Warriors lost. Commentators were speculating; have the Grizzlies figured out Steph Curry? Can he bounce back and guide the Warriors to victory in the best of 7 series?



Grizzlies grind down MVP Stephen Curry, Warriors' offense in Game 3 win

# Framing the hypothesis

Let's now frame the hypothesis.

Curry's poor 3 point percentage was shooting at below 0.20 for these two matches.  But normally he shoots at 0.40. If he shoots over 0.40 then his team wins the game 89% of the time in 2014-15.

You can frame the hypothesis as a return to form or repeat poor performance

**Golden State Warriors will return to form and win the series because Steph Curry  should be able to average a 0.40 3 point percentage after his two poor shooting games against the Memphis Grizzlies in the 2014-15 NBA playoffs.**

**Golden State Warriors will lose and bow out of the series against the Memphis Grizzlies in the 2014-15 NBA playoffs because Steph Curry will average clearly less than a 0.40 3 point percentage.**

# Deciding on the data required

- Average 3 Point Shooting percentage for all matches would be helpful to develop a baseline
- Average 3 Point Shooting percentage for all matches after he has had a poor shooting performance
- "For all matches" - just that season, or how long do you use data, is it less relevant if it's for a previous year?
- Are there any other statistics that could be useful?  Is there an issue with just playing that opponent?

# Step 2: Gather data that exists for your area of interest

We can use easily available data from [basketball-reference.com](basketball-reference.com) in this situation.

I simply took Steph Curry's game log for the 2014-15 regular season and created a .CSV file ([uploaded here](uploaded here) if you want to download it). Here's what the data looked like:

# Process for gathering the data …

Navigate your web browser to basketball-reference.com

Type Stephen Curry in the search bar and click on the associated link:

# Process for gathering the data …

Scroll down and click on 2014-15 to display data from the appropriate season

| Season | Age | Tm | Lg | Pos | G | GS | MP | FG | FGA | FG% | 3P | 3PA |
|--------|-----|-----|-----|-----|-----|-----|-------|------|------|------|------|------|
| 2009-10 | 21 | GSW | NBA | PG | 80 | 77 | 2896 | 528 | 1143 | .462 | 166 | 380 |
| 2010-11 | 22 | GSW | NBA | PG | 74 | 74 | 2489 | 505 | 1053 | .480 | 151 | 342 |
| 2011-12 | 23 | GSW | NBA | PG | 26 | 23 | 732 | 145 | 296 | .490 | 55 | 121 |
| 2012-13 | 24 | GSW | NBA | PG | 78 | 78 | 2983 | 626 | 1388 | .451 | 272 | 600 |
| 2013-14 ★ | 25 | GSW | NBA | PG | 78 | 78 | 2846 | 652 | 1383 | .471 | 261 | 615 |
| 2014-15 ★ | 26 | GSW | NBA | PG | 80 | 80 | 2613 | 653 | 1341 | .487 | 286 | 646 |
| 2015-16 ★ | 27 | GSW | NBA | PG | 50 | 50 | 1692 | 498 | 981 | .508 | 245 | 540 |
| Career | | | NBA | | 466 | 460 | 16251 | 3607 | 7585 | .476 | 1436 | 3244 |

Totals · Glossary · SHARE · Embed · CSV · Export · PRE · LINK · ?

# Process for gathering the data ...

Scroll down so that you can see all of the games played by Curry in 2014-15

**2014-15 Regular Season**  Age is Years-Days · Glossary · SHARE · Embed · CSV · Export · PRE · LINK · ?

| Rk | G | Date | Age | Tm | | Opp | | GS | MP | FG | FGA | FG% | 3P | 3PA | 3P% | FT | FTA | FT% | ORB | DRB | TRB | AST | STL | BLK | TOV | PF | PTS | GmSc | +/- |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 2014-10-29 | 26-229 | GSW | @ | SAC | W (+18) | 1 | 36:32 | 7 | 17 | .412 | 2 | 9 | .222 | 8 | 9 | .889 | 1 | 9 | 10 | 5 | 6 | 0 | 4 | 5 | 24 | 21.4 | +20 |
| 2 | 2 | 2014-11-01 | 26-232 | GSW | | LAL | W (+23) | 1 | 34:06 | 10 | 19 | .526 | 3 | 8 | .375 | 8 | 8 | 1.000 | 0 | 5 | 5 | 10 | 3 | 0 | 2 | 3 | 31 | 30.0 | +30 |
| 3 | 3 | 2014-11-02 | 26-233 | GSW | @ | POR | W (+5) | 1 | 30:50 | 6 | 18 | .333 | 1 | 5 | .200 | 8 | 8 | 1.000 | 1 | 4 | 5 | 6 | 2 | 0 | 3 | 1 | 21 | 15.5 | +12 |
| 4 | 4 | 2014-11-05 | 26-236 | GSW | | LAC | W (+17) | 1 | 33:49 | 9 | 18 | .500 | 4 | 8 | .500 | 6 | 6 | 1.000 | 0 | 6 | 6 | 7 | 1 | 0 | 5 | 4 | 28 | 20.1 | +21 |
| 5 | 5 | 2014-11-08 | 26-239 | GSW | @ | HOU | W (+11) | 1 | 39:58 | 13 | 19 | .684 | 6 | 9 | .667 | 2 | 2 | 1.000 | 1 | 8 | 9 | 5 | 4 | 0 | 5 | 2 | 34 | 30.7 | +13 |
| 6 | 6 | 2014-11-09 | 26-240 | GSW | @ | PHO | L (-12) | 1 | 34:04 | 10 | 20 | .500 | 4 | 10 | .400 | 4 | 4 | 1.000 | 0 | 2 | 2 | 10 | 5 | 0 | 10 | 5 | 28 | 18.6 | -6 |
| 7 | 7 | 2014-11-11 | 26-242 | GSW | | SAS | L (-13) | 1 | 35:46 | 7 | 18 | .389 | 0 | 7 | .000 | 2 | 2 | 1.000 | 0 | 6 | 6 | 5 | 0 | 0 | 3 | 2 | 16 | 7.7 | -18 |
| 8 | 8 | 2014-11-13 | 26-244 | GSW | | BRK | W (+8) | 1 | 31:56 | 6 | 12 | .500 | 3 | 7 | .429 | 2 | 2 | 1.000 | 1 | 2 | 3 | 5 | 0 | 0 | 3 | 5 | 17 | 10.8 | +6 |
| 9 | 9 | 2014-11-15 | 26-246 | GSW | | CHO | W (+25) | 1 | 28:18 | 8 | 15 | .533 | 3 | 6 | .500 | 0 | 1 | .000 | 2 | 3 | 5 | 9 | 1 | 0 | 1 | 3 | 19 | 18.7 | +23 |
| 10 | 10 | 2014-11-16 | 26-247 | GSW | @ | LAL | W (+21) | 1 | 29:56 | 10 | 19 | .526 | 5 | 9 | .556 | 5 | 5 | 1.000 | 0 | 4 | 4 | 15 | 1 | 0 | 3 | 0 | 30 | 30.4 | +27 |
| 11 | 11 | 2014-11-21 | 26-252 | GSW | | UTA | W (+13) | 1 | 24:30 | 3 | 9 | .333 | 2 | 6 | .333 | 0 | 0 | | 0 | 5 | 5 | 10 | 2 | 0 | 4 | 2 | 8 | 8.6 | +23 |
| 12 | 12 | 2014-11-23 | 26-254 | GSW | @ | OKC | W (+5) | 1 | 38:07 | 5 | 15 | .333 | 2 | 6 | .333 | 3 | 4 | .750 | 0 | 6 | 6 | 6 | 1 | 0 | 1 | 0 | 15 | 12.1 | +2 |
| 13 | 13 | 2014-11-25 | 26-256 | GSW | @ | MIA | W (+17) | 1 | 37:05 | 12 | 19 | .632 | 8 | 11 | .727 | 8 | 9 | .889 | 0 | 6 | 6 | 7 | 3 | 1 | 2 | 2 | 40 | 38.7 | +28 |
| 14 | 14 | 2014-11-26 | 26-257 | GSW | @ | ORL | W (+15) | 1 | 23:38 | 9 | 13 | .692 | 6 | 8 | .750 | 4 | 5 | .800 | 1 | 4 | 5 | 8 | 0 | 0 | 1 | 1 | 28 | 28.2 | +22 |
| 15 | 15 | 2014-11-28 | 26-259 | GSW | @ | CHO | W (+5) | 1 | 32:20 | 9 | 20 | .450 | 1 | 10 | .100 | 7 | 7 | 1.000 | 1 | 3 | 4 | 6 | 1 | 2 | 4 | 4 | 26 | 18.2 | +16 |
| 16 | 16 | 2014-11-30 | 26-261 | GSW | @ | DET | W (+11) | 1 | 28:22 | 5 | 9 | .556 | 1 | 3 | .333 | 5 | 6 | .833 | 1 | 2 | 3 | 10 | 1 | 0 | 1 | 2 | 16 | 18.8 | +11 |
| 17 | 17 | 2014-12-02 | 26-263 | GSW | | ORL | W (+1) | 1 | 36:27 | 8 | 17 | .471 | 3 | 9 | .333 | 3 | 3 | 1.000 | 1 | 3 | 4 | 5 | 0 | 0 | 4 | 2 | 22 | 13.6 | +8 |
| 18 | 18 | 2014-12-04 | 26-265 | GSW | | NOP | W (+27) | 1 | 30:37 | 8 | 17 | .471 | 3 | 6 | .500 | 0 | 0 | | 0 | 3 | 3 | 11 | 4 | 1 | 2 | 0 | 19 | 21.6 | +11 |

# Process for gathering the data ...

Click on the CSV option and the data should transform as follows:

# Process for gathering the data …

Highlight all of the comma separated values and copy (Ctrl-C) to clipboard

# Process for gathering the data ...

Open up a spreadsheet application and paste the data into cell A1 (Ctrl-V)

In Google Sheets

- In Cell B1, type: =arrayformula(substitute(split(subst itute(A1,",",",|")),",","), "|","")))
- Fill down Cell B1
- Highlight B1 to AE87 and Copy to clipboard (Ctrl-C)
- Edit-Paste Special-Paste values only
- Delete Column A

In Excel

- Click Data-Text to Columns
- Choose Delimited, click Next
- Choose Comma, click Next
- Click Finish

# Step 3 - Cleanse the data for use

Now we need to cleanse the data to make it usable by keeping only data relevant to our question.

- Adjust cell widths
- Ctrl-H find and replace Did Not Play with DNP
- Adjust cell widths again
- Remove repeated headers (where you see Rk/G every 20 rows)
- Name column F and H (Home, Result)
- Find which columns are needed for this comparison
    - Removed: Rk Age Tm Opp GS MP ORB DRB TRB AST STL BLK TOV PF PTS GmSc +/-

# Process for cleansing the data ...

- Changed column header with % in it: FG%, 3P% and FT% to FGP, 3PP and FTP
- Changed column header with 3 in it: 3P 3PA 3PP to ThreePM ThreePA ThreePP
- Deleted DNP data rows from table
- Saved data as Curry3PP_1415.csv

# Step 4 - Analyse the data

**Using R and R Studio**

R is a data analysis package that is powerful and R Studio is a good GUI that helps you use it.

Downloaded and installed R for windows

- https://cran.r-project.org/

Downloaded and installed RStudio for windows

- https://www.rstudio.com/products/rstudio/download/

**Using Microsoft Excel / Google Sheets**

Though it's at times harder to run multiple simulations of data in Excel than R, it's an easier tool to use and generally already installed.

All of the instructions (to the best of my knowledge) could be completed using either Excel or Sheets interchangeably.

# Process for analysing data …

In R you can create a .r file that contains code that enables you to answer questions.

Here I want to load my CSV file into data and create an array that holds Curry's average 3 Points Percentage in all games and his average when he's had a bad game (<0.20)

```
data = read.csv("Curry3PP_1415.csv", header=TRUE)
attach(data)

steph3low_pct = array()
steph3low_pct = steph3low_pct[!is.na(steph3low_pct)]
steph3avg_pct = array()
steph3avg_pct = steph3avg_pct[!is.na(steph3avg_pct)]
```

I first wanted to know what Curry's average 3 Points percentage was for every game he played during the season.

I can find this using an AVERAGE function on the entire 3 Points Percentage column

# Process for analysing data ...

I then store values in these arrays using code - most notably a for loop, then calculate and display the average (mean) for that array.

```
steph3avg_pct = c(steph3avg_pct,data$ThreePP[1])
for (i in 2:80) {
  steph3avg_pct = c(steph3avg_pct,data$ThreePP[i])
    if (data$ThreePP[i-1] <= .20)
        {steph3low_pct =
c(steph3low_pct,data$ThreePP[i])}}

mean(steph3avg_pct)
mean(steph3low_pct)
```

I then wanted to know what Curry's average 3 Points percentage was for each game AFTER he'd had a bad one.

I used a conditional format to highlight games where there were <0.20 3 point percentage

I then used an IF statement in a new column to only display a percentage value if the previous game was <0.20 3 point percentage

I then could use an AVERAGE statement on values in that new column to find out what Curry scores the match after a 0.20 game

# Process for analysing data ...

The process determines the average scores for Curry both for all games, and the game after he has a bad one:

All games: 0.435

After bad one: 0.424

The cells which contain the AVERAGE formulas should provide the average scores for Curry both for all games and the game after he has had a bad one:

All games: 0.435

After bad one: 0.424

# Step 5 - Evaluating your hypothesis …

Hypothesis - assume first one chosen:

**Golden State Warriors will return to form and win the series because Steph Curry should be able to average a 0.40 3 point percentage after his two poor shooting games against the Memphis Grizzlies in the 2014-15 NBA playoffs.**

Curry averages a 0.435 3 point percentage in his 2014-15 matches. So he is expected to return to his average and therefore proves the hypothesis true (or at least more likely)

Curry also averages a 0.424 3 point percentage after having one bad game using his 2014-15 matches. He is expected to return to this average and therefore proves the hypothesis true (or at least more likely)

# Step 5 - Evaluating your hypothesis ...

Hypothesis - assume first one chosen:

**Golden State Warriors will return to form and win the series because Steph Curry should be able to average a 0.40 3 point percentage after his two poor shooting games against the Memphis Grizzlies in the 2014-15 NBA playoffs.**

Curry also averages a 0.422 3 point percentage in his 2014-15 regular season matches against Memphis. So he is expected to return to his average and therefore proves the hypothesis true (or at least more likely)*

*This analysis can be completed in Excel using an AVERAGEIF formula

# Regression to the mean

The concept that Curry will bounce back and return to form suggests that over the long term his statistics will return to his average. This is known as "regression to the mean".

A way to question the simplicity of this analysis is by:

- checking whether his 80 games in the season is a large enough sample
- is there more reasons behind his low shooting percentage in his last two playoff games (less minutes, more shots instead of passing, more assists instead of shooting)?
- does he have a historically (> 1 season) bad record against Memphis (struggles vs one team) or in playoffs in general (struggles under pressure)?

# Step 6 - Present the Conclusion

Using data that we have analysed so far, within the scope of our investigation time period, we present the conclusion that:

**Golden State Warriors will return to form and win the series because Steph Curry  should be able to average a 0.40 3 point percentage after his two poor shooting games against the Memphis Grizzlies in the 2014-15 NBA playoffs.**

In the next 3 playoff games against the Memphis Grizzlies, Curry shoots at 0.444, 0.462 and 0.615 for 3 point percentages and the Golden State Warriors win by more than 10 points in each game.

# The rest is history

Golden State - NBA Champions 2014-15