

waRne - Putting cricket analytics in a spin

Contents

Abstract	1
Introduction	1
Why cricket needs reproducibility	2
Using waRne for teaching statistics and tidyverse	5
The Yorker	6
Easier engagement of fans	8
Dhoni Dilemma	8

Abstract

The importance of reproducibility, and the related issue of open access to data, has received a lot of recent attention. Momentum on these issues is gathering in the sports analytics community. While cricket is the second largest commercial sport in the world, unlike other popular international sports, there has been no mechanism for the public to access comprehensive statistics on players and teams. Expert commentary currently relies heavily on data that isn't made readily accessible and this produces an unnecessary barrier for the development of an inclusive sports analytics community.

Introduction

Data access is a key enabler for any analytics community. Most major sports have easy access to match statistics, for example `nbastatR` for NBA, `Lahman` for baseball, `deuce` for tennis and `nflfastR` for NFL. Through access to sports data and reproducible findings and metrics fans clubs and researchers are better able to understand the game, predict match outcomes and rate players. For example the way teams tackled 4th down decisions in the NFL has changed since (Romer 2006) seminal work. As teams have changed their 4th down decision making this allows follow up research such as (Yam and Lopez 2019) which looked at more granular data to see if this happens in practice.

By making data more accessible and more advanced metrics more accessible fans data journalism in sports has grown in recent years. For example in (Horowitz, Yurko, and Ventura 2017) has enabled EPA to enter popular discussion among fans¹².

Cricket is the second largest sport in the world. However, unfortunately there is no easy accessible way to access ball by ball data nor aggregated statistics of teams and players. Data while available on sites like `espnricinfo` are not in an easy to use form. For example, each match is listed on different webpages so hours upon hours of time would be required to copy and paste a single season, not to mention the added difficulty of linking players between games and competitions in different countries. Hence, there are significant logistical barriers for prospective fans and analysts studying the game, which stagnates understanding of cricket.

This paper describes the `waRne` package, the first to provide free and easy access to data for cricket for fans. Web scraping tools are available for fans to easily scrape the play by play commentary data on `espnricinfo`. For the first time fans can evaluate their favourite teams and players and do so in a reproducible and accessible

¹<https://twitter.com/minakimes/status/1205591497497817089>

²<https://twitter.com/KevinColePFF/status/1277974109175951368>

manner. We hope that this package can be used for fans to better understand the game, for teachers to use for interesting examples in class and for analysts in clubs who might not have access to ball by ball data.

Why cricket needs reproducibility

Data accessibility enables fans, analysts and researchers to better understand the game. Through being able to reproduce common popular metrics, visualisations and article findings.

Through being able to reproduce, fans are able to make accessible findings for others and importantly they are able to extend and grow concepts. Unfortunately what we see through leading cricket analytics providers is a track record of confusing output for fans. This can lead to lower engagement and dismissal of cricket analytics.

For example in this series of tweets we see the narrative being pushed that Steve Smith is a good player vs pace bowling unfortunately just a few months prior the same company and journalist published an article which had Steve Smith doing much worse against pace (balls above 140km/h). Unlike a similar sport baseball, fans have no easily accessible way of seeing if Steve Smith vs pace is a strength as alluded to in the original tweet, or a weakness like the same persons published online article. In comparison, fans are able to get a breakdown of Mike Trout vs fastballs from using baseballr which provides access through statcast data from baseballsavant.

Another barrier for fans to engage in analytics is the inability to look up insights around their favourite players, for cricket many fans have to take insights as standalone facts, without the ability to explore other players who might be of interest. For example after seeing this tweet fans might want to know how their favourite player goes against various types of bowling. This is in stark contrast to other major sports like the NFL whereby some of the leading analytical writers make their insights more accessible by sharing code³ and data sources.⁴⁵⁶

Without accessible data and a mechanism to get the data easily into peoples hands, how can cricket experience the growth that the NFL has throughout the past few years⁷. Through waRne, we hope to get the ball rolling on changing things for the cricket community. Fans of cricket might be interested to see if their own favourite players perform different against spin vs pace bowling, they might want to explore a different subset of players runs per over.

```
## `summarise()` regrouping output by 'bowling_style' (override with `.groups` argument)
```

³<https://www.opensourcefootball.com/posts/2020-09-09-creating-an-expected-field-goal-metric/>

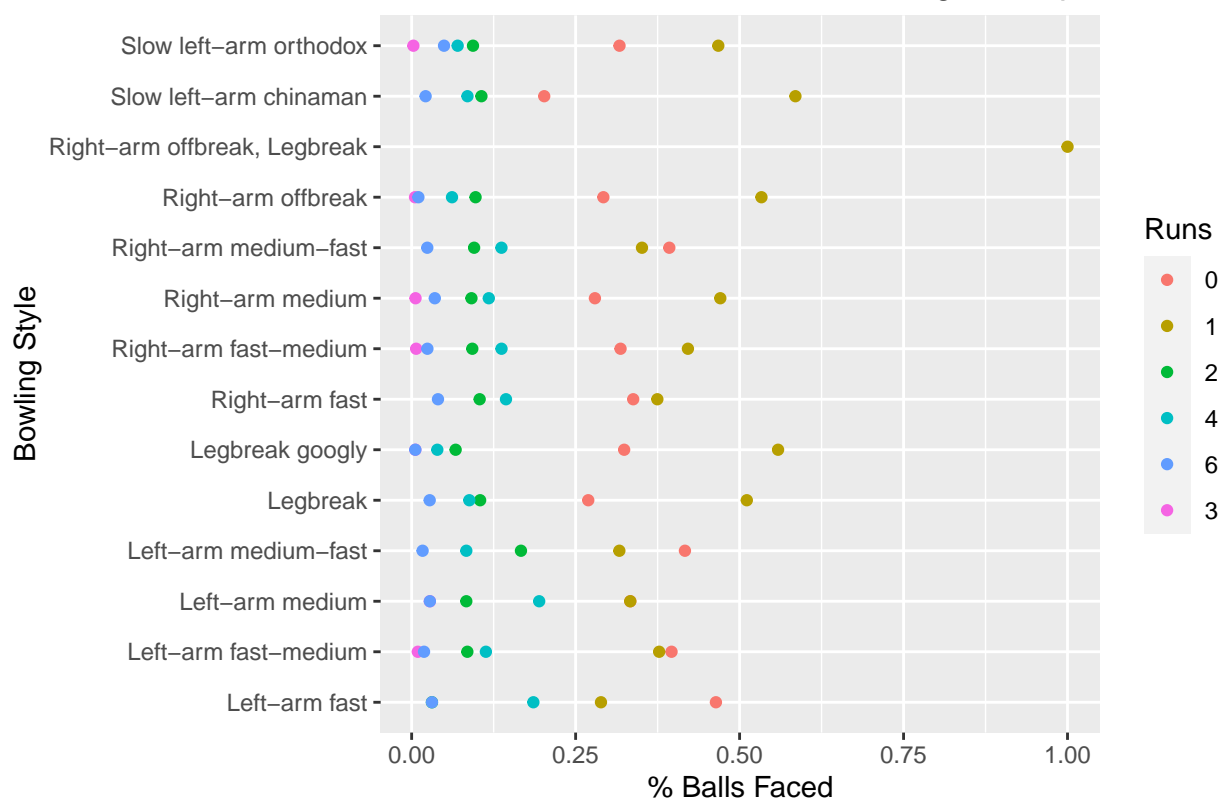
⁴<https://twitter.com/ChiefsAnalytics/status/1310758685325455360>

⁵<https://twitter.com/ChiefsAnalytics/status/1307839100896948230>

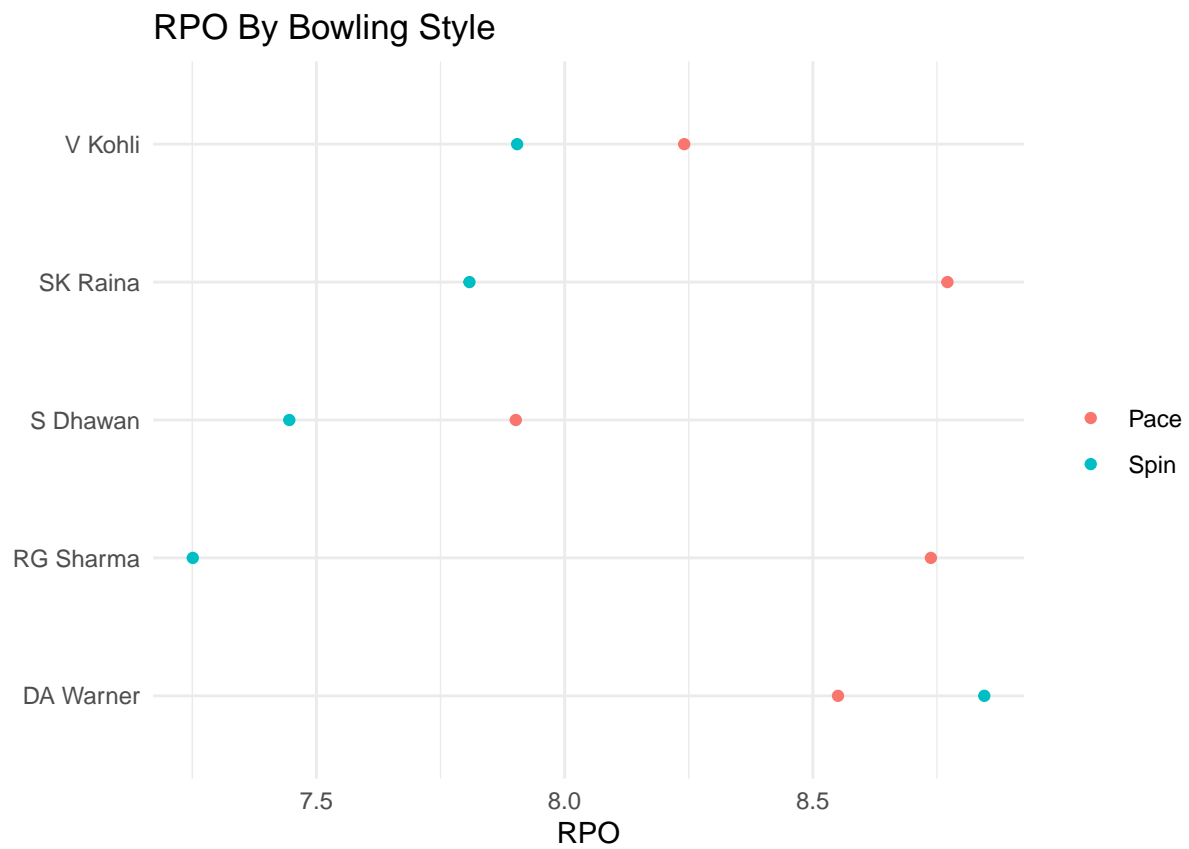
⁶<https://twitter.com/benbbaldwin/status/1311344809823657985>

⁷<https://www.samford.edu/sports-analytics/fans/2020/The-Future-of-NFL-Data-Analytics>

Steve Smith's Distribution of Runs Against Spin V Pace

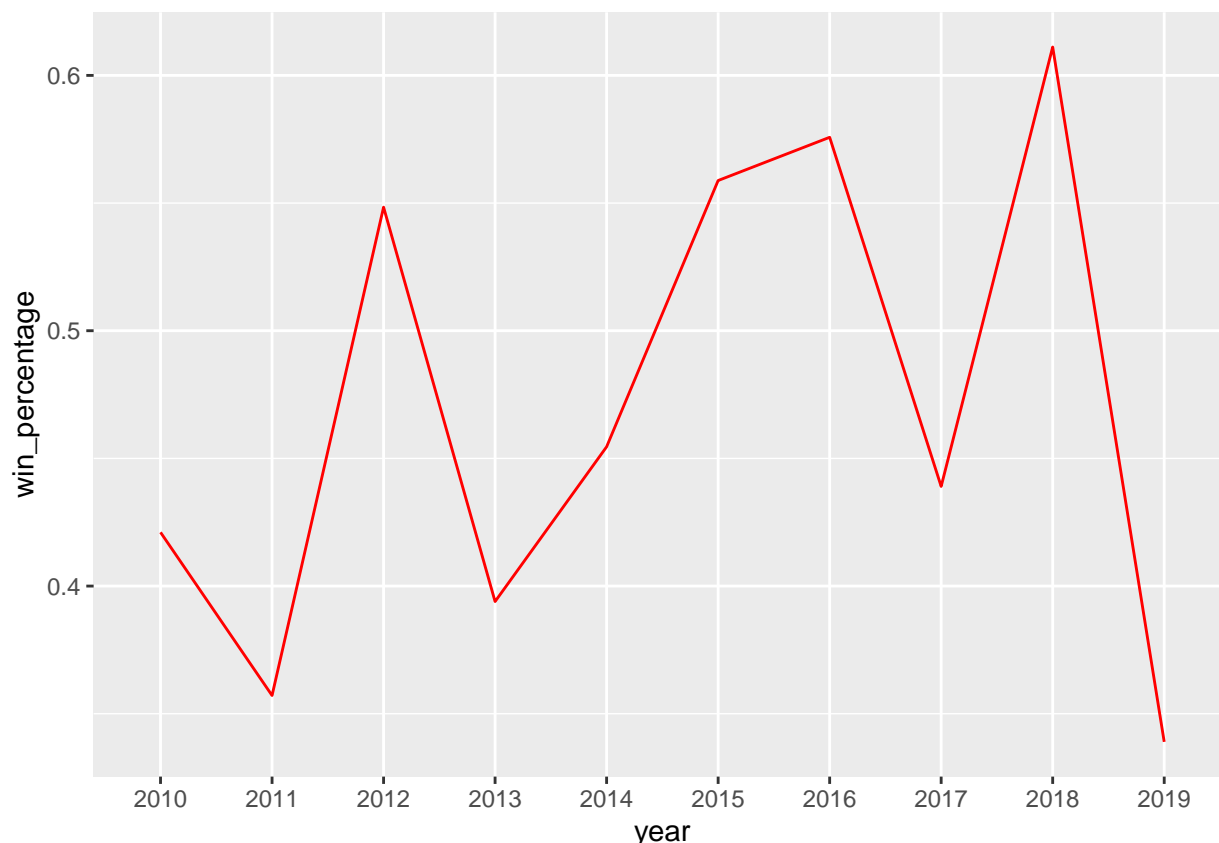


`summarise()` regrouping output by 'pace_spin' (override with `.groups` argument)



```
## Adding missing grouping variables: `match_id`
```

```
## `summarise()` regrouping output by 'batting_team_result' (override with `.groups` argument)
```



Using waRne for teaching statistics and tidyverse

We can also use waRne and other R packages as an easy way to engage students in learning statistical concepts.

For example using the above we can not only look at the answers graphically, but we can run a statistical test. A common classroom example for learning the binomial distribution is to ask if a coin is biased given a proportion of heads and tails given a sample size (McElreath 2020). Instead of asking about coins, fans of sport and cricket might be interested to ask the question; after winning the coin toss, should a team decide to bat first or bat second (Kvam and Sokol 2004).

While a coin toss is a seemingly trivial example, like in most professional sports the winner of the coin toss gets to decide what to do. In cricket, the winner of the coin toss gets to decide if they want to bat first and thus set the total that the opposition needs to pass by a single run to win, or if they want to bowl first and thus chase down the total set. Deciding what to do after winning a coin toss has proved to be a popular media piece in recent years⁸⁹¹⁰¹¹¹²¹³¹⁴

Fans and analysts of the game might want to make this decision based on a simple statistic, for example they might make it based on answer these questions instead.

⁸https://www.espnricinfo.com/story/_/id/21489056/stuart-wark-cricket-move-away-coin-toss

⁹https://www.espnricinfo.com/video/clip/_/id/23230682%5D

¹⁰https://www.espnricinfo.com/story/_/id/20429499/lehmann-backs-scrapping-toss

¹¹https://www.espnricinfo.com/story/_/id/28770755/how-much-does-losing-tosses-impact-visiting-teams

¹²<https://www.forbes.com/sites/tristanlavalette/2018/08/20/are-cricket-matches-being-decided-by-the-luck-of-a-coin-toss/#735456837eff>

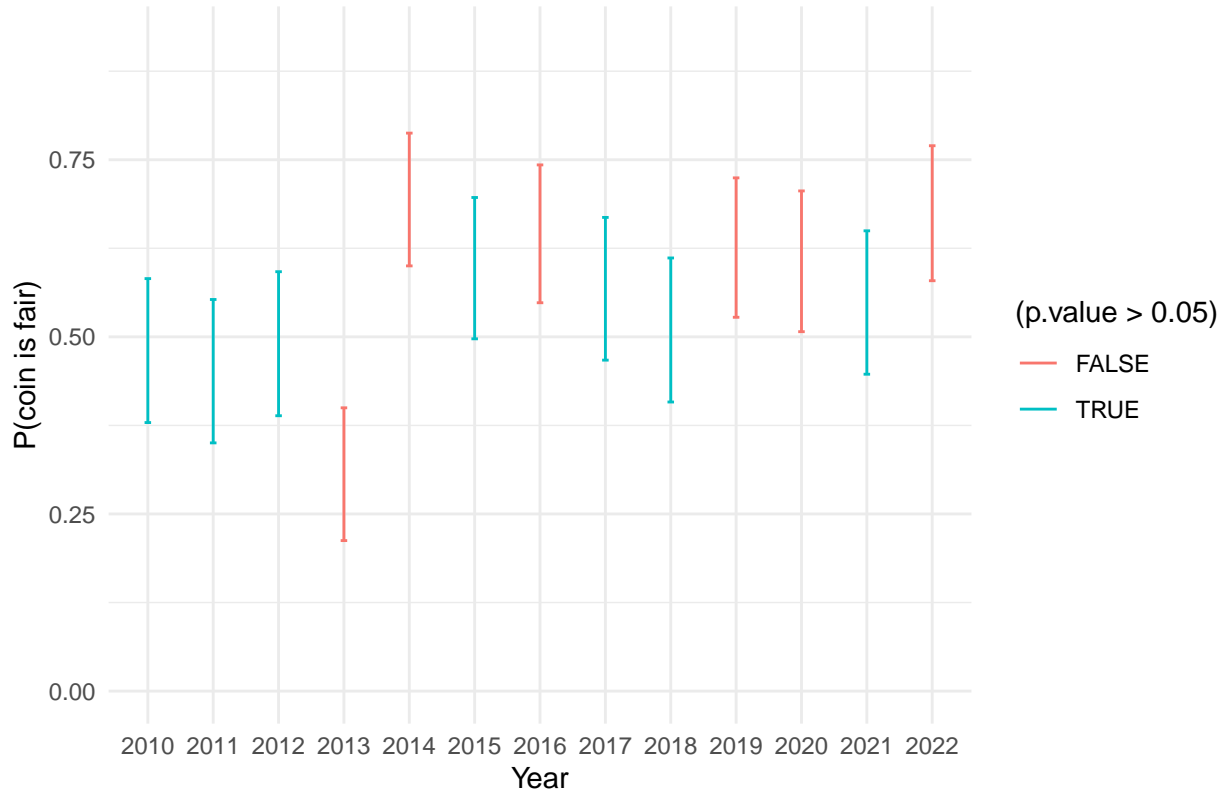
¹³<https://www.statsinsider.com.au/bbl/how-important-is-winning-the-toss-in-the-big-bash-league>

¹⁴https://www.espnricinfo.com/story/_/id/18568387/tim-wigmore-how-batting-second-become-more-fruitful-more-popular

- Do teams that bat second have a higher winning percentage than those who bat first?
- Is this consistent across leagues and levels of cricket?

```
## [1] 421 479
```

Annual coin fairness test



The Yorker

The yorker is considered the most difficult delivery to hit for a number of reason [Rob can you fill this in]. By pulling the delivery type from the play description on espnricinfo, every time a yorker was bowled we have the players who bowled and batted. We can take a look at the ability to hit yorkers for the closers we compared Dhoni to, and see how they all fare against the most difficult delivery in the game.

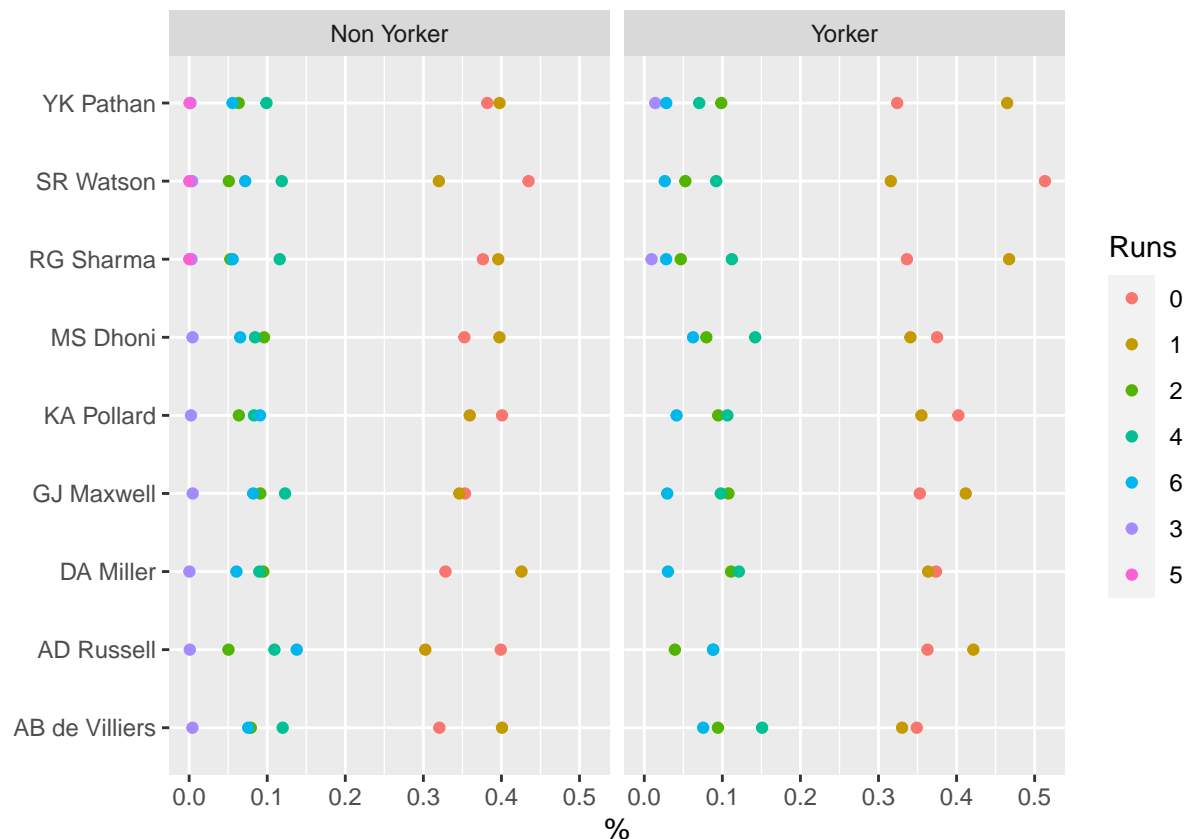
```
middle_batsman_of_interest <- c(
  230559, #pollard
  32498, #yusef pathan
  321777, #david miller
  44936, #AB De Villiers
  276298, #andre russell
  8180, #shane watson
  34102, #rohit sharma
  28081, #Dhoni
  325026 #maxwell
)

cricjam_data_big %>%
  mutate(yorker=if_else(str_detect(play_description,"york"),1,0))%>%
  select(batsman_id, earned_runs, play_description,yorker, tournament_type, batsman_name)%>%
```

```

filter(batsman_id%in% middle_batsman_of_interest) %>%
group_by(earned_runs, yorker, batsman_name)%>%
tally() %>%
filter(yorker==1) %>%
ungroup() %>%
  group_by(batsman_name) %>%
mutate(total = sum(n),
  percent = n/total,
  earned_runs = as.factor(earned_runs),
  yorker_nah = "Yorker") %>%
  bind_rows(cricjam_data_big %>%
mutate(yorker=if_else(str_detect(play_description,"york"),1,0))%>%
select(batsman_id, earned_runs, play_description,yorker, tournament_type, batsman_name)%>%
filter(batsman_id%in% middle_batsman_of_interest) %>%
group_by(earned_runs, yorker, batsman_name)%>%
tally() %>%
filter(yorker==0) %>%
ungroup() %>%
  group_by(batsman_name) %>%
mutate(total = sum(n),
  percent = n/total,
  earned_runs = as.factor(earned_runs),
  yorker_nah = "Non Yorker")) %>%
ggplot(aes(x = percent, y = batsman_name, color = earned_runs))+
geom_point() +
labs(color = "Runs",
  y = "",
  x = "%")+
facet_wrap(~yorker_nah)

```



Easier engagement of fans

To the surprise of many, there has not been a real push to engage fans in the analytics of the game of cricket. Unlike other sports which has seen a boom through the use of accessible data like Ice Hockey, NFL, Baseball and Basketball for example.

Without an easy accessible medium how can crickets version of an analytics community grow.

something something look towards how EPA has changed the way fans are engaged in NFL analytics - so maybe the change in run rate can be like EPA?

<https://twitter.com/cricvizanalyst/status/1311260989267087361>

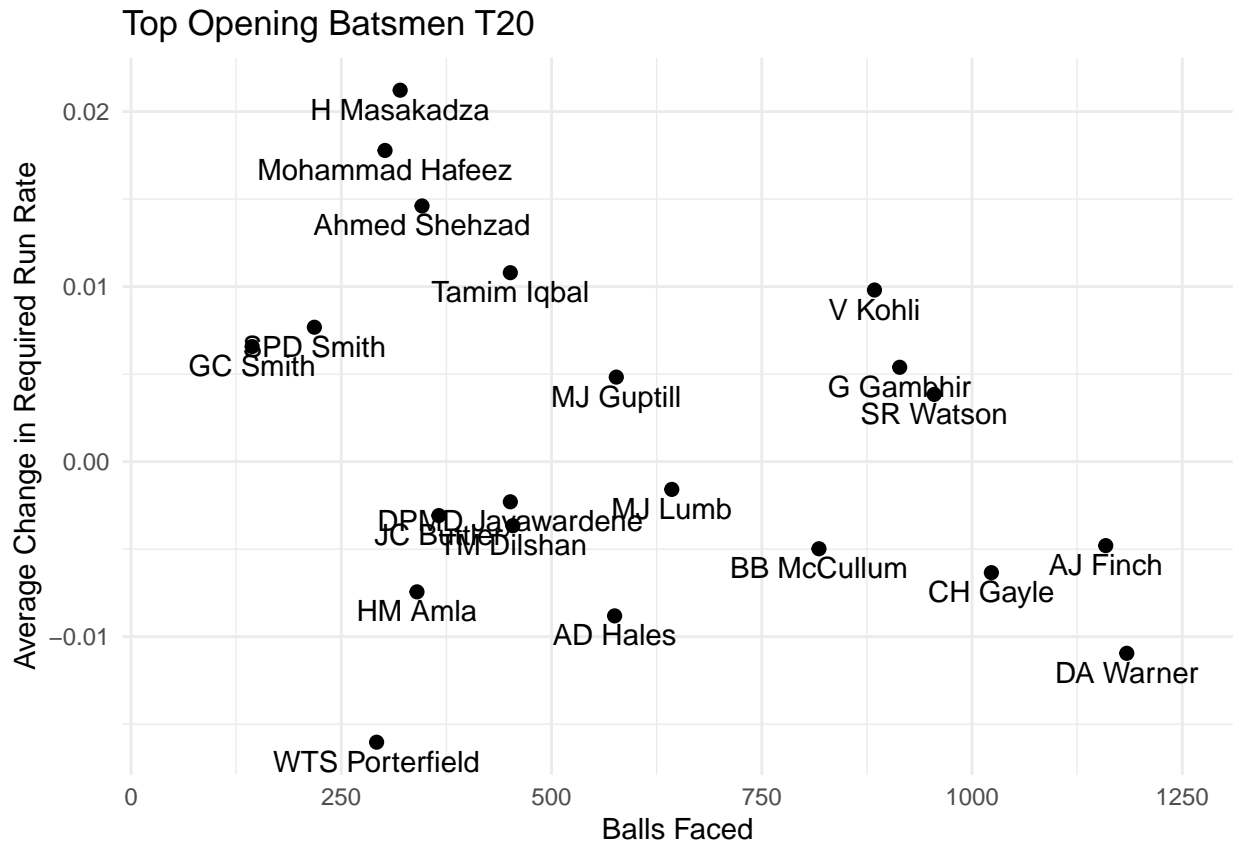
Dhoni Dilemma

With this new dataset, we allow fans to put on their analyst hats and to be able to ask themselves, what exactly makes a good closer and what exactly is the Dhoni Dilemma. To understand the Dhoni Dilemma, we need to see some interesting things to do with the change in RRR.

Cricket is an obvious sort of sport, to win in the end you have to get one run more than your opposition. A commonly used statistic during the broadcast is the required run rate, for example if you need 120 to win off 20 overs your RRR is \$6\$ runs an over. What follows is then obvious is the chasing team wins if they are able to lower the RRR that is score quicker.

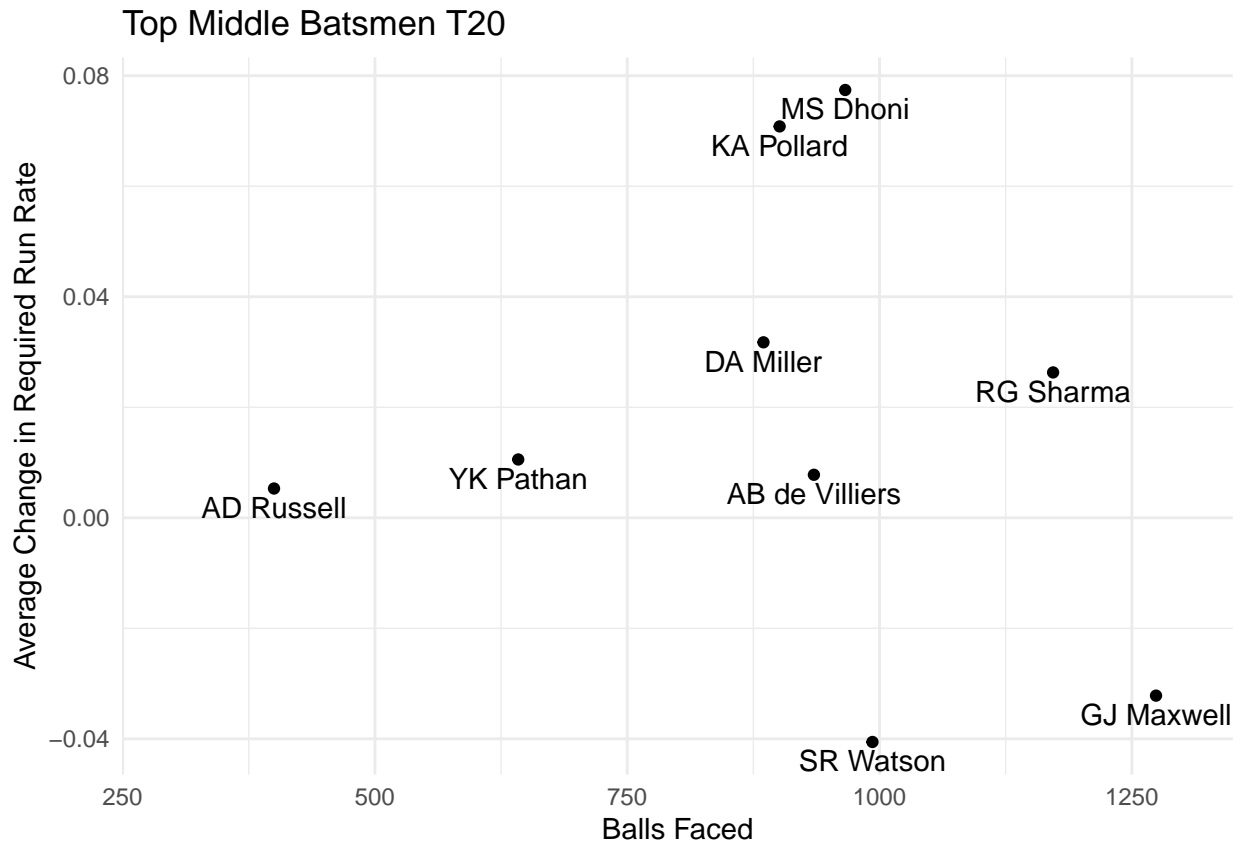
Espncricinfo did article on teams that win games get a higher proportion of their runs in the powerplay and first 10 overs. But proportion of team runs in T20 games doesn't really make sense as a batting metric as we

are not just concerned with how many runs they get, but how quickly they get them. So we introduce a new way to measure batsman change in RRR.



This graph makes sense as we can see some of the best openers in the world lower the required run rate. So what is the Dhoni Dilemma?

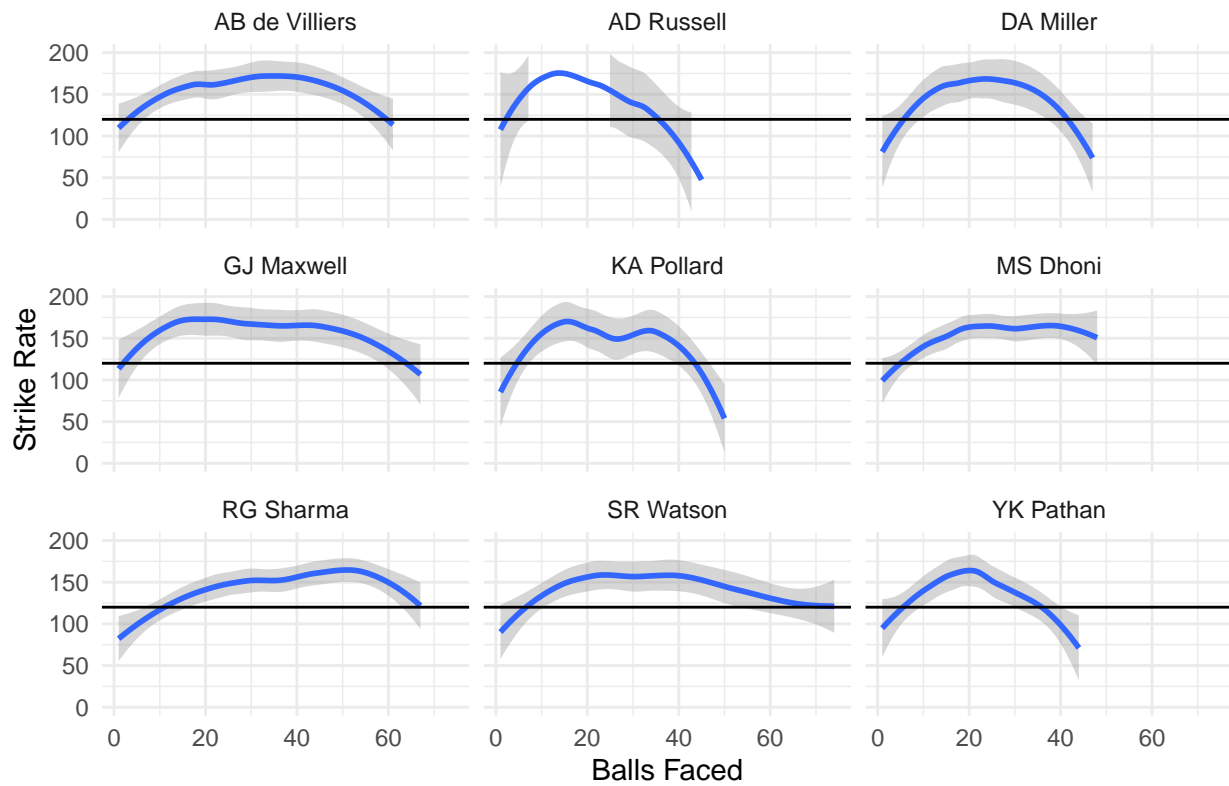
Dhoni is considered the greatest closing batsman in the world. He has pulled of many spectacular saves, and has T20 and T20I strike rates after 14 years in the game. But how does his required run rate compare to other middle order and closing batsman



Dhoni is renowned for being one of the world's best closers but we can see it seems as though in the middle he is letting the game get out of hand. How can one of the world's best closers let the run rate go up so consistently in the middle overs?

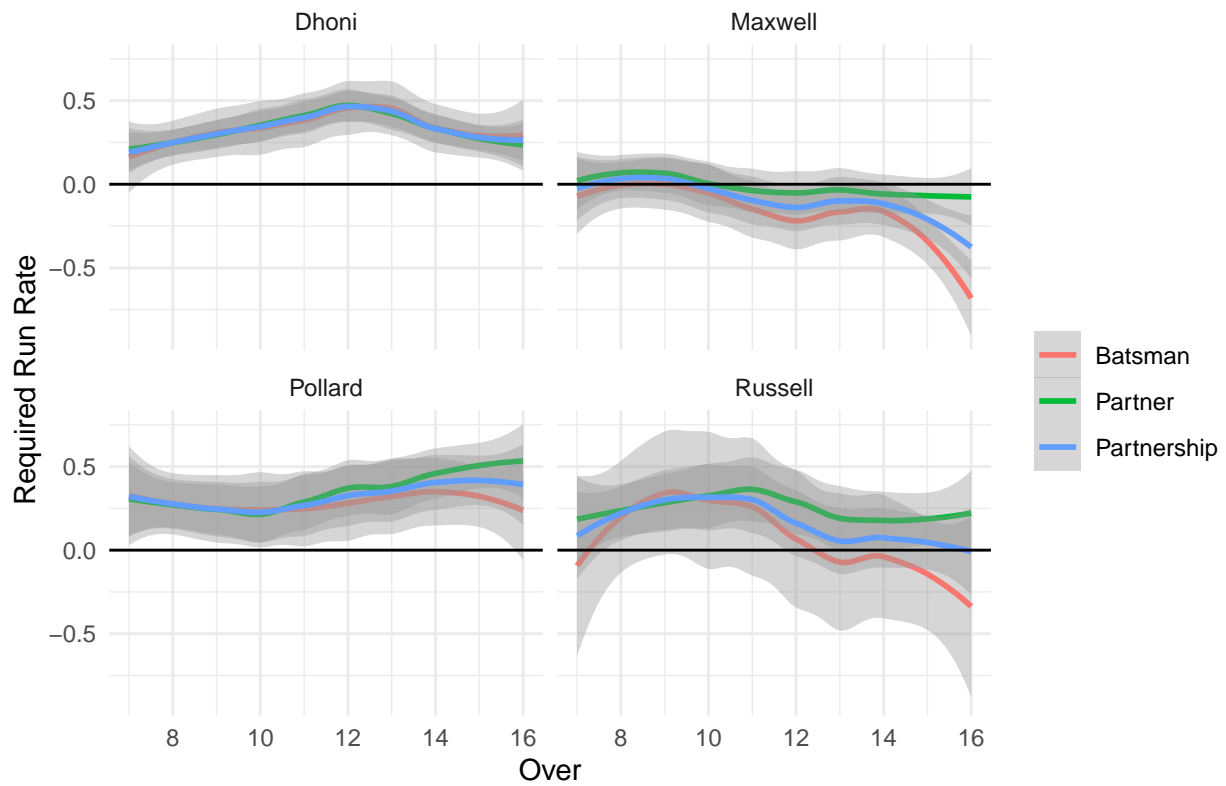
One possible explanation comes from looking at the acceleration of Dhoni compared to these other middle order batsmen. Dhoni comes in usually around overs 11/12, and from the plot below we can see that he is a slow starter, it takes several overs for him to get his strike rate up, which consistently stays high throughout the game even when other batsmen's begins to drop off.

Batsman Acceleration In Overs 11–20

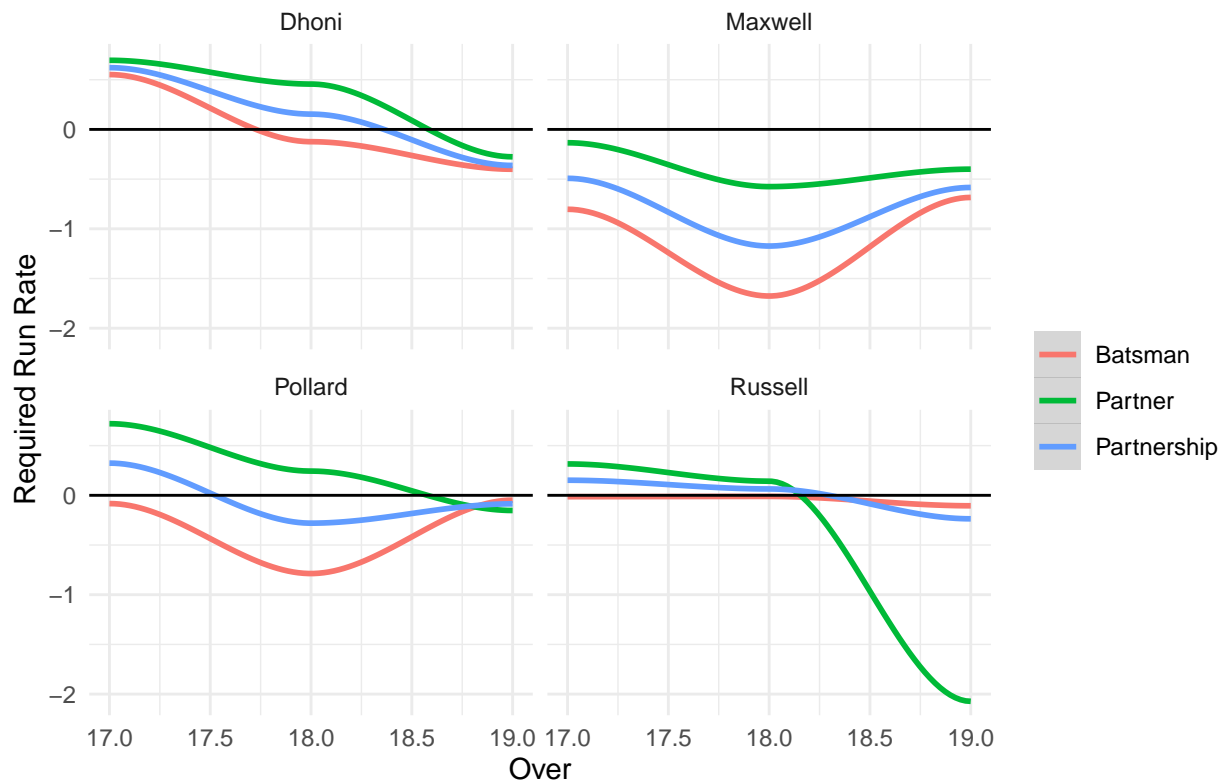


Another is the effort of Dhoni's partners in moving the game along while Dhoni is up

Change in Required Run Rate While Batting (or Partnering)

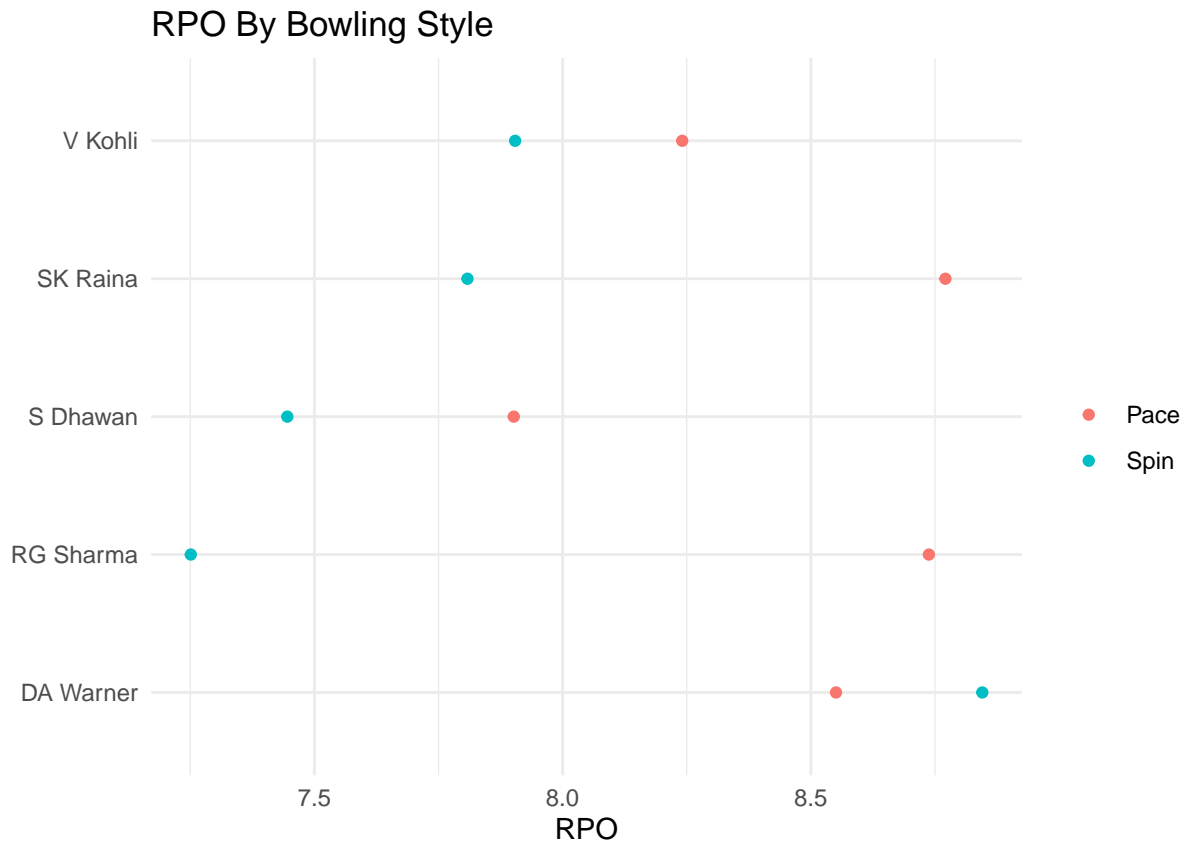


Change in Required Run Rate While Batting (or Partnering)



What we can tell is that in second inning chases, the required run rate for Dhoni is lower than our other top middle order batsmen, contrasting our initial findings. One would assume that this was due to Dhoni's partners, but it appears he still performs at an elite level

```
## `summarise()` regrouping output by 'pace_spin' (override with `.groups` argument)
```



Horowitz, Maksim, Ron Yurko, and Samuel L Ventura. 2017. "NflscrapR: Compiling the Nfl Play-by-Play Api for Easy Use in R." URL <https://Github.Com/Maksimhorowitz/nflscrapR>, R Package Version 1 (0).

Kvam, Paul H, and Joel Sokol. 2004. "Teaching Statistics with Sports Examples." *INFORMS Transactions on Education* 5 (1): 75–87.

McElreath, Richard. 2020. *Statistical Rethinking: A Bayesian Course with Examples in R and Stan*. CRC press.

Romer, David. 2006. "Do Firms Maximize? Evidence from Professional Football." *Journal of Political Economy* 114 (2): 340–65.

Yam, Derrick R, and Michael J Lopez. 2019. "What Was Lost? A Causal Estimate of Fourth down Behavior in the National Football League." *Journal of Sports Analytics* 5 (3): 153–67.