

Sports Analytics with R

Melbourne R Ladies

June 15, 2017



1 / 92

About Me

- ISEAL Research Fellow at Victoria University
- Principal Data Scientist for Tennis Australia's Game Insight Group
- Tennis Blog: on-the-t.com
- @StatsOnTheT



What It's Like to Be a Sports Statistician

- Finding sports data
- Scraping sports data
- Data wrangling
- Exploratory data analysis
- Predictive modelling
- Blogging



Finding Sports Data

Popular Types of Sports Data

- Box Scores
- Performance Metrics
- Tracking Data
- Wearable Data

Box Scores

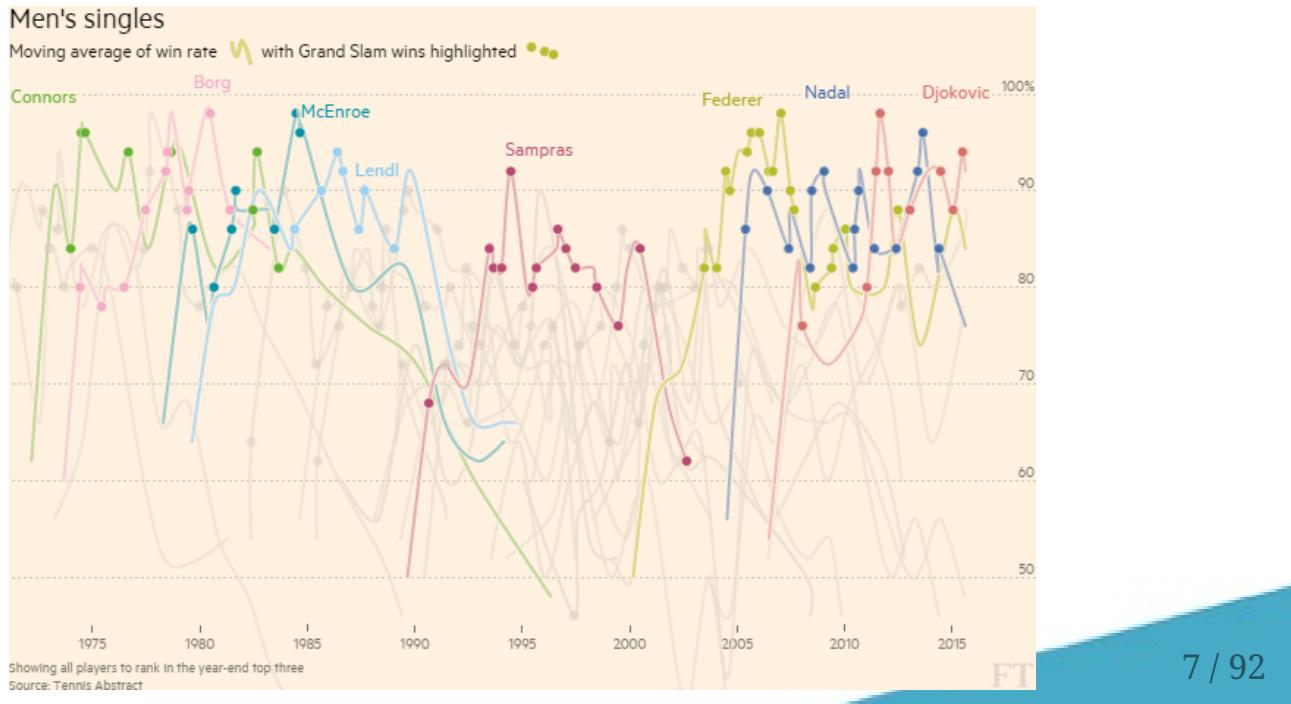
A 'box score' is a generic term for the summary statistics reported about a single sporting event.

BOSTON.	T.	R.	1B.	PO.	A.	E.	ATHLETIC.	T.	R.	1B.	PO.	A.	E.
G.Wright, s.s	6	4	4	1	5	2	Force, s. s...	5	1	2	1	3	2
Leonard, 2b.	6	3	3	4	4	3	Eggler, c. f..	5	3	3	0	0	0
O'Rourke, 1b	6	2	3	9	0	1	Fisler, r. f..	5	0	1	2	0	0
Murnan, l. f.	6	1	0	3	1	0	Meyerle, 3db	5	1	2	2	3	3
Schafer, 2d b	6	3	3	3	1	2	Sutton, 1st b.	5	1	2	10	0	0
McGinley, c.f	6	0	0	0	0	1	Coons, c....	5	1	0	1	1	3
Manning, r.f.	6	0	2	2	0	0	Hall, l. f....	5	1	3	5	0	0
Morrill, c....	6	2	2	4	1	2	Fowser, 2d b.	6	1	2	6	7	5
Josephs, p..	5	4	4	1	1	2	Knight, p....	5	2	2	0	1	2
Totals....	53	19	21	27	13	13	Totals....	46	11	17	27	15	15
Boston.....	0	1	3	3	4	1	0	2	5	19			
Athletic.....	1	0	0	0	3	3	2	/	2	0	11		
Runs earned—Boston, 4; Athletic, 5. Home-run—Hall, 1. Total bases on hits—Boston, 22; Athletic, 20. First base by errors—Boston, 8; Athletic, 5. Umpire, George White of Lowell, Mass. Time 2h. 47m.													

[1] Baseball box score from 1876.

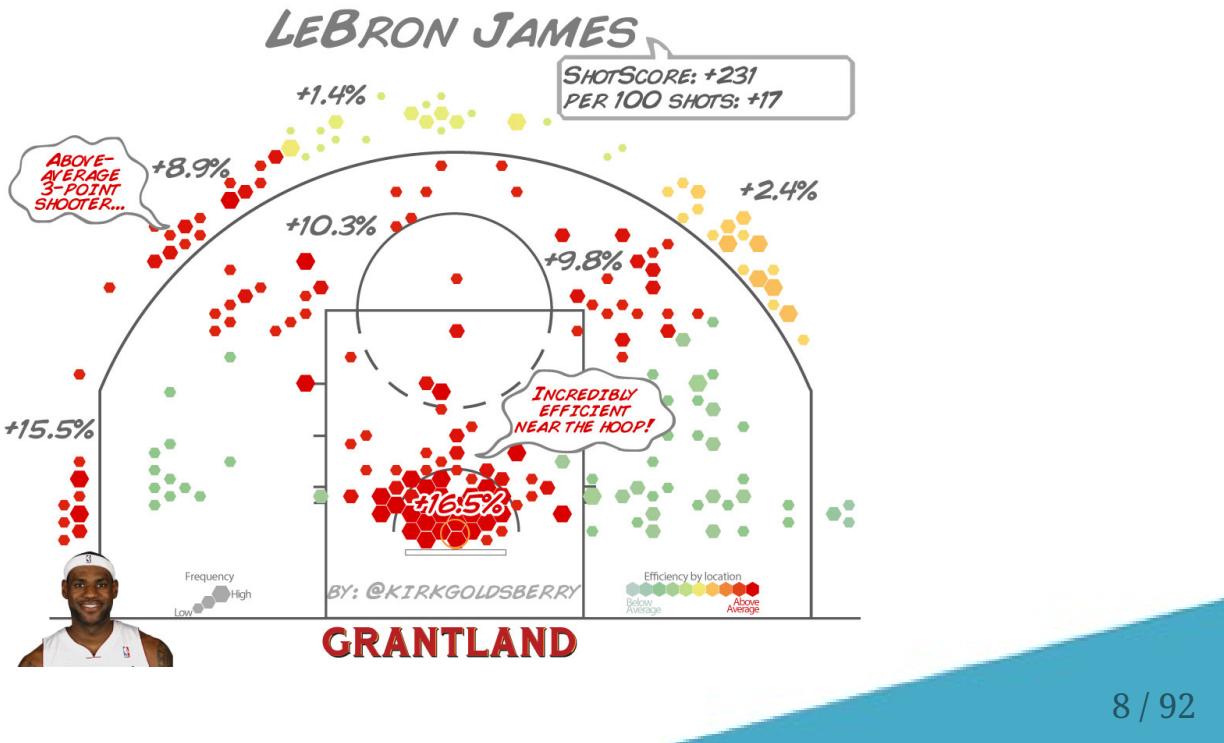
Performance Metrics

Performance metrics are usually repeated measures of team or players that are derived from historical box scores.



Tracking Data

'Tracking data' is spatio-temporal data of objects in a sporting event.



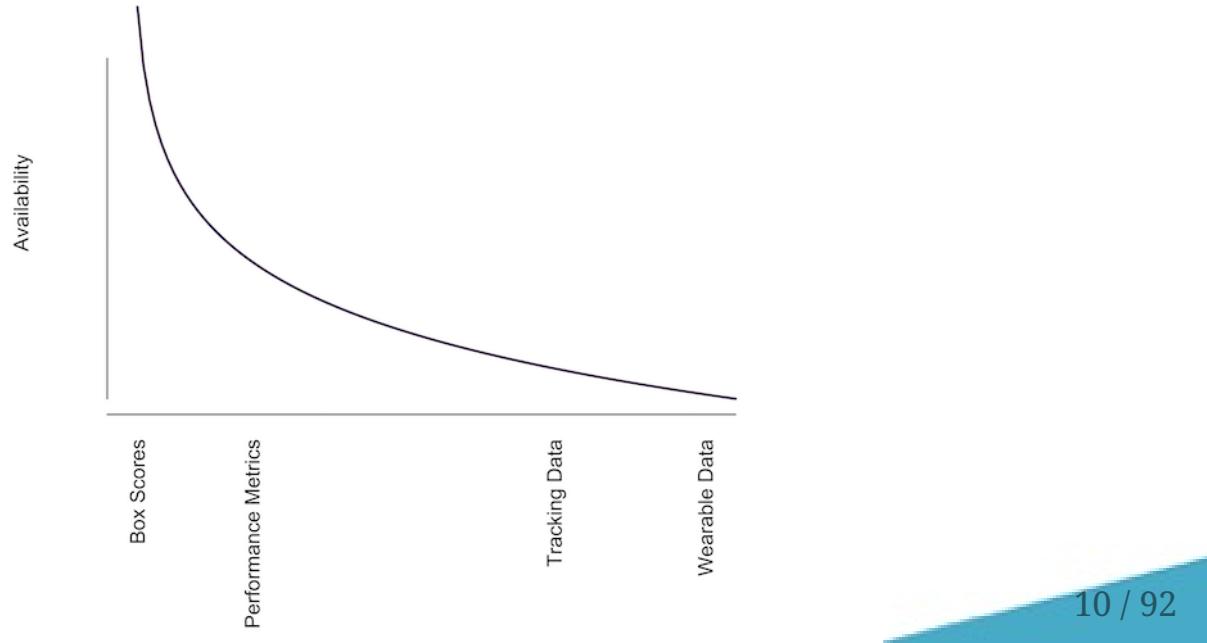
Wearables

- 'Wearables' are technologies that collect data about athletes when worn.
- Biometrics and movement are two common types of data collected.



Data Availability Curve

Sadly, not all types of sports data are equally available though the steepness of the curve depends on the particular sport.



Web Resources

opensourcesports kenpomeroycollegebasketballsit
databaseolympics cricketarchivecom
ncaafinalfour databasehockeycom whoscoredcom
baseballamerica databasegolfcom masseydata
unitedstatesfencingassociation nflstatistics
aflstatistics tennisabstract squawka
bballsportscom currrentnbastats fangraphs
footballdatacouk baseballmusings baseballlinks
cricket tennis hoopdatacom baseballcube
cricinfo nhlcom footballdata johnskiltonsbasesports
puckalytics crossfitgames apbr databasesports
retrosheet atptour thebaseballcube
soccerway collegefootballstatistics johnskiltonsbasesportsinfo
hockeyreference f1facts golf espnstatsinfo
quantockeycom sabrthesocietyforamericanbaseballresearch
fencing baseballprospectus rugbydata
finalfour basketballreference volleyballworldwide
wtatour baseballreferencecom profootballreference transfermarkt
olympicsatsportsreferencecom
databaseracingcom baseballgurucom
nhlplayersassociationstatistics masseyratings
databasebasketball washingtonpostcom
dougsteelesmlbandnbastatistics
baseballheatmaps databasefootballcom

R Resources

Sport	Package
General	stattleshipR SportsAnalytics odds.converter
Soccer	engsoccerdata
Baseball	pitchRx Lahman
Tennis	deuce
Cricket	cricketr!

More on deuce

More on deuce

- deuce is a package I created to make it easy to access large historical data on tennis matches

More on deuce

- deuce is a package I created to make it easy to access large historical data on tennis matches
- It combines data from multiple sites:
 - [flashscore.com](#)
 - [ATP Tour](#)
 - [WTA Tour](#)
 - [Tennis Abstract](#)

More on deuce

- deuce is a package I created to make it easy to access large historical data on tennis matches
- It combines data from multiple sites:
 - [flashscore.com](#)
 - ATP Tour
 - WTA Tour
 - Tennis Abstract
- Some of the data you can obtain from deuce includes:
 - Match results from Open Era (1968) to the present
 - Historical rankings
 - Player demographics
 - Point-level data for multiple years at Grand Slams
 - Shot level data from the Match Charting Project

Installing deuce

You can install deuce using devtools. It may take several minutes because of the size of the datasets being transferred.

```
library(devtools)  
devtools::install_github("skoval/deuce")
```

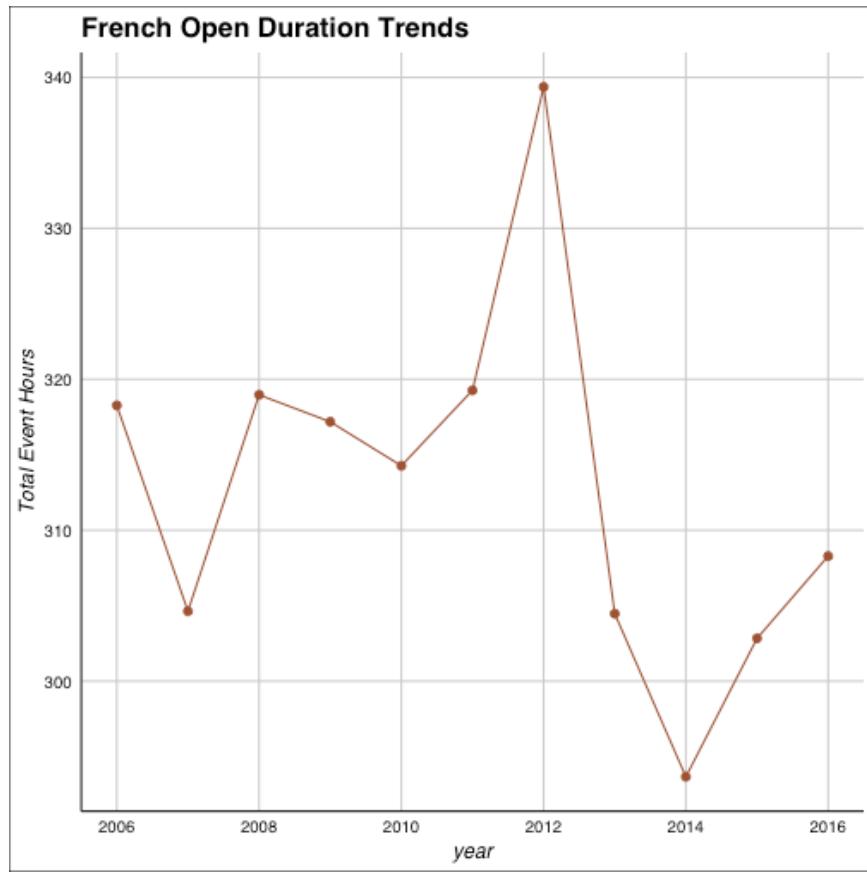
Illustrative Analysis: Match Duration

```
library(deuce)
library(dplyr)

data(atp_matches, package = "deuce")

french_minutes <- atp_matches %>%
  dplyr::filter(tourney_name == "Roland Garros",
    year >= 2006 & year < 2017)

french_minutes %>%
  group_by(year) %>%
  dplyr::summarise(
    total.time = sum(minutes, na.rm = T)
  ) %>%
  ggplot(aes(y = total.time / 60, x = year)) +
  geom_point(size = 2, col = "sienna") + geom_line(col = "sienna") +
  scale_y_continuous("Total Event Hours") +
  theme_gdocs() +
  ggtitle("French Open Duration Trends")
```



Scraping Sports Data

Ways of Getting Sports Data

- Not always fun but a necessary part of sports analysis
- There are two major ways to get data from Web:
 1. Import a file directly
 2. Extract from HTML



Example: Import Data File

- Files that can be read with `read.table` or related functions can be directly imported from a URL.
- Here we extract the most recent Australian Open match results and betting odds using `read.csv`.

```
url <- "http://www.tennis-data.co.uk/2017/ausopen.csv"  
read.csv(url)
```

Scraping from a Website

If you can't directly import data from the Web, you can still capture the data but you need to know whether it is *static* or *dynamic* data.

Scraping from a Website

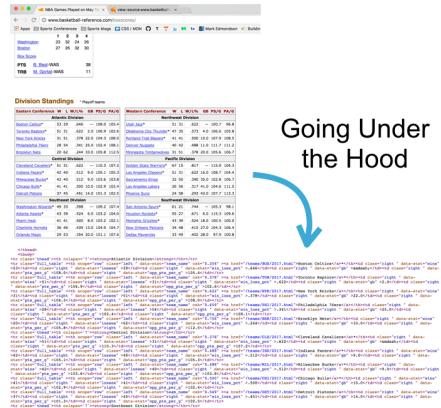
If you can't directly import data from the Web, you can still capture the data but you need to know whether it is *static* or *dynamic* data.



- Web data is *static* if you can see it in the source code.
- If you can't, the data is *dynamically* generated and you will need automation to extract it.

How Do I Know if Data is Static or Dynamic?

- You need to be able to inspect HTML and CSS
 - I recommend Google Chrome's developer tools
 - Get familiar with "View Source"
 - Try out the Selector Gadget



Example: Static or Dynamic?

Look at the source code for each of the following sites and determine whether they are examples of *static* or *dynamic* data.

Case 1. http://tennisabstract.com/reports/atp_elos_ratings.html

Case 2. <http://www.espncricinfo.com/ci/content/stats/>

Scraping Static Data with rvest

Scraping Static Data with `rvest`

- The `rvest` package is a suite of tools for scraping static Web data and putting them in easy-to-use objects (like data frames)

Scraping Static Data with `rvest`

- The `rvest` package is a suite of tools for scraping static Web data and putting them in easy-to-use objects (like data frames)
- Works with `magrittr` and allows piping commands with `%>%` operator

Scraping Static Data with `rvest`

- The `rvest` package is a suite of tools for scraping static Web data and putting them in easy-to-use objects (like data frames)
- Works with `magrittr` and allows piping commands with `%>%` operator
- Allows some browsing functionality

Scraping Static Data with `rvest`

- The `rvest` package is a suite of tools for scraping static Web data and putting them in easy-to-use objects (like data frames)
- Works with `magrittr` and allows piping commands with `%>%` operator
- Allows some browsing functionality
- Authored by Hadley Wickham

Example: Scraping Elo Ratings

```
library(rvest)

url <- "http://tennisabstract.com/reports/atp_elo_ratings.html"

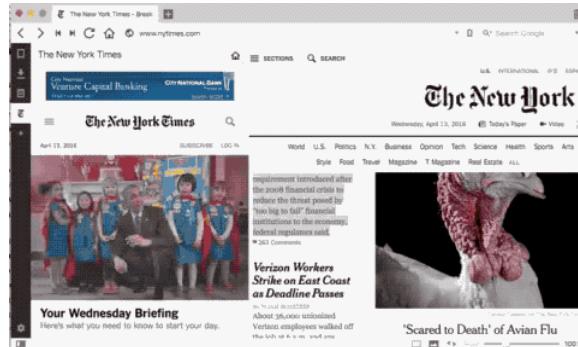
page <- read_html(url)

# Use table class to extract Elo table
elo <- page %>%
  html_nodes("table.tablesorter") %>%
  html_table()

head(elo)
```

Scraping Dynamic Data with RSelenium

- We have to automate Web browsing to get dynamic data
- *Selenium* is software that allows automated Web browsing
- **RSelenium** is a package that provides Selenium functionality in R



RSelenium: Basic Steps

1. Set the Web driver (select browser and port)
2. Find the elements with the data
3. Extract the content
4. Parse the contents

Example: Tennis Match Statistics

Consider the following match summary: [2017 Australian Open Final](#)

Example: Tennis Match Statistics

If we inspect the page, we find that these stats are dynamic data. We also find that the main table of content has the id detail.

```
<body id="top" class="tennis detailbody">
<div id="detail" class="sport-tennis"><div id="detcon">
<table class="detail">
<thead>
<tr>
<th class="header">
<div class="fleft">
<span class="flag fl_3473162"></span>ATP - SINGLES: <a href="#" onclick="window.open('/tennis/atp-singles/australian-open/'); return false;">Australian Open (Australia), hard - Final</a>
</div>
</th>
</tr>
</thead>
<tbody>
<tr>
<td class="hclean"></td>
</tr>
</tbody>
</table>
</div>
</body>
```

Using RSelenium

Below we activate the driver using a port that is not in use.

Note: You may need to activate javascript in the background for this driver to work.

```
library(RSelenium) # Load the package

# Match statistics URL
url <- "http://www.flashscore.com/match/Cj6I5iL9/#match-statistics;0"

# Establish remote driver using Chrome
remDr <-remoteDriver(port = 5556, browser = "chrome")
remDr$open(silent = TRUE)
remDr$navigate(url) # Navigate page
```

Using RSelenium

Next we extract the table of stats using the CSS `id` node.

```
# Get id element
webElem <- remDr$findElements(using = 'id', "detail")

# Use getElementText to extract the text from this element
unlist(lapply(webElem, function(x){x$getElementText()}))[[1]]

remDr$close() # Close driver when finished
```

Data Wrangling

Goal of Wrangling

The purpose of *wrangling* is to go from untidy data to a tidy row-by-column structure of our data, that has clearly named variables and valid values.

0499
stat3
111 067
101 stat1
stat2 053
131



Stat1	Stat2	Stat3
30	101	0.530
32	111	0.670
19	131	0.499

R's Arsenal for Wrangling

- `stringr` for string manipulation
- `tidyverse` a grammar for reshaping data
- `dplyr` a grammar for filtering and transforming data



Wrangling by Example

Consider this string of tennis match statistics scraped from the Web. How could we get this into a nice data structure?

```
match_stats # Example of untidy data
```

```
[1] "ALEXANDER\nZVEREV\nNOVAK\nDJOKOVIC\nMATCH STATS YTD STATS\nSERVICE STATS\n7\nAces\n1\n2\nDouble Faults\n3\n71%\n(32/45)\n1st Serve\n64%\n(43/67)\n84%\n(27/32)\n1st\nServe Points Won\n70%\n(30/43)\n69%\n(9/13)\n2nd Serve Points Won\n38%\n(9/24)\n0%\n(0/0)\nBreak Points Saved\n40%\n(2/5)\n9\nService Games Played\n10\nRETURN STATS\n30%\n(13/43)\n1st Serve Return Points Won\n16%\n(5/32)\n63%\n(15/24)\n2nd Serve Return\nPoints Won\n31%\n(4/13)\n60%\n(3/5)\nBreak Points Converted\n0%\n(0/0)\n10\nReturn\nGames Played\n9\nPOINTS STATS\n80%\n(36/45)\nReturn Points Won\n58%\n(39/67)\n42%\n(28/67)\nTotal Return Points Won\n20%\n(9/45)\n57%\n(64/112)\nTotal Points Won\n43%\n(48/112)"
```

What Steps Do We Need?

1. Identify the target structure (that is, variables and value types)
2. Split the string into the different variables
3. Extract values
4. Assign to variables in a data.frame

Goal Structure

We have a set of statistics for each player. One option is a long format with the following structure:

Statistic	Value	Player

Splitting

```
library(stringr) # Load stringr  
str_split(match_stats, "\n") # Split on return characters
```

```
## [[1]]  
## [1] "ALEXANDER"           "ZVEREV"  
## [3] "NOVAK"                "DJOKOVIC"  
## [5] "MATCH STATS YTD STATS" "SERVICE STATS"  
## [7] "7"  
## [9] "1"  
## [11] "Double Faults"  
## [13] "71%"  
## [15] "1st Serve"  
## [17] "(43/67)"  
## [19] "(27/32)"  
## [21] "70%"  
## [23] "69%"  
## [25] "2nd Serve Points Won"  
## [27] "(9/24)"  
## [29] "(0/0)"  
## [31] "40%"  
## [33] "9"  
## [35] "10"  
## [37] "11"  
## [39] "12"  
## [41] "13"  
## [43] "14"  
## [45] "15"  
## [47] "16"  
## [49] "17"  
## [51] "18"  
## [53] "19"  
## [55] "20"  
## [57] "21"  
## [59] "22"  
## [61] "23"  
## [63] "24"  
## [65] "25"  
## [67] "26"  
## [69] "27"  
## [71] "28"  
## [73] "29"  
## [75] "30"  
## [77] "31"  
## [79] "32"  
## [81] "33"  
## [83] "34"  
## [85] "35"  
## [87] "36"  
## [89] "37"  
## [91] "38"  
## [93] "39"  
## [95] "40"  
## [97] "41"  
## [99] "42"  
## [101] "43"  
## [103] "44"  
## [105] "45"  
## [107] "46"  
## [109] "47"  
## [111] "48"  
## [113] "49"  
## [115] "50"  
## [117] "51"  
## [119] "52"  
## [121] "53"  
## [123] "54"  
## [125] "55"  
## [127] "56"  
## [129] "57"  
## [131] "58"  
## [133] "59"  
## [135] "60"  
## [137] "61"  
## [139] "62"  
## [141] "63"  
## [143] "64"  
## [145] "65"  
## [147] "66"  
## [149] "67"  
## [151] "68"  
## [153] "69"  
## [155] "70"  
## [157] "71"  
## [159] "72"  
## [161] "73"  
## [163] "74"  
## [165] "75"  
## [167] "76"  
## [169] "77"  
## [171] "78"  
## [173] "79"  
## [175] "80"  
## [177] "81"  
## [179] "82"  
## [181] "83"  
## [183] "84"  
## [185] "85"  
## [187] "86"  
## [189] "87"  
## [191] "88"  
## [193] "89"  
## [195] "90"  
## [197] "91"  
## [199] "92"  
## [201] "93"  
## [203] "94"  
## [205] "95"  
## [207] "96"  
## [209] "97"  
## [211] "98"  
## [213] "99"  
## [215] "100"  
## [217] "101"  
## [219] "102"  
## [221] "103"  
## [223] "104"  
## [225] "105"  
## [227] "106"  
## [229] "107"  
## [231] "108"  
## [233] "109"  
## [235] "110"  
## [237] "111"  
## [239] "112"  
## [241] "113"  
## [243] "114"  
## [245] "115"  
## [247] "116"  
## [249] "117"  
## [251] "118"  
## [253] "119"  
## [255] "120"  
## [257] "121"  
## [259] "122"  
## [261] "123"  
## [263] "124"  
## [265] "125"  
## [267] "126"  
## [269] "127"  
## [271] "128"  
## [273] "129"  
## [275] "130"  
## [277] "131"  
## [279] "132"  
## [281] "133"  
## [283] "134"  
## [285] "135"  
## [287] "136"  
## [289] "137"  
## [291] "138"  
## [293] "139"  
## [295] "140"  
## [297] "141"  
## [299] "142"  
## [301] "143"  
## [303] "144"  
## [305] "145"  
## [307] "146"  
## [309] "147"  
## [311] "148"  
## [313] "149"  
## [315] "150"  
## [317] "151"  
## [319] "152"  
## [321] "153"  
## [323] "154"  
## [325] "155"  
## [327] "156"  
## [329] "157"  
## [331] "158"  
## [333] "159"  
## [335] "160"  
## [337] "161"  
## [339] "162"  
## [341] "163"  
## [343] "164"  
## [345] "165"  
## [347] "166"  
## [349] "167"  
## [351] "168"  
## [353] "169"  
## [355] "170"  
## [357] "171"  
## [359] "172"  
## [361] "173"  
## [363] "174"  
## [365] "175"  
## [367] "176"  
## [369] "177"  
## [371] "178"  
## [373] "179"  
## [375] "180"  
## [377] "181"  
## [379] "182"  
## [381] "183"  
## [383] "184"  
## [385] "185"  
## [387] "186"  
## [389] "187"  
## [391] "188"  
## [393] "189"  
## [395] "190"  
## [397] "191"  
## [399] "192"  
## [401] "193"  
## [403] "194"  
## [405] "195"  
## [407] "196"  
## [409] "197"  
## [411] "198"  
## [413] "199"  
## [415] "200"  
## [417] "201"  
## [419] "202"  
## [421] "203"  
## [423] "204"  
## [425] "205"  
## [427] "206"  
## [429] "207"  
## [431] "208"  
## [433] "209"  
## [435] "210"  
## [437] "211"  
## [439] "212"  
## [441] "213"  
## [443] "214"  
## [445] "215"  
## [447] "216"  
## [449] "217"  
## [451] "218"  
## [453] "219"  
## [455] "220"  
## [457] "221"  
## [459] "222"  
## [461] "223"  
## [463] "224"  
## [465] "225"  
## [467] "226"  
## [469] "227"  
## [471] "228"  
## [473] "229"  
## [475] "230"  
## [477] "231"  
## [479] "232"  
## [481] "233"  
## [483] "234"  
## [485] "235"  
## [487] "236"  
## [489] "237"  
## [491] "238"  
## [493] "239"  
## [495] "240"  
## [497] "241"  
## [499] "242"  
## [501] "243"  
## [503] "244"  
## [505] "245"  
## [507] "246"  
## [509] "247"  
## [511] "248"  
## [513] "249"  
## [515] "250"  
## [517] "251"  
## [519] "252"  
## [521] "253"  
## [523] "254"  
## [525] "255"  
## [527] "256"  
## [529] "257"  
## [531] "258"  
## [533] "259"  
## [535] "260"  
## [537] "261"  
## [539] "262"  
## [541] "263"  
## [543] "264"  
## [545] "265"  
## [547] "266"  
## [549] "267"  
## [551] "268"  
## [553] "269"  
## [555] "270"  
## [557] "271"  
## [559] "272"  
## [561] "273"  
## [563] "274"  
## [565] "275"  
## [567] "276"  
## [569] "277"  
## [571] "278"  
## [573] "279"  
## [575] "280"  
## [577] "281"  
## [579] "282"  
## [581] "283"  
## [583] "284"  
## [585] "285"  
## [587] "286"  
## [589] "287"  
## [591] "288"  
## [593] "289"  
## [595] "290"  
## [597] "291"  
## [599] "292"  
## [601] "293"  
## [603] "294"  
## [605] "295"  
## [607] "296"  
## [609] "297"  
## [611] "298"  
## [613] "299"  
## [615] "300"  
## [617] "301"  
## [619] "302"  
## [621] "303"  
## [623] "304"  
## [625] "305"  
## [627] "306"  
## [629] "307"  
## [631] "308"  
## [633] "309"  
## [635] "310"  
## [637] "311"  
## [639] "312"  
## [641] "313"  
## [643] "314"  
## [645] "315"  
## [647] "316"  
## [649] "317"  
## [651] "318"  
## [653] "319"  
## [655] "320"  
## [657] "321"  
## [659] "322"  
## [661] "323"  
## [663] "324"  
## [665] "325"  
## [667] "326"  
## [669] "327"  
## [671] "328"  
## [673] "329"  
## [675] "330"  
## [677] "331"  
## [679] "332"  
## [681] "333"  
## [683] "334"  
## [685] "335"  
## [687] "336"  
## [689] "337"  
## [691] "338"  
## [693] "339"  
## [695] "340"  
## [697] "341"  
## [699] "342"  
## [701] "343"  
## [703] "344"  
## [705] "345"  
## [707] "346"  
## [709] "347"  
## [711] "348"  
## [713] "349"  
## [715] "350"  
## [717] "351"  
## [719] "352"  
## [721] "353"  
## [723] "354"  
## [725] "355"  
## [727] "356"  
## [729] "357"  
## [731] "358"  
## [733] "359"  
## [735] "360"  
## [737] "361"  
## [739] "362"  
## [741] "363"  
## [743] "364"  
## [745] "365"  
## [747] "366"  
## [749] "367"  
## [751] "368"  
## [753] "369"  
## [755] "370"  
## [757] "371"  
## [759] "372"  
## [761] "373"  
## [763] "374"  
## [765] "375"  
## [767] "376"  
## [769] "377"  
## [771] "378"  
## [773] "379"  
## [775] "380"  
## [777] "381"  
## [779] "382"  
## [781] "383"  
## [783] "384"  
## [785] "385"  
## [787] "386"  
## [789] "387"  
## [791] "388"  
## [793] "389"  
## [795] "390"  
## [797] "391"  
## [799] "392"  
## [801] "393"  
## [803] "394"  
## [805] "395"  
## [807] "396"  
## [809] "397"  
## [811] "398"  
## [813] "399"  
## [815] "400"  
## [817] "401"  
## [819] "402"  
## [821] "403"  
## [823] "404"  
## [825] "405"  
## [827] "406"  
## [829] "407"  
## [831] "408"  
## [833] "409"  
## [835] "410"  
## [837] "411"  
## [839] "412"  
## [841] "413"  
## [843] "414"  
## [845] "415"  
## [847] "416"  
## [849] "417"  
## [851] "418"  
## [853] "419"  
## [855] "420"  
## [857] "421"  
## [859] "422"  
## [861] "423"  
## [863] "424"  
## [865] "425"  
## [867] "426"  
## [869] "427"  
## [871] "428"  
## [873] "429"  
## [875] "430"  
## [877] "431"  
## [879] "432"  
## [881] "433"  
## [883] "434"  
## [885] "435"  
## [887] "436"  
## [889] "437"  
## [891] "438"  
## [893] "439"  
## [895] "440"  
## [897] "441"  
## [899] "442"  
## [901] "443"  
## [903] "444"  
## [905] "445"  
## [907] "446"  
## [909] "447"  
## [911] "448"  
## [913] "449"  
## [915] "450"  
## [917] "451"  
## [919] "452"  
## [921] "453"  
## [923] "454"  
## [925] "455"  
## [927] "456"  
## [929] "457"  
## [931] "458"  
## [933] "459"  
## [935] "460"  
## [937] "461"  
## [939] "462"  
## [941] "463"  
## [943] "464"  
## [945] "465"  
## [947] "466"  
## [949] "467"  
## [951] "468"  
## [953] "469"  
## [955] "470"  
## [957] "471"  
## [959] "472"  
## [961] "473"  
## [963] "474"  
## [965] "475"  
## [967] "476"  
## [969] "477"  
## [971] "478"  
## [973] "479"  
## [975] "480"  
## [977] "481"  
## [979] "482"  
## [981] "483"  
## [983] "484"  
## [985] "485"  
## [987] "486"  
## [989] "487"  
## [991] "488"  
## [993] "489"  
## [995] "490"  
## [997] "491"  
## [999] "492"  
## [1001] "493"  
## [1003] "494"  
## [1005] "495"  
## [1007] "496"  
## [1009] "497"  
## [1011] "498"  
## [1013] "499"  
## [1015] "500"  
## [1017] "501"  
## [1019] "502"  
## [1021] "503"  
## [1023] "504"  
## [1025] "505"  
## [1027] "506"  
## [1029] "507"  
## [1031] "508"  
## [1033] "509"  
## [1035] "510"  
## [1037] "511"  
## [1039] "512"  
## [1041] "513"  
## [1043] "514"  
## [1045] "515"  
## [1047] "516"  
## [1049] "517"  
## [1051] "518"  
## [1053] "519"  
## [1055] "520"  
## [1057] "521"  
## [1059] "522"  
## [1061] "523"  
## [1063] "524"  
## [1065] "525"  
## [1067] "526"  
## [1069] "527"  
## [1071] "528"  
## [1073] "529"  
## [1075] "530"  
## [1077] "531"  
## [1079] "532"  
## [1081] "533"  
## [1083] "534"  
## [1085] "535"  
## [1087] "536"  
## [1089] "537"  
## [1091] "538"  
## [1093] "539"  
## [1095] "540"  
## [1097] "541"  
## [1099] "542"  
## [1101] "543"  
## [1103] "544"  
## [1105] "545"  
## [1107] "546"  
## [1109] "547"  
## [1111] "548"  
## [1113] "549"  
## [1115] "550"  
## [1117] "551"  
## [1119] "552"  
## [1121] "553"  
## [1123] "554"  
## [1125] "555"  
## [1127] "556"  
## [1129] "557"  
## [1131] "558"  
## [1133] "559"  
## [1135] "560"  
## [1137] "561"  
## [1139] "562"  
## [1141] "563"  
## [1143] "564"  
## [1145] "565"  
## [1147] "566"  
## [1149] "567"  
## [1151] "568"  
## [1153] "569"  
## [1155] "570"  
## [1157] "571"  
## [1159] "572"  
## [1161] "573"  
## [1163] "574"  
## [1165] "575"  
## [1167] "576"  
## [1169] "577"  
## [1171] "578"  
## [1173] "579"  
## [1175] "580"  
## [1177] "581"  
## [1179] "582"  
## [1181] "583"  
## [1183] "584"  
## [1185] "585"  
## [1187] "586"  
## [1189] "587"  
## [1191] "588"  
## [1193] "589"  
## [1195] "590"  
## [1197] "591"  
## [1199] "592"  
## [1201] "593"  
## [1203] "594"  
## [1205] "595"  
## [1207] "596"  
## [1209] "597"  
## [1211] "598"  
## [1213] "599"  
## [1215] "600"  
## [1217] "601"  
## [1219] "602"  
## [1221] "603"  
## [1223] "604"  
## [1225] "605"  
## [1227] "606"  
## [1229] "607"  
## [1231] "608"  
## [1233] "609"  
## [1235] "610"  
## [1237] "611"  
## [1239] "612"  
## [1241] "613"  
## [1243] "614"  
## [1245] "615"  
## [1247] "616"  
## [1249] "617"  
## [1251] "618"  
## [1253] "619"  
## [1255] "620"  
## [1257] "621"  
## [1259] "622"  
## [1261] "623"  
## [1263] "624"  
## [1265] "625"  
## [1267] "626"  
## [1269] "627"  
## [1271] "628"  
## [1273] "629"  
## [1275] "630"  
## [1277] "631"  
## [1279] "632"  
## [1281] "633"  
## [1283] "634"  
## [1285] "635"  
## [1287] "636"  
## [1289] "637"  
## [1291] "638"  
## [1293] "639"  
## [1295] "640"  
## [1297] "641"  
## [1299] "642"  
## [1301] "643"  
## [1303] "644"  
## [1305] "645"  
## [1307] "646"  
## [1309] "647"  
## [1311] "648"  
## [1313] "649"  
## [1315] "650"  
## [1317] "651"  
## [1319] "652"  
## [1321] "653"  
## [1323] "654"  
## [1325] "655"  
## [1327] "656"  
## [1329] "657"  
## [1331] "658"  
## [1333] "659"  
## [1335] "660"  
## [1337] "661"  
## [1339] "662"  
## [1341] "663"  
## [1343] "664"  
## [1345] "665"  
## [1347] "666"  
## [1349] "667"  
## [1351] "668"  
## [1353] "669"  
## [1355] "670"  
## [1357] "671"  
## [1359] "672"  
## [1361] "673"  
## [1363] "674"  
## [1365] "675"  
## [1367] "676"  
## [1369] "677"  
## [1371] "678"  
## [1373] "679"  
## [1375] "680"  
## [1377] "681"  
## [1379] "682"  
## [1381] "683"  
## [1383] "684"  
## [1385] "685"  
## [1387] "686"  
## [1389] "687"  
## [1391] "688"  
## [1393] "689"  
## [1395] "690"  
## [1397] "691"  
## [1399] "692"  
## [1401] "693"  
## [1403] "694"  
## [1405] "695"  
## [1407] "696"  
## [1409] "697"  
## [1411] "698"  
## [1413] "699"  
## [1415] "700"  
## [1417] "701"  
## [1419] "702"  
## [1421] "703"  
## [1423] "704"  
## [1425] "705"  
## [1427] "706"  
## [1429] "707"  
## [1431] "708"  
## [1433] "709"  
## [1435] "710"  
## [1437] "711"  
## [1439] "712"  
## [1441] "713"  
## [1443] "714"  
## [1445] "715"  
## [1447] "716"  
## [1449] "717"  
## [1451] "718"  
## [1453] "719"  
## [1455] "720"  
## [1457] "721"  
## [1459] "722"  
## [1461] "723"  
## [1463] "724"  
## [1465] "725"  
## [1467] "726"  
## [1469] "727"  
## [1471] "728"  
## [1473] "729"  
## [1475] "730"  
## [1477] "731"  
## [1479] "732"  
## [1481] "733"  
## [1483] "734"  
## [1485] "735"  
## [1487] "736"  
## [1489] "737"  
## [1491] "738"  
## [1493] "739"  
## [1495] "740"  
## [1497] "741"  
## [1499] "742"  
## [1501] "743"  
## [1503] "744"  
## [1505] "745"  
## [1507] "746"  
## [1509] "747"  
## [1511] "748"  
## [1513] "749"  
## [1515] "750"  
## [1517] "751"  
## [1519] "752"  
## [1521] "753"  
## [1523] "754"  
## [1525] "755"  
## [1527] "756"  
## [1529] "757"  
## [1531] "758"  
## [1533] "759"  
## [1535] "760"  
## [1537] "761"  
## [1539] "762"  
## [1541] "763"  
## [1543] "764"  
## [1545] "765"  
## [1547] "766"  
## [1549] "767"  
## [1551] "768"  
## [1553] "769"  
## [1555] "770"  
## [1557] "771"  
## [1559] "772"  
## [1561] "773"  
## [1563] "774"  
## [1565] "775"  
## [1567] "776"  
## [1569] "777"  
## [1571] "778"  
## [1573] "779"  
## [1575] "780"  
## [1577] "781"  
## [1579] "782"  
## [1581] "783"  
## [1583] "784"  
## [1585] "785"  
## [1587] "786"  
## [1589] "787"  
## [1591] "788"  
## [1593] "789"  
## [1595] "790"  
## [1597] "791"  
## [1599] "792"  
## [1601] "793"  
## [1603] "794"  
## [1605] "795"  
## [1607] "796"  
## [1609] "797"  
## [1611] "798"  
## [1613] "799"  
## [1615] "800"  
## [1617] "801"  
## [1619] "802"  
## [1621] "803"  
## [1623] "804"  
## [1625] "805"  
## [1627] "806"  
## [1629] "807"  
## [1631] "808"  
## [1633] "809"  
## [1635] "810"  
## [1637] "811"  
## [1639] "812"  
## [1641] "813"  
## [1643] "814"  
## [1645] "815"  
## [1647] "816"  
## [1649] "817"  
## [1651] "818"  
## [1653] "819"  
## [1655] "820"  
## [1657] "821"  
## [1659] "822"  
## [1661] "823"  
## [1663] "824"  
## [1665] "825"  
## [1667] "826"  
## [1669] "827"  
## [1671] "828"  
## [1673] "829"  
## [1675] "830"  
## [1677] "831"  
## [1679] "832"  
## [1681] "833"  
## [1683] "834"  
## [1685] "835"  
## [1687] "836"  
## [1689] "837"  
## [1691] "838"  
## [1693] "839"  
## [1695] "840"  
## [1697] "841"  
## [1699] "842"  
## [1701] "843"  
## [1703] "844"  
## [1705] "845"  
## [1707] "846"  
## [1709] "847"  
## [1711] "848"  
## [1713] "849"  
## [1715] "850"  
## [1717] "851"  
## [1719] "852"  
## [1721] "853"  
## [1723] "854"  
## [1725] "855"  
## [1727] "856"  
## [1729] "857"  
## [1731] "858"  
## [1733] "859"  
## [1735] "860"  
## [1737] "861"  
## [1739] "862"  
## [1741] "863"  
## [1743] "864"  
## [1745] "865"  
## [1747] "866"  
## [1749] "867"  
## [1751] "868"  
## [1753] "869"  
## [1755] "870"  
## [1757] "871"  
## [1759] "872"  
## [1761] "873"  
## [1763] "874"  
## [1765] "875"  
## [1767] "876"  
## [1769] "877"  
## [1771] "878"  
## [1773] "879"  
## [1775] "880"  
## [1777] "881"  
## [1779] "882"  
## [1781] "883"  
## [1783] "884"  
## [1785] "885"  
## [1787] "886"  
## [1789] "887"  
## [1791] "888"  
## [1793] "889"  
## [1795] "890"  
## [1797] "891"  
## [1799] "892"  
## [1801] "893"  
## [1803] "894"  
## [1805] "895"  
## [1807] "896"  
## [1809] "897"  
## [1811] "898"  
## [1813] "899"  
## [1815] "900"  
## [1817] "901"  
## [1819] "902"  
## [1821] "903"  
## [1823] "904"  
## [1825] "905"  
## [1827] "906"  
## [1829] "907"  
## [1831] "908"  
## [1833] "909"  
## [1835] "910"  
## [1837] "911"  
## [1839] "912"<
```

Group Data by Pattern

- Now that we have isolated some of the main elements of our data as a vector, we want to group data by type.
- We can use pattern-matching to separate strings by their pattern combining `stringr` and regular expressions.
- Here, we can separate the stats by using a pattern that finds elements with at least one lower-case letter

```
split <- str_split(match_stats, "\n")[[1]] # Save split vector  
pattern <- "[a-z]"  
  
stats <- str_subset(split, pattern) # Subset players and stat names
```

Using RegEx to Sort Data

We use exclusion to get all the other values

```
values <- split[  
  !str_detect(split, pattern) &  
  !str_detect(split, "[A-Z]")  
] # Get values
```

Structuring Data Frame

We notice that some stats have just counts while others have percentages and ratios. We can deal with this by flagging counts and expanding the data frame based on the condition of being a count or percentage stat.

```
counts <- stats %in% c("Aces",
  "Double Faults",
  "Service Games Played",
  "Return Games Played")

result <- data.frame(
  stat = rep(stats, ifelse(counts, 2, 4)),
  values = values
)

result
```

```
##                                     stat   values
## 1                               Aces      7
## 2                               Aces      1
## 3             Double Faults      2
## 4             Double Faults      3
## 5           1st Serve    71%
## 6           1st Serve  0.71/151
```

String Substitution

We will need to do some more tidying of the strings to get our `value` column into numeric values. String replace will be a big help. Here are some examples of removing percentage signs and parentheses using `str_replace`.

```
# We use 'all' to replace all instances
# The escapes \\ make sure () are treated as fixed
str_replace_all(values, "[\\(\\%\\)]", "")
```

```
## [1] "7"      "1"      "2"      "3"      "71"     "32/45"  "64"
## [8] "43/67"  "84"     "27/32"  "70"     "30/43"  "69"     "9/13"
## [15] "38"     "9/24"   "0"      "0/0"    "40"     "2/5"    "9"
## [22] "10"     "30"     "13/43"  "16"     "5/32"   "63"     "15/24"
## [29] "31"     "4/13"   "60"     "3/5"    "0"      "0/0"    "10"
## [36] "9"      "80"     "36/45"  "58"     "39/67"  "42"     "28/67"
## [43] "20"     "9/45"   "57"     "64/112" "43"     "48/112"
```

Reshaping a Data Frame

We frequently need to go between wide and long formats of our data. With `tidyverse` we can use `gather` to go from wide to long and `spread` to go from long to wide.

Below is an example of reshaping our `match_stats` going from long to wide.

```
library(tidyverse)

# Assign player to their stat
result <- result %>%
  group_by(stat) %>%
  dplyr::mutate(
    player = rep(c(1, 2), each = n() / 2)
  )

result %>%
  filter(stat %in% c("Aces", "Double Faults")) %>%
  tidyverse::spread(stat, values)
```

```
## # A tibble: 2 × 3
##   player    Aces `Double Faults`
## * <dbl> <fctr>          <fctr>
## 1     1      7             2
## 2     2      1             3
```

Exploring Our Data

Data Exploration

Exploratory data analysis is detective work-numerical detective work-or counting detective work-or graphical detective work...Exploratory data analysis can never be the whole story, but nothing else can serve as the foundation stone--as the first step. - John Tukey

EDA Starts By Asking Questions

- All exploration needs some direction
- Exploring aimlessly wastes time and rarely will get you where you need to be
- To avoid aimless exploration, you need to ask yourself what you are looking for? and what would be interesting to find?
- This means getting to know your variables with summary functions like: `str`, `table`, and `summary`



Exploratory Charting

Exploratory Charting

- Most of our describing and exploration in R happens with graphics

Exploratory Charting

- Most of our describing and exploration in R happens with graphics
- A powerful package for graphing is ggplot2

Exploratory Charting

- Most of our describing and exploration in R happens with graphics
- A powerful package for graphing is ggplot2
- ggplot2 provides a grammar for building univariate, bivariate, and lattice plots

Exploratory Charting

- Most of our describing and exploration in R happens with graphics
- A powerful package for graphing is ggplot2
- ggplot2 provides a grammar for building univariate, bivariate, and lattice plots
- Also, many specialty graphics packages build on ggplot2

Overview of ggplot2

To install, use `install.packages('ggplot2')`.

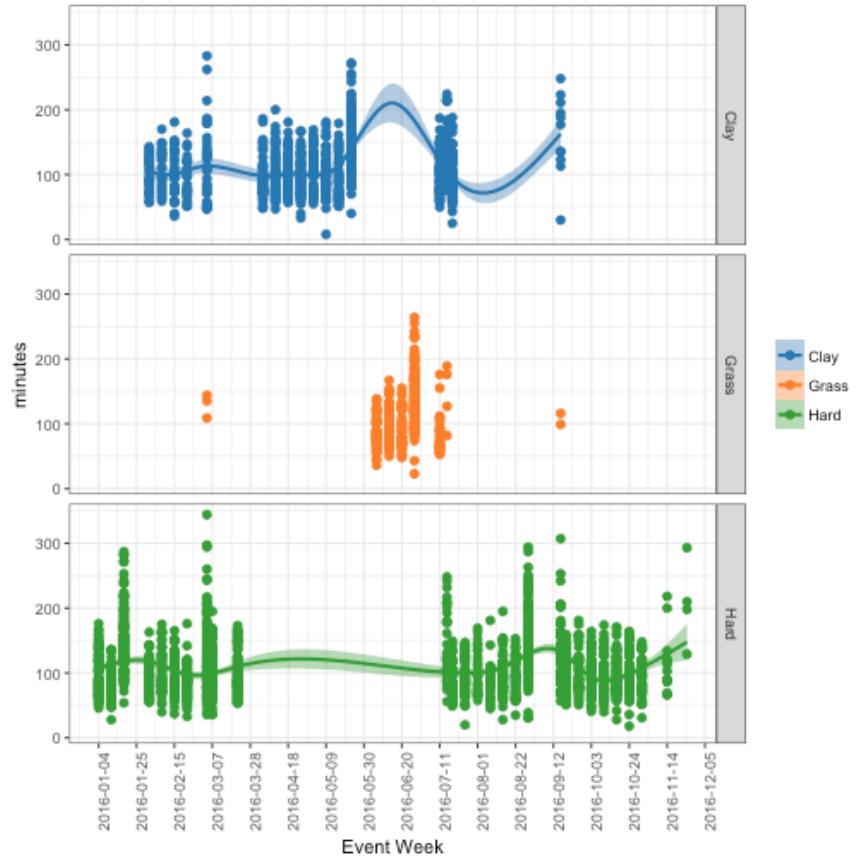
Function.Type	Description
aes	Relates variables to axes and aesthetic elements of our plot
geoms	Define how the variables will be displayed, that is, the type of chart
<hr/>	
Function.Type	Description
facet	Splits plot into rows and columns defined by a group variable
scales	Customizes the range and limits of aesthetics
themes	Further controls of the style and elements of the chart

Example: Match Duration Trends by Surface

In this example, we explore patterns in match duration by surface using a number of features in `ggplot2` and `ggthemes`.

```
data(atp_matches)
```

```
atp_matches %>%
  filter(!(minutes == 0 | minutes > 10 * 60), surface != "Carpet") %>
  ggplot(aes(y = minutes, x = tourney_date,
             col = surface, fill = surface)) +
  facet_grid(surface ~ .) + # Facet by row with no column facet
  geom_smooth() +
  geom_point(size = 2) +
  scale_color_tableau(name = ""),
  scale_fill_manual(name = "",
                    values = alpha(tableau_color_pal()(3), 0.5)) +
  scale_x_date("Event Week", date_breaks = "3 week") +
  theme_bw() +
  theme(axis.text.x = element_text(angle = 90))
```



3D Charting

- So far, we have only considered two-dimensional charts.
- For 3D data, we can use the package `plot3D`.
- But first, we need to extract the (x, y, z) coordinates for each serve along a set of time points.

Plotting Tennis Shots

In this example, we will look at some tracking data for about 60 tennis serves of ATP players. The data includes the x, y, and z locations of the flight of the serves.

```
positions[1:10,] # Serve position data
```

```
##      seconds          x          y          z id
## 1  1.532434 -11.39293  0.7152673 3.089983  1
## 2  1.532823 -11.36955  0.7105796 3.086929  1
## 3  1.533213 -11.34617  0.7058936 3.083875  1
## 4  1.533602 -11.32280  0.7012094 3.080821  1
## 5  1.533991 -11.29944  0.6965270 3.077767  1
## 6  1.534381 -11.27608  0.6918463 3.074713  1
## 7  1.534770 -11.25274  0.6871674 3.071659  1
## 8  1.535159 -11.22940  0.6824902 3.068605  1
## 9  1.535549 -11.20607  0.6778148 3.065551  1
## 10 1.535938 -11.18275  0.6731411 3.062497  1
```

Package plot3D

- The `plot3D` package provides a suite a functions for 3D plots: `points3D`, `lines3D`, `ribbon3D`, etc.
- You can install with `install.packages('plot3D')`

Plot3D Key Arguments

In addition to the basic variables x , y , z , here are some useful arguments when working with `plot3D`.

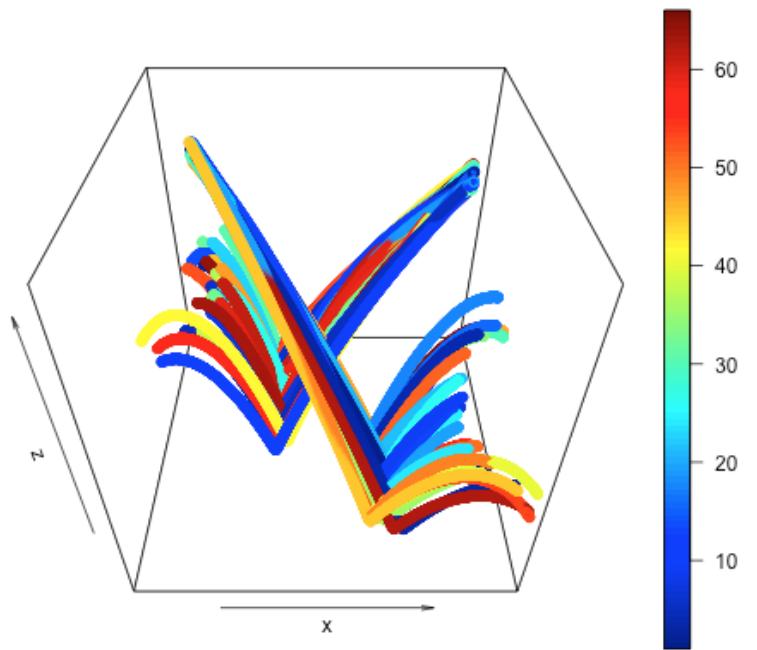
Argument	Description
colvar	Numeric group variable for coloring points
phi	Angle defining colatitude
theta	Angle defining azimuthal direction
col	Colors to use with colvar
bty	Style for the perspective box

Plotting Serves

In this chart, we plot all serves using the `points3D`. We set `theta` to 0 so that we look perpendicularly along the length of the court.

```
library(plot3D) # Load plot3D

points3D(
  x = positions$x,
  y = positions$y,
  z = positions$z,
  colvar = positions$id,
  theta = 0
)
```



Prediction Models

Sport Prediction



Predicting *who will win* is the holy grail of sports analytics and is a major research area of interest for sports statisticians.

Descriptive vs Predictive Models

Descriptive vs Predictive Models

- *Predictive modeling* usually refers to machine learning

Descriptive vs Predictive Models

- *Predictive modeling* usually refers to machine learning
- The goal of these techniques is to improve predictive performance

Descriptive vs Predictive Models

- *Predictive modeling* usually refers to machine learning
- The goal of these techniques is to improve predictive performance
- This is a very different goal from statistical models, like regression, where inference and interpretation are of primary importance

Descriptive vs Predictive Models

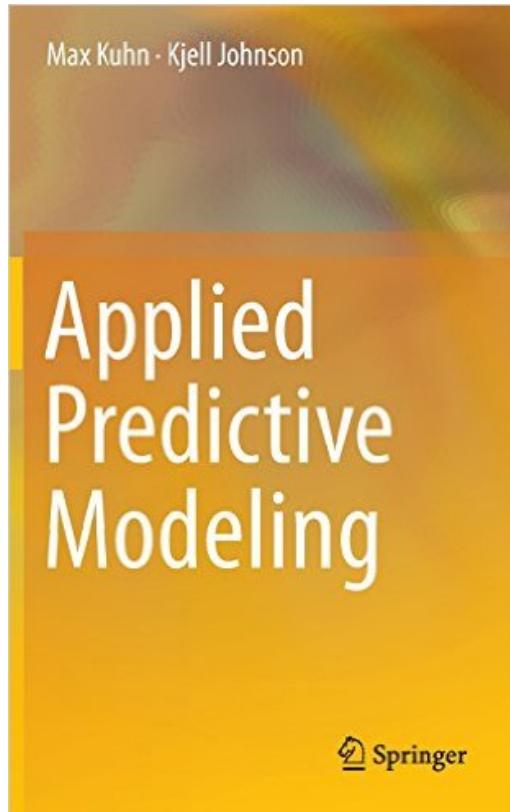
- *Predictive modeling* usually refers to machine learning
- The goal of these techniques is to improve predictive performance
- This is a very different goal from statistical models, like regression, where inference and interpretation are of primary importance
- Predictive models will often sacrifice interpretability for improved performance

Describe Before You Predict

- Because it is often a challenge to interpret the "how" of machine learning methods, it is good practice to first do some statistical modelling (e.g. `glm(y ~ x + ...)`)
- The reasons for this include:
 - Getting familiar with the interrelationships in your data
 - Identifying any issues not flagged during EDA
 - Developing some expectations for the predictive modelling results

Predictive Modelling with caret

- Once you are ready to develop your prediction model, the *caret* package, written by Max Kuhn, is a great resource for machine learning in R
- *caret* stand for Classification And REgression Training
- It includes 100+ predictive modelling methods
- It provides a unified & streamlines interface for building and evaluating models
- It allows parallel processing for faster computation
- You can install with
`install.packages('caret')`



Modelling Steps

Modelling Steps

Step 1. Decide how to spend your data

Modelling Steps

Step 1. Decide how to spend your data

Step 2. Split the data into training and test

Modelling Steps

Step 1. Decide how to spend your data

Step 2. Split the data into training and test

Step 3. Conduct pre-processing (as needed)

Modelling Steps

- Step 1. Decide how to spend your data
- Step 2. Split the data into training and test
- Step 3. Conduct pre-processing (as needed)
- Step 4. Train the model with resampling

Modelling Steps

- Step 1. Decide how to spend your data
- Step 2. Split the data into training and test
- Step 3. Conduct pre-processing (as needed)
- Step 4. Train the model with resampling
- Step 5. Evaluate the model with test data

Why Split Up the Data at All?

- Testing our data on independent samples is the strongest form of validation and evaluation
- If we trained and tested on the same data, we risk *overfitting* which is the machine-learning equivalent to a "Monday morning quarterback"
- We also resample among training for additional protection against overfitting



Using `createDataPartition`

We can use `createDataPartition` to split our data:

```
createDataPartition(y, times, p, list, ...)
```

Argument	Description
y	Outcome vector to balance sampling on
times	Number of partitions
p	Proportion of each partition allocated to training
list	Logical whether list is returned

Note: Using 70% of our data for training is typical

Example: Partitioning Data

In this example, we create one partition with 70% of our dataset allocated to training.

```
library(caret) # Load caret

# Returns matrix of indices for obs in training
train <- createDataPartition(
  y = data$outcome,
  times = 1,
  p = 0.7,
  list = F
)
```

Pre-Processing

Before we split our data, we need to pre-process our data. The pre-processing can protect against some loss in model accuracy due to scale, skew, or high correlation.

Common pre-processing steps are:

1. Centering - Give all variables a common mean of 0
2. Standardizing - Give all variables a common scale
3. Remove highly correlated variables
4. Reduce dimension (when $n \sim p$)

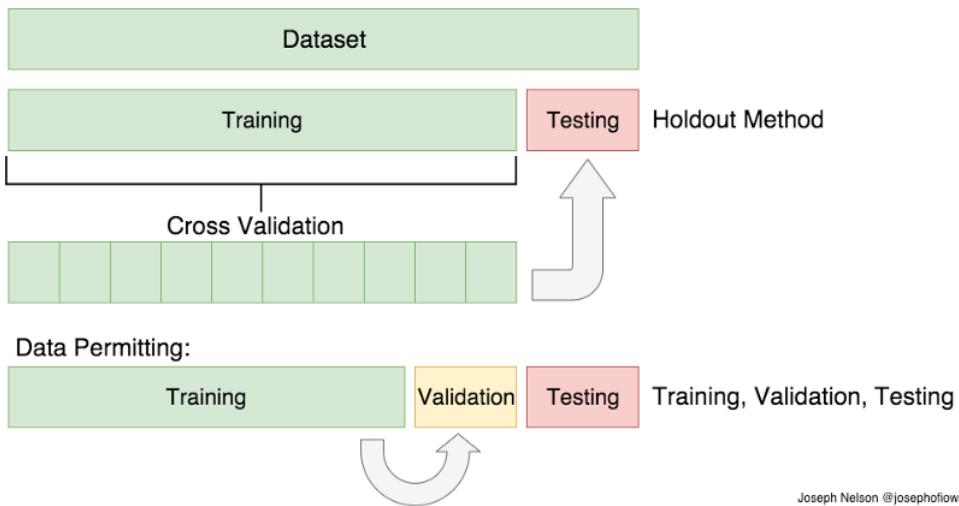
The `train` Function

The main workhouse function for model training in `caret` is `train`. Here are the main arguments you need to know to get started.

Argument	Description
form	Formula ($y \sim x$)
data	Data frame of training data
method	Character of the ML method to be used
Argument	Description
metric	Performance metric for summarizing
tuneLength	Sets granularity for tuning parameter if <code>tuneGrid</code> not specified
trControl	Control parameters for training
tuneGrid	Data frame that gives explicit range for tuning parameters

Training Control

- The `trControl` argument is a list that controls a number of aspects of training including resampling and how we summarise performance with each resample.
- The resampling is an important additional measure to protect against overfitting when training



Example: Using `trControl`

In our example, we will use `trControl` to use 5-fold cross validation and a two-class summary for our performance measures.

```
ctrlSpecs <- trainControl(  
  method = "repeatedCV",  
  repeats = 5,  
  summaryFunction = twoClassSummary, # function from caret  
  classProbs = TRUE # Needed to use twoClassSummary  
)
```

Grid Tuning

- The `tuneGrid` is a way to give a specific grid for the tuning parameters of the method
- You can use `expand.grid` to make a range of parameters
- To determine the parameters to tune and their variable names you can use `getModelInfo`

Using getModelInfo

- We can get a model of interest with `getModelInfo('model')` or get info from all models with `getModelInfo()`.
- This returns a list per model with information about the model type, parameters, grid function, etc.
- Here is an example with the `rpart` model.

```
getModelInfo('rpart')[['rpart']][1:4]
```

```
## $label
## [1] "CART"
##
## $library
## [1] "rpart"
##
## $type
## [1] "Regression"      "Classification"
##
## $parameters
##   parameter    class           label
## 1             cp numeric Complexity Parameter
```

Performance Metrics

- We can use the `metric` argument to choose our performance metric for training evaluation
- There are many metrics for evaluating classification. In general, it is best to choose one when choosing among model approaches

Metric	Description
Accuracy	Proportion of exactly correct classifications
AUC	Area under the ROC curve
Sensitivity	The true positive rate (also called 'recall')
Specificity	The true negative rate
LogLoss	Prediction-weighted loss function

Setting Performance Metric

There is no one correct performance measure. In fact, multiple should be evaluated when testing. For training, the "log loss" is good all around measure. Here is how we can set our control specs to use it.

```
ctrlSpecs <- trainControl(  
  method = "repeatedCV",  
  repeats = 5,  
  summaryFunction = mnLogLoss,  
  classProbs = TRUE  
)
```

Models

There are more than 100 models to try in `caret`! Below is just a sample of some popular kinds. The full list is available with the online [caret book](#).

Category	Description	Examples
Forest	Ensemble of multiple decision trees with bagging (bootstrap aggregation)	<code>rf</code> , <code>rfRules</code> , <code>cforest</code>
Category	Incremental building of multiple classifiers, which is a kind of correlated ensembling	<code>gbm</code> , <code>ada</code> , <code>best</code> , <code>C5.0</code>
Support Vector Machines	Collection of regression lines that try to maximally separate classes	<code>svmLinear</code> , <code>svmRadial</code>

Random Forest

Let's have a look at each category and how we could train each in caret.

Below we use the `rf` method to fit a random forest. The `tuneLength` is set to 10 to have a randomly generated grid for the forest parameters.

```
set.seed(1115)

rfFit <- train(
  outcome ~ .,
  data = trainingData,
  method = "rf",
  tuneLength = 10,
  trControl = ctrlSpecs,
  metric = "logLoss"
)
```

Evaluate the Model

- We can see the results across the different tuning parameters using:
`print` or `plot`
- The selected model can be extracted with `finalModel` and it will have all the class properties of the source method

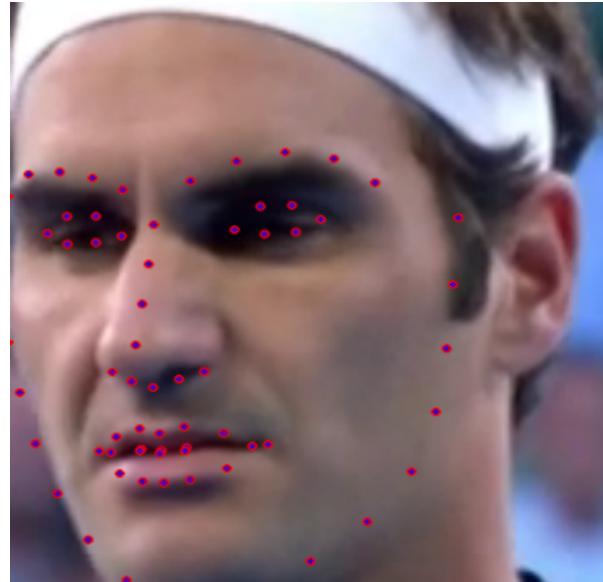
```
library(randomForest) # Class methods for RF  
importance(rfFit$finalModel) # Variable importance
```

Application: Emotion Detection

At GIG we have recently been using caret to create prediction models for player emotions based on Facial Action Units in match broadcasts.

How much does this expression indicate:

- Anxiety?
- Determination?
- Excitement?



Blogging

Why Blog?

Sharing your ideas with the world is one of the most joyful parts of sports analysis.



'On The T' Powered By Blogdown

on-the-t.com

STATS ON THE T

DEDICATED TO DATA, STATISTICS, AND TENNIS

ABOUT | **CONTACT** | **FEATURES** | **RESEARCH** | **SOFTWARE**



French Open ATP Leaders

⌚ June 11, 2017

After one shocking title winner took the ladies' crown on Saturday at the French Open, the tennis world is abuzz with speculation about whether Stan Wawrinka will be able to deny Rafael Nadal a tenth title when they meet in the men's final this afternoon? Embed from Getty Images Nadal already holds the highest record of titles won at Roland Garros, making him the definitive King of Clay. In contrast with Wawrinka, Nadal has also had the much easier course to the final, having not dropped a single set during this year's event and getting something of a break in the quarterfinal with the retirement of Pablo Carreno Busta.



French Open WTA Leaders

⌚ June 10, 2017

We are just hours away from the women's final at the 2017 French Open. After 126 matches, 25 year-old Simona Halep and 20 year-old Jelena Ostapenko are the last women's players

RECENT POSTS

- French Open ATP Leaders
- French Open WTA Leaders
- Rethinking Women's French Open Seedings
- Clay Court Performance Trends of French Open Title Men's Top 4
- Return Leaders Among ATP's Top Contenders for French Open Title
- Serve Leaders Among ATP's Top Contenders for French Open Title
- WTA Tennis Rankings Aren't Working
- Should We Be Surprised by Murray's 2017 Clay Performance So Far?

80 / 92

Getting There

- Starting a blog might seem overwhelming but, if you know R and some basic CSS, it's easy
- Producing interesting content is the tough part about blogging
- R helps to let you focus on *content* by providing easy-to-use tools that solve the technical details of going from code to the Web
- *If you can write R markdown, you can create a blog!*



Blogdown

Blogdown

- Package for writing a blog in R

Blogdown

- Package for writing a blog in R
- Create content with R markdown

Blogdown

- Package for writing a blog in R
- Create content with R markdown
- Generate site with Hugo, a static site generator

Blogdown

- Package for writing a blog in R
- Create content with R markdown
- Generate site with Hugo, a static site generator
- Deploy on Github, Netlify, or other hosting services

Blogdown

- Package for writing a blog in R
- Create content with R markdown
- Generate site with Hugo, a static site generator
- Deploy on Github, Netlify, or other hosting services
- Authored by Yihui Xie (who is at the conference!)

Getting Started

Here are the major functions you will use with `blogdown`.

Function	Description
<code>install_hugo</code>	Downloads and installs Hugo
<code>install_theme</code>	Downloads a Hugo theme from github
<code>build_site</code>	Compiles Rmd files and builds the site
<code>hugo_page</code>	Renders an Rmd file as an HTML that can be read by Hugo
<code>hugo_cmd</code>	Run Hugo commands
<code>new_content</code>	Creates new file in working directory
<code>new_site</code>	Creates environment for new site
<code>serve_site</code>	Preview working version of your site

Learning the Workflow

Yihui and Amber Thomas have a draft manual for using blogdown that will
get you started.

The screenshot shows the 'Preface' page of the 'blogdown: Creating Websites with R Markdown' book. The left sidebar contains a table of contents with sections like 'Software information and conventions', 'About the Authors', 'Yihui Xie', 'Amber Thomas', '1 Get Started' (with sub-sections 1.1 through 1.7), '2 Hugo' (with sub-sections 2.1 through 2.3), and 'Appendix'. The main content area has a header 'blogdown: Creating Websites with R Markdown', author information 'Yihui Xie & Amber Thomas', date '2017-06-13', and a 'Preface' section. The 'Preface' section includes a warning about the book being under development and a quote from Carlos Scheidegger. A note at the bottom credits Carlos's words for resonating well with the authors.

Preface

Software information and conventions

About the Authors

Yihui Xie

Amber Thomas

1 Get Started

1.1 Installation

1.2 A quick example

1.3 RStudio IDE

1.4 Global options

1.5 R Markdown vs Markdown

1.6 Other themes

1.7 A recommended workflow

2 Hugo

2.1 Static sites and Hugo

2.2 Configuration

2.2.1 TOML Syntax

2.2.2 Options

2.3 Content

blogdown: Creating Websites with R Markdown

Yihui Xie & Amber Thomas

2017-06-13

Preface

WARNING: this book is still under development. It will be updated from day to day. The blogdown package is still a beta version, so please use both the package and this book with caution until this warning is removed. We will try not to introduce breaking changes in blogdown from now on, but there is no guarantee.

In the summer of 2012, I did my internship at the AT&T Labs Research,¹ where I attended a talk given by Carlos Scheidegger (<https://cscheid.net>), and Carlos said something along the lines “if you don’t have a website nowadays, you don’t exist.” Later I paraphrased it as:

“I web, therefore I am a-spiderman.”

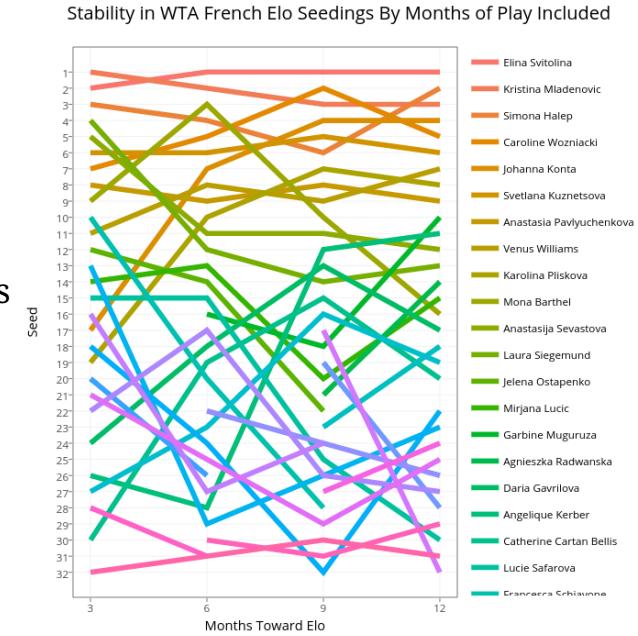
Carlos’s words resonated very well with me, although they were a little exaggerative. A well designed

Workflow Overview

1. Install `blogdown` and run `install_hugo`
2. Use `new_site` to create your site directory
3. Setup your directory with `git` and push to `github`
4. Choose your own theme choosing from themes.gohugo.io and use `install_theme`
5. Configure your site by modifying `config` file
6. Create new pages with `new_post`
7. *Deploy!*

Adding Interactive Graphics

- All sports blogs should use charts! And you can make your blog charts especially interesting by making them interactive.
- For R users who already masters of `ggplot2`, the `plotly` package is a flexible and fast way to create interactive Web graphics



Plotly

Plotly is a free online service for creating and sharing Web graphics. With the `plotly` package (by Carson Sievert) you can transform `ggplots` to interactive charts on plotly.

Getting Started with Plotly

Here are the basic steps to follow:

```
install.packages(plotly)

library(plotly)

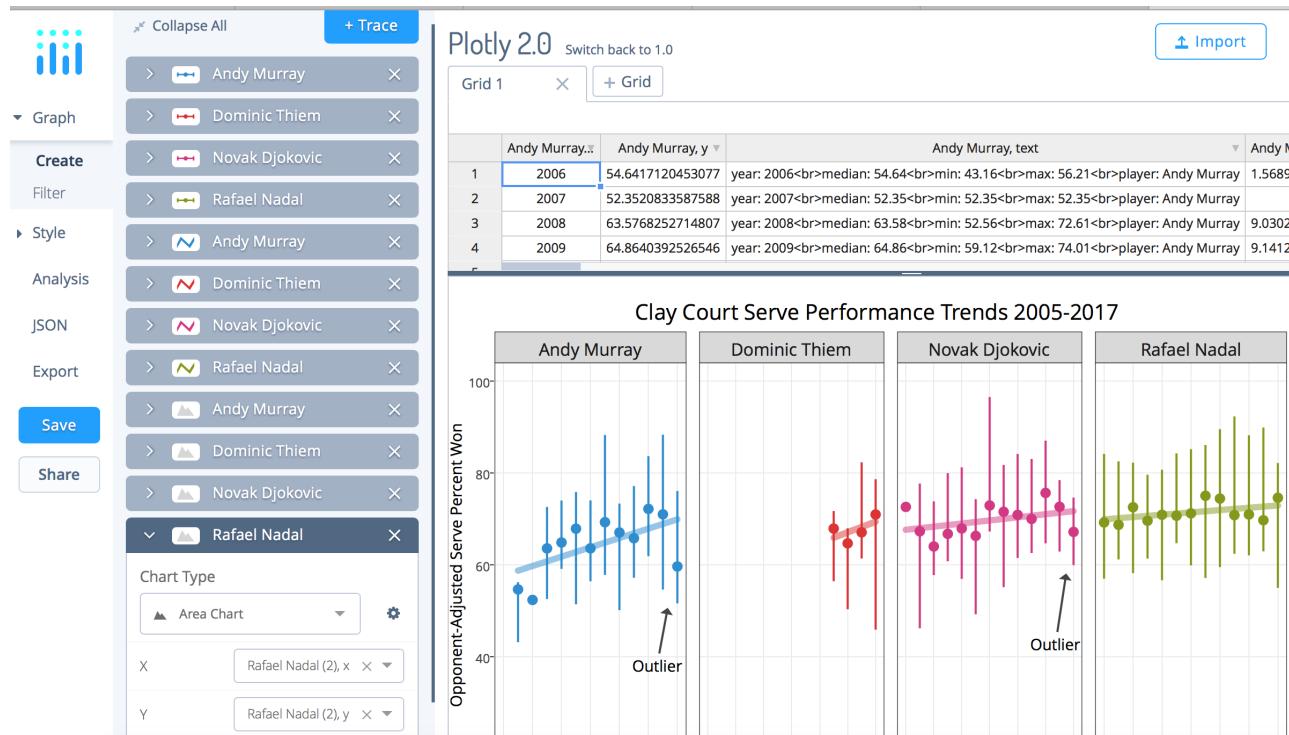
# After creating a plotly account
Sys.setenv("plotly_username"="your_plotly_username")
Sys.setenv("plotly_api_key"="your_api_key")

<CREATE GGPLOT>

plotly_POST(gg1, filename = "my-new-plot") # Publish
```

Plotly Site

Further customization can be done on the plotly site.



Plotly Embed

Under the Share link you can also find the code that you can use to embed your final plot in your R Markdown file.

The screenshot shows the 'Share' interface for a Plotly plot. At the top, there's a navigation bar with tabs: 'Link & Privacy', 'Collaborate', and 'Embed'. The 'Embed' tab is currently selected. Below the tabs, there's a section labeled 'Embed plot:' with two options: 'iframe' (which is selected) and 'html'. A large text area contains the embed code, which is highlighted with a blue background. The code is:

```
<iframe width="900" height="800" frameborder="0" scrolling="no" src="//plot.ly/~on-the-t/1217.embed"></iframe>
```

At the bottom of the interface, there's a link 'Need help? Check out our [embedding tutorial!](#)' and a page number '90 / 92'.

Summary

- Doing sports analytics requires a lot of skills
- Many tools in R can make capturing, cleaning, and analyzing data much easier
- I hope much of what you have seen in this presentation will help you make your ideas about sports data into a reality



skovalchik@tennis.com.au

@StatsOnTheT

