

# BST 270: Final Project

Dominic DiSanto

## Table of contents

<b>Set-Up</b>	<b>2</b>
<b>Preface</b>	<b>2</b>
<b>“Why Some Tennis Matches Take Forever”</b>	<b>3</b>
Table 1: Events . . . . .	3
Original . . . . .	3
Re-creation . . . . .	4
Table 2: Players . . . . .	5
Original . . . . .	5
Re-creation . . . . .	6
Summary/Analysis . . . . .	7
<b>“We Watched 906 Foul Balls To Find Out Where The Most Dangerous Ones Land”</b>	<b>8</b>
Figure 2 . . . . .	8
Original . . . . .	8
Re-creation . . . . .	9
Summary/Analysis . . . . .	10
<b>Appendix</b>	<b>11</b>
A: Session Information . . . . .	11
B: Code . . . . .	12

## Set-Up

## Preface

Tables were generated using the `kable` and `kableExtra` packages and data visualizations using `ggplot2` from the Tidyverse. This document was generated using [Quarto](#) in RStudio. Additional details are available in the Appendix (e.g. package versions, R/RStudio versions, etc.)

The `here` package has been used to load files, which *should* be robust to different operating systems, users, etc. as long as the file structure of this project (as described in the `README`) is followed. If issues persist, all specified file paths can be easily found searching for `here(` within this document and replacing with the user's preferred method of file path specification.

Each article has a brief section that displays the original table(s) or figure from the table, my recreation of the same results, and a brief summary section commenting on the reproducibility

# “Why Some Tennis Matches Take Forever”

From Carl Bialik’s article [“Why Some Tennis Matches Take Forever”](#), we will recreate his two tables (included below).

## Table 1: Events

### Original

We will first look at the original table from Bialik’s article:

Surface Speeds				
Average time added per point in men's tennis				
Fastest tournaments				
	TOURNAMENT	SURFACE	YEARS RUNNING	TIME ADDED
1	Wimbledon	Grass	1991-2014	- 2 . 98 s
2	London Olympics	Grass	2012	- 2 . 62
3	Manchester	Grass	1991-94	- 2 . 33
4	Eastbourne	Grass	2009-14	- 1 . 78
5	Birmingham	Carpet	1991	- 1 . 63
6	Queen's Club	Grass	1991-2014	- 1 . 50
7	Lyon	Hard	2009	- 1 . 35
8	Las Vegas	Hard	2006-08	- 1 . 32
9	Stockholm Masters	Carpet	1991-94	- 0 . 90
10	Nottingham	Grass	1995-2008	- 0 . 82
Slowest tournaments				
	TOURNAMENT	SURFACE	YEARS RUNNING	TIME ADDED
196	Birmingham	Clay	1994	+ 4 . 49 s
197	Oporto	Clay	1995-96	+ 4 . 69
198	Genova	Clay	1991-93	+ 4 . 71
199	Bologna	Clay	1991-98	+ 4 . 74
200	Merano	Clay	1999	+ 4 . 93
201	Viña del Mar	Clay	2009	+ 4 . 96
202	Florence	Clay	1991-94	+ 5 . 08
203	Costa do Sauipe	Clay	2004-11	+ 5 . 19
204	Maceio	Clay	1992	+ 5 . 28
205	Rio Open	Clay	2015	+ 5 . 38
Other notable tournaments				
	TOURNAMENT	SURFACE	YEARS RUNNING	TIME ADDED
21	U.S. Open	Hard	1991-2014	- 0 . 17 s
24	Australian Open	Hard	1991-2015	- 0 . 11
54	Roland Garros	Clay	1991-2014	+ 0 . 79

 FIVETHIRTYEIGHT

BASED ON DATA FROM JEFF SACKMANN

## Re-creation

The data used to recreate these tables is located in Bialik's provided `events.csv` file.

<b>Fastest tournaments</b>				
	Tournament	Surface	Years Running	Time Added
1	Wimbledon	Grass	1991-2014	-2.98s
2	London Olympics	Grass	2012	-2.62
3	Manchester	Grass	1991-1994	-2.33
4	Eastbourne	Grass	2009-2014	-1.78
5	Birmingham	Carpet	1991	-1.63
6	Queen's Club	Grass	1991-2014	-1.50
7	Lyon	Hard	2009	-1.35
8	Las Vegas	Hard	2006-2008	-1.32
9	Stockholm Masters	Carpet	1991-1994	-0.90
10	Nottingham	Grass	1995-2008	-0.82


<b>Slowest tournaments</b>				
	Tournament	Surface	Years Running	Time Added
196	Birmingham	Clay	1994	+4.49s
197	Oporto	Clay	1995-1996	+4.69
198	Genova	Clay	1991-1993	+4.71
199	Bologna	Clay	1991-1998	+4.74
200	Merano	Clay	1999	+4.93
201	Viña del Mar	Clay	2009	+4.96
202	Florence	Clay	1991-1994	+5.08
203	Costa Do Sauipe	Clay	2004-2011	+5.19
204	Maceio	Clay	1992	+5.28
205	Rio Open	Clay	2015	+5.38

<b>Other notable tournaments</b>				
	Tournament	Surface	Years Running	Time Added
21	US Open	Hard	1991-2014	-0.17s
24	Australian Open	Hard	1991-2015	-0.11
54	Roland Garros	Clay	1991-2014	+0.79

**Table 2: Players**

**Original**

Again, we first we can review the original table(s) presented in the article:

<b>Player Speeds</b>			
Average time added per point in men's tennis			
<b>Fastest players</b>			<b>Slowest players</b>
<b>1</b>	Dustin Brown	- 6 . 37 s	<b>209</b> Michael Chang + 4 . 94 s
<b>2</b>	Rohan Bopanna	- 4 . 95	<b>210</b> Joao Cunha Silva + 5 . 10
<b>3</b>	Chris Guccione	- 4 . 63	<b>211</b> Julian Knowle + 5 . 15
<b>4</b>	Benoit Paire	- 4 . 56	<b>212</b> John McEnroe + 5 . 22
<b>5</b>	Lukas Dlouhy	- 4 . 35	<b>213</b> Lucas Arnold Ker + 5 . 35
<b>6</b>	Brendan Evans	- 4 . 25	<b>214</b> T.J. Middleton + 5 . 60
<b>7</b>	Igor Sijsling	- 4 . 19	<b>215</b> Martin Stringari + 5 . 68
<b>8</b>	Lukas Rosol	- 4 . 13	<b>216</b> Rafael Nadal + 5 . 92
<b>9</b>	Alexander Kudryavtsev	- 4 . 05	<b>217</b> Nicolas Massu + 6 . 21
<b>10</b>	Sam Querrey	- 3 . 99	<b>218</b> Emanuel Couto + 6 . 35
<b>Other notable players</b>			
<b>22</b>	Goran Ivanisevic	- 3 . 15 s	
<b>36</b>	Roger Federer	- 2 . 43	
<b>124</b>	Novak Djokovic	+ 2 . 21	
<b>141</b>	Andy Murray	+ 2 . 53	
<b>191</b>	Pat Cash	+ 3 . 73	
<b>202</b>	Ivan Lendl	+ 4 . 35	
<b>203</b>	Jim Courier	+ 4 . 51	
<b>207</b>	Jimmy Connors	+ 4 . 90	
 FIVETHIRTYEIGHT		BASED ON DATA FROM JEFF SACKMANN	

## Re-creation

Now we can begin attempting to recreate this table, doing our best to mirror the format of the 538 results shown above in structure/formatting.

We generate each of the three tables included in the image above, using the `players_time.csv` table. Player rank (the integer column) was generated as noted below, player name taken from the `player` column, and the added time from the `seconds_added_per_point` column:

Fastest players			Slowest players		
1	Dustin Brown	-6.37s	209	Michael Chang	4.94s
2	Rohan Bopanna	-4.95	210	Joao Cunha Silva	+5.10
3	Chris Guccione	-4.63	211	Julian Knowle	+5.15
4	Benoit Paire	-4.56	212	John McEnroe	+5.22
5	Lukas Dlouhy	-4.35	213	Lucas Arnold Ker	+5.35
6	Brendan Evans	-4.25	214	T.J. Middleton	+5.60
7	Igor Sijsling	-4.19	215	Martin Stringari	+5.68
8	Lukas Rosol	-4.13	216	Rafael Nadal	+5.92
9	Alexander Kudryavtsev	-4.05	217	Nicolas Massu	+6.21
10	Sam Querrey	-3.99	218	Emanuel Couto	+6.35

Other notable players		
22	Goran Ivanisevic	-3.15s
36	Roger Federer	-2.43
124	Novak Djokovic	+2.21
141	Andy Murray	+2.53
191	Pat Cash	+3.73
202	Ivan Lendl	+4.35
203	Jim Courier	+4.51
207	Jimmy Connors	+4.90

We see our results mirror those in Bialik's original article exactly, and we were able to easily identify and match the data to recreate these tables from the provided GitHub repository (even without any code given by the author).

## Summary/Analysis

For Bialik's article on the pace of play in modern tennis, the data for the presented tables was immediately available and well-formatted to corroborate these results and easily output tables that are similar to those in Bialik's original article. The data used was sparse (with no additional information), which necessitated little documentation.

However considering the larger article, the two figures are presented without sufficient data to re-created. The provided data does not include any temporal data on player tendencies for the second figure, and the events data has been aggregated so that annual data is not available to create the first figure. The article also discusses regression modelling for pace of play but with no presentation of results or sharing of code to use in evaluating the model.

As an entertainment article, we obviously would not expect the same level of rigor for reproducibility as we might expect for a peer-reviewed publication or scientific article. The visualization and modelling are presented with no available data or thorough discussion of the methods to generate the results, while the tables in Bialik's article are easily reproduced from the shared materials.

## “We Watched 906 Foul Balls To Find Out Where The Most Dangerous Ones Land”

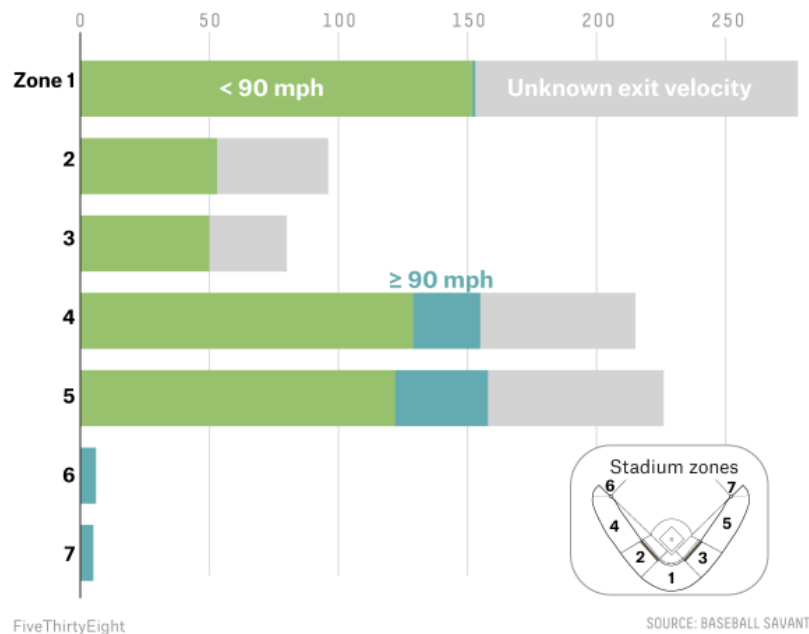
Data and analysis presented from Annette Choi’s article [We Watched 906 Fould Balls to Find Out Where the most Dangerous Ones Land](#).

Figure 2

Original

### The hardest-hit fouls seem to land in unprotected areas

Foul balls by the stadium zone they landed in and their exit velocity, among 906 fouls hit this season in the most foul-heavy day at the 10 MLB stadiums that produced the most fouls as of June 5

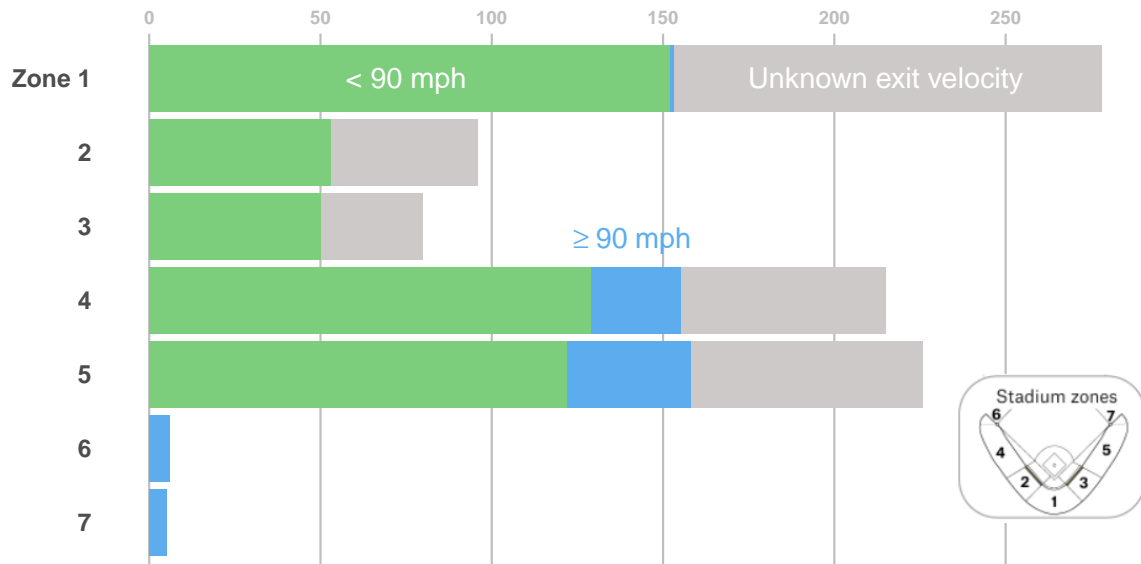




## Re-creation

### The hardest hit foul-balls seem to land in unprotected areas

Foul balls by the stadium zone they landed in and their exit velocity, among 906 foul balls hit this season in the most foul-heavy day at the 10 MLB stadiums that produced the most fouls as of June 5



## Summary/Analysis

For Annette Choi's article [We Watched 906 Foul Balls to Find Out Where the most Dangerous Ones Land](#), I recreated only her final figure, which summarizes the landing spot of foul balls into Choi's defined "zones", stratified by the exit velocity of the hit.

The data was immediately available and easy to use in generating this figure (i.e. only little wrangling necessary). The data contained three variables regarding the final "zones", with no indication in the article which was used to generate the figure (or in the analyses throughout the article). I settled on the `used_zone` variable, which uses the observed zone (`camera_zone`) of the foul ball when available and otherwise imputed a `predicted_zone` if the camera did not capture the final zone (and `camera_zone` was then missing). This imputation was discussed briefly in the article, but its use as the primary outcome variable was not mentioned explicitly in either the article or the figure.

Although I did not present the table and other figure, the table could be reproduced with the exception of the `stadium` column (although this could be inferred and created using the `matchup`'s home team). The ballpark figure could be recreated using the available data, although without shared code the method of creation of the ballpark figure is unknown and could only be approximated from the provided information.

# Appendix

## A: Session Information

R version 4.2.1 (2022-06-23 ucrt)  
Platform: x86\_64-w64-mingw32/x64 (64-bit)  
Running under: Windows 10 x64 (build 22621)

Matrix products: default

locale:

[1] LC\_COLLATE=English\_United States.utf8  
[2] LC\_CTYPE=English\_United States.utf8  
[3] LC\_MONETARY=English\_United States.utf8  
[4] LC\_NUMERIC=C  
[5] LC\_TIME=English\_United States.utf8

attached base packages:

[1] grid stats graphics grDevices utils datasets methods  
[8] base

other attached packages:

[1] png\_0.1-7 kableExtra\_1.3.4 stringr\_1.4.1 dplyr\_1.0.10  
[5] ggplot2\_3.4.0 here\_1.0.1 knitr\_1.40

loaded via a namespace (and not attached):

[1] compiler\_4.2.1 pillar\_1.8.1 tools\_4.2.1 digest\_0.6.29  
[5] viridisLite\_0.4.1 jsonlite\_1.8.2 evaluate\_0.17 lifecycle\_1.0.3  
[9] tibble\_3.1.8 gtable\_0.3.1 pkgconfig\_2.0.3 rlang\_1.0.6  
[13] cli\_3.4.1 DBI\_1.1.3 rstudioapi\_0.14 yaml\_2.3.5  
[17] xfun\_0.33 fastmap\_1.1.0 xml2\_1.3.3 httr\_1.4.4  
[21] withr\_2.5.0 systemfonts\_1.0.4 generics\_0.1.3 vctrs\_0.5.1  
[25] webshot\_0.5.4 rprojroot\_2.0.3 tidyselect\_1.2.0 svglite\_2.1.0  
[29] glue\_1.6.2 R6\_2.5.1 fansi\_1.0.3 rmarkdown\_2.17  
[33] farver\_2.1.1 magrittr\_2.0.3 scales\_1.2.1 htmltools\_0.5.3  
[37] rvest\_1.0.3 assertthat\_0.2.1 colorspace\_2.0-3 utf8\_1.2.2  
[41] stringi\_1.7.8 munsell\_0.5.0

## B: Code

```
#####  
#### Set-Up ####  
#####  
  
list.of.packages <- c("here", "ggplot2", "dplyr", "stringr", "kableExtra", "knitr")  
  
new.packages <- list.of.packages[!(list.of.packages %in% installed.packages()[,"Package"])]  
  
if(length(new.packages)>0) install.packages(new.packages)  
  
library(here) # used for file management  
library(ggplot2) # data visualization  
library(dplyr) # data wrangling  
library(stringr) # data wrangling (string detection, manipulation)  
library(kableExtra) # LaTeX/PDF table styling  
library(knitr)  
#####  
### Data Import ###  
#####  
  
events <- read.csv(here("Data", "tennis-time", "events_time.csv"), header=T) %>%  
  mutate(Rank=row_number())  
  
players <- read.csv(here("Data", "tennis-time", "players_time.csv")) %>%  
  mutate(Rank=row_number()) # adding the "Rank" variable (the integers in each table  
    # above), descending (1=fastest player, lowest `seconds_added_per_point`)  
  
#####  
### Table 1 ###  
#####  
  
# Original  
knitr::include_graphics(here("Images", "tennis-time", "bialik-tennis-time-table21.png"))  
  
# Re-creation  
  
events_10f <- events %>%  
  arrange(seconds_added_per_point) %>%  
  head(10) %>%
```

```

mutate(seconds = # reformatting as string to include the "s" for first obs
       case_when(row_number()==1 ~ paste0(sprintf("%.2f", seconds_added_per_point)
                                           , "s")
               , seconds_added_per_point<0 ~ sprintf("%.2f", seconds_added_per_point)
               , T ~ paste0("+",
                           sprintf("%.2f", seconds_added_per_point))
               )
       ) %>%
select(Rank, tournament, surface, years, seconds)

events_10s <- events %>%
mutate(Rank=row_number()) %>%
arrange(desc(seconds_added_per_point)) %>%
head(10) %>%
arrange(seconds_added_per_point) %>%
mutate(seconds = # reformatting as string to include the "s" for the first obs
       case_when(row_number()==1 ~ paste0("+",
                                           , sprintf("%.2f", seconds_added_per_point)
                                           , "s")
               , seconds_added_per_point<0 ~ sprintf("%.2f", seconds_added_per_point)
               , T ~ paste0("+", sprintf("%.2f", seconds_added_per_point))
               )
       ) %>%
select(Rank, tournament, surface, years, seconds)

events_ntbl <- events %>%
filter(str_detect(tournament, "US Op|Australian Open|Roland Garr")) %>%
mutate(seconds = # reformatting as string to include the "s" for the first obs
       case_when(row_number()==1 ~ paste0(sprintf("%.2f", seconds_added_per_point), "s"
               , seconds_added_per_point<0 ~ sprintf("%.2f", seconds_added_per_point)
               , T ~ paste0("+", sprintf("%.2f", seconds_added_per_point))
               )
       ) %>%
select(Rank, tournament, surface, years, seconds)

event_cols <- c("", "Tournament", "Surface", "Years Running", "Time Added")

events_10f %>%
kable(col.names=event_cols
      # , bottomrule=''
      , booktabs=T

```

```

    ) %>%
  kable_styling(latex_options = "HOLD_position") %>%
  add_header_above(c("Fastest tournaments"=5), align = "l"
    , bold = T
    , underline = F, line = F)
events_10s %>%
  kable(col.names=event_cols
    # , bottomrule=''
    , booktabs=T
  ) %>%
  kable_styling(latex_options = "HOLD_position") %>%
  add_header_above(c("Slowest tournaments"=5), align = "l"
    , bold = T
    , underline = F, line = F)
events_ntbl %>%
  kable(col.names=event_cols
    # , bottomrule=''
    , booktabs=T
  ) %>%
  kable_styling(latex_options = "HOLD_position") %>%
  add_header_above(c("Other notable tournaments"=5), align = "l"
    , bold = T
    , underline = F, line = F)

#####
### Table 2 ###
#####

# Original

knitr::include_graphics(here("Images", "tennis-time", "bialik-tennis-time-table1.png"))

# Re-creation

# Top 10 Fastest Players (lower/more negative seconds_added_per_point is faster pace)
t10_f <- players %>%
  arrange(seconds_added_per_point) %>% # sort ascending
  head(10) %>% # top 10
  mutate(seconds = # reformatting as above
    case_when(row_number() == 1 ~ paste0(sprintf("%0.2f", seconds_added_per_point)

```

```

                                , "s")
                                , T ~ sprintf("%.2f"
                                                , seconds_added_per_point)
                                )
                                ) %>%
select(Rank, player, seconds)

# "Top" 10 slowest players (largest values)
t10_s <- players %>%
  arrange(desc(seconds_added_per_point)) %>%
  head(10) %>%
  arrange(seconds_added_per_point) %>%
  mutate(seconds =
    case_when(row_number() == 1 ~ paste0(sprintf("%.2f", seconds_added_per_point)
                                            , "s")
              , seconds_added_per_point < 0 ~ sprintf("%.2f", seconds_added_per_point)
              , T ~ paste0("+",
                            sprintf("%.2f", seconds_added_per_point))
              )
    ) %>%
select(Rank, player, seconds)

# Specific players ("notable") taken from the full data set
tbl_pls <- players %>%
  filter(
    str_detect(player
, "Goran|Roger Fed|Novak Djok|Andy Murray|Pat Cash|Ivan Lend|Jim Courier|Jimmy Connors"))
  arrange(seconds_added_per_point) %>%
  mutate(seconds =
    case_when(row_number() == 1 ~ paste0(sprintf("%.2f", seconds_added_per_point)
                                            , "s")
              , seconds_added_per_point < 0 ~ sprintf("%.2f", seconds_added_per_point)
              , T ~ paste0("+", sprintf("%.2f", seconds_added_per_point))
              )
    ) %>%
select(Rank, player, seconds)

t10_f %>%
  kable(col.names=NULL
        # , bottomrule='')

```

```

      , booktabs=T
    ) %>%
kable_styling(latex_options = "HOLD_position") %>%
add_header_above(c("Fastest players"=3)
                  , bold=T, align="l"
                  , underline=F, line=F
                  )

tbl0_s %>%
  kable(col.names=NULL
        # , bottomrule=''
        , booktabs=T
        ) %>%
kable_styling(latex_options = "HOLD_position") %>%
add_header_above(c("Slowest players"=3)
                  , bold=T, align="l"
                  , underline=F, line=F
                  )

tbl1_pls %>%
  kable(col.names=NULL
        # , bottomrule=''
        , booktabs=T
        ) %>%
kable_styling(latex_options = "HOLD_position") %>%
add_header_above(c("Other notable players"=3)
                  , bold=T, align="l"
                  , underline=F, line=F
                  )

knitr::include_graphics(here("Images", "foul-balls", "choi-foul-0625-2-1.png"))
foul_balls <- read.csv(here("Data", "foul-balls", "foul-balls.csv"))
library(png); library(grid)

img <- readPNG(here("Images", "foul-balls", "stadium_zones.png"))
g <- rasterGrob(img, interpolate=TRUE)

foul_balls %>%
  mutate(Velocity_Cat =
         case_when(exit_velocity<90 ~ 1
                   , exit_velocity>=90 ~ 2
                   , TRUE ~ 3)
         , Zone_f = case_when(

```



```

      used_zone == 1 ~ "Zone 1"
    , TRUE ~ as.character(used_zone)
  )) %>%
mutate(fvc = factor(Velocity_Cat, levels=c(3, 2, 1))) %>%
ggplot(aes(x=reorder(as.factor(Zone_f), -used_zone), fill=fvc)) +
geom_bar() +
scale_y_continuous(breaks=seq(0, 250, 50)
                    , position="right") +
scale_fill_manual(values=c("snow3", "steelblue2", "palegreen3")) +
coord_flip() +
ylab("") + xlab("") +
# ggthemes::theme_fivethirtyeight() +
ggtitle(label="The hardest hit foul-balls seem to land in unprotected areas"
        , subtitle = "Foul balls by the stadium zone they landed in and their exit velocity")
theme(legend.position="none", plot.title = element_text(size=16)
      , panel.background = element_rect(fill = 'white', color = 'white')
      , panel.grid.major.x = element_line(colour = "gray", linewidth = 0.5)
      , axis.text.y = element_text(face="bold", size = 12)
      , axis.text.x = element_text(color="gray", face="bold")
      , axis.ticks = element_blank()
      ) +
annotation_custom(g, xmin=1, xmax=3, ymin=210, ymax=320) +
annotate("text", x = 7, y = 75, label="< 90 mph", col="white", size=5) +
annotate("text", x = 7, y = 215, label="Unknown exit velocity", col="white", size=5) +
annotate("text", x = 4.8, y = 140, label="('')>= 90 ~ 'mph'", col="steelblue2", size=5,

sessionInfo()

```