

# Q&A/DEMO SESSION OVERVIEW

- **SESSION 1 TOPICS**

- **Q&A**

- Project and Analysis
    - Data and Computation

- **Demo: Analysis in R**

- **UPCOMING SESSION TOPICS (PROJECTED)**

- **July 5:** Demo Analysis in Python
  - **July 26:** How to write a paper
  - **August 16:** Structuring and commenting code (R markdown etc.)

- **ARCHIVED VIDEOS AVAILABLE AFTER**

Q&A

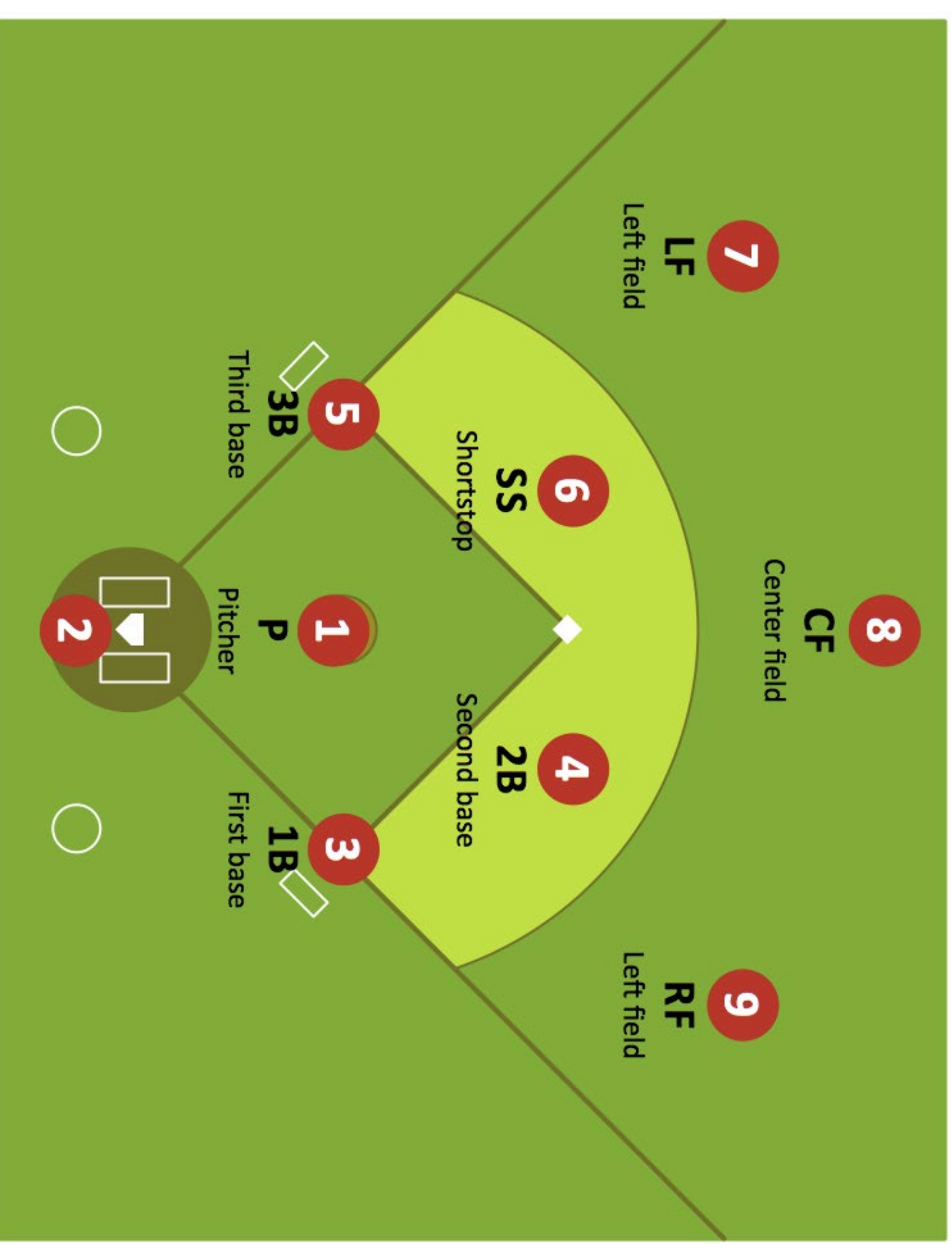


# PROJECT AND ANALYSIS

1. *Start with an idea*
2. *Turn a bug into a feature*
3. *“Frame it as a story”*
4. *Keep your audience(s) in mind*
  1. Judges
  2. Teams and other employers

# DATA AND COMPUTATION

1. Explaining *player\_position* and *event\_code*
2. What about missing data?
3. Identifying different plays
4. Do you have to use all the data?
5. Can you include other data?
6. Computer language?



# DATA AND COMPUTATION

- 1. Explaining *player\_position* and *event\_code*
- 2. What about missing data?
- 3. Identifying different plays
- 4. Do you have to use all the data?
- 5. Can you include other data?
- 6. Computer language?



## GLOSSARY

VARIABLE	TABLE	CODE	DEFINITION
player_position	game_events	1	pitcher
	player_pos (game_info)	2	catcher
		3	first baseman
		4	second baseman
		5	third baseman
		6	shortstop
		7	left field
		8	center field
		9	right field
		10	batter
		11	runner on first base
		12	runner on second base
		13	runner on third base
event_code		255	ball event with no player (e.g., ball bounce)
	game_events	1	pitch
		2	ball acquired
		3	throw (ball-in-play)
		4	ball hit into play
		5	end of play
		6	pickoff throw
		7	ball acquired - unknown field position
		8	throw (ball-in-play) - unknown field position
		9	ball deflection
		10	ball deflection off of wall
		11	home run
		16	ball bounce



# DATA AND COMPUTATION

- 1. Explaining *player\_position* and *event\_code*
- 2. What about missing data?
- 3. Identifying different plays
- 4. Do you have to use all the data?
- 5. Can you include other data?
- 6. Computer language?

game_str	play_id	at_bat	play_per_game	timestamp	player_position	event_code
1900_09_TeamKK_TeamB	154	NA	157	5799556	0	5
1900_09_TeamKK_TeamB	155	NA	158	5822725	1	1
1900_09_TeamKK_TeamB	155	NA	158	5823187	10	4
1900_09_TeamKK_TeamB	155	NA	158	5825102	6	2
1900_09_TeamKK_TeamB	155	NA	158	5826191	6	3
1900_09_TeamKK_TeamB	155	NA	158	5827412	3	2
1900_09_TeamKK_TeamB	155	NA	158	5827775	0	5
1900_09_TeamKK_TeamB	156	NA	159	5858470	1	1
1900_09_TeamKK_TeamB	156	NA	159	5858932	2	2
1900_09_TeamKK_TeamB	156	NA	159	5858932	0	5
1900_09_TeamKK_TeamB	157	NA	160	5876391	1	1
1900_09_TeamKK_TeamB	157	NA	160	5876854	2	2
1900_09_TeamKK_TeamB	157	NA	160	5876854	0	5
1900_09_TeamKK_TeamB	158	NA	161	5915205	1	1
1900_09_TeamKK_TeamB	158	NA	161	5915602	10	4
1900_09_TeamKK_TeamB	158	NA	161	5915734	0	5
1900_09_TeamKK_TeamB	159	NA	162	5945306	1	1
1900_09_TeamKK_TeamB	159	NA	162	5945768	2	2
1900_09_TeamKK_TeamB	159	NA	162	5945768	0	5
1900_09_TeamKK_TeamB	160	NA	163	5968278	1	1
1900_09_TeamKK_TeamB	160	NA	163	5968707	10	4
1900_09_TeamKK_TeamB	160	NA	163	5969268	255	16
1900_09_TeamKK_TeamB	160	NA	163	5970522	4	2
1900_09_TeamKK_TeamB	160	NA	163	5971776	4	3
1900_09_TeamKK_TeamB	160	NA	163	5972766	3	2
1900_09_TeamKK_TeamB	160	NA	163	5973261	0	5

# DEMO ANALYSIS IN R



```

1 ##### Total distance traveled vs. straight line distance to make a fly out
2
3 ### EXAMPLE: Consider fly ball routes for a single player
4
5
6
7 ### LIBRARIES
8 library(tidyverse)
9
10
11
12 ### LOAD DATA
13 game_info <- read.csv('game_info_anon_2.0.csv', header = TRUE)
14 game_events <- read.csv('game_events_anon_2.0.csv', header = TRUE)
15 team_info <- read.csv('team_info_anon_2.0.csv', header = TRUE)
16 ball_pos <- read.csv('ball_pos_anon_2.0.csv', header = TRUE)
17 player_pos <- read.csv('player_pos_anon_2.0.csv', header = TRUE)
18
19
20
21 ### CHOOSE A PLAYER
22 # Limit to center fielders
23 # Choose a player that's made a lot of fly outs
24
25 ## Using game_info, how many plays for each center fielder?
26 cf_play_count <- game_info %>% group_by(center_field) %>%
27   summarize(count=n()) %>% arrange(desc(count))
28
29 player1 <- as.numeric(cf_play_count[1,1])
30
31 ## In which games did player1 play?
32 player1_game_info <- game_info %>% filter(center_field == player1)
33 games_unique <- unique(player1_game_info$game_str)
34
35
36
37 ### FIND THE PLAYS WE WANT
38 # In which games did player1 make a fly ball catch?
39 #
40 # 1) Pull information from game_events
41 #
42 # 2) Look for sequences where a ball is put in play, then caught by the center fielder
43 #
44 #
45 #           player_position | event_code
46 #           -----
47 #           10             |      4
48 #           -----
49 #           8              |      2
50 #           -----
51 #
52 # 3) Are the catches by player1 (not the other team's center fielder)?
53 # Check using play_per_game, and keep track of play_per_game values.
54
55 ## Create a data frame containing games, play_id, and play_per_game
56 player1_plays_df <- data.frame(game_str = character(0),
57                                play_id = integer(0),
58                                play_per_game = integer(0),
59                                stringsAsFactors = FALSE)
60
61 ## Loop through player1 games in game_events
62 for(n1 in 1:length(games_unique)){
63   # for(n1 in 1:1){
64     game <- game_events %>% filter(game_str == games_unique[n1])

```



```

64 player1_game <- player1_game_info %>% filter(game_str == games_unique[n1])
65
66 ## 1) Loop through an individual game
67 for(n2 in 1:(nrow(game)-1)){
68
69   ## 2) Find catches
70   if(game$player_position[n2] == 10 & game$event_code[n2] == 4
71       & game$player_position[n2+1] == 8 & game$event_code[n2+1] == 2){
72     play_id <- game$play_id[n2]
73     play_num <- game$play_per_game[n2]
74     if(play_num %in% player1_game$play_per_game){
75       play_df <- data.frame(game_str = games_unique[n1],
76                             play_id = play_id,
77                             play_per_game = play_num)
78       player1_plays_df <- rbind(player1_plays_df, play_df)
79     }
80   }
81 }
82 }
83
84
85
86 ### -----
87 # How do we determine the start/end times for distance traveled?
88 #
89 #   - Get timestamp of hit and timestamp of catch
90 #
91 #   - Assumption: player1 starts moving as soon as the ball leaves the bat.
92 #       Start time = timestamp when event_code == 4
93 #       End time = timestamp when event_code == 2
94 #
95 # Let's see how long each ball is in the air.
96 #
97 ### (This starts with the same code as above. I'm just adding to it.)
98 #
99
100 ## Expand the data frame to include start and end times.
101 player1_plays_df <- data.frame(game_str = character(0),
102                                play_id = integer(0),
103                                play_per_game = integer(0),
104                                start_time = numeric(0),
105                                end_time = numeric(0),
106                                del_t = numeric(0),
107                                stringsAsFactors = FALSE)
108
109 ## Loop through player1 games in game_events
110 for(n1 in 1:length(games_unique)){
111   # for(n1 in 1:1){
112   game <- game_events %>% filter(game_str == games_unique[n1])
113   player1_game <- player1_game_info %>% filter(game_str == games_unique[n1])
114
115   ## 1) Loop through an individual game
116   for(n2 in 1:(nrow(game)-1)){
117
118     ## 2) Find catches
119     if(game$player_position[n2] == 10 & game$event_code[n2] == 4
120         & game$player_position[n2+1] == 8 & game$event_code[n2+1] == 2){
121       play_id <- game$play_id[n2]
122       play_num <- game$play_per_game[n2]
123       if(play_num %in% player1_game$play_per_game){
124         start_time <- game$timestamp[n2]
125         end_time <- game$timestamp[n2+1]
126         del_t <- end_time - start_time

```

```

127     play_df <- data.frame(game_str = games_unique[n1],
128                           play_id = play_id,
129                           play_per_game = play_num,
130                           start_time = start_time,
131                           end_time = end_time,
132                           del_t = del_t)
133     player1_plays_df <- rbind(player1_plays_df, play_df)
134   }
135 }
136 }
137 }
138
139
140
141 ### FIND PLAYER DISTANCE INFORMATION
142 #
143 # 1) Pull information from player_pos
144 #
145 # 2) Use info in player1_plays_df to constrain data (it's a big data frame!)
146 #
147 # 3) Add data frame columns for straight line distance and total distance traveled
148 #     STRAIGHT LINE DISTANCE: [x,y] distance between player1 location
149 #     at start and end times
150 #     TOTAL DISTANCE TRAVELED: the sum of [x,y] distances by timestep
151 #     between start and end times
152 #     ROUTE EFFICIENCY: (straight line distance)/(total distance traveled)
153 #     (N.B. - This is considered a suspect metric. We'll see why in a minute.)
154
155 ## Create a new data frame that includes distance columns
156 player1_distance_df <- data.frame(game_str = character(0),
157                                   play_id = integer(0),
158                                   play_per_game = integer(0),
159                                   start_time = numeric(0),
160                                   end_time = numeric(0),
161                                   del_t = numeric(0),
162                                   d_straight = numeric(0),
163                                   d_total = numeric(0),
164                                   route_eff = numeric(0),
165                                   stringsAsFactors = FALSE)
166
167 # NOTE: Since player_pos is big, we want to filter it down as early as possible
168 # - Use player1_plays_df to loop through player_pos
169 # - player_pos uses play_id, but not play_per_game
170
171 for(n1 in 1:nrow(player1_plays_df)){
172   # for(n1 in 1:1){
173   # n1 = 1
174   player1_pos <- player_pos %>% filter(game_str == player1_plays_df$game_str[n1]
175                                     & play_id == player1_plays_df$play_id[n1]
176                                     & player_position == 8
177                                     & timestamp >= player1_plays_df$start_time[n1]
178                                     & timestamp <= player1_plays_df$end_time[n1]
179                                     ) %>%
180     arrange(timestamp)
181
182   ## Calculate straight line distance
183   start_x <- player1_pos$field_x[1]
184   end_x <- player1_pos$field_x[nrow(player1_pos)]
185   start_y <- player1_pos$field_y[1]
186   end_y <- player1_pos$field_y[nrow(player1_pos)]
187
188   d_straight <- sqrt((end_x - start_x)^2 + (end_y - start_y)^2)
189

```

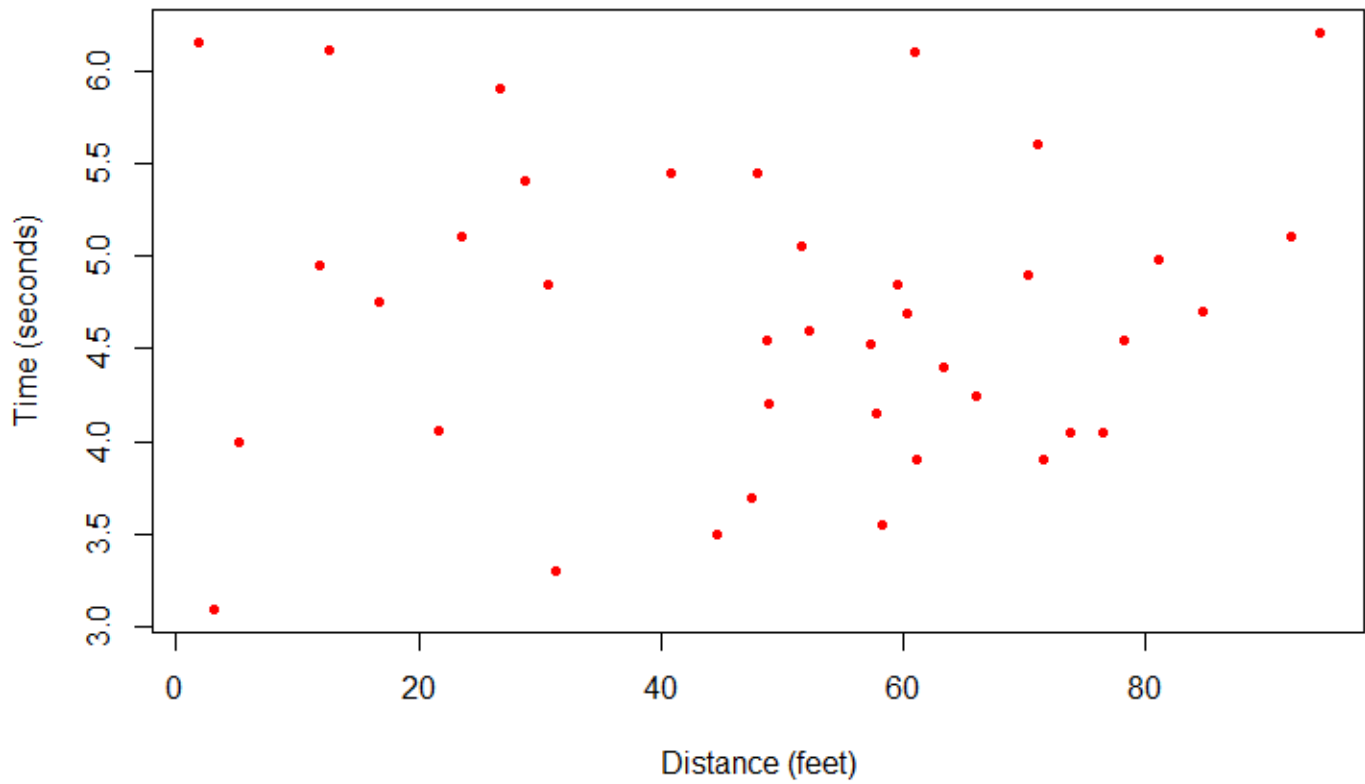
```

190     ## Calculate total distance traveled
191     d_total = 0
192     for(n2 in 1:(nrow(player1_pos)-1)){
193         start_x <- player1_pos$field_x[n2]
194         end_x <- player1_pos$field_x[n2+1]
195         start_y <- player1_pos$field_y[n2]
196         end_y <- player1_pos$field_y[n2+1]
197
198         d_step <- sqrt((end_x - start_x)^2 + (end_y - start_y)^2)
199         d_total <- d_total + d_step
200     }
201
202     ## Calculate route efficiency
203     route_eff <- d_straight/d_total
204
205     distance_df <- data.frame(game_str = player1_plays_df$game_str[n1],
206                               play_id = player1_plays_df$play_id[n1],
207                               play_per_game = player1_plays_df$play_per_game[n1],
208                               start_time = player1_plays_df$start_time[n1],
209                               end_time = player1_plays_df$end_time[n1],
210                               del_t = player1_plays_df$del_t[n1],
211                               d_straight = d_straight,
212                               d_total = d_total,
213                               route_eff = route_eff)
214     player1_distance_df <- rbind(player1_distance_df, distance_df)
215 }
216
217
218
219 ##### -----
220 ### USING VISUALS TO ANSWER QUESTIONS ABOUT PLAYER1
221
222 ## 1) How does hang time (del_t) correspond to straight line distance?
223 plot(player1_distance_df$d_straight, player1_distance_df$del_t/1000,
224       main = 'Straight Line Distance vs. Hang Time',
225       xlab = 'Distance (feet)', ylab = 'Time (seconds)',
226       col = 'red', pch = 20)
227
228 **** The farther a player has to travel to make a catch, the longer the ball has to
229 **** be in the air. Otherwise, he wouldn't be able to get to it.
230
231
232 ## 2) How does route efficiency correspond to hang time?
233 plot(player1_distance_df$del_t/1000, player1_distance_df$route_eff,
234       main = 'Hang Time vs. Route Efficiency',
235       xlab = 'Time (seconds)', ylab = 'Route Efficiency',
236       col = 'red', pch = 20)
237
238 **** If the ball is airborne longer, the player has more time to accommodate
239 **** an inefficient route and still make the catch.
240
241
242 ## 3) How does route efficiency correspond to straight line distance?
243 plot(player1_distance_df$d_straight, player1_distance_df$route_eff,
244       main = 'Straight Line Distance vs. Route Efficiency',
245       xlab = 'Distance (feet)', ylab = 'Route Efficiency',
246       col = 'red', pch = 20)
247
248 **** The farther a player has to travel to make a catch, the more efficient the route
249 **** required. He has less time to cover the distance, and if his route isn't
250 **** efficient enough, he can't make the play.
251
252

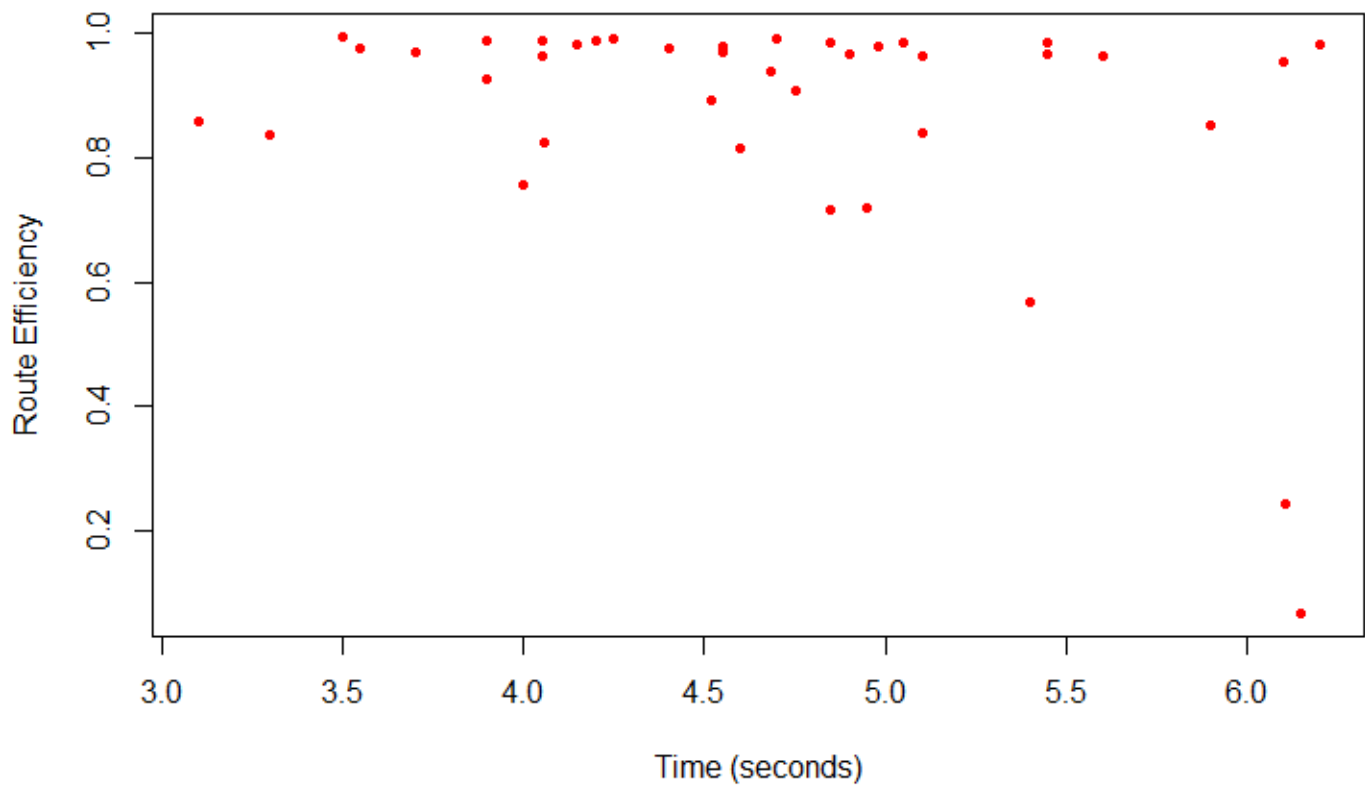
```

[illegible]

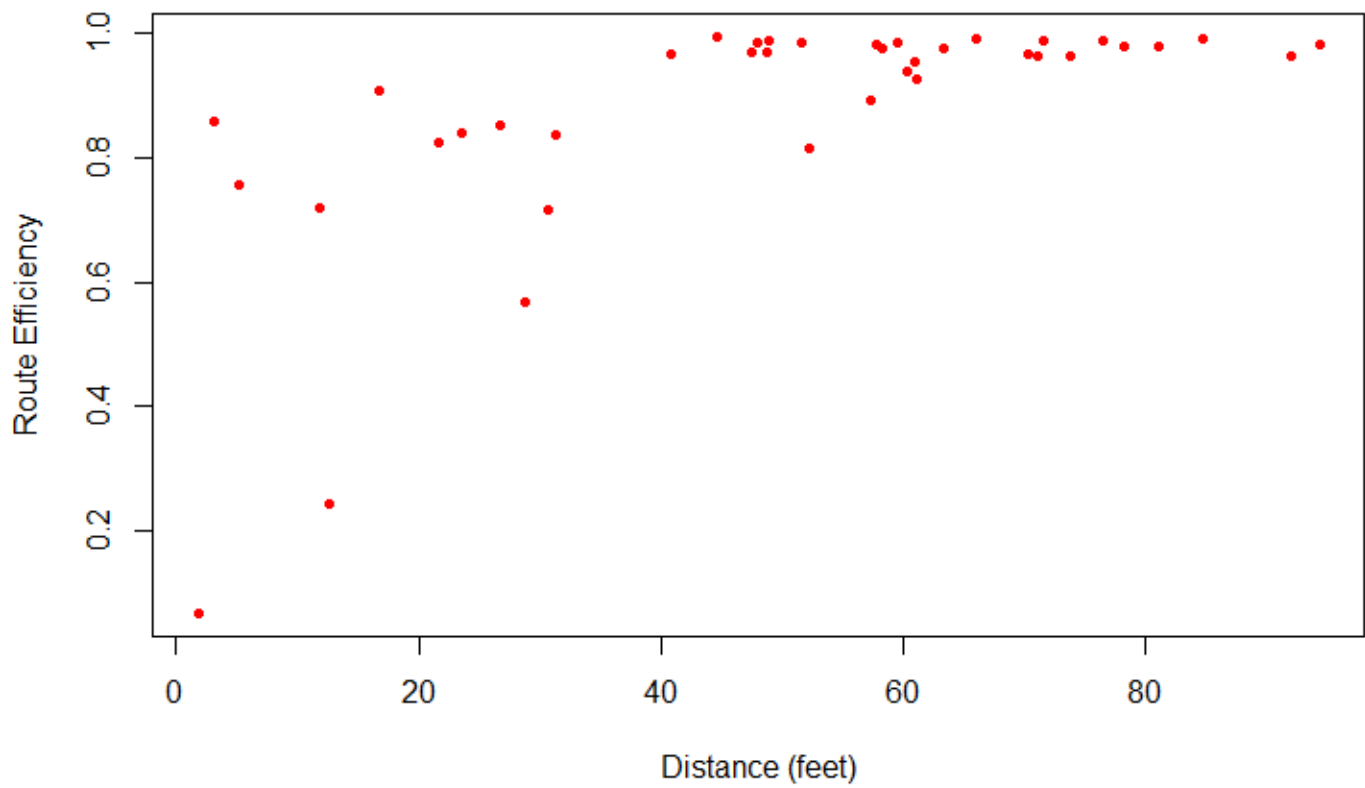
**Straight Line Distance vs. Hang Time**



**Hang Time vs. Route Efficiency**



**Straight Line Distance vs. Route Efficiency**



**Straight Line Speed vs. Route Efficiency**

