

Group 7 Final Project

Connor Beaudry, Rishi Young, James Jeffs
STAT 27815, University of Chicago

Topic and Data Collection



What are the main questions that this presentation intends to address?

- What is the effect of [Home Team Advantage](#) on scoring and match outcome?
- Does being the away or home team affect how the game is played?
- Does distance traveled by the away team effect how well they play?

Data Collection

- [worldfootballr](#) is an r package made by [Jason Zivkovic](#).
 - It contains functions for scraping data from [FBRef](#), [Transfermarkt](#), and [Understat](#).
 - These websites provide info from all major soccer leagues on match summaries, results, shooting data and more as well as player/team data.

Comparing Different League Statistics

- Let us first analyze whether there is significant difference between leagues.
- In order to do so, we will scrape data for seasons 2016-2017 through 2022-2023 for 5 major leagues: the EPL, La Liga, Ligue 1, Bundesliga, and MLS.

```
1 options(width = 80)
2 library(worldfootballR)
3
4 pull_league_data = function(countries) {
5   output = data.frame()
6   for(country in countries){
7     new = fb_match_results(
8       country = country,
9       gender = "M",
10      season_end_year = c(2017:2023),
11      tier = "1st"
12    )
13    output = rbind(output, new)
14  }
15  output = output |>
16    select(Country, Season_End_Year, Home, HomeGoals, Away, AwayGoals) |>
17    mutate(TotalGoals = HomeGoals + AwayGoals)
18
19  return(output)
20 }
21
22 League_Comparison = pull_league_data(c("ENG", "ESP", "FRA", "GER", "USA"))
```

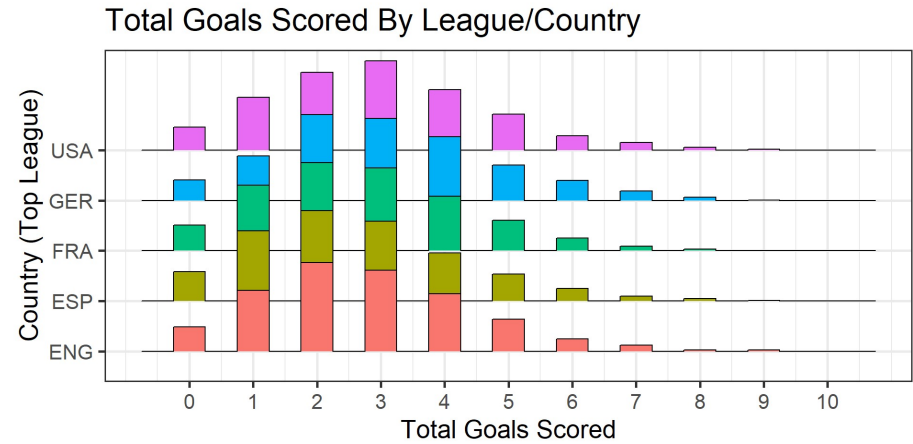
Comparing Different League Statistics

	Country	Season_End_Year	Home	HomeGoals	Away	AwayGoals	TotalGoals
412	ENG	2018	West Brom	0	Chelsea	4	4
2323	ENG	2023	Manchester City	3	Brighton	1	4
10133	GER	2023	Werder Bremen	5	Gladbach	1	6
908	ENG	2019	Everton	2	Bournemouth	0	2
2949	ESP	2017	Celta Vigo	0	Athletic Club	3	3
8154	GER	2017	Hamburger SV	2	Köln	1	3
10847	USA	2018	NYCFC	1	D.C. United	1	2
12663	USA	2023	Inter Miami	2	CF Montréal	0	2
4572	ESP	2022	Valencia	2	Mallorca	2	4
7259	FRA	2022	Paris S-G	2	Lille	1	3

Comparing Different League Statistics

Using this data, we make a couple basic visualizations:

```
1 League_Comparison |>
2 ggplot(aes(x = TotalGoals, y = Country, height = stat(dens:
3   geom_density_ridges(stat = "binline", bins = 21) +
4   guides(fill = "none") +
5   labs(x = "Total Goals Scored", y = "Country (Top League)"
6   expand_limits(y = 7) +
7   scale_x_continuous(breaks = 0:10) +
8   theme_bw(base_size = 20)
```

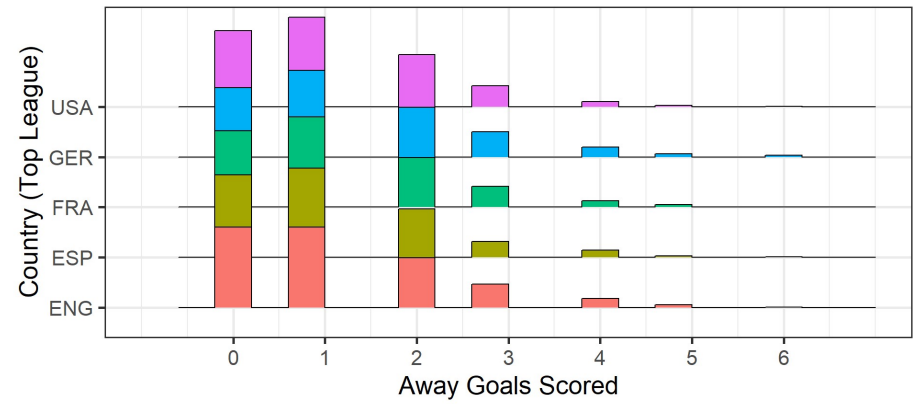


```

1 League_Comparison |>
2 ggplot(aes(x = AwayGoals, y = Country, height = stat(densit
3   geom_density_ridges(stat = "binline", bins = 21) +
4   guides(fill = "none") +
5   labs(x = "Away Goals Scored", y = "Country (Top League)",
6   expand_limits(y = 7) +
7   scale_x_continuous(breaks = 0:6, limits = c(-1, 7)) +
8   theme_bw(base_size = 20)

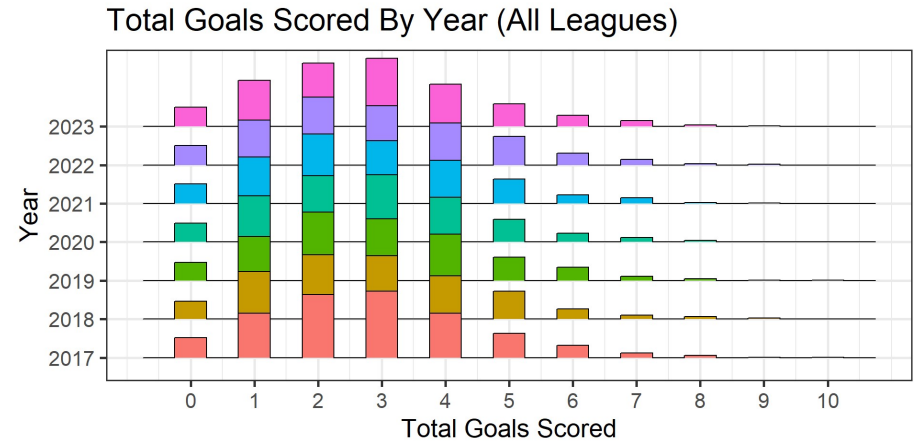
```

Away Goals Scored By League/Country



Comparing Different League Statistics

```
1 library(GGally)
2 League_Comparison |>
3   ggplot(mapping = aes(x = TotalGoals, y = factor(Season_End))) +
4   geom_density_ridges(stat = "binline", bins = 21) +
5   guides(fill = "none") +
6   labs(y = "Year", x = "Total Goals Scored", title = "Total Goals Scored By Year (All Leagues)") +
7   expand_limits(y = 9) +
8   scale_x_continuous(breaks = 0:10) +
9   theme_bw(base_size = 20)
```



The English Premier League

The English Premier League

- Making a map of The English Premier League (EPL):

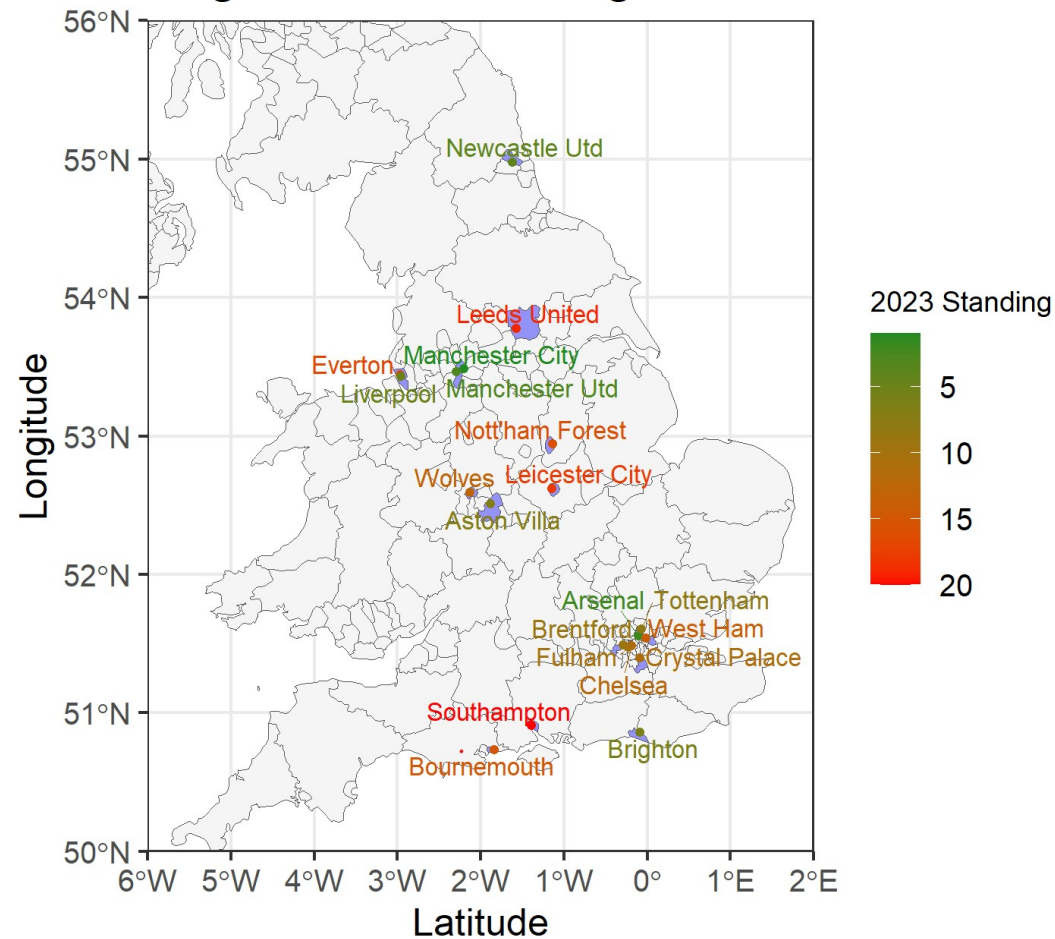
```
1 library(rnaturalearth)
2 library(rnaturalearthdata)
3 library(rnaturalearthhires)
4 library(ggmap)
5
6 locations_2023 = read_csv(file = "data/Team Stadium Locations EPL 2023.csv", show_col_types = FALSE)
7 UK_map = ne_states(country = "United Kingdom", returnclass = "sf")
8 counties = c("Newcastle upon Tyne", "Leeds", "Manchester", "Liverpool",
9             "Wolverhampton", "Nottingham", "Leicester", "Birmingham",
10            "Southampton", "Bournemouth", "Brighton and Hove",
11            "Islington", "Hounslow", "Haringey", "Hammersmith and Fulham",
12            "Newham", "Croydon")
13 UK_map_teams = filter(UK_map, name %in% counties)
14
15 ggplot(data = UK_map) +
16   geom_sf(fill = "whitesmoke") +
17   geom_sf(data = UK_map_teams, fill = "blue", alpha = 0.4) +
18   theme_bw(base_size = 30) +
19   theme(axis.title = element_text(size = 30),
20         legend.title = element_text(size = 22),
21         legend.key.size = unit(40, "pt"),
22         plot.title = element_text(size = 40)) +
23   geom_point(data = locations_2023, aes(x = Latitude,
24                                         y = Longitude,
25                                         color = Standing),
26             size = 3) +
27   geom_text_repel(data = locations_2023,
28                 aes(x = Latitude,
29                     y = Longitude,
30                     color = Standing,
31                     label = `Team Name`),
32                 size = 7, max.overlaps = 15, force = 2) +
33   scale_colour_gradient(low = "forestgreen", high = "red",
```

```
34     guide = guide_colorbar(reverse = TRUE)) +  
35     coord_sf(xlim = c(-6, 2), ylim = c(50, 56), expand = FALSE) +  
36     labs(title = "English Premier League Teams",
```

Stadium coordinate data was acquired manually from the EPL website.

The English Premier League

English Premier League Teams



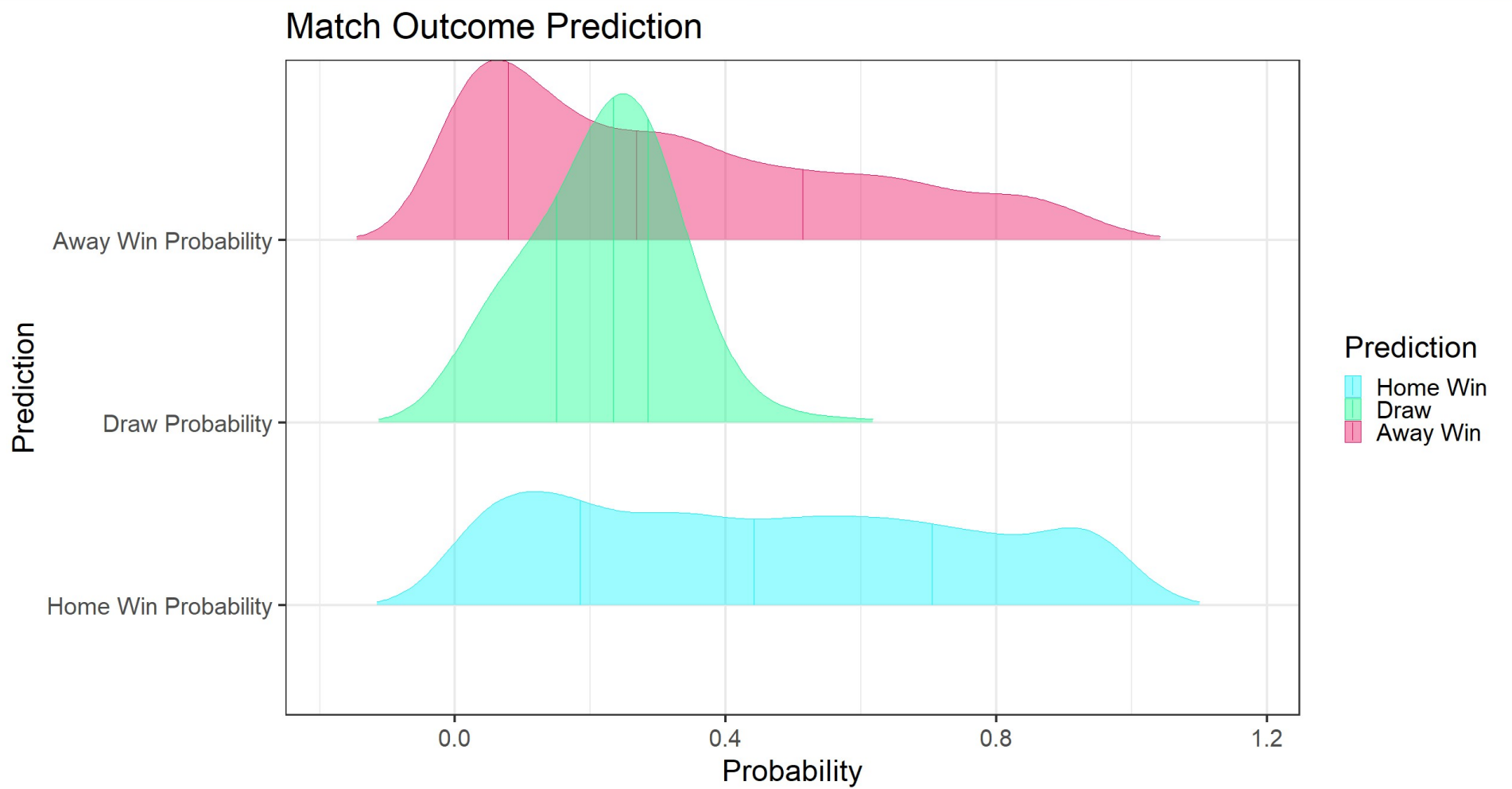
Home Team Advantage

- We will begin our analysis of the EPL 2022-23 season by determining how much of an advantage the home team gets.

```
1 library(ggthemes)
2 options(width = 100)
3 results <- understat_league_match_results(league = "EPL",
4                                           season_start_year = 2022)
5 results <- results |>
6   rename(`Home Win Probability` = forecast_win,
7          `Draw Probability` = forecast_draw,
8          `Away Win Probability` = forecast_loss)
9
10 long_data <- pivot_longer(
11   results,
12   cols = c(`Home Win Probability`,
13            `Draw Probability`,
14            `Away Win Probability`),
15   names_to = "Prediction",
16   values_to = "probability"
17 )
18
19 long_data$Prediction <- factor(
20   long_data$Prediction,
21   levels = c("Home Win Probability",
22              "Draw Probability",
23              "Away Win Probability")
24 )
```

```
1 ggplot(long_data, aes(x = probability, y = Prediction,
2                       fill = Prediction, color = Prediction)) +
3   geom_density_ridges(
4     alpha = 0.4,
5     rel_min_height = 0.01,
6     quantile_lines = TRUE,
7     quantiles = c(0.25, 0.5, 0.75)
8   ) +
9   labs(
10    title = "Match Outcome Prediction",
11    x = "Probability",
12    y = "Prediction"
13  ) +
14  scale_fill_manual(
15    values = c(`Home Win Probability` = "#04f5ff",
16               `Draw Probability` = "#00ff85",
17               `Away Win Probability` = "#e90052"),
18    labels = c("Home Win", "Draw", "Away Win")
19  ) +
20  scale_color_manual(
21    values = c(`Home Win Probability` = "#04f5ff",
22               `Draw Probability` = "#00ff85",
23               `Away Win Probability` = "#e90052"),
24    labels = c("Home Win", "Draw", "Away Win")
25  ) +
26  theme_ridges() +
27  theme(
28    plot.title = element_text(hjust = 0.5, margin = margin(10, 0, 0, 0)),
29    axis.title.x = element_text(hjust = 0.5, margin = margin(0, 0, 0, 0)),
30    axis.title.y = element_text(hjust = 0.5, margin = margin(0, 0, 0, 0)),
31    axis.text.y = element_blank(),
32    axis.ticks.y = element_blank()
```

Home Team Advantage



Away Team Travel Effect

- Data on the travel distance of Premier League clubs in the 2022-23 season was acquired from [footballteamnews](#).

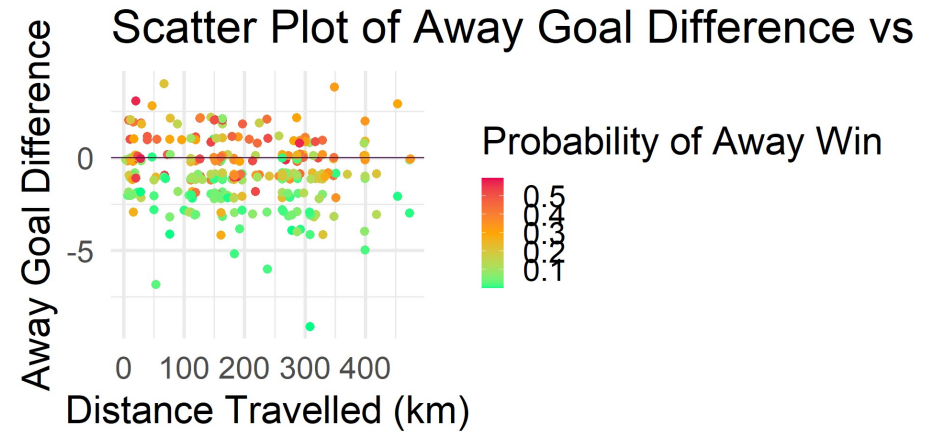
```
1 distances <- read_csv("data/England_Distances - 2023.csv")
```

```
1 find_distance <- function(home_team, away_team, distance_ma
2   teams <- colnames(distance_matrix)[2:21]
3   home_index <- which(teams == home_team)
4   away_index <- which(teams == away_team)
5   if(length(home_index) > 0 && length(away_index) > 0) {
6     return(pull(distance_matrix[home_index, away_index + 1]
7   } else {
8     return(NA)
9   }
10 }
11 modified_results <- results |>
12   rowwise() |>
13   mutate(
14     Distance_Travelled = find_distance(home_team,
15                                       away_team, distance_ma
16     Away_Goal_Diff = away_goals - home_goals
17   ) |>
18   ungroup()
```

Away Team Travel Effect

- Data was filtered to include games that did not have a high probability of the away team winning.

```
1 filtered_results <- modified_results |>
2   filter(`Away Win Probability` < 0.6)
3
4 ggplot(filtered_results, aes(x = Distance_Travelled,
5                             y = Away_Goal_Diff)) +
6   geom_jitter(aes(color = `Away Win Probability`, size = 1),
7             width = 0.2, height = 0.2) +
8   scale_color_gradient2(low = "#00ff85",
9                         mid = "orange",
10                        high = "#e90052",
11                        midpoint = 0.3,
12                        name = "Probability of Away Win") +
13   geom_hline(yintercept = 0, color = "#38003c") +
14   labs(title = "Scatter Plot of Away Goal Difference vs Distance Travelled (km)",
15        x = "Distance Travelled (km)",
16        y = "Away Goal Difference") +
17   theme_minimal(base_size = 30)
```



Away Team Effect on Shotstyle

Away Team Effect on Shotstyle

- Analyzing the effects of being the away team on shooting:

1 EPL-Shots <- understat_league_season_shots(league = "EPL", season_start_year = 2022)

	X	Y	xG	h_a	home_team	away_team	league	id	minute	result
2078	0.981	0.629	0.3262840	h	Brighton	Tottenham	EPL	491814	68	SavedShot
7361	0.978	0.428	0.0756614	a	Southampton	Manchester City	EPL	518223	40	MissedShots
6785	0.863	0.507	0.1013050	a	Arsenal	Leeds	EPL	516429	64	MissedShots
664	0.899	0.515	0.0207121	h	Leeds	Chelsea	EPL	482403	28	MissedShots
171	0.926	0.328	0.0705736	h	Manchester United	Brighton	EPL	479580	6	BlockedShot
1778	0.858	0.478	0.0583403	a	Fulham	Newcastle United	EPL	490163	4	SavedShot
4946	0.916	0.486	0.1384490	h	Arsenal	Manchester	EPL	505089	63	BlockedShot

	X	Y	xG	h_a	home_team	away_team	league	id	minute	result
						United				
3440	0.842	0.752	0.0328334	a	Tottenham	Liverpool	EPL	498841	35	BlockedShot
3549	0.866	0.626	0.0564314	h	Nottingham Forest	Crystal Palace	EPL	500158	73	MissedShots
2975	0.945	0.533	0.1269200	a	Bournemouth	Tottenham	EPL	496616	91	BlockedShot

Away Team Effect on Shotstyle

- To create a proper heat-map, we need to wrangle our data into the proper shape:

```
1 create_heatmap_data <- function(data) {  
2   x_breaks <- seq(0.65, 1, by = 0.025)  
3   y_breaks <- seq(0.35, 0.65, by = 0.025)  
4   grid <- expand_grid(xmin = head(x_breaks, -1),  
5                       xmax = tail(x_breaks, -1),  
6                       ymin = head(y_breaks, -1),  
7                       ymax = tail(y_breaks, -1))  
8   heatmap_data <- grid |>  
9     rowwise() |>  
10      mutate(shot_count = sum(data$X >= xmin & data$X < xmax  
11                             data$Y < ymax),  
12             goal_count = sum(data$X >= xmin & data$X < xmax  
13                               data$Y < ymax & data$result == "goal"),  
14             goal_percentage = ifelse(shot_count > 0,  
15                                     (goal_count / shot_count) * 100,  
16                                     0)  
17   return(heatmap_data)  
18 }
```

```
1 home_shots <- EPL_Shots |> filter(h_a == "h")  
2 away_shots <- EPL_Shots |> filter(h_a == "a")  
3 home_shots_filtered <- home_shots |>  
4   filter(X >= 0.65 & X <= 1, Y >= 0.35 & Y <= 0.65)  
5  
6 away_shots_filtered <- away_shots |>  
7   filter(X >= 0.65 & X <= 1, Y >= 0.35 & Y <= 0.65)  
8  
9 home_heatmap_data <- create_heatmap_data(home_shots_filtered)  
10 away_heatmap_data <- create_heatmap_data(away_shots_filtered)  
11 home_heatmap_data <- home_heatmap_data |>  
12   rename(home_goal_percentage = goal_percentage)  
13 away_heatmap_data <- away_heatmap_data |>  
14   rename(away_goal_percentage = goal_percentage)  
15 merged_data <- full_join(  
16   home_heatmap_data,  
17   away_heatmap_data,  
18   by = c("xmin", "xmax", "ymin", "ymax")  
19 )  
20 merged_data <- merged_data |>  
21   mutate(  
22     goal_percentage_diff = coalesce(home_goal_percentage, 0) -  
23                               coalesce(away_goal_percentage, 0)  
24   )  
25 goal_percentage_diff_data <- merged_data |>  
26   select(xmin, xmax, ymin, ymax, home_goal_percentage,  
27          away_goal_percentage, goal_percentage_diff)
```

Away Team Effect on Shotstyle

- Then we make the heat-map function for the percentage of shots that are made:

```

1 create_heatmap_plot <- function(data, title, colors, values) {
2   data <- data |>
3     mutate(
4       x = (xmin + xmax) / 2,
5       y = (ymin + ymax) / 2
6     )
7
8   ggplot(data, aes(x = x, y = y)) +
9     geom_tile(aes(fill = goal_percentage_diff), color = "white") +
10    scale_fill_gradientn(
11      colors = colors,
12      limits = c(-0.3, 0.3),
13      name = "Goal % Difference"
14    ) +
15    geom_segment(
16      aes(x = 0.82, xend = 0.82, y = 0.35, yend = 0.65),
17      color = "black", linewidth = 0.75
18    ) +
19    geom_segment(
20      aes(x = 0.94, xend = 1, y = 0.5458, yend = 0.5458),
21      color = "black", linewidth = 0.75
22    ) +
23    geom_segment(
24      aes(x = 0.94, xend = 1, y = 0.4542, yend = 0.4542),
25      color = "black", linewidth = 0.75
26    ) +
27    geom_segment(
28      aes(x = 0.94, xend = 0.94, y = 0.5458, yend = 0.4542),
29      color = "black", linewidth = 0.75
30    ) +
31    labs(title = title) +
32    theme_minimal(base_size = 25) +
33    theme(
34      panel.background = element_blank(),
35      panel.grid.major = element_blank(),
36      panel.grid.minor = element_blank(),
37      axis.text.x = element_blank()

```

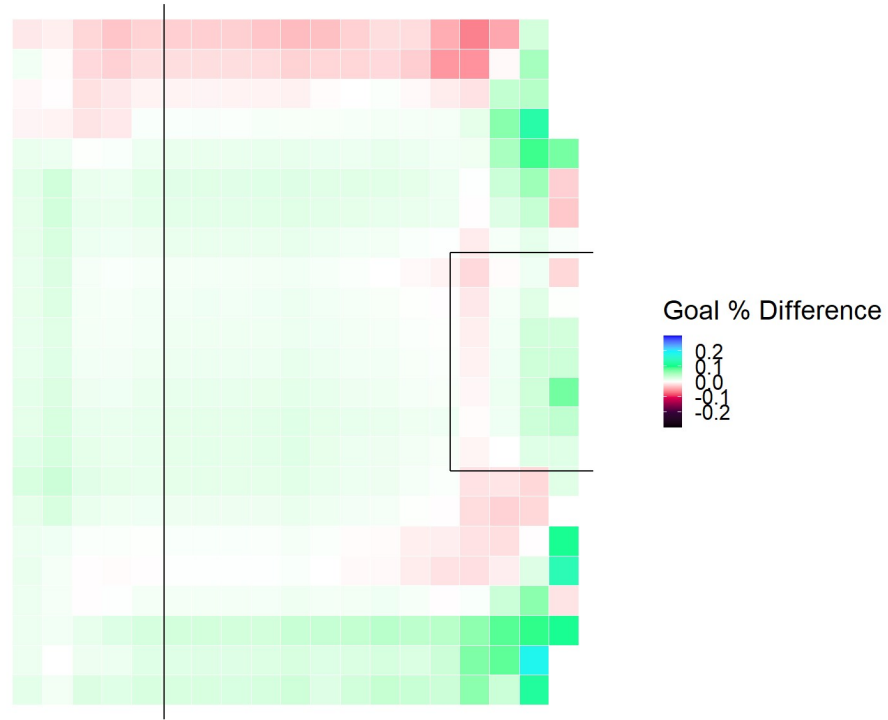
```

1 goal_percentage_diff_range <- range(goal_percentage_diff_data)
2
3 p_colors <- c("black", EPL_Colors[5], EPL_Colors[2], EPL_Colors[3],
4             EPL_Colors[4], EPL_Colors[1], "blue")
5 p_values <- c(-0.3, -0.1, -0.05, 0, 0.05, 0.1, 0.3)

```

Away Team Effect on Shotstyle

Goal Percentage Difference (Home - Away)



Away Team Effect on Shotstyle

- Using a similar function, we can produce the heat-maps for the number of shots taken:

```

1 create_shot_heatmap_data <- function(data) {
2   x_breaks <- seq(0.65, 1, by = 0.025)
3   y_breaks <- seq(0.35, 0.65, by = 0.025)
4   grid <- expand.grid(xmin = head(x_breaks, -1),
5                       xmax = tail(x_breaks, -1),
6                       ymin = head(y_breaks, -1),
7                       ymax = tail(y_breaks, -1))
8   heatmap_data <- grid |>
9     rowwise() |>
10    mutate(shot_count = sum(data$X >= xmin & data$X < xmax
11    ungroup())
12
13   return(heatmap_data)
14 }
15
16 home_shots_filtered <- home_shots |> filter(X >= 0.65 & X < 0.82)
17 away_shots_filtered <- away_shots |> filter(X >= 0.65 & X < 0.82)
18 home_shot_heatmap_data <- create_shot_heatmap_data(home_shots_filtered)
19 away_shot_heatmap_data <- create_shot_heatmap_data(away_shots_filtered)
20 mean_shot_count <- mean(home_shot_heatmap_data$shot_count)
21 sd_shot_count <- sd(home_shot_heatmap_data$shot_count)
22 shot_count_colors <- c("white", "yellow", "orange", "red")

```

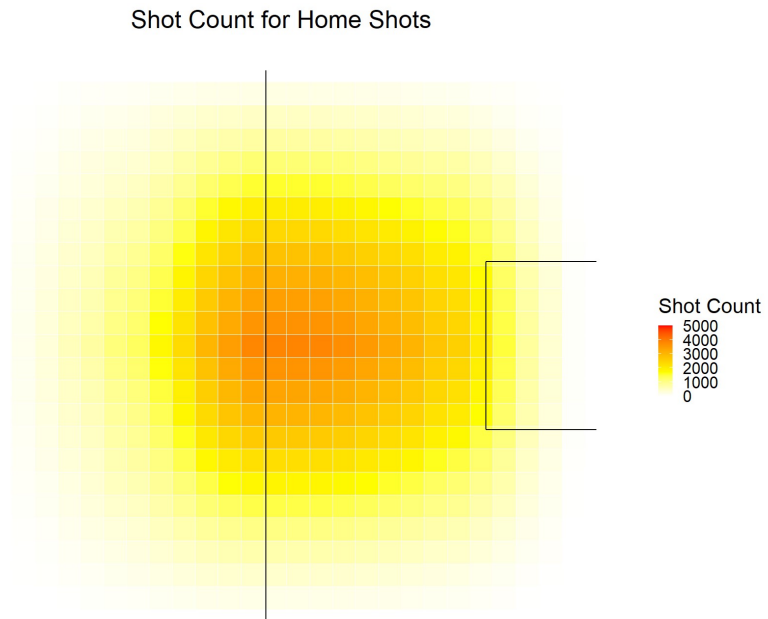
```

1 create_shot_heatmap_plot <- function(data, title, colors)
2   data <- data |>
3     mutate(x = (xmin + xmax) / 2,
4            y = (ymin + ymax) / 2)
5
6   ggplot(data, aes(x = x, y = y)) +
7     geom_tile(aes(fill = shot_count), color = "white") +
8     scale_fill_gradientn(
9       colors = colors,
10      values = scales::rescale(c(0, 1, 2, 3)),
11      limits = c(0, 5000),
12      breaks = seq(0, 5000, by = 1000),
13      name = "Shot Count"
14    ) +
15    labs(title = title) +
16    geom_segment(aes(x = 0.82, xend = 0.82, y = 0.35, yend = 0.65,
17      color = "black", size = 0.75) +
18    geom_segment(aes(x = 0.94, xend = 1, y = 0.5458, yend = 0.65,
19      color = "black", size = 0.75) +
20    geom_segment(aes(x = 0.94, xend = 1, y = 0.4542, yend = 0.65,
21      color = "black", size = 0.75) +
22    geom_segment(aes(x = 0.94, xend = 0.94, y = 0.5458, yend = 0.65,
23      color = "black", size = 0.75) +
24    coord_fixed(ratio = 1) +
25    theme_minimal(base_size = 25) +
26    theme(
27      panel.background = element_blank(),
28      panel.grid.major = element_blank(),
29      panel.grid.minor = element_blank(),
30      axis.text.x = element_blank(),
31      axis.text.y = element_blank(),
32      axis.ticks = element_blank(),
33      axis.title.x = element_blank(),
34      axis.title.y = element_blank(),
35      plot.title = element_text(hjust = 0.5)
36    )
37 }

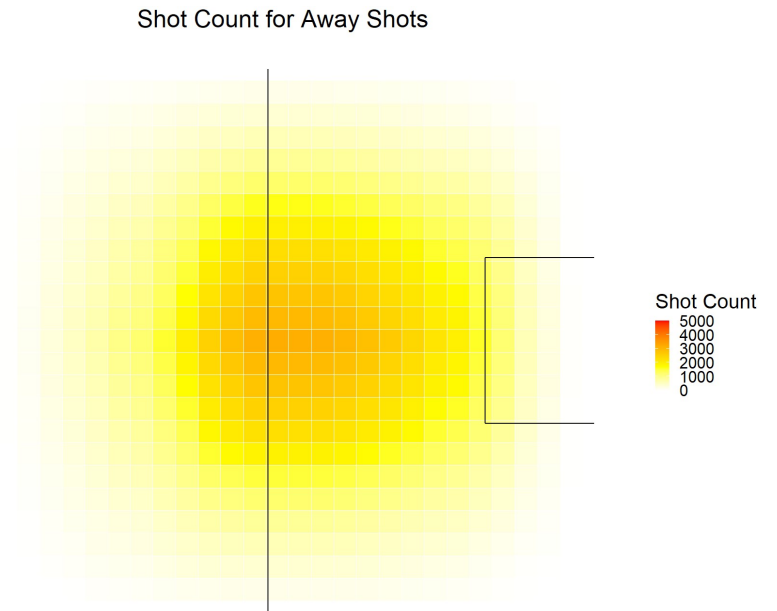
```


Away Team Effect on Shotstyle

- For home teams:



- For away teams:



Home-Away Points Breakdown

Home-Away Points Breakdown

- In soccer, a team's points over the season are calculated as **three times their number of wins plus their number of draws**.
- We will now compare the number of points earned at home games vs away games:

```
1 games <- results |>
2   select(home_team, away_team, home_goals, away_goals) |>
3   distinct()
4 unique_games <- results |>
5   select(home_team, away_team, home_goals, away_goals) |>
6   distinct()
7 unique_games <- unique_games |>
8   mutate(
9     home_points = if_else(home_goals > away_goals, 3, if_else(home_goals == away_goals, 1, 0)),
10    away_points = if_else(away_goals > home_goals, 3, if_else(away_goals == home_goals, 1, 0))
11  )
```

Home-Away Points Breakdown

Some more data manipulation:

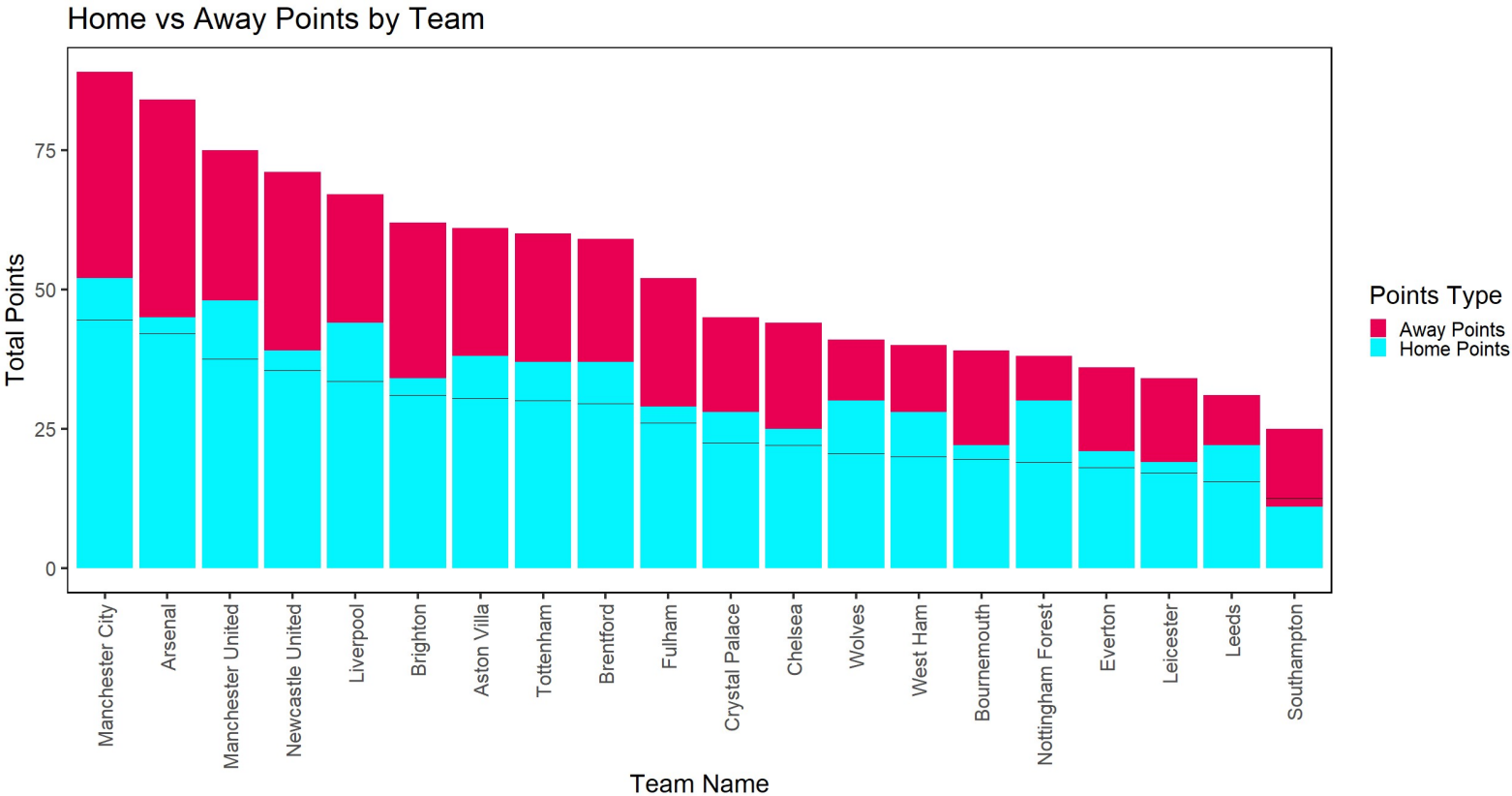
```
1 home_points_df <- unique_games |>
2   group_by(home_team) |>
3   summarize(home_points = sum(home_points, na.rm = TRUE)) |>
4   rename(team = home_team)
5 away_points_df <- unique_games |>
6   group_by(away_team) |>
7   summarize(away_points = sum(away_points, na.rm = TRUE)) |>
8   rename(team = away_team)
9 team_points_df <- full_join(home_points_df, away_points_df, by = "team")
10 team_points_df <- team_points_df |>
11   mutate(
12     home_points = replace_na(home_points, 0),
13     away_points = replace_na(away_points, 0),
14     total_points = home_points + away_points,
15     proportion = home_points / away_points,
16     half_points = total_points / 2
17   )
18 team_points_df <- team_points_df |>
19   mutate(proportion = ifelse(is.infinite(proportion), NA, proportion))
20 team_points_df <- team_points_df |>
21   arrange(desc(proportion))
```

Home-Away Points Breakdown

- The proportion of points won at home games to those won at away games:

```
1 team_points_df <- team_points_df |>
2   arrange(desc(total_points)) |>
3   mutate(team = ifelse(team == "Wolverhampton Wanderers", "Wolves", team))
4
5 team_points_df <- team_points_df |>
6   arrange(desc(total_points))
7 long_team_points_df <- team_points_df |>
8   pivot_longer(cols = c(home_points, away_points),
9               names_to = "type", values_to = "points")
10 long_team_points_df$team <- factor(long_team_points_df$team,
11                                   levels = team_points_df$team)
12
13 ggplot(long_team_points_df, aes(x = team, y = points, fill = type)) +
14   geom_bar(stat = "identity") +
15   geom_errorbar(aes(ymin = half_points, ymax = half_points),
16               width = 0.9, color = "black") +
17   scale_fill_manual(values = c("home_points" = "#04f5ff",
18                               "away_points" = "#e90052"),
19                   labels = c("home_points" = "Home Points",
20                              "away_points" = "Away Points")) +
21   labs(
22     title = "Home vs Away Points by Team",
23     x = "Team Name",
24     y = "Total Points",
25     fill = "Points Type"
26   ) +
27   theme_bw(base_size = 25) +
28   theme(axis.text.x = element_text(angle = 90, vjust = .5, hjust = 1),
29         panel.background = element_blank(),
30         panel.grid.major = element_blank(),
31         panel.grid.minor = element_blank(),
32         panel.border = element_rect(color = "black", fill = NA))
```

Home-Away Points Breakdown



Salary Analysis

Salary Data

Salary Data

```
1 premier_league_salaries = premier_league_salaries |>
2   mutate(
3     Pos = str_sub(Pos, 1, 2),
4     Pos = factor(Pos, levels = c("GK", "DF", "MF", "FW"))
5   ) |>
6   filter(Pos != "")
```

```
1 highest_earners = premier_league_salaries |>
2   group_by(Pos) |>
3   slice(which.max(AnnualWageUSD))
```

Salary Data

```
1 ggplot(premier_league_salaries, aes(x = AnnualWageUSD, y = Pos, color = Pos)) +  
2   geom_point(size = 3) +  
3   labs(  
4     title = "Premier League '23 Salaries by Position",  
5     x = "Salary (USD)",  
6     y = "Position"  
7   ) +  
8   theme_minimal() +  
9   geom_text(data = highest_earners,  
10            aes(label = paste(Player, Team, sep = "\n")),  
11              vjust = 1.5,  
12              color = "black",  
13              size = 4.5  
14            ) +  
15   guides(color = "none") +  
16   scale_x_continuous(labels = unit_format(unit = "M", scale = 1e-6)) +  
17   expand_limits(x = 38000000, y = 0) +  
18   theme_bw(base_size = 20)
```

