

# Final Report

Dylan Li, Liam Quach, Brendan Callender

## I. Introduction

The English Premier League (EPL) is the top tier of professional football (soccer) in England and is considered one of the most popular and competitive leagues in the world. The league is made up of twenty clubs (teams) that compete over a season for the Premier League title with new clubs added each year via a system of promotion and relegation. Each year, three new clubs are promoted from the second division based on the previous year's results with these promoted teams replacing the bottom three teams from the previous year's Premier League season.

Over the course of a season, each team plays a total of 38 matches, facing every other team twice—once at home and once away. Teams are rewarded points from each game as follows: 3 points for a win, 1 points for a draw, and 0 points for a loss. The team with the most points at the end of the 38-game season is crowned as the Premier League Champions.

add a little more maybe?

For our project, we are interested in exploring the following research questions:

1. What factors are associated with higher or lower point totals in the English Premier league?
2. Is spending more money in the off-season associated with earning more points the following season?
3. How do differences in expected metrics to actual metrics impact a clubs point totals?

## II. Data Source & Methods

To answer our research questions, we collected English Premier League season-level data spanning from the 2017-2018 season up to the most recently completed 2023-2024 season. Data was collected from two sites: [fbref.com](https://fbref.com) and [transfermarkt.com](https://transfermarkt.com). The data collected from fbref includes performance related metrics for each team over the season as predictors as well as point totals for each team at the end of the season for our response variable. The performance metrics include total goals scored, total goals conceded, expected goals scored, expected goals

conceded, average % possession, shooting metrics and more. The data collected from transfermarkt includes data relating to each teams expenditure and sales with respect to buying and selling players in the transfermarkt. This data includes money spent, money gained from sales, net spend, number of players bought, number of players sold and more. Money related variables are measured in thousands of euros.

Predictors relating to season totals such as goals scored and goals conceded were scaled down to per 90/ per game values for better interpretability. This was achieved by dividing these metrics by the total games played which is 38.

Table 1: Description of Dataset Variables

Variable	Role	Range of Values
Points	Response	(16, 100)
Goals/90	L1	(0.52, 2.78)
Goals Against/90	L1	(0.58, 2.74)
Average Possession of Ball (%)	L1	(35.4, 71.0)
...	...	...
Net Spend (in €1,000,000)	L1	(-118.07, 562.39)
Club Average Net Spend (in €1,000,000)	L2	(-7.72, 139.39)
Actual vs Expected Goals/90 Difference	L1	(-0.37, 0.72)
Actual vs Expected Goals/90 Against Difference	L1	(-0.37, 0.72)

See example rows of data below:

Table 2: Example Rows from Dataset

Club	Season	Pts	GF	GA	Poss	...	NetSpend	Mean_NetSpend
Chelsea	2017	70	1.63	1.0	55.6	...	65.9	139.0
Arsenal	2017	63	1.95	1.34	61.4	...	-9.55	100.0
Everton	2017	49	1.16	1.53	45.5	...	76.8	26.0

To analyze the data, we will employ multi-level regression models, also known as hierarchical linear models. This approach is well-suited for the structure of the dataset, in which we have repeat observations for different clubs over several seasons. (See figure below)

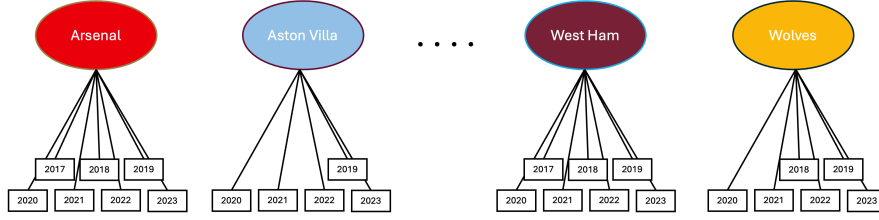


Figure 1: Multi-level Structure of Data

### III. Results

#### Exploratory Data Analysis

This section presents the exploratory data analysis conducted to understand the key relationships between variables in the dataset. This exploratory data analysis was conducted before the model fitting process to gain an initial understanding of our research questions.

*Figure 2* below shows the joint distribution of goals scored per game and goals conceded per game, colored by season point totals. We there is a strong, negative correlation between goals scored per game and goals conceded per game. This means that teams who tend to score more, also tend to conceded less as well. When considering the season point totals, we see that decreasing the number of goals conceded per game is associated with higher point totals holding goals scored constant. Additionally, increasing the number of goals scored per game is associated with higher point totals holding goals conceded constant. Lastly, we see that jointly decreasing the number of goals conceded per game and increasing the number of goals scored per game is associated with the largest increase in season point totals.

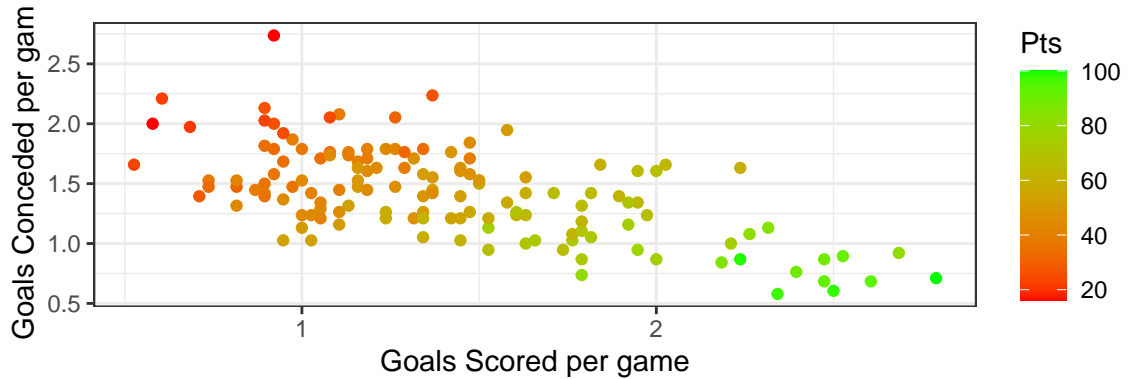


Figure 2: Impact of Goals Scored/90 and Goals Conceded/90 on Season Point Totals

*Figure 3* below shows the relationship between average % possession of the ball and season point totals. The plot shows a strong, positive association between % possession and season points with higher values for % possession associated with higher point totals. This makes sense intuitively because teams with more possession tend to have the ball more which reduces the chances of the opposing team scoring and gives your team more chances to score goals.

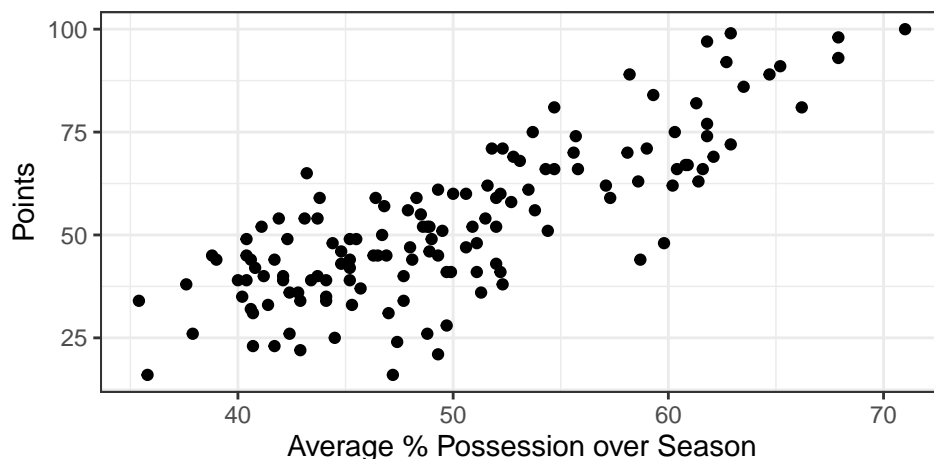


Figure 3: Season Point Totals by Season Average % Possession

Lastly, the plots below in *Figure 4* and *Figure 5* show the impact of spending on season point totals. The plot on the left shows the relationship between season point totals and the net spend of the club for individual seasons. Larger values for net spend represent a club spending more money on new players while smaller values indicate a club spending less money spent with negative values indicating a team made profit selling players in the market. From the plot, we see a weak positive association with teams who spend more money being associated with higher point totals. We also notice a major outlier in the plot with Chelsea in the 2023-2024 season. This is a valid data point and represents the season in which Chelsea had new owners invest large amounts of money into the team. This is not normal behavior for when teams get new owners and can serve as an example of how making too many changes to a team has a negative impact on performance.

In *Figure 5*, we see a much stronger positive association between net spend and points when aggregated for each club. This demonstrates that consistent investment into a team over many seasons is more strongly associated with higher point totals than just a single season of large investment. (As they say... Rome wasn't built in a day)

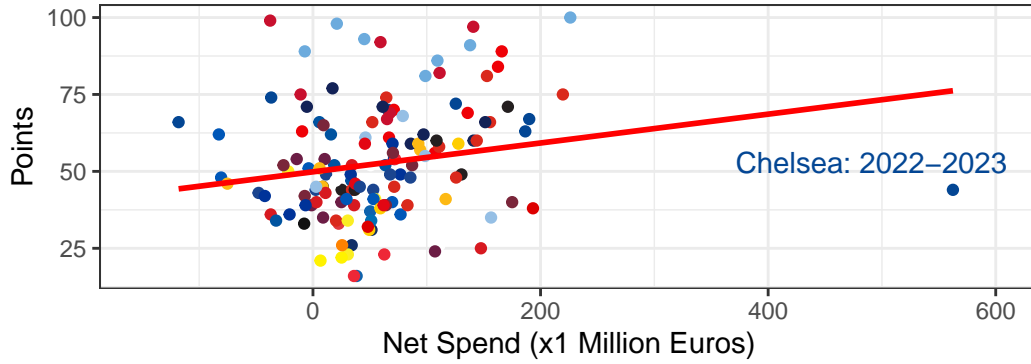


Figure 4: Season Point Totals by Single Season Net Spend

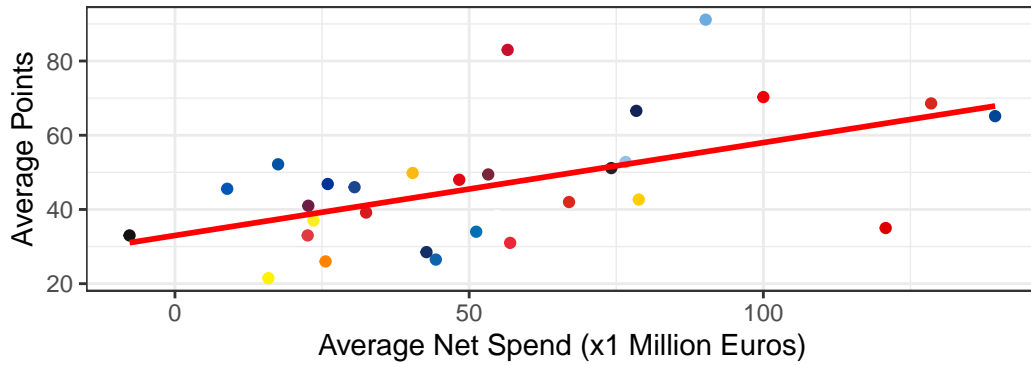


Figure 5: Club Average Points by Club Average Net Spend

## ANOVA

After performing exploratory data analysis, an initial Analysis of Variance (ANOVA) test was performed to explore whether there was significant club-to-club variability in season point totals (See results in *Table 3* below). Looking at the p-value resulting from the ANOVA, we have significant evidence that at least 2 clubs have different mean point totals. This is supported by *Figure 6* below which shows the distribution of point totals by each club. We see clubs like Manchester City have very high point totals while clubs like West Brom and Norwich City have very low point totals.

Table 3: *ANOVA for Significance of Club-to-Club Variability*

	df	SSE	MSE	F-Statistic	P-Value
Club	29	37233	1283.89	12.848	< 0.0001

	df	SSE	MSE	F-Statistic	P-Value
Residuals	110	10992	99.93		

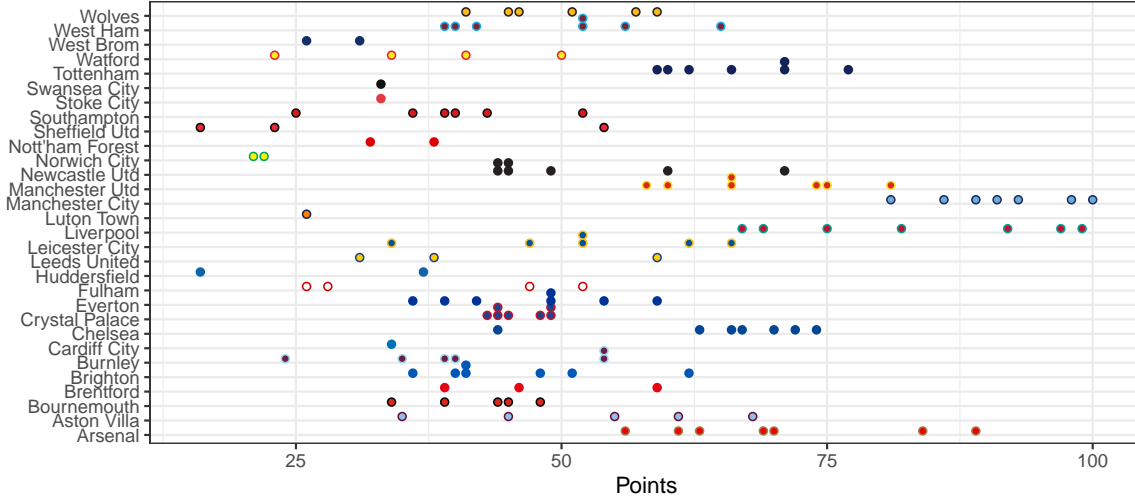


Figure 6: Distribution of Point Totals by Team

### Null Model

After finding significant club-to-club variability in the season point totals, we fit an initial null model which only includes the club as a random effect. The model can be written out as seen below:

$$Points_{ij} = \beta_{00} + u_j + \epsilon_{ij}$$

where  $u_j \sim N(0, \tau_0^2)$  and  $\epsilon_{ij} \sim N(0, \sigma^2)$

The summary of the null model output can be found below in *Table 4*. Looking at the ICC of the model, we see that approximately 72% of the variation is at the club level while 28% of the variation is within each club. This matches what we saw in *Figure 6* above with clubs point totals being similar across the seasons. Additionally, we see the model seconds the fact that there is significant club-to-club variation with the confidence interval for  $\tau_0$  not containing 0. Lastly, we see the resulting predictions from the null model in *Figure 7* below. The black points represent the mean points for each club while the red points represent the predicted points for each club. The distance from the red and black points represent the shrinkage that occurs when using a multi-level model. We can see a club like Luton Town has large shrinkage towards the mean because we only have one season of data for Luton.

Table 4: Null Model Summary

Parameter/Statistic	Estimate
$\sigma^2$	255.3
$\tau_0^2$	99.6
95% CI for $\tau_0$	(11.99, 21.18)
ICC	0.72

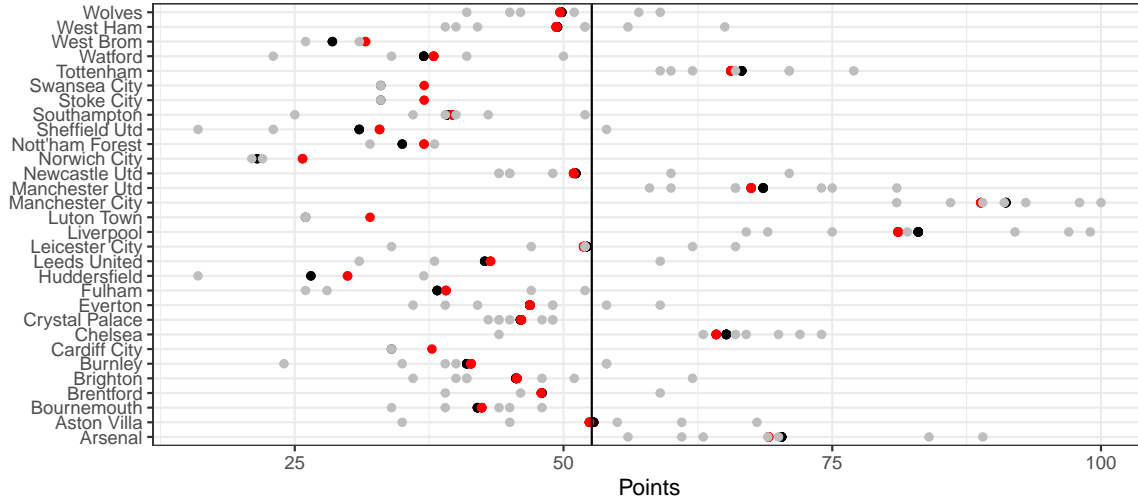


Figure 7: Null Model Predictions

### Model Fitting Process

Below shows the process of how we arrived at our final model. We first began by including goals for and goals against in our first model due to our EDA showing a strong joint association with points. We found these variables to be significant so we included them in each subsequent model. From there we continued to explore new models by adding different level 1 predictors. If a predictor was significant, it was left in the model. After exploring level 1 predictors, we added our only level 2 predictor to the model which is average net spend for each club. We found this to be significant and moved on to adding random slopes to the model.

```
model1 <- lmer(Pts ~ GF + GA + (1 | Club), data = prem)
summary(model1)
```

Add Level 1 Predictors

```
model2_1 <- lmer(Pts ~ GF + GA + xG_diff + xGA_diff + (1 | Club), data = prem)
summary(model2_1)
```

```
model2_2 <- lmer(Pts ~ GF + GA + NetSpend + (1 | Club), data = prem)
summary(model2_2)
```

```
model2_3 <- lmer(Pts ~ GF + GA + Poss + (1 | Club), data = prem)
summary(model2_3)
```

```
model2_4 <- lmer(Pts ~ GF + GA + Age + (1 | Club), data = prem)
summary(model2_4)
```

Add Level 2 Predictors

```
model2_5 <- lmer(Pts ~ GF + GA + Mean_NetSpend + (1 | Club), data = prem)
summary(model2_5)
```

```
model2_6 <- lmer(Pts ~ GF + GA + NetSpend + Mean_NetSpend + (1 | Club), data = prem)
anova(model2_5, model2_6)
```

Random Slopes

```
model3_1 <- lmer(Pts ~ GF + GA + Mean_NetSpend + (1 + GF | Club), data = prem)
anova(model2_5, model3_1)
```

```
model3_2 <- lmer(Pts ~ GF + GA + Mean_NetSpend + (1 + GA | Club), data = prem)
anova(model2_5, model3_2)
```

## Final Model

$$Points_{ij} = \beta_{00} + u_j + \beta_1(G/90)_{ij} + \beta_2(GA/90)_{ij} + \beta_3(\overline{NetSpend})_j + \epsilon_{ij}$$

$$\hat{\sigma}^2 = 19.96$$

$$\hat{\tau}_0^2 = 0$$

Level 2 Variability Explained = 100%

Level 1 Variability Explained = 92.2%

boundary (singular) fit: see help('isSingular')



```
Linear mixed model fit by REML ['lmerMod']
Formula: Pts ~ GF + GA + Mean_NetSpend + (1 | Club)
Data: prem
```

REML criterion at convergence: 815.9

Scaled residuals:

Min	1Q	Median	3Q	Max
-2.8271	-0.6048	0.1171	0.6113	3.4867

Random effects:

Groups	Name	Variance	Std.Dev.
Club	(Intercept)	0.00	0.000
	Residual	19.96	4.467

Number of obs: 140, groups: Club, 30

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	49.11412	2.95766	16.606
GF	23.03892	1.05930	21.749
GA	-21.86604	1.29747	-16.853
Mean_NetSpend	0.03120	0.01181	2.643

Correlation of Fixed Effects:

	(Intr)	GF	GA
GF	-0.769		
GA	-0.919	0.557	
Mean_NtSpnd	-0.087	-0.362	0.056

optimizer (nloptwrap) convergence code: 0 (OK)  
boundary (singular) fit: see help('isSingular')

technical writing stuff here

## Model Diagnostics

## V. Discussion

answer research questions

limitations

strengths and weaknesses

future steps

## VI. Appendix

to do: add all model code stuff

add variable labels:

### ANOVA

```
# A tibble: 2 x 6
  term      df  sumsq meansq statistic  p.value
<chr>   <int> <dbl> <dbl>    <dbl>    <dbl>
1 Club      29 37233. 1284.    12.8 1.10e-23
2 Residuals 110 10992.  99.9     NA    NA
```