# PREDICTING FOOTBALL TABLES BY A MAXIMALLY PARSIMONIOUS MODEL

KJETIL K. HAUGEN AND BRYNJULF OWREN

*Abstract.* This paper presents some useful mathematical results concerned with football table prediction. In addition, some empirical results indicate that an alternative methodology for football table prediction may produce high quality forecasts with far less resource usage than conventional methods.

## 1. INTRODUCTION

Former England international, long-time Arsenal player and present SKYSPORTS commentator Paul Merson, predicts the final Premiel League (PL) table for the 2016/2017 season in [12]. With the hindsight of time, we can check Merson's predictions compared with the true final table as indicated in Table 1.

In Table 1, the final table outcome is given in the leftmost column (*Final PL-table*), while Merson's predictions are given in the mid column (*Merson's predictions*). By defining the true (correct) final table as the consecutive integers $\{1, 2, \ldots, n^1\}$, and a table prediction as a certain permutation $P(i)^2$ of the integers $\{1, 2, \ldots, n\}$, the absolute deviations between forecasts and true values can be computed as in the rightmost column in Table 1.

If we examine Merson's tips closer, we observe that he obtained two zeros (or perfect hits) in the rightmost column in Table 1 – Chelsea as the winner, and Swansea as number 15. Furthermore, he only missed by one placement for 8 outcomes, but also missed with greater margin for instance for Bournemouth which he thought should finish at 18th place while, in fact, it ended at 9th.

The question that we will be interested in initially, is the quality of Merson's predictions. Is Merson's permutation in Table 1 a good guess? In order to attempt to answer such a question, we need to define quality. It seems reasonable to look

---

[1]$n$ is the number of teams in the league.

[2]In this example, $P(i)$ denotes Merson's predictions.

**Table 1.** Paul Merson's predictions compared to true final Premier League table..

|     | Final PL-table    | Merson's predictions | $|P(i) - i|$ |
| --- | ----------------- | -------------------- | ------------ |
| 1.  | Chelsea           | 1                    | 0            |
| 2.  | Tottenham         | 6                    | 4            |
| 3.  | Manchester City   | 2                    | 1            |
| 4.  | Liverpool         | 5                    | 1            |
| 5.  | Arsenal           | 4                    | 1            |
| 6.  | Manchester United | 3                    | 3            |
| 7.  | Everton           | 8                    | 1            |
| 8.  | Southampton       | 9                    | 1            |
| 9.  | Bournemouth       | 18                   | 9            |
| 10. | West Bromwich     | 17                   | 7            |
| 11. | West Ham          | 7                    | 4            |
| 12. | Leicester         | 11                   | 1            |
| 13. | Stoke             | 10                   | 3            |
| 14. | Crystal Palace    | 12                   | 2            |
| 15. | Swansea           | 15                   | 0            |
| 16. | Burnley           | 20                   | 4            |
| 17. | Watford           | 16                   | 1            |
| 18. | Hull              | 19                   | 1            |
| 19. | Middlesbrough     | 13                   | 6            |
| 20. | Sunderland        | 14                   | 6            |

for some function;

$$f(|P(1) - 1|, |P(2) - 2|, \ldots, |P(n) - n|)$$

which produces a single numerical value tailored for comparison. Of course, infinite possibilities exist for such a function. Fortunately, forecasting literature comes to rescue – refer for instance to [9]. The two most common measures used in similar situations are $MAE$, Mean Absolute Error or $MSE$, Mean Squared Error. With our notation, these two measures are defined as:

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |P(i) - i|, \; MSE = \frac{1}{n} \sum_{i=1}^{n} (P(i) - i)^2. \tag{1.1}$$

Although $MSE$ is more applied in statistics, preferably due to its obvious nicer mathematical properties[3], we choose to use MAE. It weighs errors equally, and it produces also an easily interpretable result; a MAE of 3 means that a prediction on average misplaces all teams by 3 places.

Given this choice, Merson's $MAE$ in Table 1 can be easily calculated as; $MAE = 2.8$. Still, we are in no position to give any statements on the quality of this $MAE$ of 2.8.

An obvious alternative way of answering our initial question, would be to gather information on other predictions, the Internet is indeed full of them ([10, 11, 13]),

---

[3]Strictly convex for instance.

calculate $MAE$ for these guesses and compare. Unfortunately, this is indeed a formidable task. As a consequence, we have chosen a slightly different path. Instead of empirical comparisons, we can investigate some basic statistical properties[4] of MAE; for instance to establish minimal ($MAE_{MIN}$), maximal ($MAE_{MAX}$) as well as the expected value for MAE ($E[MAE]$) as a simpler (or at least less time consuming) way of testing the quality of Merson's predictions. It turns out that[5] $MAE_{MIN} = 0$, $MAE_{MAX} = \frac{n}{2}$ and $E[MAE] = \frac{1}{3} \cdot \frac{n^2-1}{n}$.

The above results provide interesting information. A random table permutation (or prediction) for PL ($n = 20$) can at best produce $MAE = 0$, while at worst, it can produce $MAE = \frac{n}{2} = \frac{20}{2} = 10$. On average, a random prediction should produce $E[MAE] = \frac{1}{3} \cdot \frac{n^2-1}{n} = \frac{1}{3} \cdot \frac{20^2-1}{20} = 6.65$.

The task of guessing randomly and hit the correct table is definitely a formidable one. There are $n!$ different tables to guess, and a random guess would hence (in the case of PL, $n = 20$) have a probability of $\frac{1}{20!} = \frac{1}{2432902008176640000} \approx 4 \cdot 10^{-19}$ of hitting the correct table. Consequently, Merson's table prediction is really impressive. On average, a $MAE$ of 6.65 compared to Merson's 2.8 indicate high quality in Merson's prediction. One could of course argue that Paul Merson is an expert, and one should expect him to know this business[6]. Still, information from other countries, for instance Norway, which we will focus more on in subsequent sections, indicate that even experts may have challenges in providing tips that fit final tables.

In the next section (Section 2), we investigate some scientific attempts to produce football table forecasts. In Section 3 we argue that the trend in present research seems to be oriented in a non-parsimonious fashion, and argue why this perhaps is not a good idea. In Section 4, we discuss alternative parsimonious modelling hypotheses and test one involving goal-difference as the prime explanatory factor. Section 5 concludes, discusses and suggests further research.

## 2. The science of football table prediction

Although the Internet is "full" of football table predictions, it would be an exaggeration to say that research literature is full of serious attempts to predict the same. Still, some noteworthy exceptions exist. Three relatively recent papers by Brillinger [2–4] seem to sum up the state of the art of the area. Brillinger's touch of difference compared to other previous work seems to be that he models game outcomes in the form of Win Tie or Loss – W, T, L – directly, as opposed to other authors who uses some distributional assumptions on goal scoring frequencey, typically as seen in [8] or in [5][7].

---

[4]Under an assumption of a random guess.

[5]Refer to Appendix A for the derivation of these results as well as some other relevant statistical properties of $MAE$.

[6]Some might even argue that PL is a league with low uncertainty of outcome (a competitively imbalanced league). Hence, it is not that hard to guess final tables. Chelsea, Arsenal and Manchester United have for instance a recurring tendency to end up among the 5 best.

[7]Meeden's paper does get some interesting and perhaps unexpected criticism in [6]. Here, the whole assumption of using probability theory to model goal scoring or match outcomes is questioned by game theoretic arguments.

Almost all of the work discussed above relies on simulation to produce actual forecasts. The idea is simple. Let the computer play the games; either by drawing goal scores or W, T, L (by estimated probailistic mechanisms) for all predefined matches in the league. Register match outcomes; either by counting goals or more directly by Brillingers approach. Then, when all match outcomes are defined, the league table can be set-up. Repeating the simulation produces a new final league table, and by a large number of simulation runs, expected table placement or probabilistic table predictions can be generated. Of course, such a method opens up for updating or reestimating underlying probabilties for team quality, which then are applied if one runs a rolling horizon approach. Such rolling horizon approaches seem to be quite popular in media – refer for instance to [14].

## 3. Parsimonious forecasting

The concept of parsimony (or parameter minimzation) is both well known and well studied in time series forecasting literature. Already Box and Jenkins [1] pointed out that parsimony is desireable if forecasting accuracy is the objective. An interesting emprical test of the actual consequences of parsimony versus non-parsimomny can be found in [7].

The reason why parsimony is desirable is obvious. A model where many parameters need to be estimated generate more aggregate uncertainty than a model with fewer parameters. As a consequence, the outcome – the forecasts – tends to be more uncertain and inaccurate.

Furthermore, in many cases where causal (regression type models) are used, either alone or in combination with time series models, the causal variables will often have to be predicted in order to obtain model estimates for the target variable. And these causal variables are typically just as hard, or (perhaps) even harder, to predict reasonably correctly, than the target variable. Suppose you want to predict the number of flats sold in a certain area in London this month next year. You know that many relevant economic variables such as the UK salary level, unemployment rate and interest rates (to name just a few) affect this target variable. If a prediction model contains these variables, you need to predict them in order to predict the number of flats sold in London next year. And predicting next year's UK unemployment, salary, and interest rates is (obviously) not an easy task.

If we return to our focus – football table prediction – it should seem quite obvious that the reported methodolgy discussed in Section 2 can hardly be described as parsimonious. On the contrary, probability estimates of many teams, maybe conditional on future events like injuries, or talent logistics shoud generate much added uncertainty and it should not come as a surprise that such methods produce quite bad predictions. The fact that the team in [14] missed Greece as a potential winner of EURO 2004 may serve as an adequate example.

This said, non-parsimonious models, causal or not, have other interesting properties. They can, for instance, (far better) answer questions of the 'what if type', which in some situations are more desirable than accurate forecasts.

So, what would be a parsimonious model for football table forecasting? The answer is simple and obvious, the table itself. Either last years table, if one

predicts the final table in-between seasons, or the latest table available if the aim is to predict the final table within a rolling horizon.

Obviously, even such a simple strategy brings challenges. In most leagues, there are relegation and promotion which has the obvious effect that last season's table contains a few teams other than this season's table. Furthermore, if the tables that are to be predicted are group tables in, say, European or World Championships, there is no previous season.

Still, such problems may be solved at least if we restrict our focus to prediction after some games or rounds have been played.

4. TESTING A HYPOTHESIS OF PARSIMONIOUS FOOTBALL TABLE PREDICTION

One simple way of testing the table in a certain round $r$'s predictive power on the final table is to perform a set of linear regressions, one for each round with the final table rank as the dependent variable and table rank in round $r$ as the independent variable. Or, in our notation (the obvious $r$-subscript is omitted for simplicity):

$$i = \beta_0 + \beta_1 P(i) + \epsilon_i.$$

By calculating $R^2$ in all these regressions[8], a discrete function $R^2(r)$ is obtained. Presumably, this function will have some kind of increasing pattern (not necessarily strict), but common sense indicates that football tables change less in later than early rounds. Figure 1 shows an example from last year's Tippeliga[9] in Norway.

Looking at Figure 1, we observe our predicted pattern of non-strict positive monotonicity in $R^2(r)$. However, we also observe something else: $R^2(r)$ reaches 80% explanatory power already in round 7. That is, 80% of the final table is there, already in round 7. Surely, $R^2(r)$ drops slightly in rounds 7 to 19, but this observation indicates that our parsimonious hypothesis actually may be of relevance.

Now, is the table rank the only possible parsimonious alternative? The answer is of course no. A table contains home wins, away wins, points, goal-score to name some potential additional information. Let us focus on goal score. At the start of a season, many teams may have new players, new managers, and we may suspect that the full potential of certain good teams may not be revealed in early table rankings. Vice versa, other teams are lucky, are riding a wave and take more points than expected. These not so good teams, may have a tendency to win even matches by a single goal, but also loose other matches (against very good teams) by many goals. As such, we could suspect that goal difference (at least in earlier rounds) could perhaps bring more and better predictive information than table ranking (or points for that matter). That is, we could hope to observe patterns (of course not as smooth) similar to the "fish-form" in Figure 2.

In order to check these two hypotheses more thoroughly, We examined the Norwegian top league for some previous seasons. We stopped in 2009, as the number of teams changed to 16 (from 14) in that year, and performed regressions like those

---

[8]It is unnecessary to carry out the full regressions, as we are only interested in $R^2$. So, we calculate simply $R^2$ as the square of the correlation coefficient $r$.

[9]The name Tippeligaen has been changed to Eliteserien for this (2017) season.

$R^2$ as a function of rounds in Tippeligaen 2016 with table position in each round
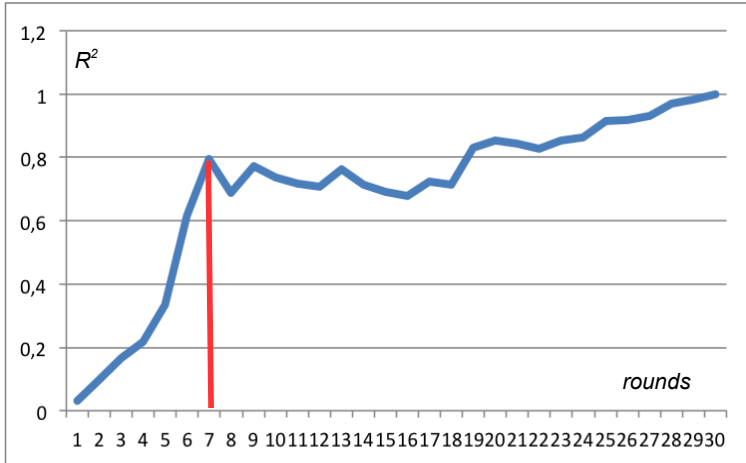as independent and final table position as dependent variables



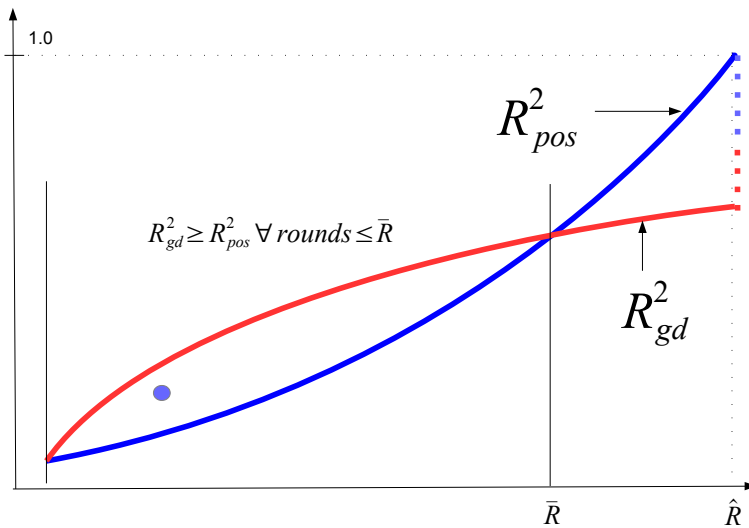**Figure 1.** An example of $R^2(r)$ from Tippeligaen 2016.



**Figure 2.** A hypothetical comparison of $R^2(r)$ for regressions with both goal difference $(R_{gd}^2)$ and table position $(R_{pos}^2)$ as independent variables.

described above with both table rank and goal difference as independent variables. Then 8 $R^2_{pos}(r)$ and 8 $R^2_{gd}(r)$ were generated. The result of this generation is shown in Figure 3.
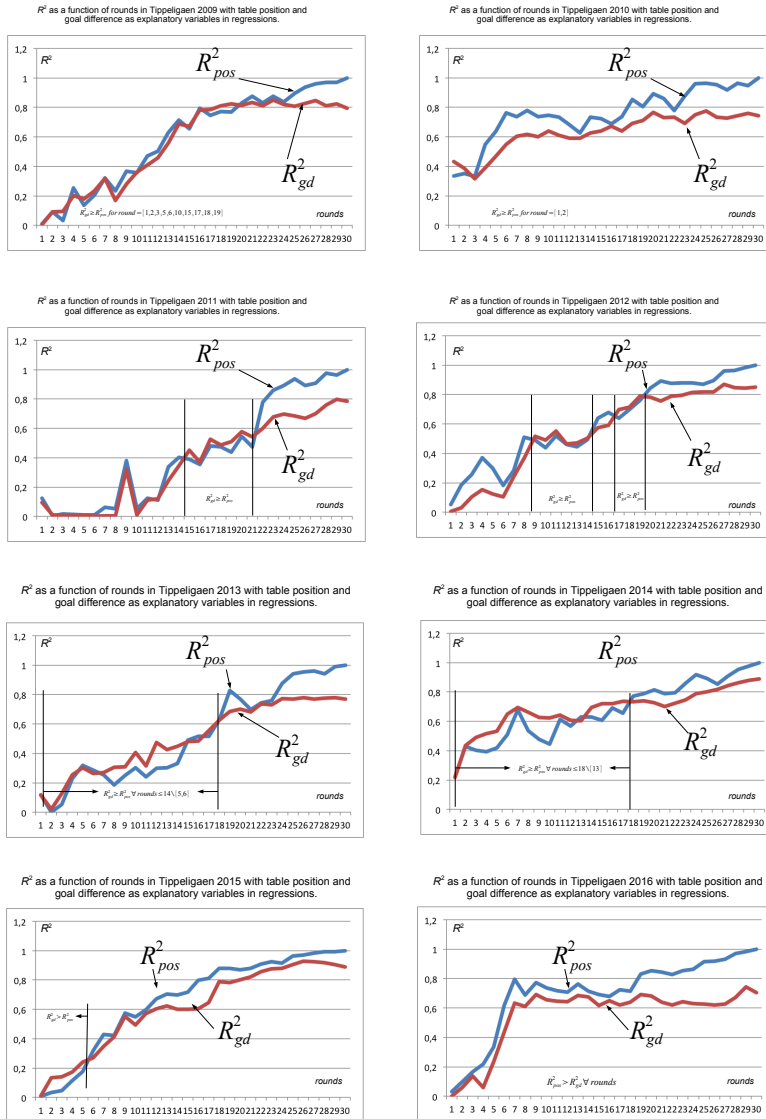


**Figure 3.** Full output from empirical analysis.

If we examine figure 3 we observe that 7 out of 8 seasons have patterns, although perhaps not visually too similar to Figure 2), where goal difference explains final table better than table rank – typically early in the season. Furthermore, around

80% explanatory power ($R^2 > 0.8$) is obtained roughly (for table rank) as early as mid-season for most cases.

Hence, an operative predictive strategy where a table prediction simply could be generated by sorting the goal differences in the early parts of the season and using the last observed table as the forecast later in the season seems reasonable.

Hence, we have demonstrated that our hypotheses are supported for these Norwegian data. Of course, our results do not say anything related to other leagues in other countries, but we do believe that similar patterns should be observable in most international leagues.

## 5. Conclusions and suggestions for further research

Apart from the fact that Paul Merson is a good predictor of the final PL-table (or at least he was before the 2016/2017 season), we have demonstrated that table rank or sometimes even better, goal difference explains major parts of final tables early. We have not (actually) checked (empirically) whether our prediction method is "better' than the existing methods in the literature. This is of course feasible, however, time consuming. What we can conclude without doubt, is that the methods applied by researchers in football table prediction are far more time and resource consuming than our methods. Simulation experiments take both coding and computing time and, if one is in doubt about whether one approach is better than the other, we could at least recommend trying our approach first.

So, is football table prediction important? Does it contribute to world welfare? Is it really necessary to spend time and resources even addressing this problem? Perhaps not. Still, most modern news agents spend a lot of time and (valuable) space on distributing such predictions each season. And, as a consequence, some real world demand seems to exist.

In any case, we have examined the problem from our perspective and even found some mathematical results we found interesting. Hopefully, our small effort may inspire other researchers to start the tedious job of empirically testing whether our approach performs better or worse than the simulation-based approaches.

## Appendix A. Statistical properties of $MAE$

### A.1. Maximal and minimal values for $MAE$

The minimal value of $MAE$ ($MAE_{MIN}$) is obvious. Even though it is unlikely (refer to Section 1), it is possible to guess correctly. In that case, $P(i) = i \forall i$, and by equations (1.1), $MAE = 0$.

Let us proceed by investigating the maximal $MAE$. We start by introducing $S(P)$ as:

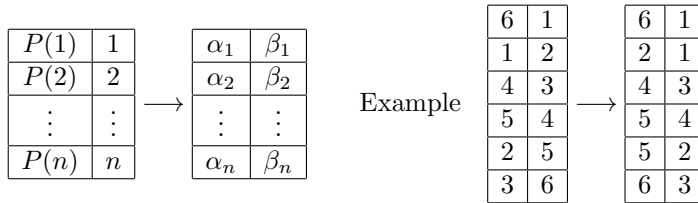$$S(P) = \sum_{i=1}^{n} |P(i) - i| \qquad (A.1)$$

then,

$$MAE(P) = \frac{1}{n} S(P). \qquad (A.2)$$

Our focus is on solving the optimization problem: $\max_P MAE(P)$. As equation (A.2) indicates only a multiplied constant difference, we might as well focus on solving:

$$\max_P S(P). \tag{A.3}$$

**Theorem A.1.** *The permutation $P = P_0 = [n, n-1, \ldots 2, 1]$ will be one (although not necessarily unique) solution to (A.3), and the value of the optimal objective (given an even numbered league[10] ) is $\frac{n}{2}$.*

*Proof.* Suppose an alternative permutation, say $P_1$ such that $S(P_1) > S(P_0)$ exists. We will show that such a permutation, $P_1$ cannot exist. In order to proceed, we introduce a little trick to handle the absolute value in equation (A.1). Let us illustrate it by an example.

| $P(1)$ | 1 |
|---|---|
| $P(2)$ | 2 |
| ⋮ | ⋮ |
| $P(n)$ | $n$ |

$\longrightarrow$

| $\alpha_1$ | $\beta_1$ |
|---|---|
| $\alpha_2$ | $\beta_2$ |
| ⋮ | ⋮ |
| $\alpha_n$ | $\beta_n$ |

Example

| 6 | 1 |
|---|---|
| 1 | 2 |
| 4 | 3 |
| 5 | 4 |
| 2 | 5 |
| 3 | 6 |

$\longrightarrow$

| 6 | 1 |
|---|---|
| 2 | 1 |
| 4 | 3 |
| 5 | 4 |
| 5 | 2 |
| 6 | 3 |

Above (on the left), a certain prediction (or permutation) $[P(1), P(2), \ldots, P(n)]$ along with the correct final table $[1, 2, \ldots n]$ is given. The trick involves a certain resort of this leftmost table into the table with $\alpha$'s and $\beta$'s, and is done as follows: $\alpha_i = \max\{P(i), i\}$ and $\beta_i = \min\{P(i), i\}$. The clue of the trick (the resort) is of course to achieve that $\alpha_i > \beta_i \forall i$, remove the absolute value sign, and hence obtain:

$$S(P) = \sum_{i=1}^{n} |P(i) - i| = \sum_{i=1}^{n} (\alpha_i - \beta_i) = \sum_{i=1}^{n} \alpha_i - \sum_{i=1}^{n} \beta_i.$$

Now, this reformulation leads to a simpler task in solving $\max_P S(P)$, as it can be done by maximising $\sum_{i=1}^{m} \alpha_i$ and minimising $\sum_{I=1}^{n} \beta_i$. The "best" that can be achieved is to get the column of $\alpha_i$'s to contain two copies of each of the numbers $m+1, \ldots, n$ such that the $\beta_i$-column contains two copies of the numbers $1, \ldots, m$ (where $m = \frac{n}{2}$). Permutations $P$ that make this happen are those where $P(i) \geq m+1$ for $i \leq m$ and $P(i) \leq m$ for $i \geq m+1$. That is, collectively, $P$ must "ship" the set $\{1, \ldots, m\}$ to $\{m+1, \ldots, n\}$ and vice versa[11]. Or, seen from the predictor's side, if the best half of the teams are predicted to end up on the bottom half of the table (and vice versa): Then, we get:

$$S^*(P) = 2\sum_{j=1}^{m}(m+j) - 2\sum_{j=1}^{m} j = 2m^2 = \frac{1}{2}n^2$$

---

[10]For obvious reasons, we focus on even numbered leagues. That is, $n = 2m$. The argument presented is, however, easily adapted to an odd $n$, in which case $P_0$ is still a solution to (A.3), but the objective value $\frac{n}{2} - \frac{1}{2n}$ is no longer an integer for $n > 1$

[11]As long as this is the case, it is immaterial how the permuation $P$ reorders each of the two sets $\{1, \ldots, m\}$ and $\{m+1, \ldots, 2m\}$, the objective value remains the same

or $MAE = \frac{1}{n}S^*(P) = \frac{n}{2}$. Luckily, the permutation $P_0$ satisfies this criteria. Reversing the order, predicting the winner to be last, number two to become second last and so on, will with necessity secure that the best $m$ are "shipped" to the last $m$ and vice versa.

$\square$

## A.2. Deriving the expression for $E[MAE]$

To derive the expression for $E[MAE]$, we assume uniform distribution of all possible tables (random guess). Let $x = (x_1, \ldots, x_n) = (P(1), \ldots, P(n))$ be a permutation. We write $S(x_1, \ldots, x_n)$ for the score and $p(x_1, \ldots, x_n) = 1/n!$ for the probability that $x$ occurs.

The expected value of $S$ is then (by definition):

$$\mathbb{E}[S] = \sum_{x \in \mathcal{S}_n} S(x_1, \ldots, x_n)p(x_1, \ldots, x_n) = \sum_{i=1}^{n} \sum_{x \in \mathcal{S}_n} |x_i - i|p(x_1, \ldots, x_n)$$

$$= \sum_{i=1}^{n} \sum_{x_i=1}^{n} |x_i - i| \sum_{y \in \mathcal{S}_{n-1}} p(y_1, \ldots, y_{i-1}, x_i, y_i, \ldots, y_{n-1}).$$

In the innermost (right) summation, we sum over all permutations$\{1, \ldots, n\}\backslash\{x_i\}$. That is, we leave out $x_i$ fom the numbers in all permutations, $\{1, \ldots, x_i - 1, x_i + 1, \ldots, n\}$. Then, the innermost summation contains all permuations of $(n-1)$ numbers. Hence, they are $(n-1)!$ in number. We get:

$$\mathbb{E}[S] = \sum_{i=1}^{n} \sum_{x_i=1}^{n} |x_i - i|\frac{(n-1)!}{n!} = \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{x_i=1}^{i-1} (i - x_i) + \sum_{x_i=i+1}^{n} (x_i - i) \right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} (f(i) + f(n + 1 - i)) = \frac{2}{n} \sum_{i=1}^{n} f(i), \quad \text{hvor } f(x) = \frac{1}{2}x(x - 1).$$

The last summation is straightforward to find and we get:

$$\mathbb{E}[S] = \frac{1}{n} \sum_{i=1}^{n} i(i - 1) = \frac{1}{3}(n^2 - 1)$$

and hence $E(MAE) = \frac{1}{3} \cdot \frac{n^2-1}{n}$. Similarly, we can also find the variance:

$$\text{Var}[S] = \frac{1}{45}(n + 1)(2n^2 + 7), \quad \text{Var}[MAE] = \frac{1}{45}\frac{(n + 1)(2n^2 + 7)}{n^2}.$$

## A.3. Some final mathematical remarks

If we return to the final paragraph of Subsection A.1, it should to be reasonably straightforward to realize that the mentioned criteria (for achieving maximal $MAE$) can be achieved in exactly $(m!)^2$ ways. (Take all $M!$ permutations of $[1, \ldots, m]$ and add $m$, for any such permutation, take all $m!$ permutations and subtract $m$.) Since there are $n! = (2m)!$ permutations totally, the probability of

obtaining the worst possible guess can be calculated as:

$$P\left(\text{Getting } MAE_{MAX} \text{ by guessing}\right) = \frac{(m!)^2}{(2m)!} = \left(\begin{array}{c} 2m \\ m \end{array}\right)^{-1} = \left(\begin{array}{c} n \\ n/2 \end{array}\right)^{-1}.$$

This outcome is unlikely. In the PL-case:

$$\left(\begin{array}{c} n \\ n/2 \end{array}\right)^{-1} = \left(\begin{array}{c} 20 \\ 10 \end{array}\right)^{-1} = \frac{1}{184756} \approx 5.4 \cdot 10^{-6}.$$

Still, far from as unlikely as guessing the correct table, as indicated in Section 1.

## References

[1] G. E. P. Box and G. M. Jenkins, *Time Series Analysis, Forecasting and Control (Revised Edition)*, Holden–Day, San Francisco, 1976.

[2] D. R. Brillinger, *Modelling some Norwegian soccer data*, in: Vijay Nair (ed.), *Advances in Statistical Modeling and Inference*, Series in Biostatistics **3**, World Scientific, Hackenshack, NJ, 2007, 3–20.

[3] D. R. Brillinger, *Modelling game outcomes of the Brazilian 2006 Series A Championship as ordinal-valued*, Brazilian Journal of Probability and Statistics **22** (2008), 80–104.

[4] D. R. Brillinger, *An analysis of Chinese Super League partial results*, Science in China, Series A **52** (2009), 1139–1156.

[5] E. Fernandez-Cantelli and G. Meeden, *An improved award system for soccer*, Chance **16** (2003), 23–29.

[6] K. K. Haugen, *An improved award system for soccer: A (game-theoretic) comment*, Chance **20** (2007), 22–24.

[7] J. Ledolter and B. Abraham, *Parsimony and its importance in time series forecasting*, Technometrics **23** (1981), 411–414.

[8] A. J. Lee, *Modeling scores in the Premier League: Is Manchester United really the best?*, Chance **10** (2009), 15–19.

[9] S. Makridakis, S. C. Wheelwright and V. E. McGee, *Forecasting, Methods and Applications*, 2nd. ed., John Wiley & Sons, 1983.

[10] J. Burn-Murdoch, *Latest predictions: Position probabilities*, Financial Times, May 2017, `https://www.ft.com/prem-predict`.

[11] B. McMahon, *Premier league predictions for 2016–2017*, Forbes/SportsMoney, August 2016, `www.forbes.com/sites/bobbymcmahon`.

[12] P. Merson, *Paul Merson predicts final Premier League table for 2016/17*, SKYSPORTS, August 2016, `www.skysports.com/football`.

[13] P. McNulty, *Premier league 2016-17: Who will finish where?*, BBC Sports, August 2016, `http://www.bbc.com/sport/football/36962933`.

[14] P. Tarrant, *Belgium to win Euro 2016? A Q&A on probabilistic predictions*, June 2016, `www.statslife.org.uk/sports`.

Kjetil K. Haugen, Faculty of Logistics, Molde University College, Specialized University in Logistics, Britveien 2, 6410 Molde, Norway
*e-mail*: `kjetil.haugen@himolde.no`

Brynjulf Owren, Department of Mathematical Sciences, NTNU, N-7034 Trondheim, Norway
*e-mail*: `brynjulf.owren@ntnu.no`