

# The Market for English Premier League Odds<sup>12</sup>

Guanhao Feng<sup>3</sup>

Nicholas Polson<sup>4</sup>

Jianeng Xu<sup>5</sup>

Booth School of Business

Booth School of Business

Department of Statistics

University of Chicago

University of Chicago

University of Chicago

Paper track: Business Sports

Paper ID: 537

## Abstract

We develop a probabilistic representation of real-time betting odds for English Premier League (EPL) soccer games and derive new empirical insights from the betting market efficiency. We use market odds data for different outcomes of win, loss, draw, and particularly the score difference. We show how a difference in Poisson processes (a.k.a. Skellam process) provides a dynamic probabilistic model for the evolution of score predictions implied by odds. As the game evolves, we use the updated market odds to re-estimate the Skellam process and to recalculate the prediction matrix of odds implied final scores. The Skellam interpretation enables us to understand the mechanism of the real-time betting market, and the evolution of odds implied volatility in the perspective of the Black-Scholes-Merton model. We verify the Skellam assumption through odds data of the win, loss, and draw for 1520 EPL games from 2012 to 2016, as well as odds data of the score difference for 18 games in the season 2016-2017. To demonstrate the flexibility of our model by showing how quickly the odds change as the score progresses, we apply it to an EPL game between Everton and West Ham in the season 2015-2016. Finally, we conclude with directions for future research.

---

<sup>1</sup> For our online latex generated version, please see <<https://arxiv.org/abs/1604.03614>>

<sup>2</sup> For a media coverage in Chicago Booth Review, please see  
<<http://review.chicagobooth.edu/economics/2016/article/real-time-gambling-odds-have-predictive-power>>

<sup>3</sup> Address: 5807 S Woodlawn Avenue, Chicago, IL 60637, USA. E-mail address: [guanhao.feng@chicagobooth.edu](mailto:guanhao.feng@chicagobooth.edu)

<sup>4</sup> Address: 5807SWoodlawnAvenue, Chicago, IL60637, USA. E-mail address: [nicholas.polson@chicagobooth.edu](mailto:nicholas.polson@chicagobooth.edu).

<sup>5</sup> Address: 5747 S Ellis Avenue, Chicago, IL 60637, USA. E-mail address: [jianeng@uchicago.edu](mailto:jianeng@uchicago.edu).

# 1. Introduction

## 1.1. The betting market for the EPL

Gambling on soccer is a global industry worth anywhere between \$700 billion and \$1 trillion a year (See “Football Betting - the Global Gambling Industry worth Billions.” BBC Sport). Spread betting, particularly fixed-odds betting, on the results of soccer matches is rapidly growing in popularity and odds are set via an online real-time betting market (Betfair, Bet365, etc.). It is any of various types of wagering on the outcome of an event where the pay-off is based on the accuracy of the wager, rather than a simple “win or lose” outcome. Traditional bookmakers, such as Ladbrokes, also offer odds on various outcomes of a match. For example, bets can be placed on the final results (win, lose, draw) as well as goals scored at half-time or full-time. Fractional odds are used in the UK, while money-line odds are favored by American bookmakers. Fractional odds of 2:1 (“two-to-one”) would imply that the bettor stands to make a \$200 profit on a \$100 stake.

A key feature of market odds is that they are updated in real time, and there is considerable interest in developing probability models for the evolution of the games score. Titman et al. (2015) assess the predictive capacity by using a model-based betting strategy to available live spread betting prices for EPL games. Stern (1994) and Polson and Stern (2015) both propose a Brownian motion model for the difference in teams scores and show how the market-based information can be used to calculate the implied volatility of a game. We build on this approach and develop a model that is tailor-made for the discrete evolution of the scores of an EPL game. Specifically, we develop a probability model based on the difference of Poisson processes (a.k.a. Skellam process, See Barndorff-Nielsen and Shephard (2012) for an introduction to Levy processes).

Another line of research, asks whether betting markets are efficient and, if not, how to exploit potential inefficiencies in the betting market. Levitt (2004) discusses the structural difference of the gambling market and financial markets. For examples, bookmakers are more skilled at game prediction than bettors and exploit bettor biases by purposely choosing prices that deviate from the market clearing price. Avery and Chevalier (1999) examine the hypothesis that sentimental bettors act like noise traders and can affect the path of prices in soccer betting markets.

In short, we provide a dynamic probabilistic representation to study the evolution of odds-implied final score prediction over the course of the game. In principle, our dynamic Skellam model can fit actual scoring data and improves by incorporating supplementary game information. However, in this paper, our primary focus is to provide a dynamic framework to interpret the real-time market odds and its prediction on final scores. The goal was to show how a two-parameter parsimonious model can be flexibly calibrated to the evolution of market odds. Relying on the betting market efficiency, we use the real-time market odds as a lens to understand the market prediction for game outcomes and provide the odds implied volatility in the perspective of the Black-Scholes-Merton model. Empirically, our Skellam model “fits” to this market of odds and provides a more efficient assessment of the outcome for the broad betting market, though any individual odds can be subject to liquidity issues.

## 1.2. Connections with Existing Work

Various probabilistic models have been proposed to predict the outcome of soccer matches motivated by the demand for assessing betting opportunities. Early models (Lee 1997) of the number of goals scored by each team use independent Poisson processes. Later models incorporate a correlation between the two scores and model the number of goals scored by each team using bivariate Poisson models (see Maher (1982) and Dixon and Coles (1997)). Our approach follows Stern (1994) by modeling the score difference (a.k.a. margin of victory), instead of modeling the number of goals and the correlation between scores directly.

Soccer scores by their very nature are discrete and not that frequent. Hence it is not adequate to apply a continuous-time stochastic model. Instead, we adapt the Poisson process and model the difference in scores via a Skellam distribution. In earlier studies, the Skellam distribution is used to model integer outcomes particularly for low-scoring sports, see Karlis and Ntzoufras (2003, 2009) and Koopman et al. (2014). Moreover, one advantage of our approach is that we can calibrate an implied volatility measure (Polson and Stern (2015)) from the market ex-ante assessment of odds for score difference.

Our approach is related to the literature of soccer gambling and market efficiency. For example, Vecer et al. (2009) estimate the scoring intensity in a soccer game using data from in-play betting markets. Dixon and Pope (2004) presents a detailed comparison of odds set by different bookmakers about the Poisson model predictions. Fitt (2009) applies the efficient portfolio theory to analyze the mispricing of cross-sectional odds. Online soccer spread bets requires bookmakers to alter the market odds dynamically to prevent arbitrage during a match and Fitt et al. (2005) models the value of online soccer spread by modeling goals and corners as Poisson processes.

We emphasize that, by incorporating the odds information, our model try to track and interpret the market participants' expectation of the game score, instead of predicting the final outcomes directly. The rest of the paper proceeds as follows. Section 2 presents our Skellam process model for tracking the difference in goals scored. We then show how to make use of an odds matrix while calibrating the model parameters. We calculate a dynamic implied prediction of any score and hence win, lose and draw outcomes, using real-time online market odds. Section 3 illustrates our methodology using an EPL game in 2015-2016 between Everton and West Ham. Finally, Section 4 discusses extensions of our basic model and concludes with directions for future research.

## 2. Skellam Process for EPL scores

Let the outcome between the two soccer teams A and B be modeled as a difference in scores,  $N(t) = N_A(t) - N_B(t)$ . We interpret  $N(t)$  as the lead-off home team A over the away team B.  $N_A(t)$  and  $N_B(t)$  denote the scores for both teams at time point  $t$  ( $0 \leq t \leq 1$ ). Negative values of  $N(t)$ , therefore, indicate that team A is behind. We assume that the game begins at time zero with  $N(0)$

= 0 and ends at time one with  $N(1)$  representing the final score difference, positive if A wins and negative if B wins.

For our analysis, we develop a probabilistic specification of the distribution of  $N(1)$  and, more generally, given  $N(t) = l$  where  $l$  is the current lead, as the game evolves. Given this probabilistic model, we can determine an implied prediction of the outcome of the whole match. For example, ex-ante  $P(N(1) > 0)$  will provide the odds of team A winning. The model can also be used for halftime betting, which is common in Europe. Half-time scores will be available for the distribution of  $N(1/2)$ , and as the game progresses we can calculate  $P(N(1/2) > 0 | N(t) = l)$  and  $P(N(1) > 0 | N(t) = l)$  where  $l$  is the current goal difference.

## 2.1. Implied Score Prediction from EPL Odds

The Skellam distribution models the difference between two independent Poisson variables, see Skellam (1946), Sellers (2012), Alzaid et al. (2010), and Barndorff-Nielsen and Shephard (2012). We now show how it can be used to model the point spread distribution in those sports with equal scored points. Karlis and Ntzoufras (2009) shows how Skellam distribution can be extended to a difference of distributions which have a specific trivariate latent variable structure.

We begin by specifying the score of home team A and away team B at time  $t$  by  $N_A(t)$  and  $N_B(t)$  respectively. Following Karlis and Ntzoufras (2003), we decompose the scores of each team as

$$\begin{cases} N_A(t) = W_A(t) + W(t) \\ N_B(t) = W_B(t) + W(t) \end{cases} \quad (1)$$

where  $W_A(t)$ ,  $W_B(t)$  and  $W(t)$  are independent processes with  $W_A(t) \sim \text{Poisson}(\lambda_{A,0}t)$ ,  $W_B(t) \sim \text{Poisson}(\lambda_{B,0}t)$ . Here  $W(t)$  is a non-negative integer-valued process to induce a correlation between the numbers of goals scored. If  $W(t)$  is a Poisson process, with  $W(t) \sim \text{Poisson}(\lambda_w t)$ , then  $N_A(t)$  and  $N_B(t)$  are both marginally Poisson with  $\text{Cov}(N_A(t), N_B(t)) = \lambda_w t$ . The advantage of analyzing the score difference,  $N(t)$ , is that it doesn't depend on the distribution of  $W(t)$ . The difference in goals scored is independent of  $W(t)$  and follows a Skellam distribution.

$$N(t) = N_A(t) - N_B(t) = W_A(t) - W_B(t) \sim \text{Skellam}(\lambda_{A,0}t, \lambda_{B,0}t) \quad (2)$$

Conditionally, we have

$$\begin{cases} W_A(1) - W_A(t) \sim \text{Poisson}(\lambda_{A,0}(1-t)) \\ W_B(1) - W_B(t) \sim \text{Poisson}(\lambda_{B,0}(1-t)) \end{cases} \quad (3)$$

Now we introduce  $N^*(1 - t)$ , the score difference of the sub-game which starts at time  $t$  and ends at time 1 and the duration is  $(1 - t)$ . By construction,  $N(1) = N(t) + N^*(1 - t)$  and by the property of Poisson process,  $N^*(1 - t)$  and  $N(t)$  are independent.

Therefore we can re-express equation (2) in terms of  $N^*(1 - t)$

$$N^*(t) = N_A^*(t) - N_B^*(t) = W_A^*(t) - W_B^*(t) \sim \text{Skellam}(\lambda_{A,t}, \lambda_{B,t}) \quad (4)$$

where  $W_A^*(1 - t) = W_A(1) - W_A(t)$  and  $W_B^*(1 - t) = W_B(1) - W_B(t)$ . The team strength is assumed to stay the same during the game, namely  $\lambda_{A,t} = \lambda_{A,0}(1 - t)$ . Later we will show how to allow for time-varying parameters. It should also be emphasized that  $\lambda_{A,t} = \lambda_{A,0}(1 - t)$  and  $\lambda_{B,t} = \lambda_{B,0}(1 - t)$  represent the expected "scoring rates" of  $W_A$  and  $W_B$  from time  $t$  to the end of the game respectively, and not of  $N_A$  and  $N_B$ . A natural interpretation of the parameters is that  $\lambda_{A,t}$  and  $\lambda_{B,t}$  reflect the "net" scoring ability of each team while the term  $W(t)$  model a common strength due to factors such as weather. We can use our model to calculate the probability of any particular score difference  $P(N(1) = x | \lambda_{A,0}, \lambda_{B,0})$  at the end of the game where the  $\lambda$ 's will be inferred from a matrix of market odds.

To derive the model implied winning probability, we use the law of total probability. The probability mass function of a Skellam random variable is the convolution of two Poisson distributions:

$$\begin{aligned} P(N(1) = x | \lambda_{A,0}, \lambda_{B,0}) &= \sum_{k=0}^{\infty} P(W_B(1) = k - x | W_A(1) = k, \lambda_{B,0}) P(W_A(1) = k | \lambda_{A,0}) \\ &= e^{-(\lambda_{A,0} + \lambda_{B,0})} \left( \frac{\lambda_{A,0}}{\lambda_{B,0}} \right)^{\frac{x}{2}} I_{|x|}(2\sqrt{\lambda_{A,0}\lambda_{B,0}}) \end{aligned} \quad (5)$$

where  $I_r(x)$  is the modified Bessel function of the first kind.

The probability of home team A winning can be easily calculated using the cumulative distribution function,

$$P(N(1) > 0 | \lambda_{A,0}, \lambda_{B,0}) = \sum_{x=1}^{\infty} P(N(1) = x | \lambda_{A,0}, \lambda_{B,0}) \quad (6)$$

In practice, we use an upper bound on the number of possible goals since the probability of an extreme score difference is always negligible. Unlike the Brownian motion model for the evolution of the outcome in a sports game (Stern (1994), Polson and Stern (2015)), the

probability of a draw in our setting is not zero. Instead,  $P(N(1) = 0 | \lambda_{A,0}, \lambda_{B,0}) > 0$  depends on the sum and product of two parameters  $\lambda_{A,0}$  and  $\lambda_{B,0}$  and thus the odds of a draw are non-zero.

For two evenly matched teams, using (5), the draw probability is a monotone decreasing function of  $\lambda$  (see Figure 1). Hence, two evenly matched teams with large  $\lambda$ 's are less likely to achieve a draw, compared with small  $\lambda$ 's.

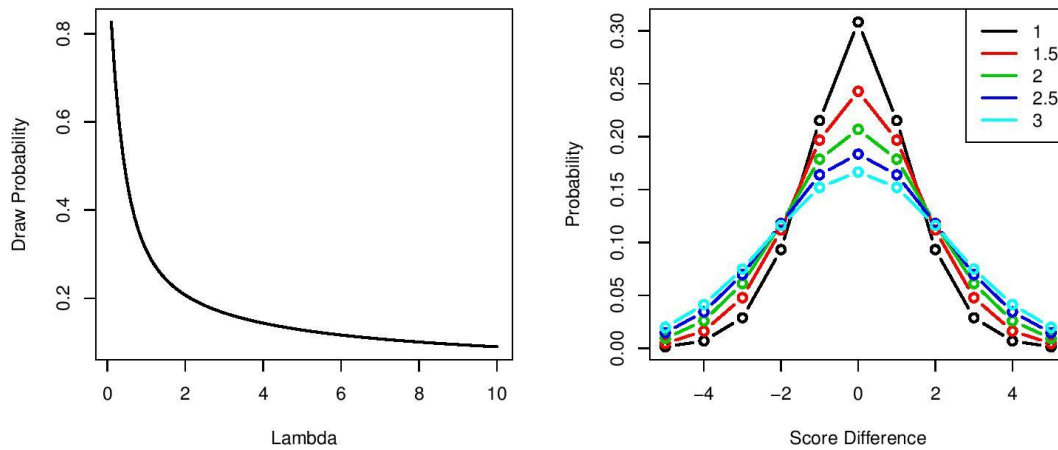


Figure 1 Left: Probability of a draw for two evenly matched teams. Right: Probability of score differences for two evenly matched teams. Lambda values are denoted by different colors.

We are mainly interested in the conditional probability of winning as the game progresses. Suppose that, the current lead at time  $t$  is  $l$  and so  $N(t) = l = N_A(t) - N_B(t)$ . With the property of Poisson process, the model updates the conditional distribution of the final score difference  $(N(1) | N(t) = l)$  by noting that  $N(1) = N(t) + N^*(1 - t)$  and that  $N(t)$  and  $N^*(1 - t)$  are independent. Specifically, conditioning on  $N(t) = l$ , we have the identity

$$N(1) = N(t) + N^*(1 - t) = l + \text{Skellam}(\lambda_{A,t}, \lambda_{B,t}).$$

From the above expression, we are now in a position to find the conditional distribution  $(N(1) = x | N(t) = l)$  for every time point  $t$  of the game given the current score. Simply put, we have the time homogeneous condition

$$P(N(1) = x | \lambda_{A,t}, \lambda_{B,t}, N(t) = l) = P(N^*(1 - t) = x - l | \lambda_{A,t}, \lambda_{B,t}) \quad (7)$$

where  $\lambda_{A,t}, \lambda_{B,t}, l$  are either given by market expectations or are known at time  $t$ .

There are two conditional probabilities of interest. First, the chances that the home team A wins are



$$P(N(1) > 0 \mid \lambda_{A,t}, \lambda_{B,t}, N(t) = l) = P(\text{Skellam}(\lambda_{A,t}, \lambda_{B,t}) > -l \mid \lambda_{A,t}, \lambda_{B,t}) \quad (8)$$

Second, the conditional probability of a draw at time  $t$  is

$$P(N(1) = 0 \mid \lambda_{A,t}, \lambda_{B,t}, N(t) = l) = P(\text{Skellam}(\lambda_{A,t}, \lambda_{B,t}) = -l \mid \lambda_{A,t}, \lambda_{B,t}) \quad (9)$$

We now turn to the calibration of our model from given market odds.

## 2.2. Market Calibration

Our information set at time  $t$ , denoted by  $I_t$ , includes the current lead  $N(t) = l$  and the market odds for  $\{\text{Win, Lose, Draw, Score}\}_t$ , where  $\text{Score}_t = \{(i - j) : i, j = 0, 1, 2, \dots\}$ . These market odds can be used to calibrate a Skellam distribution which has only two parameters  $\lambda_{A,t}$  and  $\lambda_{B,t}$ . The best fitting Skellam model with parameters  $\{\lambda_{A,t}, \lambda_{B,t}\}$  will then provide a better summary of the market's information concerning the outcome of the game than any individual market (such as win odds) as they are subject to liquidity and thinly traded effects. Suppose that the fractional odds for all possible final score outcomes are given by a bookmaker. For example, assume that we observe the final score ending with 2-1 has odds of 3/1. In this case, the bookmaker pays out three times the amount staked by the bettor if the outcome is indeed 2-1. The market implied probability makes the expected winning amount of a bet equal to 0. In this case, the implied probability  $p = 1/(1 + 3) = 1/4$  and the expected winning amount is  $\mu = -1*(1 - 1/4) + 3*(1/4) = 0$ . We denote this odds as  $\text{odds}(2, 1) = 3$  and convert all the available odds to implied probabilities, using the identity

$$P(N_A(1) = i, N_B(1) = j) = \frac{1}{1 + \text{odds}(i, j)}$$

Finally we get an odds matrix  $O$ . Its element  $o_{ij} = \text{odds}(i - 1, j - 1)$ ,  $i, j = 1, 2, 3, \dots$  for all possible combinations of final scores.

Extreme outcomes and the corresponding odds are not offered by the bookmakers. Since the probabilities are tiny, we set them equal to 0. However, the sum of the possible probabilities is still larger than 1. This phenomenon is standard in betting markets (see Dixon and Coles (1997) and Polson and Stern (2015)). The "excess" probability corresponds to a quantity known as the "market vig." For example, if the sum of all the implied probabilities is 1.1, then the expected profit of the bookmaker is 10%. To account for this phenomenon, we scale down the probabilities to make sure that the resulting sum equals to 1 before estimation. It's admitted that the scaling down step is only valid when the bookmaker gives the same market vig to all possible results. The consistency of our calibrated parameters depends on this assumption to some extent.

To determine the parameters  $\lambda_{A,t}$  and  $\lambda_{B,t}$  for the sub-game  $N^*(1 - t)$ , the odds from a bookmaker should be adjusted by  $N_A(t)$  and  $N_B(t)$ . For example, if  $N_A(0.5) = 1$ ,  $N_B(0.5) = 0$  and  $\text{odds}(2, 1) = 3$



at half time, these observations actually says that the odds for the second half score being 1-1 is 3 (the outcomes for the whole game and the first half are 2-1 and 1-0 respectively, thus the outcome for the second half is 1-1). The adjusted odds\* for  $N^*(1 - t)$  is calculated using the original odds as well as the current scores and given by

$$odds^*(i, j) = odds(i + N_A(t), j + N_B(t)) \quad (10)$$

At time  $t$  ( $0 \leq t \leq 1$ ), we calculate the implied conditional probabilities of score differences using odds information

$$P(N(1) = k | N(t) = l) \propto \sum_{i-j=k-l} \frac{1}{1 + odds^*(i, j)} \quad (11)$$

Moments of Poisson distribution make it easy to derive the moments of a Skellam random variable with parameters  $\lambda_{A,0}$  and  $\lambda_{B,0}$ . The conditional moments are given by

$$\begin{cases} E[N(1) | N(t) = l] = l + (\lambda_{A,t} - \lambda_{B,t}) \\ V[N(1) | N(t) = l] = \lambda_{A,t} + \lambda_{B,t} \end{cases} \quad (12)$$

It's not necessarily ensured that  $E^*[N(1) | N(t) = l] - l \leq V^*[N(1) | N(t) = l]$ . A simple method of moments estimate of  $\lambda$ 's, i.e. the solution to the equations

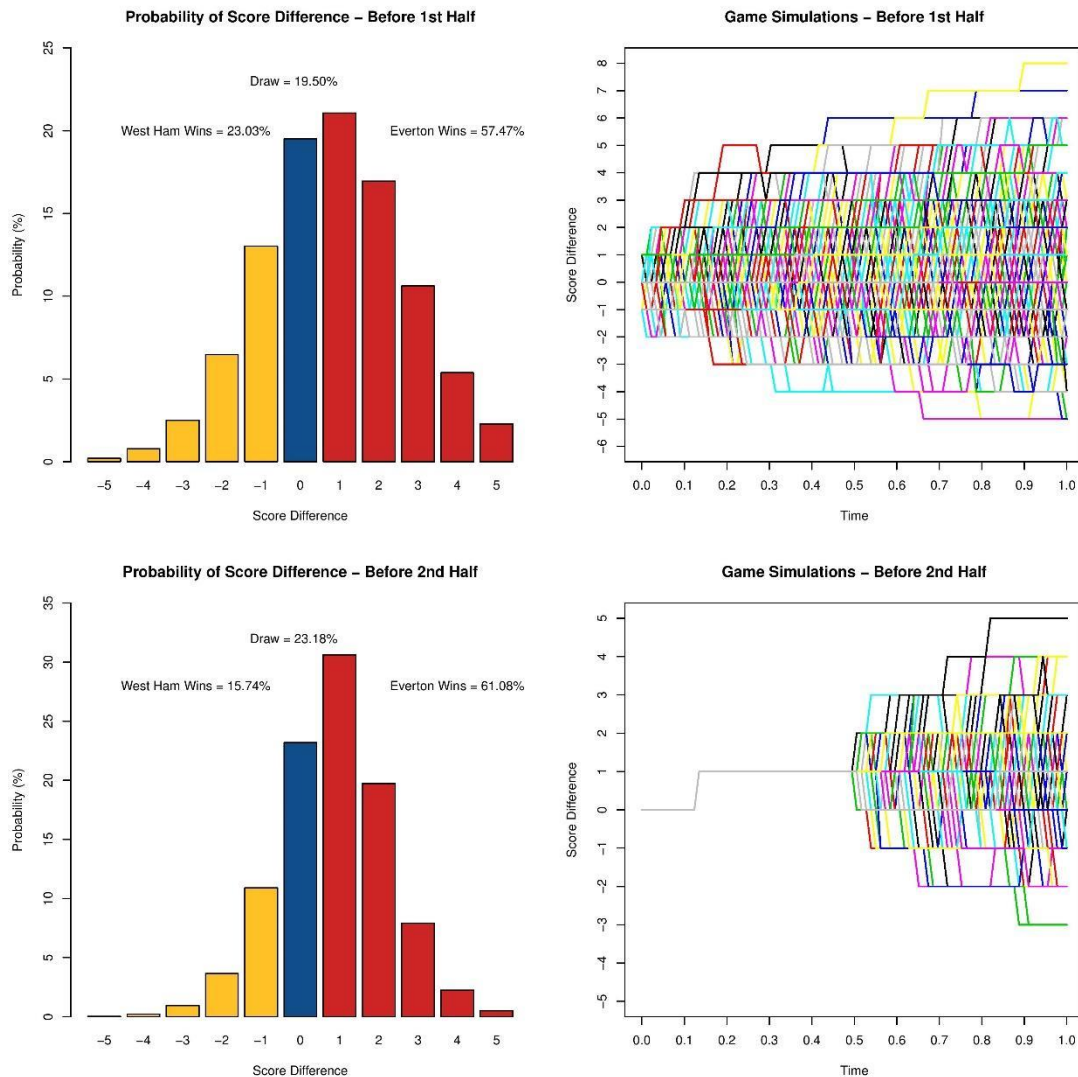
$$\begin{cases} \hat{E}[N(1) | N(t) = l] = l + (\lambda_{A,t} - \lambda_{B,t}) \\ \hat{V}[N(1) | N(t) = l] = \lambda_{A,t} + \lambda_{B,t} \end{cases} \quad (13)$$

where  $E^*$  and  $V^*$  are the expectation and variance calculated using market implied conditional probabilities, could be negative. To address this issue, we define the residuals

$$\begin{cases} D_E = E[N(1) | N(t) = l] - [l + (\lambda_{A,t} - \lambda_{B,t})] \\ D_V = V[N(1) | N(t) = l] - [\lambda_{A,t} + \lambda_{B,t}] \end{cases} \quad (14)$$

We then calibrate parameters by adding the constraints  $\lambda_{A,t} \geq 0$  and  $\lambda_{B,t} \geq 0$  and solving the following equivalent constrained optimization problem.





$$(\hat{\lambda}_{A,t}, \hat{\lambda}_{B,t}) = \arg \min \{D_E^2 + D_V^2\} \quad \text{subject to } \lambda_{A,t} \geq 0, \lambda_{B,t} \geq 0 \quad (15)$$

Figure 2 The Skellam process model for winning margin and game simulations. The top left panel shows the outcome distribution using odds data before the match starts. Each bar represents the probability of a distinct final score difference, with its color corresponding to the result of win/lose/draw. Score differences larger than 5 or smaller than -5 are not shown. The top right panel shows a set of simulated Skellam process paths for the game outcome. The bottom row has the two figures updated using odds data available at half-time.

Figure 2 illustrates a simulation evolution of an EPL game between Everton and West Ham (March 5th, 2016) with their estimated parameters. It provides a discretized version of Figure 1 in Polson and Stern (2015). The outcome probability of first half and updated second half are given in the left two panels. The top right panel illustrates a simulation-based approach to visualizing how the model works in the dynamic evolution of score difference. In the bottom left

panel, from half-time onwards, we also simulate a set of possible Monte Carlo paths to the end of the game. This illustrates the discrete nature of our Skellam process and how the scores evolve.

## 2.3. Model Diagnostics

To assess the performance our score-difference Skellam model calibration for the market odds, we have collected data from ladbrokes.com on the correct score odds of 18 EPL games (from October 15th to October 22nd, 2016) and plot the calibration result in Figure 3. The Q-Q plot of  $\log(\text{odds})$  is also shown. In average, there are 13 different outcomes per game, i.e.,  $N(1) = -6, -5, \dots, 0, \dots, 5, 6$ . In total 238 different outcomes are used. We compare our Skellam implied probabilities with the market implied probabilities for every outcome of the 18 games. If the model calibration is sufficient, all the data points should lie on the diagonal line. Figure 3 left panel demonstrates that our Skellam model is calibrated by the market odds sufficiently well, except for the underestimated draw probabilities. Karlis and Ntzoufras (2009) describe this underestimation phenomenon in a Poisson-based model for the number of goals scored. Following their approach, we apply a zero-inflated version of Skellam distribution to improve the fit on draw probabilities, namely

$$\begin{cases} \tilde{P}(N(1) = 0) = p + (1 - p) * P(N(1) = 0) \\ \tilde{P}(N(1) = x) = (1 - p) * P(N(1) = 0) \text{ if } x \neq 0 \end{cases} \quad (16)$$

Here  $0 < p < 1$  is an inflation factor and  $P^*$  denotes the inflated probabilities. We also consider another type of inflation here

$$\begin{cases} \tilde{P}(N(1) = 0) = (1 + \theta) * P(N(1) = 0) \\ \tilde{P}(N(1) = x) = (1 - \gamma) * P(N(1) = 0) \text{ if } x \neq 0 \end{cases} \quad (16)$$

where  $\theta$  is the inflation factor and  $P(N(1) = 0) = \gamma / (\gamma + \theta)$ .

Both types of inflation factors have the corresponding interpretation regarding the bookmakers' way of setting odds. With the first type of factor, the bookmakers generate two different set of probabilities, one specifically for the draw probability (namely the inflation factor  $p$ ) and the other for all the outcomes using the Skellam model. The "market vig" for all the outcomes is a constant. With the second type, the bookmakers use the Skellam model to generate the probabilities for all the outcomes. Then they apply a larger "market vig" for draws than others. Yates (1982) also point out the "collapsing" tendency in forecasting behavior, whereby the bookmakers are inclined to report forecasts of 50% when they feel they know little about the event. In Figure 3 right panel, we see that the Skellam implied  $\log(\text{odds})$  has a heavier right tail than the market implied  $\log(\text{odds})$ . This effect results from the overestimation of extreme outcomes, which in turn is due to market microstructure effect due to the market "vig".

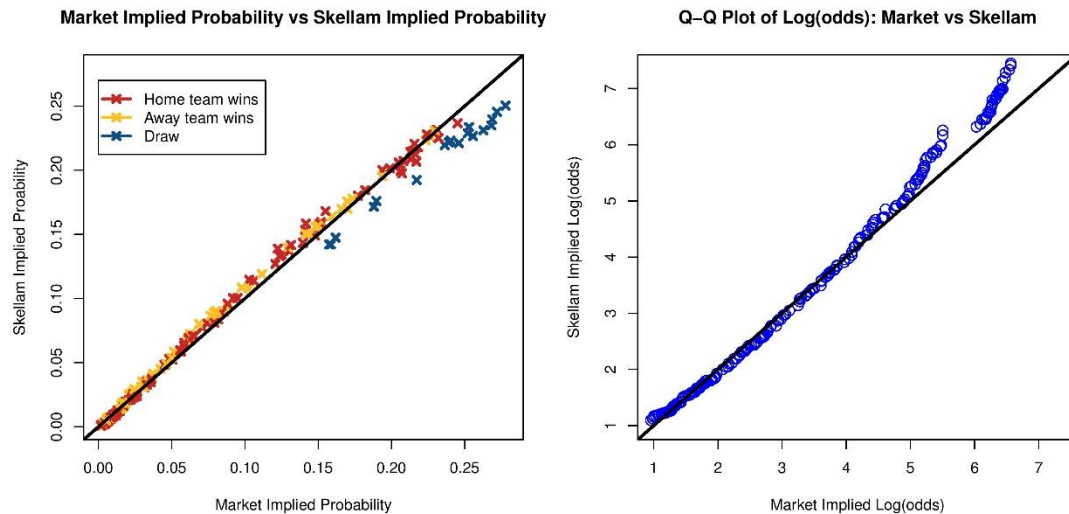


Figure 3 Left: Market implied probabilities for the score differences versus Skellam implied probabilities. Every data point represents a particular score difference; Right: Market log(odds) quantiles versus Skellam implied log(odds) quantiles. Market odds (from ladbrokes.com) of 18 games in EPL 2016-2017 are used (in average 13 score differences per game). The total number of outcomes is 238.

To assess the out-of-sample predictive ability of the Skellam model, we analyze the market (win, lose, draw) odds for 1520 EPL games (from 2012 to 2016, 380 games per season). However, the sample covariance of the end of game scores,  $N_A(1)$  and  $N_B(1)$ , is close to 0. If we assume parameters stay the same, then the estimates are  $\hat{\lambda}_{A,0} = 1.5$  and  $\hat{\lambda}_{B,0} = 1.2$ . Since the probabilities of win, lose and draw sum to 1, we only plot the market implied probabilities of win and draw. In Figure 4 left panel, the draw probability is nearly a non-linear function of the win probability. To illustrate our model, we set the value of  $\lambda_{A,0}\lambda_{B,0} = 1.5 \times 1.2 = 1.8$  and plot the curve of Skellam implied probabilities (red line). We further provide the inflated Skellam probabilities (blue line for the first type and green line for the second type). As expected, the non-inflated Skellam model (red line) underestimates the draw probabilities while the second type inflated Skellam model (green line) produces the better fit. We also group games by the market implied winning probability of home teams  $P(N(1) > 0)$ :  $(0.05, 0.1]$ ,  $(0.1, 0.15]$ ,  $\dots$ ,  $(0.8, 0.85]$ . We calculate the frequency of home team winning for each group. In Figure 4 right panel, the barplot of frequencies (x-axis is regarding scaled odds) shows that the market is efficient, i.e., the frequency is close to the corresponding market implied probability and our Skellam model is calibrated to the market outcome for this dataset.



### Skellam Model with Inflated Zero and Fix Parameter Product

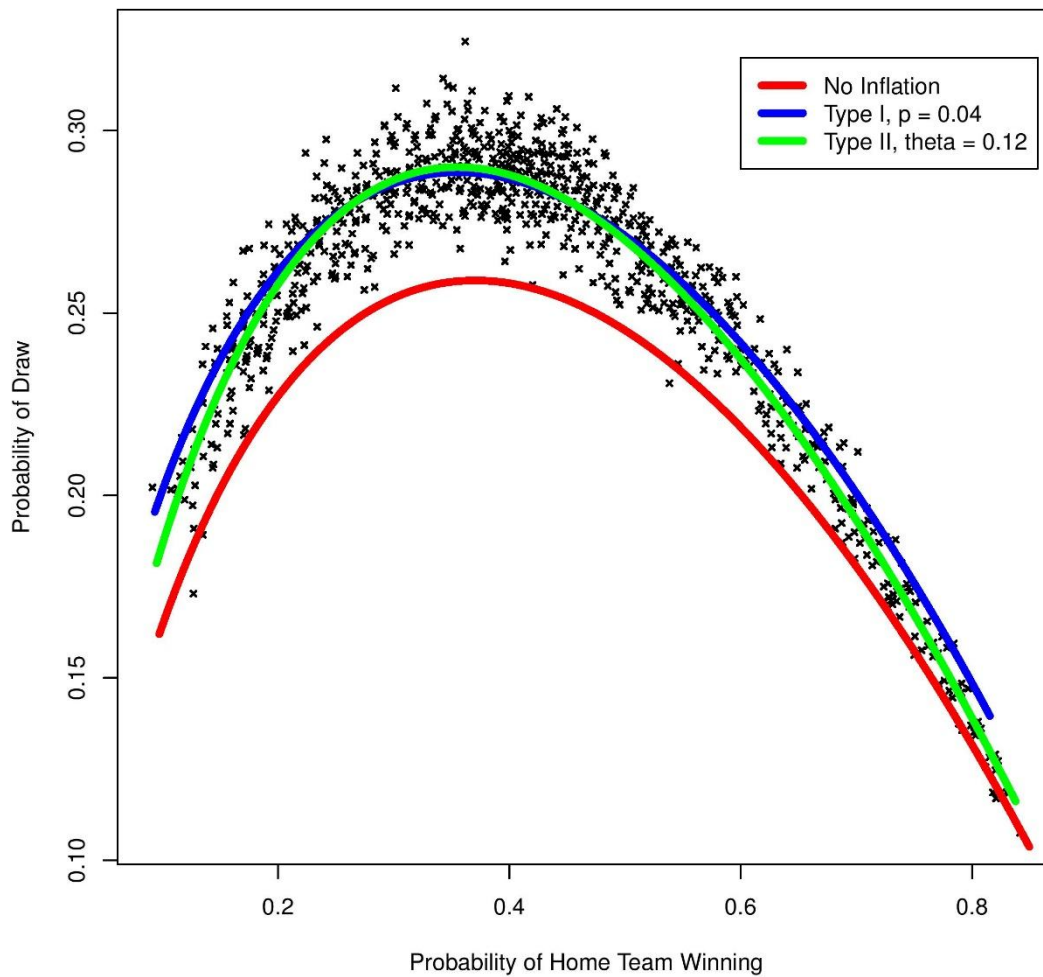


Figure 4 Left: Market implied probabilities of win and draw. The fitted curves are Skellam implied probabilities with fixed product of lambdas. Right: Market odds and result frequency of home team winning. 1520 EPL games from 2012 to 2016 are used. The dashed line represents: Frequency = Market Implied Probability.

## 2.4. Time-Varying Extension

One extension that is clearly warranted is allowing for time-varying  $\{\lambda_{A,t}, \lambda_{B,t}\}$  where the Skellam model is re-calibrated dynamically through updated market odds during the game. We use the current  $\{\lambda_{A,t}, \lambda_{B,t}\}$  to project possible results of the match in our Skellam model. Here  $\{\lambda_{A,t}, \lambda_{B,t}\}$  reveal the market expectation of scoring difference for both teams from time  $t$  to the end of the game as the game progresses. Similar to the martingale approach of Polson and Stern (2015),  $\{\lambda_{A,t}, \lambda_{B,t}\}$  reveal the best prediction of the game result. From another point of view, this approach is the same as assuming homogeneous rates for the rest of the game.

An alternative approach to time-varying  $\{\lambda_{A,t}, \lambda_{B,t}\}$  is to use a Skellam regression with conditioning information such as possession percentages, shots (on goal), corner kicks, yellow cards, red cards, etc. We would expect jumps in the  $\{\lambda_{A,t}, \lambda_{B,t}\}$  during the game when some important events happen. A typical structure takes the form

$$\begin{cases} \log(\lambda_{A,t}) = \alpha_A + \beta_A X_{A,t-1} \\ \log(\lambda_{B,t}) = \alpha_B + \beta_B X_{B,t-1} \end{cases} \quad (17)$$

estimated using standard log-linear regression.

Our approach relies on the betting market being efficient so that the updating odds should contain all information of game statistics. Using log differences as the dependent variable is another alternative with a state space evolution. Koopman et al. (2014) adopt stochastically time-varying densities in modeling the Skellam process. Barndorff-Nielsen et al. (2012) is another example of the Skellam process with different integer valued extensions in the context of high-frequency financial data. Further analysis is required, and this produces a promising area for future research.

### 3. Example: Everton vs West Ham (3/5/2016)

We collect the real-time online betting odds data from [ladbrokes.com](http://ladbrokes.com) for an EPL game between Everton and West Ham on March 5th, 2016. By collecting real-time online betting data for every 10-minute interval, we can show the evolution of betting market prediction on the final result. We do not account for the overtime for both 1st half and 2nd half of the match and focus on a 90-minute game.

#### 3.1. Implied Skellam Probabilities

Table 1 shows the raw data of odds right the game. We need to transform odds data into probabilities. For example, for the outcome 0-0, 11/1 is equivalent to a probability of 1/12. Then we can calculate the marginal probability of every score difference from -4 to 5. We neglect those extreme scores with small probabilities and rescale the sum of event probabilities to one.

| Home\Away | 0    | 1    | 2    | 3    | 4     | 5     |
|-----------|------|------|------|------|-------|-------|
| 0         | 11/1 | 12/1 | 28/1 | 66/1 | 200/1 | 450/1 |
| 1         | 13/2 | 6/1  | 14/1 | 40/1 | 100/1 | 350/1 |



|   |       |       |       |       |       |       |
|---|-------|-------|-------|-------|-------|-------|
| 2 | 7/1   | 7/1   | 14/1  | 40/1  | 125/1 | 225/1 |
| 3 | 11/1  | 11/1  | 20/1  | 50/1  | 125/1 | 275/1 |
| 4 | 22/1  | 22/1  | 40/1  | 100/1 | 250/1 | 500/1 |
| 5 | 50/1  | 50/1  | 90/1  | 150/1 | 400/1 |       |
| 6 | 100/1 | 100/1 | 200/1 | 250/1 |       |       |
| 7 | 250/1 | 275/1 | 375/1 |       |       |       |
| 8 | 325/1 | 475/1 |       |       |       |       |

Table 1 Original odds data from Ladbrokes before the game started

In Figure 5, the probabilities estimated by the model are compared with the market implied probabilities. As we see, during the course of the game, the Skellam assumption suffices to approximate market expectation of score difference distribution. This set of plots is evidence of goodness-of-fit the Skellam model.

Table 2 shows the model implied probability for the outcome of score differences before the game, compared with the market implied probability. As we see, the Skellam model appears to have longer tails. Different from independent Poisson modeling in Dixon and Coles (1997), our model is more flexible with the correlation between two teams. However, the trade-off of flexibility is that we only know the probability of score difference instead of the exact scores.

| Score difference | -4   | -3   | -2   | -1    | 0     | 1     | 2     | 3     | 4    | 5    |
|------------------|------|------|------|-------|-------|-------|-------|-------|------|------|
| Market Prob. (%) | 1.70 | 2.03 | 4.88 | 12.33 | 21.93 | 22.06 | 16.58 | 9.82  | 4.72 | 2.23 |
| Skellam Prob.(%) | 0.78 | 2.50 | 6.47 | 13.02 | 19.50 | 21.08 | 16.96 | 10.61 | 5.37 | 2.27 |

Table 2 Market implied probabilities for the score differences versus Skellam implied probabilities at different time points.



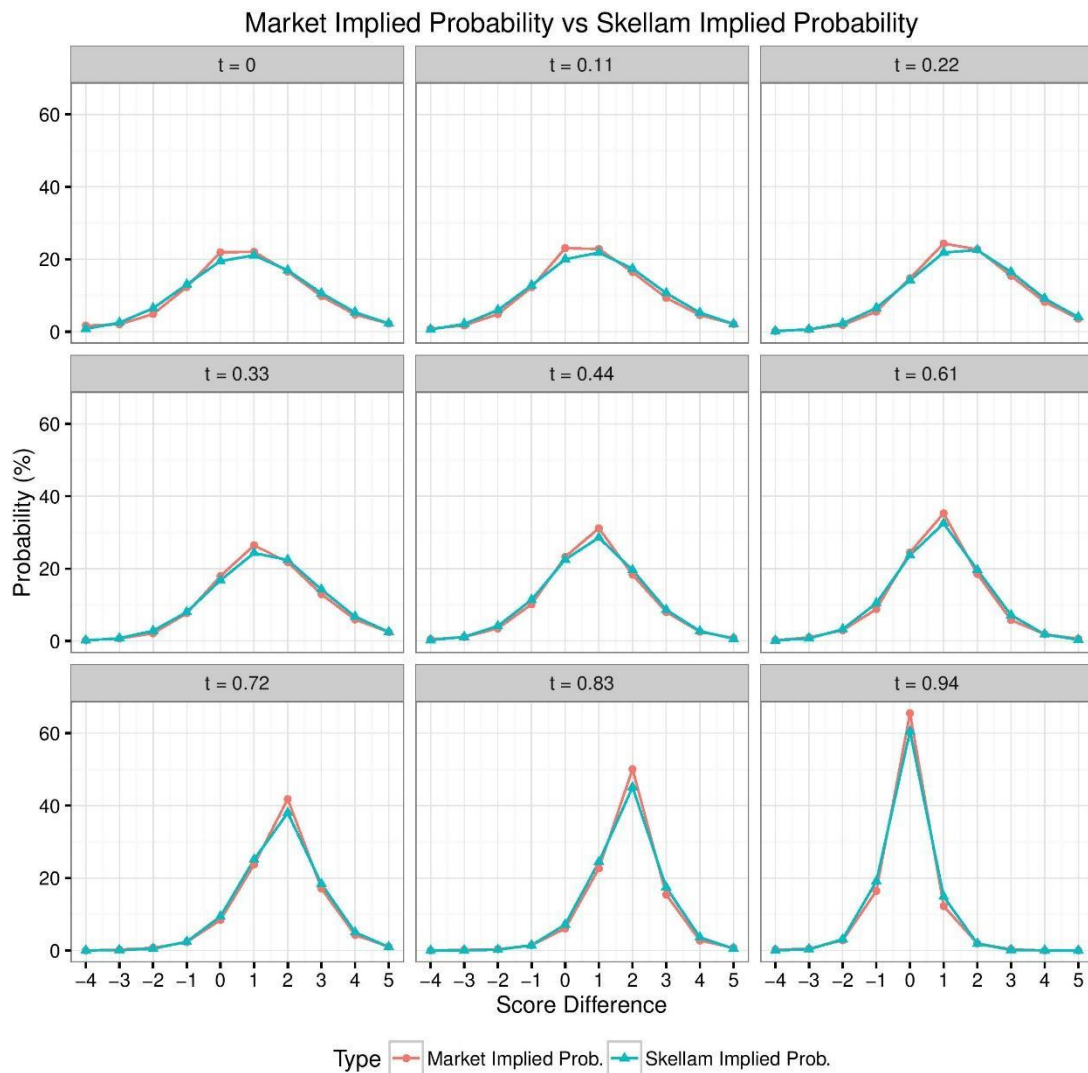


Figure 5 Market implied probabilities versus the probabilities estimated by the model at different time points, using the parameters given in Table 3.

Finally, we can plot these probability paths in Figure 6 to examine the behavior of the two teams and track the market predictions on the final result. Notably, we see the probability change of win/draw/loss for important events during the game: goals scoring and a red card penalty. In such a dramatic game, the winning probability of Everton gets raised to 90% before the first goal of West Ham in 78th minutes. The first two goals scored by West Ham in the space of 3 minutes completely reverses the probability of winning. The probability of draw gets raised to 90% until we see the last-gasp goal of West Ham that decides the game.

### 3.2. How the Market Forecast Adapts

A natural question arises to how does the market odds (win, lose, draw and actual score) adjust as the game evolves. This is similar to option pricing where Black-Scholes model uses its implied

volatility to show how market participants' beliefs change. Our Skellam model mimics its way and shows how the market forecast adapts to changing situations during the game. See Merton (1976) for references of jump models.

Our work builds on Polson and Stern (2015) who define the implied volatility of a NFL game. For an EPL game, we simply define the implied volatility as  $\sigma_{IV,t} = \sqrt{\lambda_{A,t} + \lambda_{B,t}}$ . As the market provides real-time information about  $\lambda_{A,t}$  and  $\lambda_{B,t}$ , we can dynamically estimate  $\sigma_{IV,t}$  as the game proceeds. Any goal scored is a discrete Poisson shock to the expected score difference (Skellam process) between the teams, and our odds implied volatility measure will be updated.

Figure 6 plots the path of implied volatility throughout the course of the game. Instead of a downward sloping line, we see changes in the implied volatility as critical moments occur in the game. The implied volatility path provides a visualization of the conditional variation of the market prediction for the score difference. For example, when Everton lost a player by a red card penalty at 34th minute, our estimates  $\hat{\lambda}_{A,t}$  and  $\hat{\lambda}_{B,t}$  change accordingly. There is a jump in implied volatility and our model captures the market expectation adjustment about the game prediction. The change in  $\hat{\lambda}_A$  and  $\hat{\lambda}_B$  are consistent with the findings of Vecer et al. (2009) where the scoring intensity of the penalized team drops while the scoring intensity of the opposing team increases. When a goal is scored in the 13th minute, we see the increase of  $\hat{\lambda}_{B,t}$  and the market expects that the underdog team is pressing to come back into the game, an effect that has been well-documented in the literature. Another important effect that we observe at the end of the game is that as goals are scored (in the 78th and 81st minutes), the markets expectation is that the implied volatility increases again as one might expect.

Figure 7 compares the updating implied volatility of the game with implied volatilities of fixed  $(\lambda_{A,0} + \lambda_{B,0})$ . At the beginning of the game, the red line (updating implied volatility) is under the " $(\lambda_{A,0} + \lambda_{B,0} = 4)$ "-blue line; while at the end of the game, it's above the " $(\lambda_{A,0} + \lambda_{B,0} = 8)$ "-blue line. As we expect, the value of  $(\hat{\lambda}_{A,t} + \hat{\lambda}_{B,t})/(1 - t)$  in Table 3 increases throughout the game, implying that the game became more and more intense and the market continuously updates its belief in the odds.

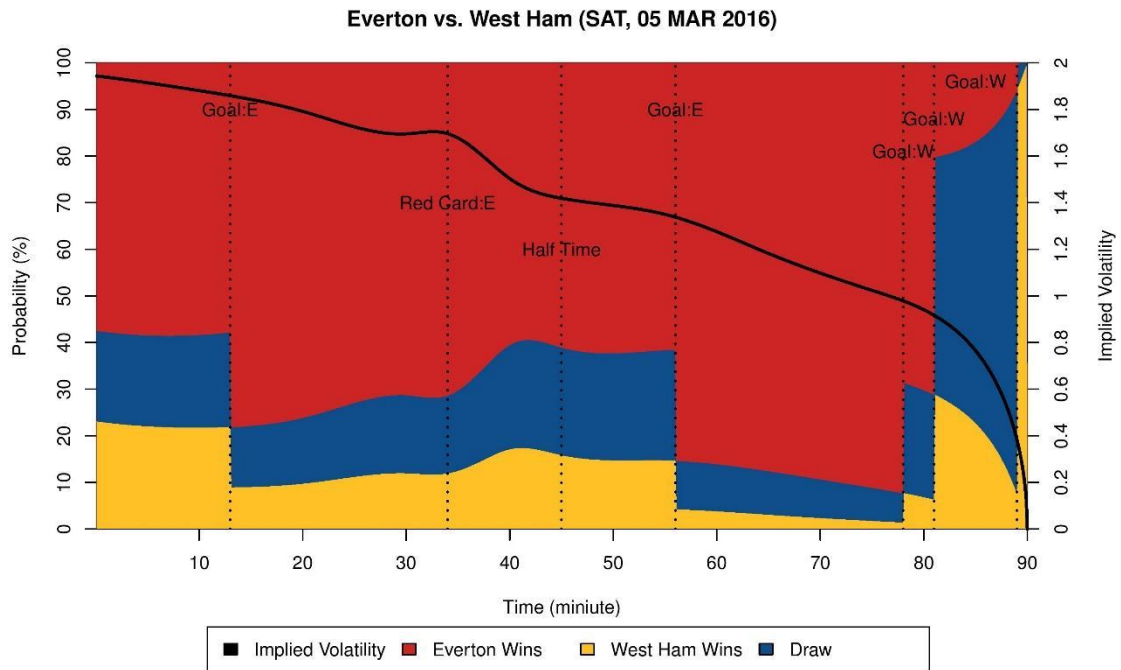


Figure 6 The betting market data for Everton and West Ham is from ladbrokes.com. Market implied probabilities (expressed as percentages) for three different results (Everton wins, West Ham wins and draw) are marked by three distinct colors, which vary dynamically as the game proceeds. The solid black line shows the evolution of the implied volatility (defined in Section 3.2). The dashed line shows significant events in the game, such as goals and red cards. Five goals in this game are 13' Everton, 56' Everton, 78' West Ham, 81' West Ham and 90' West Ham.

| $t$   | 0    | 0.11 | 0.22 | 0.33 | 0.44 | 0.50 | 0.61 | 0.72 | 0.83 | 0.94  | 1 |
|---|------|------|------|------|------|------|------|------|------|-------|---|
| $\hat{\lambda}_{A,t}/(1-t)$                         | 2.33 | 2.51 | 2.53 | 2.46 | 1.89 | 1.85 | 2.12 | 2.12 | 2.61 | 4.61  | 0 |
| $\hat{\lambda}_{B,t}/(1-t)$                         | 1.44 | 1.47 | 1.59 | 1.85 | 2.17 | 2.17 | 2.56 | 2.90 | 3.67 | 5.92  | 0 |
| $(\hat{\lambda}_{A,t} + \hat{\lambda}_{B,t})/(1-t)$ | 3.78 | 3.98 | 4.12 | 4.31 | 4.06 | 4.02 | 4.68 | 5.03 | 6.28 | 10.52 | 0 |
| $\sigma_{IV,t}$                                     | 1.94 | 1.88 | 1.79 | 1.70 | 1.50 | 1.42 | .135 | 1.18 | 1.02 | 0.76  | 0 |

Table 3 The calibrated  $\{\hat{\lambda}_{A,t}, \hat{\lambda}_{B,t}\}$  divided by  $(1 - t)$  and the implied volatility during the game.  $\{\lambda_{A,t}, \lambda_{B,t}\}$  are expected goals scored for rest of the game. The less the remaining time, the less likely to score goals. Thus  $\{\hat{\lambda}_{A,t}, \hat{\lambda}_{B,t}\}$  decrease as  $t$  increases to 1. Dividing them by  $(1 - t)$  produces an updated version of  $\lambda^0$ 's for the whole game, which are in general time-varying (but not decreasing necessarily). The calibrated  $\{\hat{\lambda}_{A,t}, \hat{\lambda}_{B,t}\}$  divided by  $(1 - t)$  and the implied volatility during the game.

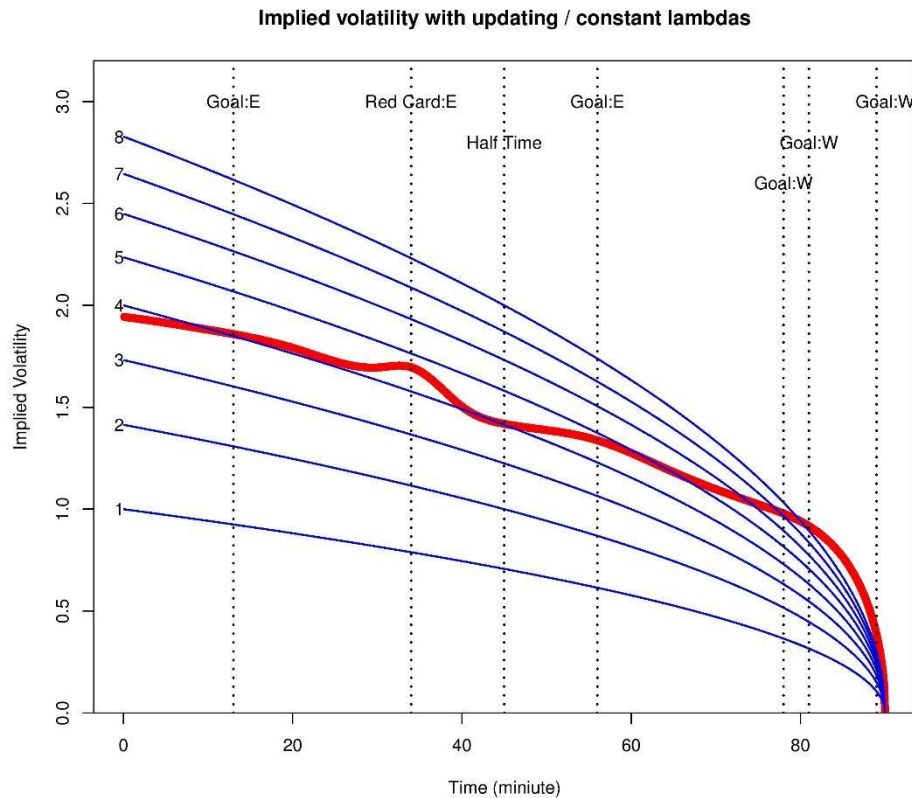


Figure 7 Red line: the path of implied volatility throughout the game. Blue lines: the path of implied volatility with constant  $\lambda_{A,0} + \lambda_{B,0}$ . Here  $(\lambda_{A,0} + \lambda_{B,0}) = 1, 2, \dots, 8$ .

## 4. Discussion

The goal of our analysis is to provide a probabilistic methodology for calibrating real-time market odds for the evolution of the score difference for a soccer game. Rather than directly using game information, we use the current odds market to calibrate a Skellam model to provide a forecast of the final result. To our knowledge, our study is the first to offer an interpretation of the betting market and to show how it reveals the market expectation of the game result through an implied volatility. One area of future research is studying the index betting. For example, a soccer game includes total goals scored in match and margin of superiority (see Jackson (1994)). The latter is the score difference in our model, and so the Skellam process directly applies.

Our Skellam model is also valid for low-scoring sports such as baseball, hockey or American football with a discrete series of scoring events. For NFL score prediction, Baker and McHale (2013) propose a point process model that performs as well as the betting market. On the one hand, our model has the advantage of implicitly considering the correlation between goals scored by both teams but on the other hand, ignores the sum of goals scored. For high-scoring sports, such as basketball, the Brownian motion adopted by Stern (1994) is more applicable.

Rosenfeld (2012) provides an extension of the model that addresses concerns of non-normality and uses a logistic distribution to estimate the relative contribution of the lead and the remaining advantage. Another avenue for future research, is to extend the Skellam model to allow for the dependent jumpiness of scores which is somewhere in between these two extremes (see Glickman and Stern (1998), Polson and Stern (2015) and Rosenfeld (2012) for further examples.)

Our model allows the researcher to test the inefficiency of EPL sports betting from a statistical arbitrage viewpoint. More importantly, we provide a probabilistic approach for calibrating dynamic market-based information. Camerer (1989) shows that the market odds are not well-calibrated and that an extreme underdog during a long losing streak is under-priced by the market. Golec and Tamarkin (1991) test the NFL and college betting markets and find bets on underdogs or home teams win more often than bets on favorites or visiting teams. Gray and Gray (1997) examine the in-sample and out-of-sample performance of different NFL betting strategies by the probit model. They find the strategy of betting on home team underdogs averages returns of over 4 percent, over commissions. In summary, a Skellam process appears to fit the dynamics of EPL soccer betting very well and produces a natural lens to view these market efficiency questions.

## Reference

- [1] Alzaid, A. A., M. A. Omaid, et al. (2010). On the poisson difference distribution inference and applications. *Bulletin of the Malaysian Mathematical Sciences Society* 8(33), 17–45.
- [2] Avery, C. and J. Chevalier (1999, October). Identifying Investor Sentiment from Price Paths: The Case of Football Betting. *The Journal of Business* 72(4), 493–521.
- [3] Baker, R. D. and I. G. McHale (2013). Forecasting exact scores in national football league games. *International Journal of Forecasting* 29(1), 122–130.
- [4] Barndorff-Nielsen, O. E., D. G. Pollard, and N. Shephard (2012). Integer-valued Levy processes and low latency financial econometrics. *Quantitative Finance* 12(4, SI), 587–605.
- [5] Barndorff-Nielsen, O. E. and N. Shephard (2012). Basics of levy processes. Technical report, Economics Group, Nuffield College, University of Oxford.
- [6] Camerer, C. F. (1989). Does the basketball market believe in the 'hot hand'? *American Economic Review* 79(1), 76–76.
- [7] Dixon, M. J. and S. G. Coles (1997, January). Modelling Association Football Scores and Inefficiencies in the Football Betting Market. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 46(2), 265–280.
- [8] Dixon, M. J. and P. F. Pope (2004, October). The Value of Statistical Forecasts in the UK Association Football Betting Market. *International Journal of Forecasting* 20(4), 697–711.
- [9] Fitt, A. D. (2009, April). Markowitz Portfolio Theory for Soccer Spread Betting. *IMA Journal of Management Mathematics* 20(2), 167–184.
- [10] Fitt, A. D., C. J. Howls, and M. Kabelka (2005, September). Valuation of Soccer Spread Bets. *Journal of the Operational Research Society* 57(8), 975–985.
- [11] Glickman, M. E. and H. S. Stern (1998, February). A State-Space Model for National Football League Scores. *Journal of the American Statistical Association* 93(441), 25–35.
- [12] Golec, J. and M. Tamarkin (1991, December). The Degree of Inefficiency in the Football Betting Market: Statistical Tests. *Journal of Financial Economics* 30(2), 311–323.
- [13] Gray, P. K. and S. F. Gray (1997, September). Testing Market Efficiency: Evidence From The NFL Sports Betting Market. *The Journal of Finance* 52(4), 1725–1737.
- [14] Jackson, D. A. (1994). Index Betting on Sports. *The Statistician* 43(2), 309.



- [15] Karlis, D. and I. Ntzoufras (2003, October). Analysis of Sports Data by Using Bivariate Poisson Models. *Journal of the Royal Statistical Society: Series D (The Statistician)* 52(3), 381–393.
- [16] Karlis, D. and I. Ntzoufras (2009, April). Bayesian Modelling of Football Outcomes: Using the Skellam's Distribution for the Goal Difference. *IMA Journal of Management Mathematics* 20(2), 133–145.
- [17] Koopman, S. J., R. Lit, and A. Lucas (2014). The dynamic skellam model with applications. Tinbergen Institute Discussion Paper 14-032/IV/DSF73.
- [18] Lee, A. J. (1997, September). Modeling Scores in the Premier League: Is Manchester United Really the Best? *Chance* 10(1), 15–19.
- [19] Levitt, S. D. (2004). Why are gambling markets organized so differently from financial markets? *Economic Journal* 114(3), 223–246.
- [20] Maher, M. J. (1982, September). Modelling Association Football Scores. *Statistica Neerlandica* 36(3), 109–118.
- [21] Merton, R. C. (1976, January). Option Pricing when Underlying Stock Returns are Discontinuous. *Journal of financial economics* 3(1-2), 125–144.
- [22] Polson, N. G. and H. S. Stern (2015). The Implied Volatility of a Sports Game. *Journal of Quantitative Analysis in Sports* 11(2), 145–153.
- [23] Rosenfeld, J. W. (2012). An in-game win probability model of the NBA. Thesis, Harvard University.
- [24] Sellers, K. F. (2012). A Distribution Describing Differences in Count Data Containing Common Dispersion Levels. *Advances and Applications in Statistical Sciences* 7(3), 35–46.
- [25] Skellam, J. G. (1946, January). The Frequency Distribution of the Difference Between Two Poisson Variates Belonging to Different Populations. *Journal of the Royal Statistical Society* 109(3).
- [26] Stern, H. S. (1994). A Brownian Motion Model for the Progress of Sports Scores. *Journal of the American Statistical Association* 89(427), 1128–1134.
- [27] Titman, A., D. Costain, P. Ridall, and K. Gregory (2015). Joint modelling of goals and bookings in association football. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 178(3), 659–683.
- [28] Vecer, J., F. Kopriva, T. Ichiba, et al. (2009). Estimating the effect of the red card in soccer: When to commit an offense in exchange for preventing a goal opportunity. *Journal of Quantitative Analysis in Sports* 5(1), 1–20.

- [29] Yates, J. F. (1982). External correspondence: Decompositions of the mean probability score.  
*Organizational Behavior and Human Performance* 30(1), 132–156.