

Real-time forecasting within soccer matches through a Bayesian lens

Chinmay Divekar^a, Soudeep Deb^a, Rishideep Roy^b

^a*Indian Institute of Management Bangalore, Bannerghatta Main Road, Bangalore, 560076, KA, India*

^b*University of Essex, School of Mathematics, Statistics and Actuarial Science, Wivenhoe Park, Colchester, CO4 3SQ, Essex, UK*

Abstract

This paper employs a Bayesian methodology to predict the results of soccer matches in real-time. Using sequential data of various events throughout the match, we utilize a multinomial probit regression in a novel framework to estimate the time-varying impact of covariates and to forecast the outcome. English Premier League data from eight seasons are used to evaluate the efficacy of our method. Different evaluation metrics establish that the proposed model outperforms potential competitors inspired by existing statistical or machine learning algorithms. Additionally, we apply robustness checks to demonstrate the model's accuracy across various scenarios.

Keywords: In-game forecasting, Ordered multinomial probit model, Soccer prediction, Bayesian method

1. Introduction

Sports and games are recreations that have attracted mankind since time immemorial. Soccer is a significant one among them. Early forms of sports involving feet and balls have been noted in many parts of the world. The game *Tsu'chu* or *Cuju* (in China), which literally stands for 'kicking the ball', is one of the earliest versions of the game from East Asia and was in practice a few millennia ago. There was a similar sport in Japan called *Kemari*, which still survives. Similarly, there have been examples of sports using balls and feet from ancient Greek and Roman civilizations, among the Native Americans, the indigenous people of Oceania, etc. There have been similar other sports recorded in Europe, particularly in the modern-day United Kingdom during the Middle Ages. With time, some of the old forms have modified themselves, some have lost out in the race, and others have taken their place. Their popularity, in terms of participation and attracting the audience, keeps the sheen of soccer. As [Walvin \(2014\)](#) notes in his book on the history of soccer, it is literally the "people's game". Spectators are especially engaged in rooting for their favourite heroes and teams and in predicting the outcome of the games. In fact, forecasting and decision-making have become intrinsic parts of sports. For avid fans, betting on their favourite teams provides just as good an experience as watching the game. The global sports betting market is estimated to increase at an approximate compound annual growth rate of 10.3% till 2032, with soccer (or association football) bringing the most attention, as stated in [European-Gaming \(2022\)](#). Arguably the most viewed sport in the world, the sports market connected to soccer is also huge.

With the increasing viewership and fandom, the concept of within-game result forecasting is turning out to be a very important aspect of soccer. The market associated with this prediction being ever so on the rise, there is a greater requirement on the precision of these predictions, which include not just the result but also the final scoreline. There are numerous techniques to forecast the outcome of a soccer match based on aggregated data at the beginning of a match, but they lack the flexibility to update the predictions based on the sequence of particular events during the match.

The fact that different types of events may occur on the soccer field with every passing minute within a match naturally motivates a Bayesian model, where one can estimate the effects of specific types of events as functions of time by taking into account the minute-by-minute data of soccer matches and subsequently use that to predict the final outcome of the game. To the best of our knowledge, there is no published work in this direction. Although many websites, such as <https://theanalyst.com/eu/> (for football), <https://fivethirtyeight.com/> (for basketball), <https://www.dimers.com/> and <https://patents.google.com/patent/US20070072679A1/en> (usually for all sports), provide live-win-probability for different matches, they lack the clarity and documentation behind their methods. In the current work, we aim to bridge that gap by developing a predictive Bayesian model that can be used for within-match forecasting in soccer. Our methodology works with the final outcome as an ordinal response variable and models the latent variable by suitably incorporating different covariates and events relevant to the game. On the one hand, the proposed approach allows us to identify the time-varying impacts of various events in soccer, while on the other, it is found to record superior predictive performance than a few other possible approaches, which are developed as extensions of existing statistical and machine learning methods.

Before delving into the proposed methodology, we find it pivotal to present a succinct account of existing literature that focuses on forecasting problems related to soccer using formal statistical methods, as well as the ones that develop in-game prediction methods for other types of sports.

1.1. Brief literature review

Accurate forecasting of results in sports has been an interest of the world for a long time. Quite naturally, due to its popularity and financial impacts, prediction problems in soccer constitute an extensively researched area. Early papers like [Maher \(1982\)](#) modelled the distribution of goals using bivariate Poisson random variables. [Rue and Salvesen \(2000\)](#) extended the bivariate modelling approach to a Bayesian setup while considering the strength of a team as a time-dependent factor. [Crowder et al. \(2002\)](#) proposed a more computationally efficient method to update the strength parameters. [Goddard \(2005\)](#) compared the forecasting performance of a goal-based model (bivariate Poisson regression model) against a result-based model (ordered probit model). [McHale and Scarf \(2007\)](#) extended the idea of dependence in a bivariate framework and utilized copulas to model the results of soccer matches. On the other hand, [Hvattum and Arntzen \(2010\)](#) worked with a rating-based system to predict the result of soccer matches, demonstrating that ELO rating-based measures can adequately incorporate past results. Another advanced technique of sparse bivariate Poisson model, along with the concepts of boosting to select appropriate covariates, was used by [Groll et al. \(2018\)](#).

Another branch of literature studies the relationship between the match results and bookmaker odds. [Forrest and Simmons \(2000\)](#) analyzed English soccer match odds provided by newspaper agencies to determine their efficiency. The authors concluded that these odds do not appropriately utilize team strength and such publicly available data to enhance their forecasting prowess. The articles by [Štrumbelj and Šikonja \(2010\)](#), [Štrumbelj \(2014\)](#) explored the quality of bookmaker odds by viewing them as probabilistic evaluations of the match results. Different aspects of the game that seemed to have an effect on the outcome were also studied. [Clarke and Norman \(1995\)](#) explored the effect of home advantage on the game result, while [Ley et al. \(2019\)](#) compared the existing team strength-based modelling approaches. [Koopman and Lit \(2019\)](#) proposed a dynamic modelling framework for the scenario when the outcome variable is goals scored, win/draw/loss and difference in goals scored. The time-dependent covariates in this method are updated as an autoregressive process. Different machine learning techniques have also been used in this regard. [Liti et al. \(2017\)](#) provide an excellent comparative discussion on the accuracy of various such

algorithms, while [Mendes-Neves and Mendes-Moreira \(2020\)](#) assess the efficacy of neural network ensemble methods in forecasting match outcomes in soccer.

Moving on to the literature of in-game forecasting, we note that since the 2000s, with the influx of data, this topic has picked pace in the domain of sports analytics. In cricket, [Easton and Uylangco \(2006\)](#) demonstrated the impact of in-game ball-by-ball events on the change in the odds of winning. In tennis, [Klaassen and Magnus \(2003\)](#) modelled the probability of the player serving to win the match and designed a way to update this probability after every point. This model was further explored by [Easton and Uylangco \(2010\)](#) to demonstrate tennis betting markets’ in-game efficiency. [Stern \(2005\)](#) forecast the winning probability in a game of basketball by modelling the difference in points scored as a Brownian process. An extension of this by incorporating the betting odds and using a Gamma process was proposed by [Song et al. \(2020\)](#). Intriguingly, even though soccer is arguably the most popular sport across the globe, we could not find any research article to tackle the problem of within-game forecasting in soccer. Our focus in this article is to address this issue through the help of a Bayesian framework.

1.2. Our contribution

In sports like basketball, baseball or cricket, the main events deciding the result (points, runs, wickets etc.) typically have a high rate of occurrence; hence a much gradual change is expected to be observed in the outcome probabilities over time. In contrast, the frequency of the corresponding variable (that is, a goal) in professional soccer is extremely low. Consequently, the fate of a soccer match is immensely volatile, which is likely to hinder the accuracy of an in-game forecaster. This article attempts to solve this problem and fill the gap in the extant literature, by proposing a Bayesian model for within-game forecasting in soccer.

The proposed methodology treats the outcome of a match with respect to the home team as an ordinal multinomial random variable (win/draw/loss). Then, utilizing appropriate time-invariant and time-dependent covariates which record the events happening in real-time in a soccer match, our model forecasts the result as the match progresses. We use a complete Bayesian framework for predicting the final outcome of the match. A latent variable with two cutoffs is used as a proxy for modelling the final outcome. We assume that this latent variable depends on different covariates pertaining to the match, as well as random errors, and this dependence is linear. The latent variable is continuously updated using Bayesian techniques, which in turn predicts the final outcome as well. As we shall discuss in detail below, it is quite a broad framework without making strong assumptions about the model. We use real-life data from the top division of English club football to demonstrate the predictive accuracy of the proposed algorithm. The model is also able to quantify the effects of specific types of events as a function of time during the progression of the match. Furthermore, different robustness checks are done to establish that the method works well across various scenarios.

The outline of the rest of the article is as follows. The methodology used for analysis is described in Section 2. We begin with our model specification in detail, followed by the Bayesian techniques used in our work. We also describe a way of evaluating our method, which includes the estimations as well as the predictive distributions. For completeness of the study, we extend a few other existing methods for conventional forecasting in soccer to within-game predictive techniques and compare their performances against our proposed model. Section 3 describes the dataset we use for our work, along with some exploratory analysis and descriptive statistics. Our results are demonstrated in Section 4 where we present a general analysis, followed by a robustness study. Next, we look at two specific case studies to understand the results better. The article ends with the summary and some important concluding remarks in Section 5.

2. Methodology

2.1. Model specification

Throughout this article, we shall use \mathbb{N} , \mathbb{Z} and \mathbb{R} to denote the sets of natural numbers, integers and real numbers, respectively. We use the shorthand notation $[K]$ for the set of all natural numbers from 1 to K , i.e. $\{k : 1 \leq k \leq K; k \in \mathbb{N}\}$. Any matrix is represented in bold capital letters, and any vector in lower-case bold letters; for example, \mathbf{a} would indicate a vector and \mathbf{A} would represent a matrix. We shall use \mathbf{I}_d to denote an identity matrix of dimension $d \times d$. Also, we shall use $\mathbf{1}(\cdot)$ to denote an indicator function. A d -dimensional normal distribution is indicated by $\mathcal{N}_d(\cdot, \cdot)$ whereas a truncated normal distribution is denoted by the shorthand $\mathcal{TN}(\mu, \sigma^2, a, b)$, where μ and σ are the mean and standard deviation of the original normal distribution, while the truncation limits are a and b , with $a < b$.

Let us formally define the framework now. We model the outcome of a game from the perspective of the team playing at home. The dependent variable is therefore denoted by an ordered multinomial random variable with three categories – loss, draw or win – correspondingly denoted by $r \in R = \{-1, 0, 1\}$, respectively. The focus of the model is on forecasting the outcome in real-time, that is, after every minute of the match. Therefore, we define a time index set $\Gamma = \{t : 1 \leq t \leq 90; t \in \mathbb{N}\}$. At the end of every minute $t \in \Gamma$, we record a set of covariates from minutes 1 through t . Hereafter, we refer to these as time-varying covariates.

To define the main model, let Y_i denote the outcome for the home team in the i^{th} match, for $i \in [n]$. The vector of the outcomes (Y_1, \dots, Y_n) will be denoted as \mathbf{y} . It is treated along the lines of [Greene \(2003\)](#), for an ordered probit model. We define a latent variable Π_i and cut-offs δ_1, δ_2 such that,

$$Y_i = \begin{cases} -1, & \text{if } \Pi_i < \delta_1 \\ 0, & \text{if } \delta_1 \leq \Pi_i \leq \delta_2 \\ 1, & \text{if } \Pi_i > \delta_2. \end{cases} \quad (2.1)$$

Here, Π_i is the value of the latent variable in the i^{th} match, which can be modelled as a function of the covariates. It is important to observe that we use 90 different models for time points $t \in \Gamma$. For the model till time t , for the i^{th} match, the response is an approximation of the latent variable till time t , $\Pi_i^{(t)}$. In this model, we use information of the events on the field till the time point $t \in \Gamma$. Note that [Goddard \(2005\)](#) modelled the ordinal outcome using probit regression with linear and additive functional forms of the covariates. [Angelini and De Angelis \(2017\)](#), [Koopman and Lit \(2019\)](#) also follow a similar specification in their papers. Motivated by these studies, we use a linear function with additive error for modelling the outcome of a soccer match in a similar setup. Mathematically, for the i^{th} match, till the time point t , $\Pi_i^{(t)}$ is modeled as

$$\Pi_i^{(t)} = \mathbf{z}_i^\top \boldsymbol{\gamma} + \sum_{k \in [K]} \mathbf{x}_{ik}^{(t)\top} \boldsymbol{\beta}_k^{(t)} + \varepsilon_i^{(t)}, \quad (2.2)$$

where \mathbf{z}_i is a $p \times 1$ vector of time-invariant covariates such as the strengths of the starting elevens for both the teams and $\boldsymbol{\gamma}$ is the corresponding vector of coefficients. On the other hand, $\boldsymbol{\beta}_k^{(t)}$ is a $2t \times 1$ vector of coefficients capturing the time-varying effect of an event of type k on the outcome of a match, whereas

$$\mathbf{x}_{ik}^{(t)} = [x_{ik,1}^H \quad x_{ik,2}^H \quad \cdots \quad x_{ik,t}^H \quad x_{ik,1}^A \quad x_{ik,2}^A \quad \cdots \quad x_{ik,t}^A]^\top, \text{ for } i \in [n],$$

is the associated vector indicating the number of occurrences for the k^{th} type of event until time

t for both teams. Here, $x_{ik,t}^H$ and $x_{ik,t}^A$ denote the corresponding covariates for the home (H) and away (A) teams, respectively. We assume that there are a total of K types of such events. For instance, one may consider goals, fouls, cards, shots-on-goal, corners etc. The variables used in this article are further elaborated in Section 3. It should be reiterated that the proposed model incorporates all time-varying events in a linear additive fashion.

Finally, we assume that all $\varepsilon_i^{(t)}$ are independent and identically distributed (iid) zero-mean Gaussian random variables, and we define $\varepsilon^{(t)} = (\varepsilon_1^{(t)}, \dots, \varepsilon_i^{(t)}, \dots, \varepsilon_n^{(t)})^\top$. Such an assumption for the error distribution is a popular choice in the current context, see [Koning \(2000\)](#) for example. [Koopman and Lit \(2019\)](#) also explored various methods to predict match results, and used Gaussian marginals since it leads to a parsimonious model with straightforward estimation procedures. Another motivation for such an assumption is that in a Bayesian setting which we use for implementation, a Normal prior for $\varepsilon^{(t)}$ induces conjugacy in the posterior distributions and is hence simpler to implement in Markov chain Monte Carlo (MCMC) methods. With this assumption, we can write

$$\varepsilon^{(t)} \sim \mathcal{N}_n(0, \sigma_y^2 \mathbf{I}_n), \quad (2.3)$$

where σ_y^2 is the error variance. Due to the structure of the latent variable in the model, we cannot estimate the cutoffs δ and the error variance σ_y^2 simultaneously (see [Higgs and Hoeting \(2010\)](#) for relevant discussions). Thus, we fix σ_y^2 at around 200 in the estimation algorithm below.

2.2. Bayesian Estimation

We use a complete Bayesian framework to estimate the model parameters and provide a probabilistic forecast for the response variable. In that regard, suitable prior specifications are necessary and we elaborate on them below. Also, we shall be using $f(\cdot)$ to denote conditional and marginal densities in the likelihood and prior, and $\pi(\cdot)$ to denote the posterior likelihoods.

A Gaussian prior for the covariates in a probit model is a well researched topic, cf. [McCulloch et al. \(2000\)](#). Such a prior has been found to work well in predicting soccer match outcomes too ([Rue and Salvesen, 2000](#)). Accordingly, we specify suitable Gaussian distributions as conjugate priors for γ and $\beta_k^{(t)}$ which result in straightforward conditional posterior distributions. In particular, for the two parameter vectors in the mean structure of the model, we consider

$$\gamma \sim \mathcal{N}_p(\mathbf{0}, \Sigma_\gamma), \quad (2.4)$$

where Σ_γ is a diagonal matrix, and

$$\beta_k^{(t)} = [\beta_{k,1}^H \quad \beta_{k,2}^H \quad \dots \quad \beta_{k,t}^H \quad \beta_{k,1}^A \quad \beta_{k,2}^A \quad \dots \quad \beta_{k,t}^A]^T \sim \mathcal{N}_{2t}(\mathbf{0}, \Sigma_k^{(t)}), \quad (2.5)$$

defined for each $k \in [K]$. Recall that γ is the vector of coefficients representing the effect of time-invariant covariates, whereas $\beta_{k,j}^H$ and $\beta_{k,j}^A$ (for $1 \leq j \leq t$) are the coefficients capturing the time-varying effects for the home (H) and away (A) teams respectively. In the above formulation, $\Sigma_k^{(t)}$ is a $2t \times 2t$ covariance matrix entailing the dependence between the effects of the same type of event with respect to various time points. We now specify the structure of $\Sigma_k^{(t)}$, assumed to be common for all events $k \in [K]$. First, we assume that there is no correlation between home (H) and away (A) team events, that is,

$$\text{Cov}(\beta_{k,t_1}^H, \beta_{k,t_2}^A) = 0 \text{ for all } t_1, t_2 \in \Gamma. \quad (2.6)$$

Next, to define the covariance structure common to both home and away coefficients, we assume that the dependence between occurrences of event k at distinct time points t_1 and t_2 is governed

by the structure

$$\text{Cov}(\beta_{k,t_1}^H, \beta_{k,t_2}^H) = \text{Cov}(\beta_{k,t_1}^A, \beta_{k,t_2}^A) = g(|t_1 - t_2|), \quad (2.7)$$

where $g : \mathbb{R} \rightarrow \mathbb{R}$ is a decreasing function. Throughout this paper, we assume an exponential function for g , i.e. $g(x) \propto e^{-x}$. We also assume that distinct events $k, k' \in [K]$ are not correlated amongst themselves.

Finally, we specify the prior distributions for the cutoffs δ_1, δ_2 . Previous works such as [Albert and Chib \(1993\)](#) and [Liddell and Kruschke \(2018\)](#) demonstrate that Gaussian priors work well for the cut-offs in a probit model. Taking motivation from these studies, we assume

$$\delta_j \sim \mathcal{N}(0, \tau^2), \quad j \in \{1, 2\}, \quad (2.8)$$

where τ is a large number to ensure high variability of the prior. For the purpose of this paper, we shall always use $\tau = 200$.

In order to explain the main steps of the Bayesian estimation procedure, one can note that the set of time-varying and time-invariant covariates can be combined into a single matrix denoted by $\mathcal{M}^{(t)}$, and the proposed model in equation (2.2) can be rewritten as

$$\Pi^{(t)} = \left(\mathcal{M}^{(t)} \right)^\top \boldsymbol{\nu}^{(t)} + \boldsymbol{\varepsilon}^{(t)}, \quad (2.9)$$

where $\Pi^{(t)}$ is the vector of outcomes of dimension $n \times 1$, and $\boldsymbol{\varepsilon}^{(t)}$ is defined the same as in eq. (2.3). Similarly, $\mathcal{M}^{(t)}$ is a matrix of dimension $(p + 2Kt) \times n$ comprising of all covariates, while $\boldsymbol{\nu}^{(t)}$ is a $(p + 2Kt) \times 1$ vector of parameters to be estimated. Since both $\boldsymbol{\gamma}$ and $\boldsymbol{\beta}_k^{(t)}$ are assumed to have Gaussian priors, we can further write

$$\boldsymbol{\nu}^{(t)} \sim \mathcal{N}_{p+2Kt}(\mathbf{0}, \boldsymbol{\Sigma}_0^{(t)}), \quad \text{where } \boldsymbol{\Sigma}_0^{(t)} = \begin{bmatrix} \boldsymbol{\Sigma}_\gamma & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_1^{(t)} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\Sigma}_K^{(t)} \end{bmatrix}_{(p+2Kt) \times (p+2Kt)}. \quad (2.10)$$

For implementation purposes, we use the concepts of Gibbs sampling in this work. Recall that the Gibbs sampling scheme is a modification of the conventional Metropolis-Hastings algorithm to obtain a sample from multivariate distributions without a closed form. Interested readers are referred to the works of [Geman and Geman \(1984\)](#) and [Gelfand \(2000\)](#) for further reading on this algorithm, which relies on the principle that iterative samples from the conditional posterior distributions will lead to a sample representative of their joint distribution. We shall now be dropping the time index t for the purposes of simplicity. We follow the Gibbs sampling procedure to obtain the joint posterior distribution of $\boldsymbol{\nu}$ and $\boldsymbol{\delta}$.

To start with, the likelihood can be written as

$$\mathcal{L}(\boldsymbol{\nu}, \boldsymbol{\delta} \mid \mathcal{M}, \Pi, \mathbf{y}) = \prod_{i=1}^n \prod_{j \in \{-1, 0, 1\}} \mathbb{P}(Y_i = j)^{\mathbf{1}(Y_i=j)}. \quad (2.11)$$

This can equivalently be expressed as

$$\mathcal{L}(\boldsymbol{\nu}, \boldsymbol{\delta} \mid \mathcal{M}, \Pi, \mathbf{y}) = \prod_{i=1}^n \mathbb{P}(\Pi_i < \delta_1)^{\mathbf{1}(\Pi_i < \delta_1)} \times \mathbb{P}(\delta_1 \leq \Pi_i \leq \delta_2)^{\mathbf{1}(\delta_1 \leq \Pi_i \leq \delta_2)} \times \mathbb{P}(\Pi_i > \delta_2)^{\mathbf{1}(\Pi_i > \delta_2)}. \quad (2.12)$$

Recall that Π_i is the latent variable in the model. Moreover, for the proposed structure, we can equivalently denote Y_i by an expression with Π_i and $\boldsymbol{\delta}$. Since we employ a Gibbs sampling procedure with the latent variable, we restrict ourselves to the use of Π and $\boldsymbol{\delta}$ for the estimation procedure instead of Y_i . For computational purposes, we treat Π_i as the parameter vector, and we compute the conditional distribution of it for the Gibbs sampling steps. It is straightforward to note that the conditional posterior of Π_i is given by:

$$\Pi_i \mid \boldsymbol{\mathcal{M}}, \boldsymbol{\nu}, \boldsymbol{\delta} \sim \mathcal{TN}(\boldsymbol{\mathcal{M}}'_i \boldsymbol{\nu}, \sigma_y^2, \delta_{j-1}, \delta_j), \quad (2.13)$$

where δ_{j-1} and δ_j are the truncation limit for an observation i dependent on the outcome Y_i . We define $\delta_0 = -\infty$, $\delta_3 = \infty$ for the sake of completeness. Next, the choice of a Gaussian prior assumed for $\boldsymbol{\nu}$ is essential and leads to a simple closed form expression for its conditional posterior, which can be written as

$$\pi(\boldsymbol{\nu} \mid \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\delta}) \propto \prod_{i=1}^n f(\Pi_i, \boldsymbol{\delta} \mid \boldsymbol{\nu}, \boldsymbol{\mathcal{M}}) f(\boldsymbol{\nu}). \quad (2.14)$$

From our assumption of a Gaussian prior distribution, we get

$$\pi(\boldsymbol{\nu} \mid \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\delta}) \propto \exp \left[-\frac{1}{2\sigma_y^2} \sum_{i=1}^n (\Pi_i - \boldsymbol{\mathcal{M}}'_i \boldsymbol{\nu})^2 \right] \exp \left[-\frac{1}{2} \boldsymbol{\nu}' \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\nu} \right]. \quad (2.15)$$

After some algebraic simplification and comparing it to a Gaussian distribution, we observe that the conditional posterior distribution of $\boldsymbol{\nu}$ can be expressed as,

$$\boldsymbol{\nu} \mid \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\delta} \sim \mathcal{N}_{p+2Kt}(\boldsymbol{\mu}_\nu, \tilde{\boldsymbol{\Sigma}}), \quad (2.16)$$

where $\tilde{\boldsymbol{\Sigma}} = (\sigma_y^{-2} \boldsymbol{\mathcal{M}}' \boldsymbol{\mathcal{M}} + \boldsymbol{\Sigma}_0^{-1})^{-1}$ and $\boldsymbol{\mu}_\nu = \sigma_y^{-2} \tilde{\boldsymbol{\Sigma}} (\boldsymbol{\mathcal{M}}' \Pi)$.

Now, to estimate $\boldsymbol{\delta}$, we refer to the work of [Dechi \(2019\)](#). The author establishes a correspondence between the multinomial categories through a transformation of the Dirichlet distribution. The Dirichlet distribution is a multivariate generalization of the Beta distribution with the support being a set of l -dimensional vectors with non-negative entries \mathcal{L} such that $\|\mathcal{L}\|_1$ is one. These can be considered as probabilities of a multinomial outcome with l categories. It is a natural conjugate for a multinomial outcome. For this paper, we assume a Dirichlet $(\alpha_1, \alpha_2, \alpha_3)$ distribution with each $\alpha_i = 1$. We choose this specification since it makes for a non-informative prior and due to the favourable convergence properties of its conditional marginals. The relationship between the Dirichlet parameters and the cutoffs δ_i is defined as,

$$\mathbb{P}(\delta_i \leq \delta < \delta_{i+1}) = p_{i+1}, \quad (2.17)$$

where p_i is the support of the distribution such that $\sum_{i=1}^3 p_i = 1$, and $p_i \in [0, 1] \ \forall i \in \{1, 2, 3\}$. Recall that we have defined $\delta_0 = -\infty$, $\delta_3 = \infty$ before. Then, the joint distribution of $\boldsymbol{\delta}$ can be written as,

$$f(\boldsymbol{\delta}) = F(\delta_1)^{\alpha_1-1} [F(\delta_2) - F(\delta_1)]^{\alpha_2-1} [1 - F(\delta_2)]^{\alpha_3-1} \prod_{j=1}^2 f(\delta_j), \quad (2.18)$$

where $F(\cdot)$ is any cumulative distribution function with \mathbb{R} being the domain of the random variable. Simple algebraic manipulation leads to the joint conditional posterior likelihood for $\boldsymbol{\delta}$ being

expressed as,

$$\pi(\boldsymbol{\delta} \mid \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\nu}) \propto f(\boldsymbol{\delta}) \prod_{i=1}^n \prod_{j=1}^3 \mathbf{1}(\delta_{j-1} < \Pi_i < \delta_j). \quad (2.19)$$

We now convert the joint likelihood into marginal by conditioning on the other δ_j 's. Because of the Gaussian prior for $\boldsymbol{\delta}$ in eq. (2.8), the conditional posterior distribution for each δ_j is simply

$$\pi(\delta_j \mid \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\nu}, \delta_{-j}) \propto [\Phi(\delta_j) - \Phi(\delta_{j-1})]^{\alpha_j-1} [\Phi(\delta_{j+1}) - \Phi(\delta_j)]^{\alpha_{j+1}-1} \phi(\delta_j) \mathbf{1}(c_{j,1} < \delta_j < c_{j,2}), \quad (2.20)$$

where $c_{j,1} = \max\{\Pi_i : Y_i = j\}$, $c_{j,2} = \min\{\Pi_i : Y_i = j+1\}$ for $j = 1, 2$, and $\Phi(\cdot)$ is the cumulative distribution function (CDF) of a univariate Gaussian distribution. For our model with three categories, we define $\Phi(\delta_0) = 0$ and $\Phi(\delta_3) = 1$ for completeness. We now identify the conditional posterior CDF for $\delta_j \mid \delta_{-j}$. In that regard, we define $\omega = (\Phi(\delta_j) - \Phi(\delta_{j-1})) / (\Phi(\delta_{j+1}) - \Phi(\delta_{j-1}))$. Letting $\Phi(\delta_{j-1}) = a$ and $\Phi(\delta_{j+1}) = b$, we can write $\delta_j = \Phi^{-1}((b-a)\omega + a)$. Utilizing the distribution of δ_j , with regard to the distribution of ω , it can be defined as $f(\omega) = \pi_\delta(\delta_j) \frac{\partial \delta_j}{\partial \omega}$, where $\pi_\delta(\cdot)$ is the unconditional posterior distribution of δ_j . By substituting the form of δ_j in eq. (2.20), we can obtain the expression

$$\pi_\delta(\Phi^{-1}((b-a)\omega + a)) \propto (\Phi(\delta_j) - a)^{\alpha_j-1} (b - \Phi(\delta_j))^{\alpha_{j+1}-1} \phi(\delta_j). \quad (2.21)$$

A simple algebraic simplification of the above equation and ignoring terms without ω leads us to,

$$\pi_\delta(\delta_j) \propto \omega^{\alpha_j-1} (1 - \omega)^{\alpha_{j+1}-1} \phi(\delta_j). \quad (2.22)$$

The derivative of ω with respect to δ_j is simply given by $\frac{\partial \omega}{\partial \delta_j} = \frac{b-a}{\phi(\delta_j)}$. Now that we have obtained simplified expressions for the terms required to identify $f(\omega)$, one can argue that

$$f(\omega) \propto \omega^{\alpha_j-1} (1 - \omega)^{\alpha_{j+1}-1} \phi(\delta_j) \times \frac{b-a}{\phi(\delta_j)}. \quad (2.23)$$

Simplifying the above expression and comparing it to a Beta distribution, we get

$$\omega \sim \text{Beta}(\alpha_j, \alpha_{j+1}). \quad (2.24)$$

In this manner, we can now incorporate the parameters of the Dirichlet distribution, namely α_j , in the distribution of δ_j through a closed-form expression for ω . Hence, the conditional posterior distribution of δ_j can be expressed as

$$\Phi(\delta_j \mid \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\nu}, \delta_{-j}) \sim [\Phi(\delta_{j+1}) - \Phi(\delta_{j-1})] \text{Beta}(\alpha_j, \alpha_{j+1}) + \Phi(\delta_{j-1}). \quad (2.25)$$

Note that $\Phi(\delta_j)$ will be truncated from below and above by $\Phi(c_{j,1})$ and $\Phi(c_{j,2})$ respectively. Let us use the shorthand notation Φ_{δ_j} to denote the conditional distribution of δ_j given the other cutoffs. Therefore, in the modelling setup above, in the case of three ordered categories we obtain the following set of conditional posterior distributions:

$$\begin{aligned} \Phi(\delta_1 \mid \delta_2, \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\nu}) &\sim \Phi(\delta_2) \text{Beta}(\alpha_1, \alpha_2) & \Phi(c_{1,1}) &\leq \Phi_{\delta_1} \leq \Phi(c_{1,2}), \\ \Phi(\delta_2 \mid \delta_1, \boldsymbol{\mathcal{M}}, \Pi, \mathbf{y}, \boldsymbol{\nu}) &\sim [1 - \Phi(\delta_1)] \text{Beta}(\alpha_2, \alpha_3) + \Phi(\delta_1) & \Phi(c_{2,1}) &\leq \Phi_{\delta_2} \leq \Phi(c_{2,2}). \end{aligned} \quad (2.26)$$

Now that we have the closed-form expressions for all the conditional posterior distributions corresponding to the unknown parameters in the model, we can implement the Gibbs sampling

algorithm to sample from the joint posterior distribution. Following the principles of this algorithm, we need to sequentially sample from the distributions of $(\Pi \mid \boldsymbol{\nu}, \boldsymbol{\delta})$, $(\boldsymbol{\nu} \mid \Pi, \boldsymbol{\delta})$ and $(\delta_j \mid \delta_{-j}, \Pi, \boldsymbol{\nu})$ until convergence. We use the GW statistic to assess the convergence, and it is explained below. After the Markov chains in the Gibbs sampler converge, we can obtain a sample from their joint distribution. To ensure the independence of the observations we incorporate thinning by drawing samples from every 10th iteration. The steps followed in the implementation of the Gibbs sampling procedure are now summarized in Algorithm 1.

Algorithm 1: Gibbs sampler for the posterior distribution in the proposed model

Input: Dataset $(\mathbf{y}, \mathcal{M})$ where, \mathbf{y} is the multinomial outcome variable and \mathcal{M} is the set of covariates.

Output: A sample of size S from the posterior distribution of the tuple $(\Pi_i, \boldsymbol{\nu}, \boldsymbol{\delta})$

Initialize: $\Pi^{(0)}$, $\boldsymbol{\nu}^{(0)}$ and $\boldsymbol{\delta}^{(0)}$. Let \mathcal{C} be the convergence criteria.

Let $m \leftarrow 1$

while \mathcal{C} not met **do**

Sample from $\Pi_i^{(m)} \mid \boldsymbol{\nu}^{(m-1)}, \boldsymbol{\delta}^{(m-1)} \sim \mathcal{TN}(\mathcal{M}_i^\top \boldsymbol{\nu}^{(m-1)}, \sigma_y^2, \delta_{j-1}, \delta_j)$

Sample from $\boldsymbol{\nu}^{(m)} \mid \Pi^{(m)}, \boldsymbol{\delta}^{(m-1)} \sim \mathcal{N}(\tilde{\boldsymbol{\nu}}^{(m)}, \tilde{\boldsymbol{\Sigma}})$, where $\tilde{\boldsymbol{\nu}}^{(m)}$ is the updated mean of the conditional posterior distribution after observing $\Pi^{(m)}$

foreach $j \in \{1, 2\}$ **do**

Sample from

$\Phi(\delta_j^{(m)}) \mid \boldsymbol{\delta}_{(-j)}^{(m-1)}, \Pi^{(m)}, \boldsymbol{\nu}^{(m)} \sim [\Phi(\delta_{j+1}^{(m-1)}) - \Phi(\delta_{j-1}^{(m)})] \text{Beta}(\alpha_j, \alpha_{j+1}) + \Phi(\delta_{j-1}^{(m)})$

with, $\Phi(\delta_j^{(m)}) \in [\Phi(c_{j,1}), \Phi(c_{j,2})]$

end

Let $m \leftarrow m + 1$

end

Discard these first M iterations (until \mathcal{C} is met) as burn-in sample.

Continue the iteration procedure given above until we reach iteration $M + S$. The requisite sample is obtained from iterations $M + 1$ to S .

In order to assess the convergence of the Gibbs sampler, we use the [Geweke et al. \(1991\)](#) statistic, hereafter abbreviated as GW statistic. In this diagnostic approach, a single chain is generated using the Gibbs sampler, and the spectral analysis is used to assess the convergence of the procedure. To elaborate, let $\{g(\theta^{(1)}), g(\theta^{(2)}), \dots, g(\theta^{(n)})\}$ be the iterations of the Gibbs sampler and $S_g(w)$ be the spectral density for the series. Then, under regularity conditions, we can write,

$$E(g(\theta)) = \bar{g}_n = \frac{1}{n} \sum_{i=1}^n g(\theta^{(i)}), \quad (2.27)$$

with the asymptotic variance being $S_g(0)/n$. After n iterations of the Gibbs sampler, the GW statistic is calculated by taking the standardized difference between the means \bar{g}_n^A and \bar{g}_n^B , which are computed based on the first n_A and the last n_B iterations respectively. As n increases, this statistic tends to a standard normal distribution ([Cowles and Carlin, 1996](#)), and that property is utilized to find evidence of the convergence.

2.3. Predictive analysis

It is important to note that the methodology developed above is primarily used to forecast the match outcome after every time-point. To explain this within-game forecasting technique, consider that the model needs to be trained on the dataset \mathcal{S}_{tr} , and let the test dataset be \mathcal{S}_{te} , with $|\mathcal{S}_{te}| = m$. We shall compare the evaluation metrics (to be elaborated below) for our model as well as for other potential models (discussed later) using the test data. Suppose, we wish to obtain the in-game predictions for the i^{th} match in \mathcal{S}_{te} . We shall compute this as a function of $t \in \Gamma$. We know that data are recorded on various events such as goals, corners, cards etc., which take place every single minute. This renders a minute-by-minute record of events which have occurred till time t . Following the same notations as before, for the training data, let the set of covariates after time t be denoted as \mathbf{M} , and the entire training data as $\mathcal{D} = \{\mathbf{M}, \mathbf{y}, \Pi\}$. Our aim is to implement the proposed model on \mathcal{D} to predict the outcome of the i^{th} match in \mathcal{S}_{te} , after t minutes have passed in the match. Let us use $\hat{\Pi}_i^{(t)}$ to denote the estimated latent variable in this regard, and define the vector $\hat{\Pi}^{(t)} = [\hat{\Pi}_1^{(t)}, \dots, \hat{\Pi}_m^{(t)}]^\top$ which furnishes the forecasts for all matches in the test dataset at time $t \in \Gamma$.

Let us use the general notation $\mathbf{M}_* = [\mathbf{M}_{1*}, \dots, \mathbf{M}_{(p+2Kt)*}]^\top$ to denote the set of covariates corresponding to the test set. Then, $\hat{\Pi}_i^{(t)}$ can be obtained from the posterior predictive distribution that estimates the probabilistic structure of the outcome variable given a new set of covariates. It incorporates the variability in the parameters by weighting the likelihood of $\hat{\Pi}^{(t)}$ by the posterior distribution of the parameters. For our model, we can write

$$\pi(\hat{\Pi}^{(t)} | \mathbf{M}_*, \mathcal{D}) = \int_{\boldsymbol{\nu}} f(\hat{\Pi}^{(t)} | \mathbf{M}_*, \mathcal{D}, \boldsymbol{\nu}) \pi(\boldsymbol{\nu} | \mathcal{D}) d\boldsymbol{\nu}. \quad (2.28)$$

Here, the distribution of the forecasts for the test data is conditional on the training data \mathcal{D} and the posterior distribution of $\boldsymbol{\nu}$. The first term in the expression is simply a likelihood taking a Gaussian form, due to eq. (2.3). The second term is the posterior distribution of $\boldsymbol{\nu}$ given by eq. (2.16). The above expression can be simplified since both are normally distributed. Thus, the posterior predictive distribution for $\hat{\Pi}_i^{(t)}$ can be simplified as

$$\pi(\hat{\Pi}_i^{(t)} | \mathbf{M}_{i*}, \mathcal{D}) \sim \mathcal{TN}(\mathbf{M}_{i*}^\top \boldsymbol{\mu}_\nu, (\sigma_y^2 + \mathbf{M}_{i*}^\top \tilde{\boldsymbol{\Sigma}} \mathbf{M}_{i*}), \delta_{j-1}, \delta_j), \quad (2.29)$$

where $\tilde{\boldsymbol{\Sigma}}$ is the posterior covariance matrix of $\boldsymbol{\nu}$, and δ_{j-1} and δ_j are the truncation limits. As stated above, the estimate of the latent variable, denoted by $\hat{\Pi}$, for the training set can be obtained from the posterior sample through the Gibbs sampler outlined in Algorithm 1. We also estimate $\tilde{\boldsymbol{\Sigma}}$, and the cutoffs $\hat{\boldsymbol{\delta}}$ from the posterior samples. Now, in order to predict the outcome, one can obtain a sample from the above posterior predictive distribution, and use that along with the cutoffs to get $\hat{Y}_i^{(t)}$ as the predicted category based on the data up to time t . Moreover, the probabilities of the different categories for the outcome variable, corresponding to the home team, can be calculated as

$$\begin{aligned} \mathbb{P}(\text{Win}) &= 1 - \Phi_{\hat{\Pi}_i}(\hat{\delta}_2), \\ \mathbb{P}(\text{Draw}) &= \Phi_{\hat{\Pi}_i}(\hat{\delta}_2) - \Phi_{\Pi_i}(\hat{\delta}_1), \\ \mathbb{P}(\text{Loss}) &= \Phi_{\hat{\Pi}_i}(\hat{\delta}_1), \end{aligned} \quad (2.30)$$

where $\Phi_{\hat{\Pi}_i}$ is the Gaussian cumulative distribution function of $\hat{\Pi}_i$. We are going to use $\hat{Y}_i^{(t)}$ and the derived probabilities in eq. (2.30) to evaluate the forecasting accuracy of the proposed methodology.

It is important to reiterate that the evaluation metrics will be computed for the vector of $\hat{\Pi}_i^{(t)}$ for every time point $t \in \Gamma$.

While the above procedure renders a point forecast for the outcome probabilities, it also gives us an opportunity to compute a credible interval for the forecasts. Due to the Bayesian framework and Algorithm 1, we have the joint posterior distribution of $\boldsymbol{\nu}$, $\boldsymbol{\delta}_1$ and $\boldsymbol{\delta}_2$ at hand. Drawing a sample from said distribution will result in a unique point estimate for $\hat{\Pi}_i^{(t)}$ and consequently $\hat{Y}_i^{(t)}$ and the predicted probability for each category by means of eq. (2.30). Repeating this procedure many times, we can generate a set of point estimates which can be used to obtain a credible interval for the predicted outcomes. For implementation of this approach, in this paper, we shall draw 1000 samples from the joint posterior distribution of $\boldsymbol{\nu}$ and $\boldsymbol{\delta}$, and obtain 1000 realizations for $\mathbb{P}(\text{Win})$, $\mathbb{P}(\text{Draw})$ and $\mathbb{P}(\text{Loss})$, as mentioned above. Then, we compute a 95% equal tailed credible interval for these probabilities. We repeat the process for each model to obtain credible intervals for every time-point. A flavour of this procedure can be found in Section 4.3, wherein we obtain the credible intervals for the predicted probabilities throughout the course of the match, for specific games.

In the main application of this article, for the completeness of the study, we are going to compare our model with a few other potential approaches. We highlight that there is no existing work on the topic of within-game forecasting in soccer, but we rely on different algorithms that have been used for predictive modelling in different capacities and extend them to develop competing approaches in our context. The first one in this regard is a version of the generalized linear model (GLM). We train the GLM in a standard probit regression framework by considering the data on the available covariates up to time t for all matches, and predict the outcome in the test set based on that. Thus, on the same lines as our proposed method, the GLM needs to be trained after every time point $t \in \Gamma$, and the forecasts are updated accordingly.

We next refer to the work of Baboota and Kaur (2019), who identified a set of features which highly influence the result of a soccer match, through extensive feature engineering and selection. The authors used support vector machine (SVM) and random forest algorithms for predicting the outcome of the match. Although not used in a temporal within-game setting, we consider these models as competitors to our model due to their flexibility in implementation. The covariates selected for modelling by the authors are similar to ours, which serve as another motivation for choosing them as competitors. Note that SVMs are very flexible supervised-learning models, usually used for classification problems. In the comparison study below, when we employ the SVM as one of the contenders, we employ two types of SVMs with respect to the kernel used for modelling, which is one of the most important hyper-parameters in this algorithm. For the first model, we assume linearly separable classes and we shall denote it as Linear SVM below. In the second case, a Gaussian radial basis function is used to incorporate non-linearity in the model. This model will be abbreviated as R-SVM hereafter. The R library `caret` is used for tuning both models.

As the fourth model in the comparative discussions, we modify the standard random forest (RF) algorithm, typically used in various classification problems related to soccer. This tree-based algorithm is widely implemented due to its tendency to incorporate dependencies between the covariates in the model. The performance of the random forest algorithm usually depends on finding the precise hyper-parameter. We employed a grid search technique for tuning the hyperparameters. We use the functions provided in the `base` library for the same. The grid search is performed for identifying the \mathcal{L}_1 and \mathcal{L}_2 regularization parameters, subsample ratio of columns when constructing each tree and hyper-parameters to control for extreme class imbalance. We search over the interval $[0,1]$ for the optimal subsample ratio and $[0,5]$ for the remaining hyper-paramters. The function `expand.grid` is used, which primarily iterates the model over all possible combination and ranges of the parameters as required (see Chambers and Hastie, 2017, for a detailed discussion). As the

set of possible features in this algorithm, akin to the other contending approaches, we use the same combination of covariates, the information being available up to time t in every match.

In order to compare the performances of different algorithms mentioned above, we first rely on the F1-score (see Ch. 7 of [Van Rijsbergen, 1979](#), for more discussions on this). It is a standard evaluation criterion of prediction accuracy for a categorical outcome variable, and is essentially the harmonic mean of sensitivity (recall) and specificity (precision). It ranges from zero to one, the latter depicting better accuracy. Below, we define this measure for our proposed model, and a similar computation will be done for the other competitors as well. To avoid notational jargon, we also remove the superscripts indicating the forecast made at time t , and it is understood throughout that the metrics are computed as a function of time during the progression of the match.

Let $\hat{\Pi}_i$ be the predicted value of the latent variable for the i^{th} match. Following the definition of the categorization, we can determine \hat{Y}_i as j if $\hat{\delta}_{j-1} \leq \hat{\Pi}_i < \hat{\delta}_j$. Then, for the j^{th} outcome category \mathcal{C}_j , the confusion matrix for the test set is reduced to a 2×2 matrix given by

$$\begin{array}{cc} & \begin{array}{c} \mathcal{C}_j \\ \mathcal{C}_{\neq j} \end{array} \\ \begin{array}{c} \mathcal{C}_j \\ \mathcal{C}_{\neq j} \end{array} & \begin{bmatrix} a_{j,11} & a_{j,12} \\ a_{j,21} & a_{j,22} \end{bmatrix}, \end{array} \quad (2.31)$$

where $a_{j,11}$ indicates the total number of matches in the test set where correct classification is made in the j^{th} category, and so on. From this confusion matrix, sensitivity and specificity are then defined as

$$\text{Sen}_j = \frac{a_{j,11}}{(a_{j,11} + a_{j,21})}, \quad \text{Spc}_j = \frac{a_{j,22}}{(a_{j,12} + a_{j,22})}. \quad (2.32)$$

Subsequently, the F1-score is computed as

$$\text{F1}_j = \frac{2(\text{Sen}_j)(\text{Spc}_j)}{\text{Sen}_j + \text{Spc}_j}. \quad (2.33)$$

Note that the F1-score is computed for each outcome category separately. This enables us to identify biases in forecasts with regard to the classes, if any. A high F1-score implies that the model can consistently forecast the correct outcome.

One of the criticisms of the F1-score is that it gives larger weight to smaller classes and favours models with similar sensitivity and specificity. It is imperative to use another criterion to evaluate the predictive accuracy of the models. Following the discussions by [Czado et al. \(2009\)](#) and [Kolassa \(2016\)](#) who pointed out the need to take into account the probability with which the forecast is made in similar problems, we are going to use the Brier score. It is a widely used scoring rule for multi-class prediction problems with mutually exclusive classes. The reader is referred to the works by [Brier et al. \(1950\)](#) and [Murphy \(1973\)](#) for the definition and related discussions on the Brier score. The scoring mechanism takes into account the probabilities of classification and compares them against the observed outcome. Recall the predicted probabilities in eq. (2.30) and let \hat{p}_{ir} be the value corresponding to the r^{th} category, for the i^{th} match in \mathcal{S}_{te} . Then, the Brier score for the test set is given by

$$\text{Brier score} = \frac{1}{m} \sum_{i=1}^m \sum_{r \in R} (\hat{p}_{ir} - \mathbf{1}(Y_i = r))^2. \quad (2.34)$$

Evidently, a lower Brier score implies a better predictive performance of the model. As mentioned before, these evaluation criteria will be reported for the competing models for every minute $t \in \Gamma$ of the match.

3. Data

3.1. Description

In this article, we use the data from English Premier League (EPL) matches from the 2008-09 season to the 2015-16 season. It is extracted from the European Soccer Database (ESD), which is available on Kaggle (link: <https://www.kaggle.com/datasets/hugomathien/soccer>).

EPL is the top division in the English soccer system. Every season, 20 teams play in a double round-robin format, and in the end, the bottom three teams are relegated to the English Football League (EFL). To maintain the 20 team format, the top three teams from EFL are promoted to play in the EPL for the next season. Such a structure results in the total number of matches recorded over 8 seasons being 3040. It is important to mention that it is not the same 20 teams playing in the EPL year on year, and consequently, the dataset has a total of 34 teams, each playing a varying number of seasons. One should note that the league matches finishing in tied scores by the end of regular time (90 minutes) do not go into extra-time or penalty shootouts. This enables an ordinal multinomial outcome with three categories. To analyze this dataset, we use a multinomial response variable to illustrate the match result for the team playing at home. As mentioned earlier, the covariates in the model are of two types: time-invariant covariates and time-varying ones. For the latter, remember that the dataset reports different types of events happening in every match, with their corresponding times of occurrence. We use eight such events in the model as time-varying covariates. These are goals, shots-on-goal, shots-off-goal, red cards, yellow cards, corners, crosses and fouls. Each event k has an associated vector of covariates as defined in eq. (2.2).

Regarding the time-invariant covariates, we consider the strength of the playing elevens for both teams. In order to define this, we rely on the ESD that records the overall rating for each player in the league, updated periodically based on their real-life performances on different parameters. These ratings consider 33 types of skills of the players, and the details can be found in the aforementioned link. Based on the information on the eleven players who start a game for each team, we compute the covariate depicting the overall strengths of the team. Specifically, the average of the ratings of the players in the starting eleven is computed based on the players' ratings at the beginning of the match. We point out that due to the absence of substitution data, the strength variable is assumed to be fixed over the course of the match, and thus, we classify it as a time-invariant covariate. Apropos to this point, note that the effects of specific opponents can also be estimated by introducing team-specific fixed effects, instead of the strength variables, in the model. We choose the latter to avoid the issues of over-fitting.

A brief discussion on the motivation behind the above choices is of the essence here. Previously, [González-Rodenas et al. \(2019\)](#) and [Gómez et al. \(2018\)](#) showed that the strength of a team based on its players and team rankings is useful in predicting the outcomes. Many other studies have used various events at an aggregate level as regressors for predicting match outcomes in soccer. [Liu et al. \(2015\)](#), for example, used a generalized linear model to identify key winning indicators from 24 different event types from World Cup data. The authors identified shots, shots-on-target, tackles, red cards and crosses as important events which influence the result of a match. An interesting finding is the negative impact of crosses, which corroborates [Vecer \(2014\)](#) who demonstrated through a multilevel Poisson regression model that crosses indeed have a negative impact on the goals scored by a team. Earlier, [Castellano et al. \(2012\)](#) used discriminant analysis to infer about various attack and defence attributes with regard to their effect on the result. Total shots and shots on target were found to have the greatest discriminatory power amongst the variables used. In more recent studies, [Ashimolowo \(2018\)](#) investigated the association of crosses, corners, free kicks and the number of shots-on-goal with the outcome of a soccer match, whereas [Červený et al.](#)

(2018) modelled the effect of a card on the goal-scoring rate of a team through a proportional hazard model.

3.2. Exploratory analysis

Before moving on to the main analysis, we find it imperative to present the descriptive statistics of the covariates used in the model. Table 1 displays the summary statistics of these variables for both home and away teams. We report the mean and standard deviation per game, along with the percentage of matches with no events of that type.

Table 1: Summary of the covariates used in the analysis. Mean and standard deviation (SD) are computed per game, whereas the zeros column indicates the percentage of matches with zero cases.

Variable	Home		Away	
	Mean (SD)	Zeros (%)	Mean (SD)	Zeros (%)
Goal	1.551 (1.312)	22.829	1.160 (1.145)	34.309
Shot-on	6.684 (3.512)	0.559	5.274 (2.936)	2.072
Shot-off	6.606 (3.091)	0.428	5.212 (2.678)	1.349
Red Card	0.065 (0.254)	93.717	0.097 (0.312)	90.822
Yellow Card	1.418 (1.170)	24.243	1.802 (1.286)	15.757
Corner	10.196 (5.810)	1.283	8.048 (5.080)	2.368
Cross	16.128 (7.609)	0.066	12.448 (6.243)	0.329
Foul	10.691 (3.546)	0	11.398 (3.687)	0
Team Strength	76.081 (3.784)	-	75.843 (3.871)	-

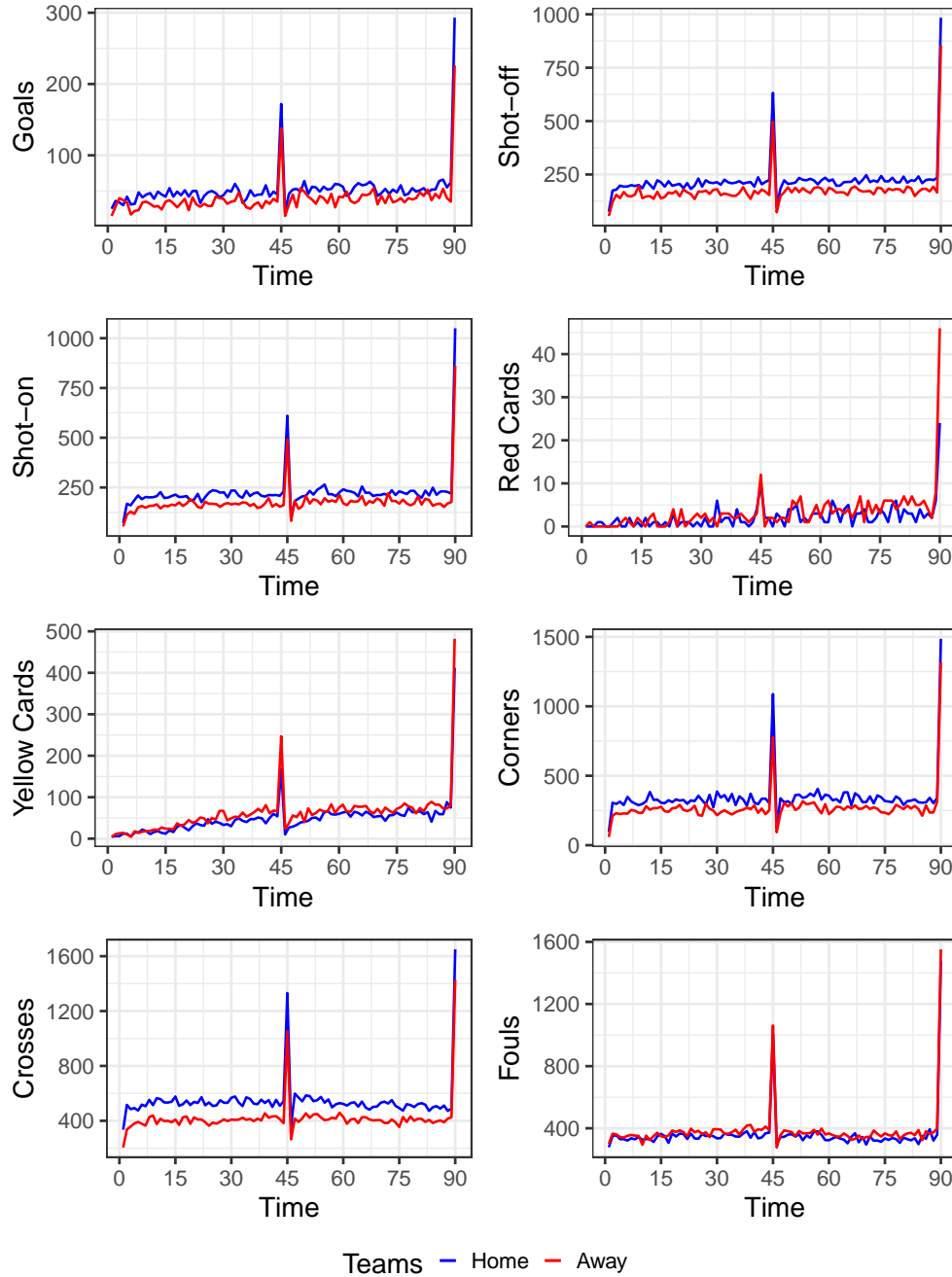
We observe that the average number of goals scored by the home teams is slightly more than that by the away teams. This home advantage can be further illustrated by the fact that a home team is held scoreless in about 23% of the matches as opposed to 34% for the away team. This difference is a result of more shots-on-goal being taken by the home team. Identical observations can be made about shots-off-goal as well. When it comes to corners, the home team is awarded two more corners than the away team on average. Average crosses per the game trend in the same direction with the home team averaging four extra crosses. However, one may note that the strength variable, as expected, does not appear to be significantly different between the home and the away teams.

Coming to the disciplinary covariates, red cards are found to be infrequent events. In EPL, a red card is shown less than once in every ten matches on average. Comparing the raw per-game numbers, we observe that the away team is 1.5 times more likely to get a red card than the home team. Although the difference in yellow cards per game for the home and away teams is much smaller, the event of no yellow cards being shown in a match happens in nearly 10% more matches for the home team as compared to the away side. In terms of the number of fouls, we observe a similar trend but to a lesser extent. One can possibly attribute these differences to refereeing bias towards the home team (Boyko et al. (2007)).

We further illustrate the nature of the time-varying covariates in the proposed model. To that end, Figure 1 provides valuable insights into the aggregate number of events over the course of a match for both teams. In these plots, the total number of events for each type, computed from the entire dataset of 3040 matches, are presented.

A striking feature for all the events is the spike observed at the 45th minute and the 90th minute. This phenomenon is related to the fact that stoppage time is provided to compensate for any delays that may have occurred during the preceding half, and during this time typically all teams tend to play with high intensity. We also want to point out that in this analysis, any event occurring

Figure 1: Aggregate number of events (computed for 3040 matches in the dataset) of different types over the course of a match for home and away teams.



in stoppage-time is reported corresponding to the last minute of the half, that is with the 45th or the 90th minute as appropriate, and therefore we are unable to assess any potential impact of the stoppage time on the outcome variable.

Among the specific covariates, the pattern of the number of cards given as a function of time appears to be interesting. We observe that the frequency of yellow cards being handed out significantly increases as a match progresses. A similar increasing trend is observed in the case of red cards, albeit at a much slower rate. The difference in means between home and away teams for red

cards can largely be attributed to the spike of the red cards given to the away team in stoppage time. However, no such trend has been found in cases of fouls, potentially hinting at the referees becoming stricter in their rulings, or a higher degree of aggressive play by both teams towards the end of the matches. We observe a systematic difference in the number of crosses and shots taken over time. This consequently results in an increased difference in the number of corners and the goals scored between the home and away teams. Interestingly, these differences are constant over time, but as we shall see in the main analysis, they are expected to have a differential impact on the final outcome of the game.

4. Results and Discussion

4.1. Primary results

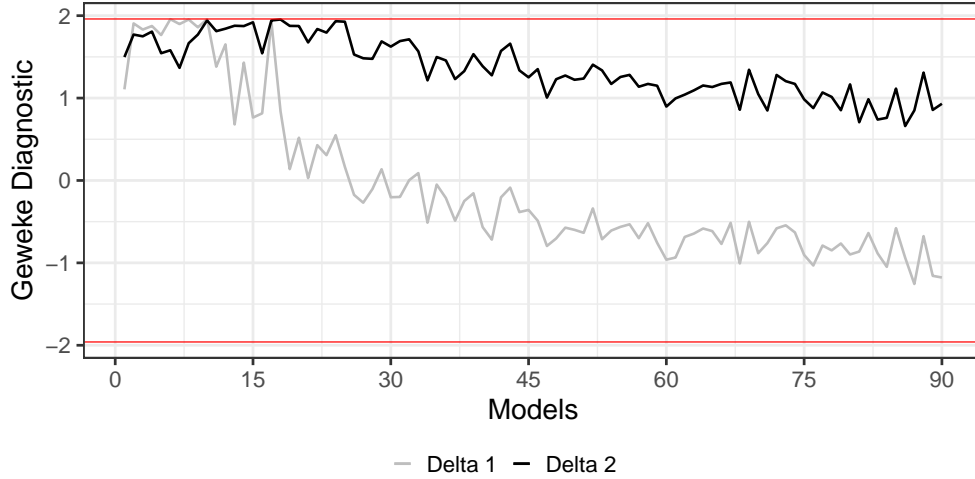
For the main analysis, the focus is on understanding the predictive accuracy of the proposed methodology. In order to assess that, we split the dataset – 90% of the data are used for training the model, while the remaining 10% are used as the test data. The results reported in this section correspond to this split, whereas in Section 4.2 we shall check the robustness of the algorithm by considering different training and test datasets. All calculations are carried out using RStudio Desktop (R version 4.3.1) using a 16-core 32GB processor. The codes for running the main algorithms are made publicly available on a GitHub repository maintained by the first author. The library `future.apply` is used for parallel computing while `e1071` and `ranger` are used for the competing models. Note that the parallelization is used to reduce the computational time for the estimation of 90 independent models, one for each time-point.

The prior specifications and the Gibbs sampling steps are detailed in Section 2. As stated there, the GW statistic is used as a diagnostic test for convergence. In order to obtain the posterior sample for the model parameters, we iterate until the GW statistic is below the threshold for all models which we classify as the burn-in period. By implementing thinning we then draw samples from the posterior after convergence. We first present the convergence results for all 90 models. A value between -1.96 and 1.96 for the GW statistic is an indicator of convergence. It is useful to note that we have specified a covariance structure which governs the posterior behaviour of $\boldsymbol{\nu}$. Hence, the convergence of the Gibbs sampler is evaluated based on the GW statistic for δ_1 and δ_2 . Figure 2 shows the GW statistic for both δ_1 and δ_2 for every model post convergence, and we can ascertain that the convergence has been reached by all models. For all these models, the number of iterations required for achieving convergence ranges between 11,000 to 200,000, with a mean number of iterations being around 30,000. It is also observed that the convergence is faster for models with more information, i.e., for models at later time points. We acknowledge that the speed of convergence is impacted by the size of the data as well as suitable prior specifications, in line with the observations in similar Bayesian estimation problem for multinomial probit models by Imai and Van Dyk (2005).

We restate that the goal is to forecast the match result in-game, that is when minute-by-minute data for the covariates are recorded in real-time. In line with this idea, we estimate the model with the time-varying covariates taken up to the t^{th} time point for each $t \in \{1, 2, \dots, 90\}$. We further forecast the match results for the test data for every such model using the mean of the posterior predictive distribution. Following the aforementioned notations, we use $\hat{\Pi}_i^{(t)}$ to obtain the predicted outcome of i^{th} game when t minutes of the match has passed, that is, the covariate information is available upto time t .

We first examine the posterior estimates of the covariate effects in the model. From the posterior samples obtained from the proposed algorithm, we find the posterior means and the credible

Figure 2: GW statistic (for the two key parameters) used to infer about convergence and decide the burn-in period for all 90 models.



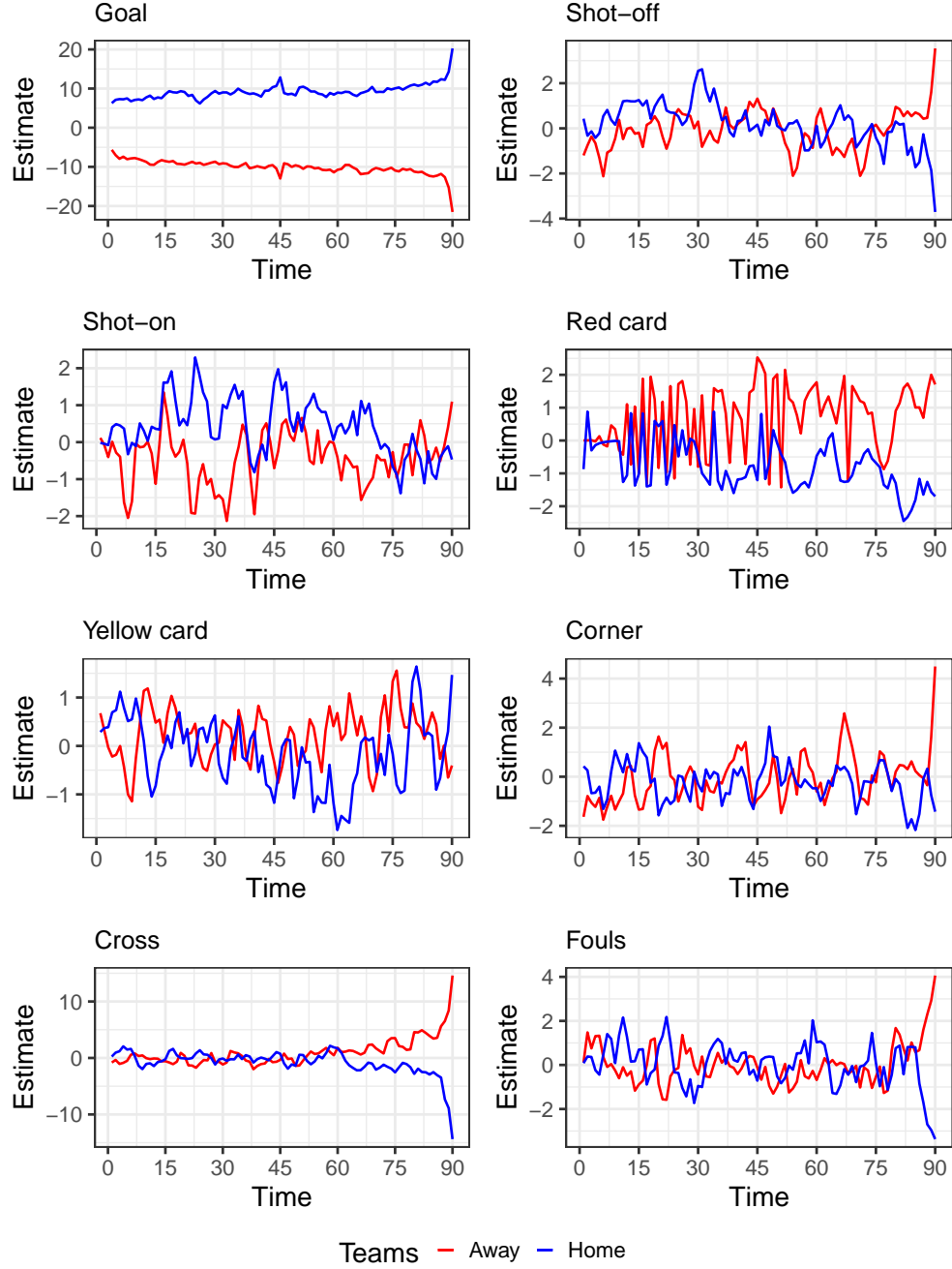
intervals for each parameter. In the case of the time-invariant covariates, the impact of the team strength is found to be substantial for both home and away teams. The coefficient of the home team is obtained to be 1.69, with the 95% credible interval being (1.17, 2.19). The same for the away team are -1.32 and $(-1.74, -0.83)$, respectively. As expected, one can see that the effect of home strength and away strength are relatively equal and of opposite signs. Next, in Figure 3, we demonstrate the temporal variation in the impact of various events on the probability of winning for the home team. A positive value implies an increase in the win probability by more occurrences of the event at that time-point, whereas a negative coefficient implies otherwise. In the figure, for ease of interpretation, all coefficients are scaled by the respective standard errors.

Goals are naturally the best indicator of the outcome of a match, as can be gleaned from the figure. We can observe an increasing effect of goals scored in the latter stages of the match as compared to the initial minutes. This phenomenon can be rationalized by the fact that a team has less time to respond when a goal is scored towards the end. The findings align with the results from [Castellano et al. \(2012\)](#), [Červený et al. \(2018\)](#) and [Rocha-Lima et al. \(2021\)](#). Next, the variables shots-on-goal and shots-off-goal exhibit marginal influence on the outcome of a soccer match. Their impact is typically found to be negligible throughout the timeline. This can be attributed to the inclusion of goals scored as an event, which might capture most of the variability explained by shots-on-goal and shots-off-goal.

When it comes to cards, both red and yellow cards are usually negatively associated with the outcome. Specifically, a red card for any home is found to have a detrimental effect on their chances of winning, making it the second most crucial event in a soccer match. It is noteworthy that red cards, akin to goals, have a significant effect on the result over the course of the match; whereas yellow cards generally do not render deciding impact. With regard to temporal fluctuations, a red card awarded to the opponent in the latter stages of the match is more advantageous for the home team as compared to the away team.

Corners, though deemed advantageous in a soccer match ([Ashimolowo \(2018\)](#)), do not prove to be as effective in forecasting the outcome. Our findings are contrary to the common notion that corners are significant factors when awarded in the dying stages of a match. This is especially true in the 90th minute and stoppage time where getting a corner does not impact the outcome in any significant manner. A contrasting behaviour is observed in the effect of a cross, which is

Figure 3: Effect of each event on the outcome of the match captured at every time-point for both teams. The estimates displayed herein are from the model with complete match data and scaled by their respective standard errors.

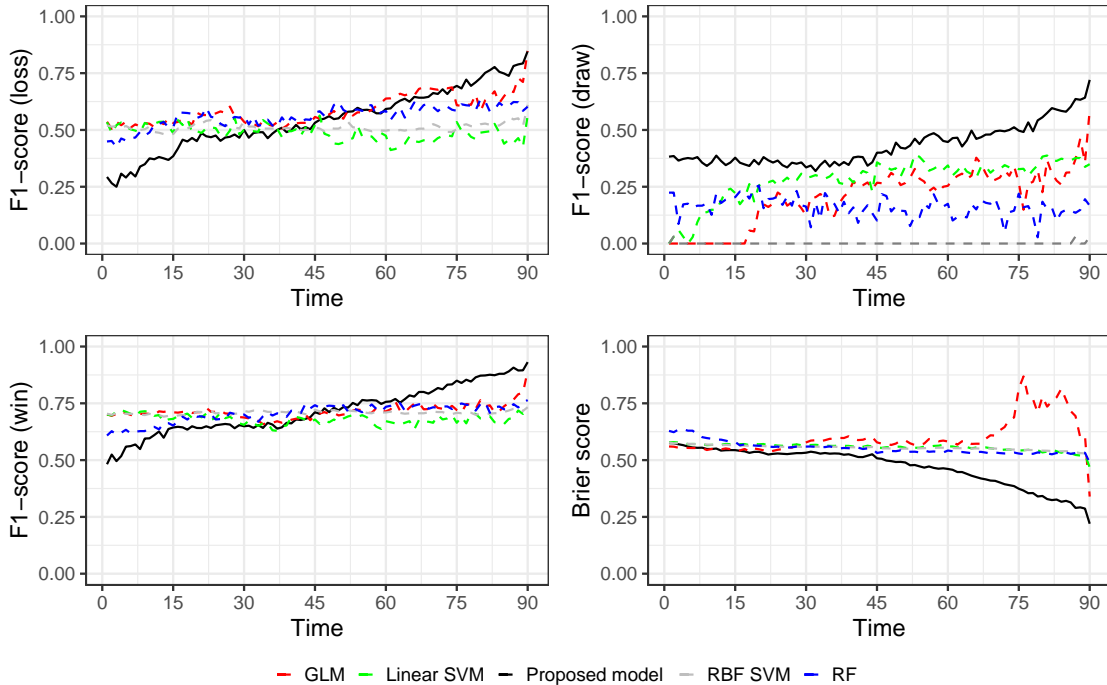


conventionally a long-range pass made towards the opponent's goal from the two sides of the field. For the majority of the match, the influence of a cross on the outcome is discovered to be negligible, but after the 75th minute, this changes. As time winds down, teams tend to be more aggressive and push for a goal. This usually leads to more long balls and crosses, and an inaccurate cross may result in a decisive counterattack on the other end, which has been hypothesized by other researchers as well (Lepschy et al., 2020, for example). Interestingly, this phenomenon is captured

by the model which estimates a significant negative impact of crosses towards the end of the match for both teams. The negative influence of crosses keeps on increasing till the end. It is imperative to recall that this observation conforms with the results of Vecer (2014) and Liu et al. (2015), as discussed earlier.

We now move on to the results of the predictive accuracy, the main focus of this article. The F1-scores and the Brier score, as defined in Section 2.3, are displayed in Figure 4. These accuracy metrics are computed as averages based on the predicted outcomes at each time point of the matches in the test set. We report the values for our method, along with the same for each of the four competing models.

Figure 4: F1-score and Brier score (averaged over the test set) for different models, with respect to their within-game forecasting accuracy as a soccer match progresses



F1-score, as discussed before, is determined for every outcome category independently. When the match result is a loss for the home team, we observe a distinctly higher F1-score for our model as the match unfolds. A similar trend is observed for the wins as well. Contrary to that, if the final outcome is a draw, the proposed method outshines its competitors throughout the time frame of the match. The linear SVM is the second-best-performing model in such cases.

To elaborate on the behaviour of the F1-scores thus obtained, one may argue that the competitor models weigh the effects of the goals in a much more severe manner as compared to the proposed model. Our model does not rely that heavily on goals and incorporates the effect of other events while forecasting. As a direct consequence of this, the proposed model performs much better than all the contender models while predicting drawn games. Also, in games where the scores are level, other covariates, such as team strength, home advantage, and player performance, are likely to have a greater impact on the forecast outcome which plays to our models' strengths. We emphasize that the aversion to relying heavily on the number of goals scored is unique to our model. This effect is further reflected upon in the case studies in Section 4.3. Moreover, due to the availability of more covariate information, it is observed that in the latter stages of the match, our model provides

uniformly better predictions than the rest.

A look at the plot for Brier score tells a similar story. We emphasize that this scoring rule incorporates the predicted probability of the outcome for model evaluation. Generally, lower values are recorded by our method, but the scores are still comparable until the end of the first half, following which the Bayesian approach stands out significantly from its competitors until the end of the match. Akin to the earlier argument, we believe that the proposed model’s performance is considerably better during the second half as it utilizes more available data in a better way. A peculiar observation can be made about the behaviour of GLM. Even though its F1-score does not display major irregularities, a significant spike in the Brier score is observed during the latter stages, likely due to its over-dependence on the number of goals. Overall, we can conclude that our model registers greater accuracy with more certainty than the contenders.

In an attempt to understand the performances in a better way, we next report the sensitivity and specificity for each model captured at 15 minute intervals in Table 2. The primary metric of interest here is sensitivity and how it interplays among various match outcomes for every model. Linear SVM outperforms all in terms of win sensitivity but has a significantly lower loss sensitivity. Furthermore, draws are not predicted by L-SVM since it performs binary classification. A behavioural pattern similar to L-SVM is observed for some of the other benchmark models. For instance, GLM lacks predictive power when it comes to draws; while RF demonstrates an inherent bias towards predicting an outcome of decisive nature. Consequently, the latter method exhibits the highest win sensitivity but the lowest draw sensitivity. Moreover, both R-SVM and GLM have decreasing win recall over time, making them less reliable at later stages of a match. Notably, all benchmark models underperform in terms of draw sensitivity, a result which gets reflected in Figure 4 as well. It is important to note that for a model to be practically usable, achieving a balance in sensitivity across all outcome categories is crucial. Compared to the other models, it is evident that the proposed approach maintains this balance most effectively. A striking observation is that this balance increases over time, with both the metrics improving along the course of a match, which is not visible in any other method.

Table 2: Sensitivity and specificity captured at specific time points, corresponding to different outcome categories, for the proposed Bayesian model (PBM), generalized linear model (GLM), linear SVM (L-SVM), SVM with radial basis functions (R-SVM) and the random forest (RF).

Outcome	Minute	Sensitivity					Specificity				
		PBM	GLM	L-SVM	R-SVM	RF	PBM	GLM	L-SVM	R-SVM	RF
Win	15th	0.603	0.849	0.897	0.732	0.720	0.766	0.500	0.354	0.614	0.551
	30th	0.589	0.746	0.925	0.651	0.801	0.778	0.607	0.392	0.690	0.551
	45th	0.637	0.712	0.931	0.646	0.877	0.823	0.696	0.399	0.766	0.550
	60th	0.678	0.719	0.904	0.630	0.863	0.873	0.728	0.405	0.772	0.551
	75th	0.781	0.733	0.897	0.603	0.883	0.918	0.722	0.411	0.804	0.538
Draw	15th	0.555	-	-	0.206	0.159	0.545	-	-	0.830	0.896
	30th	0.476	0.095	-	0.317	0.111	0.593	0.883	-	0.772	0.934
	45th	0.555	0.270	-	0.429	0.063	0.643	0.804	-	0.747	0.950
	60th	0.651	0.254	-	0.428	0.111	0.680	0.805	-	0.722	0.959
	75th	0.635	0.270	-	0.444	0.143	0.743	0.805	-	0.668	0.959
Loss	15th	0.253	0.547	0.421	0.453	0.505	0.952	0.765	0.852	0.813	0.785
	30th	0.379	0.547	0.474	0.473	0.558	0.909	0.775	0.871	0.809	0.809
	45th	0.442	0.537	0.453	0.495	0.558	0.904	0.823	0.856	0.818	0.828
	60th	0.516	0.631	0.453	0.442	0.568	0.914	0.842	0.833	0.809	0.828
	75th	0.632	0.611	0.463	0.484	0.547	0.928	0.842	0.828	0.852	0.852

To gain further insights into the model forecasts, we capture the Brier scores after every 15

minute intervals with respect to the difference in goals (home minus away) observed for the matches at the end of these intervals. We observe that the lowest Brier scores aka best accuracy were obtained for matches with a goal difference bigger than one. Naturally, an instance of zero goal difference is the most difficult to predict for. It is evident from the decreasing Brier scores that the model performance improves over the course of a match. An interesting observation is the slight increase in Brier scores from the 15th minute to the 75th for unit goal difference categories. This is likely to be the cumulative effect of time-varying covariates other than goals. This might also be an indication that goals scored in the early stages of a match provide better predictability than accounted for. It is possible that a different covariance structure than the one used in the model might be more effective in capturing this effect.

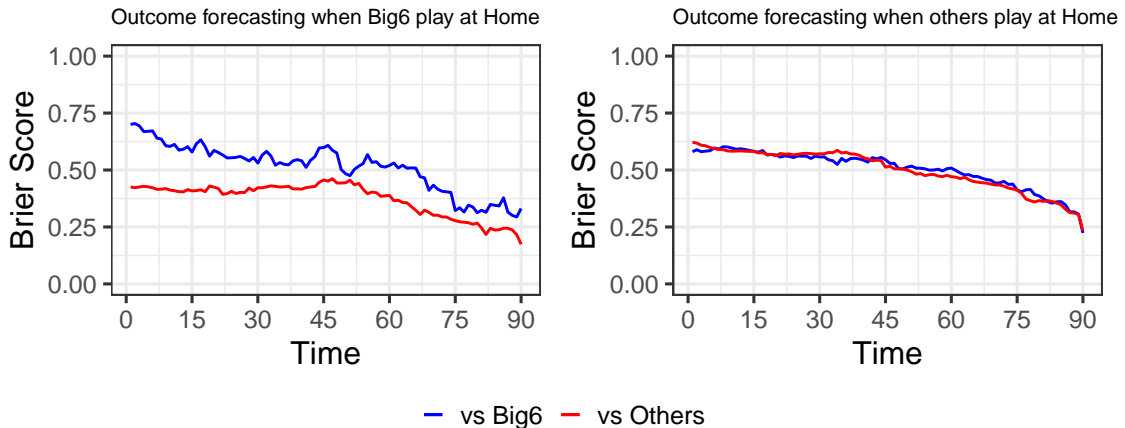
Table 3: Brier scores (with respect to goal difference) corresponding to the final outcomes based on predictions at specific time points.

Minute	Goal difference at the time of prediction						
	≤ -3	-2	-1	0	1	2	≥ 3
15th	-	0.034	0.237	1.36	0.247	0.015	-
30th	-	0.056	0.414	0.921	0.440	0.058	0.011
45th	0.004	0.135	0.411	0.768	0.439	0.111	0.033
60th	0.051	0.100	0.395	0.754	0.388	0.138	0.075
75th	0.055	0.130	0.496	0.627	0.341	0.148	0.102

4.2. Robustness of the analysis

To further support our previous results, we conduct a few robustness checks. One way to evaluate the heterogeneity in model forecasts is to compare the forecast accuracy based on the strength of a team. In that aspect, we split the set of EPL teams into two groups – ‘Big 6’ (consisting of Arsenal, Chelsea, Liverpool, Manchester City, Manchester United, Tottenham) and the rest. These ‘Big 6’ teams are known for their large payrolls and fan-bases across the world. Our objective is to identify if the predictive accuracy is different for these teams, as compared to the conventionally weaker teams. Figure 5 below shows the Brier scores for the matches corresponding to these two groups, split further according to who they played against.

Figure 5: Brier scores as a function of time when the teams are classified into two categories, namely ‘Big 6’ and others

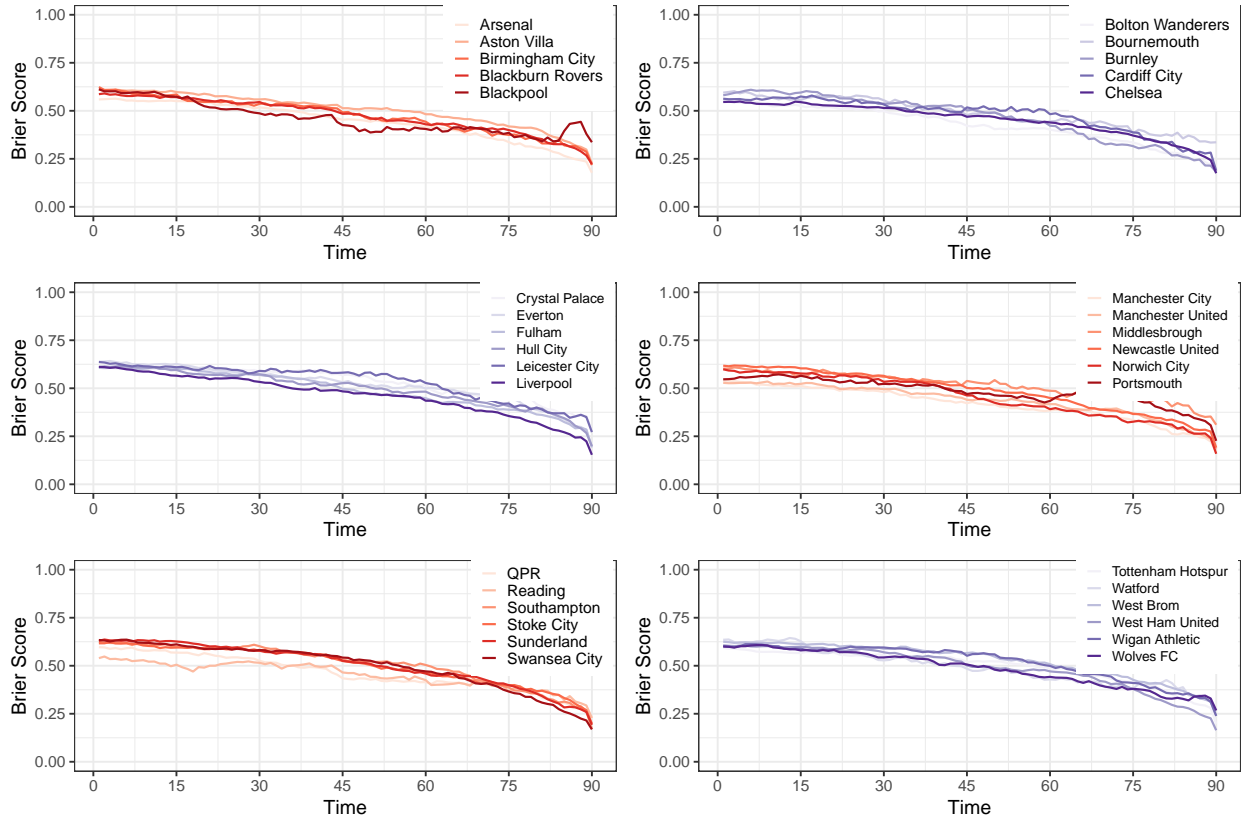


Home advantage is a commonly observed phenomenon in soccer (Staufenbiel et al. (2015)). The ‘Big 6’ clubs are the largest and most successful clubs in English soccer. One can hypothesize that

they have a stronger home advantage than smaller teams, providing better predictability of the outcomes when these clubs play at home against the others. This is reflected in the lower Brier scores in such situations, especially in case of correctly predicting the wins of the ‘Big 6’ (in fact, we observed that the brier scores of the predictions made at half-time corresponding to win is nearly four times better than the same corresponding to draw or loss of these teams). Interestingly, when two of the ‘Big 6’ clubs play against each other, the average Brier score is found to be on the higher end during the first half, suggesting that the matches are unpredictable in these circumstances. As discussed in Section 4.1 before, team strength plays a crucial role in forecasting match outcomes in the early stages. Due to the strength being comparable for both teams, the above phenomenon can be observed. For other teams playing at home, our model is robust in its accuracy irrespective of the opponents. Additionally, games involving the ‘Big 6’ clubs against other teams at their home stadiums also show significantly better forecasting accuracy in general. This suggests that home advantage is an effective performance indicator for the bigger teams, but not as much for smaller market teams. After the 60th minute mark in these matches, we note that the metric is similar for all cases, thereby establishing the proposed model’s robustness in this aspect.

Next, we look at a different perspective where we predict the outcomes of the matches for a certain team, by training the model on the dataset excluding all matches of that team. This exercise is performed for each of the 34 teams in the dataset, and we compute the average Brier scores, which are reported in Figure 6.

Figure 6: Brier scores computed after excluding a team’s games from the training data. The test data comprises entirely of matches played by the excluded team.



It must be noted that the teams have different extents of variations in their year-end positions.

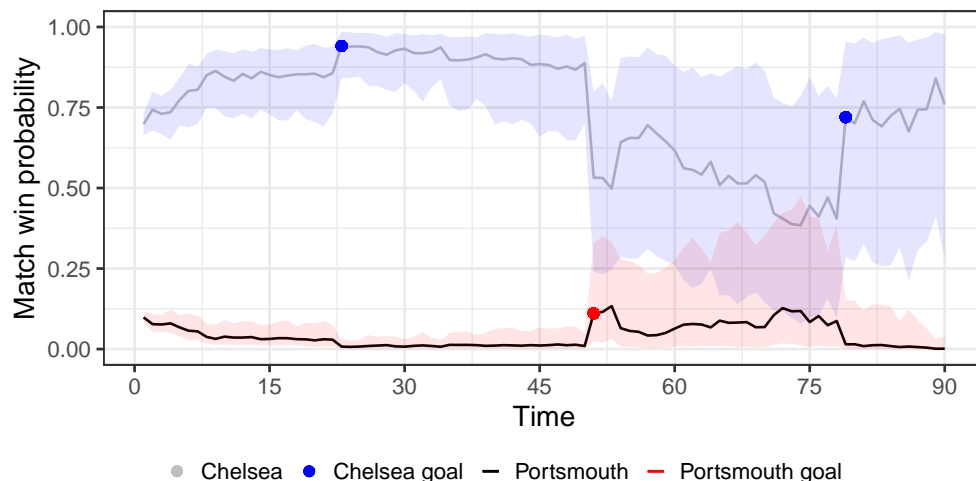
For example, Manchester United is typically at the top table in all seasons, whereas Aston Villa is a team with an average year-end position of 13 with the highest and lowest positions being 7 and 17. Some teams fluctuate in and out of the league on a yearly basis due to relegation as well. The idea is to identify how well the model performs when forecasting for all these teams, even without using its data in the training set. We believe that a robust model should be able to forecast the results of across all scenarios. The plots in Figure 6 indicate that the model registers comparable accuracy in predicting the outcomes of games for all teams. This suggests that the model effectively captures the key features of the teams and their performances, even in the presence of notable fluctuations from one year to the next. This is a significant result, as it implies that the model may be able to make reliable predictions for a wide range of teams and situations. A natural extension from a modelling perspective is the use of team-level fixed effects instead of the home-away classification we have currently used in the model. However, we observed that this gives rise to overfitting and consequently inferior results. Therefore, one will not be able to capture the effect of home advantage in such a model.

4.3. Two specific examples

As a last piece of discussion in this section, we look at the performance of the proposed approach for two specific games of different flavours. For each game, the model is trained on the data of all other matches before it, and the in-game forecasting is done on a minute-by-minute basis.

The first game we discuss was played between Chelsea and Portsmouth, at the home ground of the former. It was a back and forth match with a goal towards the end, giving Chelsea the win. This case study demonstrates the ability of our model to accurately forecast the outcome by considering a variety of factors, including team ratings and time-varying covariates to understand the flow of the game. The within-game forecasting of the win probabilities and the 95% credible intervals for the two teams is displayed in Figure 7.

Figure 7: Minute-by-minute forecast of the win probability for both teams during the match between Chelsea and Portsmouth.

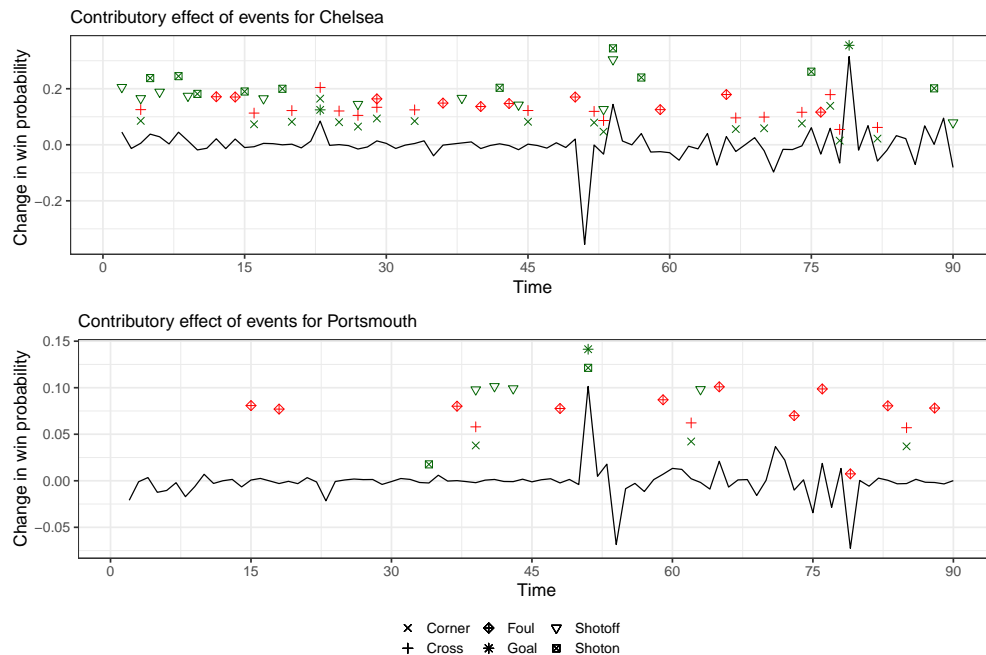


One can note that the final scoreline read 2-1 with the winning goal in the match being scored in the 79th minute. Chelsea was the heavy favourite at the start of the match, both due to being a better team (with respect to team rating) and for playing at home. The predicted win probability at the start was around 75%, and it improved for a few minutes in the beginning. A goal in the 21st minute consolidates our forecast. Later, Portsmouth equalled the scoreline just after the start

of the second half. It is indeed worth mentioning that the goals scored by the two teams are found to have differential impacts on the probabilities, potentially due to the other covariates used in the model. For instance, regardless of the even scoreline after the 50th minute, our model still considers Chelsea to be the odds on favourites to win. Chelsea is predicted to emerge victorious with 50% chance for the remaining of the match, and the result indeed goes in their favour at the end. Clearly, this case study serves as a captivating example of the importance of considering team strength in predictive modelling and the ability of our model to incorporate several factors accurately. We find it imperative to point out that the credible intervals for both teams are narrow during the early stages of the match. These get wider as the match progresses, especially after Portsmouth equalizes. This correctly reflects the uncertainty in the outcome as the scoreline is 1-1. A mix of crosses, corners and shots on goal by Chelsea resulted in large credible intervals persisting until the end of the match.

In Figure 8, we further attempt to demonstrate the probabilistic effect of the time-varying events for the match. It is to be noted that the events plotted are only for those respective teams. It is quite interesting that events other than goals have a significant impact on the opponent's win probability but a comparatively muted effect on their own. For instance, a goal by Portsmouth contributes to the sudden drop in the win probability for Chelsea, but only a marginal increase in the same for Portsmouth. On the other hand, a host of shots-on and shots-off-goals in the initial minutes enhances the win probability for Chelsea. This can have profound managerial implications. By observing the effects these time-varying covariates have on the respective team's win probabilities, the managers can make in-game adjustments to counter the same.

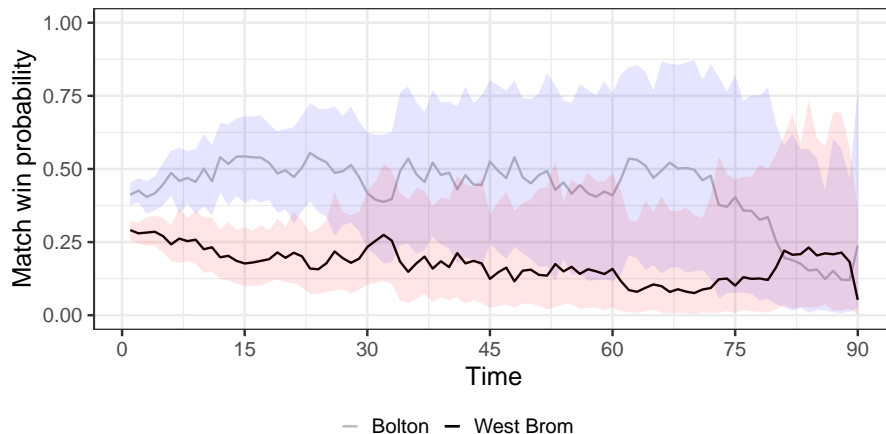
Figure 8: Contributory effect of all time-varying events on the probability of winning for Chelsea v/s Portsmouth. The first figure contains exclusively Chelsea events. The second figure contains exclusively Portsmouth events.



For completeness, we now dive deep into another case where our proposed model did not yield great results. The game was played between Bolton Wanderers and West Bromwich Albion, both teams typically being at the bottom of the league table. The score at the end of the 90th minute is 0-0, a draw. As we can see in Figures 9 and 10, at the start of the game, our model gives a slight

advantage to Bolton since it is playing at home. Then, a few shots-on-goal by Bolton pushed their win probability past 50%. This pattern persists for most parts of the match, with Bolton being considered the favourites until the 70th minute. It is only after this that the most likely predicted outcome is a draw. As we illustrate in Figure 5, home advantage is not an effective performance indicator for smaller market teams, which contributes to the inaccurate forecast for the majority of the game. This case further reflects the results obtained in Table 3 regarding the difficulty in forecasting games with an equal scoreline.

Figure 9: Minute-by-minute forecast of the win probability for both teams during the match between Bolton Wanderers vs West Bromwich Albion

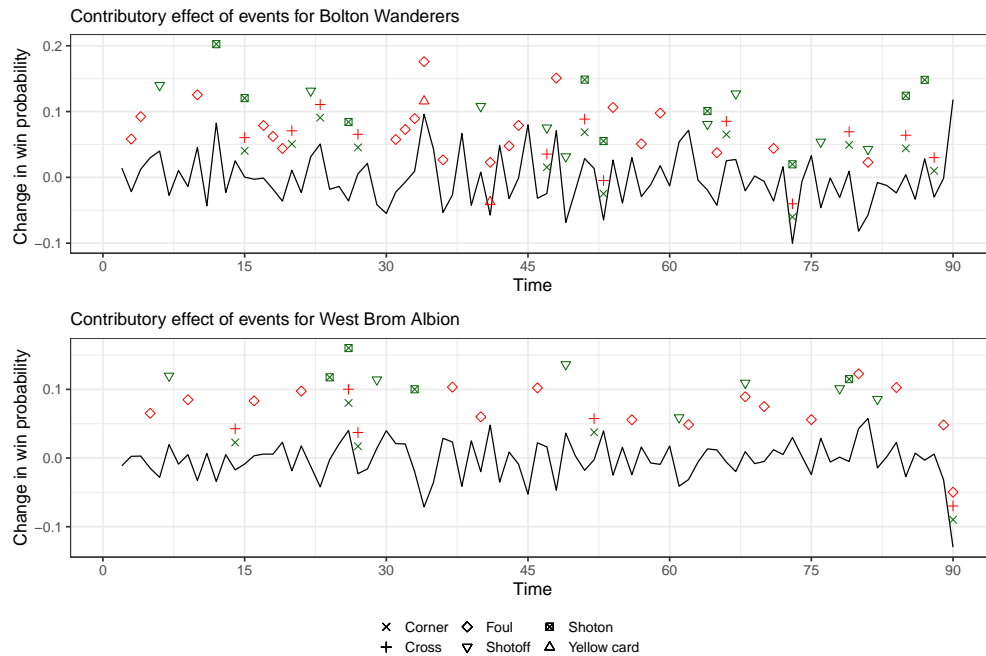


We also notice that the credible intervals get wider as the match progresses, thereby reflecting the uncertainty in the forecast. We believe that this is primarily due to the teams being equally matched. Further, the absence of goals gives us a better insight into the effect of other time-varying covariates. We can observe that shots-on-goal and fouls have a substantial effect on the probabilities in the absence of goals. A peculiar observation is the sudden rise in win probability for Bolton in the dying moments of the match. As we can identify in the bottom panel of Figure 10, a sequence of fouls and a cross by West Bromwich players towards the end may have provided this edge to the home team.

5. Conclusion

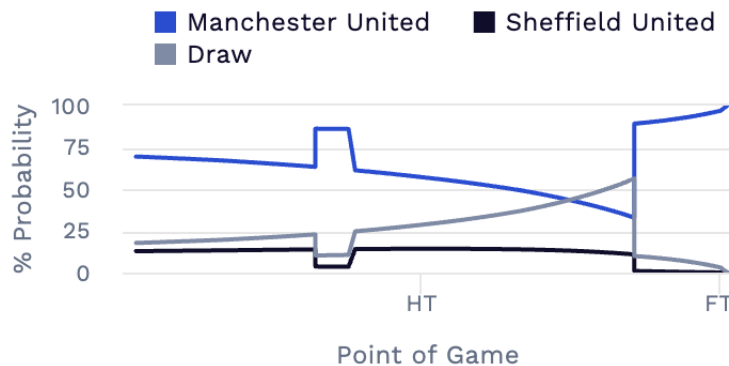
In this article, we developed a Bayesian latent variable model for analyzing and forecasting soccer match outcomes in real-time. The modelling approach, to the best of our knowledge, is the first in academic literature to furnish real-time predictions in soccer via a Bayesian technique. As mentioned earlier, although several websites offer in-game probabilities for various matches, the lack of clarity and documentation behind their methodologies affect their practical utility from a managerial standpoint. Another advantage of the proposed model over these websites is its ability to quantify the uncertainty in the probabilistic forecast at each time point. As an illustration, we present Figure 11 (sourced from [dimers.com](https://www.dimers.com)), which depicts the in-game probabilistic forecast during a recent EPL match between Manchester United and Sheffield, with the former winning by a scoreline of 2-1. A notable contrast between this and Figures 7 or 9 is the absence of credible intervals for the predicted probabilities. Moreover, noting the three jumps in the predicted probabilities in this case, it is apparent that the outputs are mainly affected by the goals scored; while our method is able to appropriately capture the effects of all types of significant events. Thus, we

Figure 10: Contributory effect of all time-varying events on the probability of winning for Bolton Wanderers vs West Bromwich Albion. The first figure contains exclusively Bolton events. The second figure contains exclusively West Brom events.



strongly believe that the proposed method, by virtue of its transparency and flexibility, can assist managers in formulating statistically sound strategies during a match. These aspects are crucial for any sort of decision making, be it from a spectator's perspective or that of a manager's.

Figure 11: Example of minute-by-minute forecast for a particular match between Manchester United and Sheffield United, sourced from dimmers.com (see <https://shorturl.at/fyILO>).



In our application, we used the data from EPL matches across eight seasons, with minute-by-minute data of various events in the matches. The main variable of interest for each game was the outcome, modelled as an ordered multivariate random variable with three categories of response. We used a latent variable and two cut-offs on it as a proxy for the final outcome and modelled this latent variable using a linear functional form. The functional form involved time-variant and time-invariant covariates as well as their corresponding coefficient, along with a random error term. In our computations, we use a Gaussian prior for the random errors, the coefficient as well as the

cut-offs, and a Dirichlet prior for the joint probabilities based on the cut-offs, to compute the step-wise conditional distributions as well as the posterior predictive distribution. It is critical to note that the Bayesian computations in our method are considerably fast due to the specific choices of the priors. However, an interesting future direction will be to relax the Gaussian assumptions and see how prior design may affect convergence in similar complex multinomial problems.

Coming to the performance in terms of forecast accuracy, we notice that our model provides insights into the effect of various events, such as corners, crosses and cards in soccer, on the outcome. Through various evaluation criteria, we find the model to be highly effective across different scenarios. We have seen that the results are robust for conventional big teams, as well as for teams with inconsistent performances. Furthermore, taking two specific examples, we demonstrate the effectiveness of the model in predicting the outcome well beforehand, a decisive goal is actually scored. Keeping that in mind, we strongly believe that the proposed methodology can be an extremely useful tool to maintain audience engagement in broadcasting soccer matches, or for the betting markets in real-time. A closely related area of research would be to understand the improvement in the win probability by potential substitutions. It can be achieved by including the player-level information in the model. Based on the within-match statistics and their deviation from the historical performances of the players, one may continuously compute the win-probability for different choices of replacements, which would in turn facilitate the coaches to devise a statistically sound substitution strategy during a match.

Continuing along the line, we would like to end the paper with a short account of other possible future extensions of this study. We note that the data did not include detailed information on the passing sequence, ball possessions, or substitutions. These aspects are expected to have a strong impact on the outcome of the game, but unfortunately, due to the lack of data, we could not add them in the current analysis. Another aspect we lose out on is team cohesion based on the interaction between individual players both in the same and opponent teams, which is also impacted by substitutions. Naturally, future implementation of the model can include these events in the framework and assess their impact on the match outcomes. Another possible improvement in the prediction accuracy can be to incorporate dependencies among different pairs of teams with respect to their previous outcomes. Similarly, one may consider the effects of the events from the home and away teams to be dependent, by relaxing the assumption in eq. (2.6). Utilizing appropriate structural forms in this context, $\Sigma_k^{(t)}$ in eq. (2.5) can be readjusted and then, it would be possible to develop a Bayesian algorithm in an identical fashion. Albeit it is expected to increase the computational burden for the Markov chains to converge, predictive ability may improve and such an idea is worth exploring in future works. Furthermore, in the current work, we do not consider self-exciting events and have considered different types of events to be independent of each other. It is not the ideal scenario, as given the occurrence of one kind of event, some other events may have higher chances of occurring. Addressing these different types of correlations in the model offers excellent scopes of future works to our method.

On a related note, another intriguing future direction is to assess the aspect of within-game momentum in soccer (see [Gauriot and Page, 2018](#); [Ötting et al., 2021](#), for some relevant discussions). In order to do this, one can extend the proposed model to a time series setting, where the response is a multivariate observation indicating the occurrence of different types of events and the regressor set includes the information of the past events in the game. This framework will allow one to estimate the momentum effect of different types of events and can subsequently lead to improved forecasting capability of the final outcome as well. Last but not the least, we recall that the assumption of Gaussian priors on the coefficients associated with the real-time events might be construed to be restrictive. We plan to address this issue in future works, by considering a more general structure,

possibly along with a semi-parametric approach to model the outcome of the game.

One must recognize that the proposed methodology can be suitably modified and adapted to various interesting research problems as well. In other sports, such as basketball or cricket, a direct extension is natural and can offer an effective tool of within-match forecasting. Not only does this have the potential to be utilized in betting markets, but it can also be used to figure out an appropriate substitution strategy in basketball where the rules allow rolling substitution. In fact, for other problems involving ordered categorical response variables, one can develop a similar technique. For instance, in sports broadcasting industry, such an algorithm can help us determine the effects of various events in retaining a customer. One may also develop a similar Bayesian method to understand the impact of different types of events and announcements on how a firm performs in an economy.

Data availability statement

The data used in this study are obtained from the publicly available European Soccer Database in Kaggle (link: <https://www.kaggle.com/datasets/hugomathien/soccer>). All codes for data extraction and for running the main algorithms are made publicly available on a GitHub repository maintained by the first author (link : <https://bit.ly/30o0SN8>).

References

- Albert, J.H., Chib, S., 1993. Bayesian analysis of binary and polychotomous response data. *Journal of the American statistical Association* 88, 669–679.
- Angelini, G., De Angelis, L., 2017. PARX model for football match predictions. *Journal of Forecasting* 36, 795–807.
- Ashimolowo, T., 2018. An Econometric Analysis of the Relationship Between Corner Kick Numbers and Football Outcomes. Ph.D. thesis. The University of North Carolina at Charlotte.
- Baboota, R., Kaur, H., 2019. Predictive analysis and modelling football results using machine learning approach for English Premier League. *International Journal of Forecasting* 35, 741–755.
- Boyko, R.H., Boyko, A.R., Boyko, M.G., 2007. Referee bias contributes to home advantage in English Premiership football. *Journal of sports sciences* 25, 1185–1194.
- Brier, G.W., et al., 1950. Verification of forecasts expressed in terms of probability. *Monthly weather review* 78, 1–3.
- Castellano, J., Casamichana, D., Lago, C., 2012. The use of match statistics that discriminate between successful and unsuccessful soccer teams. *Journal of human kinetics* 31, 139.
- Červený, J., van Ours, J.C., van Tuijl, M.A., 2018. Effects of a red card on goal-scoring in World Cup football matches. *Empirical Economics* 55, 883–903.
- Chambers, J.M., Hastie, T.J., 2017. Statistical models, in: *Statistical models* in S. Routledge, pp. 13–44.
- Clarke, S.R., Norman, J.M., 1995. Home ground advantage of individual clubs in English soccer. *Journal of the Royal Statistical Society: Series D (The Statistician)* 44, 509–521.

- Cowles, M.K., Carlin, B.P., 1996. Markov chain monte carlo convergence diagnostics: a comparative review. *Journal of the American Statistical Association* 91, 883–904.
- Crowder, M., Dixon, M., Ledford, A., Robinson, M., 2002. Dynamic modelling and prediction of English Football League matches for betting. *Journal of the Royal Statistical Society: Series D (The Statistician)* 51, 157–168.
- Czado, C., Gneiting, T., Held, L., 2009. Predictive model assessment for count data. *Biometrics* 65, 1254–1261.
- Dechi, B.O., 2019. Bayesian Analysis of Ordinal Outcomes Through Latent Variable Approach. The University of Texas at El Paso.
- Easton, S., Uylangco, K., 2010. Forecasting outcomes in tennis matches using within-match betting markets. *International Journal of Forecasting* 26, 564–575.
- Easton, S.A., Uylangco, K., 2006. An examination of in-play sports betting using one-day cricket matches. Available at SSRN 948013 .
- European-Gaming, 2022. Global Sports Betting Market to Grow At 10.3% CAGR until 2032; Football to Be the Most Sought After Sport: Fact.MR Report. URL: bit.ly/3j2jpTA.
- Forrest, D., Simmons, R., 2000. Forecasting sport: the behaviour and performance of football tipsters. *International journal of Forecasting* 16, 317–331.
- Gauriot, R., Page, L., 2018. Psychological momentum in contests: The case of scoring before half-time in football. *Journal of Economic Behavior & Organization* 149, 137–168.
- Gelfand, A.E., 2000. Gibbs sampling. *Journal of the American statistical Association* 95, 1300–1304.
- Geman, S., Geman, D., 1984. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on pattern analysis and machine intelligence* , 721–741.
- Geweke, J.F., et al., 1991. Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments. Technical Report. Federal Reserve Bank of Minneapolis.
- Goddard, J., 2005. Regression models for forecasting goals and match results in association football. *International Journal of forecasting* 21, 331–340.
- Gómez, M.Á., Mitrotasios, M., Armatas, V., Lago-Peñas, C., 2018. Analysis of playing styles according to team quality and match location in Greek professional soccer. *International Journal of Performance Analysis in Sport* 18, 986–997.
- González-Rodenas, J., Aranda-Malavés, R., Tudela-Desantes, A., Calabuig Moreno, F., Casal, C.A., Aranda, R., 2019. Effect of match location, team ranking, match status and tactical dimensions on the offensive performance in Spanish ‘La Liga’ soccer matches. *Frontiers in psychology* 10, 2089.
- Greene, W.H., 2003. *Econometric analysis*. Pearson Education India.
- Groll, A., Kneib, T., Mayr, A., Schauburger, G., 2018. On the dependency of soccer scores—a sparse bivariate poisson model for the uefa european football championship 2016. *Journal of Quantitative Analysis in Sports* 14, 65–79.

- Higgs, M.D., Hoeting, J.A., 2010. A clipped latent variable model for spatially correlated ordered categorical data. *Computational Statistics & Data Analysis* 54, 1999–2011.
- Hvattum, L.M., Arntzen, H., 2010. Using ELO ratings for match result prediction in association football. *International Journal of forecasting* 26, 460–470.
- Imai, K., Van Dyk, D.A., 2005. A bayesian analysis of the multinomial probit model using marginal data augmentation. *Journal of econometrics* 124, 311–334.
- Klaassen, F.J., Magnus, J.R., 2003. Forecasting the winner of a tennis match. *European Journal of Operational Research* 148, 257–267.
- Kolassa, S., 2016. Evaluating predictive count data distributions in retail sales forecasting. *International Journal of Forecasting* 32, 788–803.
- Koning, R.H., 2000. Balance in competition in Dutch soccer. *Journal of the Royal Statistical Society: Series D (The Statistician)* 49, 419–431.
- Koopman, S.J., Lit, R., 2019. Forecasting football match results in national league competitions using score-driven time series models. *International Journal of Forecasting* 35, 797–809.
- Lepschy, H., Wäsche, H., Woll, A., 2020. Success factors in football: an analysis of the German Bundesliga. *International Journal of Performance Analysis in Sport* 20, 150–164.
- Ley, C., Wiele, T.V.d., Eetvelde, H.V., 2019. Ranking soccer teams on the basis of their current strength: A comparison of maximum likelihood approaches. *Statistical Modelling* 19, 55–73.
- Liddell, T.M., Kruschke, J.K., 2018. Analyzing ordinal data with metric models: What could possibly go wrong? *Journal of Experimental Social Psychology* 79, 328–348.
- Liti, C., Piccialli, V., Sciandrone, M., 2017. Predicting soccer match outcome using machine learning algorithms, in: *Proceedings of MathSport International 2017 Conference*.
- Liu, H., Gomez, M.Á., Lago-Peñas, C., Sampaio, J., 2015. Match statistics related to winning in the group stage of 2014 Brazil FIFA World Cup. *Journal of sports sciences* 33, 1205–1213.
- Maher, M.J., 1982. Modelling association football scores. *Statistica Neerlandica* 36, 109–118.
- McCulloch, R.E., Polson, N.G., Rossi, P.E., 2000. A Bayesian analysis of the multinomial probit model with fully identified parameters. *Journal of econometrics* 99, 173–193.
- McHale, I., Scarf, P., 2007. Modelling soccer matches using bivariate discrete distributions with general dependence structure. *Statistica Neerlandica* 61, 432–445.
- Mendes-Neves, T., Mendes-Moreira, J., 2020. Comparing state-of-the-art neural network ensemble methods in soccer predictions, in: *International Symposium on Methodologies for Intelligent Systems*, Springer. pp. 139–149.
- Murphy, A.H., 1973. A new vector partition of the probability score. *Journal of Applied Meteorology and Climatology* 12, 595–600.
- Ötting, M., Langrock, R., Maruotti, A., 2021. A copula-based multivariate hidden Markov model for modelling momentum in football. *AStA Advances in Statistical Analysis* , 1–19.

- Rocha-Lima, E.M., Tertuliano, I.W., Fischer, C.N., 2021. The influence of crosses, shots, corner kicks and defensive movements in the results of Premier League matches. *Research, Society and Development* 10, e477101624072–e477101624072.
- Rue, H., Salvesen, O., 2000. Prediction and retrospective analysis of soccer matches in a league. *Journal of the Royal Statistical Society: Series D (The Statistician)* 49, 399–418.
- Song, K., Gao, Y., Shi, J., 2020. Making real-time predictions for NBA basketball games by combining the historical data and bookmaker’s betting line. *Physica A: Statistical Mechanics and its Applications* 547, 124411.
- Staufenbiel, K., Lobinger, B., Strauss, B., 2015. Home advantage in soccer—A matter of expectations, goal setting and tactical decisions of coaches? *Journal of sports sciences* 33, 1932–1941.
- Stern, H.S., 2005. A Brownian motion model for the progress of sports scores, in: *Anthology of Statistics in Sports*. SIAM, pp. 257–263.
- Štrumbelj, E., 2014. On determining probability forecasts from betting odds. *International journal of forecasting* 30, 934–943.
- Štrumbelj, E., Šikonja, M.R., 2010. Online bookmakers’ odds as forecasts: The case of European soccer leagues. *International Journal of Forecasting* 26, 482–488.
- Van Rijsbergen, C., 1979. Information retrieval: theory and practice, in: *Proceedings of the joint IBM/University of Newcastle upon tyne seminar on data base systems*.
- Vecer, J., 2014. Crossing in soccer has a strong negative impact on scoring: Evidence from the English Premier League the German Bundesliga and the World Cup 2014. Available at SSRN 2225728 .
- Walvin, J., 2014. *The people’s game: the history of football revisited*. Random House.