



Flexible modelling of time-varying exposures and recurrent events to analyse training load effects in team sports injuries

Lore Zumeta-Olaskoaga^{1,2} , Andreas Bender^{3,4}  and Dae-Jin Lee⁵ 

¹BCAM—Basque Center for Applied Mathematics, Bilbao, Spain

²Departamento de Matemáticas, Universidad del País Vasco UPV/EHU, Leioa, Spain

³Statistical Consulting Unit StaBLab, Ludwig-Maximilians Universität München, Munich, Germany

⁴Munich Center for Machine Learning (MCML), LMU München, Munich, Germany

⁵School of Science and Technology, IE University, Madrid, Spain

Address for correspondence: Lore Zumeta-Olaskoaga, BCAM—Basque Center for Applied Mathematics, Bilbao 48009, Spain; Departamento de Matemáticas, Universidad del País Vasco UPV/EHU, Leioa 48940, Spain.

Email: lzumeta@bcmath.org

Abstract

We present a flexible modelling approach to analyse time-varying exposures and recurrent events in team sports injuries. The approach is based on the piece-wise exponential additive mixed model where the effects of past exposures (i.e. high-intensity training loads) may accumulate over time and present complex forms of association. In order to identify a relevant time window at which past exposures have an impact on the current risk, we propose a penalty approach. We conduct a simulation study to evaluate the performance of the proposed model, under different true weight functions and different levels of heterogeneity between recurrent events. Finally, we illustrate the approach with a case study application involving an elite male football team participating in the Spanish LaLiga competition. The cohort includes time-loss injuries and external training load variables tracked by Global Positioning System devices, during the seasons 2017–2018 and 2018–2019.

Keywords: cumulative effect, football injuries, recurrent events, sports analytics, survival analysis, time-varying exposures

1 Introduction

In recent years, research on sports injuries has received much attention from a wide range of areas, including the statistical modelling field, given the ever-increasing amount of data now being collected (Casals & Finch, 2017; Sainani et al., 2021). Injuries are common, and modelling and understanding their occurrence would assist in developing tailored prevention programs. In this regard, it has been claimed that analysing changes in training load exposures over time is crucial when studying sports injury aetiology, as well as analysing recurrent injury, subsequent injury, or injury exacerbation (Nielsen et al., 2020). Athletes are repeatedly exposed to high competition demands that, in turn, increase the strain on their bodies and exposure to the risk of injury. Injury prevention depends on the athlete's capacity to tolerate repeated exposures to injury risk.

The physical load from training and competition, hereafter termed as training load, is defined as 'the cumulative stress placed on an individual from multiple training sessions and games over a period of time' (Gabbett et al., 2014). Training load can thus be applied to the athlete over varying time periods and with varying magnitudes (Soligard et al., 2016). Nowadays, vast amounts of data are gathered (e.g. training and competition time, distance, speed, power output), primarily using

Received: January 19, 2024. Accepted: October 18, 2024

© The Royal Statistical Society 2024.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Global Positioning System (GPS) devices, to quantify training load. Indeed, the study of training load is key to developing effective training plan strategies that enhance athletes' performance while also reducing their risk of injury. The relationship between training load and injury, however, remains uncertain (Griffin et al., 2020; Windt et al., 2018).

Quantifying this relationship requires the development of an etiologically plausible time-varying exposure model, which estimates how previous training affects the injury risk. The effects of past exposures may cumulate over time and exhibit complex forms of association. Moreover, the model should account for the possibility that subsequent injuries may be associated within players. To take into account the latter, specifically dependencies induced by injury recurrence, as well as the intensity and duration of past exposures, we propose a piece-wise exponential additive mixed model (PAMM) (Bender et al., 2018) with weighted cumulative exposure (WCE) cumulative effects (Sylvestre & Abrahamowicz, 2009).

PAMMs are a semi-parametric extension of the piece-wise exponential model (M. Friedman, 1982; Holford, 1980; Laird & Olivier, 1981) that allow for penalized estimation of flexible survival models (cf. Argyropoulos & Unruh, 2015; Bender et al., 2018 for a thorough overview) with a wide range of covariate effects, such as nonlinear, time-varying effects, cumulative effects, and/or random effects. This framework has also been shown to support the estimation of cumulative effects of time-varying exposure histories. The WCE-type cumulative effect suggested by Sylvestre and Abrahamowicz (2009), a weighted sum of all past exposures over a relevant time window, is a common way to address this. The weight function assigns weights to past exposures based on the time elapsed since the exposure occurred, which, ideally, is determined according to the true underlying biological mechanism. They proposed to estimate the weight function using B-spline regression.

The PAMM methodology has been applied in many recent publications. Bender et al. (2019) explore complex exposure-lag-response associations and provide a general formulation of PAMMs that includes previous approaches for cumulative effects like the WCE model and the distributed lag nonlinear model (DLNM) (Gasparrini et al., 2017) as special cases. Ramjith et al. (2022) study the PAMM framework for recurrent events analysis and show that under the assumption of proportional hazards, PAMM and the shared frailty Cox model (McGilchrist & Aisbett, 1991) are equivalent. Danieli and Abrahamowicz (2019) and Li et al. (2022) introduce approaches to model cumulative effects of time-varying exposures with competing risks, via cause-specific hazards model and subdistribution hazards model for each competing event, respectively, by incorporating separate Cox WCE models with an event-specific weight function. Recently, in the field of sports medicine, Bache-Mathiesen et al. (2022) evaluated different methods to assess the cumulative effect of training load on the risk of injury in team sports and suggest the use of DLNM.

In this work, we extend the PAMM by incorporating WCE-type cumulative effects in the recurrent events setting combined with a method to identify a relevant time window in which past exposures have an effect. We further demonstrate the practical application of this model in the field of sports medicine.

The proposed modelling framework is detailed throughout Section 2, which first introduces the PAMM framework and then focuses on how we adapt it to flexibly model time-varying exposures and recurrent events in addition to how we penalize the model for identifying a relevant window. Section 3 illustrates a real-life application of this method to assess the cumulative effects of past training exposures on the hazard of subsequent injuries in a football team. Section 4 describes the simulation study carried out to evaluate the proposed models performance and examine their properties. The final section concludes the paper with a discussion.

2 Methods

PAMMs transform a survival task to a Poisson regression task by partitioning the follow-up period into a finite number of intervals and assuming that hazards are piece-wise constant in each of these intervals. The key idea—explained in detail in the following sections—is to flexibly model the baseline hazard and other time-varying effects using penalized splines. This approach alleviates the problem of an arbitrary selection of the cut-points, that partition the follow-up time, and thus avoids over- and underfitting as well as instability issues (Bender et al., 2018). In practice, one

can simply use a relatively large number of cut-points and use spline basis functions evaluated at e.g. κ_j , the right end of each interval, and penalize the wiggleness of the estimate via penalized splines (cf. [Bender, 2018](#), for an empirical discussion on the placing of cut-points).

We denote for player l , where $l = 1, \dots, L$ and L the total number of players, Y_{i_l} as the i th recurrent time, where $i = 1, \dots, n_l$ and n_l the maximum event number that individual l has been at risk for, and C_l as the censoring time. Then, $T_{i_l} = \min Y_{i_l}, C_l$ corresponds to each follow-up time and δ_{i_l} is a binary indicator for recurrent events (i.e. subsequent injuries) which is 0 if the player has not been injured and 1 if the i th injury time, Y_{i_l} , is observed. That is, $\delta_{i_l} = \mathbf{1}_{(T_{i_l}=Y_{i_l})}$, where $\mathbf{1}$ represents the indicator function. In the following, we will use t instead of t_{i_l} , for ease of notation.

Therefore, as we proceed with the modelling of injuries, the hazard rate of the i th injury (event) of the l th player, given the player's training exposure history $z_l(t) = \{z_l(t_z) : t_z \leq t\}$, is expressed as:

$$\begin{aligned} \lambda_{i_l}(t | z_l(t), b_l) &= \lambda_0(t) \exp(g(z_l(t), t_j) + b_l) \\ &= \exp(\beta_0 + f_0(t_j) + g(z_l(t), t_j) + b_l), \end{aligned} \tag{1}$$

for all $t \in (\kappa_{j-1}, \kappa_j]$, $t > 0$ and $t_j := \kappa_j$, where κ_j are the $J + 1$ cut points defining J intervals that partition the study follow-up time $(0, t_{\max}]$, specifically, $\kappa_0 = 0, \kappa_j = t_{\max}$ and $j = 1, \dots, J - 1$.

In equation (1), the expression $\beta_0 + f_0(t_j)$ denotes the log-baseline hazard, where $f_0(t_j)$ is expressed as a smooth term of the form $\sum_{m=1}^M \gamma_{0m} B_m(t_j)$. The term $g(z_l(t), t) = \int_{\tau(t)} h(t, t_z, z_l(t_z)) dt_z$ denotes the cumulative effect of $z_l(t)$ at time t , i.e. the past exposure effects of $z_l(t)$ cumulating over a relevant time-window $\tau(t)$, resulting in a sum of weighted effects. The dependence induced by subsequent injuries is accounted for by b_l , a Gaussian random effect (i.e. frailty) associated with player l , which acts as a random intercept term for the l th player, i.e. $b_l \sim N(0, \sigma_b^2)$.

For a WCE-type effect, in equation (1), we consider time-varying exposure effects weighted by latency $t - t_z$ and linear in $z(t_z)$. That is, the contribution of covariate $z(t)$ observed at time t_z with value $z(t_z)$, is defined by $h(t, t_z, z(t_z)) := h(t - t_z)z(t_z)$, and called partial effect. Thus, the cumulative effect $g(z(t), t)$ at follow-up time t is the integral of these partial effects over exposure times t_z contained within the so-called lag-lead window, $\tau(t)$, which controls how many observations of z contribute to the cumulative effect at time t (the minimal requirement is that $t_z \leq t$).

Let $z(t) = \{z(t_z) : t_z \leq t\} = \{z(t_{z,1}), \dots, z(t_{z,Q})\}$ be the set of all registered exposure variables up to time t . Then, $g(z(t), t)$ is estimated with penalized splines (e.g. by using P-splines [Eilers & Marx, 1996, 2021](#) that penalize the differences of neighbouring basis coefficients) and with quadrature weights $\Delta_q = t_{z,q} - t_{z,q+1}$ (and $t_{z,0} = t$), the time difference between two consecutive exposure measurements, for numerical integration, as follows:

$$\int_{\tau(t)} h(t - t_z)z(t_z) dt_z \approx \sum_{q=1}^Q \tilde{\Delta}_q \tilde{h}(t - t_{z,q}) = \sum_{q=1}^Q \tilde{\Delta}_q \sum_{m=1}^M \gamma_m B_m(t - t_{z,q}) \tag{2}$$

with $\tilde{\Delta}_q = z(t_z)(t_{z,q} - t_{z,q+1})$ if $t_{z,q} \in \tau(t)$ and 0 otherwise; $B_m(\cdot)$ B-spline basis functions and γ_m the associated spline coefficients.

2.1 Penalization of the weight function

One of the challenging issues is to determine a relevant time window $\tau(t)$. Without solid prior knowledge, it can be defined as $\tau(t) = \{t_z : t_z \geq t\}$, so all past exposures, collected before actual time t , contribute to the cumulative effect $g(z(t), t)$. Yet, it is plausible that the effects of exposure variables may not be everlasting. As time passes, the effect of exposures recorded long ago may smoothly decrease to zero and eventually disappear. The exact length of the window, however, is usually unknown.

Subsequently, adapting the method by [Obermeier et al. \(2015\)](#) to the PAMM framework, we present two approaches to penalize the weight function, allowing it to transition smoothly to zero at the right end of the support interval: (i) a constrained-effect approach and (ii) a ridge penalty approach.

The smoothing parameter λ_d penalizes large differences in adjacent basis coefficients, while the regularization parameter λ_r shrinks the last basis coefficient.

Therefore, the coefficients can be estimated via penalized iteratively reweighted least squares P-IRLS (Marx & Eilers, 1998; Wood, 2017). The P-IRLS consists of iteratively updating the coefficient estimates until convergence is reached using numerical optimization methods of the restricted maximum likelihood. This is implemented in the `gam` function from `mgcv` (Wood, 2017) R package.

2.3 Software specification

The analyses of the case study and the simulation study are coded in R version 4.2.2, on a 64-bit Unix platform (x86_64 linux-gnu) computer, as well as on a high-performance cluster system. The package `msm` 1.6.9 was used to draw piece-wise exponential survival times, `pamtools` 0.5.8 and `mgcv` 1.8–41 to fit the models, `batchtools` 0.9.16 to structure, write down and submit the simulation experiment in a convenient and reproducible fashion and the package `injurytools` 1.0.1 to structure and explore the football injury data set. The code to reproduce these analyses is available at <https://github.com/lzumeta/flex-mod-training-loads-recu-injuries>.

3 Application: football injury data

3.1 Data

We apply the proposed model to observational injury data from an elite male football team that competed in LaLiga during the 2017–2018 and 2018–2019 seasons. A total of $L = 36$ players were monitored to assess their performance and health status through external training load variables (e.g. training and competition time, distance covered, speed, and heart rate, see Soligard et al., 2016) collected using tracking devices during each match and training session. These variables measure the physical exertion to which the players were exposed. A total of 72 noncontact time-loss injuries occurred among 23 players (64%) and 15 (65%, 15/23) players were reinjured during the follow-up, see Figure 1.

Our focus is on understanding the relationship between external training load and time-loss injuries (Fuller et al., 2006), particularly how the cumulative stress from multiple training sessions and matches over time impacts a player's risk of sustaining a (subsequent) injury. We specifically focus on the variables average speed per session (*Speed*) and total distance covered (*Dist*) per session, as external training load variables that represent the intensity of each session. These metrics are chosen for their clear link to physical exertion and injury risk. Average speed reflects the session's sustained intensity, while total distance indicates overall workload. Both are crucial in evaluating the cumulative stress imparted on players, potentially leading to injuries when consistently high.

3.2 Why fit a flexible time-varying exposures model?

In practical settings, a common question that practitioners face is ‘how much training load is too much among players with different characteristics?’ (Nielsen et al., 2018; Soligard et al., 2016). To address this question, researchers have studied changes in training load as a primary exposure to injury. For example, metrics such as week-to-week changes and the acute:chronic workload ratio (ACWR) have been heavily relied upon. This latter measure compares short-term (acute) training load to longer-term (chronic) load to estimate injury risk (see Section A.2 of the online supplementary material for a more comprehensive explanation). Numerous recent studies have investigated the relationship between training load and injury risk (Bache-Mathiesen et al., 2021, 2022; Carey et al., 2018; Hulin et al., 2016; Windt & Gabbett, 2017).

Before describing our modelling approach, we may first ask why a flexible time-varying exposure model might be preferred over relying on predetermined measures for calculating changes in training load. First, predefined metrics such as week-to-week changes, ACWR and its variants, or cumulative load measures, tend to simplify the complex relationships between training load and injury risk and may miss important time-varying patterns or immediate changes in load. Second, the effects of long-term exposure to training load may manifest days or even weeks later, but this time period is neither fixed nor previously known. Therefore, a method that can identify this time window for each specific context is highly valuable. Third, flexible models can capture nonlinear patterns learning directly from the data. For these reasons, we adopt a flexible time-varying exposure

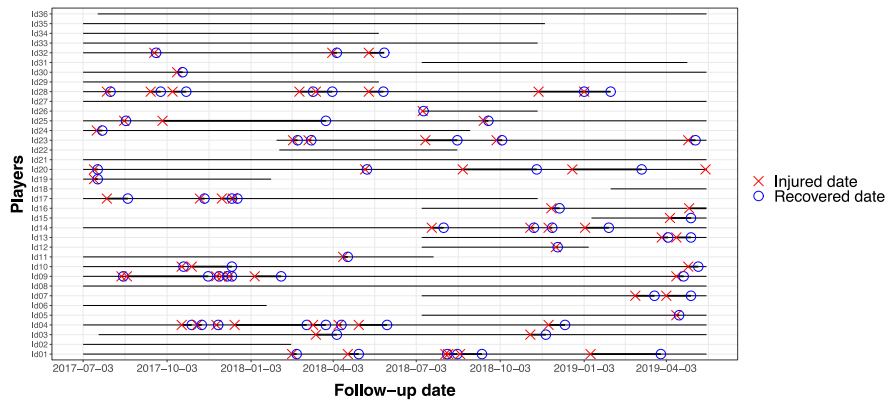


Figure 1. Timeline of the football players’ follow-up period together with the injuries they sustained. The red cross indicates the exact injury date, the blue circle the recovery date and the bold black line the duration of the noncontact time-loss injury.

model with penalization to simultaneously account for potentially nonlinear, time-varying, and cumulative effects while identifying relevant time windows.

3.3 Model specification

We consider that the unit of the follow-up time t , as well as of the exposure time t_z , to be the n th number of session (i.e. match and training sessions). The analysis time zero is defined as the first session (match or training session), from 7 July 2017, in which the player has taken part in the team. Players are followed until an injury occurs, or until they are transferred to another team, the end of the contract or the end of the study (18 May 2019), whichever occurs first. Players fully recovered from an injury are followed again until one of the previously described events occurs, see Figure 1. We consider that external training load is applied to the player over varying time periods and with varying magnitude by considering cumulative effects and we adjust for the type of session, whether training or match session, since it has been suggested as one of the primary risk factors (Bahr, 2003; Bahr et al., 2020); and account for subsequent injuries adding a random effect (Gaussian frailty). Thus, the log-hazard rate of player l of the fitted model is expressed as:

$$\log(\lambda(t | z_l(t), b_l, i)) = \beta_0 + f_0(t_j) + z_l^{\text{type session}}(t_j)\beta_1 + g_1(z_l^{\text{Speed}}(t), t) + g_2(z_l^{\text{Dist}}(t), t) + b_l, \quad (5)$$

$$\forall t \in (\kappa_{j-1}, \kappa_j], \quad t_j := \kappa_j \text{ and } b_l \sim N(0, \sigma_b),$$

where $\beta_0 + f_0(t_j)$ indicates the log-baseline hazard rate, $z_l^{\text{type session}}(t_j)$ the type of session undertaken by player l at t_j , g_1 and g_2 are nonlinear time-varying effects of the training load variables and b_l a Gaussian random intercept term associated to player l . The cumulative effects, g_1 and g_2 , are defined as $\int_{\tau_{\text{Speed}}(t)} b(t - t_z)z_l^{\text{Speed}}(t_z) dt_z$ and $\int_{\tau_{\text{Dist}}(t)} b(t - t_z)z_l^{\text{Dist}}(t_z) dt_z$, and each lag-lead window, $\tau_{\text{HRT}}(t)$ and $\tau_{\text{HSR}}(t)$, is chosen to be large enough to identify relevant past exposure effects by fitting a PAMM with a ridge penalization. Importantly, we add a minimum lag time of one session to minimize confounding by indication bias (Signorello et al., 2022), i.e. we exclude the current session to have an effect on the hazard of injury, $t > t_z$. The reason is that each session depends on the player’s physical condition and, presumably, sessions in which a player was injured had lower intensity in contrast to sessions in which he had no complaints. All smooth terms are estimated using P-Splines with second-order difference penalties.

3.4 Results

Table 1 shows the descriptive characteristics of the data, overall and by session type. The injury incidence was 4 injuries per 1,000 player-hours, and the injury burden was 90 days lost per 1,000 player-hours. The first calculates the rate at which new injury occurs (likelihood), whereas the second indicates how severe an injury is (consequences) (Bahr et al., 2020). Match injury

Table 1. Descriptive characteristics of football injury data: summary statistics related to injury and exposure variables overall and by session type

	Session type		
	Overall	Training	Match
Injury-related variables			
Injuries, n (%)	72	26 (36.1)	46 (63.9)
Days lost, n (%)	1,595	591 (37.1)	1,004 (62.9)
Total follow-up sessions, median (IQR)	220.5 (163–345)	198.5 (146–277)	37.5 (19–67)
Injury incidence, (95% CI)	4.07 (3.1–5)	1.47 (0.9–2)	2.6 (1.85–3.35)
Injury burden, (95% CI)	90.1 (85.6–94.5)	33.4 (30.4–36.1)	59.7 (53.2–60.2)
Exposure variables			
Average speed (m/s), median (IQR)	3.8 (3.24–4.72)	3.71 (3.17–4.25)	6.46 (5.96–6.88)
Total distance (m), median (IQR)	4,689 (3,586–6,122)	4,458 (3,517–5,525)	8,552 (5,138–10,022)

Note. Injury incidence and injury burden are reported per 1,000 player-hours. n = number; IQR = interquartile range; 95% CI = 95% confidence interval.

incidence and burden were higher than those during training, 2.6 vs. 1.5 injuries per 1,000 player-hours and 60 vs. 33 days lost per 1,000 player-hours, respectively.

For the PAMM with a ridge penalization modelling approach, a sufficiently large lag-lead window must be set. We selected this window based on a reasonable range suggested by domain experts. Consequently, the estimated cumulative effects are computed considering that all *Speed* and *Dist* values recorded in the last 10 sessions prior to t (i.e. approximately three weeks before)) could have an effect on the hazard of injury at time t , i.e. $\tau_1(t) = \tau_2(t) = \{t_z : t > t_z \wedge t < t_z + 11\}$. The estimated partial effects for the *Speed* and *Dist* training load variables, $\hat{h}_1(t - t_z)z_1(t_z)$ and $\hat{h}_2(t - t_z)z_2(t_z)$, are shown in Figure 2.

The results suggest that no more than seven sessions in the past are relevant for the cumulative effects of the *Speed* and *Dist* variables. Both cumulative effects are estimated to have a non-linear decaying effect on the covariate z with respect to latency and a linear effect on latency with respect to the covariate z . Regarding the estimated partial effects of both variables, the values contributing the most to the hazard are those most recently recorded, while the contribution of the values recorded longer ago diminishes. Although there is not much difference in the trend, the greater the average speed and the total distance covered in recent sessions, the greater the impact on the resulting cumulative effect. See also Figure S1 of the online supplementary material.

The Gaussian frailty term (random intercepts) that accounts for the correlation between subsequent injuries from the same player is statistically significant (p -value < 0.01), with $\hat{\sigma}_b^2 = 0.22$ the estimated variance. Players who suffered more injuries (e.g. Id04, Id28) have a higher baseline hazard of injury, as seen in Figure 3, which shows the estimated smooth baseline hazard, $\hat{\lambda}_0(t)$, together with the estimated player-specific smooth baseline hazard. With respect to the session type effect, match sessions have a higher risk of injury compared to training sessions ($\hat{\beta}_1 = 2.45$, and p -value < 0.01). See also Table S1 and Figure S3 of the online supplementary material where the estimated linear and nonlinear effects are presented.

Lastly, we evaluated different methods for capturing cumulative effects by fitting several PAMM models. Each PAMM model incorporated a different approach commonly used in the sports medicine and exercise literature for estimating the training load effect: the ACWR using rolling averages, ACWR with exponentially weighted moving averages (EWMA), and an unweighted sum of the last seven days of exposure (cf. Section A.2 of the online supplementary material). Table 2 shows the comparison of the model’s performance based on likelihood-based measures. Our model shows a better overall performance compared to the other alternatives.

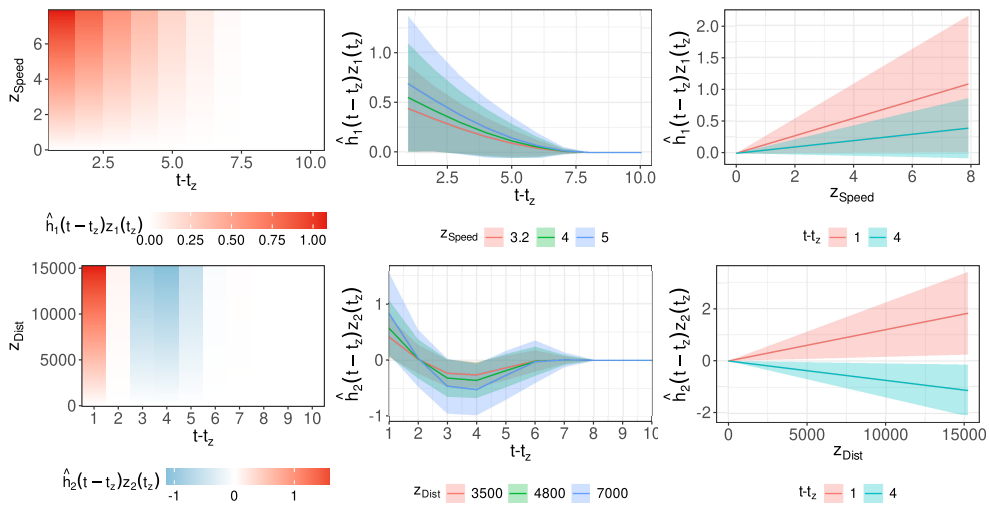


Figure 2. Top: Estimated partial effects surface (left-hand panel), one-dimensional slices through the previous surface for the covariate *Speed*, $z_1(t_z) = z_{Speed}(t_z) \in \{3.2, 4, 5\}$ (middle panel) and the latency $t - t_z \in \{1, 4\}$ (right panel) on the log-hazard scale. Bottom: the analogue for the covariate *Dist*, where the estimated one-dimensional partial effects are conditioned on the values $z_2(t_z) = z_{Dist}(t_z) \in \{3,500, 4,800, 7,000\}$ and $t - t_z \in \{1, 4\}$.

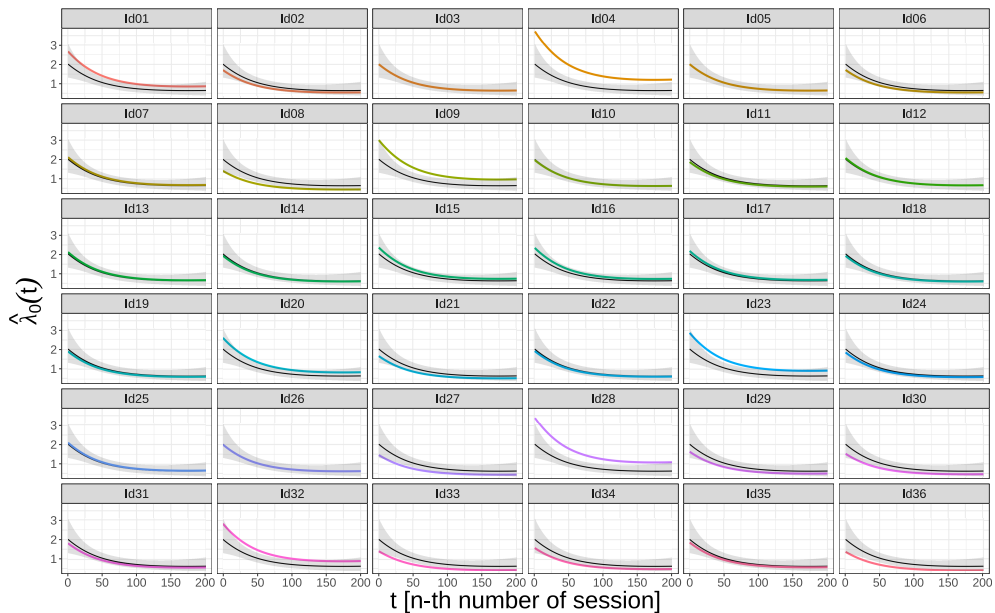


Figure 3. Estimated team smooth baseline hazard (black) and player-specific smooth baseline hazard (coloured line in each panel), together with confidence intervals (grey shadow) of the team's smooth baseline hazard estimate.

4 Simulation study

We conduct extensive simulation studies to evaluate the proposed models and to investigate their properties, aiming to validate their applicability in real-world settings. In particular, we aim to (i) assess the ability of the model to simultaneously estimate both, flexible WCE-type effects and heterogeneity resulting from recurrent events; and (ii) study the implementation of penalties on the basis coefficients to select the maximum length of the time window in which past exposures are cumulatively associated with the hazard.

Table 2. The goodness-of-fit of the fitted models, ordered from the best performance to the least, according to BIC

Model	AIC	Deviance	Deviance explained	BIC
PAMM WCE ridge model	717.21	539.58	20.57	866.41
Unweighted sum model	802.74	628.49	7.48	924.92
ACWR (rolling avg.) model	804.24	625.67	7.89	950.49
ACWR (EWMA) model	796.74	616.54	9.24	951.02

4.1 Data generation

We draw survival times from the piece-wise exponential distribution. Let n_l be the maximum event number that individual l has been at risk for, where $l = 1, \dots, L$. Then, it suffices to specify a vector of piece-wise constant hazards $\lambda = (\lambda_{1,1}, \lambda_{2,1}, \dots, \lambda_{n_1,1}, \lambda_{1,2}, \lambda_{2,2}, \dots, \lambda_{n_2,2}, \dots, \lambda_{1,L}, \lambda_{2,L}, \dots, \lambda_{n_L,L})$, in intervals defined by $J + 1$ cut-points, i.e. by the vector of interval borders $\kappa = (0 = \kappa_0, \dots, \kappa_J = t_{\max})$. That is, λ_{ij} is composed of $(\lambda_{ij,1}, \lambda_{ij,2}, \dots, \lambda_{ij,j})$, where each element λ_{ijj} is the hazard rate of i th event for individual l in the interval $j, i_l = 1, \dots, n_l$ and $j = 1, \dots, J$; and can be defined through a function of time t , current and past exposure covariates $z(t)$ and a random effect b_l , i.e. $\lambda_{ij,j}(t | z(t), b) = f(t, z(t), b) = \exp(const. + f_0(t) + \int_{[t_z: t \geq t_z]} h(t - t_z)z(t_z) dt + b_l)$, evaluated at time $t = \kappa_j$.

Then, we draw recurrent survival times from the piece-wise exponential distribution (PEXP), $t \sim \text{PEXP}(\lambda, \kappa)$, for which the algorithm is outlined in [Table S3 of the online supplementary material](#). The hazard rate vector λ is defined based on the simulation settings described in the following section. All further details on data generation are provided in [Section B of the online supplementary material](#).

4.2 Scenarios and parameter settings

Basing on real-world application data for parameter settings, we simulate $N_{\text{sim}} = 500$ times a cohort of L players, varying L in $\{20, 40, 100\}$, with exposures recorded at $t_{z,1} = 1, t_{z,2} = 2, \dots, t_{z,Q=40} = 40$ days before the time at which we model the hazard, $z_l(t) = (z_l(t_{z,1}), z_l(t_{z,2}), \dots, z_l(t_{z,Q}))$, and draw survival times from the piece-wise exponential distribution under four different true weight functions, $h(t - t_z)$, (a) *exponential decay*, (b) *bi-linear* (c) *early peak* and (d) *inverted U shapes*, each defined over a $[0, t_{z,Q}]$ interval (see the black curves in [Figure 4](#)); and under three different levels of heterogeneity between recurrent events, $\sigma_b \in \{0.05, 0.5, 1\}$, indicating very low heterogeneity, low heterogeneity and high heterogeneity, respectively (see also [Figure S4 in the online supplementary material](#)). We then fit three different PAMMs with WCE-type cumulative effects: a model with no constraint (*Uncons.*), adding a constraint (*Constr.*) and adding a ridge penalty (*Ridge*). The performance of the models is evaluated by graphical inspection of the estimated $\hat{h}(t - t_z)$ function in comparison to the true simulated $h(t - t_z)$ function; the accuracy of these $\hat{h}(t - t_z)$ WCE-type cumulative effects estimates are also evaluated via the mean RMSE, i.e. $\overline{\text{RMSE}}$, over all simulation runs, as:

$$\overline{\text{RMSE}} = \frac{1}{N_{\text{sim}}} \sum_{n=1}^{N_{\text{sim}}} \sqrt{\frac{1}{N_{t_z}} \sum_{t-t_z=0}^{40} (h(t - t_z) - \hat{h}(t - t_z)^{(n)})^2},$$

where $N_{t_z} = 41$, since $t - t_z = \{0, 1, 2, \dots, 40\}$ takes 40 + 1 number of different values.

We assess the accuracy of the standard deviation of the random effects, $\hat{\sigma}_b$, through:

$$\text{RMSE} = \sqrt{\frac{1}{N_{\text{sim}}} \sum_{n=1}^{N_{\text{sim}}} (\sigma_b - \hat{\sigma}_b^{(n)})^2}.$$

We also evaluate the rate at which the estimated confidence interval of WCE-type cumulative effects estimates contains the true estimand $h(t - t_z)$, computing the mean coverage at the $1 - \alpha$ confidence level, i.e. $\overline{\text{Coverage}}_\alpha$, as:

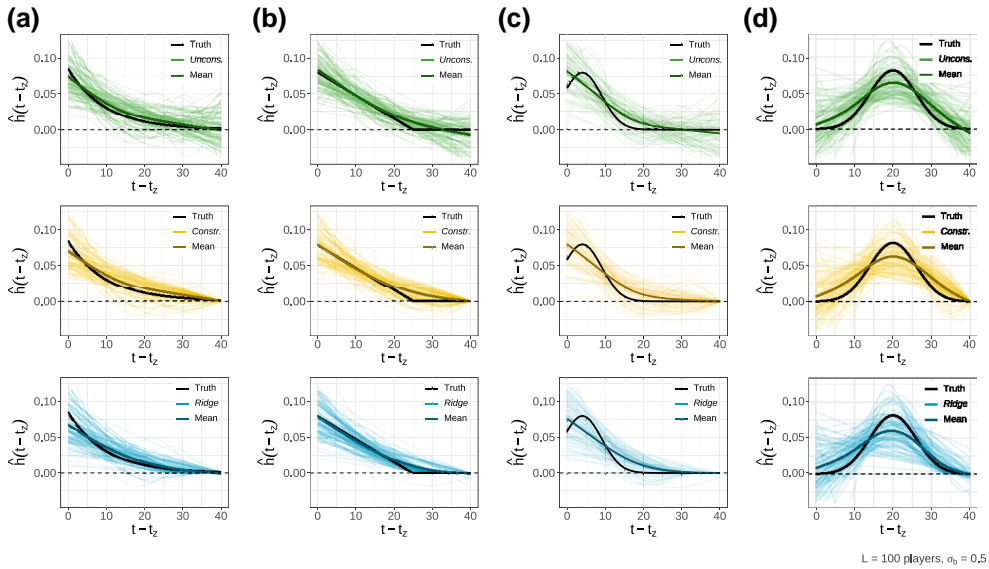


Figure 4. True vs. fitted partial effect weight function $h(t - t_z)$ for scenario $L = 100$ players and $\sigma_b = 0.5$. Rows: *Uncons.* model (top row), *Constr.* model (middle row) and *Ridge* model (bottom row). Columns: (a) *exponential decay*, (b) *bi-linear*, (c) *early peak* and (d) *inverted U*. The true shapes used for simulation are depicted as solid black lines and the mean (point-wise averages) of all simulation runs are shown as solid coloured lines. A random sample of 100 individual estimated weight functions are displayed as light-coloured curves.

$$\overline{\text{Coverage}}_\alpha = \frac{1}{N_{\text{sim}}} \sum_{n=1}^{N_{\text{sim}}} \left[\frac{1}{N_{t_z}} \sum_{t-t_z=0}^{40} \mathbb{I} \left(h(t - t_z) \in \left[\hat{h}(t - t_z)^{(n)} \mp \zeta_{1-\alpha/2} \hat{\sigma}_b^{(n)} \right] \right) \right],$$

where ζ_q is the q -quantile of the standard normal distribution and $\hat{\sigma}_b$ the standard error of the estimated $\hat{h}(t - t_z)$.

4.3 Simulation results

Figure 5 shows boxplots of the distribution of the RMSE, the 95% point-wise coverage (across all time points) and the squared error of σ_b across all simulation settings. In general, as the sample size (i.e. number of players) increases, the estimation errors for $h(t - t_z)$ and σ_b decrease and the mean point-wise coverage for $h(t - t_z)$ increases, across all models. The *Ridge* model provides the most accurate estimate for $h(t - t_z)$, outperforming the other models considered (see the first row in Figure 5). On the other hand, the estimation of the standard deviation of the random effect, σ_b , does not depend on the estimation of the shape of the weight function $h(t - t_z)$. All models provide similar estimates for σ_b , but these estimates become less accurate for higher true σ_b values (see the second row in Figure 5). The 95% coverage of $h(t - t_z)$ for shapes (a), (c) and (d) shows underfitting, specifically for models that penalize the weight function. This underfitting occurs because the true weight function in these cases has features (e.g. sharp peaks) that are not well captured by the penalized models. In addition, the method used to calculate the coverage, i.e. point-wise, may also contribute to this observed underfitting (see the third row in Figure 5). The numerical results for all simulation settings are presented in Tables S4–S7 of the online supplementary material.

Results regarding the estimation of the weight function $\hat{h}(t - t_z)$ for true weight functions (a)–(d), $L = 100$ players and $\sigma_b = 0.5$ are shown in Figure 4 (the rest of the scenarios are shown in Figures S8–S19 of the online supplementary material). In general, the model estimates capture well the underlying weight function. The models with an additional penalty (middle and bottom panels), *Constr.* and *Ridge* models, are more accurate for scenarios where the exposures that occurred relatively long ago (e.g. from the 20th lag on) have little impact on the risk at the actual time t , i.e. shapes (b) and (c). This is proved by the lower mean RMSE values in the Tables S4–S7 of the

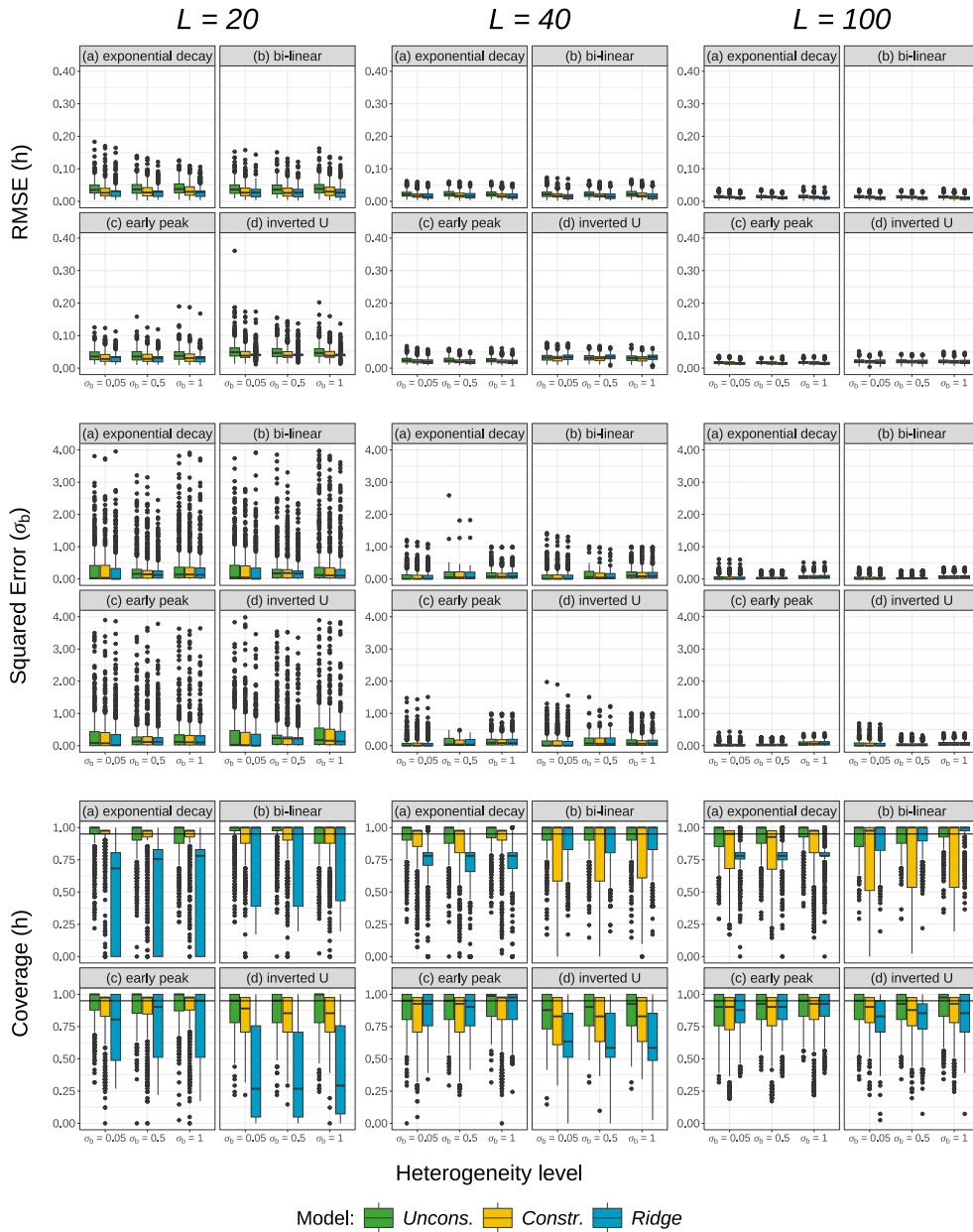


Figure 5. Summary of the simulation results. Rows: distribution of the RMSE (top row), squared error of σ_b (middle row) and 95% point-wise coverage (bottom row) across all σ_b parameters. Columns: results for sample sizes of $L = 20$ $L = 40$ and $L = 100$ players ($N_{sim} = 500$).

online supplementary material. The best estimation of $h(t - t_z)$, among the settings considered, is obtained for shape (b) *bi-linear* with ridge penalty, according to RMSE and 95% coverage, that has a mean RMSE of 0.010 and a mean coverage of 93.2% (see Tables S4–S7 of the online supplementary material).

5 Discussion

We extended and assessed the PAMM model class to the context of recurrent events with time-varying covariates, modelled as WCE-type cumulative effects. By introducing a ridge penalty to

diminish the influence of past exposures registered long ago, we presented a method to determine a relevant time window based on data. Lastly, motivated by the research question regarding the association between external training load and (subsequent) time-loss injuries, we applied the proposed methodology in the sports medicine context.

Simulations indicate that PAMM with ridge penalization is the method yielding the greatest accuracy for the partial effects estimates, $\hat{h}(t, t_z, z(t_z)) = \hat{h}(t - t_z)z(t_z)$. The additional ridge penalization of the weight function enables us to identify the relevant window $\tau(t)$ at which past exposures cumulatively affect the hazard at time t , as seen also in our application on football injury data. From a practical point of view, the presented modelling framework provides a suitable way to flexibly model training load exposures and analyse their effect on subsequent football injuries. Its flexibility in modelling the exposure window is a significant advantage, especially given the wide array of player load metrics available through GPS devices in professional sports. Being able to determine and analyse the time window in which load influences injury outcomes, without the need to fit a range of separate models (e.g. for 3, ..., 6, or 7 days back), is highly appealing to practitioners in professional sports. We propose to use wide time windows to properly determine the exposure time at which, from that time on, the estimated effects are close to zero.

Moreover, the proposed model demonstrates a better overall performance compared to other alternative measures of training load exposures used in the literature. The widely known and used acute chronic workload ratio (ACWR) (Hulin et al., 2014), and its variants (Lolli et al., 2019; Wang et al., 2020), limit to summarize past observations into a predefined unweighted metric, through a ratio of two rolling averages –last 7 days (acute load) over last 28 days (chronic load). The same applies to the exponentially weighted moving averages (EWMA) (Williams et al., 2017) metric suggested as an alternative measure of rolling averages. While EWMA more accurately accounts for the decaying nature of fitness and fatigue effects over time than rolling averages, both may fail to accurately reflect several changes in past training exposures, as well as to consider predefined time windows that are either superfluous or insufficient. Instead, our method's estimates of the cumulative effect and the relevant time window are based on the data.

Future research should evaluate negative binomial and zero-inflated models within the PAMM framework to address issues related to overdispersion and excess zeros (e.g. the low number of injuries). Although this is beyond the scope of the current study, comparing out-of-sample injury predictions from our model with other methods (e.g. regression models with ACWR-transformed training load variables and machine learning approaches like gradient boosting machines, neural networks, and random forests) would be a valuable area for further investigation. As this approach requires the estimation of many parameters and effects, which can be challenging when based on a limited number of events, another direction might involve pooling data from multiple teams or sources. This could be achieved using federated learning techniques (Archetti & Matteucci, 2023) to ensure data protection.

Simulation studies suggest that the model can recover a number of clinically plausible shapes for the true weight function under various levels of heterogeneity. Without prior knowledge about the form of association for time-varying exposures, the model proved to capture well a variety of shapes, estimating them from the data via P-splines. However, for nonsmooth effect shapes, like piece-wise constant or bi-linear, there exist other alternative methods such as adaptive splines (J. H. Friedman, 1991) or treed distributed lag nonlinear models (TDLNM) (Mork & Wilson, 2022), might be of interest.

Additionally, future research should explore the impact of the number of events per subject, held fixed in our simulation study, on model performance. It would also be worthwhile to investigate distributions other than Gaussian, such as Gamma-distributed random effects, which are commonly used in survival analysis (Balan & Putter, 2020).

By highlighting the potential value of PAMMs with WCE effects in assessing recurrent events in sports medicine, our study contributes to enriching the existing literature. We believe that this methodology would help in designing and comparing personalized training plans with insights regarding the risk of injury.

Acknowledgements

The authors thank the Medical Services of Athletic Club for data support and all the individuals who participated in the study.

Conflicts of interest: The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

This work was supported by the Basque Government through the BERC 2022–2025 program; of the Ministry of Science, Innovation and Universities through BCAM Severo Ochoa accreditation and through SEV-2017-0718 PRE2018-084007 funding; of AEI/FEDER, UE through the PID2020-115882RB-I00 with acronym ‘S3M1P4R’ and PID2023-153222OB-I00 with acronym ‘SPHERES’; and of Provincial Council of Bizkaia through the 6/12/TT/2022/00006 and acronym ‘MATH4SPORTS’.

Data availability

The application data used in this study is not publicly available due to privacy concerns. However, the simulation data and code used to generate the results are fully available and reproducible at: <https://github.com/lzumeta/flex-mod-training-loads-recu-injuries>.

Supplementary material

Supplementary material is available online at *Journal of the Royal Statistical Society: Series C*.

References

- Archetti A., Ieva F., & Matteucci M. (2023). Scaling survival analysis in healthcare with federated survival forests: A comparative study on heart failure and breast cancer genomics. *Future Generation Computer Systems*, 149, 343–358. <https://doi.org/10.1016/j.future.2023.07.036>
- Argyropoulos C., & Unruh M. L. (2015). Analysis of time to event outcomes in randomized controlled trials by Generalized Additive Models. *PLoS One*, 10(4), 1–33. <https://doi.org/10.1371/journal.pone.0123784>
- Bache-Mathiesen L. K., Andersen T. E., Dalen-Lorentsen T., Clarsen B., & Fagerland M. W. (2021). Not straight-forward: Modelling non-linearity in training load and injury research. *BMJ Open Sport & Exercise Medicine*, 7(3), e001119. <https://doi.org/10.1136/bmjsem-2021-001119>
- Bache-Mathiesen L. K., Andersen T. E., Dalen-Lorentsen T., Clarsen B., & Fagerland M. W. (2022). Assessing the cumulative effect of long-term training load on the risk of injury in team sports. *BMJ Open Sport and Exercise Medicine*, 8(2), e001342. <https://doi.org/10.1136/bmjsem-2022-001342>
- Bahr R. (2003). Risk factors for sports injuries: A methodological approach. *British Journal of Sports Medicine*, 37(5), 384–392. <https://doi.org/10.1136/bjsem.37.5.384>
- Bahr R., Clarsen B., Derman W., Dvorak J., Emery C. A., Finch C. F., Hägglund M., Junge A., Kemp S., Khan K. M., Marshall S. W., Meeuwisse W., Mountjoy M., Orchard J. W., Pluim B., Quarrie K. L., Reider B., Schweltnus M., Soligard T., ...Chamari K. (2020). International Olympic Committee consensus statement: Methods for recording and reporting of epidemiological data on injury and illness in sport 2020 (including STROBE Extension for Sport Injury and Illness Surveillance (STROBE-SIIS)). *British Journal of Sports Medicine*, 54(7), 372–389. <https://doi.org/10.1136/bjsports-2019-101969>
- Balan T. A., & Putter H. (2020). A tutorial on frailty models. *Statistical Methods in Medical Research*, 29(11), 3424–3454. <https://doi.org/10.1177/0962280220921889>
- Bender A. (2018). *Flexible modeling of time-to-event data and exposure-lag-response associations* [PhD thesis]. Ludwig Maximilian University of Munich, Germany. <https://doi.org/10.5282/edoc.22758>
- Bender A., Groll A., & Scheipl F. (2018). A generalized additive model approach to time-to-event analysis. *Statistical Modelling*, 18(3–4), 299–321. <https://doi.org/10.1177/1471082X17748083>
- Bender A., Scheipl F., Hartl W., Day A. G., & Küchenhoff H. (2019). Penalized estimation of complex, non-linear exposure-lag-response associations. *Biostatistics*, 20(2), 315–331. <https://doi.org/10.1093/biostatistics/kxy003>

- Carey D. L., Crossley K. M., Whiteley R. O. D., Mosler A., Ong K. L., Crow J., & Morris M. E. (2018). Modeling training loads and injuries: The dangers of discretization. *Medicine & Science in Sports & Exercise*, 50(11), 2267–2276. <https://doi.org/10.1249/MSS.0000000000001685>
- Casals M., & Finch C. F. (2017). Sports biostatistician: A critical member of all sports science and medicine teams for injury prevention. *British Journal of Sports Medicine*, 52(22), 1457–1461. <https://doi.org/10.1136/bjsports-2016-042211rep>
- Danieli C., & Abrahamowicz M. (2019). Competing risks modeling of cumulative effects of time-varying drug exposures. *Statistical Methods in Medical Research*, 28(1), 248–262. <https://doi.org/10.1177/0962280217720947>
- Eilers P. H., & Marx B. D. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science*, 11(2), 89–121. <https://doi.org/10.1214/ss/1038425655>
- Eilers P. H., & Marx B. D. (2021). *Practical smoothing: The joys of P-splines* (1st ed.). Cambridge University Press. <https://doi.org/10.1017/9781108610247>
- Friedman J. H. (1991). Multivariate adaptive regression splines. *Annals of Statistics*, 19(1), 1–67. <https://doi.org/10.1214/aos/1176347963>
- Friedman M. (1982). Piecewise exponential models for survival data with covariates. *Annals of Statistics*, 10(1), 101–113. <https://doi.org/10.1214/aos/1176345693>
- Fuller C. W., Colin W., Ekstrand J., Junge A., Andersen T. E., Bahr R., Dvorak J., Hägglund M., McCrory P., & Meeuwisse W. H. (2006). Consensus statement on injury definitions and data collection procedures in studies of football (soccer) injuries. *Clinical Journal of Sports Medicine*, 16(2), 97–106. <https://doi.org/10.1097/00042752-200603000-00003>
- Gabbett T. J., Whyte D. G., Hartwig T. B., Wescombe H., & Naughton G. A. (2014). The relationship between workloads, physical performance, injury and illness in adolescent male football players. *Sports Medicine*, 44(7), 989–1003. <https://doi.org/10.1007/s40279-014-0179-5>
- Gasparrini A., Scheipl F., Armstrong B., & Kenward M. G. (2017). A penalized framework for distributed lag non-linear models. *Biometrics*, 73(3), 938–948. <https://doi.org/10.1111/biom.12645>
- Griffin A., Kenny I. C., Comyns T. M., & Lyons M. (2020). The association between the acute: Chronic workload ratio and injury and its application in team sports: A systematic review. *Sports Medicine*, 50(3), 561–580. <https://doi.org/10.1007/s40279-019-01218-2>
- Holford T. R. (1980). The analysis of rates and of survivorship using log-linear models. *Biometrics*, 36(2), 299–305. <https://doi.org/10.2307/2529982>
- Hulin B. T., Gabbett T. J., Blanch P., Chapman P., Bailey D., & Orchard J. W. (2014). Spikes in acute workload are associated with increased injury risk in elite cricket fast bowlers. *British Journal of Sports Medicine*, 48, 708–712. <https://doi.org/10.1136/bjsports-2013-092524>
- Hulin B. T., Gabbett T. J., Lawson D. W., Caputi P., & Sampson J. A. (2016). The acute: Chronic workload ratio predicts injury: High chronic workload may decrease injury risk in elite rugby league players. *British Journal of Sports Medicine*, 50(4), 231–236. <https://doi.org/10.1136/bjsports-2015-094817>
- Laird N., & Olivier D. (1981). Covariance analysis of censored survival data using log-linear analysis techniques. *Journal of the American Statistical Association*, 76(374), 231–240. <https://doi.org/10.1080/01621459.1981.10477634>
- Li X., Chang C. C. H., Donohue J. M., & Krafty R. T. (2022). A competing risks regression model for the association between time-varying opioid exposure and risk of overdose. *Statistical Methods in Medical Research*, 31(6), 1013–1030. <https://doi.org/10.1177/09622802221075933>
- Lolli L., Batterham A. M., Hawkins R., Kelly D. M., Strudwick A. J., Thorpe R., Gregson W., & Atkinson G. (2019). Mathematical coupling causes spurious correlation within the conventional acute-to-chronic workload ratio calculations. *British Journal of Sports Medicine*, 53(15), 921–922. <https://doi.org/10.1136/bjsports-2017-098110>
- Marx B. D., & Eilers P. H. (1998). Direct generalized additive modeling with penalized likelihood. *Computational Statistics & Data Analysis*, 28(2), 193–209. [https://doi.org/10.1016/S0167-9473\(98\)00033-4](https://doi.org/10.1016/S0167-9473(98)00033-4)
- McGilchrist C. A., & Aisbett C. W. (1991). Regression with frailty in survival analysis. *Biometrics*, 47(2), 461–466. <https://doi.org/10.2307/2532138>
- Mork D., & Wilson A. (2022). Treed distributed lag nonlinear models. *Biostatistics*, 23(3), 754–771. <https://doi.org/10.1093/biostatistics/kxaa051>
- Nielsen R. Ø., Bertelsen M. L., Møller M., Hulme A., Windt J., Verhagen E., Mansournia M. A., Casals M., & Parner E. T. (2018). Training load and structure-specific load: Applications for sport injury causality and data analyses. *British Journal of Sports Medicine*, 52(16), 1016–1017. <https://doi.org/10.1136/bjsports-2017-097838>
- Nielsen R. Ø., Shrier I., Casals M., Nettel-Aguirre A., Møller M., Bolling C., Bittencourt N. F. N., Clarsen B., Wedderkopp N., Soligard T., Timpka T., Emery C. A., Bahr R., Jacobsson J., Whiteley R., Dahlström Ö., van Dyk N., Pluim B. M., Stamatakis E., ... Verhagen E. (2020). Statement on methods in sport injury research from the first methods matter meeting, Copenhagen, 2019. *Journal of Orthopaedic and Sports Physical Therapy*, 50(5), 226–233. <https://doi.org/10.2519/jospt.2020.9876>

- Obermeier V., Scheipl F., Heumann C., Wassermann J., & Küchenhoff H. (2015). Flexible distributed lags for modelling earthquake data. *Journal of the Royal Statistical Society. Series C: Applied Statistics*, 64(2), 395–412. <https://doi.org/10.1111/rssc.12077>
- Ramjith J., Bender A., Roes K. C. B., & Jonker M. A. (2022). Recurrent events analysis with piece-wise exponential additive mixed models'. *Statistical Modelling*, 24(3), 266–287. <https://doi.org/10.1177/1471082X221117612>
- Sainani K. L., Borg D. N., Caldwell A. R., Butson M. L., Tenan M. S., Vickers A. J., Vigotsky A. D., Warmenhoven J., Nguyen R., Lohse K. R., Knight E. J., & Bargary N. (2021). Call to increase statistical collaboration in sports science, sport and exercise medicine and sports physiotherapy. *British Journal of Sports Medicine*, 55(2), 118–122. <https://doi.org/10.1136/bjsports-2020-102607>
- Signorello L. B., McLaughlin J. K., Lipworth L., Friis S., Sørensen H. T., & Blot W. J. (2022). Confounding by indication in epidemiologic studies of commonly used analgesics. *American Journal of Therapeutics*, 9(3), 199–205. <https://doi.org/10.1097/00045391-200205000-00005>
- Soligard T., Schwelnus M., Alonso J. M., Bahr R., Clarsen B., Dijkstra H. P., Gabbett T., Gleeson M., Hägglund M., Hutchinson M. R., van Rensburg C. J., Khan K. M., Meeusen R., Orchard J. W., Pluim B. M., Raftery M., Budgett R., & Engebretsen L. (2016). How much is too much? (Part 1) International Olympic Committee consensus statement on load in sport and risk of injury. *British Journal of Sports Medicine*, 50(17), 1030–1041. <https://doi.org/10.1136/bjsports-2016-096581>
- Sylvestre M. P., & Abrahamowicz M. (2009). Flexible modeling of the cumulative effects of time-dependent exposures on the hazard. *Statistics in Medicine*, 28(27), 3437–3453. <https://doi.org/10.1002/sim.v28:27>
- Wang C., Vargas J. T., Stokes T., Steele R., & Shrier I. (2020). Analyzing activity and injury: Lessons learned from the acute: Chronic workload ratio. *Sports Medicine*, 50(7), 1243–1254. <https://doi.org/10.1007/s40279-020-01280-1>
- Williams S., West S., Cross M. J., & Stokes K. A. (2017). Better way to determine the acute: Chronic workload ratio? *British Journal of Sports Medicine*, 51(3), 209–210. <https://doi.org/10.1136/bjsports-2016-096589>
- Windt J., Ardern C. L., Gabbett T. J., Khan K. M., Cook C. E., Sporer B. C., & Zumbo B. D. (2018). Getting the most out of intensive longitudinal data: A methodological review of workload-injury studies. *BMJ Open*, 8(10), e022626. <https://doi.org/10.1136/bmjopen-2018-022626>
- Windt J., & Gabbett T. J. (2017). How do training and competition workloads relate to injury? The workload-injury aetiology model. *British Journal of Sports Medicine*, 51(5), 428–435. <https://doi.org/10.1136/bjsports-2016-096040>
- Wood S. N. (2011). Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 73(1), 3–36. <https://doi.org/10.1111/j.1467-9868.2010.00749.x>
- Wood S. N. (2017). *Generalized additive models: An introduction with R* (2nd ed.). CRC press. <https://doi.org/10.1201/9781315370279>