

Having a Ball: evaluating scoring streaks and game excitement using in-match trend estimation

Claus Thorn Ekstrøm and Andreas Kryger Jensen
Biostatistics, Institute of Public Health, University of Copenhagen
ekstrom@sund.ku.dk, aeje@sund.ku.dk

20 November, 2020

Abstract

Nu ved jeg godt nok intet om sport, men... ❁!

Keywords: APBRmetrics, Bayesian Statistics, Gaussian Processes, Sports Statistics, Trends

1 Introduction

Sports analytics receive increasing attention within statistics and not just for match prediction or betting but also for game evaluation, in-game and post-game coaching purposes, and for setting strategies and tactics in future matches.

Many popular sports such as football (soccer), basketball, boxing, table tennis, volleyball, American football, and handball involves matches between two teams or players where each team have the possibility of scoring throughout the match. Several research papers seek to predict match results (e.g., Karlis and Ntzoufras (2003); Groll et al. (2019); Gu and Saaty (2019)) or match winners for single matches (Cattelan, Varin, and Firth (2013)) in order to infer the match winner and potentially the winner of a tournament (Ekstrøm et al. (2020); Baboota and Kaur (2018)). While the overall result is highly interesting it conveys very little information about the individual development and trends throughout the match and modeling approaches that allow a finer granularity of the score difference throughout the match is needed.

The trend in score difference between the two teams is a proxy for their underlying strengths. In particular, sustained periods of time where the score difference increases suggests that one team outperforms the other whereas periods where the teams are constantly catching up to each other suggest that the team's strengths in those periods are similar. Modeling the local trend of the score difference will therefore reflect several aspects of the game in particular the team strengths and game dynamics and momentum as they develop through the match.

Figure 1 shows the development of the score difference for the final match of the playoffs in the 2020 National Basketball Association (NBA) series between Los Angeles Lakers and Miami Heat. Positive numbers indicate that LA Lakers are leading and the running score difference shows that the Lakers pulled ahead until the third quarter where Miami Heat started to keep up the scoring pace before overtaking the Lakers and reducing the lead.

In this paper we will consider the score difference between two teams as a latent Gaussian process

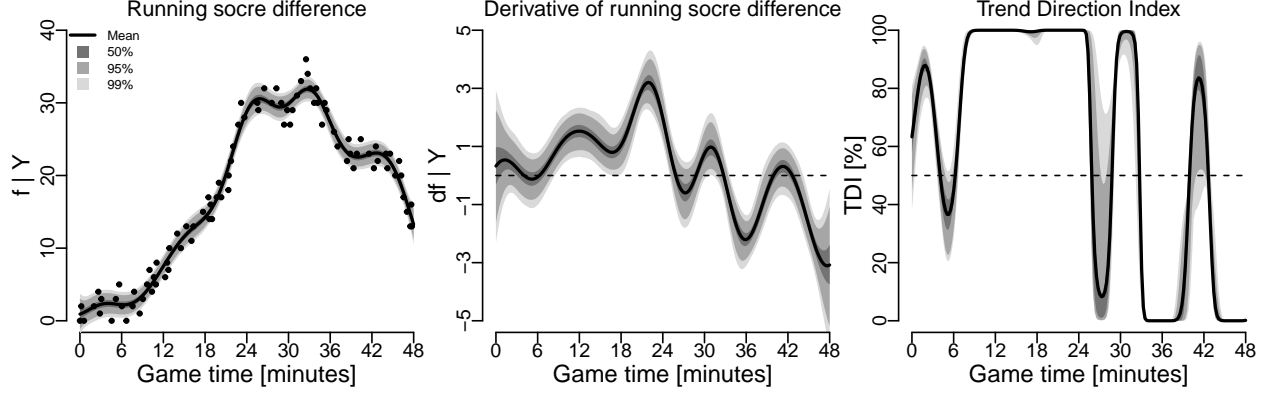


Figure 1: Los Angeles Lakers at Miami Heat. October 11, 2020 [??? elaborate ???]

and use the Trend Direction Index (TDI) from Jensen and Ekstrøm (2020b) as an accomodating measure to evaluate the local probability of the *monotonicity* of the latent process at a given time point during a match. The Trend Direction Index uses a Bayesian framework to provide a direct answer to questions such as “What is the probability that the latent process (i.e., that one team is doing better than another) is increasing at a given time-point?”. This will allow real-time evaluation of the score difference trend at the current time-point in the game and will provide post-game inference about the “hot” periods of a match where one team out-performed the other. Furthermore we will present the Excitement Trend Index (ETI) as an objective measure of spectator excitement in a given match. The ETI is defined as the expected number of times that the score difference changes monotonicity during a match. If the score difference changes monotonicity often then that echos a game where both teams frequently score that the whereas a game with a low ETI will represent a one-sided match where one team is doing consistently better than the other over sustained periods of time.

Other authors have considered using continuous processes to model the score difference of matches. Gabel and Redner (2012) shows that NBA basketball score difference is well described by a continuous-time anti-persistent random walk which suggests that a latent Gaussian process might be viable.

- Chen, Dawson, and Müller (2020)
- Chen and Fan (2018)

The paper is structured as follows. In the next section we introduce trend modeling through latent a Gaussian process and define the Trend Direction Index and Excitement Trend Index that captures the local trends in monotonicity and game excitement, respectively. In Section 3 we apply the our proposed methodology to analyse both the final match of the playoff as well as evaluating the game excitement distribution of the season by considering the ETI from all 1143 matches from the 2019–2020 National Basketball Association (NBA) season. We show how this distribution can be used to assess relative match excitement. We conclude with a discussion in Section 4. Materials to reproduce this manuscript and its analyses can be found at Jensen and Ekstrøm (2020a).

Unused stuff for now: *Philosophical question: What exactly is the sampling model for a single match? We could have a motivating figure here – e.g. a variation of Figure 1 (see its caption).*

2 Methods

Our model is based on the observed score differences D_m in a given match indexed by m . For each match we observe the random variables $\mathcal{D}_m = (t_{mi}, D_{mi})_{0 < i \leq J_m}$ where $t_{m1} < t_{m2} < t_{mi} < \dots < t_{mJ_m}$ are the ordered time points at which any teams scores, $D_{mi} = D_m(t_{mi})$ is the associated difference in scores at time t_{mi} , and J_m is the total number of scorings during the match. We use the convention that D_m is the difference in scores of the away team with respect to the home team, so that $D_m(t) > 0$ means that the away team is leading at time t .

We assume that the observed match data from a given match are noisy realizations of a latent smooth, random function defined in continuous time evaluated at the random time points where goals occur. Let d_m be the latent functions from which the realizations \mathcal{D}_m are generated. Our objective is to infer d_m and its time dynamics from \mathcal{D}_m . In pursuance of this intention we propose the following hierarchical model where d_m is a Gaussian process defined on a compact subset of the real line \mathcal{I}_m corresponding to the duration of the m 'th game, and the observed data conditional on the scoring times and the values of the latent process at these times are independently normally distributed random variables with a match specific variance. This model for a given match can be written as

$$\begin{aligned} \star_m \mid \Psi_m, \mathbf{t}_m &\sim H(\star_m \mid \Psi_m) \\ d_m(t) \mid \star_m &\sim \mathcal{GP}(\mu_{\beta_m}(t), C_{\theta_m}(s, t)) \\ D_m(t_{mi}) \mid d_m(t_{mi}), t_{mi}, \star_m &\stackrel{iid}{\sim} N(d_m(t_{mi}), \sigma_m^2) \end{aligned} \quad (1)$$

where $\star_m = (\beta_m, \theta_m, \sigma_m^2)$ is a vector of hyper-parameters governing the dynamics of the latent Gaussian process with a prior distribution H indexed by parameters Ψ_m , and $\mathbf{t}_m = (t_{m1}, \dots, t_{mJ_m})$ is the vector of time points where goals occurs in the match. The functions μ_{β_m} on \mathcal{I}_m and C_{θ_m} on $\mathcal{I}_m \times \mathcal{I}_m$ are the prior mean and covariance functions of the latent Gaussian process, and σ_m^2 is the variance characterizing the magnitudes deviations between for the observed score differences and the values of the latent process.

A Gaussian process is characterized by the multivariate joint normality of all of the joint distributions resulting from evaluating the process at any finite set of time points (Rasmussen and Williams (2006)). Specifically, for any finite set $\mathbf{t}^* \subset \mathcal{I}_m$ it follows that that the vector $d(\mathbf{t}^*) \mid \star_m$ is distributed as $N(\mu_{\beta_m}(\mathbf{t}^*), C_{\theta_m}(\mathbf{t}^*, \mathbf{t}^*))$ where $\mu_{\beta_m}(\mathbf{t}^*)$ is the vector generated by evaluating the prior mean function $\mu_{\beta_m}(t)$ at \mathbf{t}^* and $C_{\theta_m}(\mathbf{t}^*, \mathbf{t}^*)$ is the covariance matrix generated by evaluating the prior covariance function $C_{\theta_m}(s, t)$ at $\mathbf{t}^* \times \mathbf{t}^*$. Using the properties of multivariate normal distributions the posterior distribution of $d_m(\mathbf{t}^*) \mid \mathcal{D}_m, \star_m$ is also a multivariate normal distribution. This facilitates Bayesian estimation of the distribution of the latent process governing the score difference given the observed data from each match.

In addition to obtaining inference for the latent process we may also estimate its time dynamics. This follows since a Gaussian process along with its time derivatives (provided they exist) are distributed as a multivariate Gaussian process (Cramer and Leadbetter (1967)). We can therefore augment the model in Equation (1) with an additional latent structure of the first and second derivatives of d_m with respect to time as

$$\begin{bmatrix} d_m(s) \\ d'_m(t) \\ d''_m(u) \end{bmatrix} \mid \star_m \sim \mathcal{GP} \left(\begin{bmatrix} \mu_{\beta_m}(s) \\ \mu'_{\beta_m}(t) \\ \mu''_{\beta_m}(u) \end{bmatrix}, \begin{bmatrix} C_{\theta_m}(s, s') & \partial_2 C_{\theta_m}(s, t) & \partial_2^2 C_{\theta_m}(s, u) \\ \partial_1 C_{\theta_m}(t, s) & \partial_1 \partial_2 C_{\theta_m}(t, t') & \partial_1 \partial_2^2 C_{\theta_m}(t, u) \\ \partial_1^2 C_{\theta_m}(u, s) & \partial_1^2 \partial_2 C_{\theta_m}(u, t) & \partial_1^2 \partial_2^2 C_{\theta_m}(u, u') \end{bmatrix} \right) \quad (2)$$

where ' and '' denote the first and second time derivatives, ∂_j^k is the k 'th order partial derivative with respect to the j 'th variable. Combining the models in Equations (1) and (2) we therefore also have explicit expressions for the posterior distributions $d'_m | \mathcal{D}_m, \otimes_m$ and $d''_m | \mathcal{D}_m, \otimes_m$. We use the posteriors of the first and second time derivatives of the latent process to characterize the dynamical properties of each match through the Trend Direction Index (TDI) and the Excitement Trend Index (ETI).

[[[Introduce and fix notation for the posteriors and refer to supplementary for explicit expressions]]]

We define the Trend Direction Index (TDI) of a particular match m as the local posterior probability that d_m is an increasing function at any time in \mathcal{I}_m . Under our model this is equal to

$$\begin{aligned} \text{TDI}_m(t) | \otimes_m &= P(d'_m(t) > 0 | \mathcal{D}_m, \otimes_m) \\ &= \frac{1}{2} + \frac{1}{2} \text{Erf} \left(\frac{\mu_{f'}(t | \otimes_m)}{2^{1/2} \Sigma_{f'f'}(t, t | \otimes_m)^{1/2}} \right) \end{aligned} \quad (3)$$

where $\text{Erf}: x \mapsto 2\pi^{-1/2} \int_0^x \exp(-u^2) du$ is the error function and μ_{df} and Σ_{df} are the posterior mean and covariance of the trend. **[[[Write out interpretation]]]**

[Note that if we change the reference team then TDI is just 1 - TDI so it's symmetric]

To each match we assign its Excitement Trend Index, denoted by ETI_m , as a global index of excitement. The index is defined as the expected number of changes in monotonicity of the posterior distribution of d_m which is equivalent to the expected number of zero-crossings of the posterior distribution of d'_m . We hence define

$$\begin{aligned} \text{ETI}_m | \otimes_m &= \mathbb{E} \left[\# \{t \in \mathcal{I}_m : d'_m(t) = 0\} | \mathcal{D}_m, \otimes_m \right] \\ &= \int_{\mathcal{I}_m} d\text{ETI}_m(t | \otimes_m) dt \end{aligned}$$

where $d\text{ETI}_m$ is the instantaneous posterior probability of a zero-crossing of d' at any time point $t \in \mathcal{I}_m$. Under the model described in Equations (1) and (2) it can be shown that $d\text{ETI}$ is equal to

$$d\text{ETI}_m(t | \otimes_m) = \lambda(t | \otimes_m) \phi \left(\frac{\mu_{df}(t | \otimes_m)}{\Sigma_{df}(t, t | \otimes_m)^{1/2}} \right) \left(2\phi(\zeta(t | \otimes_m)) + \zeta(t | \otimes_m) \text{Erf} \left(\frac{\zeta(t | \otimes_m)}{2^{1/2}} \right) \right)$$

where $\phi: x \mapsto 2^{-1/2} \pi^{-1/2} \exp(-\frac{1}{2}x^2)$ is the standard normal density function, and λ , ω and ζ are functions defined by

$$\begin{aligned} \lambda(t | \otimes_m) &= \frac{\Sigma_{d^2f}(t, t | \otimes_m)^{1/2}}{\Sigma_{df}(t, t | \otimes_m)^{1/2}} \left(1 - \omega(t | \otimes_m)^2 \right)^{1/2} \\ \omega(t | \otimes_m) &= \frac{\Sigma_{df, d^2f}(t, t | \otimes_m)}{\Sigma_{df}(t, t | \otimes_m)^{1/2} \Sigma_{d^2f}(t, t | \otimes_m)^{1/2}} \\ \zeta(t | \otimes_m) &= \frac{\mu_{df}(t | \otimes_m) \Sigma_{d^2f}(t, t | \otimes_m)^{1/2} \omega(t) \Sigma_{df}(t, t | \otimes_m)^{-1/2} - \mu_{d^2f}(t | \otimes_m)}{\Sigma_{d^2f}(t, t | \otimes_m)^{1/2} \left(1 - \omega(t | \otimes_m)^2 \right)^{1/2}} \end{aligned}$$

A derivation of this expression can be found in the supplementary material to Jensen and Ekstrøm (2020b).

We need to argue that ETI for $S_a(t_{m_i}) - S_b(t_{m_i})$ is symmetric in a and b so that our choice of “reference group” in D_m is not important. The reason is that we look at both up- and down-crossings at 0 of df_m so the choice of sign in D_m is not relevant.

2.1 Estimation

We have implemented the model described in the previous section in Stan (Carpenter et al. 2017).

Prior mean and covariance:

$$\mu_{\beta_m}(t) = \beta_m, \quad C_{\theta_m}(s, t) = \alpha_m^2 \exp\left(-\frac{(s-t)^2}{2\rho_m^2}\right)$$

with $\theta_m = (\alpha_m, \rho_m) > 0$. [??? **Note: Infinitely differentiable sample paths. Sample path derivatives are well-defined**]

Hyper-parameters: We used independent priors on $\otimes_m = (\beta_m, \alpha_m, \rho_m, \sigma_m)$ of the form

$$H(\otimes_m | \Psi_m) = H(\beta_m | \Psi_{\beta_m})H(\alpha_m | \Psi_{\alpha_m})H(\rho_m | \Psi_{\rho_m})H(\sigma_m | \Psi_{\sigma_m})$$

where each prior is a heavy-tailed distribution with a moderate variance centered at the marginal maximum likelihood estimates. We used the following distributions

$$\beta_m \sim T(\widehat{\beta_m^{ML}}, 3, 3), \quad \alpha_m \sim T^+(\widehat{\alpha_m^{ML}}, 3, 3), \quad \rho_m \sim N^+(\widehat{\rho_m^{ML}}, 1), \quad \sigma_m \sim T^+(\widehat{\sigma_m^{ML}}, 3, 3)$$

where $T^+(\cdot, \cdot, df)$ and $N^+(\cdot, \cdot)$ denotes the location-scale half T- and normal distribution functions with df degrees of freedom.

The posterior distribution of the hyper-parameters given the observed data is then. We define $\tilde{\Theta}_m \sim P(\Theta_m | \mathbf{D}_m, \Psi_m, \mathbf{t}_m)$ hence

$$\tilde{\Theta}_m \sim \frac{G(\Theta_m | \Psi_m, \mathbf{t}_m) \int P(\mathbf{D}_m | f(\mathbf{t}_m), \Theta_m, \Psi_m, \mathbf{t}_m) dP(f_m(\mathbf{t}_m) | \Theta_m, \Psi_m, \mathbf{t}_m)}{\iint P(\mathbf{D}_m | f_m(\mathbf{t}_m), \Theta_m, \Psi_m, \mathbf{t}_m) dP(f_m(\mathbf{t}_m) | \Theta_m, \Psi_m, \mathbf{t}_m) dG(\Theta_m | \Psi_m, \mathbf{t}_m)}$$

What we estimate is then the random variable $\widehat{\text{ETI}}_m = \text{ETI}_m(\tilde{\Theta}_m)$ which can be summarized by its moments or quantiles.

For each match we ran four independent Markov chains for 25,000 iterations each with half of the iterations used for warm-up. Convergence was assessed by trace plots of the MCMC draws and the potential scale reduction factor, \hat{R} , of Gelman and Rubin (1992).

3 Application: The 2019–2020 NBA basketball season

To show the applicability of our proposed model we will apply it to data from all regular games from the 2019–2020 NBA basketball season (obtained from Sports Reference LLC (2020) and provided in Jensen and Ekström (2020a)). A lot of points are scored during a basketball match so it is easy to see the development of the score difference.

The 2019–2020 NBA season was suspended mid-March because of Covid-19 but it was resumed again in July 2020. There were a total of 1059 regular season matches. The subsequent playoffs comprised 84 matches including the final for a grand total of 1143 matches. For ease of comparison

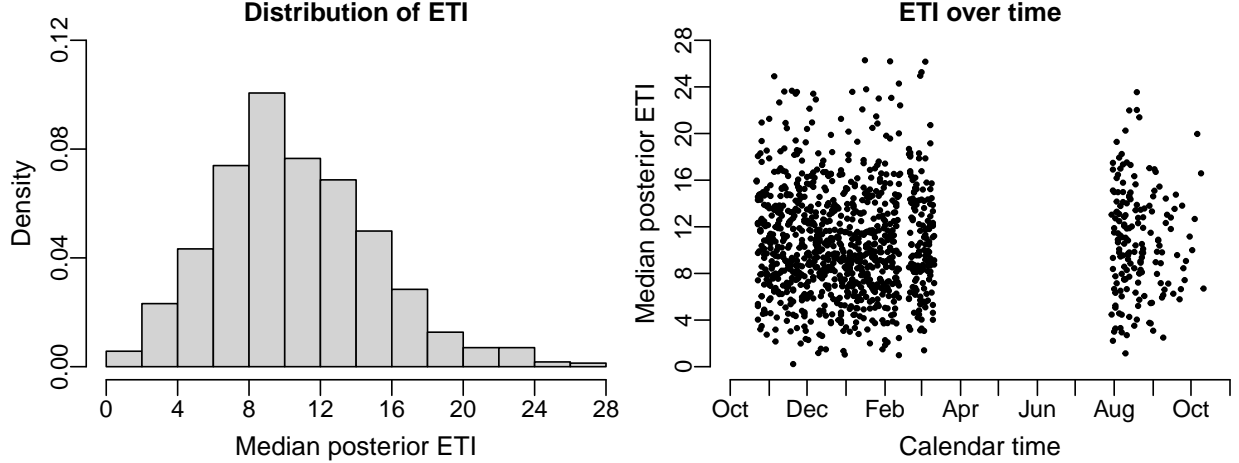


Figure 2: Histogram and kernel density estimate of the distribution of 1143 median posterior Excitement Trend Indices from the NBA 2019–2020 season (left panel) and the median posterior ETIs as a function of calendar time from October 22, 2019 to October 11, 2020 (right panel).

we are only considering the first 48 regular minutes of each match — any part of a match that goes into overtime will be disregarded. Thus, $\mathcal{I}_m = [0; 48]$ minutes without overtime.

For each of the 1143 matches during the season we fitted our model, and [???].

The left panel of Figure 2 shows the distribution of the median posterior Excitement Trend Indices estimated using our model for each of the 1143 games in the 2019–2020 NBA season. The distribution of game excitements is right skewed (skewness = 0.53) with a range of $[0.23; 26.28]$, and an median of 10.09 (mean = 10.62, SD = 4.52). This implies that the time-varying score differences of the games in the season changes monotonicity approximately 10 times during a game on average. The right panel of the same figure shows the values of the median posterior ETIs as a function of the calendar time were the matches were played. Besides illustrating the gap from the Covid-19 hiatus, the figure shows that the excitement indices are relatively evenly distributed throughout the season.

Summarizing the posterior median ETIs at the team level lead to comparable values across all 30 teams. The New Orleans Pelicans had the highest average posterior median ETI during the season with a value of 11.67 (SD = 4.91, IQR = $[3.09; 23.57]$), while the Charlotte Hornets had the lowest average posterior median ETI during the season with a value of 9.29 (SD = 3.74, IQR = $[1.96; 15.98]$). This implies that each team played games during the season that where comparable on average in terms their excitement, and the major source of variation in excitement is governed by the individual matches. The supplementary material to this paper provides summary statistics for all 30 teams.

The asymmetry of the distribution of the Excitement Trend Indices in Figure 2 conjectures the possibility of a latent, categorical variable labeling the matches to different classes. We thus fitted a heteroscedastic Gaussian mixture model to the median posterior Excitement Trend Indices on the match level, where the number of latent classes was determined by a bootstrapped likelihood ratio test using 10,000 bootstrap replicates (Scrucca et al. (2016)). Testing for 1 vs 2 latent classes gave a p-value of < 0.001 , 2 vs 3 gave a p-value of 0.027, 3 vs 4 gave a p-value of 0.008, and 4 vs 5 gave a p-value of 0.300 implying that the distribution of median posterior Excitement Trend Indices can be modeled by a 4-component Gaussian mixture distribution. The following table shows the

maximum likelihood estimates of the parameters of the Gaussian mixture model.

	class 1	class 2	class 3	class 4
π	0.24	0.22	0.28	0.26
μ	5.70	8.89	12.64	14.45
σ^2	4.44	1.87	5.15	23.31

Figure 3 shows the estimated Gaussian mixture model along with the histogram (left panel), and the proportion of matches classified into the four latent classes as determined by a maximum posterior decision rule (right panel).

[So we basically have a group of average matches (class 2), a group below average (class 1), a group above average (class 3) and then also group of 150 matches with an even large average but also a larger variance (class 4). Can we say something about these classes?]

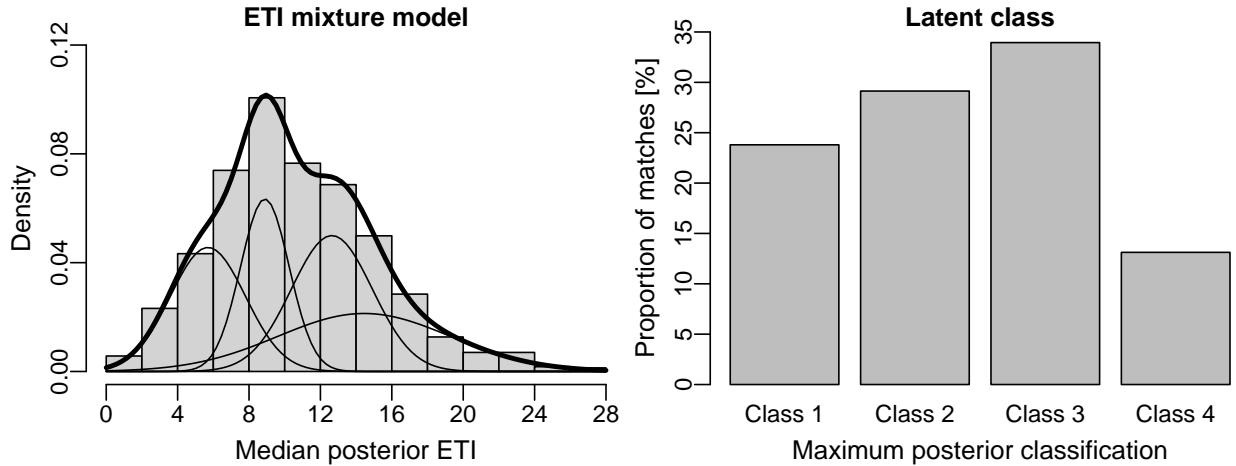


Figure 3

Figure 4 shows ...

4 Discussion

Some summarization goes on here. We have introduced etc. ...

We could define a weighted Excitement Trend Index, ETI_m^W , so that changes in monotonicity of the score differences are i) weighted higher towards the end of the game and ii) weighted lower if one team is already far away of the other team. This motivates a weighted ETI of the form

$$ETI_m^W | \otimes_m = \int_{\mathcal{I}_m} w(t, |d_m(t)|) dETI_m(t | \otimes_m) dt$$

where w is a bivariate weight function being increasing in its first variable and decreasing in its second variable. Such weight function could be constructed as a product of two kernel functions on $[0; 48] \times \mathbb{R}_{\geq 0}$ with bandwidths based on studies of psychological perception.

A different approach would be to define team-specific excitement index nested with a match. Here we would only look at the **up**-crossings at zero of df_m and we would get two excitement indices for each match (ETI_{am} , ETI_{bm}). for teams a and b . This would somehow reflect how exciting each team were in match m with respect to changing the sign of the score differences in their favor.

Acknowledgements

Bibliography

Baboota, Rahul, and Harleen Kaur. 2018. “Predictive Analysis and Modelling Football Results Using Machine Learning Approach for English Premier League.” *International Journal of Forecasting* 35 (March). <https://doi.org/10.1016/j.ijforecast.2018.01.003>.

Carpenter, Bob, Andrew Gelman, Matthew D Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. 2017. “Stan: A Probabilistic Programming Language.” *Journal of Statistical Software* 76 (1).

Cattelan, Manuela, Cristiano Varin, and David Firth. 2013. “Dynamic Bradley–Terry Modelling of Sports Tournaments.” *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 62 (1): 135–50. <https://doi.org/https://doi.org/10.1111/j.1467-9876.2012.01046.x>.

Chen, Tao, and Qingliang Fan. 2018. “A Functional Data Approach to Model Score Difference Process in Professional Basketball Games.” *Journal of Applied Statistics* 45 (1): 112–27.

Chen, Yaqing, Matthew Dawson, and Hans-Georg Müller. 2020. “Rank Dynamics for Functional Data.” *Computational Statistics & Data Analysis*, 106963.

Cramer, Harald, and M. R. Leadbetter. 1967. *Stationary and Related Stochastic Processes – Sample Function Properties and Their Applications*. John Wiley & Sons, Inc.

Ekstrøm, Claus Thorn, Hans Van Eetvelde, Christophe Ley, and Ulf Brefeld. 2020. “Evaluating One-Shot Tournament Predictions.” *Journal of Sports Analytics* Preprint (Preprint): 1–10. <https://doi.org/10.3233/JSA-200454>.

Gabel, Alan, and Sidney Redner. 2012. “Random Walk Picture of Basketball Scoring.” *Journal of Quantitative Analysis in Sports* 8 (1).

Gelman, Andrew, and Donald B. Rubin. 1992. “Inference from Iterative Simulation Using Multiple Sequences.” *Statistical Science* 7 (4): 457–72.

Groll, Andreas, Christophe Ley, Gunther Schaubberger, and Hans Van Eetvelde. 2019. “A hybrid random forest to predict soccer matches in international tournaments.” *Journal of Quantitative Analysis in Sports* 15: 271–88.

Gu, Wei, and Thomas L. Saaty. 2019. “Predicting the Outcome of a Tennis Tournament: Based on Both Data and Judgments.” *Journal of Systems Science and Systems Engineering* 28 (3): 317–43. <https://doi.org/10.1007/s11518-018-5395-3>.

Jensen, Andreas Kryger, and Claus Thorn Ekstrøm. 2020a. “GitHub Repository for Having a Ball.” 2020. <https://github.com/aejensen/Having-a-Ball>.

———. 2020b. “Quantifying the Trendiness of Trends.” *Journal of the Royal Statistical Society: Series C*.

Karlis, Dimitris, and Ioannis Ntzoufras. 2003. “Analysis of Sports Data by Using Bivariate Poisson Models.” *Journal of the Royal Statistical Society: Series D (the Statistician)* 52 (3): 381–93. <https://doi.org/10.1111/1467-9884.00366>.

Rasmussen, C. E., and C. K. I. Williams. 2006. *Gaussian Processes in Machine Learning*. MIT Press.

Scrucca, Luca, Michael Fop, T. Brendan Murphy, and Adrian E. Raftery. 2016. “mclust 5: Clustering, Classification and Density Estimation Using Gaussian Finite Mixture Models.” *The R Journal* 8 (1): 289–317. <https://doi.org/10.32614/RJ-2016-021>.

Sports Reference LLC. 2020. “Basketball Reference.” 2020. <https://www.basketball-reference.com/>.

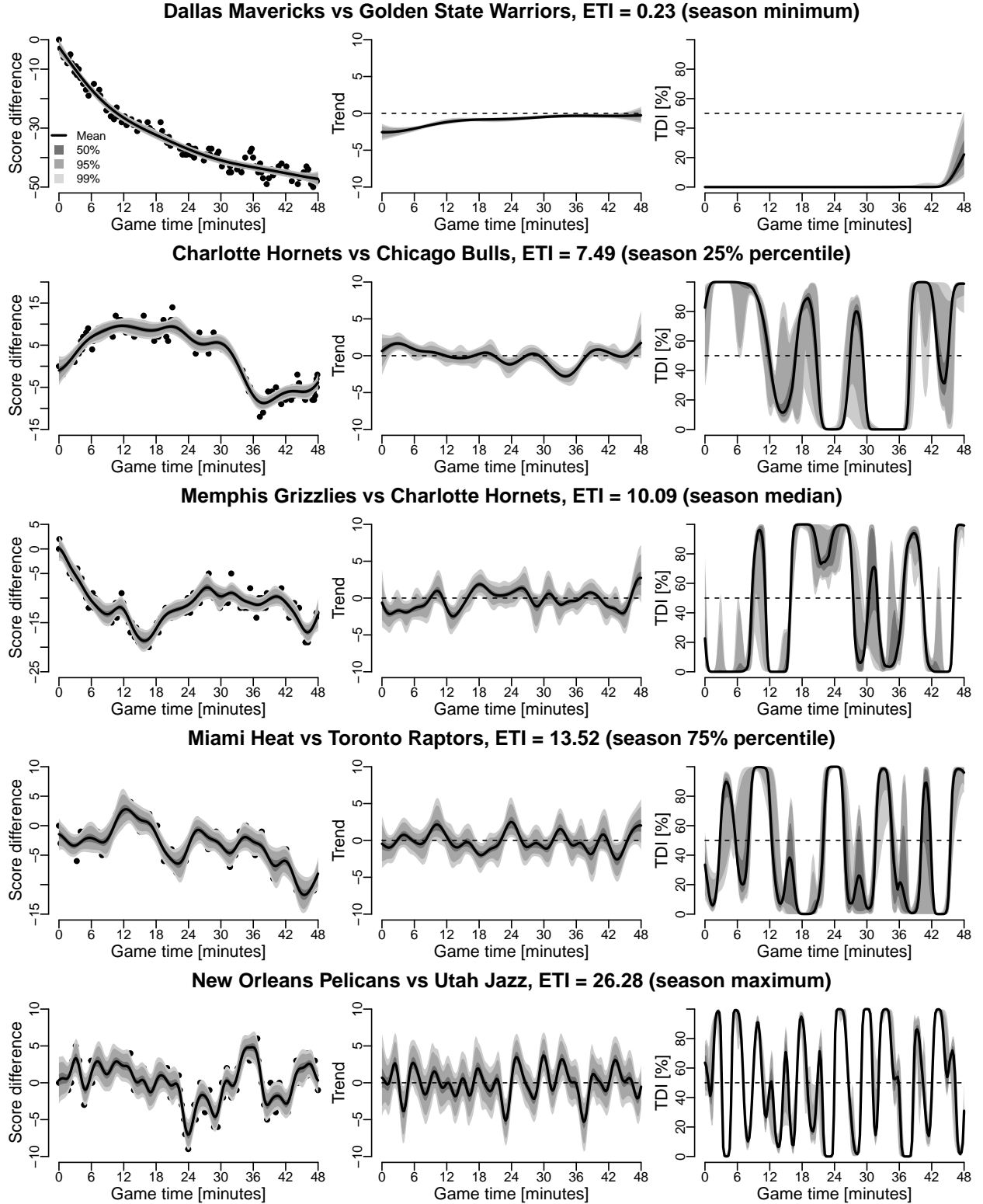


Figure 4: Observed score differences with posteriors of the latent processes (left panels), posterior trends (middle panels) and posterior Trend Direction Indices (right panels) for five games from the NBA season 2019–2020 with median posterior Excitement Trend Indices corresponding to the 0%, 25%, 50%, 75%, and 100% percentiles of the distribution of all games in the season. Gray regions depict 50%, 95%, and 99% point-wise credible intervals.