

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

14 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 increased attention within statistics and not just for match prediction or betting but also for game evaluation, coaching purposes, and for setting strategies and tactics in future matches.

Overall result of individual matches

We introduce 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.

- Excitement defineret som skift af hvem, der er i føring
- Vurdering af om et hold trender lige pt.
- Identificere hvilket “hot periods” et hold har i løbet af en kamp til efterfølgende evaluering

Reference to Quantifying the Trendiness of Trends Jensen and Ekstrøm (2020b).

We need to reference and discuss the following papers

- Chen, Dawson, and Müller (2020)
- Chen and Fan (2018)
- Gabel and Redner (2012)

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).

The paper is structured as follows. In the next we introduce the Gaussian process that can be .. and derive the TDI and ETI.

In Section 3 we apply the our proposed methodology to analyse all 1143 matches from the 2019–2020 National Basketball Association (NBA) season. We conclude with a discussion in Section 4.

Materials to reproduce this manuscript and its analyses can be found at Jensen and Ekstrøm (2020a).

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 $\mathcal{D}_m = (t_{m_i}, D_{m_i})_{0 < i \leq J_m}$ where $t_{m_1} < t_{m_2} < t_{m_i} < \dots < t_{m_{J_m}}$ is the ordered time points at which any teams scores, $D_{m_i} = D_m(t_{m_i})$ is the associated difference in scores at time t_{m_i} , and J_m is the total number of goals during the match. We use the convention that D_m is the difference in scores of the home team with respect to the the away team, so that $D_m > 0$ means that the home team is leading. [??? check that this is correct in the grønthøster]

We assume that the observed match data are a noisy realization of a latent smooth, random function defined in continuous time evaluated at random time points where goals occur. Let d_m be the latent functions from which the realizations \mathcal{D}_m are generated. We then propose the following hierarchical model where d_m is a Gaussian process and the observed data conditional on the scoring times and the latent process are independently normally distributed random variables with a match specific variance:

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

where $\Theta_m = (\beta_m, \theta_m, \sigma_m^2)$ is a vector of hyper-parameters governing the dynamics of the latent Gaussian process with a prior distribution H indexed by parameters Ψ , and $\mathbf{t}_m = (t_{m_1}, \dots, t_{m_{J_m}})$ is the vector of time points where goals occurs in the match.

By properties of Gaussian processes (see e.g., Cramer and Leadbetter (1967)) the latent process d_m and its time derivatives (assuming they exist) are distributed as a multivariate Gaussian process. We hence have that

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

where ' and '' denotes the first and second time derivatives, ∂_j^k is the k 'th order partial derivative with respect to the j 'th variable

[??? TALK ABOUT POSTERIOR HERE ???]

The Excitement Trend Index (ETI) of a particular match, denoted ETI_m , is defined as the expected number of changes in monotonicity of the score differences d_m conditional on the observed data from that match \mathcal{D}_m . This is equivalent to value of the expected number of zero-crossings of the posterior distribution of d'_m .

Formally,

$$\text{ETI}_m(\Theta_m) = \mathbb{E} [\# \{t \in \mathcal{I}_m : df_m(t) = 0\} \mid \mathbf{D}_m, \mathbf{t}_m, \Theta_m]$$

where \mathcal{I}_m is the interval of the time duration of a match i.e., $\mathcal{I}_m = [0; 48]$ minutes without overtime. The ETI is given by the integral of the local Excitement Trend Index

$$\text{ETI}_m(\Theta_m) = \int_{\mathcal{I}_m} d\text{ETI}_m(t | \Theta_m) dt$$

where $d\text{ETI}$ is the local Excitement Trend Index given by

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

and $\phi: x \mapsto 2^{-1/2} \pi^{-1/2} \exp(-\frac{1}{2}x^2)$ is the standard normal density function, $\text{Erf}: x \mapsto 2\pi^{-1/2} \int_0^x \exp(-u^2) du$ is the error function, and λ , ω and ζ are functions defined as

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

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

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.

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

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

$$G(\Theta_m | \Psi_m, \mathbf{t}_m) = G(\beta_m | \Psi_{\beta_m}) G(\alpha_m | \Psi_{\alpha_m}) G(\rho_m | \Psi_{\rho_m}) G(\sigma_m | \Psi_{\sigma_m})$$

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

$$\beta_m \sim T\left(\widehat{\beta}_m^{\text{ML}}, 3, 3\right), \quad \alpha_m \sim T^+\left(\widehat{\alpha}_m^{\text{ML}}, 3, 3\right), \quad \rho_m \sim N^+\left(\widehat{\rho}_m^{\text{ML}}, 1\right), \quad \sigma_m \sim T^+\left(\widehat{\sigma}_m^{\text{ML}}, 3, 3\right)$$

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

We use data from Sports Reference LLC (2020).

General idea: We estimate ETI_m for all matches in a given season and make a nice plot of the distribution of ETI_m . Then we can rank the matches according to increasing ETI and show the running score difference for e.g., the lowest, median and highest ranked matches. Maybe a large forest plot of ETI_m would look impressive.

Given the posteriors ETI_m we can summarize them by a posterior mean and variance and then do a **meta analysis** where we adjust for game-specific fixed effects such as number of spectators, location, stratify by season and so on.

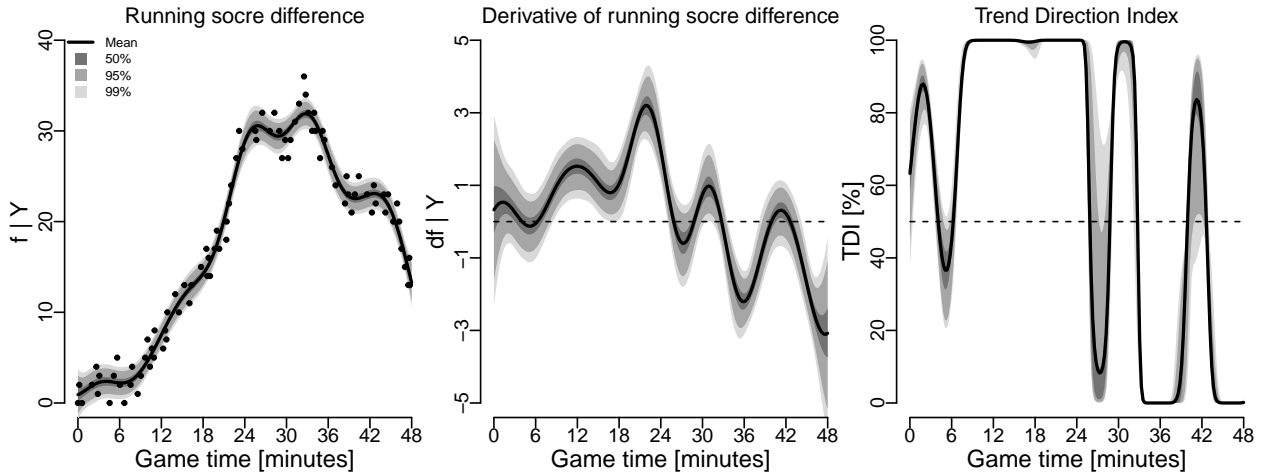


Figure 1: Los Angeles Lakers at Miami Heat. October 11, 2020. ??? Maybe this figure is superfluous in the results sections. We could instead use the first two panels as motivation in the introduction ???

4 Discussion

Some summarization goes on here. We have introduces 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

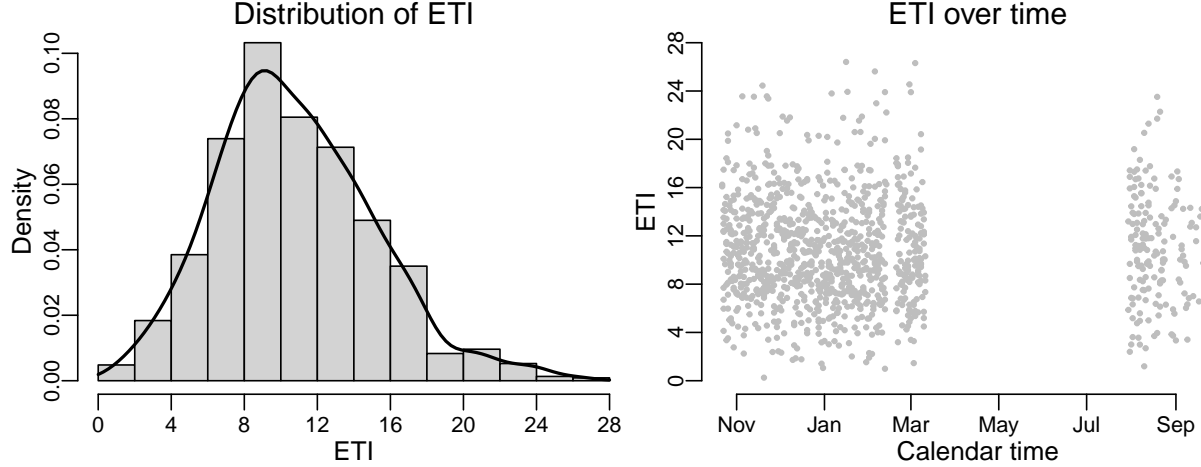


Figure 2: Caption

team is already far away of the other team. This motivates a weighted ETI of the form

$$\text{ETI}_m^W = \int_{\mathcal{I}_m} w(t, |f_m(t)|) d\text{ETI}_m(t)$$

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 ($\text{ETI}_{am}, \text{ETI}_{bm}$). 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

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

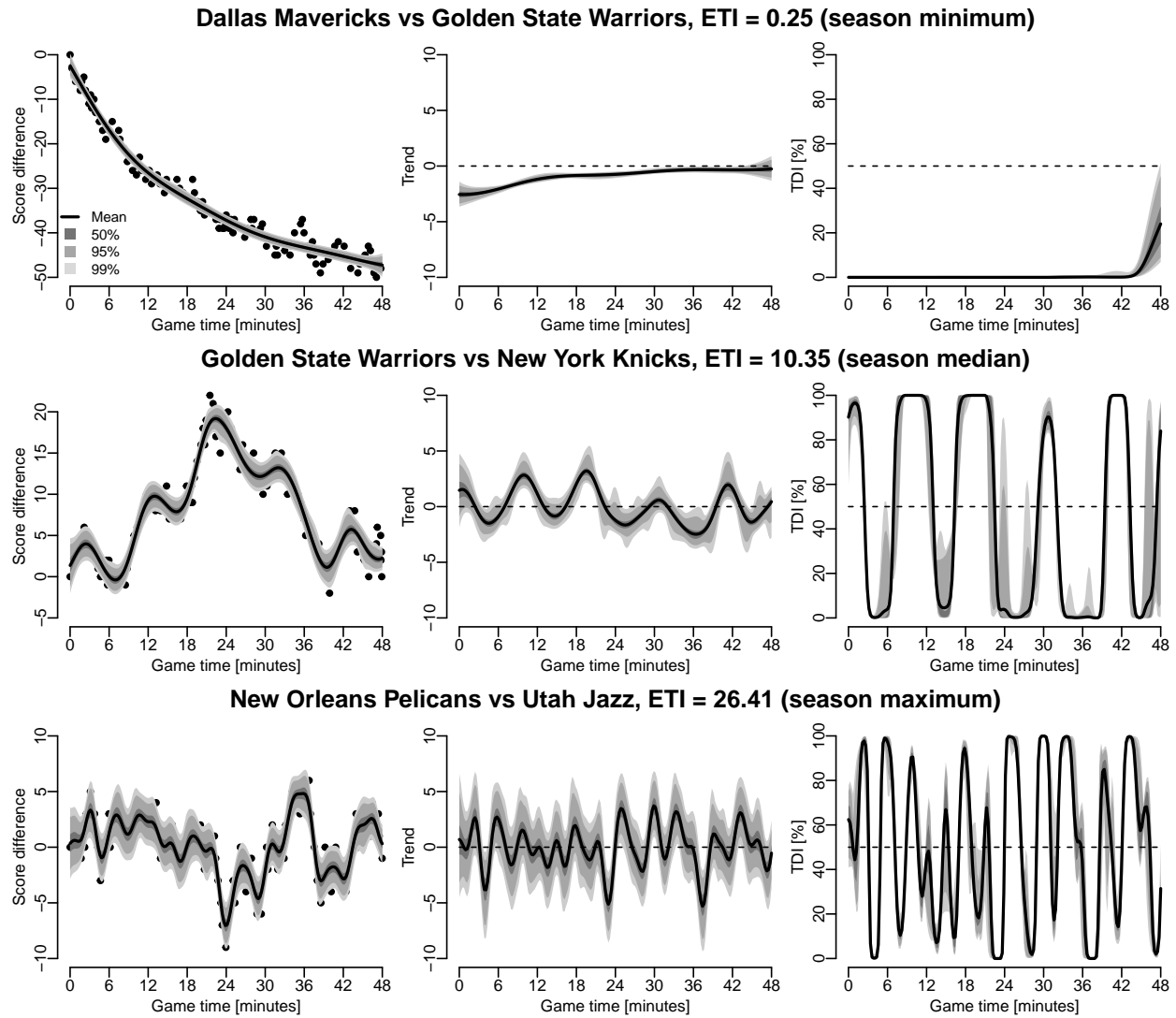


Figure 3: Caption ??? should we make this a nice full page figure and include 25% and 75% ETI games as well ???

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

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*.

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

Appendix

the joint distribution of (f, df, d^2f) conditional on \mathbf{Y} , \mathbf{t} and the hyper-parameters Θ evaluated at any finite vector \mathbf{t}^* of p time points is

$$\begin{bmatrix} f(\mathbf{t}^*) \\ df(\mathbf{t}^*) \\ d^2f(\mathbf{t}^*) \end{bmatrix} \mid \mathbf{Y}, \mathbf{t}, \Theta \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where $\boldsymbol{\mu} \in \mathbb{R}^{3p}$ is the column vector of posterior expectations and $\boldsymbol{\Sigma} \in \mathbb{R}^{3p \times 3p}$ is the joint posterior covariance matrix. Partitioning these as

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_f(\mathbf{t}^* \mid \Theta) \\ \mu_{df}(\mathbf{t}^* \mid \Theta) \\ \mu_{d^2f}(\mathbf{t}^* \mid \Theta) \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_f(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) & \Sigma_{f,df}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) & \Sigma_{f,d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) \\ \Sigma_{f,df}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta)^T & \Sigma_{df}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) & \Sigma_{df,d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) \\ \Sigma_{f,d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta)^T & \Sigma_{df,d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta)^T & \Sigma_{d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) \end{bmatrix}$$

the individual components are given by

$$\begin{aligned} \mu_f(\mathbf{t}^* \mid \Theta) &= \mu_\beta(\mathbf{t}^*) + C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} (\mathbf{Y} - \mu_\beta(\mathbf{t})) \\ \mu_{df}(\mathbf{t}^* \mid \Theta) &= d\mu_\beta(\mathbf{t}^*) + \partial_1 C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} (\mathbf{Y} - \mu_\beta(\mathbf{t})) \\ \mu_{d^2f}(\mathbf{t}^* \mid \Theta) &= d^2\mu_\beta(\mathbf{t}^*) + \partial_1^2 C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} (\mathbf{Y} - \mu_\beta(\mathbf{t})) \\ \Sigma_f(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) &= C_\theta(\mathbf{t}^*, \mathbf{t}^*) - C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} C_\theta(\mathbf{t}, \mathbf{t}^*) \\ \Sigma_{df}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) &= \partial_1 \partial_2 C_\theta(\mathbf{t}^*, \mathbf{t}^*) - \partial_1 C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} \partial_2 C_\theta(\mathbf{t}, \mathbf{t}^*) \\ \Sigma_{d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) &= \partial_1^2 \partial_2^2 C_\theta(\mathbf{t}^*, \mathbf{t}^*) - \partial_1^2 C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} \partial_2^2 C_\theta(\mathbf{t}, \mathbf{t}^*) \\ \Sigma_{f,df}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) &= \partial_2 C_\theta(\mathbf{t}^*, \mathbf{t}^*) - C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} \partial_2 C_\theta(\mathbf{t}, \mathbf{t}^*) \\ \Sigma_{f,d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) &= \partial_2^2 C_\theta(\mathbf{t}^*, \mathbf{t}^*) - C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} \partial_2^2 C_\theta(\mathbf{t}, \mathbf{t}^*) \\ \Sigma_{df,d^2f}(\mathbf{t}^*, \mathbf{t}^* \mid \Theta) &= \partial_1 \partial_2^2 C_\theta(\mathbf{t}^*, \mathbf{t}^*) - \partial_1 C_\theta(\mathbf{t}^*, \mathbf{t}) \left(C_\theta(\mathbf{t}, \mathbf{t}) + \sigma^2 I \right)^{-1} \partial_2^2 C_\theta(\mathbf{t}, \mathbf{t}^*) \end{aligned}$$