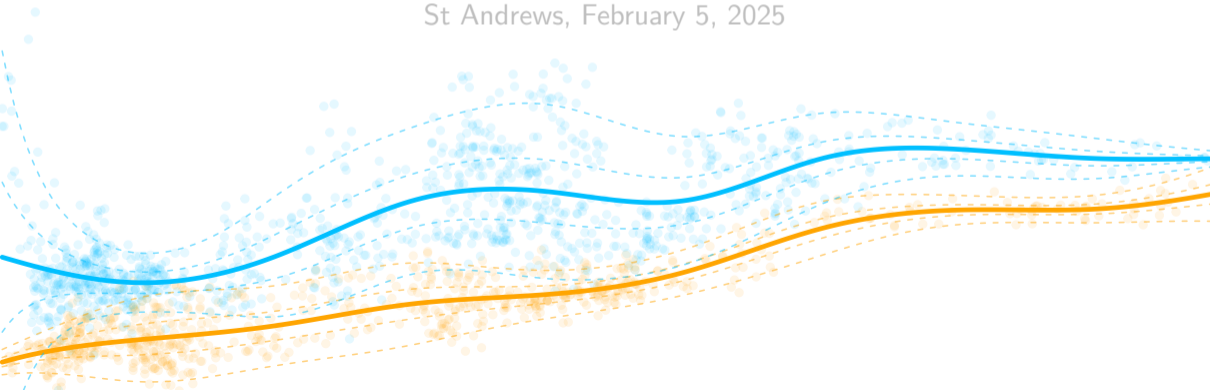


Statistics colloquium

# Efficient smoothness selection for nonparametric Markov-switching models via quasi restricted maximum likelihood

Jan-Ole Koslik

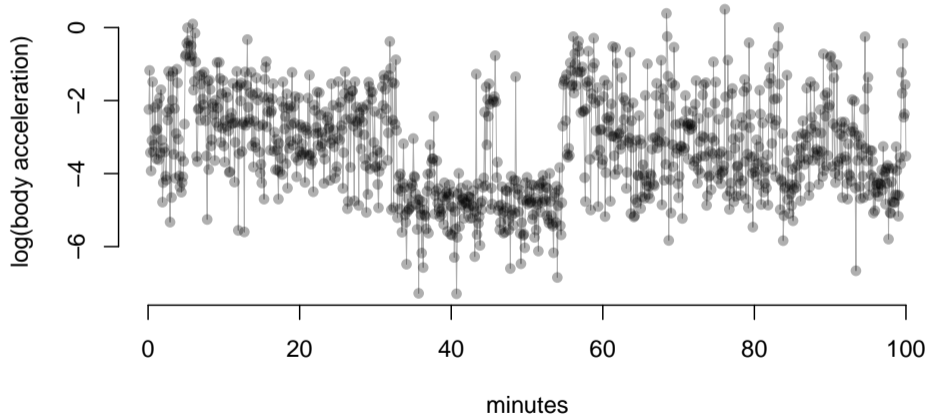
St Andrews, February 5, 2025



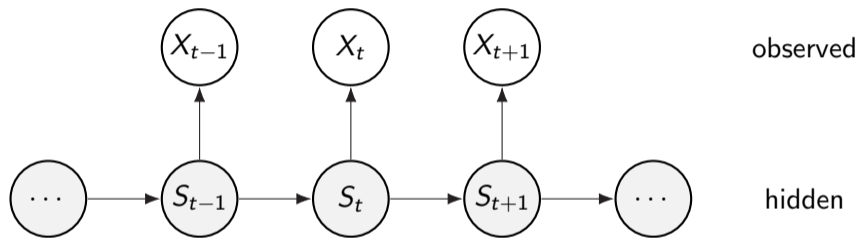


- PhD student in the **Statistics and Data Analysis Group** at Bielefeld University
- mostly working on **hidden Markov models** and their relatives

 [@olemole.bsky.social](https://bsky.app/profile/olemole.bsky.social)

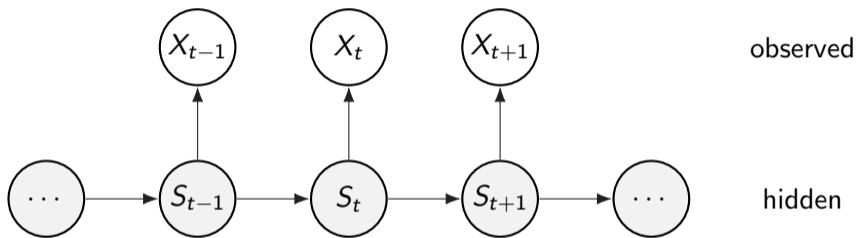


## HMM — model formulation



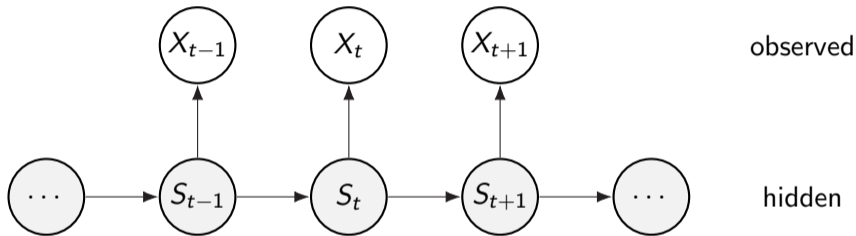
- every observation  $x_t$  is generated by one of  $N$  possible distributions  $f_1, \dots, f_N$

## HMM — model formulation

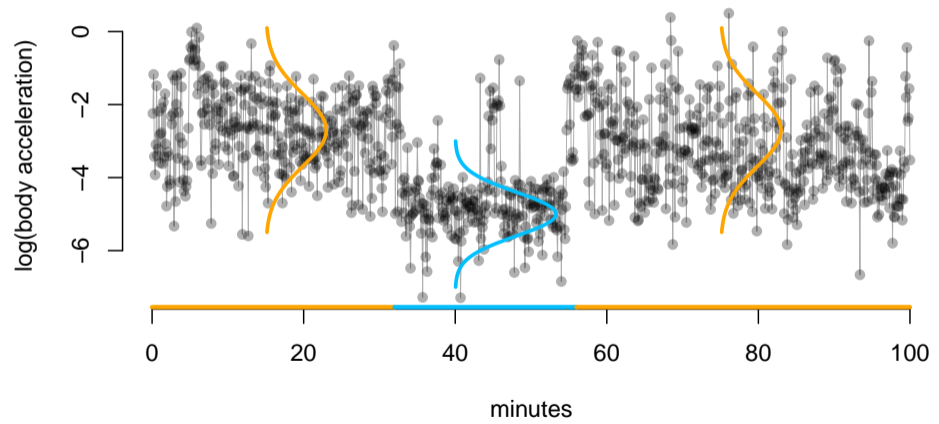


- every observation  $x_t$  is generated by one of  $N$  possible distributions  $f_1, \dots, f_N$
- hidden state process selects which distribution is active at time  $t$

## HMM — model formulation



- every observation  $x_t$  is generated by one of  $N$  possible distributions  $f_1, \dots, f_N$
- hidden state process selects which distribution is active at time  $t$
- state process is a **Markov chain**



## Reminder: Markov chains

**Markovian state process** is fully characterised by the initial distribution

$$\delta^{(1)} = (\Pr(S_1 = 1), \dots, \Pr(S_1 = N))$$

and the **transition probabilities**

$$\gamma_{ij}^{(t)} = \Pr(S_{t+1} = j \mid S_t = i),$$

which we summarise in the **transition probability matrix** (t.p.m.)

$$\mathbf{\Gamma}^{(t)} = (\gamma_{ij}^{(t)})_{i,j=1,\dots,N}.$$



## Estimating HMMs

We can efficiently calculate the HMM likelihood using the **forward algorithm**

$$\mathcal{L}(\theta) = \delta^{(1)} \mathbf{P}(x_1) \Gamma^{(1)} \mathbf{P}(x_2) \Gamma^{(2)} \cdot \dots \cdot \Gamma^{(T-1)} \mathbf{P}(x_T) \mathbf{1},$$

where  $\mathbf{P}(x_t) = \text{diag}(f_1(x_t), \dots, f_N(x_t))$ .

## Estimating HMMs

We can efficiently calculate the HMM likelihood using the **forward algorithm**

$$\mathcal{L}(\boldsymbol{\theta}) = \boldsymbol{\delta}^{(1)} \mathbf{P}(x_1) \boldsymbol{\Gamma}^{(1)} \mathbf{P}(x_2) \boldsymbol{\Gamma}^{(2)} \dots \boldsymbol{\Gamma}^{(T-1)} \mathbf{P}(x_T) \mathbf{1},$$

where  $\mathbf{P}(x_t) = \text{diag}(f_1(x_t), \dots, f_N(x_t))$ .

With some adjustments, we can also calculate the **log-likelihood**  $\ell(\boldsymbol{\theta})$  to avoid numerical underflow  $\rightarrow$  optimise in `R` using standard optimisers like `nlm()` or `optim()`.

## Estimating HMMs

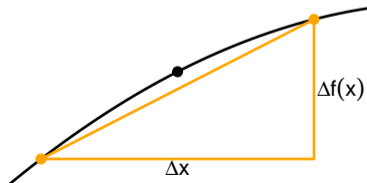
We can efficiently calculate the HMM likelihood using the **forward algorithm**

$$\mathcal{L}(\theta) = \delta^{(1)} \mathbf{P}(x_1) \Gamma^{(1)} \mathbf{P}(x_2) \Gamma^{(2)} \dots \Gamma^{(T-1)} \mathbf{P}(x_T) \mathbf{1},$$

where  $\mathbf{P}(x_t) = \text{diag}(f_1(x_t), \dots, f_N(x_t))$ .

With some adjustments, we can also calculate the **log-likelihood**  $\ell(\theta)$  to avoid numerical underflow  $\rightarrow$  optimise in `R` using standard optimisers like `nlm()` or `optim()`.

These approximate the gradient via **finite differencing**.



## Why nonparametrics?

- component distributions typically selected from **parametric** family  
→ difficult as we can't do state-specific EDA

## Why nonparametrics?

- component distributions typically selected from **parametric** family  
→ difficult as we can't do state-specific EDA
- potential covariate effects typically modelled using **linear** predictors  
→ may result in us missing interesting relationships

## Why nonparametrics?



= complicated

→ often **substantial** lack of fit

## Why nonparametrics?

Misspecification will be compensated by more states but this complicates interpretation.

# **Selecting the Number of States in Hidden Markov Models: Pragmatic Solutions Illustrated Using Animal Movement**

Jennifer POHLE, Roland LANGROCK, Floris M. van BEEST, and  
Niels Martin SCHMIDT

## Why nonparametrics?

Misspecification will be compensated by more states but this complicates interpretation.

# Selecting the Number of States in Hidden Markov Models: Pragmatic Solutions Illustrated Using Animal Movement

Jennifer POHLE, Roland LANGROCK, Floris M. van BEEST, and  
Niels Martin SCHMIDT

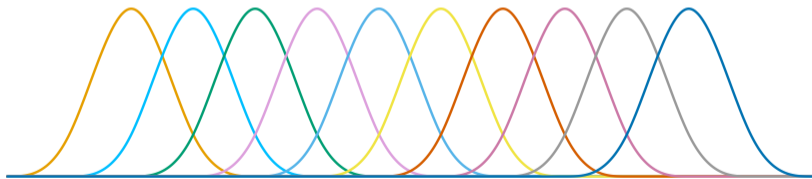
Obvious alternative: **nonparametric** approach using **penalised splines**



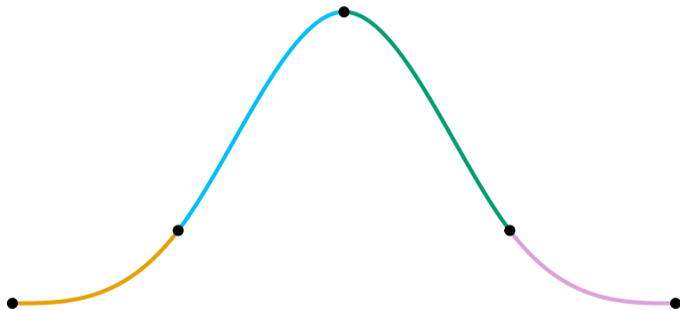
## An informal introduction to penalised splines

Idea: Perform **basis expansion** in  $x$  and represent smooth function  $s(x)$  as a linear combination of fixed basis functions  $B_k(x)$

$$s(x) = b_1 B_1(x) + b_2 B_2(x) + \dots + b_k B_k(x) = \mathbf{b}^T \mathbf{B}(x)$$

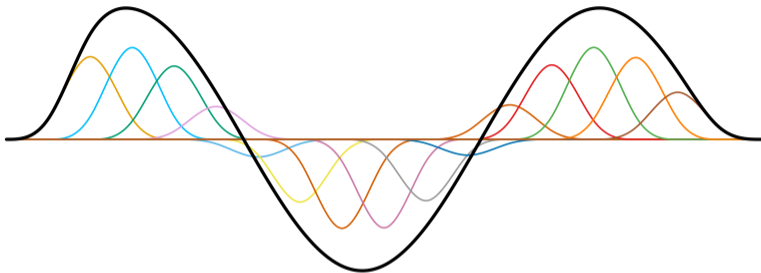


For example, when using **B-Splines**  $B_k(x)$  is a **piecewise** polynomial

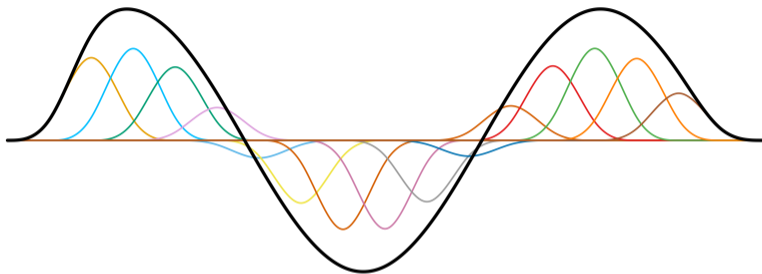


and zero outside the outer knots.

Approximate the true function with a sufficient number of basis functions:



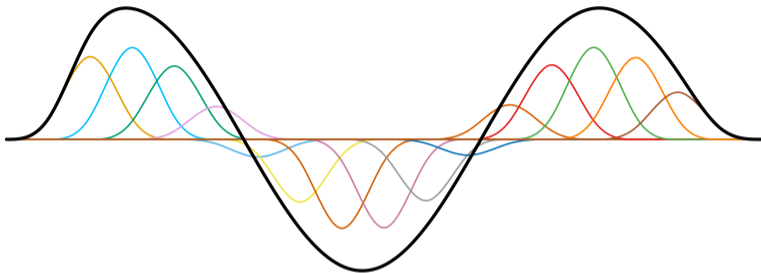
Approximate the true function with a sufficient number of basis functions:



As  $s$  should not be too **wiggly**, we add the penalty

$$\lambda \int s''(x)^2 dx$$

Approximate the true function with a sufficient number of basis functions:

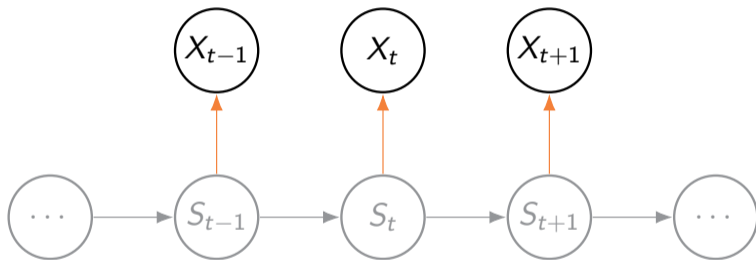


As  $s$  should not be too **wiggly**, we add the penalty

$$\lambda \int s''(x)^2 dx = \lambda \mathbf{b}^T \mathbf{S} \mathbf{b},$$

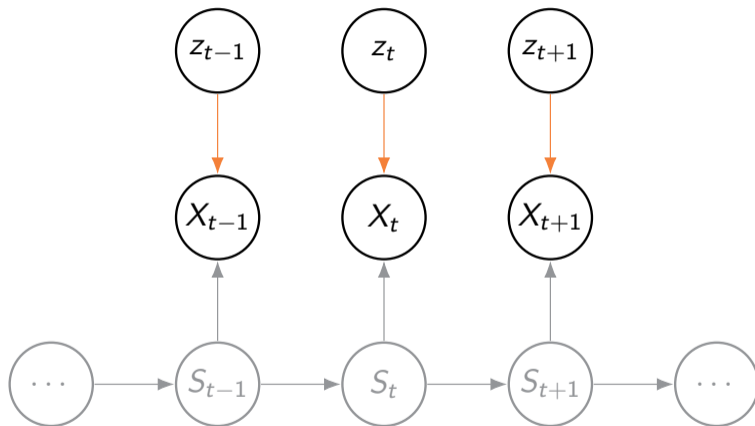
where  $\mathbf{S}$  is fixed penalty matrix with entries  $S_{ij} = \int B_i''(x) B_j''(x) dx$ .

## Nonparametric emission distributions



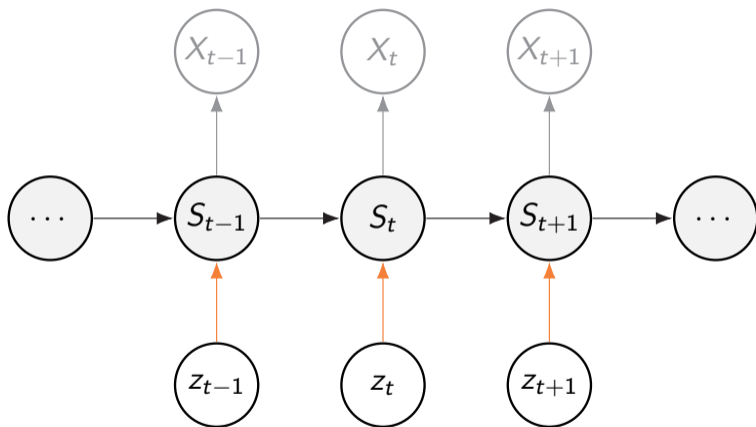
Langrock, Kneib, Sohn, and DeRuiter, 2015

## Markov-switching GAM



Langrock, Kneib, Glennie, and Michelot, 2017

## Smooth covariate effects on the state process



Feldmann et al., 2023



So this is where my spline adventure begins...

- fairly applied project:  $\gamma_{ij}^{(t)} \sim s(\text{time of day, day})$
- implementation with **fixed** penalisation straightforward



So this is where my spline adventure begins...

- fairly applied project:  $\gamma_{ij}^{(t)} \sim s(\text{time of day, day})$
- implementation with **fixed** penalisation straightforward



But how to find a good **penalty strength**??

# Smoothness selection for penalised splines

JOURNAL ARTICLE

## Nonparametric Inference in Hidden Markov Models Using P-Splines

Roland Langrock , [Thomas Kneib](#), Alexander Sohn, Stacy L. DeRuiter

*Biometrics*, Volume 71, Issue 2, June 2015, Pages 520–528,

<https://doi.org/10.1111/biom.12282>

Published: 13 January 2015 [Article history](#) ▼

## Markov-switching generalized additive models

Roland Langrock<sup>1</sup> · Thomas Kneib<sup>2</sup> · Richard Glennie<sup>3</sup> · Théo Michelot<sup>4</sup>



SPECIAL ISSUE ARTICLE |  Full Access

## Spline-based nonparametric inference in general state-switching models

Roland Langrock , Timo Adam, Vianey Leos-Barajas, Sina Mews, David L. Miller, Yannis P. Papastamatiou

## Smoothness selection for penalised splines

- cross-validation, AIC/BIC or subjective choice

## Smoothness selection for penalised splines

- cross-validation, AIC/BIC or subjective choice
- **grid search** → **curse of dimensionality**
- painfully slow and (for me) under-smoothed results

## Smoothness selection for penalised splines

- cross-validation, AIC/BIC or subjective choice
- **grid search** → **curse of dimensionality**
- painfully slow and (for me) under-smoothed results
- in general, underdeveloped, most work done for GLMMs and GAMs

## Smoothness selection for penalised splines

- cross-validation, AIC/BIC or subjective choice
- **grid search** → **curse of dimensionality**
- painfully slow and (for me) under-smoothed results
- in general, underdeveloped, most work done for GLMMs and GAMs

Ideally: define the **penalised log-likelihood** → **automatic** smoothness-selection

# Smoothness selection for penalised splines

Comput Stat (2012) 27:757–777  
DOI 10.1007/s00180-011-0289-6

---

ORIGINAL PAPER

## **Density estimation and comparison with a penalized mixture approach**

**Christian Schellhase · Göran Kauermann**

- splines as **random effects**?



# Smoothness selection for penalised splines

Comput Stat (2012) 27:757–777  
DOI 10.1007/s00180-011-0289-6

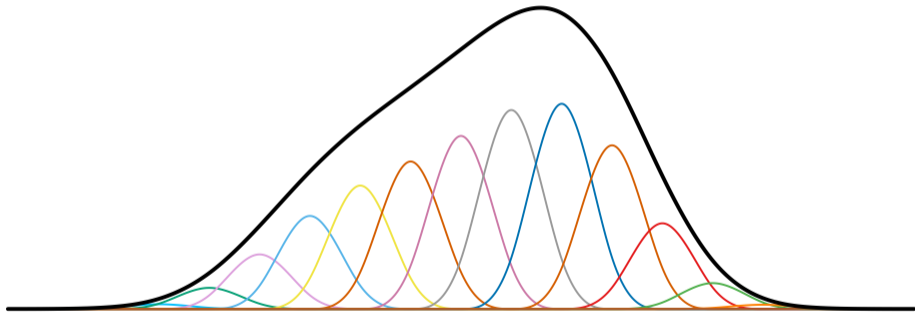
---

ORIGINAL PAPER

## Density estimation and comparison with a penalized mixture approach

Christian Schellhase · Göran Kauermann

- splines as **random effects**?
- Laplace approximation?



## Why are splines random effects?

Simple setting involving one penalised spline and no unpenalised parameters:

$$\ell_p(\mathbf{b}; \lambda) = \ell(\mathbf{b}) - \frac{1}{2} \lambda \mathbf{b}^\top \mathbf{S} \mathbf{b},$$

- coefficients  $\mathbf{b} = (b_1, \dots, b_k)$
- fixed penalty matrix  $\mathbf{S}$
- penalty strength  $\lambda$

## Why are splines random effects?

Simple setting involving one penalised spline and no unpenalised parameters:

$$\ell_p(\mathbf{b}; \lambda) = \ell(\mathbf{b}) - \frac{1}{2}\lambda\mathbf{b}^\top \mathbf{S}\mathbf{b},$$

- coefficients  $\mathbf{b} = (b_1, \dots, b_k)$
- fixed penalty matrix  $\mathbf{S}$
- penalty strength  $\lambda$

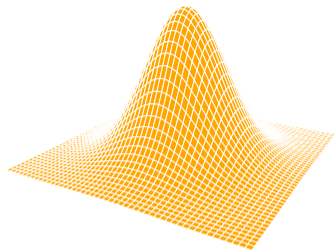
Penalised likelihood function

$$\mathcal{L}_p(\mathbf{b}; \lambda) = \mathcal{L}(\mathbf{b}) \cdot \exp\left(-\frac{1}{2}\lambda\mathbf{b}^\top \mathbf{S}\mathbf{b}\right)$$

We see that

$$\exp\left(-\frac{1}{2}\lambda\mathbf{b}^T\mathbf{S}\mathbf{b}\right)$$

is proportional to a **multivariate Gaussian density**.



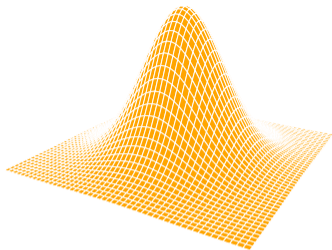
---

<sup>1</sup>ignoring potential issues with invertibility here

We see that

$$\exp\left(-\frac{1}{2}\lambda\mathbf{b}^\top\mathbf{S}\mathbf{b}\right)$$

is proportional to a **multivariate Gaussian density**.



We might as well assume  $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{S}^{-1}/\lambda)$  and add the missing **normalisation constant**<sup>1</sup>

$$\mathcal{L}_j(\mathbf{b}; \lambda) = \mathcal{L}(\mathbf{b}) \cdot (2\pi)^{-k/2} \det(\lambda\mathbf{S})^{1/2} \exp\left(-\frac{1}{2}\lambda\mathbf{b}^\top\mathbf{S}\mathbf{b}\right)$$

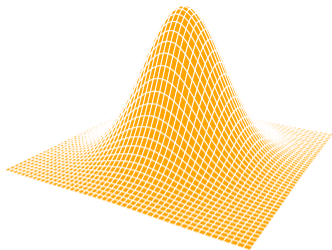
---

<sup>1</sup>ignoring potential issues with invertibility here

We see that

$$\exp\left(-\frac{1}{2}\lambda\mathbf{b}^\top\mathbf{S}\mathbf{b}\right)$$

is proportional to a **multivariate Gaussian density**.



We might as well assume  $\mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{S}^{-1}/\lambda)$  and add the missing **normalisation constant**<sup>1</sup>

$$\mathcal{L}_j(\mathbf{b}; \lambda) = \mathcal{L}(\mathbf{b}) \cdot (2\pi)^{-k/2} \det(\lambda\mathbf{S})^{1/2} \exp\left(-\frac{1}{2}\lambda\mathbf{b}^\top\mathbf{S}\mathbf{b}\right)$$

giving the joint log-likelihood

$$\ell_j(\mathbf{b}; \lambda) = \ell(\mathbf{b}) - \frac{k}{2} \log(2\pi) + \frac{1}{2} \log \det(\lambda\mathbf{S}) - \frac{1}{2} \lambda\mathbf{b}^\top\mathbf{S}\mathbf{b}$$

---

<sup>1</sup>ignoring potential issues with invertibility here

## How to estimate models with random effects?

**Joint likelihood** of the **data** and the **random effect** as a function of  $\lambda$

$$f_{\lambda}(\mathbf{x}, \mathbf{b}) = f(\mathbf{x} \mid \mathbf{b}) \cdot f_{\lambda}(\mathbf{b}),$$

with  $f(\mathbf{x} \mid \mathbf{b}) = \mathcal{L}(\mathbf{b})$  and  $f_{\lambda}(\mathbf{b}) \sim \mathcal{N}(\mathbf{0}, \mathbf{S}^{-1}/\lambda)$ .



## How to estimate models with random effects?

**Joint likelihood** of the **data** and the **random effect** as a function of  $\lambda$

$$f_{\lambda}(\mathbf{x}, \mathbf{b}) = f(\mathbf{x} | \mathbf{b}) \cdot f_{\lambda}(\mathbf{b}),$$

with  $f(\mathbf{x} | \mathbf{b}) = \mathcal{L}(\mathbf{b})$  and  $f_{\lambda}(\mathbf{b}) \sim \mathcal{N}(\mathbf{0}, \mathbf{S}^{-1}/\lambda)$ .

Marginal likelihood of the data by law of total probability

$$\mathcal{L}(\lambda) = f_{\lambda}(\mathbf{x}) = \int f_{\lambda}(\mathbf{x}, \mathbf{b}) d\mathbf{b},$$

## How to estimate models with random effects?

**Joint likelihood** of the **data** and the **random effect** as a function of  $\lambda$

$$f_{\lambda}(\mathbf{x}, \mathbf{b}) = f(\mathbf{x} | \mathbf{b}) \cdot f_{\lambda}(\mathbf{b}),$$

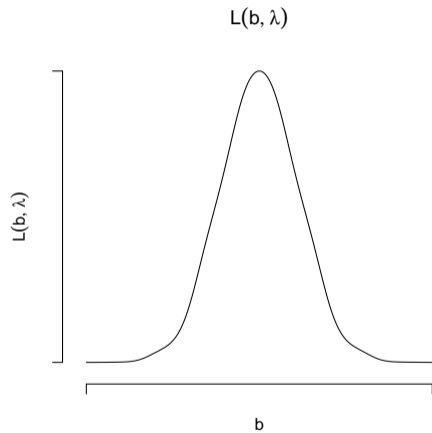
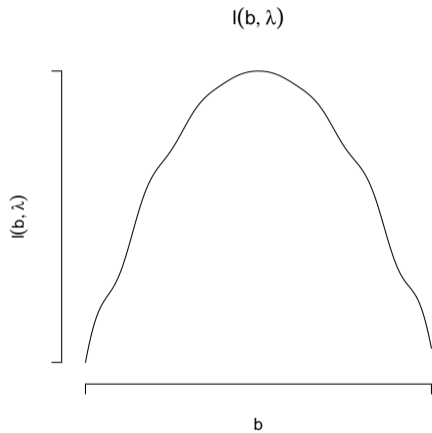
with  $f(\mathbf{x} | \mathbf{b}) = \mathcal{L}(\mathbf{b})$  and  $f_{\lambda}(\mathbf{b}) \sim \mathcal{N}(\mathbf{0}, \mathbf{S}^{-1}/\lambda)$ .

Marginal likelihood of the data by law of total probability

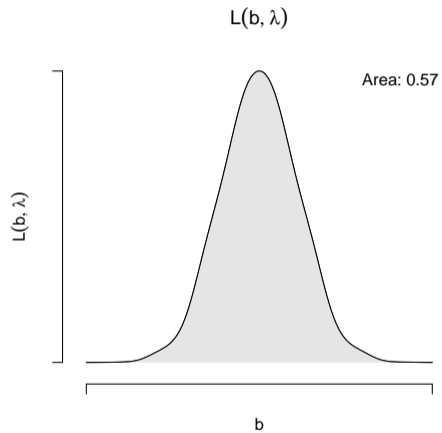
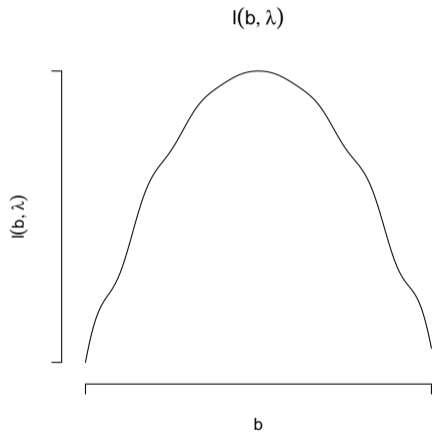
$$\mathcal{L}(\lambda) = f_{\lambda}(\mathbf{x}) = \int f_{\lambda}(\mathbf{x}, \mathbf{b}) d\mathbf{b},$$

which we would like to maximise to find the **MLE**  $\hat{\lambda}$ .

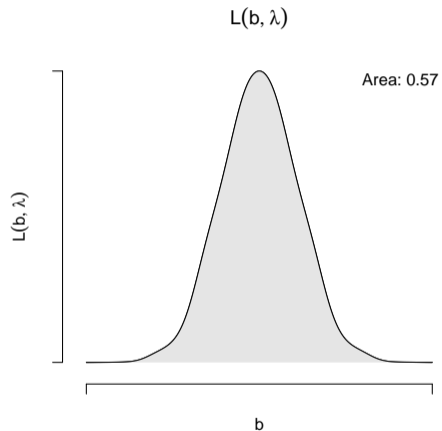
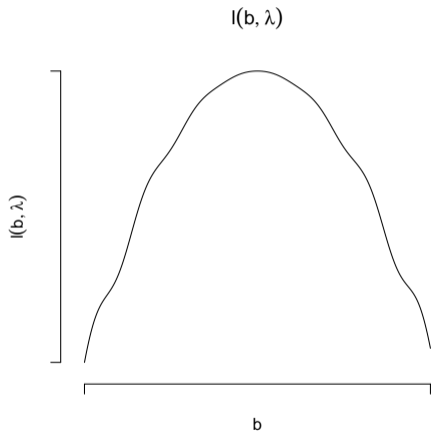
# Marginal ML



# Marginal ML



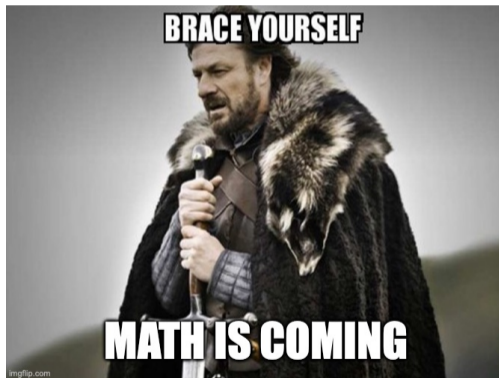
# Marginal ML



In reality, this integral is intractable!

What can we do?

What can we do?



## The Laplace approximation

- find the **mode** (in  $\mathbf{b}$ ) of the joint log-likelihood

$$\ell_j(\mathbf{b}, \lambda) = -\frac{k}{2} \log(2\pi) + \frac{1}{2} \log \det(\lambda S) + \ell(\mathbf{b}) - \frac{1}{2} \lambda \mathbf{b}^\top \mathbf{S} \mathbf{b}$$

by **penalised ML**



## The Laplace approximation

- find the **mode** (in  $\mathbf{b}$ ) of the joint log-likelihood

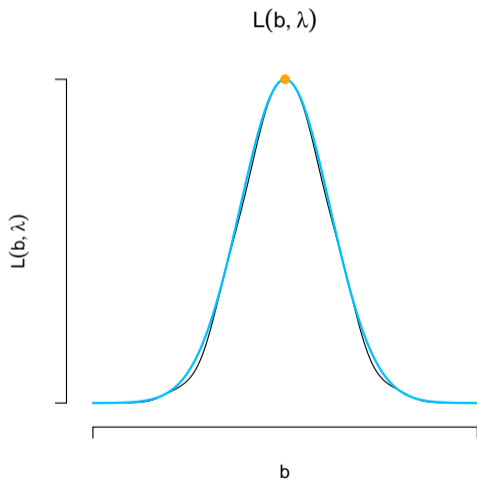
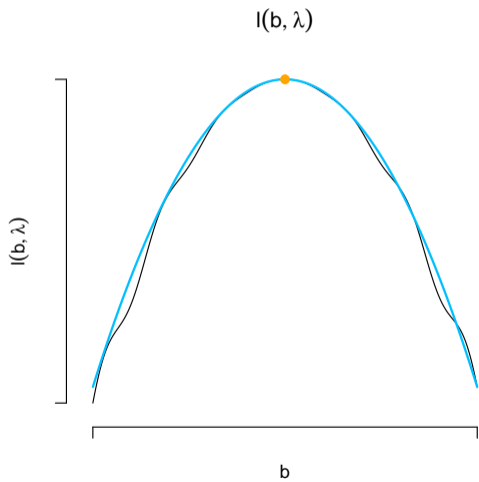
$$\ell_j(\mathbf{b}, \lambda) = -\frac{k}{2} \log(2\pi) + \frac{1}{2} \log \det(\lambda S) + \ell(\mathbf{b}) - \frac{1}{2} \lambda \mathbf{b}^\top \mathbf{S} \mathbf{b}$$

by **penalised ML**

- second-order **Taylor approximation** around the **mode**:

$$\ell_{approx}(\mathbf{b}, \lambda) = \ell_j(\hat{\mathbf{b}}, \lambda) - \frac{1}{2} (\mathbf{b} - \hat{\mathbf{b}})^\top \mathbf{J}(\lambda) (\mathbf{b} - \hat{\mathbf{b}}),$$

where  $\mathbf{J}(\lambda) = -\nabla^2 \ell_j(\hat{\mathbf{b}}, \lambda)$  is the (negative) Hessian of  $\ell_j(\mathbf{b}, \lambda)$  w.r.t.  $\mathbf{b}$  at  $\hat{\mathbf{b}}(\lambda)$ .



- now exponentiate to obtain likelihood and integrate out  $\mathbf{b}$

$$\exp(\ell_j(\hat{\mathbf{b}}, \lambda)) \cdot \int \exp\left(-\frac{1}{2}(\mathbf{b} - \hat{\mathbf{b}})^\top \mathbf{J}(\lambda)(\mathbf{b} - \hat{\mathbf{b}})\right) d\mathbf{b}$$

- now exponentiate to obtain likelihood and integrate out  $\mathbf{b}$

$$\exp(\ell_j(\hat{\mathbf{b}}, \lambda)) \cdot \int \exp\left(-\frac{1}{2}(\mathbf{b} - \hat{\mathbf{b}})^\top \mathbf{J}(\lambda)(\mathbf{b} - \hat{\mathbf{b}})\right) d\mathbf{b}$$

- right side is a Gaussian integral and equals

$$(2\pi)^{k/2} \cdot \det(\mathbf{J}(\lambda))^{-1/2}$$

- now exponentiate to obtain likelihood and integrate out  $\mathbf{b}$

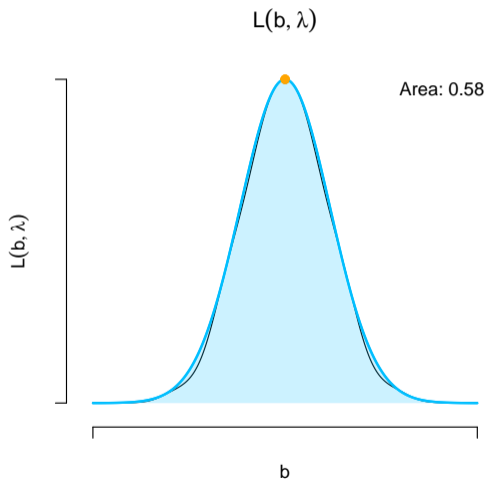
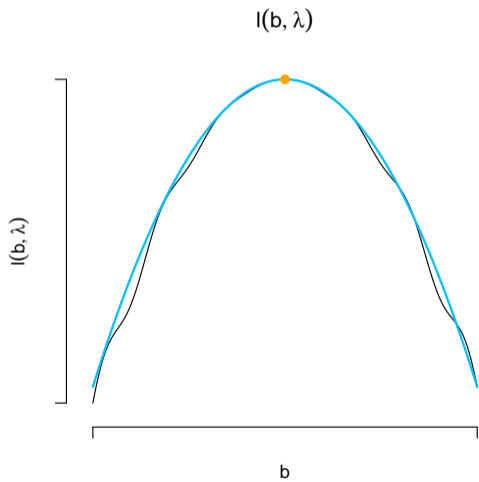
$$\exp(\ell_j(\hat{\mathbf{b}}, \lambda)) \cdot \int \exp\left(-\frac{1}{2}(\mathbf{b} - \hat{\mathbf{b}})^\top \mathbf{J}(\lambda)(\mathbf{b} - \hat{\mathbf{b}})\right) d\mathbf{b}$$

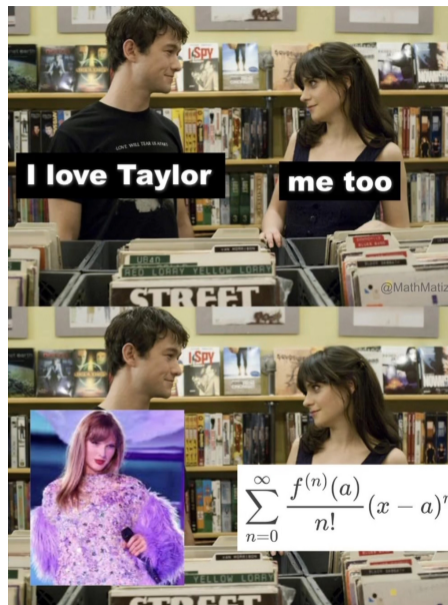
- right side is a Gaussian integral and equals

$$(2\pi)^{k/2} \cdot \det(\mathbf{J}(\lambda))^{-1/2}$$

- hence, approximate marginal log-likelihood becomes

$$\ell(\lambda) \approx \ell_j(\hat{\mathbf{b}}, \lambda) + \frac{k}{2} \log(2\pi) - \frac{1}{2} \log \det(\mathbf{J}(\lambda))$$





- typically **good approximation** because **large values** dominate the integral



- typically **good approximation** because **large values** dominate the integral
- shape of the likelihood becomes more and more **Gaussian** as  $n \rightarrow \infty$

- typically **good approximation** because **large values** dominate the integral
- shape of the likelihood becomes more and more **Gaussian** as  $n \rightarrow \infty$
- **but** each evaluation of marginal log-likelihood requires **inner optimisation** w.r.t.  **$b$**   
→ leads to nested optimisation in general

## Optimising the marginal likelihood

Schellhase and Kauermann (2012) start with the (approximate) marginal log-likelihood

$$\ell(\lambda) \approx \ell_j(\hat{\mathbf{b}}, \lambda) + \frac{k}{2} \log(2\pi) - \frac{1}{2} \log \det(\mathbf{J}(\lambda)).$$

## Optimising the marginal likelihood

Schellhase and Kauermann (2012) start with the (approximate) marginal log-likelihood

$$\ell(\lambda) \approx \ell_j(\hat{\mathbf{b}}, \lambda) + \frac{k}{2} \log(2\pi) - \frac{1}{2} \log \det(\mathbf{J}(\lambda)).$$

Writing out our joint log-likelihood, we have as our marginal log-likelihood

$$\ell(\lambda) \approx -\frac{k}{2} \log(2\pi) + \frac{1}{2} \log \det(\lambda \mathbf{S}) + \ell(\hat{\mathbf{b}}) - \frac{1}{2} \lambda \hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}} + \frac{k}{2} \log(2\pi) - \frac{1}{2} \log \det(\mathbf{J}(\lambda)),$$

## Optimising the marginal likelihood

Schellhase and Kauermann (2012) start with the (approximate) marginal log-likelihood

$$\ell(\lambda) \approx \ell_j(\hat{\mathbf{b}}, \lambda) + \frac{k}{2} \log(2\pi) - \frac{1}{2} \log \det(\mathbf{J}(\lambda)).$$

Writing out our joint log-likelihood, we have as our marginal log-likelihood

$$\ell(\lambda) \approx -\frac{k}{2} \log(2\pi) + \frac{1}{2} \log \det(\lambda \mathbf{S}) + \ell(\hat{\mathbf{b}}) - \frac{1}{2} \lambda \hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}} + \frac{k}{2} \log(2\pi) - \frac{1}{2} \log \det(\mathbf{J}(\lambda)),$$

which we can now (partially) differentiate w.r.t.  $\lambda$

## Optimising the marginal likelihood

$$\frac{\partial}{\partial \lambda} \left( \frac{1}{2} \log \det(\lambda S) \right) = \frac{\partial}{\partial \lambda} \left( \frac{1}{2} \log(\lambda^k \det(S)) \right) = \frac{\partial}{\partial \lambda} \left( \frac{k}{2} \log(\lambda) + \frac{1}{2} \log \det(S) \right) = \frac{k}{2\lambda}$$

$$\frac{\partial}{\partial \lambda} \left( -\frac{1}{2} \lambda \hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}} \right) = -\frac{1}{2} \hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}}$$

$$\frac{\partial}{\partial \lambda} \left( -\frac{1}{2} \log \det(\mathbf{J}(\lambda)) \right) = -\frac{1}{2} \text{tr}(\mathbf{J}(\lambda)^{-1} S)$$

## Optimising the marginal likelihood

Hence in total

$$\frac{\partial \ell(\lambda)}{\partial \lambda} \approx -\frac{1}{2} \hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}} + \frac{k}{2\lambda} - \frac{1}{2} \text{tr}(\mathbf{J}(\lambda)^{-1} \mathbf{S}),$$

## Optimising the marginal likelihood

Hence in total

$$\frac{\partial \ell(\lambda)}{\partial \lambda} \approx -\frac{1}{2} \hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}} + \frac{k}{2\lambda} - \frac{1}{2} \text{tr}(\mathbf{J}(\lambda)^{-1} \mathbf{S}),$$

from which we can construct the estimating equation (omitting the details)

$$\lambda = \frac{\text{tr}(\mathbf{J}(\lambda)^{-1} \mathbf{J}(\lambda = 0))}{\hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}}} = \frac{\text{dof}(\lambda)}{\hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}}}.$$



## Optimising the marginal likelihood

Hence in total

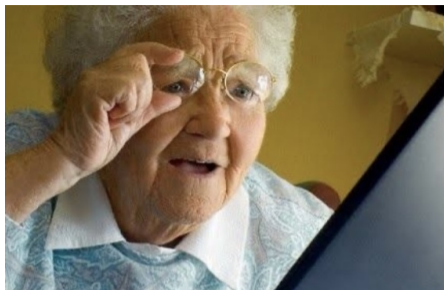
$$\frac{\partial \ell(\lambda)}{\partial \lambda} \approx -\frac{1}{2} \hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}} + \frac{k}{2\lambda} - \frac{1}{2} \text{tr}(\mathbf{J}(\lambda)^{-1} \mathbf{S}),$$

from which we can construct the estimating equation (omitting the details)

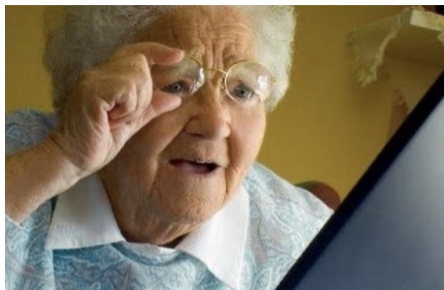
$$\lambda = \frac{\text{tr}(\mathbf{J}(\lambda)^{-1} \mathbf{J}(\lambda = 0))}{\hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}}} = \frac{\text{dof}(\lambda)}{\hat{\mathbf{b}}^\top \mathbf{S} \hat{\mathbf{b}}}.$$

which yields the iterative procedure:

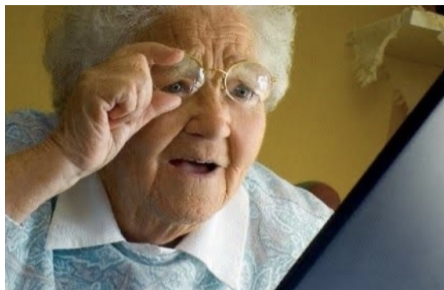
1. fit model via **penalised ML**
2. calculate **Hessian** at optimum
3. **update** penalty strength
4. repeat 1.-3. until convergence



- first try  $\rightarrow$  worked much better than CV or AIC/ BIC, requiring only a few iterations for convergence ( $\sim 20$  compared to hundreds or thousands for grid search)



- first try  $\rightarrow$  worked much better than CV or AIC/ BIC, requiring only a few iterations for convergence ( $\sim 20$  compared to hundreds or thousands for grid search)
- but my situation is more complicated:
  - fixed effects
  - multiple penalties



- first try  $\rightarrow$  worked much better than CV or AIC/ BIC, requiring only a few iterations for convergence ( $\sim 20$  compared to hundreds or thousands for grid search)
- but my situation is more complicated:
  - fixed effects
  - multiple penalties
- So why did it work??

## The full setting

$$\ell_p(\mathbf{a}, \mathbf{b}; \boldsymbol{\lambda}) = \ell(\mathbf{a}, \mathbf{b}) - \frac{1}{2} \sum_i \lambda_i \mathbf{b}_i^\top \mathbf{S} \mathbf{b}_i$$

- fixed effects  $\mathbf{a}$
- multiple random effects  $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_p)$

$$\mathbf{b}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{S}^{-1}/\lambda_i)$$

**Problem:** If we integrate out  $\mathbf{b}$ , marginal likelihood is more complicated because of  $\mathbf{a}$ .

## The full setting

$$\ell_p(\mathbf{a}, \mathbf{b}; \boldsymbol{\lambda}) = \ell(\mathbf{a}, \mathbf{b}) - \frac{1}{2} \sum_i \lambda_i \mathbf{b}_i^\top \mathbf{S} \mathbf{b}_i$$

- fixed effects  $\mathbf{a}$
- multiple random effects  $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_p)$

$$\mathbf{b}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{S}^{-1}/\lambda_i)$$

**Problem:** If we integrate out  $\mathbf{b}$ , marginal likelihood is more complicated because of  $\mathbf{a}$ .

**Solution:** Assume  $\mathbf{a} \sim \mathcal{N}(\mathbf{0}, \infty)$  and integrate out both  $\mathbf{a}$  and  $\mathbf{b} \rightarrow$  restricted maximum likelihood (REML), Laird and Ware, 1982

## Quasi REML

Then, marginal log-likelihood is a function of  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$ .

## Quasi REML

Then, marginal log-likelihood is a function of  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$ .

Partially differentiating w.r.t.  $\lambda_i$  yields very similar result:

$$\lambda_i = \frac{\text{tr}([\mathbf{J}(\boldsymbol{\lambda})^{-1}\mathbf{J}(\boldsymbol{\lambda} = 0)]_{ii})}{\hat{\mathbf{b}}_i^\top \mathbf{S} \hat{\mathbf{b}}_i} = \frac{\text{dof}(\lambda_i)}{\hat{\mathbf{b}}_i^\top \mathbf{S} \hat{\mathbf{b}}_i}$$



## Quasi REML

Then, marginal log-likelihood is a function of  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$ .

Partially differentiating w.r.t.  $\lambda_i$  yields very similar result:

$$\lambda_i = \frac{\text{tr}([\mathbf{J}(\lambda)^{-1}\mathbf{J}(\lambda = 0)]_{ii})}{\hat{\mathbf{b}}_i^\top \mathbf{S} \hat{\mathbf{b}}_i} = \frac{\text{dof}(\lambda_i)}{\hat{\mathbf{b}}_i^\top \mathbf{S} \hat{\mathbf{b}}_i}$$

Yields the same iterative procedure: model fitting via **penalised ML** and **updating** penalty strength.

## Quasi REML

Then, marginal log-likelihood is a function of  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$ .

Partially differentiating w.r.t.  $\lambda_i$  yields very similar result:

$$\lambda_i = \frac{\text{tr}([\mathbf{J}(\boldsymbol{\lambda})^{-1}\mathbf{J}(\boldsymbol{\lambda} = 0)]_{ii})}{\hat{\mathbf{b}}_i^\top \mathbf{S} \hat{\mathbf{b}}_i} = \frac{\text{dof}(\lambda_i)}{\hat{\mathbf{b}}_i^\top \mathbf{S} \hat{\mathbf{b}}_i}$$

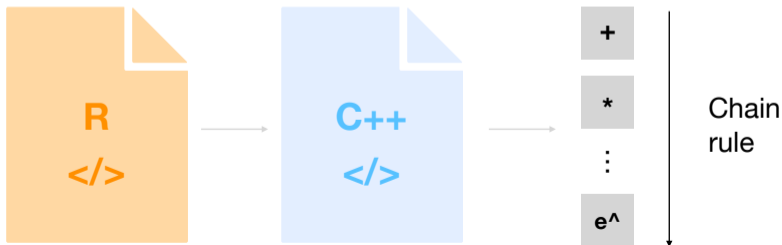
Yields the same iterative procedure: model fitting via **penalised ML** and **updating** penalty strength.

→ smoothness selection procedure that makes nonparametric HMMs feasible!

## RTMB enters the picture



Here comes `RTMB` (Kristensen, 2024) with **automatic differentiation**, natively supporting the full Laplace approximation for models written in plain `R` code



RTMB enters the picture

So did we gain anything?

## RTMB enters the picture

So did we gain anything?

- with `RTMB` (for inner optimisation), possible to implement `qrem1()` very generally

## RTMB enters the picture

So did we gain anything?

- with `RTMB` (for inner optimisation), possible to implement `qrem1()` very generally
- user only needs to specify penalised negative log-likelihood

## RTMB enters the picture

So did we gain anything?

- with `RTMB` (for inner optimisation), possible to implement `qreml()` very generally
- user only needs to specify penalised negative log-likelihood
- **qREML** + **AD** → efficiency skyrocketed!

## Practical usage

```
library(LaMa)

pnll <- function(par){ # penalised negative log-likelihood
  ... # computing the negative log-likelihood
  nll + penalty(splinePars, S, lambda)
}
```



## Practical usage

```
library(LaMa)

pnll <- function(par){ # penalised negative log-likelihood
  ... # computing the negative log-likelihood
  nll + penalty(splinePars, S, lambda)
}
```

```
par <- list(..., splinePars = matrix(0, 2, 10)) # parameter list
dat <- list(..., S = S, lambda = rep(100, 2)) # data list
```

## Practical usage

```
library(LaMa)
```

```
pnll <- function(par){ # penalised negative log-likelihood
  ... # computing the negative log-likelihood
  nll + penalty(splinePars, S, lambda)
}
```

```
par <- list(..., splinePars = matrix(0, 2, 10)) # parameter list
dat <- list(..., S = S, lambda = rep(100, 2)) # data list
```

```
mod <- qrem1(pnll, par, dat, random = "splinePars")
```

## Real-data example

## Bull sharks (Byrnes et al., 2023)

- Western Australia, extreme seasonal changes
- seven bull sharks tagged (14K observations)
- temperature, depth, and acceleration data
- response: overall dynamic body acceleration
- 2-state HMM: **low** and **high activity**



## Bull sharks (Byrnes et al., 2023)

- Western Australia, extreme seasonal changes
- seven bull sharks tagged (14K observations)
- temperature, depth, and acceleration data
- response: overall dynamic body acceleration
- 2-state HMM: **low** and **high activity**



- 
- $\text{ODBA}_k^{(t)} \mid \{S_t = i\} \sim \beta_k^{(i)} + s_k(\text{time}_t), k = 1, \dots, 7$  (ectotherms)

## Bull sharks (Byrnes et al., 2023)

- Western Australia, extreme seasonal changes
- seven bull sharks tagged (14K observations)
- temperature, depth, and acceleration data
- response: overall dynamic body acceleration
- 2-state HMM: **low** and **high activity**



- 
- $\text{ODBA}_k^{(t)} \mid \{S_t = i\} \sim \beta_k^{(i)} + s_k(\text{time}_t), k = 1, \dots, 7$  (ectotherms)
  - $\gamma_{ij}^{(t)} \sim s(\text{tod}_t) + \text{AvgTemp}_t * s(\text{tod}_t)$  (parametric in original paper)

## Bull sharks (Byrnes et al., 2023)

- Western Australia, extreme seasonal changes
- seven bull sharks tagged (14K observations)
- temperature, depth, and acceleration data
- response: overall dynamic body acceleration
- 2-state HMM: **low** and **high activity**



- 
- $\text{ODBA}_k^{(t)} \mid \{S_t = i\} \sim \beta_k^{(i)} + s_k(\text{time}_t), k = 1, \dots, 7$  (ectotherms)
  - $\gamma_{ij}^{(t)} \sim s(\text{tod}_t) + \text{AvgTemp}_t * s(\text{tod}_t)$  (parametric in original paper)
  - **11** smooths in total, **162** parameters

## Bull sharks (Byrnes et al., 2023)

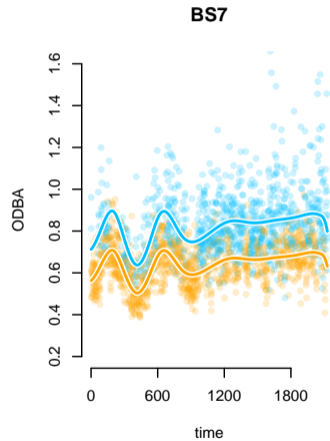
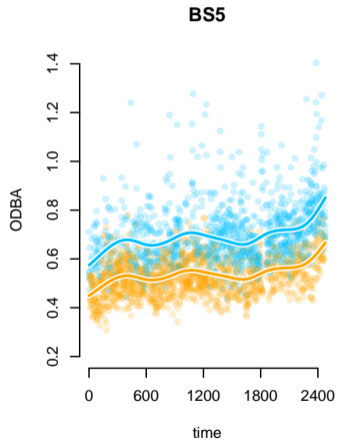
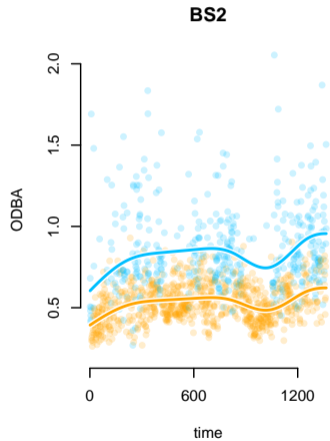
- Western Australia, extreme seasonal changes
- seven bull sharks tagged (14K observations)
- temperature, depth, and acceleration data
- response: overall dynamic body acceleration
- 2-state HMM: **low** and **high activity**



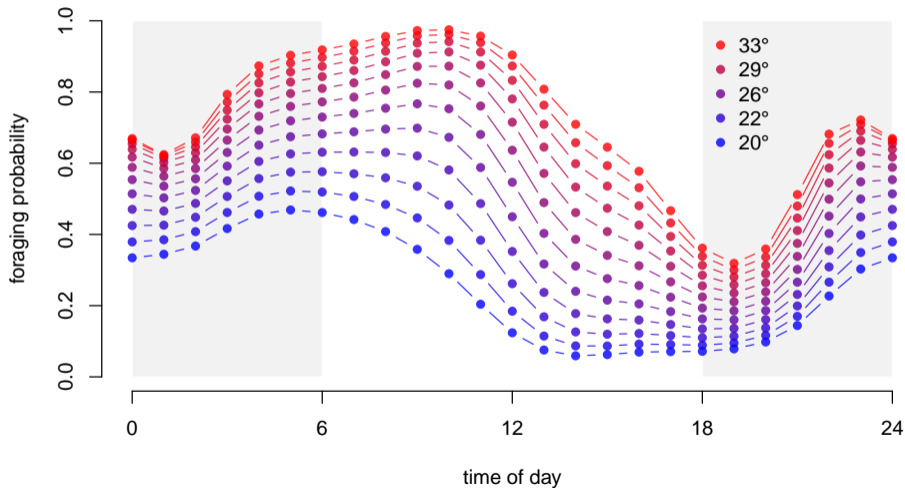
- 
- $\text{ODBA}_k^{(t)} \mid \{S_t = i\} \sim \beta_k^{(i)} + s_k(\text{time}_t), k = 1, \dots, 7$  (ectotherms)
  - $\gamma_{ij}^{(t)} \sim s(\text{tod}_t) + \text{AvgTemp}_t * s(\text{tod}_t)$  (parametric in original paper)
  - **11** smooths in total, **162** parameters
  - model fit takes  $\sim 5$  minutes (32 penalised fits until convergence)



# Bull sharks (Byrnes et al., 2023)



## Bull sharks (Byrnes et al., 2023)



## Concluding remarks

- **qREML** is (of course) much more generally applicable than just to HMMs

## Concluding remarks

- **qREML** is (of course) much more generally applicable than just to HMMs
- you can find `qrem1()` in the package `LaMa`, vignette “*Penalised splines*”

## Concluding remarks

- **qREML** is (of course) much more generally applicable than just to HMMs
- you can find `qrem1()` in the package `LaMa`, vignette “*Penalised splines*”
- models with i.i.d. random effects can be fitted using the same approach

## Concluding remarks

- **qREML** is (of course) much more generally applicable than just to HMMs
  - you can find `qrem1()` in the package `LaMa`, vignette “*Penalised splines*”
  - models with i.i.d. random effects can be fitted using the same approach
- 
- `RTMB` will become an extremely valuable tool for fitting complex models

