



THE UNIVERSITY OF
**WESTERN
AUSTRALIA**

Bayesian Regression with Functional Inequality Constraints

Joshua J. Bon with
C. Drovandi, K. Murray & B.A. Turlach

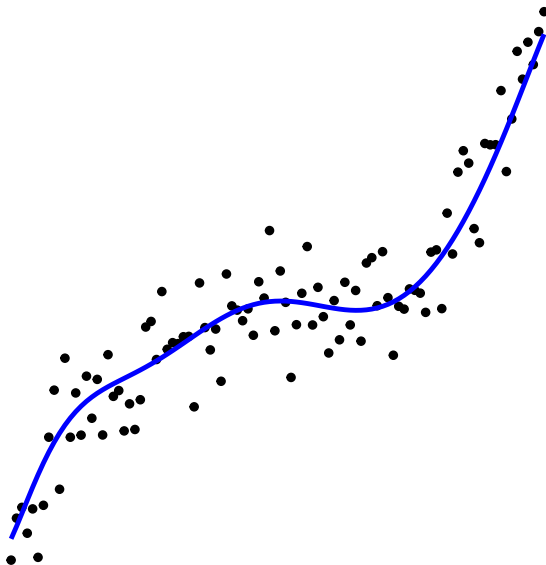
30th of November, 2017

School of Mathematics and Statistics, University of Western Australia

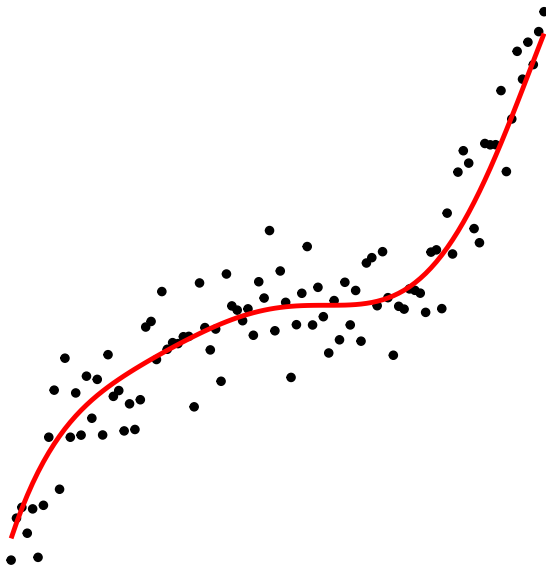
An uphill battle



An uphill battle



An uphill battle



Functional constraints

Functional constraints: Inequality constraints defined by a function that varies over an auxiliary set.

$$\beta \in \mathbb{R}^p \text{ such that}$$
$$f(\mathbf{x}; \beta) \geq c, \forall \mathbf{x} \in X \subseteq \mathbb{R}^d$$

Most often in regression problems these arise as shape constraints over certain regions, e.g. monotonic polynomials.

Functional constraints

A parameter space with functional constraints takes the form

$$\Omega = \{\beta \in \mathbb{R}^p : f(\mathbf{x}; \beta) \geq c, \forall \mathbf{x} \in X\}$$

It is not feasible to check all $\mathbf{x} \in X$ (infinite points to check).

Instead, rewrite the constraint as

$$\Omega = \left\{ \beta \in \mathbb{R}^p : \min_{\mathbf{x} \in X} f(\mathbf{x}; \beta) \geq c \right\}$$

Parameter constraints in Bayesian models

Incorporate information from the real world context into the prior probability of the parameters.

A constrained prior can be written as

$$\pi_c(\beta) \propto \pi(\beta) \times \mathbb{1}(\beta \in \Omega)$$

$\beta \notin \Omega$ are assigned a zero prior probability.

The posterior distribution is

$$\pi(\beta|y) \propto \pi(y|\beta)\pi_c(\beta)$$

Computational considerations

When estimating parameters under functional constraints (Bayesian or ML) the computational difficulty is in assessing

$$\min_{\mathbf{x} \in X} f(\mathbf{x}; \beta) \geq c$$

- When f is convex, a local minimum is the global minimum.
- Difficulties emerge when f is non-convex and $X \subseteq \mathbb{R}^d, d \geq 2$

Method: SMC + successively improving minima estimate

Sequential Monte Carlo

A subset of sequential Monte Carlo samplers¹ approximate the posterior distribution of static probabilistic models by:

- Evolving the starting distribution to the target distribution.
For example, $\pi_t \propto \pi(y|\beta)^{\phi_t} \pi(\beta)$ with increasing $\phi_t \rightarrow 1$.
- In each iteration they attempt to sample from π_t by
 1. Weighted resampling from the previous particles $\sim \pi_{t-1}$
 2. Particle mutation to increase diversity (e.g. MCMC transition kernel with invariant distribution $\sim \pi_t$)

¹P. Del Moral *et al.*, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**, 411–436 (2006).

Constrained Sequential Monte Carlo

Constrained sequential Monte Carlo samplers² move particles through a sequence of nested subsets until the constraint is satisfied. For our purposes we use:

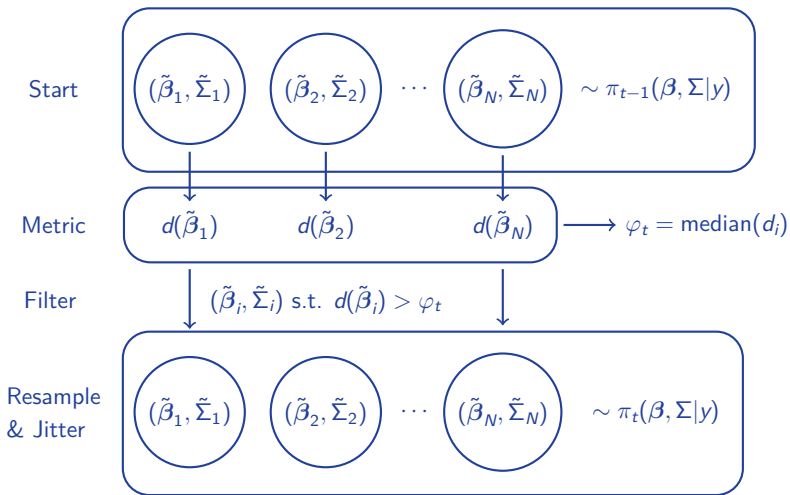
$$\pi_t(\beta|y) \propto \pi(\beta|y)\mathbb{1}(d(\beta) > \varphi_t)$$

- $\pi(\beta|y)$ is posterior distribution over the unconstrained space
- $d(\beta)$ measures “distance” away from the constrained space and $d(\beta) \geq 0 \implies \beta \in \Omega$
- $-\infty = \varphi_0 < \dots < \varphi_t < \varphi_{t+1} < \dots < \varphi_T = 0$
- Functional constraints: $d(\beta) = \min_{\mathbf{x} \in X} \{f(\mathbf{x}; \beta) - c\}$

²S. Golchi, D. A. Campbell, *Computational Statistics & Data Analysis* **97**, 98–113 (2016).

cSMC Algorithm Overview

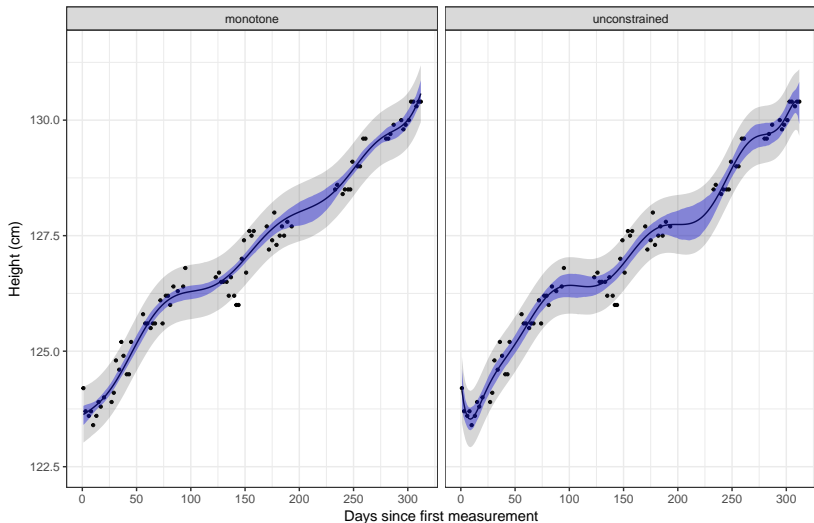
1. Generate initial particles: $\{(\tilde{\beta}_i, \tilde{\Sigma}_i)\}_{i=1}^N \sim \pi(\beta, \Sigma|y)$
2. Calc metric: $d_i = d(\tilde{\beta}_i)$, and $m = \sum_{i=1}^N \mathbb{1}(d_i \geq 0)$
3. Calc temp: $\varphi_t = \text{median}(\{d_i\}_{i=1}^N)$ or $\varphi_t = 0$ if $m > N/2$
4. Resample N particles from $\{(\tilde{\beta}_i, \tilde{\Sigma}_i)\}_{i=1}^N$ such that $d_i > \varphi_t$
5. Jitter/mutate kept particles with MCMC kernel $\sim \pi_t(\beta|y)$
6. Repeat 2–5 until $m = N$, i.e. all $\beta \in \Omega$



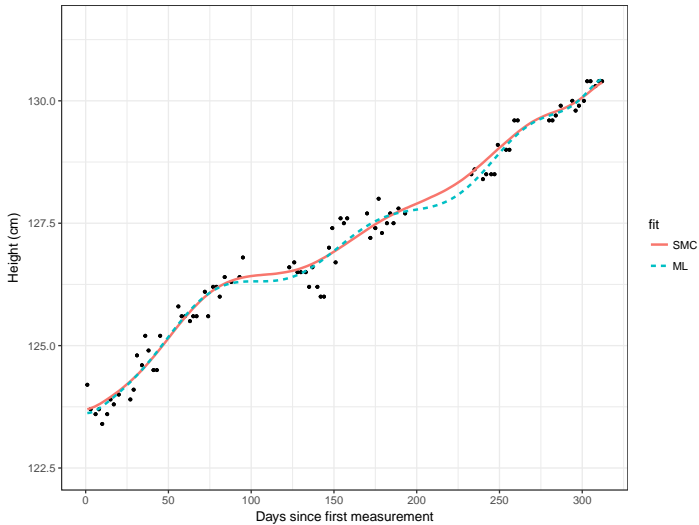
Monotonicity in one dimension:

$$d(\beta) = \min_{x \in [a, b]} f_x(x; \beta)$$

For twice-differentiable function, minimum can found at domain boundaries (if finite), or roots of $f_x(x; \beta)$.



Comparison of monotonic fit on “one child” dataset (Tuddenham and Snyder, 1954, `fda::onechild`). Fitted polynomials and 95% credible intervals for mean and posterior predictive distributions.



Maximum likelihood (ML) mean fitted values versus Bayesian mean fitted values (SMC). ML estimates from `MonoPoly` in R.

Monotonicity in two dimensions

Monotonicity in two dimensions:

$$d(\beta) = \min \{d_1(\beta), d_2(\beta)\} \text{ where } d_i(\beta) = \min_{\mathbf{x} \in X} f_{x_i}(\mathbf{x}; \beta)$$

No guarantee that f_{x_1} or f_{x_2} is convex. When f is an arbitrary polynomial this is known to be a difficult problem (many local minima and 1-dimensional boundaries).

So we approximate the distance metric:

$$d_i(\beta) \approx \text{best local min of } f_{x_i}(\mathbf{x}; \beta) \text{ from a set of starting points}$$

And attempt to improve the approximation each iteration

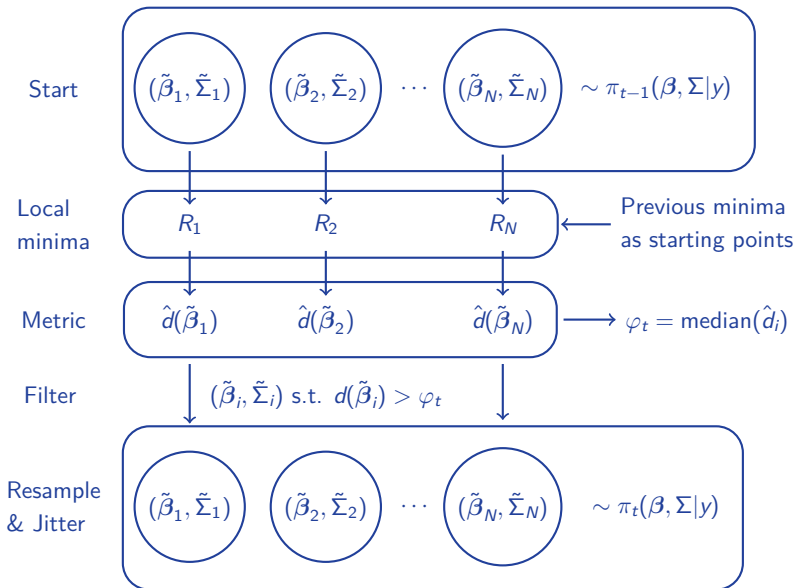
Particles now have 3 elements:

$$\{(\tilde{\beta}_i, \tilde{\Sigma}_i, R_i)\}_{i=1}^N$$

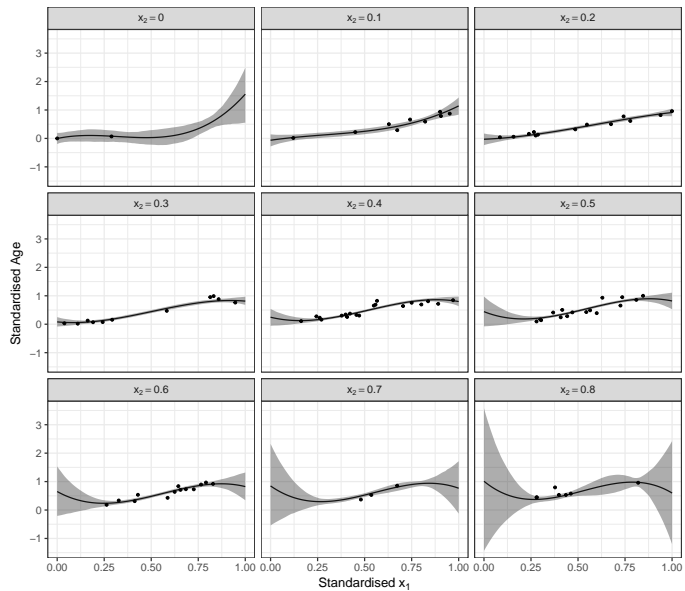
where R_i are the local minima found. Every time a particle is mutated the R_i are updated with some optimisation routine using starting points:

- The original R_i (likely to be close to new minima)
- Random sample of local minima across particles, i.e. $\{R_i\}_{i=1}^N$

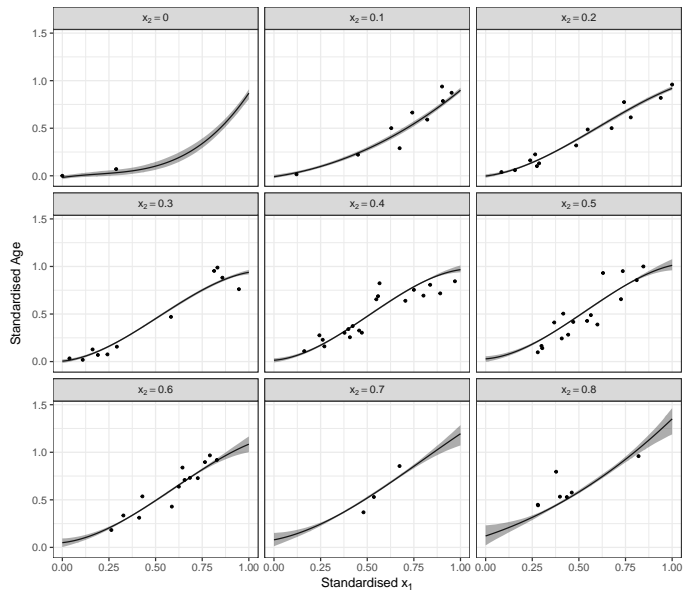
Similar to a particle-swarm algorithm. Over time a global minima (if it exists) for each particle should be found.



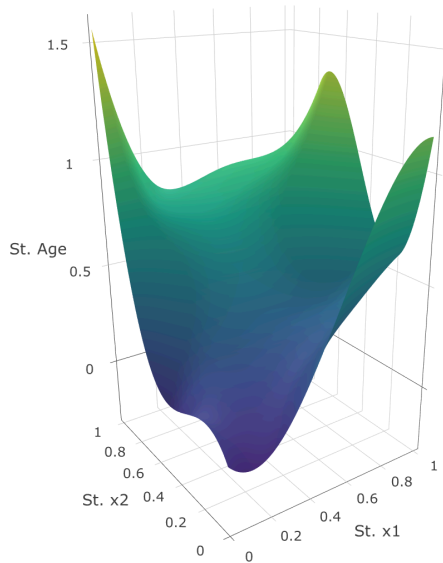
- Morphometric dataset
 - A number of morphometric measures taken from human skulls
 - Aged between 1 month and 19 years old
 - 174 individuals (93 male, 81 female)
- Aim: Predict age based on a selection of measurements
- In this example, using two measures that predict well individually for Male data.



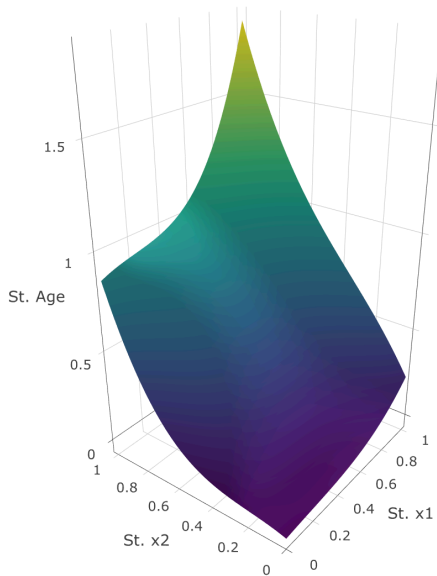
Unconstrained fit - 95% credible intervals for mean. Sliced by x_2 .



Monotone fit - 95% credible intervals for mean. Sliced by x_2 .



Unconstrained posterior mean fit.



Monotone posterior mean fit.

Conclusion

- Propose generic method for handling functional constraints in one or more dimensions
- Demonstrated with one- and two-dimensional monotonicity constraint
- What next?
 - Model selection
 - Incorporate measurement error techniques
 - Test on other datasets

Appendix

Comparison to maximum likelihood methods

ML regression generates an optimisation problem of the form

$$\max_{\beta} \log\{\pi(y|\beta)\} \text{ s.t. } \beta \in \Omega$$

again, where

$$\Omega = \{\beta \in \mathbb{R}^p : f(\mathbf{x}; \beta) \geq c, \forall \mathbf{x} \in X\}$$

This is a semi-infinite program, where X is often referred to as the index set.

Why Bayesian?

Due to the constraints, asymptotic ML theory is not applicable:

- Frequentist methods may need to rely on bootstrapping
- Whereas Bayesian computational methods can approximate entire posterior distribution

We can develop a Bayesian method to handle non-convex functional constraints with $d \geq 2$, based on:

- (Constrained) Sequential Monte Carlo
- Tracking local minima of the functional constraint

Other benefits

- No need to find optimal temperature for weights...