



Gaussian process optimisation for approximate Bayesian inference

IDA ML seminar series, Linköping University.

Johan Dahlin

johan.dahlin@liu.se

Division of Automatic Control,
Department of Electrical Engineering,
Linköping University

April 1, 2015



Linköpings universitet

■ This is ongoing/current collaborative work together with

- Dr. Fredrik Lindsten (University of Cambridge, United Kingdom)
- Prof. Thomas Schön (Uppsala University, Sweden)
- Prof. Mattias Villani (Linköping University, Sweden)

■ Aim

- Efficient approximate Bayesian parameter inference in nonlinear SSMs.
- Extending this to SSMs with intractable likelihoods.
- Inference in copula models with α -stable marginals.

■ Methods

- Gaussian process optimisation.
- Sequential Monte Carlo methods.
- Approximate Bayesian computations.

■ Contributions

- Decreased computational cost compared with popular methods.
- Interesting method for solving other costly optimisation problems.

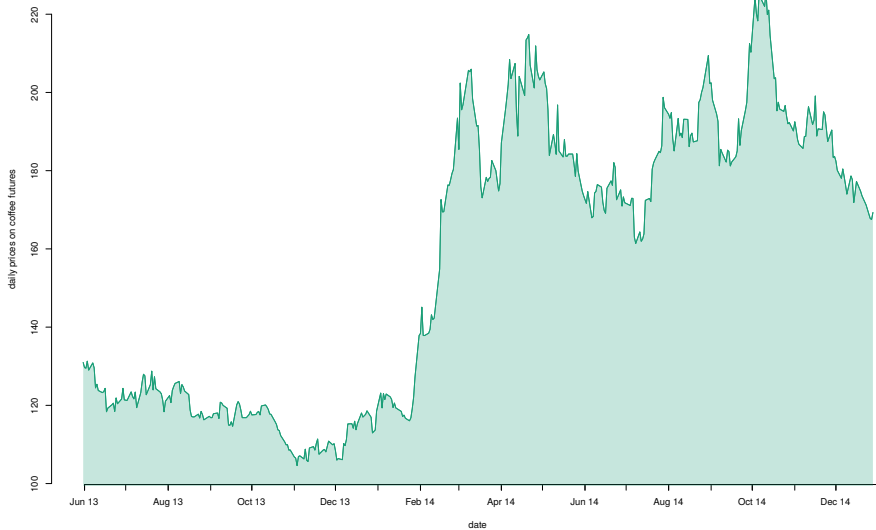
- Problem: modelling the risk in a portfolio of financial assets.
- Usually models the **log-returns**

$$y_t = \log(s_t - s_{t-1}),$$

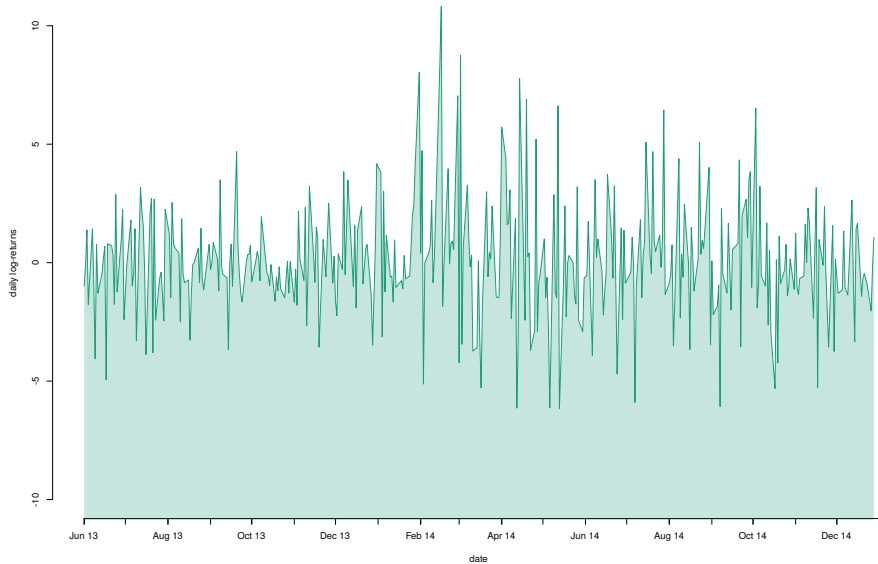
of some asset price s_t .

- **Copula model** (dependency structure and models of margins).
- Margins modelled using **stochastic volatility**.

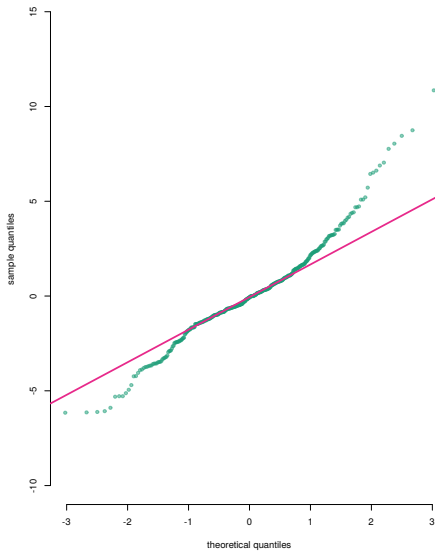
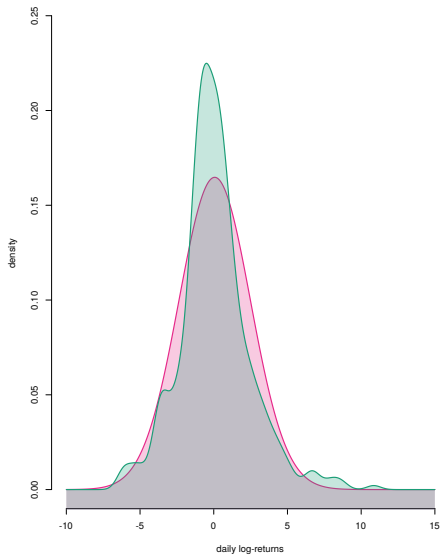
Motivating example: coffee futures



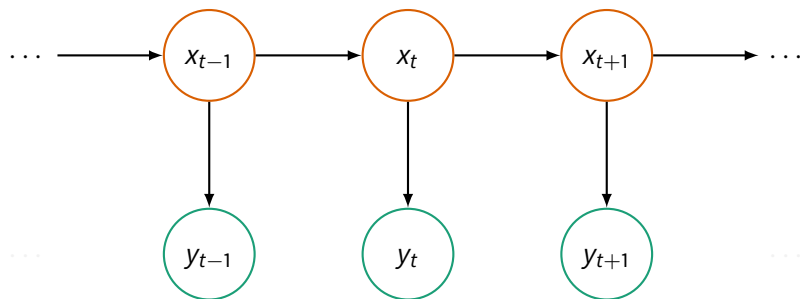
Motivating example: coffee futures



Motivating example: coffee futures



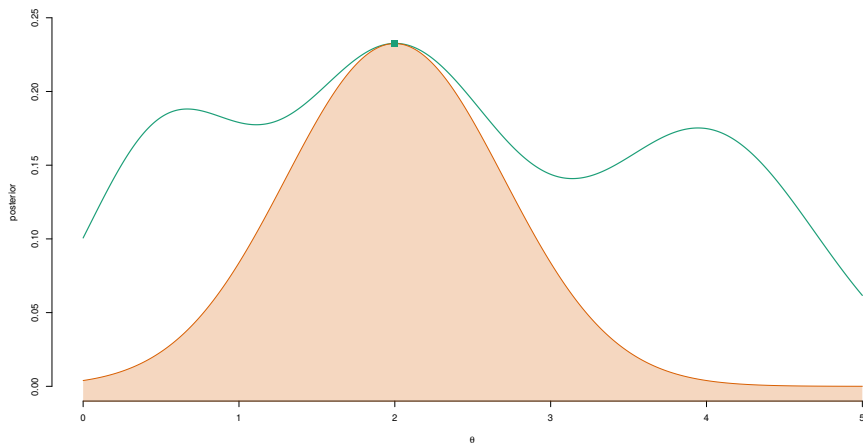
Markov chain $\{X_{0:T}, Y_{1:T}\}$ with $X_t \in \mathcal{X} = \mathbb{R}$, $Y_t \in \mathcal{Y} = \mathbb{R}$ and $t \in \mathbb{N}$.



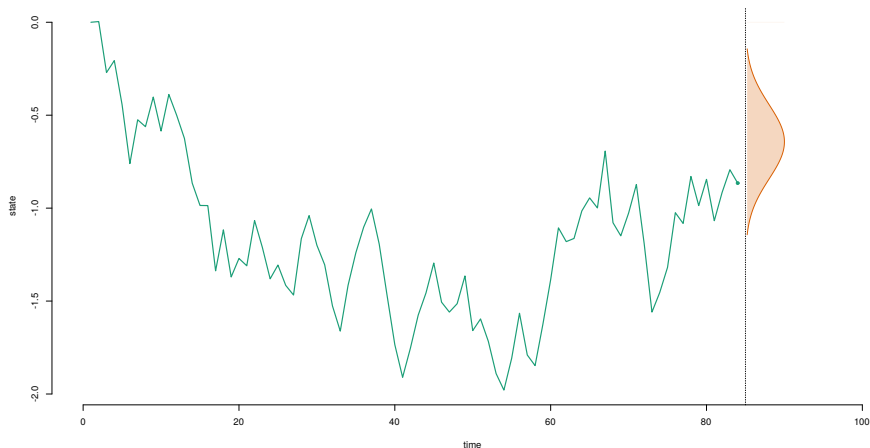
$$x_0 \sim \mu_\theta(x_0) \quad x_{t+1}|x_t \sim f_\theta(x_{t+1}|x_t), \quad y_t|x_t \sim g_\theta(y_t|x_t).$$

Example: stochastic volatility ($\theta = \{\mu, \phi, \sigma_v\}$):

$$x_{t+1}|x_t \sim \mathcal{N}\left(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v\right), \quad y_t|x_t \sim \mathcal{N}\left(y_t; 0, \exp(x_t)\right).$$



$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \left[\underbrace{\log p_{\theta}(y_{1:T}) + \log p(\theta)}_{\triangleq \log \pi(\theta)} \right]. \quad \hat{\pi}_{\text{Laplace}}(\theta) = \mathcal{N}\left(\hat{\theta}_{\text{MAP}}, \left[-\nabla^2 \log \pi(\theta) \Big|_{\theta=\hat{\theta}_{\text{MAP}}} \right]^{-1}\right).$$



filtering: $p_{\theta}(x_t|y_{1:t})$,

marginal smoothing: $p_{\theta}(x_t|y_{1:T})$,

joint smoothing: $p_{\theta}(x_{1:T}|y_{1:T})$,

- Particle-based Gaussian process optimisation
 - Sequential Monte Carlo / particle filtering.
 - Bayesian optimisation / Gaussian process optimisation.
- Extension to models with intractable likelihoods
 - α -stable processes.
 - Approximate Bayesian computations.

Based on the work presented in:

J. Dahlin and F. Lindsten, **Particle filter-based Gaussian process optimisation for parameter inference**. Proceedings of the 19th World Congress of the International Federation of Automatic Control (IFAC), Cape Town, South Africa, August 2014.

- An instance of **Bayesian optimisation**.
- *Global* optimisation method.
- Few function evaluations.
- **Gradient-free method**.
- Popular in machine learning for estimating hyperparameters.

References:

J. Mockus, V. Tiesis and A. Zilinskas, **The application of Bayesian methods for seeking the extremum**. In L.C.W. Dixon and G.P. Szegos (editors), *Toward Global optimisation*, North-Holland, 1978.

D.J. Lizotte, **Practical Bayesian optimisation**. PhD thesis, University of Alberta, 2008.

J. Snoek, H. Larochelle and R.P. Adams, **Practical Bayesian optimisation of Machine learning algorithms**. In *Advances in Neural Information Processing Systems (NIPS) 25*, Curran Associates Inc, 2012.

(i) Given iterate θ_k , estimate the log-posterior $\hat{\pi}_k \approx \pi(\theta_k)$.

Particle filtering.

(ii) Given $\{\theta_j, \hat{\pi}_j\}_{j=0}^k$, create a surrogate cost function of $\pi(\theta)$.

Gaussian process predictive distribution.

(iii) Select a new θ_{k+1} using the surrogate cost function.

Acquisition function based on the Gaussian process posterior.

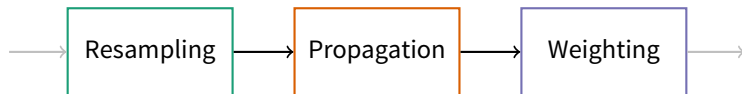
- An instance of sequential Monte Carlo (SMC) samplers.
- Estimates $\mathbb{E}[\varphi(x_t)|y_{1:t}]$ and $p_\theta(y_{1:T})$.

- Computational cost is in practice of order $\mathcal{O}(NT)$ with $N \sim T$.
- Well-understood statistical properties.
(unbiasedness, large deviation inequalities, CLTs)

References:

A. Doucet and A. Johansen, **A tutorial on particle filtering and smoothing**. In D. Crisan and B. Rozovsky (editors), The Oxford Handbook of Nonlinear Filtering. Oxford University Press, 2011.

O. Cappé, S.J. Godsill and E. Moulines, **An overview of existing methods and recent advances in sequential Monte Carlo**. In Proceedings of the IEEE 95(5), 2007.



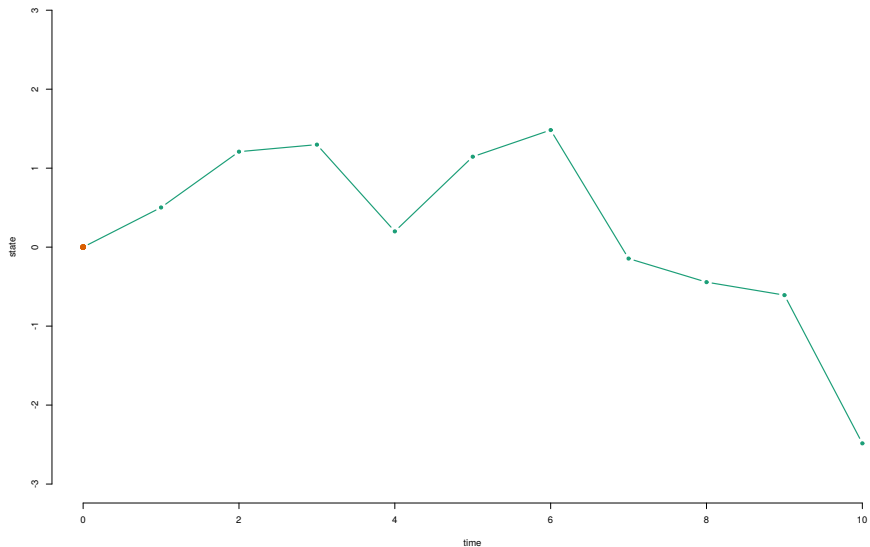
By iterating:

- **Resampling:** $\mathbb{P}(a_t^{(i)} = j) = \tilde{w}_{t-1}^{(j)}$, for $i, j = 1, \dots, N$.
- **Propagation:** $x_t^{(i)} \sim f_\theta(x_t | x_{t-1}^{(i)})$, for $i = 1, \dots, N$.
- **Weighting:** $w_t^{(i)} = g_\theta(y_t | x_t^{(i)})$, for $i = 1, \dots, N$.

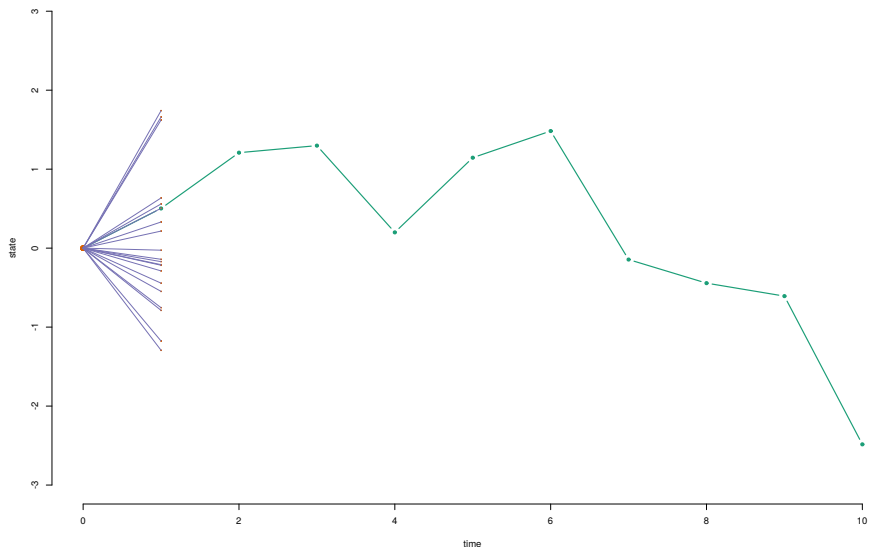
We obtain the particle system

$$\left\{ x_{0:T}^{(i)}, w_{0:T}^{(i)} \right\}_{i=1}^N.$$

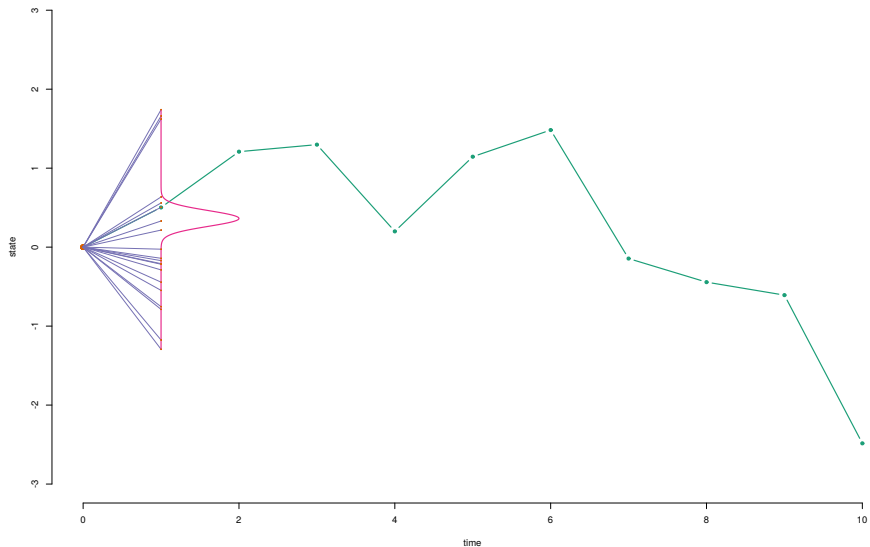
Particle filtering: animation [I/II]



Particle filtering: animation [I/II]

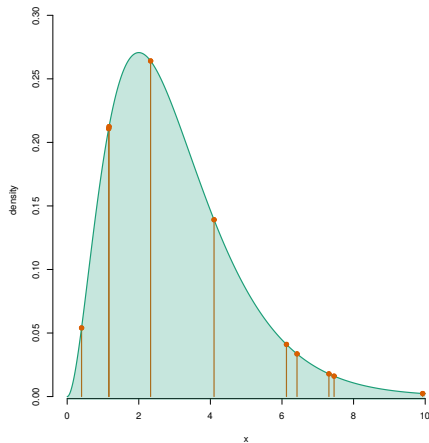
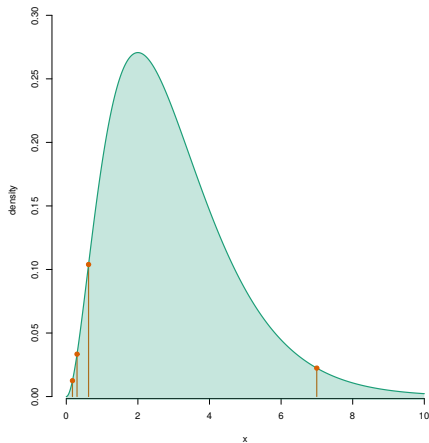


Particle filtering: animation [I/II]



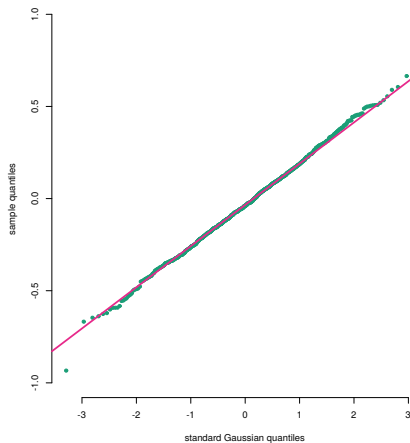
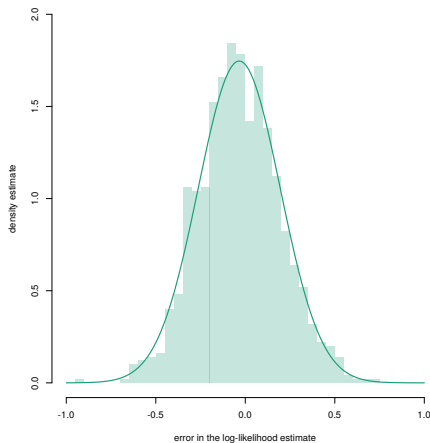
Particle filtering: animation [II/II]





$$\widehat{\varphi}_t^N \triangleq \widehat{\mathbb{E}}[\varphi(x_t) | y_{1:t}] = \sum_{i=1}^N w_t^{(i)} \varphi(x_t^{(i)}),$$

$$\sqrt{N}(\varphi_t - \widehat{\varphi}_t^N) \xrightarrow{d} \mathcal{N}(0, \sigma_t^2(\varphi)).$$



$$\log \hat{p}_\theta(y_{1:T}) = \sum_{t=1}^T \log \left(\sum_{i=1}^N w_t^{(i)} \right) - T \log N, \quad \sqrt{N} \left(\log p_\theta(y_{1:T}) - \log \hat{p}_\theta(y_{1:T}) + \frac{\sigma_\pi^2}{2N} \right) \xrightarrow{d} \mathcal{N} \left(0, \sigma_\pi^2 \right).$$

- Estimator for $\pi(\theta)$ available with Gaussian error.
- A Gaussian process is an **infinite dimensional Gaussian distribution**.
- Specified by a mean function m and covariance function (kernel) κ .
- Hyperparameters estimated using empirical Bayes.
- **Nonparametric Bayesian method**.
- Can be used for regression using a standard prior-posterior update.

References:

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

We assume a priori

$$\pi(\theta) \sim \mathcal{GP}\left(m(\theta), \kappa(\theta, \theta')\right),$$

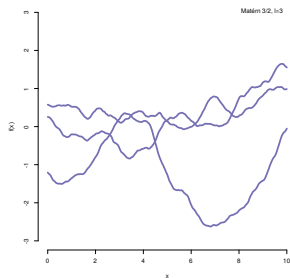
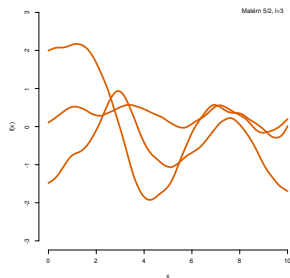
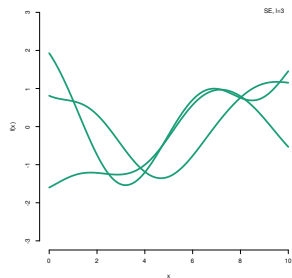
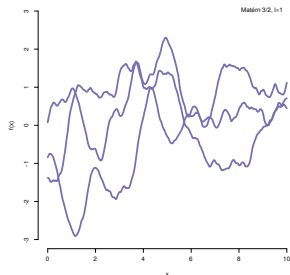
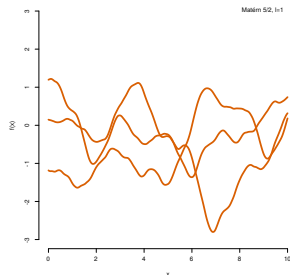
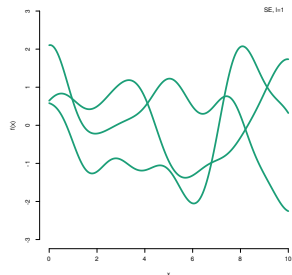
which together with the data

$$\hat{\pi}(\theta) = \pi(\theta) + \sigma_{\hat{\pi}} z_t, \quad z_t \sim \mathcal{N}(0, 1),$$

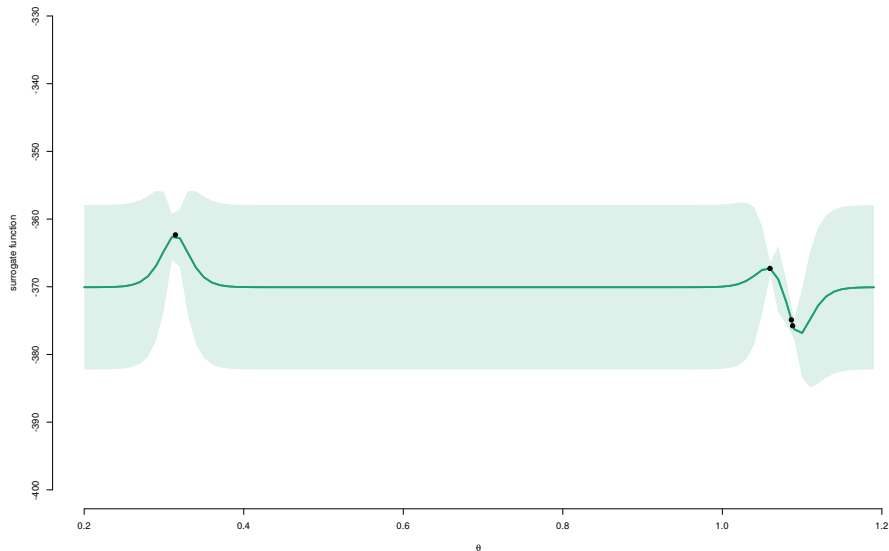
gives the **posterior predictive distribution**

$$\pi(\theta_*) | \mathcal{D}_k \sim \mathcal{N}\left(\mu(\theta_* | \mathcal{D}_k), \sigma^2(\theta_* | \mathcal{D}_k) + \sigma_{\hat{\pi}}^2\right),$$

where the current data is denoted $\mathcal{D}_k = \{\theta_j, \hat{\pi}_j\}_{j=1}^k$.



Gaussian process regression: toy example [I/II]



Gaussian process regression: toy example [II/II]



- Estimator for $\pi(\theta)$ available with Gaussian error.
- Surrogate of the log-posterior available as a Gaussian process.
- Idea: use the Gaussian process model to select θ_{k+1} .
- Balance exploration and exploitation using an acquisition rule by

$$\theta_{k+1} = \operatorname{argmax}_{\theta_{\star} \in \Theta} \text{AQ}(\theta_{\star} | \mathcal{D}_k).$$

- Cheap to evaluate $\text{AQ}(\theta_{\star} | \mathcal{D}_k)$.

References:

- J. Mockus, V. Tiesis and A. Zilinskas, **The application of Bayesian methods for seeking the extremum**. In L.C.W. Dixon and G.P. Szegso (editors), *Toward Global optimisation*, North-Holland, 1978.
- D.J. Lizotte, **Practical Bayesian optimisation**. PhD thesis, University of Alberta, 2008.
- J. Snoek, H. Larochelle and R.P. Adams, **Practical Bayesian optimisation of Machine learning algorithms**. In *Advances in Neural Information Processing Systems (NIPS) 25*, Curran Associates Inc, 2012.

Consider, the 95% upper confidence bound as the acquisition rule

$$\theta_{k+1} = \operatorname{argmax}_{\theta_* \in \Theta} \left[\mu(\theta_* | \mathcal{D}_k) + 1.96 \sqrt{\sigma^2(\theta_* | \mathcal{D}_k)} \right],$$

to determine the next iterate θ_{k+1} .

(i) Given iterate θ_k , estimate the log-posterior $\hat{\pi}_k \approx \pi(\theta_k)$.

Particle filtering.

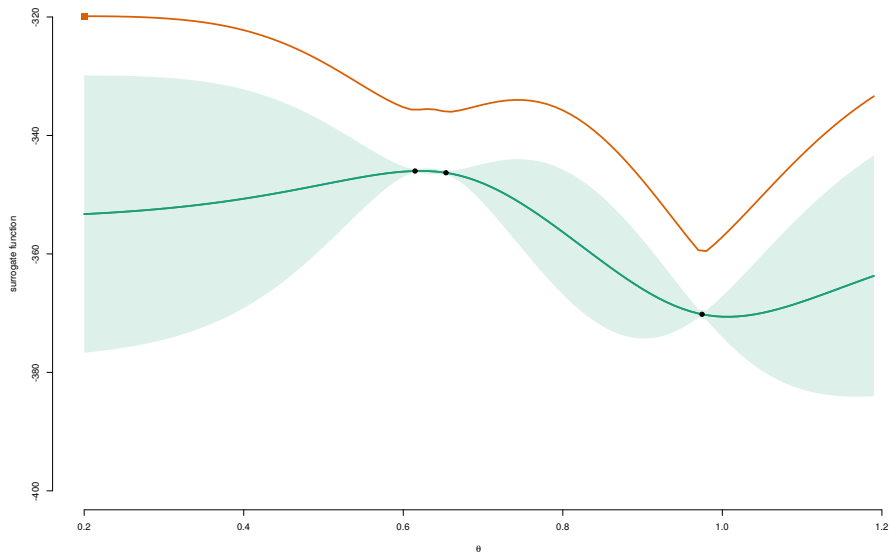
(ii) Given $\{\theta_j, \hat{\pi}_j\}_{j=0}^k$, create a surrogate cost function of $\pi(\theta)$.

Gaussian process predictive distribution.

(iii) Select a new θ_{k+1} using the surrogate cost function.

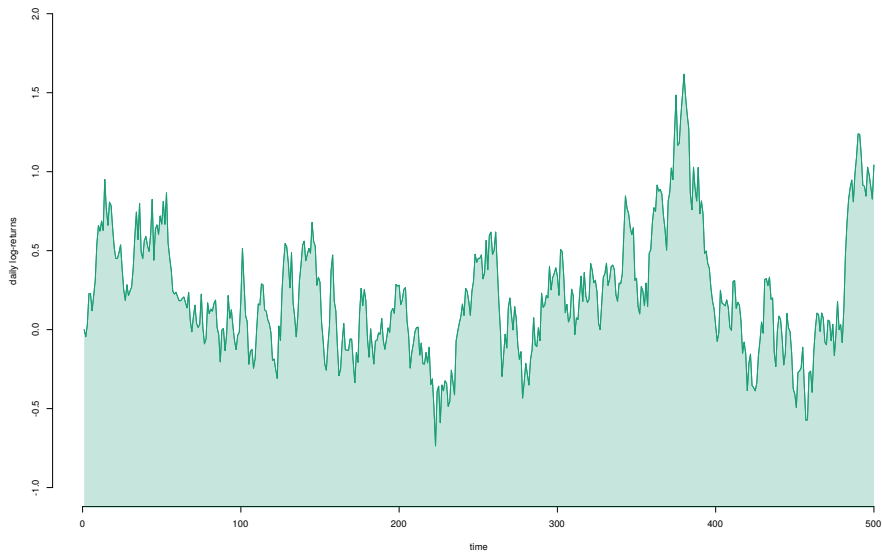
Acquisition function based on the Gaussian process posterior.

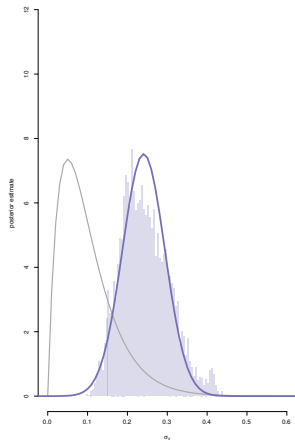
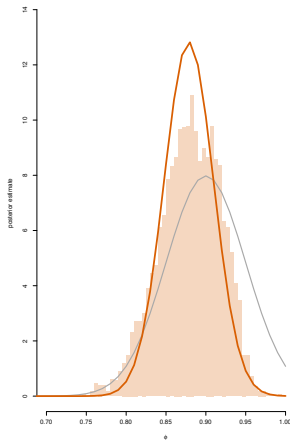
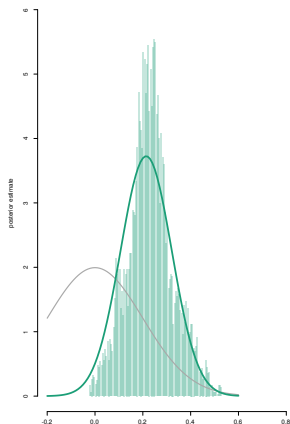
Example: Gaussian process optimisation [I/II]



Example: Gaussian process optimisation [II/II]



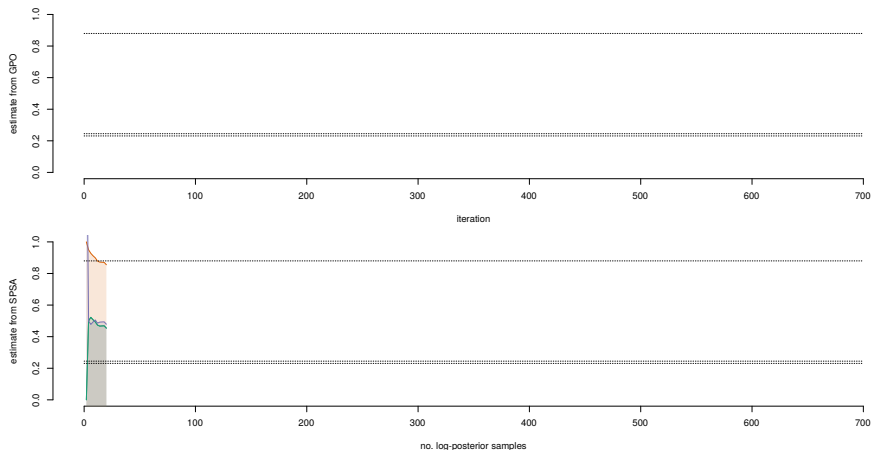




Stochastic volatility model with Gaussian returns ($\theta = \{\mu, \phi, \sigma_v\} = \{0.20, 0.96, 0.15\}$)

$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v), \quad y_t|x_t \sim \mathcal{N}(y_t; 0, \exp(x_t)).$$

Log-posterior estimates for PMH: 10 000 (histogram) and GPO: 350 (solid lines).



Stochastic volatility model with Gaussian returns ($\theta = (\mu, \phi, \sigma_v) = (0.20, 0.96, 0.15)$)

$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v), \quad y_t|x_t \sim \mathcal{N}(y_t; 0, \exp(x_t)),$$

Stochastic volatility model with Gaussian returns ($\theta = (\mu, \phi, \sigma_v) = (0.20, 0.96, 0.15)$)

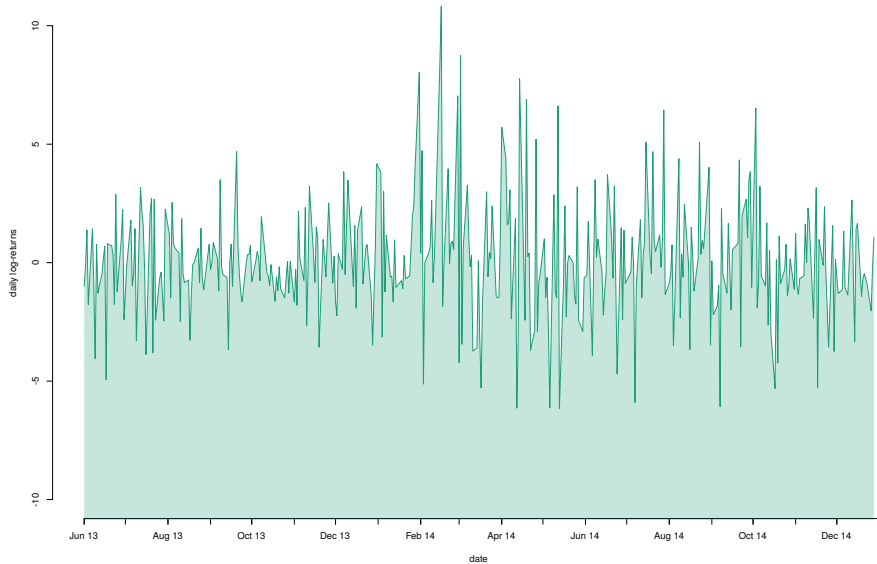
$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v), \quad y_t|x_t \sim \mathcal{N}(y_t; 0, \exp(x_t)),$$

- Particle-based Gaussian process optimisation
 - Sequential Monte Carlo / particle filtering.
 - Bayesian optimisation / Gaussian process optimisation.
- Extension to models with intractable likelihoods
 - α -stable processes.
 - Approximate Bayesian computations.

Based on the work presented in:

J. Dahlin, T. B. Schön, and M. Villani. **Approximate inference in state space models with intractable likelihoods using Gaussian process optimisation**. Technical Report LiTH-ISY-R-3075, Department of Electrical Engineering, Linköping University, Linköping, Sweden, April 2014.

Motivating example: coffee futures



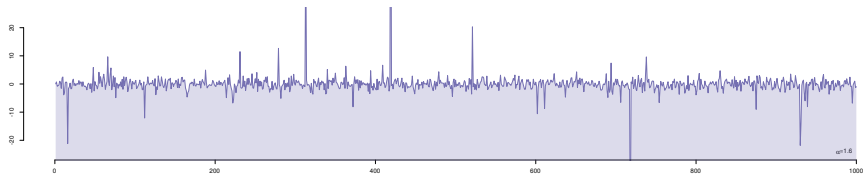
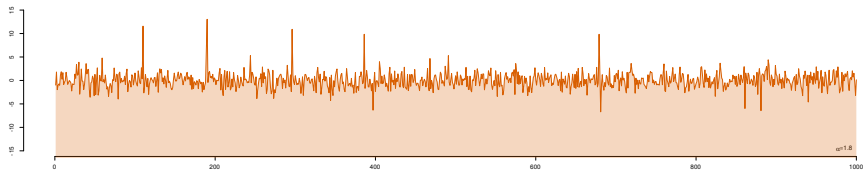
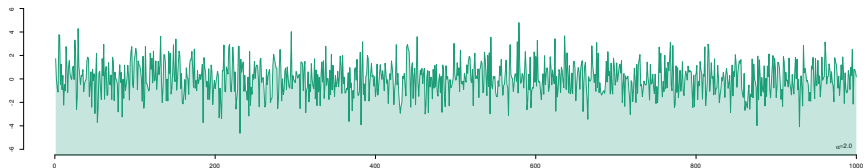
- Consider the model

$$x_{t+1}|x_t \sim \mathcal{N}\left(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v\right),$$
$$y_t|x_t \sim \mathcal{A}\left(y_t; \alpha, \exp(x_t)\right),$$

with $\theta = \{\mu, \phi, \sigma_v, \alpha\}$.

- $\mathcal{A}(\alpha, \gamma)$ denotes a symmetric α -stable distribution with zero mean and scale γ .

Stochastic volatility model with α -stable returns [II/II]



- The likelihood for this model is intractable.
- Often easy to simulate from the data distribution (likelihood).
- Idea: Simulate data from model and compare with observations.
- Augment the posterior with an auxiliary variable

$$\pi_{\epsilon}(\theta, \tilde{y}_{1:T}) = \frac{p_{\theta}(\tilde{y}_{1:T})p(\theta)\kappa_{\epsilon}(\tilde{y}_{1:T} - y_{1:T})}{\int p_{\theta'}(\tilde{y}_{1:T})p(\theta')\kappa_{\epsilon}(\tilde{y}_{1:T} - y_{1:T}) d\theta'}$$

and for a small enough ϵ , we have

$$\pi_{\epsilon}(\theta) = \int \pi_{\epsilon}(\theta, \tilde{y}_{1:T}) d\tilde{y}_{1:T}$$

References:

- S. Yildirim, S.S. Singh, T. Dean and A. Jasra. **Parameter Estimation in Hidden Markov Models with Intractable Likelihoods Using Sequential Monte Carlo**. Journal of Computational and Graphical Statistics, 2014.
- J.M., Marin, P. Pudlo, C.P. Robert, and R.J. Ryder. **Approximate Bayesian computational methods**. Statistics and Computing, 22(6), 1167-1180, 2012.

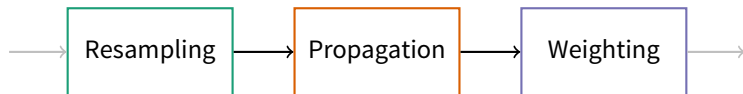
- We can simulate an α -stable r.v. by the transformation $\tilde{y}_t = \tau_\theta(v_t, x_t)$.
- Consider adding noise to the measurements

$$y_t^* = y_t + \epsilon z_t, \quad z_t \sim \mathcal{N}(0, 1).$$

- The model can then be rewritten as

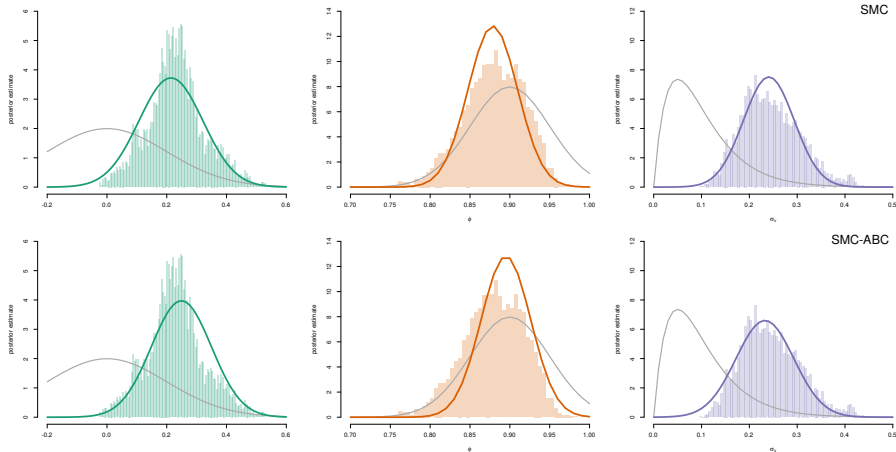
$$\begin{aligned}x_{t+1}|x_t &\sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma^2), \\v_t|x_t &\sim \nu_\theta(v_t; x_t), \\y_t^*|v_t &\sim h_{\theta, \epsilon}(y_t|v_t) = \mathcal{N}(y_t^*; \tau_\theta(v_t, x_t), \epsilon^2),\end{aligned}$$

where (x_t, v_t) is the new state vector.



- **Resampling:** $\mathbb{P}(a_t^{(i)} = j) = \tilde{w}_{t-1}^{(j)}$, for $i, j = 1, \dots, N$.
- **Propagation [1]:** $x_t^{(i)} \sim f_\theta(x_t | x_{t-1}^{(i)})$, for $i = 1, \dots, N$.
- **Propagation [2]:** $v_t^{(i)} \sim \nu_\theta(v_t | x_t^{(i)})$, for $i = 1, \dots, N$.
- **Weighting:** $w_t^{(i)} = h_{\theta, \epsilon}(y_t^* | v_t^{(i)})$, for $i = 1, \dots, N$.

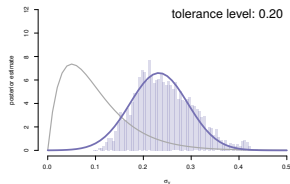
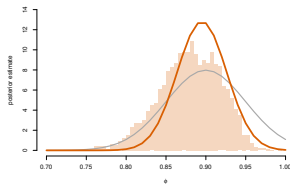
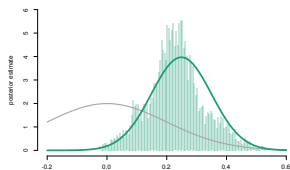
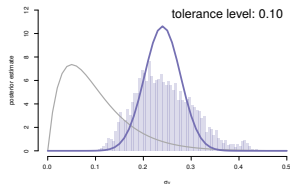
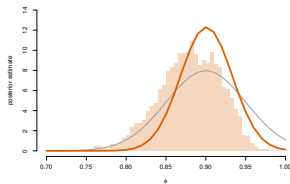
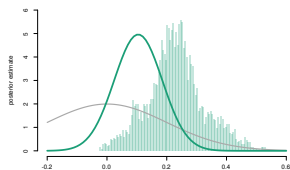
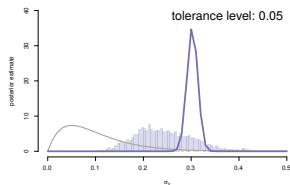
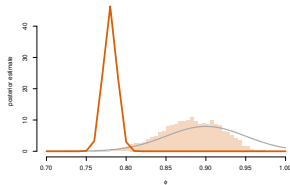
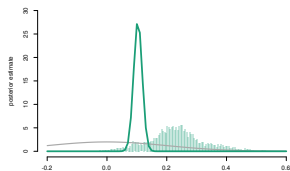
Results in an unbiased likelihood estimator (under some assumptions).



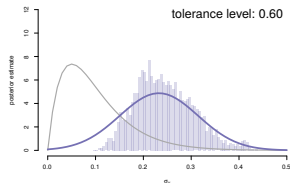
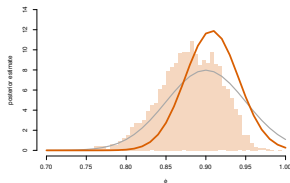
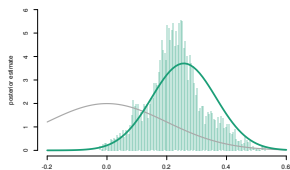
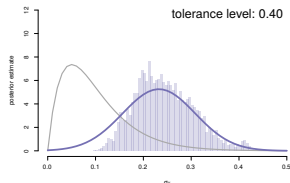
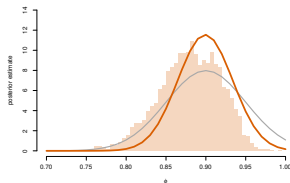
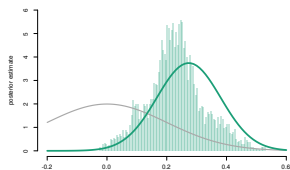
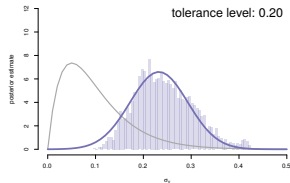
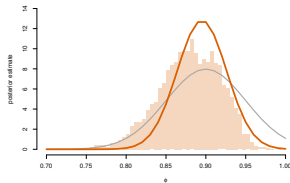
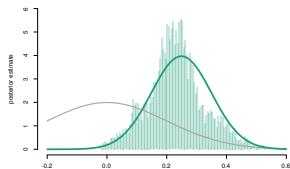
Stochastic volatility model with Gaussian returns ($\theta = \{\mu, \phi, \sigma_v\} = \{0.20, 0.96, 0.15\}$)

$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v), \quad y_t|x_t \sim \mathcal{N}(y_t; 0, \exp(x_t)).$$

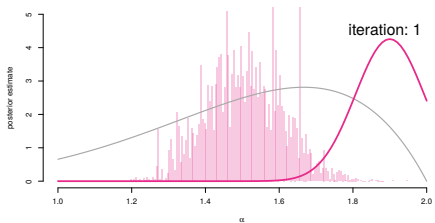
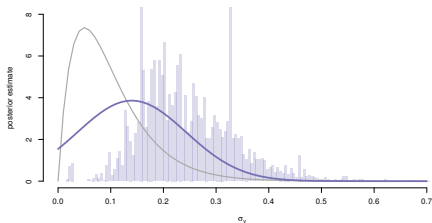
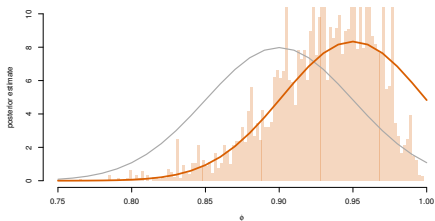
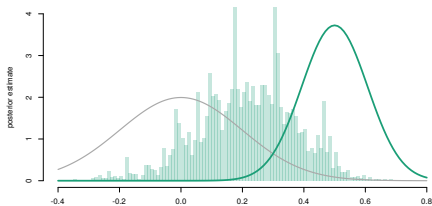
Log-posterior estimates for PMH: 10 000 (histogram) and GPO: 350 (solid lines).



Stochastic volatility: synthetic data [III/III]



Motivating example: volatility of coffee futures [I/III]



Stochastic volatility model with α -stable returns

$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v), \quad y_t|x_t \sim \mathcal{A}(y_t; \alpha, \exp(x_t)).$$

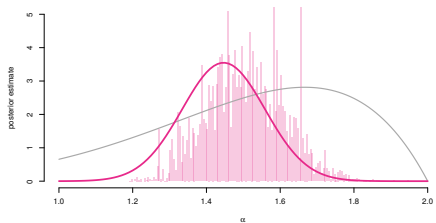
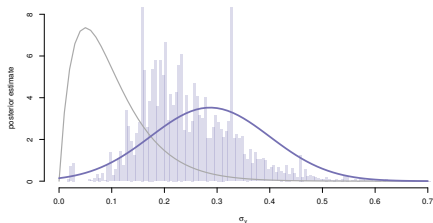
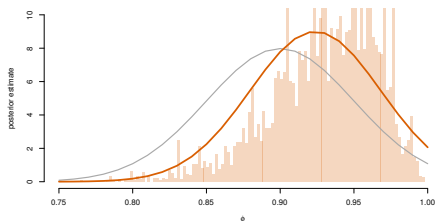
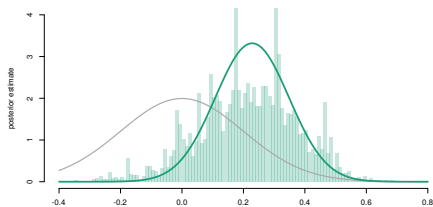
Log-posterior estimates for PMH-ABC: 30 000 (histogram) and GPO-ABC: 350 (solid lines).

Stochastic volatility model with α -stable returns

$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v), \quad y_t|x_t \sim \mathcal{A}(y_t; \alpha, \exp(x_t)).$$

Log-posterior estimates for PMH-ABC: 30 000 (histogram) and GPO-ABC: 350 (solid lines).

Motivating example: volatility of coffee futures [III/III]

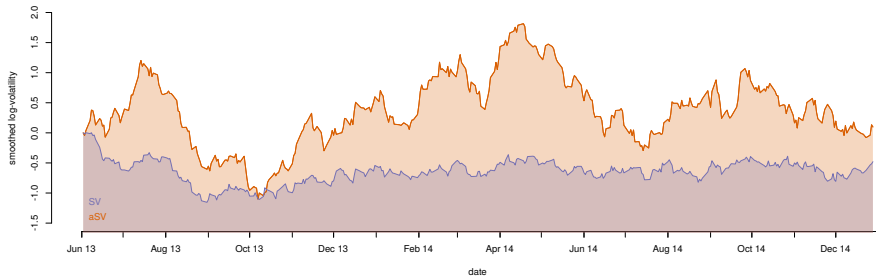
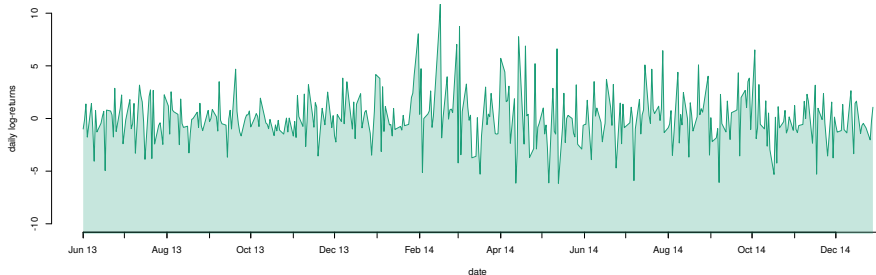


Stochastic volatility model with α -stable returns

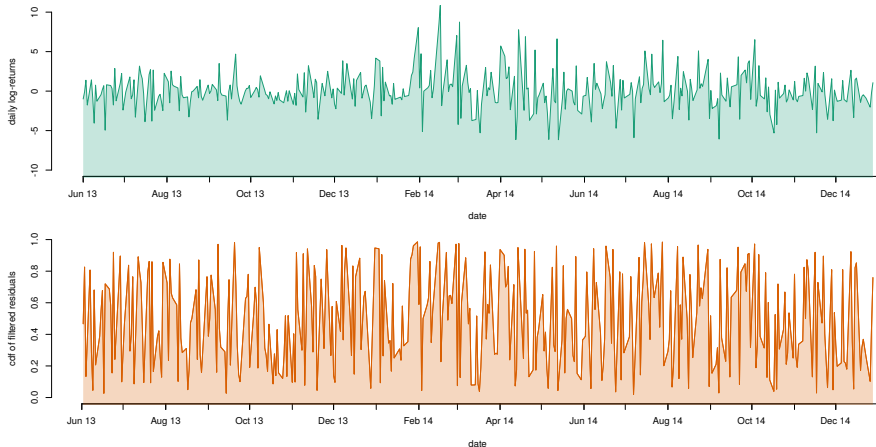
$$x_{t+1}|x_t \sim \mathcal{N}(x_{t+1}; \mu + \phi(x_t - \mu), \sigma_v), \quad y_t|x_t \sim \mathcal{A}(y_t; \alpha, \exp(x_t)).$$

Log-posterior estimates for PMH-ABC: 30 000 (histogram) and GPO-ABC: 350 (solid lines).

Motivating example: volatility of coffee futures

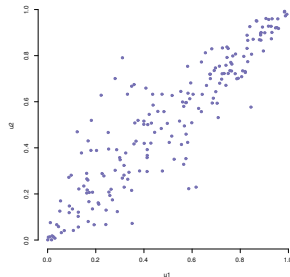
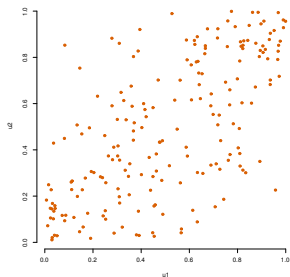
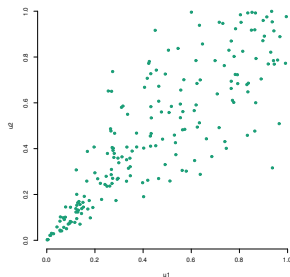
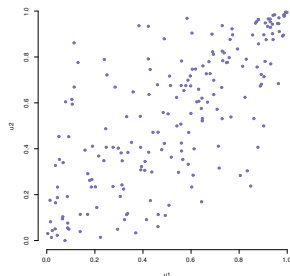
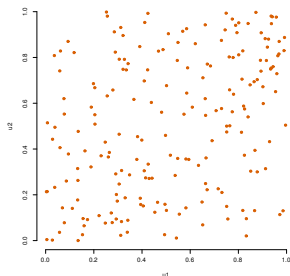
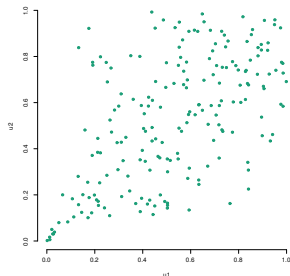


Motivating example: residuals of coffee futures



$$\hat{\epsilon}_t = y_t \exp(-\hat{x}_t/2), \quad \hat{u}_t = G_{\hat{\theta}}(\hat{\epsilon}_t).$$

Motivating example: residuals of coffee futures



■ Methods

- Particle filtering for log-likelihood estimation.
- Gaussian process surrogate model and acquisition rules.
- Approximate Bayesian computations.

■ Advantages

- Decreased computational cost compared with popular methods.
- Only makes use of *cheap* zero-order information.

■ Future work

- Bias compensation of log-likelihood estimate.
- Sparse Gaussian processes.
- Kernel design and new acquisition rules.
- Improved estimates of the Hessian.

■ Thank you for your attention!

- Questions, comments and suggestions are most welcome.
- The papers and some implementations are available for download from my homepage: <http://work.johandahlin.com/>.

This presentation is based on the work presented in:

J. Dahlin and F. Lindsten, **Particle filter-based Gaussian process optimisation for parameter inference**. Proceedings of the 19th World Congress of the International Federation of Automatic Control (IFAC), Cape Town, South Africa, August 2014.

J. Dahlin, T. B. Schön, and M. Villani. **Approximate inference in state space models with intractable likelihoods using Gaussian process optimisation**. Technical Report LiTH-ISY-R-3075, Department of Electrical Engineering, Linköping University, Linköping, Sweden, April 2014.

Consider the probability

$$\begin{aligned}\mathbb{P}(\mu(\theta_\star|\mathcal{D}_k) \geq \mu_{\max} + \epsilon) &= \mathbb{P}\left(\frac{\mu(\theta_\star|\mathcal{D}_k) - \mu_{\max} - \epsilon}{\sigma^2(\theta_\star|\mathcal{D}_k)} \geq 0\right) \\ &\triangleq \Phi(z_k(\theta_\star))\end{aligned}$$

where μ_{\max} denotes the current maximum of $\{\hat{\pi}_1, \dots, \hat{\pi}_k\}$.

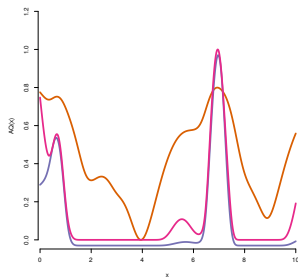
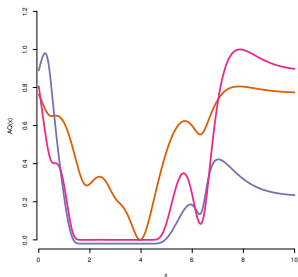
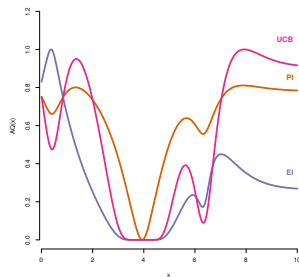
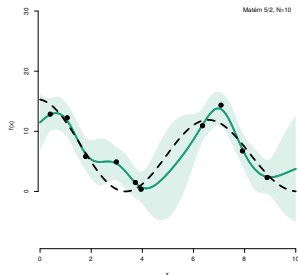
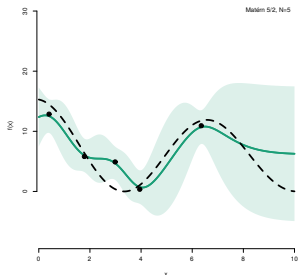
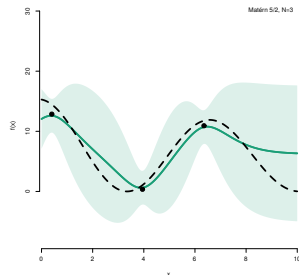
Consider the function

$$l(\theta_*) = \max \left\{ 0, \mu(\theta_* | \mathcal{D}_k) - \mu_{\max} - \epsilon \right\},$$

and its expected value with respect to the GP, given by

$$\begin{aligned} \mathbb{E} \left[l(\theta_*) | \mathcal{D}_k \right] &= \int l(\theta_*) d\Phi(z_k(\theta_*)) \\ &= \sigma^2(\theta_* | \mathcal{D}_k) \left[z_k(\theta_*) \Phi(z_k(\theta_*)) + \phi(z_k(\theta_*)) \right]. \end{aligned}$$

Acquisition rule: comparison



- (i) For $k = 1, \dots, K$
 - (a) Sample $\theta^{(k)} \sim p(\theta)$,
 - (b) Compute the weight $w^{(k)} = \pi(\theta^{(k)})$,

- (ii) Estimate the parameter posterior mean

$$\mathbb{E}[\theta | y_{1:T}] \approx \sum_{k=1}^K \frac{w^{(k)}}{\sum_{l=1}^K w^{(l)}} \theta^{(k)}.$$

- (i) For $k = 1, \dots, K$
 - (a) Sample $\theta^{(k)} \sim p(\theta)$,
 - (b) Generate $\tilde{y}_{1:T} \sim \nu_{\theta^{(k)}}$.
 - (c) Compute the weight

$$w^{(k)} = p(\theta^{(k)}) \kappa_{\epsilon} \left(y_{1:T} - \tilde{y}_{1:T} \right),$$

- (ii) Estimate the parameter posterior mean

$$\mathbb{E}[\theta | y_{1:T}] \approx \sum_{k=1}^K \frac{w^{(k)}}{\sum_{l=1}^K w^{(l)}} \theta^{(k)}.$$

