

Closed-loop automatic experimentation for Bayesian optimisation

Dave Woods

Southampton Statistical Sciences Research Institute
University of Southampton, UK
D.Woods@southampton.ac.uk
www.southampton.ac.uk/~davew

UNIVERSITY OF
Southampton

EPSRC
Engineering and Physical Sciences
Research Council

Outline

- ▶ Motivation: closed-loop optimisation of chemical processes
- ▶ Decision-theoretic sequential design of experiments for Bayesian optimisation
- ▶ Learning and down-weighting outliers
- ▶ Examples

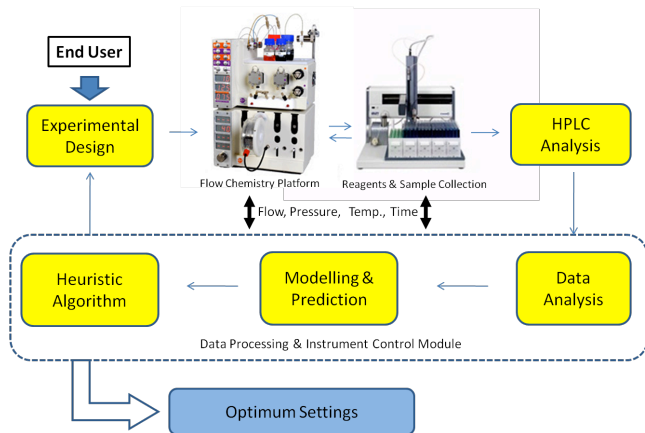
Joint work with Tim Waite (University of Manchester, UK)

Chemistry examples provided by Richard Whitby, Ian Broadwell (University of Southampton) and Xiaoping Tang (Pharmaron)

Supported by the UK Engineering and Physical Sciences Research Council

Motivation

Closed Loop Flow Chemistry Optimisation



Automated experimentation

Automation **allows** more (economically, time) efficient experimentation, 24/7 lab operation, reduction in user errors, and wider use of better designed experiments

Automation **requires** principled and computationally efficient algorithms

- ▶ Generation of (one-point-at-a-time) sequential designs
- ▶ **Robust** modelling of sequentially collected data making minimal assumptions
- ▶ Identification of possible outliers

We address these issues through

- ▶ Gaussian process regression and Bayesian updating
- ▶ Decision-theoretic sequential design for optimisation
- ▶ (Re-) weighting of observations to reflect our belief in their reliability

Bayesian optimisation

Assume the **aim** is to maximise an unknown, and expensive, function $g(\mathbf{x})$ with respect to q controllable variables $\mathbf{x} = (x_1, \dots, x_q)^T \in \mathcal{X} \subset \mathbb{R}^q$

[Here, expensive will mean an experiment is required to obtain a (noisy) evaluation of $g(\mathbf{x})$ at any \mathbf{x}]

Bayesian optimisation

- ▶ places a prior distribution on $g(\mathbf{x})$ (i.e. a statistical model)
- ▶ collects (noisy) function evaluations (data) at points chosen sequentially via an acquisition function
- ▶ updates the prior to a posterior distribution and infers the maximum of $g(\mathbf{x})$

Uncertainty in the posterior leads to an exploration/exploitation trade-off

In the statistics literature, **Expected Improvement (EI)** using **Gaussian processes** is the most popular acquisition function

Mockus (1989), Jones et al, (1998), Brochu et al., 2010

Gaussian process regression

Bayesian Gaussian process models: nonparametric modelling of a black box function $g(\mathbf{x})$

Assume we can observe $y = g(\mathbf{x}) + \varepsilon$, for $\varepsilon \sim N(0, \sigma_\varepsilon^2)$

Gaussian process prior: $g(\mathbf{x}) \sim \text{GP} \{ \mu(\mathbf{x}), \kappa(\mathbf{x}, \mathbf{x}') \}$

$$g(\mathbf{x}_1), \dots, g(\mathbf{x}_n) \sim \text{MVN}(F\boldsymbol{\beta}, \sigma_g^2 K)$$

F = model matrix for mean **trend**

K = **correlation** matrix with ij th entry $k_\lambda(\mathbf{x}_i, \mathbf{x}_j)$

Correlation function $k_\lambda(\mathbf{x}, \mathbf{x}')$ often assumed stationary and separable in \mathbf{x} , with parameter λ

Conditional on the hyperparameters, the posterior for $g(\mathbf{x})$ is also a Gaussian process; unconditional posterior inference requires numerical methods

Rasmussen & Williams (2006)

Decision-theoretic approach 1

Classic EI is for the case of non-noisy observations, $\sigma_\varepsilon^2 = 0$.

Most extensions to EI for noisy data are heuristic generalisations that do not take account of parameter updating or inference for the maximum

Huang et al. (2006), Gramacy & Polson (2011)

Problem 1. Estimation: Given data $\mathbf{y}_{1:n} = (y_1, \dots, y_n)^\top$, find a point estimate $\mathbf{x}_n^* \in \mathcal{M}_n \subseteq \mathcal{X}$ for the optimum

We assume the utility of a point \mathbf{x} is

$$u(\mathbf{x}, g) = g(\mathbf{x}),$$

with corresponding Bayes decision

$$\mathbf{x}_n^* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{M}_n} \mathbb{E}_g \{g(\mathbf{x}) \mid \mathbf{y}_{1:n}\}$$

Decision-theoretic approach 2

Now suppose one additional observation, y_{n+1} , will be taken and the posterior mean $E\{g(\mathbf{x})|\mathbf{y}_{1:(n+1)}\}$ used to locate $\mathbf{x}_{n+1}^* \in \mathcal{M}_{n+1}$

Problem 2. Design: the optimal choice of \mathbf{x}_{n+1} , the next design point, maximises

$$\begin{aligned} U(\mathbf{x}_{n+1}) &= E_{g, y_{n+1}} \{u(\mathbf{x}_{n+1}^*, g) \mid \mathbf{y}_{1:n}\} \\ &= E_{y_{n+1}} \left[\max_{\mathbf{x} \in \mathcal{M}_{n+1}} E_g \{g(\mathbf{x}) \mid \mathbf{y}_{1:(n+1)}\} \mid \mathbf{y}_{1:n} \right] \end{aligned}$$

For interpretation and presentation, we usually use the **Expected Gain in Utility (EGU)**

$$U(\mathbf{x}_{n+1}) - E_g \{g(\mathbf{x}_n^*) \mid \mathbf{y}_{1:n}\}$$

see also Osborne et al. (2010)

Two simple results

1. Expected improvement is a special case of EGU ...

... if the responses are deterministic, $y = g(\mathbf{x})$, $\mathcal{M}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, and the GP hyperparameters are known

Here, $\mathbf{x}_n^* = \operatorname{argmax}_{i=1, \dots, n} g(\mathbf{x}_i)$ and hence

$$\begin{aligned} U(\mathbf{x}_{n+1}) &= E_g \left\{ \max \left(g(\mathbf{x}_{n+1}), g(\mathbf{x}_n^*) \right) \mid \mathbf{y}_{1:n} \right\} \\ &= g(\mathbf{x}_n^*) + E_g \left\{ \max \left(g(\mathbf{x}_{n+1}) - g(\mathbf{x}_n^*), 0 \right) \mid \mathbf{y}_{1:n} \right\} \\ &= g(\mathbf{x}_n^*) + \operatorname{EI}(\mathbf{x}_{n+1}) \end{aligned}$$

2. The EGU is non-negative ...

... if $\mathcal{M}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. Then,

$$U(\mathbf{x}_{n+1}) \geq E_g \left\{ g(\mathbf{x}_n^*) \mid \mathbf{y}_{1:n} \right\}$$

Updating the GP and approximating EGU

Updating is via Sequential Monte Carlo with N_θ particle/weight pairs (θ_j, w_j) , where θ holds values of the hyperparameters and $0 < w_j \leq 1$, $\sum_j w_j = 1$

- ▶ Reweight, resample and move steps after observation of y_{n+1}
- ▶ Iterated batch importance sampling (IBIS)

Chopin (2002), Gramacy & Polson (2011), Drovandi et al. (2013)

Approximation of EGU presents two main challenges:

1. for given y_{n+1} , we must approximate the posterior for θ ;
2. the expectation must be taken with respect to the predictive distribution of y_{n+1}

Approximating EGU 1

Let $u(\mathbf{x}_{n+1}; y_{n+1}) = \max_{\mathbf{x} \in \mathcal{M}_{n+1}} \mathbb{E}_g \{g(\mathbf{x}) \mid \mathbf{y}_{1:(n+1)}\}$

Then we need to approximate

$$U(\mathbf{x}_{n+1}) = \int_{\mathbb{R}} u(\mathbf{x}_{n+1}; y_{n+1}) \pi(y_{n+1} \mid \mathbf{y}_{1:n}) dy_{n+1}$$

1. Particle approximation for the predictive distribution:

$$\hat{\pi}(y_{n+1} \mid \mathbf{y}_{1:n}) = \sum_j w_j \pi(y_{n+1} \mid \mathbf{y}_{1:n}, \theta_j)$$

2. Quadrature approximation for the integral w.r.t. y_{n+1} :

$$\hat{U}(\mathbf{x}_{n+1}) = \sum_k v_k \hat{\pi}(y_{n+1}^{(k)} \mid \mathbf{y}_{1:n}) u(\mathbf{x}_{n+1}; y_{n+1}^{(k)}),$$

with, $v_k, y_{n+1}^{(k)}$ are integration weights and abscissae

Approximating EGU 2

To estimate future maximised expected utility $u(\mathbf{x}_{n+1}; y_{n+1})$:

1. Given y_{n+1} , calculate an approximate posterior density for θ via simple reweighting, i.e. using (θ_j, \tilde{w}_j) with

$$\tilde{w}_j \propto w_j \pi(y_{n+1} \mid \mathbf{y}_{1:n}, \theta_j)$$

2. Use these weights to approximate the future posterior expectation $E_g \{g(\mathbf{x}) \mid \mathbf{y}_{1:(n+1)}\}$, giving

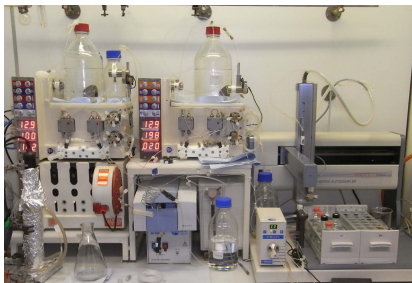
$$\hat{u}(\mathbf{x}_{n+1}; y_{n+1}) = \max_{\mathbf{x} \in \mathcal{M}_{n+1}} \left\{ \sum_j \tilde{w}_j E_g [g(\mathbf{x}) \mid \mathbf{y}_{1:(n+1)}, \theta_j] \right\}$$

Example 1

Real piece of chemistry run in Southampton Chemistry with online statistical modelling performed at Manchester

Hardware: 4-channel (A–D) Vapourtec flow reactor

Measurement via HPLC/UV



Average time for 1 run is 2 hours; practical upper bound of 30-40 runs

Variables under study

- ▶ Conc A - 0.3–1 mol dm⁻³
- ▶ Ratio B - 1–5 eq
- ▶ Residence time - 5–30 min
- ▶ Temperature - 30–180° C

Ratio D = 0.1 held fixed, Ratio C determined by Conc A

Response: product yield

Some combinations of Conc-Ratio are not physically feasible (as they would imply infeasibly low or negative flow rates)

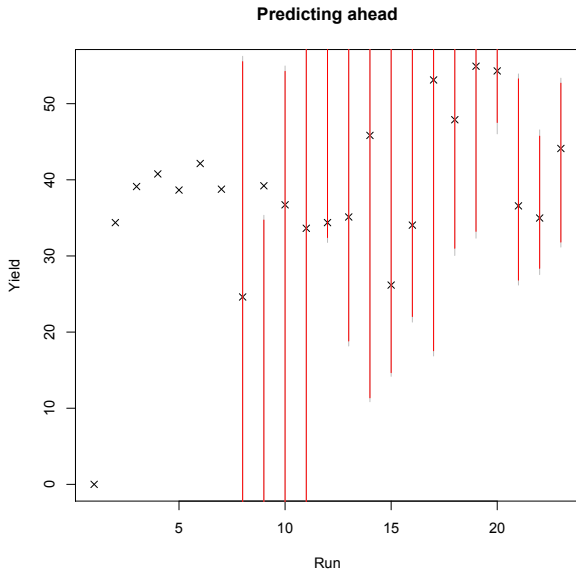
Initial design: 12 run design - 9 space-filling points (maximum-projection) + three centre points

Joseph et al. (2015)

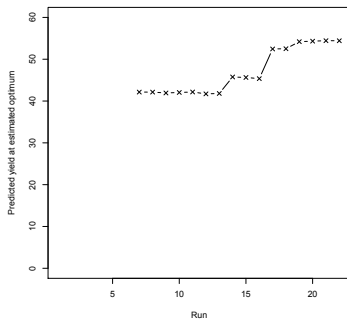
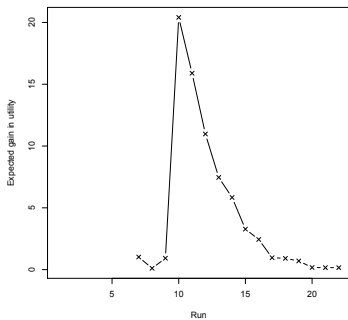
Stopping rule: $U(\mathbf{x}_{n+1}) - \mathbb{E}_g\{g(\mathbf{x}_n^*) \mid \mathbf{y}_{1:n}\} \leq 0.1$

Predicting forward using the sequential model

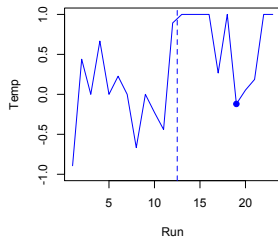
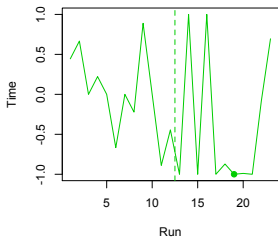
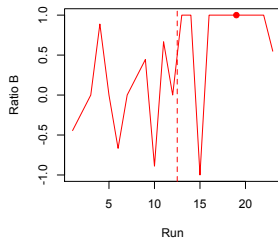
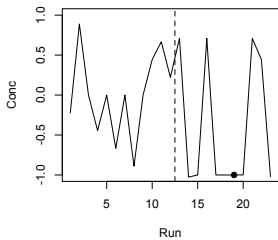
23 automated runs, with 11 chosen via EGU



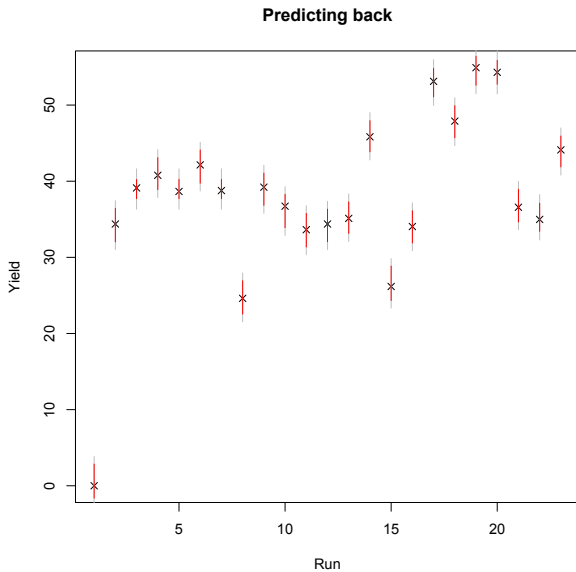
Progress of the algorithm



Design points



Predicting backwards from the final model



Robustifying the model

Outliers and unusual observations are a fact of life, and are not uncommon in chemistry experiments

Although we may wish to stop the system and investigate outlying observations manually, it is also useful to have automatic methods

To facilitate automatic outlier detection and removal, we extend the model such that

$$y_i \mid g, \sigma_\varepsilon^2, r_i \sim N(g(x_i), \sigma_\varepsilon^2 / r_i), \quad i = 1, \dots, n$$

With $r_i \sim \text{Gamma}(a_r/2, b_r/2)$, the response y_i has a non-standardised **t-distribution** conditional on g and σ_ε^2

We can use SMC and EGU with this model

cf Neal (1997)

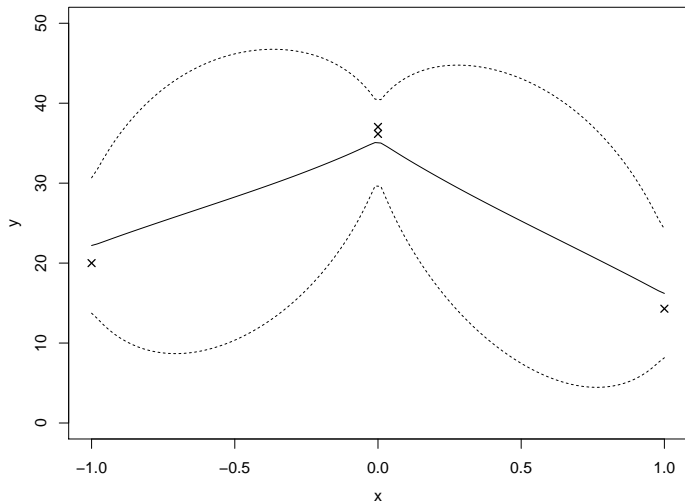
SMC for the t -response model

The parameters now consist of $(\boldsymbol{\theta}, \mathbf{r}_{1:n})$, and the dimension increases whenever we observe another data point. Thus IBIS is not applicable

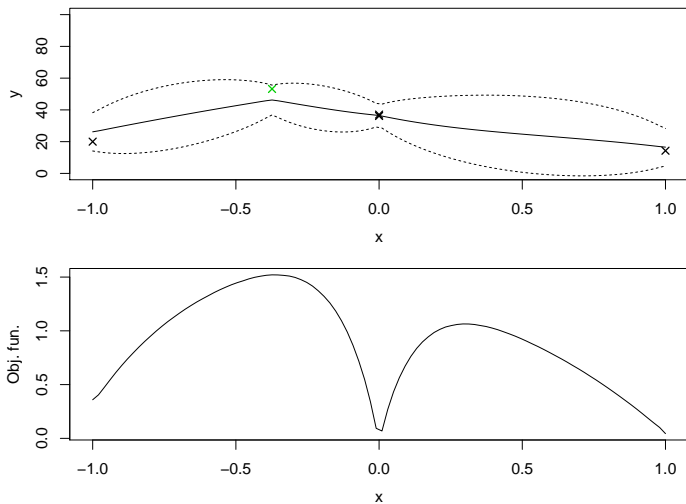
Algorithm on observing y_{n+1} ,

1. **Resample** particles $(\boldsymbol{\theta}_j, \mathbf{r}_{1:n,j})$ with resampling weights proportional to $\pi(y_{n+1} | \boldsymbol{\theta}_j, \mathbf{r}_{1:n,j}, \mathbf{y}_{1:n})$, approximated using Gauss-Laguerre quadrature
2. **Propagate** each of the resampled particles, by generating a value r_{n+1} from $\pi(r_{n+1} | \boldsymbol{\theta}, \mathbf{r}_{1:n}, \mathbf{y}_{1:(n+1)})$, via a Gibbs sampler
3. **Move** the particles using an MCMC kernel with invariant distribution $\pi(\boldsymbol{\theta}, \mathbf{r}_{1:(n+1)} | \mathbf{y}_{1:(n+1)})$

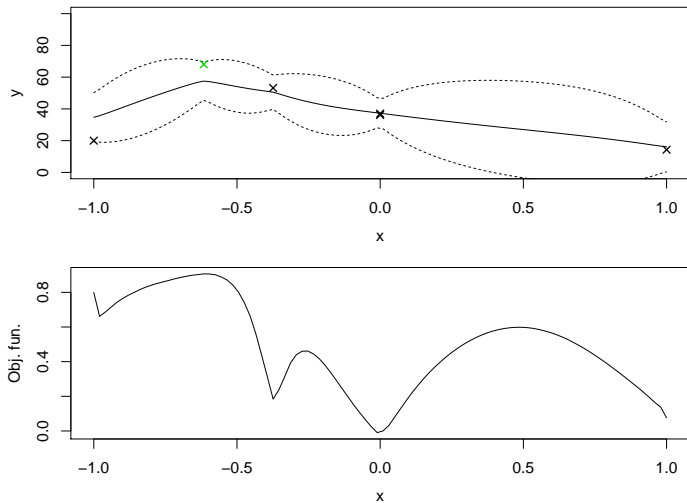
Example 2 - EGU, t -response model



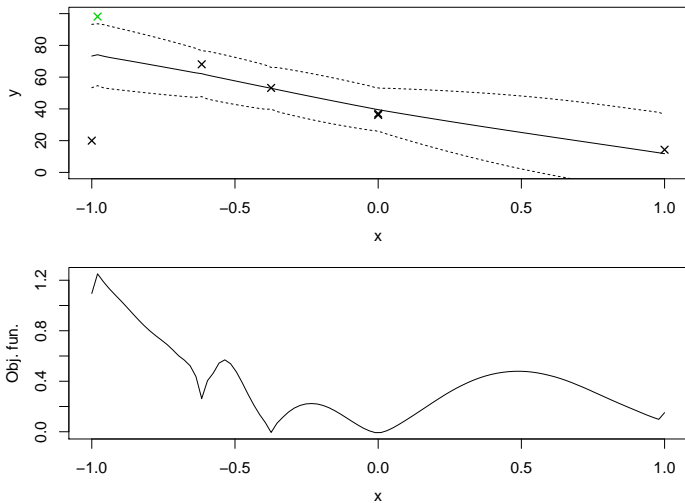
Example 2 - EGU, t -response model



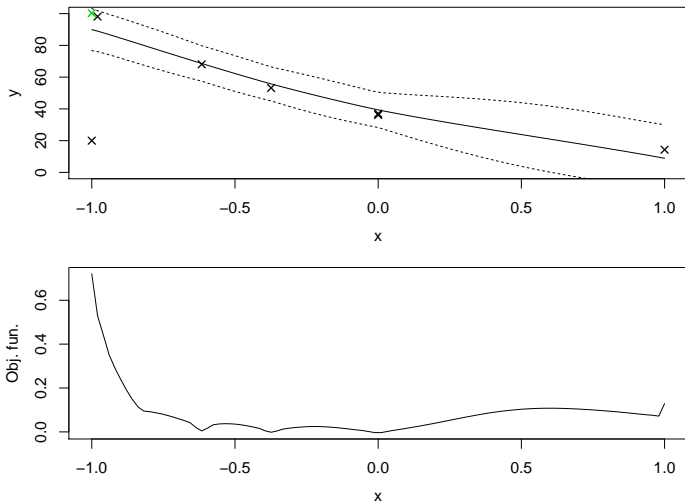
Example 2 - EGU, t -response model



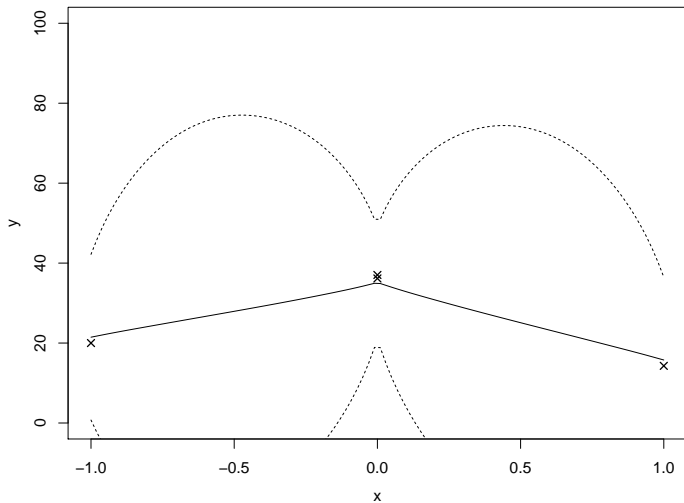
Example 2 - EGU, t -response model



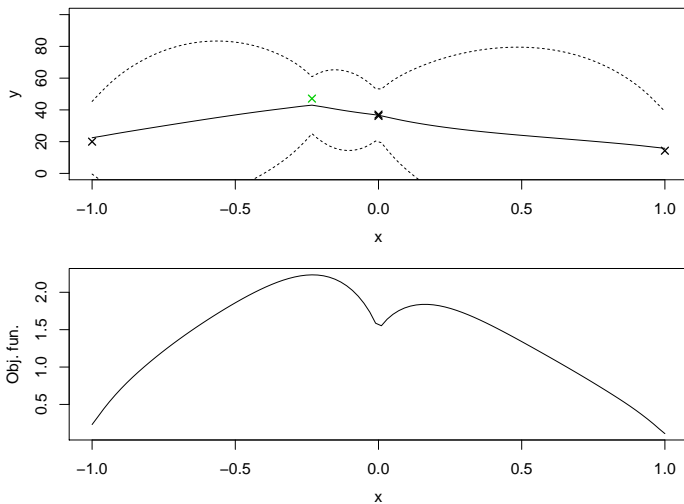
Example 2 - EGU, t -response model



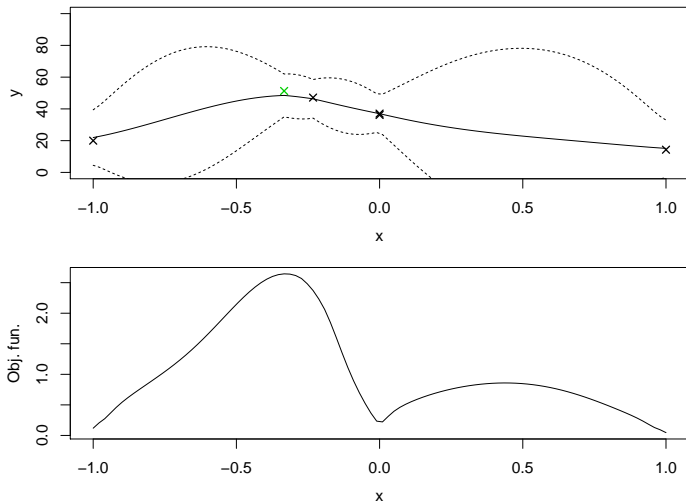
Example 2 - EI, normal-response model



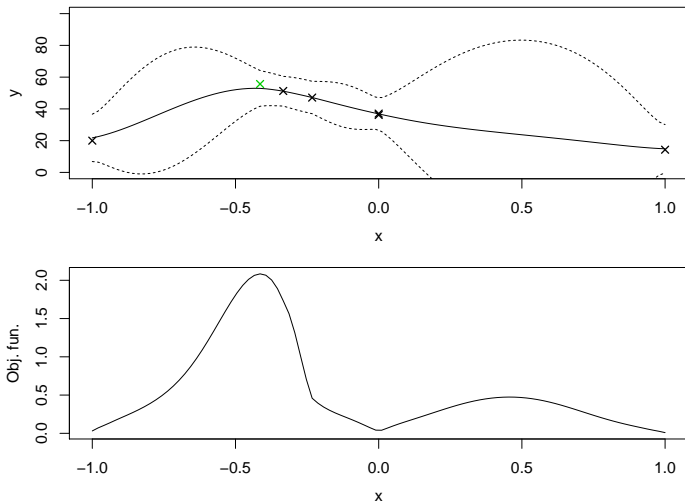
Example 2 - EI, normal-response model



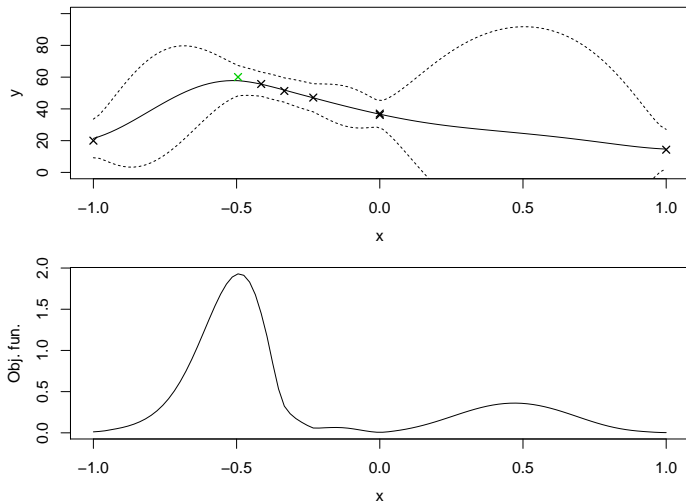
Example 2 - EI, normal-response model



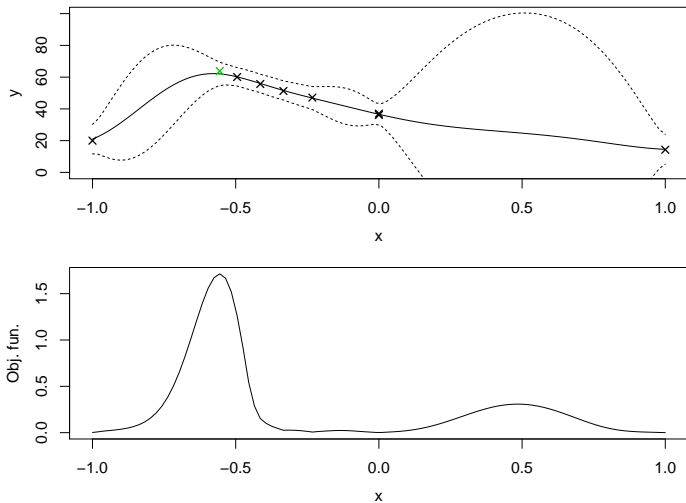
Example 2 - EI, normal-response model



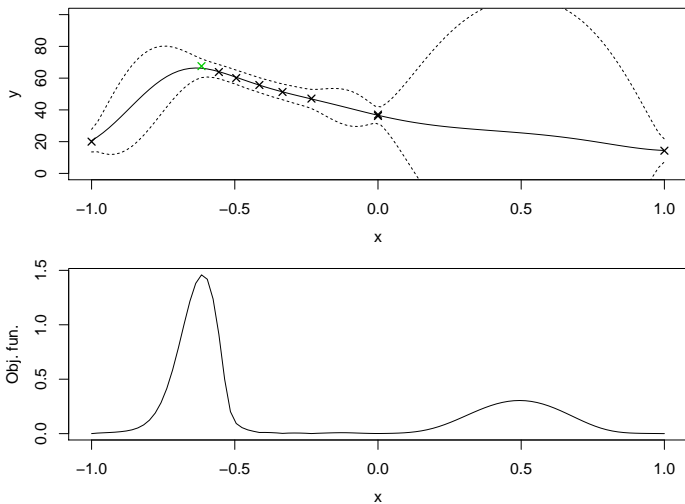
Example 2 - EI, normal-response model



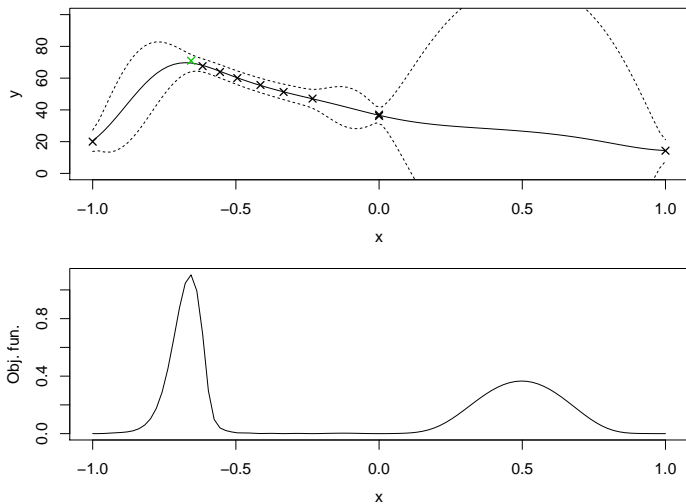
Example 2 - EI, normal-response model



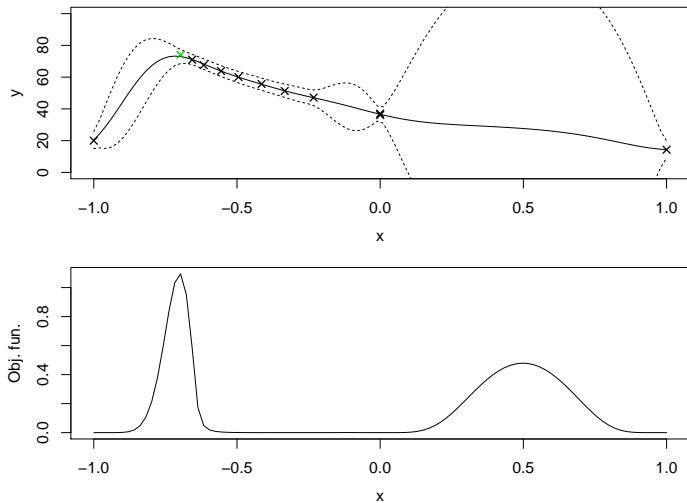
Example 2 - EI, normal-response model



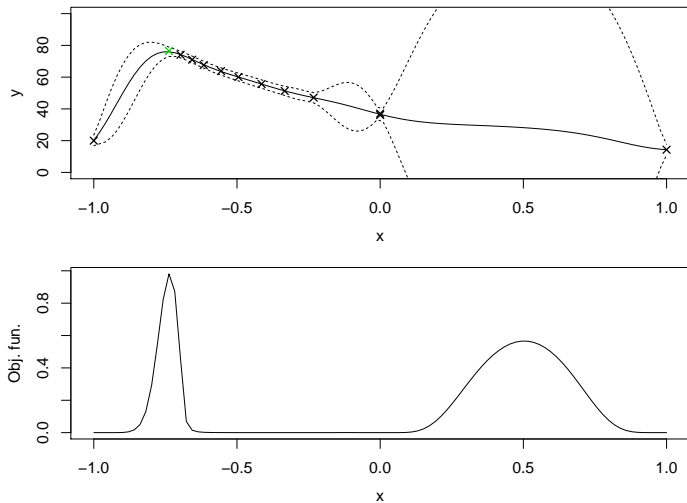
Example 2 - EI, normal-response model



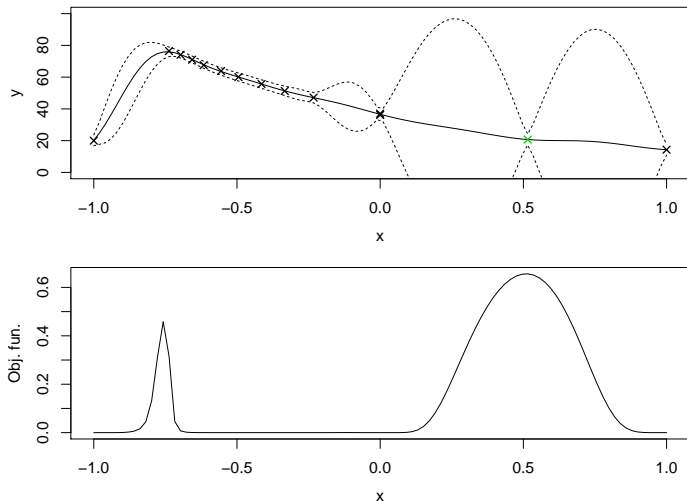
Example 2 - EI, normal-response model



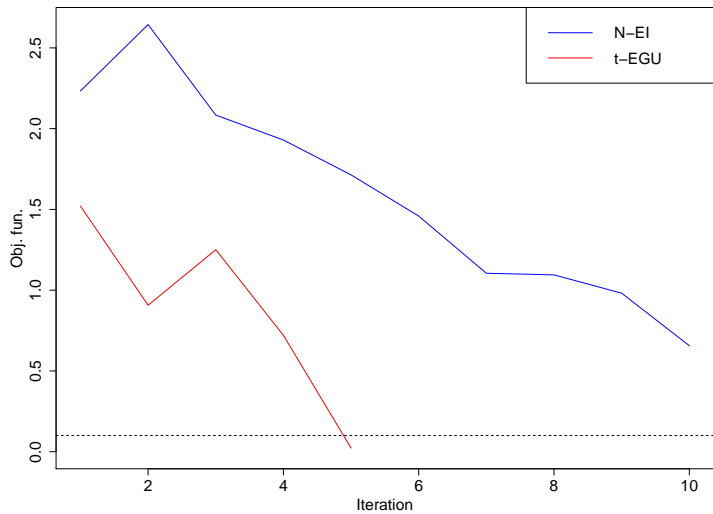
Example 2 - EI, normal-response model



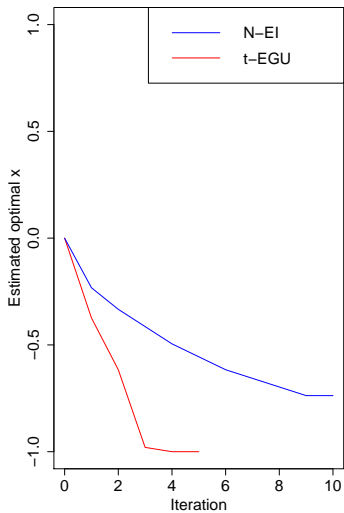
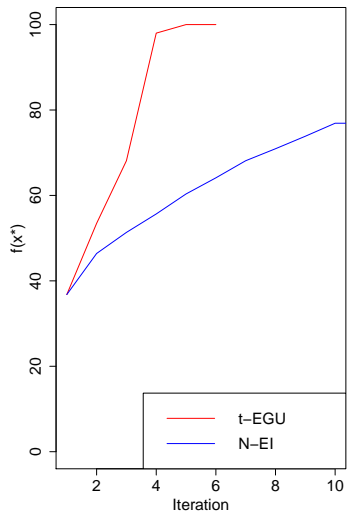
Example 2 - EI, normal-response model



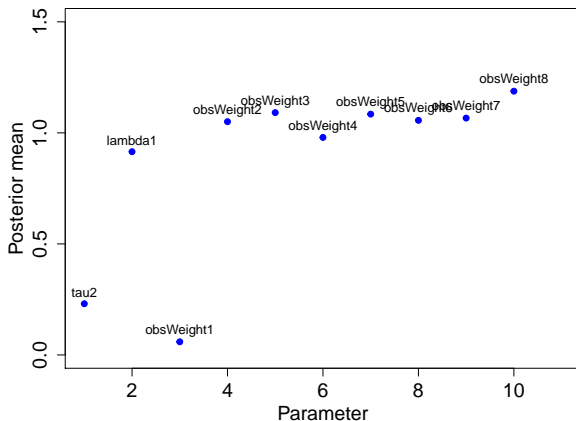
EGU and EI



Estimated maximum and its location



Mean posterior estimates from EGU algorithm



λ – correlation function hyperparameter; $\tau^2 = \sigma_\varepsilon^2 / \sigma_g^2$

Conclusions and future work

Conclusions

- ▶ Decision-theoretic design for (noisy) optimisation, with EI as a special case
- ▶ Robust modelling via combination of Gaussian process priors and weighted error variance
- ▶ Real examples via a completely automated system

Future work

- ▶ Batch sequential and non-myopic design
- ▶ More complex chemistry/faster data collection
- ▶ R package and chemistry software

References

- ▶ Brochu, E., Cora, V.M. & de Freitas, N. (2010). arXiv:1012.2599v1.
- ▶ Chopin, N. (2002). *Biometrika* 89, 539–552.
- ▶ Drovandi, C.C., McGree, J.M. & Pettitt, A.N. (2013). *CSDA*, 57, 320–335.
- ▶ Gramacy, R.B. & Polson, N.G. (2011). *Journal of Computational and Graphical Statistics*, 20, 102–118.
- ▶ Huang, D., Allen, T.T., Notz, W.I. & Zeng, N. (2006). *Journal of Global Optimization* 34, 441–466.
- ▶ Joseph, V.R., Gul, E. & Ba, S. (2015). *Biometrika*, in press.
- ▶ Jones, D.R., Schonlau, M. & Welch, W.J. (1998). *Journal of Global Optimization* 13, 455–492.
- ▶ Mockus, J. (1989). *Bayesian Approach to Global Optimization: Theory and Applications*. Kluwer.
- ▶ Neal, R.M. (1997). arXiv:9701026.
- ▶ Osbourne, M.A., Garnett, R. & Roberts, S.J. (2010). In *IEEE International Conference on Advanced Information Networking and Applications*.
- ▶ Rasmussen, C.E. & Williams, C.K.I. (2006). *Gaussian Processes for Machine Learning*, MIT Press.