

Categorical inputs, sensitivity analysis, optimization and importance tempering with **tgp** version 2, an **R** package for treed Gaussian process models

Robert B. Gramacy
Statistical Laboratory
University of Cambridge, United Kingdom
bobby@statslab.cam.ac.uk

Matthew Taddy
Booth School of Business
University of Chicago
matt.taddy@chicagobooth.edu

January 16, 2009

Abstract

This document describes the new features in version 2.0 of the **tgp** package for **R**, implementing treed Gaussian process (GP) models. The topics covered include methods for dealing with categorical inputs and excluding inputs from the tree or GP part of the model; fully Bayesian sensitivity analysis for inputs/covariates; sequential optimization of black-box functions; and a new Monte Carlo method for inference in multi-modal posterior distributions that combines simulated tempering and importance sampling. These additions extend the functionality of **tgp** across all models in the hierarchy: from Bayesian linear models, to CART, to treed Gaussian processes with jumps to the limiting linear model. It is assumed that the reader is familiar with the baseline functionality of the package, outlined in the first vignette [9].

Intended audience

This document is intended as a follow-on to the vignette for version 1.x of the **tgp** package [9]. It describes four new features that have been added to the package in version 2.x, and is divided into four essentially disjoint sections: on categorical inputs (Section 1), sensitivity analysis (Section 2), statistical optimization (Section 3), and importance tempering (Section

4). Each section will begin with a short mathematical introduction to the new feature or methodology and commence with extensive examples in R on synthetic and real data. This document has been authored in **Sweave** (try `help(Sweave)`). This means that the code quoted throughout is certified by R, and the **Stangle** command can be used to extract it. As with the first vignette, the R code in each of the sections to follow is also available as a demo in the package. Note that this tutorial was not meant to serve as an instruction manual. For more detailed documentation of the functions contained in the package, see the package help-manuals. At an R prompt, type `help(package=tgp)`. PDF documentation is also available on the world-wide-web.

<http://www.cran.r-project.org/doc/packages/tgp.pdf>

Each section starts by seeding the random number generator with `set.seed(0)`. This is done to make the results and analyses reproducible within this document, and in demo form. We recommend you try these examples with different seeds and see what happens. Usually the results will be similar, but sometimes (especially when the data (\mathbf{X}, \mathbf{Z}) is generated randomly) they may be quite different.

1 Non-real-valued, categorical and other inputs

Early versions of **tgp** worked best with real-valued inputs \mathbf{X} . While it was possible to specify ordinal, integer-valued, or even binary inputs, **tgp** would treat them the same as any other real-valued input. Two new arguments to `tgp.default.params`, and thus the ellipses `(...)` argument to the `b*` functions, provide a more natural way to model with non-real valued inputs.

Classical treed methods, such as CART [1], can cope quite naturally with categorical, binary, and ordinal, inputs. Categorical inputs can be encoded in binary, and splits can be proposed with rules such as $x_i < 1$. Once a split is made on a binary input, no further process is needed, marginally, in that dimension. Ordinal inputs can also be coded in binary, and thus treated as categorical, or treated as real-valued and handled in a default way. GP regression, however, handles such non-real-valued inputs less naturally, unless (perhaps) a custom and non-standard form of the covariance function is used. When inputs are scaled to lie in $[0, 1]$, binary-valued inputs x_i are always a constant distance apart—at the largest possible distance in the range. A separable correlation function width parameter d_i will tend to infinity (in the posterior) if the output does not vary with x_i , and will tend to zero if it is. Clearly, this functionality is more parsimoniously achieved by partitioning, e.g., using a tree. However, trees with fancy regression models at the leaves pose other problems, as discussed below.

Consider as motivation, the following modification of the Friedman data [5] (see also Section 3.5 of [9]). Augment 10 real-valued covariates in the data ($\mathbf{x} = \{x_1, x_2, \dots, x_{10}\}$) with one categorical indicator $I \in \{1, 2, 3, 4\}$ that can be encoded in binary as

$$1 \equiv (0, 0, 0) \quad 2 \equiv (0, 0, 1) \quad 3 \equiv (0, 1, 0) \quad 4 \equiv (1, 0, 0).$$

Now let the function that describes the responses (Z), observed with standard Normal noise, have a mean

$$E(Z|\mathbf{x}, I) = \begin{cases} 10 \sin(\pi x_1 x_2) & \text{if } I = 1 \\ 20(x_3 - 0.5)^2 & \text{if } I = 2 \\ 10x_4 + 5x_5 & \text{if } I = 3 \\ 10x_1 + 5x_2 + 20(x_3 - 0.5)^2 + 10 \sin(\pi x_4 x_5) & \text{if } I = 4 \end{cases} \quad (1)$$

that depends on the indicator I . Notice that when $I = 4$ the original Friedman data is recovered, but with the first five inputs in reverse order. Irrespective of I , the response depends only on $\{x_1, \dots, x_5\}$, thus combining nonlinear, linear, and irrelevant effects. When $I = 3$ the response is linear \mathbf{x} .

A new function has been included in the `tgp` package which facilitates generating random realizations from (1). Below we obtain 500 such random realizations for training purposes, and a further 1000 for testing.

```
> fb.train <- fried.bool(500)
> X <- fb.train[, 1:13]
> Z <- fb.train$Y
> fb.test <- fried.bool(1000)
> XX <- fb.test[, 1:13]
> ZZ <- fb.test$Ytrue
```

A separation into training and testing sets will be useful for later comparisons by RMSE. The names of the data frame show that the first ten columns encode \mathbf{x} and columns 11–13 encode the boolean representation of I .

```
> names(X)

[1] "X.1" "X.2" "X.3" "X.4" "X.5" "X.6" "X.7" "X.8"
[9] "X.9" "X.10" "I.1" "I.2" "I.3"
```

One, naïve approach to fitting this data would be to fit a treed GP LLM model ignoring the categorical inputs. But this model can only account for the noise, giving high RMSE, and so is not illustrated here. Clearly, the indicators must be included. One simple way to do so would be to posit a Bayesian CART model.

```
> fit1 <- bcart(X = X, Z = Z, XX = XX, verb = 0)
> rmse1 <- sqrt(mean((fit1$ZZ.mean - ZZ)^2))
> rmse1
```

```
[1] 2.983012
```

In this case the indicators are treated appropriately (as indicators), but in some sense so are the real-valued inputs as only constant models are fit at the leaves of the tree. Figure 1 shows that the tree does indeed partition on the indicators, and the other inputs, as expected.

One might expect a much better fit from a treed linear model to this data, since the response is linear in some of its inputs.

height=6, log(p)=435.729

```
graph TD; Root["X.4 <> 0.519264"] --> L1["I.2 <= 0"]; Root --> R1["I.1 <= 0"]; L1 --> L2["I.3 <= 0"]; L1 --> L3["X.4 <> 0.210015"]; R1 --> R2["I.2 <= 0"]; R1 --> R3["X.2 <> 0.502056"]; L2 --> L4["I.1 <= 0"]; L2 --> L5["4"]; L3 --> L6["5"]; L3 --> L7["6"]; R2 --> R4["I.3 <= 0"]; R2 --> R5["10"]; R3 --> R6["11"]; R3 --> R7["12"]; L4 --> L8["X.1 <> 0.499703"]; L4 --> L9["3"]; R4 --> R8["I.3 <= 0"]; R4 --> R9["9"]; L8 --> L10["1"]; L8 --> L11["2"]; R8 --> R10["7"]; R8 --> R11["8"]; L10 --> L12["0.0121  
34 obs"]; L10 --> L13["0.0063  
29 obs"]; L11 --> L14["0.0152  
62 obs"]; L5 --> L15["0.0024  
75 obs"]; L6 --> L16["0.002  
30 obs"]; L7 --> L17["0.0026  
36 obs"]; R5 --> R12["0.0059  
43 obs"]; R6 --> R13["0.0092  
38 obs"]; R7 --> R14["0.0064  
22 obs"]; R9 --> R15["0.0032  
74 obs"]; R11 --> R16["0.0057  
31 obs"]; R11 --> R17["0.0057  
26 obs"];
```

Decision tree structure and node details:

- Root: $X.4 \lessgtr 0.519264$
 - Left: $I.2 \leq 0$
 - Left: $I.3 \leq 0$
 - Left: $I.1 \leq 0$
 - Left: $X.1 \lessgtr 0.499703$
 - Left: Node 1 (34 obs, p=0.0121)
 - Right: Node 2 (29 obs, p=0.0063)
 - Right: Node 3 (62 obs, p=0.0152)
 - Right: Node 4 (75 obs, p=0.0024)
 - Right: $X.4 \lessgtr 0.210015$
 - Left: Node 5 (30 obs, p=0.002)
 - Right: Node 6 (36 obs, p=0.0026)
 - Right: $I.1 \leq 0$
 - Left: $I.2 \leq 0$
 - Left: $I.3 \leq 0$
 - Left: Node 7 (31 obs, p=0.0057)
 - Right: Node 8 (26 obs, p=0.0057)
 - Right: Node 9 (74 obs, p=0.0032)
 - Right: Node 10 (43 obs, p=0.0059)
 - Right: $X.2 \lessgtr 0.502056$
 - Left: Node 11 (38 obs, p=0.0092)
 - Right: Node 12 (22 obs, p=0.0064)

```
> fit2 <- btlm(X = X, Z = Z, XX = XX, verb = 0)
> rmse2 <- sqrt(mean((fit2$ZZ.mean - ZZ)^2))
> rmse2
```

Unfortunately, this is not the case—the RMSE obtained is similar to the one for the CART model. Figure 2 shows that the tree does indeed partition, but not on the indicator variables. When a linear model is used at the leaves of the tree the boolean indicators cannot be partitioned upon because doing so would cause the design matrix to become rank-deficient at the leaves of the tree (there would be a column of all zeros or all ones). A treed GP would have the same problem.

Let us revisit the treed LM model with `basemax` = 10.

```
> tgp.trees(fit2, "map")
```

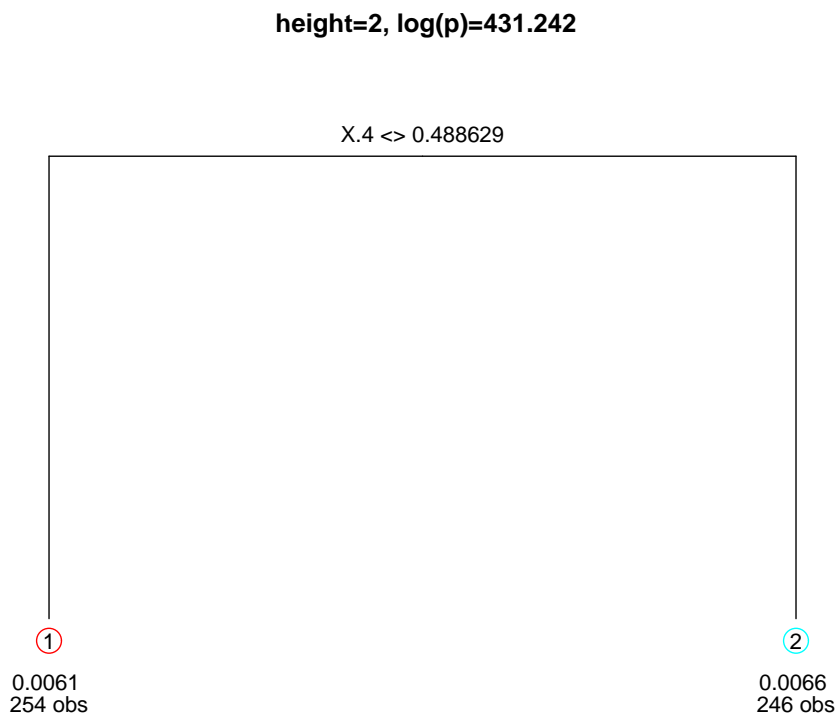


Figure 2: Diagrammatic depiction of the maximum *a posteriori* (MAP) tree for the boolean indicator version of the Friedman data in Eq. (1) using a Bayesian treed linear model.

```
> fit3 <- bt1m(X = X, Z = Z, XX = XX, basemax = 10,
+             verb = 0)
> rmse3 <- sqrt(mean((fit3$ZZ.mean - ZZ)^2))
> rmse3
```

```
[1] 1.649079
```

Figure 3 shows that the MAP tree does indeed partition on the indicators in an appropriate way—as well as on some other real-valued inputs—and the result is the lower RMSE we would expect.

A more high-powered approach would clearly be to treat all inputs as real-valued by fitting a GP at the leaves of the tree. Binary partitions are allowed on all inputs, \mathbf{X} and I , but treating the boolean indicators as real-valued in the GP is clearly inappropriate since it is known that the process does not vary smoothly over the 0 and 1 settings of the three boolean indicators representing the categorical input I .

```
> fit4 <- btgpllm(X = X, Z = Z, XX = XX, verb = 0)
> rmse4 <- sqrt(mean((fit4$ZZ.mean - ZZ)^2))
> rmse4
```

```
> tgp.trees(fit3, "map")
```

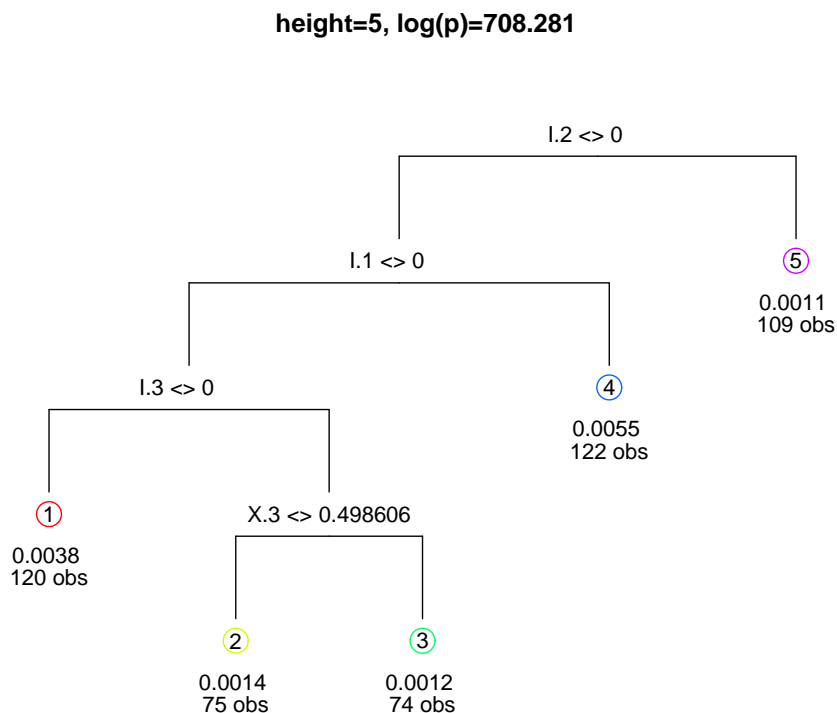


Figure 3: Diagrammatic depiction of the maximum a' posteriori (MAP) tree for the boolean indicator version of the Friedman data in Eq. (1) using a Bayesian treed linear model with the setting `basemax = 10`.

```
[1] 1.247021
```

Since the design matrices would become rank-deficient if the boolean indicators are partitioned upon, there was no partitioning in this example.

```
> fit4$gpcs
```

```

grow prune change swap
1      0      NA      NA      NA

```

Since there are large covariance matrices to invert, the MCMC inference is *very* slow. Still, the resulting fit (obtained with much patience) is better than the Bayesian CART and treed LM (with `basemax = 10`) ones, as indicated by the RMSE.

We would expect to get the best of both worlds if the setting `basemax = 10` were used when fitting the treed GP model, thus allowing partitioning on the indicators by guarding against rank deficient design matrices.

```
> fit5 <- btgpllm(X = X, Z = Z, XX = XX, basemax = 10,
+               verb = 0)
```

```
> rmse5 <- sqrt(mean((fit5$ZZ.mean - ZZ)^2))
> rmse5
```

```
[1] 0.4805412
```

And indeed this is the case.

The benefits go beyond producing full rank design matrices at the leaves of the tree. Loosely speaking, removing the boolean indicators from the GP part of the treed GP gives a more parsimonious model, without sacrificing any flexibility. The tree is able to capture all of the dependence in the response as a function of the indicator input, and the GP is the appropriate non-linear model for accounting for the remaining relationship between the real-valued inputs and outputs. We can look at the maximum *a posteriori* (MAP) tree, to see

```
> h <- fit1$post$height[which.max(fit1$post$lpst)]
> tgp.trees(fit5, "map")
```

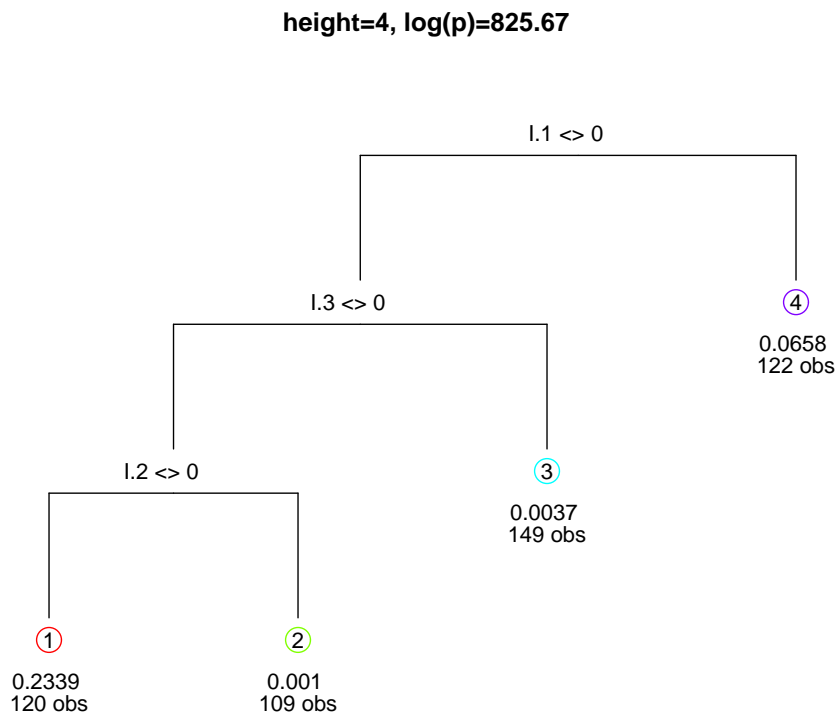


Figure 4: Diagrammatic depiction of the maximum *a posteriori* (MAP) tree for the boolean indicator version of the Friedman data in Eq. (1) using `basemax=10`.

that only (and all of) the indicators were partitioned upon in Figure 4. Further advantages to this approach include speed (a partitioned model gives smaller covariance matrices to invert) and improved mixing in the Markov chain when a separable covariance function is used. Note that using a non-separable covariance function in the presence of indicators would result in a poor fit. Good range (d) settings for the indicators would not necessarily coincide with good range settings for the real-valued inputs.

A complimentary setting, `splitmin`, allows the user to specify the first column of the inputs \mathbf{X} upon which treed partitioning is allowed. From Section 3.5 of the first `tgp` vignette [9], it was concluded that the original formulation of Friedman data was stationary, and thus treed partitioning is not required to obtain a good fit. The same would be true of the response in (1) after conditioning on the indicators. Therefore, the most parsimonious model would use `splitmin = 11`, in addition to `basemax = 10`, so that only \mathbf{X} are under the GP, and only I under the tree. Fewer viable candidate inputs for treed partitioning should yield improved mixing in the Markov chain, and thus lower RMSE.

```
> fit6 <- btgppllm(X = X, Z = Z, XX = XX, basemax = 10,
+   splitmin = 11, verb = 0)
> rmse6 <- sqrt(mean((fit6$ZZ.mean - ZZ)^2))
> rmse6
```

```
[1] 0.3844569
```

Needless to say, it is important that the input \mathbf{X} have columns which are ordered appropriately before the `basemax` and `splitmin` arguments can be properly applied. Future versions of `tgp` will have a formula-based interface to handle categorical (`factors`) and other inputs more like other R regression routines, e.g., `lm` and `glm`.

2 Analysis of sensitivity to inputs

Sensitivity Analysis (SA) is a resolving of the sources of output variability by apportioning elements of this variation to different sets of input variables. The edited volume by Saltelli et al. [25] provides a complete overview of the field. The analysis of response variability is useful in a variety of different settings. For example, when there is a large number of input variables over which an objective function is to be optimized, typically only a small subset will be influential within the confines of their uncertainty distribution. SA can be used to reduce the input space of such optimizations [31]. Other authors have used SA to assess the risk associated with dynamic factors affecting the storage of nuclear waste [12], and to investigate the uncertainty characteristics of a remote sensing model for the reflection of light by surface vegetation [18]. The `sens` function adds to `tgp` a suite of tools for global sensitivity analysis, and enables “out-of-the-box” estimation of valuable sensitivity indices for any regression relationship that may be modeled by a member of the `tgp` family.

The type of sensitivity analysis provided by `tgp` falls within the paradigm of global sensitivity analysis, wherein the variability of the response is investigated with respect to a probability distribution over the entire input space. The recent book by Saltelli et al. [26] serves as a primer on this field. Global SA is inherently a problem of statistical inference, as evidenced by the interpolation and estimation required in a study of the full range of inputs. This is in contrast with the analytical nature of local SA, which involves derivative-based investigation of the stability of the response over a small region of inputs. We will ignore local SA for the remainder of this document.

The sensitivity of a response z to a changing input \mathbf{x} is always considered in relation to a specified *uncertainty distribution*, defined by the density $u(\mathbf{x})$, and the appropriate marginal densities $u_i(x_i)$. What is represented by the uncertainty distribution changes depending upon the context. The canonical setup has that z is the response from a complicated physics or engineering simulation model, with tuning parameters \mathbf{x} , that is used to predict physical phenomena. In this situation, $u(\mathbf{x})$ represents the experimentalist’s uncertainty about real-world values of \mathbf{x} . In optimization problems, the uncertainty distribution can be used to express prior information from experimentalists or modelers on where to look for solutions. Finally, in the case of observational systems (such as air-quality or smog levels), $u(\mathbf{x})$ may be an estimate of the density governing the natural occurrence of the \mathbf{x} factors (e.g., air-pressure, temperature, wind, and cloud cover). In this setup, SA attempts to resolve the natural variability of z .

The most common notion of sensitivity is tied to the relationship between conditional and marginal variance for z . Specifically, variance-based methods decompose the variance of the objective function, with respect to the uncertainty distribution on the inputs, into variances of conditional expectations. These are a natural measure of the output association with specific sets of variables and provide a basis upon which the importance of individual inputs may be judged. The other common component of global SA is an accounting of the main effects for each input variable, $\mathbb{E}_{u_j}[z|x_j]$, which can be obtained as a by-product of the variance analysis.

Our variance-based approach to SA is a version of the method of Sobol’, wherein a deterministic objective function is decomposed into summands of functions on lower dimensional subsets of the input space. Consider the decomposition

$$f(x_1, \dots, x_d) = f_0 + \sum_{j=1}^d f_j(x_j) + \sum_{1 \leq i < j \leq d} f_{ij}(x_j, x_i) + \dots + f_{1, \dots, d}(x_1, \dots, x_d).$$

Sobol’ shows [30] that, in the case of uncorrelated input variables, a unique orthogonal decomposition of this form is built through conditional marginalization of the original function with respect to u .

It is natural to extend this setting to response functions modeled conditionally as stochastic processes. Instead of a response function $f(\mathbf{x})$, we speak of the full conditional distribution for response z given inputs \mathbf{x} . The marginal distribution for z may be defined as the marginal response distribution when \mathbf{x} is allowed to vary according to u . Global SA is then focused on the marginal conditional response distributions which arise when z has been marginalized over one subset of covariates and the complement of this subset is allowed to vary according to a marginalized uncertainty distribution. In particular, we can obtain the marginal conditional expectation $\mathbb{E}[z|\mathbf{x}_J = \{x_j : j \in J\}] = \int_{\mathbb{R}^{d-d_J}} \mathbb{E}[z|\mathbf{x}]u(\mathbf{x})d\mathbf{x}_{-J}$, where $J = \{j_1, \dots, j_{d_J}\}$ indicates a subset of input variables, and the marginal uncertainty density is given by $u_J(\mathbf{x}_J) = \int_{\mathbb{R}^{d-d_J}} u(\mathbf{x})d\{x_i : i \notin J\}$. SA concerns the variability of $\mathbb{E}[z|\mathbf{x}_J]$ with respect to changes in \mathbf{x}_J according to $u_J(\mathbf{x}_J)$.

In this new setting, if u is such that the inputs are uncorrelated, Sobol’s equation in (2) may be re-written as the variance decomposition, $\text{var}(\mathbb{E}[z|\mathbf{x}]) = \sum_{j=1}^d V_j + \sum_{1 \leq i < j \leq d} V_{ij} +$

$\dots + V_{1,\dots,d}$, where $V_j = \text{var}(\mathbb{E}[z|x_j])$, $V_{ij} = \text{var}(\mathbb{E}[z|x_i, x_j]) - V_i - V_j$, and so on. Clearly, when the inputs are correlated this identity no longer holds (although a “less-than-or-equal-to” inequality is always true). But it is useful to retain an intuitive interpretation of the V_J ’s as a portion of the overall marginal variance.

Our global SA will focus on the related sensitivity indices $S_J = V_J/\text{var}(z)$ which, as can be seen in the above equation, will sum to one over all possible J and are bounded to $[0, 1]$. These S_J ’s provide a natural measure of the *importance* of a set J of inputs and serve as the basis for an elegant analysis of sensitivity. The `sens` function allows for easy calculation of two very influential sensitivity indices associated with each input: the 1st order for the j th input variable,

$$S_j = \frac{\text{var}(\mathbb{E}[z|x_j])}{\text{var}(z)}, \quad (2)$$

and the total sensitivity for input j ,

$$T_j = \frac{\mathbb{E}[\text{var}(z|\mathbf{x}_{-j})]}{\text{var}(z)}. \quad (3)$$

The 1st order indices measure the portion of variability that is due to variation in the main effects for each input variable, while the total effect indices measure the portion of variability that is due to total variation in each input. From the identity $\mathbb{E}[\text{var}(z|\mathbf{x}_{-j})] = \text{var}(z) - \text{var}(\mathbb{E}[z|\mathbf{x}_{-j}])$, it can be seen that T_j measures the *residual* variability remaining after variability in all other inputs has been apportioned and that, for a deterministic response and uncorrelated input variables, $T_j = \sum_{J:j \in J} S_J$. This implies that the difference between T_j and S_j provides a measure of the variability in z due to interaction between input j and the other input variables. A large difference may lead the investigator to consider other sensitivity indices to determine where this interaction is most influential, and this is often a key aspect of the dimension-reduction that SA provides for optimization problems.

2.1 Monte Carlo integration for sensitivity indices

Due to the many integrals involved, estimation of the sensitivity indices is not straightforward. The influential paper by Oakley & O’Hagan [21] describes an empirical Bayes estimation procedure for the sensitivity indices; however, some variability in the indices is lost due to plug-in estimation of GP model parameters and, more worryingly, the variance ratios are only possible in the form of a ratio of expected values. Likelihood based approaches have also been proposed [33, 18]. The technique implemented in `tgp` is, in contrast, fully Bayesian and provides a complete accounting of the uncertainty involved. Briefly, at each iteration of an MCMC chain sampling from the treed GP posterior, output is predicted over a large (carefully chosen) set of input locations. Conditional on this predicted output, the sensitivity indices can be calculated via Monte Carlo integration. By conditioning on predicted response (and working as though it were the observed response), we obtain a posterior sample of the indices, incorporating variability from both the integral estimation and uncertainty about the function output. In particular, the `sens` function includes a `model`

argument which allows for SA based on any of the prediction models (the **b*** functions) in **tgp**.

Our Monte Carlo integration is based upon Saltelli's [28] efficient Latin hypercube sampling (LHS) scheme for estimation of both 1st order and total effect indices. We note that the estimation is only valid for uncorrelated inputs, such that $u(\mathbf{x}) = \prod_{j=1}^d u_j(x_j)$. The **sens** function only allows for uncertainty distributions of this type (in fact, the marginal distributions also need to be bounded), but this is a feature of nearly every "out-of-the-box" approach to SA. Studies which concern correlated inputs will inevitably require modeling for this correlation, whereas most regression models (including those in **tgp**) condition on the inputs and ignore the joint density for \mathbf{x} . Refer to the work of Saltelli & Tarantola [27] for an example of SA with correlated inputs.

We now briefly describe the integration scheme. The 2nd moment is a useful intermediate quantity in variance estimation, and we define

$$D_J = \mathbb{E} [\mathbb{E}^2 [z|\mathbf{x}_J]] = \int_{\mathbb{R}^{d_J}} \mathbb{E}^2 [z|\mathbf{x}_J] u_J(\mathbf{x}_J) d(\mathbf{x}_J).$$

Making use of an auxiliary variable,

$$\begin{aligned} D_J &= \int_{\mathbb{R}^{d_J}} \left[\int_{\mathbb{R}^{d-J}} \mathbb{E} [z|\mathbf{x}_J, \mathbf{x}_{-J}] u_{-J}(\mathbf{x}_{-J}) d\mathbf{x}_{-J} \int_{\mathbb{R}^{d-J}} \mathbb{E} [z|\mathbf{x}_J, \mathbf{x}'_{-J}] u_{-J}(\mathbf{x}'_{-J}) d\mathbf{x}'_{-J} \right] u_J(\mathbf{x}_J) d\mathbf{x}_J \\ &= \int_{\mathbb{R}^{d+d-J}} \mathbb{E} [z|\mathbf{x}_J, \mathbf{x}_{-J}] \mathbb{E} [z|\mathbf{x}_J, \mathbf{x}'_{-J}] u_{-J}(\mathbf{x}_{-J}) u_{-J}(\mathbf{x}'_{-J}) u_J(\mathbf{x}_J) d\mathbf{x} d\mathbf{x}'_{-J}. \end{aligned}$$

Thus, in the case of independent inputs,

$$D_J = \int_{\mathbb{R}^{d+d-J}} \mathbb{E} [z|\mathbf{x}] \mathbb{E} [z|\mathbf{x}_J, \mathbf{x}'_{-J}] u_{-J}(\mathbf{x}'_{-J}) u(\mathbf{x}) d\mathbf{x}'_{-J} d\mathbf{x}.$$

Note that at this point, if the inputs had been correlated, the integral would have been instead with respect to the joint density $u(\mathbf{x})u(\mathbf{x}'_{-J}|\mathbf{x}_J)$, leading to a more difficult integral estimation problem.

Recognizing that $S_j = (D_j - \mathbb{E}^2[z])/\text{var}(z)$ and $T_j = 1 - ((D_{-j} - \mathbb{E}^2[z]))/\text{var}(z)$, we need estimates of $\text{var}(z)$, $\mathbb{E}^2[z]$, and $\{(D_j, D_{-j}) : j = 1, \dots, d\}$ to calculate the sensitivity indices. Given a LHS M proportional to $u(\mathbf{x})$,

$$M = \begin{bmatrix} s_{1_1} & \cdots & s_{1_d} \\ \vdots & & \\ s_{m_1} & \cdots & s_{m_d} \end{bmatrix},$$

it is possible to estimate $\widehat{\mathbb{E}[z]} = \frac{1}{m} \sum_{k=1}^m \mathbb{E}[z|\mathbf{s}_k]$ and $\widehat{\text{var}[z]} = \frac{1}{m} \mathbb{E}^T[z|M] \mathbb{E}[z|M] - \widehat{\mathbb{E}[z]} \widehat{\mathbb{E}[z]}$, where the convenient notation $\mathbb{E}[z|M]$ is taken to mean $[\mathbb{E}[z|\mathbf{s}_1] \cdots \mathbb{E}[z|\mathbf{s}_m]]^T$. All that remains is to estimate the D 's. Define a second LHS M' proportional to u of the same size as

M and say that N_J is M' with the J columns replaced by the corresponding columns of M . Hence,

$$N_j = \begin{bmatrix} s'_{1_1} \cdots s_{1_j} \cdots s'_{1_d} \\ \vdots \\ s'_{m_1} \cdots s_{m_j} \cdots s'_{m_d} \end{bmatrix} \quad \text{and} \quad N_{-j} = \begin{bmatrix} s_{1_1} \cdots s'_{1_j} \cdots s_{1_d} \\ \vdots \\ s_{m_1} \cdots s'_{m_j} \cdots s_{m_d} \end{bmatrix}.$$

The estimates are then $\hat{D}_j = \mathbb{E}^T[z|M]\mathbb{E}[z|N_j]/(m-1)$ and $\hat{D}_{-j} = \mathbb{E}^T[z|M']\mathbb{E}[z|N_j]/(m-1) \approx \mathbb{E}^T[z|M]\mathbb{E}[z|N_{-j}]/(m-1)$. Along with the variance and expectation estimates, these can be plugged into equations for S_j and T_j in (2–3) to obtain \hat{S}_j and \hat{T}_j . Note that Saltelli recommends the use of the alternative estimate $\widehat{\mathbb{E}^2[z]} = \frac{1}{n-1}\mathbb{E}^T[z|M]\mathbb{E}[z|M']$ in calculating 1st order indices, as this brings the index closer to zero for non-influential variables. However, it has been our experience that these biased estimates can be unstable, and so **tgp** uses the standard $\widehat{\mathbb{E}^2[z]} = \mathbb{E}[z]\mathbb{E}[z]$ throughout.

The set of input locations which need to be evaluated for each calculation of the indices is $\{M, M', N_1, \dots, N_d\}$, and if m is the sample size for the Monte Carlo estimate this scheme requires $m(d+2)$ function evaluations. Hence, at each MCMC iteration of the model fitting, the $m(d+2)$ locations are drawn randomly according the LHS scheme, creating a random prediction matrix, **XX**. By allowing random draws of the input locations, the Monte Carlo error of the integral estimates will be included in the posterior variability of the indices and the posterior moments will not be dependent upon any single estimation input set. Using predicted output over this input set, a single realization of the sensitivity indices is calculated through Saltelli's scheme. At the conclusion of the MCMC, we have a representative sample from the posterior for **S** and **T**. The averages for these samples are unbiased estimates of the posterior mean, and the variability of the sample is representative of the complete uncertainty about model sensitivity.

Since a subset of the predictive locations (M and M') are actually a LHS proportional to the uncertainty distribution, we can also estimate the main effects at little extra computational cost. At each MCMC iteration, a one-dimensional nonparametric regression is fit through the scatterplot of $[s_{1_j}, \dots, s_{m_j}, s'_{1_j}, \dots, s'_{m_j}]$ vs. $[\mathbb{E}[z|M], \mathbb{E}[z|M']]$ for each of the $j = 1, \dots, d$ input variables. The resultant regression estimate provides a realization of $\mathbb{E}[z|x_j]$ over a grid of x_j values, and therefore a posterior draw of the main effect curve. Thus, at the end of the MCMC, we have not only unbiased estimates of the main effects through posterior expectation, but also a full accounting of our uncertainty about the main effect curve. This technique is not very sensitive to the method of non-parametric regression, since $2m$ will typically represent a very large sample in one-dimension. The estimation in **tgp** uses a moving average with squared distance weights and a window containing the `span*2m` nearest points (the `span` argument defaults to 0.3).

2.2 Examples

We illustrate the capabilities of the **sens** function by looking at the Friedman function considered earlier in this vignette. The function that describes the responses (Z), observed

with standard Normal noise, has mean

$$E(Z|\mathbf{x}) = 10 \sin(\pi x_1 x_2) + 20(x_3 - 0.5)^2 + 10x_4 + 5x_5. \quad (4)$$

A sensitivity analysis can be based upon any of the available regression models (e.g., `bt1m`, `bgp`, or `btgp`); we choose to specify `model=btgp1lm` for this example. The size of each LHS used in the integration scheme is specified through `nn.lhs`, such that this is equivalent to m in the above algorithm description. Thus the number of locations used for prediction—the size of the random `XX` prediction matrix—is `nn.lhs*(ncol(X)+2)`. In addition, the window for moving average estimation of the main effects is `span*2*nn.lhs` (independent of this, an `ngrid` argument with a default setting of `ngrid=100` dictates the number of grid points in each input dimension upon which main effects will be estimated).

```
> f <- friedman.1.data(250)
```

This function actually generates 10 covariates, the last five of which are completely un-influential. We'll include one of these (x_6) to show what the sensitivity analysis looks like for unrelated variables.

```
> Xf <- f[, 1:6]
> Zf <- f$Y
> sf <- sens(X = Xf, Z = Zf, nn.lhs = 600, model = bgp1lm,
+          verb = 0)
```

The progress indicators printed to the screen (for `verb > 0`) are the same as would be obtained under the specified regression `model`—`bgp1lm` in this case—so we suppress them here. All of the same options (e.g., `BTE`, `R`, etc.) apply, although if using the `trace` capabilities one should be aware that the `XX` matrix is changing throughout the MCMC. The `sens` function returns a `"tgp"`-class object, and all of the SA related material is included in the `sens` list within this object.

```
> names(sf$sens)

[1] "par"      "Xgrid"    "ZZ.mean"  "ZZ.q1"    "ZZ.q2"    "S"
[7] "T"
```

The object provides the SA parameters (`par`), the grid of locations for main effect prediction (`Xgrid`), the mean and interval estimates for these main effects (`ZZ.mean`, `ZZ.q1`, and `ZZ.q2`), and full posterior via samples of the sensitivity indices (`S` and `T`).

The plot function for `"tgp"`-class objects now provides a variety of ways to visualize the results of a sensitivity analysis. This capability is accessed by specifying `layout="sens"` in the standard plot command. By default, the mean posterior main effects are plotted next to boxplot summaries of the posterior sample for each S_j and T_j index, as in Figure 5.

A further note on the role played by `nn.lhs`: As always, the quality of the regression model estimate depends on the length of the MCMC. But now, the quality of sensitivity

```
> plot(sf, layout = "sens", legendloc = "topleft")
```

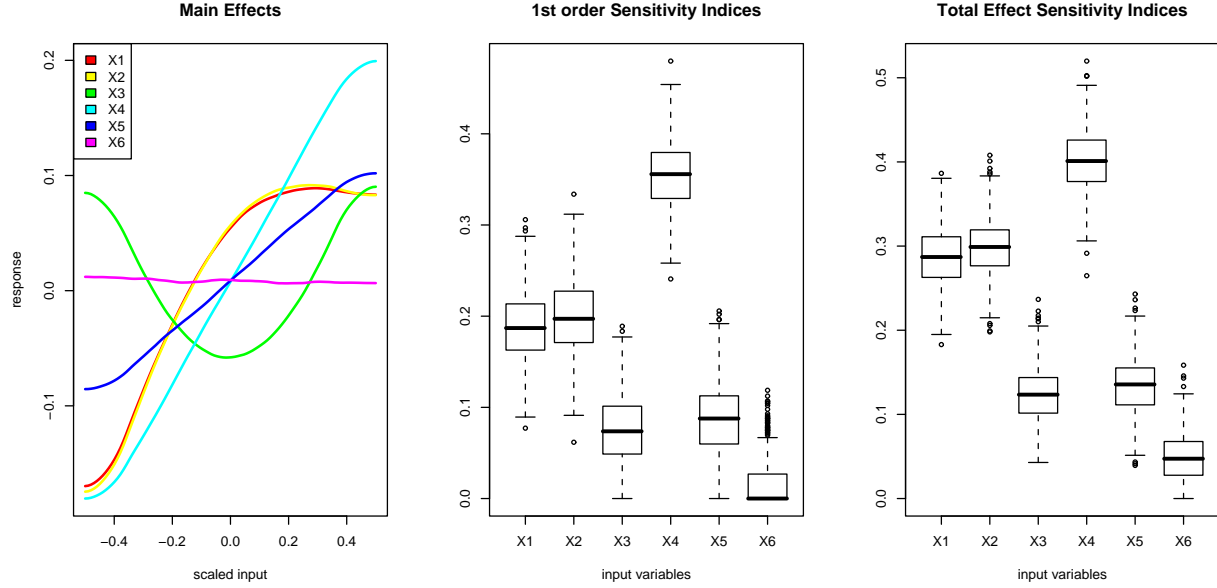


Figure 5: Full sensitivity analysis results for the Friedman function.

analysis is directly influenced by the size of the LHS used for integral approximation; as with any Monte Carlo integration scheme, the sample size (i.e., `nn.lhs`) must increase with the dimensionality of the problem. In particular, it can be seen in the estimation procedure described above that the total sensitivity indices (the T_j 's) are not forced to be non-negative. If negative values occur it is necessary to increase `nn.lhs`. In any case, the `plot.tgp` function changes any of the negative values to zero for purposes of illustration.

The `maineff` argument can be used to plot either selected main effects (Figure 6), or just the sensitivity indices (Figure 7). Note that the posterior intervals shown in these plots represent uncertainty about both the function response and the integration estimates; this full quantification of uncertainty is not presently available in any alternative SA procedures. These plots may be compared to what we know about the Friedman function (refer to Eq. (4)) to evaluate the analysis. The main effects correspond to what we would expect: sine waves for x_1 and x_2 , a parabola for x_3 , and linear effects for x_4 and x_5 . The sensitivity indices show x_1 and x_2 contributing roughly equivalent amounts of variation, while x_4 is relatively more influential than x_5 . Full effect sensitivity indices for x_3 , x_4 , and x_5 are roughly the same as the first order indices, but (due to the interaction in the Friedman function) the sensitivity indices for the total effect of x_1 and x_2 are significantly larger than the corresponding first order indices. Finally, our SA is able to determine that x_6 is unrelated to the response.

This analysis assumes the default uncertainty distribution, which is uniform over the range of input data. In other scenarios, it is useful to specify an informative $u(\mathbf{x})$. In the `sens` function, properties of u are defined through the `rect`, `shape`, and `mode` arguments. To guarantee integrability of our indices, we have restricted ourselves to bounded uncertainty

```
> par(mar = c(4, 2, 4, 2), mfrow = c(2, 3))
> plot(sf, layout = "sens", maineff = t(1:6))
```

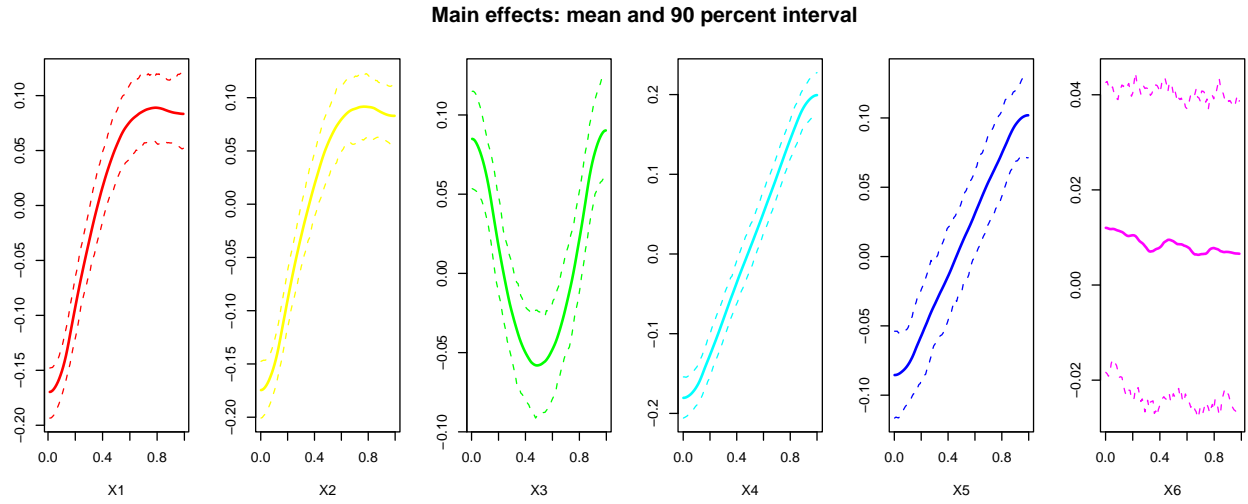


Figure 6: Freidman function main effects, with posterior 90% intervals.

```
> plot(sf, layout = "sens", maineff = FALSE)
```

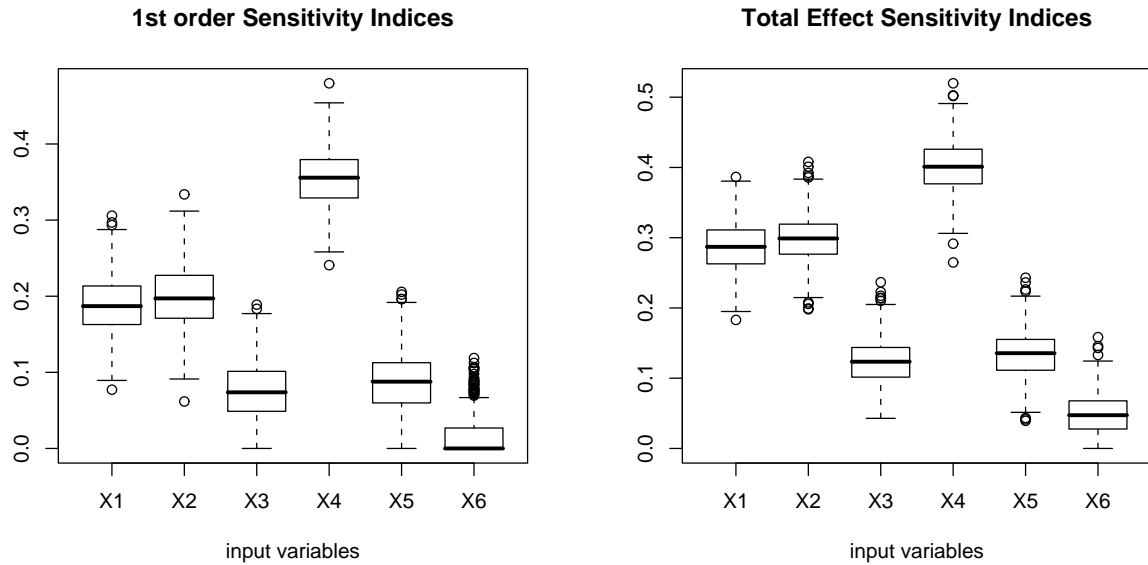


Figure 7: Sensitivity indices for the Freidman function.

distributions. Hence, `rect` defines these bounds. In particular, this defines the domain from which the LHSs are to be taken. We then use independent scaled beta distributions, parameterized by the `shape` parameter and distribution `mode`, to define an informative uncertainty distribution over this domain.

As an example of sensitivity analysis under an informative uncertainty distribution, consider the `airquality` data available with the base distribution of R. This data set contains daily readings of mean ozone in parts per billion (*Ozone*), solar radiation (*Solar.R*), wind speed (*Wind*), and maximum temperature (*Temp*) for New York City, between May 1 and September 30, 1973. Suppose that we are interested in the sensitivity of air quality to natural changes in *Solar.R*, *Wind*, and *Temp*. For convenience, we will build our uncertainty distribution while assuming independence between these inputs. Hence, for each variable, the input uncertainty distribution will be a scaled beta with `shape=2`, and `mode` equal to the data mean.

```
> X <- airquality[, 2:4]
> Z <- airquality$Ozone
> rect <- t(apply(X, 2, range, na.rm = TRUE))
> mode <- apply(X, 2, mean, na.rm = TRUE)
> shape <- rep(2, 3)
```

LHS samples from the uncertainty distribution are shown in Figure (8)

```
> Udraw <- lhs(300, rect = rect, mode = mode, shape = shape)
> par(mfrow = c(1, 3), mar = c(4, 2, 4, 2))
> for (i in 1:3) {
+   hist(Udraw[, i], breaks = 10, xlab = names(X)[i],
+       main = "", ylab = "", border = grey(0.9),
+       col = 8)
+ }
```

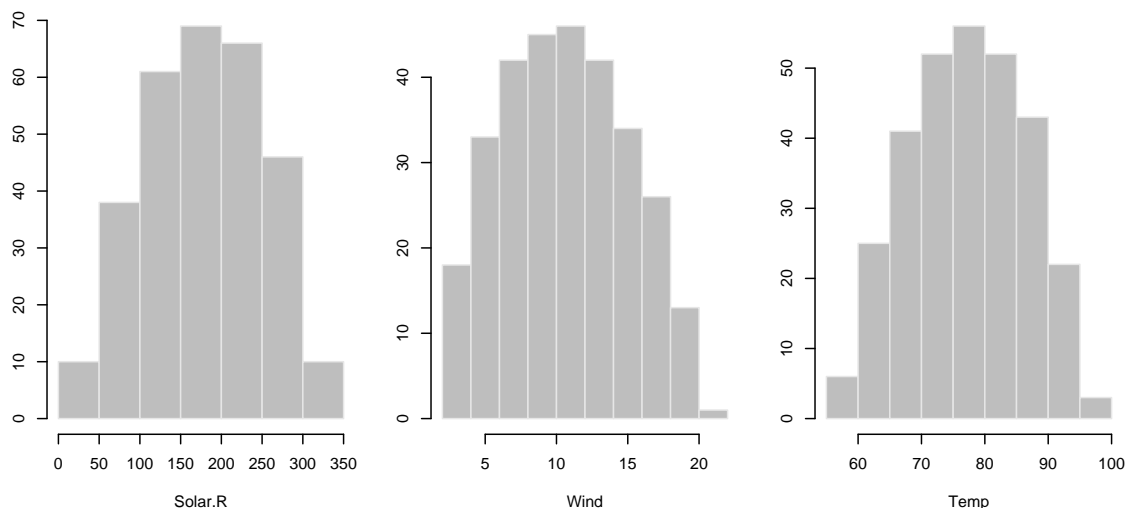


Figure 8: A sample from the marginal uncertainty distribution for the `airquality` data.

Due to missing data (discarded in the current version of `tgp`), we suppress warnings for the sensitivity analysis. We shall use the default `model=btgp`.


```
> s.air <- suppressWarnings(sens(X = X, Z = Z, nn.lhs = 300,
+   rect = rect, shape = shape, mode = mode, verb = 0))
```

Figure (9) shows the results from this analysis.

```
> plot(s.air, layout = "sens")
```

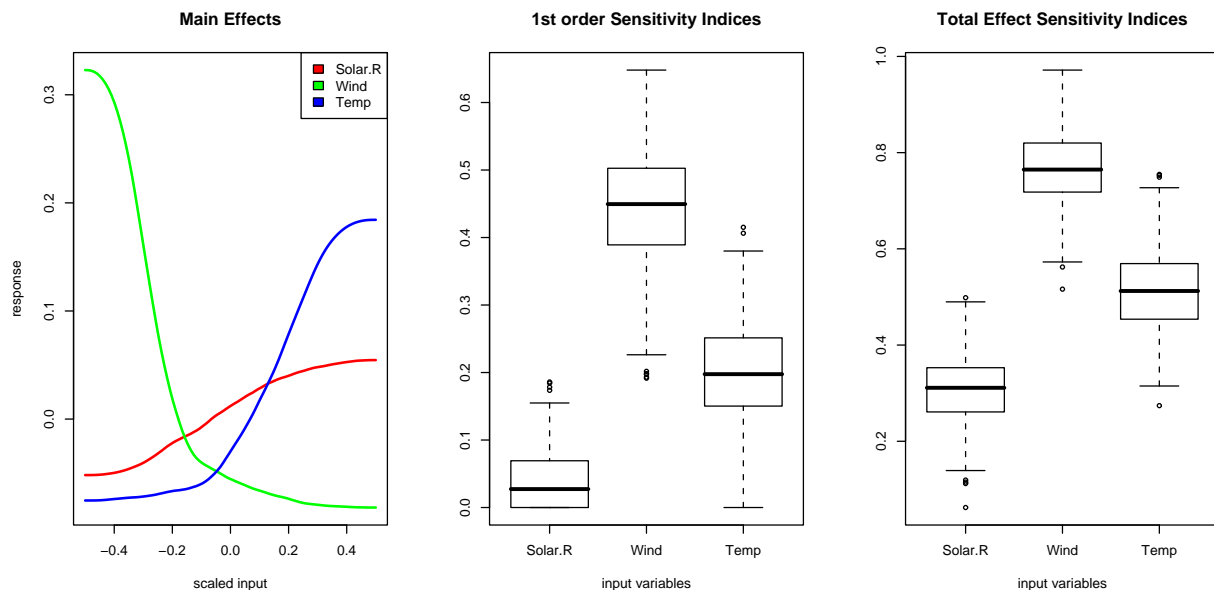


Figure 9: Sensitivity of NYC airquality to natural variation in wind, sun, and temperature.

Through use of `predict.tgp`, it is possible to quickly re-analyze with respect to a new uncertainty distribution without running new MCMC. We can, for example, look at sensitivity for air quality on only low-wind days. We thus alter the uncertainty distribution (assuming that things remain the same for the other variables)

```
> rect[2, ] <- c(0, 5)
> mode[2] <- 2
> shape[2] <- 2
```

and build a set of parameters `sens.p` with the `sens` function by setting `model=NULL`.

```
> sens.p <- suppressWarnings(sens(X = X, Z = Z, nn.lhs = 300,
+   model = NULL, rect = rect, shape = shape, mode = mode))
```

Figures (9) and (10) both show total effect indices which are much larger than the respective first order sensitivities. As one would expect, the effect on airquality is manifest largely through an interaction between variables.

Finally, it is also possible to perform SA with binary covariates, included in the regression model as described in Section 1. In this case, the uncertainty distribution is naturally characterized by a Bernoulli density. Setting `shape[i]=0` informs `sens` that the relevant variable

```

> s.air2 <- predict.tgp(s.air, BTE = c(1, 1000, 1),
+   sens.p = sens.p, verb = 0)
> plot(s.air2, layout = "sens")

```

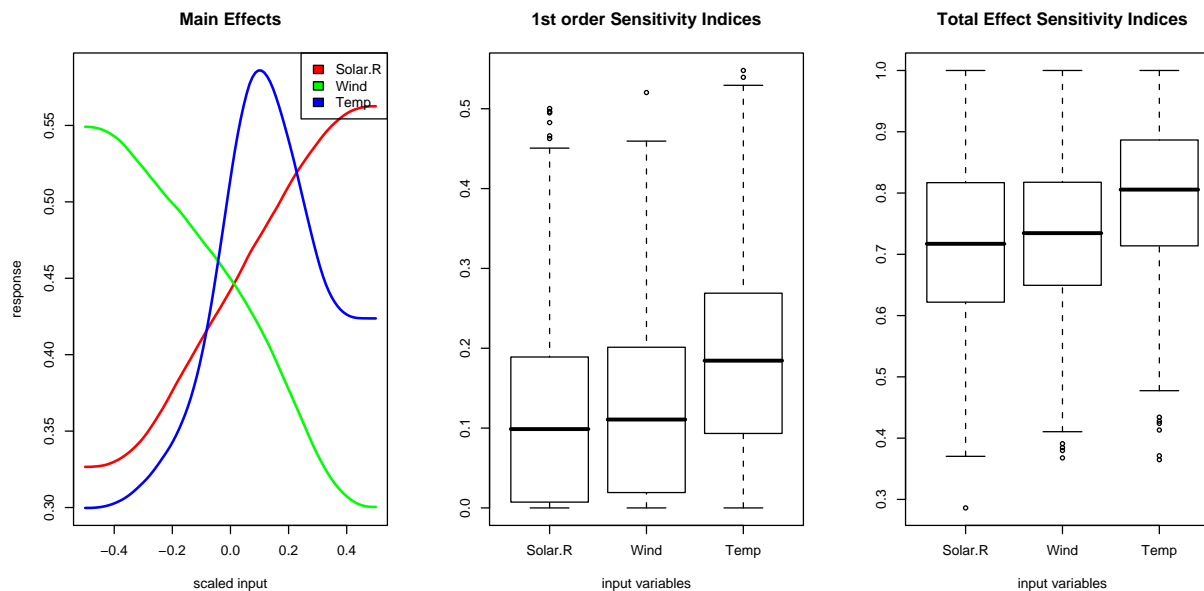


Figure 10: Air quality sensitivity on low-wind days.

is binary (perhaps encoding a categorical input as in Section 1), and that the Bernoulli uncertainty distribution should be used. In this case, the `mode[i]` parameter dictates the probability parameter for the Bernoulli, and we must have `rect[i,] = c(0,1)`. As an example, we re-analyze the original air quality data with temperature included as an indicator variable (set to one if temperature > 79 , the median, and zero otherwise).

```

> X$Temp[X$Temp > 70] <- 1
> X$Temp[X$Temp > 1] <- 0
> rect <- t(apply(X, 2, range, na.rm = TRUE))
> mode <- apply(X, 2, mean, na.rm = TRUE)
> shape <- c(2, 2, 0)
> s.air <- suppressWarnings(sens(X = X, Z = Z, nn.lhs = 300,
+   rect = rect, shape = shape, mode = mode, verb = 0,
+   basemax = 2))

```

Figure (11) shows the results from this analysis.

3 Statistical search for optimization

There has been considerable recent interest in the use of statistically generated search patterns (i.e., locations of relatively likely optima) for optimization. Such search patterns can

```
> plot(s.air, layout = "sens")
```

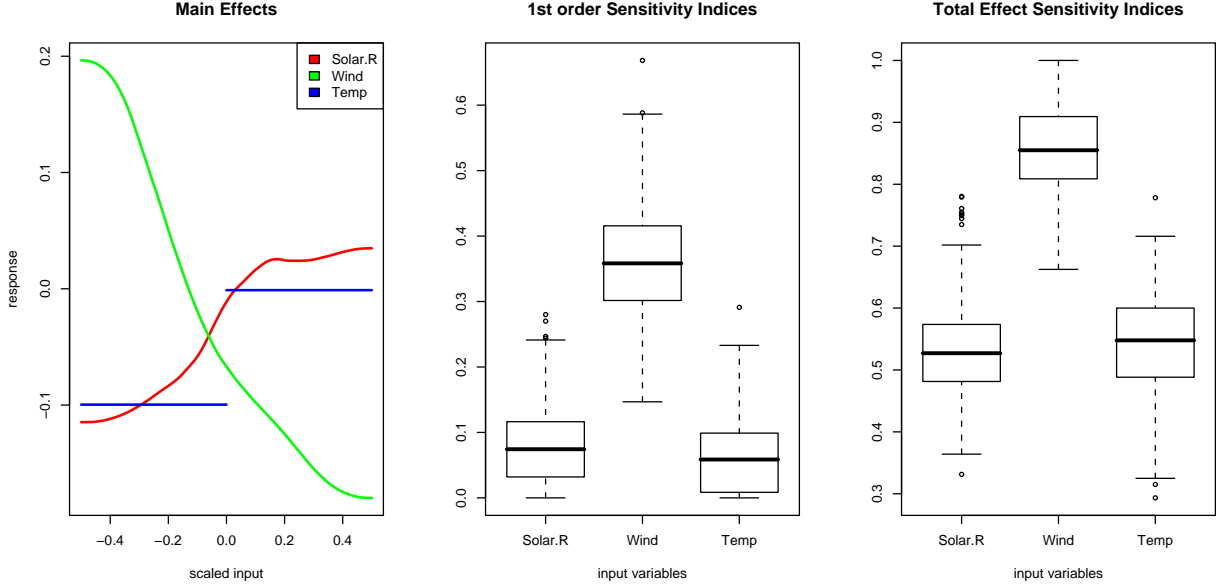


Figure 11: Sensitivity of NYC airquality to natural variation in wind, sun, and a binary temperature variable (for a threshold of 79 degrees).

be used either to direct the optimization completely (e.g., [14] or [24]) or to work in hybrid with a local pattern search optimization (as in [31]). We have implemented within **tgp** the capability to produce a list of input locations that is ordered by a measure of the potential for new optima. The procedure uses samples from the posterior predictive distribution of treed GP regression models to estimate improvement statistics and build an ordered list of search locations which maximize expected improvement. This is the approach first implemented by Taddy, et al., [31], and it provides statistically generated multi-location search patterns that are optimal over the full posterior distribution for improvement.

The single location improvement is defined $I(\mathbf{x}) = \max\{f_{\min} - f(\mathbf{x}), 0\}$, where f_{\min} is the minimum evaluated response in the search (refer to [29] for extensive discussion on general improvement statistics and initial vignette [9] for details of a base implementation in **tgp**). Thus, a high improvement corresponds to an input location that is expected to be much lower than the current minimum. The criterion is easily changed to a search for maximum values through negation of the response. The improvement is always non-negative, as points which do not turn out to be new minimum points still provide valuable information about the output surface. Thus, in the expectation, candidate locations will be rewarded for high response uncertainty (indicating a poorly explored region of the input space) as well as for low mean predicted response. Our **tgp** generated search pattern will consist of m locations that maximize (over a discrete candidate set) the expected multi-location improvement, $\mathbb{E}[I^g(\mathbf{x}_1, \dots, \mathbf{x}_m)]$, where

$$I^g(\mathbf{x}_1, \dots, \mathbf{x}_m) = (\max\{(f_{\min} - f(\mathbf{x}_1)), \dots, (f_{\min} - f(\mathbf{x}_m)), 0\})^g. \quad (5)$$

Increasing g increases the global scope of the criteria by rewarding in the expectation extra variability at \mathbf{x} . For example, $g = 0$ leads to $\mathbb{E}[I^0(\mathbf{x})] = \Pr(I(\mathbf{x}) > 0)$ (assuming the convention $0^0 = 0$), $g = 1$ yields the standard statistic, and $g = 2$ explicitly rewards the improvement variance since $\mathbb{E}[I^2(\mathbf{x})] = \text{var}[I(\mathbf{x})] + \mathbb{E}[I(\mathbf{x})]^2$.

Finding the maximum expectation of (5) is practically impossible for the full posterior distribution of $I^g(\mathbf{x}_1, \dots, \mathbf{x}_m)$, and would require conditioning on a single fit for the parameters of **tgp**. However, **tgp** just seeks to maximize over a discrete list of predictive locations. In fact, the default is to return an ordering for the entire **XX** matrix, thus defining a ranking of predictive locations by order of decreasing expected improvement. There is no restriction on the form for **XX**.¹ The structure of this scheme will dictate the form for **XX**. If it is the case that we seek simply to explore the input space and map a list of potential locations for improvement, using LHS to choose **XX** will suffice.

The discretization of decision space allows for a fast iterative solution to the optimization of $\mathbb{E}[I^g(\mathbf{x}_1, \dots, \mathbf{x}_m)]$. This begins with evaluation of the simple improvement $I^g(\tilde{\mathbf{x}}_i)$ over $\tilde{\mathbf{x}}_i \in \tilde{\mathbf{X}}$ at each of $T = \text{BTE}[2] - \text{BTE}[1]$ MCMC iterations (each corresponding to a single posterior realization of **tgp** parameters and predicted response after burn-in) to obtain the posterior sample

$$\mathcal{I} = \left\{ \begin{array}{ccc} I^g(\tilde{\mathbf{x}}_1)_1 & \dots & I^g(\tilde{\mathbf{x}}_M)_1 \\ & \vdots & \\ I^g(\tilde{\mathbf{x}}_1)_T & \dots & I^g(\tilde{\mathbf{x}}_M)_T \end{array} \right\}.$$

Recall that in **tgp** parlance, and as input to the **b*** functions: $\tilde{\mathbf{X}} \equiv \mathbf{XX}$.

We then proceed iteratively to build an *ordered* collection of m locations according to an iteratively refined improvement: Designate $\mathbf{x}_1 = \text{argmax}_{\tilde{\mathbf{x}} \in \tilde{\mathbf{X}}} \mathbb{E}[I^g(\tilde{\mathbf{x}})]$, and for $j = 2, \dots, m$, given that $\mathbf{x}_1, \dots, \mathbf{x}_{j-1}$ are already included in the collection, the next member is

$$\begin{aligned} \mathbf{x}_j &= \text{argmax}_{\tilde{\mathbf{x}} \in \tilde{\mathbf{X}}} \mathbb{E}[\max\{I^g(\mathbf{x}_1, \dots, \mathbf{x}_{j-1}), I^g(\tilde{\mathbf{x}})\}] \\ &= \text{argmax}_{\tilde{\mathbf{x}} \in \tilde{\mathbf{X}}} \mathbb{E}[(\max\{(f_{\min} - f(\mathbf{x}_1)), \dots, (f_{\min} - f(\mathbf{x}_{j-1})), (f_{\min} - f(\tilde{\mathbf{x}})), 0\})^g] \\ &= \text{argmax}_{\tilde{\mathbf{x}} \in \tilde{\mathbf{X}}} \mathbb{E}[I^g(\mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \tilde{\mathbf{x}})]. \end{aligned}$$

Thus, after each j^{th} additional point is added to the set, we have the maximum expected j -location improvement conditional on the first $j - 1$ locations. This is not necessarily the unconditionally maximal expected j -location improvement; instead, point \mathbf{x}_j is the location which will cause the greatest increase in expected improvement over the given $(j - 1)$ -location expected improvement.

The posterior sample \mathcal{I} acts as a discrete approximation to the true posterior distribution for improvement at locations within the candidate set **XX**. Based upon this approximation, iterative selection of the point set is possible without any re-fitting of the **tgp** model. Conditional on the inclusion of $\tilde{\mathbf{x}}_{i_1}, \dots, \tilde{\mathbf{x}}_{i_{l-1}}$ in the collection, a posterior sample of the l -location

¹A full optimization routine would require that the search pattern is placed within an algorithm iterating towards convergence, as in [31]. However, we concentrate here on the statistical problem of choosing the next samples optimally. We shall touch on issues of convergence in Section 3.2 where we describe a skeleton scheme for optimization extending R's internal **optim** functionality.

improvement statistics is calculated as

$$\mathcal{I}_l = \begin{Bmatrix} I^g(\tilde{\mathbf{x}}_{i_1}, \dots, \tilde{\mathbf{x}}_{i_{l-1}}, \tilde{\mathbf{x}}_1)_1 & \dots & I^g(\tilde{\mathbf{x}}_{i_1}, \dots, \tilde{\mathbf{x}}_{i_{l-1}}, \tilde{\mathbf{x}}_M)_1 \\ \vdots & & \vdots \\ I^g(\tilde{\mathbf{x}}_{i_1}, \dots, \tilde{\mathbf{x}}_{i_{l-1}}, \tilde{\mathbf{x}}_1)_T & \dots & I^g(\tilde{\mathbf{x}}_{i_1}, \dots, \tilde{\mathbf{x}}_{i_{l-1}}, \tilde{\mathbf{x}}_M)_T \end{Bmatrix},$$

where the element in the t^{th} row and j^{th} column of this matrix is calculated as $\max\{I^g(\tilde{\mathbf{x}}_{i_1}, \dots, \tilde{\mathbf{x}}_{i_{l-1}})_t, I^g(\tilde{\mathbf{x}}_j)_t\}$ and the l^{th} location included in the collection corresponds to the column of this matrix with maximum average. Since the multi-location improvement is always at least as high as the improvement at any subset of those locations, the same points will not be chosen twice for inclusion. In practice, very few iterations (about 10% of the total candidate size under the default inference and regression model(s)) through this ordering process can be performed before the iteratively updated improvement statistics become essentially zero. Increasing the number of MCMC iterations (BTE[2]–BTE[1]) can mitigate this to a large extent.²

3.1 A simple example

We shall use the Rosenbrock function to illustrate the production of an ordered collection of (possible) adaptive samples to maximize the expected improvement within **tg**p. Specifically, the two dimensional Rosenbrock function is defined as

```
> rosenbrock <- function(x) {
+   x <- matrix(x, ncol = 2)
+   100 * (x[, 1]^2 - x[, 2])^2 + (x[, 1] - 1)^2
+ }
```

and we shall bound the search space for adaptive samples to the rectangle: $-1 \leq x_i \leq 5$ for $i = 1, 2$. The single global minimum of the Rosenbrock function is at (1, 1).

```
> rosenbrock(c(1, 1))
```

```
[1] 0
```

This function involves a long steep valley with a gradually sloping floor, and is considered to be a difficult problem for local optimization routines.

We begin by drawing an LHS of 40 input locations within the bounding rectangle, and evaluating the function at these locations.

```
> rect <- cbind(c(-1, -1), c(5, 5))
> X <- lhs(40, rect)
> Z <- rosenbrock(X)
```

²Once a zero (maximal) iterative improvement is attained the rest of the ranking is essentially arbitrary, at which point **tg**p cuts off the process prematurely.

We will fit a `bgp` model to this data to predict the Rosenbrock response at unobserved (candidate) input locations in `XX`. The `improv` argument may be used to obtain an ordered list of places where we should be looking for new minima. In particular, specifying `improv=c(1,10)` will return the 10 locations which maximize the iterative multi-location expected improvement function, with $g = 1$ (i.e., Eq. (5)). Note that `improv=TRUE` is also possible, in which case `g` defaults to one and the entire list of locations is ranked. Our candidate set is just a space filling LHS design. In other situations, it may be useful to build an informative LHS design (i.e., to specify `shape` and `mode` arguments for the `lhs` function) to reflect what is already known about the location of optima.

```
> XX <- lhs(200, rect)
> rfit <- bgp(X, Z, XX, improv = c(1, 10), verb = 0)
```

Upon return, the `"tgp"`-class object `rfit` includes the matrix `improv`, which is a list of the expected single location improvement for the 200 `XX` locations, and the top 10 ranks. Note that the `ranks` for those points which are not included in the top 10 are set to `nrow(XX)=200`. Here are the top 10:

```
> cbind(rfit$improv, XX)[rfit$improv$rank <= 10, ]
```

	improv	rank	1	2
4	0.0006526866	2	1.4124348	2.15174711
20	0.0006289085	10	1.3131381	1.87566387
21	0.0006389502	6	1.6059344	2.38150541
40	0.0005311757	8	2.0701606	4.59635226
103	0.0005959594	9	0.5403690	0.01969345
110	0.0006586264	1	0.5995391	0.26334479
119	0.0006328093	3	1.6171801	2.39090831
128	0.0005580968	5	1.9180056	3.34226608
133	0.0006162374	4	-0.7459543	0.53337549
200	0.0006240170	7	2.0423112	3.99968412

This iterative algorithm may produce ranks that differ significantly from a straightforward ordering of expected improvement. This leads to a list that better explores the input space, since the expected improvement is naturally balanced against a desire to search the domain.

We plot the results with the usual function, by setting `as="improv"`, in Figure 12. The banana-shaped region of higher expected improvement corresponds to the true valley floor for the Rosenbrock function, indicating that the `bgp` model is doing a good job of prediction. Also, we note that the ordered input points are well dispersed throughout the valley—a very desirable property for adaptive sampling candidates.

It is straightforward, with `predict.tgp`, to obtain a new ordering for the more global $g=5$ (or any new g). Figure 13 shows a more diffuse expected improvement surface and a substantially different point ordering. In practice, we have found that $g=2$ provides a good compromise between local and global search.

```
> plot(rfit, as = "improv")
```

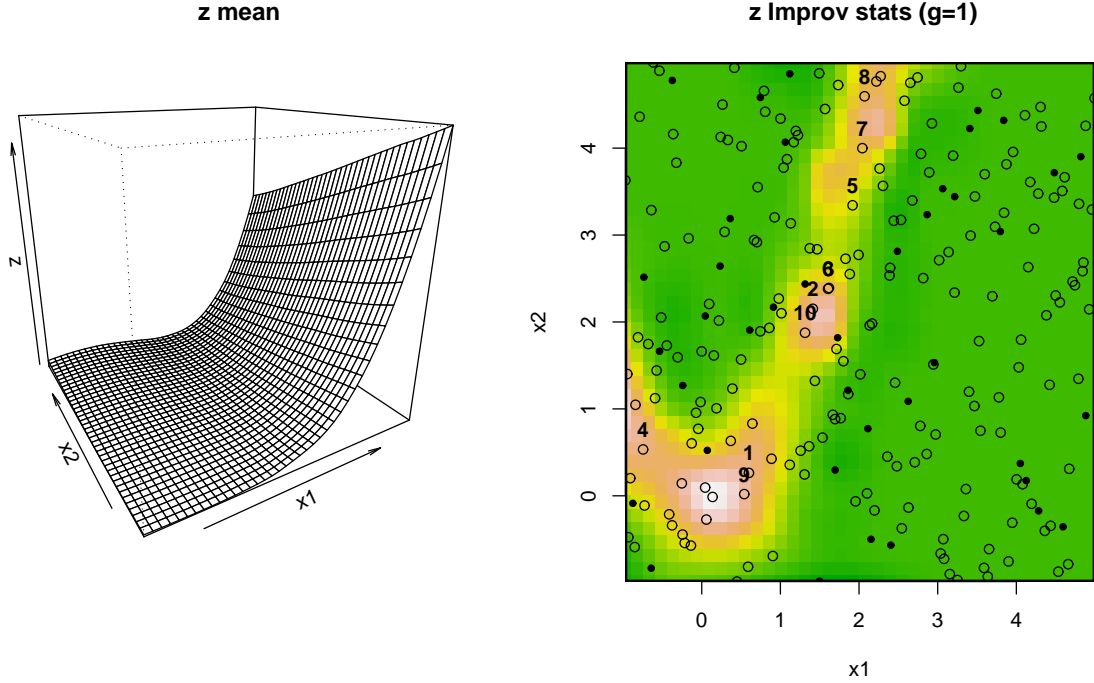


Figure 12: The *left* panel shows the mean predicted Rosenbrock function response, and on the *right* we have expected single location improvement with the top 20 points (labelled by rank) plotted on top.

3.2 A skeleton optimization scheme

The capabilities outlined above are useful in their own right, as a search list or candidate set ranked by expected improvement gain provides concrete information about potential optima. However, a full optimization framework requires that the production of these sets of search locations are nested within an iterative search scheme. The approach taken by Taddy, et al. [31], achieves this by taking the **tgp** generated sets of locations and using them to augment a local optimization search algorithm. In this way, the authors are able to achieve robust solutions which balance the convergence properties of the local methods with the global scope provided by **tgp**. Indeed, any optimization routine capable of evaluating points provided by an outside source could benefit from a **tgp** generated list of search locations.

In the absence of this sort of formal hybrid search algorithm, it is still possible to devise robust optimization algorithms based around **tgp**. A basic algorithm is as follows: first, use a LHS to explore the input space (see the **lhs** function included in **tgp**). Repeatedly fit one of the **b*** models with **improv!=FALSE** to the evaluated iterates to produce a search set, then evaluate the objective function over this search set, as described earlier. Then evaluate the objective function over the highest ranked locations in the search set. Continue until you are confident that the search has narrowed to a neighborhood around the true optimum (a good indicator of this is when all of the top-ranked points are in the same area). At this

```

> rfit2 <- predict.tgp(rfit, XX = XX, BTE = c(1, 1000,
+      1), improv = c(5, 20), verb = 0)
> plot(rfit2, layout = "as", as = "improv")

```

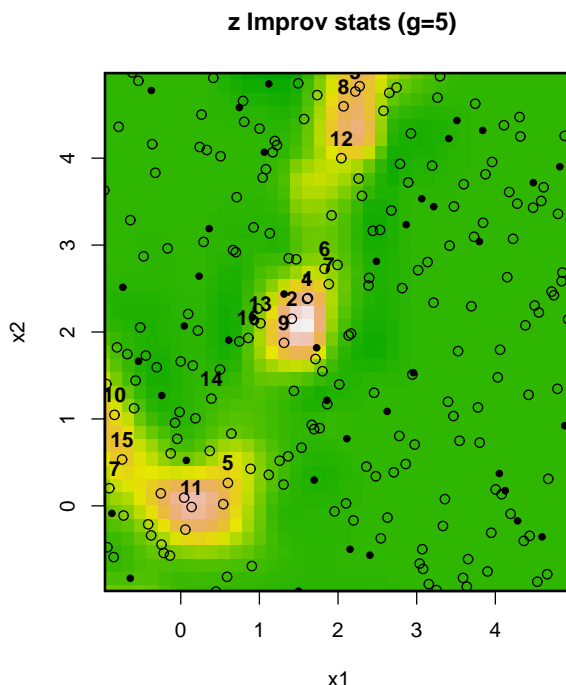


Figure 13: The expected improvement surface and top 20 ordered locations, for $g=5$.

point, the optimization may be completed by `optim`, R's general purpose local optimization algorithm in order to guarantee convergence. The `optim` routine may be initialized to the best input location (i.e. corresponding the most optimal function evaluation) found thus far by `tgp`.

Note that this approach is actually an extreme version of a template proposed by Taddy, et al. [31], where the influence of global (i.e. `tgp`) search is downweighted over time rather than cut off. In either case, a drawback to such approaches is that they do not apply when the function being optimized is deterministic. An alternative scheme is to employ both `tgp` search and a local optimization at each iteration. The idea is that a mix of local and global information is provided throughout the entire optimization, but with an added twist. Rather than apply `optim` on the stochastic function directly, which would not converge due to the noise, it can be applied on a deterministic (MAP) kriging surface provided by `tgp`. The the local optima obtained can be used to augment the candidate set of locations where the improvement statistic is gathered—which would otherwise be simple LHS. That way the search pattern produced on output is likely to have a candidate with high improvement.

To fix ideas, and for the sake of demonstration, the `tgp` package includes a skeleton function for performing a single iteration in the derivative-free optimization of noisy black-box functions. The function is called `optim.step.tgp`, and the name is intended to emphasize

that it performs a single step in an optimization by trading off local `optim`-based search of `tgp` predictive (kriging surrogate) surfaces, with the expected posterior improvement. In other words, it is loosely based on some the techniques alluded to above, but is designed to be augmented/adjusted as needed. Given N pairs of inputs and responses (\mathbf{X}, \mathbf{Z}) , `optim.step.tgp` suggests new points at which the function being optimized should be evaluated. It also returns information that can be used to assess convergence. An outline follows.

The `optim.step.tgp` function begins by constructing a set of candidate locations, either as a space filling LHS over the input space (the default) or from a treed D -optimal design, based on a previously obtained "`tgp`"-class model. R's `optim` command is used on the MAP predictive surface contained within the object to obtain an estimate of the current best guess x -location of the optimal solution. A standalone subroutine called `optim.ptgpf` is provided for this specific task, to be used within `optim.step.tgp` or otherwise. Within `optim.step.tgp`, `optim.ptgpf` is initialized with the data location currently predicted to be the best guess of the minimum. The optimal x -location found is then added into the set of candidates as it is likely that the expected improvement would be high there.

Then, a new "`tgp`"-class object is obtained by applying a `b*` function to (\mathbf{X}, \mathbf{Z}) whilst sampling from the posterior distribution of the improvement statistic. The best one, two, or several locations with highest improvement ranks are suggested for addition into the design. The values of the maximum improvement statistic are also returned in order to track progress in future iterations. The "`tgp`"-class object returned is used to construct candidates and initialize the `optim.ptgpf` function in future rounds.

To illustrate, consider the 2-d exponential data from the initial vignette [9] as our noisy function f .

```
> f <- function(x) {
+   exp2d.Z(x)$Z
+ }
```

Recall that this data is characterized by a mean value of

$$f(\mathbf{x}) = x_1 \exp(-x_1^2 - x_2^2)$$

which is observed with a small amount of Gaussian noise (with `sd` = 0.001). Elementary calculus gives that the minimum of f is obtained at $\mathbf{x} = (-\sqrt{1/2}, 0)$.

The `optim.step.tgp` function requires that the search domain be defined by a bounding rectangle, and we require an initial design to start things off. Here we shall use $[-2, 6]^2$ with an LHS design therein.

```
> rect <- rbind(c(-2, 6), c(-2, 6))
> X <- lhs(20, rect)
> Z <- f(X)
```

The following code proceeds with several rounds of sequential design towards finding the minimum of `f`.

```

> out <- progress <- NULL
> for(i in 1:20) {
+
+   ## get recommendations for the next point to sample
+   out <- optim.step.tgp(f, X=X, Z=Z, rect=rect, prev=out, verb=0)
+
+   ## add in the inputs, and newly sampled outputs
+   X <- rbind(X, out$X)
+   Z <- c(Z, f(out$X))
+
+   ## keep track of progress and best optimum
+   progress <- rbind(progress, out$progress)
+ }

```

The `progress` can be tracked through the rows of a `data.frame`, as constructed above, containing a listing of the input location of the current best guess of the minimum for each round, together with the value of the objective at that point, as well as the maximum of the improvement statistic. In addition to printing this data to the screen, plots such as the ones in Figure 14 can be valuable for assessing convergence. As can be seen in the figure, the final iteration gives an `x`-value that is very close to the correct result, and is (in some loose sense) close to convergence.

```

> out$progress[1:2]

      x1      x2
1 -0.7085659 0.00571415

```

As mentioned above, if it is known that the function evaluations are deterministic then, at any time, R's `optim` routine can be invoked—perhaps initialized by the the `x`-location in `out$progress`—and convergence to a local optimum thus guaranteed. Otherwise, `out$progress` will converge, in some sense, as long as the number of MCMC rounds used in each round, above, ($T = \text{BTE}[2] - \text{BTE}[1]$) tends to infinity. Such arguments to the `b*` functions can be set via the ellipses (`...`) arguments to `optim.step.tgp`.³ A heuristic stopping criterion can be based on the maximum improvement statistic obtained in each round as long as the candidate locations become dense in the region as $T \rightarrow \infty$. This can be adjusted by increasing the `NN` argument to `optim.step.tgp`.

The internal use of `optim` within `optim.step.tgp` on the posterior predictive (kriging surrogate) surface via `optim.ptgpf` may proceed with any of the usual method arguments. I.e.,

```

> formals(optim)$method

```

³This runs contrary to how the ellipses are used by `optim` in order to specify static arguments to `f`. If setting static arguments to `f` is required within `optim.step.tgp`, then they must be set in advance by adjusting the default arguments via `formals`.

```

> par(mfrow = c(1, 2))
> matplot(progress[, 1:2], main = "x progress", xlab = "rounds",
+         ylab = "x[,1:2]", type = "l", lwd = 2)
> plot(log(progress$improv), type = "l", main = "max log improv",
+      xlab = "rounds", ylab = "max log(improv)")

```

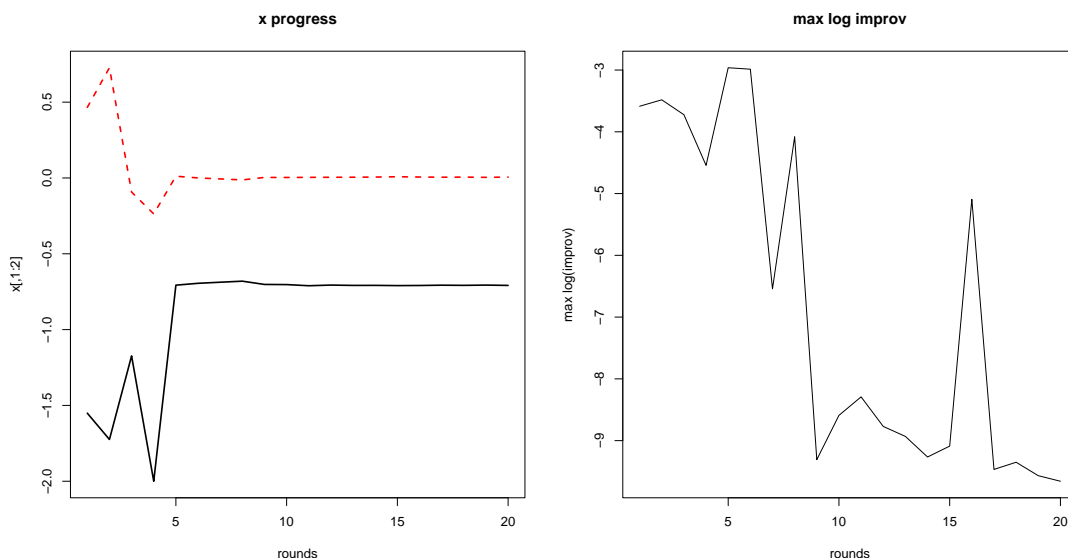


Figure 14: Progress in iterations of `optim.step.tgp` shown by tracking the \mathbf{x} -locations of the best guess of the minimum (*left*) and the logarithm of the maximum of the improvement statistics at the candidate locations (*right*)

```

c("Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN")

```

however the default ordering is switched in `optim.ptgpf` and includes one extra method.

```

> formals(optim.ptgpf)$method

```

```

c("L-BFGS-B", "Nelder-Mead", "BFGS", "CG", "SANN", "optimize")

```

Placing "L-BFGS-B" in the default position is sensible since this method enforces a rectangle of constraints as specified by `rect`. This guarantees that the additional candidate found by `optim.ptgpf` will be valid. However, the other `optim` methods generally work well despite that they do not enforce this constraint. The final method, "optimize", applies only when the inputs to `f` are 1-d. In this case, the documentation for `optim` suggests using the `optimize` function instead.

4 Importance tempering

Importance tempering (IT) is a combination of the classic method of importance sampling (IS) [16] and *simulated tempering* (ST) [8]. The ST algorithm helps obtain samples from a

multimodal density $\pi(\theta)$ where standard methods, such as Metropolis–Hastings (MH) [17, 11] and Gibbs Sampling (GS) [6], fail. It is well-known that MCMC inference in Bayesian treed methods suffers from poor mixing. For example, Chipman et al. [2, 3] recommend periodically restarting the MCMC to avoid chains becoming stuck in local modes of the posterior distribution (particularly in tree space). The `R` argument to the `b*` functions exists precisely to facilitate such restarts within the `tgpr` package.

As will be shown in our examples, ST can guard against becoming stuck in local modes of the `tgpr` posterior by encouraging better mixing *between modes* via an increase in the acceptance rate of tree modification proposals, particularly *prunes*. However, as we will see, ST suffers from inefficiency because it discards the lions share of the samples it collects. The discarded samples can be recycled if they are given appropriate importance sampling (IS) weights. These weights, if combined carefully, can be used to construct meta-estimators of expectations under the `tgpr` posterior that have much lower variance than under ST alone.

4.1 Simulated Tempering and related methods

ST is an application of the MH algorithm on the product space of parameters and inverse temperatures $k \in [0, 1]$. That is, ST uses MH to sample from the joint chain $\pi(\theta, k) \propto \pi(\theta)^k p(k)$. The inverse temperature is allowed to take on a discrete set of values $k \in \{k_1, \dots, k_m : k_1 = 1, k_i > k_{i+1} \geq 0\}$, called the *temperature ladder*. Typically, ST calls for sampling $(\theta, k)^{(t+1)}$ by first updating $\theta^{(t+1)}$ conditional on $k^{(t)}$ and (possibly) on $\theta^{(t)}$, using MH or GS. Then, for a proposed $k' \sim q(k^{(t)} \rightarrow k')$, usually giving equal probability to the nearest inverse temperatures greater and less than $k^{(t)}$, an acceptance ratio is calculated:

$$A^{(t+1)} = \frac{\pi(\theta^{(t+1)})^{k'} p(k') q(k' \rightarrow k^{(t)})}{\pi(\theta^{(t+1)})^{k^{(t)}} p(k^{(t)}) q(k^{(t)} \rightarrow k')}.$$

Finally, $k^{(t+1)}$ is determined according to the MH accept/reject rule: set $k^{(t+1)} = k'$ with probability $\alpha^{(t+1)} = \min\{1, A^{(t+1)}\}$, or $k^{(t+1)} = k^{(t)}$ otherwise. Standard theory for MH and GS gives that samples from the marginals π_{k_i} can be obtained by collecting samples $\theta^{(t)}$ where $k^{(t)} = k_i$. Samples from $\pi(\theta)$ are obtained when $k^{(t)} = 1$.

The success of ST depends crucially on the ability of the Markov chain frequently to: (a) visit high temperatures (low k) where the probability of escaping local modes is increased; (b) visit $k = 1$ to obtain samples from π . The algorithm can be tuned by: (i.) adjusting the number and location of the rungs of the temperature ladder; or (ii.) setting the pseudo-prior $p(k)$ for inverse temperature.

Geyer & Thompson [8] give ways of adjusting the spacing of the rungs of the ladder so that the ST algorithm achieves between-temperature acceptance rates of 20–40%. More recently, authors have preferred to rely on defaults, e.g.,

$$k_i = \begin{cases} (1 + \Delta_k)^{1-i} & \text{geometric spacing} \\ \{1 + \Delta_k(i-1)\}^{-1} & \text{harmonic spacing} \end{cases} \quad i = 1, \dots, m. \quad (6)$$

Motivation for such default spacings is outlined by Liu [16]. Geometric spacing, or uniform spacing of $\log(k_i)$, is also advocated by Neal [19, 20] to encourage the Markov chain to

rapidly traverse the breadth of the temperature ladder. Harmonic spacing is more often used by a related method called Metropolis coupled Markov chain Monte Carlo (MC³) [7]. Both defaults are implemented in the `tgp` package, through the provided `default.itemps` function. A new “sigmoidal” option is also implemented, as discussed below. The rate parameter $\Delta_k > 0$ can be problem specific. Rather than work with Δ_k the `default.itemps` function allows the ladder to be specified via m and the hottest temperature k_m , thus fixing Δ_k implicitly. I.e., for the geometric ladder $\Delta_k = (k_m)^{1/(1-m)} - 1$, and for the harmonic ladder $\Delta_k = \frac{(k_m)^{-1}-1}{m-1}$.

A sigmoidal ladder can provide a higher concentration of temperatures near $k = 1$ without sacrificing the other nice properties of the geometric and harmonic ladders. It is specified by first situating m indices $j_i \in \mathbb{R}$ so that $k_1 = k(j_1) = 1$ and $k_m = k(j_m) = k_m$ under

$$k(j_i) = 1.01 - \frac{1}{1 + e^{j_i}}.$$

The remaining $j_i, i = 2, \dots, (m-1)$ are spaced evenly between j_1 and j_m to fill out the ladder $k_i = k(j_i), i = 1, \dots, (m-1)$.

By way of comparison, consider generating the three different types of ladder with identical minimum inverse temperature $k_m = 0.1$, the default setting in `tgp`.

```
> geo <- default.itemps(type = "geometric")
> har <- default.itemps(type = "harmonic")
> sig <- default.itemps(type = "sigmoidal")
```

The plots in Figure 15 show the resulting inverse temperature ladders, and their logarithms. Observe how, relative to the geometric ladder, the harmonic ladder has a higher concentration of inverse temperatures near zero, whereas the sigmoidal ladder has a higher concentration near one.

Once a suitable ladder has been chosen, the `tgp` package implementation of ST follows the suggestions of Geyer & Thompson [8] in setting the pseudo-prior, starting from a uniform p_0 . First, p_0 is adjusted by *stochastic approximation*: add $c_0/[m(t+n_0)]$ to $\log p_0(k)$ for each $k_i \neq k^{(t)}$ and subtract $c_0/(t+n_0)$ from $\log p_0(k^{(t)})$ over $t = 1, \dots, B$ *burn-in* MCMC rounds sampling from the joint posterior of (θ, k) . Then, p_0 is normalized to obtain p_1 . Before subsequent runs, specified via an R `>= 2` argument, *occupation numbers* $o(k_i) = \sum_{t=1}^B 1_{\{k^{(t)}=k_i\}}$, are used update $p(k_i) \propto p_1(k_i)/o(k_i)$. Note that, in this setting, the R argument is used to update the pseudo-prior only, not to restart the Markov chain.

4.2 Importance sampling from tempered distributions

ST provides us with $\{(\theta^{(t)}, k^{(t)}) : t = 1, \dots, T\}$, where $\theta^{(t)}$ is an observation from $\pi_{k^{(t)}}$. It is convenient to write $\mathcal{T}_i = \{t : k^{(t)} = k_i\}$ for the index set of observations at the i^{th} temperature, and let $T_i = |\mathcal{T}_i|$. Let the vector of observations at the i^{th} temperature collect in $\boldsymbol{\theta}_i = (\theta_{i1}, \dots, \theta_{iT_i})$, so that $\{\theta_{ij}\}_{j=1}^{T_i} \sim \pi_{k_i}$. Each vector $\boldsymbol{\theta}_i$ can be used to construct an IS

```

> par(mfrow = c(2, 1))
> all <- cbind(geo$k, har$k, sig$k)
> matplot(all, pch = 21:23, main = "inv-temp ladders",
+         xlab = "indx", ylab = "itemp")
> legend("topright", pch = 21:23, c("geometric", "harmonic",
+         "sigmoidal"), col = 1:3)
> matplot(log(all), pch = 21:23, main = "log(inv-temp) ladders",
+         xlab = "indx", ylab = "itemp")

```

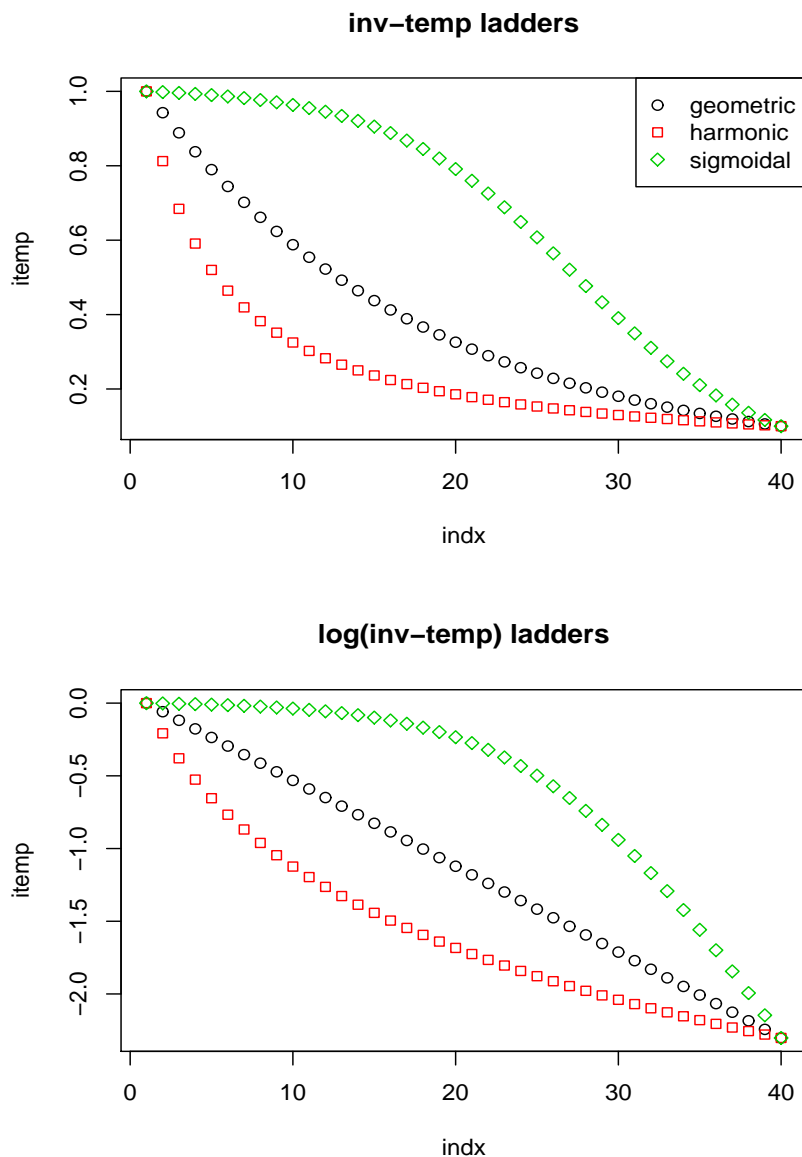


Figure 15: Three different inverse temperature ladders, each with $m = 40$ temperatures starting at $k_1 = 1$ and ending at $k_m = 0.1$

estimator of $E_\pi\{h(\theta)\}$ by setting

$$\hat{h}_i = \frac{\sum_{j=1}^{T_i} w_i(\theta_{ij})h(\theta_{ij})}{\sum_{j=1}^{T_i} w_i(\theta_{ij})} \equiv \frac{\sum_{j=1}^{T_i} w_{ij}h(\theta_{ij})}{W_i},$$

say. That is, rather than obtain one estimator from ST (at the cold temperature), we can obtain m estimators (one at each temperature) vi IS. The efficiency of each estimator, $i = 1, \dots, m$ can be measured through it's variance, but unfortunately this can be difficult to calculate in general. As a result, the notion of *effective sample size* [16] plays an important role in the study of IS estimators. Denote the vector of IS weights at the i^{th} temperature as $\mathbf{w}_i = \mathbf{w}_i(\boldsymbol{\theta}_i) = (w_i(\theta_{i1}), \dots, w_i(\theta_{iT_i}))$, where $w_i(\theta) = \pi(\theta)/\pi_{k_i}(\theta)$. The ESS of \hat{h}_i is defined by

$$\text{ESS}(\mathbf{w}_i) = \frac{T}{1 + \text{cv}^2(\mathbf{w}_i)}, \quad (7)$$

where $\text{cv}(\mathbf{w}_i)$ is the *coefficient of variation* of the weights (in the i^{th} temperature), given by

$$\text{cv}^2(\mathbf{w}_i) = \frac{\sum_{t=1}^T (w(\theta^{(t)}) - \bar{w})^2}{(T-1)\bar{w}^2}, \quad \text{where} \quad \bar{w} = T^{-1} \sum_{t=1}^T w(\theta^{(t)}).$$

In R:

```
> ESS <- function(w) {
+   mw <- mean(w)
+   cv2 <- sum((w - mw)^2)/((length(w) - 1) * mw^2)
+   ess <- length(w)/(1 + cv2)
+   return(ess)
+ }
```

This should not be confused with the concept of *effective sample size due to autocorrelation* [15] (due to serially correlated samples coming from a Markov chain as in MCMC) as implemented by the `effectiveSize` function in the `coda` package [23] for R.

Before attempting to combine m IS estimators it is fruitful backtrack briefly to obtain some perspective on the topic of applying IS with a *single* tempered proposal distribution. Jennison [13] put this idea forward more than a decade ago, though the question of how to choose the best temperature was neither posed or resolved. It is clear that larger k leads to lower variance estimators (and larger ESS), but at the expense of poorer mixing in the Markov chain. It can be shown that the optimal inverse temperature k^* for IS, in the sense of constructing a minimum variance estimator, may be significantly lower than one [10]. However, the variance of such an estimator will indeed become unbounded as $k \rightarrow 0$, just as $\text{ESS} \rightarrow 0$. Needless to say, the choice of how to best pick the best temperatures (for ST or IS) is still an open problem. But in the context of the family of tempered distributions used by ST for mixing considerations, this means that the discarded samples obtained when $k^{(t)} < 1$ may actually lead to more efficient estimators than the ones saved from the cold distribution. So ST is wasteful indeed.

However, when combining IS estimators from the multiple temperatures used in ST, the deleterious effect of the high variance ones obtained at high temperature must be mitigated. The possible strategies involved developing such a meta-estimator comprise the *importance tempering* (IT) family of methods. The idea is that small ESS will indicate high variance IS estimators which should be relegated to having only a small influence on the overall estimator.

4.3 A optimal way to combine IS estimators

It is natural to consider an overall meta-estimator of $E_\pi\{h(\theta)\}$ defined by a convex combination:

$$\hat{h}_\lambda = \sum_{i=1}^m \lambda_i \hat{h}_i, \quad \text{where} \quad 0 \leq \lambda_i \leq \sum_{i=1}^m \lambda_i = 1. \quad (8)$$

Unfortunately, if $\lambda_1, \dots, \lambda_m$ are not chosen carefully, $\text{Var}(\hat{h}_\lambda)$, can be nearly as large as the largest $\text{Var}(\hat{h}_i)$ [22], due to the considerations alluded to in Section 4.2. Notice that ST is recovered as a special case when $\lambda_1 = 1$ and $\lambda_2, \dots, \lambda_m = 0$. It may be tempting to choose $\lambda_i = W_i/W$, where $W = \sum_{i=1}^m W_i$. The resulting estimator is equivalent to

$$\hat{h} = W^{-1} \sum_{t=1}^T w(\theta^{(t)}, k^{(t)}) h(\theta^{(t)}), \quad \text{where} \quad W = \sum_{t=1}^T w(\theta^{(t)}, k^{(t)}), \quad (9)$$

and $w(\theta, k) = \pi(\theta)/\pi(\theta)^k = \pi(\theta)^{1-k}$. It can lead to a very poor estimator, even compared to ST, as will be demonstrated empirically in the examples to follow shortly.

Observe that we can equivalently write

$$\hat{h}_\lambda = \sum_{i=1}^m \sum_{j=1}^{T_i} w_{ij}^\lambda h(\theta_{ij}), \quad \text{where} \quad w_{ij}^\lambda = \lambda_i w_{ij}/W_i. \quad (10)$$

Let $\mathbf{w}^\lambda = (w_{11}^\lambda, \dots, w_{1T_1}^\lambda, w_{21}^\lambda, \dots, w_{2T_2}^\lambda, \dots, w_{m1}^\lambda, \dots, w_{mT_m}^\lambda)$. Attempting to choose $\lambda_1, \dots, \lambda_m$ to minimize $\text{Var}(\hat{h}_\lambda)$ directly can be difficult. Moreover, for the applications that we have in mind, it is important that our estimator can be constructed without knowledge of the normalizing constants of $\pi_{k_1}, \dots, \pi_{k_m}$, and without evaluating the MH transition kernels $\mathcal{K}_{\pi_{k_i}}(\cdot, \cdot)$. It is for this reason that methods like the *balance heuristic* [32], MCV [22], or population Monte Carlo (PMC) [4] cannot be applied. Instead, we seek maximize the effective sample size of \hat{h}_λ in (8), and look for an $O(T)$ operation to determine the optimal λ^* .

Among estimators of the form (8), it can be shown [10] that $\text{ESS}(\mathbf{w}^\lambda)$ is maximized by $\lambda = \lambda^*$, where, for $i = 1, \dots, m$,

$$\lambda_i^* = \frac{\ell_i}{\sum_{i=1}^m \ell_i}, \quad \text{and} \quad \ell_i = \frac{W_i^2}{\sum_{j=1}^{T_i} w_{ij}^2}.$$

The efficiency of each IS estimator \hat{h}_i can be measured through $\text{ESS}(\mathbf{w}_i)$. Intuitively, we hope that with a good choice of λ , the ESS (7) of \hat{h}_λ , would be close to the sum over i of the effective sample sizes each of \hat{h}_i . This is indeed the case for \hat{h}_{λ^*} , because it can be shown [10] that

$$\text{ESS}(\mathbf{w}^{\lambda^*}) \geq \sum_{i=1}^m \text{ESS}(\mathbf{w}_i) - \frac{1}{4} - \frac{1}{T}.$$

In practice we have found that this bound is conservative and that in fact $\text{ESS}(\mathbf{w}^{\lambda^*}) \geq \sum_{i=1}^m \text{ESS}(\mathbf{w}_i)$, as will be shown empirically in the examples that follow. Thus our optimally-combined IS estimator has a highly desirable and intuitive property in terms of its effective sample size: that the whole is greater than the sum of its parts.

$\text{ESS}(\mathbf{w}^{\lambda^*})$ depends on $\text{ESS}(\mathbf{w}_i)$ which in turn depend on the k_i . Smaller k_i will lead to better mixing in the Markov chain, but lower $\text{ESS}(\mathbf{w}_i)$. Therefore, we can expect that the geometric and sigmoidal ladders will fare better than the harmonic ones, so long as the desired improvements in mixing are achieved. In the examples to follow, we shall see that the sigmoidal ladder does indeed leader to higher $\text{ESS}(\mathbf{w}^{\lambda^*})$.

4.4 Examples

Here the IT method is shown in action for `tgp` models. IT is controlled in `b*` functions via the `items` argument: a `data.frame` coinciding with the output of the `default.items` function. The `lambda` argument to `default.items` can be used to base posterior predictive inference the other IT heuristics: ST and the naïve approach (9). Whenever the argument `m = 1` is used with `k.min != 1` the resulting estimator is constructed via tempered importance sampling at the single inverse temperature `k.min`, in the style of Jennison [13] as outlined in Section 4.2. The parameters c_0 and n_0 for stochastic approximation of the pseudo-prior can be specified as a 2-vector `c0n0` argument to `default.items`. In the examples which follow we simply use the default configuration of the IT method, adjusting only the minimum inverse temperature via the `k.min` argument.

Before delving into more involved examples, we illustrate the stages involved in a small run of importance tempering (IT) on the exponential data from Section 3.3 of [9]. The data can be obtained as:

```
> exp2d.data <- exp2d.rand()
> X <- exp2d.data$X
> Z <- exp2d.data$Z
```

Now, consider applying IT to the Bayesian treed LM with a small geometric ladder. A warning will be given if the default setting of `bprior="bflat"` is used, as this (numerically) improper prior can lead to improper posterior inference at high temperatures.

```
> its <- default.items(m = 10)
> exp.btlm <- btlm(X = X, Z = Z, bprior = "b0", R = 2,
+   items = its, pred.n = FALSE, BTE = c(1000, 3000,
+   2))
```

```

burn in: [with stoch approx (c0,n0)=(100,1000)]
**GROW** @depth 0: [1,0.45], n=(58,22)
**GROW** @depth 1: [2,0.45], n=(45,13)
**GROW** @depth 2: [1,0.2], n=(23,22)
**PRUNE** @depth 2: [1,0.25]
**GROW** @depth 1: [1,0.3], n=(12,11)
**CPRUNE** @depth 0: var=2, val=0.45->0.5, n=(60,20)
r=1000 d=[0] [0] [0]; n=(51,11,18) k=0.129155
r=2000 d=[0] [0] [0]; n=(48,14,18) k=0.359381
**PRUNE** @depth 1: [2,0.5]
**GROW** @depth 1: [2,0.45], n=(48,14)
r=3000 d=[0] [0] [0]; n=(50,13,17) k=0.278256

```

```

Sampling @ nn=0 pred locs:
**PRUNE** @depth 1: [2,0.45]
r=1000 d=[0] [0]; mh=3 n=(62,18) k=1
**GROW** @depth 1: [2,0.45], n=(48,14)
r=2000 d=[0] [0] [0]; mh=3 n=(51,11,18) k=0.278256
Grow: 2.62%, Prune: 1.19%, Change: 36.16%, Swap: 21%
finished repetition 1 of 2

```

```

burn in:
r=1000 d=[0] [0] [0]; mh=3 n=(51,11,18) k=0.16681

```

```

Sampling @ nn=0 pred locs:
r=1000 d=[0] [0] [0]; mh=3 n=(51,11,18) k=0.129155
**PRUNE** @depth 1: [2,0.5]
**GROW** @depth 1: [2,0.55], n=(51,11)
**PRUNE** @depth 1: [1,0.5]
r=2000 d=[0] [0]; mh=3 n=(62,18) k=0.129155
Grow: 1.867%, Prune: 1.241%, Change: 37.9%, Swap: 29.35%
finished repetition 2 of 2

```

```

effective sample sizes:
0: itemp=1, len=68, ess=68
1: itemp=0.774264, len=136, ess=20.6323
2: itemp=0.599484, len=214, ess=2.20532
3: itemp=0.464159, len=202, ess=1.80005
4: itemp=0.359381, len=164, ess=1.00942
5: itemp=0.278256, len=148, ess=1.42103
6: itemp=0.215443, len=202, ess=0.995485
7: itemp=0.16681, len=334, ess=1.8323

```

```

8: itemp=0.129155, len=284, ess=0.996492
9: itemp=0.1, len=248, ess=1.01059
total: len=2000, ess.sum=99.903, ess(w)=100.038
lambda-combined ess=100.038

```

Notice how the MCMC inference procedure starts with $B + T = 4000$ rounds of stochastic approximation (initial adjustment of the pseudo-prior) in place of typical (default) the $B = 1000$ burn-in rounds. Then, the first round of sampling from the posterior commences, over $T = 2000$ rounds, during which the observation counts in each temperature are tallied. The progress meter shows the current temperature the chain is in, say $k=0.629961$, after each of 1000 sampling rounds. The first repeat starts with a pseudo-prior that has been adjusted by the observation counts, which continue to be accumulated throughout the entire procedure (i.e., they are never reset). Any subsequent repeats begin after a similar (re-)adjustment.

Before finishing, the routine summarizes the sample size and effective sample sizes in each rung of the temperature ladder. The number of samples is given by `len`, and the ESS by `ess`. These quantities can also be recovered via `traces`, as shown later. The ESS of the optimal combined IT sample is the last quantity printed. This, along with the ESS and total numbers of samples in each temperature, can also be obtained via the `tgp`-class output object.

```
> exp.btlm$ess
```

```
$combined
[1] 100.0378
```

```
$each
```

	k	count	ess
1	1.0000000	68	68.0000000
2	0.7742637	136	20.6322868
3	0.5994843	214	2.2053170
4	0.4641589	202	1.8000525
5	0.3593814	164	1.0094190
6	0.2782559	148	1.4210279
7	0.2154435	202	0.9954848
8	0.1668101	334	1.8323024
9	0.1291550	284	0.9964919
10	0.1000000	248	1.0105868

4.4.1 Motorcycle accident data

Recall the motorcycle accident data of Section 3.4 of the first `tgp` vignette [9]. Consider using IT to sample from the posterior distribution of the treed GP LLM model using the geometric temperature ladder.

```
> library(MASS)
> moto.it <- btgppllm(X = mcycle[, 1], Z = mcycle[,
+   2], BTE = c(2000, 52000, 10), bprior = "b0",
+   R = 3, itemps = geo, trace = TRUE, pred.n = FALSE,
+   verb = 0)
```

Out of a total of 15600 samples from the joint chain, the resulting (optimally combined) ESS was:

```
> moto.it$ess$combined

[1] 890.3704
```

Alternatively, \mathbf{w}^{λ^*} can be extracted from the traces, and used to make the ESS calculation directly.

```
> p <- moto.it$trace$post
> ESS(p$wlambda)

[1] 890.3704
```

The unadjusted weights \mathbf{w} are also available from `trace`. We can see that the naïve choice of $\lambda_i = W_i/W$, leading to the estimator in (9), has a clearly inferior effective sample size.

```
> ESS(p$w)

[1] 97.30426
```

To see the benefit of IT over ST we can simply count the number of samples obtained when $k^{(t)} = 1$. This can be accomplished in several ways: either via the traces or through the output object.

```
> as.numeric(c(sum(p$itemp == 1), moto.it$ess$each[1,
+   2:3]))

[1] 343 343 343
```

That is, (optimal) IT gives effectively 2.6 times more samples. The naïve combination, leading to the estimator in (9), yields an estimator with an effective sample size that is 28% of the number of samples obtained under ST.

Now, we should like to compare to the MCMC samples obtained under the same model, without IT.

```
> moto.reg <- btgppllm(X = mcycle[, 1], Z = mcycle[,
+   2], BTE = c(2000, 52000, 10), R = 3, bprior = "b0",
+   trace = TRUE, pred.n = FALSE, verb = 0)
```

The easiest comparison to make is to look at the heights explored under the three chains: the regular one, the chain of heights visited at all temperatures (combined), and those obtained after applying IT via re-weighting under the optimal combination λ^* .

```
> L <- length(p$height)
> hw <- suppressWarnings(sample(p$height, L, prob = p$wlambda,
+   replace = TRUE))
> b <- hist2bar(cbind(moto.reg$trace$post$height, p$height,
+   hw))

> barplot(b, beside = TRUE, col = 1:3, xlab = "tree height",
+   ylab = "counts", main = "tree heights encountered")
> legend("topright", c("reg MCMC", "All Temps", "IT"),
+   fill = 1:3)
```

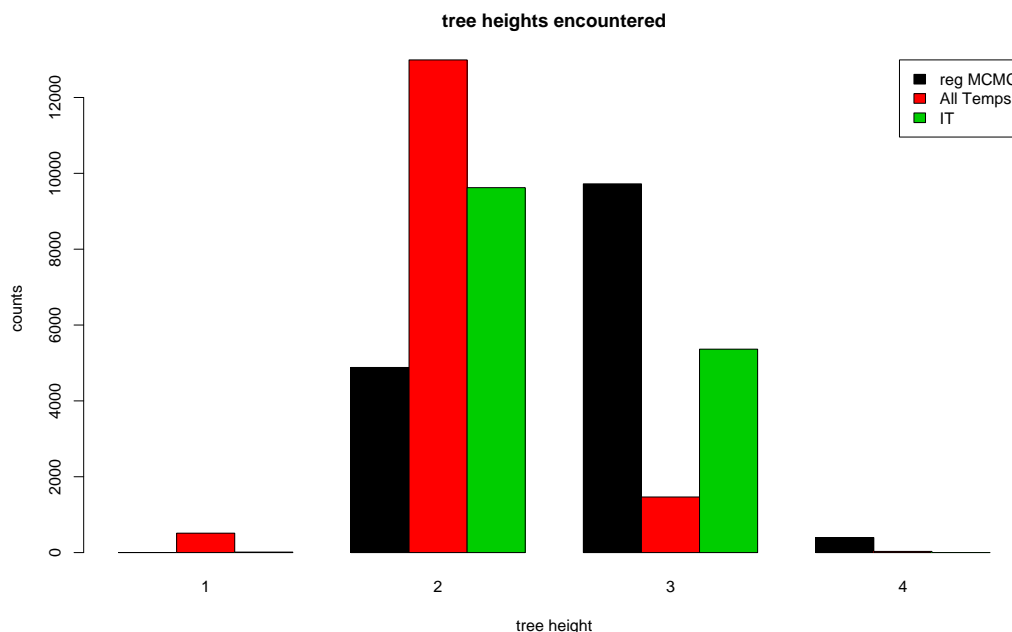


Figure 16: Barplots indicating the counts of the number of times the Markov chains (for regular MCMC, combining all temperatures in the inverse temperature ladder, and those re-weighted via IT) were in trees of various heights for the motorcycle data.

Figure 16 shows barplots indicating the count of the number of times the Markov chains were in trees of various heights after burn-in. Notice how the tempered chain (denoted “All Temps” in the figure) frequently visits trees of height one, whereas the non-tempered chain (denoted “reg MCMC”) never does. The result is that the non-tempered chain underestimates the probability of height two trees and produces a corresponding overestimate of height four trees—which are clearly not supported by the data—even visiting trees of height five. The IT estimator appropriately down-weights height one trees and provides correspondingly more realistic estimates of the probability of height two and four trees.

Whenever introducing another parameter into the model, like the inverse temperature k , it is important to check that the marginal posterior chain for that parameter is mixing well. For ST it is crucial that the chain makes rapid excursions between the cold temperature, the hottest temperatures, and visits each temperature roughly the same number of times.

```
> plot(log(moto.it$trace$post$itemp), type = "l", ylab = "log(k)",
+       xlab = "samples", main = "trace of log(k)")
```

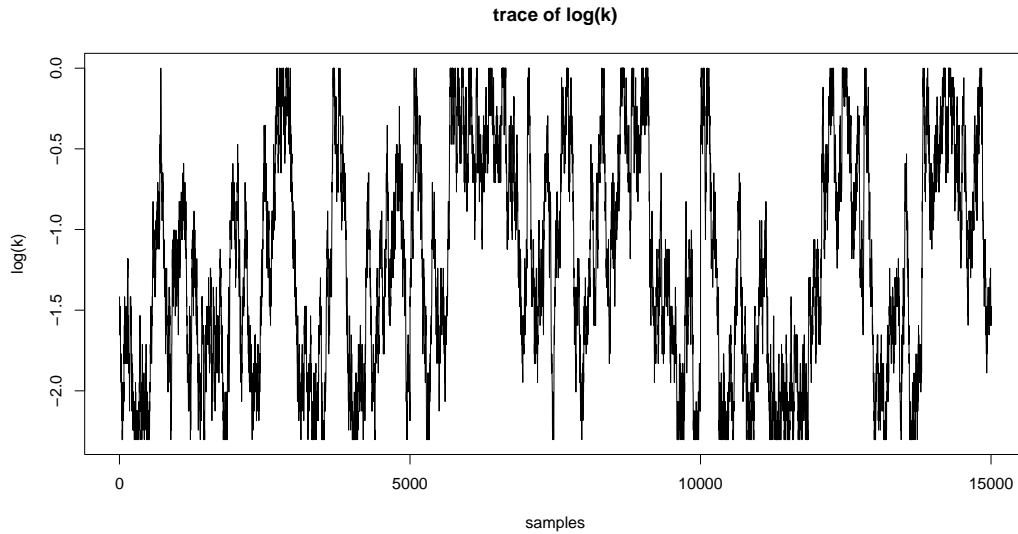


Figure 17: A trace of the MCMC samples from the marginal posterior distribution of the inverse temperature parameter, k , in the motorcycle experiment

Figure 17 shows a trace of the posterior samples for k in the motorcycle experiment. Arguably, the mixing in k -space leaves something to be desired. Since it can be very difficult to tune the pseudo-prior and MH proposal mechanism to get good mixing in k -space, it is fortunate that the IT methodology does not rely on the same mixing properties as ST does. Since samples can be obtained from the posterior distribution of the parameters of interest by re-weighting samples obtained when $k < 1$ it is only important that the chain frequently visit low temperatures to obtain good sampling, and high temperatures to obtain good mixing. The actual time spent in specific temperatures, i.e., $k = 1$ is less important. Figure 18 shows the histogram of the inverse temperatures visited in the Markov chain for the motorcycle experiment. Also plotted is a histogram of the *observation counts* in each temperature. The two histograms should have similar shape but different totals. Observation counts are tallied during every MCMC sample after burn-in, whereas the posterior samples of k are thinned (at a rate specified in `BTE[3]`). When the default `trace=FALSE` argument is used only the observation counts will be available in the `tgp`-class object, and these can be used as a surrogate for a trace of k .

The compromise IT approach obtained using the sigmoidal ladder can yield an increase in ESS.

```

> b <- itemps.barplot(moto.it, plot.it = FALSE)
> barplot(t(cbind(moto.it$items$counts, b)), col = 1:2,
+   beside = TRUE, ylab = "counts", xlab = "items",
+   main = "inv-temp observation counts")
> legend("topleft", c("observation counts", "posterior samples"),
+   fill = 1:2)

```

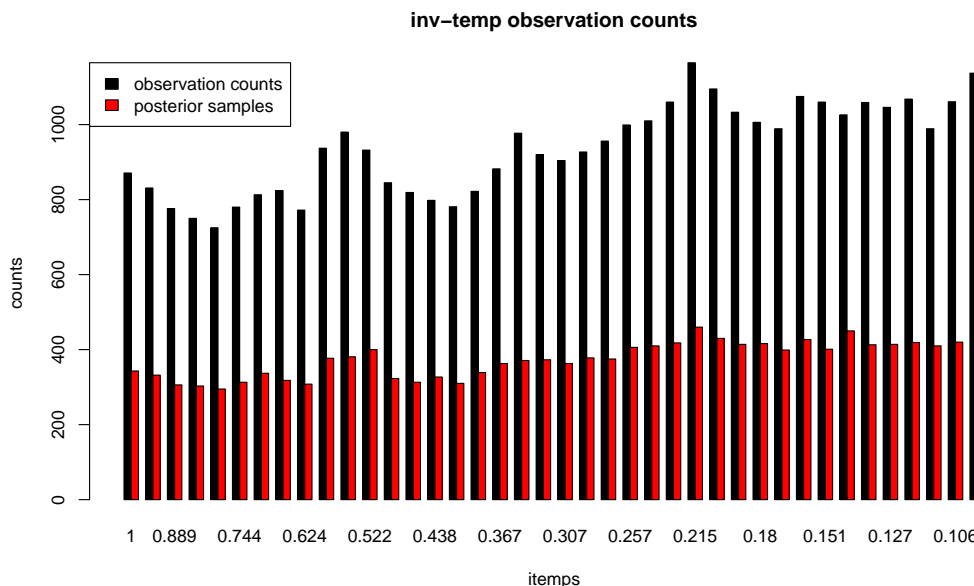


Figure 18: Comparing (thinned) samples from the posterior distribution for the inverse temperature parameter, k , (posterior samples), to the observation counts used to update the pseudo-prior, in the motorcycle experiment

```

> moto.it.sig <- btgpllm(X = mcycle[, 1], Z = mcycle[,
+   2], BTE = c(2000, 52000, 10), R = 3, bprior = "b0",
+   krige = FALSE, itemps = sig, verb = 0)

```

Compare the resulting ESS to the one given for the geometric ladder above.

```

> moto.it.sig$ess$combined

```

```

[1] 5147.576

```

Plots of the resulting predictive surface is shown in Figure 19 for comparison with those in Section 3.4 of the first `tgpl` vignette [9].

4.4.2 Synthetic 2-d Exponential Data

Recall the synthetic 2-d exponential data of Section 3.4 of the `tgpl` vignette [9], where the true response is given by

$$z(\mathbf{x}) = x_1 \exp(-x_1^2 - x_2^2).$$

```
> plot(moto.it.sig)
```

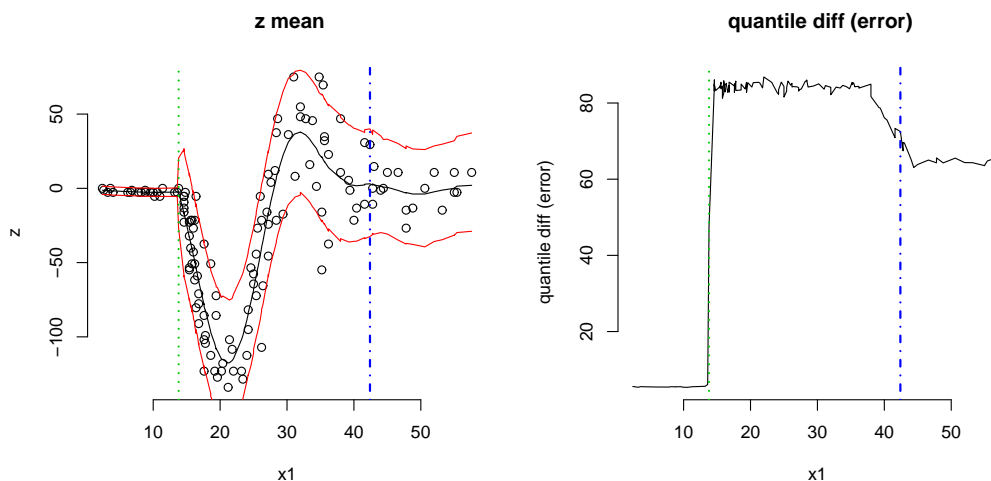


Figure 19: Posterior predictive surface for the motorcycle data, with 90% quantile errorbars, obtained under IT with the sigmoidal ladder.

Here, we will take $\mathbf{x} \in [-6, 6] \times [-6, 6]$ with a D -optimal design

```
> Xcand <- lhs(10000, rbind(c(-6, 6), c(-6, 6)))
> X <- dopt.gp(400, X = NULL, Xcand)$XX
> Z <- exp2d.Z(X)$Z
```

Consider a treed GP LLM model fit to this data using the standard RJ-MCMC.

```
> exp.reg <- btgpllm(X = X, Z = Z, BTE = c(2000, 52000,
+ 10), bprior = "b0", trace = TRUE, krige = FALSE,
+ R = 10, verb = 0)
```

Figure 20 shows the resulting posterior predictive surface. The maximum *a posteriori* (MAP) tree is drawn over the error surface in the *right-hand* plot. The height of this tree can be obtained from the `tgp`-class object.

```
> h <- exp.reg$post$height[which.max(exp.reg$post$lpst)]
> h
```

```
[1] 7
```

It is easy to see that many fewer partitions are actually necessary to separate the interesting, central, region from the surrounding flat region. Figure 21 shows a diagrammatic representation of the MAP tree. Given the apparent over-partitioning in this height 7 tree it would be surprising to find much posterior support for trees of greater height. One might indeed suspect that there are trees with fewer partitions which would have higher posterior probability, and thus guess that the Markov chain for the trees plotted in these figures possibly became stuck in a local mode of tree space while on an excursion into deeper trees.


```
> plot(exp.reg)
```

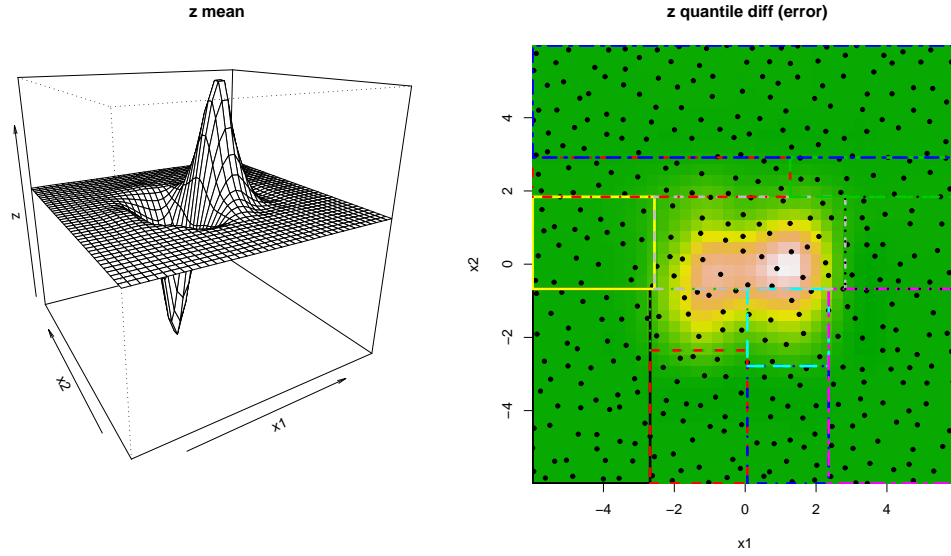


Figure 20: Posterior predictive surface for the 2-d exponential data: mean surface (*left*) and 90% quantile difference (*right*)

```
> tgp.trees(exp.reg, "map")
```

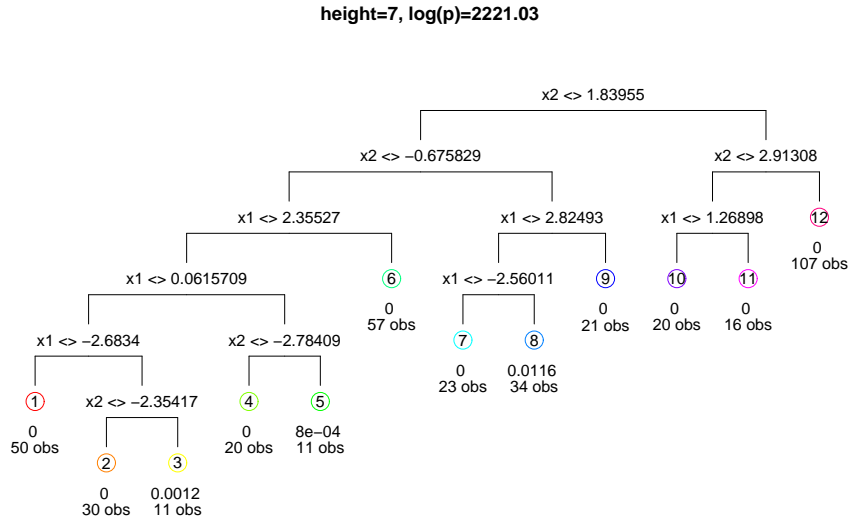


Figure 21: Diagrammatic depiction of the maximum *a posteriori* (MAP) tree for the 2-d exponential data under standard RJ-MCMC sampling

Now consider using IT. It will be important in this case to have a k_m small enough to ensure that the tree occasionally prunes back to the root. We shall therefore use a smaller k_m . Generally speaking, some pilot tuning may be necessary to choose an appropriate k_m and number of rungs m , although the defaults should give adequate performance in most

cases.

```
> its <- default.itemps(k.min = 0.02)
> exp.it <- btgpllm(X = X, Z = Z, BTE = c(2000, 52000,
+   10), bprior = "b0", trace = TRUE, krige = FALSE,
+   itemps = its, R = 10, verb = 0)
```

As expected, the tempered chain moves more rapidly throughout tree space by accepting more tree proposals. The acceptance rates of tree operations can be accessed from the `tgpl`-class object.

```
> exp.it$gpcs

      grow      prune      change      swap
1 0.0825418 0.08196492 0.8277013 0.5840494

> exp.reg$gpcs

      grow      prune      change      swap
1 0.01058038 0.007811603 0.6275904 0.3189888
```

The increased rate of *prune* operations explains how the tempered distributions helped the chain escape the local modes of deep trees.

We can quickly compare the effective sample sizes of the three possible estimators: ST, naïve IT, and optimal IT.

```
> p <- exp.it$trace$post
> data.frame(ST = sum(p$itemp == 1), nIT = ESS(p$w),
+   oIT = exp.it$ess$combined)

      ST      nIT      oIT
1 1685 1626.539 2000.642
```

Due to the thinning in the Markov chain ($\text{BTE}[3] = 10$) and the traversal between $m = 10$ temperatures in the ladder, we can be reasonably certain that the 2001 samples obtained via IT from the total of 50000 samples obtained from the posterior are far less correlated than the ones obtained via standard RJ-MCMC.

As with the motorcycle data, we can compare the tree heights visited by the two chains.

```
> L <- length(p$height)
> hw <- suppressWarnings(sample(p$height, L, prob = p$wlambda,
+   replace = TRUE))
> b <- hist2bar(cbind(exp.reg$trace$post$height, p$height,
+   hw))
```

```

> barplot(b, beside = TRUE, col = 1:3, xlab = "tree height",
+       ylab = "counts", main = "tree heights encountered")
> legend("topright", c("reg MCMC", "All Temps", "IT"),
+       fill = 1:3)

```

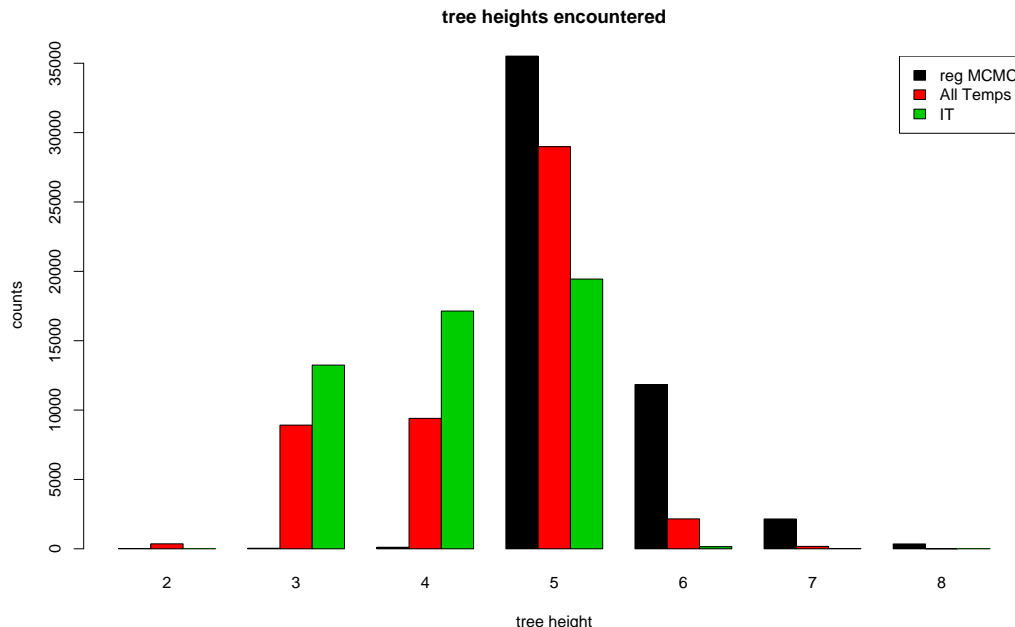


Figure 22: Barplots indicating the counts of the number of times the Markov chains (for regular MCMC, combining all temperatures in the inverse temperature ladder, and those re-weighted via IT) were in trees of various heights for the 2-d exponential data.

Figure 22 shows a barplot of \mathbf{b} , which illustrates that the tempered chain frequently visited shallow trees. IT with the optimal weights shows that the standard RJ-MCMC chain missed many trees of height three and four with considerable posterior support.

To more directly compare the mixing in tree space between the ST and tempered chains, consider the trace plots of the heights of the trees explored by the chains shown in Figure 23. Despite being restarted 10 times, the regular RJ-MCMC chain (almost) never visits trees of height less than five after burn-in and instead makes rather lengthy excursions into deeper trees, exploring a local mode in the posterior. In contrast, the tempered chain frequently prunes back to the tree root, and consequently discovers posterior modes in tree heights three and four.

To conclude, a plot of the posterior predictive surface is given in Figure 24, where the MAP tree is shown both graphically and diagrammatically.

```

> ylim <- range(p$height, exp.reg$trace$post$height)
> plot(p$height, type = "l", main = "trace of tree heights",
+      xlab = "t", ylab = "height", ylim = ylim)
> lines(exp.reg$trace$post$height, col = 2)
> legend("topright", c("tempered", "reg MCMC"), lty = c(1,
+      1), col = 1:2)

```

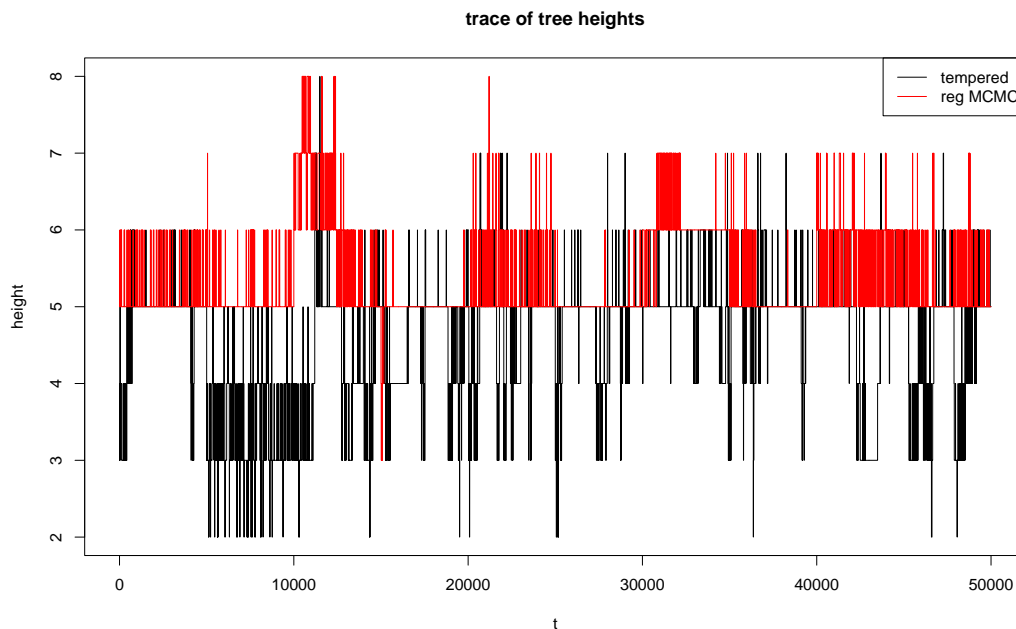


Figure 23: Traces of the tree heights obtained under the two Markov chains Markov chains (for regular MCMC, combining all temperatures in the inverse temperature ladder) on the 2-d exponential data.

References

- [1] L. Breiman, J. H. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Wadsworth, Belmont, CA, 1984.
- [2] H. A. Chipman, E. I. George, and R. E. McCulloch. Bayesian CART model search (with discussion). *Journal of the American Statistical Association*, 93:935–960, 1998.
- [3] H. A. Chipman, E. I. George, and R. E. McCulloch. Bayesian treed models. *Machine Learning*, 48:303–324, 2002.
- [4] R. Douc, A. Guillin, J.-M. Marin, and C.P. Robert. Minimum variance importance sampling via population monte carlo. Technical report, CEREMADE, Université Paris Dauphine, and CREST, INSEE, Paris, 2007.
- [5] J. H. Friedman. Multivariate adaptive regression splines. *Annals of Statistics*, 19, No. 1:1–67, March 1991.

```
> plot(exp.it)
> tgp.trees(exp.it, "map")
```

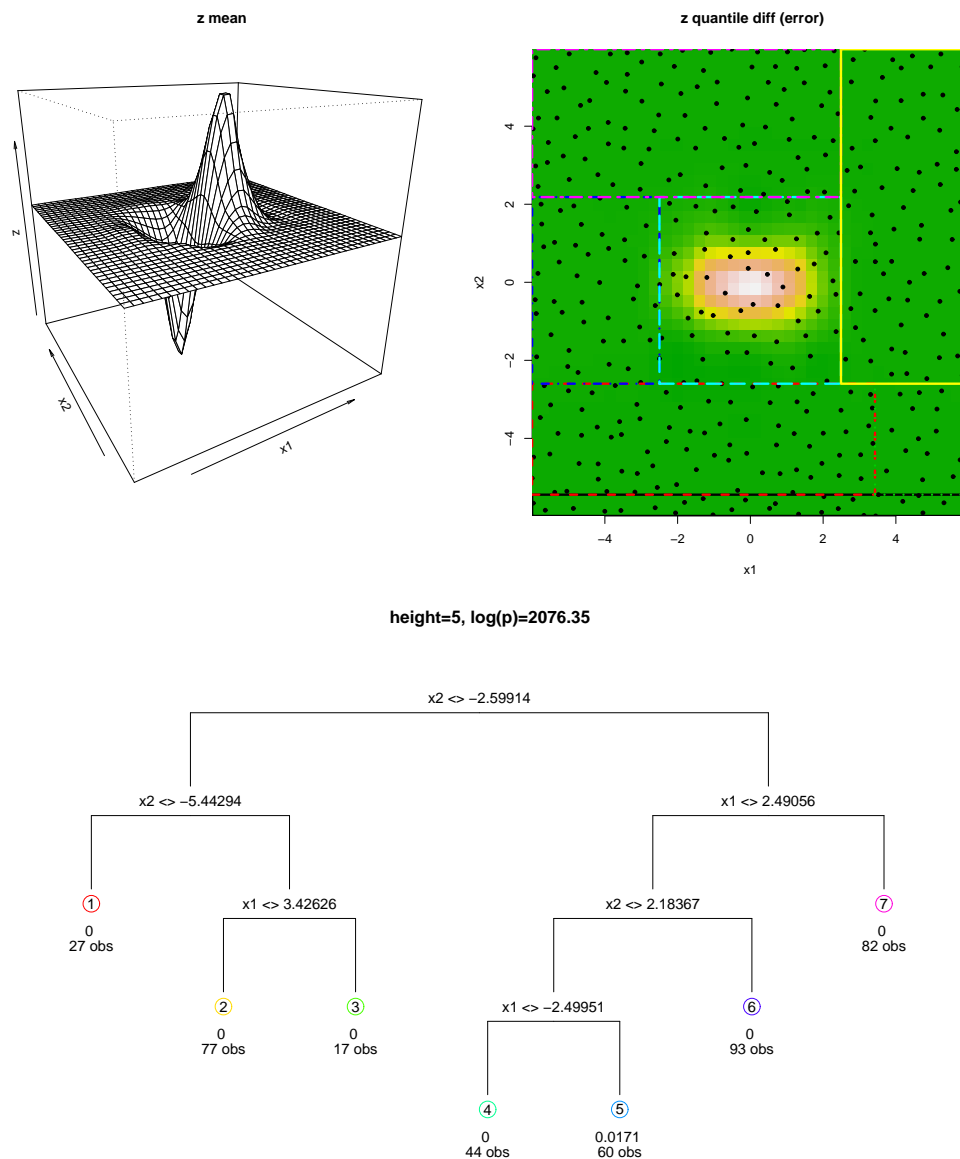


Figure 24: 2-d exponential data fit with IT. *Top*: Posterior predictive mean surface for the 2d-exponential, with the MAP tree overlayed. *Bottom*: diagrammatic representation of the MAP tree.

- [6] S. Geman and D. Geman. Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741, 1984.

- [7] C.J. Geyer. Markov chain Monte Carlo maximum likelihood. In *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*, pages 156–163, 1991.
- [8] C.J. Geyer and E.A. Thompson. Annealing Markov chain Monte Carlo with applications to ancestral inference. *Journal of the American Statistical Association*, 90:909–920, 1995.
- [9] Robert B. Gramacy. `tgp`: An R package for Bayesian nonstationary, semiparametric nonlinear regression and design by treed gaussian process models. *Journal of Statistical Software*, 19(9), 6 2007.
- [10] Robert B. Gramacy, Richard J. Samworth, and Ruth King. Importance tempering. Technical Report 0707.4242, ArXiv, 2009. to appear in *Statistics and Computing*.
- [11] W.K. Hastings. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57:97–109, 1970.
- [12] T. Homma and A. Saltelli. Importance measures in global sensitivity analysis of non-linear models. *Reliability engineering and system safety*, 52:1–17, 1996.
- [13] C. Jennison. Discussion on the meeting on the gibbs sampler and other Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society, Series B*, 55:54–56, 1993.
- [14] D.R. Jones, M. Schonlau, and W.J. Welch. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13:455–492, 1998.
- [15] Robert E. Kass, Bradley P. Carlin, Andrew Gelman, and Radford M. Neal. Markov chain monte carlo in practice: A roundtable discussion. *The American Statistician*, 52(2):93–100, May 1998.
- [16] J. S. Liu. *Monte Carlo Strategies in Scientific Computing*. Springer, New York, 2001.
- [17] N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, and R. Teller. Equations of state calculations by fast computing machine. *Journal of Chemical Physics*, 21:1087–1091, 1953.
- [18] R. D. Morris, A. Kottas, M. Taddy, R. Furfaro, and B. Ganapol. A statistical framework for the sensitivity analysis of radiative transfer models. *IEEE Transactions on Geoscience and Remote Sensing*, To appear.
- [19] Radford M. Neal. Sampling from multimodal distributions using tempered transition. *Statistics and Computing*, 6:353–366, 1996.
- [20] Radford M. Neal. Annealed importance sampling. *Statistics and Computing*, 11:125–129, 2001.
- [21] J.E. Oakley and A. O’Hagan. Probabilistic sensitivity analysis of complex models: a Bayesian approach. *Journal of the Royal Statistical Society Series B*, 66:751–769, 2004.

- [22] Art Owen and Yi Zhou. Safe and effective importance sampling. *Journal of the American Statistical Association*, 95(449):135–143, March 2000.
- [23] Martyn Plummer, Nicky Best, Kate Cowles, and Karen Vines. *coda: Output analysis and diagnostics for MCMC*, 2008. R package version 0.13-3.
- [24] Rommel G. Regis and Christine A. Shoemaker. Improved strategies for radial basis function methods for global optimization. *J. of Global Optimization*, 37(1):113–135, 2007.
- [25] A. Saltelli, K. Chan, and E.M. Scott, editors. *Sensitivity Analysis*. John Wiley and Sons, 2000.
- [26] A. Saltelli, M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana, and S. Tarantola. *Global Sensitivity Analysis: The Primer*. John Wiley & Sons, 2008.
- [27] A. Saltelli and S. Tarantola. On the relative importance of input factors in mathematical models: safety assessment for nuclear waste disposal. *Journal of the American Statistical Association*, 97:702–709, 2002.
- [28] Andrea Saltelli. Making best use of model evaluations to compute sensitivity indices. *Computer Physics Communications*, 145:280–297, 2002.
- [29] M. Schonlau, D.R. Jones, and W.J. Welch. Global versus local search in constrained optimization of computer models. In *New Developments and applications in experimental design*, number 34 in IMS Lecture Notes - Monograph Series, pages 11–25. IMS, 1998.
- [30] I. M. Sobol’. Sensitivity analysis for nonlinear mathematical models. *Mathematical Modeling and Computational Experiment*, 1:407–414, 1993.
- [31] Matthew Taddy, Herbert K. H. Lee, Genetha A. Gray, and Joshua D. Griffin. Bayesian guided pattern search for robust local optimization. *Technometrics*, To appear.
- [32] Eric Veach and Leonidas J. Guibas. Optimally combining sampling techniques for monte carlo rendering. In *SIGGRAPH ’95 Conference Proceedings*, pages 419–428, Reading, MA, 1995. Addison–Wesley.
- [33] W. J. Welch, R. J. Buck, J. Sacks, H. P. Wynn, T.J Mitchell, and M. D. Morris. Screening, predicting, and computer experiment. *Technometrics*, 34:15–25, 1992.