

---

# Probabilistic Programming with Gaussian Process Memoization

---

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

This paper describes the *Gaussian process memoizer*, a probabilistic programming technique that uses Gaussian processes to provide a statistical alternative to memorization. Memoizing a target procedure results in a self-caching wrapper that remembers previously computed values. Gaussian process memoization additionally produces a statistical emulator based on a Gaussian process whose predictions automatically improve whenever a new value of the target procedure becomes available. This paper also introduces an efficient implementation, named `gpmem`, that can use kernels given by a broad class of probabilistic programs. The flexibility of `gpmem` is illustrated via three applications: (i) GP regression with hierarchical hyper-parameter learning, (ii) Bayesian structure learning via compositional kernels generated by a probabilistic grammar, and (iii) a bandit formulation of Bayesian optimization with automatic inference and action selection. All applications share a single 50-line Python library and require fewer than 20 lines of probabilistic code each.

## 1 Introduction

Gaussian Processes (GP) are widely used tools in statistics (Barry, 1986), machine learning (Neal, 1995; Williams and Barber, 1998; Kuss and Rasmussen, 2005; Rasmussen and Williams, 2006; Damianou and Lawrence, 2013), robotics (Ferris et al., 2006), computer vision (Kemmler et al., 2013), and scientific computation (Kennedy and O’Hagan, 2001; Schneider et al., 2008; Kwan et al., 2013). They are also central to probabilistic numerics, an emerging effort to develop more computationally efficient numerical procedures, and to Bayesian optimization, a family of meta-optimization techniques that are widely used to tune parameters for deep learning algorithms (Snoek et al., 2012; Gelbart et al., 2014). They have even seen use in artificial intelligence; for example, they provide the key technology behind a project that produces qualitative natural language descriptions of time series (Duvenaud et al., 2013; Lloyd et al., 2014).

This paper describes Gaussian process memoization, a technique for integrating GPs into a probabilistic programming language, and demonstrates its utility by re-implementing and extending state-of-the-art applications of the GP. Memoization, typically implemented by a procedure called `mem()`, is a classic higher-order programming technique in which a procedure is augmented with an input-output cache that is checked each time before the function is invoked. This prevents unnecessary recomputation, potentially saving time at the cost of increased storage requirements. Gaussian process memoization, implemented by the `gpmem()` procedure, generalizes this idea to include a statistical emulator that uses previously computed values as data in a statistical model that can cheaply forecast probable outputs. The covariance function for the Gaussian process is also allowed to be an arbitrary probabilistic program.

This paper presents three applications of `gpmem`: (i) a replication of (Neal, 1997) results on outlier rejection via hyper-parameter inference; (ii) a fully Bayesian extension to the Automated Statis-

054 tician project; and (iii) an implementation of Bayesian optimization via Thompson sampling. The  
 055 first application can in principle be replicated in several other probabilistic languages embedding the  
 056 proposal that is described in this paper. The remaining two applications rely on distinctive capabili-  
 057 ties of Venture: support for fully Bayesian structure learning and language constructs for inference  
 058 programming. All applications share a single 50-line Python library and require fewer than 20 lines  
 059 of probabilistic code each.  
 060

## 061 2 Background on Gaussian Processes

063 Let the data be pairs of real-valued scalars  $\{(x_i, y_i)\}_{i=1}^n$  (complete data will be denoted by column  
 064 vectors  $\mathbf{x}, \mathbf{y}$ ). GPs present a non-parametric way to express prior knowledge on the space of possible  
 065 functions  $f$  modeling a regression relationship. Formally, a GP is an infinite-dimensional extension  
 066 of the multivariate Gaussian distribution. Often one assumes the values  $\mathbf{y}$  are noisily measured, that  
 067 is, one only sees the values of  $\mathbf{y}_{\text{noisy}} = \mathbf{y} + \mathbf{w}$  where  $\mathbf{w}$  is Gaussian white noise with variance  $\sigma_{\text{noise}}^2$ .  
 068 In that case, the log-likelihood of a GP is

$$069 \log p(\mathbf{y}_{\text{noisy}} | \mathbf{x}) = -\frac{1}{2}\mathbf{y}^\top(\Sigma + \sigma_{\text{noise}}^2\mathbf{I})^{-1}\mathbf{y} - \frac{1}{2}\log|\Sigma + \sigma_{\text{noise}}^2\mathbf{I}| - \frac{n}{2}\log 2\pi \quad (1)$$

070 where  $n$  is the number of data points. Both log-likelihood and predictive posterior can be computed  
 071 efficiently in a Venture SP with an algorithm that resorts to Cholesky factorization(Rasmussen and  
 072 Williams, 2006, chap. 2) resulting in a computational complexity of  $\mathcal{O}(n^3)$  in the number of data  
 073 points.  
 074

075 The covariance function (or kernel) of a GP governs high-level properties of the observed data such  
 076 as linearity, periodicity and smoothness. It comes with few free parameters that we call hyper-  
 077 parameters. Adjusting these results in minor changes, for example with regards to when two data  
 078 points are treated similar. More drastically different covariance functions are achieved by changing  
 079 the structure of the covariance function itself. Note that covariance function structures are compo-  
 080 sitional: adding or multiplying two valid covariance functions results in another valid covariance  
 081 function.  
 082

083 Venture includes the primitive `make_gp`, which takes as arguments a unary function `mean` and a  
 084 binary (symmetric, positive-semidefinite) function `cov` and produces a function `g` distributed as a  
 085 Gaussian process with the supplied mean and covariance. For example, a function  $g \sim \mathcal{GP}(0, \text{SE})$ ,  
 086 where `SE` is a squared-exponential covariance

$$087 \text{SE}(x, x') = \sigma^2 \exp\left(\frac{(x - x')^2}{2\ell}\right)$$

088 with  $\sigma = 1$  and  $\ell = 1$ , can be instantiated as follows:  
 089

```
090 assume zero = make_const_func( 0.0 )
091 assume se = make_squaredexp( 1.0, 1.0 )
092 assume g = make_gp( zero, se )
```

093 There are two ways two view `g` as a “random function.” In the first view, the `assume` directive that  
 094 instantiates `g` does not use any randomness—only the subsequent calls to `g` do—and coherence  
 095 constraints are upheld by the interpreter by keeping track of which evaluations of `g` exist in the current  
 096 execution trace. Namely, if the current trace contains evaluations of `g` at the points  $x_1, \dots, x_N$  with  
 097 return values  $y_1, \dots, y_N$ , then the next evaluation of `g` (say, jointly at the points  $x_{N+1}, \dots, x_{N+n}$ )  
 098 will be distributed according to the joint conditional distribution  
 099

$$100 P((g x_{N+1}), \dots, (g x_{N+n}) | (g x_i) = y_i \text{ for } i = 1, \dots, N).$$

101 In the second view, `g` is a randomly chosen deterministic function, chosen from the space of all  
 102 deterministic real-valued functions; in this view, the `assume` directive contains *all* the randomness,  
 103 and subsequent invocations of `g` are deterministic. The first view is procedural and is faithful to the  
 104 computation that occurs behind the scenes in Venture. The second view is declarative and is faithful  
 105 to notations like “ $g \sim P(g)$ ” which are often used in mathematical treatments. Because a model  
 106 program could make arbitrarily many calls to `g`, and the joint distribution on the return values of the  
 107

108 calls could have arbitrarily high entropy, it is not computationally possible in finite time to choose  
 109 the entire function  $g$  all at once as in the second view. Thus, it stands to reason that any computationally  
 110 implementable notion of “nonparametric random functions” must involve incremental random  
 111 choices in one way or another, and Gaussian processes in Venture are no exception.  
 112

## 113 2.1 The `gpmem` construct

114  
 115 Memoization is the practice of storing previously computed values of a function so that future calls  
 116 with the same inputs can be evaluated by lookup rather than recomputation. Although memoization  
 117 does not change the semantics of a deterministic program, it does change that of a stochastic  
 118 program (Goodman et al., 2008). The authors provide an intuitive example: let  $f$  be a function  
 119 that flips a coin and return “head” or “tails”. The probability that two calls of  $f$  are equivalent is  
 120 0.5. However, if the function call is memoized, it is 1. In fact, there is an infinite range of possible  
 121 caching policies (specifications of when to use a stored value and when to recompute), each potentially  
 122 having a different semantics. Any particular caching policy can be understood by random  
 123 world semantics (Poole, 1993; Sato, 1995) over the stochastic program: each possible world corre-  
 124 sponds to a mapping from function input sequence to function output sequence (McAllester et al.,  
 125 2008). In Venture, these possible worlds are first-class objects, and correspond to the *probabilistic*  
 126 *execution traces* (Mansinghka et al., 2014).

127 To transfer this idea to probabilistic programming, we now introduce a language construct called a  
 128 *statistical memoizer*. Suppose we have a function  $f$  which can be evaluated but we wish to learn  
 129 about the behavior of  $f$  using as few evaluations as possible. The statistical memoizer, which here  
 130 we give the name `gpmem`, was motivated by this purpose. It produces two outputs:  
 131

$$f \xrightarrow{\text{gpmem}} (f_{\text{probe}}, f_{\text{emu}}).$$

132 The function  $f_{\text{probe}}$  calls  $f$  and stores the output in a memo table, just as traditional memoization  
 133 does. The function  $f_{\text{emu}}$  is an online statistical emulator which uses the memo table as its training  
 134 data. A fully Bayesian emulator, modelling the true function  $f$  as a random function  $f \sim P(f)$ ,  
 135 would satisfy

$$(f_{\text{emu}} x_1 \dots x_k) \sim P(f(x_1), \dots, f(x_k) | f(x) = (f x) \text{ for each } x \text{ in memo table}).$$

136 Different implementations of the statistical memoizer can have different prior distributions  $P(f)$ ; in  
 137 this paper, we deploy a Gaussian process prior (implemented as `gpmem` below). Note that we require  
 138 the ability to sample  $f_{\text{emu}}$  jointly at multiple inputs because the values of  $f(x_1), \dots, f(x_k)$  will in  
 139 general be dependent.  
 140

141 In Venture, we initialize `gpmem` as follows:

```
142
143     assume K = make_squaredexp (sf, 1)
144     assume (f_compute f_emu) = gpmem( f, K))
```

145 Adding a single line to the program, such as

```
146
147     infer mh( quote( parameters), one, 50),
```

148 allows us perform Bayesian inference over the parameters  $sf$  and  $l$ . Such inference can be performed  
 149 without the refactoring and elaboration needed if one would describe the above declaratively as in  
 150 traditional statistics notation. In contrast to traditional statistics notation, probabilistic programs  
 151 are written procedurally; reasoning declaratively about the dynamics of a fundamentally procedural  
 152 inference algorithm is often unwieldy due to the absence of programming constructs such as loops  
 153 and mutable state.  
 154

155 We implement `gpmem` by memoizing a target procedure in a wrapper that remembers previously  
 156 computed values. This comes with interesting implications: from the standpoint of computation,  
 157 a data set of the form  $\{(x_i, y_i)\}$  can be thought of as a function  $y = f_{\text{look-up}}(x)$ , where  $f_{\text{look-up}}$  is  
 158 restricted to only allow evaluation at a specific set of inputs  $x$ .

159 Modelling the data set with a GP then amounts to trying to learn a smooth function  $f_{\text{emu}}$  (“emu”  
 160 stands for “emulator”) which extends  $f$  to its full domain. We can then incorporate observations  
 161 in two different ways: we either are either told that the at a point  $x$  the value is  $y$ :

```
162     observe f_emu ( x ) = y
```

163  
164 Or we express this as  
165

```
166     predict f_compute ( x )
```

167  
168 The second expression has at least two benefits: (i) readability (in some cases), and (ii) amenability  
169 to active learning. As to (i), the statistical code of creating a Gaussian process is replaced with a  
170 memoization-like idiom, which will be more familiar to programmers. As to (ii), when using gpmem,  
171 it is quite easy to decide incrementally which data point to sample next: for example, the loop from  
172  $x[1]$  to  $x[n]$  could be replaced by a loop in which the next index  $i$  is chosen by a supplied decision  
173 rule. In this way, we could use gpmem to perform online learning using only a subset of the available  
174 data.

175 We illustrate the use of gpmem in this context in a tutorial in Fig. 1.  
176

### 177 3 Applications 178

179 gpmem is an elegant linguistic framework for function learning-related tasks. The technique allows  
180 language constructs from programming to help express models which would be cumbersome to  
181 express in statistics notation. We will now illustrate this with three example applications.  
182

#### 183 3.1 Nonlinear regression in the presence of outliers 184

185 The probability of the hyper-parameters of a GP as defined above and given covariance function  
186 structure  $\mathbf{K}$  is:

$$187 P(\boldsymbol{\theta} | \mathbf{D}, \mathbf{K}) = \frac{P(\mathbf{D} | \boldsymbol{\theta}, \mathbf{K})P(\boldsymbol{\theta} | \mathbf{K})}{P(\mathbf{D} | \mathbf{K})}. \quad (2)$$

188

189 Let the  $\mathbf{K}$  be the sum of a smoothing and a white noise (WN) kernel. For this case, Neal (1997)  
190 suggested the problem of outliers in data as a use-case for a hierarchical Bayesian treatment of  
191 Gaussian processes<sup>1</sup>. The work suggests a hierarchical system of hyper-parameterization. Here, we  
192 draw hyper-parameters from a  $\Gamma$  distributions:

$$193 \ell^{(t)} \sim \Gamma(\alpha_1, \beta_1), \sigma^{(t)} \sim \Gamma(\alpha_2, \beta_2) \quad (3)$$

194 and in turn sample the  $\alpha$  and  $\beta$  from  $\Gamma$  distributions as well:

$$195 \alpha_1^{(t)} \sim \Gamma(\alpha_\alpha^1, \beta_\alpha^1), \alpha_2^{(t)} \sim \Gamma(\alpha_\alpha^2, \beta_\alpha^2), \dots \quad (4)$$

196

197 One can represent this kind of model using gpmem (Fig. 2). Neal provides a custom inference  
198 algorithm setting and evaluates it using the following synthetic data problem. Let  $f$  be the underlying  
199 function that generates the data:

$$200 f(x) = 0.3 + 0.4x + 0.5 \sin(2.7x) + \frac{1.1}{(1+x^2)} + \eta \quad \text{with } \eta \sim \mathcal{N}(0, \sigma) \quad (5)$$

201

202 We synthetically generate outliers by setting  $\sigma = 0.1$  in 95% of the cases and to  $\sigma = 1$  in the  
203 remaining cases. gpmem can capture the true underlying function within only 100 MH steps on  
204 the hyper-parameters to get a good approximation for their posterior. Note that Neal devises an  
205 additional noise model and performs a large number of Hybrid-Monte Carlo and Gibbs steps.

#### 206 3.2 Discovering qualitative structure from time series data 207

208 Inductive learning of symbolic expression for continuous-valued time series data is a hard task which  
209 has recently been tackled using a greedy search over the approximate posterior of the possible kernel  
210 compositions for GPs (Duvenaud et al., 2013; Lloyd et al., 2014)<sup>2</sup>.

211 With gpmem we can provide a fully Bayesian treatment of this, previously unavailable, using a stochastic  
212 grammar (see Fig. 3).

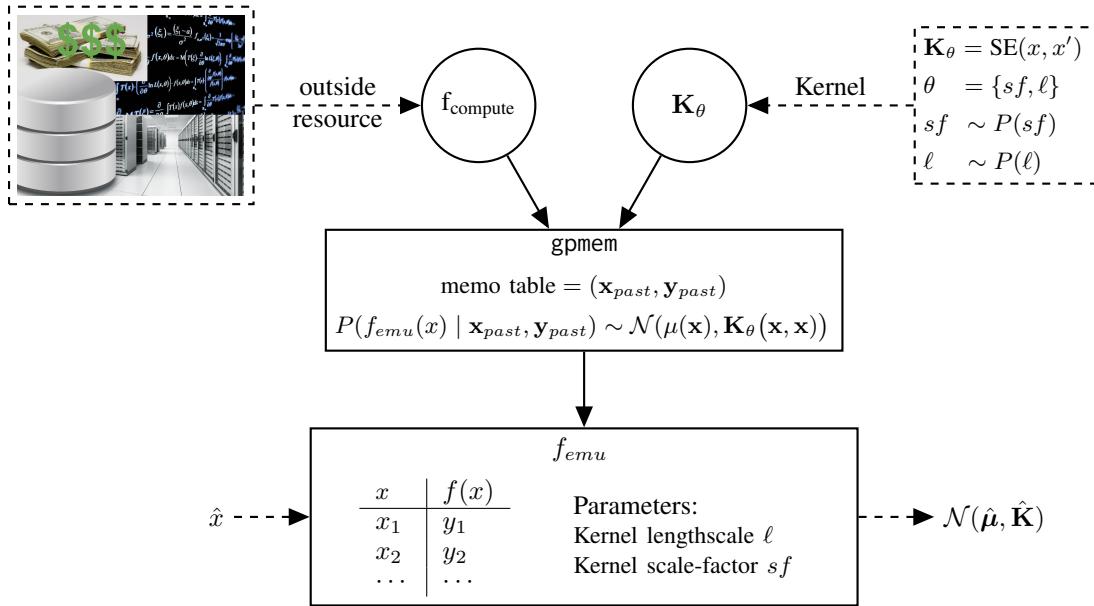
213  
214 <sup>1</sup>In (Neal, 1997) the sum of an SE plus a constant kernel is used. We keep the WN kernel for illustrative  
215 purposes.

<sup>2</sup><http://www.automaticstatistician.com/>

```

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

```



```

define f = proc( x ) {
    exp(-0.1*abs(x-2))) *
    10* cos(0.4*x) + 0.2
}
assume (f_compute f_emu) = gpmem( f, K)
sample f_emu( array( -20, ..., 20))

predict f_compute( 12.6)
sample f_emu( array( -20, ..., 20))

predict f_compute( -6.4)
sample f_emu( array( -20, ..., 20))

observe f_emu( -3.1) = 2.60
observe f_emu( 7.8) = -7.60
observe f_emu( 0.0) = 10.19

sample f_emu( array( -20, ..., 20))

infer mh(quote(hyper-parameter), one, 50)

sample f_emu( array( -20, ..., 20))

```

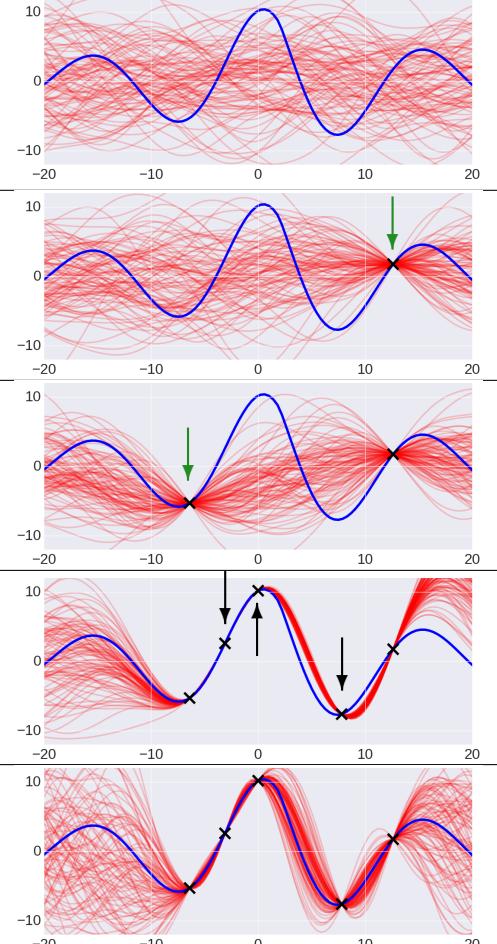


Figure 1: gpmem tutorial. The top shows a schematic of gpmem.  $f_{\text{compute}}$  probes an outside resource. This can be expensive (top left). Every probe is memoized and improves the GP-based emulator. Below the schematic we see a movie of the evolution of gpmem's state of believe of the world given certain Venture directives.

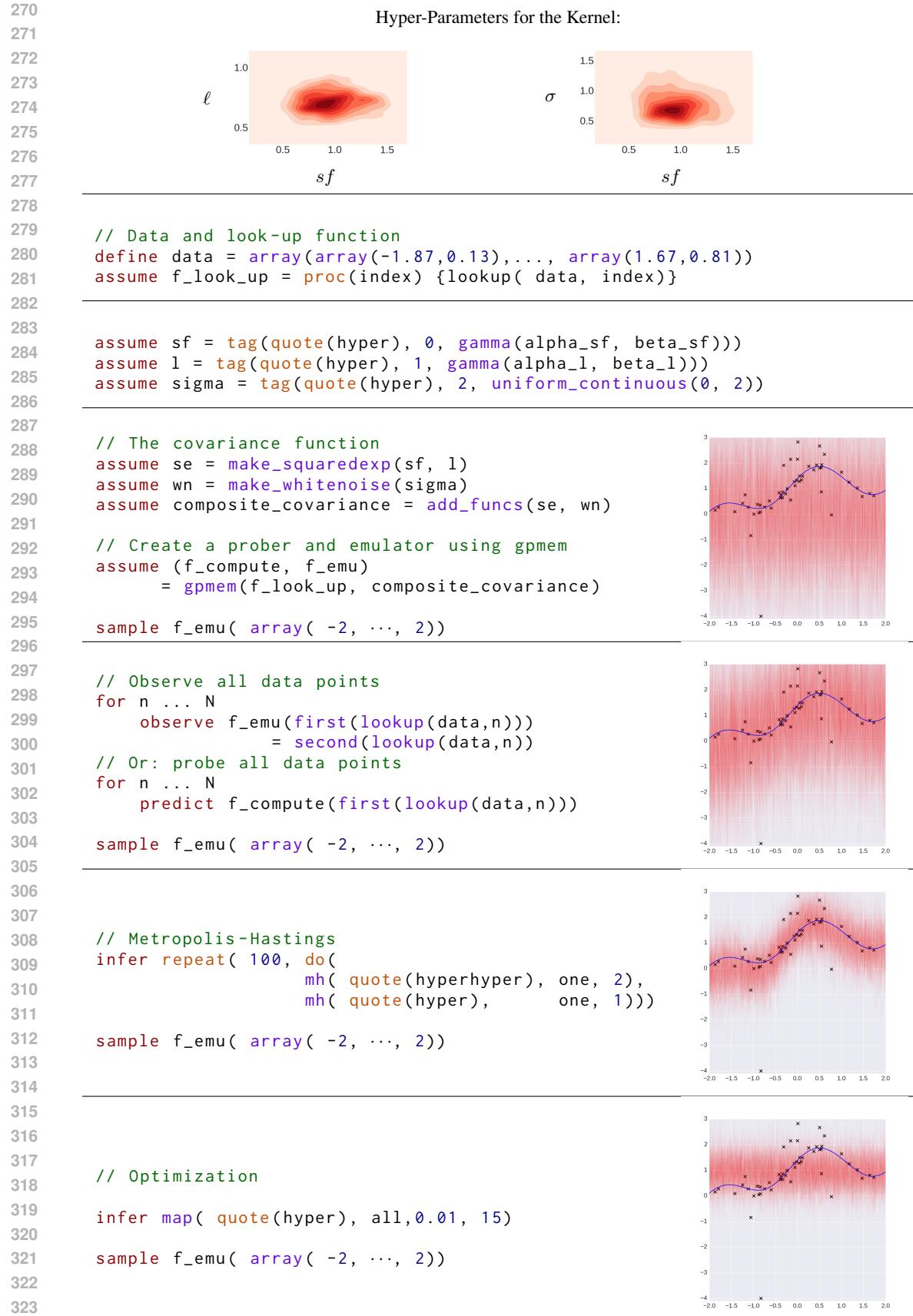
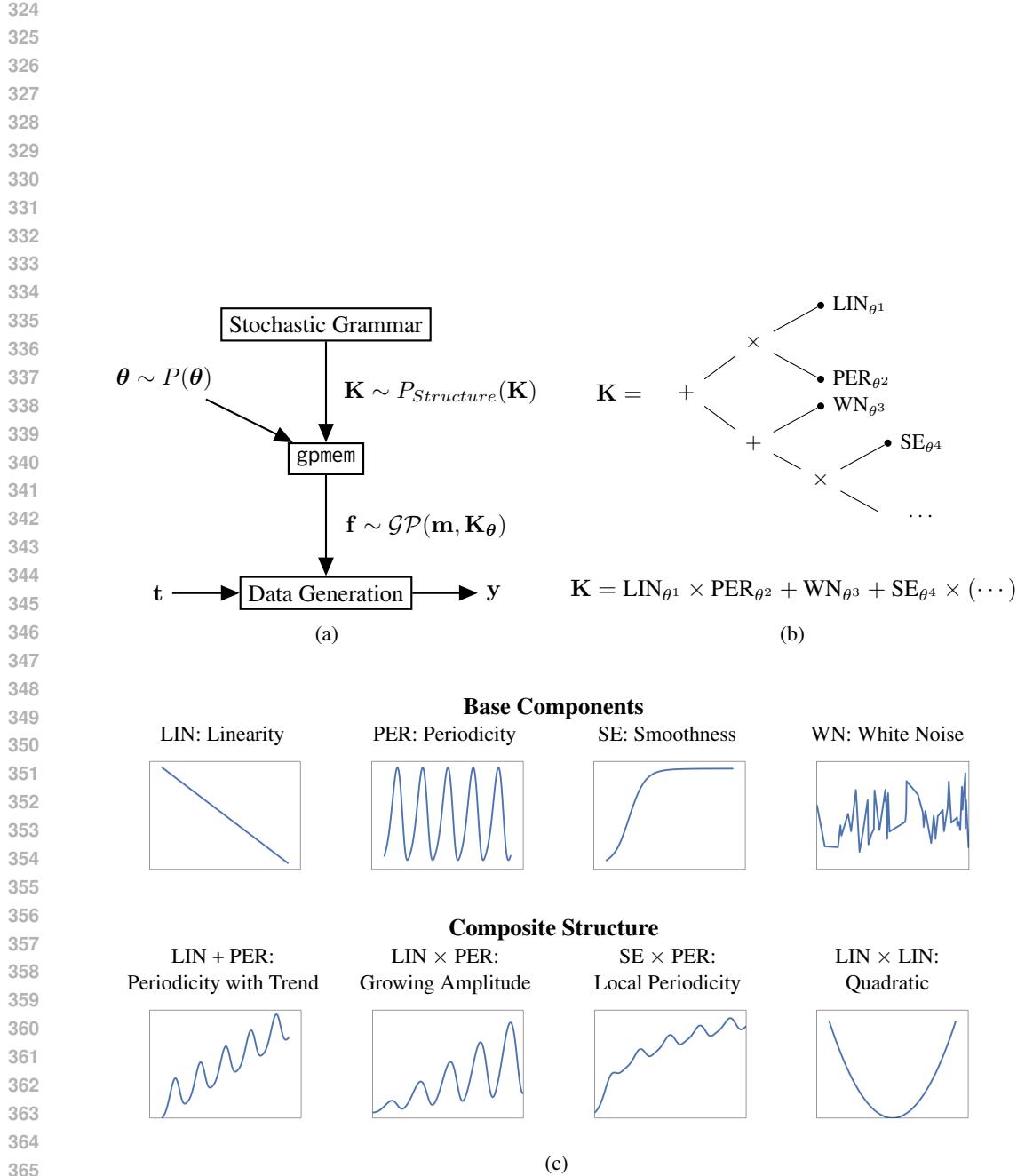


Figure 2: Regression with outliers and hierarchical prior structure.



366      Figure 3: (a) Graphical description of Bayesian GP structure learning. (b) Composite structure. (c)  
 367      The natural language interpretation of the structure.  
 368  
 369  
 370  
 371  
 372  
 373  
 374  
 375  
 376  
 377

378 We deploy a probabilistic context free grammar for our prior on structures. An input of non-  
 379 composite kernels (base kernels) is supplied to generate a posterior distributions of composite struc-  
 380 ture to express local and global aspects of the data.  
 381

382 We approximate the following intractable integrals of the expectation for the prediction:

$$383 \quad \mathbb{E}[y^* | x^*, \mathbf{D}, \mathbf{K}] = \iint f(x^*, \boldsymbol{\theta}, \mathbf{K}) P(\boldsymbol{\theta} | \mathbf{D}, \mathbf{K}) P(\mathbf{K} | \boldsymbol{\Omega}, s, n) d\boldsymbol{\theta} d\mathbf{K}. \quad (6)$$

386 This is done by sampling from the posterior probability distribution of the hyper-parameters and the  
 387 possible kernel:  
 388

$$389 \quad y^* \approx \frac{1}{T} \sum_{t=1}^T f(x^* | \boldsymbol{\theta}^{(t)}, \mathbf{K}^{(t)}). \quad (7)$$

391 In order to provide the sampling of the kernel, we introduce a stochastic process that simulates the  
 392 grammar for algebraic expressions of covariance function algebra:  
 393

$$394 \quad \mathbf{K}^{(t)} \sim P(\mathbf{K} | \boldsymbol{\Omega}, s, n) \quad (8)$$

395 Here, we start with the set of given base kernels and draw a random subset. For this subset of size  
 396  $n$ , we sample a set of possible operators  $\boldsymbol{\Omega}$  combining base kernels. The marginal probability of a  
 397 composite structure

$$398 \quad P(\mathbf{K} | \boldsymbol{\Omega}, s, n) = P(\boldsymbol{\Omega} | s, n) \times P(s | n) \times P(n), \quad (9)$$

400 is characterized by the prior  $P(n)$  on the number of base kernels used, the probability of a uniformly  
 401 chosen subset of the set of  $n$  possible covariance functions

$$403 \quad P(s | n) = \frac{n!}{|s|!}, \quad (10)$$

405 and the probability of sampling a global or a local structure, which is given by a binomial distribu-  
 406 tion:  
 407

$$408 \quad P(\boldsymbol{\Omega} | s, n) = \binom{n}{r} p_{+ \times}^k (1 - p_{+ \times})^{n-k}. \quad (11)$$

409 Many equivalent covariance structures can be sampled due to covariance function algebra and equiv-  
 410 alent representations with different parameterization (Lloyd et al., 2014). To inspect the posterior  
 411 of these equivalent structures we convert each kernel expression into a sum of products and subse-  
 412 quently simplify. All base kernels can be found in Appendix A, rules for this simplification can be  
 413 found in appendix B. The code for learning of kernel structure is as follows:  
 414

415  
 416  
 417  
 418  
 419  
 420  
 421  
 422  
 423  
 424  
 425  
 426  
 427  
 428  
 429  
 430  
 431

```

432
433 // GRAMMAR FOR KERNEL STRUCTURE
434 1 assume kernels = list(se, wn, lin, per, rq) // defined as above
435
436 // prior on the number of kernels
437 2 assume p_number_k = uniform_structure(n)
438 3 assume subset_kernels = tag(quote(grammar), 0,
439                                subset(kernels, p_number_k))
440
441 // kernel composition
442 5 assume composition = proc(l) {
443 6   if (size(l) <= 1)
444 7     { first(l) }
445 8   else { if (bernoulli())
446 9     { add_funcs(first(l), composition(rest(l))) }
44710     else { mult_funcs(first(l), composition(rest(l))) }
44811   }
44912 }
450
45113 assume K = tag(quote(grammar), 1, composition(subset_kernels))
452
45314 assume (f_compute f_emu) = gpmem(f_look_up, K)
454
455 // Probe all data points
45615 for n ... N
45716   predict f_compute(get_data_xs(n))
458
459 // PERFORMING INFERENCE
46017 infer repeat(2000, do(
46118   mh(quote(grammar), one, 1),
46219     for kernel in K:
46320       mh(quote(parameters_kernel), one, 1)))
464
465
466
467
468
469
470
```

462 We defined the space of covariance structures in a way that allows us to produce results coherent with  
463 work presented in Automatic Statistician. For example, for the airline data set describing monthly  
464 totals of international airline passengers (Box et al., 1997, according to Duvenaud et al., 2013). Our  
465 most frequent sample is identical with the highest scoring result reported in previous work using a  
466 search-and-score method (Duvenaud et al., 2013) for the CO<sub>2</sub> data set (see Rasmussen and Williams,  
467 2006 for a description) and the predictive capability is comparable. However, the components factor  
468 in a different way due to different parameterization of the individual base kernels. We see that the  
469 most probable alternatives for a structural description both recover the data dynamics (Fig. 9 for the  
470 airline data set).

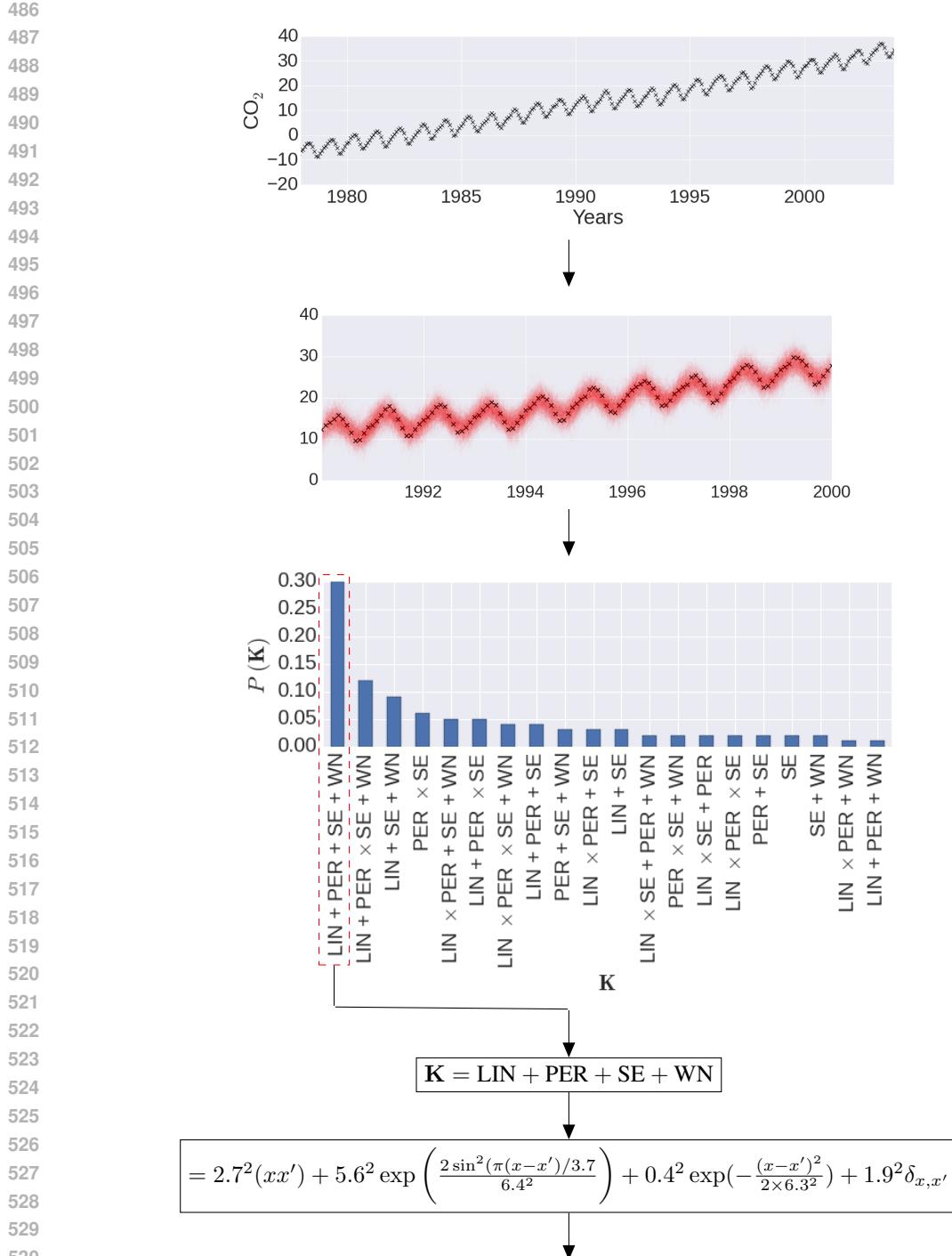
471 Confident about our results, we can now query the data for certain structures being present. We  
472 illustrate this using the Mauna Loa data used in previous work on automated kernel discovery (Du-  
473 venaud et al., 2013). We assume a relatively simple hypothesis space consisting of only four kernels,  
474 a linear, a smoothing, a periodic and a white noise kernel. In this experiment, we resort to the white  
475 noise kernel instead RQ (similar to (Lloyd et al., 2014)). We can now run the algorithm, compute a  
476 posterior of structures (see Fig. 4). We can also query this posterior distribution for the marginal of  
477 certain simple structures to occur. We demonstrate this in Fig. 5

478

### 479 3.3 Bayesian optimization

480

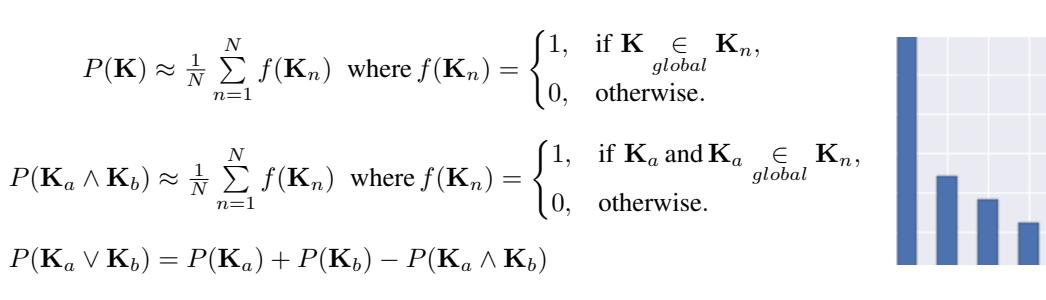
481 We introduce Thompson sampling, the basic solution strategy underlying the Bayesian optimization  
482 with gpmem. Thompson sampling (Thompson 1933) is a widely used Bayesian framework for ad-  
483 dressing the trade-off between exploration and exploitation in multi-armed (or continuum-armed)  
484 bandit problems. We cast the multi-armed bandit problem as a one-state Markov decision process,  
485 and describe how Thompson sampling can be used to choose actions for that Markov Decision  
Processes (MDP).



**Qualitative Interpretation:** The posterior peaks at a kernel structure with four additive components. Additive components hold globally, that is there are no higher level, qualitative aspects of the data that vary with the input space. The additive components are as follows: (i) a linearly increasing function or trend; (ii) a periodic function; (iii) a smooth function; and (iv) white noise.

Figure 4: Posterior of structure and qualitative, human interpretable reading. We take the raw data (top), compute a posterior distribution on structures (bar plot). We take the peak of this distribution ( $\text{LIN} + \text{PER} + \text{SE} + \text{WN}$ ) and show its human readable interpretation (left of bar plot). Below the bar plot show one sample of this structure with corresponding parameters. On the bottom, we sample from a GP with this kernel-hyper-parameter combination.

540  
541  
542  
543  
544  
545



555  
556  
557

What is the probability of a trend, a recurring pattern **and** noise in the data?

$$P((\text{LIN} \vee \text{LIN} \times \text{SE}) \wedge (\text{PER} \vee \text{PER} \times \text{SE} \vee \text{PER} \times \text{LIN}) \wedge (\text{WN} \vee \text{LIN} \times \text{WN})) = 0.36$$

560

561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574

Is there a trend?  
 $P(\text{LIN} \vee \text{LIN} \times \text{SE}) = 0.65$

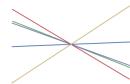


575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

Is there noise?  
 $P(\text{WN} \vee \text{LIN} \times \text{WN}) = 0.75$



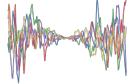
A linear trend?  
 $P(\text{LIN}) = 0.63$



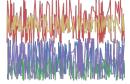
A smooth trend?  
 $P(\text{LIN} \times \text{SE}) = 0.02$



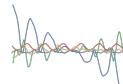
Heteroskedastic noise?  
 $P(\text{LIN} \times \text{WN}) = 0$



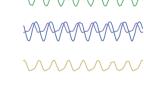
White noise?  
 $P(\text{WN}) = 0.75$



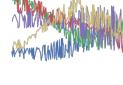
Is there repeating structure?  
 $P(\text{PER} \vee \text{PER} \times \text{SE} \vee \text{PER} \times \text{LIN}) = 0.73$



PER = 0.32



PER × SE = 0.34



PER × LIN = 0.07



Figure 5: We can query the data if some logical statements are probable to be true, for example, is it true that there is a trend?

594 Here, an agent is to take a sequence of actions  $a_1, a_2, \dots$  from a (possibly infinite) set of possible  
 595 actions  $\mathcal{A}$ . After each action, a reward  $r \in \mathbb{R}$  is received, according to an unknown conditional  
 596 distribution  $P_{\text{true}}(r | a)$ . The agent's goal is to maximize the total reward received for all actions in  
 597 an online manner. In Thompson sampling, the agent accomplishes this by placing a prior distribution  
 598  $P(\vartheta)$  on the possible "contexts"  $\vartheta \in \Theta$ . Here a context is a believed model of the conditional  
 599 distributions  $\{P(r | a)\}_{a \in \mathcal{A}}$ , or at least, a believed statistic of these conditional distributions which  
 600 is sufficient for deciding an action  $a$ . If actions are chosen so as to maximize expected reward, then  
 601 one such sufficient statistic is the believed conditional mean  $V(a | \vartheta) = \mathbb{E}[r | a; \vartheta]$ , which can be  
 602 viewed as a believed value function. For consistency with what follows, we will assume our context  
 603  $\vartheta$  takes the form  $(\theta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$  where  $\mathbf{a}_{\text{past}}$  is the vector of past actions,  $\mathbf{r}_{\text{past}}$  is the vector of their  
 604 rewards, and  $\theta$  (the "semicontext") contains any other information that is included in the context.  
 605

In this setup, Thompson sampling has the following steps:

---

**Algorithm 1** Thompson sampling.

---

609 Repeat as long as desired:

- 610 1. **Sample.** Sample a semicontext  $\theta \sim P(\theta)$ .
  - 611 2. **Search (and act).** Choose an action  $a \in \mathcal{A}$  which (approximately) maximizes  $V(a | \vartheta) =$   
 612      $\mathbb{E}[r | a; \vartheta] = \mathbb{E}[r | a; \theta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}]$ .
  - 613 3. **Update.** Let  $r_{\text{true}}$  be the reward received for action  $a$ . Update the believed distribution on  
 614      $\theta$ , i.e.,  $P(\theta) \leftarrow P_{\text{new}}(\theta)$  where  $P_{\text{new}}(\theta) = P(\theta | a \mapsto r_{\text{true}})$ .
- 

617 Note that when  $\mathbb{E}[r | a; \vartheta]$  (under the sampled value of  $\theta$  for some points  $a$ ) is far from the true value  
 618  $\mathbb{E}_{P_{\text{true}}}[r | a]$ , the chosen action  $a$  may be far from optimal, but the information gained by probing  
 619 action  $a$  will improve the belief  $\vartheta$ . This amounts to "exploration." When  $\mathbb{E}[r | a; \vartheta]$  is close to the  
 620 true value except at points  $a$  for which  $\mathbb{E}[r | a; \vartheta]$  is low, exploration will be less likely to occur, but  
 621 the chosen actions  $a$  will tend to receive high rewards. This amounts to "exploitation." The trade-off  
 622 between exploration and exploitation is illustrated in Figure 6. Roughly speaking, exploration will  
 623 happen until the context  $\vartheta$  is reasonably sure that the unexplored actions are probably not optimal,  
 624 at which time the Thompson sampler will exploit by choosing actions in regions it knows to have  
 625 high value.

626 Typically, when Thompson sampling is implemented, the search over contexts  $\vartheta \in \Theta$  is limited  
 627 by the choice of representation. In traditional programming environments,  $\theta$  often consists of a few  
 628 numerical parameters for a family of distributions of a fixed functional form. With work, a mixture of  
 629 a few functional forms is possible; but without probabilistic programming machinery, implementing  
 630 a rich context space  $\Theta$  would be an unworkably large technical burden. In a probabilistic programming  
 631 language, however, the representation of heterogeneously structured or infinite-dimensional context  
 632 spaces is quite natural. Any computable model of the conditional distributions  $\{P(r | a)\}_{a \in \mathcal{A}}$  can  
 633 be represented as a stochastic procedure  $(\lambda(a) \dots)$ . Thus, for computational Thompson sampling,  
 634 the most general context space  $\hat{\Theta}$  is the space of program texts. Any other context space  $\Theta$  has a  
 635 natural embedding as a subset of  $\hat{\Theta}$ .

**A Mathematical Specification**

637 Our mathematical specification assert the following properties:

- 640 • The regression function has a Gaussian process prior.
- 641 • The actions  $a_1, a_2, \dots \in \mathcal{A}$  are chosen by a Metropolis-like search strategy with Gaussian  
 642     drift proposals.
- 643 • The hyperparameters of the Gaussian process are inferred using Metropolis–Hastings sam-  
 644     pling after each action.

645 In this version of Thompson sampling, the contexts  $\vartheta$  are Gaussian processes over the action space  
 646  $\mathcal{A} = [-20, 20] \subseteq \mathbb{R}$ . That is,

$$V \sim \mathcal{GP}(\mu, K),$$

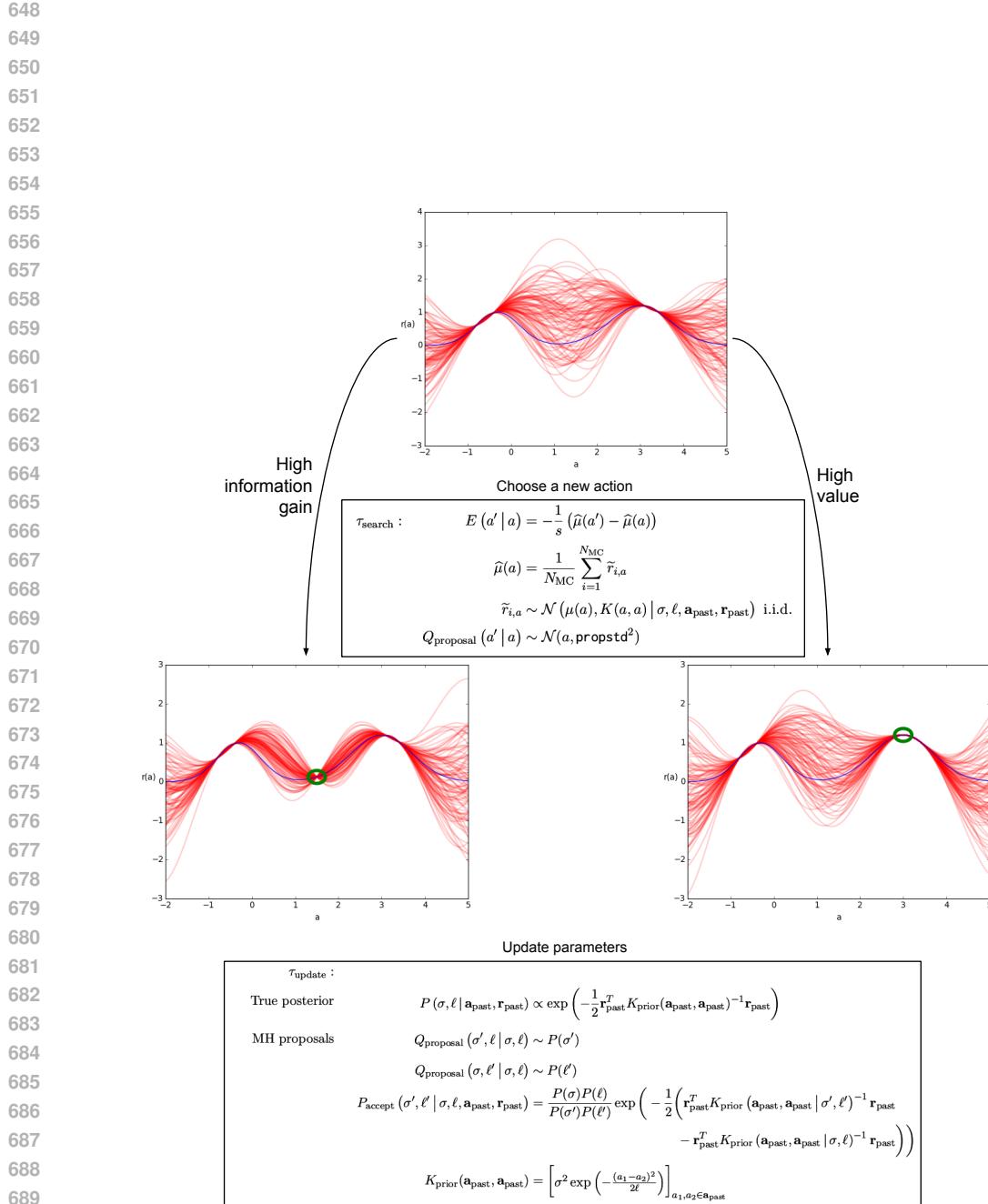


Figure 6: Two possible actions (in green) for an iteration of Thompson sampling. The believed distribution on the value function  $V$  is depicted in red. In this example, the true reward function is deterministic, and is drawn in blue. The action on the right receives a high reward, while the action on the left receives a low reward but greatly improves the accuracy of the believed distribution on  $V$ . The transition operators  $\tau_{\text{search}}$  and  $\tau_{\text{update}}$  are described in Section 3.3.

where the mean  $\mu$  is a computable function  $\mathcal{A} \rightarrow \mathbb{R}$  and the covariance  $K$  is a computable (symmetric, positive-semidefinite) function  $\mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}$ . This represents a Gaussian process  $\{R_a\}_{a \in \mathcal{A}}$ , where  $R_a$  represents the reward for action  $a$ . Computationally, we represent a context as a data structure

$$\vartheta = (\theta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}) = (\mu_{\text{prior}}, K_{\text{prior}}, \eta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}),$$

where  $\mu_{\text{prior}}$  is a procedure to be used as the prior mean function. w.l.o.g. we set  $\mu_{\text{prior}} \equiv 0$ .  $K_{\text{prior}}$  is a procedure to be used as the prior covariance function, parameterized by  $\eta$ .

The posterior mean and covariance for such a context  $\vartheta$  are gotten by the usual conditioning formulas (assuming, for ease of exposition as above, that the prior mean is zero):<sup>3</sup>

$$\begin{aligned} \mu(\mathbf{a}) &= \mu(\mathbf{a} | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}) \\ &= K_{\text{prior}}(\mathbf{a}, \mathbf{a}_{\text{past}}) K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}})^{-1} \mathbf{r}_{\text{past}} \\ K(\mathbf{a}, \mathbf{a}) &= K(\mathbf{a}, \mathbf{a} | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}) \\ &= K_{\text{prior}}(\mathbf{a}, \mathbf{a}) - K_{\text{prior}}(\mathbf{a}, \mathbf{a}_{\text{past}}) K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}})^{-1} K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}). \end{aligned}$$

Note that the context space  $\Theta$  is not a finite-dimensional parametric family, since the vectors  $\mathbf{a}_{\text{past}}$  and  $\mathbf{r}_{\text{past}}$  grow as more samples are taken.  $\Theta$  is, however, representable as a computational procedure together with parameters and past samples, as we do in the representation  $\vartheta = (\mu_{\text{prior}}, K_{\text{prior}}, \eta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$ .

We combine the Update and Sample steps of Algorithm 1 by running a Metropolis–Hastings (MH) sampler whose stationary distribution is the posterior  $P(\theta | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$ . The functional forms of  $\mu_{\text{prior}}$  and  $K_{\text{prior}}$  are fixed in our case, so inference is only done over the parameters  $\eta = \{\sigma, \ell\}$ ; hence we equivalently write  $P(\sigma, \ell | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$  for the stationary distribution. We make MH proposals to one variable at a time, using the prior as proposal distribution:

$$Q_{\text{proposal}}(\sigma', \ell' | \sigma, \ell) = P(\sigma')$$

and

$$Q_{\text{proposal}}(\sigma, \ell' | \sigma, \ell) = P(\ell').$$

The MH acceptance probability for such a proposal is

$$P_{\text{accept}}(\sigma', \ell' | \sigma, \ell) = \min \left\{ 1, \frac{Q_{\text{proposal}}(\sigma, \ell | \sigma', \ell')}{Q_{\text{proposal}}(\sigma', \ell' | \sigma, \ell)} \cdot \frac{P(\mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}} | \sigma', \ell')}{P(\mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}} | \sigma, \ell)} \right\}$$

Because the priors on  $\sigma$  and  $\ell$  are uniform in our case, the term involving  $Q_{\text{proposal}}$  equals 1 and we have simply

$$\begin{aligned} P_{\text{accept}}(\sigma', \ell' | \sigma, \ell) &= \min \left\{ 1, \frac{P(\mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}} | \sigma', \ell')}{P(\mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}} | \sigma, \ell)} \right\} \\ &= \min \left\{ 1, \exp \left( -\frac{1}{2} \left( \mathbf{r}_{\text{past}}^T K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}} | \sigma', \ell')^{-1} \mathbf{r}_{\text{past}} \right. \right. \right. \\ &\quad \left. \left. \left. - \mathbf{r}_{\text{past}}^T K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}} | \sigma, \ell)^{-1} \mathbf{r}_{\text{past}} \right) \right) \right\}. \end{aligned}$$

The proposal and acceptance/rejection process described above define a transition operator  $\tau_{\text{update}}$  which is iterated a specified number of times; the resulting state of the MH Markov chain is taken as the sampled semicontext  $\theta$  in Step 1 of Algorithm 1.

For Step 2 (Search) of Thompson sampling, we explore the action space using an MH-like transition operator  $\tau_{\text{search}}$ . As in MH, each iteration of  $\tau_{\text{search}}$  produces a proposal which is either accepted or rejected, and the state of this Markov chain after a specified number of steps is the new action  $a$ . The Markov chain's initial state is the most recent action, and the proposal distribution is Gaussian drift:

$$Q_{\text{proposal}}(a' | a) \sim \mathcal{N}(a, \text{propstd}^2),$$

---

<sup>3</sup>Here, for vectors  $\mathbf{a} = (a_i)_{i=1}^n$  and  $\mathbf{a}' = (a'_i)_{i=1}^{n'}$ ,  $\mu(\mathbf{a})$  denotes the vector  $(\mu(a_i))_{i=1}^n$  and  $K(\mathbf{a}, \mathbf{a}')$  denotes the matrix  $[K(a_i, a'_j)]_{1 \leq i \leq n, 1 \leq j \leq n'}$ .

756 where the drift width `propstd` is specified ahead of time. The acceptance probability of such a  
 757 proposal is

$$P_{\text{accept}}(a' | a) = \min \{1, \exp(-E(a' | a))\},$$

759 where the energy function  $E(\bullet | a)$  is given by a Monte Carlo estimate of the difference in value  
 760 from the current action:

$$E(a' | a) = -\frac{1}{s} (\hat{\mu}(a') - \hat{\mu}(a))$$

763 where

$$\hat{\mu}(a) = \frac{1}{N_{\text{avg}}} \sum_{i=1}^{N_{\text{avg}}} \tilde{r}_{i,a}$$

767 and

$$\tilde{r}_{i,a} \sim \mathcal{N}(\mu(a), K(a, a))$$

769 and  $\{\tilde{r}_{i,a}\}_{i=1}^{N_{\text{avg}}}$  are i.i.d. for a fixed  $a$ . Here the temperature parameter  $s \geq 0$  and the population size  
 770  $N_{\text{avg}}$  are specified ahead of time. Proposals of estimated value higher than that of the current action  
 771 are always accepted, while proposals of estimated value lower than that of the current action are  
 772 accepted with a probability that decays exponentially with respect to the difference in value. The  
 773 rate of the decay is determined by the temperature parameter  $s$ , where high temperature corresponds  
 774 to generous acceptance probabilities. For  $s = 0$ , all proposals of lower value are rejected; for  
 775  $s = \infty$ , all proposals are accepted. For points  $a$  at which the posterior mean  $\mu(a)$  is low but the  
 776 posterior variance  $K(a, a)$  is high, it is possible (especially when  $N_{\text{avg}}$  is small) to draw a “wild”  
 777 value of  $\hat{\mu}(a)$ , resulting in a favorable acceptance probability.

778 Indeed, taking an action  $a$  with low estimated value but high uncertainty serves the useful function  
 779 of improving the accuracy of the estimated value function at points near  $a$  (see Figure 6).<sup>4,5</sup> We see  
 780 a complete probabilistic program with `gpmem` implementing Bayesian optimization with Thompson  
 781 Sampling below (Listing 1).

782  
 783  
 784  
 785  
 786  
 787  
 788  
 789  
 790  
 791  
 792  
 793  
 794  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802  
 803  
 804  
 805  
 806  
 807

---

<sup>4</sup>At least, this is true when we use a smoothing prior covariance function such as the squared exponential.

<sup>5</sup>For this reason, we consider the sensitivity of  $\hat{\mu}$  to uncertainty to be a desirable property; indeed, this is why we use  $\hat{\mu}$  rather than the exact posterior mean  $\mu$ .

```

810
811
812
813
814 1 assume sf = tag(quote(hyper), 0, uniform_continuous(0, 10))
815 2 assume l = tag(quote(hyper), 1, uniform_continuous(0, 10))
816 3 assume se = make_squaredexp(sf, 1)
817 4 assume blackbox_f = get_bayesopt_blackbox()
818 5 assume (f_compute, f_emulate) = gpmem(blackbox_f, se)

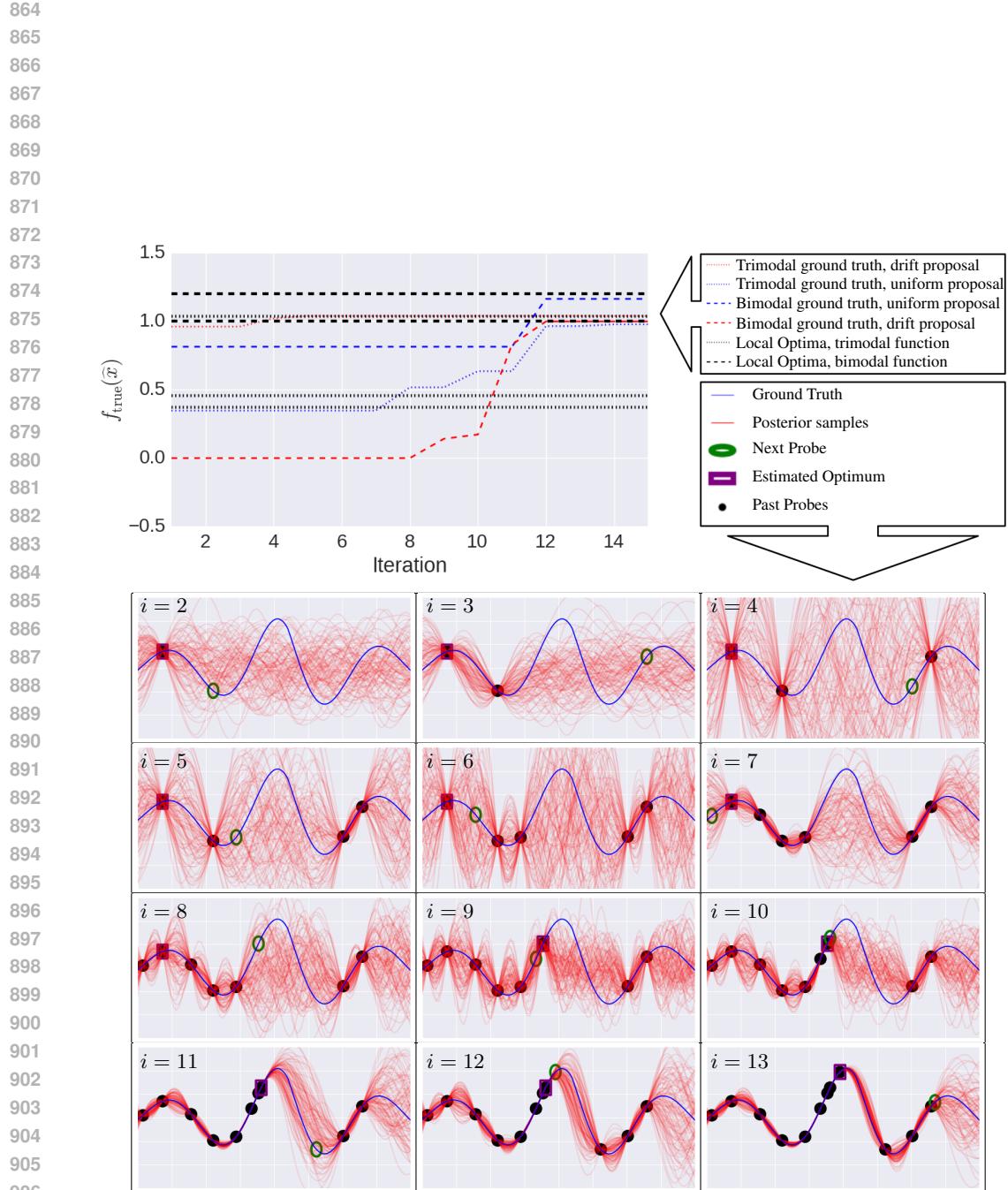
819 // A naive estimate of the argmax of the given function
820 define mc_argmax = proc(func) {
821 6   candidate_xs = mapv(proc(i) {uniform_continuous(-20, 20)},
822 7     arange(20));
823 8   candidate_ys = mapv(func, candidate_xs);
824 9   lookup(candidate_xs, argmax_of_array(candidate_ys))
825 10 };
826
827 // Shortcut to sample the emulator at a single point without packing
828 // and unpacking arrays
829 define emulate_pointwise = proc(x) {
830 12   run(sample(lookup(f_emulate(array(unquote(x))), 0)))
831 13 };
832
833 // Main inference loop
834 infer repeat(15, do(pass,
835 15   // Probe V at the point mc_argmax(emulate_pointwise)
836   predict(f_compute(unquote(mc_argmax(emulate_pointwise)))),
837   // Infer hyperparameters
838   mh(quote(hyper), one, 50)));
839
```

## 4 Discussion

We provided gpmem, an elegant linguistic framework for function learning-related tasks such as Bayesian optimization and GP kernel structure learning. We highlighted how gpmem overcomes shortcomings of the notations currently used in statistics, and how language constructs from programming allow the expression of models which would be cumbersome (prohibitively so, in some cases) to express in statistics notation. We evaluated our contribution on a range of hard problems for state-of-the-art Bayesian nonparametrics.

```

840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
```



907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Figure 7: Top: the estimated optimum over time. Blue and Red represent optimization with uniform and Gaussian drift proposals. Black lines indicate the local optima of the true functions. Bottom: a sequence of actions. Depicted are iterations 7-12 with uniform proposals.

918      **Appendix**

919  
920      **A Covariance Functions**  
921

922  
923       $SE = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$       (12)  
924

925       $LIN = \sigma^2(xx')$       (13)  
926

927       $C = \sigma^2$       (14)  
928

929       $WN = \sigma^2 \delta_{x,x'}$       (15)  
930

931       $RQ = \sigma^2 \left(1 + \frac{(x - x')^2}{2\alpha\ell^2}\right)^{-\alpha}$       (16)  
932

933       $PER = \sigma^2 \exp\left(\frac{2\sin^2(\pi(x - x')/p)}{\ell^2}\right).$       (17)

934      **B Covariance Simplification**  
935

936       $SE \times SE \rightarrow SE$   
937       $\{SE, PER, C, WN\} \times WN \rightarrow WN$   
938       $LIN + LIN \rightarrow LIN$   
939       $\{SE, PER, C, WN, LIN\} \times C \rightarrow \{SE, PER, C, WN, LIN\}$

940      Rule 1 is derived as follows:  
941

942       $\sigma_c^2 \exp\left(-\frac{(x - x')^2}{2\ell_c^2}\right) = \sigma_a^2 \exp\left(-\frac{(x - x')^2}{2\ell_a^2}\right) \times \sigma_b^2 \exp\left(-\frac{(x - x')^2}{2\ell_b^2}\right)$   
943  
944       $= \sigma_c^2 \exp\left(-\frac{(x - x')^2}{2\ell_a^2}\right) \times \exp\left(-\frac{(x - x')^2}{2\ell_b^2}\right)$       (18)  
945  
946       $= \sigma_c^2 \exp\left(-\frac{(x - x')^2}{2\ell_a^2} - \frac{(x - x')^2}{2\ell_b^2}\right)$   
947  
948       $= \sigma_c^2 \exp\left(-\frac{(x - x')^2}{2\ell_c^2}\right)$   
949  
950

951      For stationary kernels that only depend on the lag vector between  $x$  and  $x'$  it holds that multiplying  
952      such a kernel with a WN kernel we get another WN kernel (Rule 2). Take for example the SE kernel:  
953

954       $\sigma_a^2 \exp\left(-\frac{(x - x')^2}{2\ell_c^2}\right) \times \sigma_b \delta_{x,x'} = \sigma_a \sigma_b \delta_{x,x'}$       (19)  
955  
956

957      Rule 3 is derived as follows:  
958

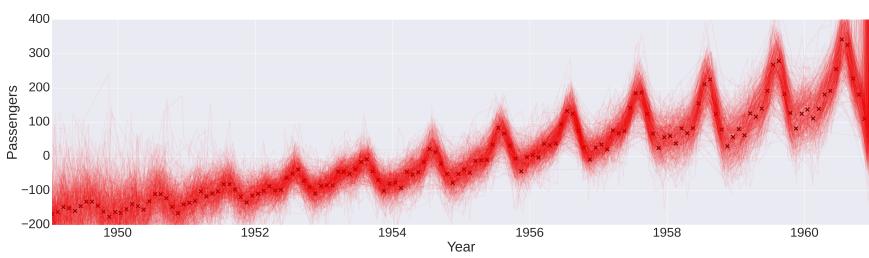
959       $\theta_c(x \times x') = \theta_a(x \times x') + \theta_b(x \times x')$       (20)  
960

961      Multiplying any kernel with a constant obviously changes only the scale parameter of a kernel (Rule  
962      4).  
963

964      **C Additional Structure Learning Results**  
965

966      Below we depict additional results for the airline data set using a hypothesis space of LIN, PER, SE  
967      and RQ covariance functions as base kernels for the composition. Fig 8 shows the training data and  
968      posterior samples drawn from the GP with posterior composite structure.  
969

970      Fig. 9 shows why it advantageous to be Bayesian here. We see two possible hypotheses for the  
971      composite structure of  $\mathbf{K}$ .  
972



972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
Figure 8: Data (black x) and posterior samples (red) for the airline data set.



Figure 9: We see two possible hypotheses for the composite structure of  $\mathbf{K}$ . (a) Most frequent sample drawn from the posterior on structure. We have found two global components. First, a smooth trend ( $\text{LIN} \times \text{SE}$ ) with a non-linear increasing slope. Second, a periodic component with increasing variation and noise. (b) Second most frequent sample drawn from the posterior on structure. We found one global component. It is comprised of local changes that are periodic and with changing variation.

## References

- Barry, D. (1986). Nonparametric bayesian regression. *The Annals of Statistics*, 14(3):934–953.
- Box, G. E. P., Jenkins, G. M., and Reinsel, G. C. (1997). *Time series analysis: forecasting and control*.
- Damianou, A. and Lawrence, N. (2013). Deep gaussian processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 207–215.
- Duvenaud, D., Lloyd, J. R., Grosse, R., Tenenbaum, J., and Ghahramani, Z. (2013). Structure discovery in nonparametric regression through compositional kernel search. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1166–1174.
- Ferris, B., Haehnel, D., and Fox, D. (2006). Gaussian processes for signal strength-based location estimation. In *Proceedings of the Conference on Robotics Science and Systems*. Citeseer.
- Gelbart, M. A., Snoek, J., and Adams, R. P. (2014). Bayesian optimization with unknown constraints. *arXiv preprint arXiv:1403.5607*.
- Goodman, N. D .and Mansinghka, V. K., Roy, D., Bonawitz, K., and Tenenbaum, J. (2008). Church: A language for generative models. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 220–229.
- Kemmler, M., Rodner, E., Wacker, E., and Denzler, J. (2013). One-class classification with gaussian processes. *Pattern Recognition*, 46(12):3507–3518.
- Kennedy, M. C. and O'Hagan, A. (2001). Bayesian calibration of computer models. *Journal of the Royal Statistical Society. Series B, Statistical Methodology*, pages 425–464.
- Kuss, M. and Rasmussen, C. E. (2005). Assessing approximate inference for binary gaussian process classification. *The Journal of Machine Learning Research*, 6:1679–1704.
- Kwan, J., Bhattacharya, S., Heitmann, K., and Habib, S. (2013). Cosmic emulation: The concentration-mass relation for wcdm universes. *The Astrophysical Journal*, 768(2):123.
- Lloyd, J. R., Duvenaud, D., Grosse, R., Tenenbaum, J., and Ghahramani, Z. (2014). Automatic construction and natural-language description of nonparametric regression models. In *Proceedings of the Conference on Artificial Intelligence (AAAI)*.
- Mansinghka, V. K., Selsam, D., and Perov, Y. (2014). Venture: a higher-order probabilistic programming platform with programmable inference. *arXiv preprint arXiv:1404.0099*.

- 1026 McAllester, D., Milch, B., and Goodman, N. D. (2008). Random-world semantics and syntactic  
1027 independence for expressive languages. Technical report.  
1028
- 1029 Neal, R. M. (1995). *Bayesian Learning for Neural Networks*. PhD thesis, University of Toronto.  
1030
- 1031 Neal, R. M. (1997). Monte carlo implementation of gaussian process models for bayesian regression  
and classification. *arXiv preprint physics/9701026*.
- 1032 Poole, D. (1993). Probabilistic horn abduction and bayesian networks. *Artificial Intelligence*,  
1033 64(1):81–129.
- 1034 Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning (Adap-*  
1035 *tive Computation and Machine Learning)*. The MIT Press.
- 1036 Sato, T. (1995). A statistical learning method for logic programs with distribution semantics. In *In*  
1037 *Proceedings of the International Conference on Logic Programming*. Citeseer.
- 1038 Schneider, M. D., Knox, L., Habib, S., Heitmann, K., Higdon, D., and Nakhleh, C. (2008). Simula-  
1039 tions and cosmological inference: A statistical model for power spectra means and covariances.  
1040 *Physical Review D*, 78(6):063529.
- 1041 Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical bayesian optimization of machine  
1042 learning algorithms. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2951–  
1043 2959.
- 1044 Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view  
1045 of the evidence of two samples. *Biometrika*, pages 285–294.
- 1046 Williams, C. K. I. and Barber, D. (1998). Bayesian classification with gaussian processes. *IEEE*  
1047 *Transactions on Pattern Analysis and Machine Intelligence*, 20(12):1342–1351.
- 1048
- 1049
- 1050
- 1051
- 1052
- 1053
- 1054
- 1055
- 1056
- 1057
- 1058
- 1059
- 1060
- 1061
- 1062
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079