
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 results (Neal, 1997) on outlier rejection via hyper-parameter inference; (ii) a fully Bayesian extension to the Automated Statis-

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
tician project; and (iii) an implementation of Bayesian optimization via Thompson sampling. The first application can in principle be replicated in several other probabilistic languages embedding the proposal that is described in this paper. The remaining two applications rely on distinctive capabilities of Venture: support for fully Bayesian structure learning and language constructs for inference programming. All applications share a single 50-line Python library and require fewer than 20 lines of probabilistic code each.

2 Background on Gaussian Processes

Let the data be pairs of real-valued scalars $\{(x_i, y_i)\}_{i=1}^n$ (complete data will be denoted by column vectors \mathbf{x}, \mathbf{y}). GPs present a non-parametric way to express prior knowledge on the space of possible functions f modeling a regression relationship. Formally, a GP is an infinite-dimensional extension of the multivariate Gaussian distribution. The joint prior distribution of training and test output is defined as

$$\begin{bmatrix} \mathbf{f} \\ \hat{\mathbf{f}} \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} K(\mathbf{x}, \mathbf{x}) & K(\mathbf{x}, \hat{\mathbf{x}}) \\ K(\hat{\mathbf{x}}, \mathbf{x}) & K(\hat{\mathbf{x}}, \hat{\mathbf{x}}) \end{bmatrix}\right) \quad (1)$$

We sample from the predictive posterior by applying Bayes rule to the multivariate Gaussian that is determined by the GP:

$$\hat{\mathbf{f}} | \hat{\mathbf{x}}, \mathbf{x}, \mathbf{f} \sim \mathcal{N}(K(\hat{\mathbf{x}}, \mathbf{x})K(\mathbf{x}, \mathbf{x})^{-1}\mathbf{f}, K(\hat{\mathbf{x}}, \hat{\mathbf{x}}) - K(\hat{\mathbf{x}}, \mathbf{x})K(\mathbf{x}, \mathbf{x})^{-1}K(\mathbf{x}, \hat{\mathbf{x}})) \quad (2)$$

$$= \mathcal{N}(\hat{\mu}, \hat{\mathbf{K}}) \quad (3)$$

Often one assumes the values \mathbf{y} are noisily measured, that is, one only sees the values of $\mathbf{y}_{\text{noisy}} = \mathbf{f} + \mathbf{w}$ where \mathbf{w} is Gaussian white noise with variance σ_{noise}^2 . This noise term can be absorbed into the covariance matrix $\mathbf{K}(\mathbf{x}, \mathbf{x})$ which in the following, we will write as \mathbf{K} for readability. The log-likelihood of a GP can then be written as:

$$\log P(\mathbf{f} | \mathbf{x}) = -\frac{1}{2}\mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{n}{2} \log 2\pi \quad (4)$$

where n is the number of data points. Both log-likelihood and predictive posterior can be computed efficiently in a Venture SP with an algorithm that resorts to Cholesky factorization(Rasmussen and Williams, 2006, chap. 2) where we compute (4) as

$$\log(P(\mathbf{f} | \mathbf{x})) := -\frac{1}{2}\mathbf{f}^\top \boldsymbol{\alpha} - \sum_i \log \mathbf{L}_{ii} - \frac{n}{2} \log 2\pi \quad (5)$$

where

$$\mathbf{L} := \text{cholesky}(\mathbf{K}) \quad (6)$$

and

$$\boldsymbol{\alpha} := \mathbf{L}^\top \backslash (\mathbf{L} \backslash \mathbf{f}). \quad (7)$$

This results in a computational complexity of $\mathcal{O}(n^3)$ in the number of data points for sampling with a complexity of $n^3/6$ for (6) an $n^2/2$ for (7).

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

Venture includes the primitive `make_gp`, which takes as arguments a unary function `mean` and a binary (symmetric, positive-semidefinite) function `cov` and produces a function `g` distributed as a Gaussian process with the supplied mean and covariance. For example, a function `g ~ GP(0, SE)`, where `SE` is a squared-exponential covariance

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

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

```

108 assume zero = make_const_func( 0.0)
109 assume se = make_squaredexp( 1.0, 1.0)
110 assume g = make_gp( zero, se)
111

```

112 There are two ways to view g as a “random function.” In the first view, the `assume` directive that
 113 instantiates g does not use any randomness—only the subsequent calls to g do—and coherence con-
 114 straints are upheld by the interpreter by keeping track of which evaluations of g exist in the current
 115 execution trace. Namely, if the current trace contains evaluations of g at the points x_1, \dots, x_N with
 116 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})
 117 will be distributed according to the joint conditional distribution

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

118 In the second view, g is a randomly chosen deterministic function, chosen from the space of all
 119 deterministic real-valued functions; in this view, the `assume` directive contains *all* the randomness,
 120 and subsequent invocations of g are deterministic. The first view is procedural and is faithful to the
 121 computation that occurs behind the scenes in Venture. The second view is declarative and is faithful
 122 to notations like “ $g \sim P(g)$ ” which are often used in mathematical treatments. Because a model
 123 program could make arbitrarily many calls to g , and the joint distribution on the return values of the
 124 calls could have arbitrarily high entropy, it is not computationally possible in finite time to choose
 125 the entire function g all at once as in the second view. Thus, it stands to reason that any computa-
 126 tionally implementable notion of “nonparametric random functions” must involve incremental random
 127 choices in one way or another, and Gaussian processes in Venture are no exception.

129 2.1 The `gpmem` construct

131 Memoization is the practice of storing previously computed values of a function so that future calls
 132 with the same inputs can be evaluated by lookup rather than recomputation. Although memoiza-
 133 tion does not change the semantics of a deterministic program, it does change that of a stochastic
 134 program (Goodman et al., 2008). The authors provide an intuitive example: let f be a function
 135 that flips a coin and return “head” or “tails”. The probability that two calls of f are equivalent is
 136 0.5. However, if the function call is memoized, it is 1. In fact, there is an infinite range of possible
 137 caching policies (specifications of when to use a stored value and when to recompute), each poten-
 138 tially having a different semantics. Any particular caching policy can be understood by random
 139 world semantics (Poole, 1993; Sato, 1995) over the stochastic program: each possible world corre-
 140 sponds to a mapping from function input sequence to function output sequence (McAllester et al.,
 141 2008). In Venture, these possible worlds are first-class objects, and correspond to the *probabilistic*
 142 *execution traces* (*traces* (Mansinghka et al., 2014)).

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

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

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

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

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

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

```

156
157
158
159
160
161      assume K = make_squaredexp (sf, 1)
162      assume (f_compute f_emu) = gpmem( f, K))

```

```

162 Adding a single line to the program, such as
163
164
165     infer mh( quote( parameters), one, 50),
166
167 allows us perform Bayesian inference over the parameters  $s_f$  and  $l$ . Such inference can be performed
168 without the refactoring and elaboration needed if one would describe the above declaratively as in
169 traditional statistics notation. In contrast to traditional statistics notation, probabilistic programs
170 are written procedurally; reasoning declaratively about the dynamics of a fundamentally procedural
171 inference algorithm is often unwieldy due to the absence of programming constructs such as loops
172 and mutable state.
173
174 We implement gpmem by memoizing a target procedure in a wrapper that remembers previously
175 computed values. This comes with interesting implications: from the standpoint of computation,
176 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
177 restricted to only allow evaluation at a specific set of inputs  $x$ .
178
179 Modelling the data set with a GP then amounts to trying to learn a smooth function  $f_{\text{emu}}$  (“emu”
180 stands for “emulator”) which extends  $f$  to its full domain. We can then the incorporate observations
181 in two different ways: we either are either told that the at a point  $x$  the value is  $y$ :
182
183
184     observe f_emu ( x ) = y
185
186 Or we express this as
187
188
189     predict f_compute ( x )
190
191 The second expression has at least two benefits: (i) readability (in some cases), and (ii) amenability
192 to active learning. As to (i), the statistical code of creating a Gaussian process is replaced with a
193 memoization-like idiom, which will be more familiar to programmers. As to (ii), when using gpmem,
194 it is quite easy to decide incrementally which data point to sample next: for example, the loop from
195  $x[1]$  to  $x[n]$  could be replaced by a loop in which the next index  $i$  is chosen by a supplied decision
196 rule. In this way, we could use gpmem to perform online learning using only a subset of the available
197 data.
198
199 We illustrate the use of gpmem in this context in a tutorial in Fig. 1.
200
201
202
203 

### 3 Applications


204
205 gpmem is an elegant linguistic framework for function learning-related tasks. The technique allows
206 language constructs from programming to help express models which would be cumbersome to
207 express in statistics notation. We will now illustrate this with three example applications.
208
209
210 

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


211
212 The probability of the hyper-parameters of a GP as defined above and given covariance function
213 structure  $\mathbf{K}$  is:
214

$$P(\boldsymbol{\theta} | \mathbf{D}, \mathbf{K}) = \frac{P(\mathbf{D} | \boldsymbol{\theta}, \mathbf{K})P(\boldsymbol{\theta} | \mathbf{K})}{P(\mathbf{D} | \mathbf{K})}. \quad (8)$$

215
216 Let the  $\mathbf{K}$  be the sum of a smoothing and a white noise (WN) kernel. For this case, Neal (1997)
217 suggested the problem of outliers in data as a use-case for a hierarchical Bayesian treatment of
218 Gaussian processes1. The work suggests a hierarchical system of hyper-parameterization. Here, we
219 draw hyper-parameters from a  $\Gamma$  distributions:
220

$$\ell^{(t)} \sim \Gamma(\alpha_1, \beta_1), \sigma^{(t)} \sim \Gamma(\alpha_2, \beta_2) \quad (9)$$

221


---

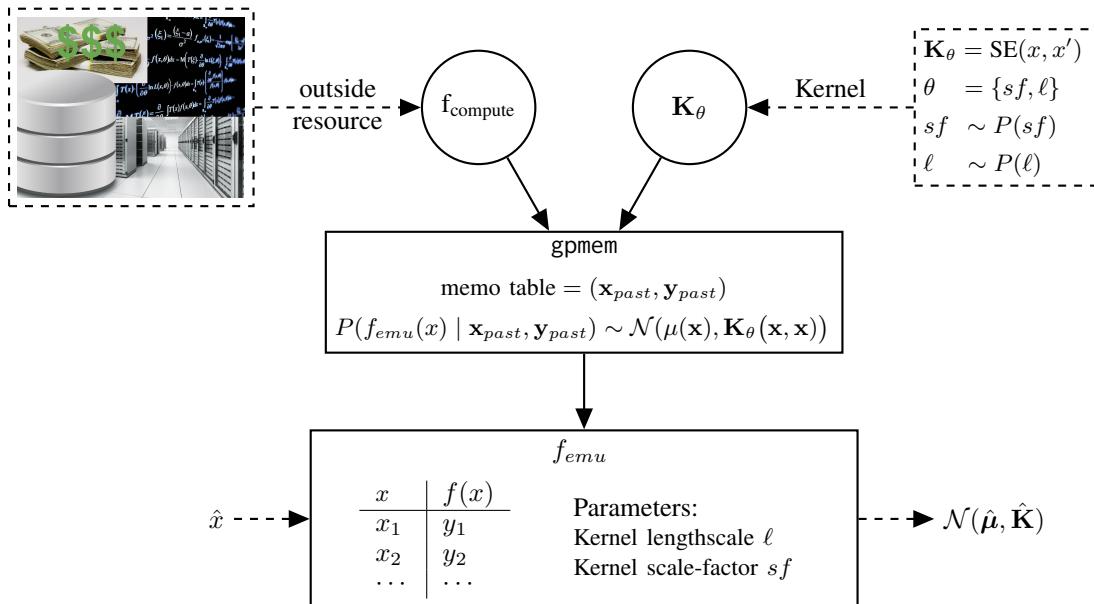

222 1In (Neal, 1997) the sum of an SE plus a constant kernel is used. We keep the WN kernel for illustrative
223 purposes.

```

```

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))

```

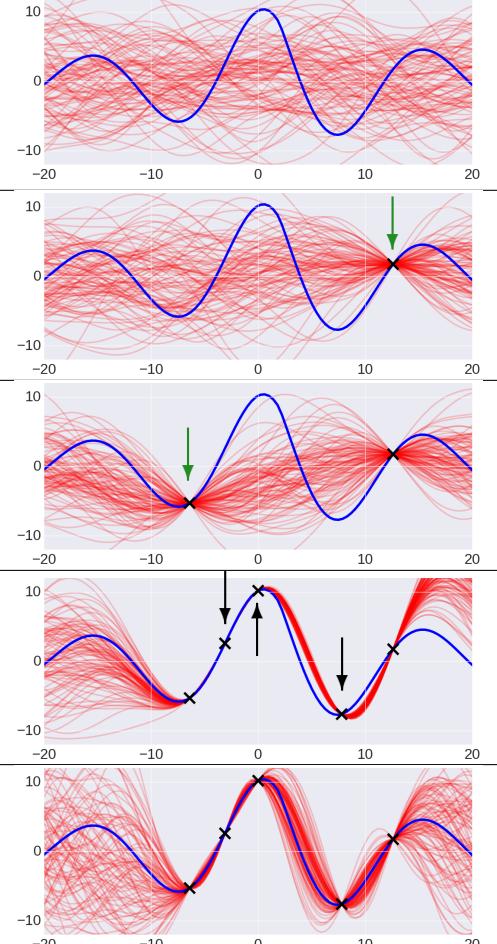


Figure 1: gpmem tutorial. The top shows a schematic of gpmem. f_{compute} probes an outside resource. This can be expensive (top left). Every probe is memoized and improves the GP-based emulator. Below the schematic we see a movie of the evolution of gpmem's state of believe of the world given certain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

and in turn sample the α and β from Γ distributions as well:

$$\alpha_1^{(t)} \sim \Gamma(\alpha_\alpha^1, \beta_\alpha^1), \alpha_2^{(t)} \sim \Gamma(\alpha_\alpha^2, \beta_\alpha^2), \dots \quad (10)$$

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

$$f(x) = 0.3 + 0.4x + 0.5 \sin(2.7x) + \frac{1.1}{(1+x^2)} + \eta \quad \text{with } \eta \sim \mathcal{N}(0, \sigma) \quad (11)$$

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

3.2 Discovering qualitative structure from time series data

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

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

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

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

$$\mathbb{E}[y^* | x^*, \mathbf{D}, \mathbf{K}] = \iint f(x^*, \boldsymbol{\theta}, \mathbf{K}) P(\boldsymbol{\theta} | \mathbf{D}, \mathbf{K}) P(\mathbf{K} | \boldsymbol{\Omega}, s, n) d\boldsymbol{\theta} d\mathbf{K}. \quad (12)$$

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

$$y^* \approx \frac{1}{T} \sum_{t=1}^T f(x^* | \boldsymbol{\theta}^{(t)}, \mathbf{K}^{(t)}). \quad (13)$$

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

$$\mathbf{K}^{(t)} \sim P(\mathbf{K} | \boldsymbol{\Omega}, s, n) \quad (14)$$

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

$$P(\mathbf{K} | \boldsymbol{\Omega}, s, n) = P(\boldsymbol{\Omega} | s, n) \times P(s | n) \times P(n), \quad (15)$$

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

$$P(s | n) = \frac{n!}{|s|!}, \quad (16)$$

and the probability of sampling a global or a local structure, which is given by a binomial distribution:

$$P(\boldsymbol{\Omega} | s, n) = \binom{n}{r} p_{+ \times}^k (1 - p_{+ \times})^{n-k}. \quad (17)$$

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

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

```

324
325
326
327
328
329
330
331
332
333 // Data and look-up function
334 define data = array(array(-1.87,0.13),..., array(1.67,0.81))
335 assume f_look_up = proc(index) {lookup( data, index)}
336
337
338 assume sf = tag(quote(hyper), 0, gamma(alpha_sf, beta_sf)))
339 assume l = tag(quote(hyper), 1, gamma(alpha_l, beta_l)))
340 assume sigma = tag(quote(hyper), 2, uniform_continuous(0, 2))
341
342 // The covariance function
343 assume se = make_squaredexp(sf, l)
344 assume wn = make_whitenoise(sigma)
345 assume composite_covariance = add_funcs(se, wn)
346
347 // Create a prober and emulator using gpmem
348 assume (f_compute, f_emu)
349 = gpmem(f_look_up, composite_covariance)
350
351 sample f_emu( array( -2, ..., 2))
352
353 // Observe all data points
354 for n ... N
355 observe f_emu(first(lookup(data,n)))
356 = second(lookup(data,n))
357 // Or: probe all data points
358 for n ... N
359 predict f_compute(first(lookup(data,n)))
360
361 sample f_emu( array( -2, ..., 2))
362
363 // Metropolis-Hastings
364 infer repeat( 100, do(
365 mh( quote(hyperhyper), one, 2),
366 mh( quote(hyper), one, 1)))
367
368 sample f_emu( array( -2, ..., 2))
369
370
371 // Optimization
372
373 infer map( quote(hyper), all, 0.01, 15)
374
375 sample f_emu( array( -2, ..., 2))
376
377
```

Hyper-Parameters for the Kernel:

Figure 2: Regression with outliers and hierarchical prior structure.

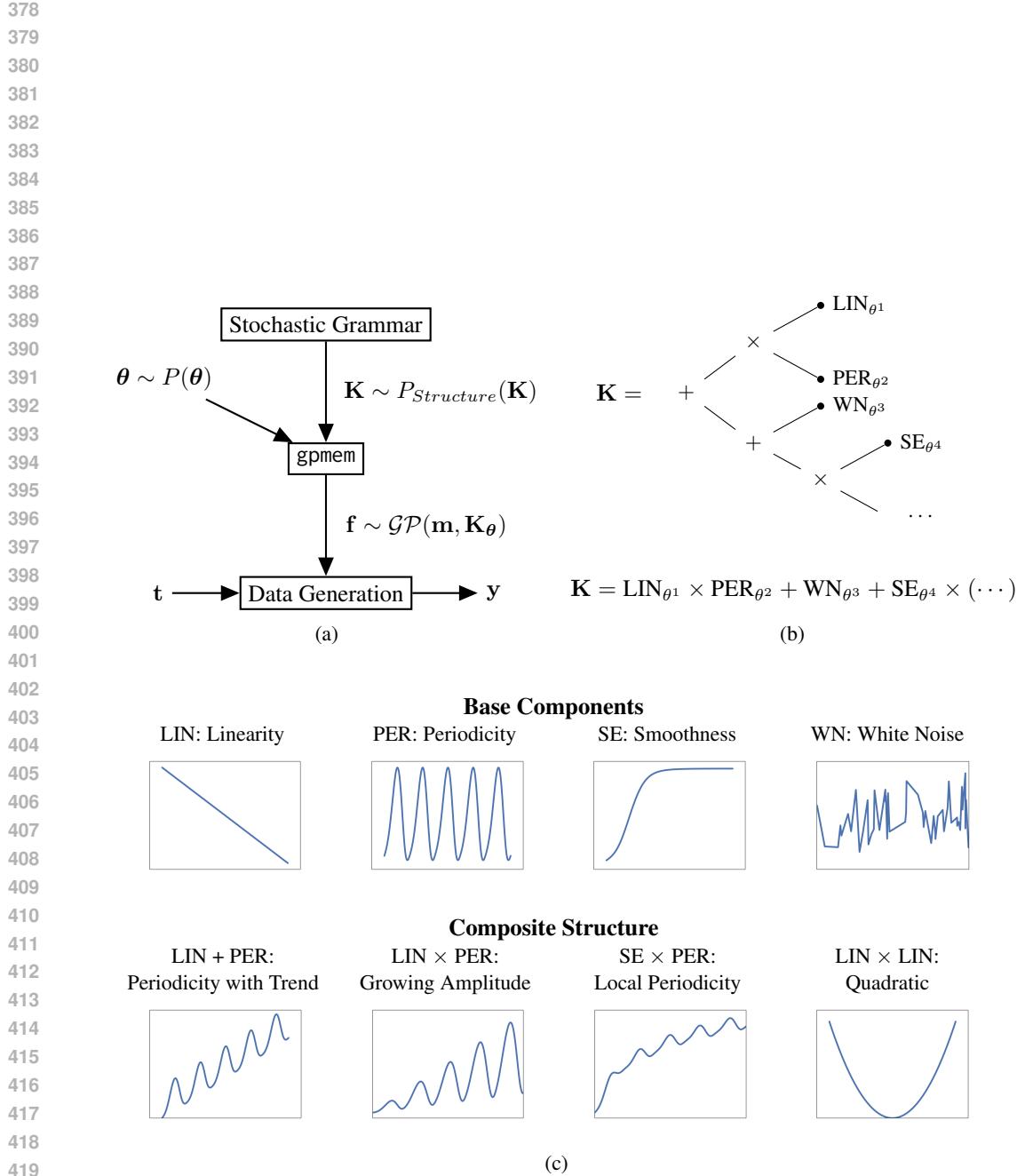


Figure 3: (a) Graphical description of Bayesian GP structure learning. (b) Composite structure. (c) The natural language interpretation of the structure.

```

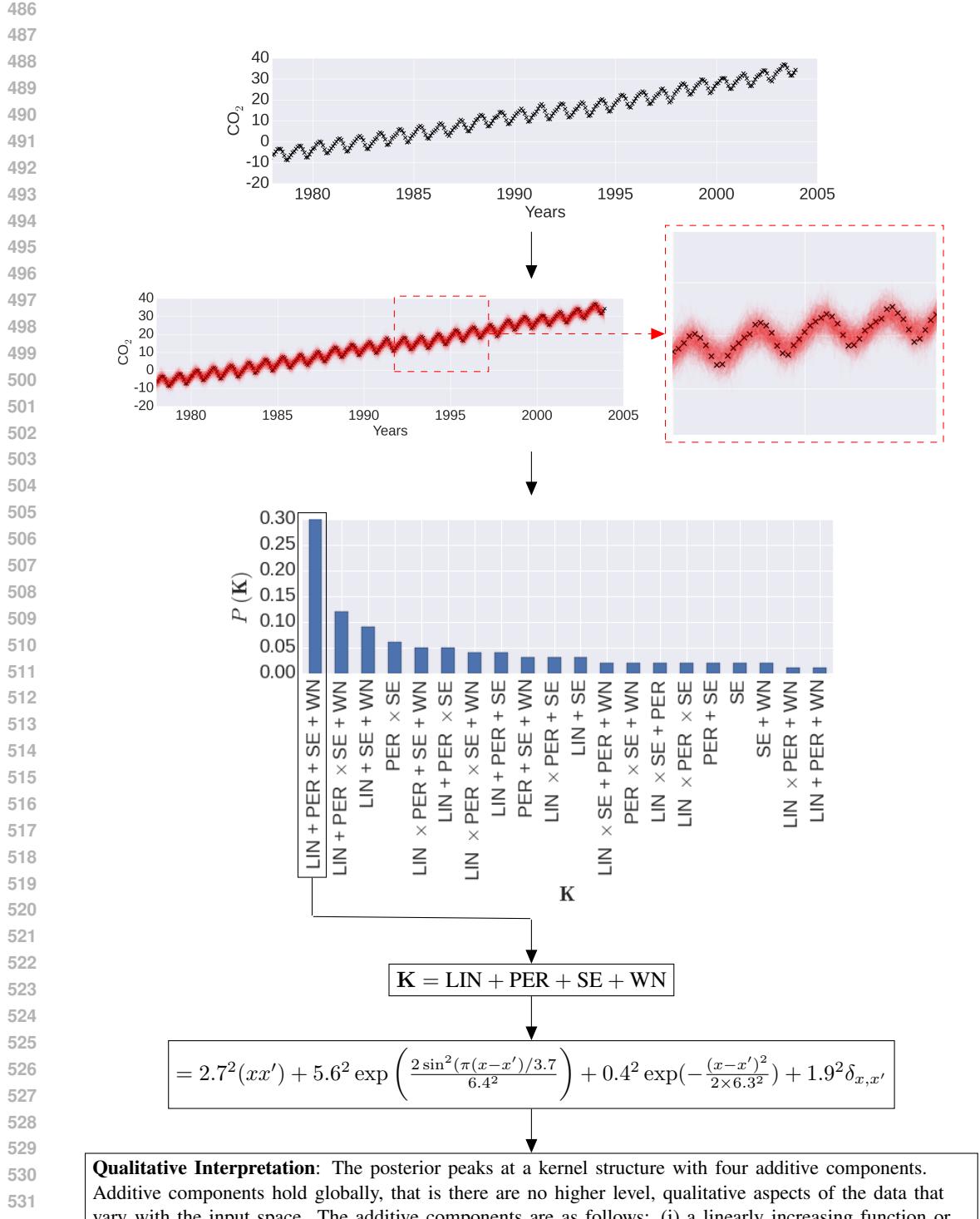
432
433 // GRAMMAR FOR KERNEL STRUCTURE
434 1 assume kernels = list(se, wn, lin, per, rq) // defined as above
435
436 // prior on the number of kernels
437 2 assume p_number_k = uniform_structure(n)
438 3 assume subset_kernels = tag(quote(grammar), 0,
439                                subset(kernels, 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
```

We defined the space of covariance structures in a way that allows us to produce results coherent with work presented in Automatic Statistician. For example, for the airline data set describing monthly totals of international airline passengers (Box et al., 1997, according to Duvenaud et al., 2013). Our most frequent sample is identical with the highest scoring result reported in previous work using a search-and-score method (Duvenaud et al., 2013) for the CO₂ data set (see Rasmussen and Williams, 2006 for a description) and the predictive capability is comparable. However, the components factor in a different way due to different parameterization of the individual base kernels. We see that the most probable alternatives for a structural description both recover the data dynamics (Fig. ?? for the airline data set).

Confident about our results, we can now query the data for certain structures being present. We illustrate this using the Mauna Loa data used in previous work on automated kernel discovery (Duvenaud et al., 2013). We assume a relatively simple hypothesis space consisting of only four kernels, a linear, a smoothing, a periodic and a white noise kernel. In this experiment, we resort to the white noise kernel instead RQ (similar to (Lloyd et al., 2014)). We can now run the algorithm, compute a posterior of structures (see Fig. 4 and 5). We can also query this posterior distribution for the marginal of certain simple structures to occur. We demonstrate this in Fig. 6

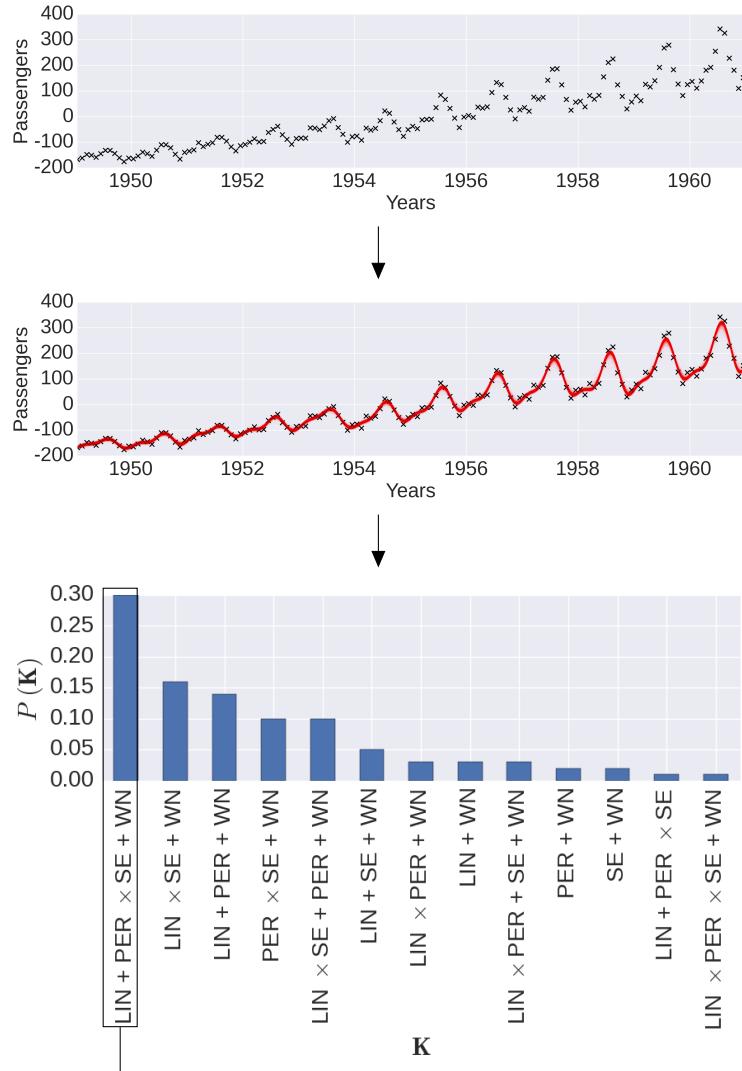
3.3 Bayesian optimization

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



534
535
536
537
538
539

Figure 4: Posterior of structure and qualitative, human interpretable reading. We take the raw data (top), compute a posterior distribution on structures (red samples and bar plot). We take the peak of this distribution ($\text{LIN} + \text{PER} + \text{SE} + \text{WN}$) with the sampled parameters used to generate the samples for the second plot from the top and write it in functional form with parameters. We depict the human readable interpretation of the equation on the bottom.



$$= 7.47^2(xx') + \left(0.27^2 \exp\left(-\frac{(x-x')^2}{2 \times 4.63^2}\right) \times 7.34^2 \exp\left(\frac{2 \sin^2(\pi(x-x')/4.4)}{4.55^2}\right) \right) + 2.93^2 \delta_{x,x'}$$

Qualitative Interpretation: The posterior peaks at a kernel structure with three 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 approximate periodic function; and (iv) white noise.

Figure 5: Posterior of structure and qualitative, human interpretable reading. We take the raw data (top), compute a posterior distribution on structures (red samples and bar plot). We take the peak of this distribution ($\text{LIN} + \text{PER} \times \text{SE} + \text{WN}$) with the sampled parameters used to generate the samples for the second plot from the top and write it in functional form with parameters. We depict the human readable interpretation of the equation on the bottom.

$$P(\mathbf{K}) \approx \frac{1}{N} \sum_{n=1}^N f(\mathbf{K}_n) \text{ where } f(\mathbf{K}_n) = \begin{cases} 1, & \text{if } \mathbf{K}_n \in \mathbf{K}_{global}, \\ 0, & \text{otherwise.} \end{cases}$$

$$P(\mathbf{K}_a \wedge \mathbf{K}_b) \approx \frac{1}{N} \sum_{n=1}^N f(\mathbf{K}_n) \text{ where } f(\mathbf{K}_n) = \begin{cases} 1, & \text{if } \mathbf{K}_a \text{ and } \mathbf{K}_b \in \text{global } \mathbf{K}_n, \\ 0, & \text{otherwise.} \end{cases}$$

$$P(\mathbf{K}_a \vee \mathbf{K}_b) = P(\mathbf{K}_a) + P(\mathbf{K}_b) - P(\mathbf{K}_a \wedge \mathbf{K}_b)$$



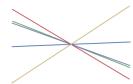
What is the probability of a trend, a recurring pattern **and** noise in the data?

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

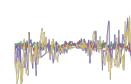
Is there a trend?
 $P(\text{LIN} \vee \text{LIN} \times \text{SE}) = 0.65$



A linear trend? $P(\text{LIN}) = 0.63$ A smooth trend?
 $P(\text{LIN} \times \text{SE}) = 0.02$



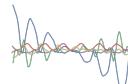
Is there noise?



Heteroskedastic noise? White noise?
 $P(\text{LIN} \times \text{WN}) = 0$ $P(\text{WN}) = 0.75$



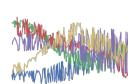
Is there repeating structure?
 $P(\text{PER} \vee \text{PER} \times \text{SE} \vee \text{PER} \times \text{LJN}) = 0.73$



PER = 0.32



PER X



$\times \text{LIN} = 0.07$



Figure 6: We can query the data if some logical statements are probable to be true, for example, is it true that there is a trend?

648
649 Here, an agent is to take a sequence of actions a_1, a_2, \dots from a (possibly infinite) set of possible
650 actions \mathcal{A} . After each action, a reward $r \in \mathbb{R}$ is received, according to an unknown conditional
651 distribution $P_{\text{true}}(r | a)$. The agent's goal is to maximize the total reward received for all actions in
652 an online manner. In Thompson sampling, the agent accomplishes this by placing a prior distribution
653 $P(\vartheta)$ on the possible "contexts" $\vartheta \in \Theta$. Here a context is a believed model of the conditional
654 distributions $\{P(r | a)\}_{a \in \mathcal{A}}$, or at least, a believed statistic of these conditional distributions which
655 is sufficient for deciding an action a . If actions are chosen so as to maximize expected reward, then
656 one such sufficient statistic is the believed conditional mean $V(a | \vartheta) = \mathbb{E}[r | a; \vartheta]$, which can be
657 viewed as a believed value function. For consistency with what follows, we will assume our context
658 ϑ takes the form $(\theta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$ where \mathbf{a}_{past} is the vector of past actions, \mathbf{r}_{past} is the vector of their
659 rewards, and θ (the "semicontext") contains any other information that is included in the context.

660 In this setup, Thompson sampling has the following steps:

661 **Algorithm 1** Thompson sampling.

663 Repeat as long as desired:

- 664 1. **Sample.** Sample a semicontext $\theta \sim P(\theta)$.
 - 665 2. **Search (and act).** Choose an action $a \in \mathcal{A}$ which (approximately) maximizes $V(a | \vartheta) =$
666 $\mathbb{E}[r | a; \vartheta] = \mathbb{E}[r | a; \theta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}]$.
 - 667 3. **Update.** Let r_{true} be the reward received for action a . Update the believed distribution on
668 θ , i.e., $P(\theta) \leftarrow P_{\text{new}}(\theta)$ where $P_{\text{new}}(\theta) = P(\theta | a \mapsto r_{\text{true}})$.
-

671 Note that when $\mathbb{E}[r | a; \vartheta]$ (under the sampled value of θ for some points a) is far from the true value
672 $\mathbb{E}_{P_{\text{true}}}[r | a]$, the chosen action a may be far from optimal, but the information gained by probing
673 action a will improve the belief ϑ . This amounts to "exploration." When $\mathbb{E}[r | a; \vartheta]$ is close to the
674 true value except at points a for which $\mathbb{E}[r | a; \vartheta]$ is low, exploration will be less likely to occur, but
675 the chosen actions a will tend to receive high rewards. This amounts to "exploitation." The trade-off
676 between exploration and exploitation is illustrated in Figure 7. Roughly speaking, exploration will
677 happen until the context ϑ is reasonably sure that the unexplored actions are probably not optimal,
678 at which time the Thompson sampler will exploit by choosing actions in regions it knows to have
679 high value.

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

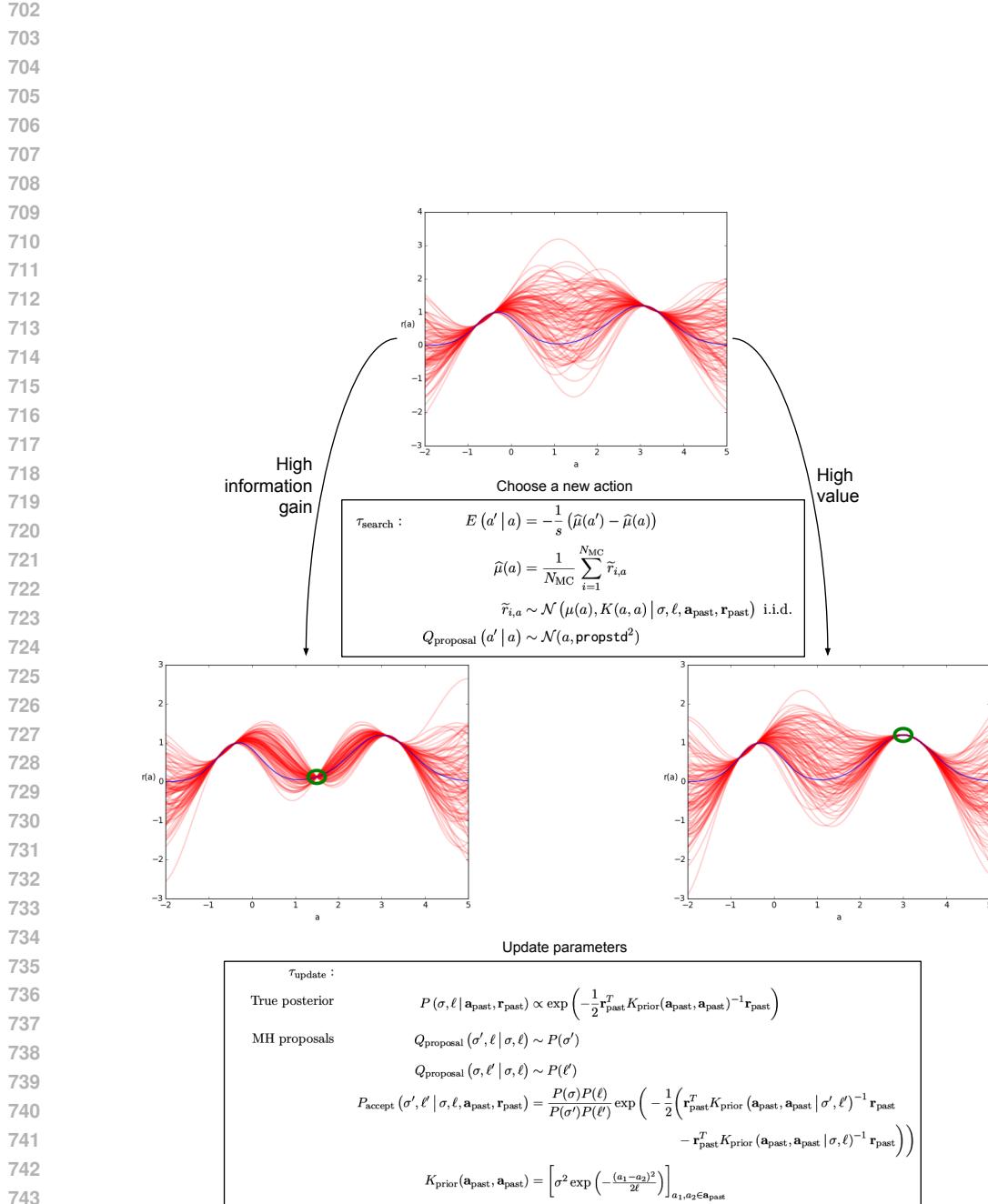
690 **A Mathematical Specification**

691 Our mathematical specification assert the following properties:

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

697 In this version of Thompson sampling, the contexts ϑ are Gaussian processes over the action space
698 $\mathcal{A} = [-20, 20] \subseteq \mathbb{R}$. That is,

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



745 Figure 7: Two possible actions (in green) for an iteration of Thompson sampling. The believed
 746 distribution on the value function V is depicted in red. In this example, the true reward function is
 747 deterministic, and is drawn in blue. The action on the right receives a high reward, while the action
 748 on the left receives a low reward but greatly improves the accuracy of the believed distribution on
 749 V . The transition operators τ_{search} and τ_{update} are described in Section 3.3.

750
 751
 752
 753
 754
 755

756 where the mean μ 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 $\{R_a\}_{a \in \mathcal{A}}$,
 757 where R_a represents the reward for action a . Computationally, we represent a context as a data
 758 structure
 759

$$760 \quad \vartheta = (\theta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}) = (\mu_{\text{prior}}, K_{\text{prior}}, \eta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}),$$

761 where μ_{prior} is a procedure to be used as the prior mean function. w.l.o.g. we set $\mu_{\text{prior}} \equiv 0$. K_{prior} is
 762 a procedure to be used as the prior covariance function, parameterized by η .
 763

764 The posterior mean and covariance for such a context ϑ are gotten by the usual conditioning formulas
 765 (assuming, for ease of exposition as above, that the prior mean is zero):³
 766

$$\begin{aligned} 767 \quad \mu(\mathbf{a}) &= \mu(\mathbf{a} | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}) \\ 768 \quad &= K_{\text{prior}}(\mathbf{a}, \mathbf{a}_{\text{past}}) K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}})^{-1} \mathbf{r}_{\text{past}} \\ 769 \quad K(\mathbf{a}, \mathbf{a}) &= K(\mathbf{a}, \mathbf{a} | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}}) \\ 770 \quad &= K_{\text{prior}}(\mathbf{a}, \mathbf{a}) - K_{\text{prior}}(\mathbf{a}, \mathbf{a}_{\text{past}}) K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}})^{-1} K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}). \end{aligned}$$

771 Note that the context space Θ is not a finite-dimensional parametric family, since the vectors
 772 \mathbf{a}_{past} and \mathbf{r}_{past} grow as more samples are taken. Θ is, however, representable as a computational
 773 procedure together with parameters and past samples, as we do in the representation $\vartheta =$
 774 $(\mu_{\text{prior}}, K_{\text{prior}}, \eta, \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$.
 775

776 We combine the Update and Sample steps of Algorithm 1 by running a Metropolis–Hastings (MH)
 777 sampler whose stationary distribution is the posterior $P(\theta | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$. The functional forms of
 778 μ_{prior} and K_{prior} are fixed in our case, so inference is only done over the parameters $\eta = \{\sigma, \ell\}$; hence
 779 we equivalently write $P(\sigma, \ell | \mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}})$ for the stationary distribution. We make MH proposals to
 780 one variable at a time, using the prior as proposal distribution:
 781

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

782 and
 783

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

784 The MH acceptance probability for such a proposal is
 785

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

786 Because the priors on σ and ℓ are uniform in our case, the term involving Q_{proposal} equals 1 and we
 787 have simply
 788

$$\begin{aligned} 789 \quad P_{\text{accept}}(\sigma', \ell' | \sigma, \ell) &= \min \left\{ 1, \frac{P(\mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}} | \sigma', \ell')}{P(\mathbf{a}_{\text{past}}, \mathbf{r}_{\text{past}} | \sigma, \ell)} \right\} \\ 790 \quad &= \min \left\{ 1, \exp \left(-\frac{1}{2} \left(\mathbf{r}_{\text{past}}^T K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}} | \sigma', \ell')^{-1} \mathbf{r}_{\text{past}} \right. \right. \right. \\ 791 \quad &\quad \left. \left. \left. - \mathbf{r}_{\text{past}}^T K_{\text{prior}}(\mathbf{a}_{\text{past}}, \mathbf{a}_{\text{past}} | \sigma, \ell)^{-1} \mathbf{r}_{\text{past}} \right) \right) \right\}. \end{aligned}$$

792 The proposal and acceptance/rejection process described above define a transition operator τ_{update}
 793 which is iterated a specified number of times; the resulting state of the MH Markov chain is taken
 794 as the sampled semicontext θ in Step 1 of Algorithm 1.
 795

796 For Step 2 (Search) of Thompson sampling, we explore the action space using an MH-like transition
 797 operator τ_{search} . As in MH, each iteration of τ_{search} produces a proposal which is either accepted or
 798 rejected, and the state of this Markov chain after a specified number of steps is the new action a .
 799 The Markov chain's initial state is the most recent action, and the proposal distribution is Gaussian
 800 drift:
 801

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

802 ³Here, 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
 803 the matrix $[K(a_i, a'_j)]_{1 \leq i \leq n, 1 \leq j \leq n'}$.

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

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

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

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

815 where

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

816 and

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

817 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
 818 N_{avg} are specified ahead of time. Proposals of estimated value higher than that of the current action
 819 are always accepted, while proposals of estimated value lower than that of the current action are
 820 accepted with a probability that decays exponentially with respect to the difference in value. The
 821 rate of the decay is determined by the temperature parameter s , where high temperature corresponds
 822 to generous acceptance probabilities. For $s = 0$, all proposals of lower value are rejected; for
 823 $s = \infty$, all proposals are accepted. For points a at which the posterior mean $\mu(a)$ is low but the
 824 posterior variance $K(a, a)$ is high, it is possible (especially when N_{avg} is small) to draw a “wild”
 825 value of $\hat{\mu}(a)$, resulting in a favorable acceptance probability.

826 Indeed, taking an action a with low estimated value but high uncertainty serves the useful function
 827 of improving the accuracy of the estimated value function at points near a (see Figure 7).^{4,5} We see
 828 a complete probabilistic program with `gpmem` implementing Bayesian optimization with Thompson
 829 Sampling below (Listing 1).

830
 831
 832
 833
 834
 835
 836
 837
 838
 839
 840
 841
 842
 843
 844
 845
 846
 847
 848
 849
 850
 851
 852
 853
 854
 855
 856
 857
 858
 859
 860
 861

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

⁵For 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 μ .

```

864
865 Listing 1: Bayesian optimization using gpmem
866
867
868 1 assume sf = tag(quote(hyper), 0, uniform_continuous(0, 10))
869 2 assume l = tag(quote(hyper), 1, uniform_continuous(0, 10))
870 3 assume se = make_squaredexp(sf, l)
871 4 assume blackbox_f = get_bayesopt_blackbox()
872 5 assume (f_compute, f_emulate) = gpmem(blackbox_f, se)
873
874 // A naive estimate of the argmax of the given function
875 define mc_argmax = proc(func) {
876     candidate_xs = mapv(proc(i) {uniform_continuous(-20, 20)},
877                          arange(20));
878     candidate_ys = mapv(func, candidate_xs);
879     lookup(candidate_xs, argmax_of_array(candidate_ys))
880 };
881
882 // Shortcut to sample the emulator at a single point without packing
883 // and unpacking arrays
884 define emulate_pointwise = proc(x) {
885     run(sample(lookup(f_emulate(array(unquote(x))), 0)))
886 };
887
888 // Main inference loop
889 infer repeat(15, do(pass,
890     // Probe V at the point mc_argmax(emulate_pointwise)
891     predict(f_compute(unquote(mc_argmax(emulate_pointwise)))),
892     // Infer hyperparameters
893     mh(quote(hyper), one, 50)));

```

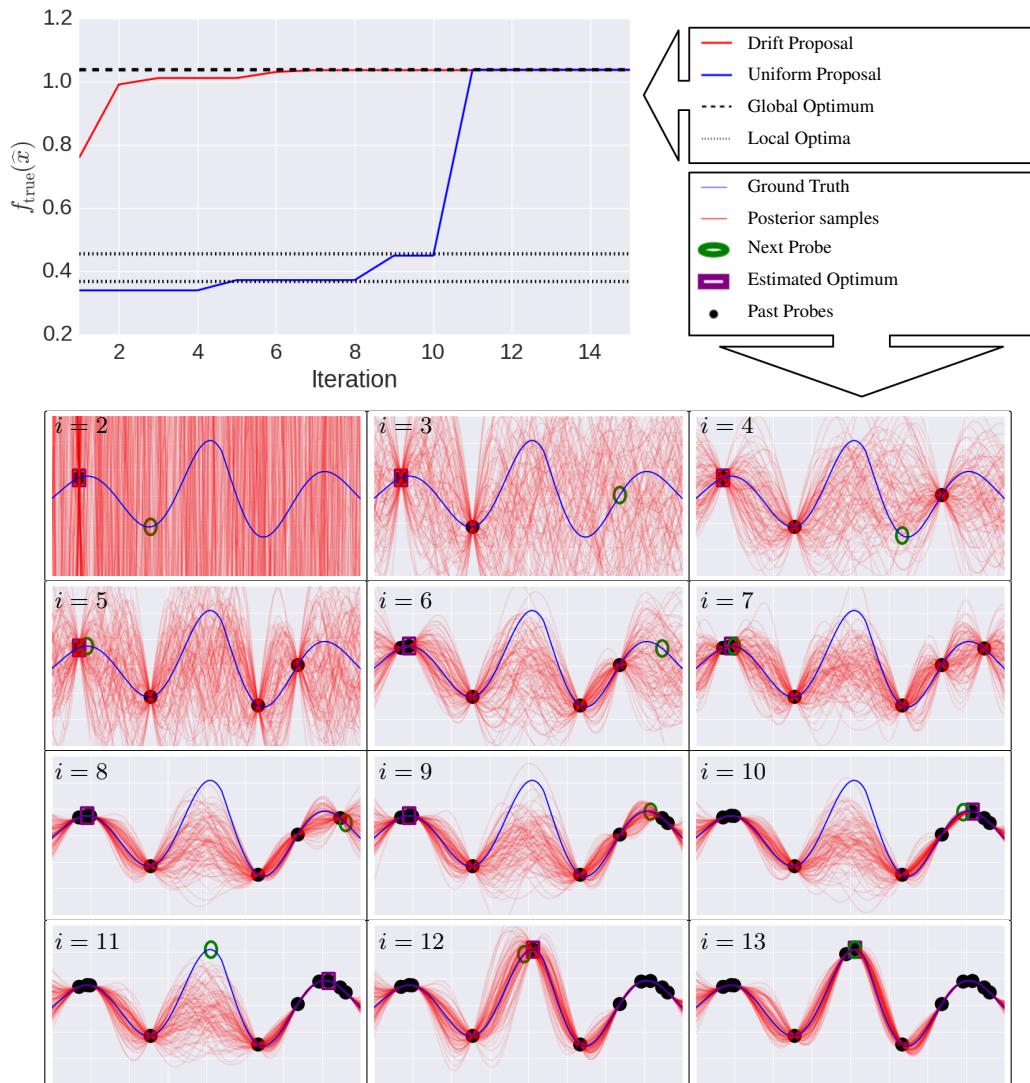
4 Discussion

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.

```

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

```



Appendix

A Covariance Functions

$$\text{SE} = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right) \quad (18)$$

$$\text{LIN} = \sigma^2(xx') \quad (19)$$

$$C = \sigma^2 \quad (20)$$

$$\text{WN} = \sigma^2 \delta_{x,x'} \quad (21)$$

$$RQ = \sigma^2 \left(1 + \frac{(x - x')^2}{2\alpha\ell^2} \right)^{-\alpha} \quad (22)$$

$$\text{PER} = \sigma^2 \exp\left(\frac{2 \sin^2(\pi(x - x')/p)}{\ell^2}\right). \quad (23)$$

B Covariance Simplification

$$\begin{array}{ll}
 \text{SE} \times \text{SE} & \rightarrow \text{SE} \\
 \{\text{SE}, \text{PER}, \text{C}, \text{WN}\} \times \text{WN} & \rightarrow \text{WN} \\
 \text{LIN} + \text{LIN} & \rightarrow \text{LIN} \\
 \{\text{SE}, \text{PER}, \text{C}, \text{WN}, \text{LIN}\} \times \text{C} & \rightarrow \{\text{SE}, \text{PER}, \text{C}, \text{WN}, \text{LIN}\}
 \end{array}$$

Rule 1 is derived as follows:

$$\begin{aligned}
\sigma_c^2 \exp\left(-\frac{(x-x')^2}{2\ell_c^2}\right) &= \sigma_a^2 \exp\left(-\frac{(x-x')^2}{2\ell_a^2}\right) \times \sigma_b^2 \exp\left(-\frac{(x-x')^2}{2\ell_b^2}\right) \\
&= \sigma_c^2 \exp\left(-\frac{(x-x')^2}{2\ell_a^2}\right) \times \exp\left(-\frac{(x-x')^2}{2\ell_b^2}\right) \\
&= \sigma_c^2 \exp\left(-\frac{(x-x')^2}{2\ell_a^2} - \frac{(x-x')^2}{2\ell_b^2}\right) \\
&= \sigma_c^2 \exp\left(-\frac{(x-x')^2}{2\ell_c^2}\right)
\end{aligned} \tag{24}$$

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

$$\sigma_a^2 \exp\left(-\frac{(x-x')^2}{2\ell_z^2}\right) \times \sigma_b \delta_{x,x'} = \sigma_a \sigma_b \delta_{x,x'} \quad (25)$$

Rule 3 is derived as follows:

$$\theta_c(x \times x') = \theta_a(x \times x') + \theta_b(x \times x') \quad (26)$$

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

1026 **References**
1027

- 1028 Barry, D. (1986). Nonparametric bayesian regression. *The Annals of Statistics*, 14(3):934–953.
1029 Box, G. E. P., Jenkins, G. M., and Reinsel, G. C. (1997). *Time series analysis: forecasting and*
1030 *control*.
1031 Damianou, A. and Lawrence, N. (2013). Deep gaussian processes. In *Proceedings of the Interna-*
1032 *tional Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 207–215.
1033 Duvenaud, D., Lloyd, J. R., Grosse, R., Tenenbaum, J., and Ghahramani, Z. (2013). Structure
1034 discovery in nonparametric regression through compositional kernel search. In *Proceedings of*
1035 *the International Conference on Machine Learning (ICML)*, pages 1166–1174.
1036 Ferris, B., Haehnel, D., and Fox, D. (2006). Gaussian processes for signal strength-based location
1037 estimation. In *Proceedings of the Conference on Robotics Science and Systems*. Citeseer.
1038 Gelbart, M. A., Snoek, J., and Adams, R. P. (2014). Bayesian optimization with unknown con-
1039 straints. *arXiv preprint arXiv:1403.5607*.
1040 Goodman, N. D .and Mansinghka, V. K., Roy, D., Bonawitz, K., and Tenenbaum, J. (2008). Church:
1041 A language for generative models. In *Proceedings of the Conference on Uncertainty in Artificial*
1042 *Intelligence (UAI)*, pages 220–229.
1043 Kemmler, M., Rodner, E., Wacker, E., and Denzler, J. (2013). One-class classification with gaussian
1044 processes. *Pattern Recognition*, 46(12):3507–3518.
1045 Kennedy, M. C. and O’Hagan, A. (2001). Bayesian calibration of computer models. *Journal of the*
1046 *Royal Statistical Society. Series B, Statistical Methodology*, pages 425–464.
1047 Kuss, M. and Rasmussen, C. E. (2005). Assessing approximate inference for binary gaussian process
1048 classification. *The Journal of Machine Learning Research*, 6:1679–1704.
1049 Kwan, J., Bhattacharya, S., Heitmann, K., and Habib, S. (2013). Cosmic emulation: The
1050 concentration-mass relation for wcdm universes. *The Astrophysical Journal*, 768(2):123.
1051 Lloyd, J. R., Duvenaud, D., Grosse, R., Tenenbaum, J., and Ghahramani, Z. (2014). Automatic con-
1052 struction and natural-language description of nonparametric regression models. In *Proceedings*
1053 *of the Conference on Artificial Intelligence (AAAI)*.
1054 Mansinghka, V. K., Selsam, D., and Perov, Y. (2014). Venture: a higher-order probabilistic pro-
1055 gramming platform with programmable inference. *arXiv preprint arXiv:1404.0099*.
1056 McAllester, D., Milch, B., and Goodman, N. D. (2008). Random-world semantics and syntactic
1057 independence for expressive languages. Technical report.
1058 Neal, R. M. (1995). *Bayesian Learning for Neural Networks*. PhD thesis, University of Toronto.
1059 Neal, R. M. (1997). Monte carlo implementation of gaussian process models for bayesian regression
1060 and classification. *arXiv preprint physics/9701026*.
1061 Poole, D. (1993). Probabilistic horn abduction and bayesian networks. *Artificial Intelligence*,
1062 64(1):81–129.
1063 Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning (Adap-*
1064 *tive Computation and Machine Learning)*. The MIT Press.
1065 Sato, T. (1995). A statistical learning method for logic programs with distribution semantics. In *In*
1066 *Proceedings of the International Conference on Logic Programming*. Citeseer.
1067 Schneider, M. D., Knox, L., Habib, S., Heitmann, K., Higdon, D., and Nakhleh, C. (2008). Simula-
1068 tions and cosmological inference: A statistical model for power spectra means and covariances.
1069 *Physical Review D*, 78(6):063529.
1070 Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical bayesian optimization of machine
1071 learning algorithms. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2951–
1072 2959.
1073 Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view
1074 of the evidence of two samples. *Biometrika*, pages 285–294.
1075 Williams, C. K. I. and Barber, D. (1998). Bayesian classification with gaussian processes. *IEEE*
1076 *Transactions on Pattern Analysis and Machine Intelligence*, 20(12):1342–1351.