

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

Gaussian Processes with Probabilistic Programming

Anonymous Author(s)

Affiliation

Address

email

Abstract

We introduce Venture GPs, a way to formulate Gaussian Processes with probabilistic programming. Gaussian Processes are flexible non-parametric models that can be applied to a broad class of problems. We represent Gaussian Processes in Venture, a Turing complete probabilistic programming language. Venture provides a compositional language with a generalized inference engine that builds on a stochastic procedure interface. This stochastic procedure interface specifies and encapsulates primitive random variables analogously conditional probability tables in Bayesian Networks. The programming language is extended with a set of stochastic processes that allow a user to formulate Gaussian Process models and to perform numerically stable inference over them while hiding the linear algebra needed for this inside the language. We show how we can extend the non-parametric model to incorporate hierarchical causal priors on model structure and hyper-parameters with only a few lines of code. We also show state-of-the-art applications of Gaussian Processes in this framework, namely structure discovery of high-level properties of Gaussian Processes, Bayesian Optimization and hyper-parameter inference. We evaluate the performance of the programs with synthetic and real world data.

1 Introduction

Probabilistic programming could be revolutionary for machine intelligence due to universal inference engines and the rapid prototyping for novel models (Ghahramani, 2015). Probabilistic programming languages aim to provide a formal language to specify probabilistic models in the style of computer programming and can represent any computable probability distribution as a program. In this work, we will introduce new features of Venture, a recently developed probabilistic programming language and the first probabilistic programming language suitable for general purpose use (Mansinghka et al., 2014). Venture comes with scalable performance on hard problems and with a general purpose inference engine. The inference engine is based on Markov Chain Monte Carlo (MCMC) methods (for an introduction, see Andrieu et al. (2003)). MCMC lends itself models with complex structures such as probabilistic programs or hierarchical Bayesian non-parametric models since they can provide a vehicle to express otherwise intractable integrals necessary for a fully Bayesian representation. MCMC is scalable, often distributable and also compositional. That is, one can arbitrarily chain MCMC kernels to infer over several hierarchically connected or nested models as they will emerge in probabilistic programming.

One very powerful model yet unseen in probabilistic programming languages are Gaussian Processes (GPs). GPs are gaining increasing attention for representing unknown functions by posterior probability distributions in various fields such as machine learning, signal processing, computer vision and bio-medical data analysis. In the following, we will present Gaussian Processes as a novel feature for probabilistic programming languages. Our contribution is threefold:

- we introduce a new stochastic process for GPs in a probabilistic programming language;

- we show how one can solve hard problems of state-of-the-art machine learning with only a few lines of Venture code; and
- we introduce an additional stochastic process that samples from a probabilistic context free grammar for GP covariance structure generation.

The paper is structured as follows, we will first provide some background on probabilistic programming in Venture and GPs. We will then elaborate on our new stochastic processes. Finally, we will show how we can apply those on problems of hyper-parameter inference, structure discovery for Gaussian Processes and Bayesian Optimization including experiments with real world and synthetic data.

2 Background

2.1 Venture

Venture is a compositional language for custom inference strategies that comes with a Scheme- and Java-Script-like front-end syntax. It's implementation is based on three concepts:

1. stochastic procedure interfaces that specify and encapsulate random variables, analogously to conditional probability tables in a Bayesian network;
2. probabilistic execution traces that represent execution histories and capture conditional dependencies; and
3. scaffolds that partition execution histories and factor global inference problems into sub-problems.

These building blocks provide a powerful way to represent probability distributions; some of which cannot be expressed with density functions. For the purpose of this work the most important Venture directives that operate on these building blocks to understand are ASSUME, OBSERVE, SAMPLE and INFER. ASSUME induces a hypothesis space for (probabilistic) models including random variables by binding the result of an expression to a symbol. SAMPLE simulates a model expression and returns a value. OBSERVE adds constraints to model expressions. INFER instructions incorporate observations and cause Venture to find a hypothesis that is probable given the data.

INFER is most commonly done by deploying the Metropolis-Hastings algorithm (MH) (Metropolis et al., 1953). Many algorithms used in the MCMC world can be interpreted as special cases of MH (Andrieu et al., 2003). We can outline the MH algorithm as follows. For T steps we sample x^* from a proposal distribution q :

$$x^* \sim q(x^* | x^{(t)}) \quad (1)$$

which we accept ($x^{t+1} \leftarrow x^*$) with ratio:

$$\alpha = \min \left\{ 1, \frac{p(x^*)q(x^t | x^*)}{p(x^{(t)})q(x^* | x^t)} \right\} \quad (2)$$

Venture implements an MH transition operator for probabilistic execution traces.

2.2 Gaussian Processes

In the following, we will introduce GP related theory and notations. We will exclusively work on two variable regression problems. Let the data be 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 that we assume to have generated the data. f is assumed latent and the GP prior is given by a multivariate Gaussian $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(x_i, x'_i))$, where $m(\mathbf{x})$ is a function of the mean of all functions that map to y_i at x_i and $k(x_i, x'_i)$ is a kernel or covariance function that summarizes the covariance of all functions that map to y_i at x_i . We can absorb the mean function into the covariance function so without loss of generality we can set the mean to zero. The marginal likelihood can be expressed as:

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}, \mathbf{x}) p(\mathbf{f}|\mathbf{x}) d\mathbf{f} \quad (3)$$

108 where the prior is Gaussian $\mathbf{f}|\mathbf{x} \sim \mathcal{N}(0, k(\mathbf{x}, \mathbf{x}'))$. For a zero mean Gaussian Process this results in
 109 a Gaussian posterior $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$ with mean:
 110

$$\boldsymbol{\mu} = \mathbf{K}(\mathbf{x}, \mathbf{x}^*) \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*)^{-1} \mathbf{y} \quad (4)$$

112 and covariance
 113

$$\Sigma = \mathbf{K}(\mathbf{x}, \mathbf{x}) + \mathbf{K}(\mathbf{x}, \mathbf{x}^*) \mathbf{K}(\mathbf{x}^*, \mathbf{x}^*)^{-1} \mathbf{K}(\mathbf{x}^*, \mathbf{x}). \quad (5)$$

115 where \mathbf{K} is a covariance function. We can compute both the log-likelihood and and the predictive
 116 posterior in a Venture SP efficiently by resorting to an algorithm that resorts to Cholesky factoriza-
 117 tion(Rasmussen and Williams, 2006, chap. 2) resulting in a computational complexity of $\mathcal{O}(n^3)$ in
 118 the number of data-points.

119 The covariance function covers general high-level properties of the observed data such as linear-
 120 ity, periodicity and smoothness. The most widely used type of covariance function is the squared
 121 exponential covariance function:

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right) \quad (6)$$

125 where σ and ℓ are hyper-parameters. σ is a scaling factor and ℓ is the typical length-scale. Smaller
 126 variations can be achieved by adapting these hyper-parameters.
 127

128 3 Venture GPs

130 Below, we see how we can express simple GP smoothing with a few lines
 131 of Venture code while allowing users to custom design covariance func-
 132 tions. Throughout the paper, we will use the Scheme-like front-end syntax.
 133

134 Listing 1: GP Smoothing

```
135 [ASSUME 1 1] ∈ {hyper-parameters}
136 [ASSUME sf 2] ∈ {hyper-parameters}
137
138  $k(x, x') := \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$ 
139
140 [ASSUME f VentureFunction(k, σ, ℓ) ]
141 [ASSUME SE make-se (apply-function f 1 sf) ]
142 [ASSUME (make-gp 0 SE) ]
143
144 [SAMPLE GP (array 1 2 3)] % Prior
145
146 [OBSERVE GP D]
147
148 [SAMPLE GP (array 1 2 3)]
149
150 [INFER (MH {hyper-parameters} one 100) ]
151
152 [SAMPLE GP (array 1 2 3)] % Posterior
```

153 The first two lines depict the hyper-parameters. We tag both of them to belong to the set {hyper-
 154 parameters}. Every member of this set belongs to the same inference scope. This scope controls the
 155 application of the inference procedure used. In this paper, we use MH throughout. Each scope is
 156 further subdivided into blocks that allow to do block-proposals. In the following we omit the block
 157 notation for readability, since we always choose the block of a certain scope at random.

158 The ASSUME directives describe the assumptions we make for the GP model, we assume the hyper-
 159 parameters 1 and sf (corresponding to ℓ, σ) to be 1 and 2. The squared exponential covariance
 160 function can be defined outside the Venture code with foreign conventional programming languages,
 161 e.g. Python. In that way, the user can define custom covariance functions without being restricted to
 the most common ones. We then integrate the foreign function into Venture as VentureFunction. In

162 the next line this function is associated with the hyper-parameters. Finally, we assume a Gaussian
 163 Process SP with a zero mean and the previously assumed squared exponential covariance function.
 164

165 In the case where hyper-parameters are unknown they can be found deterministically by optimizing
 166 the marginal likelihood using a gradient based optimizer. Non-deterministic, Bayesian representa-
 167 tions of this case are also known (Neal, 1997). Extending the program described in listing 1 for a
 168 Bayesian treatment of hyper-parameters is simple using the build in stochastic procedure that simu-
 169 lates drawing samples from a gamma distribution:

170 Listing 2: Bayesian GP Smoothing
 171

```
172 [ASSUME 1 (gamma 1 3)]
173 [ASSUME sf (gamma 1 2)]
174  $k(x, x') := \sigma^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$ 
175
176 [ASSUME f VentureFunction ( $k, \sigma, \ell$ ) ]
177 [ASSUME SE make-se (apply-function f 1 sf) ]
178 [ASSUME (make-gp 0 SE ) ]
179
```

180 Larger variations are achieved by changing the type of the covariance function structure. A different
 181 type could be a linear covariance function:
 182

$$k(x, x') = \sigma^2(x - \ell)(x' - \ell). \quad (7)$$

184 Note that covariance function structures are compositional. We can add covariance functions if we
 185 want to model globally valid structures
 186

$$k_3(x, x') = k_1(x, x') + k_2(x, x') \quad (8)$$

188 and we can multiply covariance functions if the data is best explained by local structure
 189

$$k_4(x, x') = k_1(x, x') \times k_2(x, x'); \quad (9)$$

192 both, k_3 and k_4 are valid covariance function structures. This leads to an infinite space of possible
 193 structures that could potentially explain the observed data best (e.g. Fig. 1). In the following, we
 194 will refer to covariance functions that are not composite as base covariance functions. Note that this
 195 form of composition can be easily expressed in Venture, for example if one wishes to add a linear
 196 and a periodic kernel:
 197

198 Listing 3: LIN × PER
 199

```
200 [ASSUME 1 (gamma 1 3)]
201 [ASSUME sf (gamma 1 2)]
202 [ASSUME a (gamma 2 2)]
203
204  $k_{LIN}(x, x') = \sigma_1^2(x - \ell)(x' - \ell)$ 
205  $k_{PER}(x, x') := \sigma_2^2 \exp\left(-\frac{2 \sin^2(\pi(x-x')/p)}{p^2}\right)$ 
206
207 [ASSUME fLIN VentureFunction ( $k_{LIN}, \sigma_1$ ) ]
208 [ASSUME fPER VentureFunction ( $k_{PER}, \sigma_2, \ell, p$ ) ]
209 [ASSUME LIN (make-LIN (apply-function fLIN a)) ]
210 [ASSUME PER (make-PER (apply-function fPER 1 sf)) ]
211 [ASSUME (make-gp 0 (function-times LIN PER)) ]
```

212 Knowledge about the composite nature of covariance functions is not new, however, until recently,
 213 the choice and the composition of covariance functions were done ad-hoc. The Automated Statisti-
 214 cian Project came up with an approximate search over the possible space of kernel structures (Du-
 215 venaud et al., 2013; Lloyd et al., 2014). However, a fully Bayesian treatment of this was not done
 before.

216 **3.1 A Bayesian interpretation**

218 In the following, we will explore a Bayesian representation of GP. The probability of the hyper-
 219 parameters of a GP with assumptions as above and given covariance function structure \mathbf{K} can be
 220 described as:

221
$$P(\boldsymbol{\theta} | \mathbf{D}, \mathbf{K}) = \frac{P(\mathbf{D} | \boldsymbol{\theta}, \mathbf{K})P(\boldsymbol{\theta} | \mathbf{K})}{P(\mathbf{D} | \mathbf{K})}. \quad (10)$$

223 Neal suggested the treatment of outliers as a use-case for a Bayesian treatment of Gaussian pro-
 224 cesses (1997). He evaluates his MCMC setting using the following synthetic data problem. Let f
 225 be the underlying function that generates the data:

227
$$f(x) = 0.3 + 0.4x + 0.5 \sin(2.7x) + \frac{1.1}{(1+x^2)} + \eta \quad \text{with } \eta \sim \mathcal{N}(0, \sigma) \quad (11)$$

229 We synthetically generate outliers by setting $\sigma = 0.1$ in 95% of the case and to $\sigma = 1$ in the
 230 remaining cases. Venture GPs can capture the true underlying function within only 100 MH steps
 231 (see Fig. 2). Note that Neal devices an additional noise model and performs large numbe of Hybrid-
 232 Monte Carlo and Gibbs steps.

234 **3.2 Structure Learning**

236 The case where the covariance structure is not given is even more interesting. Our probabilistic
 237 programming based MCMC framework approximates the following intractable integrals of the ex-
 238 pectation for the prediction:

239
$$\mathbb{E}[y^* | x^*, D, \mathbf{K}_\Omega^s] = \iint f(x^*, \boldsymbol{\theta}, \mathbf{K}) P(\boldsymbol{\theta} | \mathbf{D}, \mathbf{K}) P(\mathbf{K} | \boldsymbol{\Omega}, s, n) d\boldsymbol{\theta} d\mathbf{K}. \quad (12)$$

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

244
$$y^* \approx \frac{1}{T} \sum_{t=1}^T f(x^* | \boldsymbol{\theta}^{(t)}, \mathbf{K}^{(t)}). \quad (13)$$

247 In order to provide the sampling of the kernel, we introduce a stochastic process to the SP that
 248 simulates the grammar for algebraic expressions of kernel algebra. Here, we start with a set of
 249 possible kernels and draw a random subset. For this subset of size n , we sample a set of possible
 250 operators that operate on the base kernels.

251 The marginal probability of a kernel structure which allows us to sample is characterized by the
 252 probability of a uniformly chosen subset of the set of n possible covariance functions times the
 253 probability of sampling a global or a local structure which is given by a binomial distribution:

255
$$P(\mathbf{K} | \boldsymbol{\Omega}, s, n) = P(\boldsymbol{\Omega} | s, n) \times P(s | n) \times P(n), \quad (14)$$

256 with

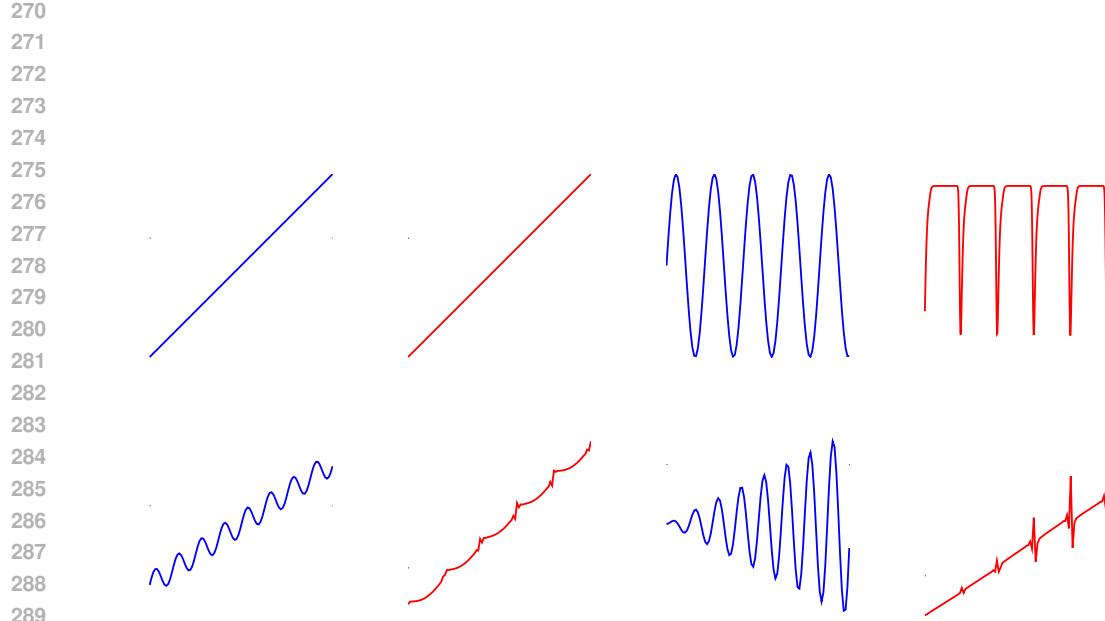
258
$$P(\boldsymbol{\Omega} | s, n) = \binom{n}{r} p_{+ \times}^k (1 - p_{+ \times})^{n-k} \quad (15)$$

260 and

261
$$P(s | n) = \frac{n!}{|s|!} \quad (16)$$

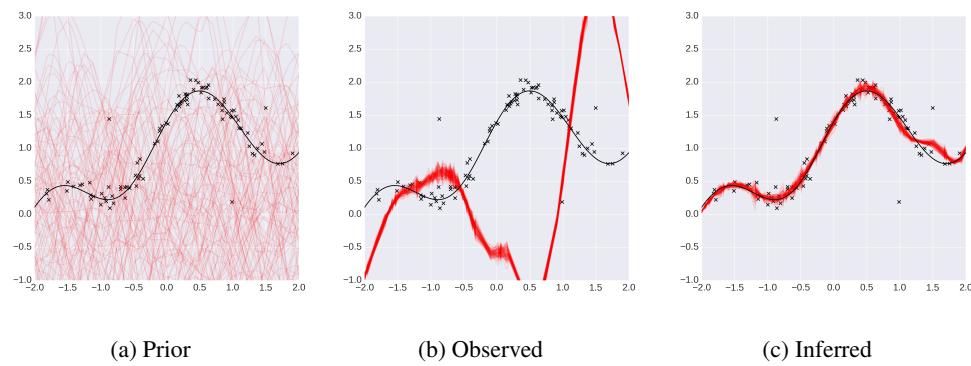
263 where $P(n)$ is a prior on the number of base kernels used. It is possible to also assign a prior for the
 264 probability to sample global or local priors, however, we have assigned complete uncertainty to this
 265 with the binomial $p = 0.5$.

266 Many equivalent covariance structures can be sampled due to covariance function algebra
 267 and equivalent representations with different parameterization (Lloyd et al., 2014). Cer-
 268 tain covariance functions can differ in terms of the hyper-parameterization but can be
 269 absorbed into a single covariance function with a different parameterization. To in-
 spect the posterior of these equivalent structures we convert each kernel expression into



290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315

Figure 1: Composition of covariance functions (blue, left) and samples from the distribution of curves they can produce (red, right).



316
317
318
319
320
321
322
323

Figure 2: Running a Venture GP on Neal's example for MCMC showing the prior, after having observed the data and after performing inference on the hyper-parameters. Note how the GP is choosing outliers to smooth instead of essential data before inference takes place.

324 a sum of products and subsequently simplify expressions using the following grammar:
325

326 Listing 4: Grammar to simplify expressions

327 SE × SE	→ SE
328 {SE, PER, C, WN} × WN	→ WN
329 LIN + LIN	→ LIN
330 {SE, PER, C, WN, LIN} × C	→ {SE, PER, C, WN, LIN}

332 For reproducing results from the Automated Statistician Project in a Bayesian fashion we first define
333 a prior on the hypothesis space. Note that, as in the implementation of the Automated Statistician,
334 we upper-bound the complexity of the space of covariance functions we want to explore. We also
335 put vague priors on hyper-parameters.

337 Listing 5: Venture Code for Bayesian GP Structure Learning

```
338 [ASSUME S (array K1,K2,...,Kn)] // (defined as above)
339 [ASSUME pn (uniform_structure n)]
340 [ASSUME S (array K1,K2,...,Kn)]
341 [ASSUME K* (grammar S pn)]
342 [ASSUME GP (make-gp 0 K*)]
343
344 [OBSERVE GP D]
345
346 [INFER (REPEAT 2000 (DO
347   (MH 10 pn one 1)
348   (MH 10 K* one 1)
349   (MH 10 {hyper-parameters} one 10))]
```

350 We defined the space of covariance structures in a way allowing us to reproduce results for covariance
351 function structure learning as in the Automated Statistician. This lead to coherent results, for
352 example for the airline data set. We will elaborate the result using a sample from the posterior (Fig.
353 3). The sample is identical with the highest scoring result reported in previous work using a search-
354 and-score method (Duvenaud et al., 2013) and the predictive capability is comparable. However, the
355 components factor in a different way due to different parameterization of the individual base kernels.

356 We further investigated the quality of our stochastic processes by running a leave one out cross-
357 validation to gain confidence on the posterior. This resulted in 545 independent runs of the Markov
358 chain that produced a coherent posterior: our Bayesian interpretation of GP structure and GPs pro-
359 duced a posterior of structures that is in line with previous results on this data set (Duvenaud et al.,
360 2013; see Fig. 4).

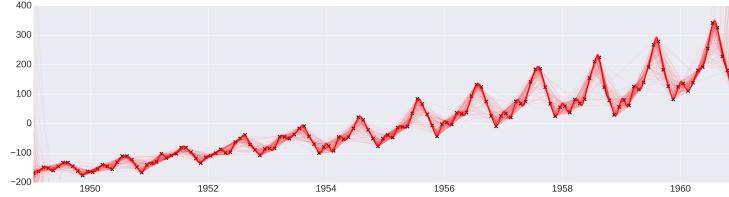
361 We found the final sample of multiple runs to be most informative. This kind of Markov Chain
362 seems to produce samples that are highly auto-correlated.

365 4 Bayesian Optimization

366 Bayesian Optimization poses the problem of finding the global maximum of an unknown function
367 as a hierarchical decision problem (Ghahramani, 2015). Evaluating the actual function can be
368 very expensive. For example, finding the best configuration for the learning algorithm of a large
369 convolutional neural network implies expensive function evaluations to compare a potentially infinite
370 number of configurations. Another common example is the example of data acquisition. For
371 problems with large amounts of data available it may be interested to chose certain informative data-
372 points to evaluate a model on. In continuous domains, many Bayesian Optimization methods deploy
373 GPs (e.g. Snoek et al., 2012).

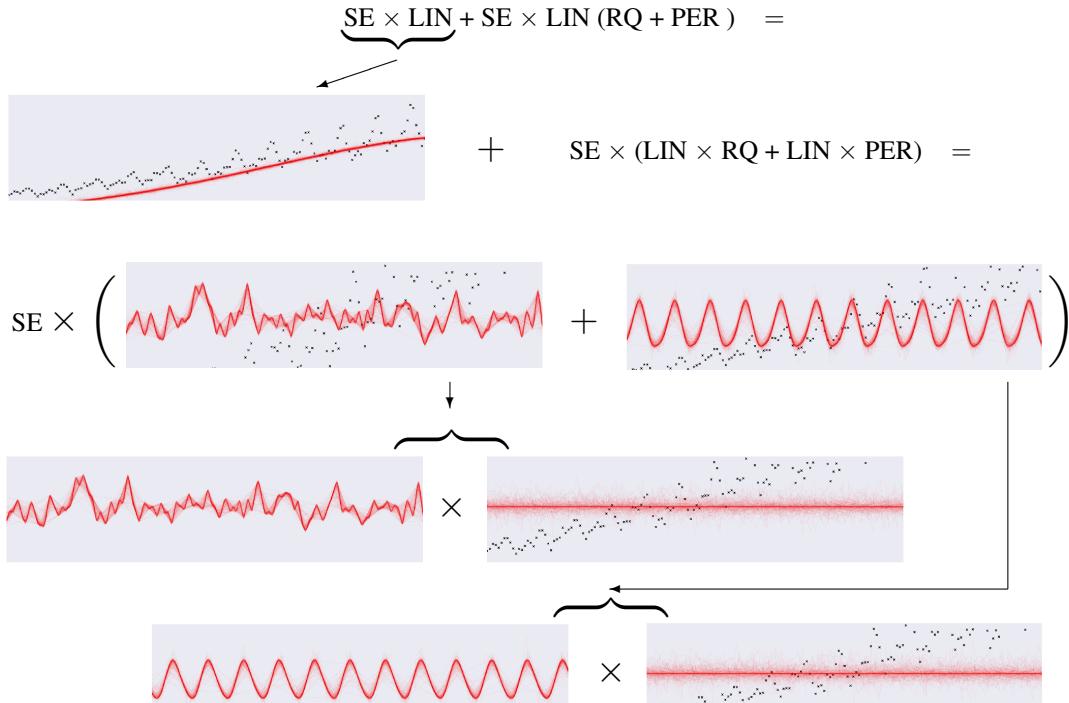
374 The hierarchical nature of Bayesian Optimization makes it an ideal application for GPs in Venture.
375 The following Bayesian Optimization scheme is closely related to Thompson Sampling Thompson
376 (1933). Thompson Sampling is a general framework to solve exploration-exploitation problems that

378
379
380
381
382
383
384
385
386
387
388
389
390



(a) The predictive posterior on given the full grammar structure.

391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417



(b) Compositional Structure

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

Figure 3: a) We see the predictive posterior as an result 1000 nested MH steps on the airline data set. b) depicts a decomposition of this posterior for the structures sampled by Venture. RQ is the rational quadratic covariance function. Note that although the overall result is in line with what Duvenaud et al. (2013) report a slightly different composition is implied by the sampled parameters for the structure. The first line shows the global trend and denotes the rest of the structure that is shown above. In the second line, the see the periodic component on the right hand side. The left hand side denotes short term deviations both multiplied by a smoothing kernel. The third and fourth lines denote how we reach the second line: both periodic and rational quadratic covariance functions are multiplied by a line with slope zero.

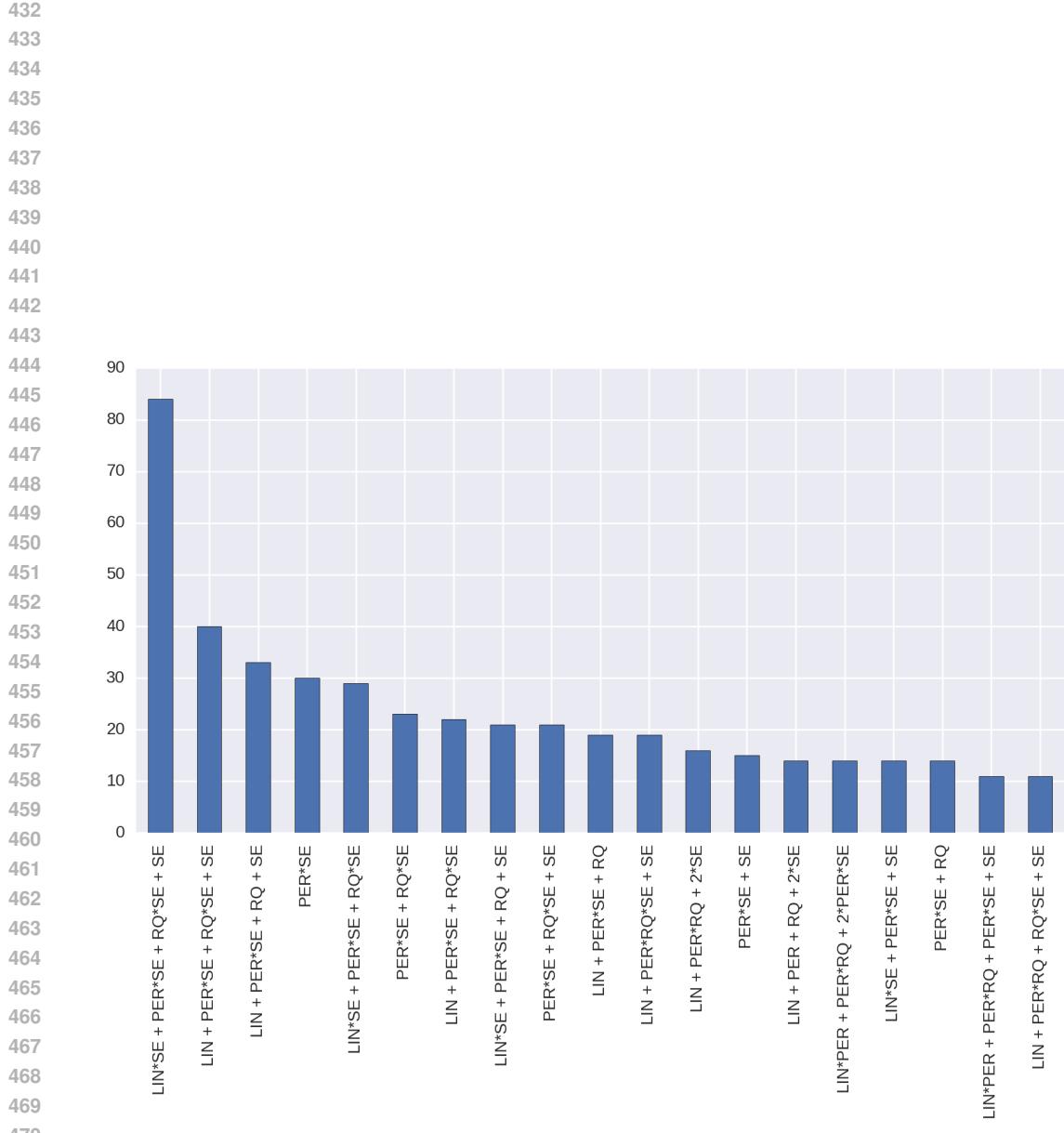


Figure 4: Posterior on structure of the CO₂ data. We have cut the tail of the distribution for space reasons since the number of possible structures is large. We see the final sample of the each of the 545 chains with 2000 nested steps each. Note that Duvenaud et al. (2013) report LIN × SE + PER × SE + RQ × SE.

486 applies to our notion of Bayesian Optimization. We sample a probe from the posterior
487

488 $\hat{\theta} \sim P(\hat{\theta}|\vec{x})$ (17)

489 and then optimize the expected reward by choosing an action a
490

491 $a_i = \arg \max_a \mathbb{E}_{P(\cdot|\theta)}[r(world_\theta(a))]$ (18)
492

493 $x_i \sim world_\theta(a_i)$ (19)
494

495 References

- 496 Andrieu, C., De Freitas, N., Doucet, A., and Jordan, M. I. (2003). An introduction to mcmc for
497 machine learning. *Machine learning*, 50(1-2):5–43.
- 498 Duvenaud, D., Lloyd, J. R., Grosse, R., Tenenbaum, J., and Ghahramani, Z. (2013). Structure
499 discovery in nonparametric regression through compositional kernel search. In *Proceedings of
500 the 30th International Conference on Machine Learning (ICML-13)*, pages 1166–1174.
- 501 Ghahramani, Z. (2015). Probabilistic machine learning and artificial intelligence. *Nature*,
502 521(7553):452–459.
- 503 Lloyd, J. R., Duvenaud, D., Grosse, R., Tenenbaum, J., and Ghahramani, Z. (2014). Automatic
504 construction and natural-language description of nonparametric regression models. In *Twenty-
505 Eighth AAAI Conference on Artificial Intelligence*.
- 506 Mansinghka, V. K., Selsam, D., and Perov, Y. (2014). Venture: a higher-order probabilistic pro-
507 gramming platform with programmable inference. *arXiv preprint arXiv:1404.0099*.
- 508 Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953). Equation
509 of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6):1087–
510 1092.
- 511 Neal, R. M. (1997). Monte carlo implementation of gaussian process models for bayesian regression
512 and classification. *arXiv preprint physics/9701026*.
- 513 Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning (Adap-
514 tive Computation and Machine Learning)*. The MIT Press.
- 515 Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical bayesian optimization of machine
516 learning algorithms. In *Advances in Neural Information Processing Systems*, pages 2951–2959.
- 517 Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view
518 of the evidence of two samples. *Biometrika*, pages 285–294.
- 519

520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539