

How to Train Your Program

A probabilistic programming pattern for Bayesian learning from data

DAVID TOLPIN, Ben-Gurion University of the Negev & PUB+, Israel

We present a Bayesian approach to machine learning with probabilistic programs. In our approach, training on available data is implemented as inference on a hierarchical model. The posterior distribution of model parameters is then used to *stochastically condition* a complementary model, such that inference on new data yields the same posterior distribution of latent parameters corresponding to the new data as inference on a hierarchical model on the combination of both previously available and new data, at a lower computation cost. We frame the approach as a design pattern of probabilistic programming referred to herein as ‘stump and fungus’, and illustrate realization of the pattern on a didactic case study.

1 INTRODUCTION

The ultimate Bayesian approach to learning from data is embodied by hierarchical models [Gelman et al. 2013; Goodman et al. 2016; McElreath 2020]. In a hierarchical model, each observation or a group of observations y_i corresponding to a single item in the data set is conditioned on a parameter θ_i , and all parameters are conditioned on a hyperparameter τ :

$$\begin{aligned}\tau &\sim H \\ \theta_i &\sim D(\tau) \\ y_i &\sim F(\theta_i)\end{aligned}\tag{1}$$

A hierarchical model can be thought of as a way of inferring, or ‘learning’, the prior of each θ_i from all observations in the data set. Consider the following example problem: multiple boxes are randomly filled by K marbles from a bag containing a mixture of blue and white marbles. We are presented by a few draws with replacement from each of the boxes, y_{ij} being the j th draw from the i th box; our goal is to infer the number of blue marbles θ_i in each box. Intuitively, since the boxes are filled from the same bag, the posterior distribution of θ_i should account both for draws from the i th box and, indirectly, for draws from all other boxes. This is formalized by the following hierarchical model:

$$\begin{aligned}\tau &\sim \text{Beta}(1, 1) \\ \theta_i &\sim \text{Binomial}(K, \tau) \\ y_{ij} &\sim \text{Bernoulli}\left(\frac{\theta_i}{K}\right)\end{aligned}\tag{2}$$

Model (2) learns from the data in the sense that inference for each box is influenced by draws from all boxes. However, learning from *training data* to improve inference on *future data* with a hierarchical model is computationally inefficient — if a new box is presented, one has to add observations of the new box to the previously available data and re-run inference on the extended data set. Inference performance can be improved by employing data subsampling [Bardenet et al. 2014, 2017; Korattikara et al. 2014; Maclaurin and Adams 2014; Quiroz et al. 2018], but the whole training data set still needs to be kept and made accessible to the inference algorithm. A hierarchical model cannot ‘compress’, or summarize, training data for efficient inference on future observations.

An alternative approach, known as *empirical Bayes* [Casella 1985; Robbins 1951, 1992], consists in adjusting the hyperprior based on the training data; e.g. by fixing τ at a likely value, or by replacing H in (1) with a suitable approximation of the posterior distribution of τ . However, empirical Bayes

is not Bayesian. While practical efficiency of empirical Bayes was demonstrated in a number of settings [Robbins 1992], in other settings empirical Bayes may result in a critically misspecified model and overconfident or biased inference outcomes.

In this work, we propose an approach to learning from data in probabilistic programs which is both Bayesian in nature and computationally efficient. First, we state the problem of learning from data in the context of Bayesian generative models (Section 2). Then, we introduce the approach and discuss its implementation (Section 3). As an illustration, we apply the approach to a didactic example based on a case study of tumor incidence in rats (Section 4). Finally, we conclude with a discussion of related work and further research (Section 5).

2 PROBLEM: LEARNING FROM DATA

The challenge we tackle here is re-using inference outcomes on the training data set for inference on new data. Formally, population \mathcal{Y} is a set of sets $y_i \in Y$ of observations $y_{ij} \in y_i$. Members of each y_i are assumed to be drawn from a known distribution F with unobserved parameter θ_i , $y_{ij} \sim F(\theta_i)$. θ_i are assumed to be drawn from a common distribution H . Our goal is to devise a scheme that, given a subset $Y \subset \mathcal{Y}$, the *training set*, infers the posterior distribution of $\theta_k|Y, y_k$ for any $y_k \in \mathcal{Y}$ in a shorter amortized time than running inference on a hierarchical model $Y \cup \{y_k\}$. By amortized time we mean here average time per y_k , $k \in 1 : K$ as $K \rightarrow \infty$.

In other words, we look for a scheme that works in two stages. At the first stage, inference is performed on the training set Y only. At the second stage, the inference outcome of the first stage is used, together with y_k , to infer $\theta_k|Y, y_k$. We anticipate a scheme that ‘compresses’ the training set at the first stage, resulting in a shorter running time of the second stage. Such scheme bears similarity to the conventional machine learning paradigm: an expensive computation on the training data results in shorter running times on new data.

3 MAIN IDEA: STUMP AND FUNGUS

In quest of devising such a scheme, we make two observations which eventually help us arrive at a satisfactory solution:

- (1) In Bayesian modelling, information about data is usually conveyed through conditioning of the model on various aspects of the data.
- (2) In a hierarchical model, influence of the i th group of observations on the hyperparameters τ and, consequently, on other groups, passes exclusively through the group parameters θ_i .

If, instead of conditioning on training data y_i , we could condition on parameters θ_i corresponding to the training data, then we could perform inference on new data item y_k at a lower time and space complexity. Continuing the well known analogy between a hierarchical model and a tree, with the hyperparameter τ at the root and observations y_{ij} in the leaves, we can liken a model which receives all θ_i of the training data and new data item y_k as a *stump* (the hierarchical model with the trunk cut off just after the hyperparameters) and a *fungus* growing on the stump – the new data item. The problem is, of course, that we infer **distributions**, rather than fixed values, of θ_i , and the model must be, somewhat unconventionally, conditioned on the distributions of θ_i .



However, a recently introduced notion of stochastic conditioning [Tolpin et al. 2021] makes conditioning on distributions of θ_i possible, both theoretically and in the practical case when the

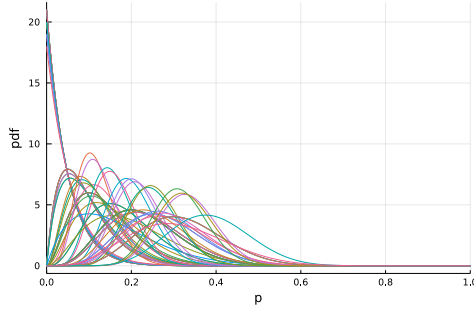


Fig. 1. Separate model

posteriors of θ_i are approximated using Monte Carlo samples. Moreover, conditioning the model both *stochastically* on the posterior distributions of θ_i on training data and *deterministically* on new data y_k yields the same posterior distribution of θ_k as inference on the full hierarchical model. Based on this, we propose the ‘stump-and-fungus’ pattern for learning from data in probabilistic programs:

- Training is accomplished through inference on a hierarchical model, in the usual way.
- Training outcomes are summarized as a collection of samples $\tilde{\theta}$, representing the mixture distribution of θ_i of all groups.
- For inference on new data item y , a stump-and-fungus model is employed:

$$\begin{array}{c}
 \tilde{\theta} \sim \text{Hierarchical}(Y) \\
 \hline
 \tau \sim H \\
 \tilde{\theta}, \theta | \tau \sim D(\tau) \\
 y | \theta \sim F(\theta)
 \end{array} \tag{3}$$

Although two models — hierarchical and stump-and-fungus — are involved in the pattern, the models are in fact two roles fulfilled by the same generative model, combining stochastic conditioning on training data and deterministic conditioning on new data (consisting potentially of multiple data items). This preserves a common practice in machine learning in which the same model is used for both training and inference.

4 CASE STUDY: TUMOR INCIDENCE IN RATS

In this case study, based on [Tarone 1982] and discussed in [Gelman et al. 2013, Chapter 5], data on tumor incidence in rats in $N = 70$ laboratory experiments is used to infer tumor incidence based on outcomes of yet another experiment. A different number of rats n_i was involved in each experiment, and the number of tumor cases y_i was reported. The posteriors of independent models for all experiments are shown in Figure 1. A schoolbook solution for the problem is to perform inference on a hierarchical model:

$$\begin{array}{c}
 \alpha, \beta \sim \text{Uniform}(0, \infty) \\
 p_i | \alpha, \beta \sim \text{Beta}(\alpha, \beta) \\
 y_i | p_i \sim \text{Binomial}(n_i, p_i)
 \end{array} \tag{4}$$

Inference on Model (4) can be performed efficiently thanks to summarization of n observations from Bernoulli(p) as a single observation from Binomial(n, p). In general however, the use of a

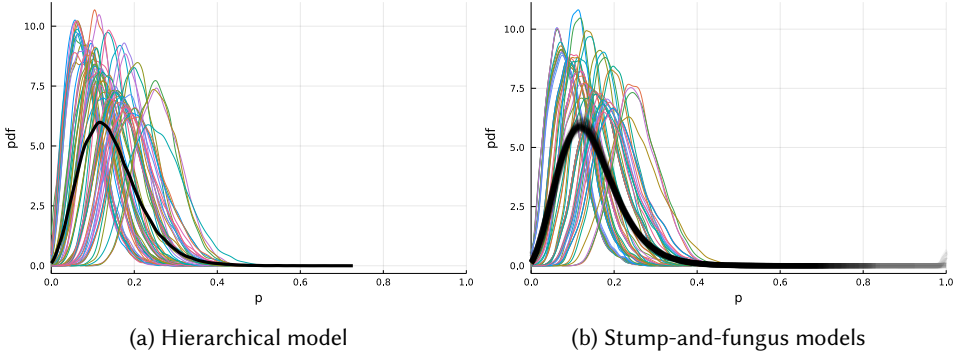


Fig. 2. Hierarchical vs. stump-and-fungus model. The posteriors obtained via either method are the same, except for small discrepancies apparently caused by finite sample size approximation. Colorful lines are posterior distributions of p for each of 71 experiments. Black lines are posterior distributions of p for a future experiment.

hierarchical model would require carrying all observations of all previous experiments for learned inference on findings of a new experiment.

The stump-and-fungus pattern is straightforwardly applicable to the problem:

$$\begin{aligned}
 & \frac{p_{-k} \sim \text{Hierarchical}(n_{-k}, y_{-k})}{\alpha, \beta \sim \text{Uniform}(0, \infty)} \\
 & p | \alpha, \beta \sim \text{Beta}(\alpha, \beta) \\
 & y_k | p_k \sim \text{Binomial}(n_k, p_k)
 \end{aligned} \tag{5}$$

In (5), p_{-k} are stochastically observable from the posterior of hierarchical model (4) conditioned on all but the k th experiment. Figure 2 shows the posterior distributions for p inferred on the hierarchical model (Figure 2a) and through $N + 1$ applications of stump-and-fungus (similar to leave-one-out cross-validation). The Infergo [Tolpin 2019] source code of the model is provided in the appendix. The inference was performed with HMC [Neal 2011] on the hierarchical model (4) and *stochastic gradient* HMC [Chen et al. 2014] on the stump-and-fungus model (5). 1000 samples were used to visualize the posterior. One can see that the posteriors obtained via either method appear to be the same, except for small discrepancies apparently caused by finite sample size approximation. The data and code for the case study are available at <https://bitbucket.org/dtolpin/h2typ-studies>.

5 DISCUSSION

We presented a probabilistic programming pattern for Bayesian learning from data. The importance of learning from data is well appreciated in probabilistic programming. Along with empirical Bayes, applicable to probabilistic programming as well as to Bayesian generative models in general, probabilistic-programming specific approaches were proposed. One possibility is to induce a probabilistic program suited for a particular data set [Hwang et al. 2011; Liang et al. 2010; Perov 2018; Perov and Wood 2014]. A related but different research direction is inference compilation [Baydin et al. 2019; Le et al. 2017], where the cost of inference is amortized through learning proposal distributions from data. Another line of research is concerned by speeding up inference algorithms by tuning them based on training data [Eslami et al. 2014; Mansinghka et al. 2018]. Our approach to learning from data in probabilistic programs is different in that it does not require any particular

implementation of probabilistic programming to be used, nor introspection into the structure of probabilistic programs or inference algorithms. Instead, the approach uses inference in ubiquitously adopted hierarchical models for training, and conditioning on observations for incorporation of training outcomes in inference.

ACKNOWLEDGMENTS

We thank PUB+ for supporting development of *Infergo*.

REFERENCES

- Rémi Bardenet, Arnaud Doucet, and Chris Holmes. 2014. Towards scaling up Markov chain Monte Carlo: an adaptive subsampling approach. In *Proceedings of the 31st International Conference on Machine Learning*. 405–413.
- Rémi Bardenet, Arnaud Doucet, and Chris Holmes. 2017. On Markov chain Monte Carlo methods for tall data. *Journal of Machine Learning Research* 18, 47 (2017), 1–43.
- Atilim Gunes Baydin, Lei Shao, W. Bhimji, L. Heinrich, Lawrence Meadows, Jialin Liu, Andreas Munk, Saeid Naderiparizi, Bradley Gram-Hansen, Gilles Louppe, Mingfei Ma, X. Zhao, P. Torr, V. Lee, K. Cranmer, Prabhat, and Frank D. Wood. 2019. Etalumis: bringing probabilistic programming to scientific simulators at scale. *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* (2019).
- George Casella. 1985. An Introduction to Empirical Bayes Data Analysis. *The American Statistician* 39, 2 (1985), 83–87.
- Tianqi Chen, Emily B. Fox, and Carlos Guestrin. 2014. Stochastic Gradient Hamiltonian Monte Carlo. In *Proceedings of the 31st International Conference on Machine Learning*. 1683–1691.
- Ali Eslami, Daniel Tarlow, Pushmeet Kohli, and John Winn. 2014. Just-In-Time Learning for Fast and Flexible Inference. In *NIPS'14 Proceedings of the 27th International Conference on Neural Information Processing Systems (nips'14 proceedings of the 27th international conference on neural information processing systems ed.)*. MIT Press Cambridge, 154–162.
- A. Gelman, J.B. Carlin, H.S. Stern, D.B. Dunson, A. Vehtari, and D.B. Rubin. 2013. *Bayesian Data Analysis*. CRC Press.
- Noah D Goodman, Joshua B. Tenenbaum, and the ProbMods contributors. 2016. *Probabilistic Models of Cognition* (second ed.). <http://probmods.org/v2> electronic; retrieved: 2019-4-29.
- I. Hwang, Andreas Stuhlmüller, and Noah D. Goodman. 2011. Inducing Probabilistic Programs by Bayesian Program Merging. *ArXiv abs/1110.5667* (2011).
- Anoop Korattikara, Yutian Chen, and Max Welling. 2014. Austerity in MCMC Land: Cutting the Metropolis-Hastings Budget. In *Proceedings of the 31st International Conference on Machine Learning*. 181–189.
- Tuan Anh Le, Atilim Gunes Baydin, and Frank Wood. 2017. Inference Compilation and Universal Probabilistic Programming. In *Proceedings of the 34th International Conference on Machine Learning*. 1338–1348.
- Percy Liang, Michael I. Jordan, and Dan Klein. 2010. Learning Programs: A Hierarchical Bayesian Approach. In *Proceedings of the 27th International Conference on Machine Learning*. Omnipress, 639–646.
- Dougal Maclaurin and Ryan P. Adams. 2014. Firefly Monte Carlo: Exact MCMC with subsets of data. In *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence*. 543–552.
- Vikash K. Mansinghka, Ulrich Schaehtle, Shivam Handa, Alexey Radul, Yutian Chen, and Martin Rinard. 2018. Probabilistic Programming with Programmable Inference. *SIGPLAN Not.* 53, 4 (June 2018), 603–616.
- Richard McElreath. 2020. *Statistical Rethinking* (2nd ed.). CRC Press.
- Radford M. Neal. 2011. MCMC using Hamiltonian dynamics. Chapter 5 of the *Handbook of Markov Chain Monte Carlo*.
- Yura Perov. 2018. Inference Over Programs That Make Predictions. *arXiv:arXiv:1810.01190*
- Yura N. Perov and Frank D. Wood. 2014. Learning Probabilistic Programs. *arXiv:arXiv:1407.2646*
- Matias Quiroz, Mattias Villani, Robert Kohn, Minh-Ngoc Tran, and Khue-Dung Dang. 2018. Subsampling MCMC — an Introduction for the Survey Statistician. *Sankhya A* 80, 1 (01 Dec 2018), 33–69.
- Herbert Robbins. 1951. Asymptotically Subminimax Solutions of Compound Statistical Decision Problems. In *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*. University of California Press, Berkeley, Calif., 131–149. <https://euclid-web1.library.cornell.edu/euclid.bsmsp/1200500224>
- Herbert Robbins. 1992. *An Empirical Bayes Approach to Statistics*. Springer New York, New York, NY, 388–394.
- Robert Tarone. 1982. The Use of Historical Control Information in Testing for a Trend in Proportions. *Biometrics* 38, 1 (1982), 215–220.
- David Tolpin. 2019. Deployable Probabilistic Programming. In *Proceedings of the 2019 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software*. 1–16.
- David Tolpin, Yuan Zhou, Tom Rainforth, and Hongseok Yang. 2021. Probabilistic Programs with Stochastic Conditioning. *arXiv:arXiv:2010.00282*

A INFERGO STUMP-AND-FUNGUS MODEL

```

1  package model
2
3  import (
4      . "bitbucket.org/dtolpin/infergo/dist"
5      "bitbucket.org/dtolpin/infergo/mathx"
6      "math"
7  )
8
9  // Model is shared by both the tree and the fungus. When K is 0,
10 // the model reduces to a conventional hierarchical model.
11 // Otherwise, K is amount of evidence flowing from the tree
12 // to the fungus, with stochastic conditioning fed through P.
13 type Model struct {
14     P      <-chan float64
15     K      int
16     Y, N []float64
17 }
18
19 func (m *Model) Observe(x []float64) (lp float64) {
20     // Hyperparameters
21     alpha := math.Exp(x[0])
22     beta := math.Exp(x[1])
23     lp += x[0] + x[1]
24
25     // Parameters
26     x = x[2:]
27     // x is a vector of odds, and is regularized through
28     // imposing the Cauchy prior.
29     lp += Cauchy.Logps(0, 10, x...)
30     p := make([]float64, len(x))
31     for i := range p {
32         p[i] = mathx.Sigm(x[i])
33         lp += mathx.LogDSigm(x[i])
34     }
35
36     // Parameters given hyperparameters
37     lp += Beta.Logps(alpha, beta, p...)
38     // Stochastic conditioning given hyperparameters
39     for i := 0; i != m.K; i++ {
40         lp += Beta.Logp(alpha, beta, <-m.P)
41     }
42
43     // Observations given parameters
44     for i := range m.Y {
45         lp += Binomial.Logp(m.N[i], p[i], m.Y[i])
46     }
47
48     return lp
49 }

```