Patterns of Scalable Bayesian Inference

Elaine Angelino^{*} UC Berkeley elaine@eecs.berkeley.edu Matthew James Johnson* Harvard University mattjj@csail.mit.edu

Ryan P. Adams Harvard University and Twitter rpa@seas.harvard.edu

*Authors contributed equally

Contents

1	Introduction						
	1.1	Why be Bayesian with big data?	3				
	1.2	The fidelity of approximate integration	5				
	1.3	Outline	5				
2	Background						
	2.1	Exponential families	8				
	2.2	Markov Chain Monte Carlo inference	12				
		2.2.1 Bias and variance of estimators	12				
		2.2.2 Monte Carlo estimates from indepdent samples	13				
		2.2.3 Markov chains	15				
		2.2.4 Markov chain Monte Carlo (MCMC)	16				
		2.2.5 Metropolis-Hastings (MH) sampling	22				
		2.2.6 Gibbs sampling	24				
	2.3	Mean field variational inference	25				
	2.4	Stochastic gradient optimization	29				
3	мс	MC with data subsets	32				
	3.1 Factoring the joint density						
		3.2.1 $$ An approximate MH test based on a data subset $$.	34				

		3.2.2	Approximate MH with an adaptive stopping rule .	35				
		3.2.3	Using a <i>t</i> -statistic hypothesis test	37				
		3.2.4	Using concentration inequalities	39				
		3.2.5	Error bounds on the stationary distribution	41				
	3.3	Sub-se	electing data via a lower bound on the likelihood	44				
	3.4	Stocha	astic gradients of the log joint density	46				
	3.5	Summ	ary	49				
	3.6	Discus	sion	51				
4	Para	allel and	d distributed MCMC	54				
	4.1	Paralle	lizing standard MCMC algorithms	55				
		4.1.1	Conditional independence and graph structure	55				
		4.1.2	Speculative execution and prefetching	56				
	4.2	Defini	ng new data-parallel dynamics	58				
		4.2.1	Aggregating from subposteriors	60				
			Embarrassingly parallel consensus of subposteriors .	61				
			Weighted averaging of subposterior samples	63				
			Subposterior density estimation	64				
			Weierstrass samplers	67				
		4.2.2	Hogwild Gibbs	72				
			Defining Hogwild Gibbs variants	73				
			Theoretical analysis	75				
	4.3	Summ	ary	77				
	4.4	Discus	sion	81				
5	Scaling variational mean field algorithms 83							
	5.1	Stocha	astic optimization and variational inference	84				
		5.1.1	SVI for complete-data conjugate models	85				
		5.1.2	Stochastic gradients with general nonconjugate					
			models	89				
		5.1.3	Exploiting reparameterization for some nonconju-					
		_	gate models	92				
	5.2	Strean	ning variational Bayes (SVB)	94				
	5.3	Summary						
	5.4	Discus	sion	98				

6	Challenges and questions	99
Re	eferences	107

iv

Abstract

Datasets are growing not just in size but in complexity, creating a demand for rich models and quantification of uncertainty. Bayesian methods are an excellent fit for this demand, but scaling Bayesian inference is a challenge. In response to this challenge, there has been considerable recent work based on varying assumptions about model structure, underlying computational resources, and the importance of asymptotic correctness. As a result, there is a zoo of ideas with few clear overarching principles.

In this paper, we seek to identify unifying principles, patterns, and intuitions for scaling Bayesian inference. We review existing work on utilizing modern computing resources with both MCMC and variational approximation techniques. From this taxonomy of ideas, we characterize the general principles that have proven successful for designing scalable inference procedures and comment on the path forward.

1

Introduction

We have entered a new era of scientific discovery, in which computational insights are being integrated with large-scale statistical data analysis to enable researchers to ask both grander and more subtle questions about our natural world. This viewpoint asserts that we need not be limited to the narrow hypotheses that can be framed by traditional small-scale analysis techniques. Supporting new kinds of data-driven queries, however, requires that new methods be developed for statistical inference that can *scale up* along multiple axes — more samples, more dimensions, and greater model complexity — as well as *scale out* by taking advantage of modern parallel compute environments.

There are a variety of methodological frameworks for performing statistical inference, e.g., performing estimation and evaluating hypotheses; here we are concerned with the Bayesian formalism. In the Bayesian setting, queries about structure in data are framed as interrogations of the posterior distribution over parameters, missing data, and other unknowns; these unobserved quantities are treated as random variables. By conditioning on the data, the Bayesian hopes to not only perform point estimation, but also to understand the uncertainties associated with those estimates.

1.1. Why be Bayesian with big data?

Accounting for uncertainty is central to Bayesian analysis, and so the computations associated with most common tasks – e.g., estimation, prediction, evaluation of hypotheses – are typically integrations. In some situations, it is possible to perform such integrations exactly, either by taking advantage of conjugate structure in the prior-likelihood pair, or by using dynamic programming when the dependencies between random variables are appropriately simple. Unfortunately, most real-world analysis problems are not amenable to these exact inference procedures and so most of the interest in Bayesian computation focuses on better methods of approximate inference.

There are two dominant paradigms for approximate inference in Bayesian models: Monte Carlo sampling methods and variational approximations. The Monte Carlo approach observes that integrations performed to query posterior distributions can be framed as expectations, and thus estimated with samples; such samples are most often generated via simulation from carefully designed Markov chains. Variational inference seeks to compute these integrals by approximating the posterior distribution with a more tractable alternative, where identification of the best approximation can then be performed using powerful optimization techniques.

In this paper, we examine how these techniques can be scaled up to larger problems and scaled out across parallel computational resources. This is not intended to be an exhaustive survey of a rapidly-evolving area of research; rather, we seek to identify the main ideas and themes that are emerging in this area, and articulate what we believe are some of the significant open questions and challenges.

1.1 Why be Bayesian with big data?

The Bayesian paradigm is fundamentally about integration: integration computes posterior estimates and measures of uncertainty, eliminates nuisance variables or missing data, and averages models to compute predictions or perform model comparison. While some statistical methods, such as MAP estimation, can be described from a Bayesian perspective, in which case the prior might serves as a regularizer in an optimization problem, such methods are not inherently or exclusively Bayesian. Posterior integration is the distinguishing characteristic of Bayesian statistics, and so a defense of Bayesian ideas in the big data regime rests on the utility of integration.

But from a classical perspective, the big data setting might seem to be precisely where integration isn't important: as the dataset grows, shouldn't the posterior distribution concentrate towards a point mass? If big data means we end up making predictions with such concentrated posteriors, why not focus on point estimation and avoid the specification of priors and the burden of approximate integration?

These objections certainly apply to settings where the number of parameters is small and fixed ("tall data"). However, many models of interest have many parameters ("wide data"), or indeed have a number of parameters that grows along with the amount of data.

For example, an Internet company making inferences about its users' viewing and buying habits may have terabytes of data in total but only a few observations for its newest customers, the ones most important to impress with personalized recommendations. Moreover, it may wish to adapt its model in an online way as data arrive, a task that benefits from calibrated posterior uncertainties [Stern et al., 2009]. As another example, consider a healthcare company. As its dataset grows, it might hope to make more detailed and complex inferences about populations while also making careful predictions with calibrated uncertainty for each patient, even in the presence of massive missing data [Lawrence, 2015]. These scaling issues also arise in astronomy, where hundreds of billions of light sources, such as stars, galaxies, and quasars, each have latent variables that must be estimated from very weak observations, and are coupled in a large hierarchical model [Regier et al., 2015]. In Microsoft Bing's sponsored search advertising, predictive probabilities inform the pricing in the keyword auction mechanism. This problem nevertheless must be solved at scale, with tens of millions of impressions per hour [Graepel et al., 2010].

These are the regimes where big data can be small [Lawrence, 2015] and the number and complexity of statistical hypotheses grows with the data. The Bayesian methods we survey in this paper may provide

solutions to these challenges.

1.2 The fidelity of approximate integration

Bayesian inference may be important in some modern big data regimes, but exact integration in general is computationally out of reach. While decades of research in Bayesian inference in both statistics and machine learning have produced many powerful approximate inference algorithms, the big data setting poses some new challenges. Iterative algorithms that read the entire dataset before making each update become prohibitively expensive. Sequential computation that cannot leverage parallel and distributed computing resources is at a significant and growing disadvantage. Insisting on zero asymptotic bias from Monte Carlo estimates of expectations may leave us swamped in errors from high variance [Korattikara et al., 2014] or transient bias.

These challenges, and the tradeoffs that may be necessary to address them, can be viewed in terms of how accurate the integration in our approximate inference algorithms must be. Markov chain Monte Carlo (MCMC) algorithms that admit the exact posterior as a stationary distribution may be the gold standard for generically estimating posterior expectations, but if standard MCMC algorithms become intractable in the big data regime we must find alternatives and understand their tradeoffs. Indeed, someone using Bayesian methods for machine learning may be less constrained than a classical Bayesian statistician: if the ultimate goal is to form predictions that perform well according to a specific loss function, computational gains at the expense of the internal posterior representation may be worthwhile. The methods studied here cover a range of such approximate integration tradeoffs.

1.3 Outline

The remainder of this review is organized as five chapters. In Chapter 2, we provide relevant background material on exponential families, MCMC inference, mean field variational inference, and stochastic gradient optimization. The next three chapters survey recent algorithmic ideas for scaling Bayesian inference, highlighting theoretical results where possible. Each of these central technical chapters ends with a summary and discussion, identifying emergent themes and patterns as well as open questions. Chapters 3 and 4 focus on MCMC algorithms, which are inherently serial and often slow to converge; the algorithms in the first of these use various forms of data subsampling to scale up serial MCMC and in the second use a diverse array of strategies to scale out on parallel resources. In Chapter 5 we discuss two recent techniques for scaling variational mean field algorithms. Both process data in minibatches: the first applies stochastic gradient optimization methods and the second is based on incremental posterior updating. Finally, in Chapter 6 we provide an overarching discussion of the ideas we survey, focusing on challenges and open questions in large-scale Bayesian inference.

2

Background

In this chapter we summarize background material on which the ideas in subsequent chapters are based. This chapter also serves to fix some common notation. Throughout the chapter, we avoid measure-theoretic definitions and instead assume that any density exists with respect either to Lebesgue measure or counting measure, depending on its context.

First, we cover some relevant aspects of exponential families. Second, we cover the foundations of Markov chain Monte Carlo (MCMC) algorithms, which are the workhorses of Bayesian statistics and are common in Bayesian machine learning. Indeed, the algorithms discussed in Chapters 3 and 4 are either MCMC algorithms or aim to approximate MCMC algorithms. Next, we describe the basics of mean field variational inference and stochastic gradient optimization, both of which are used extensively in Chapter 5. Finally, we close the chapter with notes on computational architectures and useful notions for measuring performance.

2.1 Exponential families

Exponential families of densities play a key role in Bayesian analysis and many practical Bayesian methods. In particular, likelihoods that are exponential families yield natural conjugate prior families, which provide analytical and computational advantages in both MCMC and variational inference algorithms. Exponential families are also particularly relevant in the context of large datasets: in a precise sense, they are the only families of densities which admit a finite-dimensional sufficient statistic. Thus only exponential families allow arbitrarily large amounts of data to be summarized with a fixed-size description.

In this section we give basic definitions, notation, and results concerning exponential families. For perspectives from convex analysis see Wainwright and Jordan [2008], and for perspectives from differential geometry see Amari and Nagaoka [2007].

Exponential families are defined in terms of densities with respect to some underlying σ -finite measure, which we denote ν .

Definition 2.1 (Exponential family). We say a parameterized family of densities $\{p(\cdot | \theta) : \theta \in \Theta\}$ is an *exponential family* if each density can be written as

$$p(x|\theta) = h(x) \exp\{\langle \eta(\theta), t(x) \rangle - \log Z(\eta(\theta))\}$$
(2.1)

where $\langle \cdot, \cdot \rangle$ is an inner product on a finite-dimensional real vector space. We call $\eta(\theta)$ the *natural parameter* vector, t(x) the *statistic* vector, $h(\cdot)$ the *base density*, and

$$\log Z(\eta) \triangleq \log \int e^{\langle \eta, t(x) \rangle} h(x) \nu(dx)$$
(2.2)

the log partition function.

We restrict our attention to families for which the support of the density does not depend on θ . When $\eta(\theta) = \theta$ we say the family is written in *natural parameters* or *natural coordinates*, which we denote by writing $p(x|\eta)$. We say a family is *regular* if Θ is open, and *minimal* if there is no nonzero *a* such that $\langle a, t(x) \rangle$ is equal to a constant (ν -a.e.).

2.1. Exponential families

The statistic t is *sufficient* in the sense of the Fisher-Neyman Factorization Theorem [Keener, 2010, Theorem 3.6] by construction

$$p(x|\theta) \propto h(x) \exp\{\langle \eta(\theta), t(x) \},\$$

and hence t(x) contains all the information about x that is relevant for the parameter θ . In the context of Bayesian analysis, in which θ is a random variable, this definition of sufficiency is equivalent to the conditional independence statement $\theta \perp X \mid t(X)$. The Koopman-Pitman-Darmois Theorem shows that exponential families are the only families which provide this powerful summarization property, under some mild regularity conditions [Hipp, 1974].

Exponential families have many convenient analytical and computational properties. In particular, differentiating the log partition function $\log Z$ generates cumulants:

Proposition 2.1 (Mean mapping and cumulants). For a regular exponential family of densities of the form (2.1) with $X \sim p(\cdot | \eta)$, we have $\nabla \log Z : \Theta \to \mathcal{M}$ and

$$\nabla \log Z(\eta) = \mathbb{E}[t(X)] \tag{2.3}$$

and writing $\mu \triangleq \mathbb{E}[t(X)]$ we have

$$\nabla^2 \log Z(\eta) = \mathbb{E}[t(X)t(X)^{\mathsf{T}}] - \mu\mu^{\mathsf{T}}.$$
(2.4)

More generally, the moment generating function of t(X) can be written

$$\mathcal{M}_{t(X)}(s) \triangleq \mathbb{E}[e^{\langle s, t(X) \rangle}] = e^{\log Z(\eta+s) - \log Z(\eta)}.$$
 (2.5)

and so derivatives of $\log Z$ give *cumulants* of t(X), where the first cumulant is the mean and the second and third cumulants are the second and third central moments, respectively.

Proof. To show 2.3, we write

$$\nabla_{\eta} \log Z(\eta) = \nabla_{\eta} \log \int e^{\langle \eta, t(x) \rangle} h(x) \nu(dx)$$
(2.6)

$$= \frac{1}{\int e^{\langle \eta, t(x) \rangle} h(x)\nu(dx)} \int t(x) e^{\langle \eta, t(x) \rangle} h(x)\nu(dx) \qquad (2.7)$$

$$= \int t(x)p(x|\eta)\nu(dx)$$
(2.8)

$$=\mathbb{E}[t(X)].$$
(2.9)

To derive the form of the moment generating function, we write

$$\mathbb{E}[e^{\langle s,t(X)\rangle}] = \int e^{\langle s,t(x)\rangle} p(x)\nu(dx)$$
(2.10)

$$= \int e^{\langle s,t(x)\rangle} e^{\langle \eta,t(x)\rangle - \log Z(\eta)} h(x)\nu(dx)$$
(2.11)
$$\log Z(\eta+s) - \log Z(\eta)$$
(2.12)

$$=e^{i0g\,2\,(\eta+0)-i0g\,2\,(\eta)}.$$
(2.12)

For members of an exponential family, many quantities can be expressed generically in terms of the natural parameter, expected statistics under that parameter, and the log partition function.

Proposition 2.2 (Score and Fisher information). For a regular exponential family in natural coordinates, with $X \sim p(\cdot | \eta)$ and $\mu(\eta) \triangleq \mathbb{E}[t(X)]$ we have

1. When the family is regular, the *score* with respect to the natural parameter is

$$v(x,\eta) \triangleq \nabla_{\eta} \log p(x|\eta) = t(x) - \mu(\eta) \tag{2.13}$$

2. When the family is regular, the *Fisher information* with respect to the natural parameter is

$$\mathcal{I}(\eta) \triangleq \mathbb{E}[v(X,\eta)v(X,\eta)^{\mathsf{T}}] = \nabla^2 \log Z(\eta).$$
 (2.14)

Proof. Each follows from (2.1) and Proposition 2.1.

Below, we define a notion of conjugacy for pairs of families of distributions. Conjugate families are especially useful for Bayesian analysis and algorithms.

Definition 2.2. A parameterized (not necessarily exponential) family of densities $\mathcal{F} = \{p(\cdot|\alpha) : \alpha \in \mathcal{A}\}$ is *conjugate* for a likelihood function $p(x|\cdot)$ if for every density $p(\cdot|\alpha)$ in \mathcal{F} the posterior distribution

$$p(\theta|\alpha') \propto p(\theta|\alpha)p(x|\theta)$$
 (2.15)

also belongs to \mathcal{F} , for some $\alpha' = \alpha'(x, \alpha)$ that may depend on x and α .

2.1. Exponential families

Conjugate pairs are particularly useful in Bayesian analysis because if we have a prior family $p(\theta|\alpha)$ and we observe data generated according to a likelihood $p(x|\theta)$ then the posterior $p(\theta|x, \alpha)$ is in the same family as the prior. In the context of Bayesian updating, we call α the hyperparameter and α' the posterior hyperparameter.

Given a regular exponential family likelihood, we can always define a conjugate prior, as shown in the next proposition.

Proposition 2.3. Given a regular exponential family

$$p_{X|\theta}(x|\theta) = h_X(x) \exp\{\langle \eta_X(\theta), t_X(x) \rangle - \log Z_X(\eta_X(\theta))\}$$
(2.16)

$$= h_X(x) \exp\{\langle (\eta_X(\theta), -\log Z_X(\eta(\theta))), (t_X(x), 1) \rangle\}$$
(2.17)

then if we define the statistic $t_{\theta}(\theta) \triangleq (\eta_X(\theta), -\log Z_X(\eta(\theta)))$ and an exponential family of densities with respect to that statistic as

$$p_{\theta|\alpha}(\theta|\alpha) = h_{\theta}(\theta) \exp\{\langle \eta_{\theta}(\alpha), t_{\theta}(\theta) \rangle - \log Z_{\theta}(\eta_{\theta}(\alpha))\}$$
(2.18)

then the pair $(p_{\theta|\alpha}, p_{X|\theta})$ is a conjugate pair of families with

$$p(\theta|\alpha)p(x|\theta) \propto h_{\theta}(\theta) \exp\{\langle \eta_{\theta}(\alpha) + (t_X(x), 1), t_{\theta}(\theta) \rangle\}$$
(2.19)

and hence we can write the posterior hyperparameter as

$$\alpha' = \eta_{\theta}^{-1}(\eta_{\theta}(\alpha) + (t_X(x), 1)).$$
(2.20)

When the prior family is parameterized with its natural parameter, we have $\eta' = \eta + (t_X(x), 1)$.

As a consequence of Proposition 2.3, if the prior family is written with natural parameters and we generate data $\{x_i\}_{i=1}^n$ according to the model

$$\theta \sim p_{\theta|\eta}(\ \cdot \ |\eta) \tag{2.21}$$

$$x_i | \theta \stackrel{\text{iid}}{\sim} p_{X|\theta}(\cdot | \theta) \quad i = 1, 2, \dots, n,$$
 (2.22)

where the notation $x_i \stackrel{\text{iid}}{\sim} p(\cdot)$ denotes that the random variables x_i are independently and identically distributed, then $p(\theta|\{x_i\}_{i=1}^n, \eta)$ has posterior hyperparameter $\eta' = \eta + (\sum_{i=1}^n t(x_i), n)$. Therefore any tractable computations in the prior, such as simulation or computing expectations, are shared by the posterior. Furthermore, for inferences about θ , the entire dataset can be summarized by the statistic $(\sum_{i=1}^n t(x_i), n)$.

2.2 Markov Chain Monte Carlo inference

Markov chain Monte Carlo (MCMC) is a class of algorithms for estimating expectations with respect to distributions. These distributions may be intractable, such as most posterior distributions arising in Bayesian inference. Given a target distribution, a standard MCMC algorithm proceeds by simulating an ergodic random walk that admits the target distribution as its stationary distribution. As we develop in the following subsections, by collecting samples from the simulated trajectory and forming Monte Carlo estimates, expectations of many functions can be approximated to arbitrary accuracy. Thus MCMC is employed when samples or expectations from a distribution cannot be obtained directly, as is often the case with complex, high-dimensional systems arising across disciplines, such as estimating bulk material properties from molecular dynamics physics simulations or performing inference in Bayesian probabilistic models.

In this section, we first review the two underlying ideas behind MCMC algorithms: Monte Carlo methods and Markov chains. First we define the bias and variance of estimators. Next, we introduce Monte Carlo estimators based on independent and identically distributed samples. We then describe how Monte Carlo estimates can be formed using mutually dependent samples generated by a Markov chain simulation. Finally, we introduce two general MCMC algorithms commonly applied to Bayesian posterior inference, the Metropolis-Hastings and Gibbs sampling algorithms. Our exposition here mostly follows the standard treatment, such as in Brooks et al. [2011, Chapter 1], Geyer [1992], and Robert and Casella [2004].

2.2.1 Bias and variance of estimators

Notions of bias and variance are fundamental to understanding and comparing estimator performance, and much of our discussion of MCMC methods is framed in these terms.

Consider using a scalar-valued random variable $\hat{\theta}$ to estimate a fixed scalar quantity of interest θ . The bias and variance of the estimator $\hat{\theta}$

are defined as

$$\operatorname{Bias}[\hat{\theta}] = \mathbb{E}[\hat{\theta} - \theta] \tag{2.23}$$

$$\operatorname{Var}[\hat{\theta}] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^{2}].$$
(2.24)

The mean squared error $\mathbb{E}[(\hat{\theta} - \theta)^2]$ can be decomposed in terms of the variance and the square of the bias:

$$\mathbb{E}[(\hat{\theta} - \theta)^2] = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2]$$
(2.25)

$$= \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2] + (\mathbb{E}[\hat{\theta}] - \theta)^2$$
(2.26)

$$= \operatorname{Var}[\hat{\theta}] + \operatorname{Bias}^{2}[\hat{\theta}] \tag{2.27}$$

This decomposition provides a basic language for evaluating estimators and thinking about tradeoffs. Among unbiased estimators, those with lower variance are generally preferrable. However, when an unbiased estimator has high variance, a biased estimator that achieves low variance can have a lower overall mean squared error.

As we describe in the following sections, a substantial amount of the study of Bayesian statistical computation has focused on algorithms that produce asymptotically unbiased estimates of posterior expectations, in which the bias due to initialization is transient and is washed out relatively quickly. In this setting, the error is typically considered to be dominated by the variance term, which can be made as small as desired by increasing computation time without bound. When computation becomes expensive as in the big data setting, errors under a realistic computational budget may in fact be dominated by variance, as observed by Korattikara et al. [2014], or, as we argue in Chapter 6, transient bias. Several of the new algorithms we examine in Chapters 3 and 4 aim to adjust this tradeoff by allowing some asymptotic bias while effectively reducing the variance and transient bias contributions through more efficient computation.

2.2.2 Monte Carlo estimates from indepdent samples

Let X be a random variable with $\mathbb{E}[X] = \mu < \infty$, and let $(X_i : i \in \mathbb{N})$ be a sequence of i.i.d. random variables each with the same distribution

as X. The Strong Law of Large Numbers (LLN) states that the sample average converges almost surely to the expectation μ as $n \to \infty$:

$$\mathbb{P}\left(\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} X_i = \mu\right) = 1.$$
(2.28)

This convergence immediately suggests the Monte Carlo method: to approximate the expectation of X, which to compute exactly may involve an intractable integral, one can use i.i.d. samples and compute a sample average. In addition, because for any measurable function f the sequence $(f(X_i) : i \in \mathbb{N})$ is also a sequence of i.i.d. random variables, we can form the Monte Carlo estimate

$$\mathbb{E}[f(X)] \approx \frac{1}{n} \sum_{i=1}^{n} f(X_i).$$
(2.29)

Monte Carlo estimates of this form are unbiased by construction, and so the quality of a Monte Carlo estimate can be evaluated in terms of its variance as a function of the number of samples n, which in turn can be understood with the Central Limit Theorem (CLT), at least in the asymptotic regime. If X is real-valued and has finite variance $\mathbb{E}[(X - \mu)^2] = \sigma^2 < \infty$, then the CLT states that the deviation $\frac{1}{n} \sum_{i=1}^n X_i - \mu$, rescaled appropriately, converges in distribution and is asymptotically normal:

$$\lim_{n \to \infty} \mathbb{P}\left(\frac{1}{\sqrt{n}} \sum_{i=1}^{n} (X_i - \mu) < \alpha\right) = \mathbb{P}(Z < \alpha)$$
(2.30)

where $Z \sim \mathcal{N}(0, \sigma^2)$. In particular, as *n* grows, the standard deviation of the sample average $\frac{1}{n} \sum_{i=1}^{n} X_i - \mu$ converges to zero at an asymptotic rate proportional to $\frac{1}{\sqrt{n}}$. More generally, for any real-valued measurable function *f*, the Monte Carlo standard error (MCSE) in the estimate (2.29) asymptotically scales as $\frac{1}{\sqrt{n}}$ regardless of the dimension of *X*.

Monte Carlo estimators effectively reduce the problem of computing expectations to the problem of generating samples. However, the preceding statements require the samples used in the Monte Carlo estimate to be independent, and independent samples can be computationally difficult to generate. Instead of relying on independent samples, Markov chain Monte Carlo algorithms compute estimates using mutually dependent samples generated by simulating a Markov chain.

2.2.3 Markov chains

Let \mathcal{X} be a discrete or continuous state space and let $x, x' \in \mathcal{X}$ denote states. A time-homogeneous Markov chain is a discrete-time stochastic process $(X_t : t \in \mathbb{N})$ governed by a transition operator $T(x \to x')$ that specifies the probability density of transitioning to a state x' from a given state x:

$$\mathbb{P}(X_{t+1} \in A \mid X_t = x) = \int_A T(x \to x') \, dx' \quad \forall t \in \mathbb{N}$$
(2.31)

for all measurable sets A. A Markov chain is memoryless in the sense that its future behavior depends only on the current state and is independent of its past history.

Given an initial density $\pi_0(x)$ for X_0 , a Markov chain evolves this density from one time point to the next through iterative application of the transition operator. We write the application of the transition operator to a density π_0 to yield a new density π_1 as

$$\pi_1(x') = (\pi_0 T)(x') = \int_{\mathcal{X}} T(x \to x') \pi_0(x) \, dx. \tag{2.32}$$

Writing T^t to denote t repeated applications of the transition operator T, the density of X_t induced by π_0 and T is then given by $\pi_t = \pi_0 T^t$.

Markov chain simulation follows this iterative definition by iteratively sampling the next state using the current state and the transition operator. That is, after first sampling X_0 from $\pi_0(\cdot)$, Markov chain simulation proceeds at time step t by sampling X_{t+1} according to the density $T(x_t \to \cdot)$ induced by the fixed sample x_t .

We are interested in Markov chains that converge in total variation to a unique stationary density $\pi(x)$ in the sense that

$$\lim_{t \to \infty} \|\pi_t - \pi\|_{\rm TV} = 0 \tag{2.33}$$

for any initial distribution π_0 , where $\|\cdot\|_{\text{TV}}$ denotes the total variation norm on densities:

$$\|p - q\|_{\rm TV} = \frac{1}{2} \int_{\mathcal{X}} |p(x) - q(x)| \, dx.$$
 (2.34)

For a transition operator $T(x \to x')$ to admit $\pi(x)$ as a stationary density, its application must leave $\pi(x)$ invariant:

$$\pi = \pi T. \tag{2.35}$$

For a discussion of general conditions that guarantee a Markov chain converges to a unique stationary distribution, i.e., that the chain is ergodic, see Meyn and Tweedie [2009].

In some cases it is easy to show that a transition operator has a particular unique stationary distribution. In particular, it is clear that π is the unique stationary distribution when a transition operator $T(x \to x')$ is *reversible* with respect to π , i.e., it satisfies the detailed balance condition with respect to a density $\pi(x)$,

$$T(x \to x')\pi(x) = T(x' \to x)\pi(x') \quad \forall x, x' \in \mathcal{X},$$
(2.36)

which is a pointwise condition over $\mathcal{X} \times \mathcal{X}$. Integrating over x on both sides gives:

$$\int_{\mathcal{X}} T(x \to x') \pi(x) \, dx = \int_{\mathcal{X}} T(x' \to x) \pi(x') \, dx$$
$$= \pi(x') \int_{\mathcal{X}} T(x' \to x) \, dx$$
$$= \pi(x'),$$

which is precisely the required condition from (2.35). We can interpret (2.36) as stating that, for a reversible Markov chain starting from its stationary distribution, any transition $x \to x'$ is equilibrated by the corresponding reverse transition $x' \to x$. Many MCMC methods are based on deriving reversible transition operators.

For a thorough introduction to Markov chains, see Robert and Casella [2004, Chapter 6] and Meyn and Tweedie [2009].

2.2.4 Markov chain Monte Carlo (MCMC)

Markov chain Monte Carlo (MCMC) methods simulate a Markov chain for which the stationary distribution is equal to a target distribution of interest, and use the simulated samples to form Monte Carlo estimates of expectations. That is, consider simulating a Markov chain

with unique stationary density $\pi(x)$, as in Section 2.2.3, and collecting its trajectory into a set of samples $\{X_i\}_{i=1}^n$. These collected samples can be used to form a Monte Carlo estimate for a function f of a random variable X with density $\pi(x)$ via

$$\mathbb{E}[f(X)] = \int_{\mathcal{X}} f(x)\pi(x) \, dx \approx \frac{1}{n} \sum_{i=1}^{n} f(X_i). \tag{2.37}$$

Even though this Markov chain Monte Carlo estimate is not constructed from independent samples, it can asymptotically satisfy analogs of the Law of Large Numbers (LLN) and Central Limit Theorem (CLT) that were used to justify ordinary Monte Carlo methods in Section 2.2.2. We sketch these important results here.

The MCMC analog of the LLN states that

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} f(X_i) = \int_{\mathcal{X}} f(x) \, \pi(x) \, dx \quad (a.s.)$$
(2.38)

for all functions f that are absolutely integrable with respect to π , i.e. all $f : \mathcal{X} \to \mathbb{R}$ that satisfy $\int_{\mathcal{X}} |f(x)| \pi(x) dx < \infty$. To quantify the asymptotic variance of MCMC estimates, the analog of the CLT must take into account both the Markov dependency structure among the samples used in the Monte Carlo estimate and also the initial state in which the chain was started. However, under mild conditions on both the Markov chain and the function f, the sample average for any initial distribution π_0 is asymptotically normal in distribution (with appropriate scaling):

$$\lim_{n \to \infty} \mathbb{P}\left(\frac{1}{\sqrt{n}} \sum_{n=1}^{n} (f(X_i) - \mu) < \alpha\right) = \mathbb{P}(Z < \alpha), \tag{2.39}$$

$$Z \sim \mathcal{N}\left(0, \sigma^2\right),$$
 (2.40)

$$\sigma^2 = \operatorname{Var}_{\pi}[f(X_0)] + 2\sum_{t=1}^{\infty} \operatorname{Cov}_{\pi}[f(X_0), f(X_t)]$$
(2.41)

where $\mu = \int_{\mathcal{X}} f(x) \pi(x) dx$ and where Var_{π} and Cov_{π} denote the variance and covariance operators with the chain (X_i) initialized in stationarity with $\pi_0 = \pi$. Thus standard error in the MCMC estimate also

scales asymptotically as $\frac{1}{\sqrt{n}}$, with a constant that depends on the autocovariance function of the stationary version of the chain. See Meyn and Tweedie [2009, chapter 17] and Robert and Casella [2004, section 6.7] for precise statements of both the LLN and CLT for Markov chain Monte Carlo estimates and for conditions on the Markov chain which guarantee that these theorems hold.

These results show that the asymptotic behavior of MCMC estimates of the form (2.37) is generally comparable to that of ordinary Monte Carlo estimates as discussed in Section 2.2.2. However, in the non-asymptotic regime MCMC estimates differ from ordinary Monte Carlo estimates in an important respect: there is a transient bias due to initializing the Markov chain out of stationarity. That is, the initial distribution π_0 from which the first iterate is sampled is generally not the chain's stationary distribution π , since if it were then ordinary Monte Carlo could be performed directly. While the marginal distribution of each Markov chain iterate converges to the stationary distribution, the effects of initialization on the initial iterates of the chain contribute an error term to Eq. (2.37) in the form of a transient bias.

This transient bias does not factor into the asymptotic behavior described by the MCMC analogs of the LLN and the CLT; asymptotically, it decreases at a rate of at least $\mathcal{O}(\frac{1}{n})$ and is hence dominated by the Monte Carlo standard error which decreases only at rate $\mathcal{O}(\frac{1}{\sqrt{n}})$. However, its effects can be significant in practice, especially in machine learning. Whenever a sampled chain seems "unmixed" because its iterates are too dependent on the initialization, errors in MCMC estimates are dominated by this transient bias.

The simulation in Figure 2.1 illustrates these error terms in MCMC estimates and how they can behave as more Markov chain samples are collected. The LLN and CLT for MCMC describe the regime on the far right of the plot: the total error can be driven arbitrarily small because the MCMC estimates are asymptotically unbiased, and the total error is asymptotically dominated by the Monte Carlo standard error. However, before reaching the asymptotic regime, the error is often dominated by the transient initialization bias. Several of the new methods we survey can be understood as attempts to alter the traditional MCMC tradeoffs,



Figure 2.1: A simulation illustrating error terms in MCMC estimator (2.37) as a function of the number of Markov chain iterations (log scale). The marginal distributions of the Markov chain iterates converge to the target distribution (top panel), while the errors in MCMC estimates due to transient bias and Monte Carlo standard error are eventually driven arbitrarily small at rates of $\mathcal{O}(\frac{1}{n})$ and $\mathcal{O}(\frac{1}{\sqrt{n}})$, respectively (bottom panel).

as we discuss further in Chapter 6.

Transient bias can be traded off against Monte Carlo standard error by choosing different subsets of Markov chain samples in the MCMC estimator. As an extreme choice, instead of using the MCMC estimator (2.37) with the full set of Markov chain samples $\{X_i\}_{i=1}^n$, transient bias can be minimized by forming estimates using only the last Markov chain sample:

$$\mathbb{E}[f(X)] \approx f(X_n). \tag{2.42}$$

However, this choice of MCMC estimator maximizes the Monte Carlo standard error, which asymptotically cannot be decreased below the posterior variance of the estimand. A practical choice is to form MCMC estimates using the last $\lceil n/2 \rceil$ Monte Carlo samples, resulting in an estimator

$$\mathbb{E}[f(X)] \approx \frac{1}{\lceil n/2 \rceil} \sum_{i=\lfloor n/2 \rfloor}^{n} f(X_i).$$
(2.43)

With this choice, once the marginal distribution of the Markov chain iterates approaches the stationary distribution the error due to transient bias is reduced at up to exponential rates. See Figure 2.2 for an illustration. With any choice of MCMC estimator, transient bias can be asymptotically decreased at least as fast as $\mathcal{O}(\frac{1}{n})$, and potentially much faster, while MCSE can decrease only as fast as $\mathcal{O}(\frac{1}{\sqrt{n}})$.

Using these ideas, MCMC algorithms provide a general means for estimating posterior expectations of interest: first construct an algorithm to simulate an ergodic Markov chain that admits the intended posterior density as its stationary distribution, and then simply run the simulation, collect samples, and form Monte Carlo estimates from the samples. The task then is to design an algorithm to simulate from such a Markov chain with the intended stationary distribution. In the following sections, we briefly review two canonical procedures for constructing such algorithms: Metropolis-Hastings and Gibbs sampling. For a thorough treatment, see Robert and Casella [2004] and Brooks et al. [2011, Chapter 1].



Figure 2.2: A simulation illustrating error terms in MCMC estimator (2.43) as a function of the number of Markov chain iterations (log scale). Because the first half of the Markov chain samples are not used in the estimate, the error due to transient bias is reduced much more quickly than in Figure 2.1 at the cost of shifting up the standard error curve.

Algorithm 1 Metropolis-Hastings for posterior sampling

Input: Initial state θ_0 , number of iterations T, joint density $p(\theta, \mathbf{x})$, proposal density $q(\theta' \mid \theta)$ **Output:** Samples $\theta_1, \ldots, \theta_T$ for t in 0, ..., T - 1 do $\theta' \sim q(\theta' \mid \theta_t)$ \triangleright Generate proposal $\alpha(\theta, \theta') \leftarrow \min\left(1, \frac{p(\theta', \mathbf{x})q(\theta_t \mid \theta')}{p(\theta_t, \mathbf{x})q(\theta' \mid \theta_t)}\right)$ ▷ Acceptance probability $u \sim \text{Unif}(0, 1)$ \triangleright Set stochastic threshold if $\alpha(\theta, \theta') > u$ then $\theta_{t+1} \leftarrow \theta'$ \triangleright Accept proposal else $\theta_{t+1} \leftarrow \theta_t$ \triangleright Reject proposal

2.2.5 Metropolis-Hastings (MH) sampling

In the context of Bayesian posterior inference, the Metropolis-Hastings (MH) algorithm simulates a reversible Markov chain over a state space Θ that admits the posterior density $p(\theta | x)$ as its stationary distribution. The algorithm depends on a user-specified proposal density, $q(\theta'|\theta)$, which can be evaluated numerically and sampled from efficiently, and also requires that the joint density $p(\theta, x)$ can be evaluated (up to proportionality). The MH algorithm then generates a sequence of states $\theta_1, \ldots, \theta_T \in \Theta$ according to Algorithm 1.

In each iteration, a proposal for the next state θ' is drawn from the proposal distribution, conditioned on the current state θ . The proposal is stochastically accepted with probability given by the *acceptance probability*,

$$\alpha(\theta, \theta') = \min\left(1, \frac{p(\theta', x)q(\theta \mid \theta')}{p(\theta, x)q(\theta' \mid \theta)}\right), \qquad (2.44)$$

via comparison to a random variate u drawn uniformly from the interval [0, 1]. If $u < \alpha(\theta, \theta')$, then the next state is set to the proposal, otherwise, the proposal is rejected and the next state is set to the current state. MH is a generalization of the *Metropolis algorithm* [Metropolis et al., 1953], which requires the proposal distribution to be symmetric, *i.e.*, $q(\theta' | \theta) = q(\theta | \theta')$, in which case the acceptance probability is

2.2. Markov Chain Monte Carlo inference

simply

$$\min\left(1, \frac{p(\theta', x)}{p(\theta, x)}\right). \tag{2.45}$$

Hastings [1970] later relaxed this by showing that the proposal distribution could be arbitrary.

One can show that the stationary distribution is indeed $p(\theta | x)$ by showing that the MH transition operator satisfies detailed balance (2.36). The MH transition operator density is a two-component mixture corresponding to the 'accept' event and the 'reject' event:

$$T(\theta \to \theta') = \alpha(\theta, \theta')q(\theta' \mid x) + (1 - \beta(\theta))\delta_{\theta}(\theta')$$
(2.46)

$$\beta(\theta) = \int_{\Theta} \alpha(\theta, \theta') q(\theta' \mid \theta) \ d\theta'.$$
(2.47)

To show detailed balance, it suffices to show the two balance conditions

$$\alpha(\theta, \theta')q(\theta' \mid \theta)p(\theta \mid x) = \alpha(\theta', \theta)q(\theta \mid \theta')p(\theta' \mid x)$$
(2.48)

$$(1 - \beta(\theta))\delta_{\theta}(\theta')p(\theta \mid x) = (1 - \beta(\theta'))\delta_{\theta'}(\theta)p(\theta' \mid x).$$
(2.49)

To show (2.48) we write

$$\begin{aligned} \alpha(\theta, \theta')q(\theta' \mid \theta)p(\theta \mid x) &= \min\left(1, \frac{p(\theta', x)q(\theta \mid \theta')}{p(\theta, x)q(\theta' \mid \theta)}\right)q(\theta' \mid \theta)p(\theta \mid x) \\ &= \min\left(q(\theta' \mid \theta)p(\theta \mid x), p(\theta', x)q(\theta \mid \theta')\frac{p(\theta \mid x)}{p(\theta, x)}\right) \\ &= \min\left(q(\theta' \mid \theta)p(\theta \mid x), p(\theta' \mid x)q(\theta \mid \theta')\right) \\ &= \min\left(q(\theta' \mid \theta)p(\theta, x)\frac{p(\theta' \mid x)}{p(\theta', x)}, p(\theta' \mid x)q(\theta \mid \theta')\right) \\ &= \min\left(1, \frac{p(\theta, x)q(\theta' \mid \theta)}{p(\theta', x)q(\theta \mid \theta')}\right)q(\theta \mid \theta')p(\theta' \mid x). \end{aligned}$$

$$(2.50)$$

To show (2.49), we need to verify that

$$(1 - \beta(\theta))p(\theta \,|\, x) = (1 - \beta(\theta'))p(\theta' \,|\, x), \qquad (2.51)$$

and we can use the same manipulation as in (2.50) under the integral

Algorithm 2 Gibbs sampling

Input: X Markov on graph G with nodes $\{1, 2, ..., N\}$, Markov blankets $MB_G(i)$ and subroutines to sample $X_i | X_{MB_G(i)}$ for each $i \in V$ **Output:** Samples $\{\hat{x}^{(t)}\}$

Initialize
$$x = (x_1, x_2, \dots, x_N)$$

for $t = 1, 2, \dots$ do
for $i = 1, 2, \dots, N$ do
 $x_i \leftarrow \text{ sample } X_i \mid X_{\text{MB}_G(i)} = x_{\text{MB}_G(i)}$
 $\hat{x}^{(t)} \leftarrow (x_1, x_2, \dots, x_N)$

sign:

$$(1 - \beta(\theta))p(\theta \mid x) = \left(1 - \int \alpha(\theta, \theta')q(\theta' \mid \theta) \ d\theta'\right)p(\theta \mid x)$$
$$= \left(1 - \int \alpha(\theta', \theta)q(\theta \mid \theta') \ d\theta\right)p(\theta' \mid x)$$
$$= (1 - \beta(\theta'))p(\theta' \mid x).$$
(2.52)

See Robert and Casella [2004, Section 7.3] for a more detailed treatment of the Metropolis-Hastings algorithm.

2.2.6 Gibbs sampling

Given a collection of n random variables $X = \{X_i : i \in [n]\}$, the Gibbs sampling algorithm iteratively samples each variable conditioned on the sampled values of the others. When the random variables are Markov on a graph G = (V, E), the conditioning can be reduced to each variable's respective Markov blanket, as in Algorithm 2. In the context of Bayesian inference, the posterior of interest may correspond to conditioning on some subset of the random variables, fixing them to observed values.

A variant of the *systematic scan* of Algorithm 2, in which nodes are traversed in a fixed order for each outer iteration, is the *random scan*, in which nodes are traversed according to a random permutation sampled for each outer iteration. An advantage of the random scan (and other variants) is that the chain becomes reversible and therefore

simpler to analyze [Robert and Casella, 2004, Section 10.1.2]. With the conditional independencies implied by a graph, some sampling steps may be performed in parallel.

The Gibbs sampling algorithm can be analyzed as a special case of the Metropolis-Hastings algorithm, where the proposal distribution is based on the conditional distributions and the acceptance probability is always one. If the Markov chain produced by a Gibbs sampling algorithm is ergodic, then the stationary distribution is the target distribution of X [Robert and Casella, 2004, Theorem 10.6]. The Markov chain for a Gibbs sampler can fail to be ergodic if, for example, the support of the target distribution is disconnected [Robert and Casella, 2004, Example 10.7]. A sufficient condition for Gibbs sampling to be ergodic is that all conditional densities exist and are positive everywhere [Robert and Casella, 2004, Theorem 10.8].

For a more detailed treatment of Gibbs sampling theory, see Robert and Casella [2004, Chapters 6 and 10].

2.3 Mean field variational inference

In mean field, and variational inference more generally, the task is to approximate an intractable distribution, such as a complex posterior, with a distribution from a tractable family so that the posterior can be efficiently interrogated for estimations of interest. In this section we define the mean field optimization problem and derive the standard coordinate optimization algorithm. We also give some basic results on the relationship between mean field and both graphical model and exponential family structure. For concreteness and simpler notation, we work mostly with undirected graphical models; the results extend immediately to directed models.

Mean field inference makes use of several densities and distributions, and so we use a subscript notation for expectations to clarify the measure used in the integration when it cannot easily be inferred from context. Given a function f and a random variable X with range \mathcal{X} and density p with respect to a base measure ν , we write the expectation

of f as

$$\mathbb{E}_{p(X)}\left[f(X)\right] = \int_{\mathcal{X}} f(x)p(x)\nu(dx).$$
(2.53)

Proposition 2.4 (Mean field variational inequality). For a probability density p with respect to a base measure ν of the form

$$p(x) = \frac{1}{Z}\bar{p}(x)$$
 with $Z \triangleq \int \bar{p}(x)\nu(dx),$ (2.54)

where \bar{p} is the unnormalized density, for all densities q with respect to ν we have

$$\log Z = \mathcal{L}[q] + \mathrm{KL}(q||p) \ge \mathcal{L}[q]$$
(2.55)

where

$$\mathcal{L}[q] \triangleq \mathbb{E}_{q(X)} \left[\log \frac{\bar{p}(X)}{q(X)} \right] = \mathbb{E}_{q(X)} \left[\log \bar{p}(X) \right] + \mathbb{H}[q]$$
(2.56)

$$\mathrm{KL}(q\|p) \triangleq \mathbb{E}_{q(X)}\left[\log\frac{q(X)}{p(X)}\right].$$
(2.57)

Here, $\mathbb{H}[q]$ is the differential entropy of q.

Proof. To show the equality, with $X \sim q$ we write

$$\mathcal{L}[q] + \mathrm{KL}(q||p) = \mathbb{E}_{q(X)} \left[\frac{\bar{p}(X)}{q(X)} \right] + \mathbb{E}_{q(X)} \left[\log \frac{q(X)}{p(X)} \right]$$
(2.58)

$$= \mathbb{E}_{q(X)} \left[\log \frac{\bar{p}(X)}{p(X)} \right]$$
(2.59)

$$= \log Z. \tag{2.60}$$

The inequality follows from the property $\text{KL}(q||p) \geq 0$, known as Gibbs's inequality, which follows from Jensen's inequality and the fact that the logarithm is concave:

$$-\operatorname{KL}(q\|p) = \mathbb{E}_{q(X)}\left[\log\frac{q(X)}{p(X)}\right] \le \log\int q(x)\frac{p(x)}{q(x)}\nu(dx) = 0 \quad (2.61)$$

with equality if and only if q = p (ν -a.e.).

We call the negative log of \bar{p} in (2.54) the *energy* and $\mathcal{L}[q]$ the variational lower bound, and say $\mathcal{L}[q]$ decomposes into the entropy minus the average energy as in (2.56). For two densities q and p with respect

2.3. Mean field variational inference

to the same base measure, KL(q||p) is the Kullback-Leibler divergence from q to p, used as a measure of dissimilarity between pairs of densities [Amari and Nagaoka, 2007].

The variational inequality given in Proposition 2.4 is useful in inference because if we wish to approximate an intractable p with a tractable q by minimizing KL(q||p), we can equivalently choose q to maximize $\mathcal{L}[q]$, which is possible to evaluate since it does not include the partition function Z.

In the context of Bayesian inference, p is usually an intractable posterior distribution of the form $p(\theta|x,\alpha)$, \bar{p} is the unnormalized joint distribution $\bar{p}(\theta) = p(\theta|\alpha)p(x|\theta)$, and Z is the marginal likelihood $p(x|\alpha) = \int p(x|\theta)p(\theta|\alpha)\nu(d\theta)$, which plays a central role in Bayesian model selection and the minimum description length (MDL) criterion [MacKay, 2002, Chapter 28] [Hastie et al., 2001, Chapter 7].

Given that graphical model structure can affect the complexity of probabilistic inference [Koller and Friedman, 2009] it is natural to consider families q that factor according to tractable graphs.

Definition 2.3 (Mean field variational inference). Let p be the density with respect to ν for a collection of random variables $X = (X_i : i \in V)$, and let

$$\mathcal{Q} \triangleq \{q: q(x) \propto \prod_{C \in \mathcal{C}} q_C(x_C)\}$$
(2.62)

be a family of densities with respect to ν that factorize according to a graph G = (V, E) with \mathcal{C} being the set of maximal cliques of G. Then the mean field optimization problem is

$$q^* = \operatorname*{arg\,max}_{q \in \mathcal{Q}} \mathcal{L}[q] \tag{2.63}$$

where $\mathcal{L}[q]$ is defined as in (2.56).

Note that this optimization problem is not in general convex¹ and so one can only expect to find a local optimum of the objective [Wainwright and Jordan, 2008]. However, when the model distribution is an exponential family the objective is concave in each q_C individually and hence an optimization procedure that updates each factor in turn, while

¹In the sense of maximizing a concave objective over a convex set.

holding the rest constant, will converge to a local optimum [Wainwright and Jordan, 2008] [Bishop, 2006, Section 10.1.1] [Murphy, 2012, Section 22.3]. We call such a coordinate ascent procedure on (2.63) a mean field algorithm.

For approximating families in a factored form, we can derive a generic update to be used in a mean field algorithm.

Proposition 2.5 (Mean field update). Given a mean field objective as in Definition 2.3, the optimal update to a factor q_A fixing the other factors defined by $q_A^* = \arg \max_{q_A} \mathcal{L}[q]$ is

$$q_A^*(x_A) \propto \exp\{\mathbb{E}[\log \bar{p}(x_A, X_{A^c})]\}$$
(2.64)

where the expectation is over $X_{A^c} \sim q_{A^c}$ with

$$q_{A^c}(x_{A^c}) \propto \prod_{C \in \mathcal{C} \setminus A} q_C(x_C) .$$
 (2.65)

Proof. Dropping terms constant with respect to q_A , we write

$$q_A^* = \underset{q_A}{\operatorname{arg\,min}} \operatorname{KL}(q \| p) \tag{2.66}$$

$$= \underset{q_A}{\operatorname{arg\,min}} \mathbb{E}_{q_A} \left[\log q_A(X_A) \right] + \mathbb{E}_{q_A} \left[\mathbb{E}_{q_{A^c}} \left[\log \bar{p}(X) \right] \right]$$
(2.67)

$$= \operatorname*{arg\,min}_{q_A} \operatorname{KL}(q_A \| \widetilde{p}_A) \tag{2.68}$$

where $\tilde{p}_A(x_A) \propto \exp\{\mathbb{E}_{q_{A^c}}[\log \bar{p}(x_A, X_{A^c})]\}$. Therefore, we achieve the unique (ν -a.e.) minimum by setting $q_A = \tilde{p}_A$.

Finally, we note the simple form of updates for exponential family conjugate pairs.

Proposition 2.6 (Mean field and conjugacy). If x_i appears in \bar{p} only in an exponential family conjugate pair (p_1, p_2) where

$$p_1(x_i|x_{\pi_G(i)}) \propto \exp\{\langle \eta(x_{\pi_G(i)}), t(x_i) \rangle\}$$
 (2.69)

$$p_2(x_{c_G(i)}|x_i) = \exp\{\langle t(x_i), (t(x_{c_G(i)}), 1) \rangle\}$$
(2.70)

then the optimal factor $q_i(x_i)$ is in the prior family with natural parameter

$$\widetilde{\eta} \triangleq \mathbb{E}_q[\eta(X_{\pi_G(i)})] + \mathbb{E}_q[(t(X_{c_G(i)}), 1)].$$
(2.71)

Proof. The result follows from substituting (2.69) and (2.70) into (2.64).

See Wainwright and Jordan [2008, Chapter 5] for a convex analysis perspective on mean field algorithms in graphical models composed of exponential families.

2.4 Stochastic gradient optimization

In this section we briefly review some basic ideas in stochastic gradient optimization. In particular, the basic algorithm we use in this paper is given in Algorithm 3 and sufficient conditions for its convergence to a local extreme point are given in Theorem 2.1.

Given a dataset $\bar{y} = \{\bar{y}^{(k)}\}_{k=1}^{K}$, where each $\bar{y}^{(k)}$ is a data *minibatch*, consider the optimization problem

$$\phi^* = \operatorname*{arg\,max}_{\phi} f(\phi) \tag{2.72}$$

where the objective function f decomposes according to

$$f(\phi) = \sum_{k=1}^{K} g(\phi, \bar{y}^{(k)}).$$
(2.73)

In the context of variational Bayesian inference, the objective f may be a variational lower bound on the model evidence and ϕ may be the parameters of the variational family. In MAP inference, f may be proportional to the posterior density and ϕ may be its parameters.

Using the decomposition of f, we can compute unbiased Monte Carlo estimates of its gradient. In particular, if the random index \hat{k} is sampled from $\{1, 2, \ldots, K\}$, denoting the probability of sampling index k as $p_k > 0$, we have

$$\nabla_{\phi} f(\phi) = \sum_{k=1}^{K} p_k \frac{1}{p_k} \nabla_{\phi} g(\phi, \bar{y}^{(k)}) = \mathbb{E}_{\hat{k}} \left[\frac{1}{p_{\hat{k}}} \nabla_{\phi} g(\phi, \bar{y}^{(\hat{k})}) \right].$$
(2.74)

Thus by considering a Monte Carlo approximation to the expectation over \hat{k} , we can generate stochastic approximate gradients of the objective f using only a single $\bar{y}^{(k)}$ at a time.

Algorithm 3 Stochastic gradient ascent

Input: $f : \mathbb{R}^n \to \mathbb{R}$ of the form (2.74), sequences $\rho^{(t)}$ and $G^{(t)}$ Initialize $\phi^{(0)} \in \mathbb{R}^n$ **for** $t = 0, 1, 2, \dots$ **do** $\hat{k}^{(t)} \leftarrow$ sample index k with probability p_k , for $k = 1, 2, \dots, K$ $\phi^{(t+1)} \leftarrow \phi^{(t)} + \rho^{(t)} \frac{1}{p_{\hat{k}}} G^{(t)} \nabla_{\phi} g(\phi^{(t)}, \bar{y}^{(\hat{k}^{(t)})})$

A stochastic gradient ascent algorithm uses these approximate gradients to perform updates and find a local optimum to the optimization problem. At each iteration, such an algorithm samples a data minibatch, computes a gradient with respect to that minibatch, and takes a step in that direction. In particular, for a sequence of stepsizes $\rho^{(t)}$ and a sequence of positive definite matrices $G^{(t)}$, a typical stochastic gradient ascent algorithm is given in Algorithm 3.

Stochastic gradient algorithms have very general convergence guarantees, requiring only weak conditions on the step size sequence and even the accuracy of the gradients themselves. We summarize a common set of sufficient conditions in Theorem 2.1. Proofs of this result, along with more general versions, can be found in Bertsekas and Tsitsiklis [1989] and Bottou [1998]. Note also that while the construction here has assumed that the stochasticity in the gradients arises only from randomly subsampling a finite sum, more general versions allow for other sources of stochasticity, typically requiring only bounded variance (and even allowing biased gradients) Bertsekas and Tsitsiklis [1989, Section 7.8].

Theorem 2.1. Given a function $f : \mathbb{R}^n \to \mathbb{R}$ of the form (2.73), if

- 1. there exists a constant C_0 such that $f(\phi) \leq C_0$ for all $\phi \in \mathbb{R}^n$,
- 2. there exists a constant C_1 such that

$$\|\nabla f(\phi) - \nabla f(\phi')\|_2 \le C_1 \|\phi - \phi'\|_2 \quad \forall \phi, \phi' \in \mathbb{R}^n,$$

3. there are positive constants C_2 and C_3 such that

$$\forall t \ C_2 I \prec G^{(t)} \prec C_3 I,$$

2.4. Stochastic gradient optimization

4. and the stepsize sequence $\rho^{(t)}$ satisfies

$$\sum_{t=0}^{\infty} \rho^{(t)} = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} (\rho^{(t)})^2 < \infty,$$

then Algorithm 3 converges to a local stationary point in the sense that

$$\liminf_{t \to \infty} \|\nabla f(\phi^{(t)})\| = 0 \tag{2.75}$$

with probability 1.

While stochastic optimization theory provides convergence guarantees, there is no general theory to analyze rates of convergence for nonconvex problems such as those that commonly arise in posterior inference. Indeed, the empirical rate of convergence often depends strongly on the variance of the stochastic gradient updates and on the choice of step size sequence. There are automatic methods to tune or adapt the sequence of stepsizes [Snoek et al., 2012, Ranganath et al., 2013], though we do not discuss them here. To make a single-pass algorithm, the minibatches can be sampled without replacement.

3

MCMC with data subsets

In MCMC sampling for Bayesian inference, the task is to simulate a Markov chain that admits as its stationary distribution the posterior distribution of interest. While there are many standard procedures for constructing and simulating from such Markov chains, when the dataset is large many of these algorithms' updates become computationally expensive. This growth in complexity naturally suggests the question of whether there are MCMC procedures that can generate approximate posterior samples without using the full dataset in each update. In this chapter, we focus on recent MCMC sampling schemes that scale Bayesian inference by operating on only subsets of data at a time.

3.1 Factoring the joint density

In most Bayesian inference problems, the fundamental object of interest is the posterior density, which for fixed data is proportional to the product of the prior and the likelihood:

$$\pi(\theta \,|\, \mathbf{x}) \propto \pi(\theta, \mathbf{x}) = \pi_0(\theta) \pi(\mathbf{x} \,|\, \theta). \tag{3.1}$$

In this survey we are often concerned with posteriors where the data $\mathbf{x} = \{x_n\}_{n=1}^N$ are conditionally independent given the model pa-
rameters θ , and hence the likelihood can be decomposed into a product of terms:

$$\pi(\theta \mid \mathbf{x}) \propto \pi_0(\theta) \pi(\mathbf{x} \mid \theta) = \pi_0(\theta) \prod_{n=1}^N \pi(x_n \mid \theta).$$
(3.2)

When N is large, this factorization can be exploited to construct MCMC algorithms in which the updates depend only on subsets of the data.

In particular, we can use subsets of data to form an unbiased Monte Carlo estimate of the log likelihood and consequently the log joint density. The log likelihood is a sum of terms:

$$\log \pi(\mathbf{x} \mid \theta) = \sum_{n=1}^{N} \log \pi(x_n \mid \theta), \qquad (3.3)$$

and we can approximate this sum using a random subset of m < N terms

$$\log \pi(\mathbf{x} \mid \theta) \approx \frac{N}{m} \sum_{n=1}^{m} \log \pi(x_n^* \mid \theta), \qquad (3.4)$$

where $\{x_n^*\}_{n=1}^m$ is a uniformly random subset of $\{x_n\}_{n=1}^N$. This approximation is an unbiased estimator and yields an unbiased estimate of the log joint density:

$$\log \pi(\theta)\pi(\mathbf{x} \mid \theta) \approx \log \pi_0(\theta) + \frac{N}{m} \sum_{n=1}^m \log \pi(x_n^* \mid \theta).$$
(3.5)

Several of the methods reviewed in this chapter exploit this estimator to perform MCMC updates.

3.2 Adaptive subsampling for Metropolis–Hastings

In traditional Metropolis–Hastings (MH), we evaluate the joint density to decide whether to accept or reject a proposal. As noted by Korattikara et al. [2014], because the value of the joint density depends on the full dataset, when N is large this is an unappealing amount of computation to reach a binary decision. In this section, we survey ideas for using approximate MH tests that depend on only a subset of the full dataset. The resulting approximate MCMC algorithms proceed in each iteration by reading only as much data as required to satisfy some estimated error tolerance.

While there are several variations, the common idea is to model the probability that the outcome of such an approximate MH test differs from the exact MH test. This probability model allows us to construct an approximate MCMC sampler, outlined in Section 3.2.2, where the user specifies some tolerance for the error in an MH test and the amount of data evaluated is controlled by an *adaptive stopping rule*. Different models for the MH test error lead to different stopping rules. Korattikara et al. [2014] use a normal model to construct a *t*statistic hypothesis test, which we describe in Section 3.2.3. Bardenet et al. [2014] instead use concentration inequalities, which we describe in Section 3.2.4. Given an error model and resulting stopping rule, both schemes rely on an MH test based on a Monte Carlo estimate of the log joint density, which we summarize in Section 3.2.1. Our notation in this section follows Bardenet et al. [2014].

Bardenet et al. [2014] observe that similar ideas have been developed both in the context of simulated annealing¹ by the operations research community [Bulgak and Sanders, 1988, Alkhamis et al., 1999, Wang and Zhang, 2006], and in the context of MCMC inference for factor graphs [Singh et al., 2012].

3.2.1 An approximate MH test based on a data subset

In the Metropolis–Hastings algorithm (§2.2.5), the proposal is stochastically accepted when

$$\frac{\pi(\theta' \mid \mathbf{x})q(\theta \mid \theta')}{\pi(\theta \mid \mathbf{x})q(\theta' \mid \theta)} > u, \tag{3.6}$$

where $u \sim \text{Unif}(0, 1)$. Rearranging and using log probabilities gives

$$\log\left[\frac{\pi(\mathbf{x}\mid\boldsymbol{\theta}')}{\pi(\mathbf{x}\mid\boldsymbol{\theta})}\right] > \log\left[u\frac{q(\boldsymbol{\theta}'\mid\boldsymbol{\theta})\pi_{0}(\boldsymbol{\theta})}{q(\boldsymbol{\theta}\mid\boldsymbol{\theta}')\pi_{0}(\boldsymbol{\theta}')}\right].$$
(3.7)

¹Simulated annealing is a stochastic optimization heuristic that is operationally similar to MH.

Scaling both sides by 1/N gives an equivalent threshold,

$$\Lambda(\theta, \theta') > \psi(u, \theta, \theta'), \tag{3.8}$$

where on the left, $\Lambda(\theta, \theta')$ is the average log likelihood ratio,

$$\Lambda(\theta, \theta') = \frac{1}{N} \sum_{n=1}^{N} \log\left[\frac{\pi(x_n \mid \theta')}{\pi(x_n \mid \theta)}\right] \equiv \frac{1}{N} \sum_{n=1}^{N} \ell_n, \qquad (3.9)$$

where

$$\ell_n = \log \pi(x_n \mid \theta') - \log \pi(x_n \mid \theta), \qquad (3.10)$$

and on the right,

$$\psi(u,\theta,\theta') = \frac{1}{N} \log \left[u \frac{q(\theta' \mid \theta) \pi_0(\theta)}{q(\theta \mid \theta') \pi_0(\theta')} \right].$$
(3.11)

We can form an approximate threshold by subsampling the ℓ_n . Let $\{\ell_n^*\}_{n=1}^m$ be a subsample of size m < N, without replacement, from $\{\ell_n\}_{n=1}^N$. This gives the following approximate test:

$$\hat{\Lambda}_m(\theta, \theta') > \psi(u, \theta, \theta'), \qquad (3.12)$$

where

$$\hat{\Lambda}_m(\theta, \theta') = \frac{1}{m} \sum_{n=1}^m \log\left[\frac{\pi(x_n^* \mid \theta')}{\pi(x_n^* \mid \theta)}\right] \equiv \frac{1}{m} \sum_{n=1}^m \ell_n^*.$$
(3.13)

This subsampled average log likelihood ratio $\hat{\Lambda}_m(\theta, \theta')$ is an unbiased estimate of the average log likelihood ratio $\Lambda(\theta, \theta')$. However, an error is made in the event that the approximate test (3.12) disagrees with the exact test (3.8), and the probability of such an error event depends on the distribution of $\hat{\Lambda}_m(\theta, \theta')$ and not just its mean.

Note that because the proposal θ' is usually a small perturbation of θ , we expect $\log \pi(x_n | \theta')$ to be similar to $\log \pi(x_n | \theta)$. In this case, we expect the log likelihood ratios ℓ_n have a smaller variance compared to the variance of $\log \pi(x_n | \theta)$ across data terms.

3.2.2 Approximate MH with an adaptive stopping rule

A nested sequence of data subsets, sampled without replacement, that converges to the complete dataset gives us a sequence of approximate

Algorithm 4 Approximate MH with an adaptive stopping rule

Input: Initial state θ_0 , number of iterations T, data $\mathbf{x} = \{x_n\}_{n=1}^N$, posterior $\pi(\theta \mid \mathbf{x})$, proposal $q(\theta' \mid \theta)$ **Output:** Samples $\theta_1, \ldots, \theta_T$ for t in $0, \ldots, T-1$ do $\theta' \sim q(\theta' \mid \theta_t)$ ▷ Generate proposal $u \sim \text{Unif}(0, 1)$ ▷ Draw random number $u \sim \text{Omi}(0, 1)$ $\psi(u, \theta, \theta') \leftarrow \frac{1}{N} \log \left[u \frac{q(\theta' \mid \theta) \pi_0(\theta)}{q(\theta \mid \theta') \pi_0(\theta')} \right]$ $\hat{\Lambda}(\theta, \theta') \leftarrow \text{AvgLogLikeRatioEstimate}(\theta, \theta', \psi(u, \theta, \theta'))$ if $\hat{\Lambda}(\theta, \theta') > \psi(u, \theta, \theta')$ then ▷ Approximate MH test $\theta_{t+1} \leftarrow \theta'$ \triangleright Accept proposal else $\theta_{t+1} \leftarrow \theta_t$ \triangleright Reject proposal

MH tests that converges to the exact MH test. Modeling the error of such an approximate MH test gives us a mechanism for designing an approximate MH algorithm in which, at each iteration, we incrementally read more data until an adaptive stopping rule informs us that our error is less than some user-specified tolerance. Algorithm 4 outlines this approach. The function AVGLOGLIKERATIOESTIMATE computes $\hat{\Lambda}(\theta, \theta')$ according to an adaptive stopping rule that depends on an error model, *i.e.*, a way to approximate or bound the probability that the approximate outcome disagrees with the full-data outcome:

$$\mathbb{P}\left|\left(\left(\hat{\Lambda}_m(\theta,\theta') > \psi(u,\theta,\theta')\right) \neq \left(\left(\Lambda(\theta,\theta') > \psi(u,\theta,\theta')\right)\right)\right|.$$
 (3.14)

We describe two possible error models in Sections 3.2.3 and 3.2.4.

A practical issue with adaptive subsampling is choosing the sizes of the data subsets. One approach, taken by Korattikara et al. [2014], is to use a fixed batch size b and read b more data points at a time. Bardenet et al. [2014] instead geometrically increase the total subsample size, and also discuss connections between adaptive stopping rules and related ideas such as bandit problems, racing algorithms and boosting.

3.2.3 Using a *t*-statistic hypothesis test

Korattikara et al. [2014] propose an approximate MH acceptance probability that uses a parametric test of significance as its error model. By assuming a normal model for the log likelihood estimate $\hat{\Lambda}(\theta, \theta')$, a *t*-statistic hypothesis test then provides an estimate of whether the approximate outcome agrees with the full-data outcome, *i.e.*, the expression in Equation (3.14). This leads to an adaptive framework as in Section 3.2.2 where, at each iteration, the data are processed incrementally until the *t*-test satisfies some user-specified tolerance ϵ .

Let us model the ℓ_n as i.i.d. from a normal distribution with bounded variance σ^2 :

$$\ell_n \sim \mathcal{N}(\mu, \sigma^2) \,. \tag{3.15}$$

The mean estimate $\hat{\mu}_m$ for μ based on the subset of size m is equal to $\hat{\Lambda}_m(\theta, \theta')$:

$$\hat{\mu}_m = \hat{\Lambda}_m(\theta, \theta') = \frac{1}{m} \sum_{n=1}^m \ell_n^*.$$
 (3.16)

The error estimate $\hat{\sigma}_m$ for σ may be derived from s_m/\sqrt{m} , where s_m is the empirical standard deviation of the *m* subsampled ℓ_n terms, *i.e.*,

$$s_m = \sqrt{\frac{m}{m-1} \left(\hat{\Lambda}_m^2(\theta, \theta') - \hat{\Lambda}_m(\theta, \theta')^2 \right)}, \qquad (3.17)$$

where

$$\hat{\Lambda}_m^2(\theta, \theta') = \frac{1}{m} \sum_{n=1}^m (\ell_n^*)^2.$$
(3.18)

To obtain a confidence interval, we multiply this estimate by the finite population correction, giving:

$$\hat{\sigma}_m = \frac{s_m}{\sqrt{m}} \sqrt{\frac{N-m}{N-1}} \,. \tag{3.19}$$

If m is large enough for the CLT to hold, the test statistic

$$t = \frac{\hat{\Lambda}_m(\theta, \theta') - \psi(u, \theta, \theta')}{\hat{\sigma}_m}$$
(3.20)

Algorithm 5 Estimate of the average log likelihood ratio. The adaptive stopping rule uses a *t*-statistic hypothesis test.

Parameters: batch size b, user-defined error tolerance ϵ function AVGLOGLIKERATIOESTIMATE $(\theta, \theta', \psi(u, \theta, \theta'))$ $m, \hat{\Lambda}(\theta, \theta'), \hat{\Lambda}^2(\theta, \theta') \leftarrow 0, 0, 0$ while True do $c \leftarrow \min(b, N - m)$ $\hat{\Lambda}(\theta, \theta') \leftarrow \frac{1}{m+c} \left(m\hat{\Lambda}(\theta, \theta') + \sum_{n=m+1}^{m+c} \log \frac{\pi(x_n \mid \theta')}{\pi(x_n \mid \theta)} \right)$ $\hat{\Lambda}^2(\theta, \theta') \leftarrow \frac{1}{m+c} \left(m\hat{\Lambda}^2(\theta, \theta') + \sum_{n=m+1}^{m+c} \left[\log \frac{\pi(x_n \mid \theta')}{\pi(x_n \mid \theta)} \right]^2 \right)$ $m \leftarrow m + c$ $s \leftarrow \sqrt{\frac{m}{m-1}} \left(\hat{\Lambda}^2(\theta, \theta') - \hat{\Lambda}(\theta, \theta')^2 \right)$ $\hat{\sigma} \leftarrow \frac{s}{\sqrt{m}} \sqrt{\frac{N-m}{N-1}}$ $\rho \leftarrow 1 - \phi_{m-1} \left(\left| \frac{\hat{\Lambda}(\theta, \theta') - \psi(u, \theta, \theta')}{\hat{\sigma}} \right| \right)$ if $\rho > \epsilon$ or m = N then return $\hat{\Lambda}(\theta, \theta')$

follows a Student's *t*-distribution with m-1 degrees of freedom when $\Lambda(\theta, \theta') = \psi(u, \theta, \theta')$. The tail probability for |t| then gives the probability that the approximate and actual outcomes agree, and thus

$$\rho = 1 - \phi_{m-1}(|t|) \tag{3.21}$$

is the probability that they disagree, where $\phi_{m-1}(\cdot)$ is the CDF of the Student's *t*-distribution with m-1 degrees of freedom. The *t*-test thus gives an adaptive stopping rule, *i.e.*, for any user-provided tolerance $\epsilon \geq 0$, we can incrementally increase m until $\rho \leq \epsilon$. We illustrate this approach in Algorithm 5.

3.2.4 Using concentration inequalities

Bardenet et al. [2014] propose an adaptive subsampling method that is mechanically similar to using a t-test but instead uses concentration inequalities. In addition to a bound on the error (of the approximate acceptance probability) that is local to each iteration, concentration bounds yield a bound on the total variation distance between the approximate and true stationary distributions.

As in Section 3.2.3, we evaluate an approximate MH threshold based on a data subset of size m, given in Equation (3.12). We bound the probability that the approximate binary outcome is incorrect via *concentration inequalities* that characterize the quality of $\hat{\Lambda}_m(\theta, \theta')$ as an estimate for $\Lambda(\theta, \theta')$. Such a concentration inequality is a probabilistic statement that, for $\delta_m \in (0, 1)$ and some constant c_m ,

$$\mathbb{P}\left(\left|\hat{\Lambda}_m(\theta, \theta') - \Lambda(\theta, \theta')\right| \le c_m\right) \ge 1 - \delta_m.$$
(3.22)

For example, in Hoeffding's inequality without replacement [Serfling, 1974]

$$c_m = C_{\theta,\theta'} \sqrt{\frac{2}{m} \left(1 - \frac{m-1}{N}\right) \log\left(\frac{2}{\delta_m}\right)}$$
(3.23)

where

$$C_{\theta,\theta'} = \max_{1 \le n \le N} \left| \log \pi(x_n \mid \theta') - \log \pi(x_n \mid \theta) \right| = \max_{1 \le n \le N} |\ell_n|, \qquad (3.24)$$

using ℓ_n as in Equation (3.10). Alternatively, if the empirical standard deviation s_m of the *m* subsampled ℓ_n^* terms is small, then the empirical Bernstein bound,

$$c_m = s_m \sqrt{\frac{2\log(3/\delta_m)}{m}} + \frac{6C_{\theta,\theta'}\log(3/\delta_m)}{m},$$
 (3.25)

is tighter [Audibert et al., 2009], where s_m is given in Equation (3.17). While $C_{\theta,\theta'}$ can be obtained via all the ℓ_n , this is precisely the computation we want to avoid. Therefore, the user must provide an estimate of $C_{\theta,\theta'}$.

Bardenet et al. [2014] use a concentration bound to construct an adaptive stopping rule based on a strategy called empirical Bernstein



Figure 3.1: Reproduction of Figure 2 from Bardenet et al. [2014]. If $|\hat{\Lambda}_m(\theta, \theta') - \psi(u, \theta, \theta')| > c_m$, then the adaptive stopping rule using a concentration bound is satisfied and we use the approximate MH test based on $\hat{\Lambda}_m(\theta, \theta')$.

stopping [Mnih et al., 2008]. Let c_m be a concentration bound as in Equation (3.23) or (3.25) and let δ_m be the associated error. This concentration bound states that $|\hat{\Lambda}_m(\theta, \theta') - \Lambda(\theta, \theta')| \leq c_m$ with probability $1 - \delta_m$. If $|\hat{\Lambda}_m(\theta, \theta') - \psi(u, \theta, \theta')| > c_m$, then the approximate MH test agrees with the exact MH test with probability $1 - \delta_m$. We reproduce a helpful illustration of this scenario from Bardenet et al. [2014] in Figure 3.1. If instead $|\hat{\Lambda}_m(\theta, \theta') - \psi(u, \theta, \theta')| \leq c_m$, then we want to increase *m* until this is no longer the case. Let *M* be the *stopping time*, *i.e.*, the number of data points evaluated using this criterion,

$$M = \min\left(N, \inf_{m \ge 1} \left| \hat{\Lambda}_m(\theta, \theta') - \psi(u, \theta, \theta') \right| > c_m \right).$$
(3.26)

We can set δ_m according to a user-defined parameter $\epsilon \in (0, 1)$ so that ϵ gives an upper bound on the error of the approximate acceptance probability. Let p > 1 and set

$$\delta_m = \frac{p-1}{pm^p}\epsilon, \quad \text{thus} \quad \sum_{m\geq 1}\delta_m \leq \epsilon.$$
 (3.27)

A union bound argument gives

$$\mathbb{P}\left(\bigcap_{m\geq 1}\left\{\left|\hat{\Lambda}_m(\theta,\theta') - \Lambda(\theta,\theta')\right| \le c_m\right\}\right) \ge 1 - \epsilon, \qquad (3.28)$$

under sampling without replacement. Hence, with probability $1 - \epsilon$, the approximate MH test based on $\hat{\Lambda}_M(\theta, \theta')$ agrees with the exact MH test. In other words, the stopping rule for computing $\hat{\Lambda}_m(\theta, \theta')$ in Algorithm 4 is satisfied once we observe $|\hat{\Lambda}_m(\theta, \theta') - \psi(u, \theta, \theta')| > c_m$.

40

We illustrate this approach in Algorithm 6, using Hoeffding's inequality without replacement.

In their actual implementation, Bardenet et al. [2014] modify δ_m to reflect the number of batches processed instead of the subsample size m. For example, suppose we use the concentration bound in Equation (3.23), *i.e.*, Hoeffding's inequality without replacement. Then after processing a subsample of size m in k batches, the adaptive stopping rule checks whether $|\hat{\Lambda}_m(\theta, \theta') - \psi(u, \theta, \theta')| > c_m$, where

$$c_m = C_{\theta,\theta'} \sqrt{\frac{2}{m} \left(1 - \frac{m-1}{N}\right) \log\left(\frac{2}{\delta_k}\right)}$$
(3.29)

and

$$\delta_k = \frac{p-1}{pk^p}\epsilon. \tag{3.30}$$

Also, as mentioned in Section 3.2.2, Bardenet et al. [2014] geometrically increase the subsample size by a factor γ . In their experiments, they use the empirical Bernstein-Serfling bound [Bardenet and Maillard, 2015]. For the hyperparameters, they set p = 2, $\gamma = 2$, and $\epsilon = 0.01$, and remark that they empirically found their algorithm to be robust to the choice of ϵ .

3.2.5 Error bounds on the stationary distribution

In this and the next subsection, we reproduce some theoretical results from Korattikara et al. [2014] and Bardenet et al. [2014]. After setting up some notation, we emphasize the most general aspects of these results, which apply to pairs of transition kernels whose differences are bounded, and thus are not specific to adaptive subsampling procedures. The central theorem is an upper bound on the difference between the stationary distributions of such pairs of kernels in the case of Metropolis–Hastings. Its proof depends on the ability to bound the difference in the acceptance probabilities, at each iteration, of the two MH transition kernels.

Preliminaries and notation. Let P and Q be probability measures (distributions) with Radon–Nikodym derivatives (densities) f_P and f_Q ,

Algorithm 6 Estimate of the average log likelihood ratio. The adaptive stopping rule uses Hoeffding's inequality without replacement.

Parameters: batch size *b*, user-defined error tolerance ϵ , estimate of $C_{\theta,\theta'} = \max_n |\ell_n|, p > 1$ **function** AVGLOGLIKERATIOESTIMATE $(\theta, \theta', \psi(u, \theta, \theta'))$ *m*, $\hat{\Lambda}(\theta, \theta') \leftarrow 0, 0$ **while** True **do** $c \leftarrow \min(b, N - M)$ $\hat{\Lambda}(\theta, \theta') \leftarrow \frac{1}{m+c} \left(m\hat{\Lambda}(\theta, \theta') + \sum_{n=m+1}^{m+c} \log \frac{\pi(x_n \mid \theta')}{\pi(x_n \mid \theta)} \right)$ $m \leftarrow m + c$ $\delta \leftarrow \frac{p-1}{pm^p} \epsilon$ $c \leftarrow C_{\theta,\theta'} \sqrt{\frac{2}{m} \left(1 - \frac{m-1}{N}\right) \log \left(\frac{2}{\delta}\right)}$ **if** $\left| \hat{\Lambda}(\theta, \theta') - \psi(u, \theta, \theta') \right| > c$ **or** m = N **then return** $\hat{\Lambda}(\theta, \theta')$

respectively, and absolutely continuous with respect to measure ν . The *total variation distance* between P and Q is

$$\|P - Q\|_{\mathrm{TV}} \equiv \frac{1}{2} \int_{\theta \in \Theta} d\nu(\theta) |f_P(\theta) - f_Q(\theta)|.$$
 (3.31)

For any transition kernel T, let T^k denote the kernel obtained via k iterations of T. Let T denote a transition kernel with stationary distribution $\pi(\theta | \mathbf{x})$. Let \tilde{T} denote an approximation to T, with stationary distribution $\tilde{\pi}$. When T is a MH transition kernel, let $q(\theta' | \theta)$ denote its proposal, and let $\alpha(\theta, \theta')$ denote its acceptance probability, given current and proposed states θ and θ' , respectively, *i.e.*,

$$\alpha(\theta, \theta') = \min\left(1, \frac{\pi(\theta')q(\theta \mid \theta')}{\pi(\theta)q(\theta' \mid \theta)}\right).$$
(3.32)

In this case, \tilde{T} is an approximate MH transition kernel with the same proposal $q(\theta' | \theta)$, and let $\tilde{\alpha}(\theta, \theta')$ denote its acceptance probability. Throughout this section, we specify when \tilde{T} is constructed from T

via an adaptive stopping rule; some of the results are more general. Let

$$\mathcal{E}(\theta, \theta') = \tilde{\alpha}(\theta, \theta') - \alpha(\theta, \theta') \tag{3.33}$$

be the acceptance probability error of the approximate MH test, with respect to the exact test. Finally, let

$$\mathcal{E}_{\max} = \sup_{\theta, \theta'} |\mathcal{E}(\theta, \theta')| \tag{3.34}$$

be the worst case absolute acceptance probability error.

Theoretical results. The theorem below provides an upper bound on the total variation distance between the stationary distributions of T and \tilde{T} ; the bound is linear in \mathcal{E}_{max} .

Theorem 3.1 (Total variation bound under uniform geometric ergodicity [Bardenet et al., 2014]). Let T be uniformly geometrically ergodic, *i.e.*, there exists an integer $h < \infty$, probability measure ν on $(\Theta, \mathcal{B}(\Theta))$, and constant $\lambda \in [0, 1)$ such that for all $\theta \in \Theta$ and $B \in \mathcal{B}(\Theta)$,

$$T^{h}(\theta, B) \ge (1 - \lambda)\nu(B), \qquad (3.35)$$

and thus there exists a constant $A < \infty$ such that for all $\theta \in \Theta$ and k > 0,

$$\|T^{k}(\theta, \cdot) - \pi\|_{\mathrm{TV}} \le A\lambda^{\lfloor k/h \rfloor}.$$
(3.36)

It follows that there exists a constant $C < \infty$ such that for all $\theta \in \Theta$ and k > 0,

$$\|\tilde{T}^{k}(\theta,\cdot) - \tilde{\pi}\|_{\mathrm{TV}} \le C \left(1 - (1-\epsilon)^{h}(1-\lambda)\right)^{\lfloor k/h \rfloor}.$$
 (3.37)

Moreover,

$$\|\pi - \tilde{\pi}\|_{\mathrm{TV}} \le \frac{Ah\mathcal{E}_{\mathrm{max}}}{1 - \lambda}.$$
(3.38)

The upper bound in Equation (3.38) depends on the worst case acceptance probability error. For adaptive subsampling schemes, this depends on the choice of adaptive procedure. We briefly outline a proof from Korattikara et al. [2014] of a similar theorem that exploits a stronger assumption on T. Specifically, assume T satisfies the contraction condition,

$$\|PT - \pi\|_{\rm TV} \le \eta \|P - \pi\|_{\rm TV},\tag{3.39}$$

for all probability distributions P and some constant $\eta \in [0, 1)$. We can combine the contraction condition with a bound on the *one-step error* between T and \tilde{T} , defined as $\|P\tilde{T} - PT\|_{\text{TV}}$, to bound $\|\pi - \tilde{\pi}\|_{\text{TV}}$. Note that this result does not require T to be a MH kernel.

For approximate MH with an adaptive stopping rule, \mathcal{E}_{max} , the maximum acceptance probability error, gives an upper bound on the one-step error. Korattikara et al. [2014] show how to calculate an upper bound on \mathcal{E}_{max} when using a *t*-test. Using concentration inequalities leads to a simpler bound: by construction, the user-defined error tolerance, ϵ , directly gives an upper bound on \mathcal{E}_{max} [Bardenet et al., 2014].

Finally, we note that an adaptive subsampling schemes using a concentration inequality enables an upper bound on the stopping time [Bardenet et al., 2014].

3.3 Sub-selecting data via a lower bound on the likelihood

Maclaurin and Adams [2014] introduce *Firefly Monte Carlo* (FlyMC), an auxiliary variable MCMC sampling procedure that operates on only subsets of data in each iteration. At each iteration, the algorithm dynamically selects what data to evaluate based on the random indicators included in the Markov chain state. In addition, it generates samples from the exact target posterior rather than an approximation. However, FlyMC requires a lower bound on the likelihood with a particular "collapsible" structure (essentially an exponential family lower bound) and is therefore not as generally applicable. The algorithm's performance depends on the tightness of the bound; it can achieve impressive gains in performance when model structure allows.

FlyMC samples from an augmented posterior that eliminates potentially many likelihood factors. Define

$$L_n(\theta) = p(x_n \mid \theta) \tag{3.40}$$

3.3. Sub-selecting data via a lower bound on the likelihood

and let $B_n(\theta)$ be a strictly positive lower bound on $L_n(\theta)$, *i.e.*, $0 < B_n(\theta) \leq L_n(\theta)$. For each datum, we introduce a binary auxiliary variable $z_n \in \{0, 1\}$ conditionally distributed according to a Bernoulli distribution,

$$p(z_n \mid x_n, \theta) = \left[\frac{L_n(\theta) - B_n(\theta)}{L_n(\theta)}\right]^{z_n} \left[\frac{B_n(\theta)}{L_n(\theta)}\right]^{1-z_n}, \quad (3.41)$$

where the z_n are independent for different n. When the bound is tight, *i.e.*, $B_n(\theta) = L_n(\theta)$, then $z_n = 0$ with probability 1. More generally, a tighter bound results in a higher probability that $z_n = 0$. Augmenting the density with $\mathbf{z} = \{z_n\}_{n=1}^N$ gives:

$$\tilde{\pi}(\theta, \mathbf{z} \mid \mathbf{x}) \propto \pi(\theta \mid \mathbf{x}) p(\mathbf{z} \mid \mathbf{x}, \theta)$$

$$= \pi_0(\theta) \prod_{n=1}^N \pi(x_n \mid \theta) p(z_n \mid x_n, \theta).$$
(3.42)

Using Equations (3.40) and (3.41), we can now write:

$$\tilde{\pi}(\theta, \mathbf{z} \mid \mathbf{x}) \propto \pi_0(\theta) \prod_{n=1}^N L_n(\theta) \left[\frac{L_n(\theta) - B_n(\theta)}{L_n(\theta)} \right]^{z_n} \left[\frac{B_n(\theta)}{L_n(\theta)} \right]^{1-z_n}$$
$$= \pi_0(\theta) \prod_{n=1}^N (L_n(\theta) - B_n(\theta))^{z_n} B_n(\theta)^{1-z_n}$$
$$= \pi_0(\theta) \prod_{n:z_n=1} (L_n(\theta) - B_n(\theta)) \prod_{n:z_n=0} B_n(\theta).$$
(3.43)

Thus for any fixed configuration of \mathbf{z} we can evaluate the joint density using only the likelihood terms $L_n(\theta)$ where $z_n = 1$ and the bound values $B_n(\theta)$ for each n = 1, 2, ..., N.

While Equation (3.43) still involves a product of N terms, if the product of the bound terms $\prod_{n:z_n=0} B_n(\theta)$ can be evaluated without reading each corresponding data point then the joint density can be evaluated reading only the data x_n for which $z_n = 1$. In particular, if the form of $B_n(\theta)$ is an exponential family density, then the product $\prod_{n:z_n=0} B_n(\theta)$ can be evaluated using only a finite-dimensional sufficient statistic for the data $\{x_n : z_n = 0\}$. Thus by exploiting lower bounds in the exponential family, FlyMC can reduce the amount of

data required at each iteration of the algorithm while maintaining the exact posterior as its stationary distribution. Maclaurin and Adams [2014] show an application of this methodology to Bayesian logistic regression.

FlyMC presents three main challenges. The first is constructing a collapsible lower bound, such as an exponential family, that is sufficiently tight. The second is designing an efficient implementation. Maclaurin and Adams [2014] discuss these issues and, in particular, design a cache-like data structure for managing the relationship between the N indicator values and the data. Finally, it is likely that the inclusion of these auxiliary variables slows the mixing of the Markov chain, but Maclaurin and Adams [2014] only provide empirical evidence that this effect is small relative to the computational savings from using data subsets.

3.4 Stochastic gradients of the log joint density

In this section, we review recent efforts to develop MCMC algorithms inspired by stochastic optimization techniques. This is motivated by the existence of, first, MCMC algorithms that can be thought of as the sampling analogues of optimization algorithms, and second, scalable stochastic versions of these optimization algorithms.

Traditional gradient ascent or descent performs optimization by iteratively computing and following a local gradient [Dennis and Schnabel, 1983]. In Bayesian MAP inference, the objective function is typically a log joint density and the update rule for gradient ascent is given by

$$\theta_{t+1} = \theta_t + \frac{\epsilon_t}{2} \left(\nabla \log \pi(\theta_t, \mathbf{x}) \right)$$
(3.44)

for $t = 1, ..., \infty$. As discussed in Section 2.4, *stochastic* gradient descent (SGD) is simple modification of gradient descent that exploits situations where the objective function decomposes into a sum of many terms. While the traditional gradient descent update depends on all the

data, i.e.,

$$\theta_{t+1} = \theta_t + \frac{\epsilon_t}{2} \left(\nabla \log \pi_0(\theta_t) + \sum_{n=1}^N \nabla \log \pi(x_n \,|\, \theta_t) \right), \qquad (3.45)$$

SGD forms an update based on only a data subset,

$$\theta_{t+1} = \theta_t + \frac{\epsilon_t}{2} \left(\nabla \log \pi_0(\theta_t) + \frac{N}{m} \sum_{n=1}^m \nabla \log \pi(x_n \,|\, \theta_t) \right). \tag{3.46}$$

The iterates converge to a local extreme point of the log joint density in the sense that $\lim_{t\to\infty} \nabla \log \pi(\theta_t | \mathbf{x}) = 0$ if the step size sequence $\{\epsilon_t\}_{t=1}^{\infty}$ satisfies

$$\sum_{t=1}^{\infty} \epsilon_t = \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \epsilon_t^2 < \infty.$$
(3.47)

A common choice of step size sequence is $\epsilon_t = \alpha(\beta + t)^{-\gamma}$ for some $\beta > 0$ and $\gamma \in (0.5, 1]$.

Welling and Teh [2011] propose stochastic gradient Langevin dynamics (SGLD), an approximate MCMC procedure that combines SGD with a simple kind of Langevin dynamics (Langevin Monte Carlo) [Neal, 1994]. They extend the Metropolis-adjusted Langevin algorithm (MALA) that uses noisy gradient steps to generate proposals for a Metropolis-Hastings chain [Roberts and Tweedie, 1996]. At iteration t, the MH proposal is

$$\theta' = \theta_t + \frac{\epsilon}{2} \left(\nabla \log \pi(\theta_t, \mathbf{x}) \right) + \eta_t, \qquad (3.48)$$

where the injected noise $\eta_t \sim \mathcal{N}(0, \epsilon)$ is Gaussian. Notice that the scale of the noise is $\sqrt{\epsilon}$, *i.e.*, is constant and set by the gradient step size parameter. The MALA proposal is thus a stochastic gradient step, constructed by adding noise to a step in the direction of the gradient.

SGLD modifies the Langevin dynamics in Equation (3.48) by using stochastic gradients based on data subsets, as in Equation (3.46), and requiring that the step size parameter satisfy Equation (3.47). Thus, at iteration t, the proposal is

$$\theta' = \theta_t + \frac{\epsilon_t}{2} \left(\nabla \log \pi_0(\theta_t) + \frac{N}{m} \sum_{n=1}^m \nabla \log \pi(x_n \mid \theta_t) \right) + \eta_t, \quad (3.49)$$

Algorithm 7 Stochastic gradient Langevin dynamics (SGLD).

Input: Initial state θ_0 , number of iterations T, data **x**, grad log prior $\nabla \log \pi_0(\theta)$, grad log likelihood $\nabla \log \pi(x \mid \theta)$, batch size m, step size tuning parameters $(e.q., \alpha, \beta, \gamma)$ **Output:** Samples $\theta_1, \ldots, \theta_T$ J = N/mfor τ in $0, \ldots, T/J - 1$ do $\mathbf{x} \leftarrow \text{Permute}(\mathbf{x})$ \triangleright For sampling without replacement for k in 0, ..., J - 1 do $t = \tau J + k$ $\epsilon_t \leftarrow \alpha(\beta + t)^{-\gamma}$ \triangleright Example step size $\eta_t \sim \mathcal{N}(0, \epsilon_t)$ ▷ Draw noise to inject $\theta' \leftarrow \theta_t + \frac{\epsilon_t}{2} \left(\nabla \log \pi_0(\theta_t) + \frac{N}{m} \sum_{n=km+1}^{km+m} \nabla \log \pi(x_n \,|\, \theta_t) \right) + \eta_t$ \triangleright Accept proposal with probability 1

where $\eta_t \sim \mathcal{N}(0, \epsilon_t)$. Notice that the injected noise decays with the gradient step size parameter, but at a slower rate. Specifically, if ϵ_t decays as $t^{-\gamma}$, then η_t decays as $t^{-\gamma/2}$. As in MALA, the SGLD proposal is a stochastic gradient step, where the noise comes from subsampling as well as the injected noise.

An actual Metropolis–Hastings algorithm would accept or reject the proposal in Equation (3.49) by evaluating the full (log) joint density at θ' and θ_t , but this is precisely the computation we wish to avoid. Welling and Teh [2011] observe that as $\epsilon_t \to 0$, $\theta' \to \theta_t$ in both Equations (3.48) and (3.49). In this limit, the probability of accepting the proposal converges to 1, but the chain stops completely. The authors suggest that ϵ_t can be decayed to a value that is large enough for efficient sampling, yet small enough for the acceptance probability to essentially be 1. These assumptions lead to a scheme where $\epsilon_t > \epsilon_{\infty} > 0$, for all t, and all proposals are accepted, therefore the acceptance probability is never evaluated. We show this scheme in Algorithm 7. Without the stochastic MH acceptance step, however, asymptotic samples are no longer guaranteed to represent the target distribution.

3.5. Summary

	Adaptive subsampling	FlyMC	SGLD
Approach	Approximate MH test	Auxiliary variables	Optimization plus noise
Requirements	Error model, $e.g.$, t -test	Likelihood lower bound	Gradients, <i>i.e.</i> , $\nabla \log \pi(\theta, \mathbf{x})$
Data access pattern	Mini-batches	Random	Mini-batches
Hyperparameters	Batch size, error tolerance per iteration	None	Batch size, error tolerance, annealing schedule
Asymptotic bias	Bounded TV	None	Bounded weak error

Table 3.1: Summary of recent MCMC methods for Bayesian inference that operate on data subsets. Error refers to the total variation distance between the stationary distribution of the Markov chain and the target posterior distribution.

In more recent work, Patterson and Teh [2013] apply SGLD to *Rie-mann manifold Langevin dynamics* [Girolami and Calderhead, 2011] and Chen et al. [2014] combine the idea of SGD with *Hamiltonian Monte Carlo* (HMC), an improved generalization of Langevin dynamics [Neal, 1994, 2010]. Finally, we note that all the methods in this section require gradient information that might not be readily computable.

3.5 Summary

In this chapter, we have surveyed three recent approaches to scaling MCMC that operate on subsets of data. Below and in Table 3.1, we summarize and compare adaptive subsampling approaches (§3.2), FlyMC (§3.3), and SGLD (§3.4) along several axes.

Approaches. Adaptive subsampling approaches replace the Metropolis–Hastings (MH) test, a function of all the data, with an approximate test that depends on only a subset. FlyMC is an auxiliary variable method that stochastically replaces likelihood computations with a collapsible lower bound. Stochastic gradient Langevin dynamics (SGLD) replaces gradients in a Metropolis-adjusted Langevin algorithm (MALA) with stochastic gradients based on data subsets and eliminates the Metropolis–Hastings test.

Generality, requirements, and assumptions. Each of the methods exploits assumptions or additional problem structure. Adaptive subsampling methods require an error model that accurately represents the probability that an approximate MH test will disagree with the exact MH test. A normal model [Korattikara et al., 2014] or concentration bounds [Bardenet et al., 2014] represent natural choices; under certain conditions, tighter concentration bounds may apply. FlyMC requires a strictly positive collapsible lower bound on the likelihood, essentially an exponential family lower bound, which may not in general be available. SGLD requires the log gradients of the prior and likelihood.

Data access patterns. While all the methods use subsets of data, their access patterns differ. Adaptive subsampling and SGLD require randomization to avoid issues of bias due to data order, but this randomization can be achieved by permuting the data before each pass and hence these algorithms allow data access that is mostly sequential. In contrast, FlyMC operates on random subsets of data determined by the Markov chain itself, leading to a random access pattern. However, subsets from one iteration to the next tend to be correlated, and motivate implementation details such as the proposed cache data structure.

Hyperparameters. FlyMC does not introduce additional hyperparameters that require tuning. Both adaptive subsampling methods and SGLD introduce hyperparameters that can significantly affect performance. Both are mini-batch methods, and thus have the batch size as a tuning parameter. In adaptive subsampling methods, the stopping criterion is evaluated potentially more than once before it is satisfied. This motivates schemes that geometrically increase the amount of data processed whenever the stopping criterion is not satisfied, which introduces additional hyperparameters. Adaptive subsampling methods additionally provide a single tuning parameter that allows the user to control the error at each iteration. Finally, since these adaptive methods define an approximate MH test, they implicitly also require that the user specify a proposal distribution. For SGLD, the user must specify an annealing schedule for the step size parameter; in particular, it should

converge to a small positive value so that the injected noise term dominates, while not being too large compared to the scale of the posterior distribution.

Error. FlyMC is exact in the sense that the target posterior distribution is a marginal of its augmented state space. The adaptive subsampling approaches and SGLD are approximate methods in that neither has a stationary distribution equal to the target posterior. The adaptive subsampling approaches bound the error of the MH test at each iteration, and for MH transition kernels with uniform ergodicity this one-step error bound leads to an upper bound on the total variation distance between the approximate stationary distribution and the target posterior distribution. The theoretical analysis of SGLD is less clear [Sato and Nakagawa, 2014].

3.6 Discussion

Data subsets. The methods surveyed in this chapter achieve computational gains by using data subsets in place of an entire dataset of interest. The adaptive subsampling algorithms (§3.2) are more successful when a small subsample leads to an accurate estimator for the exact MH test's accept/reject decision. Intuitively, such an estimator is easier to construct when the log posterior values at the proposed and current states are significantly different. This tends to be true far away from the mode(s) of the posterior, e.g., in the tails of a distribution that decay exponentially fast, compared to the area around a mode, which is locally more flat. Thus, these algorithms tend to evaluate more data when the chain is in the vicinity of a mode, and less data when the chain is far away (which tends to be the case for an arbitrary initial condition). SGLD (§3.4) exhibits somewhat related behavior. Recall that SGLD behaves more like SGD when the update rule is dominated by the gradient term, which tends to be true during the initial execution phase. Similar to SGD, the chain progresses toward a mode at a rate that depends on the accuracy of the stochastic gradients. For a log posterior target, stochastic gradients tend to be more accurate estimators of true gradients far away from the mode(s). In contrast, the MAP-tuned version of FlyMC (§3.3) requires the fewest data evaluations when the chain is close to the MAP, since by design, the lower likelihood bounds are tightest there. Meanwhile, the untuned version of FlyMC tends to exhibit the opposite behavior.

Adaptive proposal distributions. The Metropolis-Hastings algorithm requires the user to specify a proposal distribution. Fixing proposal distribution can be problematic, because the behavior of MH is sensitive to the proposal distribution and can furthermore change as the chain converges. A common solution, employed *e.g.*, by Bardenet et al. [2014], is to use an adaptive MH scheme [Haario et al., 2001, Andrieu and Moulines, 2006]. These algorithms tune the proposal distribution during execution, using information from the samples as they are generated, in a way that provably converges asymptotically. Often, it is desirable for the proposal distribution to be close to the target. This motivates adaptive schemes that fit a distribution to the observed samples and use this fitted model as the proposal distribution. For example, a simple online procedure can update the mean μ and covariance Σ of a multidimensional Gaussian model as follows:

$$\mu_{t+1} = \mu_t + \gamma_{t+1}(\theta_{t+1} - \mu_t) \quad t \ge 0 \Sigma_{t+1} = \Sigma_k + \gamma_{t+1}((\theta_{t+1} - \mu_t)(\theta_{t+1} - \mu_t)^\top - \Sigma_t),$$

where t indexes the MH iterations and γ_{t+1} controls the speed with which the adaptation vanishes. An appropriate choice is $\gamma_t = t^{-\alpha}$ for $\alpha \in [1/2, 1)$. The tutorial by Andrieu and Thoms [2008] provides a review of this and other, more sophisticated, adaptive MH algorithms.

Combining methods. The subsampling-based methods in this chapter are conceptually modular, and some may be combined. For example, it might be of interest to consider a 'tunable' version of FlyMC that achieves even greater computational efficiency at the cost of its original exactness. For example, we might use an adaptive subsampling scheme (§3.2) to evaluate only a subset of terms in Equation (3.43); this subset would need to represent terms corresponding to both possible values of z_n . As another example, Korattikara et al. [2014] suggest

3.6. Discussion

using adaptive subsampling as a way to 'fix up' SGLD. Recall that the original SGLD algorithm completely eliminates the MH test and blindly accepts all proposals, in order to avoid evaluating the full posterior. A reasonable compromise is to instead evaluate a fraction of the data within the adaptive subsampling framework, since this bounds the per-iteration error.

Unbiased likelihood estimators. The estimator in Equation (3.4) based on a data subset is an unbiased estimator for the log likelihood; to be explicit,

$$\exp\left\{\frac{N}{m}\sum_{n=1}^{m}\log\pi(x_n^*\,|\,\theta)\right\}$$
(3.50)

is not an unbiased estimate of the likelihood. While it is possible to transform Equation (3.50) into an unbiased likelihood estimate, *e.g.*, using a Poisson estimator [Wagner, 1987, Papaspiliopoulos, 2009, Fearnhead et al., 2010], it is not necessarily non-negative, which is a requirement to incorporate the estimator into a Metropolis-Hastings algorithm. In general, we cannot derive estimators that are both unbiased and nonnegative [Jacob and Thiery, 2015, Lyne et al., 2015]. *Pseudo-marginal MCMC* algorithms,² first introduced by Lin et al. [2000], rely on non-negative unbiased likelihood estimators to construct unbiased MCMC procedures [Andrieu and Roberts, 2009]. In this context, methods for constructing unbiased non-negative likelihood estimators include importance sampling [Beaumont, 2003] and particle filters [Andrieu et al., 2010, Doucet et al., 2015].

²Pseudo-marginal MCMC is also known as *exact-approximate sampling*.

4

Parallel and distributed MCMC

MCMC procedures that take advantage of parallel computing resources form another broad approach to scaling Bayesian inference. Because the computational requirements of inference often scale with the amount of data involved, and because large datasets may not even fit on a single machine, these approaches often focus on data parallelism. In this chapter we consider several approaches to scaling MCMC by exploiting parallel computation, either by adapting classical MCMC algorithms or by defining new simulation dynamics that are inherently parallel.

One way to use parallel computing resources is to run multiple sequential MCMC algorithms at once. However, running identical chains in parallel does not reduce the transient bias in MCMC estimates of posterior expectations, though it would reduce their variance. Instead of using parallel computation only to collect more MCMC samples and thus reduce only estimator variance without improving transient bias, it is often preferable to use computational resources to speed up the simulation of the chain itself. Section 4.1 surveys several methods that use parallel computation to speed up the execution of MCMC procedures, including both basic methods and more recent ideas.

Alternatively, instead of adapting serial MCMC procedures to ex-

ploit parallel resources, another approach is to design new approximate algorithms that are inherently parallel. Section 4.2 summarizes some recent ideas for simulations that can be executed in a data-parallel manner and have their results aggregated or corrected to represent posterior samples.

4.1 Parallelizing standard MCMC algorithms

An advantage to parallelizing standard MCMC algorithms is that they retain their theoretical guarantees and analyses. Indeed, a common goal is to produce identical samples under serial and parallel execution, so that parallel resources enable speedups without introducing new approximations. This section first summarizes some basic opportunities for parallelism in MCMC and then surveys the speculative execution framework for MH.

4.1.1 Conditional independence and graph structure

The MH algorithm has a straightforward opportunity for parallelism. In particular, if the target posterior can be written as

$$\pi(\theta \,|\, \mathbf{x}) \propto \pi_0(\theta) \pi(\mathbf{x} \,|\, \theta) = \pi_0(\theta) \prod_{n=1}^N \pi(x_n \,|\, \theta), \tag{4.1}$$

then when the number of likelihood terms N is large it may be beneficial to parallize the evaluation of the product of likelihoods. The communication between processors is limited to transmitting the value of the parameter and the scalar values of likelihood products. This basic parallelization, which naturally fits in a bulk synchronous parallel (BSP) computational model, exploits conditional independence in the probabilistic model, namely that the data are independent given the parameter.

Gibbs sampling algorithms can exploit more fine-grained conditional independence structure, and are thus a natural fit for graphical models which express such structure. Given a graphical model and a corresponding graph coloring with K colors that partitions the set of random variables into K groups, the random variables in each color



Figure 4.1: Graphical models and graph colorings can expose opportunities for parallelism in Gibbs samplers. (a) In this directed graphical model for a discrete mixture, each label (red) can be sampled in parallel conditioned on the parameters (blue) and data (gray) and similarly each parameter can be resampled in parallel conditioned on the labels. (b) This undirected grid has a classical "red-black" coloring, emphasizing that the variables corresponding to red nodes can be resampled in parallel given the values of black nodes and vice-versa.

group can be resampled in parallel while conditioning on the values in the other K-1 groups [Gonzalez et al., 2011]. Thus graphical models provide a natural perspective on opportunities for parallelism. See Figure 4.1 for some examples.

These opportunities for parallelism, while powerful in some cases, are limited by the fact that they require frequent global synchronization and communication. Indeed, at each iteration it is often the case that every element of the dataset is read by some processor and many processors must mutually communicate. The methods we survey in the remainder of this chapter aim to mitigate these limitations by adjusting the allocation of parallel resources or by reducing communication.

4.1.2 Speculative execution and prefetching

Another class of parallel MCMC algorithms uses speculative parallel execution to accelerate individual chains. This idea is called *prefetching* in some of the literature and appears to have received only limited



Figure 4.2: Metropolis-Hastings conceptualized as a binary tree. Nodes at depth d correspond to iteration t + d, where the root is at depth 0, and branching to the right/left indicates that the proposal is accepted/rejected. Each subscript is a sequence, of length d, of 0's and 1's, corresponding to the history of rejected and accepted proposals with respect to the root.

attention.

As shown in Algorithm 1, the body of a MH implementation is a loop containing a single conditional statement and two associated branches. We can thus view the possible execution paths as a binary tree, illustrated in Figure 4.2. The vanilla version of parallel prefetching speculatively evaluates all paths in this binary tree on parallel processors [Brockwell, 2006]. The sampled path will be exactly one of these, so with J processors this approach achieves a speedup of $\log_2 J$ with respect to single core execution, ignoring communication and bookkeeping overheads.

Naïve prefetching can be improved by observing that the two branches in Algorithm 1 are not taken with equal probability. For typical algorithm tunings, the reject branch tends to be more probable; a classic result for the optimal MH acceptance rate in the Gaussian case is 0.234 [Roberts et al., 1997], so prefetching scheduling policies can be built around the expectation of rejection. Angelino et al. [2014] provides a thorough review of these strategies.

Parallel predictive prefetching makes more efficient use of parallel

resources by dynamically predicting the outcome of each MH test [Angelino et al., 2014]. In the case of Bayesian inference, these predictions can be constructed in the same manner as the approximate MH algorithms based on subsets of data, as discussed in Section 3.2.2. Furthermore, these predictions can be made in the context of an error model, *e.g.*, with the concentration inequalities used by Bardenet et al. [2014]. This yields a straightforward and rational mechanism for allocating parallel cores to computations most likely to fall along the true execution path.

Algorithms 8 and 9 sketch pseudocode for an implementation of parallel predictive prefetching that follows a master-worker pattern. See Angelino [2014] for a formal description of the algorithm and implementation details.

4.2 Defining new data-parallel dynamics

In this section we survey two ideas for performing inference using new data-parallel dynamics. These algorithms define new dynamics in the sense that their iterates do not form ergodic Markov chains which admit the posterior distribution as an invariant distribution, and thus they do not qualify as classical MCMC schemes. Instead, while some of the updates in these algorithms resemble standard MCMC updates, the overall dynamics are designed to exploit parallel and distributed computation. A unifying theme of these new methods is to perform local computation on data while controlling the amount of global synchronization or communication.

One such family of ideas involves the definition of *subposteriors*, defined using only subsets of the full dataset. Inference in the subposteriors can be performed in parallel, and the results are then globally aggregated into an approximate representation of the full posterior. Because the synchronization and communication costs—as well as the approximation quality—are determined by the aggregation step, several such aggregation procedures have been proposed. In Section 4.2.1 we summarize some of these proposals.

Another class of data-parallel dynamics does not define indepen-

Algorithm 8 Parallel predictive prefetching master process		
repeat		
Receive message from worker j		
if worker j wants work then		
Find highest utility node ρ in tree with work left to do		
Send worker j the computational state of ρ		
else if message contains state θ_{ρ} at proposal node ρ then		
Record state θ_{ρ} at ρ		
else if message contains update at ρ then		
Update estimate of $\pi(\theta_{\rho} \mathbf{x})$ at ρ		
for node α in { ρ and its descendants} do		
Update utility of α		
if utility of α below threshold and worker k at α then		
Send worker k message to stop current computation		
if posterior computation at ρ and its parent complete then		
Know definitively whether to accept or reject θ_{ρ}		
Delete subtree corresponding to branch not taken		
if node ρ is the root's child then		
repeat		
Trim old root so that new root points to child		
Output state at root, the next state in the chain		
until posterior computation at root's child incomplete		
until master has output T Metropolis–Hastings chain states		
Terminate all worker processes		

dent subposteriors but instead, motivated by Gibbs sampling, focuses on simulating from local conditional distributions with out-of-date information. In standard Gibbs sampling, updates can be parallelized in models with conditional independence structure (Section 4.1), but without such structure the Gibbs updates may depend on the full dataset and all latent variables, and thus must be performed sequentially. These sequential updates can be especially expensive with large or distributed datasets. A natural approximation to consider is to run the same local Gibbs updates in parallel with out-of-date global infor-

Algorithm 9 Parallel predictive prefetching worker process
repeat
Send master request for work
Receive work assignment at node ρ from master
if the corresponding state θ_{ρ} has not yet been generated then
Generate proposal θ_{ρ}
repeat
Advance the computation of $\pi(\theta_{\rho} \mathbf{x})$
Send update at ρ to master
\mathbf{if} receive message to stop current computation \mathbf{then}
break
until computation of $\pi(\theta_{\rho} \mathbf{x})$ is complete
until terminated by master

mation and only infrequent communication. While such a procedure loses the theoretical guarantees provided by standard Gibbs sampling analysis, some empirical and theoretical results are promising. We refer to this broad class of methods as *Hogwild* Gibbs algorithms, and we survey some particular algorithms and analyses in Section 4.2.2.

4.2.1 Aggregating from subposteriors

Suppose we want to divide the evaluation of the posterior across J parallel cores. We can divide the data into J partition elements, $\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(J)}$, also called *shards*, and factor the posterior into J corresponding *subposteriors*, as

$$\pi(\theta \mid \mathbf{x}) = \prod_{j=1}^{J} \pi^{(j)}(\theta \mid \mathbf{x}^{(j)}), \qquad (4.2)$$

where

$$\pi^{(j)}(\theta \,|\, \mathbf{x}^{(j)}) = \pi_0(\theta)^{1/J} \prod_{x \in \mathbf{x}^{(j)}} \pi(x \,|\, \theta), \quad j = 1, \dots, J.$$
(4.3)

The contribution from the original prior is down-weighted so that the posterior is equal to the product of the J subposteriors,

4.2. Defining new data-parallel dynamics

i.e., $\pi(\theta | \mathbf{x}) = \prod_{j=1}^{J} \pi^{(j)}(\theta | \mathbf{x}^{(j)})$. Note that a subposterior is not the same as the posterior formed from the corresponding partition, *i.e.*,

$$\pi^{(j)}(\theta \,|\, \mathbf{x}^{(j)}) \neq \pi(\theta \,|\, \mathbf{x}^{(j)}) = \pi_0(\theta) \prod_{x \in \mathbf{x}^{(j)}} \pi(x \,|\, \theta). \tag{4.4}$$

Embarrassingly parallel consensus of subposteriors

Once a large dataset has been partitioned across multiple machines, a natural alternative is to try running MCMC inference on each partition element separately and in parallel. This yields samples from each subposterior in Equation 4.3, but there is no obvious choice for how to combine them in a coherent fashion to form approximate samples of the full posterior. In this section, we survey various proposals for forming such a *consensus* solution from the subposterior samples. Algorithm 10 outlines the structure of consensus strategies for embarrassingly parallel posterior sampling. This terminology, used by Huang and Gelman [2005] and Scott et al. [2013], invokes related notions of consensus, notably those that have existed for decades in the optimization literature on data-parallel algorithms in decentralized or distributed settings. We discuss this topic briefly in Section 4.2.1.

Below, we present two recent consensus strategies for combining subposterior samples, through weighted averaging and density estimation, respectively. The earlier report by Huang and Gelman [2005] proposes four consensus strategies, based either on normal approximations or importance resampling; the authors focus on Gibbs sampling for hierarchical models and do not evaluate any actual parallel implementations. Another consensus strategy is the recently proposed *variational consensus Monte Carlo* (VCMC) algorithm, which casts the consensus problem within a variational Bayes framework [Rabinovich et al., 2015].

Throughout this section, Gaussian densities provide a useful reference point and motivate some of the consensus strategies. Consider the jointly Gaussian model

$$\theta \sim \mathcal{N}(0, \Sigma_0) \tag{4.5}$$

$$\mathbf{x}^{(j)} \mid \boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma}_j). \tag{4.6}$$

Algorithm 10 Embarrassingly parallel consensus of subposteriors

Input: Initial state θ_0 , number of samples T, data partitions $\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(J)}$, subposteriors $\pi^{(1)}(\theta \mid \mathbf{x}^{(1)}), \ldots, \pi^{(J)}(\theta \mid \mathbf{x}^{(J)})$ **Output:** Approximate samples $\hat{\theta}_1, \ldots, \hat{\theta}_T$ **for** $j = 1, 2, \ldots, J$ in parallel **do** Initialize $\theta_{j,0}$ **for** $t = 1, 2, \ldots, T$ **do** Simulate MCMC sample $\theta_{j,t}$ from subposterior $\pi^{(j)}(\theta \mid \mathbf{x}^{(j)})$ Collect $\theta_{j,1}, \ldots, \theta_{j,T}$ $\hat{\theta}_1, \ldots, \hat{\theta}_T \leftarrow \text{CONSENSUSSAMPLES}(\{\theta_{j,1}, \ldots, \theta_{j,T}\}_{j=1}^J)$

The joint density is:

$$\begin{split} p(\theta, \mathbf{x}) &= p(\theta) \prod_{j=1}^{J} p(\mathbf{x}^{(j)} \mid \theta) \\ &\propto \exp\left\{-\frac{1}{2} \theta^{\top} \Sigma_{0}^{-1} \theta\right\} \prod_{j=1}^{J} \exp\left\{-\frac{1}{2} \left(\mathbf{x}^{(j)} - \theta\right)^{\top} \Sigma_{J}^{-1} \left(\mathbf{x}^{(j)} - \theta\right)\right\} \\ &\propto \exp\left\{-\frac{1}{2} \theta^{\top} \left(\Sigma_{0}^{-1} + \sum_{j=1}^{J} \Sigma_{j}^{-1}\right) \theta + \left(\sum_{j=1}^{J} \Sigma_{j}^{-1} \mathbf{x}^{(j)}\right)^{\top} \theta\right\}. \end{split}$$

Thus the posterior is Gaussian:

$$\theta \mid \mathbf{x} \sim \mathcal{N}(\mu, \Sigma),$$
 (4.7)

where

$$\Sigma = \left(\Sigma_0^{-1} + \sum_{j=1}^J \Sigma_j^{-1}\right)^{-1}$$
(4.8)

$$\mu = \Sigma \left(\sum_{j=1}^{J} \Sigma_j^{-1} \mathbf{x}^{(j)} \right).$$
(4.9)

To arrive at an expression for the subposteriors, we begin by factoring the joint distribution into an appropriate product:

$$p(\theta, \mathbf{x}) \propto \prod_{j=1}^{J} f_j(\theta),$$
 (4.10)

where

$$f_{j}(\theta) = p(\theta)^{1/J} p(\mathbf{x}^{(j)} | \theta)$$

= $\exp\left\{-\frac{1}{2}\theta^{\top}(\Sigma_{0}^{-1}/J)\theta\right\} \exp\left\{-\frac{1}{2}\left(\mathbf{x}^{(j)} - \theta\right)^{\top}\Sigma_{j}^{-1}\left(\mathbf{x}^{(j)} - \theta\right)\right\}$
 $\propto \exp\left\{-\frac{1}{2}\theta^{\top}\left(\Sigma_{0}^{-1}/J + \Sigma_{j}^{-1}\right)\theta + \left(\Sigma_{j}^{-1}\mathbf{x}^{(j)}\right)^{\top}\theta\right\}.$

Thus the subposteriors are also Gaussian:

$$\theta_j \sim \mathcal{N}\left(\tilde{\mu}_j, \tilde{\Sigma}_j\right) \propto f_j(\theta)$$
(4.11)

where

$$\tilde{\Sigma}_j = \left(\Sigma_0^{-1}/J + \Sigma_j^{-1}\right)^{-1} \tag{4.12}$$

$$\tilde{\mu}_j = \left(\Sigma_0^{-1}/J + \Sigma_j^{-1}\right)^{-1} \left(\Sigma_j^{-1} \mathbf{x}^{(j)}\right).$$
(4.13)

Weighted averaging of subposterior samples

One approach is to combine the subposterior samples via weighted averaging [Scott et al., 2013]. For simplicity, we assume that we obtain Tsamples in \mathbb{R}^d from each subposterior, and let $\{\theta_{j,t}\}_{t=1}^T$ denote the samples from the *j*th subposterior. The goal is to construct T consensus posterior samples $\{\hat{\theta}_t\}_{t=1}^T$, that (approximately) represent the full posterior, from the JT subposterior samples, where each $\hat{\theta}_t$ combines subposterior samples $\{\theta_{j,t}\}_{j=1}^J$. We associate with each subposterior *j* a matrix $W_j \in \mathbb{R}^{d \times d}$ and assume that each consensus posterior sample is a weighted¹ average:

$$\hat{\theta}_t = \sum_{j=1}^J W_j \theta_{j,t} \,. \tag{4.14}$$

The challenge now is to design an appropriate set of weights.

Following Scott et al. [2013], we consider the special case of Gaussian subposteriors, as in Equation 4.11. Our presentation is slightly different, as we also account for the effect of having a prior. We also drop the subscript t from our notation for simplicity. Let $\{\theta_j\}_{j=1}^J$ be

¹Our notation differs slightly from that of Scott et al. [2013] in that our weights W_j are normalized.

a set of draws from the J subposteriors. Each θ_j is an independent Gaussian and thus $\hat{\theta} = \sum_{j=1}^{J} W_j \theta_j$ is Gaussian. From Equation 4.13, its mean is

$$\mathbb{E}[\hat{\theta}] = \sum_{j=1}^{J} W_j \mathbb{E}[\theta_j] = \sum_{j=1}^{J} W_j \tilde{\mu}_j$$
$$= \sum_{j=1}^{J} W_j \left(\Sigma_0^{-1} / J + \Sigma_j^{-1} \right)^{-1} \left(\Sigma_j^{-1} \mathbf{x}^{(j)} \right). \quad (4.15)$$

Thus, if we choose

$$W_j = \Sigma \left(\Sigma_0^{-1} / J + \Sigma_j^{-1} \right) = \left(\Sigma_0^{-1} + \sum_{j=1}^J \Sigma_j^{-1} \right)^{-1} \left(\Sigma_0^{-1} / J + \Sigma_j^{-1} \right)$$
(4.16)

where Σ is the posterior covariance in Equation 4.8, then

$$\mathbb{E}[\hat{\theta}] = \Sigma\left(\sum_{j=1}^{J} \Sigma_j^{-1} \mathbf{x}^{(j)}\right) = \mu, \qquad (4.17)$$

where μ is the posterior mean in Equation 4.9. A similar calculation shows that $\text{Cov}(\hat{\theta}) = \Sigma$.

Thus for the Gaussian model, $\hat{\theta}$ is distributed according to the posterior distribution, indicating that Equation 4.16 gives the appropriate weights. Each weight matrix W_j is a function of Σ_0 , the prior covariance, and the subposterior covariances $\{\Sigma_j\}_{j=1}^J$. We can form a Monte Carlo estimate of each Σ_j using the empirical sample covariance $\bar{\Sigma}_j$. Algorithm 11 summarizes this consensus approach with weighted averaging. While this weighting is optimal in the Gaussian setting, Scott et al. [2013] shows it to be effective in some non-Gaussian models. Scott et al. [2013] also suggests weighting each dimension of a sample θ by the reciprocal of its marginal posterior variance, effectively restricting the weight matrices W_j to be diagonal.

Subposterior density estimation

Another consensus strategy relies on density estimation [Neiswanger et al., 2014]. First, use the subposterior samples to separately fit a

Algorithm 11 Consensus of subposteriors via weighted averaging.

Parameters: Prior covariance Σ_0 function CONSENSUSSAMPLES $(\{\theta_{j,1}, \ldots, \theta_{j,T}\}_{j=1}^J)$ for $j = 1, 2, \ldots, J$ do $\overline{\Sigma}_j \leftarrow$ Sample covariance of $\{\theta_{j,1}, \ldots, \theta_{j,T}\}$ $\Sigma \leftarrow \left(\Sigma_0^{-1} + \sum_{j=1}^J \overline{\Sigma}_j^{-1}\right)^{-1}$ for $j = 1, 2, \ldots, J$ do $W_j \leftarrow \Sigma \left(\Sigma_0^{-1}/J + \overline{\Sigma}_j^{-1}\right)$ \triangleright Compute weight matrices for $t = 1, 2, \ldots, T$ do $\hat{\theta}_t \leftarrow \sum_{j=1}^J W_j \theta_{j,t}$ \triangleright Compute weighted averages return $\hat{\theta}_1, \ldots, \hat{\theta}_T$

density estimator, $\tilde{\pi}^{(j)}(\theta \mid \mathbf{x}^{(j)})$, to each subposterior. The product of these density estimators then represents a density estimator for the full posterior target, *i.e.*,

$$\pi(\theta \,|\, \mathbf{x}) \approx \tilde{\pi}(\theta \,|\, \mathbf{x}) = \prod_{j=1}^{J} \tilde{\pi}^{(j)}(\theta \,|\, \mathbf{x}^{(j)}) \,. \tag{4.18}$$

Finally, one can sample from this posterior density estimator using MCMC; ideally, this density is straightforward to obtain and sample. In general, however, density estimation can yield complex models that are not amenable to efficient sampling.

Neiswanger et al. [2014] explore three density estimation approaches of various complexities. Their first approach assumes a parametric model and is therefore approximate. Specifically, they fit a Gaussian to each set of subposterior samples, yielding

$$\tilde{\pi}(\theta \,|\, \mathbf{x}) = \prod_{j=1}^{J} \mathcal{N}(\bar{\mu}_j, \bar{\Sigma}_j), \qquad (4.19)$$

where $\bar{\mu}_j$ and Σ_j are the empirical mean and covariance, respectively, of the samples from the *j*th subposterior. This product of Gaussians

simplifies to a single Gaussian $\mathcal{N}(\hat{\mu}_J, \hat{\Sigma}_J)$, where

$$\hat{\Sigma}_J = \left(\sum_{j=1}^J \bar{\Sigma}_j^{-1}\right)^{-1} \tag{4.20}$$

$$\hat{\mu}_J = \hat{\Sigma}_J \left(\sum_{j=1}^J \bar{\Sigma}_j^{-1} \bar{\mu}_j \right). \tag{4.21}$$

These parameters are straightforward to compute and the overall density estimate can be sampled with reasonable efficiency and even in parallel, if desired. Algorithm 12 summarizes this consensus strategy based on fits to Gaussians.

In the case when the model is jointly Gaussian, the parametric density estimator we form is $\mathcal{N}(\hat{\mu}_J, \hat{\Sigma}_J)$, with $\hat{\mu}_J$ and $\hat{\Sigma}_J$ given in Equations 4.21 and 4.20, respectively. In this special case, the estimator exactly represents the Gaussian posterior. However, recall that we could have instead written the exact posterior directly as $\mathcal{N}(\mu, \Sigma)$, where μ and Σ are in Equations 4.9 and 4.8, respectively. Thus, computing the exact posterior is more or less as expensive as computing the density estimator, *i.e.*, *J* local matrix inversions (or corresponding linear system solves).

4.2. Defining new data-parallel dynamics

The second approach proposed by Neiswanger et al. [2014] is to use a nonparametric kernel density estimate (KDE) for each subposterior. Suppose we obtain T samples $\{\theta_{j,t}\}_{t=1}^{T}$ from the *j*th subposterior, then its KDE with bandwidth parameter *h* has the following functional form:

$$\tilde{\pi}^{(j)}(\theta \,|\, \mathbf{x}^{(j)}) = \frac{1}{T} \sum_{t=1}^{T} \frac{1}{h^d} K\left(\frac{\|\theta - \theta_{j,t}\|}{h}\right), \tag{4.22}$$

i.e., the KDE is a mixture of T kernels, each centered at one of the samples. If we use T samples from each subposterior, then the density estimator for the full posterior is a complicated function with T^J terms, since it the a product of J such mixtures, and is therefore very challenging to sample from. Neiswanger et al. [2014] use a Gaussian KDE for each subposterior, and from this derive a density estimator for the full posterior that is a mixture of T^J Gaussians with unnormalized mixture weights. They also consider a third, semi-parametric approach to density estimation given by the product of a parametric (Gaussian) model and a nonparametric (Gaussian KDE) correction. As the number of samples $T \to \infty$, the nonparametric and semi-parametric density estimates exactly represent the subposterior densities and are therefore asymptotically exact. Unfortunately, their complex mixture representations grow exponentially in size, rendering them somewhat unwieldy in practice.

Weierstrass samplers

The consensus strategies surveyed so far are embarrassingly parallel. These methods obtain samples from each subposterior independently and in parallel, and from these attempt to construct samples that (approximately) represent the posterior post-hoc. The methods in this section proceed similarly, but introduce some amount of information sharing between the parallel samplers. This communication pattern is reminiscent of the alternating direction method of multipliers (ADMM) algorithm for data-parallel convex optimization; for a detailed treatment of ADMM, see the review by Boyd et al. [2011].

Weierstrass samplers [Wang and Dunson, 2013] are named for the

Weierstrass transform:

$$W_h f(\theta) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}h} \exp\left\{-\frac{(\theta-\xi)^2}{2h^2}\right\} f(\xi)d\xi, \qquad (4.23)$$

which was introduced by Weierstrass [1885]. The transformed function $W_h f(\theta)$ is the convolution of a one-dimensional function $f(\theta)$ with a Gaussian density of standard deviation h, and so converges pointwise to $f(\theta)$ as $h \to 0$,

$$\lim_{h \to 0} W_h f(\theta) = \int_{-\infty}^{\infty} \delta(\theta - \xi) f(\xi) d\xi = f(\theta),$$

where $\delta(\tau)$ is the Dirac delta function. For h > 0, $W_h f(\theta)$ can be thought of as a smoothed approximation to $f(\theta)$. Equivalently, if $f(\theta)$ is the density of a random variable θ , then $W_h f(\theta)$ is the density of a noisy measurement of θ , where the noise is an additive Gaussian with zero mean and standard deviation h.

Wang and Dunson [2013] analyzes a more general class of Weierstrass transforms by defining a multivariate version and also allowing non-Gaussian kernels:

$$W_h^{(K)}f(\theta_1,\ldots,\theta_d) = \int_{-\infty}^{\infty} f(\xi_1,\ldots,\xi_d) \prod_{i=1}^d h_i^{-1} K_i\left(\frac{\theta_i - \xi_i}{h_i}\right) d\xi_i.$$

For simplicity, we restrict our attention to the one-dimensional Weierstrass transform.

Weierstrass samplers use Weierstrass transforms on subposterior densities to define an augmented model. Let $f_j(\theta)$ denote the *j*-th subposterior,

$$f_{j}(\theta) = \pi^{(j)}(\theta \,|\, \mathbf{x}^{(j)}) = \pi_{0}(\theta)^{1/J} \prod_{x \in \mathbf{x}^{(j)}} \pi(x \,|\, \theta), \qquad (4.24)$$

68
so that the full posterior can be approximated as

$$\pi(\theta \mid \mathbf{x}) \propto \prod_{j=1}^{J} f_j(\theta) \approx \prod_{j=1}^{J} W_h f_j(\theta)$$
$$= \prod_{j=1}^{J} \int \frac{1}{\sqrt{2\pi}h} \exp\left\{-\frac{(\theta - \xi_j)^2}{2h^2}\right\} f_j(\xi_j) d\xi_j$$
$$\propto \int \prod_{j=1}^{J} \exp\left\{-\frac{(\theta - \xi_j)^2}{2h^2}\right\} f_j(\xi_j) d\xi_j \,. \tag{4.25}$$

The integrand of (4.25) defines the joint density of an augmented model that includes the $\xi = \{\xi_j\}_{j=1}^J$ as auxiliary variables:

$$\pi_h(\theta, \xi \,|\, \mathbf{x}) \propto \prod_{j=1}^J \exp\left\{-\frac{(\theta - \xi_j)^2}{2h^2}\right\} f_j(\xi_j).$$
(4.26)

The posterior of interest can then be approximated by the marginal distribution of θ in the augmented model,

$$\int \pi_h(\theta, \xi \,|\, \mathbf{x}) d\xi \approx \pi(\theta \,|\, \mathbf{x}) \tag{4.27}$$

with pointwise equality in the limit as $h \to 0$. Thus by running MCMC in the augmented model, producing Markov chain samples of both θ and ξ , we can generate approximate samples of the posterior. Furthermore, the augmented model is more amenable to parallelization due to its conditional independence structure: conditioned on θ , the subposterior parameters ξ are rendered independent.

The same augmented model construction can be motivated without explicit reference to the Weierstrass transform of densities. Consider the factor graph model of the posterior in Figure 4.3a, which represents the definition of the posterior in terms of subposterior factors,

$$\pi(\theta \,|\, \mathbf{x}) \propto \prod_{j=1}^{J} f_j(\theta). \tag{4.28}$$

This model can be equivalently expressed as a model where each subposterior depends on an exact local copy of θ . That is, writing ξ_j as



(a) Factor graph for $\pi(\theta | \mathbf{x})$ in terms of subposterior factors $f_j(\theta)$.



(b) Factor graph for the Weierstrass augmented model $\pi(\theta, \xi | \mathbf{x})$.

Figure 4.3: Factor graphs defining the augmented model of the Weierstrass sampler.

4.2. Defining new data-parallel dynamics

the local copy of θ for subposterior j, the posterior is the marginal of a new augmented model given by

$$\pi(\theta, \xi \,|\, \mathbf{x}) \propto \prod_{j=1}^{J} f_j(\xi_j) \delta(\xi_j - \theta) \,. \tag{4.29}$$

This new model can be represented by the factor graph in Figure 4.3b, with potentials $\psi(\xi_j, \theta) = \delta(\xi_j - \theta)$. Finally, rather than taking the ξ_j to be exact local copies of θ , we can instead relax them to be noisy Gaussian measurements of θ :

$$\pi_h(\theta, \xi \,|\, \mathbf{x}) \propto \prod_{j=1}^J f_j(\xi_j) \psi_h(\xi_j, \theta) \tag{4.30}$$

$$\psi_h(\xi_j, \theta) = \exp\left\{-\frac{(\theta - \xi_j)^2}{2h^2}\right\}.$$
(4.31)

Thus the potentials $\psi_h(\xi_j, \theta)$ enforce some consistency across the noisy local copies of the parameter but allow them to be decoupled, where the amount of decoupling depends on h. With smaller values of h the approximate model is more accurate, but the local copies are more coupled and hence sampling in the augmented model is less efficient.

We can construct a Gibbs sampler for the joint distribution $\pi(\theta, \xi | \mathbf{x})$ in Equation 4.25 by alternately sampling from $p(\theta | \xi)$ and $p(\xi_j | \theta, \mathbf{x}^{(j)})$, for j = 1, ..., J. It follows from Equation 4.25 that

$$p(\theta \mid \xi_1, \dots, \xi_j, \mathbf{x}) \propto \prod_{j=1}^J \exp\left\{-\frac{(\theta^2 - 2\theta\xi_j)}{2h^2}\right\}.$$
 (4.32)

Rearranging terms gives

$$p(\theta \mid \xi_1, \dots, \xi_j, \mathbf{x}) \propto \exp\left\{-\frac{(\theta - \bar{\xi})^2}{2h^2/J}\right\},$$
 (4.33)

where $\bar{\xi} = J^{-1} \sum_{j=1}^{J} \xi_j$. The remaining Gibbs updates follow from Equation 4.25, which directly yields

$$p(\xi_j \mid \theta, \mathbf{x}^{(j)}) \propto \exp\left\{-\frac{(\theta - \xi_j)^2}{2h^2}\right\} f_j(\xi_j), \quad j = 1, \dots, J.$$
(4.34)

Algorithm 13 Weierstrass Gibbs sampling. For simplicity, $\theta \in \mathbb{R}$.

Input: Initial state θ_0 , number of samples T, data partitions $\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(J)}$, subposteriors $f_1(\theta), \ldots, f_J(\theta)$, tuning parameter h **Output:** Samples $\theta_1, \ldots, \theta_T$ Initialize θ_0 **for** $t = 0, 1, \ldots, T - 1$ **do** Send θ_t to each processor **for** $j = 1, 2, \ldots, J$ in parallel **do** $\xi_{j,t+1} \sim p(\xi_{j,t+1} | \theta_t, \mathbf{x}^{(j)}) \propto \mathcal{N}(\xi_{j,t+1} | \theta_t, h^2) f_j(\xi_{j,t+1})$ Collect $\xi_{1,t+1}, \ldots, \xi_{J,t+1}$ $\bar{\xi}_{t+1} = \frac{1}{J} \sum_{j=1}^{J} \xi_{j,t+1}$ $\theta_{t+1} \sim \mathcal{N}(\theta_{t+1} | \bar{\xi}_{t+1}, h^2/J)$

This Gibbs sampler allows for parallelism but requires communication at every round. A straightforward parallel implementation, shown in Algorithm 13, generates the updates for ξ_1, \ldots, ξ_J in parallel, but the update for θ depends on the most recent values of all the ξ_j . Wang and Dunson [2013] describes an approximate variant of the full Gibbs procedure that avoids frequent communication by only occasionally updating θ . In other efforts to exploit parallelism while avoiding communication, the authors propose alternate Weierstrass samplers based on importance sampling and rejection sampling.

4.2.2 Hogwild Gibbs

Instead of designing new data-parallel algorithms from scratch, another approach is to take an existing MCMC algorithm and execute its updates in parallel at the expense of accuracy or theoretical guarantees. In particular, Hogwild Gibbs algorithms take a Gibbs sampling algorithm (§2.2.6) with interdependent sequential updates (*e.g.*, due to collapsed parameters or lack of graphical model structure) and simply run the updates in parallel anyway, using only occasional communication and out-of-date (stale) information from other processors. Because these strategies take existing algorithms and let the updates run 'hogwild' in the spirit of Hogwild! stochastic gradient descent in convex optimization [Recht et al., 2011], we refer to these methods as Hogwild Gibbs.

Similar approaches have a long history. Indeed, Gonzalez et al. [2011] attributes a version of this strategy, Synchronous Gibbs, to the original Gibbs sampling paper [Geman and Geman, 1984]. However, these strategies have seen renewed interest, particularly due to extensive empirical work on Approximate Distributed Latent Dirichlet Allocation (AD-LDA) [Newman et al., 2007, 2009, Asuncion et al., 2008, Liu et al., 2011, Ihler and Newman, 2012], which showed that running collapsed Gibbs sampling updates in parallel allowed for near-perfect parallelism without a loss in predictive likelihood performance. With the growing challenge of scaling MCMC both to not only big datasets but also big models, it is increasingly important to understand when and how these approaches may be useful.

In this section, we first define some variations of Hogwild Gibbs based on examples in the literature. Next, we survey the empirical results and summarize the current state of theoretical understanding.

Defining Hogwild Gibbs variants

Here we define some Hogwild Gibbs methods and related schemes, such as the stale synchronous parameter server. In particular, we consider bulk-synchronous parallel and asynchronous variations. We also fix some notation used for the remainder of the section.

For all of the Hogwild Gibbs algorithms, as with standard Gibbs sampling, we are given a collection of n random variables, $\{x_i : i \in [n]\}$ where $[n] \triangleq \{1, 2, \ldots, n\}$, and we assume that we can sample from the conditional distributions $x_i | x_{\neg i}$, where $x_{\neg i}$ denotes $\{x_j : j \neq i\}$. For the Hogwild Gibbs algorithms, we also assume we have K processors, each of which is assigned a set of variables on which to perform MCMC updates. We represent an assignment of variables to processors by fixing a partition $\{\mathcal{I}_1, \mathcal{I}_1, \ldots, \mathcal{I}_K\}$ of [n], so that the kth processor performs updates on the state values indexed by \mathcal{I}_k .

Algorithm 14 Bulk-synchronous parallel (BSP) Hogwild Gibbs

```
Input: Joint distribution over x = (x_1, ..., x_n), partition \{\mathcal{I}_1, ..., \mathcal{I}_K\}
of \{1, 2, ..., n\}, iteration schedule q(t, k)
Initialize \bar{x}^{(1)}
for t = 1, 2, ..., K in parallel do
\bar{x}_{\mathcal{I}_k}^{(t+1)} \leftarrow \text{LOCALGIBBS}(\bar{x}^{(t)}, \mathcal{I}_k, q(t, k))
Synchronize
function LOCALGIBBS(\bar{x}, \mathcal{I}, q)
for j = 1, 2, ..., q do
for i \in \mathcal{I} in order do
\bar{x}_i \leftarrow \text{sample } x_i | x_{\neg i} = \bar{x}_{\neg i}
return \bar{x}
```

Bulk-synchronous parallel Hogwild Gibbs

A bulk-synchronous parallel (BSP) Hogwild Gibbs algorithm assigns variables to processors and alternates between performing parallel processor-local updates and global synchronization steps. During epoch t, the kth processor performs q(t, k) MCMC updates, such as Gibbs updates, on the variables $\{x_i : i \in \mathcal{I}_k\}$ without communicating with the other processors; in particular, these updates are computed using out-of-date values for all $\{x_j : j \notin \mathcal{I}_k\}$. After all processors have completed their local updates, all processors communicate the updated state values in a global synchronization step and the system advances to the next epoch. We summarize this Hogwild Gibbs variant in Algorithm 14, in which the local MCMC updates are taken to be Gibbs updates.

Several special cases of the BSP Hogwild Gibbs scheme have been of interest. The Synchronous Gibbs scheme of Gonzalez et al. [2011] associates one variable with each processor, so that $|\mathcal{I}_k| = 1$ for each $k = 1, 2, \ldots, K$ (in which case we may take q = 1 since no local iterations are needed with a single variable). One may also consider the case where the partition is arbitrary and q is very large, in which case the local MCMC iterations may converge and exact block samples are

74

drawn on each processor using old statistics from other processors for each outer iteration. Finally, note that setting K = 1 and q(t, k) = 1reduces to standard Gibbs sampling on a single processor.

Asynchronous Hogwild Gibbs

Another Hogwild Gibbs pattern involves performing updates asynchronously. That is, processors might communicate only by sending messages to one another instead of by a global synchronization. Versions of this Hogwild Gibbs pattern has proven effective both for collapsed latent Dirichlet allocation topic model inference [Asuncion et al., 2008], and for Indian Buffet Process inference [Doshi-Velez et al., 2009]. A version was also explored in the Gibbs sampler of the Stale Synchronous Parameter (SSP) server of Ho et al. [2013], which placed an upper bound on the staleness of the entries of the state vector on each processor.

There are many possible communication strategies in the asynchronous setting, and so we follow a version of the random communication strategy employed by Asuncion et al. [2008]. In this approach, after performing some number of local updates, a processor sends its updated state information to a set of randomly-chosen processors and receives updates from other processors. The processor then updates its state representation and performs another round of local updates. A version of this asynchronous Hogwild Gibbs strategy is summarized in Algorithm 15.

Theoretical analysis

Despite its empirical successes, theoretical understanding of Hogwild Gibbs algorithms is limited. There are two settings in which some analysis has been offered: first, in a variant of AD-LDA, *i.e.*, Hogwild Gibbs applied to Latent Dirichlet Allocation models, and second in the jointly Gaussian case.

The work of Ihler and Newman [2012] provides some understanding of the effectiveness of a variant of AD-LDA by bounding in terms of run-time quantities the one-step error probability induced by pro-

Algorithm 15 Asynchronous Hogwild Gibbs

Initialize $\bar{x}^{(1)}$ for each processor k = 1, 2, ..., K in parallel do for t = 1, 2, ... do $\bar{x}_{\mathcal{I}_k}^{(t+1)} \leftarrow \text{LOCALGIBBS}(\bar{x}^{(t)}, \mathcal{I}_k, q(t, k))$ Send $\bar{x}_{\mathcal{I}_k}^{(t+1)}$ to K' randomly-chosen processors for each $k' \neq k$ do if update $\bar{x}_{\mathcal{I}_{k'}}$ received from processor k' then $\bar{x}_{\mathcal{I}_{k'}}^{(t+1)} \leftarrow \bar{x}_{\mathcal{I}_{k'}}$ else $\bar{x}_{\mathcal{I}_{k'}}^{(t+1)} \leftarrow \bar{x}_{\mathcal{I}_{k'}}^{(t)}$

ceeding with sampling steps in parallel, thereby allowing an AD-LDA user to inspect the computed error bound after inference [Ihler and Newman, 2012, Section 4.2]. In experiments, the authors empirically demonstrate very small upper bounds on these one-step error probabilities, e.g., a value of their parameter $\varepsilon = 10^{-4}$ meaning that at least 99.99% of samples are expected to be drawn just as if they were sampled sequentially. However, this per-sample error does not necessarily provide a direct understanding of the effectiveness of the overall algorithm because errors might accumulate over sampling steps; indeed, understanding this potential error accumulation is of critical importance in iterative systems. Furthermore, the bound is in terms of empirical runtime quantities, and thus it does not provide guidance on which other models the Hogwild strategy may be effective. Ihler and Newman [2012, Section 4.3] also provides approximate scaling analysis by estimating the order of the one-step bound in terms of a Gaussian approximation and some distributional assumptions.

The jointly Gaussian case is more tractable for analysis [Johnson et al., 2013, Johnson, 2014]. In particular, Johnson [2014, Theorem 7.6.6] shows that for the BSP Hogwild Gibbs process to be stable, *i.e.*, to form an ergodic Markov chain and have a well-defined stationary distribution, for any variable partition and any iteration schedule it suffices for the model's joint Gaussian precision matrix to satisfy a gen-

4.3. Summary

eralized diagonal dominance condition. Because the precision matrix contains the coefficients of the log potentials in a Gaussian graphical model, the diagonal dominance condition captures the intuition that Hogwild Gibbs should be stable when variables do not interact too strongly. Johnson [2014, Proposition 7.6.8] gives a more refined condition for the case where the number of processor-local Gibbs iterations is large.

When a bulk-synchronous parallel Gaussian Hogwild Gibbs process defines an ergodic Markov chain and has a stationary distribution, Johnson [2014, Chapter 7] also provides an understanding of how that stationary distribution relates to the model distribution. Because both the model distribution and the Hogwild Gibbs process stationary distribution are Gaussian, accuracy can be measured in terms of the mean vector and covariance matrix. Johnson [2014, Proposition 7.6.1] shows that the mean of a stable Gaussian Hogwild Gibbs process is always correct. Johnson [2014, Propositions 7.7.2 and 7.7.3] identify a tradeoff in the accuracy of the process covariance matrix as a function of the number of processor-local Gibbs iterations: at least when the processor interactions are sufficiently weak, more processor-local iterations between synchronization steps increase the accuracy of the covariances among variables within each processor but decrease the accuracy of the covariances between variables on different processors. Johnson [2014, Proposition 7.7.4] also gives a more refined error bound as well as an inexpensive way to correct covariance estimates for the case where the number of processor-local Gibbs iterations is large.

4.3 Summary

Many ideas for parallelizing MCMC have been proposed, exhibiting many tradeoffs. These ideas vary in generality, in faithfulness to the posterior, and in the parallel computation architectures for which they are best suited. Here we summarize the surveyed methods, emphasizing their relative strengths on these criteria. See Table 4.1 for an overview.

	Parallel density evaluation (§4.1.1)	Prefetching $($ §4.1.2 $)$	Consensus (§4.2.1)	Weierstrass ($§4.2.1$)	Hogwild Gibbs ($\$4.2.2$)	
Requirements	Conditional independence	None	Approximate factorization	Approximate factorization	Weak dependencies across processors	
Parallel model	BSP	Speculative execution	MapReduce	BSP	BSP and asynchronous mes- sage passing variants	
Communication	Each iteration	Master scheduling	Once	Tuneable	Tuneable	
Design choices	None	Scheduling policy	Data partition, consensus algorithm	Data partition, synchronization frequency, Weierstrass h	Data partition, communication frequency	
Computational overheads	None	Scheduling, bookkeeping	Consensus step	Auxiliary variable sampling	None	
Approximation error	None	None	Depends on number of proces- sors and consensus algorithm	Depends on number of processors and Weierstrass h	Depends on number of proces- sors and staleness	

 Table 4.1: Summary of recent approaches to parallel MCMC.

4.3. Summary

Simulating independent Markov chains Independent instances of serial MCMC algorithms can be run in an embarrassingly parallel manner, requiring only minimal communication between processors to ensure distinct initializations and to collect samples. This approach can reduce Monte Carlo variance by increasing the number of samples collected in any time budget, achieving an ideal parallel speedup, but does nothing to accelerate the warm-up period of the chains during which the transient bias is eliminated (see Section 2.2.4 and Chapter 6). That is, using parallel resources to run independent chains does nothing to improve mixing unless there is some mechanism for information sharing as in Nishihara et al. [2014]. In addition, running an independent MCMC chain on each processor requires each processor to access the full dataset, which may be problematic for especially large datasets. These considerations motivate both subposterior methods and Hogwild Gibbs.

Direct parallelization of standard updates Some MCMC algorithms applied to models with particular structure allow for straightforward parallel implementation. In particular, when the likelihood is factorized across data points, the computation of the Metropolis–Hastings acceptance probability can be parallelized. This strategy lends itself to a bulk-synchronous parallel (BSP) computational model. Parallelizing MH in this way yields exact MCMC updates and can be effective at reducing the mixing time required by serial MH, but it requires a simple likelihood function and its implementation requires frequent synchronization and communication, mitigating parallel speedups unless the likelihood function is very expensive.

Gibbs sampling also presents an opportunity for direct parallelization for particular graphical model structures. In particular, given a graph coloring of the graphical model, variables corresponding to nodes assigned a particular color are conditionally mutually independent and can be updated in parallel without communication. However, frequent synchronization and significant communication can be required to transmit sampled values to neighbors after each update. Relaxing both the strict conditional independence requirements and synchronization requirements motivates Hogwild Gibbs.

Prefetching and speculative execution The prefetching algorithms studied in Section 4.1.2 use speculative execution to transform traditional (serial) Metropolis–Hastings into a parallel algorithm without incurring approximate updates or requiring any model structure. The implementation naturally follows a master-worker pattern, where the master allocates (possibly speculative) computational work, such as proposal generation or (partial) density evaluation, to worker processors. Ignoring overheads, basic prefetching algorithms achieve at least logarithmic speedup in the number of processors available. More sophisticated scheduling by the master, such as predictive prefetching [Angelino et al., 2014], can increase speedup significantly. While this method is very general and yields the same iterates as serial MH, the speedup can be limited.

Subposterior consensus and Weierstrass samplers Subposterior methods, such as the consensus Monte Carlo algorithms and the Weierstrass samplers of Section 4.2.1, allow for data parallelism and minimal communication because each subposterior Markov chain can be allocated to a processor and simulation can proceed independently. Communication is required only for final sample aggregation in consensus Monte Carlo or the periodic resampling of the global parameter in the Weierstrass sampler. In consensus Monte Carlo, the quality of the approximate inference depends on both the effectiveness of the aggregation strategy and the extent to which dependencies in the posterior can be factorized into subposteriors. The Weierstrass samplers directly trade off approximation quality and the amount of decoupling between subposteriors.

The consensus Monte Carlo approach originated at Google [Scott et al., 2013] and naturally fits the MapReduce programming model, allowing it to be executed on large computational clusters. Recent work has extended consensus Monte Carlo and provides tools for designing simple consensus strategies [Rabinovich et al., 2015], but the generality and approximation quality of subposterior methods remain unclear.

4.4. Discussion

The Weierstrass sampler fits well into a BSP model.

Hogwild Gibbs Hogwild Gibbs of Section 4.2.2 also allows for data parallelism but avoids factorizing the posterior as in consensus Monte Carlo or instantiating coupled copies of a global parameter as in the Weierstrass sampler. Instead, processor-local sampling steps (such as local Gibbs updates) are performed with each processor treating other processors' states as fixed at stale values; processors can communicate updated states less frequently, either via synchronous or asynchronous communication. Hogwild Gibbs variants span a range of parallel computation paradigms from fully synchronous BSP to fully asynchronous message-passing. While Hogwild Gibbs has proven effective in practice for several models, its applicability and approximation tradeoffs remain unclear.

4.4 Discussion

The ideas surveyed in this chapter suggest several challenges and questions.

More parallelism means less accuracy The new data-parallel methods surveyed here, namely consensus Monte Carlo, the Weierstrass samplers, and Hogwild Gibbs, do not generate samples that are asymptotically distributed according to the target posterior. Instead, each generates samples that are asymptotically distributed according to a distribution that is meant to approximate the target posterior. While the nature of these approximations differ, each has a tradeoff between parallelism and accuracy: increasing parallelism by using more processors decreases the faithfulness of the asymptotic posterior approximation. This tradeoff may be inherent to most data-parallel MCMC schemes, though parallel predictive prefetching strategies do not suffer the same drawback.

How to split up data? The performance of data-parallel methods may be significantly affected by the data partitioning that assigns data subsets to processors. In the case of Hogwild Gibbs, it is probably best to choose a data partition that minimizes the strength of cross-processor dependence. Similarly, in the case of subposterior methods, some factorizations may be more effective than others. Since data paritioning is likely to have significant practical effects for all of these methods, it may be fruitful to develop and analyze general heuristics for assigning data to processors.

Analysis of approximation quality and tradeoffs The works surveyed in this chapter introduce several alternative approximations. While each is well motivated, it is unclear how to choose the most appropriate method for a given model, or how to think about and compare the various approximations and tradeoffs. A more unified perspective is necessary, through empirical comparison or through analyzing their application to simple models that are tractable for analysis.

82

5

Scaling variational mean field algorithms

Variational inference is a standard paradigm for posterior inference in Bayesian models. Because variational methods pose inference as an optimization problem, ideas in scalable optimization can in principle yield scalable posterior inference algorithms. In this chapter, we consider such scalable algorithms mainly in the context of mean field variational inference, which is often called variational Bayes.

These scalable variational inference algorithms can be compared to the algorithms of Chapters 3 and 4 in the same way that variational methods are usually compared to MCMC. That is, because inference is typically performed in a family of distributions that does not include the exact posterior, it can be said that variational methods do not fully instantiate the Bayesian computation that MCMC methods do (at least, when given unbounded computation time). Indeed, MAP inference, in which the posterior is represented only as a single atom, is an extreme case of variational inference. More generally, mean field variational families typically provide only unimodal approximations, and additionally cannot represent some posterior correlations near particular modes. As a result, MCMC methods can provide better performance even when the Markov chain only explores a single mode in a reasonable number of iterations.

Despite these potential shortcomings, variational inference is widely used in machine learning because the computational advantage over MCMC can be significant. This computational advantage is particularly salient in the context of scaling inference to large datasets. The big data context may also inform the relative cost of performing inference in a constrained variational family rather than attempting to represent the posterior exactly: when the posterior is concentrated, a variational approximation may suffice. While such questions may ultimately need to be explored empirically on a case-by-case basis, the scalable variational inference methods surveyed in this chapter provide the tools for such an exploration.

In this chapter we summarize two patterns of scalable variational inference. First, in Section 5.1, we discuss the application of stochastic gradient optimization methods to mean field variational inference problems. Second, in Section 5.2, we describe an alternative approach that instead leverages the idea of incremental posterior updating to develop an inference algorithm with minibatch-based updates.

5.1 Stochastic optimization and variational inference

Stochastic gradient optimization is a powerful tool for scaling optimization algorithms to large datasets, and it has been applied to mean field variational inference problems to great effect. While many traditional algorithms for optimizing mean field objective functions, including both gradient-based and coordinate optimization methods, require re-reading the entire dataset in each iteration, the stochastic gradient framework allows each update to be computed with respect to minibatches of the dataset while providing very general asymptotic convergence guarantees.

In this section we first summarize the stochastic variational inference (SVI) framework of Hoffman et al. [2013], which applies to models with complete-data conjugacy. Next, we discuss alternatives and extensions which can handle more general models at the cost of updates with greater variance and, hence, slower convergence.

84



Figure 5.1: Prototypical graphical model for stochastic variational inference (SVI). The global latent variables are represented by ϕ and the local latent variables by $z^{(k)}$.

5.1.1 SVI for complete-data conjugate models

This section follows the development in Hoffman et al. [2013]. It depends on results from stochastic gradient optimization theory; see Section 2.4 for a review. For notational simplicity we consider each minibatch to consist of only a single observation; the generalization to minibatches of arbitrary sizes is immediate.

Many common probabilistic models are hierarchical: they can be written in terms of *global* latent variables (or parameters), *local* latent variables, and observations. That is, many models can be written as

$$p(\phi, z, y) = p(\phi) \prod_{k=1}^{K} p(z^{(k)} | \phi) p(y^{(k)} | z^{(k)}, \phi)$$
(5.1)

where ϕ denotes global latent variables, $z = \{z^{(k)}\}_{k=1}^{K}$ denotes local latent variables, and $y = \{y^{(k)}\}_{k=1}^{K}$ denotes observations. See Figure 5.1 for a graphical model. Given such a class of models, the mean field variational inference problem is to approximate the posterior $p(\phi, z \mid \bar{y})$ for fixed data \bar{y} with a distribution of the form $q(\phi)q(z) = q(\phi)\prod_{k}q(z^{(k)})$ by finding a local minimum of the KL divergence from the approximating distribution to the posterior or, equivalently, finding a local maximum of the marginal likelihood lower bound

$$\mathcal{L}[q(\phi)q(z)] \triangleq \mathbb{E}_{q(\phi)q(z)} \left[\log \frac{p(\phi, z, \bar{y})}{q(\phi)q(z)} \right] \le \log p(\bar{y}).$$
(5.2)

Hoffman et al. [2013] develops a stochastic gradient ascent algorithm for such models that leverages complete-data conjugacy. Gradients of \mathcal{L} with respect to the parameters of $q(\phi)$ have a convenient form if we assume the prior $p(\phi)$ and each complete-data likelihood $p(z^{(k)}, y^{(k)} | \phi)$ are a conjugate pair of exponential family densities. That is, if we have

$$\log p(\phi) = \langle \eta_{\phi}, t_{\phi}(\phi) \rangle - \log Z_{\phi}(\eta_{\phi})$$
(5.3)

$$\log p(z^{(k)}, y^{(k)} | \phi) = \langle \eta_{zy}(\phi), \ t_{zy}(z^{(k)}, y^{(k)}) \rangle - \log Z_{zy}(\eta_{zy}(\phi))$$
(5.4)

then conjugacy identifies the statistic of the prior with the natural parameter and log partition function of the likelihood via

$$t_{\phi}(\phi) = (\eta_{zy}(\phi), -\log Z_{zy}(\eta_{zy}(\phi))),$$
 (5.5)

so that

$$p(\phi, z^{(k)}, \bar{y}^{(k)}) \propto \exp\left\{ \langle \eta_{\phi} + (t_{zy}(z^{(k)}, \bar{y}^{(k)}), 1), t_{\phi}(\phi) \rangle \right\}.$$
 (5.6)

Conjugacy implies the optimal variational factor $q(\phi)$ has the same form as the prior; that is, without loss of generality we can write $q(\phi)$ in the same form as (5.3),

$$q(\phi) = \exp\left\{\langle \widetilde{\eta}_{\phi}, t_{\phi}(\phi) \rangle - \log Z_{\phi}(\widetilde{\eta}_{\phi})\right\}, \qquad (5.7)$$

for some variational parameter $\tilde{\eta}_{\phi}$.

Given this conjugacy structure, we can find a simple expression for the gradient of \mathcal{L} with respect to the global variational parameter $\tilde{\eta}_{\phi}$, optimizing out the local variational factor q(z). That is, we write the variational objective over global parameters as

$$\mathcal{L}(\widetilde{\eta}_{\phi}) = \max_{q(z)} \mathcal{L}[q(\phi)q(z)].$$
(5.8)

Writing the optimal parameters of q(z) as $\tilde{\eta}_z^*$, note that when q(z) is partially optimized to a stationary point of \mathcal{L} , so that $\frac{\partial \mathcal{L}}{\partial \tilde{\eta}_z^*} = 0$ at $\tilde{\eta}_z^*$, the chain rule implies that the gradient with respect to the global variational parameters simplifies:

$$\frac{\partial \mathcal{L}}{\partial \tilde{\eta}_{\phi}}(\tilde{\eta}_{\phi}) = \frac{\partial \mathcal{L}}{\partial \tilde{\eta}_{\phi}}(\tilde{\eta}_{\phi}, \tilde{\eta}_{z}^{*}) + \frac{\partial \mathcal{L}}{\partial \tilde{\eta}_{z}^{*}} \frac{\partial \tilde{\eta}_{z}^{*}}{\partial \tilde{\eta}_{\phi}}(\tilde{\eta}_{\phi}, \tilde{\eta}_{z}^{*})$$
(5.9)

$$= \frac{\partial \mathcal{L}}{\partial \tilde{\eta}_{\phi}} (\tilde{\eta}_{\phi}, \tilde{\eta}_{z}^{*}) .$$
(5.10)

86

5.1. Stochastic optimization and variational inference

Because the optimal local factor q(z) can be computed with local mean field updates for a fixed value of the global variational parameter $\tilde{\eta}_{\phi}$, we need only find an expression for the gradient $\nabla_{\tilde{\eta}_{\phi}} \mathcal{L}(\tilde{\eta}_{\phi})$ in terms of the optimized local factors.

To find an expression for the gradient $\nabla_{\tilde{\eta}_{\phi}} \mathcal{L}(\tilde{\phi}_{\phi})$ that exploits conjugacy structure, using (5.6) we can substitute

$$p(\phi, z, \bar{y}) \propto \exp\left\{ \langle \eta_{\phi} + \sum_{k} (t_{zy}(z^{(k)}, \bar{y}^{(k)}), 1), t_{\phi}(\phi) \rangle \right\},$$
 (5.11)

into the definition of \mathcal{L} in (5.2). Using the optimal form of $q(\phi)$, we have

$$\mathcal{L}(\tilde{\eta}_{\phi}) = \mathbb{E}_{q(\phi)q(z)} \left[\langle \eta_{\phi} + \sum_{k} t_{zy}(z^{(k)}, \bar{y}^{(k)}) - \tilde{\eta}_{\phi}, t_{\phi}(\phi) \rangle \right] + \log Z_{\phi}(\tilde{\eta}_{\phi}) + \text{const.}$$
(5.12)
$$= \langle \eta_{\phi} + \sum_{k} \mathbb{E}_{q(z^{(k)})}[t_{zy}(z^{(k)}, \bar{y}^{(k)})] - \tilde{\eta}_{\phi}, \mathbb{E}_{q(\phi)}[t_{\phi}(\phi)] \rangle + \log Z_{\phi}(\tilde{\eta}_{\phi}) + \text{const},$$
(5.13)

where the constant does not depend on $\tilde{\eta}_{\phi}$. Using the identity for natural exponential families that

$$\nabla \log Z_{\phi}(\tilde{\eta}_{\phi}) = \mathbb{E}_{q(\phi)}[t_{\phi}(\phi)], \qquad (5.14)$$

we can write the same expression as

$$\mathcal{L}(\tilde{\eta}_{\phi}) = \langle \eta_{\phi} + \sum_{k} \mathbb{E}_{q(z^{(k)})}[t_{zy}(z^{(k)}, \bar{y}^{(k)})] - \tilde{\eta}_{\phi}, \ \nabla \log Z_{\phi}(\tilde{\eta}_{\phi}) \rangle + \log Z_{\phi}(\tilde{\eta}_{\phi}) + \text{const.}$$
(5.15)

Thus we can compute the gradient of $\mathcal{L}(\tilde{\eta}_{\phi})$ with respect to the global variational parameters $\tilde{\eta}_{\phi}$ as

$$\nabla_{\tilde{\eta}_{\phi}} \mathcal{L}(\tilde{\eta}_{\phi}) = \langle \nabla^2 \log Z_{\phi}(\tilde{\eta}_{\phi}), \eta_{\phi} + \sum_k \mathbb{E}_{q(z^{(k)})}[t_{zy}(z^{(k)}, \bar{y}^{(k)})] - \tilde{\eta}_{\phi} \rangle$$
$$- \nabla \log Z_{\phi}(\tilde{\eta}_{\phi}) + \nabla \log Z_{\phi}(\tilde{\eta}_{\phi}) \qquad (5.16)$$
$$= \langle \nabla^2 \log Z_{\phi}(\tilde{\eta}_{\phi}), \eta_{\phi} + \sum_k \mathbb{E}_{q(z^{(k)})}[t_{zy}(z^{(k)}, \bar{y}^{(k)})] - \tilde{\eta}_{\phi} \rangle$$

where the first two terms come from applying the product rule.

The matrix $\nabla^2 \log Z_{\phi}(\tilde{\eta}_{\phi})$ is the Fisher information of the variational family, since

$$-\mathbb{E}_{q(\phi)}\left[\nabla_{\widetilde{\eta}_{\phi}}^{2}\log q(\phi)\right] = \nabla^{2}\log Z_{\phi}(\widetilde{\eta}_{\phi}).$$
(5.17)

In the context of stochastic gradient ascent, we can cancel the multiplication by the matrix $\nabla^2 \log Z_{\phi}(\tilde{\eta}_{\phi})$ simply by choosing the sequence of positive definite matrices in Algorithm 3 to be $G^{(t)} \triangleq \nabla^2 \log Z_{\phi}(\tilde{\eta}_{\phi}^{(t)})^{-1}$. This choice yields a stochastic *natural* gradient ascent algorithm [Amari and Nagaoka, 2007], where the updates are stochastic approximations to the natural gradient

$$\widetilde{\nabla}_{\widetilde{\eta}_{\phi}} \mathcal{L} = \eta_{\phi} + \sum_{k} \mathbb{E}_{q(z^{(k)})}[t_{zy}(z^{(k)}, \bar{y}^{(k)})] - \widetilde{\eta}_{\phi}.$$
(5.18)

Natural gradients effectively include a second-order quasi-Newton correction for local curvature in the variational family, making the updates invariant to reparameterization of the variational family and thus often improving performance of the algorithm. More importantly, at least for the case of complete-data conjugate families considered here, natural gradient steps are in fact easier to compute than 'flat' gradient steps in either the natural parameterization or moment parameterization of the variational family $q(\phi)$.

Therefore a stochastic natural gradient ascent algorithm on the global variational parameter $\tilde{\eta}_{\phi}$ proceeds at iteration t by sampling a minibatch $\bar{y}^{(k)}$ and taking a step of some size $\rho^{(t)}$ in an approximate natural gradient direction via

$$\widetilde{\eta}_{\phi} \leftarrow (1 - \rho^{(t)})\widetilde{\eta}_{\phi} + \rho^{(t)} \left(\eta_{\phi} + K \mathbb{E}_{q(z^{(k)})}[t(z^{(k)}, \overline{y}^{(k)})] \right)$$
(5.19)

where we have assumed the minibatches are of equal size to simplify notation. The local variational factor $q(z^{(k)})$ is computed using a local mean field update on the data minibatch and the global variational factor. That is, if $q(z^{(k)})$ is not further factorized in the mean field approximation, it is computed according to

$$q(z^{(k)}) \propto \exp\left\{\mathbb{E}_{q(\phi)}[\log p(z^{(k)} | \phi)p(\bar{y}^{(k)} | z^{(k)}, \phi)]\right\}.$$
(5.20)

We summarize the general SVI algorithm in Algorithm 16.

88

Al	gorithm	16	Stochastic	Variational	Inference ((SVI)	
----	---------	----	------------	-------------	-------------	-------	--

Initialize global variational parameter $\tilde{\eta}_{\phi}^{(0)}$ **for** $t = 0, 1, 2, \dots$ **do** $\hat{k} \leftarrow \text{sample index } k \text{ with probability } p_k > 0, \text{ for } k = 1, 2, \dots, K$ $q(z^{(\hat{k})}) \leftarrow \text{LOCALMEANFIELD}(\tilde{\eta}^{(t)}, \bar{y}^{(\hat{k})})$ $\tilde{\eta}_{\phi}^{(t+1)} \leftarrow (1 - \rho^{(t)}) \tilde{\eta}_{\phi}^{(t)} + \rho^{(t)} \left(\eta_{\phi} + \frac{1}{p_{\hat{k}}} \mathbb{E}_{q(z^{(\hat{k})})}\left[t(z^{(\hat{k})}, \bar{y}^{(\hat{k})})\right]\right)$

5.1.2 Stochastic gradients with general nonconjugate models

The development of SVI in the preceding section assumes that $p(\phi)$ and $p(z, y | \phi)$ are a conjugate pair of exponential families. This assumption led to a particularly convenient form for the natural gradient of the mean field variational objective and hence an efficient stochastic gradient ascent algorithm. However, when models do not have this conjugacy structure, more general algorithms are required.

In this section we review Black Box Variational Inference (BBVI), which is a stochastic gradient algorithm for variational inference that can be applied at scale [Ranganath et al., 2014]. The "black box" name suggests its generality: while the stochastic variational inference of Section 5.1.1 requires particular model structure, BBVI only requires that the model's log joint distribution can be evaluated. It also makes few demands of the variational family, since it only requires that the family can be sampled and that the gradient of its log joint with respect to the variational parameters can be computed efficiently. With these minimal requirements, BBVI is not only useful in the big-data setting but also a tool for handling nonconjugate variational inference more generally. Because BBVI uses Monte Carlo approximation to compute stochastic gradient updates, it fits naturally into a stochastic gradient optimization framework, and hence it has the additional benefit of yielding a scalable algorithm simply by adding minibatch sampling to its updates at the cost of increasing their variance. In this subsection we review the general BBVI algorithm and then compare it to the SVI algorithm of Section 5.1.1. For a review of Monte Carlo estimation, see Section 2.2.2.

Scaling variational mean field algorithms

We consider a general model $p(\theta, y) = p(\theta) \prod_{k=1}^{K} p(y^{(k)} | \theta)$ including parameters θ and observations $y = \{y^{(k)}\}_{k=1}^{K}$ divided into K minibatches. The distribution of interest is the posterior $p(\theta | y)$ and we write the variational family as $q(\theta) = q(\theta | \tilde{\eta}_{\theta})$, where we suppress the particular mean field factorization structure of $q(\theta)$ from the notation. The mean field variational lower bound is then

$$\mathcal{L} = \mathbb{E}_{q(\theta)} \left[\log \frac{p(\theta, y)}{q(\theta)} \right].$$
 (5.21)

Taking the gradient with respect to the variational parameter $\tilde{\eta}_{\theta}$ and expanding the expectation into an integral, we have

$$\nabla_{\widetilde{\eta}_{\theta}} \mathcal{L} = \nabla_{\widetilde{\eta}_{\theta}} \int q(\theta) \log \frac{p(\theta, y)}{q(\theta)} d\theta$$
(5.22)

$$= \int \nabla_{\widetilde{\eta}_{\theta}} \left[\log \frac{p(\theta, y)}{q(\theta)} \right] q(\theta) d\theta + \int \log \frac{p(\theta, y)}{q(\theta)} \nabla_{\widetilde{\eta}_{\theta}} q(\theta) d\theta \quad (5.23)$$

where we have moved the gradient into the integrand and applied the product rule to yield two terms. The first term is identically zero:

$$\int \nabla_{\widetilde{\eta}_{\theta}} \left[\log \frac{p(\theta, y)}{q(\theta)} \right] q(\theta) d\theta = -\int \frac{1}{q(\theta)} \nabla_{\widetilde{\eta}_{\theta}} \left[q(\theta) \right] q(\theta) d\theta \qquad (5.24)$$

$$= -\int \nabla_{\widetilde{\eta}_{\theta}} q(\theta) d\theta \tag{5.25}$$

$$= -\nabla_{\widetilde{\eta}_{\theta}} \int q(\theta) d\theta = 0 \tag{5.26}$$

where we have used $\nabla_{\tilde{\eta}_{\theta}} \log p(\theta, y) = 0$. To write the second term of (5.23) in a form that allows convenient Monte Carlo approximation, we first note the identity

$$\nabla_{\widetilde{\eta}_{\theta}} \log q(\theta) = \frac{\nabla_{\widetilde{\eta}_{\theta}} q(\theta)}{q(\theta)} \implies \nabla_{\widetilde{\eta}_{\theta}} q(\theta) = q(\theta) \nabla_{\widetilde{\eta}_{\theta}} \log q(\theta) \qquad (5.27)$$

and hence we can write the second term of (5.23) as

$$\int \log \frac{p(\theta, y)}{q(\theta)} \nabla_{\widetilde{\eta}_{\theta}} q(\theta) d\theta = \int \log \frac{p(\theta, y)}{q(\theta)} \nabla_{\widetilde{\eta}_{\theta}} \left[\log q(\theta) \right] q(\theta) d\theta \quad (5.28)$$

$$= \mathbb{E}_{q(\theta)} \left[\log \frac{p(\theta, y)}{q(\theta)} \nabla_{\widetilde{\eta}_{\theta}} \log q(\theta) \right]$$
(5.29)

$$\approx \frac{1}{|\mathcal{S}|} \sum_{\hat{\theta} \in S} \log \frac{p(\hat{\theta}, y)}{q(\hat{\theta})} \nabla_{\widetilde{\eta}_{\theta}} \log q(\hat{\theta})$$
(5.30)

5.1. Stochastic optimization and variational inference

where in the final line we have written the expectation as a Monte Carlo estimate using a set of samples S, where $\hat{\theta} \stackrel{\text{iid}}{\sim} q(\theta)$ for $\hat{\theta} \in S$. Notice that the gradient is written as a weighted sum of gradients of the variational log density with respect to the variational parameters, where the weights depend on the model log joint density.

The BBVI algorithm uses the Monte Carlo estimate (5.30) to compute stochastic gradient updates. This gradient estimator is also known as the score function estimator [Kleijnen and Rubinstein, 1996, Gelman and Meng, 1998]. The variance of these updates, and hence the convergence of the overall stochastic gradient algorithm, depends both on the sizes of the gradients of the variational log density and on the variance of $q(\theta)$. Large variance in the gradient estimates can lead to very slow optimization, and so Ranganath et al. [2014] proposes and evaluates two variance reduction schemes, including a control variate method as well as a Rao-Blackwellization method that can exploit factorization structure in the variational family.

To provide a scalable version of BBVI, gradients can be further approximated by subsampling minibatches of data. That is, using $\log p(\theta, y) = \log p(\theta) + \sum_{k=1}^{K} \log p(y^{(k)} | \theta)$ we write (5.29) and (5.30) as

$$\nabla_{\widetilde{\eta}_{\theta}} \mathcal{L} = \mathbb{E}_{\hat{k}} \left[\mathbb{E}_{q(\theta)} \left[\left(\log \frac{p(\theta)}{q(\theta)} + K \log p(y^{(\hat{k})} \mid \hat{\theta}) \right) \nabla_{\widetilde{\eta}_{\theta}} \log q(\theta) \right] \right]$$
(5.31)

$$\approx \frac{1}{|\mathcal{S}|} \sum_{\hat{\theta} \in \mathcal{S}} \left(\log \frac{p(\theta)}{q(\hat{\theta})} + K \mathbb{E}_{\hat{k}} \left[\log p(y^{(\hat{k})} \mid \hat{\theta}) \right] \right) \nabla_{\tilde{\eta}_{\theta}} \log q(\hat{\theta}) \quad (5.32)$$

with the minibatch index \hat{k} distributed uniformly over $\{1, 2, \ldots, K\}$ and the minibatches are assumed to be the same size for simpler notation. This subsampling over minibatches further increases the variance of the updates and thus may further limit the rate of convergence of the algorithm. We summarize this version of the BBVI algorithm in Algorithm 17.

It is instructive to compare the fully general BBVI algorithm applied to hierarchical models to the SVI algorithm of Section 5.1.1; this comparison not only shows the benefits of exploiting conjugacy structure but also suggests a potential Rao-Blackwellization scheme. Tak-

Algorithm 17 Minibatch Black-Box Variational Inference (BBVI)

Initialize $\tilde{\eta}_{\theta}^{(0)}$ for t = 0, 1, 2, ... do $\mathcal{S} \leftarrow \{\hat{\theta}_s\}$ where $\hat{\theta}_s \sim q(\cdot \mid \tilde{\eta}_{\theta}^{(t)})$ $\hat{k} \sim \text{Uniform}(\{1, 2, ..., K\})$ $\tilde{\eta}_{\theta}^{(t+1)} \leftarrow \tilde{\eta}_{\theta}^{(t)} + \frac{1}{|\mathcal{S}|} \sum_{\hat{\theta} \in \mathcal{S}} \left(\log \frac{p(\hat{\theta})}{q(\hat{\theta})} + K \log p(y^{(\hat{k})} \mid \hat{\theta})\right) \nabla_{\tilde{\eta}_{\theta}} \log q(\hat{\theta})$

ing $\theta = (\phi, z)$ and $q(\theta) = q(\phi)q(z)$ and starting from (5.22) and (5.29), we can write the gradient as

$$\nabla_{\tilde{\eta}_{\theta}} \mathcal{L} = \mathbb{E}_{q(\phi)q(z)} \left[\log \frac{p(\phi, z, y)}{q(\phi)q(z)} \nabla_{\tilde{\eta}_{\phi}} \log q(\phi)q(z) \right]$$
(5.33)
$$\approx \frac{1}{|S|} \sum_{\hat{\phi} \in S} \left(\mathbb{E}_{q(z)} \log \frac{p(\hat{\phi}, z, y)}{q(\hat{\phi})} - \mathbb{E}_{q(z)} \log q(z) \right) \nabla_{\tilde{\eta}_{\phi}} \log q(\hat{\phi})$$
(5.34)

where S is a set of samples with $\hat{\phi} \stackrel{\text{iid}}{\sim} q(\phi)$ for $\hat{\phi} \in S$. Thus if the entropy of the local variational distribution q(z) and the expectations with respect to q(z) of the log density $\log p(\hat{\phi}, z, y)$ can be computed without resorting to Monte Carlo estimation, then the resulting update would likely have a lower variance than the BBVI update that requires sampling over both $q(\phi)$ and q(z).

This comparison also makes clear the advantages of exploiting conjugacy in SVI: when the updates of Section 5.1.1 can be used, neither $q(\phi)$ nor q(z) needs to be sampled. Furthermore, while BBVI uses stochastic gradients in its updates, the SVI algorithm of Section 5.1.1 uses stochastic natural gradients, adapting to the local curvature of the variational family. Computing stochastic natural gradients in BBVI would require both computing the Fisher information matrix of the variational family and solving a linear system with it.

5.1.3 Exploiting reparameterization for some nonconjugate models

While the score function estimator developed for BBVI in Section 5.1.2 is sufficiently general to handle essentially any model, some noncon-

jugate models admit convenient stochastic gradient estimators that can have lower variance. In particular, in settings where the latent variables are continuous (or any discrete latent variables that can be marginalized efficiently) samples from some variational distributions can be reparameterized in a way that enables an alternative stochastic gradient estimator. This technique is related to non-centered reparameterizations [Papaspiliopoulos et al., 2007] and has recently been called the reparameterization trick [Kingma and Welling, 2014, Rezende et al., 2014].

The reparameterization trick applies when samples $\hat{\theta} \sim q(\theta)$, where $q(\theta)$ has parameter $\tilde{\eta}_{\theta}$, can be written as

$$\hat{\theta} = f(\tilde{\eta}_{\theta}, \epsilon) \tag{5.35}$$

where $\epsilon \sim p(\epsilon)$ is a random variable with a distribution $p(\epsilon)$ that does not depend on $\tilde{\eta}_{\theta}$ and where $\nabla_{\tilde{\eta}_{\theta}} f(\tilde{\eta}_{\theta}, \epsilon)$ can be computed efficiently for almost every value of ϵ . In this case, we can compute stochastic estimates of the gradient of the variational objective by first writing a Monte Carlo approximation of the objective function itself:

$$\mathcal{L} = \mathbb{E}_{q(\theta)} \left[\log \frac{p(\theta, y)}{q(\theta)} \right] \approx \frac{1}{|S|} \sum_{\hat{\epsilon} \in S} \log \frac{p(f(\tilde{\eta}_{\theta}, \hat{\epsilon}), y)}{q(f(\tilde{\eta}_{\theta}, \hat{\epsilon}))}$$
(5.36)

where $\hat{\epsilon} \stackrel{\text{iid}}{\sim} p(\epsilon)$ for each $\hat{\epsilon} \in S$. Alternatively, when the variational entropy term $\mathbb{E}_{q(\theta)} \log q(\theta)$ can be computed efficiently, *e.g.*, if the variational distribution is a Gaussian, only the energy term needs to be approximated via Monte Carlo:

$$\mathcal{L} \approx -\mathbb{E}_{q(\theta)} \log q(\theta) + \frac{1}{|S|} \sum_{\hat{\epsilon} \in S} \log p(f(\tilde{\eta}_{\theta}, \hat{\epsilon}), y).$$
(5.37)

This Monte Carlo approximation is a differentiable unbiased estimate of \mathcal{L} as a function of the variational parameter $\tilde{\eta}_{\theta}$, and so we can form a Monte Carlo estimate of the gradient of the variational objective simply by differentiating it:

$$\nabla_{\widetilde{\eta}_{\theta}} \mathcal{L} \approx -\nabla_{\widetilde{\eta}_{\theta}} \mathbb{E}_{q(\theta)} \log q(\theta) + \frac{1}{|S|} \sum_{\hat{\epsilon} \in S} \nabla_{\widetilde{\eta}_{\theta}} \log p(f(\widetilde{\eta}_{\theta}, \hat{\epsilon}), y).$$
(5.38)

This estimator often has lower variance than the fully general score function estimator [Kingma and Welling, 2014] and can be easier to compute.

5.2 Streaming variational Bayes (SVB)

Streaming variational Bayes (SVB) provides an alternative framework in which to derive minibatch-based scalable variational inference [Broderick et al., 2013]. While the methods of Section 5.1 generally apply stochastic gradient optimization algorithms to a fixed variational mean field objective, SVB instead considers the streaming data setting, in which case there may be no fixed dataset size and hence no fixed variational objective. To handle streaming data, the SVB approach is based on the classical idea of Bayesian updating, in which a posterior is updated to reflect new data as they become available. This sequence of posteriors is approximated by a sequence of variational models, and each variational model is computed from the previous variational model via an incremental update on new data.

More concretely, given a prior $p(\theta)$ over a parameter θ and a (possibly infinite) sequence of data minibatches $y^{(1)}, y^{(2)}, \ldots$, each distributed independently according to a likelihood distribution $p(y^{(k)} | \theta)$, we consider the sequence of posteriors

$$p(\theta | y^{(1)}, \dots, y^{(t)}), \quad t = 1, 2, \dots$$
 (5.39)

Given an approximation updating algorithm \mathcal{A} one can compute a corresponding sequence of approximations

$$p(\theta | y^{(1)}, \cdots, y^{(t)}) \approx q_t(\theta) = \mathcal{A}\left(y^{(t)}, q_{t-1}(\theta)\right), \quad t = 1, 2, \dots$$
 (5.40)

with $q_0(\theta)p(\theta)$. This sequential updating view naturally suggests an online or one-pass algorithm in which the update (5.40) is applied successively to each of a sequence of minibatches.

A sequence of such updates may also exploit parallel or distributed computing resources. For example, the sequence of approximations may

5.2. Streaming variational Bayes (SVB)

be computed as

$$p(\theta \mid y^{(1)}, \cdots, y^{(K_t)}) \approx q_t(\theta)$$
(5.41)

$$= q_{t-1}(\theta) \left(\prod_{k=K_{t-1}+1}^{K_t} \mathcal{A}\left(y^{(k)}, q_{t-1}(\theta)\right) q_{t-1}(\theta)^{-1} \right)$$
(5.42)

where $K_{t-1} + 1, K_{t-1} + 2, \ldots, K_t$ indexes a set of data minibatches for which each update is computed in parallel before being combined in the final update from $q_{t-1}(\theta)$ to $q_t(\theta)$.

This combination of partial results is especially appealing when the prior $p(\theta)$ and the family of approximating distributions $q(\theta)$ are in the same exponential family,

$$p(\theta) \propto \exp\{\langle \eta, t(\theta) \rangle\} \quad q_0(\theta) \propto \exp\{\langle \tilde{\eta}_0, t(\theta) \rangle\}$$
 (5.43)

$$q_t(\theta) = \mathcal{A}(y^{(k)}, q_{t-1}(\theta)) \propto \exp\left\{\langle \tilde{\eta}_t, \ t(\theta) \rangle\right\}$$
(5.44)

for a prior natural parameter η and a sequence of variational parameters $\tilde{\eta}_t$. In the exponential family case, the updates (5.42) can be written

$$p(\theta \mid y^{(1)}, \cdots, y^{(K_t)}) \approx q_t(\theta) \propto \exp\left\{\langle \tilde{\eta}_t, t(\theta) \rangle\right\}$$

$$= \exp\left\{\langle \tilde{\eta}_{t-1} + \sum_k (\tilde{\eta}_k - \tilde{\eta}_{t-1}), t(\theta) \rangle\right\}$$
(5.45)
(5.46)

where we may take the algorithm \mathcal{A} to return an updated natural parameter, $\tilde{\eta}_k = \mathcal{A}(y^{(k)}, \tilde{\eta}_t)$.

Finally, similar updates can be performed in an asynchronous distributed master-worker setting. Each worker can process a minibatch and send the corresponding natural parameter increment to a master process, which updates the global variational parameter and transmits back the updated variational parameter along with a new data minibatch. In symbols, we can write that a worker operating on minibatch $y^{(k)}$ for some minibatch index k computes the update increment $\Delta \tilde{\eta}_k$ according to

$$\Delta \widetilde{\eta}_k = \mathcal{A}(y^{(k)}, \widetilde{\eta}_{\tau(k)}) \tag{5.47}$$

Algorithm 18 Streaming Variational Bayes (SVB)

Initialize $\tilde{\eta}_0$
for each worker $p = 1, 2, \ldots, P$ do
Send task $(y^{(p)}, \tilde{\eta}_0)$ to worker p
${\bf as}$ workers send updates ${\bf do}$
Receive update $\Delta \tilde{\eta}_k$ from worker p
$\widetilde{\eta}_{t+1} \leftarrow \widetilde{\eta}_t + \Delta \widetilde{\eta}_k$
Retrieve new data minibatch index k^\prime
Send new task $(y^{(k')}, \tilde{\eta}_{t+1})$ to worker p

Algorithm 19 SVB Worker Process	
repeat	
Receive task $(y^{(k)}, \tilde{\eta}_t)$ from master	
$\Delta \widetilde{\eta}_k \leftarrow \mathcal{A}(y^{(k)}, \widetilde{\eta}_t) - \widetilde{\eta}_t$	
Send update $\Delta \tilde{\eta}_k$ to master	
until no tasks remain	

where $\tau(k)$ is the index of the global variational parameter used in the worker's computation. Upon receiving an update, the master updates its global variational parameter synchronously according to

$$\widetilde{\eta}_{t+1} = \widetilde{\eta}_t + \Delta \widetilde{\eta}_k. \tag{5.48}$$

We summarize a version of this process in Algorithms 18 and 19.

A related algorithm, which we do not detail here, is Memoized Variational Inference (MVI) [Hughes and Sudderth, 2013, Hughes et al., 2015]. While this algorithm is designed for the fixed dataset setting rather than the streaming setting, the updates can be similar to those of SVB. In particular, MVI optimizes the mean field objective in the conjugate exponential family setting using the mean field coordinate descent algorithm but with an atypical update order, in which only some local variational factors are updated at a time. This update order enables minibatch-based updating but, unlike the stochastic gradient algorithms, does not optimize out the other local variational factors not included in the minibatch and instead leaves them fixed.

5.3. Summary

Streaming variational inference algorithms similar to SVB have also recently been studied in some Bayesian nonparametric mixture models [Tank et al., 2015].

5.3 Summary

Stochastic gradients vs. streaming. The methods of Section 5.1 apply stochastic optimization to variational mean field inference objectives. In optimization literature and practice, stochastic gradient methods have a large body of both theoretical and empirical support, and so such methods offer a compelling framework for scalable inference. The streaming ideas surveyed in Section 5.2 are less well understood, but by treating the streaming setting, rather than the setting of a large fixed-size dataset, they may extend the reach of Bayesian modeling and inference.

Minibatching. All of this chapter's scalable approaches to mean field variational inference are based on processing minibatches of data. These algorithms arrive at this data access pattern via two routes: the first applies stochastic gradient optimization to mean field variational inference (§5.1) and the second considers the streaming data setting (§5.2). SVI (§5.1.1) and BBVI (§5.1.2) optimize the variational objective by replacing full gradient updates with stochastic gradient updates. In both SVI and BBVI, these approximate gradients arise from randomly sampling data minibatches, while in BBVI there is additional stochasticity due to the Monte Carlo approximation required to handle nonconjugate structure. In contrast to SVI and BBVI, SVB (§5.2) processes data minibatches to drive incremental posterior updates, constructing a sequence of approximate posterior distributions that correspond to classical sequential Bayesian updating without having a single fixed objective to optimize.

Generality, requirements, and assumptions. As with most approaches to scaling MCMC samplers for Bayesian inference, these minibatch-based variational inference methods depend on model struc-

ture. In SVI and scalable BBVI, minibatches map to terms in a factorization of the joint probability. In SVB, minibatches map to a sequence of likelihoods to be incorporated into the variational posterior. Some of these methods further depend on and exploit exponential family and conjugacy structure. SVI is based on complete-data conjugacy, while BBVI was specifically developed for nonconjugate models. SVB is a general framework, but in the conjugate exponential family case the updates can be written in terms of simple updates to natural parameters. A direction for future research might be to develop new methods based on identifying and exploiting some 'middle ground' between the structural requirements of SVI and BBVI, or similarly of SVB with and without exponential family structure.

5.4 Discussion

The minibatch-based variational inference methods Parallel variants. developed in this chapter suggest parallel and asynchronous variants. In the case of SVB, distributed and asynchronous versions, such as the master-worker pattern depicted by Algorithms 18 and 19, have been empirically studied Broderick et al. [2013]. However, we lack theoretical understanding about these procedures, and it is unclear how to define and track notions of convergence or stability. Methods based on stochastic gradients, such as SVI, can naturally be extended to exploit parallel and asynchronous (or "Hogwild") variants of stochastic gradient ascent. In such parallel settings, these optimization-based techniques benefit from powerful gradient convergence results [Bertsekas and Tsitsiklis, 1989, Section 7.8, though tuning such algorithms is still a challenge. Other parallel versions of these ideas and algorithms have also been developed in Campbell and How [2014] and Campbell et al. [2015].

6

Challenges and questions

In this review, we have examined a variety of different views on scaling Bayesian inference up to large datasets and greater model complexity and out to parallel compute resources. Several different themes have emerged, from techniques that exploit subsets of data for computational savings to proposals for distributing inference computations across multiple machines. Progress is being made, but there remain significant open questions and outstanding challenges to be tackled as this research programme moves forward.

Trading off errors in MCMC One of the key insights underpinning much of the recent work on scaling Bayesian inference can be framed in terms of a kind of bias-variance tradeoff. Traditional MCMC theory provides asymptotically unbiased estimators for which the error can eventually be driven arbitrarily small. However, in practice, under limited computational budgets the error can be significant. This error has two components: transient bias, in which the samples produced are too dependent on the Markov chain's initialization, and Monte Carlo standard error, in which the samples collected may be too few or too highly correlated to produce good estimates.



Figure 6.1: A simulation illustrating the error terms in traditional MCMC estimators as a function of wall-clock time (log scale). The marginal distributions of the Markov chain iterates converge to the target distribution (top panel), while the errors in MCMC estimates due to transient bias and Monte Carlo standard error are driven arbitrarily small.

Figure 6.1 illustrates the error regimes and tradeoffs in traditional MCMC.¹ Asymptotic analysis describes the regime on the right of the plot, after the sampler has mixed sufficiently well. In this regime, the marginal distribution of each sample is essentially equal to the target distribution, and the transient bias from initialization, which affects only the early samples in the Monte Carlo sum, is washed out rapidly at least at a $\mathcal{O}(\frac{1}{n})$ rate. The dominant source of error is due to Monte Carlo standard error, which diminishes only at a $\mathcal{O}(\frac{1}{\sqrt{n}})$ rate.

However, machine learning practitioners using MCMC often find themselves in another regime: in the middle of the plot, the error is decreasing but dominated instead by the transient bias. The challenge in practice is often to get through this regime, or even to get into it at all. When the underlying Markov chain does not mix sufficiently well or when the transitions cannot be computed sufficiently quickly, getting to this regime may be practically infeasible for a realistic computational budget.

Several of the new MCMC techniques we have studied aim to address this challenge. In particular, the parallel predictive prefetching method of Section 4.1.2 accelerates this phase of MCMC without affecting the stationary distribution. Other methods instead introduce approximate transition operators that can be executed more efficiently. For example, the adaptive subsampling methods of Section 3.2 and the stochastic gradient sampler of Section 3.4 can execute updates more efficiently by operating only on data subsets, while the Weierstrass and Hogwild Gibbs samplers of Sections 4.2.1 and 4.2.2, respectively, execute more quickly by leveraging data parallelism. These transition operators are approximate in that they do not admit the exact target distribution as a stationary distribution: instead, the stationary distribution is only intended to be close to the target. Framed in terms of Monte Carlo estimates, these approximations effectively accelerate the execution of the chain at the cost of introducing an asymptotic bias. Figure 6.2 illustrates this new tradeoff.

Allowing some asymptotic bias to reduce transient bias or even Monte Carlo variance is likely to enable MCMC inference at a new

¹See also Section 2.2.4



Figure 6.2: A simulation illustrating the new tradeoffs in some proposed scalable MCMC methods. Compare to Figure 6.1. As a function of wall-clock time (log scale), the Markov chain iterations execute more than an order of magnitude faster, and hence the marginal distributions of the Markov chain iterates converge to the stationary distribution more quickly; however, because the stationary distribution is not the target distribution, an asymptotic bias remains (top panel). Correspondingly, MCMC estimator error, particularly the transient bias, can be driven to a small value more rapidly, but there is an error floor due to the introduction of the asymptotic bias (bottom panel).

scale. However, both the amount of asymptotic bias introduced by these methods and the ways in which it depends on model and algorithm parameters remain unclear. More theoretical understanding and empirical study is necessary to guide machine learning practice.

Scaling limits of Bayesian inference Scalability in the context of Bayesian inference is ultimately about spending computational resources to better interrogate posterior distributions. It is therefore important to consider whether there are fundamental limits to what can be achieved by, *e.g.*, spending more money on Amazon EC2, for either faster computers or more of them.

In parallel systems, linear scaling is ideal: twice as much computational power yields twice as much useful work. Unfortunately, even if this lofty parallel speedup goal is achieved, the asymptotic picture for MCMC is dim: in the asymptotic regime, doubling the number of samples collected can only reduce the Monte Carlo standard error by a factor of $\sqrt{2}$. This scaling means that there are diminishing returns to purchasing additional computational resources, even if those resources provide linear speedup in terms of accelerating the execution of the MCMC algorithm.

Interestingly, variational methods may not suffer from such intrinsic limits. In particular, the stochastic gradient variational inference methods surveyed in Section 5.1 can utilize optimization methods that, at least in smooth convex settings, can converge at least at $\mathcal{O}(\frac{1}{n})$ rates [Bubeck, 2015]. When these convergence properties are maintained for achieving local minima in non-convex problems, applying additional computational resources would not inherently suffer from the problem of diminishing marginal returns.

Measuring performance With all the ideas surveyed here, one thing is clear: there are many alternatives for how to scale Bayesian inference. How should we compare these alternative algorithms? Can we tell when any of these algorithms work well in an absolute sense?

One standard approach for evaluating MCMC procedures is to define a set of scalar-valued test functions (or estimands of interest) and compute effective sample size [Gelman et al., 2014, Section 11.5] as a function of wall-clock time. However, in complex models designing an appropriately comprehensive set of test functions may be difficult. Furthermore, many such measures require the Markov chain to mix and do not account for any asyptotic bias [Gorham and Mackey, 2015], hence limiting their applicability to measuring the performance of many of the new inference methods studied here.

To confront these challenges, one recently-proposed approach [Gorham and Mackey, 2015] draws on Stein's method, classically used as an analytical tool, to design an efficiently-computable measure of discrepancy between a target distribution and a set of samples. A natural measure of discrepancy between a target density p(x) and a (weighted) sample distribution q(x), where $q(x) = \sum_{i=1}^{n} w_i \delta_{x_i}(x)$ for some set of samples $\{x_i\}_{i=1}^{n}$ and weights $\{w_i\}_{i=1}^{n}$, is to consider their largest absolute difference across a large class of test functions:

$$d_{\mathcal{H}}(q,p) = \sup_{h \in \mathcal{H}} |\mathbb{E}_q h(X) - \mathbb{E}_p h(X)|$$
(6.1)

where \mathcal{H} is the class of test functions. While expectations with respect to the target density p may be difficult to compute, by designing \mathcal{H} such that $\mathbb{E}_p h(X) = 0$ for every $h \in \mathcal{H}$, we need only compute expectations with respect to the sample distribution q. To meet this requirement, instead of designing \mathcal{H} directly, we can instead choose \mathcal{H} to be the image of another function class \mathcal{G} under an operator \mathcal{T}_p that may depend on p, so that $\mathcal{H} = \mathcal{T}_p \mathcal{G}$ and the requirement becomes $\mathbb{E}_p(\mathcal{T}_p g)(x) = 0$ and the discrepancy measure becomes

$$d_{\mathcal{T}_p\mathcal{G}}(q,p) = \sup_{g \in \mathcal{G}} |\mathbb{E}_q(\mathcal{T}_p g)(X)|.$$
(6.2)

Such operators \mathcal{T}_p can be designed using infinitessimal generators from continuous-time ergodic Markov processes, and Gorham and Mackey [2015] suggest using the operator

$$(\mathcal{T}_p g)(x) \triangleq \langle g(x), \nabla \log p(x) \rangle + \langle \nabla, \nabla g(x) \rangle$$
(6.3)

which requires computing only the gradient of the target log density. Furthermore, while the optimization in (6.2) is infinite-dimensional in general and might have infinitely many smoothness constraints from \mathcal{G} ,
Gorham and Mackey [2015] shows that for the sample distribution q the test function g need only be evaluated at the finitely-many sample points $\{x_i\}_{i=1}^n$ and that only a small number of constraints must be enforced. This new performance metric does not require assumptions on whether the samples are generated from an unbiased, stationary Markov chain, and so it may provide clear ways to compare across a broad spectrum sampling-based approximate inference algorithms.

Another recently-proposed approach attempts to estimate or bound the KL divergence from an algorithm's approximate posterior representation to the true posterior, at least when applied to synthetic data. This approach, called bidirectional Monte Carlo (BDMC) [Grosse et al., 2015], can be applied to measure the performance of both variational mean field algorithms as well as annealed importance sampling (AIS) and sequential Monte Carlo (SMC) algorithms. By rearranging the variational identity (2.55), we can write the KL divergence KL(q||p) from an approximating distribution $q(z, \theta)$ to a target posterior $p(z, \theta | \bar{y})$ in terms of the log marginal likelihood log $p(\bar{y})$ and an expectation with respect to $q(z, \theta)$:

$$\operatorname{KL}(q\|p) = \log p(\bar{y}) - \mathbb{E}_{q(z,\theta)} \left[\log \frac{p(z,\theta \mid \bar{y})}{q(z,\theta)} \right].$$
(6.4)

Because the expectation can be readily computed in a mean field setting or stochastically lower-bounded when using AIS [Grosse et al., 2015, Section 4.1], with a stochastic upper bound on $\log p(\bar{y})$ we can use (6.4) to compute a stochastic upper bound on the KL divergence $\mathrm{KL}(q||p)$. BDMC provides a method to compute such stochastic upper bounds on $\log p(\bar{y})$ for synthetic datasets \bar{y} , and so may enable new performance metrics that apply to both sampling-based algorithms as well as variational mean field algorithms. However, while MCMC transition operators are used to construct AIS algorithms, BDMC does not directly apply to evaluating the performance of such transition operators in standard MCMC inference.

Developing performance metrics and evaluation procedures is critical to making progress. As observed in Grosse et al. [2015],

In many application areas of machine learning, especially supervised learning, benchmark datasets have spurred rapid progress in developing new algorithms and clever refinements to existing algorithms. [...] So far, the lack of quantitative performance evaluations in marginal likelihood estimation, and in sampling-based inference more generally, has left us fumbling around in the dark.

By developing better ways to measure the performance of these Bayesian inference algorithms, we will be much better equipped to compare, improve, and extend them.

Acknowledgements

This work was funded in part by NSF IIS-1421780 and the Alfred P. Sloan Foundation. E.A. is supported by the Miller Institute for Basic Research in Science, University of California, Berkeley. M.J. is supported by a fellowship from the Harvard/MIT Joint Grants program.

- Talal M. Alkhamis, Mohamed A. Ahmed, and Vu Kim Tuan. Simulated annealing for discrete optimization with estimation. *European Journal of Operational Research*, 116(3):530–544, 1999.
- Shun-ichi Amari and Hiroshi Nagaoka. *Methods of Information Geometry*. American Mathematical Society, 2007.
- Christophe Andrieu and Eric Moulines. On the ergodicity properties of some adaptive MCMC algorithms. *The Annals of Applied Probability*, 16(3): 1462–1505, 2006.
- Christophe Andrieu and Gareth O. Roberts. The pseudo-marginal approach for efficient Monte Carlo computations. *Annals of Statistics*, pages 697–725, 2009.
- Christophe Andrieu and Johannes Thoms. A tutorial on adaptive MCMC. Statistics and Computing, 18(4):343–373, 2008.
- Christophe Andrieu, Arnaud Doucet, and Roman Holenstein. Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society Series B*, 72(3):269–342, 2010.
- Elaine Angelino. Accelerating Markov chain Monte Carlo via parallel predictive prefetching. PhD thesis, School of Engineering and Applied Sciences, Harvard University, 2014.
- Elaine Angelino, Eddie Kohler, Amos Waterland, Margo Seltzer, and Ryan P. Adams. Accelerating MCMC via parallel predictive prefetching. In 30th Conference on Uncertainty in Artificial Intelligence, 2014.

- Arthur U. Asuncion, Padhraic Smyth, and Max Welling. Asynchronous distributed learning of topic models. In Advances in Neural Information Processing Systems 21, pages 81–88, 2008.
- Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Explorationexploitation tradeoff using variance estimates in multi-armed bandits. *The*oretical Computer Science, 410(19):1876–1902, 2009.
- Rémi Bardenet and Odalric-Ambrym Maillard. Concentration inequalities for sampling without replacement. *Bernoulli*, 20(3):1361–1385, 2015.
- Rémi Bardenet, Arnaud Doucet, and Chris Holmes. Towards scaling up Markov chain Monte Carlo: An adaptive subsampling approach. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.
- Mark A. Beaumont. Estimation of population growth or decline in genetically monitored populations. *Genetics*, 164(3):1139–60, 2003.
- Dimitri P. Bertsekas and John N. Tsitsiklis. Parallel and Distributed Computation: Numerical Methods. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1989.
- Christopher M. Bishop. Pattern Recognition and Machine Learning (Information Science and Statistics). Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- Léon Bottou. On-line learning and stochastic approximations. In David Saad, editor, On-line Learning in Neural Networks, pages 9–42. Cambridge University Press, New York, NY, USA, 1998.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends in Machine Learning, 3(1):1–122, January 2011.
- A. E. Brockwell. Parallel Markov chain Monte Carlo simulation by prefetching. Journal of Computational and Graphical Statistics, 15(1):246–261, March 2006.
- Tamara Broderick, Nicholas Boyd, Andre Wibisono, Ashia C. Wilson, and Michael I. Jordan. Streaming variational Bayes. In Advances in Neural Information Processing Systems 26, pages 1727–1735, 2013.
- Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. Handbook of Markov Chain Monte Carlo. Chapman & Hall/CRC Handbooks of Modern Statistical Methods. CRC press, 2011.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. Foundations and Trends in Machine Learning, 8(3-4), 2015.

- Akif Asil Bulgak and Jerry L. Sanders. Integrating a modified simulated annealing algorithm with the simulation of a manufacturing system to optimize buffer sizes in automatic assembly systems. In *Proceedings of the 20th Conference on Winter Simulation*, pages 684–690, New York, NY, USA, 1988. ACM.
- Trevor Campbell and Jonathan P. How. Approximate decentralized Bayesian inference. In Uncertainty in Artificial Intelligence, 2014.
- Trevor Campbell, Julian Straub, John W. Fisher III, and Jonathan P. How. Streaming, massively parallel variational inference for Bayesian nonparametrics. In Advances in Neural Information Processing Systems 28, 2015.
- Tianqi Chen, Emily B. Fox, and Carlos Guestrin. Stochastic gradient Hamiltonian Monte Carlo. In *Proceedings of the 31st International Conference* on Machine Learning, June 2014.
- J. E. Dennis, Jr. and Robert B. Schnabel. Numerical Methods for Unconstrained Optimization and Nonlinear Equations. Prentice-Hall Series in Computational Mathematics, 1983.
- Finale Doshi-Velez, David A. Knowles, Shakir Mohamed, and Zoubin Ghahramani. Large scale nonparametric Bayesian inference: Data parallelisation in the Indian buffet process. In Advances in Neural Information Processing Systems 22, pages 1294–1302, 2009.
- Arnaud Doucet, Michael Pitt, Robert Kohn, and George Deligiannidis. Efficient implementation of Markov chain Monte Carlo when using an unbiased likelihood estimator. *Biometrika*, 102(2):295–313, 2015.
- Paul Fearnhead, Omiros Papaspiliopoulos, Gareth O. Roberts, and Andrew Stuart. Random-weight particle filtering of continuous time processes. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 72 (4):497–512, 2010.
- Andrew Gelman and Xiao-Li Meng. Simulating normalizing constants: From importance sampling to bridge sampling to path sampling. *Statistical science*, pages 163–185, 1998.
- Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. *Bayesian data analysis*. Chapman and Hall/CRC, 3rd edition, 2014.
- Stuart Geman and Donald Geman. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. Pattern Analysis and Machine Intelligence, IEEE Transactions on, pages 721–741, 1984.
- Charles J. Geyer. Practical Markov chain Monte Carlo. Statistical Science, pages 473–483, 1992.

- Mark Girolami and Ben Calderhead. Riemann manifold Langevin and Hamiltonian Monte Carlo methods. Journal of the Royal Statistical Society: Series B (Statistical Methodology) (With Discussion), 73:123 – 214, 03/2011 2011.
- Joseph Gonzalez, Yucheng Low, Arthur Gretton, and Carlos Guestrin. Parallel Gibbs sampling: From colored fields to thin junction trees. In *International Conference on Artificial Intelligence and Statistics*, pages 324–332, 2011.
- Jackson Gorham and Lester Mackey. Measuring sample quality with stein's method. In Advances in Neural Information Processing Systems, pages 226–234, 2015.
- Thore Graepel, Joaquin Quñonero Candela, Thomas Borchert, and Ralf Herbrich. Web-scale Bayesian click-through rate prediction for sponsored search advertising in Microsoft's Bing search engine. In *Proceedings of the* 27th International Conference on Machine Learning, pages 13–20, 2010.
- Roger B Grosse, Zoubin Ghahramani, and Ryan P Adams. Sandwiching the marginal likelihood using bidirectional Monte Carlo. arXiv preprint arXiv:1511.02543, 2015.
- Heikki Haario, Eero Saksman, and Johanna Tamminen. An adaptive Metropolis algorithm. *Bernoulli*, 7(2):223–242, 04 2001.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. The Elements of Statistical Learning. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- W. K. Hastings. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, 57(1):97–109, April 1970.
- Christian Hipp. Sufficient statistics and exponential families. Ann. Statist., 2 (6):1283–1292, 11 1974.
- Qirong Ho, James Cipar, Henggang Cui, Seunghak Lee, Jin Kyu Kim, Phillip B. Gibbons, Gregory R. Ganger, Garth Gibson, and Eric P. Xing. More effective distributed ML via a stale synchronous parallel parameter server. In Advances in Neural Information Processing Systems 26, 2013.
- Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference. Journal of Machine Learning Research, 14(1): 1303–1347, May 2013.
- Zaiying Huang and Andrew Gelman. Sampling for Bayesian computation with large datasets. Technical report, Columbia University, 2005.
- Michael C Hughes and Erik Sudderth. Memoized online variational inference for Dirichlet process mixture models. In Advances in Neural Information Processing Systems 26, pages 1133–1141, 2013.

- Michael C Hughes, Dae II Kim, and Erik B Sudderth. Reliable and scalable variational inference for the hierarchical Dirichlet process. In *Proceedings* of the Eighteenth International Conference on Artificial Intelligence and Statistics, pages 370–378, 2015.
- Alexander Ihler and David Newman. Understanding errors in approximate distributed latent Dirichlet allocation. *IEEE Transactions on Knowledge* and Data Engineering, 24(5):952–960, 2012.
- Pierre E. Jacob and Alexandre H. Thiery. On nonnegative unbiased estimators. Annals of Statistics, 43(2):769–784, 04 2015.
- Matthew Johnson, James Saunderson, and Alan Willsky. Analyzing Hogwild parallel Gaussian Gibbs sampling. In Advances in Neural Information Processing Systems 26, pages 2715–2723, 2013.
- Matthew James Johnson. Bayesian Time Series Models and Scalable Inference. PhD thesis, Massachusetts Institute of Technology, 2014.
- Robert W. Keener. *Theoretical Statistics: Topics for a Core Course*. Springer Texts in Statistics. Springer New York, 2010.
- Diederik P. Kingma and Max Welling. Auto-encoding variational Bayes. International Conference on Learning Representations, 2014.
- Jack PC Kleijnen and Reuven Y Rubinstein. Optimization and sensitivity analysis of computer simulation models by the score function method. *European Journal of Operational Research*, 88(3):413–427, 1996.
- Daphne Koller and Nir Friedman. Probabilistic Graphical Models: Principles and Techniques. MIT Press, 2009.
- Anoop Korattikara, Yutian Chen, and Max Welling. Austerity in MCMC land: Cutting the Metropolis-Hastings budget. In *Proceedings of the 31th International Conference on Machine Learning*, volume 32, pages 181–189, 2014.
- Neil D. Lawrence. Modelling in the context of massively missing data, 2015. URL http://staffwww.dcs.shef.ac.uk/people/N.Lawrence/ talks/missingdata_tuebingen15.pdf.
- L. Lin, K. F. Liu, and J. Sloan. A noisy Monte Carlo algorithm. *Physical Review D*, 61:074505, Mar 2000.
- Zhiyuan Liu, Yuzhou Zhang, Edward Y. Chang, and Maosong Sun. PLDA+: Parallel latent Dirichlet allocation with data placement and pipeline processing. ACM Transactions on Intelligent Systems and Technology, 2(3): 26:1–26:18, May 2011.

- Anne-Marie Lyne, Mark Girolami, Yves Atchaé, Heiko Strathmann, and Daniel Simpson. On Russian roulette estimates for Bayesian inference with doubly-intractable likelihoods. *Statistical Science*, 30(4):443–467, 11 2015.
- David J. C. MacKay. Information Theory, Inference & Learning Algorithms. Cambridge University Press, New York, NY, USA, 2002.
- Dougal Maclaurin and Ryan P. Adams. Firefly Monte Carlo: Exact MCMC with subsets of data. In *30th Conference on Uncertainty in Artificial Intelligence*, 2014.
- Nicholas Metropolis, Arianna W. Rosenbluth, Marshall N. Rosenbluth, Augusta H. Teller, and Edward Teller. Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6):1087–1092, 1953.
- Sean Meyn and Richard L. Tweedie. *Markov Chains and Stochastic Stability*. Cambridge University Press, New York, NY, USA, 2nd edition, 2009.
- Volodymyr Mnih, Csaba Szepesvári, and Jean-Yves Audibert. Empirical Bernstein stopping. In Proceedings of the 25th International Conference on Machine Learning, pages 672–679, New York, NY, USA, 2008. ACM.
- Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- Radford M. Neal. An improved acceptance procedure for the hybrid Monte Carlo algorithm. J. Comput. Phys., 111(1):194–203, March 1994.
- Radford M. Neal. MCMC using Hamiltonian dynamics. In Handbook of Markov Chain Monte Carlo, Chapman & Hall/CRC Handbooks of Modern Statistical Methods, pages 113–162. CRC Press, 2010.
- Willie Neiswanger, Chong Wang, and Eric Xing. Asymptotically exact, embarrassingly parallel MCMC. In 30th Conference on Uncertainty in Artificial Intelligence, 2014.
- David Newman, Padhraic Smyth, Max Welling, and Arthur U. Asuncion. Distributed inference for latent Dirichlet allocation. In Advances in Neural Information Processing Systems 20, pages 1081–1088, 2007.
- David Newman, Arthur U. Asuncion, Padhraic Smyth, and Max Welling. Distributed algorithms for topic models. *Journal of Machine Learning Re*search, 10:1801–1828, 2009.
- Robert Nishihara, Iain Murray, and Ryan P. Adams. Parallel MCMC with generalized elliptical slice sampling. *Journal of Machine Learning Research*, 15:2087–2112, 2014.

- Omiros Papaspiliopoulos. A methodological framework for Monte Carlo probabilistic inference for diffusion processes. Technical report, Centre for Research in Statistical Methodology, University of Warwick, June 2009.
- Omiros Papaspiliopoulos, Gareth O Roberts, and Martin Sköld. A general framework for the parametrization of hierarchical models. *Statistical Science*, pages 59–73, 2007.
- Sam Patterson and Yee Whye Teh. Stochastic gradient Riemannian Langevin dynamics on the probability simplex. In Advances in Neural Information Processing Systems 26, 2013.
- Maxim Rabinovich, Elaine Angelino, and Michael I. Jordan. Variational Consensus Monte Carlo. In Advances in Neural Information Processing Systems 28, 2015.
- Rajesh Ranganath, Chong Wang, David M. Blei, and Eric P. Xing. An adaptive learning rate for stochastic variational inference. In *International Conference on Machine Learning*, volume 28, pages 298–306, 2013.
- Rajesh Ranganath, Sean Gerrish, and David M. Blei. Black box variational inference. In International Conference on Artificial Intelligence and Statistics, pages 814–822, 2014.
- Benjamin Recht, Christopher Re, Stephen J. Wright, and Feng Niu. Hogwild!: A lock-free approach to parallelizing stochastic gradient descent. In Advances in Neural Information Processing Systems 24, pages 693–701, 2011.
- Jeffrey Regier, Andrew Miller, Jon McAuliffe, Ryan P. Adams, Matt Hoffman, Dustin Lang, David Schlegel, and Prabhat. Celeste: Variational inference for a generative model of astronomical images. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 2095–2103, 2015.
- Danilo J Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning*, pages 1278–1286, 2014.
- Christian P. Robert and George Casella. *Monte Carlo Statistical Methods* (Springer Texts in Statistics). Springer-Verlag New York, Inc., 2004.
- Gareth O. Roberts and Richard L. Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, 2(4): 341–363, 12 1996.
- Gareth O. Roberts, Andrew Gelman, and Walter R. Gilks. Weak convergence and optimal scaling of random walk Metropolis algorithms. *Annals* of Applied Probability, 7:110–120, 1997.

- Issei Sato and Hiroshi Nakagawa. Approximation analysis of stochastic gradient Langevin dynamics by using Fokker-Planck equation and Ito process. In *International Conference on Machine Learning*, pages 982–990, 2014.
- Steven L. Scott, Alexander W. Blocker, and Fernando V. Bonassi. Bayes and big data: The consensus Monte Carlo algorithm. In *Bayes 250*, 2013.
- R. J. Serfling. Probability inequalities for the sum in sampling without replacement. *The Annals of Statistics*, 2(1):39–48, 1974.
- Sameer Singh, Michael L. Wick, and Andrew McCallum. Monte Carlo MCMC: Efficient inference by approximate sampling. In Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, pages 1104–1113, 2012.
- Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine learning algorithms. In Advances in Neural Information Processing Systems 25, pages 2951–2959. 2012.
- David H. Stern, Ralf Herbrich, and Thore Graepel. Matchbox: Large scale online Bayesian recommendations. In Proceedings of the 18th International Conference on World Wide Web, pages 111–120, 2009.
- Alex Tank, Nicholas Foti, and Emily Fox. Streaming variational inference for Bayesian nonparametric mixture models. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, pages 968–976, 2015.
- Wolfgang Wagner. Unbiased Monte Carlo evaluation of certain functional integrals. *Journal of Computational Physics*, 71(1):21–33, 1987.
- Martin J. Wainwright and Michael I. Jordan. Graphical models, exponential families, and variational inference. Foundations and Trends in Machine Learning, 1(1-2):1–305, November 2008.
- Ling Wang and Liang Zhang. Stochastic optimization using simulated annealing with hypothesis test. *Applied Mathematics and Computation*, 174 (2):1329–1342, 2006.
- Xiangyu Wang and David B. Dunson. Parallel MCMC via Weierstrass sampler. arXiv preprint arXiv:1312.4605, 2013.
- Karl Weierstrass. Über die analytische darstellbarkeit sogenannter willkrlicher functionen einer reellen vernderlichen. Sitzungsberichte der Königlich Preuischen Akademie der Wissenschaften zu Berlin, 1885. (II). Erste Mitteilung (part 1) pp. 633–639, Zweite Mitteilung (part 2) pp. 789–805.
- Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th International Conference* on Machine Learning, 2011.