

Neural Information **Processing Systems** Foundation

Distributed Bayesian Posterior Sampling via Moment Sharing Minjie Xu¹, Balaji Lakshminarayanan², Yee Whye Teh³, Jun Zhu¹, and Bo Zhang¹ ¹Dept. of CST, Tsinghua University; ²Gatsby Unit, University College London; ³Dept. of Stats, University of Oxford

CONTRIBUTIONS

We proposed a **distributed MCMC inference** algorithm for large scale **Bayesian posterior simulation** that

- scales to "big data" and compute nodes,
- converges fast and is exempt from a final combination stage,
- supports flexible distributed schemes: synchronous, asynchronous decentralized, etc.,
- incurs low communication costs,
- and achieves high approximation accuracy.

We empirically studied the performance of our algorithm and compared it with the state-of-the-art, on Bayesian logistic regression and sparse linear regression with a spike-and-slab prior.

MOTIVATION

Bayesian inference with "Big Data" $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^{N}$ $p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}, \boldsymbol{\theta}) = p_0(\boldsymbol{\theta}) \cdot \prod_{i=1}^N p(\mathbf{x}_i|\boldsymbol{\theta})$

Variational inference: use $q(\boldsymbol{\theta}) \in \mathcal{Q}$.

$$\mathcal{L}(q) = \mathbb{E}_q[\log p(\mathcal{D}, \theta)] + H(q)$$
 (evidence lower bound)

• MCMC sampling: use $\Theta = \{ \boldsymbol{\theta}^{(t)} \}_{t=1}^T$.

 $A(\boldsymbol{\theta}^{(t)} \to \boldsymbol{\theta}') = \min\left(1, \frac{p(\boldsymbol{\theta}'|\mathcal{D})}{p(\boldsymbol{\theta}^{(t)}|\mathcal{D})} \frac{q(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}')}{q(\boldsymbol{\theta}'|\boldsymbol{\theta}^{(t)})}\right)$ (acceptance rate)

In both cases, the cost of one iteration (sample) is O(N).

To scale up, people generally resort to two types of approaches, namely stochastic approximation (Welling'11, Hoffman'13, etc.) and distributed algorithms (Scott'13, Neiswanger'14, etc.)



SAMPLING VIA MOMENT SHARING (SMS)

Key idea: to encourage mutual awareness and agreement among the local samplers so as to improve inference quality, by enforcing sharing across them a small number of moment statistics of the local posteriors.

| Local posterior: $ilde{p}_i(oldsymbol{ heta}|\mathcal{D}_i) \propto q(oldsymbol{ heta};\eta_{-i}) p(\mathcal{D}_i|oldsymbol{ heta}).$ 1. each local sampler independently draws samples from it; 2. moments of interest: $\mathbb{E}_{\tilde{p}_i(\boldsymbol{\theta}|\mathcal{D}_i)}[S(\boldsymbol{\theta})]$ for some sufficient statistics function $S(\boldsymbol{\theta})$; 3. the effective local prior $q(\theta; \eta_{-i})$ is assumed a member of the exponential family with sufficient statistics $S(\boldsymbol{\theta})$ and natural parameter η_{-i} ; 4. moment sharing: find η_{-i} so that $\mathbb{E}_{\tilde{p}_i(\boldsymbol{\theta}|\mathcal{D}_i)}[S(\boldsymbol{\theta})] = \mu$ for some shared μ . And that's where expectation propagation comes into play.

- 1. approximate each $p(\mathcal{D}_i|\boldsymbol{\theta})$ with $q(\boldsymbol{\theta};\eta_i)$ (use $q(\boldsymbol{\theta};\eta_0)$ for $p_0(\boldsymbol{\theta})$)
- 2. posterior is thus approximated with $q(\theta; \nu)$, where $\nu = \eta_0 + \sum_{i=1}^m \eta_i$
- 3. iteratively solve for each η_i as

 $\operatorname{argmin}_{n_i} \operatorname{KL}(q(\boldsymbol{\theta}; \eta_{-i}) p(\mathcal{D}_i | \boldsymbol{\theta}) || q(\boldsymbol{\theta}; \boldsymbol{\theta})$

which equates matching the moments of the two arguments in $KL(\cdot \| \cdot)$.

Combining the above two ideas, we come up with our SMS algorithm.

EP estimated posterior $q(\theta; \nu)$.

Algorithm 1: UPDATE (ON NODE *i*)

Input: current global natural parameter ν .

- **Output**: updated local natural parameter η_i .
- l compute $\eta_{-i} \leftarrow \nu \eta_i$
- where $\tilde{p}_i(\boldsymbol{\theta}|\mathcal{D}) \propto q(\boldsymbol{\theta}; \eta_{-i}) p(\mathcal{D}_i|\boldsymbol{\theta})$
- **3** compute the empirical moment parameter $\mu_i \leftarrow \frac{1}{T} \sum_{t=1}^T S(\boldsymbol{\theta}_i^{(t)})$
- 4 transform μ_i into its natural parameter $\nu_i \leftarrow \nabla_\mu F^*(\mu_i)$
- **5** update $\eta_i \leftarrow \nu_i \eta_{-i}$
- 6 return η_i

On master node: update $u = \eta_0 + \sum \eta_i$

Multivariate Gaussian family: motivated by the Bernstein-von Mises Theorem for big \mathcal{D}_i

1. sufficient statistics: $S(\boldsymbol{\theta}) = (\boldsymbol{\theta}, \boldsymbol{\theta} \boldsymbol{\theta}^{\top})$ 2. moment parameter: $\mu_i = \frac{1}{T} \sum_{t=1}^{T} \boldsymbol{\theta}_i^{(t)}$,

3. natural parameter: $u_i=(\Omega_i\mu_i,\Omega_i)$ wh

is an unbiased estimate of the precision m

Expectation propagation (Minka'01): variational approximation via moment matching

$$(\eta_{-i})q(\boldsymbol{\theta};\eta_i)) \qquad (\eta_{-i}=\nu-\eta_i=\eta_0+\sum_{j\neq i}\eta_j)$$

The following shows one iteration in the SMS algorithm. Upon convergence, $\nu_i = \nu$, $\mu_i = \mu$, and SMS outputs the collection of samples $\{\Theta_i\}_{i=1}^m$ from the last iteration as well as the



$$\begin{split} & \Sigma_{i} = \frac{1}{T-1} \sum_{t=1}^{T} (\boldsymbol{\theta}_{i}^{(t)} - \mu_{i}) (\boldsymbol{\theta}_{i}^{(t)} - \mu_{i}) \\ & \text{ere} \\ & = \frac{T-d-2}{T-1} \Sigma_{i}^{-1} \\ & \text{atrix.} \end{split}$$

EMPIRICAL RESULTS

base sampler: NUTS (Hoffman'14), burnin = 400, thinning = 2.



Figure 2: Convergence of local posterior means on a smaller Bayesian logistic regression dataset (N = 1000, d = 5). The x-axis indicates the number of likelihood evaluations, with vertical lines denoted EP iteration numbers. The y-axis indicates the estimated posterior means (dimensions indicated by different colours). We show ground truth with solid horizontal lines, the EP estimated mean with asterisks, and local sample estimated means dots connected with dash lines.



Figure 3: Errors (log-scale) against the cumulative number of samples drawn on all nodes (kTm). We tested two random splits of the dataset (hence 2 curves for each algorithm). Each complete EP iteration is highlighted by a vertical grid line.



(a) Approximate KL-divergence

(b) Approximate KL-divergence (c) Approximate KL-divergence

Figure 4: Cross comparison with different numbers of nodes. Note that the x-axes have different meanings. In figure (a), it is the cumulative number of samples drawn locally on each node (kT). For the asynchronous SMS(a), we only plot every m iterations so as to mimic the behaviour of SMS(s)for a more direct comparison. In figure (b) however, it is the cumulative number of likelihood evaluations on each node (kTN/m), which more accurately reflect computation time.





Figure 5: Results on Boston housing dataset for Bayesian sparse linear regression model with spike and slab prior. The x-axis plots the number of data points per node (equals the number of likelihood evaluations per sample) times the cumulative number of samples drawn per node, which is a surrogate for the computation times of the methods. The y-axis plots the ground truth (solid), local sample estimated means (dashed) and EP estimated mean (asterisks) at every iteration.

Algorithm costs		
Algorithm	Local sampling (per node)	Communication (overall)
SCOT	KT samples $\boldsymbol{\theta} \sim p_i(\boldsymbol{\theta})$	$\{\Theta_i\}_{i=1}^m (1 \cdot O(mdKT))$
NEIS(p)	KT samples $\boldsymbol{\theta} \sim p_i(\boldsymbol{\theta})$	$\{\eta_i\}_{i=1}^m (1 \cdot O(m S(\boldsymbol{\theta})))$
NEIS(n)	KT samples $\boldsymbol{\theta} \sim p_i(\boldsymbol{\theta})$	$\{\Theta_i\}_{i=1}^m (1 \cdot O(mdKT))$
WANG	K' iterations, each time KT/K'	all the local samples
	samples $oldsymbol{t} \sim \mathcal{N}(oldsymbol{t} oldsymbol{ heta}, \mathbb{H}) p_i(oldsymbol{t})$	$(K' \cdot O(mdKT/K'))$
SMS	K iterations, each time T samples	$\{\eta_i\}_{i=1}^m$ for K iterations
	$\boldsymbol{\theta} \sim q(\boldsymbol{\theta}; \eta_{-i}) p(\mathcal{D}_i \boldsymbol{\theta})$	$(K \cdot O(m S(\boldsymbol{\theta})))$
where $\pi(\boldsymbol{\theta}) \propto \pi(\boldsymbol{\mathcal{D}} \mid \boldsymbol{\theta}) \pi(\boldsymbol{\theta})^{1/m} \prod_{k=1}^{m} \operatorname{dis}_{\boldsymbol{\theta}} \pi(\boldsymbol{\theta}) = h$ and $ \mathcal{L}(\boldsymbol{\theta}) = O(d^2)$ for		

where $p_i(\theta) \propto p(D_i|\theta)p_0(\theta)^{1/m}$, $\mathbb{H} = \operatorname{diag}(h_1, \ldots, h_d)$ and $|S(\theta)| = O(d^2)$ for multivariate Gaussians

REFERENCES

- [Welling'11] Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient Langevin dynamics n Proceedings of the 28th International Conference on Machine Learning (ICML-11), pages 681-688, 2011.
- . [Hoffman'13] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. The Journal of Machine Learning Research, 14(1):1303-1347, 2013.

3. [Scott'13] Steven L Scott, Alexander W Blocker, Fernando V Bonassi, Hugh A Chipman, Edward I George, and Robert E McCulloch. Bayes and big data: The consensus Monte Carlo algorithm. *EFaBBayes 250 conference*, 16, 2013.

. [Neiswanger'14] Willie Neiswanger, Chong Wang, and Eric Xing. Asymptotically exact, embarrassingly parallel MCMC. In Proceedings of the 30th International Conference on Uncertainty in Artificial Intelligence (UAI-14), pages 623-632, 2014.

5. [Minka'01] Thomas P Minka. A family of algorithms for approximate Bayesian inference. PhD thesis, Aassachusetts Institute of Technology, 2001.

5. [Hoffman'14] Matthew D Hoffman and Andrew Gelman. The No-U-Turn sampler: Adaptively setting path engths in Hamiltonian Monte Carlo. Journal of Machine Learning Research, 15:1593-1623, 2014.

ACKNOWLEDGEMENT

We thank Willie Neiswanger for sharing his implementation of NEIS(n), and Michalis K Titsias for sharing the code used in spike-and-slab sampling. MX, JZ and BZ gratefully acknowledge funding from the National Basic Research Program of China (No. 2013CB329403) and National NSF of China (Nos. 61322308, 61332007). BL gratefully acknowledges generous funding from the Gatsby charitable foundation. YWT gratefully acknowledges EPSRC for research funding through grant EP/K009362/1.