

A versatile distributed MCMC algorithm for large scale inverse problems

P.-A. THOUVENIN^{*}, A. REPETTI^{†‡}, P. CHAINAIS^{*}

^{*}University of Lille, CNRS, Centrale Lille, UMR 9189 CRISTAL, France

[†]Department of Actuarial Mathematics & Stats., Heriot-Watt University, U. K.

[‡]Institute of Sensors, Signals and Systems, Heriot-Watt University, U. K.

EUSIPCO 2022, Aug. 29 – Sep. 2, Belgrade

AUGUST 30, 2022

Work supported by the ANR-20-CHIA-0031-01 and ANR-16-IDEX-0004 projects.

Inverse problems

► Observation model

$$\mathbf{y} = \mathcal{D}(\mathbf{A}\mathbf{x}),$$

$\mathbf{y} \in \mathbb{R}^M$	observations
$\mathbf{x} \in \mathbb{R}^N$	unknown parameters (image, ...)
$\mathbf{A} \in \mathbb{R}^{M \times N}$	measurement operator
$\mathcal{D}: \mathbb{R}^M \rightarrow \mathbb{R}^M$	noise

► **Objective:** find estimate of \mathbf{x} from \mathbf{y}

$$\pi(\mathbf{x} \mid \mathbf{y}) \propto \exp\left(-f_{\mathbf{y}}(\mathbf{A}\mathbf{x}) - g(\mathbf{B}\mathbf{x})\right),$$

$f_{\mathbf{y}}: \mathbb{R}^M \rightarrow]-\infty, +\infty]$	data-fidelity
$g: \mathbb{R}^P \rightarrow]-\infty, +\infty]$, $\mathbf{B} \in \mathbb{R}^{P \times N}$	prior

Challenges

Typical applications (imaging, ...)

- ▶ large number of parameters (N large, from 10^3 to 10^8)
- ▶ large datasets (M large, $M \approx N$)

Scaling with the dimensions of the problem

- ▶ distribute **parameters** and **observations**
- ▶ decompose inference over multiple *workers*
- ▶ limit communication bottlenecks

⇒ Exploit the structure of the problem

⇒ Algorithm amenable to parallel implementation?

Distributed state-of-the-art approaches

Optimization literature

- ▶ *Splitting* methods: large class of parallelizable algorithms: ADMM (Boyd et al. 2011), primal-dual (Chambolle et al. 2011), ...
- ✗ MAP estimator only
- ✓ Distributed implementation (client-server, SPMD)

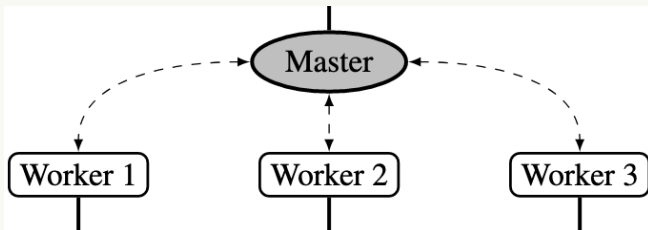
MCMC methods

- ▶ *Asymptotically exact data augmentation* (AXDA) approach (Vono et al. 2021; Rendell et al. 2021)
Divide to conquer strategy (splitting) *via* data augmentation
- ✓ Estimator (MMSE, MAP, ...) + credibility interval
- ⚠ Distributed implementation (**limited to client-server so far**)

Client-server vs SPMD distributed architecture

Client-server

- ▶ each *client* executes a specific task assigned by the *server*
- ▶ information from clients aggregated periodically on the server
- ✗ communications bottleneck (copies across workers, ...)
- ✗ load balancing issues (if drastically different tasks)

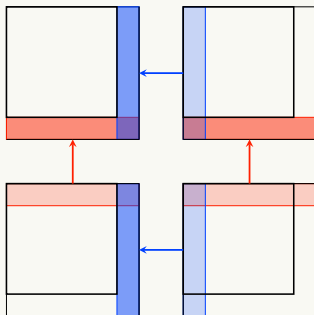


Client-server communications

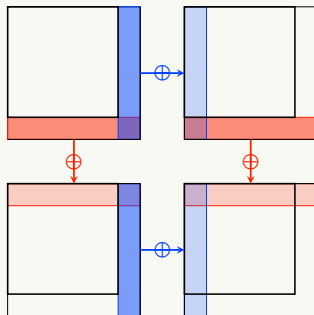
Client-server vs SPMD distributed architecture

Single Program Multiple Data (SPMD) (Darema 2001)

- ▶ all workers execute the same task
- ✓ communications between a worker and its neighbours
- ✓ each worker responsible for a chunk of y and x (**data locality**)



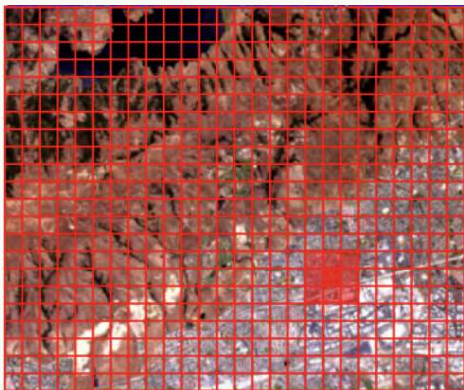
SPMD: exchanging borders



SPMD: aggregating borders

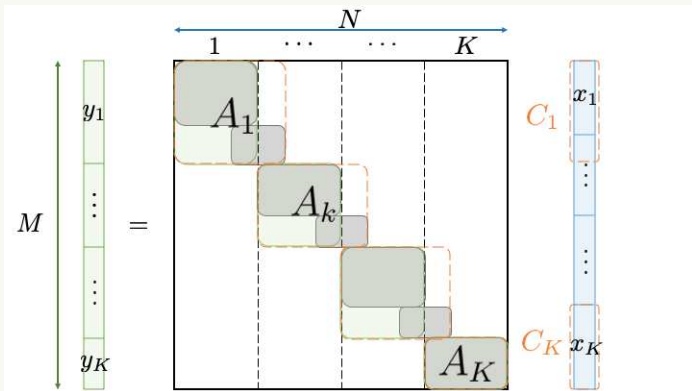
Problem structure

Idea: many problems involve **localized operators** A and B
e.g., inpainting mask, convolution, Laplacian, gradient, ...



Block-sparse model

Assumption: A and B are block-sparse (**localized structure**)



- ▶ **Overlap** between contiguous blocks $\ll \lfloor N/K \rfloor$
- ▶ $(A_k)_{1 \leq k \leq K}$ (dashed orange), $x = (x_k)_{1 \leq k \leq K}$, $x_k \in \mathbb{R}^{N_k}$

Model (I)

Assumption: \mathcal{D} block-separable: $\mathcal{D}(z) = \left(\mathcal{D}_k(z_k) \right)_{1 \leq k \leq K}$

e.g., additive white Gaussian noise, Poisson noise, ...

$$\mathbf{y} = \mathcal{D}(\mathbf{A}\mathbf{x}) = (\mathbf{y}_k)_{1 \leq k \leq K}, \text{ where } (\forall k) \mathbf{y}_k = \mathcal{D}_k(\mathbf{A}_k \mathbf{C}_k \mathbf{x})$$

$\mathbf{C}_k \in \mathbb{R}^{\tilde{N}_k \times N}$ selection operator, with $N \leq \sum_k \tilde{N}_k$
(allows **overlap** in pixel selection)

$\mathbf{A}_k \in \mathbb{R}^{M_k \times \tilde{N}_k}$ local operator

- ▶ **Additively separable** data-fidelity term

$$f_{\mathbf{y}}(\mathbf{A}\mathbf{x}) = \sum_{k=1}^K f_{\mathbf{y}_k}(\mathbf{A}_k \mathbf{C}_k \mathbf{x}), \quad f_{\mathbf{y}_k}: \mathbb{R}^{M_k} \rightarrow]-\infty, +\infty]$$

Model (II)

- ▶ **Additively separable** data-fidelity term

$$f_{\mathbf{y}}(\mathbf{A}\mathbf{x}) = \sum_{k=1}^K f_{\mathbf{y}_k}(\mathbf{A}_k \mathbf{C}_k \mathbf{x}), \quad f_{\mathbf{y}_k}: \mathbb{R}^{M_k} \rightarrow]-\infty, +\infty]$$

- ▶ **g additively separable** (e.g., ℓ_1 -norm (Laplace prior), $\ell_{1,2}$ -norm)

$$g(\mathbf{B}\mathbf{x}) = \sum_{k=1}^K g_k(\mathbf{B}_k \mathbf{D}_k \mathbf{x}) \quad g_k: \mathbb{R}^{P_k} \rightarrow]-\infty, +\infty]$$

$\mathbf{D}_k \in \mathbb{R}^{\bar{N}_k \times N}$ selection operator, $\mathbf{B}_k \in \mathbb{R}^{P_k \times \bar{N}_k}$

- ✓ Problem amenable to SPMD architecture

Algorithm structure compatible with model structure?

Proposed approach

Model of interest

- 1 Block sparse structure for matrices
- 2 Additively separable functions

Towards a SPMD sampler

- 1 Leverage AXDA (Vono et al. 2019b)
 - ▶ **splitting** via approximate data augmentation
- 2 Gibbs sampler
- 3 Langevin kernels to handle non-trivial conditional distributions (e.g., PSGLA (Salim et al. 2020), MYULA (Durmus et al. 2018))
- 4 Parallel implementation (SPMD)

Illustration on image inpainting

- ▶ Observation model: additive white Gaussian noise

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}, \text{ with } \mathbf{w} \sim \mathcal{N}(\mathbf{0}_M, \sigma^2 \mathbf{I}_{M \times M})$$

- ▶ Data fidelity term

$$f_{\mathbf{y}_k}(\mathbf{A}_k \mathbf{C}_k \mathbf{x}) = \frac{1}{2\sigma^2} \|\mathbf{A}_k \mathbf{x}_k - \mathbf{y}_k\|_2^2, \quad \mathbf{A}_k : \text{local inpainting mask}$$

- ▶ Total variation (TV) prior: \mathbf{B} discrete gradient (Chambolle 2004),

$$g_k(\mathbf{B}_k \mathbf{D}_k \mathbf{x}) = \tau \|\mathbf{B}_k \mathbf{D}_k \mathbf{x}\|_{1,2}, \quad \tau > 0.$$

Application of AXDA

- ▶ Approximate posterior distribution $\tilde{\pi}_\alpha$
 - ▶ $\tilde{\pi}_\alpha \rightarrow \pi$ when $\alpha \rightarrow 0$
 - ▶ splitting variable \mathbf{z} , with $p(\mathbf{z} | \alpha) \propto \exp(-\phi_\alpha(\mathbf{z}))$, ϕ_α separable
- ▶ (Optional) additional exact augmentation, variable \mathbf{u}
 - ▶ improves mixing of the chain (Vono et al. 2019a)

$$\tilde{\pi}_\alpha(\mathbf{x}, \mathbf{z}, \mathbf{u} | \mathbf{y}, \beta) \propto \exp\left(-\sum_{k=1}^K h_k(\mathbf{x}, \mathbf{z}_k, \mathbf{u}_k; \alpha, \beta)\right),$$

$$h_k(\mathbf{x}, \mathbf{z}_k, \mathbf{u}_k; \alpha, \beta) = f_{y_k}(\mathbf{A}_k \mathbf{x}_k) + \tau \|\mathbf{z}_k\|_{2,1} + \frac{1}{2\beta} \|\mathbf{u}_k\|_2^2$$

$$+ \frac{1}{2\alpha} \|\mathbf{B}_k \mathbf{D}_k \mathbf{x} - \mathbf{z}_k + \mathbf{u}_k\|_2^2.$$

$$\underbrace{\hspace{10em}}_{\phi_{k,\alpha}(\mathbf{z}_k - \mathbf{u}_k)}$$

Proposed SPMD sampler

For $t \in \{0, \dots, N_{\text{MC}} - 1\}$, on worker k , each sample generated as

// Update \mathbf{x}_k with PSGLA kernel (Salim et al. 2020)

Communications induced by \mathbf{D}_k to compute $\nabla_{\mathbf{x}} h_k$

$$\mathbf{x}_k^{(t+1)} = \mathbf{x}_k^{(t)} - \gamma \nabla_{\mathbf{x}} h_k(\mathbf{x}^{(t)}, \mathbf{z}_k^{(t)}, \mathbf{u}_k^{(t)}; \alpha, \beta) + \sqrt{2\gamma} \xi_k,$$

with $\xi_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{N_k \times N_k})$

Communications to compute $\mathbf{B}_k \mathbf{D}_k \mathbf{x}^{(t+1)}$

// Update \mathbf{z} with PSGLA kernel (Salim et al. 2020)

$$\mathbf{z}_k^{(t+1)} = \text{prox}_{\eta g_k} \left(\mathbf{z}_k^{(t)} - \frac{\eta}{\alpha} (\mathbf{z}_k^{(t)} + \mathbf{B}_k \mathbf{D}_k \mathbf{x}^{(t+1)} - \mathbf{u}_k^{(t)}) + \sqrt{2\eta} \zeta_k \right),$$

with $\zeta_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{M_k \times M_k})$

// Sample \mathbf{u} from its full conditional

$$\mathbf{u}_k^{(t+1)} \mid \mathbf{x}^{(t+1)}, \mathbf{z}_k^{(t+1)} \sim \mathcal{N} \left(\frac{\nu}{\alpha} (\mathbf{z}_k^{(t+1)} - \mathbf{B}_k \mathbf{D}_k \mathbf{x}^{(t+1)}), \nu \mathbf{I}_{P_k \times P_k} \right)$$

Experiment

Simulation settings

- ▶ $M = \lfloor 0.6N \rfloor$ observations, σ^2 such that SNR = 40 dB
- ▶ $N_{MC} = 10^4$ samples, $N_{bi} = 5 \times 10^3$ burn-in, $(\alpha, \beta, \tau) = (9, 1, 0.2)$
- ▶ **Strong scaling** (fixed problem size, increasing K)

Parallel setting

- ▶ HPC computer grid from University of Lille¹
- ▶ Single node: two 2.1 GHz, 18-core, Intel Xeon E5-2695 v4 series processors (36 CPU cores in total)
- ▶ Parallelization: mpi4py (Dalcin et al. 2021) library
(1 worker = 1 process on a CPU core)

¹<https://hpc.univ-lille.fr/>

Strong scaling experiment

K	SNR (MMSE)	SNR (MAP)	Time per iter. ($\times 10^{-3}$ s)	Speedup	Runtime (s)
1 (Vono et al. 2019a)	23.33	22.45	65.56 (2.08)	0.19	262.20
1	23.45	22.95	12.21 (0.63)	1.00	61.04
2	23.46	22.88	6.07 (0.42)	2.01	30.37
4	23.48	22.88	3.50 (0.21)	3.49	17.50
8	23.44	22.86	1.93 (0.77)	6.33	9.63
16	23.48	22.90	1.08 (2.35)	11.30	5.38

Table 1: Strong scaling experiment results.

Experiment results



(a) Ground truth



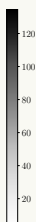
(b) MMSE (Vono 2019)



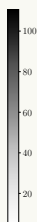
(c) MMSE (prop.)



(a) Obs.



(b) CI (Vono 2019)



(c) CI (prop.)

Conclusions and perspectives

Conclusions: SPMD-distributed sampler (PSGLA within Gibbs)

- ✓ quality comparable to (Vono et al. 2019a);
- ✓ lower runtime (distribution flexibility);
- ✓ strong scaling behaviour.

Conclusions and perspectives

Conclusions: SPMD-distributed sampler (PSGLA within Gibbs)

- ✓ quality comparable to (Vono et al. 2019a);
- ✓ lower runtime (distribution flexibility);
- ✓ strong scaling behaviour.

Perspectives

- infer AXDA parameters α, β ;
- more general applications
 - ↪ inverse problems on hyper-graphs;
 - ↪ assessment on multiple nodes;
- extension to handle asynchronous communications.

Conclusions and perspectives

Conclusions: SPMD-distributed sampler (PSGLA within Gibbs)

- ✓ quality comparable to (Vono et al. 2019a);
- ✓ lower runtime (distribution flexibility);
- ✓ strong scaling behaviour.

Perspectives

- infer AXDA parameters α, β ;
- more general applications
 - ↪ inverse problems on hyper-graphs;
 - ↪ assessment on multiple nodes;
- extension to handle asynchronous communications.

Thank you for your attention.

A versatile distributed MCMC algorithm for large scale inverse problems

P.-A. THOUVENIN^{*}, A. REPETTI^{†‡}, P. CHAINAIS^{*}

^{*}University of Lille, CNRS, Centrale Lille, UMR 9189 CRISTAL, France

[†]Department of Actuarial Mathematics & Stats., Heriot-Watt University, U. K.

[‡]Institute of Sensors, Signals and Systems, Heriot-Watt University, U. K.

EUSIPCO 2022, Aug. 29 – Sep. 2, Belgrade

AUGUST 30, 2022

Work supported by the ANR-20-CHIA-0031-01 and ANR-16-IDEX-0004 projects.

Backup slides

- ▶ References

References I

- Boyd, Stephen et al. (2011). “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers”. In: *Foundations and Trends® in Machine Learning* 3.1, pp. 1–122.
- Chambolle, Antonin (Jan. 2004). “An algorithm for total variation minimization and applications”. In: *J. Math. Imag. Vision* 20.1, pp. 89–97.
- Chambolle, Antonin and Thomas Pock (May 2011). “A First-Order Primal-Dual Algorithm for Convex Problems with Applications to Imaging”. In: *J. Math. Imag. Vision* 40.1, pp. 120–145.
- Dalcin, Lisandro and Yao-Lung L. Fang (2021). “mpi4py: Status Update After 12 Years of Development”. In: *IEEE Comput. Sci. Eng.* 23.4, pp. 47–54.

References II

- Darema, Frederica (2001). “The SPMD Model: Past, Present and Future”. In: *Recent Advances in Parallel Virtual Machine and Message Passing Interface*. Ed. by Yiannis Cotronis and Jack Dongarra. Berlin, Heidelberg, pp. 1–1.
- Durmus, Alain, Eric Moulines, and Marcelo Pereyra (2018). “Efficient Bayesian computation by proximal Markov chain Monte Carlo: when Langevin meets Moreau”. In: *SIAM J. Imaging Sci.* 11.1, pp. 473–506.
- Rendell, Lewis J. et al. (2021). “Global Consensus Monte Carlo”. In: *J. Comput. and Graph. Stat.* 30.2, pp. 249–259.
- Salim, Adil and Peter Richtàrik (2020). “Primal Dual Interpretation of the Proximal Stochastic Gradient Langevin Algorithm”. In: *Adv. in Neural Information Processing Systems*. Vol. 33, pp. 3786–3796.

References III

- Vono, M., N. Dobigeon, and P. Chainais (Mar. 2019a). “Split-and-augmented Gibbs sampler - Application to large-scale inference problems”. In: *IEEE Trans. Signal Process.* 67.6, pp. 1648–1661.
- Vono, Maxime, Nicolas Dobigeon, and Pierre Chainais (2019b). “Asymptotically exact data augmentation: models, properties and algorithms”. In: arXiv preprint.
- (2021). “Asymptotically Exact Data Augmentation: Models, Properties, and Algorithms”. In: *J. Comput. and Graph. Stat.* 30.2, pp. 335–348.