

On Stochastic Variance Reduction for Penalised PET Reconstruction

Zeljko Kereta¹, Robert Twyman², Simon Arridge¹,
Kris Thielemans², Bangti Jin¹

¹Computer Science Department, University College London

²Institute of Nuclear Medicine, University College London

Advanced Image Reconstruction Methods

12 September 2022

Introduction

$$A\mathbf{f} + \mathbf{r} = E[\mathbf{g}]$$

- Penalised PET reconstruction is concerned with

$$\mathbf{f}_{\text{map}} = \underset{\mathbf{f}}{\operatorname{argmax}} f(\mathbf{f}) := L(\mathbf{f}) + R(\mathbf{f})g$$

log likelihood ← L(f)
← R(f)g
← penalty
← penalty strength

- Iterative optimisation is widely used in the reconstruction
- Using all the data can be costly and slow - **ordered subsets** (mini batching) dramatically improve convergence in early iterations

Motivation

- Standard subset methods experience loss of convergence towards the maximising solution
- This often leads to limit cycle behaviour and significant variability between successive image updates
- Convergence can be achieved by a relaxed step size sequence
- Alternatively, we can utilise variance reduction algorithms

Preconditioned Gradient Ascent

- Update equation via diagonally preconditioned gradient ascent methods

$$\mathbf{f}^{k+1} = \mathbf{P}_0 \mathbf{f}^k + \underbrace{\alpha_k}_{\text{stepsize}} \underbrace{D_{t_k}(\mathbf{f}^k)}_{\text{diagonal, non-negative preconditioner}} \underbrace{\beta_{t_k}(\mathbf{f}^k)}_{\text{gradient estimator}}$$

- \mathbf{P}_0 is a non-negativity projection; t_k is a subset index

Example: SGD/OSEM

- $r_{t_k}(\mathbf{f}) = \mathbf{A}_t^T \left(\frac{\mathbf{g}_m}{\mathbf{A}_t \mathbf{f} + \mathbf{r}_m} \right) \quad 1) \quad \overline{M}^r R(\mathbf{f}), \quad D_t(\mathbf{f}) = \text{diag} \left[\frac{\mathbf{f}}{\mathbf{A}_t^T \mathbf{1}} \right]$

Stochastic Variance Reduction

Common features

- Gradient estimator $\hat{g}_{t_k}(\mathbf{f}^k)$ is computed by using both a new subset gradient $r_{t_k}(\mathbf{f}^k)$ and a history of stored gradients
- Each update has approximately the same computational cost as one subset gradient update, but the memory is increased
- Use of a constant step-size

Stochastic Variance Reduction

- **SAGA** -replace stored subset gradient with the current one

$$\rho_{t_k}(\mathbf{f}^k) = M(r_{t_k}(\mathbf{f}^k) - \mathbf{g}_{t_k}) + \sum_{t=1}^{N_S} \mathbf{g}_t$$

update $\mathbf{g}_{t_k} = r_{t_k}(\mathbf{f}^k)$; and for $t \notin t_k$ keep \mathbf{g}_t

- **SVRG** -periodically recompute subset gradient history at current image

$$\rho_{t_k}(\mathbf{f}^k) = M(r_{t_k}(\mathbf{f}^k) - r_{t_k}(\mathbf{f}^{\text{anc}})) + \mathcal{G}$$

If $k \bmod N_S = 0$ set $\mathbf{f}^{\text{anc}} = \mathbf{f}^k$ and update $\mathcal{G} = r_{t_k}(\mathbf{f}^{\text{anc}})$

Convergence for SAGA and SVRG

Theorem

Let $d \in \mathbb{R}_{>0}^N$, denote by $L = \max_{t \in \{1, \dots, N_s\}} L_t$ where L_t is the Lipschitz constant of sub-objective gradients $\nabla_t(f)$ and by $d_{\max} = \max_i d_i$, and assume $\arg\max_{f \in \mathcal{C}} (f) \neq \emptyset$.

Taking $\frac{1}{3Ld_{\max}^{1/2}}$ and $D_t(f_{\text{saga}}^k) = \text{diag}(d)$ in the SAGA algorithm we have

$\|f_{\text{saga}}^k - f^*\| \leq \frac{1}{3Ld_{\max}^{1/2}}$ and $f_{\text{saga}}^k \rightarrow f^*$ almost surely.

Taking $\frac{1}{4Ld_{\max}^{1/2}(N_s+2)}$ and $D_t(f_{\text{svrg}}^k) = \text{diag}(d)$ in the SVRG algorithm we have

$\|f_{\text{svrg}}^k - f^*\| \leq \frac{1}{4Ld_{\max}^{1/2}(N_s+2)}$ almost surely and $E[\|f_{\text{svrg}}^k - f^*\|^2] = O(1/k)$.

NB: Subset gradients ∇_t are in general not Lipschitz. But, the assumptions are satisfied in physically realistic cases (everywhere non-zero backgrounds $r > 0$) or with the use of a modified log-likelihood

Experiments - XCAT Phantom

- GE Discovery 690 Scanner
- XCAT torso phantom with a 1cm diameter, 1cm long hot lesion (2.5:1 contrast)
- Data binned into 288 projection angles
- Image reconstruction algorithms based on STIR
- BSREM used for comparison

50 million event simulated data

Figure: Photon emission distribution

50 million event simulated data

Figure: Attenuation distribution cm-1

50 million event simulated data

Figure: Initial OSEM image

Experiments - Algorithm Setup

- Relative difference prior

$$R(x) = \prod_{i=1}^{N_v} \prod_{j \in N_i} p \frac{(f_i - f_j)^2}{(f_i + f_j) + 2j f_i - f_j}$$

- Preconditioner $D(f) = \text{diag}(\frac{f+}{A>1})$, with ϵ small
- Algorithms warm started with 1 epoch of OSEM (24 subsets)
- ROI drawn over the lung lesion, and we compute

$$\text{ROI Percentage Error} = 100\% \frac{k}{\hat{\Lambda}}$$

Experiments - Investigations

Investigate the impact of

- Choice of preconditioner: variable $D(f^k)$, frozen at initial point $D(f_0)$, after 5 epochs $D(f_{5M})$ or after 10 epochs $D(f_{10M})$
- Constant ($\eta_k = 1$), and decaying ($\eta_k = \frac{1}{k=M+1}$) stepsizes
- The number of subsets
- Subsets are pre-binned and then randomly sampled - other methodologies were tested but are not shown

Number of subsets

Figure: BSREM

Number of subsets

Figure: SAGA

Number of subsets

Figure: SVRG

Role of the preconditioner

Figure: SAGA

Role of the preconditioner

Figure: SVRG

Role of the preconditioner

Figure: ROI outside the torso

Variance Reduction in MAP-EM

- Traditionally, PET was optimised through an MLEM (maximum-likelihood expectation maximisation) algorithm

$$f^{k+1} = \underset{f \geq 0}{\operatorname{argmax}} E_{G|g;f}[\log p(G|f)]$$

- The issue with (explicit) maximisation with general priors (MAP-EM) is that the gradients are not spatially independent, and thus there is no closed-form maximiser

Online EM

- Write MAP-EM as

$$f^{k+1} = \underset{f \geq 0}{\operatorname{argmax}} \left[\log(f)^{\top} T(f^k) + \sum_{m=1}^M a_m^{\top} f R(f) \right];$$

depend on f

- Here

$$T(f^k) = \frac{1}{N_s} \sum_{t=1}^{N_s} t(f^k)$$

full conditional statistic

where

$$t(f) = N_s \operatorname{diag}(f) (rL_t(f) + A_t^{\top} \mathbf{1})$$

subset conditional statistic

Stochastic EM

Instead of $T(f^k)$ we compute

- SEM

$$\mathbf{p}^{k+1} = (1 - \alpha_k) \mathbf{p}^k + \alpha_k t_k(\mathbf{p}^k)$$

- SVREM

$$\mathbf{p}^{k+1} = (1 - \alpha_k) \mathbf{p}^k + \alpha_k t_k(\mathbf{p}^k) + \alpha_k t_k(\mathbf{p}^{\text{anc}}) + \mathbf{s}^{\text{anc}}$$

If $k \bmod N_s = 0$; set $\mathbf{f}^{\text{anc}} = \mathbf{f}^k$ and update $\mathbf{s}^{\text{anc}} = T(\mathbf{f}^{\text{anc}})$

- SAGAEM

$$\mathbf{p}^{k+1} = (1 - \alpha_k) \mathbf{p}^k + \alpha_k t_k(\mathbf{p}^k) + \frac{1}{N_s} \sum_{t=1}^{N_s} \mathbf{s}_t$$

Draw $t_k \sim \mathcal{U}[N_s]$; set $\mathbf{s}_{t_k} = t_k(\mathbf{p}^k)$; keep the rest intact

Stochastic EM

Instead of $T(f^k)$ we compute

■ SEM

$$\mathbf{p}^{k+1} = (1 - \alpha_k) \mathbf{p}^k + \alpha_k \mathbf{t}_k(\mathbf{p}^k)$$

■ SVREM

$$\mathbf{p}^{k+1} = (1 - \alpha_k) \mathbf{p}^k + \alpha_k \mathbf{t}_k(\mathbf{p}^k) + \alpha_k \mathbf{t}_k(\mathbf{p}^{\text{anc}}) + \mathbf{s}^{\text{anc}}$$

If $k \bmod N_s = 0$; set $\mathbf{f}^{\text{anc}} = \mathbf{f}^k$ and update $\mathbf{s}^{\text{anc}} = T(\mathbf{f}^{\text{anc}})$

■ SAGAEM

$$\mathbf{p}^{k+1} = (1 - \alpha_k) \mathbf{p}^k + \alpha_k \mathbf{t}_k(\mathbf{p}^k) + \frac{1}{N_s} \sum_{t=1}^{N_s} \mathbf{s}_t$$

Draw $t_k \sim [N_s]$; set $\mathbf{s}_{t_k} = \mathbf{t}_{t_k}(\mathbf{p}^k)$; keep the rest intact

Separable surrogates

- We consider (standard) priors of the form

$$R(f) = \frac{1}{2} \sum_{n=1}^N \sum_{j=1}^{2N_n} w_{nj} f_n f_j$$

and use parabolic surrogates

$$b^k(f_n; f_j) = (f_n^k f_j^k) f_n \frac{f_n^k + f_j^k}{2}^2 + f_j \frac{f_n^k + f_j^k}{2}^2 ;$$

where $(f) = \frac{(f)}{f}$

- The surrogate M-step for MAP-SEM/SVREM/SAGAEM is given by

$$f^{k+1} = \underset{f \geq 0}{\operatorname{argmax}} \log(f)^T p^{k+1} - \sum_{m=1}^M a_m^T f - \mathcal{R}(f; f^k)$$

where

$$\mathcal{R}(f; f^k) = \frac{1}{2} \sum_{n=1}^N \sum_{j \in 2N_n} w_{nj} \wedge^k f_n; f_j$$

- Explicit maximiser (root of the gradient is a quadratic polynomial with a single non-negative solution)

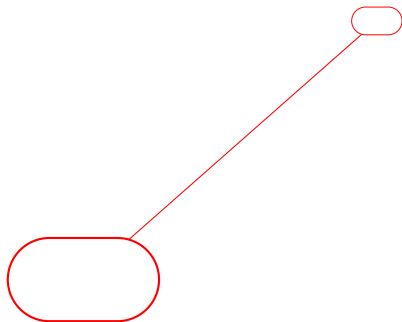
Admissible potentials

	$\phi(t)$	$\phi(x)$	$\psi(x)$
quadratic	$\frac{x^2}{2}$	x	1
log cosh	$\rho^2 \log \cosh(x)$	$\tanh(x)$	$\frac{\tanh(x)}{x}$
hyperbola	$\frac{1}{1+(x)^2}$	$\rho \frac{x}{1+(x)^2}$	$\rho \frac{1}{1+(x)^2}$

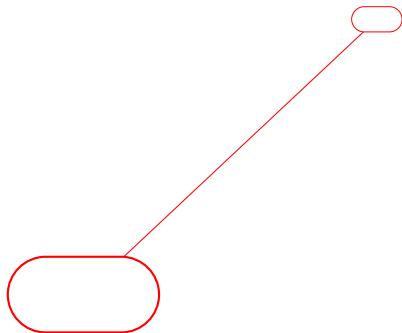
Experiments - XCAT Phantom

- XCAT torso phantom; 280 view scanner
- log coshprior with hand selected λ and penalty strength
- Initialised with 5 epochs of OSEM
- Sinogram data pre-binned as OS. A subset index is then sampled at random in each iteration

Objective Value - 40 Subsets



Objective Value - 40 Subsets



SVREM Reconstruction Progression

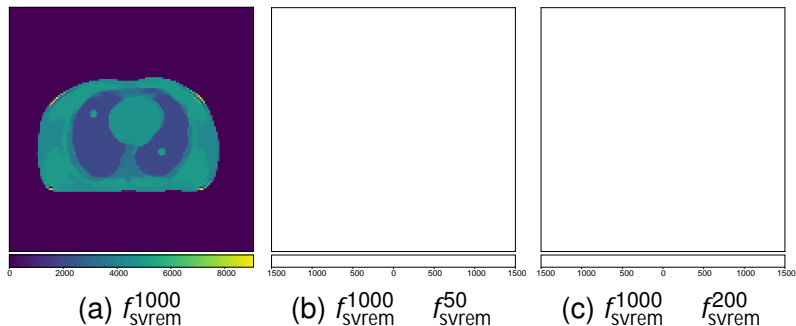


Figure: (a) SVREM reconstruction after 1000, and (b)-(c) pixel-wise differences of SVREM reconstructions after 200 and 50 epochs.