# Stochastic EM methods with Variance Reduction for Penalised PET Reconstructions

23 November 2021

# Introduction

$$\mathbf{Af} + \mathbf{r} = \mathbb{E}[\mathbf{g}]$$

- Iterative methods are widely used in PET reconstruction

- EM-ML[1] and its variants are particularly prevalent

$$\mathbf{f}^{k+1} = \underset{\mathbf{f} \geq \mathbf{0}}{\mathrm{argmax}} \, \mathbb{E}_{\mathbf{G}|\mathbf{g}, \mathbf{f}^k}[\log p(\mathbf{G}|\mathbf{f})]$$

image ←     measured data ←     complete data →

- Explicit solution in each step

- Ordered subset (OS) methods improve the convergence in early iterations

---

[1] Shepp and Vardi, '82

# Introduction

- Iterative methods are widely used in PET reconstruction

- EM and its variants are particularly prevalent

$$\mathbf{f}^{k+1} = \underset{\mathbf{f} \geq \mathbf{0}}{\arg\max} \, \mathbb{E}_{\mathbf{G}_t | \mathbf{g}_t, \mathbf{f}^k} [\log p(\mathbf{G}_t | \mathbf{f})]$$

→ subset index

- Explicit solution in each step

- Ordered subset (OS) methods improve the convergence in early iterations

UCL CENTRE FOR
INVERSE
PROBLEMS

Two common issues with OSEM methods:

- Loss of convergence towards the maximising solution
    - Instead we enter a limit cycle behaviour

- Problems if there is a penalty (MAP-EM)

$$\mathbf{f}_{\mathrm{map}} = \underset{\mathbf{f} \geq \mathbf{0}}{\mathrm{argmax}}\{\Phi(\mathbf{f}) := \mathcal{L}(\mathbf{f}) - \beta \, \mathcal{R}(\mathbf{f})\}$$

log likelihood ←     penalty →

penalty strength →

- Maximisation is no longer analytical and thus further approximations are needed

Two common issues with OSEM methods:

- Loss of convergence towards the maximising solution

    - Instead we enter a limit cycle behaviour

- Problems if there is a penalty (MAP-EM)

$$\mathbf{f}_{\text{map}}^{k+1} = \underset{\mathbf{f} \geq \mathbf{0}}{\operatorname{argmax}} \{ \underbrace{\mathbb{E}_{\mathbf{G}|\mathbf{g},\mathbf{f}^k}[\log p(\mathbf{G}|\mathbf{f})]}_{\text{conditional expectation}} - \underbrace{\beta}_{} \underbrace{\mathcal{R}(\mathbf{f})}_{} \}$$

conditional expectation

penalty

penalty strength

- Maximisation is no longer analytical and thus further approximations are needed

- Alternative: optimise using gradient ascent based methods (as discussed by Robbie)

- Instead we consider stochastic EM algorithms for MAP-EM which

    - Uses OS and **exponentially moving average** of the expected statistic
    - Employs **separable parabolic surrogates**[1] for the prior

---

[1] de Pierro '95, de Pierro and Yamagishi '01, Erdogan and Fessler '98, etc.

# Online EM [2]

- Write MAP-EM as

depend on $\mathbf{f}$

$$\mathbf{f}^{k+1} = \underset{\mathbf{f} \geq \mathbf{0}}{\operatorname{argmax}} \left\{ \boxed{\log(\mathbf{f})^\top} s(\mathbf{f}^k) - \boxed{\sum_{m=1}^{M} a_m^\top \mathbf{f} - \beta \mathcal{R}(\mathbf{f})} \right\},$$

- Here

full conditional statistic

$$\boxed{s(\mathbf{f}^k)} = \mathbb{E}_{\mathbf{G}|\mathbf{g},\mathbf{f}^k} \left[ \log\left( \sum_{m=1}^{M} g_{mn} \right)_{n=1}^{N} \right] = \frac{1}{N_s} \sum_{t=1}^{N_s} \tau_t(\mathbf{f}^k)$$

where

$$\boxed{\tau_t(\mathbf{f})} = N_s f \odot \left( \nabla \mathcal{L}_t(\mathbf{f}) + A_t^\top \mathbf{1} \right)$$

subset conditional statistic

# Stochastic EM

Instead of $s(\mathbf{f}^k)$ we compute [2, 1, 3]

- **SEM**

$$\widehat{s}^{k+1} = (1 - \alpha_k)\widehat{s}^k + \alpha_k\, \tau_{t_k}(\widehat{\mathbf{f}}_{\text{sem}}^k)$$

- **SVREM**

$$\widehat{s}^{k+1} = (1 - \alpha)\widehat{s}^k + \alpha\left(\tau_{t_k}(\widehat{\mathbf{f}}_{\text{svrem}}^k) - \tau_{t_k}(\widehat{\mathbf{f}}^{\text{anc}}) + s^{\text{anc}}\right)$$

If $k \bmod \eta N_s = 0$, set $\mathbf{f}^{\text{anc}} = \mathbf{f}_{\text{svrem}}^k$ and update $s^{\text{anc}} = s(\mathbf{f}^{\text{anc}})$

- **SAGAEM**

$$\widehat{s}^{k+1} = (1 - \alpha)\widehat{s}^k + \alpha\left(\tau_{t_k}(\widehat{\mathbf{f}}_{\text{sagaem}}^k) - \mathfrak{s}_{t_k} + \frac{1}{N_s}\sum_{t=1}^{N_s}\mathfrak{s}_t\right)$$

Draw $\tilde{t}_k \in [N_s]$, set $\mathfrak{s}_{\tilde{t}_k} = \tau_{t_k}(\widehat{\mathbf{f}}_{\text{sagaem}}^k)$, keep the rest intact

# Stochastic EM

Instead of $s(\mathbf{f}^k)$ we compute [2, 1, 3]

- **SEM**

$$\widehat{s}^{k+1} = (1 - \alpha_k)\widehat{s}^k + \alpha_k \, \tau_{t_k}(\widehat{\mathbf{f}}^k_{\text{sem}})$$

- **SVREM**

$$\widehat{s}^{k+1} = (1 - \alpha)\widehat{s}^k + \alpha \left(\tau_{t_k}(\widehat{\mathbf{f}}^k_{\text{svrem}}) - \tau_{t_k}(\widehat{\mathbf{f}}^{\text{anc}}) + s^{\text{anc}}\right)$$

If $k \bmod \eta N_s = 0$, set $\mathbf{f}^{\text{anc}} = \mathbf{f}^k_{\text{svrem}}$ and update $s^{\text{anc}} = s(\mathbf{f}^{\text{anc}})$

- **SAGAEM**

$$\widehat{s}^{k+1} = (1 - \alpha)\widehat{s}^k + \alpha \left(\tau_{t_k}(\widehat{\mathbf{f}}^k_{\text{sagaem}}) - \mathfrak{s}_{t_k} + \frac{1}{N_s}\sum_{t=1}^{N_s}\mathfrak{s}_t\right)$$

Draw $\tilde{t}_k \in [N_s]$, set $\mathfrak{s}_{\tilde{t}_k} = \tau_{\tilde{t}_k}(\widehat{\mathbf{f}}^k_{\text{sagaem}})$, keep the rest intact

# Separable surrogates

- We consider (standard) priors of the form

$$\mathcal{R}(\mathbf{f}) = \frac{1}{2} \sum_{n=1}^{N} \sum_{j \in \mathcal{N}_n} w_{nj} \, \boxed{\rho(f_n - f_j)}$$

  → smooth, non decreasing function of $|f_n - f_j|$

- The issue with (explicit) maximisation with general priors is that the gradients are not spatially independent

- Instead of $\rho$ use a parabolic surrogate [4]

$$\widehat{\rho}^k(f_n; f_j) = \gamma_\rho(f_n^k - f_j^k)\Big( \big(f_n - \tfrac{f_n^k + f_j^k}{2}\big)^2 + \big(f_j - \tfrac{f_n^k + f_j^k}{2}\big)^2 \Big),$$

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

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

$$\mathbf{f}^{k+1} = \underset{\mathbf{f} \geq \mathbf{0}}{\operatorname{argmax}} \left\{ \log(\mathbf{f})^\top \widehat{s}^{k+1} - \sum_{m=1}^M a_m^\top \mathbf{f} - \beta \widehat{\mathcal{R}}(\mathbf{f}; \mathbf{f}^k) \right\}$$

where

$$\widehat{\mathcal{R}}(\mathbf{f}; \mathbf{f}^k) = \frac{1}{2} \sum_{n=1}^N \sum_{j \in \mathcal{N}_n} w_{nj} \, \hat{\rho}^k (f_n; f_j)$$

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

# Explicit maximiser

Let $d_{nj} := w_{nj}\gamma_\rho(f_n^k - f_j^k)$ and

$$a_n = \widehat{s}_n^k, \quad b_n = \beta \sum_{j \in \mathcal{N}_n} d_{nj},$$

$$c_n = \beta f_n^k \sum_{j \in \mathcal{N}_n} d_{nj} + \beta \sum_{j \in \mathcal{N}_n} d_{nj} f_j^k - \sum_{m=1}^{M} a_{mn}$$

Then

$$f_n^{k+1} = \frac{c_n + \sqrt{c_n^2 + 8a_n b_n}}{4b_n}$$

UCL CENTRE FOR
**INVERSE
PROBLEMS**

# Admissible potentials

| | $\rho(t)$ | $\rho'(\mathsf{x})$ | $\gamma_\rho(\mathsf{x})$ |
|---|---|---|---|
| **quadratic** | $\frac{\mathsf{x}^2}{2}$ | $\mathsf{x}$ | $1$ |
| **log cosh** | $\delta^2 \log \cosh(\mathsf{x}/\delta)$ | $\delta \tanh(\mathsf{x}/\delta)$ | $\delta \frac{\tanh(\mathsf{x}/\delta)}{\mathsf{x}}$ |
| **hyperbola** | $\delta\left(\sqrt{1 + (\mathsf{x}/\delta)^2} - 1\right)$ | $\frac{\mathsf{x}}{\sqrt{1+(\mathsf{x}/\delta)^2}}$ | $\frac{1}{\sqrt{1+(\mathsf{x}/\delta)^2}}$ |

UCL CENTRE FOR
**INVERSE
PROBLEMS**

# Convergence for SAGA and SVRG

As a reminder, variance reduced gradient ascent methods obey

$$\mathbf{f}^{k+1} = \mathbf{P}_{\geq 0}\Big(\mathbf{f}^k + \alpha \mathbf{D}_k(\mathbf{f}^k)\tilde{\nabla}_k\Big)$$

---

**Theorem**

*Let $d \in \mathbb{R}^N_{>0}$, denote by $L = \max_{t \in N_s} L_t$ where $L_t$ is the Lipschitz constant of sub-objective gradients $\widetilde{\Phi}_t(\mathbf{f})$ and by $d_{\max} = \|d\|_\infty$, and assume $\mathrm{argmax}_{\mathbf{f} \geq \mathbf{0}} \Phi(\mathbf{f}) \neq \emptyset$. Taking $\alpha \leq \frac{1}{3Ld_{\max}^{1/2}}$ and $\mathbf{D}_t(\mathbf{f}^k_{\mathrm{saga}}) = \mathrm{diag}(d)$ in the SAGA algorithm we have $\widetilde{\Phi}(\mathbf{f}^k_{\mathrm{saga}}) \to \Phi(\mathbf{f}^\star)$ and $\mathbf{f}^k_{\mathrm{saga}} \to \mathbf{f}^\star$ almost surely. Taking $\alpha \leq \frac{1}{4Ld_{\max}^{1/2}(\eta N_s + 2)}$ and $\mathbf{D}_t(\mathbf{f}^k_{\mathrm{svrg}}) = \mathrm{diag}(d)$ in the SVRG algorithm we have $\mathbf{f}^k_{\mathrm{svrg}} \to \mathbf{f}^\star$ almost surely and $\mathbb{E}[\Phi(\mathbf{f}^\star) - \widetilde{\Phi}(\mathbf{f}^{k\eta N_s}_{\mathrm{svrg}})] = \mathcal{O}(1/k)$.*

---

UCL CENTRE FOR
**INVERSE
PROBLEMS**

# Convergence

- Subset gradients $\Phi_t$ are in general not Lipschitz

- But, the assumptions are satisfied in physically realistic cases (everywhere non-zero backgrounds $\mathbf{r} > \mathbf{0}$) or with the used of a modified log-likelihood[2]

- Under current theory SVREM and SAGAEM will also converge under certain (but somewhat stronger) Lipschitz assumptions
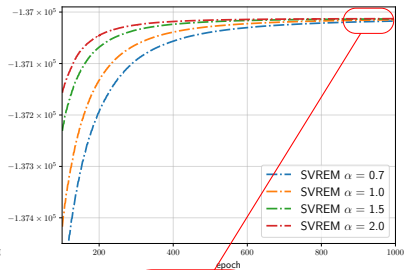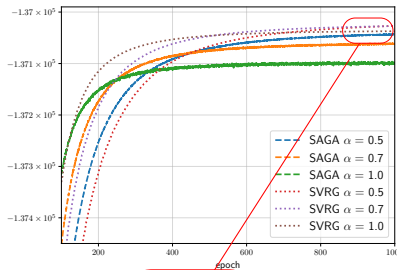
---

[2]Ahn and Fessler, '03

UCL CENTRE FOR
INVERSE
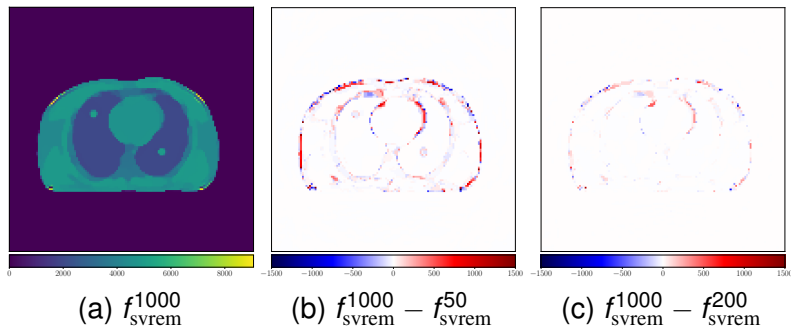PROBLEMS

# Experiments - XCAT Phantom

- XCAT torso phantom; 280 view scanner

- $\log \cosh$ prior with hand selected $\delta$ and penalty strength $\beta$

- 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

# SVREM Reconstruction Progression



(a) $f_{\text{svrem}}^{1000}$     (b) $f_{\text{svrem}}^{1000} - f_{\text{svrem}}^{50}$     (c) $f_{\text{svrem}}^{1000} - f_{\text{svrem}}^{200}$

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

UCL CENTRE FOR
**INVERSE
PROBLEMS**

# Quick and easy way heuristics to accelearate the convergence

- Using *SVRG without the outer loop* (and adjusting $\eta$)

- Nonlinear acceleration through extrapolation
    - Improves performance drastically on simple data
    - Inconsistent on more realistic data

- Nesterov, etc.

UCL CENTRE FOR
**INVƎRSE
PROBLEMS**

J. Chen, J. Zhu, Y. W. Teh, and T. Zhang,
Stochastic Expectation Maximization with Variance Reduction,
*NeurIPS*, (2018).

O. Cappé and E. Moulines,
Online Expectation-Maximization Algorithm for Latent Data Models,
*Journal of the Royal Statistical Society: Series B*, **71(3)** (2009).

B. Karimi, H.-T. Wai, E. Moulines, and M. Lavielle,
On the Global Convergence of (Fast) Incremental Expectation Maximization
Methods,
*NeurIPS*, (2019).

J.-H. Chang, J. M. M. Anderson, and J. R. Votaw,
Regularized Image Reconstruction Algorithms for Positron Emission Tomography,

*IEEE Trans. Med. Imag.*, **23(9)** (2004).