



Computational Strategies for Bayesian Inversion with Conditionally Gaussian Sparsity Priors

SIAM CSE23

3rd March 2023



Jonathan Lindblom[†], Jan Glaubitz, and Anne Gelb

Outline

1. Introduction
2. Hierarchical model
3. Data fusion
4. Despeckling
5. Data augmentation

Problem setup

We would like to reconstruct an unknown signal $\mathbf{x} \in \mathbb{R}^n$ given an indirect noisy observation

$$\mathbf{y} = \mathbf{F}\mathbf{x} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\text{noise}}).$$

Here $\mathbf{y}, \boldsymbol{\epsilon} \in \mathbb{R}^m$ with linear forward operator $\mathbf{F} \in \mathbb{R}^{m \times n}$ and SPD noise covariance $\boldsymbol{\Sigma}_{\text{noise}} \in \mathbb{R}^{m \times m}$.

Problem setup

We would like to reconstruct an unknown signal $\mathbf{x} \in \mathbb{R}^n$ given an indirect noisy observation

$$\mathbf{y} = \mathbf{F}\mathbf{x} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\text{noise}}).$$

Here $\mathbf{y}, \boldsymbol{\epsilon} \in \mathbb{R}^m$ with linear forward operator $\mathbf{F} \in \mathbb{R}^{m \times n}$ and SPD noise covariance $\boldsymbol{\Sigma}_{\text{noise}} \in \mathbb{R}^{m \times m}$.

Assumptions:

1. \mathbf{x} is (relatively) sparse under some linear transformation $\mathbf{R} \in \mathbb{R}^{k \times n}$
2. We permit $\ker(\mathbf{R}) \neq \{\mathbf{0}\}$, but require $\ker(\mathbf{F}) \cap \ker(\mathbf{R}) = \{\mathbf{0}\}$
3. We do not know $\boldsymbol{\Sigma}_{\text{noise}}$ a priori but believe it has a simple parametric form

Conditionally Gaussian sparsity priors

We could opt to work with a sparsity prior directly.

Alternatively, many sparsity priors admit a representation as a conditionally Gaussian hierarchical prior, in which case of the form

$$\beta_j \stackrel{iid}{\sim} \text{Hyper-prior}$$

$$\boldsymbol{x} | \boldsymbol{\beta} \sim \mathcal{N} \left(\mathbf{0}, \left(\boldsymbol{R}^T \boldsymbol{D}_{\boldsymbol{\beta}}^{-1} \boldsymbol{R} \right)^{-1} \right)$$

$$\boldsymbol{y} | \boldsymbol{x} \sim \mathcal{N} (\boldsymbol{F}\boldsymbol{x}, \boldsymbol{\Sigma}_{\text{noise}})$$

Conditionally Gaussian sparsity priors

ℓ_1 -regularization:

$$\boldsymbol{x}^* = \arg \min_{\boldsymbol{x}} \|\boldsymbol{F}\boldsymbol{x} - \boldsymbol{y}\|_{\Sigma_{\text{noise}}}^2 + \lambda \|\boldsymbol{R}\boldsymbol{x}\|_1$$

via conditionally Gaussian representation [1] :

$$(\boldsymbol{x}^*, \boldsymbol{\beta}^*) \stackrel{\eta \rightarrow 0}{=} \arg \min_{\boldsymbol{x}, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{F}\boldsymbol{x} - \boldsymbol{y}\|_{\Sigma_{\text{noise}}}^2 + \frac{1}{2} \|\boldsymbol{R}\boldsymbol{x}\|_{D_{\boldsymbol{\beta}}^{-1}}^2 + \frac{1}{2} \sum_{j=1}^k \log \beta_j - \log \pi(\boldsymbol{\beta})$$

¹Calvetti, Somersalo and Strang, "Hierachical Bayesian models and sparsity: ℓ_2 -magic", 2019.

Negative log posterior

$$\pi_{\text{GG}}(\theta | r, s, \vartheta) \propto \theta^{rs-1} \exp \left\{ - \left(\frac{\theta}{\vartheta} \right)^r \right\} \mathbf{1}_{>0}(\theta)$$

$$E(\boldsymbol{x}, \boldsymbol{\beta}) = \frac{1}{2} \|\boldsymbol{F}\boldsymbol{x} - \boldsymbol{y}\|_{\boldsymbol{\Sigma}_{\text{noise}}^{-1}}^2 + \frac{1}{2} \|\boldsymbol{R}\boldsymbol{x}\|_{\boldsymbol{D}_{\boldsymbol{\beta}}^{-1}}^2 - \eta \sum_{j=1}^n \log \beta_j + \sum_{j=1}^n \left(\frac{\beta_j}{\vartheta} \right)^r$$

Theorem [2]

Let $\eta := rs - \frac{3}{2}$ and $\ker(\boldsymbol{R}) = \{\mathbf{0}\}$. Then $E(\boldsymbol{x}, \boldsymbol{\beta})$ is globally convex if $r \geq 1$ and $\eta > 0$, and is locally convex at $(\boldsymbol{x}, \boldsymbol{\beta})$ if either (i) $0 < r < 1$ and $\eta > 0$ or (ii) $r < 0$ and $s > 0$, and

$$\beta_j < \vartheta \left(\eta / (r|r-1|) \right)^{1/r}, \quad j = 1, \dots, k.$$

²Calvetti et al., "Sparse reconstructions from few noisy data: analysis of hierarchical Bayesian models with generalized gamma hyperpriors", 2020.

MAP estimation via block coordinate descent (IAS)

IAS Algorithm

Choose initial $\boldsymbol{x}_0, \boldsymbol{\beta}_0$.

for $k = 1, \dots, n_{\text{maxits}}$ do

$$\boldsymbol{\beta}_k \leftarrow \arg \max_{\boldsymbol{\beta}} \pi(\boldsymbol{x}_{k-1}, \boldsymbol{\beta} | \boldsymbol{y})$$

$$\boldsymbol{x}_k \leftarrow \arg \max_{\boldsymbol{x}} \pi(\boldsymbol{x}, \boldsymbol{\beta}_k | \boldsymbol{y})$$

end for

$$\arg \max_{\boldsymbol{x}} \pi(\boldsymbol{x}, \boldsymbol{\beta} | \boldsymbol{y}) = \left(\boldsymbol{F}^T \boldsymbol{\Sigma}_{\text{noise}}^{-1} \boldsymbol{F} + \boldsymbol{R}^T \boldsymbol{D}_{\boldsymbol{\beta}}^{-1} \boldsymbol{R} \right)^{-1} \boldsymbol{F}^T \boldsymbol{\Sigma}_{\text{noise}}^{-1} \boldsymbol{y}$$

Hierarchical model

Consider $\Sigma_{\text{noise}} = \alpha \mathbf{I}_m$. Our first modification is to consider the model

$$\begin{aligned}\alpha &\sim \text{GG}(r_1, s_1, \vartheta_1), \\ \beta_j &\stackrel{iid}{\sim} \text{GG}(r_2, s_2, \vartheta_2), \quad j = 1, \dots, k, \\ \mathbf{x} | \boldsymbol{\beta} &\sim \mathcal{N} \left(\mathbf{0}, \left(\mathbf{R}^T \mathbf{D}_{\boldsymbol{\beta}}^{-1} \mathbf{R} \right)^{-1} \right), \\ \mathbf{y} | \mathbf{x}, \alpha &\sim \mathcal{N} (\mathbf{F}\mathbf{x}, \alpha \mathbf{I}_m).\end{aligned}$$

Negative log posterior

The negative log posterior is given by $E(\mathbf{x}, \alpha, \boldsymbol{\beta}) = E_1(\mathbf{x}, \alpha) + E_2(\mathbf{x}, \boldsymbol{\beta})$, where

$$E_1(\mathbf{x}, \alpha) = \frac{1}{2\alpha} \|\mathbf{F}\mathbf{x} - \mathbf{y}\|_2^2 - \eta_1 \log \frac{\alpha}{\vartheta_1} + \left(\frac{\alpha}{\vartheta_1} \right)^{r_1}$$

$$E_2(\mathbf{x}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{R}\mathbf{x}\|_{\mathbf{D}_{\boldsymbol{\beta}}^{-1}}^2 - \eta_2 \sum_{j=1}^k \log \frac{\beta_j}{\vartheta_2} + \sum_{j=1}^k \left(\frac{\beta_j}{\vartheta_2} \right)^{r_2}$$

Here $\eta_2 := r_2 s_2 - \frac{3}{2}$, and $\eta_1 := r_1 s_1 - \frac{m+2}{2}$.

Convexity in the variance parameterization

Convexity result for $E_1(\mathbf{x}, \alpha)$

Same result in terms of (r, s, η) , η depends on m (size of measurement vector).

Convexity result for $E_2(\mathbf{x}, \beta)$

Same result, but we have shown explicitly that it still holds for $\ker(\mathbf{R}) \neq \{\mathbf{0}\}$.

Main observation:

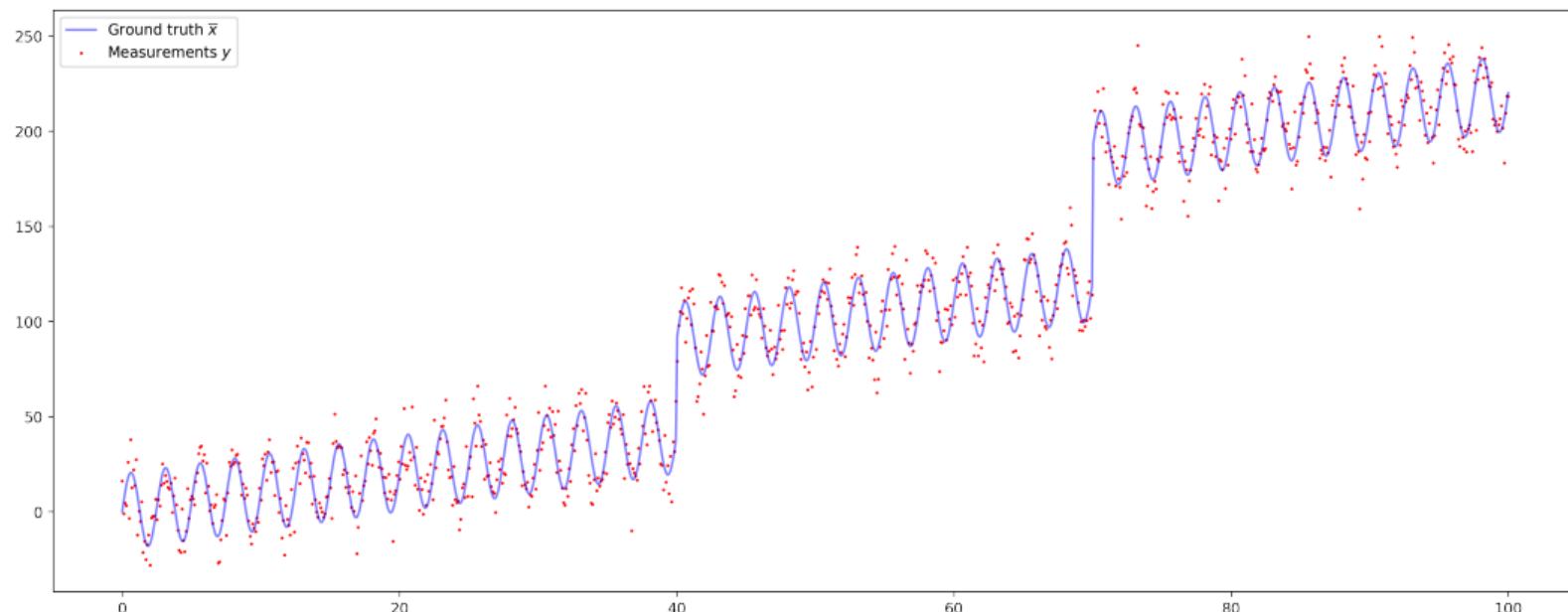
$$\nabla_{\mathbf{x}} \nabla_{\beta} E_2(\mathbf{x}, \beta) = -\mathbf{R}^T \mathbf{D}_{\beta}^{-1} \text{diag} \left(\mathbf{D}_{\beta}^{-1} \mathbf{R} \mathbf{x} \right)$$

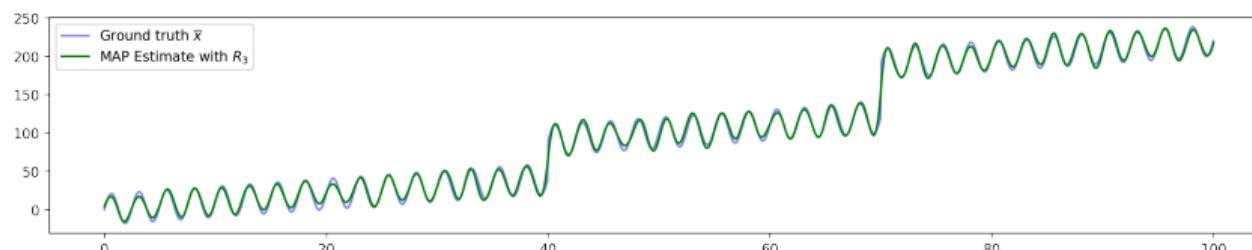
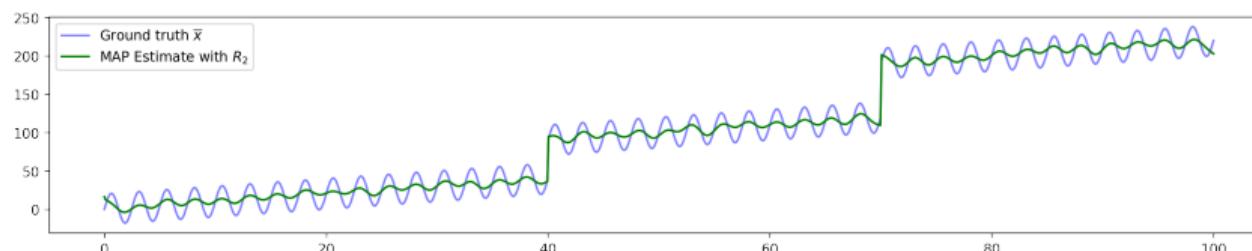
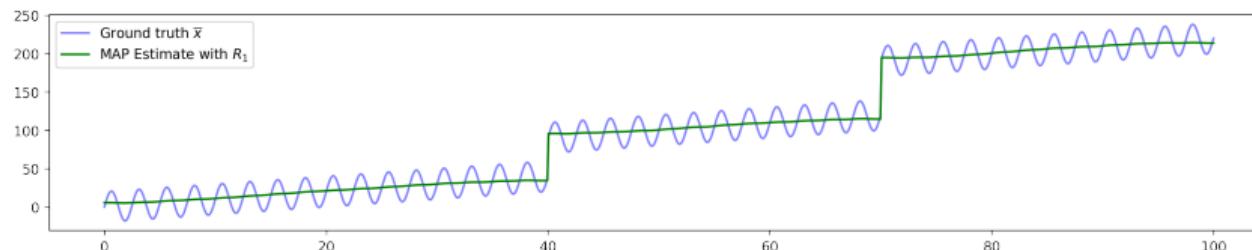
Uninformative hyper-priors

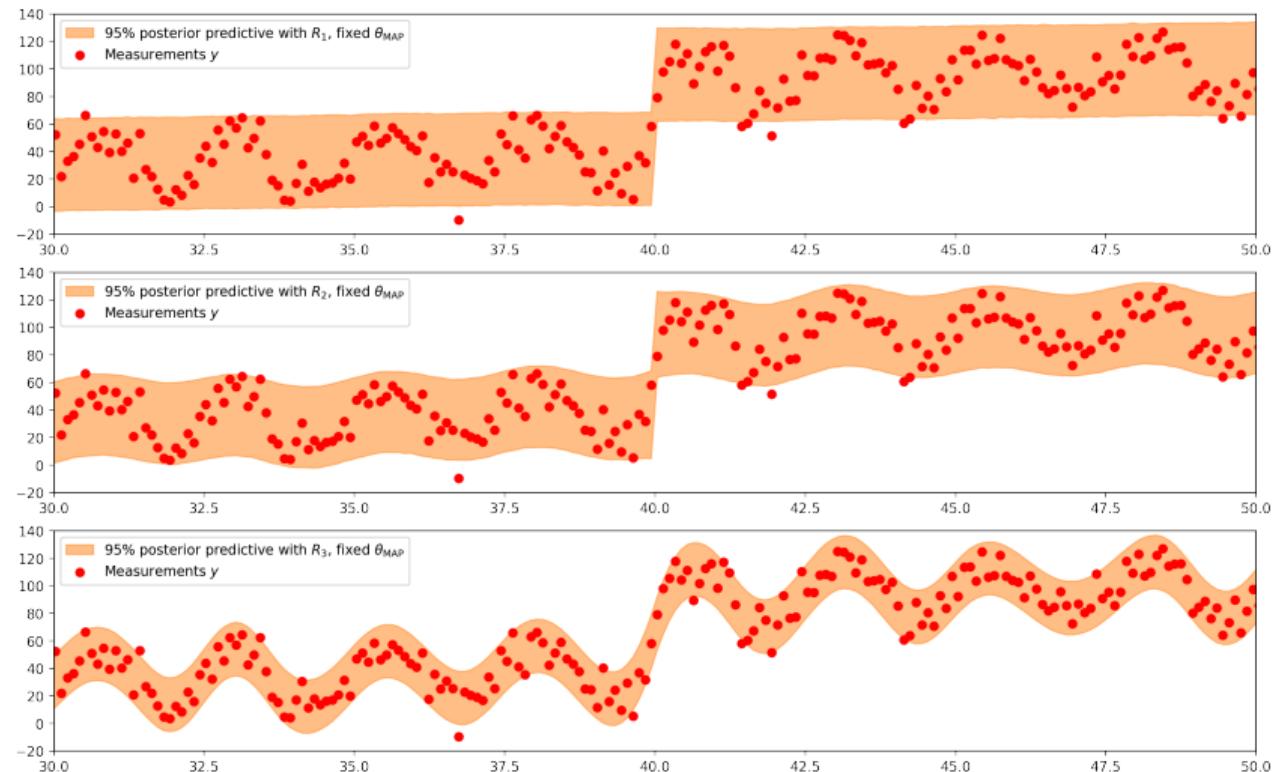
- The choice $r_1 = -1, s_1 = 1, \vartheta_1 = 10^{-4}$ corresponds to an “uninformative” inverse Gamma variance hyper-prior
- In literature and in practice, accurately recovers true noise variance
- Stationary points satisfy discrepancy principle w.r.t. inferred noise variance
- If a stationary point (\mathbf{x}^*, α^*) satisfies the discrepancy principle $\|\mathbf{F}\mathbf{x}^* - \mathbf{y}\|_2^2 \approx m\alpha_{\text{true}}$, then $\alpha^* \approx \alpha_{\text{true}}$
- However, our sufficient condition for local convexity in this case is that

$$\alpha \leq \frac{4 \times 10^{-4}}{m + 4}$$

- Condition is not useful?







Precision parameterization

What changes if α and β denote inverse variances (precisions)?

Lemma

$\theta \sim \text{GG}(r, s, \vartheta) \Leftrightarrow \frac{1}{\theta} \sim \text{GG}(-r, s, \vartheta^{-1})$, so we can easily convert between the two parameterizations.

But in general the convexity and the MAP estimate will depend on the parameterization.

Negative log posterior

The negative log posterior is given by $E(\mathbf{x}, \alpha, \boldsymbol{\beta}) = E_3(\mathbf{x}, \alpha) + E_4(\mathbf{x}, \boldsymbol{\beta})$, where

$$E_3(\mathbf{x}, \alpha) = \frac{\alpha}{2} \|\mathbf{F}\mathbf{x} - \mathbf{y}\|_2^2 - \eta_3 \log \frac{\alpha}{\vartheta_3} + \left(\frac{\alpha}{\vartheta_3} \right)^{r_3},$$

$$E_4(\mathbf{x}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{R}\mathbf{x}\|_{\mathbf{D}_{\boldsymbol{\beta}}}^2 - \eta_4 \sum_{j=1}^k \log \frac{\beta_j}{\vartheta_4} + \sum_{j=1}^k \left(\frac{\beta_j}{\vartheta_4} \right)^{r_4}.$$

Here $\eta_3 = r_3 s_3 + \frac{m-2}{2}$, $\eta_4 = r_4 s_4 - \frac{1}{2}$.

Convexity in the precision parameterization

Convexity result for $E_3(\mathbf{x}, \alpha)$

$$\eta_3\alpha^{-1} + \frac{r_3(r_3 - 1)}{\vartheta_3^{r_3}}\alpha^{r_3-1} - \|\mathbf{F}\mathbf{x} - \mathbf{y}\|_2^2 \geq 0$$

Convexity result for $E_4(\mathbf{x}, \beta)$

$$\eta_4\beta_j^{-1} + \frac{r_4(r_4 - 1)}{\vartheta_4^{r_4}}\beta_j^{r_4-1} - [\mathbf{R}\mathbf{x}]_j^2 \geq 0, \quad j = 1, \dots, k$$

Uninformative hyper-priors

- The choice $r_3 = 1, s_3 = 1, \vartheta_3 = 10^4$ corresponds to an “uninformative” Gamma precision hyper-prior
- From our convexity condition, local convexity only guaranteed in this case if

$$\frac{m}{2} \geq \alpha \|\mathbf{F}\mathbf{x} - \mathbf{y}\|_2^2$$

- If \mathbf{x}^* approximately satisfies the discrepancy principle $\|\mathbf{F}\mathbf{x}^* - \mathbf{y}\|_2^2 \approx \frac{m}{\alpha_{\text{true}}}$, then this becomes

$$\frac{1}{\alpha} \geq \frac{2}{\alpha_{\text{true}}}$$

IAS in the precision parameterization

- The \boldsymbol{x} -update is still given by the solution to a linear system
- What about the α and β updates?

Lemma

If $r < 0$ and $\eta < 0$, or if $r > 0$ and $\eta > 0$, for $z \geq 0$ and $\xi > 0$ the unique solution to

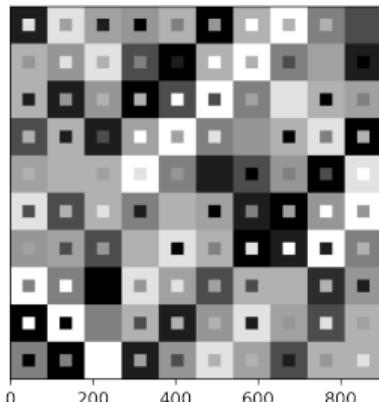
$$\frac{1}{2}z^2 - \eta\xi^{-1} + r\xi^{r-1} = 0$$

is given by the solution to the initial value problem

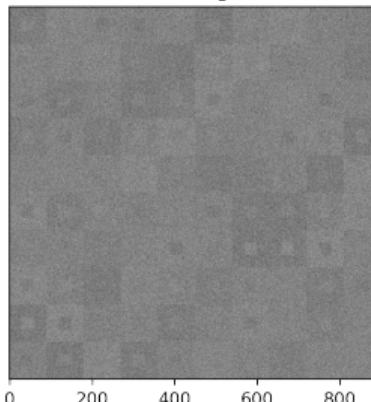
$$\psi'(t) = -\frac{2t\psi(t)}{2r^2\psi(t)^{r-1} + t^2}, \quad \psi(0) = (\eta/r)^{1/r}$$

evaluated at $t = |z|$.

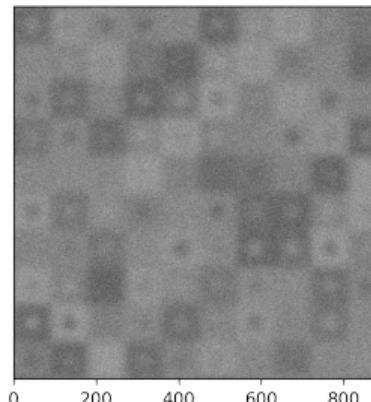
Ground truth



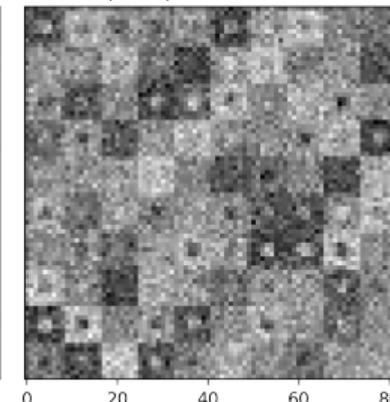
Point-wise (high noise)



Blurred (medium noise)



Upsampled (low noise)



Conditioned on any hyper-parameters θ , the posterior density looks like

$$\pi(\mathbf{x} \mid \mathbf{y}^{1:p}, \theta) \propto \pi_1(\mathbf{x}) \times \cdots \times \pi_K(\mathbf{x}), \quad \pi_j(\mathbf{x}) \propto \exp \left\{ -\frac{1}{2} \|\mathbf{F}_j \mathbf{x} - \mathbf{y}_j\|_{Q_j}^2 \right\}.$$

Conditioned on any hyper-parameters θ , the posterior density looks like

$$\pi(\mathbf{x} | \mathbf{y}^{1:p}, \theta) \propto \pi_1(\mathbf{x}) \times \cdots \times \pi_K(\mathbf{x}), \quad \pi_j(\mathbf{x}) \propto \exp \left\{ -\frac{1}{2} \|\mathbf{F}_j \mathbf{x} - \mathbf{y}_j\|_{\mathbf{Q}_j}^2 \right\}.$$

Then $\pi(\mathbf{x} | \mathbf{y}^{1:p}, \theta)$ is the Gaussian $\mathcal{N}(\boldsymbol{\mu}, \mathbf{Q}^{-1})$ with

$$\begin{aligned} \mathbf{Q} &= \sum_{j=1}^K \mathbf{F}_j^T \mathbf{Q}_j \mathbf{F}_j, \\ \boldsymbol{\mu} &= \mathbf{Q}^{-1} \left(\sum_{j=1}^K \mathbf{F}_j^T \mathbf{Q}_j \mathbf{y}_j \right). \end{aligned}$$

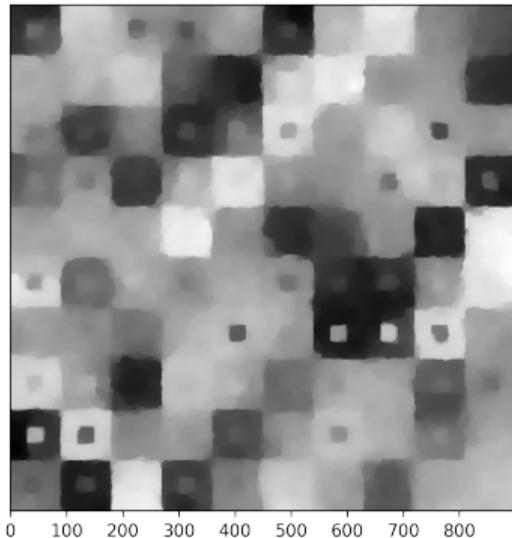
IAS algorithm for data fusion

Input Data $\{\mathbf{y}^{(j)}\}_{j=1}^p$, forward operators $\mathbf{F}^{(j)}$, sparsifying transformation \mathbf{R} , hyper-prior parameters, initialization $(\mathbf{x}_0, \{\alpha_0\}_{j=1}^p, \boldsymbol{\beta}_0)$, number of iterations N .

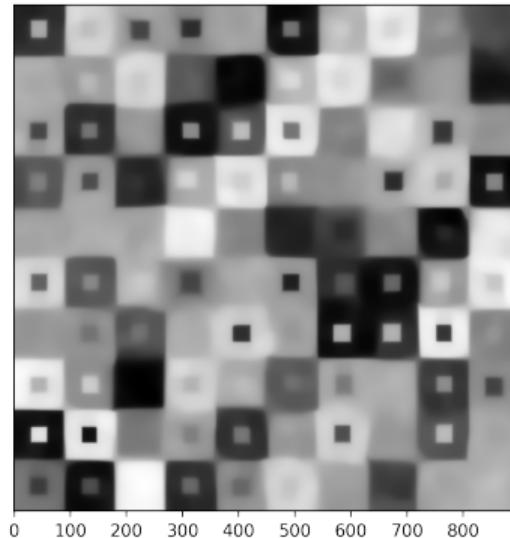
Output Approximation to the MAP estimate.

- 1 For $k = 1, \dots, N$ do
 - i. Set $\alpha_k^{(j)} = \arg \max_{\alpha^{(j)}} \pi(\mathbf{x}_{k-1}, \alpha_{k-1}^{(1)}, \dots, \alpha_k^{(j)}, \dots, \alpha_{k-1}^{(p)}, \boldsymbol{\beta}_{k-1} | \mathbf{y})$ for each $j = 1, \dots, p$.
 - ii. Set $\boldsymbol{\beta}_k = \arg \max_{\boldsymbol{\beta}} \pi(\mathbf{x}_{k-1}, \alpha_k^{(1)}, \dots, \alpha_k^{(p)}, \boldsymbol{\beta} | \mathbf{y})$.
 - iii. Set $\mathbf{x}_k = \arg \max_{\mathbf{x}} \pi(\mathbf{x}, \alpha_k^{(1)}, \dots, \alpha_k^{(p)}, \boldsymbol{\beta}_k | \mathbf{y})$.
 - 2 Return $(\mathbf{x}_N, \{\alpha_N^{(j)}\}_{j=1}^p, \boldsymbol{\beta}_N)$.
-

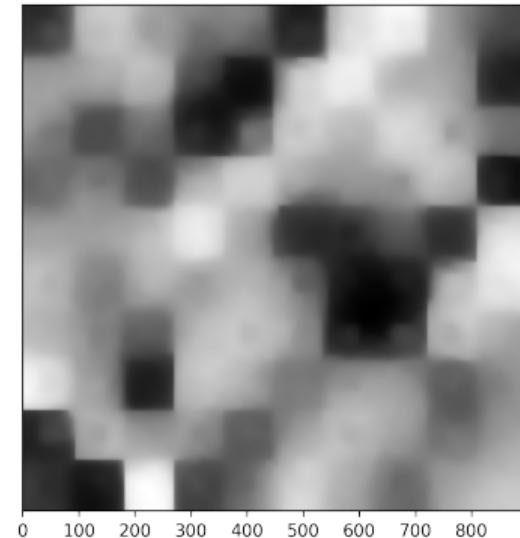
From point-wise observation



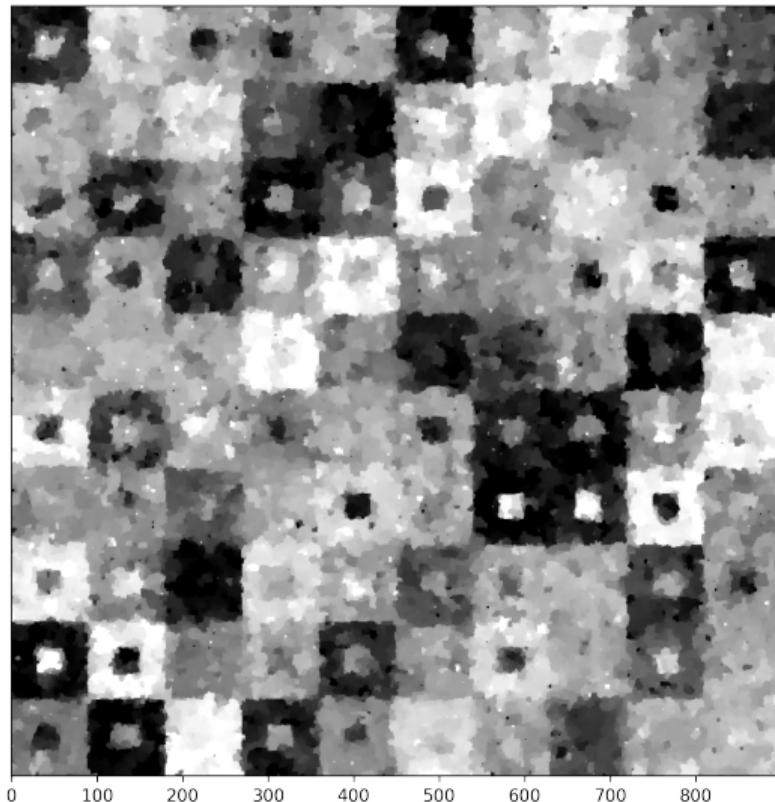
From blurred observation



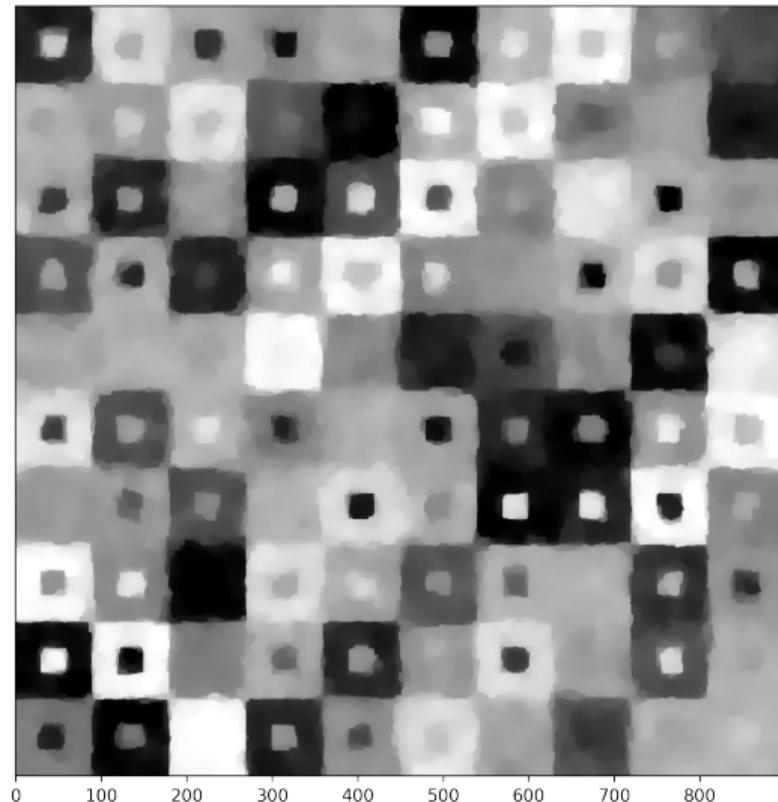
From up-sampled observation



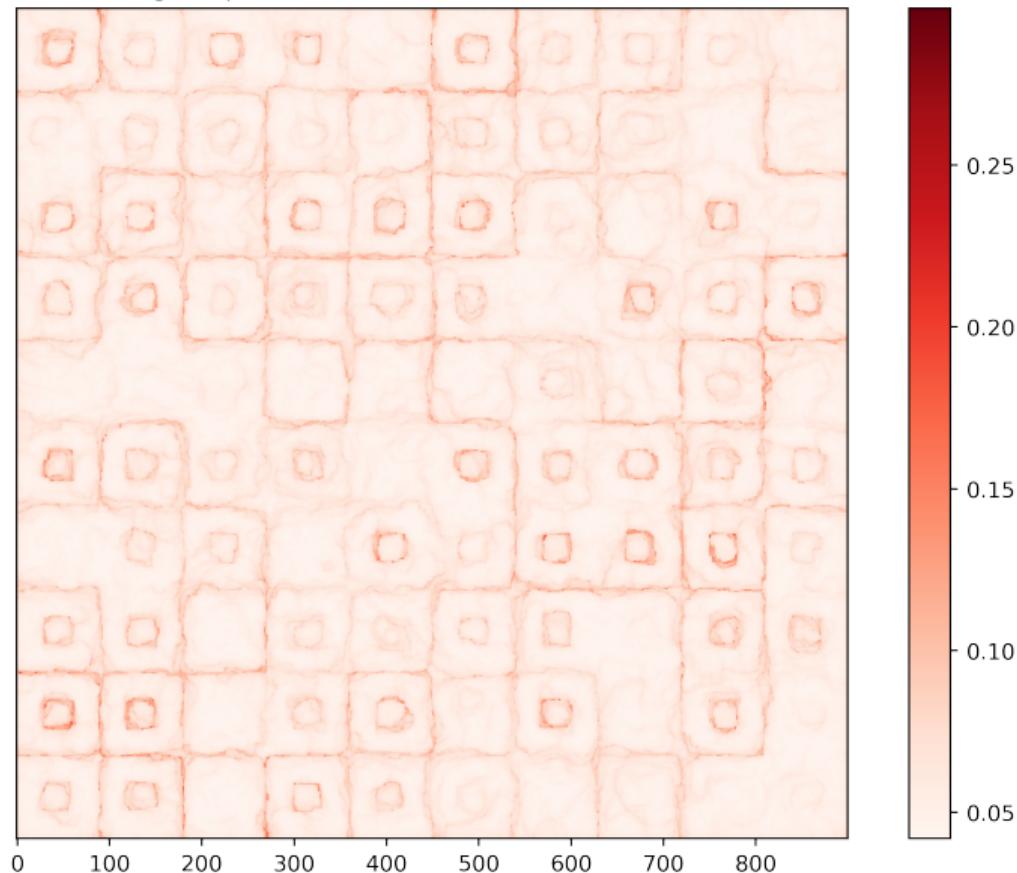
Stacked reconstruction

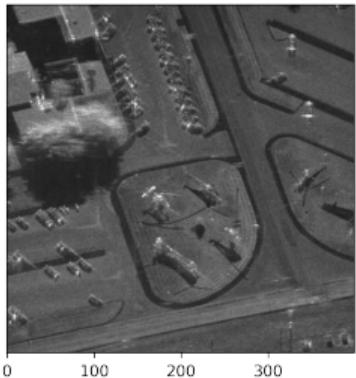
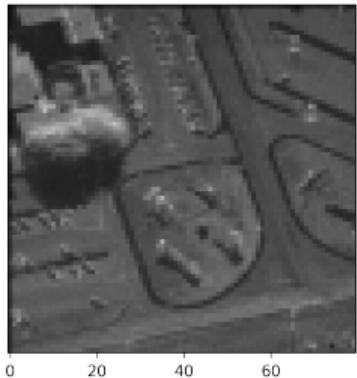
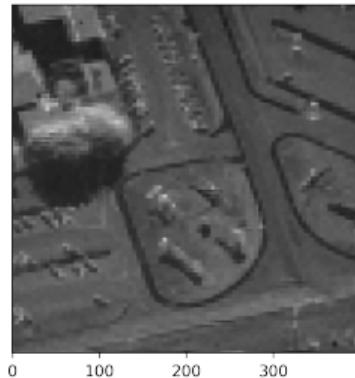


Reconstruction with three variances



Reconstruction	MSE	SSIM
From point-wise	0.0223	0.7518
From blurred	0.0113	0.7763
From up-sampled	0.0478	0.7101
Stacked	0.0157	0.6339
With three learned variances	0.0073*	0.8162*

Marginal pixel-wise standard deviation, fixed θ_{MAP} 

$\bar{x}^{(1)}$  $y^{(1)}$ IAS reconstruction, initialized at $\bar{x}^{(1)}$  $\bar{x}^{(2)}$  $y^{(2)}$ IAS reconstruction, initialized at $\bar{x}^{(2)}$ 

Marginal penalty

In the variance parameterization, it is known that we can pick (r, s, ϑ) such that the local minimizers agree with those corresponding to ℓ_p -norm ($0 < p < 2$) and Student- t ($\nu = 2$) regularization.

In the precision parameterization, we have shown we can do the same to achieve Cauchy regularization.

Despeckling

- Multiplicative gamma-distributed noise
- Logarithmic transformation yields a model with log-concave Fisher-Tippett likelihood with

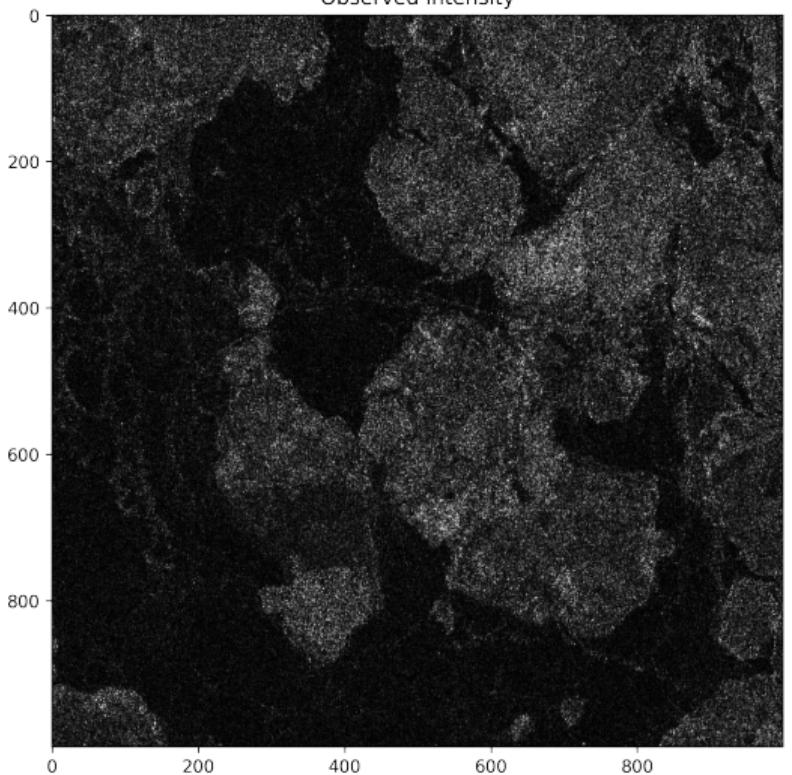
$$-\log \pi(\mathbf{y} | \mathbf{x}) = L \left(\sum_{i=1}^n \exp\{\mathbf{y}_i - \mathbf{x}_i\} + \mathbf{x}_i - \mathbf{y}_i \right)$$

- Coupled with hierarchical sparsity prior, only modification to IAS is solution to nonlinear equation via Newton's method or other

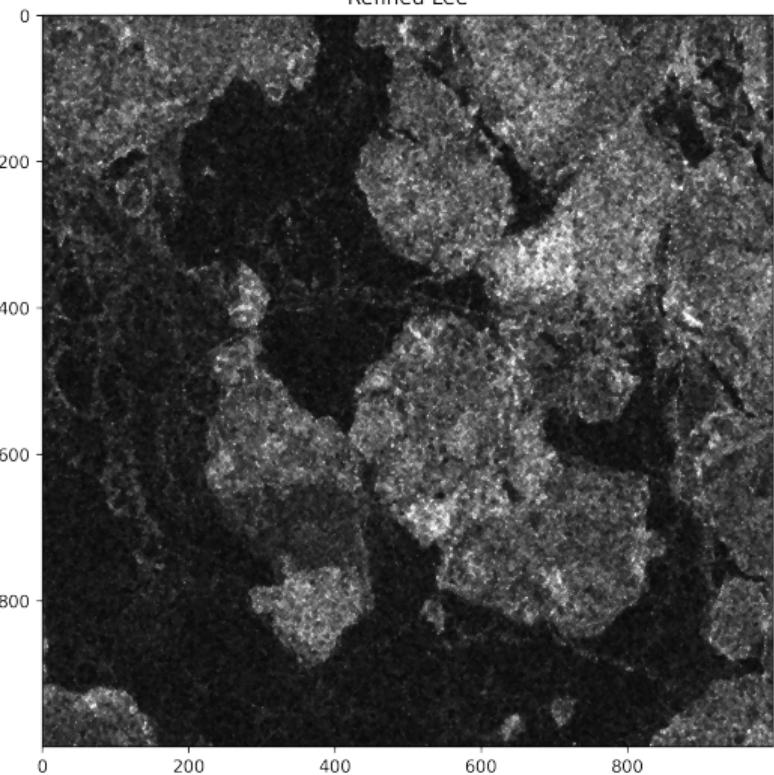


Sentinel-1, courtesy of ESA

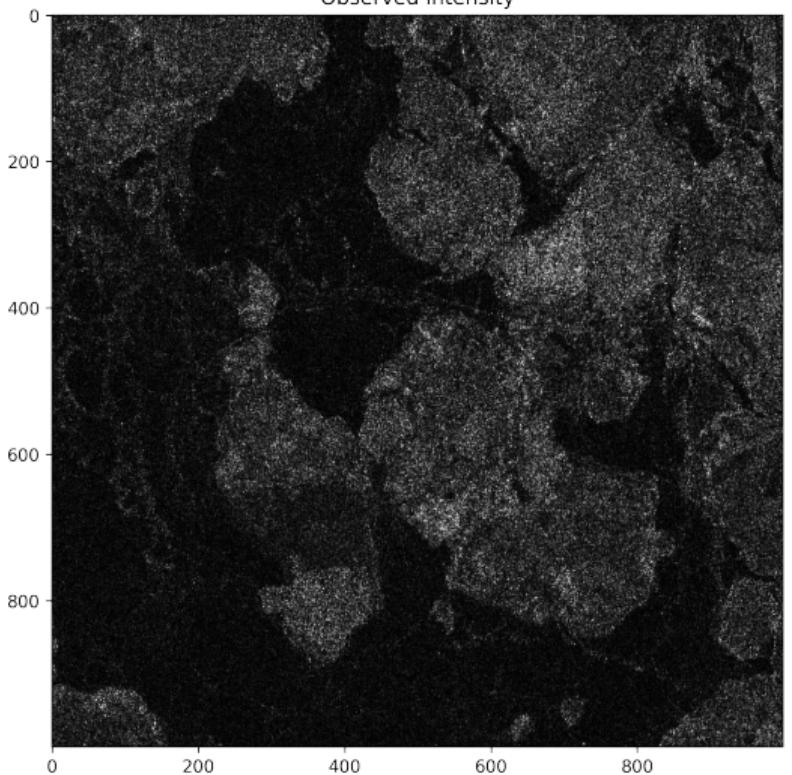
Observed intensity



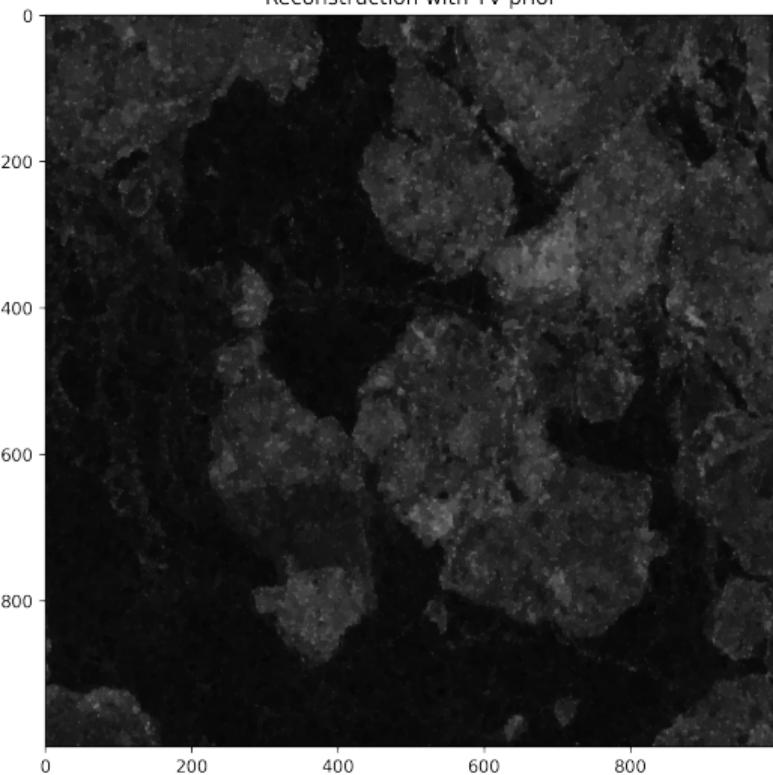
Refined Lee



Observed intensity



Reconstruction with TV prior



Data augmentation

- When using a Gaussian prior $\pi(\mathbf{x}) \propto \exp\{-\frac{\lambda}{2}\|\mathbf{R}\mathbf{x}\|_2^2\}$, we are able to make use of a common diagonalization of $\mathbf{F}^T \mathbf{F}$ and $\mathbf{R}^T \mathbf{R}$ to directly compute an \mathbf{x} -update, e.g., if both block-circulant with circulant blocks (BCCB) then

$$\mathbf{F}^T \mathbf{F} = \mathcal{F}^H \boldsymbol{\Lambda}_F \mathcal{F}, \quad \mathbf{R}^T \mathbf{R} = \mathcal{F}^H \boldsymbol{\Lambda}_R \mathcal{F}$$

³Marnissi et al., "A Data Augmentation Approach for Sampling Gaussian Models in High Dimension", 2019.

Data augmentation

- When using a Gaussian prior $\pi(\mathbf{x}) \propto \exp\{-\frac{\lambda}{2}\|\mathbf{R}\mathbf{x}\|_2^2\}$, we are able to make use of a common diagonalization of $\mathbf{F}^T \mathbf{F}$ and $\mathbf{R}^T \mathbf{R}$ to directly compute an \mathbf{x} -update, e.g., if both block-circulant with circulant blocks (BCCB) then

$$\mathbf{F}^T \mathbf{F} = \mathcal{F}^H \boldsymbol{\Lambda}_F \mathcal{F}, \quad \mathbf{R}^T \mathbf{R} = \mathcal{F}^H \boldsymbol{\Lambda}_R \mathcal{F}$$

- When applying IAS with a sparse prior, \mathbf{D}_{β}^{-1} "gets-in-the-way" since

$$\mathbf{R}^T \mathbf{D}_{\beta}^{-1} \mathbf{R} = ?$$

³Marnissi et al., "A Data Augmentation Approach for Sampling Gaussian Models in High Dimension", 2019.

Data augmentation

- When using a Gaussian prior $\pi(\mathbf{x}) \propto \exp\{-\frac{\lambda}{2}\|\mathbf{R}\mathbf{x}\|_2^2\}$, we are able to make use of a common diagonalization of $\mathbf{F}^T \mathbf{F}$ and $\mathbf{R}^T \mathbf{R}$ to directly compute an \mathbf{x} -update, e.g., if both block-circulant with circulant blocks (BCCB) then

$$\mathbf{F}^T \mathbf{F} = \mathcal{F}^H \boldsymbol{\Lambda}_F \mathcal{F}, \quad \mathbf{R}^T \mathbf{R} = \mathcal{F}^H \boldsymbol{\Lambda}_R \mathcal{F}$$

- When applying IAS with a sparse prior, \mathbf{D}_β^{-1} "gets-in-the-way" since

$$\mathbf{R}^T \mathbf{D}_\beta^{-1} \mathbf{R} = ?$$

- We have made use of recent data augmentation techniques³ to get around this

³Marnissi et al., "A Data Augmentation Approach for Sampling Gaussian Models in High Dimension", 2019.

$$\boldsymbol{v} | \boldsymbol{x}, \boldsymbol{\beta} \sim \mathcal{N} \left(\boldsymbol{\mu}_v, \boldsymbol{Q}_v^{-1} \right), \quad \boldsymbol{Q}_v^{-1} = \frac{1}{\lambda(\boldsymbol{\beta})} \boldsymbol{I} - \boldsymbol{D}_{\boldsymbol{\beta}}^{-1}, \quad \boldsymbol{\mu}_v = \boldsymbol{Q}^{-1} \boldsymbol{R} \boldsymbol{x},$$

$$\Rightarrow \boldsymbol{x} | \boldsymbol{v}, \alpha, \boldsymbol{\beta}, \boldsymbol{y} \sim \mathcal{N} \left(\boldsymbol{\mu}, \boldsymbol{Q}^{-1} \right), \quad \boldsymbol{Q} = \frac{1}{\alpha} \boldsymbol{F}^T \boldsymbol{F} + \frac{1}{\lambda} \boldsymbol{R}^T \boldsymbol{R}, \quad \boldsymbol{\mu} = \boldsymbol{Q}^{-1} \left(\frac{1}{\alpha} \boldsymbol{F}^T \boldsymbol{y} + \boldsymbol{R}^T \boldsymbol{v} \right).$$

PCG-IAS Algorithm

Choose initial $\boldsymbol{x}_0, \alpha_0, \boldsymbol{\beta}_0, \boldsymbol{v}_0$.

for $k = 1, \dots, n_{\text{maxits}}$ do

$$\alpha_k \leftarrow \arg \max_{\alpha} \pi(\alpha | \boldsymbol{x}_{k-1}, \boldsymbol{\beta}_{k-1}, \boldsymbol{y}).$$

$$\boldsymbol{\beta}_k \leftarrow \arg \max_{\boldsymbol{\beta}} \pi(\boldsymbol{\beta} | \alpha_k, \boldsymbol{x}_{k-1}, \boldsymbol{y}).$$

$$\boldsymbol{v}_k \leftarrow \arg \max_{\boldsymbol{v}} \pi(\boldsymbol{v} | \alpha_k, \boldsymbol{\beta}_k, \boldsymbol{x}_{k-1}, \boldsymbol{y})$$

$$\boldsymbol{x}_k \leftarrow \arg \max_{\boldsymbol{x}} \pi(\boldsymbol{x} | \alpha_k, \boldsymbol{\beta}_k, \boldsymbol{v}_k, \boldsymbol{y}).$$

end for
