

Hyper-parameters Estimation in Bayesian Models

Guohui Song

Old Dominion University

Joint work with Feng Yu (UT El Paso) and Lixin Shen (Syracuse University)

Linear Inverse Problems

- Consider a linear inverse problem: $\mathbf{y} = \mathbf{F}\mathbf{x} + \boldsymbol{\epsilon}$
- It is usually ill-posed.
 - We need some extra/prior information on \mathbf{x} .

Linear Inverse Problems

Consider a linear inverse problem: $\mathbf{y} = \mathbf{F}\mathbf{x} + \boldsymbol{\epsilon}$

↑ observations ↗ unknown parameters

- Ⓐ It is usually ill-posed.
- Ⓐ We need some extra/prior information on \mathbf{x} .

Regularization approach: impose a regularization/penalty term $R(\mathbf{x})$

$$\arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \lambda R(\mathbf{x})$$

- Ridge regression: $R(\mathbf{x}) = \|\mathbf{x}\|^2$
- Lasso: $R(\mathbf{x}) = \|\mathbf{x}\|_1$
- Group Lasso: $R(\mathbf{x}) = \sum_{g=1}^G \|\mathbf{x}_g\|_2$

Bayesian Models

- Bayesian approach: introduce a prior distribution on \boldsymbol{x}
 - likelihood: $p(\boldsymbol{y}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{y}|\mathbf{F}\boldsymbol{x}, \sigma^2\mathbf{I})$
 - prior: $p(\boldsymbol{x})$
 - compute the posterior distribution $p(\boldsymbol{x}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x})}{p(\boldsymbol{y})}$

Bayesian Models

Bayesian approach: introduce a prior distribution on \boldsymbol{x}

⊖ likelihood: $p(\boldsymbol{y}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{y}|\mathbf{F}\boldsymbol{x}, \sigma^2\mathbf{I})$

⊖ prior: $p(\boldsymbol{x})$

⊖ compute the posterior distribution $p(\boldsymbol{x}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x})}{p(\boldsymbol{y})}$

MAP estimate

$$\arg \max_{\boldsymbol{x}} p(\boldsymbol{x}|\boldsymbol{y}) = \arg \max_{\boldsymbol{x}} p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x}) = \arg \min_{\boldsymbol{x}} \frac{1}{2\sigma^2} \|\boldsymbol{y} - \mathbf{F}\boldsymbol{x}\|^2 - \log p(\boldsymbol{x})$$

○ Equivalent to the regularization approaches for appropriate choices of the prior.

Bayesian Models

Bayesian approach: introduce a prior distribution on \boldsymbol{x}

⊖ likelihood: $p(\boldsymbol{y}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{y}|\mathbf{F}\boldsymbol{x}, \sigma^2\mathbf{I})$

⊖ prior: $p(\boldsymbol{x})$

⊖ compute the posterior distribution $p(\boldsymbol{x}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x})}{p(\boldsymbol{y})}$

MAP estimate

$$\arg \max_{\boldsymbol{x}} p(\boldsymbol{x}|\boldsymbol{y}) = \arg \max_{\boldsymbol{x}} p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x}) = \arg \min_{\boldsymbol{x}} \frac{1}{2\sigma^2} \|\boldsymbol{y} - \mathbf{F}\boldsymbol{x}\|^2 - \log p(\boldsymbol{x})$$

⊖ Equivalent to the regularization approaches for appropriate choices of the prior.

Uncertainty quantification

⊖ The posterior distribution $p(\boldsymbol{x}|\boldsymbol{y})$ provides a full distribution of \boldsymbol{x} .

Bayesian Models

Bayesian approach: introduce a prior distribution on \boldsymbol{x}

⌚ likelihood: $p(\boldsymbol{y}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{y}|\mathbf{F}\boldsymbol{x}, \sigma^2\mathbf{I})$

⌚ prior: $p(\boldsymbol{x})$

⌚ compute the posterior distribution $p(\boldsymbol{x}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x})}{p(\boldsymbol{y})}$

MAP estimate

$$\arg \max_{\boldsymbol{x}} p(\boldsymbol{x}|\boldsymbol{y}) = \arg \max_{\boldsymbol{x}} p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x}) = \arg \min_{\boldsymbol{x}} \frac{1}{2\sigma^2} \|\boldsymbol{y} - \mathbf{F}\boldsymbol{x}\|^2 - \log p(\boldsymbol{x})$$

⌚ Equivalent to the regularization approaches for appropriate choices of the prior.

Uncertainty quantification

⌚ The posterior distribution $p(\boldsymbol{x}|\boldsymbol{y})$ provides a full distribution of \boldsymbol{x} .

Bayesian Models

Bayesian approach: introduce a prior distribution on \boldsymbol{x}

- ⊖ likelihood: $p(\boldsymbol{y}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{y}|\mathbf{F}\boldsymbol{x}, \sigma^2\mathbf{I})$
- ⊖ prior: $p(\boldsymbol{x})$
- ⊖ compute the posterior distribution $p(\boldsymbol{x}|\boldsymbol{y}) = \frac{p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x})}{p(\boldsymbol{y})}$

Challenges

- The computation of $p(\boldsymbol{y}) = \int p(\boldsymbol{y}|\boldsymbol{x})p(\boldsymbol{x})d\boldsymbol{x}$ is usually intractable.
- The choice of the prior distribution is crucial.

Hierarchical Bayesian Learning

- Consider a hierarchical prior (scale mixtures of Gaussians)

- $p(\mathbf{x}|\boldsymbol{\gamma}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \text{diag}(\gamma_1, \dots, \gamma_n))$

hyper-parameters

- impose a hyper-prior on $\boldsymbol{\gamma}$

Hierarchical Bayesian Learning

- ✓ Consider a hierarchical prior (scale mixtures of Gaussians)

- ⌚ $p(\mathbf{x}|\boldsymbol{\gamma}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \text{diag}(\gamma_1, \dots, \gamma_n))$

hyper-parameters

- ⌚ impose a hyper-prior on $\boldsymbol{\gamma}$

- Benefits

- Scale mixture of Gaussians contains all symmetric stable distributions, Laplace distributions, logistic distributions, and exponential power distributions.
 - The conditional posterior $p(\mathbf{x}|\mathbf{y}, \boldsymbol{\gamma})$ is analytically tractable.
 - Individualized hyper-parameters encode sparsity.

Hierarchical Bayesian Learning

Consider a hierarchical prior (scale mixtures of Gaussians)

Ⓐ $p(\mathbf{x}|\boldsymbol{\gamma}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \text{diag}(\gamma_1, \dots, \gamma_n))$

hyper-parameters

Ⓐ impose a hyper-prior on $\boldsymbol{\gamma}$

Benefits

Ⓐ Scale mixture of Gaussians contains all symmetric stable distributions, Laplace distributions, logistic distributions, and exponential power distributions.

Ⓐ The conditional posterior $p(\mathbf{x}|\mathbf{y}, \boldsymbol{\gamma})$ is analytically tractable.

Ⓐ Individualized hyper-parameters encode sparsity.

Challenge: how to estimate the hyper-parameters $\boldsymbol{\gamma}$?

○ The dimension of hyper-parameters $\boldsymbol{\gamma}$ scales with \mathbf{x} .

○ It is not easy to choose the hyper-prior.

Hyper-parameters Estimation

- Learning from data (Evidence approach)

$$\arg \max_{\gamma} p(\mathbf{y}|\gamma) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\gamma)d\mathbf{x}$$

Hyper-parameters Estimation

- Learning from data (Evidence approach)

$$\arg \max_{\gamma} p(\mathbf{y}|\gamma) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\gamma)d\mathbf{x}$$

- Equivalent to

$$\arg \min_{\gamma} \mathbf{y}^T (\mathbf{S}(\gamma))^{-1} \mathbf{y} + \log \det(\mathbf{S}(\gamma)), \quad \mathbf{S}(\gamma) = \sigma^2 \mathbf{I} + \mathbf{F} \text{diag}(\gamma) \mathbf{F}^T$$

- It is a high-dimensional non-convex optimization problem
- Existing methods: EM [Dempster et.al. 1977], MacKay [MacKay, 1992], Convex Bound-ing [Wipf and Nagarajan, 2009], etc.

Hyper-parameters Estimation

- Learning from data (Evidence approach)

$$\arg \max_{\gamma} p(\mathbf{y}|\gamma) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\gamma)d\mathbf{x}$$

- Equivalent to

$$\arg \min_{\gamma} \mathbf{y}^T (\mathbf{S}(\gamma))^{-1} \mathbf{y} + \log \det(\mathbf{S}(\gamma)), \quad \mathbf{S}(\gamma) = \sigma^2 \mathbf{I} + \mathbf{F} \text{diag}(\gamma) \mathbf{F}^T$$

- It is a high-dimensional non-convex optimization problem
- Existing methods: EM [Dempster et.al. 1977], MacKay [MacKay, 1992], Convex Bound-ing [Wipf and Nagarajan, 2009], etc.
- Goal: develop more efficient algorithms

Alternating Minimization

- The objective function:

$$\min_{\gamma} \mathbf{y}^T (\mathbf{S}(\gamma))^{-1} \mathbf{y} + \log \det(\mathbf{S}(\gamma)), \quad \mathbf{S}(\gamma) = \sigma^2 \mathbf{I} + \mathbf{F} \text{diag}(\gamma) \mathbf{F}^T$$

Alternating Minimization

- ✓ The objective function:

$$\min_{\gamma} \underline{\mathbf{y}^T (\mathbf{S}(\gamma))^{-1} \mathbf{y}} + \log \det(\mathbf{S}(\gamma)), \quad \mathbf{S}(\gamma) = \sigma^2 \mathbf{I} + \mathbf{F} \text{diag}(\gamma) \mathbf{F}^T$$

- Rewrite the first term: $\underline{\min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^T \text{diag}(\gamma^{-1}) \mathbf{x}}$

Alternating Minimization

- ✓ The objective function:

$$\min_{\gamma} \underline{\mathbf{y}^T (\mathbf{S}(\gamma))^{-1} \mathbf{y}} + \log \det(\mathbf{S}(\gamma)), \quad \mathbf{S}(\gamma) = \sigma^2 \mathbf{I} + \mathbf{F} \text{diag}(\gamma) \mathbf{F}^T$$

- ✓ Rewrite the first term: $\underline{\min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^T \text{diag}(\gamma^{-1}) \mathbf{x}}$

- Reformulate the objective function:

$$\min_{\gamma} \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^T \text{diag}(\gamma^{-1}) \mathbf{x} + \log \det(\mathbf{S}(\gamma))$$

Alternating Minimization

- Reformulate the objective function:

$$\min_{\gamma} \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^\top \text{diag}(\boldsymbol{\gamma}^{-1})\mathbf{x} + \log \det(\mathbf{S}(\boldsymbol{\gamma}))$$

- Alternating Minimization:

Alternating Minimization

- ✓ Reformulate the objective function:

$$\min_{\gamma} \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^\top \text{diag}(\boldsymbol{\gamma}^{-1}) \mathbf{x} + \log \det(\mathbf{S}(\boldsymbol{\gamma}))$$

- ✓ Alternating Minimization:

- \mathbf{x} -update: $\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^\top \text{diag}((\boldsymbol{\gamma}^{(k)})^{-1}) \mathbf{x}$

Alternating Minimization

- ✓ Reformulate the objective function:

$$\min_{\gamma} \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathsf{F}\mathbf{x}\|^2 + \mathbf{x}^\top \text{diag}(\boldsymbol{\gamma}^{-1}) \mathbf{x} + \log \det(\mathsf{S}(\boldsymbol{\gamma}))$$

- ✓ Alternating Minimization:

- \mathbf{x} -update: $\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathsf{F}\mathbf{x}\|^2 + \mathbf{x}^\top \text{diag}((\boldsymbol{\gamma}^{(k)})^{-1}) \mathbf{x}$

- $\boldsymbol{\gamma}$ -update: $\boldsymbol{\gamma}^{(k+1)} = \arg \min_{\boldsymbol{\gamma}} (\mathbf{x}^{(k+1)})^\top \text{diag}(\boldsymbol{\gamma}^{-1}) \mathbf{x}^{(k+1)} + \log \det(\mathsf{S}(\boldsymbol{\gamma}))$

Alternating Minimization

- ✓ Reformulate the objective function:

$$\min_{\gamma} \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^\top \text{diag}(\boldsymbol{\gamma}^{-1}) \mathbf{x} + \log \det(\mathbf{S}(\boldsymbol{\gamma}))$$

- ✓ Alternating Minimization:

- ⊖ \mathbf{x} -update: $\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x}} \sigma^{-2} \|\mathbf{y} - \mathbf{F}\mathbf{x}\|^2 + \mathbf{x}^\top \text{diag}((\boldsymbol{\gamma}^{(k)})^{-1}) \mathbf{x}$
- ⊖ $\boldsymbol{\gamma}$ -update: $\boldsymbol{\gamma}^{(k+1)} = \arg \min_{\boldsymbol{\gamma}} (\mathbf{x}^{(k+1)})^\top \text{diag}(\boldsymbol{\gamma}^{-1}) \mathbf{x}^{(k+1)} + \log \det(\mathbf{S}(\boldsymbol{\gamma}))$

- The challenge comes from the log-determinant term.

- Find a separable surrogate for it.
- Use first-order Taylor expansion.

Alternating Minimization with Linearization

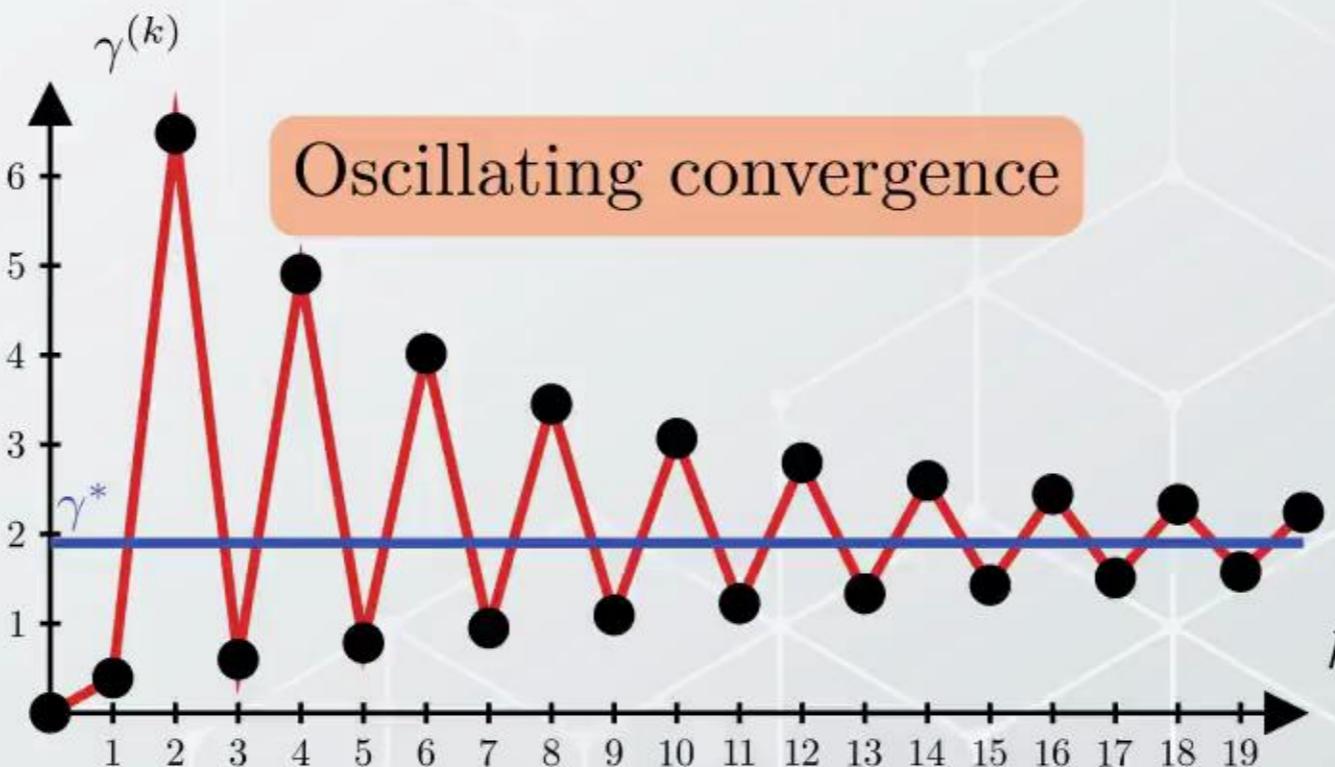
- Consider a linear approximation of $\log \det(\mathbf{S}(\boldsymbol{\gamma}))$
 - We could update each γ_i separately.
 - We have a closed-form solution of each $\gamma_i^{(k+1)}$ in terms of $\gamma_i^{(k)}$

Alternating Minimization with Linearization

- ☒ Consider a linear approximation of $\log \det(\mathbf{S}(\boldsymbol{\gamma}))$
 - ⌚ We could update each γ_i separately.
 - ⌚ We have a closed-form solution of each $\gamma_i^{(k+1)}$ in terms of $\gamma_i^{(k)}$
 - ⌚ We prove the convergence of the algorithm.
 - ⌚ The convergence might not be fast enough in some cases.

Alternating Minimization with Linearization

- ✓ Consider a linear approximation of $\log \det(\mathbf{S}(\boldsymbol{\gamma}))$
 - ⌚ We could update each γ_i separately.
 - ⌚ We have a closed-form solution of each $\gamma_i^{(k+1)}$ in terms of $\gamma_i^{(k)}$
 - ⌚ We prove the convergence of the algorithm.
 - ⌚ The convergence might not be fast enough in some cases.



Alternating Minimization with Quadratic Surrogates

Alternating Minimization with Quadratic Surrogates

- Consider a quadratic surrogate for $\log \det(\mathbf{S}(\boldsymbol{\gamma}))$
 - Add a penalty term $\tau \|\boldsymbol{\gamma} - \boldsymbol{\gamma}^{(k)}\|^2$ to the linear approximation.

Alternating Minimization with Quadratic Surrogates

- ☒ Consider a quadratic surrogate for $\log \det(\mathbf{S}(\boldsymbol{\gamma}))$
 - Ⓐ Add a penalty term $\tau \|\boldsymbol{\gamma} - \boldsymbol{\gamma}^{(k)}\|^2$ to the linear approximation.
 - Update each γ_i separately through minimizing a quadratic surrogate.
 - Have a closed-form solution of each $\gamma_i^{(k+1)}$ in terms of $\gamma_i^{(k)}$

Alternating Minimization with Quadratic Surrogates

- ✓ Consider a quadratic surrogate for $\log \det(\mathbf{S}(\boldsymbol{\gamma}))$
 - ⌚ Add a penalty term $\tau \|\boldsymbol{\gamma} - \boldsymbol{\gamma}^{(k)}\|^2$ to the linear approximation.
 - ⌚ Update each γ_i separately through minimizing a quadratic surrogate.
 - ⌚ Have a closed-form solution of each $\gamma_i^{(k+1)}$ in terms of $\gamma_i^{(k)}$

Theorem [Feng, Shen, and S., 2024]

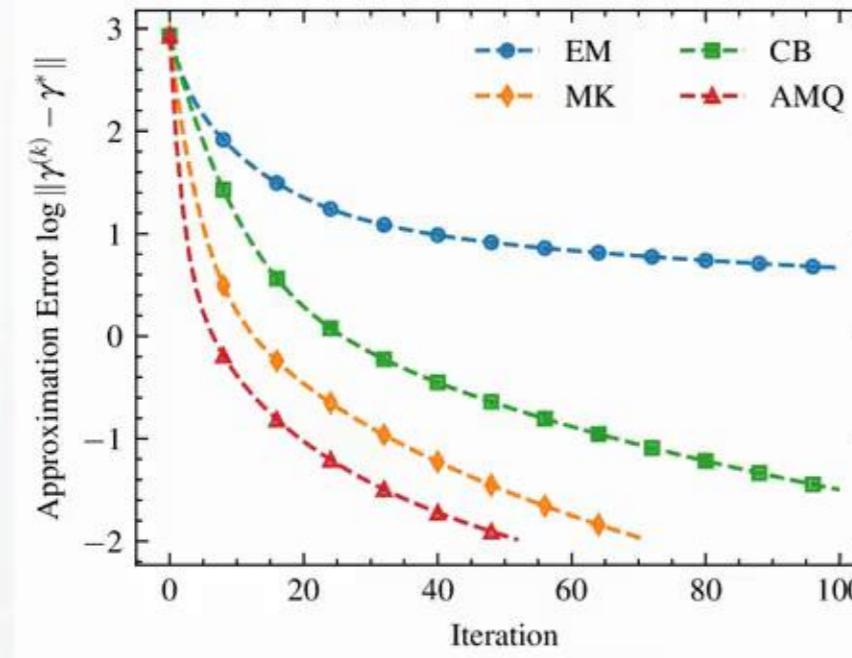
If τ_k is properly chosen, the proposed algorithm converges to a stationary point $\boldsymbol{\gamma}^*$ of the objective function $L(\boldsymbol{\gamma})$. That is, for each i , either $\gamma_i^* = 0$ or $\frac{\partial L(\boldsymbol{\gamma}^*)}{\partial \gamma_i} = 0$.

Numerical Results: Denoising

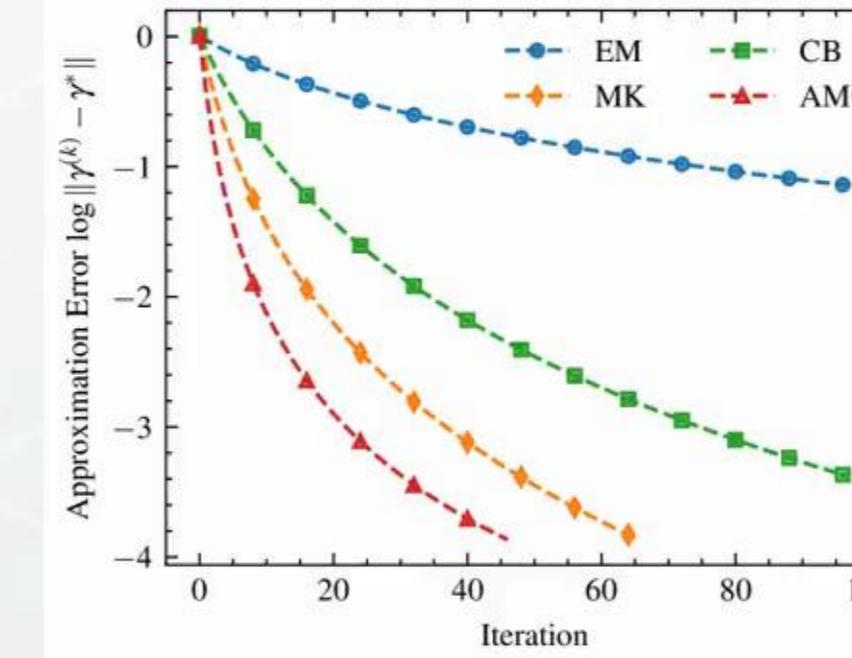
Numerical Results: Denoising

$s = 20\%$

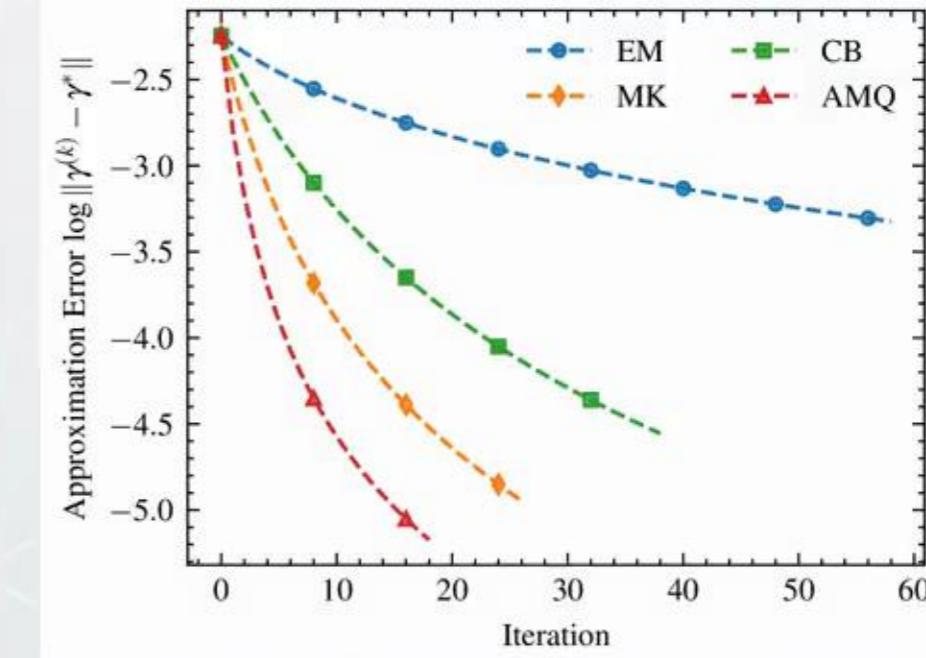
$$\sigma^2 = 10$$



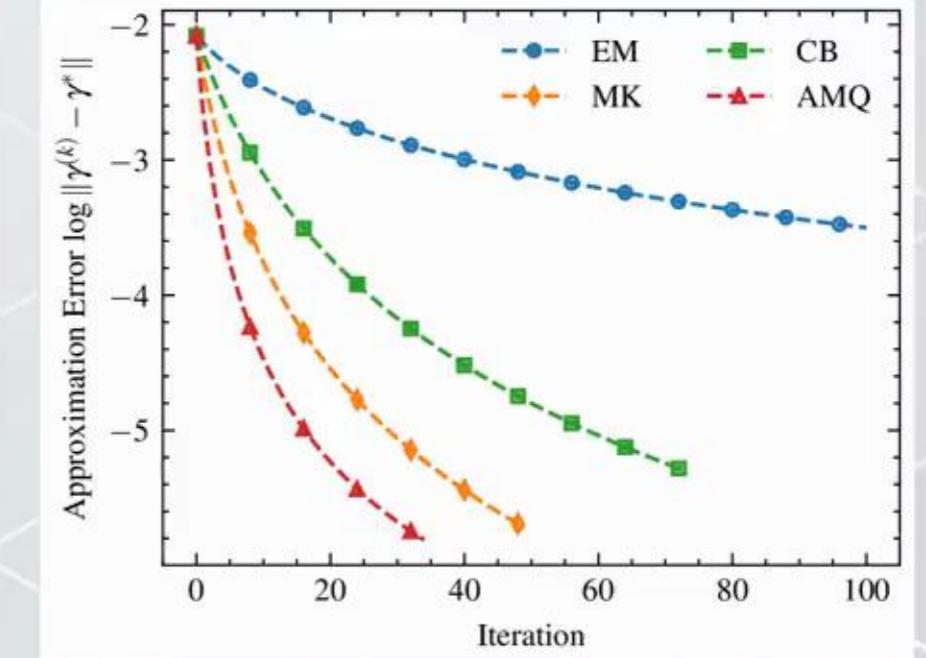
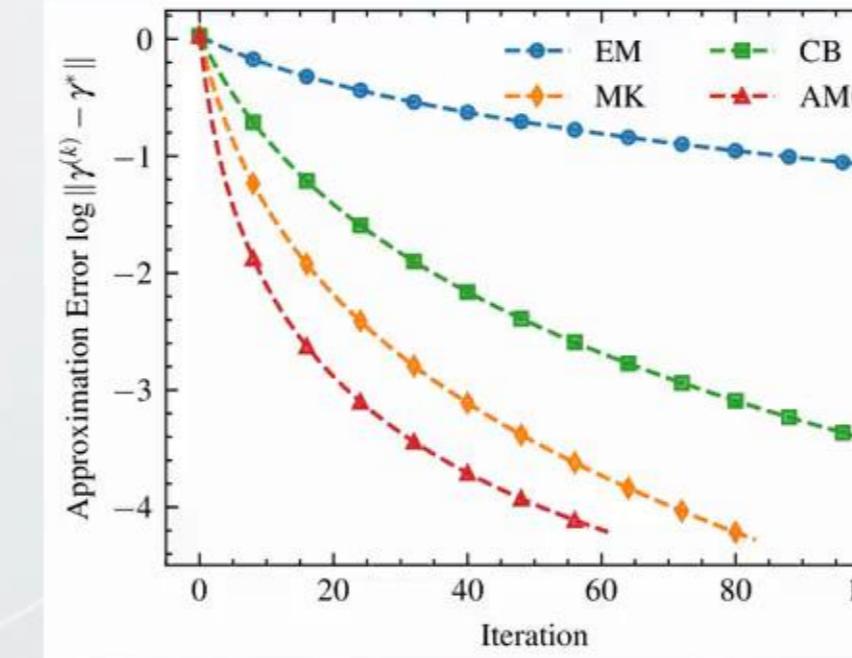
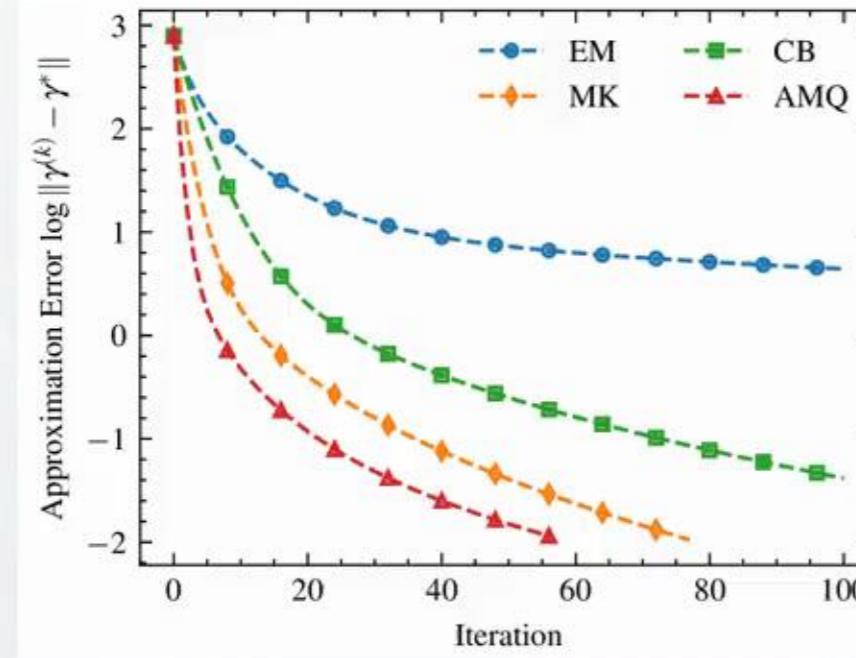
$$\sigma^2 = 1$$



$$\sigma^2 = 0.1$$



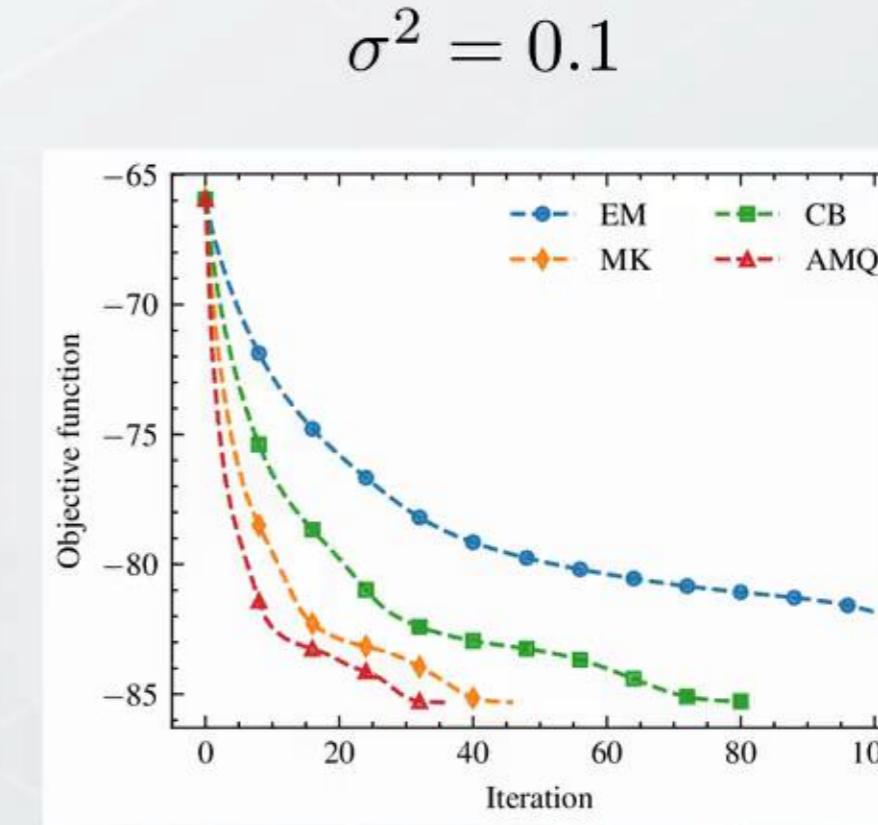
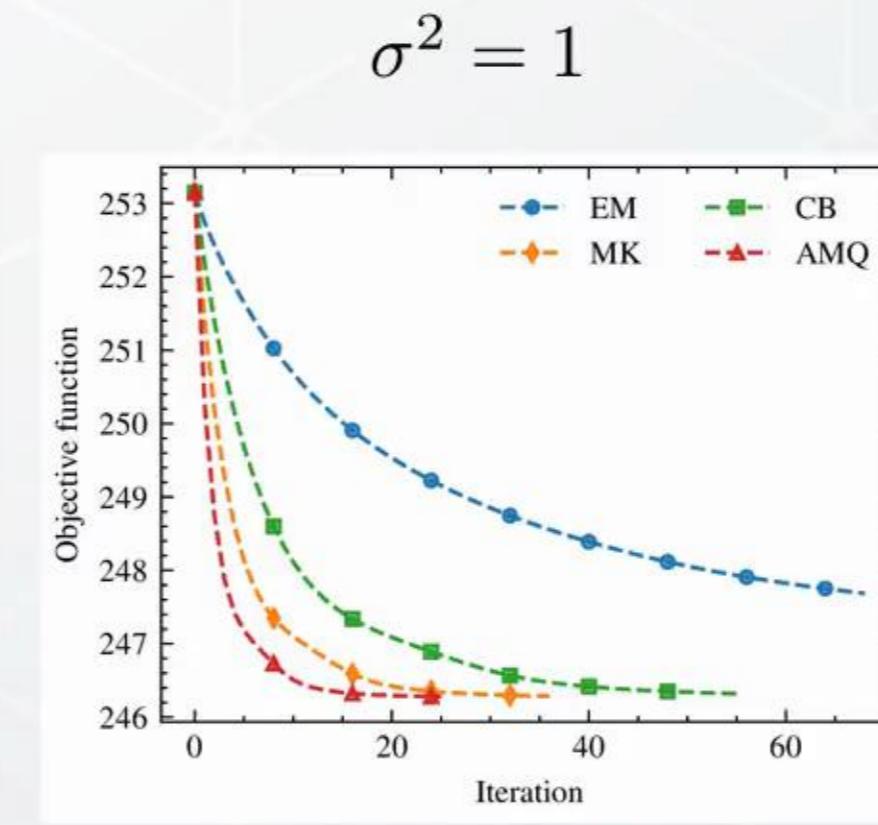
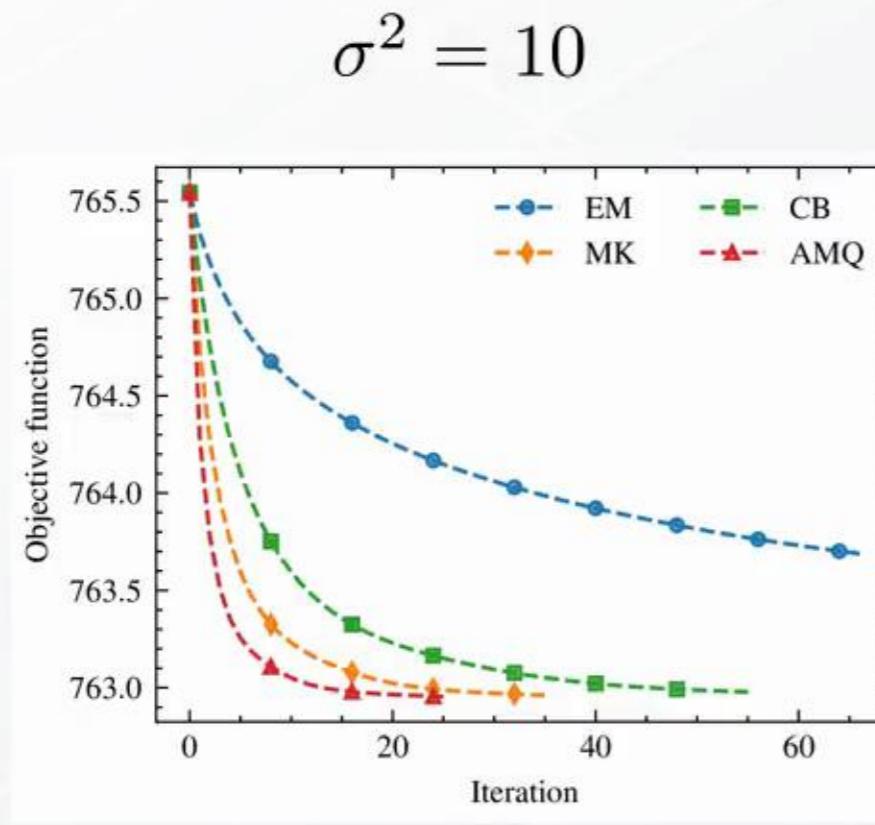
$s = 90\%$



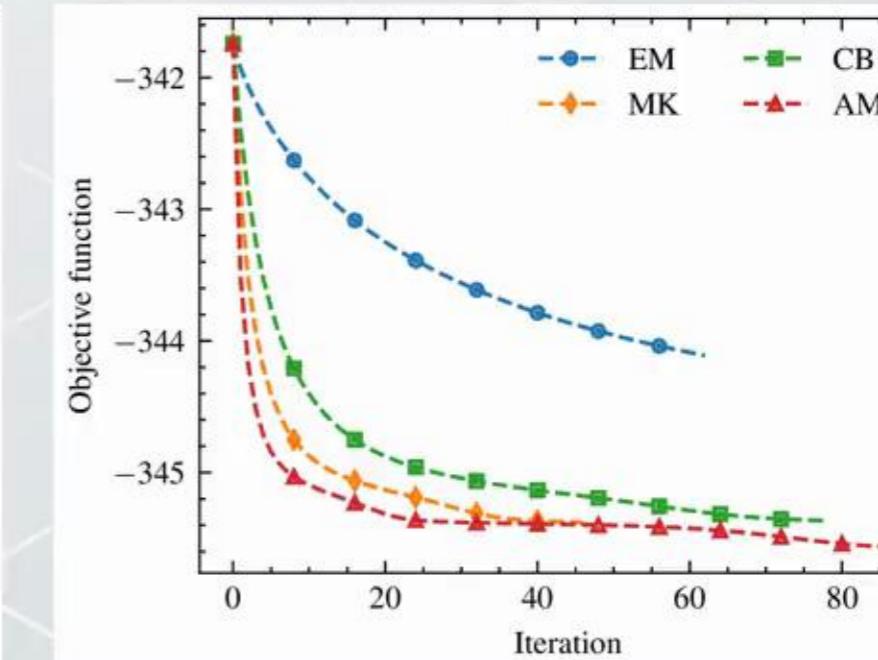
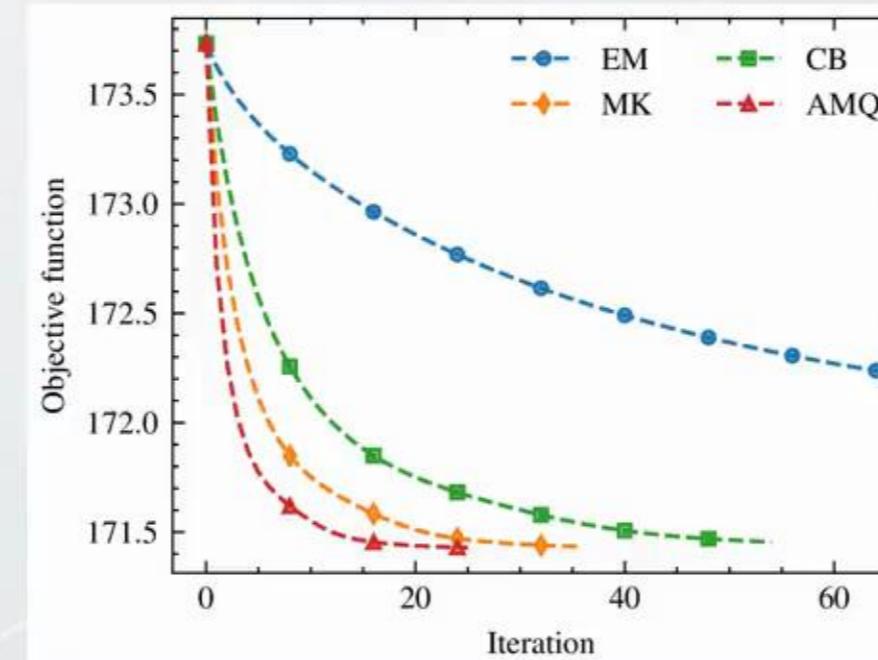
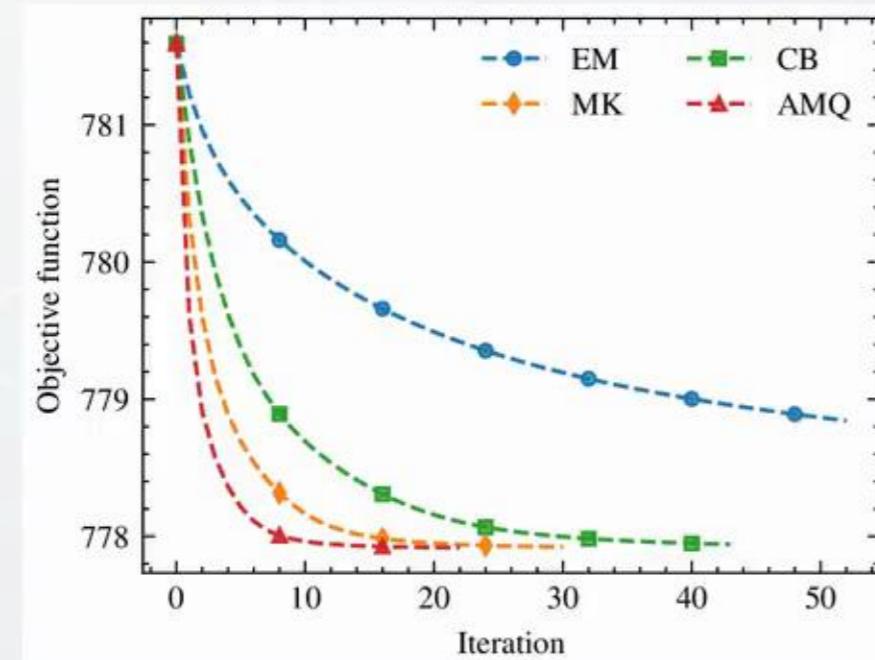
Numerical Results: Partial DCT

Numerical Results: Partial DCT

$s = 20\%$



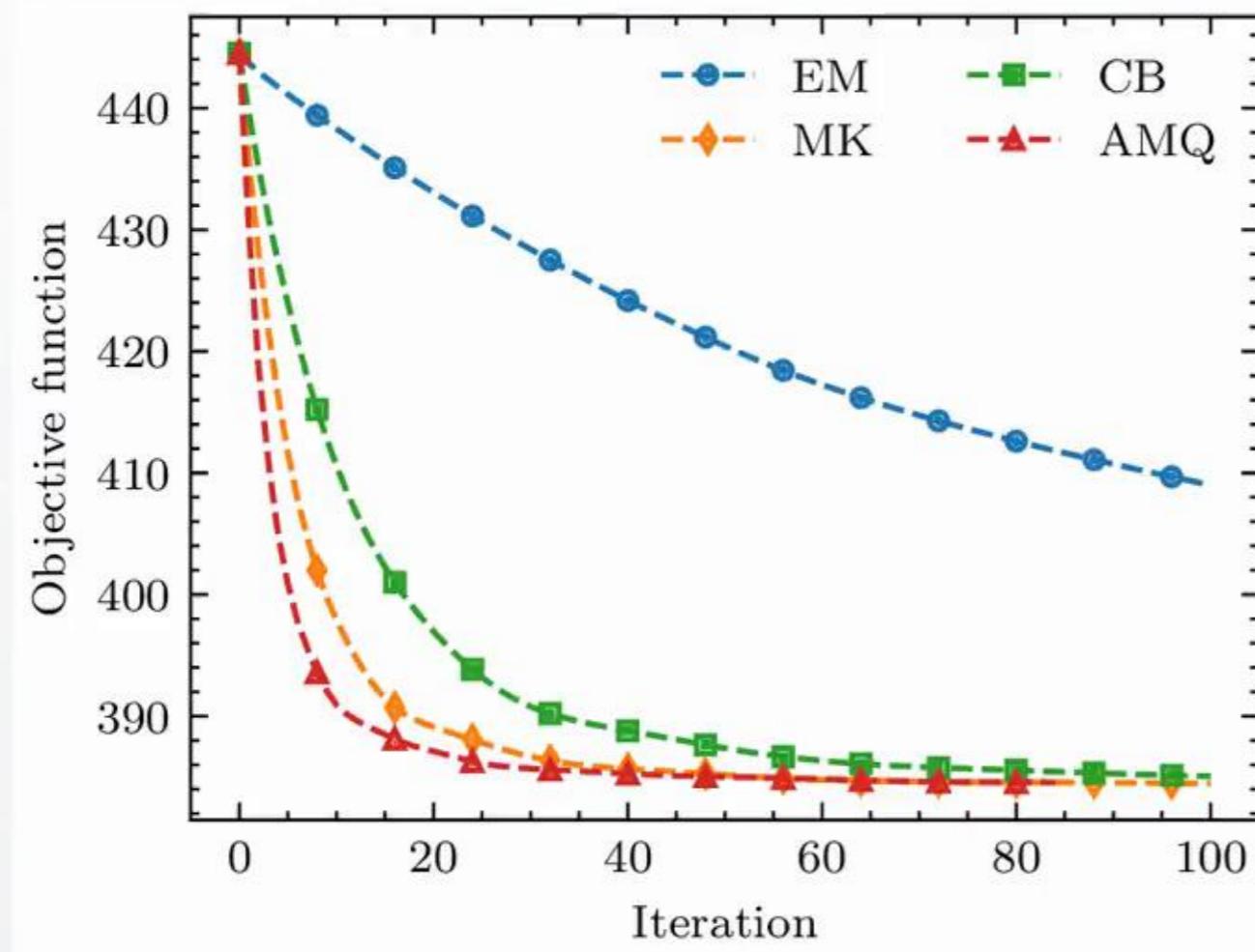
$s = 90\%$



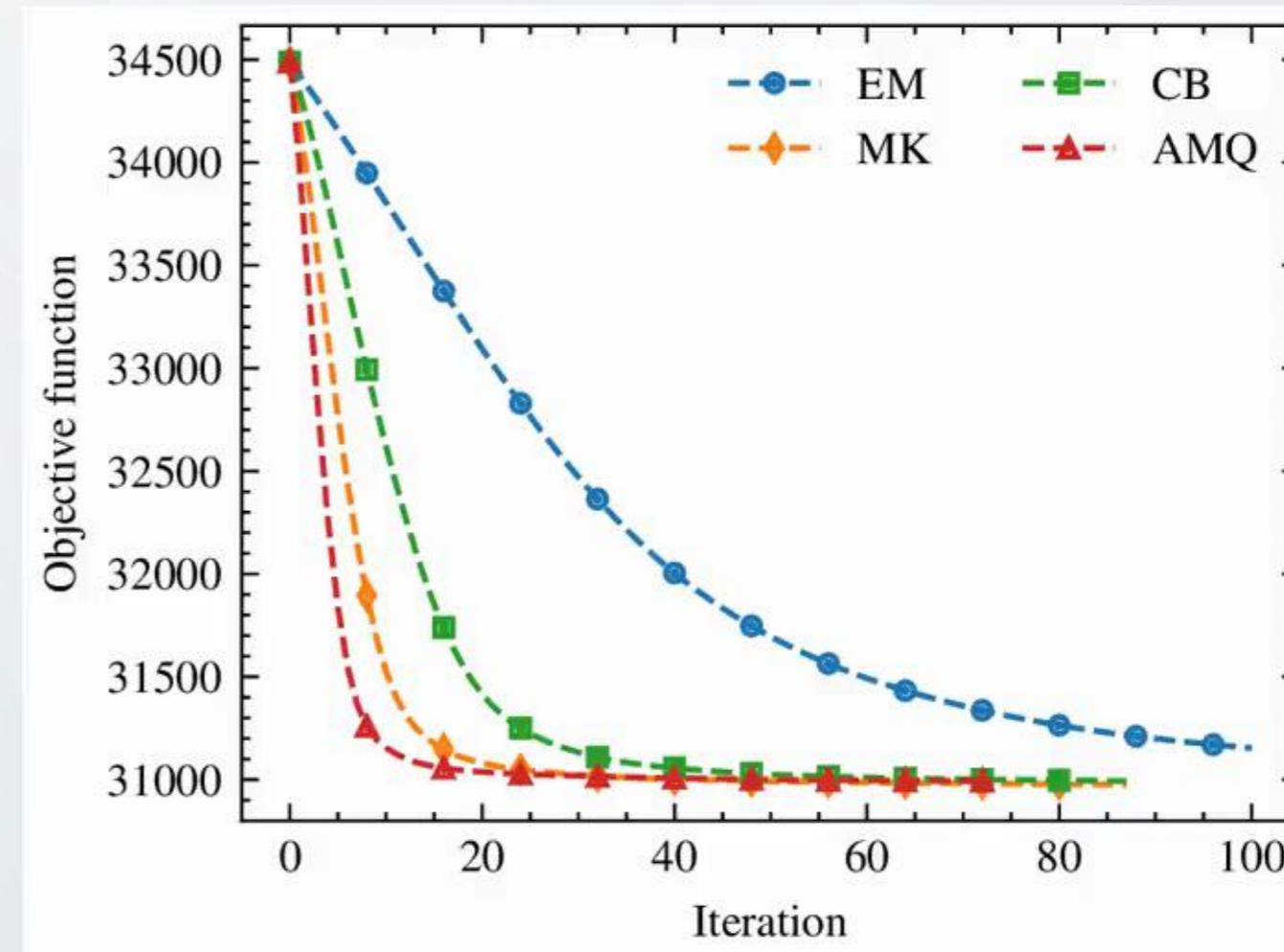
Numerical Results: SAR and EEG

Numerical Results: SAR and EEG

EEG



SAR



Summary

□ Comments

- Bayesian models offer uncertainty quantification.
- Hierarchical Bayesian models provide flexibility and scalability in incorporating prior information.
- We propose a new algorithm for hyper-parameters estimation in hierarchical (linear) Bayesian models.

Summary

Comments

- Ⓐ Bayesian models offer uncertainty quantification.
- Ⓐ Hierarchical Bayesian models provide flexibility and scalability in incorporating prior information.
- Ⓐ We propose a new algorithm for hyper-parameters estimation in hierarchical (linear) Bayesian models.

Issues

- Non-Gaussian distributions
- Non-linear models
- Challenge: more efficient sampling (Monte Carlo) algorithms

Thank you for your attention!

??

Guohui Song, gsong@odu.edu