

On approximate message passing and Expectation-Propagation for scalable inference

Yoann Altmann

Joint work with A. Perelli, D. Yao, S. McLaughlin and M. Davies

*RAEng Research Fellow
School of Engineering and Physical Sciences
Heriot-Watt University*

Edinburgh, February, 20th 2019

Motivations

- Bayesian estimation for complex models
 - Hierarchical
 - Discrete/continuous variables
 - Large scale (imaging)
 - Model selection
- Simulation methods
 - Flexible but not (yet) scalable
- More efficient methods needed

Potential solutions

- Point estimation
 - MAP: convex/non-convex optimization
 - Limited uncertainty quantification
- Approximate methods
 - Proximal MCMC
 - Variational Bayes (VB) methods
 - Approximate message passing (AMP/EP)

Bayesian modeling

- Likelihood (observation model)

$$f(y|x)$$

- (Hierarchical) prior model

$$f(x, \theta) = f(x|\theta)f(\theta)$$

Bayesian modeling

- Exact model

$$f(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) = \frac{f(\mathbf{y}|\mathbf{x})f(\mathbf{x}|\boldsymbol{\theta})f(\boldsymbol{\theta})}{f(\mathbf{y})}$$

- Approximating distribution

$$f(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) \approx q(\mathbf{x}, \boldsymbol{\theta})$$

- Proximal MCMC: Moreau approx.
- VB/EP: Divergence-based
- $f(\mathbf{y})$: evidence (model selection)

Approximate methods

- Variational Bayes

Often $f(\mathbf{x}, \theta | \mathbf{y}) \approx q(\mathbf{x}, \theta) = q(\mathbf{x})q(\theta)$

$$\min_{q(\mathbf{x}), q(\theta)} KL(q(\mathbf{x})q(\theta) || f(\mathbf{x}, \theta | \mathbf{y}))$$

- Expectation Propagation

$$f(\mathbf{y}, \mathbf{x}, \theta) = f(\mathbf{y}|\mathbf{x})f(\mathbf{x}|\theta)f(\theta)$$
$$q(\mathbf{x}, \theta) \propto q_0(\mathbf{x})q_1(\mathbf{x}, \theta)q_0(\theta)$$

$$\min_{Z, q(\mathbf{x}, \theta)} KL(f(\mathbf{y}, \mathbf{x}, \theta) || Zq(\mathbf{x}, \theta))$$
$$Z \approx f(\mathbf{y})$$

EP: preserves better marginals

Iterative minimization

$$f(\mathbf{y}, \mathbf{x}, \theta) = f(\mathbf{y}|\mathbf{x})f(\mathbf{x}|\theta)f(\theta)$$
$$q(\mathbf{x}, \theta) \propto q_0(\mathbf{x})q_1(\mathbf{x}, \theta)q_0(\theta)$$

- $\min_{q_0(\mathbf{x})} KL(f(\mathbf{y}|\mathbf{x})q_1(\mathbf{x}, \theta)q_0(\theta) || q_0(\mathbf{x})q_1(\mathbf{x}, \theta)q_0(\theta))$
- $\min_{q_1(\mathbf{x}, \theta)} KL(q_0(\mathbf{x})f(\mathbf{x}|\theta)q_0(\theta) || q_0(\mathbf{x})q_1(\mathbf{x}, \theta)q_0(\theta))$
- $\min_{q_0(\theta)} KL(q_0(\mathbf{x})q_1(\mathbf{x}, \theta)f(\theta) || q_0(\mathbf{x})q_1(\mathbf{x}, \theta)q_0(\theta))$

... until convergence

Divergence minimization

- Difficult in general
 - Gaussian approx.: moment matching
 - Exponential families: natural parameters
- More difficult with multivariate distributions
- Approximating graph/distributions
 - Tradeoff accuracy/complexity
 - Sequential updates
 - Possibility to add constraints

Example

$$y = Ax + e \Rightarrow f(y|Ax) = \prod_{m=1}^M f(y_m | a_m^T x)$$

$$f(x) = \prod_{n=1}^N f_n(x_n)$$

- A : blur, dictionary, ...
- Arbitrary observation noise (Gaussian, Poisson, ...)
- Separable prior model

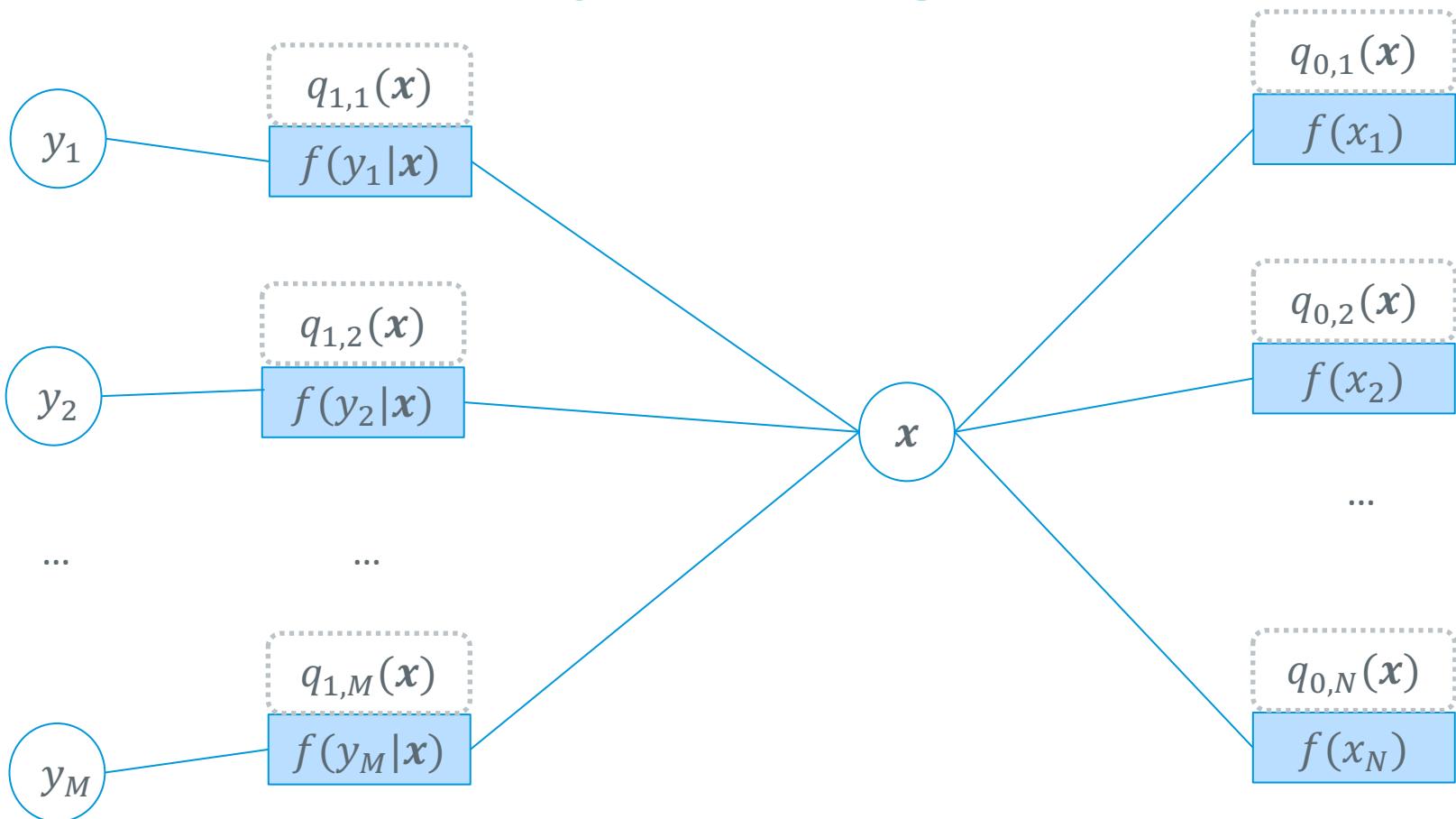
Classical EP scheme

$$f(\mathbf{y}, \mathbf{x}) = \left[\prod_{m=1}^M f(y_m | \mathbf{a}_m^T \mathbf{x}) \right] \left[\prod_{n=1}^N f_n(x_n) \right]$$

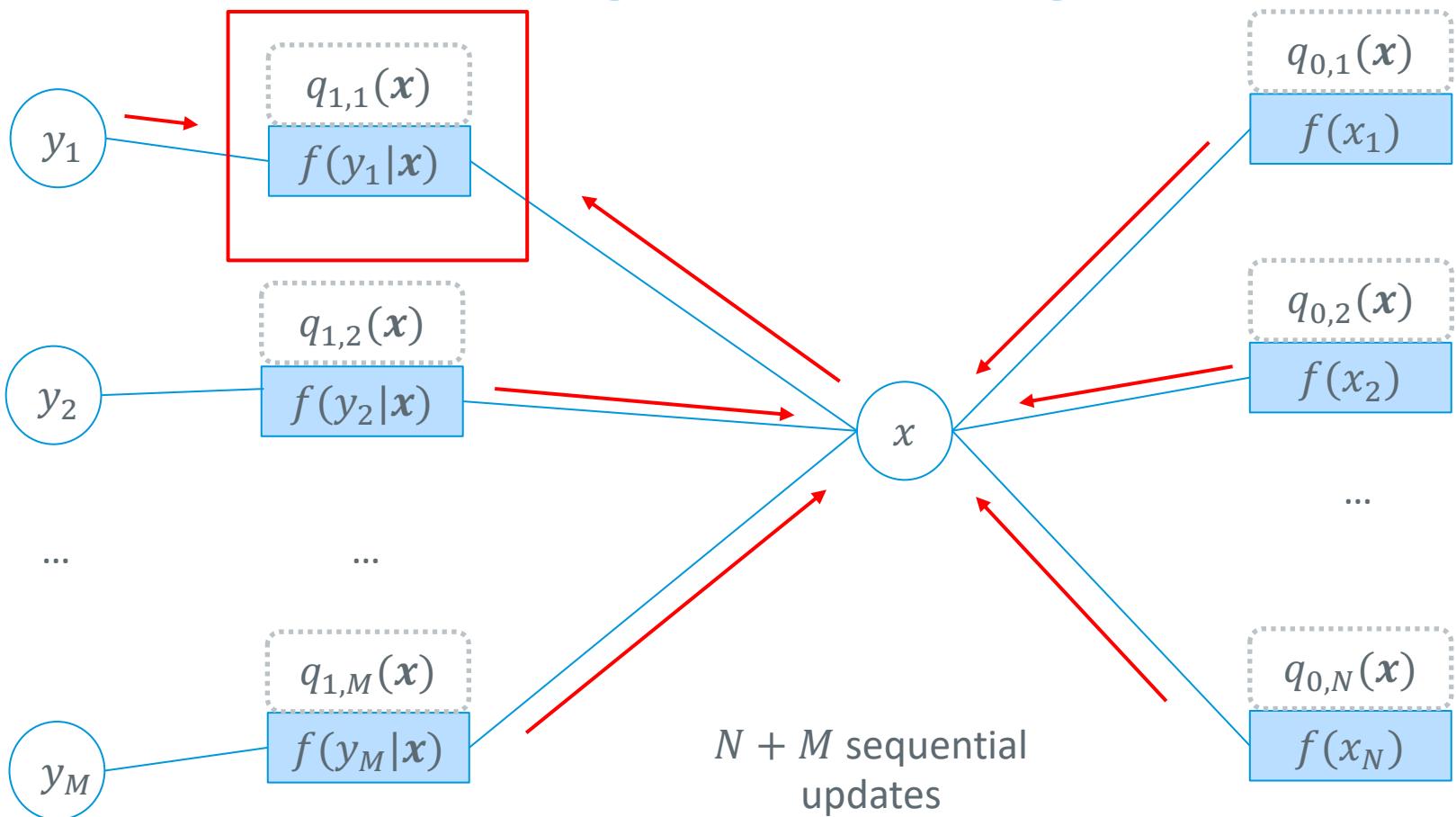
$$q(\mathbf{x}) \propto \left[\prod_{m=1}^M q_{1,m}(\mathbf{x}) \right] \left[\prod_{n=1}^N q_{0,n}(\mathbf{x}) \right]$$

$N + M$ factors in the approximation

Classical Bayesian graph



EP as message passing



EP with compact graph



- Simple models: 2 updates
 - $(q_1(x), q_1(x))$ isotropic cov.: AMP
 - $(q_1(x), q_1(x))$ diagonal cov.: VAMP
 - $f(x)$ implicit: D-AMP
- Messages potentially difficult to compute...

Extended model

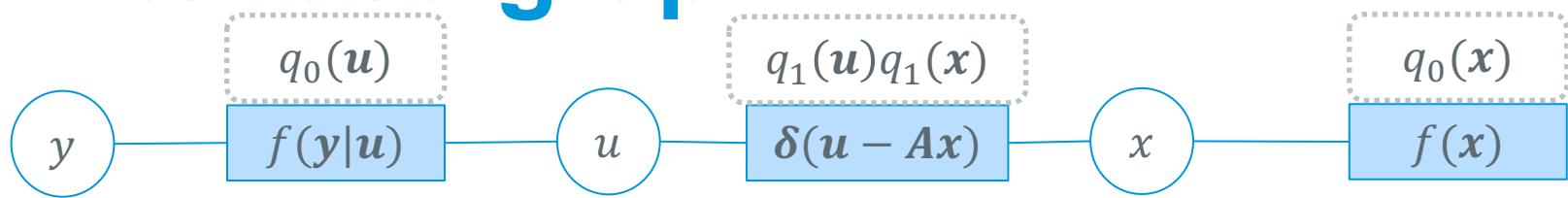
$$f(\mathbf{y}|\mathbf{u}) = \prod_{m=1}^M f(y_m|u_m)$$

$$f(\mathbf{u}|\mathbf{x}) = \prod_{m=1}^M f(u_m|\mathbf{x}) = \prod_{m=1}^M \delta(u_m - a_m^T \mathbf{x})$$

$$f(\mathbf{x}) = \prod_{n=1}^N f_n(x_n)$$

$$f(\mathbf{y}|\mathbf{x}) = \int f(\mathbf{y}|\mathbf{u})f(\mathbf{u}|\mathbf{x})du$$

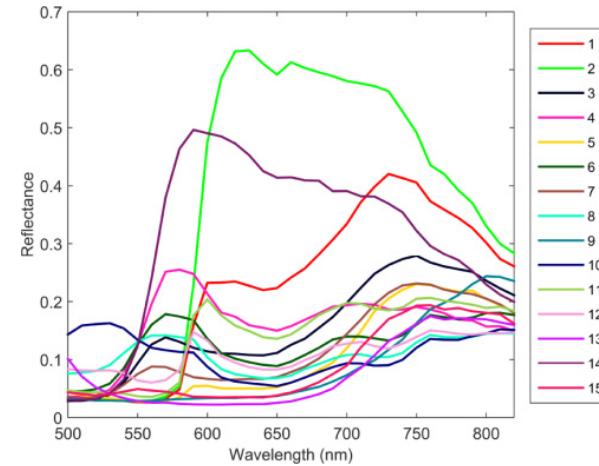
Extended graph



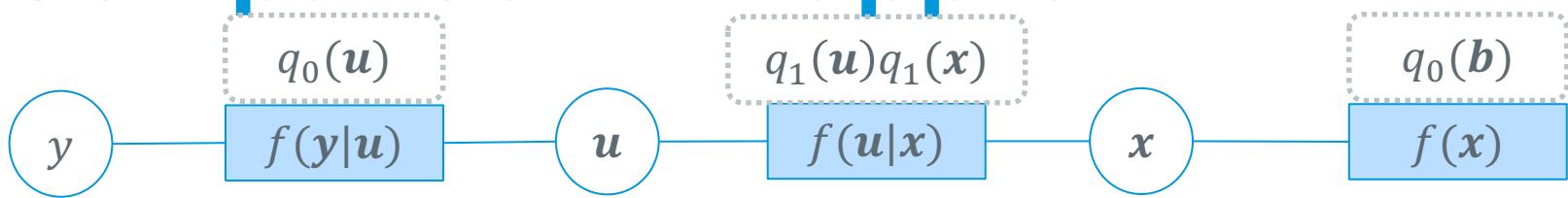
- Flexibility in choosing $q_0(\mathbf{u})$, $q_0(\mathbf{x})$ and $q_1(\mathbf{u}, \mathbf{x})$
 - Partial correlation
- Here (\mathbf{u}, \mathbf{x}) a posteriori independent
 - But covariance for \mathbf{x} and/or \mathbf{u}
- $N \ll M$: regression ; $M \ll N$: CS
- Isotropic/diagonal messages: AMP/VAMP

Example: regression with Poisson noise

- Spectral unmixing in photon-starved regime
- $M = 33$ spectral bands
- $N = 15$ sources (polymer clay)

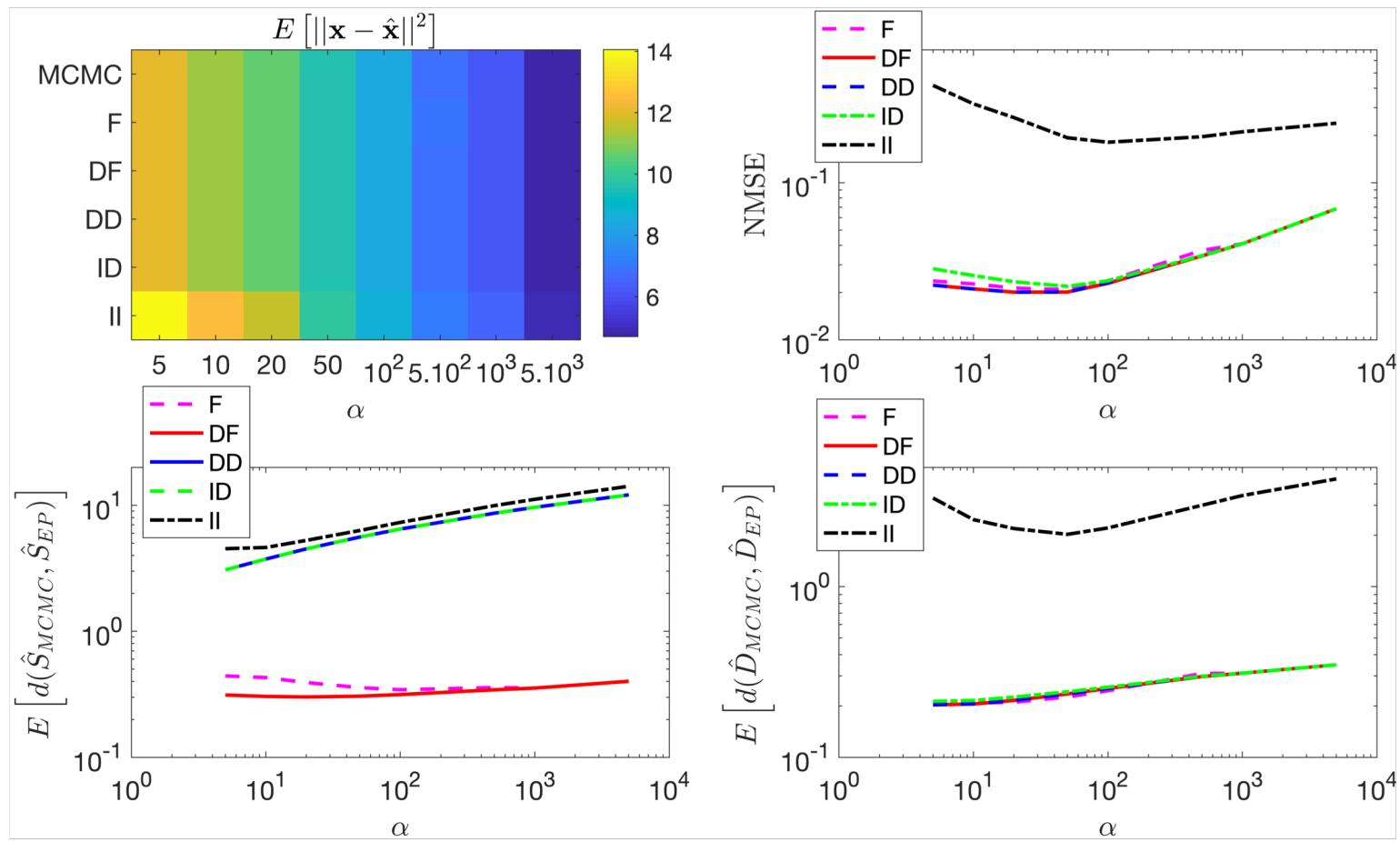


Comparison EP approx.

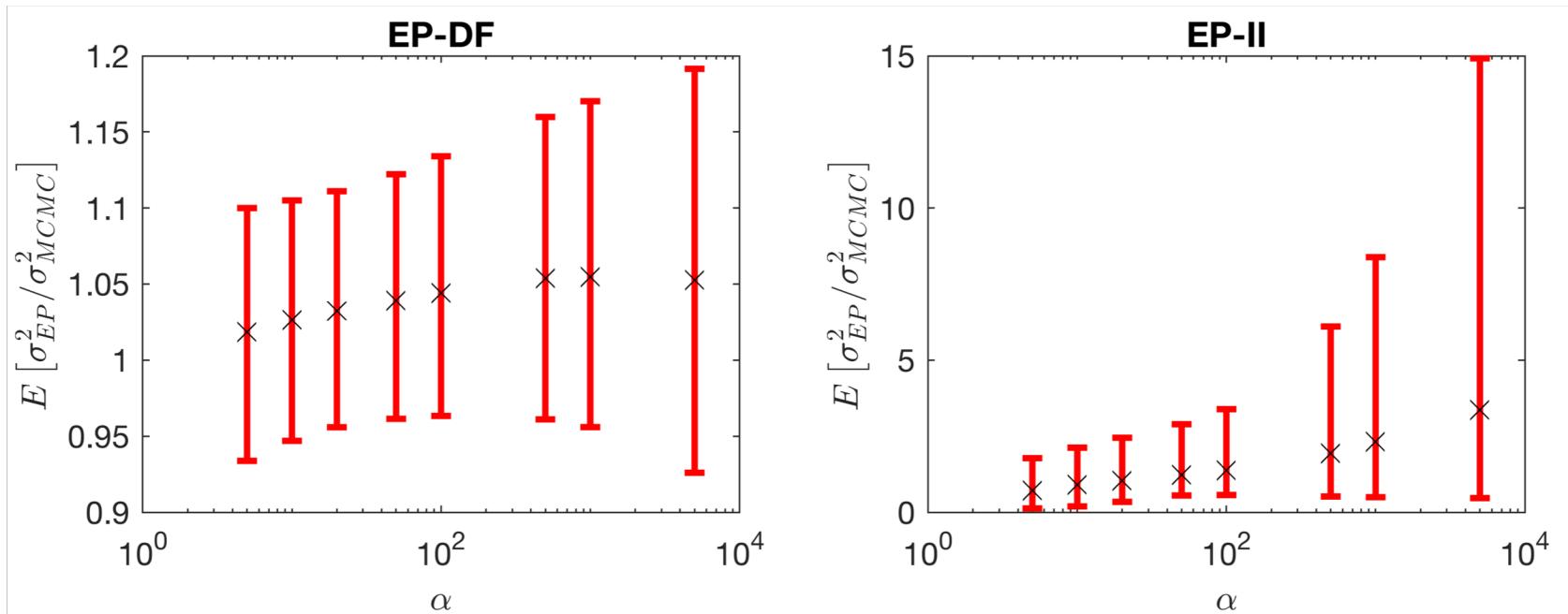


- Exponential priors $f(x)$
- Gaussian approx. only
- Different choices of covariance matrices
- $q_0(u)$: diag. (D) $q_1(u)$: iso./diag.
- $q_0(x)$: diag. $q_1(x)$: iso./diag./full (I/D/F)
- Positivity not enforced in $q(x)$

Results



Results



EP-DF only slightly overestimates the marginal variances

Conclusion

- Approximate message passing
 - Scalable/distributed
 - Flexible
 - Fast convergence
- Good estimation performance
- Approx. evidence: by-product
- But can be difficult to implement

Ongoing work

- Convergence properties of EP
 - Constraints seem to help
- Approximations
 - for non-Gaussian noise
 - Hard constraints (e.g., positivity)
- EP within simulation schemes
- Application to imaging (Denoising-AMP)