

# Low Complexity Static and Dynamic Sparse Bayesian Learning Combining BP, VB and EP Message Passing

Christo Kurisummoottil Thomas, Dirk Slock

kurisumm@eurecom.fr, slock@eurecom.fr



State of the Art

Novel Contributions

## Why Combined BP/MF/EP ?

- ▶ Dynamic autoregressive SBL (DAR-SBL): a case of joint Kalman filtering (KF) with a linear time-invariant diagonal state-space model, and parameter estimation, which can be considered an instance of nonlinear filtering.
- ▶ The authors in [RieglerTIT2013] propose a message passing (MP) approach for inferring the posteriors combining belief propagation (BP) and mean field (MF) approximations. As usual, they apply BP to detection (discrete parameters) only and MF to all estimation (continuous parameters).
- ▶ The advantages of the MF approach are that it always admits a convergent implementation while BP yields a good approximation of the posterior marginals if the factor graph has no cycles.
- ▶ Expectation Propagation (EP)-Gaussian approx., intractable distributions.

## Sparse Bayesian Learning (SBL)

- ▶ Dynamic SBL (AR(1) state space model):

State Update:  $\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \mathbf{w}_t$ ,  $\mathbf{w}_t \sim \mathcal{CN}(\mathbf{0}, \mathbf{\Lambda}^{-1})$

Observation:  $\mathbf{y}_t = \mathbf{A}^{(t)}\mathbf{x}_t + \mathbf{v}_t$ ,  $\mathbf{v}_t \sim \mathcal{CN}(\mathbf{0}, \frac{1}{\gamma}\mathbf{I})$ .



- ▶ SBL for compressed sensing does not exactly sparsify  $\mathbf{x}_t$  but works well in the case of relatively limited data (eg underdetermined) in which case estimation emphasis is given to large unknowns while small unknowns get mostly ignored.
- ▶ Gamma prior to the precision of the state  $\mathbf{x}_t$  and the innovation sequences  $\mathbf{w}_t$  (both with same support), allowing to sparsify the components of  $\mathbf{x}_t$ .  
 $\mathbf{\Lambda}^{-1} = \mathbf{\Gamma}(\mathbf{I} - \mathbf{F}\mathbf{F}^H) = \text{diag}(\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_M})$ ,  $\boldsymbol{\theta} = \{\mathbf{x}, \mathbf{\Lambda}, \gamma, \mathbf{F}\}$

## Combined BP and MF Approximation

The fixed points of the BP algorithm are shown to be the stationary points of the Bethe free energy (BFE) [RieglerTIT2013]. However, for the MF approximation in variational Bayes (VB) [BealThesis2003], the approximate posteriors are shown to be converging to a local minimum of the MF free energy which is an approximation of the BFE.

$$F_{BP, MF} = \sum_{a \in \mathcal{A}_{BP}} \sum_{\theta_a} q_a(\theta_a) \ln \frac{q_a(\theta_a)}{f_a(\theta_a)} - \sum_{a \in \mathcal{A}_{MF}} \sum_{i \in \mathcal{N}(a)} \prod_{\theta_i} q_i(\theta_i) \ln f_a(\theta_a) - \sum_{i \in \mathcal{I}} (|\mathcal{N}_{BP}(i)| - 1) \sum_{\theta_i} q_i(\theta_i) \ln q_i(\theta_i).$$

$$\text{Constraints: } \sum_{\theta_i} q_i(\theta_i) = 1, \forall i \in \mathcal{I}_{MF} \setminus \mathcal{I}_{BP}, \sum_{\theta_a} q_a(\theta_a) = 1,$$

$$\forall a \in \mathcal{A}_{BP}, q_i(\theta_i) = \sum_{\theta_a} q_a(\theta_a), \forall a \in \mathcal{A}_{BP}, i \in \mathcal{N}(a).$$

## Message Passing Expressions

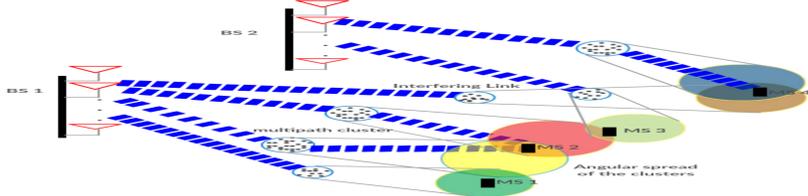
The fixed point equations corresponding to the constrained optimization of BFE can be written as follows [RieglerTIT2013], (a-factor nodes, i-variable nodes)

$$q_i(\theta_i) = z_i \prod_{a \in \mathcal{N}_{BP}(i)} m_{a \rightarrow i}^{BP}(\theta_i) \prod_{a \in \mathcal{N}_{MF}(i)} m_{a \rightarrow i}^{MF}(\theta_i), \quad n_{i \rightarrow a}(\theta_i) = \prod_{a \in \mathcal{N}_{BP}(i) \setminus a} m_{a \rightarrow i}(\theta_i) \prod_{a \in \mathcal{N}_{MF}(i)} m_{a \rightarrow i}(\theta_i),$$

$$m_{a \rightarrow i}^{MF}(\theta_i) = \exp(\langle \ln f_a(\theta_a) \rangle_{\prod_{j \in \mathcal{N}(a) \setminus i} n_{j \rightarrow a}(\theta_j)}), \quad m_{a \rightarrow i}^{BP}(\theta_i) = \left( \int \prod_{j \in \mathcal{N}(a) \setminus i} n_{j \rightarrow a}(\theta_j) f_a(\theta_a) \prod_{j \neq i} d\theta_j \right).$$

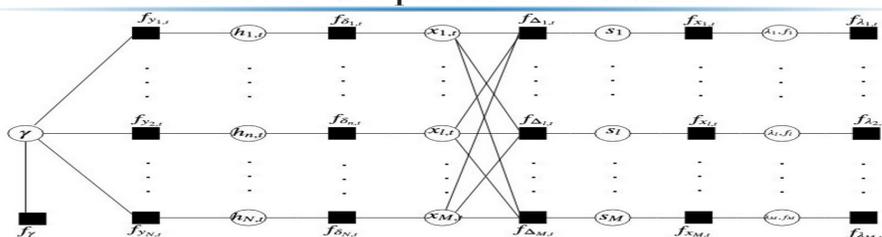
## Applications

- ▶ Massive MIMO channel estimation:



- ▶ Bayesian adaptive filtering.

## Factor Graph for DAR-SBL



All the messages (beliefs or continuous pdfs) passed between the nodes in the above factor graph can be shown to be Gaussian [TanLiTIT2013], parameterized by the mean and variance of the beliefs. With the hard constraints, the equivalent observation model:

$$y_n - \sum_{l \neq n} A_{n,l} \tilde{x}_{l,n} = A_{n,n} x_n + \sum_{l \neq n} A_{n,l} \tilde{x}_{l,n} + v_n, \text{ where,}$$

$$\tilde{x}_{l,n} \sim \mathcal{CN}(0, \nu_{l,n}), \text{ and } m_{f_{n \rightarrow x}} \propto \mathcal{CN}(x; \hat{x}_{n,l}, \nu_{n,l}),$$

## Combined BP-MF-EP DAR SBL

Initialization  $\hat{f}_{l|0}, \hat{\lambda}_{l|0} = \frac{a}{b}, \hat{\gamma}_0 = \frac{c}{d}, \hat{x}_{l,0|0} = 0, \sigma_{l,0|0}^2 = 0, \forall l$ . Define  $\boldsymbol{\Sigma}_{t-1|t-1} = \text{diag}(\sigma_{l,t|t-1}^2)$ . for  $t = 1 : T$  do

**Prediction Stage:**

- ▶ Compute  $\hat{x}_{l,t|t-1}, \sigma_{l,t|t-1}^2$  using EP.

**Filtering Stage (BP for  $\hat{x}_{l,t|t}, \sigma_{l,t|t}^2$ ):**

- ▶ Compute  $\hat{x}_{n,t}^{(t)}, \nu_{n,t}^{(t)}$  and update  $\hat{x}_{l,t|t}, \sigma_{l,t|t}^2$ .
- ▶ Compute  $\nu_{l,n}^{(t)}, \hat{x}_{l,n}^{(t)}$ .
- ▶ Continue steps 1) to 2) until convergence.

**Smoothing Stage:**

**Initialization:**  $\boldsymbol{\Sigma}_{t-1|t-1}^{(0)} = \boldsymbol{\Sigma}_{t-1|t-1}, \hat{\mathbf{x}}_{t-1|t-1}^{(0)} = \hat{\mathbf{x}}_{t-1|t-1}$ . Define  $\mathbf{B}^{(t)} = \mathbf{F}^T \mathbf{A}^{(t)T} \tilde{\mathbf{R}}_t^{-1} \mathbf{A}^{(t)} \mathbf{F} + \boldsymbol{\Sigma}_{t-1|t-1}^{(0)}$ ,  $\mathbf{h}_t = \mathbf{F}^T \mathbf{A}^{(t)T} \tilde{\mathbf{R}}_t^{-1} \mathbf{y}_t$ ,  $\tilde{\mathbf{R}}_t = \mathbf{A}^{(t)} \mathbf{\Lambda}^{-1} \mathbf{A}^{(t)H} + \frac{1}{\gamma} \mathbf{I}$ .

- ▶  $P_{i,j} = \frac{-B_{i,j}^{(0)}}{B_{i,i}^{(0)} + \sum_{k \in \mathcal{N}(i) \setminus j} P_{k,i} \mu_{k,i}}$ ,  $\mu_{i,j} = (h_{i,t} + \sum_{k \in \mathcal{N}(i) \setminus j} P_{k,i} \mu_{k,i}), \forall i, j$ .
- ▶  $\sigma_{i,t|t-1}^2 = B_{i,i}^{(0)} + \sum_{k \in \mathcal{N}(i)} P_{k,i} \hat{x}_{k,t-1|t-1} = \sigma_{i,t-1|t-1}^2 (h_{i,t} + \sum_{k \in \mathcal{N}(i)} P_{k,i} \mu_{k,i})$
- ▶  $\boldsymbol{\Sigma}_{t-1|t}^{(i)} = (\hat{\mathbf{F}}_t^H \mathbf{A}^{(i)H} \tilde{\mathbf{R}}_t^{-1} \mathbf{A}^{(i)} \hat{\mathbf{F}}_t + \text{diag}(\mathbf{A}^{(i)H} \tilde{\mathbf{R}}_t^{-1} \mathbf{A}^{(i)}) \boldsymbol{\Sigma}_{\mathbf{F}_t} + \boldsymbol{\Sigma}_{t-1|t}^{(i-1)})$ ,  $\tilde{\mathbf{R}}_t = \mathbf{A}^{(t)} \mathbf{\Lambda}^{-1} \mathbf{A}^{(t)H} + \frac{1}{\gamma} \mathbf{I}$ .
- ▶  $\hat{\mathbf{x}}_{t-1|t}^{(i)} = \hat{\boldsymbol{\Sigma}}_{t-1|t}^{(i)} (\hat{\boldsymbol{\Sigma}}_{t-1|t-1}^{-1} \hat{\mathbf{x}}_{t-1|t-1}^{(i-1)} + \hat{\mathbf{F}}_t^H \mathbf{A}^{(i)H} \tilde{\mathbf{R}}_t^{-1} \mathbf{y}_t)$ .

**Estimation of hyperparameters (Define:  $x'_{k,t} = x_{k,t} - f_k x_{k,t-1}, \zeta_t = \beta \zeta_{t-1} + (1 - \beta) \langle \|\mathbf{y}_t - \mathbf{A}^{(t)} \mathbf{x}_t\|^2 \rangle$ ):**

- ▶ Compute  $\hat{f}_{l|t}, \sigma_{f_{l|t}}^2$  using MF rule,  $\hat{\gamma}_t = \frac{c+N}{(\zeta_t+d)}$  and  $\lambda_{l|t} = \frac{(a+1)}{(\langle |x'_{k,t}|^2 \rangle_{t+b})}$ .

## Optimal Partitioning of BP/MF nodes

mCRB refers to mismatched CRB (CRB under model misspecification) [RichmondTSP15].

**Theorem:** If the parameter partitioning in VB is such that the different parameter blocks are decoupled at the level of Fisher Information Matrix (FIM), then VB is not suboptimal in terms of (mismatched) Cramer-Rao Bound. If a finer partitioning granularity is used (such as up to scalar level as in MF), then VB becomes quite suboptimal, which can be alleviated by using BP instead.

$$mCRB_{BP} = \text{blkdiag}(CRB) = \text{blkdiag}(FIM^{-1}),$$

$$mCRB_{VB} = (\text{blkdiag}(FIM))^{-1},$$

$$\text{So, } mCRB_{BP} = mCRB_{VB} \text{ if } FIM = \text{blkdiag}(FIM).$$

**Hence:** BP may also improve parameter estimation.

## Identifiability

Non-singularity of FIM  $\implies$  local identifiability.

**Lemma:** The AR(1) model parameters require (at least lag 1) smoothing for identifiability.

For the AR(1) parameters, we obtain the FIM submatrix  $\mathbf{J}_t = \mathbf{D} - \mathbf{D}(\mathbf{D} + \mathbf{J}_{t-1})^{-1} \mathbf{D}$ , if  $\mathbf{J}_{-1} = 0$ , then  $\mathbf{J}_t = 0, \forall t \geq 0$ . For  $\mathbf{x}_t$ ,  $\mathbf{J}_{\mathbf{x},t} = \mathbf{\Lambda} + \gamma \mathbf{A}^{(t)H} \mathbf{A}^{(t)} + \mathbf{\Lambda} \mathbf{F} (\mathbf{F} \mathbf{\Lambda} \mathbf{F}^H + \mathbf{J}_{\mathbf{x},t-1})^{-1} \mathbf{\Lambda} \mathbf{F}^H$ , diagonal  $\mathbf{J}_{\mathbf{x},t} \implies$  MF is sufficient for  $\mathbf{x}_t$  at prediction stage.

## Smoothing

$\mathbf{y}_t = \mathbf{A}^{(t)} \mathbf{F} \mathbf{x}_{t-1} + \tilde{\mathbf{v}}_t$ , where  $\tilde{\mathbf{v}}_t = \mathbf{A}^{(t)} \mathbf{w}_{t-1} + \mathbf{v}_t$ , Define  $\mathbf{D} = (\mathbf{I} - \mathbf{F}\mathbf{F}^H)^{-1}$ ,

$$\mathbf{J}_{\mathbf{F},t} = \mathbf{\Gamma} \text{diag}(\mathbf{A}^{(t)H} \tilde{\mathbf{R}}_t^{-1} \mathbf{A}^{(t)}) + \mathbf{D} - \mathbf{D}(\mathbf{D} + \mathbf{J}_{\mathbf{F},t})^{-1} \mathbf{D}, \quad \mathbf{J}_{p,t} = \begin{bmatrix} \mathbf{J}_{\mathbf{\Lambda},t} & \mathbf{J}_{\mathbf{\Lambda},\gamma,t} \\ \mathbf{J}_{\mathbf{\Lambda},\gamma,t} & \mathbf{J}_{\gamma,t} \end{bmatrix}$$

We obtain  $\mathbf{J}_{\mathbf{x},t} = \mathbf{F}^T \mathbf{A}^{(t)H} \tilde{\mathbf{R}}_t^{-1} \mathbf{A}^{(t)} \mathbf{F} + \mathbf{\Lambda} - \mathbf{\Lambda} \mathbf{F} (\mathbf{F} \mathbf{\Lambda} \mathbf{F}^H + \mathbf{J}_{\mathbf{x},t-1})^{-1} \mathbf{\Lambda} \mathbf{F}^H$ , which is a full matrix.

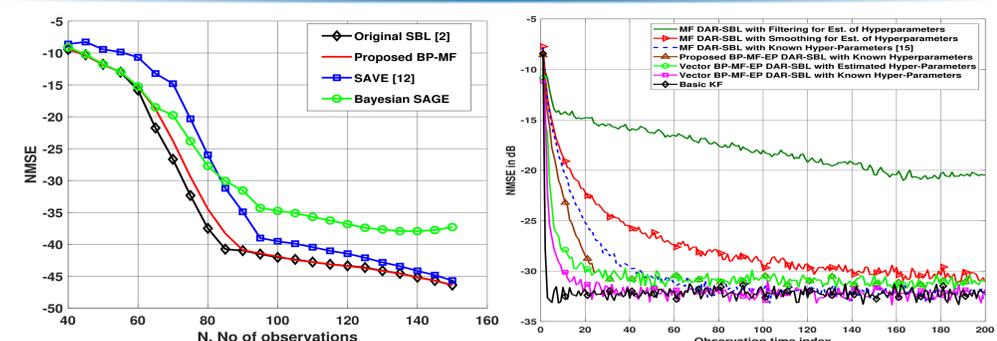
**Conclusions:**

- ▶ BP required for smoothing of  $\mathbf{x}_t, \mathbf{x}_{t-1|t}$ .
- ▶ Lag 1 smoothing is sufficient for AR(1) parameters.

**Corollary**

- ▶ For the smoothing stage, an optimal partitioning is to apply BP for estimation of the state vector  $\mathbf{x}$ ,  $\hat{\mathbf{x}}_{t-1|t}$  and MF for  $\mathbf{F}$ .

## Numerical Results



Static SBL

DAR-SBL

## Conclusion and Future Work

- ▶ Introduced a fast SBL algorithm called BP-MF-EP DAR-SBL, which uses a combination of BP, VB and EP techniques to better approximate the posteriors of the sparse state  $\mathbf{x}$  and parameters and track a time varying  $\mathbf{x}$ .
- ▶ Extension to dictionary learning ongoing.