

SPACE ALTERNATING VARIATIONAL ESTIMATION AND KRONECKER STRUCTURED DICTIONARY LEARNING

Christo Kurisummoottil Thomas, Dirk Slock

kurisumm@eurecom.fr, slock@eurecom.fr



Motivation

Why SAVED-KS ?

- ▶ VB: **analytical approximations to the posterior distributions** of interest even when exact inference of these distributions is intractable.
- ▶ Propose a novel fast algorithm called **space alternating variational estimation with Kronecker structured dictionary learning (SAVED-KS)**, which is a version of VB(-SBL) pushed to the scalar level.
- ▶ The component-wise approach of SAVE compared to SBL renders it less likely to get stuck in bad local optima and its inherent damping (more cautious progression) also leads to typically **faster convergence of the non-convex optimization process**
- ▶ Unstructured KS dictionary matrices learning

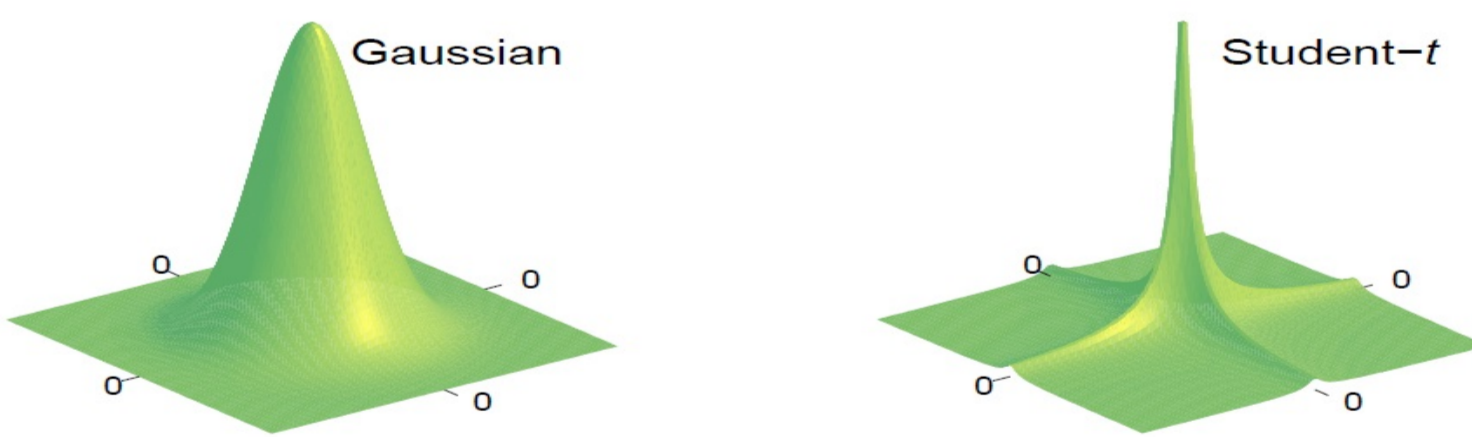
Sparse Bayesian Learning

- ▶ Bayesian Compressed Sensing: 2-layer hierarchical prior for \mathbf{x} as in [Tipping:JMLR01, WipfRao:TSP04], inducing sparsity for \mathbf{x} .

$$p(x_i|\alpha_i) = \mathcal{N}(0, \alpha_i^{-1}), p(\alpha_i/a, b) = \Gamma^{-1}(a)b^a \alpha_i^{a-1} e^{-b\alpha_i}$$

\Rightarrow sparsifying Student-t marginal

$$p(x_i) = \frac{b^a \Gamma(a+\frac{1}{2})}{(2\pi)^{\frac{1}{2}} \Gamma(a)} (b + x_i^2/2)^{-(a+\frac{1}{2})}$$



- ▶ SBL which is geared towards compressed sensing assuming sparse unknown vectors. It's just that the SBL approach works well in the case of relatively limited data (for a given state-space dimension) in which case estimation emphasis is given to large unknowns and small unknowns get more or less ignored.
- ▶ We apply the (Gamma) prior to the precision of the state x , allowing to sparsify the components of x .

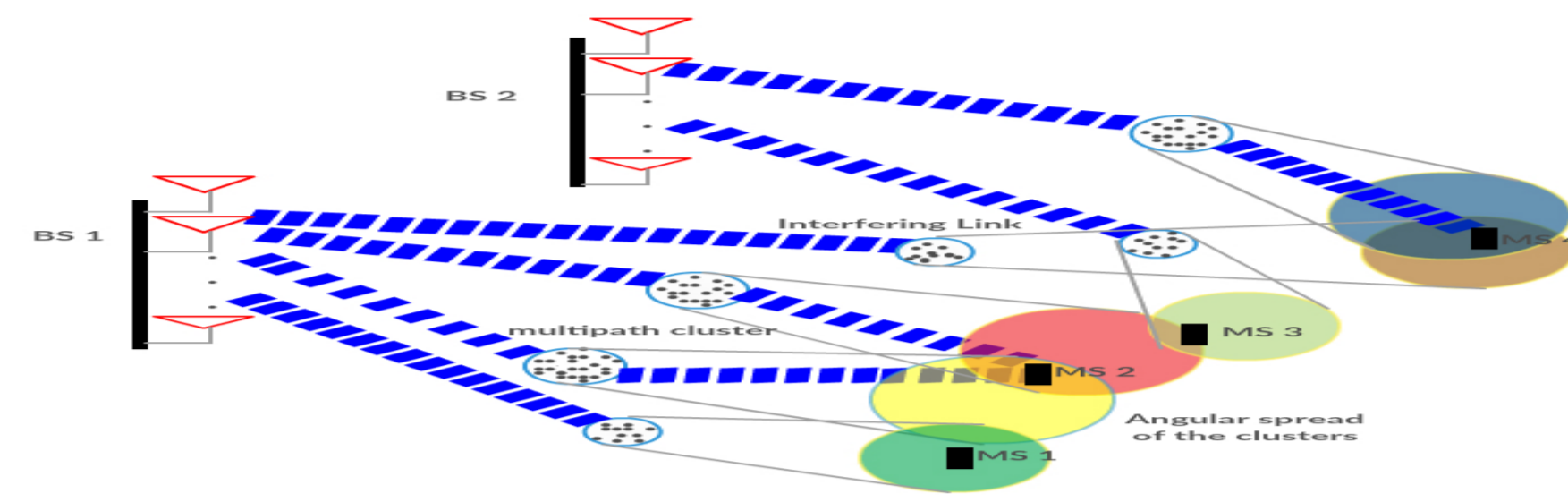
Application: Massive MIMO Channel Estimation

We get for the matrix impulse response of a time-varying frequency-selective MIMO channel $\mathbf{H}(t, \tau)$,

$$\mathbf{H}(t, \tau) = \sum_{i=1}^{N_p} A_i(t) e^{j2\pi f_i t} \mathbf{h}_r(\phi_i) \mathbf{h}_t^T(\psi_i) p(\tau - \tau_i)$$

with N_p (specular) pathwise contributions where

Massive MIMO Channel Estimation



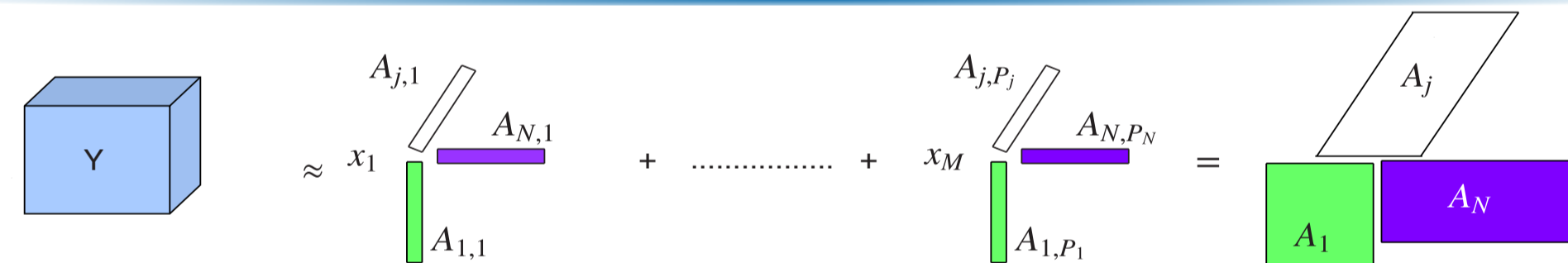
- ▶ A_i : complex attenuation, f_i : Doppler shift
- ▶ ψ_i : AoD (azimuth, elevation, polar), ϕ_i : AoA (azimuth, elevation, polarization)
- ▶ τ_i : path delay (ToA), $\mathbf{h}_t(\cdot)$, $\mathbf{h}_r(\cdot)$: $N_t/N_r \times 1$ Tx/Rx antenna array response, $p(\cdot)$: pulse shape (Tx filter)

The channel impulse response \mathbf{H} has per path a rank one contribution in four dimensions (Tx and Rx spatial multi-antenna dimensions, delay spread and Doppler spread). Hence, going to the frequency domain, we get

$$\text{vec}(\mathbf{H}(1:t, f_1:f_2)) = \sum_{i=1}^{N_p} A_i \mathbf{h}_t(\psi_i) \otimes \mathbf{h}_r(\phi_i) \otimes \mathbf{v}_f(\tau_i) \otimes \mathbf{v}_t(f_i)$$

where $\mathbf{v}_f(\cdot)$, $\mathbf{v}_t(\cdot)$ are appropriate Vandermonde vectors (possibly subsampled in the case of $\mathbf{v}_f(\cdot)$). Hence we get a sum of rank one $4D$ tensors. \mathbf{h}_r , \mathbf{h}_t : Kronecker structure in the case of polarization or in the case of $2D$ antenna arrays with separable structure [Sidiropoulos:icassp18].

System Model



Let Y_{i_1, \dots, i_N} represents the $i_1 i_2 \dots i_N$ th element of the tensor and $\mathbf{y} = [y_{1,1, \dots, 1}, y_{1,1, \dots, 2}, \dots, y_{1,1,2, \dots, I_N}]^T$, then it can be verified that [Sidiropoulos:TSP17],

$$\mathbf{y} = (\mathbf{A}_1 \otimes \mathbf{A}_2 \dots \otimes \mathbf{A}_N) \mathbf{x} + \mathbf{w}, \mathbf{w} \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$$

$$\text{Matrix Unfolding: } \mathbf{Y}^{(n)} = \mathbf{A}_n \mathbf{X}^{(n)} (\mathbf{A}_N \otimes \dots \otimes \mathbf{A}_{n+1} \otimes \mathbf{A}_{n-1} \dots \otimes \mathbf{A}_1)^T$$

$$q(\mathbf{x}, \boldsymbol{\alpha}, \gamma, \mathbf{A}) = q_\gamma(\gamma) \prod_{i=1}^M q_{x_i}(x_i) \prod_{i=1}^M q_{\alpha_i}(\alpha_i) \prod_{i=1}^M \prod_{j=1}^M q_{a_{j,i}}(\mathbf{a}_{j,i})$$

Variational Bayesian Inference

- ▶ VB compute the factors q by minimizing the Kullback-Leibler distance between the true posterior distribution $p(\mathbf{x}, \boldsymbol{\alpha}, \gamma, \mathbf{A}|\mathbf{y})$ and the $q(\mathbf{x}, \boldsymbol{\alpha}, \gamma, \mathbf{A})$.
- ▶ Equivalent to maximizing the evidence lower bound (ELBO) [Tzikas:SPMag08], $\boldsymbol{\theta} = \{\mathbf{x}, \boldsymbol{\alpha}, \gamma, \mathbf{A}\}$.

$$\ln p(\mathbf{y}) = L(q) + KLD_{VB}, \text{ where,}$$

$$L(q) = \int q(\boldsymbol{\theta}) \ln \frac{p(\mathbf{y}, \boldsymbol{\theta})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}, KLD_{VB} = - \int q(\boldsymbol{\theta}) \ln \frac{p(\boldsymbol{\theta}|\mathbf{y})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}$$

$$\ln(q_i(\theta_i)) = \langle \ln p(\mathbf{y}, \boldsymbol{\theta}) \rangle_{k \neq i} + c_i$$

SAVED-KS Equations

Joint Distribution:

$$\ln p(\mathbf{y}, \boldsymbol{\theta}) = N \ln \gamma - \gamma \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \sum_{i=1}^M (\ln \alpha_i - \alpha_i |x_i|^2) + \sum_{i=1}^M ((a-1) \ln \alpha_i + a \ln b - b \alpha_i) + (c-1) \ln \gamma + c \ln d - d \gamma + \text{constants}$$

Gaussian q for x_i or for $\mathbf{a}_{j,i}$ (Multivariate):

$$\sigma_i^2 = \frac{1}{\langle \gamma \rangle \prod_{j=1}^N \langle \|\mathbf{A}_{j,p_{ji}}\|^2 \rangle + \langle \alpha_i \rangle}, \mathbf{C}_i = \left(\bigotimes_{j=1}^N \mathbf{A}_{j,p_{ji}} \right)$$

$$\hat{x}_i = \sigma_i^2 (\langle \mathbf{C}_i^H \mathbf{y} \rangle - \langle (\mathbf{C}_i^H \mathbf{C}_i) \rangle \langle \mathbf{x}_i \rangle) / \langle \gamma \rangle$$

$$\hat{\mathbf{a}}_{j,i} = (\mathbf{b}_j)_T, \mathbf{b}_j = (\mathbf{Y}^{(j)} \langle \mathbf{X}^{(j)} \rangle \langle \bigotimes_{k=N, k \neq j}^N \mathbf{A}_k \rangle^T) / \langle \gamma \rangle$$

$$\Upsilon_{j,i} = \beta_{j,i} \mathbf{I}, \beta_{j,i} = \text{tr} \left\{ \left(\bigotimes_{k=N, k \neq j}^N \langle \mathbf{A}_k^T \mathbf{A}_k \rangle \right) \langle \mathbf{X}^{(j)H} \mathbf{X}^{(j)} \rangle \right\}$$

Gamma q for hyper-parameters

$$\langle \alpha_i \rangle = \frac{a+\frac{1}{2}}{\langle |x_i|^2 \rangle + b}, \text{ where } \langle |x_i|^2 \rangle = |\hat{x}_i|^2 + \sigma_i^2$$

$$\langle \gamma \rangle = \frac{c+\frac{N}{2}}{\langle \|\mathbf{y} - \left(\bigotimes_{j=1}^N \mathbf{A}_j \right) \mathbf{x}\|^2 \rangle + d}$$

Joint VB for KS Matrices

Complex Matrix Normal Distribution

$$\mathbf{M}_j = \hat{\mathbf{A}}_j = \langle \gamma \rangle \mathbf{B}_j \boldsymbol{\Psi}_j$$

$$\boldsymbol{\Psi}_j = \left(\langle \gamma \rangle \mathbf{X} \bigotimes_{k=1, k \neq j}^N \langle \mathbf{A}_k^T \mathbf{A}_k \rangle \mathbf{X}^H \right)^{-1}, \mathbf{X} = \text{diag}(\mathbf{x})$$

\mathbf{B}_j is with the first row of $(\mathbf{Y}^{(j)} \langle \bigotimes_{k=1, k \neq j}^N \mathbf{A}_k \rangle^* \rangle \langle \mathbf{X}^H \rangle$ removed.

SAVED-KS SBL Algorithm

Given: $\mathbf{y}, \mathbf{A}, M, N$.

Initialization: a, b, c, d are taken to be very low, on the order of 10^{-10} . $\alpha_i^0 = a/b, \forall i, \gamma^0 = c/d$ and $\sigma_i^{2,0} = \frac{1}{\|\mathbf{A}_i\|^2 \gamma^0 + \alpha_i^0}, \mathbf{x}^0 = \mathbf{0}$.

At iteration $t+1$,

- ▶ Update $\sigma_i^{2,t+1}, \hat{x}_i^{t+1}, \forall i$ from using \mathbf{x}_{i-}^{t+1} and \mathbf{x}_{i+}^t .
- ▶ Update $\hat{\mathbf{A}}_{j,i}, \forall i, j$ or $\mathbf{A}_j, \forall j$.
- ▶ Compute $\langle x_i^{2,t+1} \rangle$ and update α_i^t .
- ▶ Update the noise variance, γ^{t+1} .
- ▶ Continue steps 1 – 4 till convergence of the algorithm.

Comparison to SotA

- ▶ **Lowering Complexity:** No matrix inversions compared to standard SBL and ALS.
- ▶ **Improving Convergence** compared to standard ALS.

Identifiability

- ▶ The local identifiability (upto permutation ambiguity) of the KS DL is ensured if the FIM is non-singular.

$$\mathbf{J}(\boldsymbol{\theta}, \mathbf{x}) = [\mathbf{J}(\boldsymbol{\theta}) \mathbf{J}(\mathbf{x})], \mathbf{J}(\boldsymbol{\theta}) = [\mathbf{J}(\boldsymbol{\theta}_1) \dots \mathbf{J}(\boldsymbol{\theta}_N)]$$

where, $\mathbf{J}(\boldsymbol{\theta}_j) = \mathbf{F}(\mathbf{x})(\boldsymbol{\theta}_1 \otimes \dots \otimes \mathbf{I}_{I_j P_j} \dots \otimes \boldsymbol{\theta}_N)$,

$$\mathbf{J}(\mathbf{x}) = [\mathbf{F}_1(\bigotimes_{j=1}^N \boldsymbol{\theta}_j), \dots, \mathbf{F}_M(\bigotimes_{j=1}^N \boldsymbol{\theta}_j)]$$

$$\text{FIM} = \begin{bmatrix} E(\gamma) \mathbf{J}(\boldsymbol{\theta})^H \mathbf{J}(\boldsymbol{\theta}) & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & E(\gamma) \mathbf{J}(\mathbf{x})^H \mathbf{J}(\mathbf{x}) + E(\boldsymbol{\Gamma}) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & a E(\boldsymbol{\Gamma}^{-2}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & N' E(\gamma^{-2}) \end{bmatrix}$$

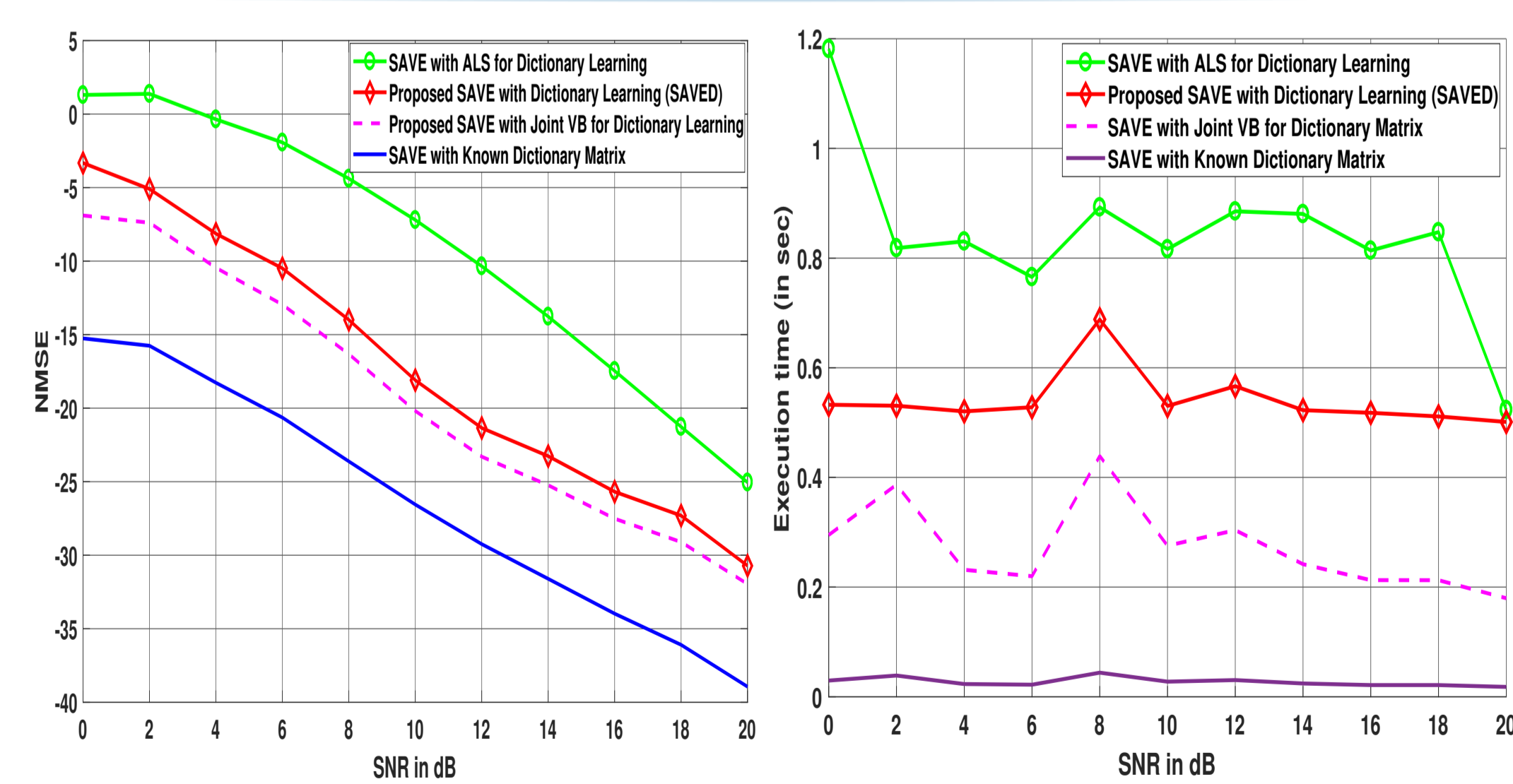
- ▶ For the FIM analysis (with known support of \mathbf{x}), then $E(\gamma) \mathbf{J}(\mathbf{x})^H \mathbf{J}(\mathbf{x}) + E(\boldsymbol{\Gamma})$ and $a E(\boldsymbol{\Gamma}^{-2})$ becomes invertible if $\prod_{j=1}^N I_j > K$. Assuming

$$\prod_{j=1}^N I_j > \sum_{j=1}^N (I_j - 1) P_j, \text{ i.e. no. of degrees of freedom in the dictionary } \langle \prod_{j=1}^N I_j, \text{ FIM is non-singular.}$$

- ▶ Possibility of FIM singularity even under single measurement vector case.
- ▶ Mixture of $P (< N)$ Vandermonde matrix factors and non-parametric KS factors: The identifiability conditions can be restated as,

$$\prod_{j=1}^N I_j > \sum_{j=1}^P P_j + \sum_{j=P+1}^N (I_j - 1) P_j$$

Numerical Results



$I_1 = 4, I_2 = 10, I_3 = 4$. Convergence behaviour (20 dB).

Conclusion and Future Work

- ▶ Convex combination of structured and unstructured KS factor matrices, For eg, DoA response closeness to the vandermonde.
- ▶ Asymptotic performance analysis, mismatched CRBs.