

EE 4715 Array Processing

10. Factor Analysis

June 2022

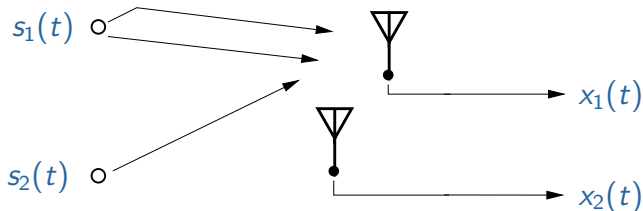
Introduction: Subspace estimation

Classical data model

Consider Q sources $\mathbf{s}[n]$, a $P \times Q$ mixing matrix \mathbf{A} , and white noise $\mathbf{n}[n]$:

$$\mathbf{x}[n] = \mathbf{A}\mathbf{s}[n] + \mathbf{n}[n]$$

$$\mathbf{R} := E[\mathbf{x}\mathbf{x}^H] = \mathbf{A}\Sigma_s\mathbf{A}^H + \sigma^2\mathbf{I}$$



Introduction: Subspace estimation

$$R = A\Sigma_s A^H + \sigma^2 I$$

Eigenvalue decompositions are used to estimate the column span of A :

- **No noise:**

$$R =: U\Lambda U^H = \begin{bmatrix} U_0 & U_1 \end{bmatrix} \begin{bmatrix} \Lambda_0 & \\ & 0 \end{bmatrix} \begin{bmatrix} U_0^H \\ U_1^H \end{bmatrix}$$

- **With additive white noise:**

$$R = \begin{bmatrix} U_0 & U_1 \end{bmatrix} \begin{bmatrix} \Lambda_0 + \sigma^2 I & \\ & \sigma^2 I \end{bmatrix} \begin{bmatrix} U_0^H \\ U_1^H \end{bmatrix}$$

In both cases, U_0 is a basis for the column span of A . This is used in subspace-based techniques, e.g. in MUSIC.

Subspace estimation

- **With spatially colored noise:** $R_n = D$ (some diagonal)

$$R = A\Sigma_s A^H + D$$

The eigenvalue decomposition does not give the correct subspace. We need to do *prewhitening*: work with $R_n^{-1/2}x[n]$.

$$\begin{aligned} R_n^{-1/2} R R_n^{-1/2} &= (R_n^{-1/2} A) \Sigma_s (A^H R_n^{-1/2}) + I \\ &= \begin{bmatrix} U_0 & U_1 \end{bmatrix} \begin{bmatrix} \Lambda_0 + I & \\ & I \end{bmatrix} \begin{bmatrix} U_0^H \\ U_1^H \end{bmatrix} \end{aligned}$$

After prewhitening, U_0 is a basis for the column span of $R_n^{-1/2} A$.

This can be done if R_n is known (requires calibration). What if it isn't?

Factor Analysis

Data model

We work with the model

$$R = AA^H + D \quad [\text{low rank plus diagonal}]$$

The goal in FA is, given R , estimate A and D .

FA is an old problem (1920?), usually for real factors. It is used in many domains where data has arbitrary scalings.

Uniqueness

- A is unique only up to a rotation Q : if A fits the model, then also AQ . This is solved by placing constraints on A .

The subspace $\text{ran}(A)$ is invariant anyway.

Factor Analysis

If $Q \geq P$ then the problem is not identifiable: we can set $D = 0$. What is the maximum Q (number of sources?)

Counting number of equations and number of unknowns:

- For a complex R , we have P^2 real-valued “observations”
- Number of real-valued unknowns in A is $2PQ$, in D is P
- Need Q^2 constraints to make A unique (e.g., constraining $A^H D^{-1} A$ to be diagonal and the first row of A real)

Total degrees of freedom:

$$s = P^2 - (2PQ + P - Q^2) = (P - Q)^2 - P$$

Identifiability requires $s > 0$ (more equations than unknowns)

$$s > 0 \Leftrightarrow Q < P - \sqrt{P}$$

FA algorithms

Ad hoc algorithm

In practice we do not have \mathbf{R} , but an estimate $\hat{\mathbf{R}}$.

We can formulate Factor Analysis as a covariance fitting problem (least squares):

$$\min_{\mathbf{A}, \mathbf{D}} \|\hat{\mathbf{R}} - \mathbf{A}\mathbf{A}^H - \mathbf{D}\|_F^2$$

Solve using Alternating Least Squares: estimate \mathbf{A} , then estimate \mathbf{D} , etc.

This converges *very* slowly.

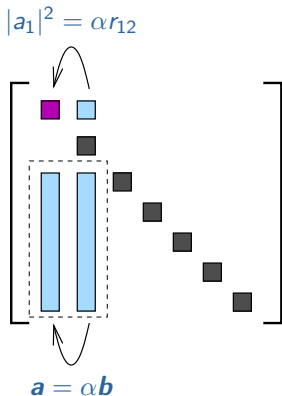
FA algorithms

Closed-form approximate solution

We can also view this as a matrix extension problem: replace the main diagonal of \hat{R} such that it becomes low rank ($= \mathbf{A}\mathbf{A}^H$).

- Rank-1 factor model:

$$R' = R - D = \mathbf{a}\mathbf{a}^H =$$



Property: each submatrix not involving the main diagonal is rank-1.
The ratios can be used to estimate the diagonal entries of $\mathbf{a}\mathbf{a}^H$.

Intermezzo

Recap: the Kronecker product

For a matrix, $\text{vec}(\cdot)$ denotes the stacking of the columns of a matrix into a vector.

For two matrices \mathbf{A} and \mathbf{B} , the Kronecker product is defined as

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1N}\mathbf{B} \\ \vdots & & \vdots \\ a_{M1}\mathbf{B} & \cdots & a_{MN}\mathbf{B} \end{bmatrix},$$

Some properties:

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) &= \mathbf{AC} \otimes \mathbf{BD} \\ \text{vec}(\mathbf{ab}^H) &= \mathbf{b}^* \otimes \mathbf{a} \\ \text{vec}(\mathbf{ABC}) &= (\mathbf{C}^T \otimes \mathbf{A})\text{vec}(\mathbf{B}) \end{aligned}$$

Intermezzo

The variance of a sample covariance matrix

Consider an i.i.d. zero mean random sequence \mathbf{x}_n .

Let $\mathbf{R} = E[\mathbf{x}_n \mathbf{x}_n^H]$. We estimate \mathbf{R} from N data samples as

$$\hat{\mathbf{R}} = \frac{1}{N} \sum_{n=0}^{N-1} \mathbf{x}_n \mathbf{x}_n^H$$

Then $E[\hat{\mathbf{R}}] = \mathbf{R}$: unbiased. But what is the variance of this estimate?

To define this, consider $\mathbf{r} = \text{vec}(\mathbf{R})$, $\hat{\mathbf{r}} = \text{vec}(\hat{\mathbf{R}})$. Then

$$\hat{\mathbf{r}} = \mathbf{r} + \mathbf{e}$$

where $\hat{\mathbf{r}}$ is a random variable, \mathbf{r} is its expected value (mean), and \mathbf{e} zero mean “finite sample noise”. We define

$$\text{cov}[\hat{\mathbf{R}}] =: \text{cov}[\hat{\mathbf{r}}] = E[\mathbf{e} \mathbf{e}^H]$$

Intermezzo

The variance of a sample covariance matrix

Recall

$$\hat{\mathbf{R}} = \frac{1}{N} \sum \mathbf{x}[n]\mathbf{x}[n]^H \Rightarrow \hat{\mathbf{r}} = \frac{1}{N} \sum \mathbf{x}_n^* \otimes \mathbf{x}_n$$

Use the fact that \mathbf{x}_i is independent of \mathbf{x}_j for $i \neq j$ to derive

$$\begin{aligned} \text{cov}[\hat{\mathbf{R}}] &= \mathbb{E} \left[\left(\frac{1}{N} \sum (\mathbf{x}_i^* \otimes \mathbf{x}_i) - \mathbb{E}[\mathbf{x}_i^* \otimes \mathbf{x}_i] \right) \left(\frac{1}{N} \sum (\mathbf{x}_j^* \otimes \mathbf{x}_j) - \mathbb{E}[\mathbf{x}_j^* \otimes \mathbf{x}_j] \right)^H \right] \\ &= \frac{1}{N^2} \sum \sum \mathbb{E} [(\mathbf{x}_i^* \otimes \mathbf{x}_i - \mathbb{E}[\mathbf{x}_i^* \otimes \mathbf{x}_i])(\mathbf{x}_j^* \otimes \mathbf{x}_j - \mathbb{E}[\mathbf{x}_j^* \otimes \mathbf{x}_j])^H] \\ &= \frac{1}{N^2} \sum \mathbb{E} [(\mathbf{x}_i^* \otimes \mathbf{x}_i - \mathbb{E}[\mathbf{x}_i^* \otimes \mathbf{x}_i])(\mathbf{x}_i^* \otimes \mathbf{x}_i - \mathbb{E}[\mathbf{x}_i^* \otimes \mathbf{x}_i])^H] \\ &= \frac{1}{N} (\mathbb{E}[(\mathbf{x}^* \otimes \mathbf{x})(\mathbf{x}^* \otimes \mathbf{x})^H] - \mathbb{E}[\mathbf{x}^* \otimes \mathbf{x}]\mathbb{E}[\mathbf{x}^* \otimes \mathbf{x}]^H) \\ &= \frac{1}{N} \mathbf{C} \end{aligned}$$

where $\mathbf{C} = \mathbb{E}[(\mathbf{x}_k^* \otimes \mathbf{x}_k)(\mathbf{x}_k^* \otimes \mathbf{x}_k)^H] - \mathbb{E}[\mathbf{x}_k^* \otimes \mathbf{x}_k]\mathbb{E}[\mathbf{x}_k^* \otimes \mathbf{x}_k]^H$

Intermezzo

The variance of a sample covariance matrix

The covariance of $\hat{\mathbf{R}}$ involves fourth-order correlations. These can often be described in simpler terms using cumulants.

For the special case where the entries x_i of \mathbf{x} are zero-mean and **jointly Gaussian** distributed, it is known that (for arbitrary indices $a, b, c, d = 0, \dots, M - 1$)

$$E[x_a x_b^* x_c x_d^*] = E[x_a x_b^*]E[x_c x_d^*] + E[x_a x_d^*]E[x_b^* x_c] + E[x_a x_c]E[x_b^* x_d^*]$$

“Proper” (or circularly symmetric) complex variables are such that $E[\mathbf{x}\mathbf{x}^T] = 0$. In this case, the last term vanishes.

Thus, for Gaussians, higher-order moments reduce to a function of second-order moments.

Intermezzo

The variance of a sample covariance matrix

Stacking in a matrix with row-index $a + M b$ and column-index $c + M d$, we can write this expression compactly as

$$E[(\mathbf{x}^* \otimes \mathbf{x})(\mathbf{x}^* \otimes \mathbf{x})^H] = E[\mathbf{x}^* \otimes \mathbf{x}]E[\mathbf{x}^* \otimes \mathbf{x}]^H + E[\mathbf{x}^* \mathbf{x}^{*H}] \otimes E[\mathbf{x}\mathbf{x}^H] \\ + E[(\mathbf{x}^* \otimes \mathbf{1})(\mathbf{1} \otimes \mathbf{x})^H] \odot E[(\mathbf{1} \otimes \mathbf{x})(\mathbf{x}^* \otimes \mathbf{1})^H]$$

For proper complex variables, the last term vanishes, and thus, for **zero mean proper complex-valued Gaussian random variables**,

$$\text{cov}[\hat{\mathbf{R}}] = \frac{1}{N} \mathbf{R}^* \otimes \mathbf{R}$$

FA algorithms

Parametrization

We parametrize $\mathbf{R} = \mathbf{R}(\boldsymbol{\theta}) = \mathbf{A}\mathbf{A}^H + \mathbf{D}$ using parameters

$$\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\theta}_A \\ \boldsymbol{\theta}_{A^*} \\ \boldsymbol{\theta}_D \end{bmatrix} = \begin{bmatrix} \text{vec}(\mathbf{A}) \\ \text{vec}(\mathbf{A}^*) \\ \text{diag}(\mathbf{D}) \end{bmatrix}$$

These are redundant, but the constraints can be implemented later (after estimating $\boldsymbol{\theta}$).

The parameters are **complex** and we use **Wirtinger calculus** to differentiate. Essentially: treat $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^*$ as independent.

The Jacobian $\mathbf{J}(\boldsymbol{\theta})$ for $\mathbf{R}(\boldsymbol{\theta})$ is given by

$$\mathbf{J} = \frac{\partial \text{vec}(\mathbf{R})}{\partial \boldsymbol{\theta}^T} = [\mathbf{J}_A, \quad \mathbf{J}_{A^*}, \quad \mathbf{J}_D] = [\mathbf{A}^* \otimes \mathbf{I}, \quad (\mathbf{I} \otimes \mathbf{A})\mathbf{K}, \quad \mathbf{I} \circ \mathbf{I}]$$

where \circ denotes the Khatri-Rao product (column-wise Kronecker product), and $\text{vec}(\mathbf{A}^T) = \mathbf{K}\text{vec}(\mathbf{A})$.

FA algorithms

Non-Linear Weighted Least Squares

Define $\mathbf{r} = \text{vec}(\mathbf{R})$, then we can formulate a covariance matching problem

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \|\mathbf{W}^{1/2}[\hat{\mathbf{r}} - \mathbf{r}(\boldsymbol{\theta})]\|^2 = \arg \min_{\boldsymbol{\theta}} [\hat{\mathbf{r}} - \mathbf{r}(\boldsymbol{\theta})]^H \mathbf{W} [\hat{\mathbf{r}} - \mathbf{r}(\boldsymbol{\theta})]$$

where \mathbf{W} is a weighting matrix. The optimum weight is the inverse of the covariance matrix of $\hat{\mathbf{r}}$, for proper Gaussian-distributed data estimated as

$$\mathbf{W} = \hat{\mathbf{R}}^{-T} \otimes \hat{\mathbf{R}}^{-1}$$

NOte that we use $\hat{\mathbf{R}}$ instead of \mathbf{R} . Asymptotically (for large N), this converges to the ML solution.

FA algorithms

Gauss-Newton for solving Non-Linear WLS

The Gauss-Newton technique for solving nonlinear WLS is the iteration

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + \mu^{(k)} \boldsymbol{\delta}$$

where the direction of descent $\boldsymbol{\delta}$ is found by solving

$$\mathbf{B}(\boldsymbol{\theta}^{(k)}) \boldsymbol{\delta} = \mathbf{g}(\boldsymbol{\theta}^{(k)})$$

$\mathbf{g}(\boldsymbol{\theta})$ is the gradient of the cost function, and $\mathbf{B}(\boldsymbol{\theta})$ is the Gramian of the Jacobians (approximating the Hessian):

$$\mathbf{B}(\boldsymbol{\theta}) = \mathbf{J}^H(\boldsymbol{\theta}) \mathbf{W} \mathbf{J}(\boldsymbol{\theta}), \quad \mathbf{g}(\boldsymbol{\theta}) = \mathbf{J}^H(\boldsymbol{\theta}) \mathbf{W} [\hat{\mathbf{r}} - \mathbf{r}(\boldsymbol{\theta})]$$

FA algorithms

Closed-form solution for direction of descent

$B(\theta)$ can be very large: hard to construct and invert.

A complicated derivation shows how we can solve for δ_D inside δ in closed form. Define

$$\begin{aligned}\tilde{W} &= \hat{R}^{-1} - \hat{R}^{-1}A(A^H\hat{R}^{-1}A)^{-1}A^H\hat{R}^{-1} & [\tilde{W}A = 0] \\ \tilde{B}_D &= J_D^H \left(\tilde{W}^T \otimes \tilde{W} \right) J_D = \tilde{W}^T \odot \tilde{W} & [J_D = I \circ I] \\ \tilde{g}_D &= J_D^H \left(\tilde{W}^T \otimes \tilde{W} \right) \text{vec}[\hat{R} - R(\theta)]\end{aligned}$$

Then the computation of

$$\delta = \begin{bmatrix} \text{vec}(\Delta_A) \\ \text{vec}(\Delta_{A^*}) \\ \delta_D \end{bmatrix}$$

reduces to the computation of δ_D from $\tilde{B}_D\delta_D = \tilde{g}_D$, while

$$\Delta_A = \frac{1}{2}(I + \hat{R}\tilde{W})(\hat{R} - R(\theta) - \text{diag}(\delta_D))\hat{R}^{-1}A(A^H\hat{R}^{-1}A)^{-1}$$

FA algorithms

Alternating WLS method

If we take step size $\mu = 1$, then the closed-form result simplifies.

[...]

The result is the Alternating Weighted Least Squares (AWLS) algorithm (cf the Ad-Hoc algorithm), where $\mathbf{W} = \hat{\mathbf{R}}^{-1}$:

$$\mathbf{U}\mathbf{\Lambda}\mathbf{U}^H := \mathbf{D}_{(k)}^{-1/2} \hat{\mathbf{R}} \mathbf{D}_{(k)}^{-1/2} \quad [\text{EVD; prewhitening of } \hat{\mathbf{R}}]$$

$$\mathbf{A}_{(k+1)} := \mathbf{D}_{(k)}^{1/2} \mathbf{U}_0 (\mathbf{\Lambda}_0 - \mathbf{I})^{1/2}$$

$$\tilde{\mathbf{W}} := \mathbf{W} - \mathbf{W} \mathbf{A}_{(k+1)} (\mathbf{A}_{(k+1)}^H \mathbf{W} \mathbf{A}_{(k+1)})^{-1} \mathbf{A}_{(k+1)}^H \mathbf{W}$$

$$\mathbf{d}_{(k+1)} := \left[\tilde{\mathbf{W}}^T \odot \tilde{\mathbf{W}} \right]^{-1} \text{vecdiag}(\tilde{\mathbf{W}})$$

$$\mathbf{D}_{(k+1)} := \text{diag}(\mathbf{d}_{(k+1)})$$

Extensions

Extended FA

Let \mathbf{M} with $m_{ij} \in \{0, 1\}$ be a “masking matrix”. Then

$$\mathbf{R} = \mathbf{A}\mathbf{A}^H + \mathbf{\Psi}, \quad \mathbf{\Psi} = \mathbf{M} \odot \mathbf{\Psi}$$

This generalizes \mathbf{D} to a more general structure.

This can be further generalized to linear models of the form $\text{vec}(\mathbf{\Psi}) = \mathbf{G}\boldsymbol{\theta}_\psi$, where \mathbf{G} is a fixed basis (e.g. selected columns of a Fourier matrix to model spatially lowpass noise).

Joint FA

Suppose we have multiple “snapshots” \mathbf{R}_m , with different \mathbf{A}_m but a common noise covariance matrix,

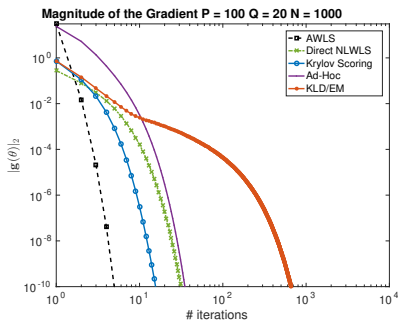
$$\mathbf{R}_m = \mathbf{A}_m\mathbf{A}_m^H + \mathbf{D}, \quad m = 1, \dots, M$$

Simulations

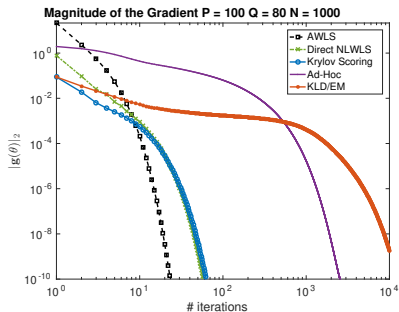
Convergence

Convergence speed for $P = 100$ sensors, $N = 1000$ samples

$Q = 20$ sources:



$Q = 80$ sources:

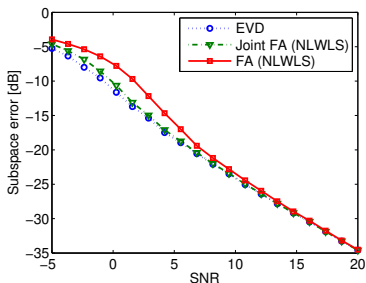


Simulations

Subspace Estimation Performance

Subspace error as function of SNR for $\Psi = I$ (white noise).

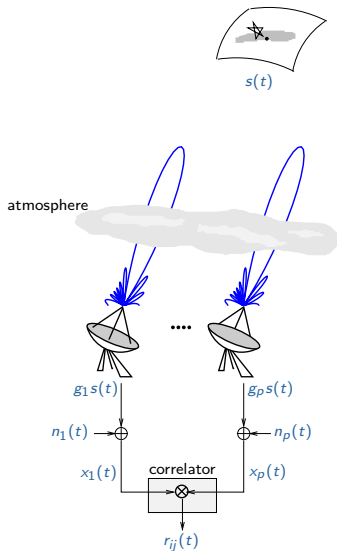
$$Q_m = 2, P = 5, M = 5, N = 100$$



For low SNR, there is some drop in performance in particular for non-joint processing. For higher SNR, there is no performance penalty.

Applications - 1 (FA)

Gain and noise power calibration at the Westerbork Radio Telescope

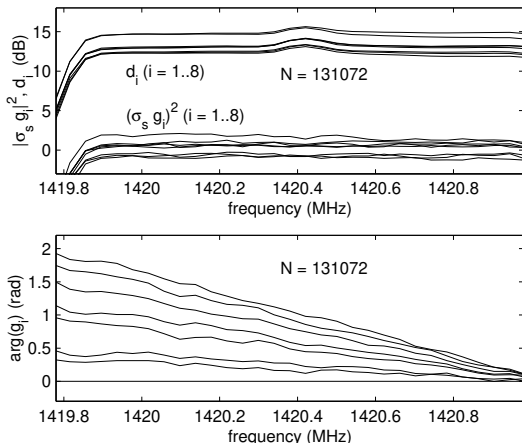


- General calibration strategy: observe a single 'strong' sky source
- Stack the p antenna signals in a vector \mathbf{x} and analyze $\mathbf{R} = \mathbb{E}[\mathbf{x}\mathbf{x}^H]$ with model

$$\mathbf{R} = \mathbf{g}\sigma_s^2\mathbf{g}^H + \mathbf{D}$$

- Need to estimate the complex gains g_i and noise powers in \mathbf{D} .

Applications - 1 (FA)



Estimates of (a) gain magnitude, noise powers, and (b) gain phase, as function of frequency, around $f = 1420$ MHz. Based on observation of the astronomical source 3C48. Frequencies processed independently.

Applications - 2 (JFA)

Narrowband interference subspace estimation at the Westerbork Radio Telescope

For each frequency f , we have a sequence of short-term correlation matrices $\hat{\mathbf{R}}_{k,f}$ (averaged to 10 ms) with model

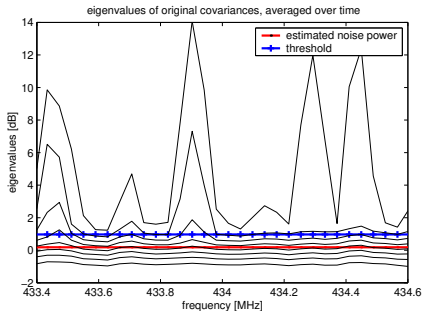
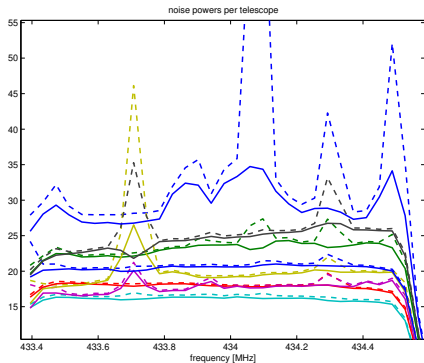
$$\mathbf{R}_{k,f} = \mathbf{A}_{k,f} \mathbf{A}_{k,f}^H + \mathbf{D}_f, \quad k = 1, 2, \dots, 1000$$

- $\mathbf{A}_{k,f}$: interference subspace at time k , frequency f (32 frequency bins)
- \mathbf{D}_f : noise power estimate at a single frequency – to be calibrated

For each frequency, we use JFA to estimate the factors and noise powers.

Applications - 2 (JFA)

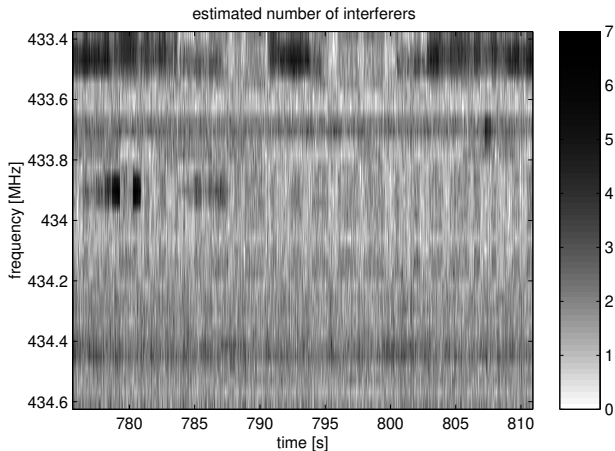
Amateur AM radio broadcast interference at 434 MHz, both continuous and intermittent, recorded at the WSRT pointed at 3C48:



(a) estimated noise powers, (b) averaged eigenvalues (after whitening)

Applications - 2 (JFA)

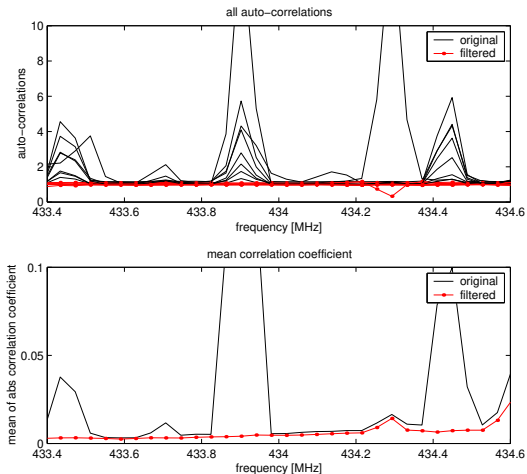
Estimated number of interferers



Applications - 2 (JFA)

Correlation spectra after interference cancellation

The interference subspaces are projected out, and the results are further averaged over time.



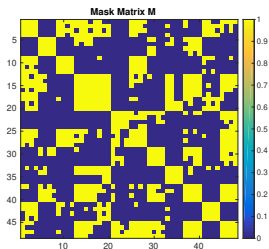
Applications - 3 (EFA)

Suppression of the Milky Way at a LOFAR station

Using a LOFAR LBA station, we observe two strong sources, Cas A and Cyg A, and the Milky Way (cloud-like emission).

The Milky way is modeled by R_n with mask M incorporating all baselines shorter than 4 wavelengths.

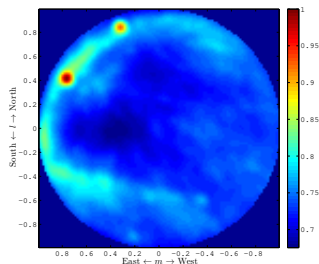
After estimating R_n using Extended Factor Analysis, the residual $\hat{R}_0 = \hat{R} - R_n$ is imaged.



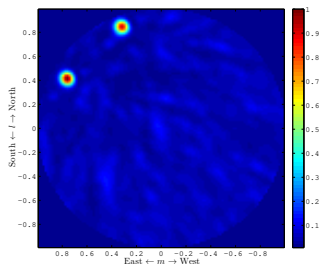
Applications - 3 (EFA)

Total sky image for 1 station using classical DFT beamforming (“dirty image”) combining data from 24 156kHz subbands distributed between 45.3 and 67.3 MHz and 10 seconds of integration per channel.

DFT Imaging without FA



DFT Imaging with Extended FA

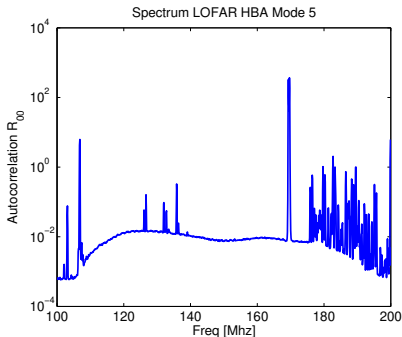


Next, it is possible to calibrate the gains of the individual antennas using the (known) intensities of the point sources.

Applications - 4

Spatial Filtering at a LOFAR station

Data from the LOFAR station RS409 in HBA mode 5 (100-200 MHz), tracking the strong astronomical source Cyg A.



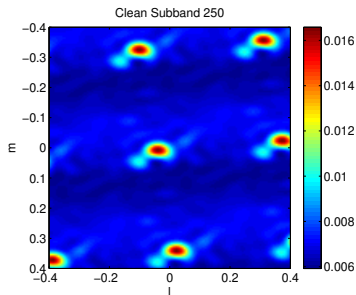
Above 174 MHz, the spectrum is heavily contaminated by wideband DAB transmissions.

Applications - 4

Snapshot image of the sky

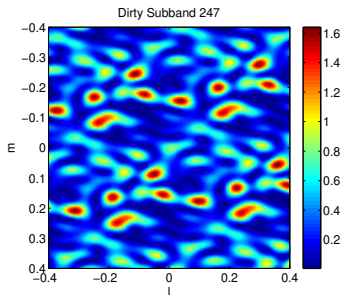
Uncontaminated image

subband 250 at 175.59 MHz



RFI-contaminated data

subband 247 at 175.88 MHz



Applications - 4

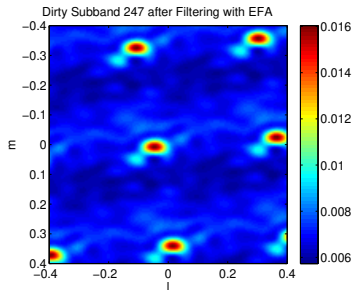
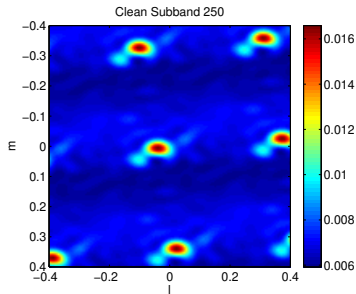
For spatial filtering, we take 6 of the 46 receiving elements as reference array. The covariance model is

$$\mathbf{R}_k = \mathbf{A}_k \mathbf{A}_k^H + \Psi = \mathbf{A}_k \mathbf{A}_k^H + \left[\begin{array}{c|c} \Psi_{00} & 0 \\ \hline 0 & \Sigma_1 \end{array} \right], \quad \mathbf{M} = \left[\begin{array}{c|c} \mathbf{1}\mathbf{1}^T & 0 \\ \hline 0 & \mathbf{I} \end{array} \right]$$

Ψ_{00} contains the astronomical covariances of interest. The mask \mathbf{M} is used to avoid a subspace model for these. \mathbf{A}_k is the subspace of the interference.

Applications - 4

EFA is used to estimate the subspace of the interference. This is then used as spatial filter on the 40 remaining antennas (primary array).



Summary

Factor Analysis is viewed as an extension of the eigenvalue decomposition for nonwhite noise.

We proposed:

- New algorithms based on Gauss-Newton iterations
- Extensions of FA to multiple matrices and more general noise models

The algorithms show reliable and efficient convergence, feasible for moderately large problem sizes ($P = 100$ sensors).

Even if the noise is white, the performance penalty with respect to EVD is minor.