

# *Nonparametric* Independent Component Analysis for Nonstationary Mixtures

Nima Reyhani<sup>1</sup> and Peter Bickel<sup>2</sup>

<sup>1</sup>Helsinki University of Technology

<sup>2</sup>UC Berkeley, Department of Statistics

September 7, 2009

# Independent Component Analysis-set up

- ▶ ICA assumption
  - ▶  $S = (s_1, \dots, s_m)^T$ , where  $s_i$  are independent r.v.s
  - ▶ assume model

$$X = AS, \quad A \in \mathbb{R}^{m \times m} \quad X, S \in \mathbb{R}^m$$

for  $X$  observable and  $S$  latent r.v.s.

- ▶ Question: Given observations  $X_1, \dots, X_n$  estimate the mixing matrix  $A$

Identifiability condition: (for iid) At most one of sources are Normal.

Algorithms:

- ▶ fastICA: Non-Gaussianity
- ▶ Samarov and Tsybakov: Empirical Fisher Information-like matrix.

# ICA with non-iid sources

- ▶ Assume  $s_i(\cdot), i = 1, \dots, m$ , is a time series and  $s_i \perp s_j, i \neq j$ , i.e. temporal structure.
- ▶ Let's define time-lagged covariance matrices (they are the same in iid case)

$$C_i := \mathbb{E}\{X(t)X(t + \tau_i)\}$$

$$\begin{aligned} X(t) &= AS(t) \Rightarrow \mathbb{E}\{X(t)X(t + \tau_i)^T\} \\ &= AS(t)S(t + \tau_i)^T A^T = A\Lambda(\tau_i)A^T \\ [\Lambda(\tau_i)]_{j,j} &= [\Lambda_i]_{j,j} = \mathbb{E}\{s_j(t)s_j(t + \tau)\}. \end{aligned}$$

- ▶ in stationary case,  $\|\Lambda_i - \Lambda_j\|_F \neq 0$  for some  $i, j \Rightarrow$  contrast
- ▶ Estimating  $A$  is equivalent to find a matrix simultaneously diagonalizes  $C_i$  and  $C_j$

# Block Stationary ICA

Block Stationary:

take  $[0, T] = [0, t_1) \cup \dots \cup [t_L, T]$  and assume  $s_i(t) \forall i$  stationary time series within each interval.

Sources of contrast:

e.g. **temporal fixed structure within each block** and different statistical properties between blocks. Also, we usually assume **non-Gaussianity**.

# Block-Stationary ICA

- ▶ Solution proposed by Pham & Cardoso:
- ▶ log-likelihood of source vector:

$$\begin{aligned}C_{ML} &= -\frac{1}{2} \sum \frac{s_i^2(t)}{\sigma_i^2(t)} + \log(2\pi\sigma_i^2(t)) = \\ &= -\frac{1}{2} \text{tr}[\Sigma^{-2}(t)S(t)S(t)^T] - \frac{1}{2} \log \det[2\pi\Sigma^2(t)] = \\ &= -\frac{1}{2} \text{tr}[\Sigma^{-2}(t)A^{-1}X(t)X(t)^{-T}A^{-T}] - \frac{1}{2} \log \det[2\pi\Sigma^2(t)] \\ \Sigma^2(t) &= \text{diag}[\sigma_1^2(t), \dots, \sigma_K^2(t)]\end{aligned}$$

assume  $\sigma_i^2(t) = \sigma_{i,l}^2$  and set  $\Sigma_l^2 = \text{diag}(\sigma_{1,l}^2, \dots, \sigma_{K,l}^2)$

- ▶ With some calculations:

$$C_{ML} = \frac{1}{2} \sum_{l=1}^L w_l [\text{tr}(R_l^{-1} \hat{R}_l) - \log \det(R_l^{-1} \hat{R}_l)] + K, \quad R_l = A \Sigma_l^2 A^T$$

$$C_{ML} = \sum_{l=1}^L \text{off}(A^{-1} \hat{R}_l A^{-T}), \quad \hat{R}_l = \frac{1}{\#T_l} \sum X_l(t) X_l(t)^T$$

- ▶ *off* measures deviation from diagonality

# Proposed Algorithm

Let's consider

$$M^l(\gamma_l) := \mathbb{E}_p\{\gamma_l(X_l)\nabla p(X_l)\nabla^T p(X_l)\} = \sum_{i,j=1}^m M_{i,j}^l \beta_i \beta_j^T$$

$$M_{ij}^l = \det(B)^2 \mathbb{E} \left\{ \gamma_l(X_l) \prod_{m \neq i} p_m(X_l^T \beta_m) \prod_{k \neq j} p_k(X_l^T \beta_k) p'_i(X_l^T \beta_i) p'_j(X_l^T \beta_j) \right\}$$

$\beta_i$  is the  $i$ -th row vector of  $B = A^{-1}$ ,  $u_i^l := X_l^T \beta_i$

$$M_{ij}^l = m_{ij}^l \int \gamma_l(X) p_i^2(u_i^l) p_j^2(u_j^l) p'_i(u_i^l) p'_j(u_j^l) du_i du_j$$

$$\int p_j^2(u) p'_j(u) du = 0, \forall j, \quad M_{ij}^l = 0, i \neq j \quad M^l = B \text{diag}(M_{ii}^l) B^T \Rightarrow$$

Replace the covariance matrix with the Fisher-like information matrix.

# Empirical estimates

Assume  $\gamma \equiv 1$ ,

$$\hat{M}' = \frac{1}{\#T_l} \sum_{i \in T_l} \nabla \hat{p}_{-i}(X_i^l) \nabla^T \hat{p}_{-i}(X_i^l)$$

$$\frac{\partial \hat{p}_{-i}}{\partial x_s}(X_i^l) = \frac{1}{(\#T_l - 1)h^{m+1}} \sum_{j=1, j \neq i} K_1 \left( \frac{X_{js}^l - X_{is}^l}{h} \right) \prod_{k=1, k \neq s} K \left( \frac{X_{jk}^l - X_{ik}^l}{h} \right),$$

# Joint Diagonalization

- ▶ Suppose matrices  $C_1, \dots, C_n$ , have a common diagonalizing matrix  $A$

$$C_i = A\Lambda_i A^T,$$

where  $\Lambda_i$  diagonal matrix.

- ▶ We are looking for the common diagonalizing matrix  $A$ .

$$\operatorname{argmin}_{A \in \mathbb{R}^{m \times m}, \Lambda_i} \sum_{i=1}^n \|C_i - A\Lambda_i A^T\|_F$$

for  $\|\cdot\|_F$  Frobenius-norm.

- ▶ Optimization technique: Block coordinate descent or Alternating Projection/Least Squares



# Conclusion

- Constructing Fisher-like information matrix is more expensive but we get
  - ▶ more robustness to the noise
  - ▶ faster algorithmic convergence
  - ▶ more efficient than using covariance matrices
- Non-stationary is not always bad :-)

Thank you