

Lien entre méthodes de séparation de sources, ACP, ACI et principes d'apprentissage

Ali MOHAMMAD-DJAFARI

Laboratoire des Signaux et Systèmes

UMR 08506 du CNRS-Supélec-UPS

Supélec, Plateau de Moulon

91192 Gif-sur-Yvette, FRANCE.

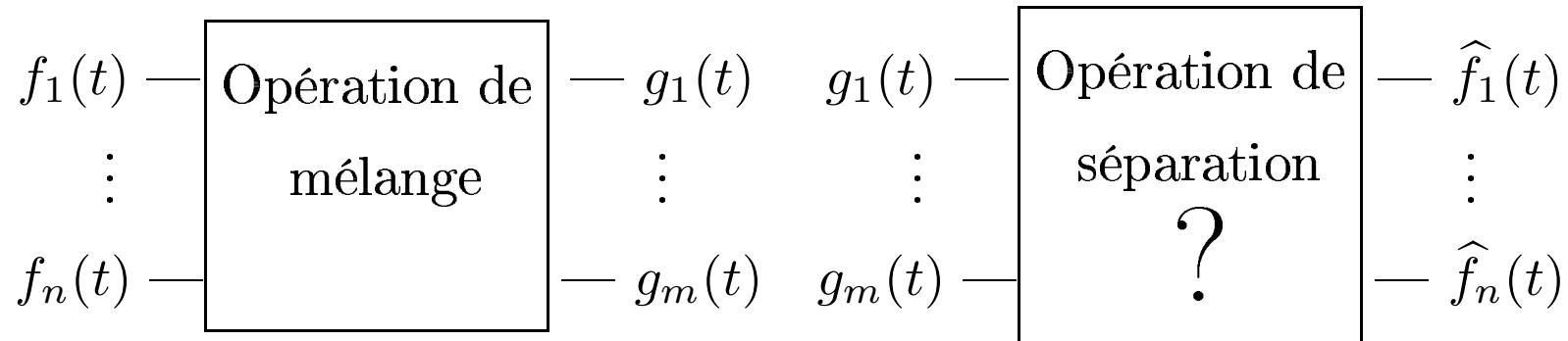
`djafari@lss.supelec.fr`

`http ://www.lss.supelec.fr`

`http ://djafari.free.fr`

INTRODUCTION

Problème de la séparation de source :



– **Mélange linéaire :**

$$g(t) = \int A(t, t') f(t') dt'$$

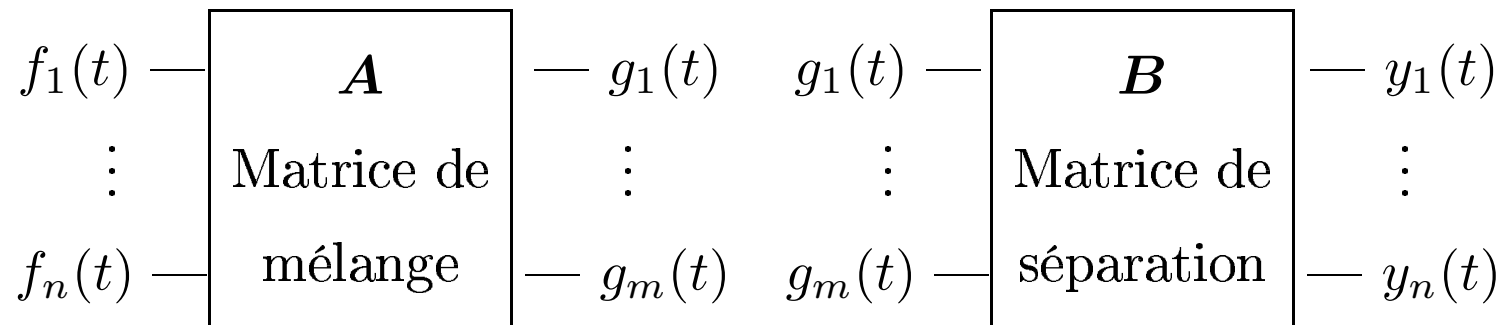
– **Mélange convolutive :**

$$g(t) = \int A(t - t') f(t') dt'$$

– **Mélange instantanée :**

$$g(t) = A f(t)$$

Mélange convolutive \longrightarrow **Déconvolution multi canaux**
Mélange instantané \longrightarrow **Séparation de sources**



– **Indétermination principale** : $\longrightarrow \mathbf{B} = \mathbf{\Sigma} \mathbf{\Lambda} \mathbf{A}^{-1}$

où $\mathbf{\Sigma}$ est une matrice de permutation et $\mathbf{\Lambda}$ une matrice diagonale.

– Hypothèses fondamentales : $f_1(t), \dots, f_n(t)$ sont **non corrélés (ACP)** ou **independants (ACI)**.

– Méthodes classiques : Infomax, Contrast, Ordres supérieures, Maximum de vraisemblance, Approche bayésienne

ANALYSE EN COMPOSANTS PRINCIPALES (ACP)

– Principale hypothèse :

$f_1(t), \dots, f_n(t)$ sont non corrélés entre eux et blanches.

$$\mathbf{R}_{ff} = \mathbb{E} \mathbf{f} \mathbf{f}^t = \Lambda$$

– Algorithme : $\mathbf{R}_{gg} = \mathbb{E} \mathbf{g} \mathbf{g}^t = \mathbb{E} \mathbf{A} \mathbf{f} \mathbf{f}^t \mathbf{A}^t = \mathbf{A} \mathbf{R}_{ff} \mathbf{A}^t = \mathbf{A} \Lambda \mathbf{A}^t$

– Estimation de la matrice de covariance

$$[\mathbf{R}_{gg}]_{kl} = \frac{1}{T} \sum_t g_k(t) g_l(t)$$

– Décomposition en valeurs singulières

$$\mathbf{R}_{gg} = \mathbf{A} \Lambda \mathbf{A}^t \longrightarrow \hat{\mathbf{f}}(t) = (\Lambda^+)^{1/2} \mathbf{A}^t \mathbf{g}(t)$$

– \mathbf{A} est déterminée à une rotation et à un facteur d'échelle près :

$$\mathbf{A} \longrightarrow \mathbf{A} \Theta \longrightarrow \mathbf{A} \Theta \Lambda \Theta^t \mathbf{A}^t = \mathbf{A} \Lambda \mathbf{A}^t \quad \text{avec } \Theta \text{ une matrice orthogonale}$$

ANALYSE EN COMPOSANTS INDÉPENDANTS (ACI)

- Imposer une structure pour l'estimateur :

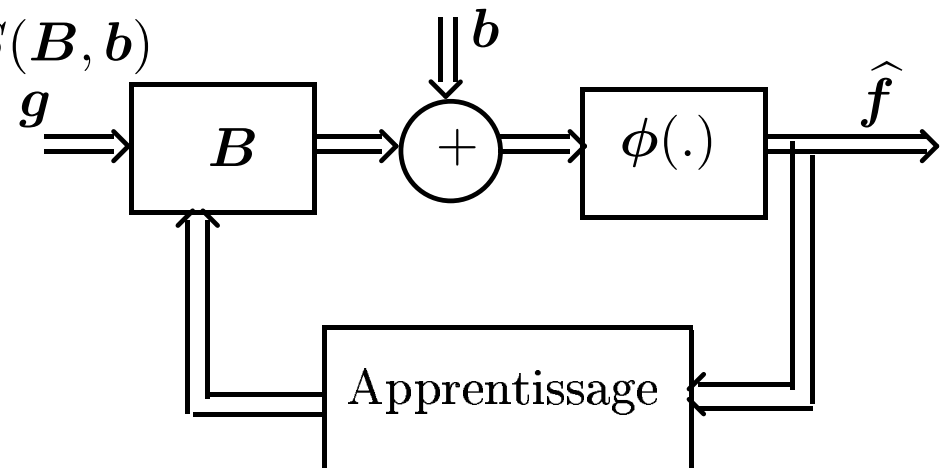
$$\hat{f}_i = \phi_i([\mathbf{B}\mathbf{g} + \mathbf{b}]_i) \longrightarrow \hat{\mathbf{f}} = \phi(\mathbf{B}\mathbf{g} + \mathbf{b})$$

- Choisir un critère de mesure de l'indépendance des composants de $\hat{\mathbf{f}}$:

$$S = - \sum_i p_i(\hat{f}_i) \ln p_i(\hat{f}_i) = - \sum_i p_i(\phi_i([\mathbf{B}\mathbf{g}]_i + b_i)) \ln p_i(\phi_i([\mathbf{B}\mathbf{g}]_i + b_i))$$

- Optimiser S : $(\hat{\mathbf{B}}, \hat{\mathbf{b}}) = \arg \max (\mathbf{B}, \mathbf{b}) S(\mathbf{B}, \mathbf{b})$

- Structure d'optimisation dans un réseau de neurons



MÉTHODES BASÉES SUR DES FONCTIONS DE CONTRASTS

Définir une fonction de contraste $c(\mathbf{y}) = c(\mathbf{B}\mathbf{g})$ qui aura pour minimiseur la matrice de séparation \mathbf{B} .

Exemple :

$$c(\mathbf{B}) = KL \left(p(\mathbf{y}) : \prod_i p_i(y_i) \right) = \int p(\mathbf{y}) \ln \frac{p(\mathbf{y})}{\prod_i p_i(y_i)} d\mathbf{y}$$

MÉTHODES BASÉES SUR LES STATISTIQUES D'ORDRE SUPÉRIEURS À DEUX

Approximation de $p_i(y_i)$ par ses cumulants

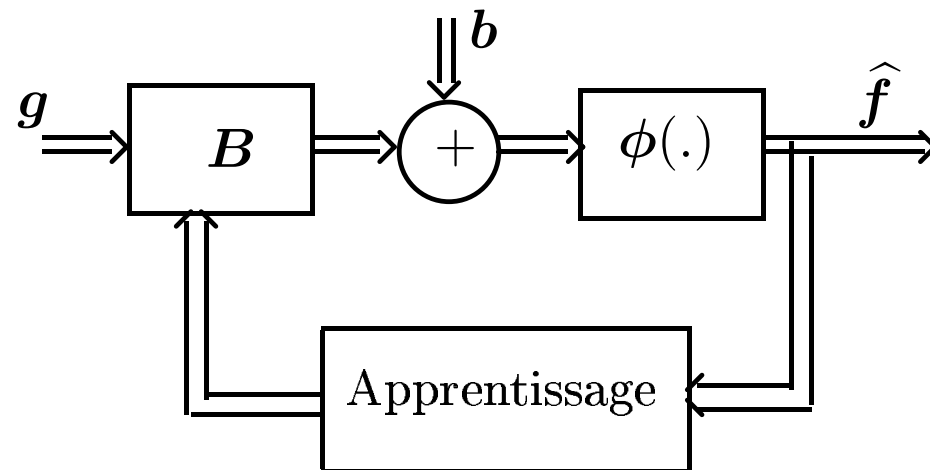
Développement en série de Taylor de la fonction caractéristique :

$$\frac{\partial c(\mathbf{B})}{\partial \mathbf{B}} \propto \mathbb{E} \frac{p'_i(y_i)}{p_i(y_i)} \propto \text{cumulants de } (y_i)$$

PRINCIPALES LIMITATIONS DES MÉTHODES CLASSIQUES

- La plus part de ces méthodes ne prennent pas en compte explicitement le bruit (les erreurs de la modélisation et le bruit de mesure) ;
- La plus part de ces méthodes supposent que la matrice de mélange est carrée et inversible.
- La plus part de ces méthodes supposent que les sources sont i.i.d. →
Baluchement avant ACP ou ACI
- La plus part de ces méthodes imposent *a priori* une structure :

$$\hat{f} = \phi(Bg + b)$$



APPROCHE BAYÉSIENNE

$$\mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) + \boldsymbol{\epsilon}(t), \quad t = 1, \dots, T$$

- Information sur $\boldsymbol{\epsilon}$ $\longrightarrow p(\mathbf{g}_{1..T} | \mathbf{A}, \mathbf{f}_{1..T})$
- Information sur $\mathbf{f}(t)$ $\longrightarrow p(\mathbf{f}_{1..T})$
- Information sur \mathbf{A} $\longrightarrow p(\mathbf{A})$
- Bayes :

$$p(\mathbf{A}, \mathbf{f}_{1..T} | \mathbf{g}_{1..T}) \propto p(\mathbf{g}_{1..T} | \mathbf{A}, \mathbf{f}_{1..T}) p(\mathbf{f}_{1..T}) p(\mathbf{A})$$

- Définir un estimateur utilisant :

$$p(\mathbf{A}, \mathbf{f}_{1..T} | \mathbf{g}_{1..T}), \quad p(\mathbf{A} | \mathbf{g}_{1..T}) \quad \text{ou} \quad p(\mathbf{f}_{1..T} | \mathbf{g}_{1..T})$$

Exemples de $p(\mathbf{A})$:

$$p(\mathbf{A}) \propto \exp[-\lambda \|\mathbf{A}\|^2], \quad p(\mathbf{A}) \propto \exp\left[-\frac{1}{2\sigma_a^2} \|\mathbf{I} - \mathbf{A}\|^2\right]$$

CAS D'UN MODÈLE EXACT, \mathbf{A} INVERSIBLE ET LES SOURCES IID

– \mathbf{A} inversible

$$\mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) = \mathbf{B}^{-1}\mathbf{f}(t), \quad t = 1, \dots, T$$

– Sources iid : $\mathbf{g} = \mathbf{A}\mathbf{f}$

$$p(\mathbf{f}) = \prod_i p_i(f_i) \longrightarrow p(\mathbf{g}|\mathbf{B}) = |\det(\mathbf{B})| \prod_i p_i([\mathbf{B}\mathbf{x}]_i)$$

$$\ln p(\mathbf{g}_{1..T}|\mathbf{B}) = \ln |\det(\mathbf{B})|^T + \sum_t \sum_i p_i(y_i(t)), \quad \text{avec } \mathbf{y}(t) = \mathbf{B}\mathbf{g}(t)$$

$$J(\mathbf{B}) = -\ln p(\mathbf{B}|\mathbf{g}_{1..T}) = -T \ln |\det(\mathbf{B})| - \sum_t \sum_i \ln p_i(y_i(t)) + \ln p(\mathbf{B}) + cte.$$

MAP ou Maximum de vraisemblance :

$$\frac{\partial J(\mathbf{B})}{\partial \mathbf{B}} = - \sum_t \mathbf{H}(\mathbf{y}(t)) \quad \text{avec } \mathbf{H}(\mathbf{y}) = \frac{\partial}{\partial \mathbf{B}} \left[\sum_i \ln p_i(y_i) + \ln |\det(\mathbf{B})| + \ln p(\mathbf{B}) \right]$$

Cas particulier : $p(\mathbf{B})$ uniforme \longrightarrow Maximum de vraisemblance

$$\mathbf{H}(\mathbf{y}) = \phi(\mathbf{y}) \mathbf{y}^\dagger - \mathbf{I},$$

avec

$$\phi_i(y_i) = -\frac{p'_i(y_i)}{p_i(y_i)}.$$

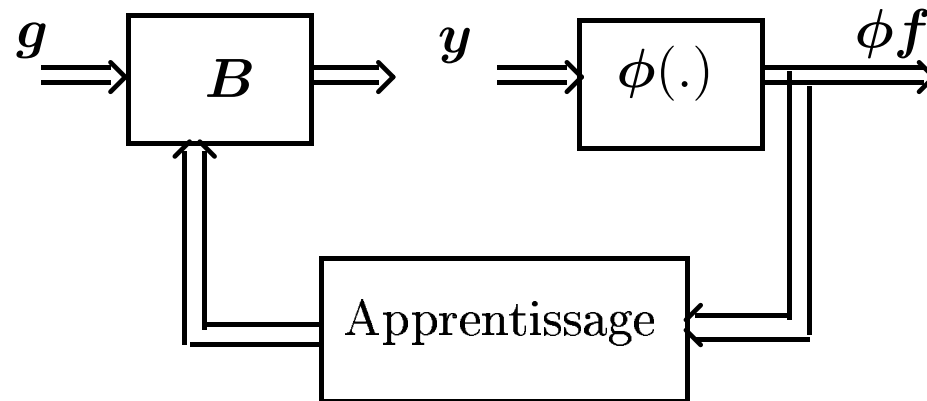
Gauss	$p(z) \propto \exp[-\alpha z^2]$	$\phi(z) = 2\alpha z$
Laplace	$p(z) \propto \exp[-\alpha z]$	$\phi(z) = \alpha \text{sign}(z)$
Cauchy	$p(z) \propto \frac{1}{1 + (z/\alpha)^2}$	$\phi(z) = \frac{2z/\alpha^2}{1+(z/\alpha)^2}$
sous gaussien	$p(z) \propto \exp\left[-\frac{1}{2}z^2\right] \text{sech}^2(z)$	$\phi(z) = z + \tanh(z)$
mélange de gaussiennes	$p(z) \propto \exp\left[-\frac{1}{2}(z - \alpha)^2\right] + \exp\left[-\frac{1}{2}(z + \alpha)^2\right]$	$\phi(z) = \alpha z - \alpha \tanh(\alpha z)$

LIEN AVEC L'APPRENTISSAGE ET RN

$$H(\mathbf{y}) = \frac{\partial}{\partial \mathbf{B}} \left[\sum_i \ln p_i(y_i) - \ln |\det(\mathbf{B})| \right].$$

Optimisation par un algorithme du descente :

$$\Delta \mathbf{B} \propto H(\mathbf{y}) = [\mathbf{I} - \phi(\mathbf{y})\mathbf{y}^t] \mathbf{B}$$



PRISE EN COMPTE DES ERREURS

$$\mathbf{g}(t) = \mathbf{A} \mathbf{f}(t) + \boldsymbol{\epsilon}(t), \quad t = 1, \dots, T.$$

$$\ln p(\boldsymbol{\epsilon}(1), \dots, \boldsymbol{\epsilon}(T)) = \sum_t \sum_i \ln p_i(\epsilon_i(t)).$$

– Vraisemblance :

$$\ln p(\mathbf{g}_{1..T} | \mathbf{A}, \mathbf{f}_{1..T}) = \sum_t \sum_i q_i(g_i(t) - [\mathbf{A}\mathbf{f}]_i(t))$$

avec $q_i(\cdot) = \ln p_i(\cdot)$.

– A posteriori :

$$\begin{aligned} \ln p(\mathbf{A}, \mathbf{f}_{1..T} | \mathbf{g}_{1..T}) &= \ln p(\mathbf{g}_{1..T} | \mathbf{A}, \mathbf{f}_{1..T}) + \ln p(\mathbf{f}_{1..T}) + \ln p(\mathbf{A}) + cte \\ &= \sum_t \sum_i q_i(g_i(t) - [\mathbf{A}\mathbf{f}]_i(t)) + \ln p(\mathbf{f}_{1..T}) + \ln p(\mathbf{A}) + cte. \end{aligned}$$

Trois directions (dépendant de l'application) :

- Intégrer hors du problème \mathbf{A} pour obtenir $p(\mathbf{f}_{1..T}|\mathbf{g}_{1..T})$ et estimer $\mathbf{f}_{1..T}$.

Par exemple

$$\hat{\mathbf{f}}_{1..T} = \arg \max \mathbf{f}_{1..T} p(\mathbf{f}_{1..T}|\mathbf{g}_{1..T})$$

- Intégrer hors du problème $\mathbf{f}_{1..T}$ pour obtenir $p(\mathbf{A}|\mathbf{g}_{1..T})$ et estimer \mathbf{A} .

Par exemple

$$\hat{\mathbf{A}} = \arg \max \mathbf{A} p(\mathbf{A}|\mathbf{g}_{1..T}),$$

et ensuite estimer \mathbf{f} par $\hat{\mathbf{f}} = \hat{\mathbf{A}}\mathbf{g}$;

- Estimation jointe de $\mathbf{f}_{1..T}$ et de \mathbf{A} en utilisant $p(\mathbf{A}, \mathbf{f}_{1..T}|\mathbf{g}_{1..T})$.

Par exemple

$$\begin{cases} \hat{\mathbf{f}}_{1..T}^{(k)} &= \arg \max \mathbf{f}_{1..T} p(\hat{\mathbf{A}}^{(k-1)}, \mathbf{f}_{1..T}|\mathbf{g}_{1..T}) \\ \hat{\mathbf{A}}^{(k)} &= \arg \max \mathbf{A} p(\mathbf{A}, \hat{\mathbf{f}}_{1..T}^{(k-1)}|\mathbf{g}_{1..T}) \end{cases}$$

SOURCES INDÉPENDANTES ET BLANCHES

$$\ln p(\mathbf{f}_{1..T}) = \sum_t \sum_j r_j(f_j(t))$$

$$p(\mathbf{A}) \propto \exp \left[-\frac{1}{2\sigma_a^2} \sum_k \sum_l a_{kl}^2 \right],$$

$$\ln p(\mathbf{A}, \mathbf{f}_{1..T} | \mathbf{g}_{1..T}) = \sum_t \sum_i q_i (g_i(t) - y_i(t)) + \sum_t \sum_j r_j(f_j(t)) + \frac{1}{2\sigma_a^2} \sum_k \sum_l a_{kl}^2 + cte.$$

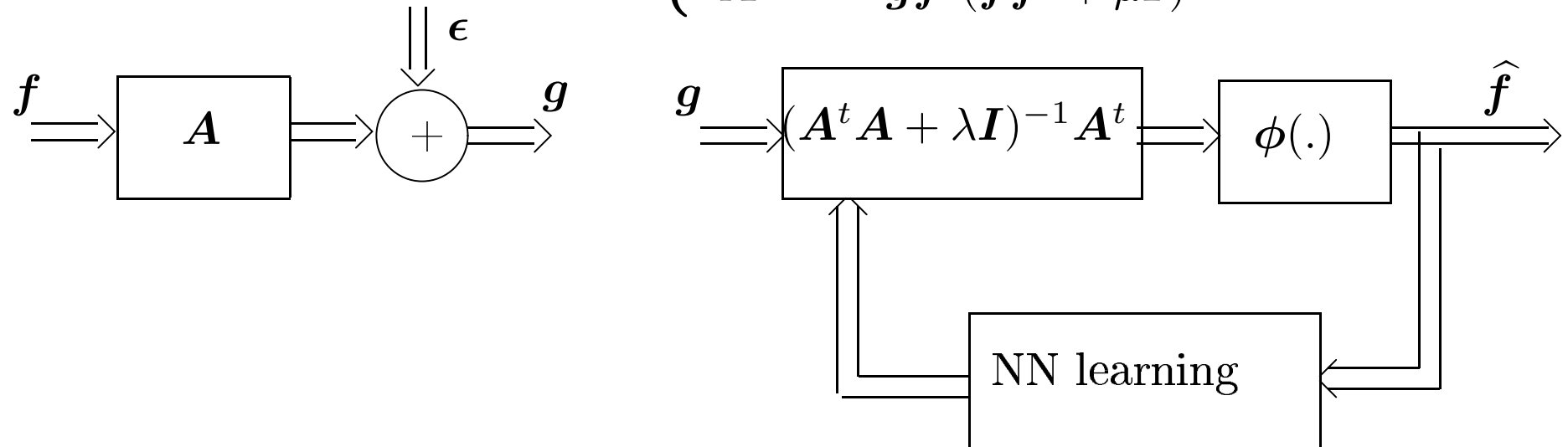
$$\left\{ \begin{array}{l} \hat{\mathbf{f}}^{(k)} = \arg \max_{\mathbf{f}} \sum_i q_i (g_i - y_i) + \sum_j r_j(f_j) \\ \hat{\mathbf{A}}^{(k)} = \arg \max_{\mathbf{A}} \sum_i q_i (g_i - y_i) + \frac{1}{2\sigma_a^2} \sum_k \sum_l a_{kl}^2 \end{array} \right.$$

- **Lois gaussiennes** pour le bruit et les sources :

$$\begin{cases} \mathbf{f} &= (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^t \mathbf{g} \\ \mathbf{A} &= \mathbf{g} \mathbf{f}^t (\mathbf{f} \mathbf{f}^t + \mu \mathbf{I})^{-1} \end{cases}$$

avec $\lambda = \sigma_n^2 / \sigma_s^2$ et $\mu = \sigma_n^2 / \sigma_a^2$.

- **Lois non gaussienne pour f** :
- $$\begin{cases} \mathbf{y} &= (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^t \mathbf{g} \\ \mathbf{f} &= \phi(\mathbf{y}) \\ \mathbf{A} &= \mathbf{g} \mathbf{f}^t (\mathbf{f} \mathbf{f}^t + \mu \mathbf{I})^{-1} \end{cases}$$



SOURCES DÉPENDANTES

$$\mathbf{f} = (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^t \mathbf{g} \quad \longrightarrow \quad \mathbf{f} = (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{D}^t \mathbf{D})^{-1} \mathbf{A}^t \mathbf{g}$$

SOURCES INDÉPENDANTES MAIS COLORÉES TEMPORELLEMENT

Principale difficulté : Modélisation de $p(f_j(1), \dots, f_j(T))$

– **Modèles markoviens :**

$$\ln p(f_j(1), \dots, f_j(T)) = \sum_t \ln p(f_j(t) | f_j(t-1))$$

– **Modèle de Gauss-Markov :**

$$\mathbf{f}(t) = (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{I})^{-1} [\text{diag}[\lambda_1, \dots, \lambda_n] \mathbf{f}(t-1) + \mathbf{A}^t \mathbf{g}(t)]$$

MARGINALISATION

$$\mathbf{g} = \mathbf{A}\mathbf{f} + \epsilon$$

$$p(\mathbf{A}, \mathbf{f}|\mathbf{g}) \propto \exp \left[-\frac{1}{2\sigma_n^2} J(\mathbf{A}, \mathbf{f}) \right] \quad \text{avec} \quad J(\mathbf{A}, \mathbf{f}) = \|\mathbf{g} - \mathbf{A}\mathbf{f}\|^2 + \lambda\phi(\mathbf{f}) + \mu\psi(\mathbf{A})$$

$$p(\mathbf{A}|\mathbf{g}) = \int p(\mathbf{A}, \mathbf{f}|\mathbf{g}) \, d\mathbf{f}$$

Approximation du second ordre :

$$-\ln p(\mathbf{A}|\mathbf{g}) \propto -\ln \left| \det \left(\widehat{\mathbf{P}}_s^{-1} \right) \right| - J(\mathbf{A}, \widehat{\mathbf{f}})$$

$$\widehat{\mathbf{f}} = (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^t \mathbf{g} \quad \text{and} \quad \widehat{\mathbf{P}}_s = (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{I})^{-1};$$

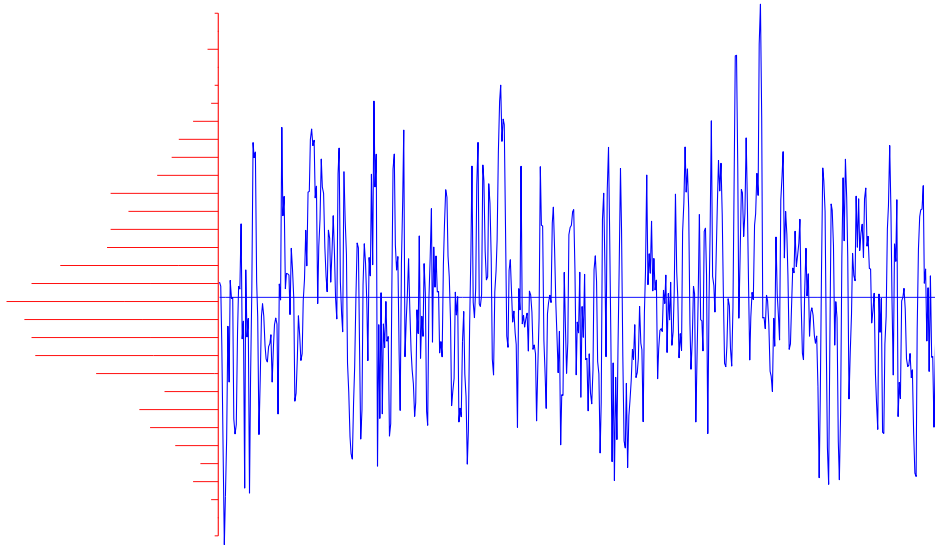
→ On obtient le même schéma algorithmique de RN avec :

$$\Delta \mathbf{A} \propto \mathbf{A}^t (\mathbf{A}^t \mathbf{A} + \lambda \mathbf{I})^{-1} + \mathbf{g}\mathbf{f} + \mu\psi'(\mathbf{A})$$

Algorithme Espérance-Maximisation (EM) :

Données complètes (\mathbf{g}, \mathbf{f}) et incomplètes \mathbf{g} .

CHOIX D'UNE LOI DE PROBABILITÉ A PRIORI

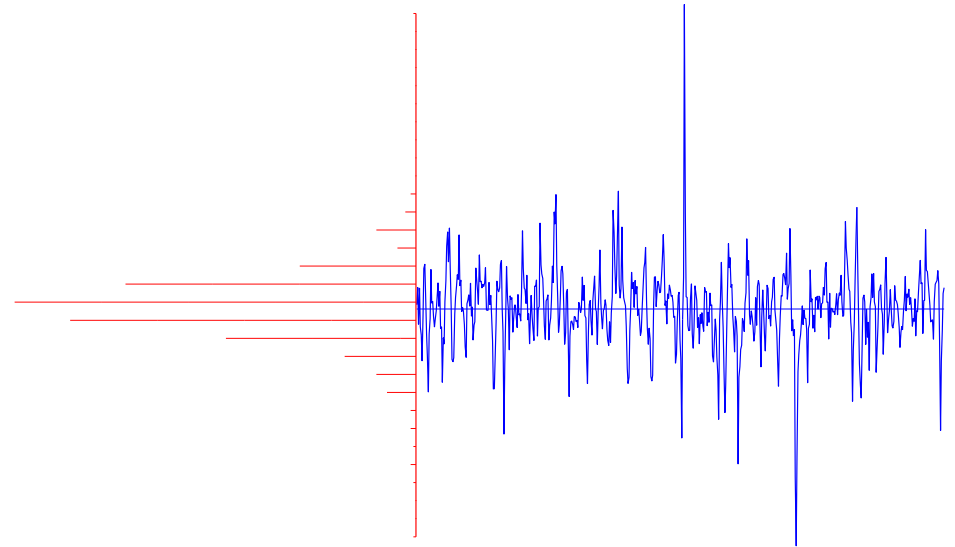


Gaussienne (Gauss-Markov)

$$p(f_j | f_{j-1}) = \mathcal{N}(f_{j-1}, \sigma_f^2)$$

$$p(\mathbf{f}) \propto \exp \left[-\alpha \sum_j |f_j - f_{j-1}|^2 \right]$$

$$p(\mathbf{f}) \propto \exp \left[-\alpha \sum_{r \in \mathcal{R}} |f(\mathbf{r}) - \beta \sum_{s \in \mathcal{V}(\mathbf{r})} f(\mathbf{s})|^2 \right]$$



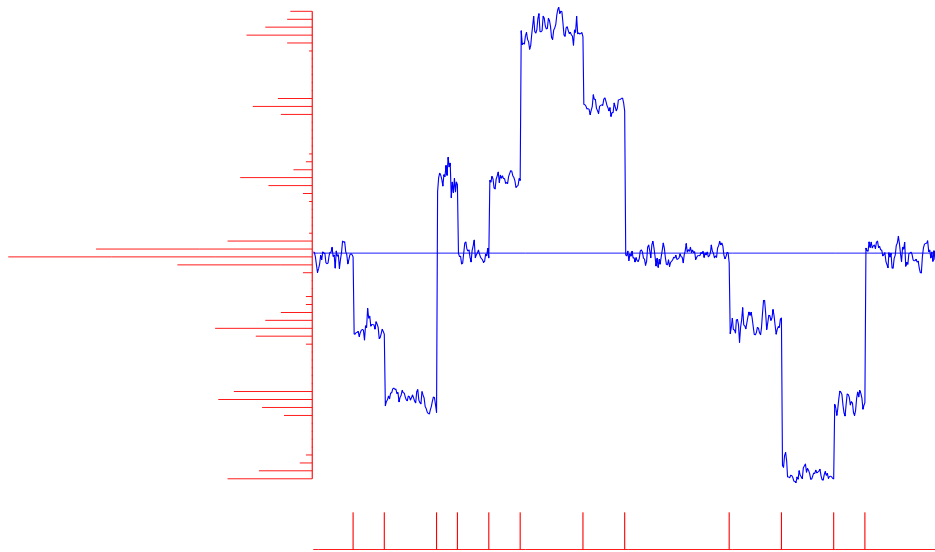
Gaussienne généralisée

$$p(f_j | f_{j-1}) \propto \exp [-\alpha |f_j - f_{j-1}|^p]$$

$$p(\mathbf{f}) \propto \exp \left[-\alpha \sum_j |f_j - f_{j-1}|^p \right]$$

$$p(\mathbf{f}) \propto \exp \left[-\alpha \sum_{r \in \mathcal{R}} |f(\mathbf{r}) - \beta \sum_{s \in \mathcal{V}(\mathbf{r})} f(\mathbf{s})|^p \right]$$

MODÈLES A PRIORI HIÉRARCHIQUES AVEC CHAMPS CACHÉS (HMRF)

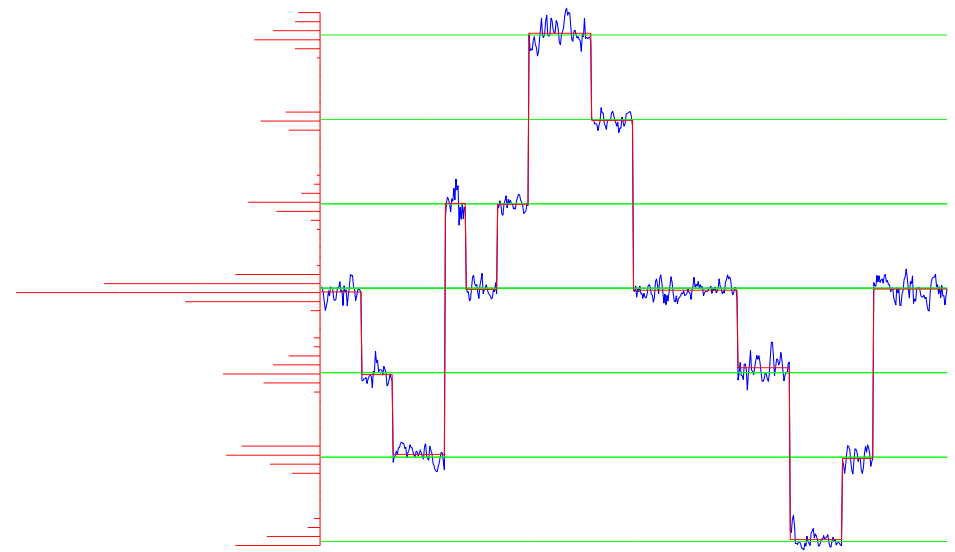


Gauss-Markov + processus de ligne

$$p(\mathbf{f}_j | q_j, \mathbf{f}_{j-1}) = \mathcal{N}((1 - q_j)\mathbf{f}_{j-1}, \sigma_f^2)$$

$$p(\mathbf{f} | \mathbf{q}) \propto \exp \left[-\alpha \sum_j (1 - q_j) |\mathbf{f}_j - \mathbf{f}_{j-1}|^2 \right]$$

$$P(q_j = 1) = \alpha, P(q_j = 0) = 1 - \alpha$$



Mélange de gaussiennes ou GM+régions

$$p(\mathbf{f}_j | z_j = k) = \mathcal{N}(m_k, \sigma_k^2)$$

$$p(\mathbf{f} | \mathbf{z}) \propto \exp \left[-\alpha \sum_k \sum_{j \in \mathcal{R}_k} (\mathbf{f}_j - m_k / \sigma_k)^2 \right]$$

$$p(\mathbf{z}) \propto \exp \left[\alpha \sum_j \sum_{i \in \mathcal{V}(j)} \delta(z_j - z_i) \right]$$

MODÉLISATION MÉLANGE DE GAUSSIENNES ET VARIABLES CACHÉES

$$p(f_j | z_j = k) = \mathcal{N}(m_{jk}, \sigma_{jk}^2) \longrightarrow p(f_j) = \sum_{k=1}^{K_j} \alpha_{jk} \mathcal{N}(m_{jk}, \sigma_{jk}^2)$$

– **Variable cachée** : $z_j \in \mathcal{Z}_j = (1, \dots, K_j)$ avec $\alpha_{jk} = p(z_j = k)$

$$p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}) = \prod_j \prod_{z_j \in \mathcal{Z}_j} \mathcal{N}(m_{jk}, \sigma_{jk}^2), \quad \boldsymbol{\theta} = \{(m_{jk}, \sigma_{jk}^2)\}$$

– Loi *a priori*
$$p(\mathbf{f} | \boldsymbol{\theta}) = \sum_{\mathbf{z} \in \mathcal{Z}} p(\mathbf{z}) p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta})$$

– Loi *a posteriori*
$$p(\mathbf{f} | \mathbf{g}, \boldsymbol{\theta}) = \sum_{\mathbf{z} \in \mathcal{Z}} p(\mathbf{z} | \mathbf{g}, \boldsymbol{\theta}) p(\mathbf{f} | \mathbf{g}, \mathbf{z}, \boldsymbol{\theta})$$

– **Algorithme globale de l'estimation** :

- Estimer \mathbf{z} utilisant $p(\mathbf{z} | \mathbf{g}, \boldsymbol{\theta})$;
- Estimer \mathbf{f} utilisant $p(\mathbf{f} | \mathbf{g}, \mathbf{z}, \boldsymbol{\theta})$
- Estimer les hyperparamètres $\boldsymbol{\theta}$ utilisant $p(\boldsymbol{\theta} | \mathbf{g}, \mathbf{z}, \mathbf{f})$

CONCLUSIONS

- Il y a un lien entre les problèmes de séparation de sources et la notion de l'apprentissage.
- Les méthodes classiques supposent \mathbf{A} inversible et ne prennent pas en compte explicitement les incertitudes.
- Approche bayésienne permet de combiner l'information dans les données et l'information *a priori*, et prends en compte explicitement les incertitudes.
- Dans les méthodes classique la structure du schéma de l'apprentissage est fixée d'avance. Dans l'approche bayésienne, cette structure dépend de la modélisation du bruit et des signaux d'entrée.
- En général, le calcul bayésien est coûteux :
nécessité de faire des approximations,
par exemple l'approximation en champs moyens (Mean field approximation)