

Teoria da Informação

Charles Casimiro Cavalcante

`charles@gtel.ufc.br`

Grupo de Pesquisa em Telecomunicações Sem Fio – GTEL
Programa de Pós-Graduação em Engenharia de Teleinformática
Universidade Federal do Ceará – UFC
<http://www.gtel.ufc.br/~charles>

“A principal função de um sistema de comunicação é reproduzir, exatamente ou de forma aproximada, uma informação proveniente de outro ponto diferente.”

Claude Shannon, 1948

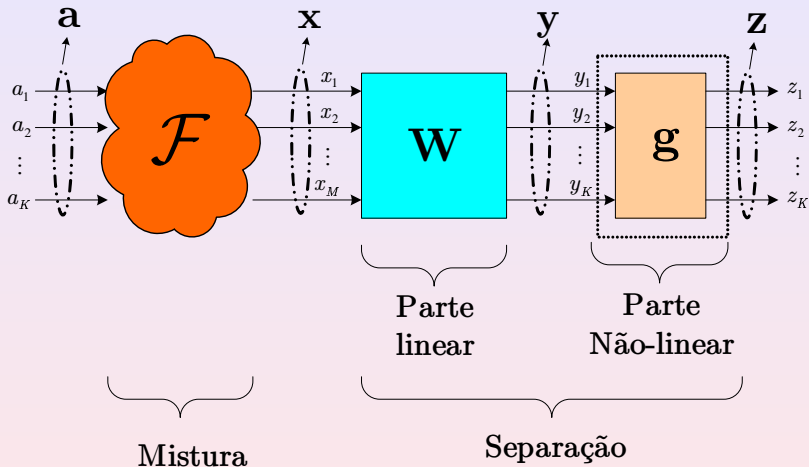
Conteúdo do curso

- 1 Revisão de probabilidade
- 2 Informação e Entropia
- 3 Codificação de fontes
- 4 Codificação e capacidade de canal
- 5 Complexidade de Kolmogorov
- 6 Funções de otimização
- 7 *Independent Component Analysis*

Parte VII

Independent Component Analysis

Modelo geral



Modelo

$$\mathbf{x}(n) = \mathcal{F}(\mathbf{a}(n), \mathbf{v}(n), n) \leftarrow \text{Mapeamento}$$

$$\mathbf{a}(n) = \begin{bmatrix} a_1(n) & a_2(n) & \cdots & a_K(n) \end{bmatrix}^T \leftarrow K \text{ fontes}$$

$$\mathbf{v}(n) = \begin{bmatrix} v_1(n) & v_2(n) & \cdots & v_V(n) \end{bmatrix}^T \leftarrow V \text{ sinais de ruído}$$

$$\mathbf{x}(n) = \begin{bmatrix} x_1(n) & x_2(n) & \cdots & x_M(n) \end{bmatrix}^T \leftarrow M \text{ sensores.}$$

Considerações usuais

- \mathcal{F} é linear e invariante no tempo
- Fontes mutuamente independentes e independentes do ruído
- $V = M$
- $M \geq K$ (mais sensores que fontes no mínimo)

Mistura

$$\mathbf{x}(n) = \mathbf{H}\mathbf{a}(n) + \mathbf{v}(n)$$

Separação

$$\mathbf{y}(n) = \mathbf{W}^H \mathbf{x}(n) = \hat{\mathbf{a}}(n)$$

Características

- 1 Indeterminação em relação a permutação e escalonamento

$$\mathbf{y}(n) = \mathbf{P}\mathbf{D}\mathbf{a}(n)$$

\mathbf{P} é uma matriz de permutação de ordem $K \times K$ e \mathbf{D} é uma matriz diagonal e inversível de ordem $K \times K$.

- 2 Possível inserção de não-linearidade após \mathbf{W}

Equações: 2 fontes e 2 sensores

$$\mathbf{x} = \mathbf{H}\mathbf{a}$$

$$\begin{cases} x_1 = h_{11}a_1 + h_{12}a_2 \\ x_2 = h_{21}a_1 + h_{22}a_2 \end{cases}$$

Duas (K) variáveis e duas (M) equações!

Separando as fontes...

- Caso sem ruído: *separar é possível se* $\mathbf{W} = \mathbf{H}^{-1}$

Questão 1

Como identificar \mathbf{H} para projetar \mathbf{W} ?

Questão 2

Quais (e quantas) estatísticas são necessárias para prover a separação?

Resposta para questão 1: Como identificar \mathbf{H} ?

- Observando a matriz de autocorrelação do vetor de sinais recebidos:

$$\mathbf{R}_x = \mathbb{E} \{ \mathbf{x}(n) \mathbf{x}^T(n) \} = \mathbf{H} \mathbf{R}_a \mathbf{H}^T = \mathbf{H} \mathbf{H}^T,$$

- Observar que $\mathbf{H} \mathbf{Q}^T$, em que \mathbf{Q} é uma matriz ortogonal também soluciona equação
- Conclui-se que: $\mathbf{H} = \mathbf{R}_x^{\frac{1}{2}}$
- Extração de raiz quadrada de matrizes: através de decomposição em valores singulares (SVD, Singular Value Decomposition)

Identificando \mathbf{H} : possibilidades

- Escrevendo

$$\mathbf{H} = \mathbf{U}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{V}^T,$$

em que \mathbf{U} e \mathbf{V} são matrizes retangulares de ordem $K \times M$, tais que $\mathbf{U}\mathbf{U}^T = \mathbf{V}\mathbf{V}^T = \mathbf{I}_M$ e $\mathbf{U}^T\mathbf{U} = \mathbf{V}^T\mathbf{V} = \mathbf{I}_K$

- Então...

$$\begin{aligned}\mathbf{R}_x &= \mathbf{U}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{V}\mathbf{V}^T\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{U}^T \\ &= \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T.\end{aligned}$$

- Resultado: através da matriz de autocorrelação (estatística de ordem 2) só é possível identificar as matrizes \mathbf{U} e $\mathbf{\Lambda}$
- Estatísticas de ordem 2 (SOS, Second Order Statistics) não resolvem o problema por completo
- **Resposta parcial da questão 2: estatísticas necessárias para separação!**

Máximo possível com SOS

Processamento

Projeção dos dados na direção da inversa da matriz de autocorrelação: $\mathbf{T} = \mathbf{R}_x^{-\frac{1}{2}}$

- Assim, tem-se o seguinte conjunto de dados:

$$\bar{\mathbf{x}}(n) = \mathbf{T}\mathbf{x}(n)$$

de tal forma que

$$\begin{aligned}\mathbf{R}_{\bar{x}} &= \mathbf{T}\mathbf{R}_x\mathbf{T}^T \\ &= \mathbf{U}\mathbf{\Lambda}^{-\frac{1}{2}}\mathbf{U}^T\mathbf{U}\mathbf{\Lambda}\mathbf{U}^T\mathbf{U}^T\mathbf{\Lambda}^{-\frac{1}{2}}\mathbf{U} \\ &= \mathbf{I}_K\end{aligned}$$

- Branqueamento!** - em separação cega de fontes é denominado de *esferatização*

- Projeção dos dados na direção dos principais autovetores de \mathbf{R}_x
- Técnica de análise por componentes principais (PCA, Principal Component Analysis): projeção dos dados nas direções definidas pelos principais autovetores de \mathbf{H}
- Matriz de branqueamento, ou *transformação de Mahalanobis*, é escrita como

$$\mathbf{T} = \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}_x^T.$$

- Redução do problema para uma matriz de mistura ortogonal

$$\bar{\mathbf{x}}(n) = \mathbf{T} \mathbf{H} \mathbf{a}(n) = \mathbf{Q} \mathbf{a}(n)$$

em que $\mathbf{Q} \mathbf{Q}^T = \mathbf{I}$.

Encontrando a matriz \mathbf{V}

Considerar um atraso arbitrário ℓ para o qual, não haja duas fontes, i e j , com a mesma autocorrelação

$$\mathbb{E} \{a_i(n)a_i(n-\ell)\} \neq \mathbb{E} \{a_j(n)a_j(n-\ell)\} \quad \forall i \neq j,$$

Diferentes densidades espectrais!

$$\begin{aligned} \mathbf{R}_{\bar{x}}(k) &= \mathbb{E} \{ \bar{\mathbf{x}}(n) \bar{\mathbf{x}}^T(n-\ell) \} \\ &= \mathbf{T} \mathbf{H} \mathbf{R}_a(\ell) \mathbf{H}^T \\ &= \underbrace{\mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T}_{\mathbf{T}} \underbrace{\mathbf{U} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{V}^T}_{\mathbf{H}} \mathbf{R}_a(\ell) \underbrace{\mathbf{V} \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{U}^T}_{\mathbf{H}^T} \underbrace{\mathbf{U} \mathbf{\Lambda}^{-\frac{1}{2}}}_{\mathbf{T}^T} \\ &= \mathbf{V}^T \mathbf{R}_a(\ell) \mathbf{V}. \end{aligned}$$

- Fontes independentes $\Rightarrow \mathbf{R}_a(\ell)$ diagonal
- Autovalores distintos: fontes com diferentes autocorrelações
- Conclusão: *uso somente de SOS permite separar fontes com espectros diferentes!*
- Algoritmo clássico: *Algorithm for Multiple Unknown Signals Extraction* (AMUSE)
- Na prática, fontes com espectros similares (embora diferentes) não são separáveis
- Impossibilidade de separar fonte brancas ou i.i.d.

Como separar fontes com mesmo espectro?

Proposição de uma nova técnica

- Não procura-se estimar \mathbf{H} para projetar \mathbf{W}
- Fator decisivo: suposição de independência das fontes!
- Projetar \mathbf{W} de maneira que sejam obtidas fontes *o mais independentes possível* na saída do dispositivo de separação
- *Independent Component Analysis (ICA)*
- Restrição: no máximo uma fonte pode ser gaussiana
 - Teorema Central do Limite: mistura de gaussianas é gaussiana!
- Interpretação: ponto chave é a não-gaussianidade das fontes.

Questão

Como medir a não-gaussianidade das fontes e utilizar o fato de que as mesmas são estatisticamente independentes para separá-las?

Independência estatística

$$p_y(\mathbf{y}) \triangleq \prod_{i=1}^K p_{y_i}(y_i)$$

$$\mathbb{E}\{y_1 \cdot y_2 \cdots y_K\} = \mathbb{E}\{y_1\} \cdot \mathbb{E}\{y_2\} \cdots \mathbb{E}\{y_K\}$$

Descorrelação estatística

$$\mathbb{E}\{y_1 \cdot y_2 \cdots y_K\} - \mathbb{E}\{y_1\} \cdot \mathbb{E}\{y_2\} \cdots \mathbb{E}\{y_K\} = 0.$$

Independência \Rightarrow Descorrelação
 \nLeftarrow

Entropia e informação mútua

$$\mathcal{H}(\mathbf{x}) \triangleq -\mathbb{E} \{ \ln [p_x(\mathbf{x})] \} = - \int_{-\infty}^{\infty} p_x(\mathbf{x}) \cdot \ln [p_x(\mathbf{x})] d\mathbf{x}$$

entropia diferencial

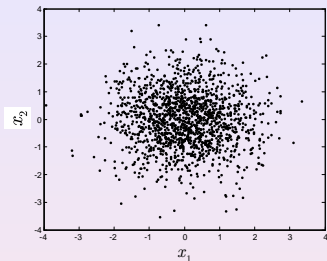
$$\begin{aligned} \mathcal{H}(\mathbf{x}|\mathbf{y}) &= -\mathbb{E} \{ \ln [p_{x|y}(\mathbf{x}|\mathbf{y})] \} \\ &= \int p_{x,y}(\mathbf{x}, \mathbf{y}) \cdot \ln [p_{x|y}(\mathbf{x}|\mathbf{y})] d\mathbf{x}d\mathbf{y} \end{aligned}$$

entropia condicional

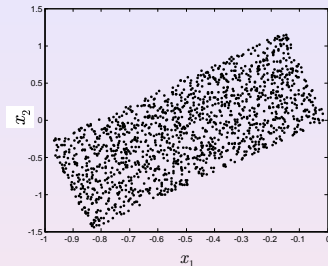
$$\mathcal{I}(\mathbf{x}, \mathbf{y}) = \mathcal{H}(\mathbf{x}) - \mathcal{H}(\mathbf{x}|\mathbf{y})$$

informação mútua entre \mathbf{x} e \mathbf{y}

E se as fontes forem gaussianas?



Mistura de fontes gaussianas



Mistura de fontes uniformes

Não há direções preferenciais!

Divergência de Kullback-Leibler (KLD)

$$D(p_x(\mathbf{x})||g_x(\mathbf{x})) \triangleq \int_{-\infty}^{\infty} p_x(\mathbf{x}) \cdot \ln \left[\frac{p_x(\mathbf{x})}{g_x(\mathbf{x})} \right] dx$$

em que $p(x)$ e $g(x)$ são duas funções estritamente positivas

Propriedades:

- ① é sempre de valor positivo ou zero; KLD é zero para o caso específico de $p_x(\mathbf{x}) = g_x(\mathbf{x})$.
- ② é invariante com relação as seguintes mudanças nas componentes do vetor \mathbf{x} ;
 - permutação de ordem
 - escalonamento de amplitude
 - transformação monotônica não-linear

Usando a KLD...

$$\begin{aligned}\mathcal{I}(\mathbf{x}, \mathbf{y}) &= \int p_{x,y}(\mathbf{x}, \mathbf{y}) \cdot \ln \left[\frac{p_{x,y}(\mathbf{x}, \mathbf{y})}{p_x(\mathbf{x}) p_y(\mathbf{y})} \right] d\mathbf{x} d\mathbf{y} \\ &= D(p_{x,y}(\mathbf{x}, \mathbf{y}) \| p_x(\mathbf{x}) p_y(\mathbf{y}))\end{aligned}$$

ou ainda

$$\begin{aligned}\mathcal{I}(\mathbf{y}) &= D(p_y(\mathbf{y}) \| \tilde{p}_y(\mathbf{y})) \\ &= \int \cdots \int p_y(y_1, \dots, y_K) \cdot \ln \left[\frac{p_y(y_1, \dots, y_K)}{\prod_{i=1}^K p_{y_i}(y_i)} \right] dy_1 \cdots dy_K\end{aligned}$$

Finalmente...

$$\mathcal{I}(\mathbf{y}) = -\mathcal{H}(\mathbf{y}) + \sum_{i=1}^K \mathcal{H}(y_i),$$

o que significa que minimizar a informação mútua entre os componentes do vetor \mathbf{y} significa tornar a entropia de \mathbf{y} o mais próximo possível da soma de suas entropias marginais.

Negentropia - I

- Medida de “não-gaussianidade” baseado na medida de entropia diferencial
- Diferença entre a entropia da v.a. \mathbf{y} e a entropia de uma v.a. \mathbf{y}^G de distribuição gaussiana e com os mesmos momentos de ordem um e dois (média e variância) de \mathbf{y}

$$N_G(\mathbf{y}) \triangleq \mathcal{H}(\mathbf{y}^G) - \mathcal{H}(\mathbf{y}).$$

ou

$$N_G(\mathbf{y}) \triangleq D(p_{\mathbf{y}}(\mathbf{y}) \| p_{\mathbf{y}^G}(\mathbf{y}))$$

- **Problema: requer conhecimento ou estimativa da densidade de probabilidade de \mathbf{y} !**

Negentropia - II

- Possibilidade de expressar estatísticas necessárias para separação
- Utilizando KLD e informação mútua pode-se escrever:

$$\mathcal{I}(\mathbf{y}) = \mathcal{I}(\mathbf{y}^G) + \left(N_G(\mathbf{y}) - \sum_{i=1}^K N_G(y_i) \right).$$

- 1 Primeiro termo utiliza SOS
- 2 Segundo termo mede não-gaussianidade do sinal

Kurtosis

- Cumulante de quarta ordem

$$\mathcal{K}\{y\} \triangleq \mathbb{E}\{y^4\} - 3 \cdot (\mathbb{E}\{y^2\})^2.$$

- Faixa de valores:

- ▶ Distribuição gaussiana: $\mathcal{K}\{y\} = 0$
- ▶ Distribuição sub-gaussiana: $\mathcal{K}\{y\} \leq 0$
- ▶ Distribuição super-gaussiana: $\mathcal{K}\{y\} \geq 0$

- Propriedades

$$\mathcal{K}\{y_1 + y_2\} = \mathcal{K}\{y_1\} + \mathcal{K}\{y_2\}$$

$$\mathcal{K}\{\alpha \cdot y\} = \alpha^4 \cdot \mathcal{K}\{y\}$$

Funções de contraste

- Uma função $\Psi(\cdot)$, no espaço de K fdps (distintas ou não) é dita ser um *contraste* se respeita as seguintes condições:

- 1 $\Psi(p_y)$ é invariante a permutações:

$$\Psi(p_{Py}) = \Psi(p_y) \text{ para qualquer matriz de permutação } \mathbf{P}$$

- 2 $\Psi(p_y)$ é invariante a mudanças de escala:

$$\Psi(p_{Dy}) = \Psi(p_y) \text{ para qualquer matriz diagonal } \mathbf{D}$$

- 3 Se \mathbf{y} possui componentes independentes, então:

$$\Psi(p_{Wy}) \leq \Psi(p_y) \text{ para qualquer matriz inversível } \mathbf{W}$$

Algumas funções de contraste

$$\Psi_{\text{ICA}}(p_y) = -\mathcal{I}(\mathbf{y})$$

ou ainda em relação à matriz de separação

$$\Psi_{\text{ICA}}(\mathbf{W}) = \ln [|\det(\mathbf{W})|] - \mathbb{E} \left\{ \ln \left[\prod_{i=1}^K p_{y_i}(y_i) \right] \right\}$$

Aproximações da negentropia

- Necessário para evitar estimações da fdp
- Momentos de ordem superior (clássica)

$$N_G(\mathbf{y}) \approx \frac{1}{12} [\mathbb{E} \{\mathbf{y}^3\}]^2 + \frac{1}{48} [\mathcal{K} \{\mathbf{y}\}]^2,$$

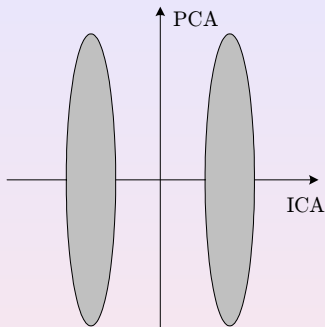
- Classe de aproximações (*FastICA*)

$$N_G(\mathbf{y}) \approx \sum_{i=1}^p \varrho_i \cdot [\mathbb{E} \{g_i(\mathbf{y})\} - \mathbb{E} \{g_i(\nu)\}]^2$$

$$g_1(u) = \frac{1}{a_1} \log [\cosh(a_1 u)] \text{ e}$$

$$g_2(u) = -\exp\left(-\frac{u^2}{2}\right)$$

Diferença ICA \times PCA



- 1 PCA: projeta numa dimensão menor preservando *máxima variância* dos dados
- 2 ICA: projeta numa dimensão menor preservando *estrutura* dos dados

MaxEnt/InfoMax - I

- Separação é obtida quando estimativas das fontes são independentes $\Rightarrow I(\mathbf{y}) = 0$

$$\mathcal{I}(\mathbf{y}, \mathbf{x}) = \mathcal{H}(\mathbf{y}) - \mathcal{H}(\mathbf{y}|\mathbf{x})$$

- Mapeamento é determinístico, logo $\mathcal{H}(\mathbf{y}|\mathbf{x}) = 0$

$$\mathcal{I}(\mathbf{y}, \mathbf{x}) = \mathcal{H}(\mathbf{y})$$

- Entropia de \mathbf{y} não é limitada, então

$$\mathbf{z} = \mathbf{g}(\mathbf{W}\mathbf{x})$$

MaxEnt/InfoMax - II

- Funções $g_i(\cdot)$ monotonicamente crescentes, limitadas de tal forma que $g_i(-\infty) = 0$ e $g_i(\infty) = 1$
- Se $g_i(\cdot)$ for igual à função de distribuição cumulativa (fdc) da i -ésima fonte

$$p_z(\mathbf{z}) = U[0, 1], \quad \text{para}$$

- Adaptação

$$\Delta \mathbf{W} \propto (\mathbf{W}^{-T})^{-1} - 2 \cdot \tanh(\mathbf{W}\mathbf{x}) \mathbf{x}^T,$$

- Função contraste

$$\max_{\mathbf{W}} \left(\Psi_{\text{InfoMax}}(\mathbf{W}) \triangleq \ln[|\det(\mathbf{W})|] - \mathbb{E} \left\{ \ln \left[\prod_{i=1}^K g'_i(y_i) \right] \right\} \right).$$

Máxima verossimilhança - I

- Entropia

$$\mathcal{H}(\mathbf{z}) = - \int p_z(\mathbf{z}) \ln \left[\frac{p_z(\mathbf{z})}{\prod_{i=1}^K U(z_i)} \right] d\mathbf{z} = -D(p_z(\mathbf{z}) \| U_N(\mathbf{z}))$$

- Modelo paramétrico Θ

$$\mathcal{L}_Q(\Theta) = \frac{1}{Q} \ln \left[\prod_{q=1}^Q p_x(\mathbf{x}(q) | \Theta) \right] = \frac{1}{Q} \sum_{q=1}^Q \ln [p_x(\mathbf{x}(q) | \Theta)] .$$

Máxima verossimilhança - II

- Lei dos grandes números

$$\mathcal{L}_Q(\Theta) \xrightarrow{Q \rightarrow \infty} \mathcal{L}(\mathbf{W}) \triangleq \int p_y(\mathbf{y}|\mathbf{W}) \ln [p_{\tilde{\mathbf{a}}}(\mathbf{y})] d\mathbf{y} + \ln [|\det(\mathbf{W})|] .$$

- Usando KLD

$$\begin{aligned} \mathcal{L}(\mathbf{W}) &= -D(p_x(\mathbf{y}|\mathbf{W}) \| p_{\tilde{\mathbf{a}}}(\mathbf{y})) - \mathcal{H}(\mathbf{y}|\mathbf{W}) + \ln [|\det(\mathbf{W})|] \\ &= -D(p_x(\mathbf{y}|\mathbf{W}) \| p_{\tilde{\mathbf{a}}}(\mathbf{y})) - \mathcal{H}(\mathbf{x}) . \end{aligned}$$

- Função contraste MV

$$\Psi_{\text{MV}}(\mathbf{W}) = -D(\mathbf{W}\mathbf{H}\mathbf{a} \| \tilde{\mathbf{a}}) .$$

Máxima verossimilhança - III

- Interpretação da MV

$$-\Psi_{MV}(\mathbf{W}) = -\Psi_{ICA}(\mathbf{W}) + \sum_{i=1}^K D(\tilde{z}_i \| \tilde{a}_i)$$

ou seja

$$\left(\begin{array}{c} \text{Desvio} \\ \text{total} \end{array} \right) = \left(\begin{array}{c} \text{Desvio da} \\ \text{independência} \end{array} \right) + \left(\begin{array}{c} \text{Desvio} \\ \text{marginal} \end{array} \right).$$

e ainda

$$\Psi_{ICA}(\mathbf{W}) = \max_{\tilde{\mathbf{a}}} (\Psi_{MV}(\mathbf{W})).$$

Critério “universal”

$$J_{\text{BSS}}(\mathbf{W}) = \ln [|\det|(\mathbf{W})] - \ln \left[\prod_{i=1}^K \phi_i(y_i) \right]$$

em que $\mathbf{y} = \mathbf{W}\mathbf{x}$ e, idealmente, as funções $\phi_i(\cdot)$ são as fdps das fontes.

Questão

Como encontrar ou estimar as fdps das fontes?

Resposta possível

Expansões polinomiais através dos cumulantes.