

Machine Intelligence 2

2.2 ICA: The Infomax method

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

SS 2017

Statistical Independence & Infomax

$$I(X, Y) = \underset{\text{PG/DKL divergence}}{KL[P_j(X, Y); P_i(X, Y)]} = \underset{\text{Entropy}}{H(Y) - H(Y|X)}$$

Statistical independence

Scenario Data

observations: $\underline{\mathbf{x}} \in \mathbb{R}^N$ with distribution $P_{\underline{\mathbf{x}}}(\underline{\mathbf{x}})$

estimated sources: $\hat{\underline{\mathbf{s}}} = \underline{\mathbf{W}} \cdot \underline{\mathbf{x}}$ Model Parameter: Matrix $\underline{\mathbf{W}}$

$P_{\underline{\mathbf{s}}}(\hat{\underline{\mathbf{s}}})$: family of true (unknown) densities, parametrized by $\underline{\mathbf{W}}$

$\leadsto P_{\underline{\mathbf{s}}}(\hat{\underline{\mathbf{s}}}) = P_{\underline{\mathbf{s}}}(\underline{\mathbf{W}} \cdot \underline{\mathbf{x}})$

Model selection

$$\hat{P}_{\underline{\mathbf{s}}}(\hat{\underline{\mathbf{s}}}) = \prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i) \leftarrow \text{assumption: statistical independence}$$

Performance measure

$$D_{\text{KL}} = D_{\text{KL}}[P_{\underline{\mathbf{s}}}(\hat{\underline{\mathbf{s}}}), \hat{P}_{\underline{\mathbf{s}}}(\hat{\underline{\mathbf{s}}})] = \int d\underline{\mathbf{s}} P_{\underline{\mathbf{s}}}(\underline{\mathbf{s}}) \ln \frac{P_{\underline{\mathbf{s}}}(\underline{\mathbf{s}})}{\prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)} \stackrel{!}{=} \min_{\underline{\mathbf{W}}}$$

Equally distributed sources



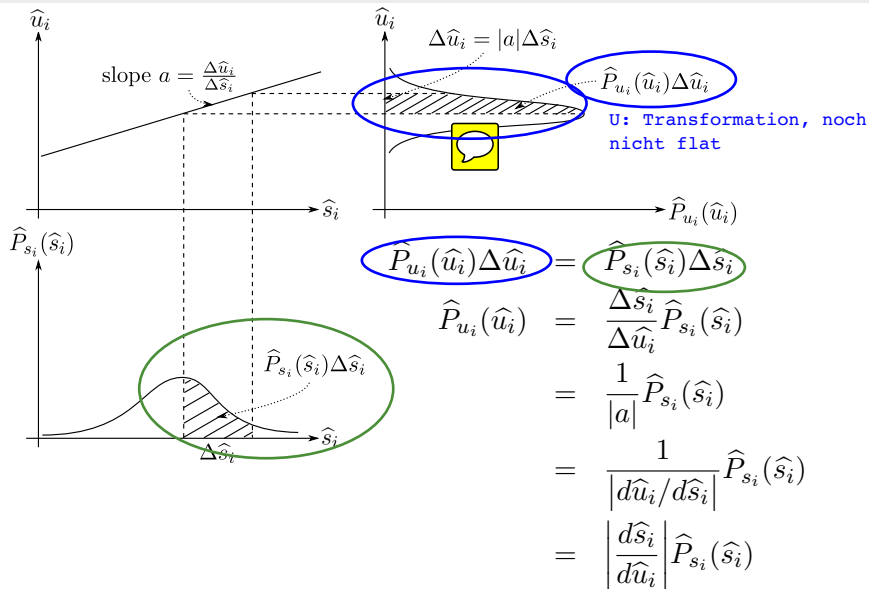
non-linear transformations: $\hat{u}_i = \hat{f}_i(\hat{s}_i)$, such that $\hat{P}_{u_i}(\hat{u}_i) = \text{const.}$

conservation of probability

$$\hat{P}_{u_i}(\hat{u}_i)d\hat{u}_i = \hat{P}_{s_i}(\hat{s}_i)d\hat{s}_i$$



Conservation of probability



Equally distributed sources

non-linear transformations: $\hat{u}_i = \hat{f}_i(\hat{s}_i)$, such that $\hat{P}_{u_i}(\hat{u}_i) = \text{const.}$

conservation of probability



$$\hat{P}_{u_i}(\hat{u}_i) d\hat{u}_i = \hat{P}_{s_i}(\hat{s}_i) d\hat{s}_i$$

$$\hat{P}_{u_i}(\hat{u}_i) = \left| \frac{d\hat{s}_i}{d\hat{u}_i} \right| \hat{P}_{s_i}(\hat{s}_i) = \frac{1}{|\hat{f}_i'(\hat{s}_i)|} \hat{P}_{s_i}(\hat{s}_i) \stackrel{!}{=} 1$$

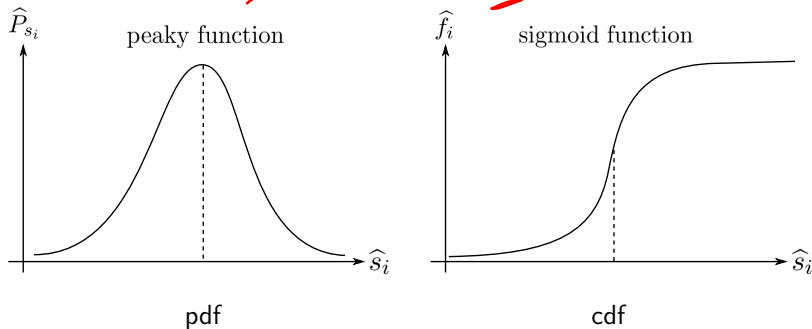
It follows:

$$|\hat{f}_i'(\hat{s}_i)| = \hat{P}_{s_i}(\hat{s}_i) \Rightarrow \hat{f}_i(\hat{s}_i) = \int_{-\infty}^{\hat{s}_i} dy \hat{P}_{s_i}(y)$$

\hat{f}_i : cumulative density function (cdf) of $\hat{P}_{s_i}(\hat{s}_i)$

Equally distributed sources

Make it flat



Flat is good bc no density required?

The Infomax principle

Statistical independence:



S. Aufschrieb

$$D_{\text{KL}} = \int d\hat{\underline{\mathbf{s}}} P_{\underline{\mathbf{s}}}(\hat{\underline{\mathbf{s}}}) \ln \frac{P_{\underline{\mathbf{s}}}(\hat{\underline{\mathbf{s}}})}{\prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)} \stackrel{!}{=} \min$$

Transformation:

Nonlinear Transformation of \mathbf{s}

$$\hat{u}_i = \hat{f}_i(\underbrace{\mathbf{W} \cdot \mathbf{x}}_{\hat{\underline{\mathbf{s}}}})$$

see blackboard

The Infomax principle (Bell, Sejnowski, 1995)

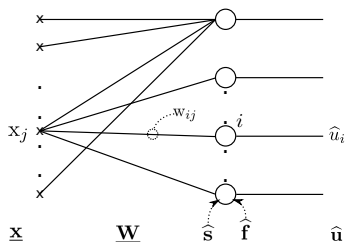
Max entropy

$$H = - \int d\hat{\underline{\mathbf{u}}} P_{\underline{\mathbf{u}}}(\hat{\underline{\mathbf{u}}}) \ln P_{\underline{\mathbf{u}}}(\hat{\underline{\mathbf{u}}}) \stackrel{!}{=} \max_{\underline{\mathbf{W}}} \quad \text{}$$



Perceptron implementation

Architecture



perceptron: $\hat{u}_i = \hat{f}_i\left(\sum_j w_{ij}x_j\right)$ observations: $\underline{x}^{(\alpha)} \in \mathbb{R}^N, \alpha = 1, \dots, p$

Cost function

$$H = - \int d\underline{\hat{\mathbf{u}}} P_{\underline{\mathbf{u}}}(\underline{\hat{\mathbf{u}}}) \ln P_{\underline{\mathbf{u}}}(\underline{\hat{\mathbf{u}}}) \stackrel{!}{=} \max_{\underline{\mathbf{W}}} \quad \text{💬}$$

Empirical Risk Minimization

$$H = - \int d\underline{\hat{\mathbf{u}}} P_{\underline{\mathbf{u}}}(\underline{\hat{\mathbf{u}}}) \ln P_{\underline{\mathbf{u}}}(\underline{\hat{\mathbf{u}}}) \stackrel{!}{=} \max_{\underline{\mathbf{w}}} \quad$$

see blackboard

Empirical Risk Minimization

$$H = - \int d\hat{\underline{\mathbf{u}}} P_{\underline{\mathbf{u}}}(\hat{\underline{\mathbf{u}}}) \ln P_{\underline{\mathbf{u}}}(\hat{\underline{\mathbf{u}}}) \stackrel{!}{=} \max_{\underline{\mathbf{W}}} \quad$$

Model selection: maximize E^G



$$E^G = \ln |\det \underline{\mathbf{W}}| + \int d\underline{\mathbf{x}} P_{\underline{\mathbf{x}}}(\underline{\mathbf{x}}) \left\{ \sum_{l=1}^N \ln \tilde{f}_l \left(\sum_{k=1}^N w_{lk} x_k \right) \right\}$$

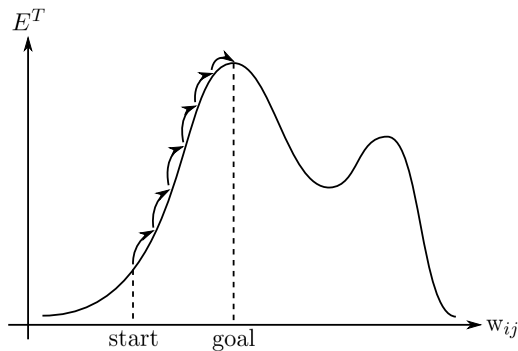
ERM principle: mathematical expectation $E^G \longrightarrow$ empirical average E^T

Approximation weil nicht berechenbar

$$E^T = \ln |\det \underline{\mathbf{W}}| + \frac{1}{p} \sum_{\alpha=1}^p \sum_{l=1}^N \ln \tilde{f}_l \left(\sum_{k=1}^N w_{lk} x_k^{(\alpha)} \right)$$

Gradient based optimization

Gradient Ascent



$$\Delta w_{ij} = \underbrace{\varepsilon}_{\text{learning step}} \cdot \frac{\partial E^T}{\partial w_{ij}}$$



Gradient ascent on the training cost.

Gradient based optimization

Batch Learning: $\Delta w_{ij} = \frac{\partial E^T}{\partial w_{ij}} = \frac{\varepsilon}{p} \sum_{\alpha} \frac{\partial e^{\alpha}}{\partial w_{ij}}$



On-line learning: $\Delta w_{ij} = \eta \frac{\partial e^{\alpha}}{\partial w_{ij}}$ and time-dependent (\searrow) learning rate η

Gradient based optimization

Batch Learning: $\Delta w_{ij} = \frac{\partial E^T}{\partial w_{ij}} = \frac{\varepsilon}{p} \sum_{\alpha} \frac{\partial e^{\alpha}}{\partial w_{ij}}$

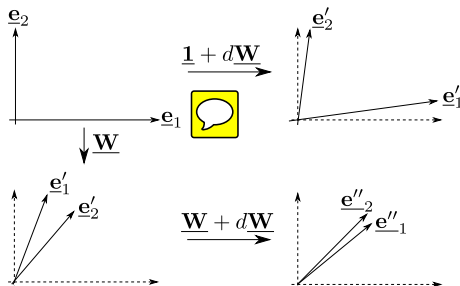
On-line learning: $\Delta w_{ij} = \eta \frac{\partial e^{\alpha}}{\partial w_{ij}}$ and time-dependent (\searrow) learning rate η

$$e^{(\alpha)} = \ln |\det \underline{\mathbf{W}}| + \sum_{l=1}^N \ln \hat{f}_l' \left(\sum_{k=1}^N w_{lk} x_k^{(\alpha)} \right)$$

$$e'^{(\alpha)} = \frac{\partial e^{(\alpha)}}{\partial w_{ij}} = \underbrace{(\underline{\mathbf{W}}^{-1})_{ji}}_{\text{costly computation}} + \frac{\hat{f}_i'' \left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)} \right)}{\hat{f}_i' \left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)} \right)} \cdot x_j^{(\alpha)} \quad \text{Learning Step}$$

Natural Gradient

linear transformations: $d\underline{\mathbf{W}}, \underline{\mathbf{W}}$



Normalized step size:

$$d\underline{\mathbf{Z}} = \underbrace{d\underline{\mathbf{W}}}_{\text{then do } d\underline{\mathbf{W}}} \cdot \underbrace{\underline{\mathbf{W}}^{-1}}_{\text{transform back to } \underline{\mathbf{1}}}$$

\Rightarrow make learning steps "comparable"

Natural Gradient

Taylor expansion of e ($e^{(\alpha)}$, but α suppressed in the following):

$$\begin{aligned} e(\underline{\mathbf{W}} + d\underline{\mathbf{W}}) &= e(\underline{\mathbf{W}}) + \nabla e(\underline{\mathbf{W}}) d\underline{\mathbf{W}} \\ &\stackrel{d\underline{\mathbf{W}} = \varepsilon \underline{\mathbf{D}}_{\mathbf{w}}}{=} e(\underline{\mathbf{W}}) + \varepsilon [\nabla e(\underline{\mathbf{W}})]^T \cdot \underline{\mathbf{D}}_{\mathbf{w}} \end{aligned}$$

learning step:

$$\begin{aligned} d\underline{\mathbf{Z}} &= d\underline{\mathbf{W}} \cdot \underline{\mathbf{W}}^{-1} \\ &= \varepsilon \underline{\mathbf{D}}_{\mathbf{w}} \cdot \underline{\mathbf{W}}^{-1} \end{aligned}$$

direction of steepest ascent under normalized step-size:

$$\begin{aligned} [\nabla e(\underline{\mathbf{W}})] \underline{\mathbf{D}}_{\mathbf{w}} &\stackrel{!}{=} \max_{\underline{\mathbf{D}}_{\mathbf{w}}} \\ \left(\underline{\mathbf{D}}_{\mathbf{w}} \cdot \underline{\mathbf{W}}^{-1} \right)^2 &\stackrel{!}{=} 1 \end{aligned}$$

Natural Gradient

Solution using Lagrange multipliers:

$$\sum_{i,j=1}^N \frac{\partial e}{\partial \mathbf{W}_{ij}} (\underline{\mathbf{D}}_{\mathbf{W}})_{ij} - \lambda \sum_{i,j,k,l=1}^N (\underline{\mathbf{D}}_{\mathbf{W}})_{ij} (\underline{\mathbf{W}}^{-1})_{jl} (\underline{\mathbf{D}}_{\mathbf{W}})_{ik} (\underline{\mathbf{W}}^{-1})_{kl} \stackrel{!}{=} \max_{\underline{\mathbf{D}}_{\mathbf{W}}} \quad$$

Taking the derivative w.r.t. $(\underline{\mathbf{D}})_{ps}$ and setting to zero yields:

$$\begin{aligned} \frac{\partial e}{\partial (\underline{\mathbf{D}}_{\mathbf{W}})_{ps}} &= \frac{\partial e}{\partial (\underline{\mathbf{W}})_{ps}} - \lambda \sum_{k,l=1}^N (\underline{\mathbf{W}}^{-1})_{sl} (\underline{\mathbf{D}}_{\mathbf{W}})_{pk} (\underline{\mathbf{W}}^{-1})_{kl} \\ &\quad - \lambda \sum_{i,j=1}^N (\underline{\mathbf{D}}_{\mathbf{W}})_{pj} (\underline{\mathbf{W}}^{-1})_{jl} (\underline{\mathbf{W}}^{-1})_{sl} \\ &= \frac{\partial e}{\partial (\underline{\mathbf{W}})_{ps}} - 2\lambda \sum_{k,l=1}^N (\underline{\mathbf{D}}_{\mathbf{W}})_{pk} (\underline{\mathbf{W}}^{-1})_{kl} (\underline{\mathbf{W}}^{-1})_{sl} \stackrel{!}{=} 0 \\ \frac{\partial e}{\partial \underline{\mathbf{W}}} &= 2\lambda \underline{\mathbf{D}}_{\mathbf{W}} \underline{\mathbf{W}}^{-1} (\underline{\mathbf{W}}^{-1})^T \\ \underline{\mathbf{D}}_{\mathbf{W}} &= \frac{1}{2\lambda} \frac{\partial e}{\partial \underline{\mathbf{W}}} \underline{\mathbf{W}}^T \underline{\mathbf{W}} \end{aligned}$$

Natural Gradient

$$e(\underline{\mathbf{W}} + d\underline{\mathbf{W}}) = e(\underline{\mathbf{W}}) + \varepsilon [\nabla e(\underline{\mathbf{W}})] \cdot \underline{\mathbf{D}}_{\mathbf{W}}$$

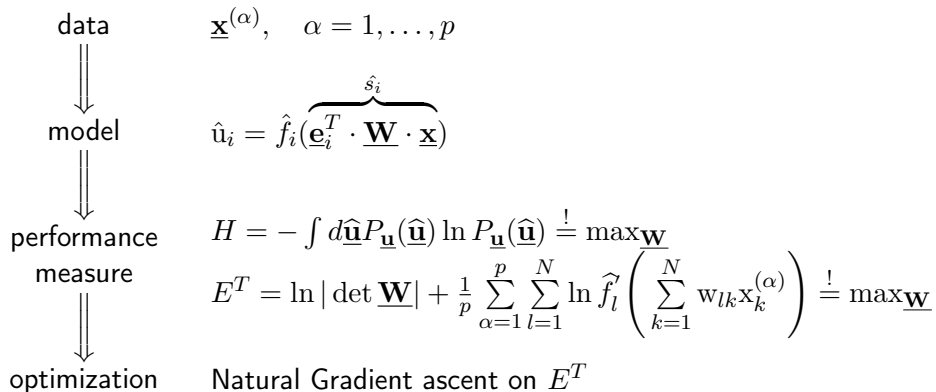
Inserting the optimal direction for "natural" gradient ascent yields:

$$\Delta \underline{\mathbf{W}} = \varepsilon \overbrace{\frac{\partial e}{\partial \underline{\mathbf{W}}}}^{\text{"original" gradient}} \underbrace{\underline{\mathbf{W}}^T \underline{\mathbf{W}}}_{\text{normalization of step size}}$$

i.e.

$$\Delta w_{ij} = \varepsilon \sum_{l=1}^N \left\{ \delta_{il} + \frac{\hat{f}_i'' \left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)} \right)}{\hat{f}_i' \left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)} \right)} \sum_{k=1}^N w_{lk} x_k^{(\alpha)} \right\} w_{lj}$$

Summary: The Infomax method



Practical aspects: Source amplitudes

Problem: Undetermined source amplitudes \leadsto convergence problems

① Bell-Sejnowski solution:

$$\Delta w_{ii} = 0 \text{ and } w_{ii} = 1 \text{ for all } i$$

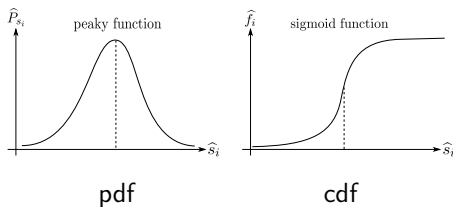
② Amari solution: Learning steps are always orthogonal to subspace of equivalent unmixing matrices.

$$\Delta w_{ij} = \varepsilon \frac{\hat{f}_i'' \left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)} \right)}{\hat{f}_i' \left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)} \right)} \sum_{l \neq i}^N \left(\sum_{k=1}^N w_{lk} x_k^{(\alpha)} \right) w_{lj}$$

Practical aspects: Choice of \hat{f}_i

Problem: True distribution (of u and s) and its cumulative distribution function is **unknown**.

Solution: Choose peaky PDF \leadsto sigmoid CDF



typical choice:

$$\hat{f}_{(y)} = \frac{1}{1 + \exp(-y)} \Rightarrow \frac{\hat{f}_{(y)}''}{\hat{f}_{(y)}} = 1 - 2\hat{f}_{(y)}$$

Observation: ICA is fairly robust against details of the choice of \hat{f} !

Practical aspects: Choice of \hat{f}_i

Natural Gradient (batch):

$$\Delta \mathbf{w}_{ij} = \varepsilon \sum_{l=1}^N \left\{ \delta_{il} + \frac{1}{p} \sum_{\alpha=1}^p \frac{\hat{f}_i'' \left(\sum_{k=1}^N \mathbf{w}_{ik} \mathbf{x}_k^{(\alpha)} \right)}{\hat{f}_i' \left(\sum_{k=1}^N \mathbf{w}_{ik} \mathbf{x}_k^{(\alpha)} \right)} \sum_{k=1}^N \mathbf{w}_{lk} \mathbf{x}_k^{(\alpha)} \right\} \mathbf{w}_{lj}$$

Stationary state:

$$\Delta \underline{\mathbf{w}}_{ij} \stackrel{!}{=} 0 \quad \Rightarrow \quad \delta_{il} \stackrel{!}{=} - \frac{1}{p} \sum_{\alpha=1}^p \underbrace{\frac{\hat{f}_i'' \left(\sum_{k=1}^N \mathbf{w}_{ik} \mathbf{x}_k^{(\alpha)} \right)}{\hat{f}_i' \left(\sum_{k=1}^N \mathbf{w}_{ik} \mathbf{x}_k^{(\alpha)} \right)}}_{\varphi_i(\hat{s}_i^{(\alpha)})} \cdot \underbrace{\sum_{k=1}^N \mathbf{w}_{lk} \mathbf{x}_k^{(\alpha)}}_{\hat{s}_l^{(\alpha)}}$$

Practical aspects: Choice of \hat{f}_i

$$-\frac{1}{p} \sum_{\alpha=1}^p \varphi_i(\hat{s}_i^{(\alpha)}) \hat{s}_i^{(\alpha)} \stackrel{!}{=} \delta_{il}$$

Ansatz: $\hat{s}_i = \lambda_i s_i \rightarrow$ estimated \sim true source signals

$i = l$: Through proper choice of λ_i we can always fulfill:

$$-\frac{1}{p} \sum_{\alpha=1}^p \varphi_i(\hat{s}_i^{(\alpha)}) \lambda_i s_i^{(\alpha)} \stackrel{!}{=} 1$$

$i \neq l$: Limit of large number of observations:

$$\begin{aligned} \frac{1}{p} \sum_{\alpha=1}^p \varphi_i(\hat{s}_i^{(\alpha)}) \lambda_l s_l^{(\alpha)} &\rightarrow \left\langle \varphi_i(\hat{s}_i^{(\alpha)}) \lambda_l s_l^{(\alpha)} \right\rangle_{P_{\underline{s}}} \\ \left\langle \varphi_i(\lambda_i s_i) \lambda_l s_l \right\rangle &\underbrace{=}_{\text{statistical independence}} \left\langle \varphi_i(\lambda_i s_i) \right\rangle \left\langle \lambda_l s_l \right\rangle \stackrel{!}{=} 0 \end{aligned}$$

Can always be fulfilled if data is centered: $\langle s_l \rangle = 0$

Practical aspects: Choice of \hat{f}_i

True (independent) source signals are always a fixed point of the natural gradient ascent \rightarrow independent of choice of \hat{f}_i

However: if \hat{f}_i deviates too strongly from its true shape, the fixed point may become unstable:

\Rightarrow if in doubt (and enough data available):

\leadsto make a parametrized ansatz for \hat{f}_i

\leadsto estimate parameters in addition to \mathbf{w}