

# Chapter 1

## Independent Component Analysis(ICA)

ICA tries to find original signals  $\hat{S}$ . from mixed signal  $X$ , called Blind Source Separation problem.

$$X = AS \quad \rightarrow \quad \hat{S} = WX$$

where

- $X \in \mathbb{R}^{N \times p}$
- $A, W \in \mathbb{R}^{N \times N}$
- $S, \hat{S} \in \mathbb{R}^{N \times p}$
- $N$  is no. sources
- $p$  is no. observations in the source

One could say that ICA tries to find  $W \approx A^{-1}$ .

### 1.1 ICA vs PCA

The goal of PCA is to find directions that don't correlate with each other( orthogonal projections ). However, ICA finds sources that are independent to each other. The right most figure of Fig 1.1 shows that even correlation  $C_{12}, C_{21}$  is zero, but  $x_1, x_2$  are dependent to each other( knowing one of them can tell position of the other ).

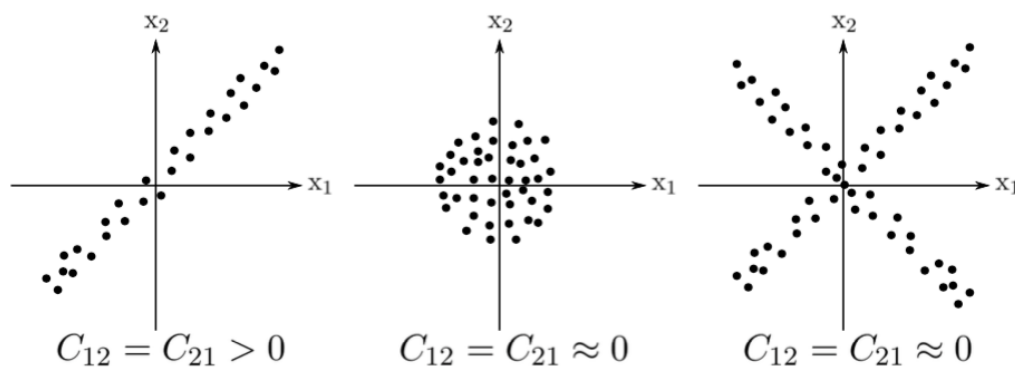


Figure 1.1: Decorrelation vs Independent (Drawn from MI2 Slides)

### 1.2 ICA Limitations

#### 1.2.1 Permutation of Sources

$\hat{s}_i$  that found that ICA is not necessary the same as the original source  $s_i$ .

### 1.2.2 Discovered Amplitude

### 1.2.3 Gaussian Sources

The underly theory behind ICA is Central Limit Theorem in which it implies that sum of random variables will approach Gaussian distribution. Thus, if the original sources come from Gaussian distribution, we won't have any clue to separate the sources.

## 1.3 ICA Approaches

An ICA algorithm can be derived using different cost functions using assumption that

$$\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}}) = \prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)$$

and  $P_{\mathbf{s}}(\hat{\mathbf{s}}) = P_{\mathbf{s}}(W\mathbf{x})$

### 1.3.1 Informax

From Kullback-Leibler divergence, we know that

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

Thus, we can formulate the problem as

$$\min_W D_{KL}(P_{\mathbf{s}}(\hat{\mathbf{s}})||\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}}))$$

However, deriving cost from  $D_{KL}$  directly would not work as it require probability density estimation. Hence, we shall use a transformation :

$$\hat{u}_i = \hat{f}(\hat{s}_i) \quad \text{s.t.} \quad \hat{P}_{u_i}(\hat{u}_i) = \text{const.}$$

From Conservation of Probability, we know that

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

Then

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

This implies that  $\hat{f}(\hat{s}_i)$  is CDF of  $\hat{P}_{s_i}(\hat{s}_i)$ . From this knowledge, we can derive  $D_{KL}$  :

$$\begin{aligned} D_{KL}(P_{\mathbf{s}}(\hat{\mathbf{s}})||\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}})) &= \int_{-\infty}^{\infty} d\hat{\mathbf{s}} P_{\mathbf{s}}(\hat{\mathbf{s}}) \log \frac{P_{\mathbf{s}}(\hat{\mathbf{s}})}{\prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)} \\ &= \int_{-\infty}^{\infty} d\hat{\mathbf{s}} P_{\mathbf{s}}(\hat{\mathbf{s}}) \log \frac{P_{\mathbf{s}}(\hat{\mathbf{s}})}{\prod_{i=1}^N \hat{P}_{s_i}(\hat{s}_i)} \frac{\prod_{i=1}^N 1/\hat{f}'}{\prod_{i=1}^N 1/\hat{f}'} \\ &= \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log \frac{P_{\mathbf{u}}(\hat{\mathbf{u}})}{\prod_{i=1}^N \hat{P}_{u_i}(\hat{u}_i)} \\ &= \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log P_{\mathbf{u}}(\hat{\mathbf{u}}) - \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log \prod_{i=1}^N \hat{P}_{u_i}(\hat{u}_i) \end{aligned}$$

Hence,

$$\begin{aligned}
\min D_{KL}(P_{\mathbf{s}}(\hat{\mathbf{s}})||\hat{P}_{\mathbf{s}}(\hat{\mathbf{s}})) &\approx \min \int_{-\infty}^{\infty} d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log P_{\mathbf{u}}(\hat{\mathbf{u}}) \\
&\approx \min -H(u) \\
&\approx \max H(u)
\end{aligned}
\tag{Informax Principle}$$

### Learning via Neural Network and Empirical Risk Minimization

From Conservation of Probability, we know that

$$\begin{aligned}
P(\hat{\mathbf{u}})d\hat{\mathbf{u}} &= P(\hat{\mathbf{s}})d\hat{\mathbf{s}} \\
&= P(\mathbf{x})d\mathbf{x}
\end{aligned}$$

and,

$$\begin{aligned}
\hat{\mathbf{u}} &= \hat{f}(\hat{\mathbf{s}}) \\
&= \hat{f}(W\mathbf{x})
\end{aligned}$$

Thus,

$$\begin{aligned}
P(\hat{\mathbf{u}}) &= P(\mathbf{x})d\mathbf{x} \\
&= \left| \frac{d\mathbf{x}}{d\hat{\mathbf{u}}} \right| P(\mathbf{x}) \\
&= \frac{1}{|M|} P(\mathbf{x})
\end{aligned}$$

where  $|M|$  is functional determinant and

$$\begin{aligned}
|M| &= \left| \frac{\partial \hat{\mathbf{u}}}{\partial \mathbf{x}} \right| \\
&= |W| \prod_{l=1}^N \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k \right)
\end{aligned}$$

Hence,

$$\begin{aligned}
H(\mathbf{u}) &= - \int d\hat{\mathbf{u}} P_{\mathbf{u}}(\hat{\mathbf{u}}) \log P_{\mathbf{u}}(\hat{\mathbf{u}}) \\
&= - \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log \frac{P_{\mathbf{x}}(\hat{\mathbf{x}})}{|M|} \\
&= - \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log P_{\mathbf{x}} + \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log |M|
\end{aligned}$$

Because  $-\int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log P_{\mathbf{x}}$  is a constant if we optimize over  $w$ , therefore, the cost function  $E^G$  is reduced to

$$\begin{aligned}
E^G &= \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log |M| \\
&= \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log |W| + \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \log \prod_{l=1}^N \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k \right) \\
&= \log |W| + \int d\mathbf{x} P_{\mathbf{x}}(\mathbf{x}) \left\{ \sum_{l=1}^N \log \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k \right) \right\}
\end{aligned}$$

Using Empirical Risk Minimization, we convert expectation to average thus

$$\begin{aligned}
E^T &= \log |W| + \frac{1}{p} \sum_{\alpha=1}^p \sum_{l=1}^N \log \hat{f}' \left( \sum_{k=1}^N w_{lk} x_k^{(\alpha)} \right) \\
&\stackrel{!}{=} \max_W
\end{aligned}$$

We can learn  $W$  using gradient ascent :

$$W \leftarrow W + \eta \nabla_W E^T$$

For individual cost of each observation  $e^{(\alpha)}$ ,

$$\frac{\partial e^{(\alpha)}}{\partial w_{ij}} = W_{ji}^{-1} + \frac{\hat{f}''\left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)}\right)}{\hat{f}'\left(\sum_{k=1}^N w_{ik} x_k^{(\alpha)}\right)} x_j^{(\alpha)}$$

Or

$$\frac{\partial}{\partial W} e^{(\alpha)} = (W^{-1})^T + \psi(Wx^{(\alpha)})(x^{(\alpha)})^T$$

One could observe that there is  $W^{-1}$  in the computation which causes performance.

### Learning via Natural Gradient

Due to the fact that  $W$  space is not orthogonal coordinate, hence gradient doest point to the direction of steepest ascent. Thus, we need to transform the space back to comparable space before computing gradient. This yields

$$dZ = dW \cdot W^{-1}$$

Using Taylor expansion of  $e^{(\alpha)}$  or  $e$  for short, around  $W$ , we have :

$$\begin{aligned} e(W + dW) &= e(W) + \nabla_W e^T dW \\ &= e(W) + \eta \nabla_W e^T z_w \end{aligned}$$

where  $z_w = \eta dw$

To find the direction of steepest ascent, we have to

$$\nabla_W e^T z_w \stackrel{!}{=} \max \quad \text{s.t.} \quad (z_w \cdot W^{-1})^2 \stackrel{!}{=} 1$$

Using Lagrange multiplier we have

$$\mathcal{L}(z_w, \lambda) = \nabla_W e^T z_w - \lambda((z_w \cdot W^{-1})^2 - 1)$$

Taking derivative wrt to  $z_w$  and set to zero, we have

$$\begin{aligned} 0 &= \nabla_W e - 2\lambda(z_w \cdot W^{-1})(W^{-1})^T \\ \nabla_W e &= 2\lambda(z_w \cdot W^{-1})(W^{-1})^T \\ z_w &= \frac{1}{2\lambda} \nabla_W e W^T W \end{aligned}$$

Hence the direction for **Natural Gradient** is

$$\Delta W = \eta \nabla_W e W^T W$$

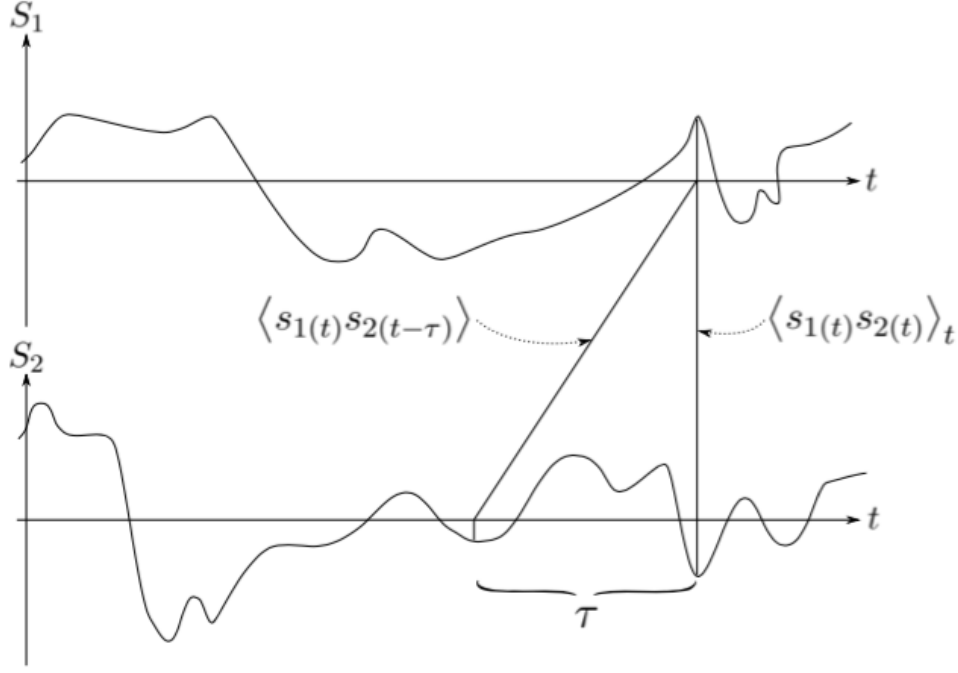
where  $\nabla_W e$  is original gradient.

Putting things together, we have

$$\Delta W = \eta \left( \mathbb{I} + \psi(Wx)(Wx)^T \right) W$$

Commonly, people use the sigmoid function as CDF :

$$f(x) = \frac{1}{1 + \exp(-x)} \quad \rightarrow \quad \psi(x) = \frac{f''(x)}{f'(x)} = 1 - 2f(x)$$

Figure 1.2: Autocorrelation with time  $\tau$  shifted (Drawn from MI2 Slides)

### 1.3.2 Second-order Approach

Due to the fact that, higher moment measurements are sensitive to noise, especially when we have small set of samples, hence it might yield unreliable result. Moreover, signals, such as EEG, fMRI data, might correlate to each other in temporal pattern. Because most of the time we know the length of the signal, we can exploit this fact by using autocorrelation.

Formally, the motivation of 2nd order approach is to find unmixing matrix  $W$  that cross-correlation vanishes.

#### Source Separation with 2 shifts

1. Whitening data

$$\mathbf{x}_c = \mathbf{x} - \frac{1}{p} \sum_{\alpha} \mathbf{x}^{(\alpha)}$$

Then compute eigen value decomposition of  $C_x^{(0)}$  where  $(C_x)_{ij} = \langle x_i, x_j \rangle$ . Hence,

$$\begin{aligned} \mathbf{u} &= \Lambda_0^{-1} \cdot E_0 \mathbf{x} \\ &= M_0 \mathbf{x} \end{aligned}$$

where  $\Lambda_0^{-1}$  is a diagonal of inverse eigen values and  $E_0$  is matrix of eigen vectors.

2. Orthogonal Transformation

Assume :  $\mathbf{s} = B\mathbf{u}$ . We know that

$$\begin{aligned} \langle s_i s_j \rangle &= \delta_{ij} \\ &= \sum_{k,l=1}^N B_{ik} \langle u_k u_l \rangle B_{l,j}^T \\ &= \sum_{k,l=1}^N B_{ik} B_{l,j}^T \quad (\langle u_k u_l \rangle = 1 \text{ from Whitening}) \end{aligned}$$

Hence,

$$BB^T = B^T B = \mathbb{I} \rightsquigarrow \text{Orthogonal Transformation}$$

3. Find a candidate of  $B$  using time-shifted cross-correlation matrix. Apply eigen value decomposition on  $C_u^{(\tau)}$  where

$$(C_u^{(\tau)})_{ij} = \langle u_i^{(t)} u_j^{(t-\tau)} \rangle$$

This results in a matrix of eigen vectors  $E_\tau$  whose  $E_\tau^T E_\tau = \mathbb{I}$ . Putting things together, we have

$$\hat{\mathbf{s}} = E_\tau \Lambda_0^{-1} E_0 \mathbf{x}$$

### Noise robust algorithms in general case

Given a set of  $T$  zero mean observations :  $\mathbf{x}^{(t)}$ . We compute joint cross-correlation metric  $C_x^{(\tau)}$  as :

$$(C_x^{(\tau)})_{ij} = \frac{1}{T} \sum_{t=0}^{T-1} x_i^{(t)} x_j^{(t-\tau)}$$

To find  $N \times N$  matrix  $W$  that diagonalizes  $C^{(\tau)}$ ,  $\tau = 0, 1, \dots$ , we use sum of off-diagonal entries of  $W C^{(\tau)} W^T$  as a cost function.

1. QDIAG Algorithm

$$E_W^T = \sum_{\tau} \alpha_{\tau} \sum_{i \neq j} \left( W C^{(\tau)} W^T \right)_{ij}^2$$

where  $\alpha_{\tau}$  is a weighting factor and the optimization problem is  $E^T \stackrel{!}{=} \min$  s.t.

$$\forall i, \quad \left( W C^{(0)} W^T \right)_{ii} = 1$$

2. FFDIAG Algorithm

$$E_W^T = \sum_{\tau} \sum_{i \neq j} \left( W C^{(\tau)} W^T \right)_{ij}^2$$

with a constraint that  $W$  is a non-singular matrix.

### 1.3.3 Maximize Non-Gaussianity

According to Central Limit Theorem, the sum of random variables is more Gaussian than the original sources. Hence, we can leverage this knowledge by finding a unmixing matrix that minimizes Gaussianity of the unmixed sources yielding interesting directions (projection pursuit). There are 2 ways to measure Gaussianity, namely Kurtosis (4th moment) and Negentropy.

#### Kurtosis

$$\text{kurt}(x) = \langle x^4 \rangle - 3 \underbrace{(\langle x^2 \rangle)^2}_{\substack{1 \text{ for} \\ \text{sphered data}}}$$

- $\text{kurt}(x) = 0$  for Gaussian distribution
- $\text{kurt}(x) > 0$  for **Super**-Gaussian : long-tail, e.g. Laplace Distribution
- $\text{kurt}(x) < 0$  for **Sub**-Gaussian : constant distribution.

For independent variables  $x$  and  $y$  we have :

$$\begin{aligned} \text{kurt}(x + y) &= \text{kurt}(x) + \text{kurt}(y) \\ \text{kurt}(ax) &= a^4 \text{kurt}(x) \end{aligned}$$

Because,

$$\hat{\mathbf{s}} = W \mathbf{x} = W A \mathbf{s} = b^T \tilde{A} \mathbf{s} = b^T \mathbf{u}$$

Hence, the optimization problem can be written as

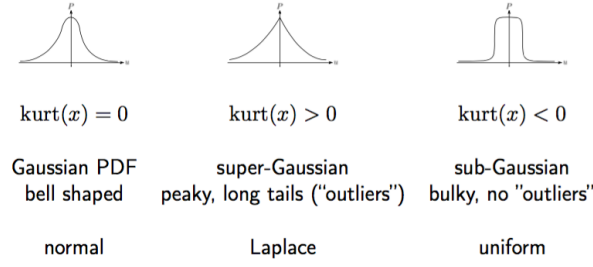


Figure 1.3: Distribution and Kurtosis (Drawn from MI2 Slides)

$$\max_{\mathbf{b}} |\text{kurt}(\mathbf{b}^T \mathbf{u})| \quad \text{s.t.} \quad |\tilde{A}^T \mathbf{b}| = |\mathbf{b}| = 1$$

We have :

$$\begin{aligned} \frac{\partial}{\partial \mathbf{b}} |\text{kurt}(\mathbf{b}^T \mathbf{u})| &= 4 \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \left( \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle - 3\mathbf{b} \langle (\mathbf{b}^T \mathbf{u})^2 \rangle \right) \\ &\stackrel{!}{=} 0 \quad \text{s.t.} \quad |\mathbf{b}| = 1 \end{aligned}$$

Due to the fact that

$$\begin{aligned} 3\mathbf{b} \langle (\mathbf{b}^T \mathbf{u})^2 \rangle &= 3\mathbf{b} \langle ((\mathbf{b}^T \mathbf{u})(\mathbf{b}^T \mathbf{u})^T) \rangle \\ &= 3\mathbf{b} \langle ((\mathbf{b}^T \mathbf{u})(\mathbf{u}^T \mathbf{b})) \rangle \\ &= 3\mathbf{b} \mathbf{b}^T \langle \mathbf{u} \mathbf{u}^T \rangle \mathbf{b} \\ &= 3\mathbf{b} \mathbf{b}^T \mathbb{I} \mathbf{b} \\ &= 3\mathbf{b} |\mathbf{b}|^2 \\ &= 3\mathbf{b} \end{aligned} \quad (|\mathbf{b}|^2 = 1)$$

which changes only direction of  $\mathbf{b}$ , we can drop it and hence the problem is reduced to

$$\begin{aligned} 4 \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle &\stackrel{!}{=} 0 \\ \epsilon \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle &\approx 0 \end{aligned}$$

Therefore for **batch** learning, we have

$$\begin{aligned} \Delta \mathbf{b} &= \epsilon \text{sgn}[\text{kurt}(\mathbf{b}^T \mathbf{u})] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle \\ \mathbf{b} &\leftarrow \mathbf{b} / |\mathbf{b}| \end{aligned}$$

where  $\text{kurt}$  and  $\langle \cdot \rangle$  are replaced with corresponding empirical method and  $\mathbf{b}_0$  is a arbitrary vector with unit length.

For **online** approach, we use running average of kurt  $\gamma$  where  $\gamma_0 = 0$ .

$$\begin{aligned} \Delta \mathbf{b} &= \epsilon \text{sgn}[\gamma] \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle \\ \Delta \gamma &= \eta [(\mathbf{b}^T \mathbf{u})^4 - 3 - \gamma] \\ \mathbf{b} &\leftarrow \mathbf{b} / |\mathbf{b}| \end{aligned}$$

At the equilibrium point of gradient ascent, one could observe that

$$\begin{aligned} \mathbf{b} &\propto \Delta \mathbf{b} \\ &\rightsquigarrow \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle - 3\mathbf{b} \end{aligned}$$

This yields fixed-point algorithm ( **Kurtosis-based fastICA** )

$$\begin{aligned} \mathbf{b} &\leftarrow \langle \mathbf{u}(\mathbf{b}^T \mathbf{u})^3 \rangle - 3\mathbf{b} \\ \mathbf{b} &\leftarrow \mathbf{b} / |\mathbf{b}| \end{aligned}$$

For whitened data  $\mathbf{u}^{(\alpha)}, \alpha = 1, \dots, p$ , we have

$$\begin{aligned}\mathbf{b} &\leftarrow \frac{1}{p} \sum_{\alpha=1}^p \mathbf{u}^{(\alpha)} (\mathbf{b}^T \mathbf{u}^{(\alpha)})^3 - 3\mathbf{b} \\ \mathbf{b} &\leftarrow \mathbf{b}/|\mathbf{b}|\end{aligned}$$

One remark is that **kurtosis** tends to be very sensitive to outliers.

#### 1.3.4 Negentropy