

# Chapter 4

# Independent

# Component Analysis

Part II: Algorithms



# ICA definition

- Given  $n$  observations of  $m$  random variables in matrix  $\mathbf{X}$ , find  $n$  observations of  $m$  independent components in  $\mathbf{S}$  and  $m$ -by- $m$  invertible mixing matrix  $\mathbf{A}$  s.t.  $\mathbf{X} = \mathbf{SA}$ 
  - Components are statistically independent
  - At most one is Gaussian
  - We can assume  $\mathbf{A}$  is orthogonal (by whitening  $\mathbf{X}$ )

# Maximal non- Gaussian

# Central limit theorem

- Average of i.i.d. variables converges to normal distribution
  - $\sqrt{n} \left( \left( \frac{1}{n} \sum_{i=1}^n X_i \right) - \mu \right) \xrightarrow{d} N(0, \sigma^2)$  as  $n \rightarrow \infty$
- Hence  $(X_1 + X_2)/2$  is “more Gaussian” than  $X_1$  or  $X_2$  alone
  - For i.i.d. zero-centered non-Gaussian  $X_1$  and  $X_2$
- Hence, we can try to find components  $s$  that are “maximally non-Gaussian”

# Re-writing ICA

- Recall, in ICA  $\mathbf{x} = \mathbf{sA} \Leftrightarrow \mathbf{s} = \mathbf{xA}^{-1}$ 
  - Hence,  $s_j$  is a linear combination of  $x_i$
- Approximate  $s_j \approx y = \mathbf{xb}^T$  ( $\mathbf{b}$  to be determined)
  - Now  $y = \mathbf{sAb}^T$  so  $y$  is a lin. comb. of  $\mathbf{s}$
  - Let  $\mathbf{q}^T = \mathbf{Ab}^T$  and write  $y = \mathbf{xb}^T = \mathbf{sq}^T$

# More re-writings

- Now  $s_j \approx y = \mathbf{x}\mathbf{b}^T = \mathbf{s}\mathbf{q}^T$
- If  $\mathbf{b}^T$  is a column of  $\mathbf{A}^{-1}$ ,  $s_j = y$  and  $q_j = 1$  and  $\mathbf{q}$  is 0 elsewhere
- CLT:  $\mathbf{s}\mathbf{q}^T$  is least Gaussian when  $\mathbf{q}$  looks correct
  - We don't know  $\mathbf{s}$ , so we can't vary  $\mathbf{q}$
  - But we can vary  $\mathbf{b}$  and study  $\mathbf{x}\mathbf{b}^T$
- **Approach:** find  $\mathbf{b}$  s.t.  $\mathbf{x}\mathbf{b}^T$  is least Gaussian

# Kurtosis

- One way to measure how Gaussian a random variable is is its **kurtosis**
- $\text{kurt}(y) = E[(y - \mu)^4] - 3(E[(y - \mu)^2])^2$ 
  - $E[y] = \mu$
  - Normalized version of the fourth central moment  $E[(y - \mu)^4]$
- If  $y \sim N(\mu, \sigma^2)$ ,  $\text{kurt}(y) = 0$ , most other distributions have non-zero kurtosis (positive or negative)

# Computing with kurtosis

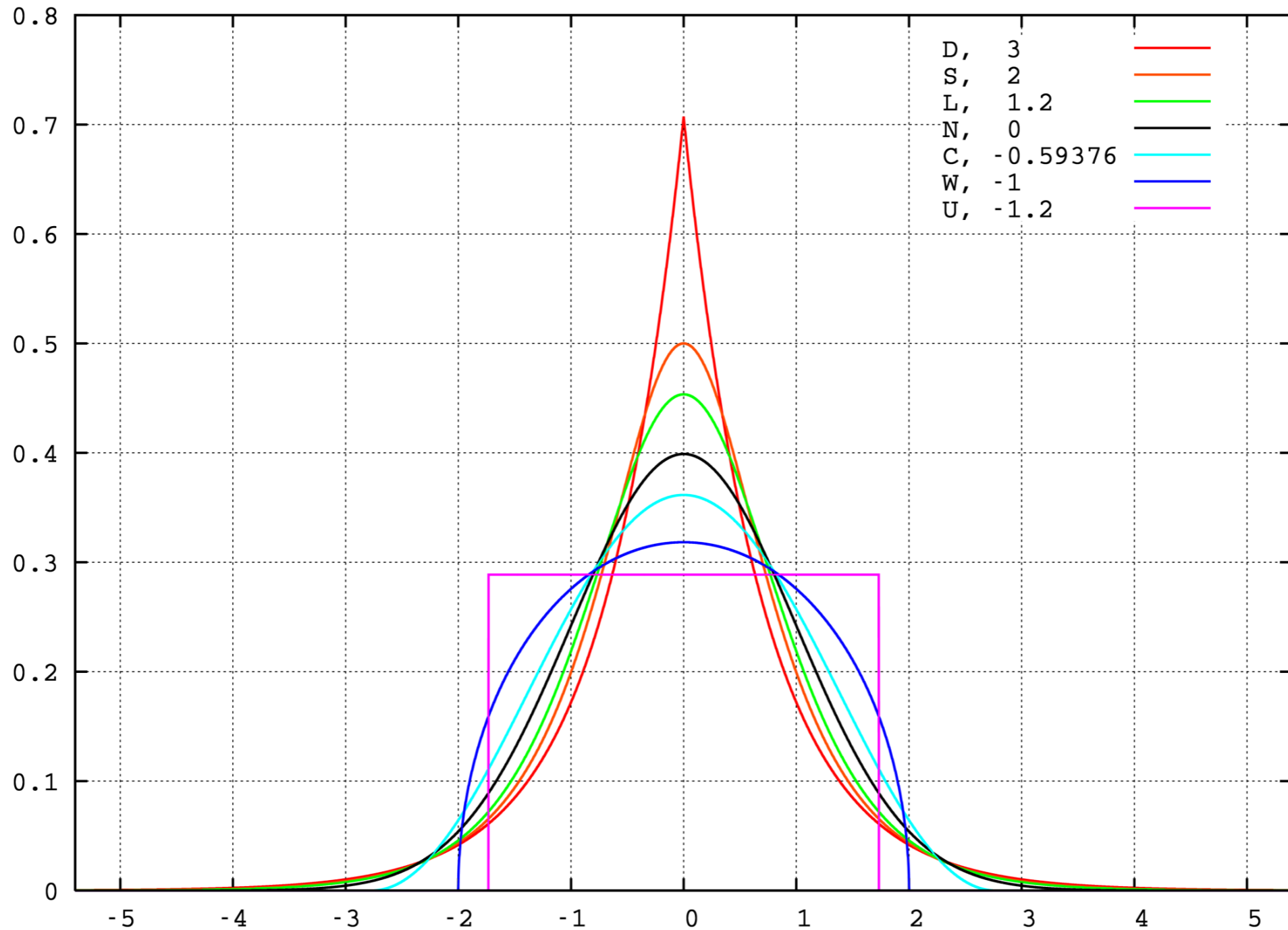
- If  $x$  and  $y$  are independent random variables:
  - $\text{kurt}(x + y) = \text{kurt}(x) + \text{kurt}(y)$
  - Homework
- If  $\alpha$  is a constant:
  - $\text{kurt}(\alpha x) = \alpha^4 \text{kurt}(x)$
  - $E[(\alpha x)^4] - 3(E[(\alpha x)^2])^2 = \alpha^4 E[x^4] - \alpha^4 3(E[x^2])^2$



# Sub- and super-Gaussian distributions

- Distributions with negative kurtosis are **sub-Gaussian** (or **platykurtic**)
  - Flatter than Gaussian
- Distributions with positive kurtosis are **super-Gaussian** (or **leptokurtic**)
  - Spikier than Gaussian

# Examples



[https://en.wikipedia.org/wiki/Kurtosis#/media/File:Standard\\_symmetric\\_pdfs.png](https://en.wikipedia.org/wiki/Kurtosis#/media/File:Standard_symmetric_pdfs.png)

# Back to optimization

- Recall: with two components

$$y = \mathbf{x}\mathbf{b}^T = \mathbf{s}\mathbf{q}^T = q_1s_1 + q_2s_2$$

- $s_i$  have unit variance
- We want to find  $\pm\mathbf{b} = \operatorname{argmax} |\operatorname{kurt}(\mathbf{x}\mathbf{b}^T)|$ 
  - We can't determine the sign
- We want  $y$  to be either  $s_1$  or  $s_2$ , hence
$$E[y^2] = q_1^2 + q_2^2 = 1$$

# Whitening, again

- Generally,  $\|\mathbf{q}\|^2 = 1$
- Recall:  $\mathbf{Z} = \mathbf{U} = \mathbf{XV}\boldsymbol{\Sigma}^{-1}$  is the whitened version of  $\mathbf{X}$
- Target becomes  $\pm\mathbf{w} = \operatorname{argmax} |\operatorname{kurt}(\mathbf{z}\mathbf{w}^T)|$
- Now  $\|\mathbf{q}\|_2^2 = (\mathbf{w}\mathbf{U}^T)(\mathbf{U}\mathbf{w}^T) = \|\mathbf{w}\|_2^2$ 
  - Hence we have constraint  $\|\mathbf{w}\|^2 = 1$

# Gradient-based algorithm

- Gradient is

$$\frac{\partial |\text{kurt}(\mathbf{z}\mathbf{w}^T)|}{\partial \mathbf{w}} = 4 \text{sign}(\text{kurt}(\mathbf{z}\mathbf{w}^T))(E[(\mathbf{z}\mathbf{w}^T)^3 \mathbf{z}] - 3\mathbf{w} \|\mathbf{w}\|_2^2)$$

- $E[(\mathbf{z}\mathbf{w}^T)^2] = \|\mathbf{w}\|^2$  for whitened data
- We can optimize this using standard gradient methods
- To satisfy the constraint  $\|\mathbf{w}\|^2 = 1$ , we divide  $\mathbf{w}$  with its norm after every update

# FastICA for one IC

- Noticing that  $\|\mathbf{w}\|^2 = 1$  by constraint and taking infinite step update, we get  
$$\mathbf{w} \leftarrow E[(\mathbf{z}\mathbf{w}^T)^3\mathbf{z}] - 3\mathbf{w}$$
- Again set  $\mathbf{w} \leftarrow \mathbf{w}/\|\mathbf{w}\|$  after every update
- Expectation has naturally to be estimated
- No theoretical guarantees but works in practice

# Multiple components

- So far we have found only one component
  - To find more, remember that vectors  $\mathbf{w}_i$  are orthogonal (whitening!)
- General idea:
  - Find one vector  $\mathbf{w}$
  - Find second that is orthogonal to the first one
  - Find third that is orthogonal to the two previous ones, etc.

# Symmetric orthogonalization

- We can compute  $\mathbf{w}_i$ s in parallel
  - Update  $\mathbf{w}_i$ s independently
  - Run orthogonalization after every update step
    - $\mathbf{W} \leftarrow (\mathbf{W}\mathbf{W}^T)^{-1/2}\mathbf{W}$
- Iterate until convergence



# Maximum Likelihood

# Maximum-likelihood algorithms

- **Idea:** We are given observations  $\mathbf{X}$  that are drawn from some parameterized family of distributions  $D(\Theta)$ 
  - The **likelihood** of  $\mathbf{X}$  given  $\Theta$ ,  $L(\Theta; \mathbf{X}) = p_D(\mathbf{X}; \Theta)$ , where  $p_D(\cdot; \Theta)$  is the probability density function of  $D$  with parameters  $\Theta$
- In **maximum-likelihood estimation** (MLE) we try to find  $\Theta$  that maximizes the likelihood given  $\mathbf{X}$

# ICA as MLE

- If  $p_{\mathbf{x}}(\mathbf{x})$  is the pdf of  $\mathbf{x} = \mathbf{s}\mathbf{A}$ , then

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{s}}(\mathbf{s}) |\det \mathbf{B}| = |\det \mathbf{B}| \prod_i p_i(s_i) = |\det \mathbf{B}| \prod_i p_i(\mathbf{x}\mathbf{b}_i^T)$$

- $\mathbf{B} = \mathbf{A}^{-1}$

- For  $t$  observations  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$  the log-likelihood of  $\mathbf{B}$  given  $\mathbf{X}$  is

$$\log L(\mathbf{B}; \mathbf{X}) = \sum_{t=1}^T \sum_{i=1}^m \log p_i(\mathbf{x}_t \mathbf{b}_i^T) + T \log |\det \mathbf{B}|$$

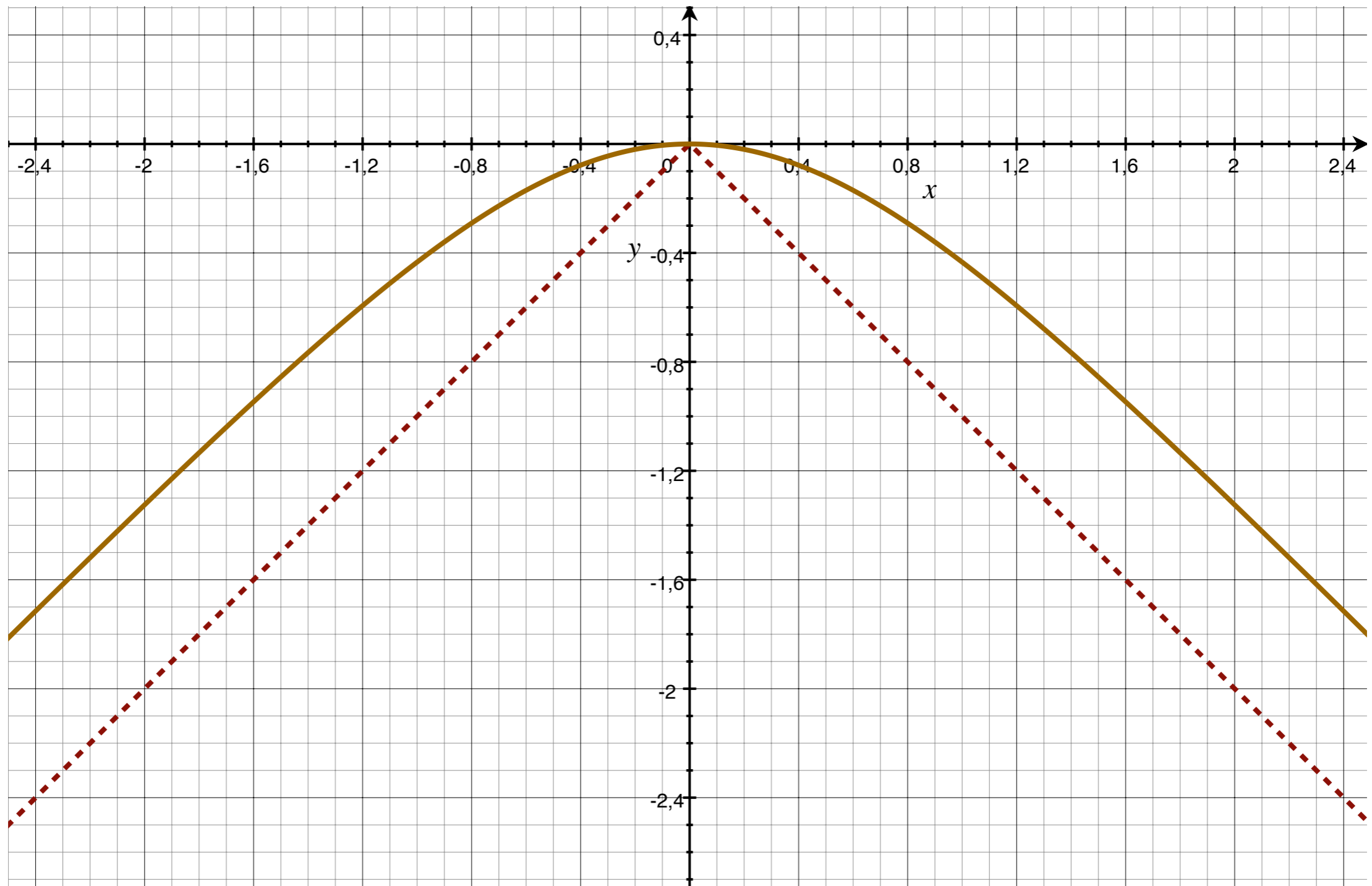
# Problems with MLE

- The likelihood is expressed as a function of  **$B$**
- But we also need to estimate the pdfs  $p_i()$ 
  - Non-parametric problem, infinite number of different pdfs
- Very hard problem...

# If we know the pdfs

- Sometimes we know the pdfs of the components
  - We only need to estimate their parameters and  **$B$**
- Sometimes we know only that the pdfs are super-Gaussian (for example)
  - We can use  $\log p_i(s_i) = -\log \cosh(s_i)$ 
    - Requires normalization

$$-\log \cosh(x) \approx -|x|$$



# Nothing on the pdfs is known

- We might not know whether the pdfs of the components are sub- or super-Gaussian
  - It is enough to estimate which one they are!
- For super-Gaussian,  
$$\log p_i^+(s_i) = \alpha_1 - 2\log \cosh(s_i)$$
- For sub-Gaussian,  
$$\log p_i^-(s_i) = \alpha_2 - (s_i^2/2 - \log \cosh(s_i))$$

$\alpha_i$  are only needed to make these logs of pdfs – not in optimization

# Log-likelihood gradient

- The gradient is  $\frac{\partial \log L}{\partial \mathbf{B}} = (\mathbf{B}^T)^{-1} + \sum_{t=1}^T \mathbf{g}(\mathbf{x}_t \mathbf{B}^T)^T \mathbf{x}_t$ 
  - Here  $\mathbf{g}(\mathbf{y}) = (g_i(y_i))_{i=1}^n$  with  $g_i(y_i) = (\log p_i(y_i))' = p_i'(y_i)/p_i(y_i)$
- This gives us  $\mathbf{B} \leftarrow \mathbf{B} + \delta((\mathbf{B}^T)^{-1} + \sum_t \mathbf{g}(\mathbf{x}_t \mathbf{B}^T)^T \mathbf{x}_t)$
- Multiplying from right with  $\mathbf{B}^T \mathbf{B}$  and defining  $\mathbf{y}_t = \mathbf{x}_t \mathbf{B}^T$  gives  $\mathbf{B} \leftarrow \mathbf{B} + \delta(\mathbf{I} + \sum_t \mathbf{g}(\mathbf{y}_t)^T \mathbf{y}_t) \mathbf{B}$ 
  - So-called **infomax** algorithm



# Setting $g()$

- We compute  $E[-\tanh(s_i)s_i + (1 - \tanh(s_i)^2)]$ 
  - If positive, set  $g(y) = -2\tanh(y)$
  - If negative (or zero), set  $g(y) = \tanh(y) - y$
- Use current estimates of  $s_i$

# Putting it all together

- Start with random  $\mathbf{B}$  and  $\gamma$ , choose learning rates  $\delta$  and  $\delta_\gamma$
- Iterate until convergence
  - $\mathbf{y} \leftarrow \mathbf{B}\mathbf{x}$  and normalize  $\mathbf{y}$  to unit variance
  - $\gamma_i \leftarrow (1 - \delta_\gamma)\gamma_{i-1} + \delta_\gamma E[-\tanh(y_i)y_i + (1 - \tanh(y_i)^2)]$ 
    - if  $\gamma_i > 0$ , use super-Gaussian  $g$ ; o/w sub-Gaussian  $g$
  - $\mathbf{B} \leftarrow \mathbf{B} + \delta(\mathbf{I} + \sum_t \mathbf{g}(\mathbf{y}_t)^T \mathbf{y}_t)\mathbf{B}$

# ICA summary

- ICA can recover independent source signals
  - if they are non-Gaussian
- Does not reduce rank
- Many applications, special case of blind source separation
- Standard algorithmic technique is to maximize non-Gaussianity of the recovered components

# ICA literature

- Hyvärinen & Oja (2000): *Independent Component Analysis: Algorithms and Applications*. Neural networks 13(4), 411–430
- Hyvärinen (2013): *Independent component analysis: recent advances*. Phil. Trans. R. Soc. A 371:20110534