

# The Development of NonGaussian Component Analysis

Elmar Diederichs<sup>2</sup>

joint work with Anatoli Juditsky<sup>3</sup> and Vladimir Spokoiny<sup>1</sup>

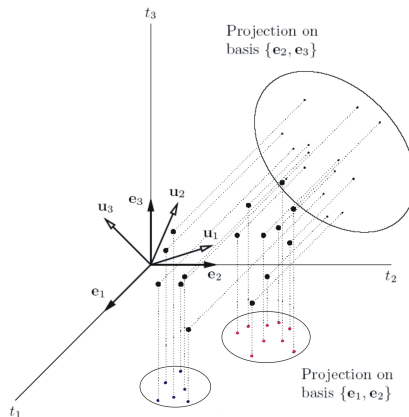
<sup>1</sup>Weierstrass-Institute and Humbolt University, Berlin

<sup>2</sup>Moscow Institute for Physics and Technology, Moscow

<sup>3</sup>U. Joseph Fourier, Grenoble

International Workshop on Statistical Learning  
Moscow, 26.-27.05.2013

# Structural Analysis



# Design of Unsupervised Feature Extraction

Data  $X_1, \dots, X_n \in \mathbb{R}^d$  i.i.d.,  $d$  large. For simplicity let  $\mathbf{E}[X_i] = 0$  for all  $i$ .

**Problem:** most of all random projections  $X^\top \omega$  are almost approximately normal

**Approach:** **Gaussian component** of the data is entropy-maximizing and hence **uninformative** (noise). Project the data on the **non-Gaussian components**.

**Requirements for an acceptable statistical method:**

- i) No apriori knowledge about the data density is used.
- ii) No **dependency on the magnitude of second moments** of Gaussian and non-Gaussian components as found e.g. in PCA or **unrealistic assumptions on the data density** as found e.g. in ICA.

# The Semi-Parametric Model

**Semiparametric structural assumption:**

$$\rho(x) = \phi_{\mu=0, \Sigma}(x) q(Tx) \quad (1)$$

This links pure Gaussian Analysis (PCA) and pure NonGaussian Analysis (ICA).

$q : \mathbb{R}^m \rightarrow \mathbb{R}$  smooth nonlinear function,  $m \leq d$

$T : \mathbb{R}^d \rightarrow \mathbb{R}^m$  linear operator with  $\mathcal{I} = \text{Ker}(T)^\perp$ .

$\mathcal{I}$ : linear **target subspace** of the non-Gaussian components.

**goal:** Estimate a projector without estimating the model parameter.

**interpretation:** (1) lead to the stationary data model  $X = Z + \zeta$  where  $\zeta$  represents independent Gaussian noise components and  $Z$  the signal.

# Recovery of Target Space

## Lemma

Assume that  $\rho(x)$  follows the semiparametric assumption with  $\mathbf{E}[X] = 0$ . Then for  $\psi(x) \in \mathcal{C}^1(\mathbb{R}^d, \mathbb{R})$  and

$$\beta(\psi) \stackrel{\text{def}}{=} \mathbf{E}[\nabla \psi(x)] \quad (2)$$

there is  $\beta \in \mathcal{I}$  such that there is a uniform error bound

$$\|\beta(\psi) - \beta\|_2 \leq \left\| \Sigma^{-1} \mathbf{E}[x\psi(x)] \right\|_2 \quad (3)$$

of  $\text{dist}(\beta(\psi), \mathcal{I})$ . Moreover if  $\mathbf{E}[x\psi(x)] = 0$ , then  $\beta(\psi) \in \mathcal{I}$ .

# Unsupervised Feature Extraction Using Projections

- 1) linear approach:  $\psi(x) = \sum_{\ell}^L c_{\ell} h_{\ell}(x)$
- 2) test functions:  $h_{\ell}(x) \stackrel{\text{def}}{=} h(\omega_{\ell}^{\top} x) e^{-\lambda \|x\|^2/2}$
- 3) define:

$$\gamma_{\ell} \stackrel{\text{def}}{=} E[X h_{\ell}(X)], \quad \eta_{\ell} \stackrel{\text{def}}{=} E[\nabla h_{\ell}(X)],$$

and let  $\hat{\gamma}_{\ell}$  and  $\hat{\eta}_{\ell}$  be their "empirical counterparts" that can be computed, such that for a set  $A$  of probability at least  $1 - \epsilon$  it holds  $\max_{|c|_1 \leq 1} \sum_{\ell} |\hat{\eta}_{\ell} - \eta_{\ell}|_2 \leq \delta_N$  and  $\max_{|c|_1 \leq 1} \sum_{\ell} |\hat{\gamma}_{\ell} - \gamma_{\ell}|_2 \leq \nu_N$ .

## Lemma

Let  $h$  be bounded and continuously differentiable. For a fixed constant  $C = C(h)$ , it holds

$$E \sup_{\omega \in \mathcal{B}_d} |\hat{\gamma}_{\omega} - \gamma_{\omega}|^2 + |\hat{\eta}_{\omega} - \eta_{\omega}|^2 \leq CN^{-1/2} \sqrt{d}$$

# NonGaussian Component Analysis

Approach of Blanchard et al.(JMLR, 2007):

Consider  $\psi(x) = h(x) - \alpha^T x$  and select  $\alpha$  s.t.

$$E[X\psi(X)] = E[Xh(X)] - EXX^T\alpha = 0.$$

in order to suppress the noise when estimating elements from  $\mathcal{I}$ .

Then  $\beta(\psi) \stackrel{\text{def}}{=} E[\nabla\psi(x)]$  leads to:

$$\hat{\beta}_\ell = \hat{\eta}_\ell - \hat{\Sigma}^{-1}\hat{\gamma}_\ell.$$

**drawbacks:** requires to compute and study  $\hat{\Sigma}^{-1}$ ,  $m$  cannot be estimated.

# Sparse NonGaussian Component Analysis

Approach of Diederichs et al. (IEEE Trans. Inf. Theo, 2009):

Set  $\psi(x) = \sum_{\ell} c_{\ell} h_{\ell}(x)$  and consider the convex projection problem

$$\hat{c} = \arg \min_c \left\{ \|\xi - \sum_{\ell} c_{\ell} \hat{\eta}_{\ell}\|_2 \mid \sum_{\ell} c_{\ell} \hat{\gamma}_{\ell} = 0, \|c\|_1 \leq 1 \right\}$$

and define  $\hat{\beta}_{\ell} = \sum_{\ell} \hat{c}_{\ell} \hat{\eta}_{\ell}$ . Then under some regularity conditions there is  $\mathcal{C} \stackrel{\text{def}}{=} \{ \|c\|_1 \leq 1, \sum_{\ell} c_{\ell} \hat{\gamma}_{\ell} \}$

$$\|(\mathbf{1}_d - \Pi^*)\hat{\beta}\|_2 \leq \sqrt{d}C(d, N^{-1/2})(1 + \|\Sigma^{-1}\|_2)$$

**drawbacks:** choice of informative probe vector  $\xi$ ,  $m$  cannot be estimated

Both approaches require the solution of the **Reduced Rank Regression** problem: given  $m$ , recover  $\mathcal{I}$  or the projector  $\Pi_{\mathcal{I}}$  from  $\hat{\beta}_1, \dots, \hat{\beta}_L$ .



## Dimension Reduction Step: the RRR problem

Suppose to be given the vectors  $\hat{\beta}_1, \dots, \hat{\beta}_L$  such that

$$\|(\mathbf{1}_d - \Pi_{\mathcal{I}})\hat{\beta}_\ell\| \leq \epsilon$$

where  $\Pi_{\mathcal{I}}$  is a projector on a  $m$ -dimensional subspace.

PCA solution:

$$\hat{\mathcal{I}} = \arg \min_{\dim(\mathcal{I})=m} \sum_{\ell} \|(\mathbf{1}_d - \Pi_{\mathcal{I}})\hat{\beta}_\ell\|^2 = \langle \text{first } m \text{ eigenvectors of } \sum_{\ell} \hat{\beta}_\ell \hat{\beta}_\ell^T \rangle.$$

Requires that  $\lambda_m(\sum_{\ell} \beta_{\ell} \beta_{\ell}^T) \geq L\epsilon^2$ . Works poorly if most of the  $\hat{\beta}_{\ell}$ 's are **non-informative**.

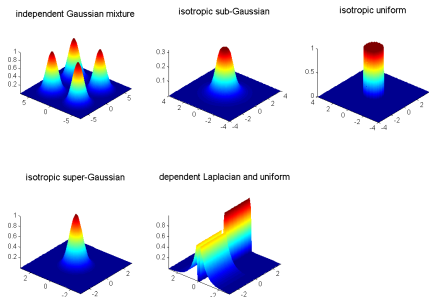
**Rounding Ellipsoid:** (see Yu.Nesterov, 2004) Define the convex set

$$\mathcal{S} := \{\hat{\beta}_1, -\hat{\beta}_1, \hat{\beta}_2, -\hat{\beta}_2, \dots\}.$$

For  $\mathcal{S}$  a centered ellipsoid  $\mathcal{E}$  of minimum volume that encloses  $\mathcal{S}$  always exists and recovering  $\mathcal{I}$  from the principal axis of  $\mathcal{E}$  comes with the accuracy

$$\|\Pi_{\mathcal{I}} - \Pi_{\hat{\mathcal{I}}}\|_2^2 \leq C(\lambda_m, \epsilon^2)d\sqrt{d}.$$

# Numerical Illustration of Progress



**Figure :** (A) NonGaussian test densities: 2 d independent Gaussian mixtures, (B) 2 d isotropic super-Gaussian, (C) 2 d isotropic uniform and (D) dependent 1 d Laplacian with additive 1 d uniform with  $N = 1000$  respectively.

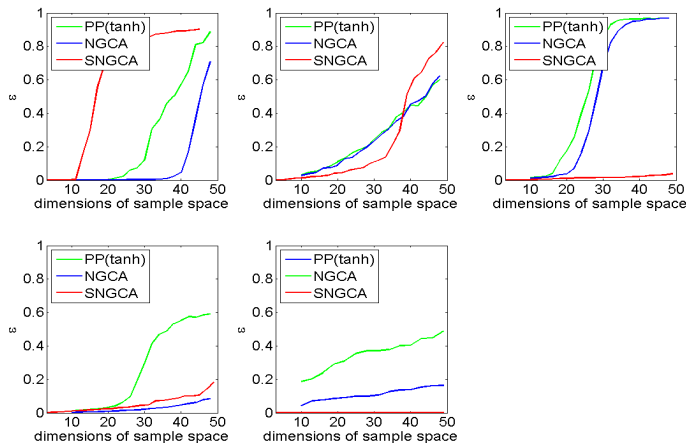
# Error Measure

The closeness of the subspace  $\mathcal{I}$  and its estimate  $\hat{\mathcal{I}}$  can be measured by the error function

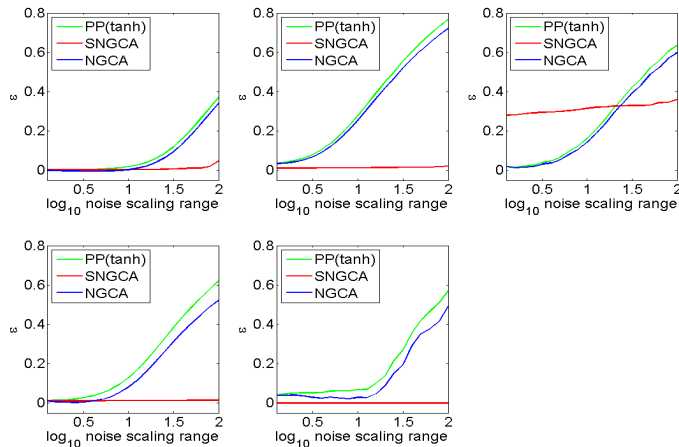
$$\mathcal{E}(\hat{\mathcal{I}}, \mathcal{I}) = \frac{1}{m} \sum_{i=1}^m \|(\mathbf{1}_d - \Pi) \mathbf{v}_i\|_2^2 \quad (4)$$

where  $\Pi$  denotes the orthogonal projection onto  $\hat{\mathcal{I}}$ ,  $\{\mathbf{v}_i\}_{i=1}^m$  is an orthonormal basis of  $\mathcal{I}$  and  $I$  denotes the identity matrix.

# Estimation Error for Increasing Dimensionality



# Comparison of Methods Cont'd: Noise



# Design of Semidefinite Structural Analysis

first step of approach of Juditsky et al. (JML, 2012): Choose randomly a family of directions  $\{\omega_\ell\}$ ,  $\ell = 1, \dots, L$  and compute

$$\hat{\gamma}_\ell \stackrel{\text{def}}{=} \mathbf{E}_N[Xh_\ell(X)], \quad \hat{\eta}_\ell \stackrel{\text{def}}{=} \mathbf{E}_N[\nabla h_\ell(X)]$$

second step: avoid any probe vectors and the RRR problem by solving the semidefinite problem

$$\Pi^* = \min_{\Pi} \max_c \left\{ \left\| (I - \Pi)Uc \right\|_2^2 \mid \begin{array}{l} \Pi \text{ is a projector on a} \\ m\text{-dimensional subspace of } \mathbb{R}^d \\ c \in \mathbb{R}^L, \ Gc = 0 \end{array} \right\}. \quad (5)$$

where  $U = [\eta_1, \dots, \eta_L] \in \mathbb{R}^{d \times L}$ ,  $G = [\gamma_1, \dots, \gamma_L] \in \mathbb{R}^{d \times L}$  and  $\Pi^*$  is the Euclidean projector on  $\mathcal{I}$ .

# Design of Semiparametric Structural Analysis cont'd

**third step:** structural adaptation idea (Hristache, Juditsky, Polzehl and Spokoiny, 2003):

use the estimated projector  $\hat{\Pi}_{k-1}$  as a prior information for the directional sampling to improve the quality of estimation at iteration  $k$  of SD-NGCA.

This leads to a **sequential procedure**: alternate two steps

- a) sample some directions  $\omega_\ell$  from the space spanned by the  $m$  principal directions of  $\hat{\Pi}_{k-1}$
- b) estimate the structure  $\hat{\Pi}_k$

This ensures that a certain fraction of  $\hat{\gamma}_\ell$  and  $\hat{\eta}_\ell$  is informative.

# Relaxation of the hard problem in SD-NGCA

**idea:** drop constraints to get convexity and solve an approximating problem

- i) Use:  $\|(I - \Pi)\hat{U}c\|_2^2 = \text{Tr} [\hat{U}(I - \Pi)\hat{U}cc^T]$ .
- ii) **Linearization:** positive semidefinite matrix  $X = cc^T$  with  $\text{rank}X = 1$  as "new variable".
- iii) Set  $|X|_1 \stackrel{\text{def}}{=} \sum_{i,j=1}^L |X_{ij}|$  and transform  $\|\hat{G}c\|_2 \leq \delta$  into  $\text{Tr}[\hat{G}X\hat{G}] \leq \varrho^2$ .
- iv) **Drop** the non-convex constraints  $\text{rank}X = 1$  and  $\text{rank}\Pi = m$ .

Then we arrive at the **relaxed** and constrained saddle point problem:

$$\min_P \max_X \left\{ \text{Tr} [\hat{U}(I - P)\hat{U}X] \mid \begin{array}{l} 0 \preceq P \preceq I, \text{Tr}[P] = m, \\ X \succeq 0, |X|_1 \leq 1, \text{Tr}[\hat{G}X\hat{G}] \leq \varrho^2 \end{array} \right\}. \quad (6)$$

where  $|X|_1 := \sum_{ij} |X_{ij}|$ . We call  $P$  - defined as above - a subprojector.



# Accuracy of the estimated projector

## Lemma

Let  $\hat{P}$  be an *optimal solution of the relaxed SDP* and assume that

- i)  $\Pi^*$  on  $\mathcal{I}$  is a convex combination of rank-one matrices  $Ucc^T U^T$
- ii)  $c$  satisfies  $Gc = 0$  and  $\|c\|_1 \leq 1$ .

Then it holds of  $\hat{\Pi}$ , spanned by the first  $m$  eigenvectors of  $\hat{P}$ , with probability  $\geq (1 - \epsilon)$ ;

$$\begin{aligned} \|(\mathbf{1}_d - \hat{\Pi})Uc\|_2 &\leq C_1 \sqrt{m+1} ((\varrho + \nu_N) \lambda_{\min}^{-1}(\Sigma) + 2\delta_N) \\ \text{Tr}[(\mathbf{1}_d - \hat{\Pi})\Pi^*] &\leq C_2 [(\varrho + \nu_N) \lambda_{\min}^{-1}(\Sigma) + 2\delta_N]^2 \\ \|\hat{\Pi} - \Pi^*\|_{Frob}^2 &\leq C_3 (m+1) [(\varrho + \nu_N) \lambda_{\min}^{-1}(\Sigma) + 2\delta_N]^2 \end{aligned}$$

where  $C_i = C_i(h)$  with  $i = 1, 2, 3$  does not depend on  $d$  or  $L$ .

# Linear Constraints

Observe that  $\widehat{G}^T \widehat{G} = \Gamma \Lambda \Gamma^T$  and  $X$  are symmetric and positive. Hence:

$$\text{Tr}(\widehat{G}^T \widehat{G} X) = 0 \quad \Rightarrow \quad X = QZQ^T \quad (7)$$

where  $Z \in \mathcal{S}^{L-d}$  and  $Q \in \mathcal{S}^{L \times (L-d)}$  is a submatrix of columns of  $\Gamma$  corresponding to the vanishing eigenvalues of  $\widehat{G}^T \widehat{G}$ .

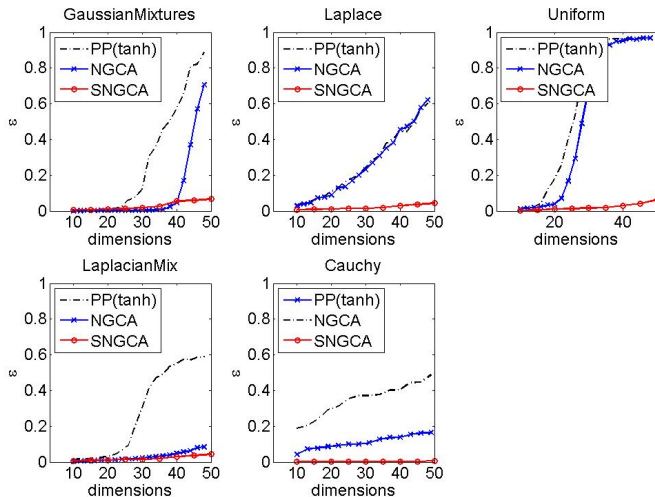
Let  $V = \widehat{G}Q$ . Then we get a **regularized** and hence **unconstrained convex reformulation** of the relaxed problem:

$$\min_{\Pi, W} \left[ \max_{Z \in \mathcal{Z}, Y} \text{Tr}[V^T (I - \Pi_{\widehat{\mathcal{I}}}) V Z] + \text{Tr}[W(QZQ^T - Y)] \right] \quad (8)$$

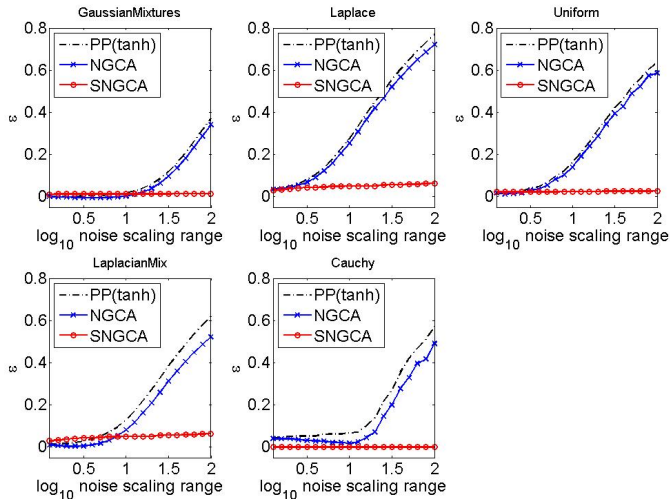
where  $Z \in \mathcal{Z}$  and  $\mathcal{Z} := \{Z \in \mathcal{S}_{L-d} \mid Z \succeq 0, \text{Tr}(Z) \leq 1\}$ .

The latter problem can be solved using a gradient-type method with complexity  $\mathcal{O}(d \log d)$  and  $\mathcal{O}(k^{-1})$  iterations (Nesterov 2007).

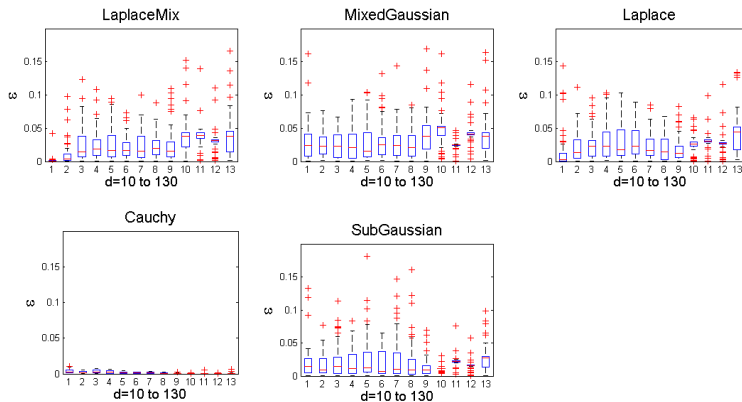
# Estimation Error for Increasing Dimensionality



# Comparison of Methods Cont'd: Noise



# Numerical Performance



# Final Slide

Thank you for your attention!

# Dual Extrapolation Algorithm:

Consider the linear matrix game

$$\min_{x \in \mathcal{A}} \max_{y \in \mathcal{B}} \langle x, Ay \rangle + \langle a, x \rangle + \langle b, y \rangle. \quad (9)$$

with  $\mathcal{A} \subset \mathbb{E}^n$  and  $\mathcal{B} \subset \mathbb{E}^m$  be closed and convex sets.

**motivation:** size  $L^2 \sim 10^6$  of the variable  $X$  rules out the possibility of using state-of-the-art interior point methods

**idea:** use a subgradient method to get low complexity

## Scheme of Dual Extrapolation Algorithm

- a) vector field of descend-ascend directions:  $F(z) = (-A^T y - a, Ax + b)$
- b)  $z = (x, y)$ ,  $z_k, z_k^+ \in \mathcal{A} \times \mathcal{B}$  and  $s_k \in E^*$  at the  $k$ -th iteration
- c) minimizer of a distance-generating function over  $\mathcal{A} \times \mathcal{B}$ :  $\bar{z}$

update:

$$\begin{aligned} z_{k+1} &= T(\bar{z}, s_k), \\ z_{k+1}^+ &= T(z_{k+1}, \lambda_k F(z_{k+1})), \\ s_{k+1} &= s_k + \lambda_k F(z_{k+1}^+) \end{aligned}$$

where  $\lambda_k > 0$  is the current (adaptively chosen) stepsize.

approximate solution:

$$\hat{z}_{k+1} = \frac{1}{k+1} \sum_{i=1}^{k+1} z_i^+.$$



# Choice of Prox-Transform

- a) **distance-generating function**  $d(z) = d_x(x) + d_y(y)$ :

$\alpha$ -strongly convex and differentiable on  $\mathcal{A} \times \mathcal{B}$

- b) **prox-function**  $V$  on  $\mathcal{A} \times \mathcal{B}$  in  $z_0 = (x_0, y_0)$ :

$$V(z_0, z) \stackrel{\text{def}}{=} d(z) - d(z_0) - \langle \nabla d(z_0), z - z_0 \rangle.$$

- c) **prox-transform**  $T(z_0, s)$  of  $s = (s_x, s_y)$

$$T(z_0, s) \stackrel{\text{def}}{=} \arg \min_{z \in \mathcal{A} \times \mathcal{B}} [\langle s, z - z_0 \rangle - V(z_0, z)].$$