# Constructing $\psi$ Satisfying (A) <sup>92</sup> $\int x\psi(x)p(x)dx = 0$ (A)

- $\begin{array}{c} \blacksquare h(\boldsymbol{x}) : \text{Smooth non-linear function from} \\ \mathbb{R}^d \text{ to } \mathbb{R} \end{array} \end{array}$
- Then following  $\psi(x)$  satisfies (A):

$$\boldsymbol{\psi}(\boldsymbol{x}) = \boldsymbol{x}^{\top} \int \boldsymbol{x} h(\boldsymbol{x}) p(\boldsymbol{x}) d\boldsymbol{x} - h(\boldsymbol{x})$$

# Constructing $\beta$ from h

Then  $\beta$  is given by

$$\boldsymbol{\beta}(h) = \int g(\boldsymbol{x}) p(\boldsymbol{x}) d\boldsymbol{x}$$
$$g(\boldsymbol{x}) = \boldsymbol{x} h(\boldsymbol{x}) - \nabla h(\boldsymbol{x})$$

93

Empirical approximation:

$$\widehat{\boldsymbol{\beta}}(h) = \frac{1}{n} \sum_{i=1}^{n} g(\boldsymbol{x}_i)$$

Note:  $\widehat{\beta}(h)$  only approximately belongs to the non-Gaussian space.

# Identifying Non-Gaussian Subspace

- Each function  $h(\mathbf{x})$  yields a vector  $\widehat{\boldsymbol{\beta}}(h)$
- Prepare a set of different non-linear functions:  $\{h_i(\boldsymbol{x})\}_{i=1}^p$
- Calculate corresponding vectors:  $\{\widehat{\beta}_i\}_{i=1}^p$

$$\widehat{\boldsymbol{\beta}}_i = \frac{1}{n} \sum_{j=1}^n g_i(\boldsymbol{x}_j)$$

 $g_i(\boldsymbol{x}) = \boldsymbol{x}h_i(\boldsymbol{x}) - \nabla h_i(\boldsymbol{x})$ 

Identifying Non-Gaussian Subspace (cont.)

- All  $\{\widehat{\beta}_i\}_{i=1}^p$  approximately belong to the non-Gaussian subspace.
- The non-Gaussian subspace may be estimated by "principal subspace" of  $\{\widehat{\beta}_i\}_{i=1}^p$ .
- We apply PCA to  $\{\widehat{\beta}_i\}_{i=1}^p$  and extract leading *m* directions.

#### 96 **Examples of Non-Linear Functions** For a random vector w. • $h(x) = \sin(\langle w, x \rangle)$ • $h(x) = \cos(\langle w, x \rangle)$ • $h(x) = \langle w, x \rangle^3 \exp(-\langle w, x \rangle^2)$ • $h(x) = \tanh(\langle w, x \rangle)$ Corresponding derivatives are • $\nabla h(\boldsymbol{x}) = \boldsymbol{w} \cos(\langle \boldsymbol{w}, \boldsymbol{x} \rangle)$ • $\nabla h(\boldsymbol{x}) = -\boldsymbol{w}\sin(\langle \boldsymbol{w}, \boldsymbol{x} \rangle)$ • $\nabla h(\boldsymbol{x}) = 3\boldsymbol{w}\langle \boldsymbol{w}, \boldsymbol{x} \rangle^2 \exp(-\langle \boldsymbol{w}, \boldsymbol{x} \rangle^2)$ $-2\boldsymbol{w}\langle \boldsymbol{w}, \boldsymbol{x} \rangle^4 \exp(-\langle \boldsymbol{w}, \boldsymbol{x} \rangle^2)$ • $\nabla h(\boldsymbol{x}) = \boldsymbol{w}(1 - \tanh(\langle \boldsymbol{w}, \boldsymbol{x} \rangle))$

# Norm of $\widehat{\boldsymbol{\beta}}$

$$\widehat{\boldsymbol{\beta}}_i = \frac{1}{n} \sum_{j=1}^n g_i(\boldsymbol{x}_j) \quad g_i(\boldsymbol{x}) = \boldsymbol{x} h_i(\boldsymbol{x}) - \nabla h_i(\boldsymbol{x})$$

Derivative is a linear operation.
Mapping h → β̂ is therefore linear.
Norm of β̂ can be arbitrary by rescaling h.
For example, h(x) and 2h(x) give the same direction but length is different.

$$\widehat{\boldsymbol{\beta}}(h) = \widehat{\boldsymbol{\beta}}(2h) - \mathcal{R}(\boldsymbol{T}^{\top})$$

## Normalization

- In PCA, long vectors are more "powerful" than short ones.
- In order to have better estimate of the non-Gaussian subspace by PCA,  $\{\hat{\beta}_i\}_{i=1}^p$ should be reasonably normalized.
- Normalization should be carried out such that accurate  $\hat{\beta}_i$  has large norm.



# Normalization (cont.)

This may be achieved by normalizing  $\{\hat{\beta}_i\}_{i=1}^p$  such that the standard deviation  $\sqrt{\varepsilon_i}$  is equivalent for all  $\{\hat{\beta}_i\}_{i=1}^p$ .

$$\varepsilon_{i} = \mathbb{E}_{\{\boldsymbol{x}_{i}\}_{i=1}^{n}} \|\widehat{\boldsymbol{\beta}}_{i} - \mathbb{E}_{\{\boldsymbol{x}_{i}\}_{i=1}^{n}}\widehat{\boldsymbol{\beta}}_{i}\|$$



However,  $\varepsilon_i$  is inaccessible.

# **Empirical Approximation**

100

Variance is expressed as

$$\varepsilon_i = \frac{1}{n} \mathbb{E}_{\boldsymbol{x}} ||g_i(\boldsymbol{x})||^2 - \frac{1}{n} ||\mathbb{E}_{\boldsymbol{x}} g_i(\boldsymbol{x})||^2$$

Empirical approximation:

$$\widehat{\varepsilon}_{i} = \frac{1}{n^{2}} \sum_{j=1}^{n} ||g_{i}(\boldsymbol{x}_{j})||^{2} - \frac{1}{n} ||\frac{1}{n} \sum_{j=1}^{n} g_{i}(\boldsymbol{x}_{j})||^{2}$$

# Algorithm of Non-Gaussian <sup>101</sup> Component Analysis

Prepare *p* different non-linear functions:  ${h_i(x)}_{i=1}^p$ Calculate  $\hat{\boldsymbol{\beta}}_i = \frac{1}{n} \sum_{j=1}^n g_i(\boldsymbol{x}_j)$  $g_i(\boldsymbol{x}) = \boldsymbol{x} h_i(\boldsymbol{x}) - \nabla h_i(\boldsymbol{x})$ **Normalize**  $\{\widehat{\boldsymbol{\beta}}_i\}_{i=1}^p$  :  $\widehat{\boldsymbol{\beta}}_i \leftarrow \widehat{\boldsymbol{\beta}}_i / \sqrt{\widehat{\varepsilon}_i}$  $\widehat{\varepsilon}_{i} = \frac{1}{n^{2}} \sum_{j=1}^{n} \|g_{i}(\boldsymbol{x}_{j})\|^{2} - \frac{1}{n} \|\frac{1}{n} \sum_{j=1}^{n} g_{i}(\boldsymbol{x}_{j})\|^{2}$ Apply PCA to  $\{\widehat{\beta}_i\}_{i=1}^p$  and extract leading m components.

### Examples

p = 100 $h_i(\boldsymbol{x}) = \sin(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle)$  $\{\widehat{\boldsymbol{eta}}_i\}_{i=1}^p$  $\{x_i\}_{i=1}^n$ 

### Examples (cont.)

$$p = 100$$
$$h_i(x) = \sin(\langle w_i, x \rangle)$$



Examples (cont.)

#### With outliers, NGCA with sin works well.



## Homework

$$\begin{array}{l} \boldsymbol{x} = \boldsymbol{s} + \boldsymbol{\epsilon} \quad \boldsymbol{s} \sim q(\boldsymbol{s}) \quad \boldsymbol{\epsilon} \sim \phi_{\boldsymbol{0},\boldsymbol{\Gamma}}(\boldsymbol{\epsilon}) \\ \boldsymbol{P} : \text{Projection onto } \boldsymbol{\Gamma}^{-1} \mathcal{S} \\ \boldsymbol{T} : \quad m \times d \text{ matrix such that } \boldsymbol{P} = \boldsymbol{T}^{\top} \boldsymbol{T} \\ \quad m = \dim(\boldsymbol{\Gamma}^{-1} \mathcal{S}) \\ \end{array} \\ \textbf{Prove } p(\boldsymbol{x}) = f(\boldsymbol{T} \boldsymbol{x}) \phi_{\boldsymbol{0},\boldsymbol{\Gamma}}(\boldsymbol{x}) \\ \quad f(\boldsymbol{z}) = g(\boldsymbol{T}^{\top} \boldsymbol{z}) \\ \qquad g(\boldsymbol{x}) = \int q(\boldsymbol{s}) e^{-\frac{1}{2} \langle \boldsymbol{\Gamma}^{-1} \boldsymbol{s}, \boldsymbol{s} \rangle} e^{\langle \boldsymbol{\Gamma}^{-1} \boldsymbol{s}, \boldsymbol{x} \rangle} d\boldsymbol{s} \end{array}$$

## Homework (cont.)

# Prove $\varepsilon_i = \mathbb{E}_{\{\boldsymbol{x}_i\}_{i=1}^n} \|\widehat{\boldsymbol{\beta}}_i - \mathbb{E}_{\{\boldsymbol{x}_i\}_{i=1}^n} \widehat{\boldsymbol{\beta}}_i\|^2$ is expressed as

$$\varepsilon_i = \frac{1}{n} \mathbb{E}_{\boldsymbol{x}} ||g_i(\boldsymbol{x})||^2 - \frac{1}{n} ||\mathbb{E}_{\boldsymbol{x}} g_i(\boldsymbol{x})||^2$$

$$\widehat{\boldsymbol{\beta}}_i = \frac{1}{n} \sum_{j=1}^n g(\boldsymbol{x}_j)$$

# Linear Dimensionality Reduction<sup>108</sup> Methods: Summary

Method	Advantages	Disadvantages
PCA	Good data description	Structure can be missed
	Analytic solution available	
	No tuning parameter	
LPP	Local structure preservation	Tuning parameters included
	Analytic solution available	
PP	Interesting structure discovery	No analytic solution
		NG measure prefixed
NGCA	Interesting structure discovery	No analytic solution

# Data with Curved Structures <sup>109</sup>



If the data cloud is bent, any linear methods fail to find the curved structure.



Limitation of linear method!

## Suggestion

- Read the following article for the next class:
- B. Schölkopf, A. Smola and K.-R. Müller: Nonlinear Component Analysis as a Kernel Eigenvalue Problem, *Neural Computation*, 10(5), 1299-1319, 1998.