Advanced Data Analysis: Projection Pursuit (2)

Masashi Sugiyama (Computer Science)

W8E-406, <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi

# Projection Pursuit 209

Find the most non-Gaussian direction.

Original formulation: maximize distance of kurtosis from 3

$$oldsymbol{\psi} = rgmax_{oldsymbol{b} \in \mathbb{R}^d} \left( rac{1}{n} \sum_{i=1}^n \langle oldsymbol{b}, \widetilde{oldsymbol{x}}_i 
angle^4 - 3 
ight)^2$$

subject to 
$$\|\boldsymbol{b}\| = 1$$

Gradient ascent algorithm

$$\boldsymbol{b} \longleftarrow \boldsymbol{b} + \varepsilon \left( \frac{1}{n} \sum_{i=1}^{n} \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^4 - 3 \right) \frac{1}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^3$$

•  $\boldsymbol{b} \longleftarrow \boldsymbol{b} / \| \boldsymbol{b} \|$ 

## Drawbacks of Gradient Method<sup>10</sup>

Choice of  $\varepsilon$  affects speed of convergence.

- If  $\varepsilon$  is small: Slow convergence
- If  $\varepsilon$  is large: Fast but less accurate
- Appropriately choosing  $\varepsilon$  is not easy in practice.
- Demonstrations:
  - demo(1): appropriate  $\varepsilon$
  - demo(2): small  $\varepsilon$
  - demo(3): large  $\varepsilon$

## Alternative Formulation <sup>211</sup>

#### Maximize or minimize kurtosis

• 
$$\psi_{max} = \operatorname*{argmax}_{\boldsymbol{b} \in \mathbb{R}^d} \left[ \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^4 \right]$$
 subject to  $\|\boldsymbol{b}\|^2 = 1$   
•  $\psi_{min} = \operatorname*{argmin}_{\boldsymbol{b} \in \mathbb{R}^d} \left[ \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^4 \right]$  subject to  $\|\boldsymbol{b}\|^2 = 1$ 

•  $\psi$  is given by  $\psi_{max}$  or  $\psi_{min}$  .

## Lagrangian

In either minimization or maximization case, Lagrangian is given by

$$L(\boldsymbol{b}, \lambda) = \frac{1}{n} \sum_{i=1}^{n} \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^4 + \lambda(\|\boldsymbol{b}\|^2 - 1)$$

Stationary (necessary) condition:

$$\frac{\partial L}{\partial \boldsymbol{b}} = \frac{4}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^3 + 2\lambda \boldsymbol{b} = \boldsymbol{0}$$

We want to find b such that

$$\frac{\partial L}{\partial \boldsymbol{b}} = \boldsymbol{0}$$



## Newton Method (Multi-Dim.) <sup>214</sup>

Problem: Find **b** such that  $f(\mathbf{b}) = \mathbf{0}$ 

$$oldsymbol{b}_{k+1} \longleftarrow oldsymbol{b}_k - \left( rac{\partial f}{\partial oldsymbol{b}} igg|_{oldsymbol{b} = oldsymbol{b}_k} 
ight)^{-1} f(oldsymbol{b}_k)$$

#### Note:

- f(b) is a d-dimensional vector.
- $\frac{\partial f}{\partial b}$  is a *d*-dimensional matrix.

## Newton-Based PP Method <sup>215</sup>

In the current setting,

$$f(\boldsymbol{b}) = rac{4}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i 
angle^3 + 2\lambda \boldsymbol{b}$$

$$\frac{\partial f}{\partial \boldsymbol{b}} = \frac{12}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_{i} \widetilde{\boldsymbol{x}}_{i}^{\top} \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_{i} \rangle^{2} + 2\lambda \boldsymbol{I}_{d}$$

Drawbacks:

- Calculating inverse  $\left(\frac{\partial f}{\partial b}\right)^{-1}$  in each step is computationally demanding.
- $\lambda$  is unknown.

## Approximation 216

$$\frac{1}{n}\sum_{i=1}^{n}\widetilde{x}_{i}\widetilde{x}_{i}^{\top}\langle \boldsymbol{b},\widetilde{x}_{i}\rangle^{2} \approx \left(\frac{1}{n}\sum_{i=1}^{n}\widetilde{x}_{i}\widetilde{x}_{i}^{\top}\right)\left(\frac{1}{n}\sum_{i=1}^{n}\langle \boldsymbol{b},\widetilde{x}_{i}\rangle^{2}\right) = \boldsymbol{I}_{d}$$
$$\frac{1}{n}\sum_{i=1}^{n}\widetilde{x}_{i}\widetilde{x}_{i}^{\top} = \boldsymbol{I}_{d} \quad \|\boldsymbol{b}\| = 1$$
$$\text{Then}$$
$$\frac{\partial f}{\partial \boldsymbol{b}} = \frac{12}{n}\sum_{i=1}^{n}\widetilde{x}_{i}\widetilde{x}_{i}^{\top}\langle \boldsymbol{b},\widetilde{x}_{i}\rangle^{2} + 2\lambda\boldsymbol{I}_{d}$$
$$\approx (12+2\lambda)\boldsymbol{I}_{d}$$

Calculating inverse is easy!

## Approximation (cont.) 217

$$f(\boldsymbol{b}) = rac{4}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^3 + 2\lambda \boldsymbol{b}$$

$$\frac{\partial f}{\partial \boldsymbol{b}} \approx \left(12 + 2\lambda\right) \boldsymbol{I}_d$$

Approximate updating rule is given by

$$oldsymbol{b} \longleftarrow rac{2}{6+\lambda} \left( 3oldsymbol{b} - rac{1}{n} \sum_{i=1}^n \widetilde{oldsymbol{x}}_i \langle oldsymbol{b}, \widetilde{oldsymbol{x}}_i 
angle^3 
ight)$$

**b** is later normalized, so the scaling factor can be dropped:  $b = 2b = 1 \sum_{n=1}^{n} \sum_{n=1}^{\infty} \frac{1}{2} \sum_{n=1}^{\infty} \frac{1}{2$ 

$$oldsymbol{b} \longleftarrow 3oldsymbol{b} - rac{1}{n}\sum_{i=1}^n \widetilde{oldsymbol{x}}_i \langle oldsymbol{b}, \widetilde{oldsymbol{x}}_i 
angle^3$$

The update rule does not depend on  $\lambda$ !

## Approximate Newton-Based <sup>218</sup> PP Method

Problem to be solved:

$$f(\boldsymbol{b}) = \boldsymbol{0}$$
 subject to  $\|\boldsymbol{b}\|^2 = 1$ 

Repeat until convergence:

• Update **b** by approximate Newton method to satisfy the stationary point condition  $\partial L/\partial b = 0$ :

$$oldsymbol{b} \longleftarrow 3oldsymbol{b} - rac{1}{n}\sum_{i=1}^n \widetilde{oldsymbol{x}}_i \langle oldsymbol{b}, \widetilde{oldsymbol{x}}_i 
angle^3$$

• Modify  $\boldsymbol{b}$  to satisfy  $\|\boldsymbol{b}\| = 1$ :

 $oldsymbol{b} \longleftarrow oldsymbol{b} / \|oldsymbol{b}\|$ 



Demonstrations:

- demo(1): Gradient ascent with appropriate  $\varepsilon$
- demo(4): Approximate Newton

#### Approximate Newton

- is much faster than gradient ascent.
- does not include any tuning parameter!

## **Outliers**



Outliers: Irregular large values
 If a Gaussian component contains outliers, its non-Gaussianity becomes very large since kurtosis contains 4th power.







A single outlier can totally corrupt the result.
 Influence of outliers needs to be reduced!

## General Non-Gaussian Measures

For some function G(s), we define a general non-Gaussian measure by

$$\frac{1}{n}\sum_{i=1}^{n}G(\langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle)$$

•  $G(s) = s^4$  corresponds to Kurtosis.

To suppress the effect of outliers, using a "gentler" function would be appropriate.

#### General Non-Gaussian Measures

Examples of smooth functions:

• 
$$G(s) = \log \cosh(s)$$
  
•  $G(s) = -\exp(-s^2/2)$ 



# Approximate Newton Procedur<sup>24</sup>

Approximate Newton procedure for centered and sphered data:

• Update **b** to satisfy the stationary-point condition:

$$\begin{split} g(s) &= G'(s) \\ b \longleftarrow \frac{1}{n} b \sum_{i=1}^{n} g'(\langle b, \widetilde{x}_i \rangle) - \frac{1}{n} \sum_{i=1}^{n} \widetilde{x}_i g(\langle b, \widetilde{x}_i \rangle) \\ & \text{(Homework)} \\ \bullet \text{ Modify } b \text{ to satisfy } \|b\| = 1 \text{ :} \end{split}$$

 $oldsymbol{b} \longleftarrow oldsymbol{b} / \|oldsymbol{b}\|$ 

#### **Derivatives**

Derivatives:

• 
$$(s^4)' = 4s^3$$
  
 $(4s^3)' = 12s^2$ 

- $(\log \cosh(s))' = \tanh(s)$  $(\tanh(s))' = 1 - \tanh^2(s)$
- $(-\exp(-s^2/2))' = s \exp(-s^2/2)$  $(s \exp(-s^2/2))' = (1 - s^2) \exp(-s^2/2)$



Approximate Newton with Kurtosis:

$$g(s) = 4s^3$$



226

×

Approximate Newton with log(cosh):

$$g(s) = \tanh(s)$$

Approximate Newton with log(cosh) is robust against outliers!

# Extracting Several Non-Gaussian Directions

227

Running the algorithm many times from different initial points may give different non-Gaussian directions.

However, this is not computationally efficient.

Another idea: Find orthogonal directions

This is achieved by modifying the direction as

$$oldsymbol{b} \longleftarrow oldsymbol{b} - \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{\psi}_i \ \mathbf{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{b}, oldsymbol{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{\psi}_i 
angle oldsymbol{b} oldsymbol{b} oldsymbol{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{b}, oldsymbol{b} = \sum_{i=1}^{k-1} \langle oldsymbol{b}, oldsymbol{b} = \sum_$$

# **Full Algorithm**

Center and sphere samples:  $\widetilde{X} = (XH^2X)^{-\frac{1}{2}}XH$ 

For 
$$k = 1, 2, ..., m$$

• Repeat until convergence:

$$\mathbf{b} \longleftarrow \frac{1}{n} \mathbf{b} \sum_{i=1}^{n} g'(\langle \mathbf{b}, \widetilde{\mathbf{x}}_i \rangle) - \frac{1}{n} \sum_{i=1}^{n} \widetilde{\mathbf{x}}_i g(\langle \mathbf{b}, \widetilde{\mathbf{x}}_i \rangle)$$

$$\mathbf{b} \longleftrightarrow \mathbf{b} - \sum_{i=1}^{k-1} \langle \mathbf{b}, \psi_i \rangle \psi_i$$

$$\mathbf{b} \longleftrightarrow \mathbf{b} / \| \mathbf{b} \|$$

$$\mathbf{X} = (\mathbf{x}_1 | \mathbf{x}_2 | \cdots | \mathbf{x}_n)$$

$$\mathbf{X} = (\mathbf{x}_1 | \mathbf{x}_2 | \cdots | \mathbf{x}_n)$$

$$\mathbf{H} - \mathbf{I} - \frac{1}{1}$$

 $\mathcal{N}$ 

• 
$$\psi_k = b$$
  
Embed the data  $x$  by  
 $\overline{z} = B_{PP}(x - \frac{1}{2}X\mathbf{1}_n)$ 

$$oldsymbol{H} = oldsymbol{I}_n - rac{1}{n} oldsymbol{1}_{n imes n}$$
  
 $oldsymbol{I}_n$ : *n*-dimensional identity matrix  
 $oldsymbol{1}_{n imes n}$ : *n imes n* matrix with all ones  
 $oldsymbol{1}_n$ : *n*-dimensional vector with all ones  
 $oldsymbol{B}_{PP} = (oldsymbol{\psi}_1 | oldsymbol{\psi}_2 | \cdots | oldsymbol{\psi}_m)^{ op}$ 

 $\langle \Psi_1 | \Psi_2 |$ 

 $\mathbf{D}PP$ 

228

## Homework

229

 Implement approximate Newton-based PP method with general non-Gaussianity measure and reproduce the 2-dimensional examples with an outlier shown in the class. You may create similar (or more interesting) data sets by yourself.

http://sugiyama-www.cs.titech.ac.jp/~sugi/data/DataAnalysis



## Homework (cont.)

230

2. Prove that the approximate Newton updating rule is given by

$$oldsymbol{b} \longleftarrow rac{1}{n}oldsymbol{b}\sum_{i=1}^n g'(\langle oldsymbol{b}, \widetilde{oldsymbol{x}}_i 
angle) - rac{1}{n}\sum_{i=1}^n \widetilde{oldsymbol{x}}_i g(\langle oldsymbol{b}, \widetilde{oldsymbol{x}}_i 
angle)$$

under the following approximation:

$$\frac{1}{n}\sum_{i=1}^{n}\widetilde{\boldsymbol{x}}_{i}\widetilde{\boldsymbol{x}}_{i}^{\top}g'(\langle \boldsymbol{b},\widetilde{\boldsymbol{x}}_{i}\rangle) \approx \frac{1}{n}\sum_{i=1}^{n}g'(\langle \boldsymbol{b},\widetilde{\boldsymbol{x}}_{i}\rangle)\boldsymbol{I}_{d}$$

#### Schedule

June 25<sup>th</sup>: Projection Pursuit (2)

- Application Deadline to Mini-Conference
- July 2<sup>nd</sup>: Independent Component Analysis
- July 9<sup>th</sup>: Preparation for Mini-Conference
- July 16<sup>th</sup>: Mini-Conference Day 1
- July 23<sup>rd</sup>: Mini-Conference Day 2