Advanced Data Analysis: Projection Pursuit (2)

Masashi Sugiyama (Computer Science)

W8E-505, <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi

Drawbacks of Gradient Method¹³

Choice of ε affects speed of convergence.

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

Alternative Formulation ²¹⁴

Original formulation: maximize distance from 3

$$\boldsymbol{\psi} = \operatorname*{argmax}_{\boldsymbol{b} \in \mathbb{R}^d} \left(\frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^4 - 3 \right)^2$$

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

Alternative formulation: maximize or minimize kurtosis

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

ert ψ is given by ψ_{max} or ψ_{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 point (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.) ²¹⁷

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

$$\boldsymbol{b}_{k+1} \leftarrow \boldsymbol{b}_k - \left(\frac{\partial f}{\partial \boldsymbol{b}} \bigg|_{\boldsymbol{b} = \boldsymbol{b}_k} \right)^{-1} f(\boldsymbol{b}_k)$$

Note:

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

Newton-Based PP Method ²¹⁸

In the current setting,

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

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

Drawbacks:

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

Approximation 219

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

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

Calculating inverse is easy!

Approximation (cont.) 220

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

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

Approximate updating rule is given by

$$\boldsymbol{b} \longleftarrow rac{1}{12+2\lambda} \left(12\boldsymbol{b} - rac{4}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^3 \right)$$

b is later normalized, thus the scaling factor can be dropped: $b \leftarrow 3b = \frac{1}{2} \sum_{n=1}^{n} \widetilde{x} \cdot b = \widetilde{x} \cdot \sqrt{3}$

$$\boldsymbol{b} \longleftarrow 3\boldsymbol{b} - \frac{1}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^3$$

The update rule does not depend on λ !

Approximate Newton-Based ²²¹ PP Method

Problem to be solved:

$$f(b) = 0$$
 subject to $||b||^2 = 1$

Repeat until convergence:

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

$$\boldsymbol{b} \longleftarrow 3\boldsymbol{b} - \frac{1}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \langle \boldsymbol{b}, \widetilde{\boldsymbol{x}}_i \rangle^3$$

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

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



Demonstrations:

- demo(1): Gradient ascent with appropriate ε
- 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 deemphasized!

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 Procedur²⁷

Approximate Newton procedure for centered and sphered data:

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

$$g(s) = G'(s)$$

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

$$\bullet Modify b \text{ to satisfy } ||b|| = 1:$$

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

Derivatives

228

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

 $g(s) = 4s^3$



229

×

Approximate Newton with log(cosh)

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

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

Extracting Several Non-Gaussian Directions

- Running the algorithm many times from different initial points may give different non-Gaussian directions.
- However, this might not be efficient.
- Another idea: Find orthogonal directions
- This is achieved modifying the direction by

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



230

Full Algorithm

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

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

• Repeat until convergence:

$$oldsymbol{\psi}_k = oldsymbol{b}$$

Embed the data x by

$$\overline{\boldsymbol{z}} = \boldsymbol{B}_{PP}(\boldsymbol{x} - \frac{1}{n}\boldsymbol{X}\boldsymbol{1}_n)$$

$$egin{aligned} \widetilde{oldsymbol{X}} &= (\widetilde{oldsymbol{x}}_1 | \widetilde{oldsymbol{x}}_2 | \cdots | \widetilde{oldsymbol{x}}_n) \ oldsymbol{X} &= (oldsymbol{x}_1 | oldsymbol{x}_2 | \cdots | oldsymbol{x}_n) \ oldsymbol{H} &= oldsymbol{I}_n - rac{1}{n} oldsymbol{1}_{n imes n} \end{aligned}$$

231

 \boldsymbol{I}_n : *n*-dimensional identity matrix $\boldsymbol{1}_{n \times n}$: $n \times n$ matrix with all ones $\boldsymbol{1}_n$: *n*-dimensional vector with all ones $\boldsymbol{B}_{PP} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \cdots | \boldsymbol{\psi}_m)^{\top}$

232

Mini-Conference on Data Analysis: Program

July 14th

- Tatsuya Shigemura
- Youhei Namiki
- Rosset Matthieu
- Supaporn Spanurattana
- Liang Zheng
- Nakamura Satoru
- Jacob Montiel
- Sangeeta Biswas
- Satoshi Gomori
- Chativit Prayoonsri

July 21st

- Jaak Simm
- Swit Phuvipadawat
- Takeshi Motoda
- Wisnu Ananta Kusuma
- Nakamasa Inoue
- Shintaro Matsui
- Felipe Gomez
- Anders Lind
- Cetinkaya Ahmet

Talk: 7 minutes, Q&A: 2 minutes

If You Are ...

233

eager to do homework, try the following two problems.

Homework

234

 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.)

235

2. Prove that approximate Newton updating rule is given by

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

under the following approximation:

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