Drawbacks of Gradient Method⁵¹

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

$$\boldsymbol{\psi} = \operatorname*{argmax}_{\boldsymbol{b} \perp \{\boldsymbol{\psi}_i\}_{i=1}^{k-1}} \left(\frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^4 - 3 \right)^2 \text{ subject to } \|\boldsymbol{b}\| = 1$$

 ψ is given by ψ_{max} or ψ_{min} :

$$\boldsymbol{\psi}_{max} = \underset{\boldsymbol{b} \perp \{\boldsymbol{\psi}_i\}_{i=1}^{k-1}}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^{n} \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^4 \quad \text{subject to } \|\boldsymbol{b}\|^2 = 1$$

$$\boldsymbol{\psi}_{min} = \operatorname*{argmin}_{\boldsymbol{b} \perp \{\boldsymbol{\psi}_i\}_{i=1}^{k-1}} \frac{1}{n} \sum_{i=1}^{n} \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^4 \text{ subject to } \|\boldsymbol{b}\|^2 = 1$$

Lagrangian

- For the moment, ignore $b \perp \{\psi_i\}_{i=1}^{k-1}$.
- In either minimization or maximization case, Lagrangian is given by

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

Stationary point (necessary condition):

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

We want to find b which satisfies

$$f(\boldsymbol{b}) = \frac{\partial L}{\partial \boldsymbol{b}} = \mathbf{0}$$



Newton Method (Multi-Dim.)

Find **b** s.t.
$$f(\mathbf{b}) = \mathbf{0}$$

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

Note:

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

Newton Approach

f(b) = 0 subject to $||b||^2 = 1$ and $b \perp \{\psi_i\}_{i=1}^{k-1}$

Repeat the following until convergence:

• Update **b** by Newton method to satisfy the stationary point condition $\frac{\partial L}{\partial \mathbf{b}} = \mathbf{0}$:

$$\boldsymbol{b} \longleftarrow \boldsymbol{b} - \left(\frac{\partial f}{\partial \boldsymbol{b}}\right)^{-1} f(\boldsymbol{b})$$

• Modify \boldsymbol{b} to satisfy $\boldsymbol{b} \perp \{\boldsymbol{\psi}_i\}_{i=1}^{k-1}$: $\boldsymbol{b} \longleftarrow \boldsymbol{b} - \sum_{i=1}^{k-1} \langle \boldsymbol{b}, \boldsymbol{\psi}_i \rangle \boldsymbol{\psi}_i$

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

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

Newton Approach (cont.)

$$f(\boldsymbol{b}) = \frac{4}{n} \sum_{i=1}^{n} \boldsymbol{x}_i \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^3 + 2\lambda \boldsymbol{b}$$
$$\frac{\partial f}{\partial \boldsymbol{b}} = \frac{12}{n} \sum_{i=1}^{n} \boldsymbol{x}_i \boldsymbol{x}_i^{\top} \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^2 + 2\lambda \boldsymbol{I}_d$$

However,

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



59

$$f(\mathbf{b}) = \frac{4}{n} \sum_{i=1}^{n} \mathbf{x}_i \langle \mathbf{b}, \mathbf{x}_i \rangle^3 + 2\lambda \mathbf{b} \qquad \frac{\partial f}{\partial \mathbf{b}} \approx (12 + 2\lambda) \mathbf{I}_d$$

Approximate updating rule is

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

b is normalized so we can ignore constant:

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

 λ can be removed!



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.





Demonstrations:

- demo(4): Approximate Newton (without outlier)
- demo(5): Approximate Newton (with single outlier)

A single outlier can totally corrupt the result.
 Influence of outliers should be deemphasized!

General Non-Gaussian Measures

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

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

 $G(s) = s^4$ 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 Procedure⁶⁶

- Approximate Newton procedure for centered and sphered data:
 - Update *b* to satisfy the stationary-point condition:

$$\begin{aligned} \mathbf{b} &\leftarrow \frac{1}{n} \mathbf{b} \sum_{i=1}^{n} g'(\langle \mathbf{b}, \mathbf{x}_i \rangle) - \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i g(\langle \mathbf{b}, \mathbf{x}_i \rangle) \\ \bullet &\text{ Modify } \mathbf{b} \text{ to satisfy } \mathbf{b} \bot \{ \mathbf{\psi}_i \}_{i=1}^{k-1} : \qquad g(s) = G'(s) \\ &\mathbf{b} \leftarrow \mathbf{b} - \sum_{i=1}^{k-1} \langle \mathbf{b}, \mathbf{\psi}_i \rangle \mathbf{\psi}_i \end{aligned}$$

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

 $egin{array}{c} egin{array}{c} egin{array}$

Derivatives

Derivatives:

•
$$(s^4)' = 4s^3$$

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

•
$$(\log \cosh(s))' = \tanh(s)$$

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

Examples

Demonstrations:

- demo(5): Approximate Newton with Kurtosis ${\it g}(s)=4s^3$
- demo(6): Approximate Newton with log(cosh) $\label{eq:g} \begin{array}{l} g(s) = \tanh(s) \end{array}$

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

Homework

Prove that approximate Newton updating rule is given by

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

under the following approximation:

$$\frac{1}{n}\sum_{i=1}^{n} \boldsymbol{x}_{i}\boldsymbol{x}_{i}^{\mathsf{T}}g'(\langle \boldsymbol{b},\boldsymbol{x}_{i}\rangle) \approx \frac{1}{n}\sum_{i=1}^{n}g'(\langle \boldsymbol{b},\boldsymbol{x}_{i}\rangle)\boldsymbol{I}_{d}$$

Homework (cont.)

Implement projection pursuit algorithms.

Apply the algorithms to your data set and extract two non-Gaussian directions in your data.

Project your data onto the twodimensional subspace and find something interesting (data mining!).