Effect of Data Rescaling

Rescaling of the data affects the results of LPP and PCA, because Euclidean distance depends on the scale of data.

$$egin{aligned} m{B}_{PCA} &= rgmin_{m{B}\in\mathbb{R}^{m imes d}} \left[\sum_{i=1}^n \|m{B}^ op m{B} m{x}_i - m{x}_i\|^2
ight] \ & ext{ subject to } m{B}m{B}^ op = m{I}_m \ & ext{B}_{LPP} &= rgmin_{m{B}\in\mathbb{R}^{m imes d}} \left[\sum_{i,j=1}^n \|m{B} m{x}_i - m{B} m{x}_j\|^2 m{W}_{i,j}
ight], \ & ext{ subject to } m{B} m{X} m{D} m{X}^ op m{B}^ op = m{I}_m \end{aligned}$$

Invariance under Data Rescaling³⁶

- To be invariant under data rescaling, embedding criterion should not depend on the scale of the data.
- Idea: Use data distributions, rather than distances.
- Suppose data samples are i.i.d. random variables.

$$oldsymbol{x}_i \stackrel{i.i.d.}{\sim} p(oldsymbol{x})$$

Gaussian Distribution

Gaussian distribution: Probability density function is given by

$$\phi_{\boldsymbol{\theta},\boldsymbol{\Gamma}}(\boldsymbol{x}) = \frac{1}{(2\pi)^{\frac{d}{2}} |\boldsymbol{\Gamma}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\theta})^{\top}\boldsymbol{\Gamma}^{-1}(\boldsymbol{x}-\boldsymbol{\theta})\right)$$

θ, Γ :Mean, covariance
 $\mathbb{E}[x] = θ$ $\mathbb{E}[(x - θ)(x - θ)^{\top}] = Γ$



Interesting Directions for Data Visualization

What distribution is interesting to visualize?

If data follows the Gaussian distribution, samples are spherically distributed.

Visualizing spherically distributed samples is not so interesting.
 What about "non-Gaussian" data?



Non-Gaussian data look interesting:
 We want to project the data so that it has non-Gaussian distributions

Uniform (sharp edge)



Projection Pursuit

Idea: Iteratively find non-Gaussian directions in the data

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

Find the most non-Gaussian direction in data:

$$\psi_k = \underset{\boldsymbol{b} \perp \{\psi_i\}_{i=1}^{k-1}}{\operatorname{argmax}} J_{PP}(\boldsymbol{b})$$

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

PP embedding of a sample x':

$$oldsymbol{z}' = oldsymbol{B}_{PP}oldsymbol{x}' \qquad oldsymbol{B}_{PP} = (oldsymbol{\psi}_1|oldsymbol{\psi}_2|\cdots|oldsymbol{\psi}_m)^{ op}$$

Kurtosis

PP needs a non-Gaussianity measure. Kurtosis: $\mathbb{R}[(e - \mathbb{R}[e])^4]$

$$\beta_4 = \frac{\mathbb{E}[(s - \mathbb{E}[s])^4]}{(\mathbb{E}[(s - \mathbb{E}[s])^2])^2} \quad (>0)$$
$$s = \langle \boldsymbol{b}, \boldsymbol{x} \rangle$$

If tail of distribution is



Kurtosis (cont.)

- $\beta_4 = 3$: Gaussian distribution
- $\beta_4 < 3$: Sub-Gaussian distribution
- $\beta_4 > 3$: Super-Gaussian distribution



43 **Kurtosis-Based Non-Gaussianity Measure** Non-Gaussianity is strong if $(\beta_4 - 3)^2$ is large. In practice, we use empirical approximation: $J_{PP}(\mathbf{b}) = \left(\frac{\frac{1}{n}\sum_{i=1}^{n}[(s_{i}-\overline{s})^{4}]}{(\frac{1}{n}\sum_{i=1}^{n}[(s_{i}-\overline{s})^{2}])^{2}} - 3\right)^{2} \qquad s_{i} = \langle \mathbf{b}, \mathbf{x}_{i} \rangle$ $\overline{s} = \frac{1}{n}\sum_{i=1}^{n}s_{i}$

There is no known method for analytically solving the optimization problem.

$$\boldsymbol{\psi} = \operatorname*{argmax}_{\boldsymbol{b} \perp \{\boldsymbol{\psi}_i\}_{i=1}^{k-1}} J_{PP}(\boldsymbol{b}) \text{ subject to } \|\boldsymbol{b}\| = 1$$

We resort to numerical methods.

Gradient Ascent Approach

(0)

Repeat the following until convergence:

• Update \boldsymbol{b} to increase J_{PP} :

$$\boldsymbol{b} \longleftarrow \boldsymbol{b} + \varepsilon \frac{\partial J_{PP}}{\partial \boldsymbol{b}} \quad (\varepsilon >$$

• Modify $m{b}$ to satisfy $m{b}ot\{\psi_i\}_{i=1}^{k-1}$:

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



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



Data Centering and Sphering ⁴⁵

Centering:

$$oldsymbol{x}_i \longleftarrow oldsymbol{x}_i - \overline{oldsymbol{x}}_i$$

$$\overline{\boldsymbol{x}} = \frac{1}{n} \sum_{j=1}^{n} \boldsymbol{x}_j$$

Sphering (or pre-whitening):

$$x_i \longleftarrow C^{-rac{1}{2}} x_i$$



Covariance matrix for sphered data:

$$C = I_d$$



Gradient for Sphered Data

For centered and sphered data, gradient is given by

$$\frac{\partial J_{PP}}{\partial \boldsymbol{b}} = 2\left(\frac{1}{n}\sum_{i=1}^{n} \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^4 - 3\right) \left(\frac{4}{n}\sum_{i=1}^{n} \boldsymbol{x}_i \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^3\right)$$

Updating rule is

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

• $\boldsymbol{b} \leftarrow \boldsymbol{b} - \sum_{i=1}^{k-1} \langle \boldsymbol{b}, \boldsymbol{\psi}_i \rangle \boldsymbol{\psi}_i$

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

Examples

$$d = 2, m = 1, n = 1000$$







Examples (cont.)

$$d = 2, m = 1, n = 1000$$





Homework

- Prove the followings for centered and sphered data:
 - Covariance matrix is given by

$$\frac{1}{n}\sum_{j=1}^{n}(\boldsymbol{x}_{j}-\overline{\boldsymbol{x}})(\boldsymbol{x}_{j}-\overline{\boldsymbol{x}})^{\top}=\boldsymbol{I}_{d}$$

• J_{PP} under $\|\boldsymbol{b}\| = 1$ is given by

$$J_{PP}(\boldsymbol{b}) = \left(\frac{1}{n}\sum_{i=1}^{n} \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^4 - 3\right)^2$$

• Gradient $\partial J_{PP}/\partial b$ is given by

$$\frac{\partial J_{PP}}{\partial \boldsymbol{b}} = 2\left(\frac{1}{n}\sum_{i=1}^{n} \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^4 - 3\right) \left(\frac{4}{n}\sum_{i=1}^{n} \boldsymbol{x}_i \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle^3\right)$$