# Advanced Data Analysis: Projection Pursuit (1)

Masashi Sugiyama (Computer Science)

W8E-505,  sugi@cs.titech.ac.jp

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

# I.i.d. Samples

- Independent and identically distributed (i.i.d.) samples

$$\boldsymbol{x}_i \overset{i.i.d.}{\sim} P(\boldsymbol{x})$$

  - Independent: joint probability is a product of each probability

  $$P(\boldsymbol{x}_i, \boldsymbol{x}_j) = P(\boldsymbol{x}_i)P(\boldsymbol{x}_j)$$

  - Identically distributed: each variable follow the identical distribution:

  $$\boldsymbol{x}_i \sim P(\boldsymbol{x})$$

# Gaussian Distribution

■ **Gaussian distribution**: Probability density function is given by

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

■ $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ :Mean, covariance

$$\mathbb{E}[\boldsymbol{x}] = \boldsymbol{\mu}$$

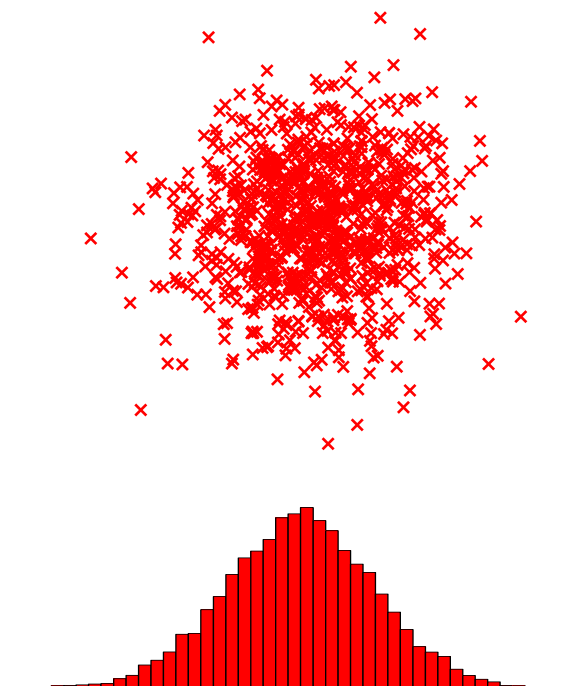$$\mathbb{E}[(\boldsymbol{x}-\boldsymbol{\mu})(\boldsymbol{x}-\boldsymbol{\mu})^{\top}] = \boldsymbol{\Sigma}$$

■ When one-dimensional,

$$\phi_{\mu,\sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

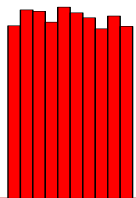# Interesting Directions for Data Visualization

- Which distribution is interesting to visualize?

- If data follows the Gaussian distribution, samples are <span style="color:red">spherically</span> distributed.

- Visualizing spherically distributed samples is not so interesting.

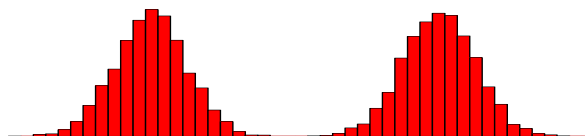- What about "<span style="color:red">non-Gaussian</span>" data?

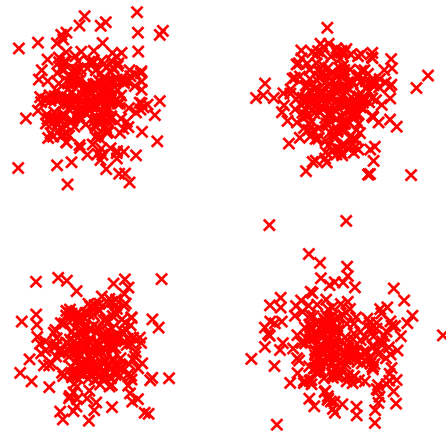# Non-Gaussian Distributed Data

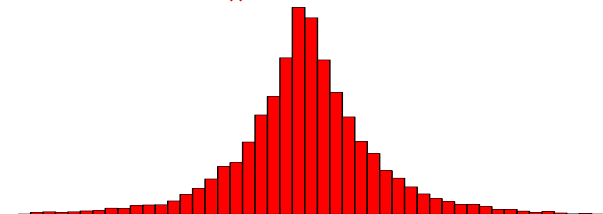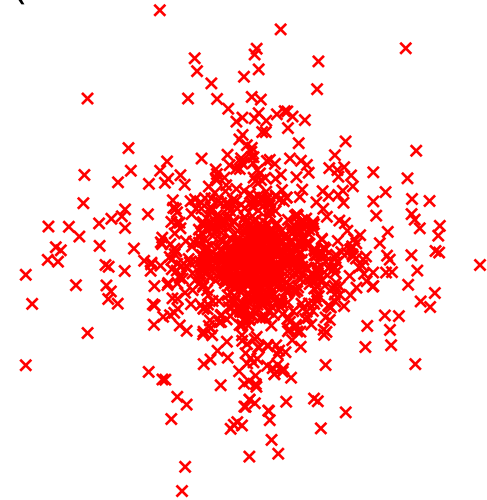- Non-Gaussian data look more interesting than Gaussian:

Uniform
(sharp edge)

Gaussian mixture
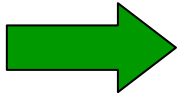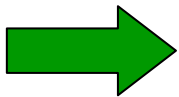(cluster structure)

Laplacian
(existence of outliers)

# Projection Pursuit

- Idea: Find the most non-Gaussian direction in the data

- For this purpose, we need a criterion to measure non-Gaussianity of data as a function of the direction.

# Kurtosis

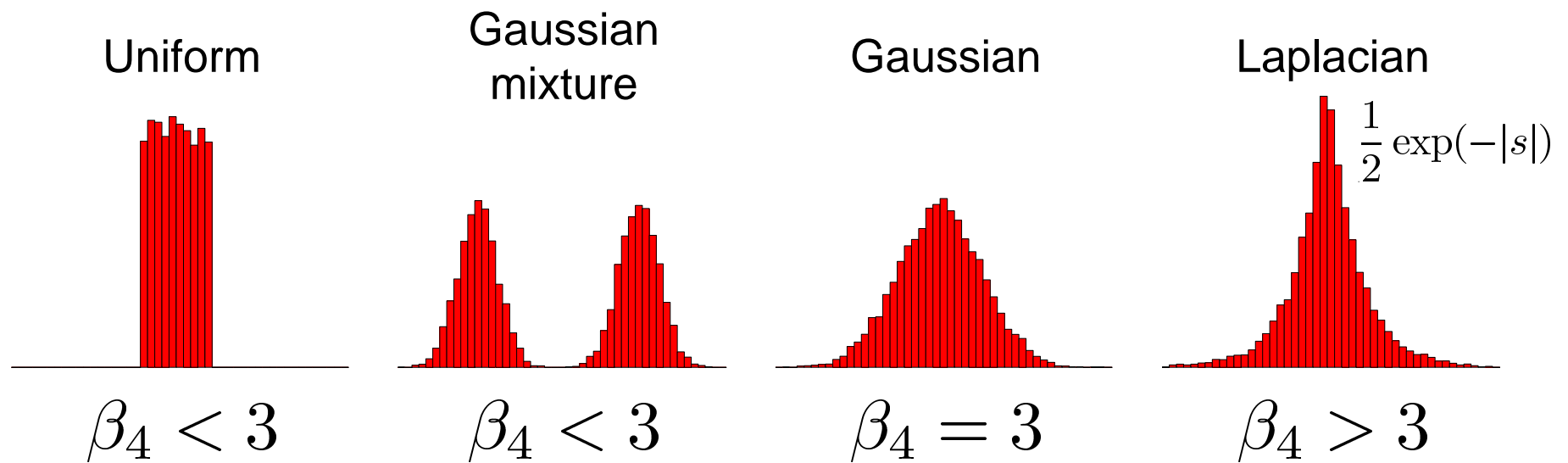■ **Kurtosis** for a one-dimensional random variable $s$:

$$\beta_4 = \frac{\mathbb{E}[(s - \mathbb{E}[s])^4]}{(\mathbb{E}[(s - \mathbb{E}[s])^2])^2} \quad (> 0)$$

■ Kurtosis measures the "sharpness" of the distributions.

■ If tail of distribution is

- Heavy ⟹ $\beta_4$ is large
- Light ⟹ $\beta_4$ is small

# Kurtosis (cont.)

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



Uniform

Gaussian mixture

Gaussian

Laplacian

$\frac{1}{2}\exp(-|s|)$

$\beta_4 < 3$        $\beta_4 < 3$        $\beta_4 = 3$        $\beta_4 > 3$

# Kurtosis-Based Non-Gaussianity Measure

$$\beta_4 = \frac{\mathbb{E}[(s - \mathbb{E}[s])^4]}{(\mathbb{E}[(s - \mathbb{E}[s])^2])^2}$$

- Non-Gaussianity is strong if $(\beta_4 - 3)^2$ is large.

- Non-Gaussianity of the data for a direction $b$ can be measured by letting $s = \langle b, x \rangle$ and $\|b\| = 1$ .

# PP Criterion

- In practice, we use empirical approximation:

$$J_{PP}(\boldsymbol{b}) = \left( \frac{\frac{1}{n}\sum_{i=1}^{n}(s_i - \overline{s})^4}{(\frac{1}{n}\sum_{i=1}^{n}(s_i - \overline{s})^2)^2} - 3 \right)^2$$

$$s_i = \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle$$

$$\overline{s} = \frac{1}{n}\sum_{i=1}^{n} s_i$$

- PP criterion:

$$\boldsymbol{\psi} = \underset{\boldsymbol{b} \in \mathbb{R}^d}{\arg\max}\, J_{PP}(\boldsymbol{b})$$

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

- There is no known method for analytically solving this optimization problem.

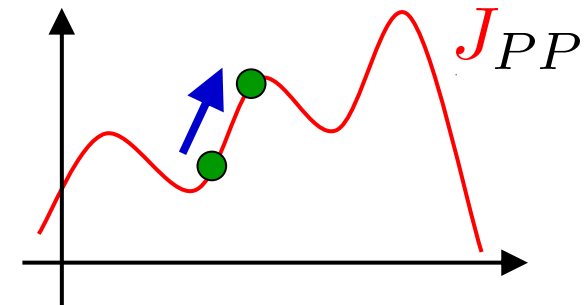- We resort to numerical methods.

# Gradient Ascent Approach

■ Repeat until convergence:
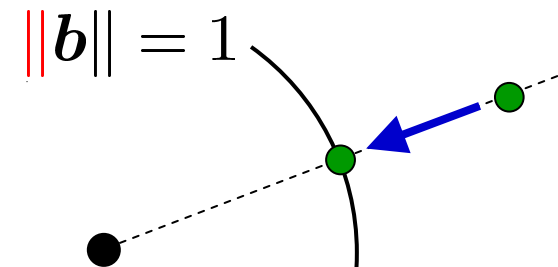
● Update $\boldsymbol{b}$ to increase $J_{PP}$ :

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

$$(\varepsilon > 0)$$

$J_{PP}$

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

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

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

- Centering:

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

- Sphering (or pre-whitening):

$$\widetilde{\boldsymbol{x}}_i = \left(\frac{1}{n}\sum_{i=1}^{n}\overline{\boldsymbol{x}}_i\overline{\boldsymbol{x}}_i^{\top}\right)^{-\frac{1}{2}}\overline{\boldsymbol{x}}_i$$

- In matrix,

$$\widetilde{\boldsymbol{X}} = (\frac{1}{n}\boldsymbol{X}\boldsymbol{H}^2\boldsymbol{X}^{\top})^{-\frac{1}{2}}\boldsymbol{X}\boldsymbol{H}$$

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

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

$$\boldsymbol{H} = \boldsymbol{I}_n - \frac{1}{n}\boldsymbol{1}_{n\times n}$$

$\boldsymbol{I}_n$: $n$-dimensional identity matrix

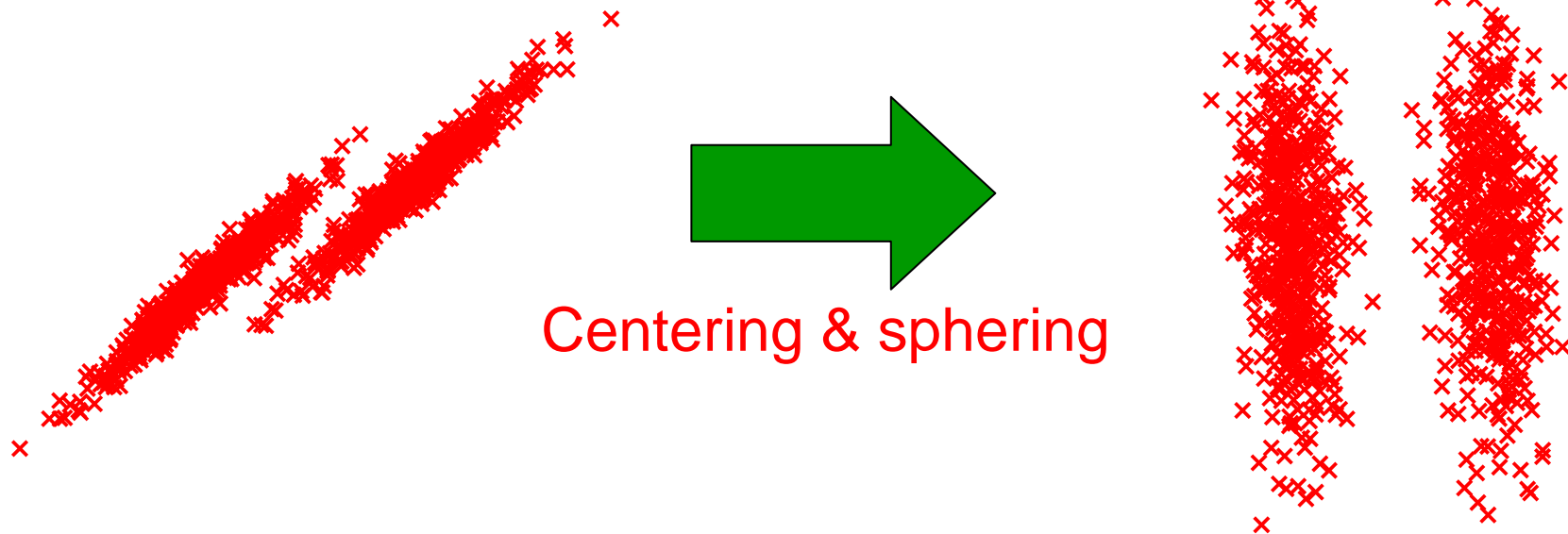$\boldsymbol{1}_{n\times n}$: $n\times n$ matrix with all ones

# Data Centering and Sphering

■ By centering and sphering, covariance matrix becomes identity:

$$\frac{1}{n} \sum_{i=1}^{n} \widetilde{\boldsymbol{x}}_i \widetilde{\boldsymbol{x}}_i^{\top} = \boldsymbol{I}_d$$

Homework: Prove it!



Centering & sphering

# Simplification for Sphered Data

■ For centered and sphered samples $\{\widetilde{\boldsymbol{x}}_i\}_{i=1}^n$ ,

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

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

■ Gradient update rule is

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

■ Don't forget normalization: $\boldsymbol{b} \longleftarrow \boldsymbol{b}/\|\boldsymbol{b}\|$

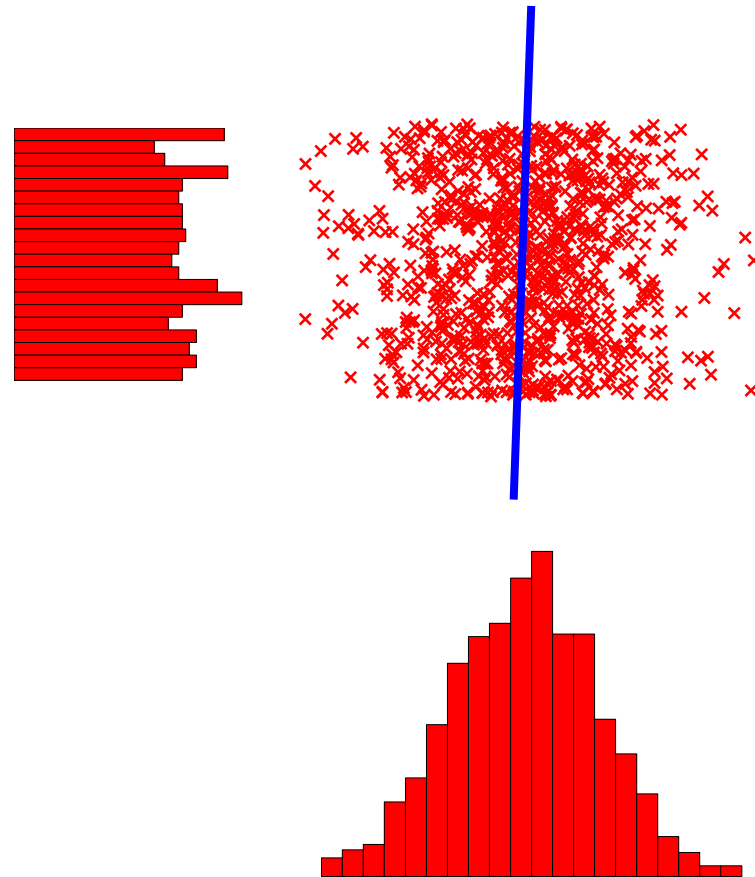■ Homework: Prove them!

# Examples

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

- $\boldsymbol{x} = \begin{pmatrix} s \\ t \end{pmatrix}$

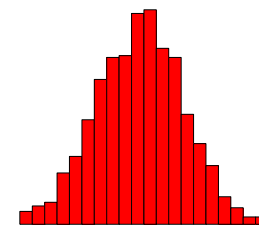- $s \sim N(0, 1)$
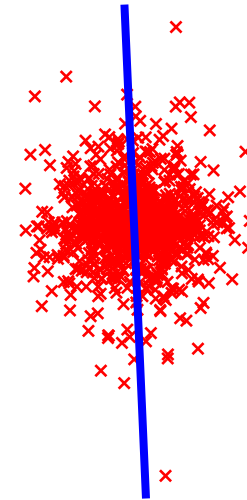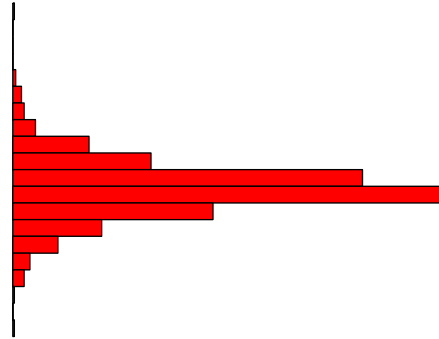- $t \sim U(-\sqrt{3}, \sqrt{3})$

# Examples (cont.)

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

- $\boldsymbol{x} = \begin{pmatrix} s \\ t \end{pmatrix}$

- $s \sim N(0, 1)$
- $t \sim Lap(0, 1)$

# Notification of Final Assignment

- **Data Analysis**: Apply dimensionality reduction or clustering techniques to your own data set and "mine" something interesting!

# Mini-Conference on Data Analysis

- At the end of the semester, we have a mini-conference on data analysis.

- Some of the students may present their data analysis results.

- Those who give a talk at the conference will have very good grades!

# Schedule

- June 23rd: Regular lecture (projection pursuit 1)
- June 30th: Regular lecture (projection pursuit 2)
- July 7th: Preparation for the mini-conference (no lecture)
- July 14th: Mini-conference (day 1)
- July 21st: Mini-conference (day 2 if necessary)

# Mini-Conference on Data Analysis

- Application procedure: On June 23rd, just say to me "I want to give a talk!".

- Presentation: approx. 10 min (?)
  - Description of your data
  - Methods to be used
  - Outcome

- Slides should be in English.

- Better to speak in English, but Japanese may also be allowed (perhaps your friends will provide simultaneous translation!).
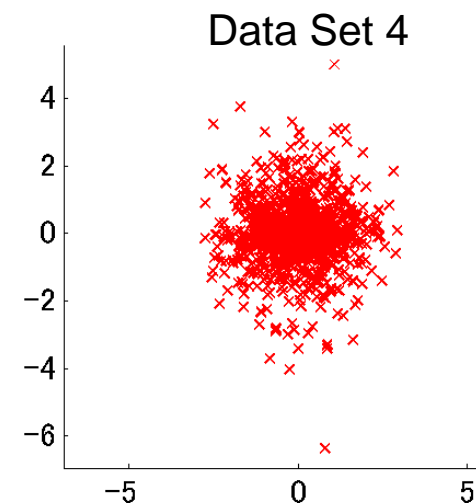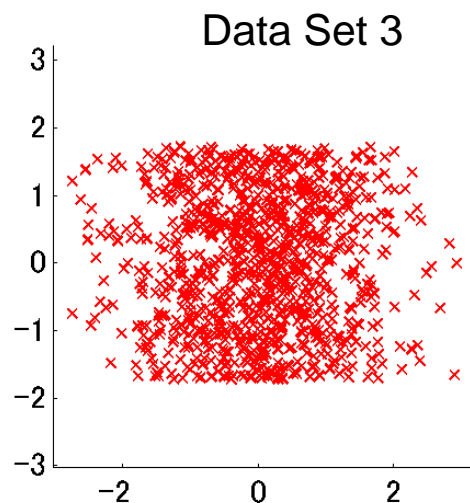
# If You Are ...

- eager to do homework, try the following two problems.

# Homework

1. Implement PP and reproduce the 2-dimensional examples shown in the class.

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



Data Set 3

Data Set 4

   You may create similar (and more interesting) data sets by yourself.

# Homework (cont.)

2.  Prove the following for centered and sphered samples $\{\widetilde{\boldsymbol{x}}_i\}_{i=1}^n$ :

   A) Covariance matrix is given by

   $$\frac{1}{n}\sum_{i=1}^{n}\widetilde{\boldsymbol{x}}_i\widetilde{\boldsymbol{x}}_i^\top = \boldsymbol{I}_d$$

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

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

   C) Gradient $\partial J_{PP}/\partial\boldsymbol{b}$ is given by

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