# Dimensionality Reduction

Makoto Yamada
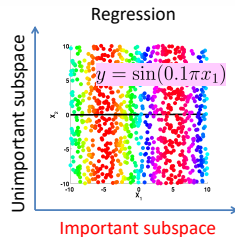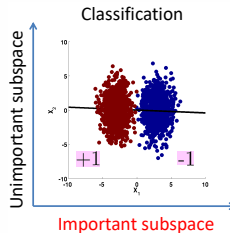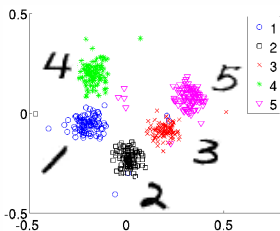
`myamada@i.kyoto-u.ac.jp`

Kyoto University

# Dimensionality Reduction

Dimensionality reduction is a method to reduce the dimensionality of data.

- Feature selection is a dimensionality reduction method. Select a set of $m$ features among $d$ features ($m < d$).
- We use feature selection for interpretation.
- We use dimensionality reduction to compress data, to visualize data, etc.

# Problem Formulation

Dimension reduction (DR) is to find a low-dimensional mapping $f : \mathbb{R}^d \to \mathbb{R}^m$ $(d > m)$ $(\boldsymbol{x} \in \mathbb{R}^d)$

- It is useful for data visualization.
- Keep the original information as much as possible
- The DR outputs the combination of features.
- Linear dimension reduction $\boldsymbol{z} = \boldsymbol{U}^\top \boldsymbol{x}$ $(\boldsymbol{U} \in \mathbb{R}^{d \times m})$.

$$m \left[ \boldsymbol{z} \right] = m \left[ \boldsymbol{U}^\top \right] \underbrace{}_{d} \left[ \boldsymbol{x} \right] d$$

- Nonlinear dimension reduction $\boldsymbol{z} = \boldsymbol{g}(\boldsymbol{x})$. For example, deep learning model: $\boldsymbol{g}(\boldsymbol{x}) = \boldsymbol{\sigma}(\boldsymbol{W}_1(\boldsymbol{\sigma}(\boldsymbol{W}_2)))$
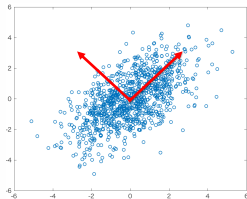
# Principal Component Analysis

The principal component analysis (PCA) is given as:

$$\widehat{U} = \underset{U^\top U = I}{\operatorname{argmax}} \ \operatorname{tr}(U^\top R U),$$

where $R = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^\top \in \mathbb{R}^{d \times d}$ (we assume $\mathbb{E}[x] = 0$) is the covariance matrix.

Find a direction that maximizes the variance. For 1d case i.e., $u \in \mathbb{R}^d$, $\operatorname{tr}(u^\top R u) = \frac{1}{n} \sum_{i=1}^{n} (u^\top x_i)^2$ and $\mathbb{E}[u^\top x] = 0$.

# Obtain the first principal component

To obtain the first principal component:

$$\underset{\boldsymbol{u}^\top \boldsymbol{u}=1.}{\text{argmax}} \quad \boldsymbol{u}^\top \boldsymbol{R} \boldsymbol{u} = \underset{\boldsymbol{u}}{\text{argmax}} \frac{\boldsymbol{u}^\top \boldsymbol{R} \boldsymbol{u}}{\|\boldsymbol{u}\|_2^2},$$

where $\frac{\boldsymbol{u}}{\|\boldsymbol{u}\|_2}$ is a unit vector and $\frac{\boldsymbol{u}^\top \boldsymbol{R} \boldsymbol{u}}{\|\boldsymbol{u}\|_2^2}$ is called as the Rayleigh quotient.

Using the Lagrange multiplier $\lambda$ to find a critical point:

$$L(\boldsymbol{u}) = \boldsymbol{u}^\top \boldsymbol{R} \boldsymbol{u} - \lambda(\boldsymbol{u}^\top \boldsymbol{u} - 1)$$

To take the derivative with respect to $\boldsymbol{u}$, we have

$$\frac{\partial L(\boldsymbol{u})}{\partial \boldsymbol{u}} = 2\boldsymbol{R}\boldsymbol{u} - 2\lambda\boldsymbol{u} = \boldsymbol{0} \rightarrow \boldsymbol{R}\boldsymbol{u} = \lambda\boldsymbol{u}.$$

This is an eigenvalue decomposition problem where $\lambda$ is the eigenvalue and $\boldsymbol{u}$ is the eigenvector. Variance is $\boldsymbol{u}^\top \boldsymbol{R} \boldsymbol{u} = \lambda$.

# PCA with eigenvalue decomposition

PCA can be solved by using eigenvalue decomposition of the covariance matrix $\boldsymbol{R}$!

The eigenvalue decomposition of covariance matrix $\boldsymbol{R} \in \mathbb{R}^{d \times d}$:

$$\boldsymbol{R} = \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^\top \text{ or } \boldsymbol{U}^\top \boldsymbol{R}\boldsymbol{U} = \boldsymbol{\Lambda}$$

where

- $\boldsymbol{\Lambda} = \mathsf{diag}(\lambda_1, \lambda_2, \ldots, \lambda_d) \in \mathbb{R}^{d \times d}$, $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$. If $\boldsymbol{R}$ is a positive definite matrix $\lambda_d \geq 0$.
- $\boldsymbol{U} \in \mathbb{R}^{d \times d}$ is an orthogonal matrix $\boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{U}\boldsymbol{U}^\top = \boldsymbol{I}_d$
- $\mathsf{tr}(\boldsymbol{U}^\top \boldsymbol{R}\boldsymbol{U}) = \mathsf{tr}(\boldsymbol{U}^\top \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^\top \boldsymbol{U}) = \sum_{i=1}^{d} \lambda_i$.

# Relationship to Linear Auto-encoder (1/2)

Assume that $\mathbb{E}[\boldsymbol{x}] = \boldsymbol{0}$. Then, consider the following linear Auto-encoder problem:

$$\widehat{\boldsymbol{U}} = \underset{\boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{I}}{\operatorname{argmin}} \quad \frac{1}{n} \sum_{i=1}^{n} \|\boldsymbol{x}_i - \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{x}_i\|_2^2,$$

The loss function term can be written as

$$\frac{1}{n} \sum_{i=1}^{n} \|\boldsymbol{x}_i - \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{x}_i\|_2^2 = \frac{1}{n} \sum_{i=1}^{n} \left( \boldsymbol{x}_i^\top \boldsymbol{x}_i - 2\boldsymbol{x}_i^\top \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{x}_i + \boldsymbol{x}_i^\top \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{x}_i \right)$$

$$\propto -\frac{1}{n} \sum_{i=1}^{n} \left( \boldsymbol{x}_i^\top \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{x}_i \right) \quad (\boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{I})$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \left( \operatorname{tr}(\boldsymbol{U}^\top \boldsymbol{x}_i \boldsymbol{x}_i^\top \boldsymbol{U}) \right) \quad (\operatorname{tr}(\boldsymbol{A}\boldsymbol{B}) = \operatorname{tr}(\boldsymbol{B}\boldsymbol{A}))$$

$$= -\operatorname{tr}(\boldsymbol{U}^\top \boldsymbol{R}\boldsymbol{U}), \quad (\boldsymbol{R} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{x}_i \boldsymbol{x}_i^\top)$$

The minimization problem can be written as the maximization problem:

$$\operatorname*{argmin}_{\boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{I}} \frac{1}{n} \sum_{i=1}^{n} \|\boldsymbol{x}_i - \boldsymbol{U}\boldsymbol{U}^\top \boldsymbol{x}_i\|_2^2, \leftrightarrow \operatorname*{argmax}_{\boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{I}} \operatorname{tr}(\boldsymbol{U}^\top \boldsymbol{R} \boldsymbol{U})$$

Thus, PCA is related to the linear Auto-encoder.

# Nonlinear Auto-encoder

We consider the following Auto-encoder problem:

$$\widehat{\boldsymbol{\Theta}} = \underset{\boldsymbol{\Theta}}{\operatorname{argmin}} \quad \frac{1}{n} \sum_{i=1}^{n} \| \boldsymbol{x}_i - \boldsymbol{f_x}(\boldsymbol{g_x}(\boldsymbol{x}_i)) \|_2^2,$$

The nonlinear auto-encoder can be illustrated as

# Stochastic Neighbor Embedding (SNE)

The asymmetric probability $p_{ij}$ that $i$-th sample would pick $j$-th sample as its neighbor:

$$p_{ij} = \frac{\exp(-d_{ij}^2)}{\sum_{k \neq i} \exp(-d_{ki}^2)}, \quad d_{ij}^2 = \frac{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2}{2\sigma_i^2},$$

where $\sigma_i$ is a tuning parameter.

The model:

$$q_{ij} = \frac{\exp(-\|\boldsymbol{y}_i - \boldsymbol{y}_j\|_2^2)}{\sum_{k \neq i} \exp(-\|\boldsymbol{y}_k - \boldsymbol{y}_i\|_2^2)}$$

Optimization:

$$\widehat{\boldsymbol{y}}_1, \ldots, \widehat{\boldsymbol{y}}_n = \underset{\boldsymbol{y}_1, \ldots, \boldsymbol{y}_n}{\operatorname{argmin}} \quad \sum_{i=1}^{n} \sum_{j=1}^{n} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

# Symmetric SNE

The symmetric probability $p_{ij}$ that $i$-th sample would pick $j$-th sample as its neighbor:

$$p_{ij} = \frac{\exp(-d_{ij}^2)}{\sum_{k \neq l} \exp(-d_{kl}^2)}, \quad d_{ij}^2 = \frac{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2}{2\sigma^2},$$

where $\sigma$ is a tuning parameter.

The model:

$$q_{ij} = \frac{\exp(-\|\boldsymbol{y}_i - \boldsymbol{y}_j\|_2^2)}{\sum_{k \neq l} \exp(-\|\boldsymbol{y}_k - \boldsymbol{y}_l\|_2^2)}$$

Optimization:

$$\widehat{\boldsymbol{y}}_1, \ldots, \widehat{\boldsymbol{y}}_n = \operatorname*{argmin}_{\boldsymbol{y}_1, \ldots, \boldsymbol{y}_n} \quad \sum_{i=1}^{n} \sum_{j=1}^{n} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

# t-Stochastic Neighbor Embedding (t-SNE)

The asymmetric probability $p_{ij}$ that $i$-th sample would pick $j$-th sample as its neighbor:

$$p_{ij} = \frac{\exp(-d_{ij}^2)}{\sum_{k \neq l} \exp(-d_{ik}^2)}, \quad d_{ij}^2 = \frac{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2}{2\sigma^2},$$

where $\sigma$ is a tuning parameter.
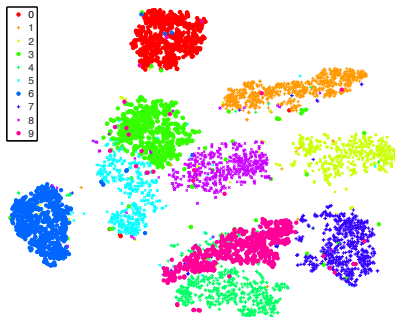
The model (Cauchy distribution):

$$q_{ij} = \frac{(1 + \|\boldsymbol{y}_i - \boldsymbol{y}_j\|_2^2)^{-1}}{\sum_{k \neq l}(1 + \|\boldsymbol{y}_k - \boldsymbol{y}_l\|_2^2)^{-1}}$$

Optimization:

$$\widehat{\boldsymbol{y}}_1, \ldots, \widehat{\boldsymbol{y}}_n = \operatorname*{argmin}_{\boldsymbol{y}_1, \ldots, \boldsymbol{y}_n} \quad \sum_{i=1}^{n} \sum_{j=1}^{n} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

Image taken from [1]



(a) Visualization by t-SNE.

t-SNE is heavily used in biology data such as the expression data.

# Multi-modal Dimensionality Reduction

PCA and auto-encoders are for uni-modal input (i.e., only image or only text).

How to do dimensionality reduction for multi-modal data (i.e., image and text)?

We have $(\boldsymbol{x}, \boldsymbol{y})$, where $\boldsymbol{x} \in \mathbb{R}^{d_x}$ and $\boldsymbol{y} \in \mathbb{R}^{d_y}$.

- Linear dimension reduction $\boldsymbol{z}_x = \boldsymbol{U}^\top \boldsymbol{x}$ and $\boldsymbol{z}_y = \boldsymbol{V}^\top \boldsymbol{y}$. $\boldsymbol{U} \in \mathbb{R}^{d_x \times m}$ and $\boldsymbol{V} \in \mathbb{R}^{d_y \times m}$.

- Nonlinear dimension reduction $\boldsymbol{z}_x = \boldsymbol{g_x}(\boldsymbol{x})$ and $\boldsymbol{z}_y = \boldsymbol{g_y}(\boldsymbol{y})$.

# Canonical Correlation Analysis (1/3)

Canonical Correlation Analysis (CCA) is to find dimensionality reduction that maximize the similarity between $\boldsymbol{z}_x = \boldsymbol{U}^\top \boldsymbol{x}$ and $\boldsymbol{z}_y = \boldsymbol{V}^\top \boldsymbol{y}$.

Assume that $\mathbb{E}[\boldsymbol{x}] = \boldsymbol{0}$ and $\mathbb{E}[\boldsymbol{y}] = \boldsymbol{0}$.

$$\mathsf{Corr}(X, Y) = \frac{\frac{1}{n} \sum_{i=1}^n \boldsymbol{z}_{x,i}^\top \boldsymbol{z}_{y,i}}{\sqrt{\frac{1}{n} \sum_{i=1}^n \boldsymbol{z}_{x,i}^\top \boldsymbol{z}_{x,i}} \sqrt{\frac{1}{n} \sum_{i=1}^n \boldsymbol{z}_{y,i}^\top \boldsymbol{z}_{y,i}}}$$

$$\frac{1}{n} \sum_{i=1}^n \boldsymbol{z}_{x,i}^\top \boldsymbol{z}_{y,i} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}_i^\top \boldsymbol{U} \boldsymbol{V}^\top \boldsymbol{y}_i$$
$$= \mathsf{tr}(\boldsymbol{U}^\top \boldsymbol{R}_{xy} \boldsymbol{V})$$

where $\boldsymbol{R}_{xy} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}_i \boldsymbol{y}_i^\top \in \mathbb{R}^{d_x \times d_y}$.

# Canonical Correlation Analysis (2/3)

The optimization problem of CCA is given as

$$\widehat{U}, \widehat{V} = \underset{U, V}{\operatorname{argmax}} \quad \operatorname{tr}(U^\top R_{xy} V),$$

$$\text{s.t.} \quad U^\top R_{xx} U = I, V^\top R_{yy} V = I,$$

where $R_{xx} = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^\top$ and $R_{yy} = \frac{1}{n} \sum_{i=1}^{n} y_i y_i^\top$.

Then, CCA can be written as

$$\max_{U, V} \quad \operatorname{tr}\left( \begin{bmatrix} U^\top & V^\top \end{bmatrix} \begin{bmatrix} O & R_{xy} \\ R_{xy}^\top & O \end{bmatrix} \begin{bmatrix} U \\ V \end{bmatrix} \right),$$

$$\text{s.t.} \quad \begin{bmatrix} U^\top & V^\top \end{bmatrix} \begin{bmatrix} R_{xx} & O \\ O & R_{yy} \end{bmatrix} \begin{bmatrix} U \\ V \end{bmatrix} = I,$$

This is a generalized eigenvalue decomposition (GEV) problem.

Let us transform the variables as

$$
\left[ \begin{array}{c} \bar{U} \\ \bar{V} \end{array} \right] = \left[ \begin{array}{cc} R_{xx}^{1/2} & O \\ O & R_{yy}^{1/2} \end{array} \right] \left[ \begin{array}{c} U \\ V \end{array} \right]
$$

we can rewrite the CCA optimization problem as

$$
\max_{\bar{U}, \bar{V}} \ \frac{1}{2} \mathrm{tr} \left( \left[ \begin{array}{cc} \bar{U}^{\top} & \bar{V}^{\top} \end{array} \right] \left[ \begin{array}{cc} O & R_{xx}^{-1/2} R_{xy} R_{yy}^{-1/2} \\ (R_{xx}^{-1/2} R_{xy} R_{yy}^{-1/2})^{\top} & O \end{array} \right] \left[ \begin{array}{c} \bar{U} \\ \bar{V} \end{array} \right] \right),
$$

$$
\text{s.t.} \ \left[ \begin{array}{cc} \bar{U}^{\top} & \bar{V}^{\top} \end{array} \right] \left[ \begin{array}{c} \bar{U} \\ \bar{V} \end{array} \right] = I,
$$

Thus, we can solve the CCA problem by using eigenvalue decomposition!

# Other dimensionality reduction methods

- Fisher Discriminant Analysis (FDA)
- Independent Component Analysis (ICA)
- Sufficient Dimensionality Reduction (SDR)
- Locally Linear Embedding (LLE)
- etc.

# Summary of today's lecture

- Dimensionality reduction (Reduce the dimensionality of features).
- Feature selection is to interpret features, while dimensionality reduction is to reduce the dimensionality (for compression and visualization).
- Principal Component Analysis (PCA)
- Canonical Correlation Analysis (Multi-modal data)

[1] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.