

# 3 カーネルの実際

## -3.2 カーネル正準相関分析

10NM733X 林華

# 3.2.1 正準相関分析の復習-1

□ 正準相関分析(CCA)：2種類の多変量データX,Yに依存関係を調べる線形手法である。

•X は p 次元のデータを n 個集めた  $n \times p$  のデータ行列 .

•Y は q 次元のデータを n 個集めた  $n \times q$  のデータ行列 .  $p \leq q$  としておく .

共分散行列：

• $S_{XX} = 1/n X^{\sim T} X^{\sim}$

• $S_{YY} = 1/n Y^{\sim T} Y^{\sim}$

• $S_{XY} = 1/n X^{\sim T} Y^{\sim}$

(ただし,  $X^{\sim}$  と  $Y^{\sim}$  は, それぞれ X と Y から, それぞれの列方向の平均値を, 各要素から引いた平均偏差行列 . )

p 次元のベクトル a と q 次元のベクトル b を導入 . これらのベクトルと,  $X^{\sim}$  や  $Y^{\sim}$  とを一次結合したベクトルの相関を考える .

$$\rho = \max_{a \in R^m, b \in R^l} \frac{a^T S_{XY} b}{\sqrt{a^T S_{XX} a} \sqrt{b^T S_{YY} b}}$$

この式を最大化する a と b を, 第1正準相関ベクトルと呼び,  $a_1$  と  $b_1$  で表す .

第k 正準相関ベクトル( $k=1, \dots, p$ )は,  $a_i^T S_{XX} a_i = 1, b_i^T S_{YY} b_i = 1$ を満たすものの中で,  $\rho$ を最大化するように求める .

## 3.2.1 正準相関分析の復習-2

ここで $V_{XY}, V_{XX}, V_{YY}$ は標本分散共分散行列で、上の最大化問題は

$$\max a^T V_{XY} b \quad \text{subject to} \quad a^T V_{XX} a = b^T V_{YY} b = 1$$

という形に書き直せる．さらにLagrange乗数法を適用することにより、以下の一般化固有値問題の最大固有値に一致

$$\begin{pmatrix} \mathbf{O} & V_{XY} \\ V_{YX} & \mathbf{O} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \rho \begin{pmatrix} V_{XX} & \mathbf{O} \\ \mathbf{O} & V_{YY} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix}$$

## 3.2.2 カーネルCCA-1

- 独立同分布データ  $X, Y$
- 対応する再生核ヒルベルト空間  $H_X, H_Y$
- 特徴写像  $\Phi_X, \Phi_Y$ 、 $\tilde{\Phi}_X(X_i)$  は標本平均を差し引いた特徴ベクトル

正規相関を与える方向ベクトル  $f \in H_X, g \in H_Y$  と正準関数  $\rho$  は次と定義

$$\rho = \max_{f \in H_X, g \in H_Y} \frac{\sum_{i=1}^N \left\langle f, \tilde{\Phi}_X(X_i) \right\rangle_{H_X} \left\langle g, \tilde{\Phi}_Y(Y_i) \right\rangle_{H_Y}}{\sqrt{\sum_{i=1}^N \left\langle f, \tilde{\Phi}_X(X_i) \right\rangle_{H_X}^2} \sqrt{\sum_{i=1}^N \left\langle g, \tilde{\Phi}_Y(Y_i) \right\rangle_{H_Y}^2}}$$

ここで上式を最大にする  $f \in H_X, g \in H_Y$  は下式の形と考えればよい。

$$f = \sum_{i=1}^N \alpha_i \tilde{\Phi}_X(X_i), \quad g = \sum_{i=1}^N \beta_i \tilde{\Phi}_Y(Y_i)$$

$$\rho = \max_{\alpha \in R^N, \beta \in R^N} \frac{\alpha^T \tilde{K}_X \tilde{K}_Y \beta}{\sqrt{\alpha^T \tilde{K}_X^2 \alpha} \sqrt{\beta^T \tilde{K}_Y^2 \beta}}$$

## 3.2.2 カーネルCCA-2

問題：  $u = \tilde{K}_X \alpha, v = \tilde{K}_Y \beta$  とすれば、

$$\rho = \max_{u \in R(K_X), v \in R(K_Y)} \frac{u^T v}{\|u\| \|v\|}$$

対策：正則化係数  $\varepsilon_N$  を導入して正則化

$$\|f\|^2 = \alpha^T \tilde{K}_X \alpha, \|g\|^2 = \beta^T \tilde{K}_Y \beta \quad \text{により}$$

$$\rho = \max_{f \in H_x, g \in H_y} \frac{\sum_{i=1}^N \left\langle f, \tilde{\Phi}_x(X_i) \right\rangle_{H_x} \left\langle g, \tilde{\Phi}_y(Y_i) \right\rangle_{H_y}}{\sqrt{\sum_{i=1}^N \left\langle f, \tilde{\Phi}_x(X_i) \right\rangle_{H_x}^2 + N\varepsilon_N \|f\|^2} \sqrt{\sum_{i=1}^N \left\langle g, \tilde{\Phi}_y(Y_i) \right\rangle_{H_y}^2 + N\varepsilon_N \|g\|^2}}$$

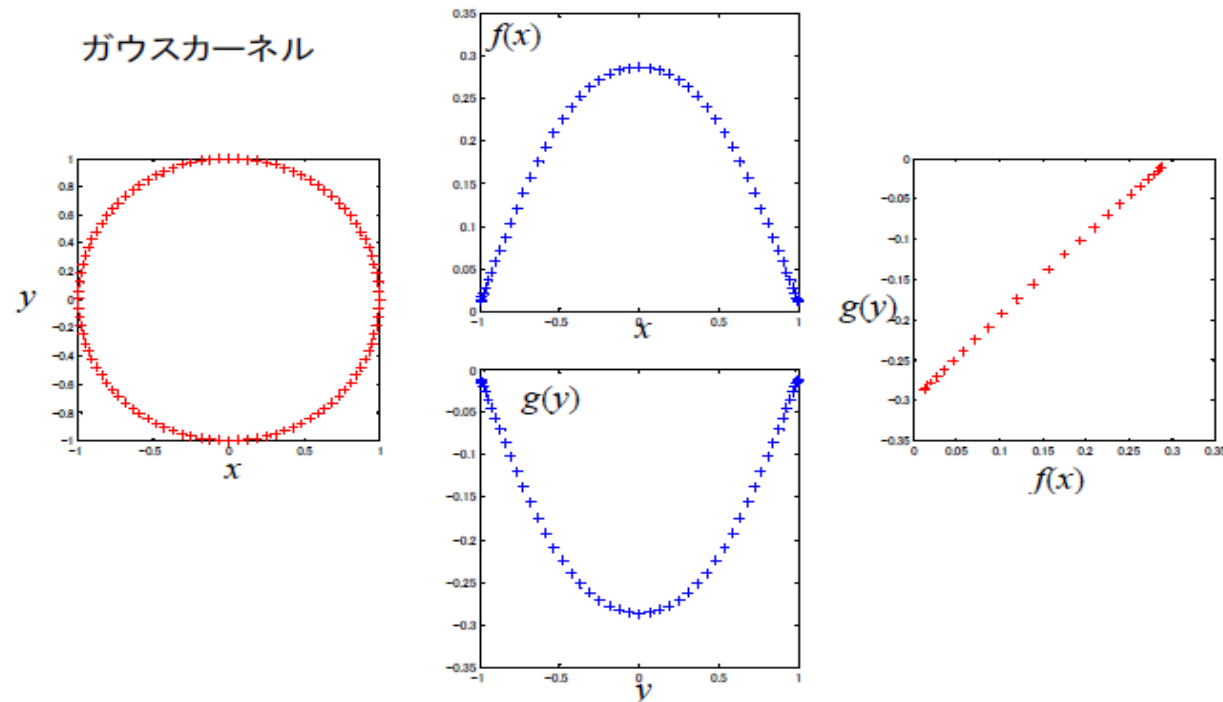
$$\rho = \max_{\alpha \in R^N, \beta \in R^N} \frac{\alpha^T \tilde{K}_X \tilde{K}_Y \beta}{\sqrt{\alpha^T \tilde{K}_X^2 \alpha + N\varepsilon_N \tilde{K}_X \alpha} \sqrt{\beta^T \tilde{K}_Y^2 \beta + N\varepsilon_N \tilde{K}_Y \beta}}$$

$$\begin{pmatrix} O & \tilde{K}_X \tilde{K}_Y \\ \tilde{K}_Y \tilde{K}_X & O \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \rho \begin{pmatrix} \tilde{K}_X^2 + N\varepsilon_N \tilde{K}_X & O \\ O & \tilde{K}_Y^2 + N\varepsilon_N \tilde{K}_Y \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

# 3.2.3 カーネルCCAの応用例-1

- $X, Y$ : 円上にある100点からなる1次元データ
- $k(x, y) = \exp(-1/2(x-y)^2)$
- $\varepsilon = 1/10^5$

- カーネルCCAの実験例



## 3.2.3 カーネルCCAの応用例-2

- X: image 色情報データ, Y: text 単語の頻度データ
- それぞれ400データ、中の200データがテスト用

- データにカーネルCCAを実施し、大きいd個の固有値に対応する固有ベクトル $f_1, \dots, f_d, g_1, \dots, g_d$ を求める

- 各画像 $X_i$ に対するd次元部分空間への正写像を計算  $\xi_i = \left( \left\langle \tilde{\Phi}_x(X_i), f_a \right\rangle_{H_x} \right)_{a=1}^d \in R^d$

- 新たなテキストデータ $Y_{new}$ が与えられたとき、その特徴を計算

$$\zeta = \left( \left\langle \tilde{\Phi}_y(Y_{new}), g_a \right\rangle_{H_x} \right)_{a=1}^d \in R^d$$

- テキストに最も適した画像データとして  $\arg \max_i \xi^T \zeta$  を与える  $X_i$ を出力する