

- パターン認識の手順
  - 特徴抽出: 特徴ベクトル(パターン)の定義
  - 学習パターンを用いて特徴空間を各クラスに分割
- 問題点
  - 特徴ベクトルの各成分間のスケールリング, 正規化
  - 特徴空間の次元数, 次元の呪い  
特徴の数を増やせば情報量が増え, 認識率も向上するという誤解
    - \* 特徴の数を増やせば増やすほど相関の高い特徴の組が混入する可能性が高まる
    - \* 統計計算に要する計算量は少なくとも次元のべき乗のオーダーになる. 計算量の爆発.
    - \* 有限個のパターンから識別器を設計する際, 次元を高くしていくと誤識別率がかえって上昇する. ヒューズ現象
- 情報量正規化
- 次元削減
  - KL展開: パターン全体の分布をもっともよく近似する部分空間
  - 線形判別法: 各クラスのパターン分布の分離度を最大にする部分空間

- 元特徴ベクトル:  $\mathbf{x}$  , 変換後のベクトル:  $\mathbf{y}$  , 変換行列:  $A$

$$\mathbf{y} = A\mathbf{x}$$

- パターン:  $\mathbf{x}_p = (x_{p1}, x_{p2}, \dots, x_{pd})^T$

- 変換行列:  $A = \text{diag}(a_1, a_2, \dots, a_d)$

- 変換後のベクトル:  $\mathbf{y}(y_{pi} = a_i x_{pi})$

- $p$  番目のパターンとその他のパターンの平均距離

$$r_p^2 = \frac{1}{n-1} \sum_{q=1}^n \sum_{j=1}^d (y_{pj} - y_{qj})^2$$

- 全パターンの平均距離

$$\begin{aligned} R^2 &= \frac{1}{n} \sum_{p=1}^n r_p^2 = \frac{1}{n(n-1)} \sum_{p=1}^n \sum_{q=1}^n \sum_{j=1}^d a_j^2 (x_{pj} - x_{qj})^2 \\ &= \frac{n}{n-1} \sum_{j=1}^d a_j^2 \frac{1}{n} \sum_{p=1}^n \sum_{q=1}^n \left( \frac{1}{n} x_{pj}^2 - \frac{2}{n} x_{pj} x_{qj} + \frac{1}{n} x_{qj}^2 \right) \\ &= \frac{n}{n-1} \sum_{j=1}^d a_j^2 \left( \frac{1}{n} \sum_{q=1}^n \frac{1}{n} \sum_{p=1}^n x_{pj}^2 \right. \\ &\quad \left. - 2 \frac{1}{n} \sum_{p=1}^n x_{pj} \frac{1}{n} \sum_{q=1}^n x_{qj} + \frac{1}{n} \sum_{p=1}^n \frac{1}{n} \sum_{p=1}^n x_{qj}^2 \right) \\ &= \frac{2n}{n-1} \sum_{j=1}^d a_j^2 (\bar{x}_j^2 - \bar{x}_j^2) \\ &= 2 \sum_{j=1}^d a_j^2 \sigma_j^2 \end{aligned}$$

- 特徴の分散

$$\sigma_j^2 = \frac{1}{n-1} \sum_{p=1}^n (x_{pj} - \bar{x}_j)^2 = \frac{n}{n-1} \sum_{p=1}^n (\bar{x}_j^2 - \bar{x}_j^2)$$

# パターン相互の距離の最大化

---

- 制約条件 (特徴空間の体積不変)

$$\prod_{j=1}^d a_j = 1$$

- Lagrange 未定乗数法

$$L = 2 \sum_{j=1}^d a_j^2 \sigma_j^2 - \lambda \left( \prod_{j=1}^d a_j - 1 \right)$$

$$\frac{\partial L}{\partial a_j} = 4a_j \sigma_j^2 - \lambda \prod_{k \neq j} a_k = 0$$

$$a_j = \frac{\sqrt{\lambda}}{2\sigma_j}, \quad \lambda = 4 \left( \prod_{j=1}^d \sigma_j \right)^{2/d}$$

- 最適係数

$$a_j = \frac{1}{\sigma_j} \left( \prod_{k=1}^d \sigma_k \right)^{1/d}$$

各特徴軸を標準偏差で正規化し，平均の回りの分散 (すなわちパターンの広がり) を均一化

- Karhunen-Loève展開: 線形空間における特徴ベクトルの分布をもっともよく近似する部分空間を求める方法
  - 分散最大基準
  - 平均二乗誤差最小基準
- 手法は多変量解析(主成分分析)
- $d$ 次元の特徴量ベクトル  $\mathbf{x}$  から  $\bar{d}$  ( $< d$ )次元の特徴  $\mathbf{y}$  への変換  $A$

$$\mathbf{y} = A\mathbf{x}, \quad A = (\mathbf{u}_1, \dots, \mathbf{u}_{\bar{d}}), \quad A^T A = I$$

- 分散最大基準

$$\begin{aligned}\tilde{\sigma}^2(A) &= \frac{1}{n} \sum_{\mathbf{y}} (\mathbf{y} - \tilde{\mathbf{m}})^T (\mathbf{y} - \tilde{\mathbf{m}}) \quad (\tilde{\mathbf{m}} = A^T \mathbf{m}) \\ &= \text{tr}(A^T \Sigma A) \quad (\Sigma = \frac{1}{n} \sum_{\mathbf{x}} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T)\end{aligned}$$

$$\max\{\tilde{\sigma}^2(A)\} = \max\{\text{tr}(A^T \Sigma A)\} = \max\{\Lambda\} = \sum_{i=1}^{\bar{d}} \lambda_i$$

- 平均二乗誤差最小

$$\begin{aligned}\varepsilon^2(A) &= \frac{1}{n} \sum (A\mathbf{y} - \mathbf{x})^T (A\mathbf{y} - \mathbf{x}) \\ &= \text{tr}R - \text{tr}(A^T R A) \quad (R = \frac{1}{n} \sum \mathbf{x}\mathbf{x}^T)\end{aligned}$$

$$\min\{\varepsilon^2(A)\} = \text{tr}R - \sum_{i=1}^{\bar{d}} \lambda_i$$

# 線形判別法 (Fisher's discriminant method)

- 特徴空間上の2クラスのパターンの分布からこの2クラスを識別するのに最適な1次元軸を求める手法(クラス内変動・クラス内分布比最大)

- 変動行列:  $S_i$

$$S_i = \sum_{\mathbf{x} \in \mathcal{X}_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T$$

- クラス内変動行列:  $S_W$

$$S_W = S_1 + S_2 = \sum_{i=1,2} \sum_{\mathbf{x} \in \mathcal{X}_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T$$

- クラス間変動行列:  $S_B$

$$S_B = \sum_{i=1,2} n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T = \frac{n_1 n_2}{n} (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$$

- $d$ 次元特徴空間から1次元空間への変換:  $A \quad y = A\mathbf{x}$

- 変換後のクラス内変動行列:  $\tilde{S}_W$ , クラス間変動行列:  $\tilde{S}_B$

$$\tilde{S}_W = \tilde{S}_1 + \tilde{S}_2 = \sum_{i=1,2} \sum_{y \in \mathcal{Y}_i} (y - \tilde{m}_i)^2 = n_1 \tilde{\sigma}_1^2 + n_2 \tilde{\sigma}_2^2 = A^T S_W A$$

$$\tilde{S}_B = \sum_{i=1,2} n_i (\tilde{m}_i - \tilde{m})^2 = \frac{n_1 n_2}{n} (\tilde{m}_1 - \tilde{m}_2)^2 = A^T S_B A$$

- フィッシャーの評価基準(クラス間変動のクラス内変動に対する比)

$$J_S(A) = \frac{\tilde{S}_B}{\tilde{S}_W} = \frac{n_1 n_2 (\tilde{m}_1 - \tilde{m}_2)^2}{n (n_1 \tilde{\sigma}_1^2 + n_2 \tilde{\sigma}_2^2)} = \frac{A^T S_B A}{A^T S_W A}$$

# 評価基準最大化

---

- $J_S(A)$  の最大化:  $\tilde{S}_W = I$  の条件下で  $\tilde{S}_B$  の最大化
- Lagrange 乗数

$$J(A) = A^T S_B A - \lambda(A^T S_W A - I)$$

$$\frac{\partial J(A)}{\partial A} = S_B A - \lambda S_W A = 0$$

- $S_W$  が正則なら

$$(S_W^{-1} S_B - \lambda I) A = 0$$

- $S_W^{-1} S_B$  の最大固有値を  $\lambda_1$  とすると

$$\max\{J_S(A)\} = \lambda_1$$

$J_S$  を最大にする  $A$  は最大固有値に対応する固有ベクトル

# KL 展開と線形判別法

- KL 展開

- 特徴ベクトルの分布全体の持つ情報をなるべく最大限に反映できるように特徴空間の次元を削減する方法
- クラスの識別に有効かどうかはわからない
- 目的は表現のための (圧縮のための) 次元削減

- 線形判別法

- クラスの分布の分離度を考慮した空間の変換方法
- 判別のための次元削減

- それにも関わらず KL 展開が用いられる理由

- 文字認識や音声認識などの高度な認識には，高次元の特徴ベクトルが必要であり，次元の呪いから逃れるためには次元削減が不可欠である．
- はじめに選ばれた (人間が選んだ) 特徴には相関を持つ特徴の組が含まれている可能性がある．

