

EE 7403 LEC 12 Feature Extraction/Dimension Reduction v. Machine Learning

1. Intro

Pattern Recognition is a series of processes that reduce the dimensionality and the variation of samples for the same class and keep their discrimination for different classes.

⇒ Dimensionality Reduction & Discriminative info extraction.

feature extraction & dimensionality reduction

- human expert knowledge: fingerprint (minutia points)
- image local structure: corners, blobs, interesting points
- image global structure
- Machine Learning from training database

2. Feature Extraction Based on Image Global Structure

Transform image $f(x, y)$ to feature $g(u, v)$

$$g(u, v) = T\{f(x, y)\}$$

$$= \iint w(u, v, x, y) f(x, y) dx dy \quad \text{linear transform}$$

$$\Rightarrow \sum_1^h \sum_1^w w(u, v, x, y) f(x, y) \quad \text{for digital image}$$

$$\Rightarrow \sum_1^h \sum_1^w e^{-2\pi j(ux+vy)} f(x, y) \quad \text{Fourier transform}$$

$$\Rightarrow \sum_1^h \sum_1^w x^u y^v f(x, y) \quad \text{moments computing}$$

$$\Rightarrow f = W^T x \quad \text{vector-matrix representation}$$

Polar Complex Exponential Transform (PCET)

$$g(u, v) = \frac{1}{\pi} \int_0^{2\pi} \int_0^1 e^{-j(u^2 r^2 \cos^2 \theta + v^2 r^2 \sin^2 \theta)} f(r, \theta) dr d\theta$$

3. Principal Component Analysis *unsupervised*

Given q n -dim training samples:

$$x_1, x_2, \dots, x_q$$

$$\mu = \frac{1}{q} \sum_{i=1}^q x_i$$

$$\tilde{x}_i = x_i - \mu \quad X = [\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_q]$$

use 1-dim to represent \tilde{x}_i best

$$a_i = \phi^T \tilde{x}_i \quad (\phi, \|\phi\|^2 = \phi^T \phi = 1)$$

reduce computational complexity

best info maybe x important for

data representation but \downarrow for discriminating possibly.

EE7403 LEC 12

the best ϕ makes reconstruction error minimum

$$\epsilon^2 = \sum_{i=1}^q \|\tilde{x}_i - a_i \phi\|^2 \Rightarrow \text{minimum}$$

$$\begin{aligned} \Rightarrow \epsilon^2 &= \sum \|\tilde{x}_i - a_i \phi\|^2 = \sum (\tilde{x}_i - a_i \phi)^T (\tilde{x}_i - a_i \phi) = \sum (\tilde{x}_i^T \tilde{x}_i - 2a_i \phi^T \tilde{x}_i + a_i^2 \phi^T \phi) \\ &= \sum (\tilde{x}_i^T \tilde{x}_i - 2a_i^2 + a_i^2) = \sum (\tilde{x}_i^T \tilde{x}_i - a_i^2) \\ &= \sum \|\tilde{x}_i\|^2 - \sum a_i^2 = \sum \|\tilde{x}_i\|^2 - \sum (\phi^T \tilde{x}_i)(\phi^T \tilde{x}_i)^T \\ &= \sum \|\tilde{x}_i\|^2 - \phi^T \left(\sum \tilde{x}_i \tilde{x}_i^T \right) \phi \end{aligned}$$

the sample covariance matrix of all training data

$$S^t = \frac{1}{q} \sum_{i=1}^q (x_i - \mu)(x_i - \mu)^T = \frac{1}{q} \sum_{i=1}^q \tilde{x}_i \tilde{x}_i^T = \frac{1}{q} X X^T$$

S^t is total scatter matrix

the rank of S^t is $\min(q-1, n)$ or most

to minimize $\epsilon^2 = \sum \|\tilde{x}_i\|^2 - q \phi^T S^t \phi$

use Lagrange optimization

$$f(\phi, \lambda) = \phi^T S^t \phi - \lambda (\phi^T \phi - 1) \Rightarrow \text{maximum}$$

$$\frac{\partial f}{\partial \phi} = 2 S^t \phi - 2 \lambda \phi = 0$$

$$S^t \phi = \lambda \phi$$

\therefore the solution is the the eigenvalues and eigenvectors of S^t . (and eigenvectors ϕ satisfies $\phi^T \phi = 1$)

notes:

require $f(x, y)$ reaches maximum when $g(x, y) = c$

$$L(x, y, \lambda) = f(x, y) + \lambda (g(x, y) - c)$$

whose pole values contains pole values of $f(x, y)$

$$\Rightarrow \nabla L = 0$$

$$\epsilon^2 = \sum \|\tilde{x}_i\|^2 - q \phi^T S^t \phi = \sum \|\tilde{x}_i\|^2 - q \phi^T \lambda \phi = \sum \|\tilde{x}_i\|^2 - \lambda q$$

Reducing the x_i into lower-dim m -dim y_i by

$$y_i = \phi^T (x_i - \mu) \quad \phi = [\phi_1, \phi_2, \dots, \phi_m] \quad m < n$$

m largest eigenvalues of S^t

reconstruct

$$\hat{x}_i = \phi y_i + \mu$$

the reconstruction error $E[\|\Delta\|^2] = E[\Delta^T \Delta] = \sum_{k=m+1}^n \lambda_k$ (λ_k are λ started in descending order)

4. Eigen decomposition

eigenvalue and eigenvector :

$$\Sigma \phi_i = \lambda_i \phi_i \quad i=1, 2, \dots, n$$

if Σ is symmetric eigenvectors corresponding to the distinct eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ are orthogonal

Take the unit length of $\phi_1, \phi_2, \dots, \phi_n$

$$\phi_i^T \phi_j = \begin{cases} 1, & \text{if } i=j \\ 0, & \text{if } i \neq j \end{cases}$$

(orthogonal + unit length = orthonormal)

let Φ be the orthonormal matrix formed by eigenvectors

$$\Phi = [\phi_1, \phi_2, \dots, \phi_n]$$

$$\text{Obviously } \Phi^T \Phi = I \quad \therefore \Phi^T = \Phi^T \therefore \Phi \Phi^T = I$$

let Λ be a diagonal matrix $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$

from $\Sigma \phi_i = \lambda_i \phi_i$

$$\Rightarrow \Sigma \Phi = \Phi \Lambda \quad \therefore \boxed{\Sigma = \Phi \Lambda \Phi^T} \quad \text{or} \quad \boxed{\Lambda = \Phi^T \Sigma \Phi}$$

$$\Lambda = \Phi^T \Sigma \Phi = \Phi^T X X^T \Phi \xrightarrow{Y = \Phi^T X} Y Y^T = \Sigma_y = \Lambda$$

Problem of PCA:

① the lost info less important for representing data could be critical for discriminating

② PCA projects data vertically instead of horizontally.

while maximizes the variation

keeps irrelevant info to discriminate.

EE7403 LEC12.

5. Linear Discriminant Analysis (LDA)

properties to determine:

maximize separation between projected class mean

minimize projected within class scatter (variance)

$$y = \phi^T(x - \mu_i)$$

Given q n -dim samples of classes c

$$x_1, x_2, \dots, x_q$$

the num of samples in class w_j is q_j , $j=1, 2, \dots, c$.

$$\Sigma_j = \frac{1}{q_j} \sum_{x_i \in w_j} (x_i - \mu_j)(x_i - \mu_j)^T, \text{ where } \mu_j = \frac{1}{q_j} \sum_{x_i \in w_j} x_i$$

within-class scatter matrix is:

$$S^w = \sum_{j=1}^c \frac{q_j}{q} \Sigma_j$$

between-class scatter matrix

$$S^b = \sum_{j=1}^c \frac{q_j}{q} (\mu_j - \mu)(\mu_j - \mu)^T$$

$$\mu = \frac{1}{q} \sum_{i=1}^q x_i = \sum_{j=1}^c \frac{q_j}{q} \mu_j$$

$$\text{LDA: } \min \text{trace}[\phi^T S^w \phi] \quad \max \text{trace}[\phi^T S^b \phi]$$

$$\Rightarrow \max \text{trace}[\phi^T S^{w-1} S^b \phi] = \sum_{k=1}^m \lambda_k^{b/w}$$

$$S^t = S^w + S^b$$

Obviously, LDA extracts much more discriminative features than PCA

the rank of S^w is $\min[q-c, n]$ at most. (mostly $q-c \ll n$)

Problems: S^w is singular and its inverse does not exist.