

Note:

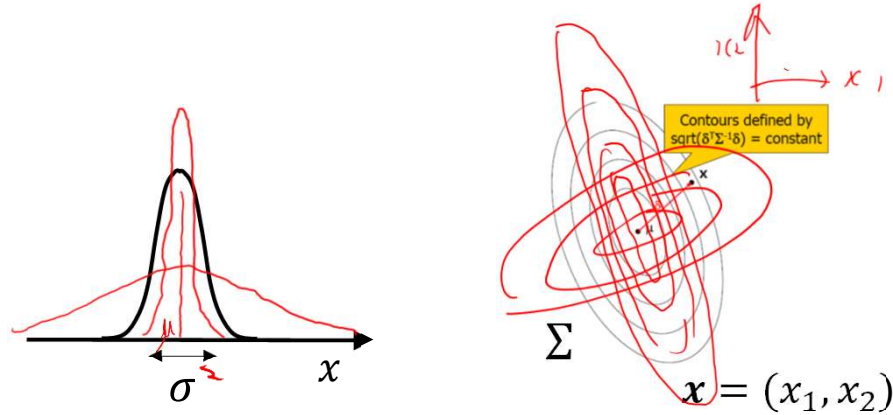
- Homework #4
 - On Lecture 8-9, **Due in two weeks**
 - Accessible on *Canvas* now.
 - Submit the PDF file to “Submission for HW #4” link by 4/15 4pm. No late submission. Strictly applied.
 - **Do not procrastinate this homework which could be time-consuming**
- Fast Prototyping Exercise #2 on Mean Shift continues (second session).
 - <https://bidal.sfsu.edu/~kazokada/csc872/PD2.pdf>

Statistical Modeling

PF: Gaussian Mixture Model

CSC 872
Pattern Analysis and Machine Intelligence

Let's be OK with Multivariate Gaussian



Variance controls the shape

Covariance controls the shape

3

Understanding Variance: Review

- **Univariate Domain:**
 - Given a random scalar variable X
 - **True Mean Definition:** $E[X] = \mu$
 - **True Variance Definition:** $Var[X] = E[(X-\mu)^2] = E[X^2] - (E[X])^2 = \sigma^2$

- **MLE of $P(X)$ as a Gaussian Distribution**

- Given a sample x_1, \dots, x_N drawn from a Gaussian

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

- **MLE of mean** is sample mean *unbiased*
- **MLE of variance** is sample variance *biased*

$$\mu^{mle} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\sigma_{mle}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu^{mle})^2$$

4

Covariance: Definition

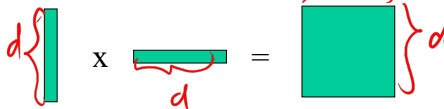
- **Multivariate Domain:**

- Given a random **column-vector** variable X 

- **True Mean Definition:**

- **True Covariance Definition:**

$$E[X] = \mu \quad \left. \vphantom{\mu} \right\} d$$

$$\text{Cov}[X] = E[(X-\mu)(X-\mu)^T] = \Sigma \quad \left. \vphantom{\Sigma} \right\} d$$


CSC872: PAMI – Kazunori Okada (C) 2025

5

5

Covariance: MLE for Gaussian

- **Multivariate Domain:**

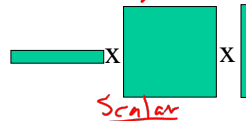
- **MLE of P(X) as a Gaussian Distribution**

- Given a sample x_1, \dots, x_N drawn from a Gaussian

$$N(x|\mu, \Sigma) = \frac{1}{|2\pi\Sigma|^{d/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

quadratic form

$$\mu^{mle} = \frac{1}{N} \sum_{n=1}^N x_n$$



biased!

$$\Sigma^{mle} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu^{mle})(x_n - \mu^{mle})^T$$

(1/(N-1)) unbiased!

CSC872: PAMI – Kazunori Okada (C) 2025

6

6

Understanding Covariance

- Symmetric & Square
 - Transpose of A is the same as A
- Positive semi-definite (non-negative definite)
 - Eigen values of A are all positive or zero
 - Quadratic function $x^T A x$ is positive or zero for all x
 - The power of exponent in the multivariate Gaussian is always negative!
- Ellipsoidal shape and Cov
 - Eigen vectors = Ellipsoidal axes
 - Eigen values = Ellipsoidal axis length

$$\Sigma = \Sigma^T$$

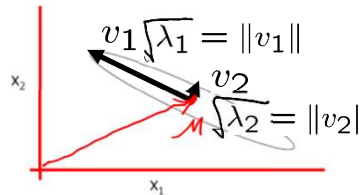
Positive definite

all positive

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{pmatrix} \quad \Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \dots & \sigma_{1m} \\ \sigma_{12} & \sigma_{22}^2 & \dots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \dots & \sigma_{mm}^2 \end{pmatrix}$$

$$v_1, \dots, v_n$$

$$\lambda_1, \dots, \lambda_n$$



CSC872: PAMI – Kazunori Okada (C) 2025

7

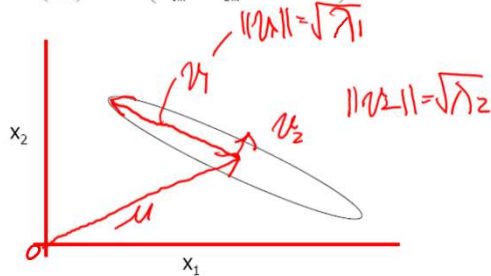
7

Understanding Covariance

- General Gaussian: Fully-valued Covariance
- Any oriented ellipsoidal shape

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{pmatrix} \quad \Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \dots & \sigma_{1m} \\ \sigma_{12} & \sigma_{22}^2 & \dots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \dots & \sigma_{mm}^2 \end{pmatrix}$$

Eigen decomp.



CSC872: PAMI – Kazunori Okada (C) 2025

8

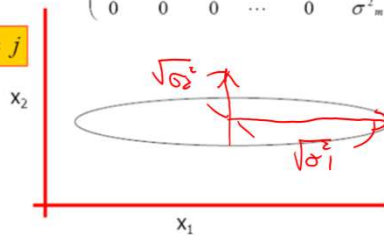
8

Understanding Covariance

- **Axis-Aligned Gaussian: Diagonal Covariance**
- Any axis-aligned ellipsoidal shape
- Every X_i are independent to each other

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma^2_m \end{pmatrix}$$

$X_i \perp X_j$ for $i \neq j$



CSC872: PAMI – Kazunori Okada (C) 2025

9

9

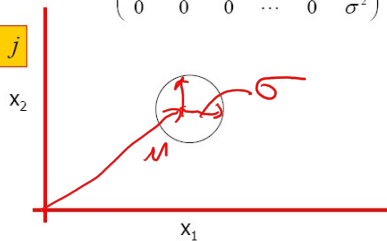
Understanding Covariance

- **Spherical Gaussian: $\Sigma = \sigma^2 I$**
- Spherical shape
- Independent & identical size

$$\Sigma = \sigma^2 I$$

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_m \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma^2 \end{pmatrix}$$

$X_i \perp X_j$ for $i \neq j$



CSC872: PAMI – Kazunori Okada (C) 2025

10

10

Two Statistical Modeling Approaches

- **Non-Parametric: *Histogram & KDE***
 - Yes: **Flexible**, accurately describe arbitrary distributions
 - No: **High Time and Space Complexity**
- **Parametric: *MLE & MAP***
 - Yes: **Low Time and Space Complexity**
 - No: **Rigid**, it may not be accurate
- **Any flexible but economic parametric model???**

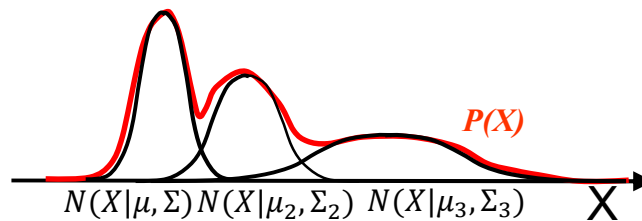
CSC872: PAMI – Kazunori Okada (C) 2025

11

11

YES

- **Gaussian Mixture Model**
 - Suppose $P(X)$ takes a form of a **weighted sum** of K different **Gaussian components**



$$P(X) = \sum_{k=1}^K \pi_k N(X|\mu_k, \Sigma_k)$$

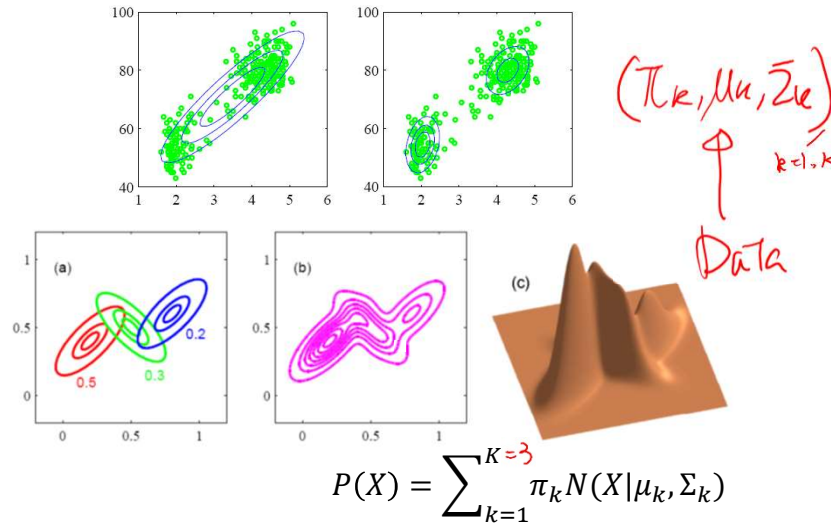
π_k : mixing weights $\sum_{k=1}^K \pi_k = 1$

CSC872: PAMI – Kazunori Okada (C) 2025

12

12

In 2D



CSC872: PAMI – Kazunori Okada (C) 2025

13

13

PF: Mixture Model

- More general form of a Mixture Model
- A mixture (=weighted sum) of arbitrary component distributions parameterized by θ_k

$$P(X) = \sum_{k=1}^N \pi_k P(X|\theta_k)$$

Data \rightarrow $\{\pi_1, \dots, \pi_N\}$
 $\{\theta_1, \dots, \theta_N\}$

CSC872: PAMI – Kazunori Okada (C) 2025

14

14

Sampling from GMM

Gaussian Mixture Model

- Sampling: ^{random} The Inverse Problem of Modeling

- Draw a set of data points from a known prob. dist.

- Sampling from GMM $P(X)$ is two-step!

- 1) Pick one of Gaussians according to $\pi = \pi_1, \dots, \pi_K$

- 2) Generate $x \sim$ the chosen Gaussian component

- Latent variable $Z \in \{1, \dots, K\}$

- A random variable that picks one of the Gaussians!

- $\pi_k = P(Z=k|\pi)$

$$P(X) = \sum_{k=1}^K P(Z = k|\pi) N(X|\mu_k, \Sigma_k)$$

{ μ_k, Σ_k, π_k }

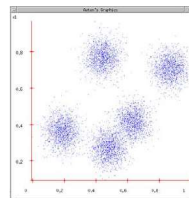
CSC872: PAMI – Kazunori Okada (C) 2025

15

15

STOP: Problem

- You have N data points
- You know they all come from K Gaussian Mixture
- Let me ask you. Can I get MLE of the μ 's?
- No problem!
- MLE of Gaussian is sample mean! so I just need to compute K sample means
- Oh there's one thing.
- None of the data are labeled. I don't know which Gaussian each point is from...
- Oh oh. You cannot do the MLE...



CSC872: PAMI – Kazunori Okada (C) 2025

16

16

Problem: More Formally

- Latent variable Z is **not observable**
- Data is unlabeled >>> **Unsupervised Learning**
- Data is labeled then >>> **Supervised Learning (Later)**

- Incomplete data likelihood $f(\theta) = \log P(x_1, \dots, x_N | \theta)$ $z_1 \sim z_N$
- Unknown z_1, \dots, z_N and θ $= \sum_n \log P(x_n | \theta)$
- Good old MLE recipe $= \sum_n \log \sum_k P(x_n, z_k | \theta)$ complete data likelihood

$$\frac{\partial f(\theta)}{\partial \theta} = 0 \quad \frac{\partial f(\theta)}{\partial z} = 0$$

$$= \sum_n \log \sum_k P(z_k | \pi) N(x_n | z_k, \theta)$$

- Well, it yields non-linear eq. **You cannot solve them**

CSC872: PAMI – Kazunori Okada (C) 2025

17

17

Expectation Maximization Algorithm

- Goal:** solve MLE problem of $f(\theta)$ *iteratively*
- Basic Idea
 - MLE of incomplete likelihood is difficult due to the unknown labels Z so...
- First find best label Z guess for each data point X**
 - E-step:** expectation
 - Expected value of label Z is computed, solving this problem probabilistically
- With the guessed Z, you can find MLE of θ with X**
 - M-step:** maximization
 - Optimize complete likelihood instead of incomplete likelihood
- Iterate these two steps (EM-algorithm) *probabilistically* mean shift
- You can prove that this converges to a nearest local optimum/peak!
(without step-size tuning like Gradient D) η

CSC872: PAMI – Kazunori Okada (C) 2025

18

18

Cost function for EM

- Given observed X and unobserved Z
- Complete log data likelihood

$$\begin{aligned} \log P(X, Z | \theta) &= \log \prod_n P(z_n | \pi) P(x_n | z_n, \theta) \\ &= \sum_n \log P(z_n | \pi) N(x_n | z_n, \theta) \end{aligned}$$

- Auxiliary function $Q(\theta | \theta')$ that you are going to optimize

$$Q(\theta | \theta') = E_{Z|X, \theta'} [\log P(X, Z | \theta)]$$

Marginalizing!

PS: EM Algorithm for GMM

- Given $X = x_1, \dots, x_n, \dots, x_N$, $\theta = (\{\pi_k\}, \{\mu_k\}, \{\Sigma_k\})$, K
- E-step: calculate $P(z_n | x_n, \theta_{old})$ for each x_n

$$P(z_n = k | x_n, \theta) = \frac{N(x_n | \mu_k, \Sigma_k) \pi_k}{\sum_{k'} N(x_n | \mu_{k'}, \Sigma_{k'}) \pi_{k'}} = r_{nk}$$

Sum = 1

- M-step: replace current θ_{old} by solving

$$\theta_{new} \leftarrow \operatorname{argmax}_{\theta'} Q(\theta' | \theta_{old})$$

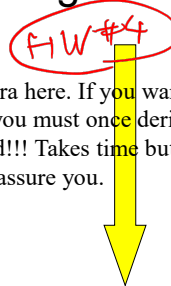
$$\pi_k^{new} = \frac{1}{N} \sum_n r_{nk}$$

$$\mu_k^{new} = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

$$\Sigma_k^{new} = \frac{\sum_n r_{nk} (x_n - \mu_k^{new})(x_n - \mu_k^{new})^T}{\sum_n r_{nk}}$$

Rather Simple !

Lots of algebra here. If you want to master EM, you must once derive these by hand!!! Takes time but rewarding, I assure you.



Relation to Clustering

DS Data Science
DM Data Mining
Data Engineering
↓ Unsupervised Learning

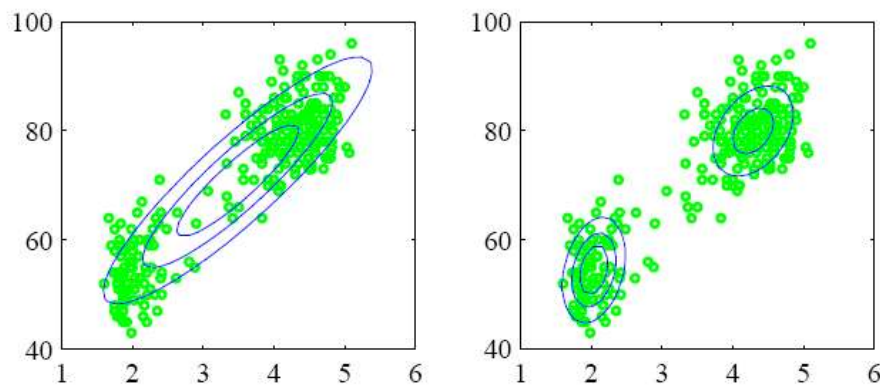
- It is a parametric clustering!!!
- Bayesian way of doing clustering
 - After running GMM-EM on your data, you have a Gaussian prob. distribution for each cluster
 - Any given point, you can do MAP
 - Which cluster label k maximizes $P(z_n=k|x_n, \theta)$?
- Let's compare this with mean shift
- Better?: yes, faster and more economical
- Worse?: yes, well both can get stuck at local optima but really you need to know K beforehand for EM

CSC872: PAMI – Kazunori Okada (C) 2025

21

21

Number of Clusters...



CSC872: PAMI – Kazunori Okada (C) 2025

22

22

PS: K-Mean Clustering

- What about other clustering methods I learned in your other favorite courses???
- **Key Observation:** Notice that EM is not really computing the value for Z at each step. E-step: calculate $P(z_n|x_n, \theta_{old})$
- It only computes a probability distribution of Z
- This is called **Soft-Assignment**
- Replace it with **Hard-Assignment?**
- At each iteration, you do the **MAP** and get Z value *Sample mean*
- Then use regular **MLE** formulae to get parameters *Sample variance*
- **It is the K-Mean Clustering!!! (with $\Sigma = \sigma^2 I, \sigma \rightarrow 0$)**
- **Moral: GMM-EM is a generalization of K-Mean**

CSC872: PAMI – Kazunori Okada (C) 2025

23

23

Curse of Dimensionality

- **Here is A MAJOR problem!!!**
- Linear increase in dimension of domain yields
- **Exponential increase in volume**
- Examples
 - Joint distribution of random variables with 10 attributes $\mathbf{x}=(x_1, \dots, x_{10})$
 - 2 variables $P(x,y)$: $10^2 = 100$ possible combinations
 - 10 variables $P(x,y,z,a,\dots,g)$: $10^{10} = 10$ billion possible combinations!!!
- Examples
 - 1D time series analysis to 3D RGB color analysis
 - Bioinformatics: a study on a few markers to a full set of genes

CSC872: PAMI – Kazunori Okada (C) 2025

24

24

The deal is that...

- **Because of this, solving a real-world problem is difficult!!!**
 - **Cannot model** joint distribution in high-dimension
 - **Cannot sample** sufficiently from a joint distribution
 - Takes **a lot of time** to compute
 - Takes **a lot of space** to keep in memory
- We learned that **Parametric Modeling** makes things better
- But you may **not have enough number of samples** to accurately estimate the covariance
- Oh NO... **What we gonna do???**

CSC872: PAMI – Kazunori Okada (C) 2025

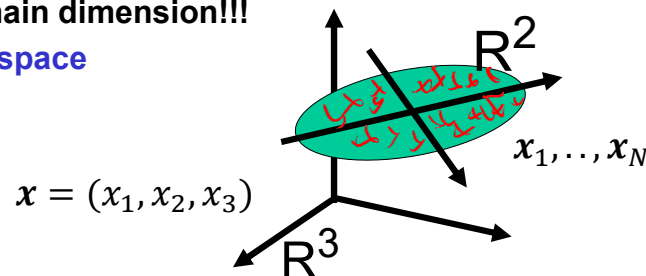
25

25

Dimensionality Reduction



- **A SERIOUSLY GOOD NEWS OF TODAY**
 - **You can describe entire information of your data by using much smaller number of variables**
 - **Intrinsic Dimensionality**
 - **Variance of your data is (most of time) confined to a space whose dimensionality is lower than your domain dimension!!!**
 - **Subspace**



CSC872: PAMI – Kazunori Okada (C) 2025

26

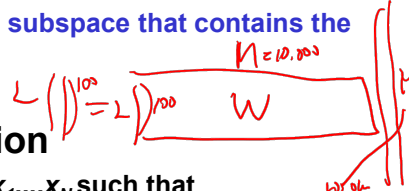
26

Dimensionality Reduction



- **Example**

- Consider a 100x100 pixel image
- We take it as 10,000 dimensional vector $\in R^{10000}$
- Take 100 images
- We have 100 points in R^{10000}
- What is the dimensionality of the subspace that contains the 100 points?
- Typically $\sim 100 \lll 10,000$!!!



- **Linear Dimension Reduction**

- Find a linear transform W given x_1, \dots, x_N such that

$$y = Wx$$

$100 = L < M = 10,000$

$$x \in R^M = 10,000$$

$$y \in R^L = 100$$

$$W \in R^{L \times M} = 100 \times 10,000$$

CSC872: PAMI – Kazunori Okada (C) 2025

27

27

PS: Principal Component Analysis

- **PCA:** Find a linear subspace with a lower-dimension that captures most data variance of N d -variate samples

- **Steps (Pseudo Code)**

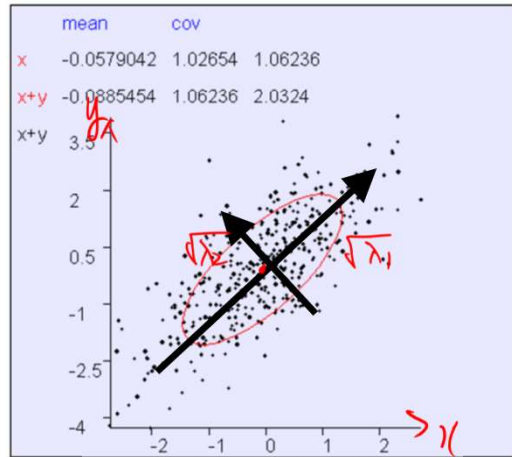
- 1) Compute $d \times d$ **sample covariance** matrix
- 2) Solve its **eigen-value problem**
- 3) Resulting a set of eigen-values and corresponding eigen-vectors
- 4) **Sort the eigen-vectors** according to the eigen-values
- 5) **Choose top-K** eigen-vectors with highest eigen-values
- 6) Set each row of W by the K eigen-vectors (called principal components)

CSC872: PAMI – Kazunori Okada (C) 2025

28

28

Principal Components in 2D



Eigenvalues = σ^2

PCs are perpendicular to each other

PCs are independent to each other

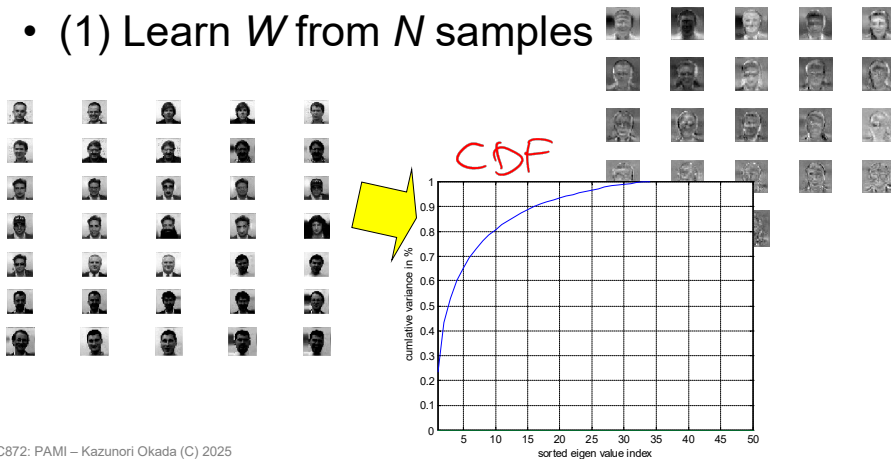
CSC872: PAMI – Kazunori Okada (C) 2025

29

29

Eigenface: Learning

- Use PCA to extract economic feature for describing facial images
- (1) Learn W from N samples



CSC872: PAMI – Kazunori Okada (C) 2025

30

30

PCA: Principal Component Anal.

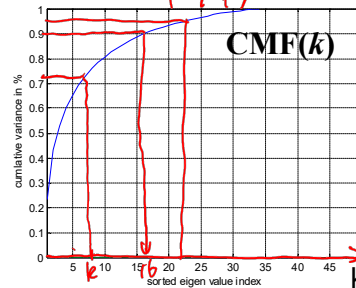
- Steps

- 1) Collect Training Images (must be aligned)
- 2) Vectorize the Images: $X = \{x_1, \dots, x_N\}$
- 3) Construct Covariance Matrix: $C = XX^T$
- 4) Solve Eigenvalue Problem: $Cv_i = \lambda_i v_i$
- 5) Select Top Eigenvectors $W = \{v_1, \dots, v_K\}^T$

$$\lambda_1 > \lambda_2 > \dots > \lambda_k > \dots > \lambda_M$$

$$CMF(k) = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^M \lambda_i}$$

k=1 k=2 k=3
 λ_1 $\lambda_1 + \lambda_2$ $\lambda_1 + \lambda_2 + \lambda_3$ $\lambda_1 + \dots + \lambda_M$



CSC872: PAMI – Kazunori Okada (C) 2025

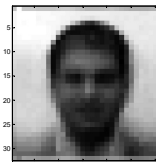
31

Eigenface: Feature Extraction

- Use PCA to extract economic feature for describing facial images
- (2) Extract a feature f of a face x using W

$$f = W(x - \mu)$$

$$\mu = \frac{1}{N} \sum_{n=1}^N x_n$$



CSC872: PAMI – Kazunori Okada (C) 2025

32

32

Face Recognition with Eigen Face

- Preprocess:
 - Prepare DB of N known face by using W and μ
- Recognition: Given an input face x , which entry of the known person DB is closest to the input?
- Steps (Pseudo Code)
 - 1) Extract the feature of input using W and μ
 - 2) Compute the Euclidean distance of the input feature to the DB features.
 - 3) Find the DB entry with the smallest Euclidean distance to the input
 - 4) Output the index of the best match entry
 - 5) Show the image of the best match entry



CSC872: PAMI – Kazunori Okada (C) 2025

33

33

Summary

- Gaussian Mixture Model
 - Gaussian Mixture Model
 - EM Algorithm
 - K-Mean Clustering
 - Curse of Dimensionality
 - Dimensionality Reduction
 - Principal Component Analysis
 - Eigenface and Face Recognition
- Next
 - General Framework for Machine Learning
 - Pattern Classification Examples

CSC872: PAMI – Kazunori Okada (C) 2025

34

34