# EE5907: Pattern Recognition

(Second Half)

# &

# EE5026: Machine Learning for Data Analytics

## Dr. WANG Si

Email: si.wang@nus.edu.sg

Office: E1a-04-01

# Outlines

- Pattern Representation Learning
  - Unsupervised Representation Learning (week 7)
  - Supervised Representation Learning (week 8)
  - Unified Framework: Graph Embedding (week 8)
- Patter Recognition Models
  - Clustering and Applications (week 9)
  - Gaussian Mixture Model and Boosting (week 10)
  - Support Vector Machines (week 11)
- Deep Learning (week 12)
- Revision and Q&A (week 13)

# Continuous Assessments

(20%)                    (40%)
❑ EE5907 CA2 & EE5026 CA1: 9 Oct. 2025 - 14 Nov. 2025 1800 SGT (Thu. Week 8 - Fri. Week 13)

Policy on late submission:

- Up to 12 hours late: 80% of your earned points

- Up to 1 day late: 50% of your earned points

- Up to 2 days late: 30% of your earned points
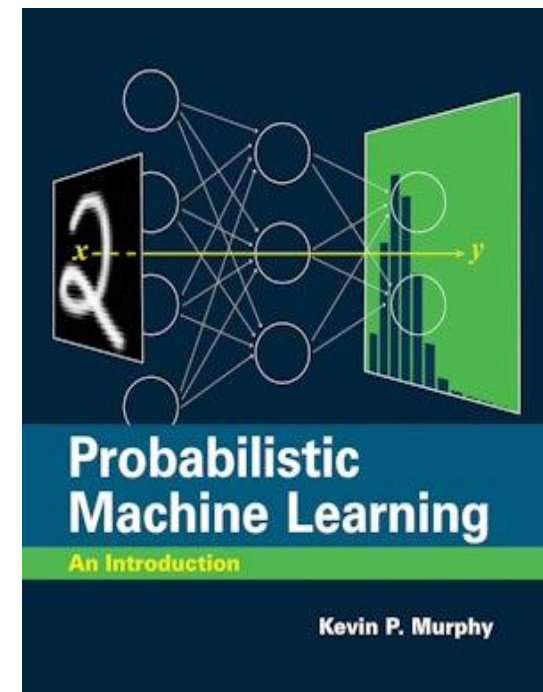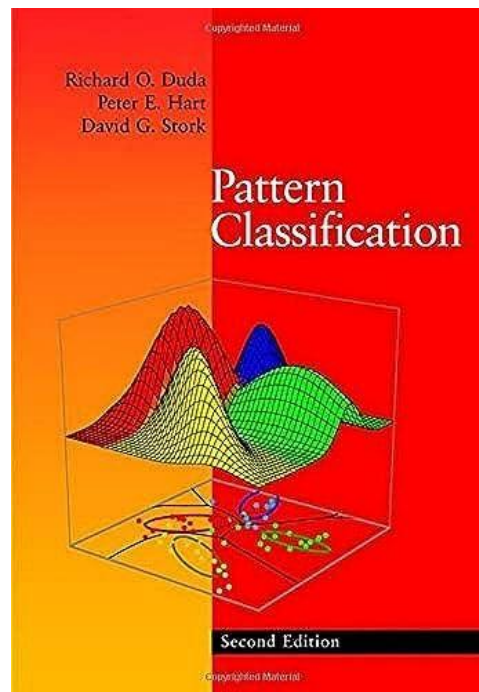
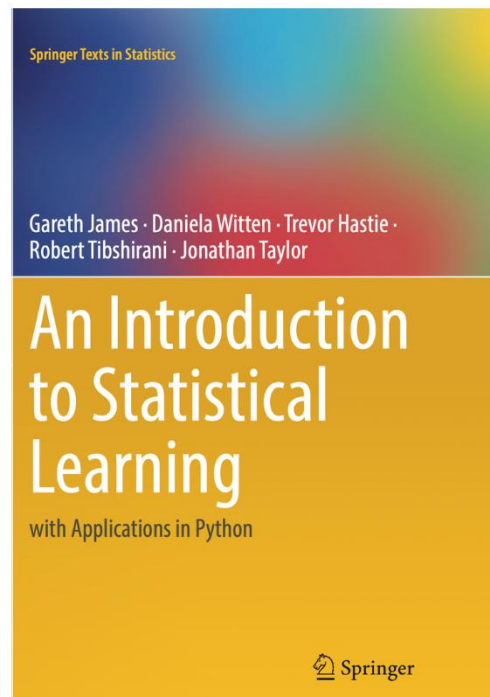- More than 2 days late: No marks (0%)

❑ EE5907 & EE5026 Final Exam (60%): Closed book with one-page double-sided A4-sized cheat sheet

# Textbooks and References

## (no fixed textbook)

- Books
  - Gareth, J., et al. "**An Introduction to Statistical Learning: with Applications in Python**." (2023)
  - Richard O. Duda, et al. "**Pattern Classification.**" (2001, Wiley)
  - Murphy, Kevin P., "**Probabilistic Machine Learning: An Introduction.**" (2022, MIT press)

# Pattern Representation Learning

➢ We will discuss:

  ➢ What is pattern representation learning

  ➢ Principal Component Analysis (PCA)

  ➢ Non-negative Matrix Factorization (NMF)

➢ At the end of this lecture, you should be able to:

  ➢ Know what pattern representation learning is

  ➢ Understand rationales behind PCA and NMF

  ➢ Perform PCA and NMF

# What are we doing with Pattern Recognition?

Raw Data

Data Preprocessing

Extract Informative Pattern Representation
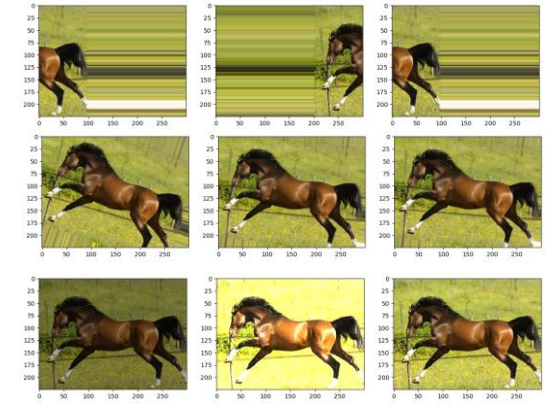
PR Model (after training)

Recognition Results



Grayscale | Normalization | Augmentation

|  | Face | Neck | Legs |
|---|---|---|---|
| Horse | Long | Long | Long |
| Cat | Round | Short | Short |

**Unsupervised**
- K-means
- Hierarchical clustering
- Gaussian Mixture Model

**Supervised**
- Boosting
- Support Vector Machine

Deep Neural Network

Cluster assignment          Classification / Regression

6

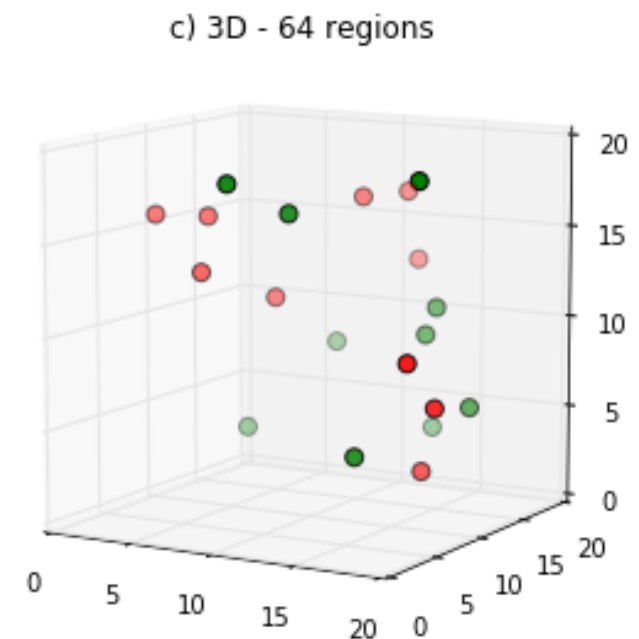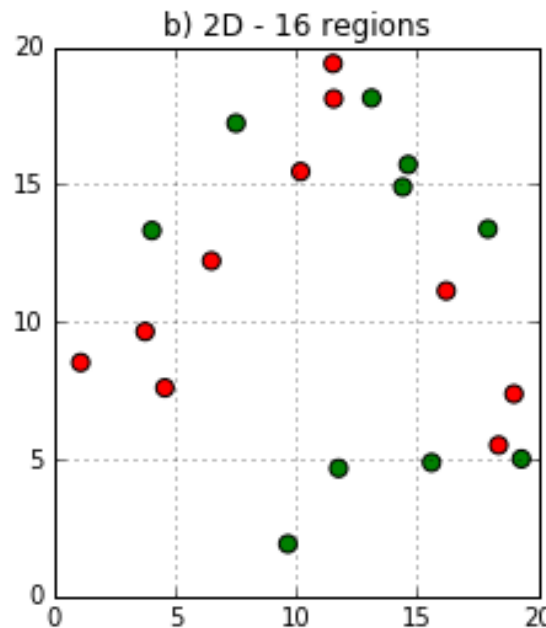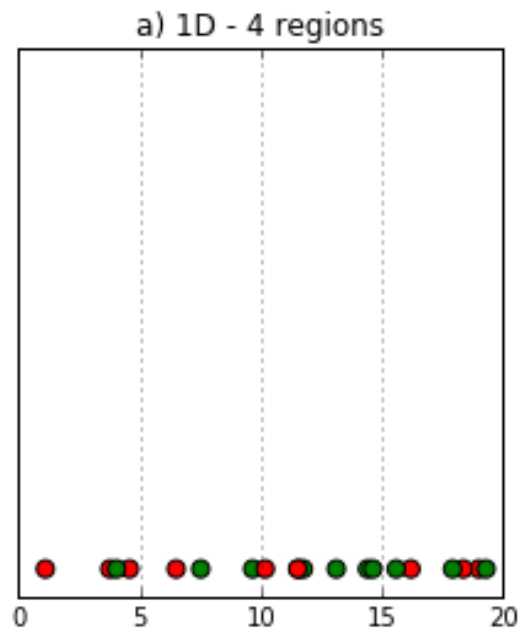# Representation Learning / Extraction

- Representation learning refers to <span style="color:red">transforming</span> the raw/preprocessed data into another (possibly lower-dimensional) space.

- Such that, in the new space, one can perform pattern recognition <span style="color:red">more easily</span>.

- Criterion for good representation in different problem settings:
  - **Unsupervised**: minimize information loss (no class information)
  - **Supervised**: maximize discrimination (with class information)

|  | Face | Neck | Legs |
|------|------|------|------|
| **Horse** | Long | Long | Long |
| **Cat** | Round | Short | Short |

# Why Feature Extraction?

- Many traditional pattern recognition techniques may not be effective for high-dimensional data
    - Curse of dimensionality



a) 1D - 4 regions
b) 2D - 16 regions
c) 3D - 64 regions

# Why Feature Extraction?

- Patterns may have small <span style="color:red">intrinsic</span> dimension



Intrinsic dimensions → (50-100 dimensions)

400x400x3=480000 dimensions

# Why Feature Extraction?

- Visualization: projecting high-dimensional data onto 2D or 3D planes

- Data compression: efficient storage and retrieval

- Noise removal: positive effect on testing accuracy

# Applications of Feature Extraction

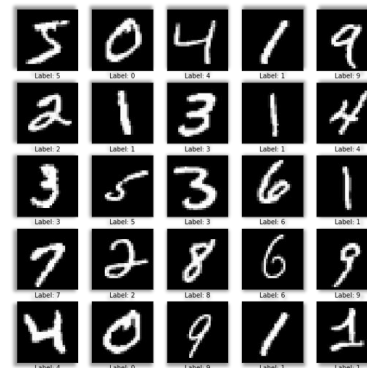- Face recognition

- Handwritten digit recognition

- Text mining

- Image retrieval

- Protein classification

Face Images

Text

Proteins

Handwritten digit

Image database

# Unsupervised Feature Learning
# **Principal Component Analysis**

# Principal Component Analysis (PCA)

- Probably the most widely-used and well-known multivariate analysis method.

- Introduced by Pearson (1901)

- First applied in ecology by Goodall (1954) under the name "factor analysis".



Karl Pearson

Pearson, K. (1901). "On Lines and Planes of Closest Fit to Systems of Points in Space" Philosophical Magazine **2** (11): 559–572.

# What is Principal Component Analysis?

- Principal component analysis (PCA)

  – Reduce the dimensionality of a collection of observations by projecting onto a new smaller set of variables, called principal components (PCs), while preserving as much information as possible.

- What is PC?

  – Capture the big (principal) variability in the data

  – Ordered by captured variations from largest to smallest (e.g., 1st PC captures the largest variability)

  – PCs are uncorrelated (orthogonal to each other)

# Geometric Picture of Principal Components

# Geometric Picture of Principal Components

A

B

Which one is the 1st PC?

# Geometric Picture of Principal Components



**How to find this line?**

- The 1$^{st}$ PC $Z_1$ is parallel to <u>the line</u> minimum-distance fitted to the data points

- The 2$^{nd}$ PC $Z_2$ is parallel to <u>the line</u> minimum-distance fitted to the data points after $z_1$ has been taken into account

# Geometric Picture of Principal Components

https://people.duke.edu/~hpgavin/SystemID/CourseNotes/TotalLeastSquares.pdf

# Algebraic Definition of PCs

Given a sample set of $n$ observations on a vector of $d$ variables

$$\{x_1, x_2, \ldots, x_n\} \subset \mathbb{R}^d$$

Define the first principal component ($PC_1$) to be the normalized linear projection vector $\mathbf{z}_1 = [z_{1,1}, \ldots, z_{d,1}]^T$, such that

$$y_1 = \mathbf{z}_1^T \mathbf{x}$$

$var[y_1]$ is maximum, where $y_1$ is the projection along 1st PC axis

# Algebraic Derivation of PCs

To find $\mathbf{z}_1$ that maximizes $\text{var}[y_1]$, subject to $\mathbf{z}_1^T\mathbf{z}_1 = 1$

$$\underset{\mathbf{z}_1}{\arg\max} \, \text{var}[y_1] \, , s.t., \mathbf{z}_1^T\mathbf{z}_1 = 1 \qquad \bar{\mathbf{x}} = \frac{1}{n}\sum_{i=1}^{n}\mathbf{x}_i$$

$$\text{var}[y_1] = \mathbb{E}[(y_1 - \overline{y_1})^2] = \frac{1}{n}\sum_{i=1}^{n}(\mathbf{z}_1^T\mathbf{x}_i - \mathbf{z}_1^T\bar{\mathbf{x}})^2 = \mathbf{z}_1^T(\frac{1}{n}\sum_{i=1}^{n}(\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T)\mathbf{z}_1$$

Covariance matrix **S**

$$\begin{bmatrix} \text{Cov}(X_1, X_1) & \cdots & \text{Cov}(X_1, X_d) \\ \vdots & \ddots & \vdots \\ \text{Cov}(X_d, X_1) & \cdots & \text{Cov}(X_d, X_d) \end{bmatrix}$$

$$f = \mathbf{z}_1^T\mathbf{S}\mathbf{z}_1, \, g = \mathbf{z}_1^T\mathbf{z}_1 - 1$$

Let $\lambda$ be a Lagrange multiplier,

$\text{Cov}(A,B) = \mathbb{E}[(A - \mathbb{E}[A])(B - \mathbb{E}(B))]$

$$\nabla_{\mathbf{z}_1}(\mathbf{z}_1^T\mathbf{S}\mathbf{z}_1) = \lambda\nabla_{\mathbf{z}_1}(\mathbf{z}_1^T\mathbf{z}_1 - 1)$$

$$\mathbf{S}\mathbf{z}_1 = \lambda\mathbf{z}_1$$

Therefore, $\mathbf{z}_1$ is an eigenvector of **S,**
corresponding to the <u>largest</u> eigenvalue $\lambda = \lambda_1$

Why?

Lagrange Multipliers:
Given *f* to be optimized, subject to the constraint *g*=0, then, solve: $\nabla f = \lambda\nabla g$

Eigenvector:
direction unchanged after applying a linear transformation **S**

# Algebraic Derivation of PCs

$$\underset{\mathbf{z}_1}{\operatorname{argmax}} \operatorname{var}[y_1], s.t., \mathbf{z}_1^T \mathbf{z}_1 = 1$$

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i$$

$$\operatorname{var}[y_1] = \mathbb{E}[(y_1 - \overline{y_1})^2] = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{z}_1^T \mathbf{x}_i - \mathbf{z}_1^T \bar{\mathbf{x}})^2 = \mathbf{z}_1^T (\underbrace{\frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T}) \mathbf{z}_1$$

Covariance matrix **S**

$f = \mathbf{z}_1^T \mathbf{S} \mathbf{z}_1, g = \mathbf{z}_1^T \mathbf{z}_1 - 1$

Let $\lambda$ be a Lagrange multiplier,

$$\nabla_{\mathbf{z}_1}(\mathbf{z}_1^T \mathbf{S} \mathbf{z}_1) = \lambda \nabla_{\mathbf{z}_1}(\mathbf{z}_1^T \mathbf{z}_1 - 1)$$

$$\mathbf{S} \mathbf{z}_1 = \lambda \mathbf{z}_1$$

Therefore, $\mathbf{z}_1$ is an eigenvector of **S**, corresponding to the <u>largest</u> eigenvalue $\lambda = \lambda_1$

$$\operatorname{var}[y_1] = \mathbf{z}_1^T \mathbf{S} \mathbf{z}_1 = \mathbf{z}_1^T \lambda \mathbf{z}_1 = \lambda$$

Therefore, $\lambda$ is largest eigenvalue in order to maximize $\operatorname{var}[y_1]$

Lagrange Multipliers:
Given *f* to be optimized, subject to the constraint *g*=0, then, solve: $\nabla f = \lambda \nabla g$

Eigenvector:
direction unchanged after applying a linear transformation *S*

# Algebraic Derivation of PCs

Similarly, $\mathbf{z}_2$ is also an eigenvector of $\mathbf{S}$,

whose eigenvalue $\lambda = \lambda_2$ is the second largest

In general,

$$\text{var}[y_k] = \mathbf{z}_k^T \mathbf{S} \mathbf{z}_k = \lambda_k$$

Are obtained PCs orthogonal?

- The $k$th largest eigenvalue of $\mathbf{S}$ is the variance of the $k$th PC projections.

- The $k$th PC $\mathbf{z}_k$ captures the $k$th greatest variation in the samples

# Algebraic Derivation of PCs

Proof: Obtained PCs are orthogonal to each other

$$Sz_1 = \lambda_1 z_1, \qquad Sz_2 = \lambda_2 z_2, \ldots \ldots$$

$$z_2^T S z_1 = \lambda_1 z_2^T z_1 \qquad z_1^T S z_2 = \lambda_2 z_1^T z_2$$

$$(z_1^T S z_2)^T = \lambda_2 (z_1^T z_2)^T$$

$$z_2^T S^T z_1 = \lambda_2 z_2^T z_1$$

$$\Downarrow S^T = S$$

$$z_2^T S z_1 = \lambda_2 z_2^T z_1$$

$$\lambda_1 z_2^T z_1 = \lambda_2 z_2^T z_1 \quad \Rightarrow \quad z_2^T z_1 = 0$$

# Compute PCs

- Main steps for computing PCs and Projections on PCs

  - Calculate the covariance matrix $\mathbf{S}$.

  - Compute its eigenvectors: $\mathbf{z}_i$

  $$(\mathbf{S} - \lambda \mathbf{I}_d)\mathbf{z}_i = 0 \quad \textcircled{1}$$

  Solve $|\mathbf{S} - \lambda \mathbf{I}_d| = 0$, Get $\lambda$

  Substitute each $\lambda$ to $\textcircled{1}$, Get $\mathbf{z}_i$

  - The first $p$ eigenvectors $\{\mathbf{z}_i\}_{i=1}^{p}$ form the $p$ PCs

  - The transformation matrix $\mathbf{Z}$ consists of the $p$ PCs:

  $$\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_p],$$

  $$\mathbf{y} = \mathbf{Z}^T \mathbf{x}$$

$$A \dot{=} [c_1, c_2 \cdots c_d]$$

$$z_1 c_1 + z_2 c_2 + \cdots z_d c_d = 0$$

24

# Examples: Compute PCs

Given the following dataset, compute the first PC and its corresponding projection for each data sample.

| Observations | Feature 1 | Feature 2 |
|---|---|---|
| 1 | 1 | 2 |
| 2 | 3 | 4 |
| 3 | 5 | 6 |
| 4 | 7 | 8 |

# Examples: Compute PCs

## Step 1: Compute Covariance Matrix $\mathbf{S}$

$$\mathbf{S} = \frac{1}{n}\sum_{i=1}^{n}(\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

$$\bar{\mathbf{x}} = \frac{1}{n}\sum_{i=1}^{n}\mathbf{x}_i =$$

$$\frac{1}{4}\left(\begin{bmatrix}1\\2\end{bmatrix}+\begin{bmatrix}3\\4\end{bmatrix}+\begin{bmatrix}5\\6\end{bmatrix}+\begin{bmatrix}7\\8\end{bmatrix}\right) = \begin{bmatrix}4\\5\end{bmatrix}$$

n=4, d=2

| Observations | Feature 1 | Feature 2 |
|---|---|---|
| 1 | 1 | 2 |
| 2 | 3 | 4 |
| 3 | 5 | 6 |
| 4 | 7 | 8 |

$$= \frac{1}{4}\left(\left(\begin{bmatrix}1\\2\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)\left(\begin{bmatrix}1\\2\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)^T + \left(\begin{bmatrix}3\\4\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)\left(\begin{bmatrix}3\\4\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)^T + \right.$$

$$\left.\left(\begin{bmatrix}5\\6\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)\left(\begin{bmatrix}5\\6\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)^T + \left(\begin{bmatrix}7\\8\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)\left(\begin{bmatrix}7\\8\end{bmatrix} - \begin{bmatrix}4\\5\end{bmatrix}\right)^T\right) = \begin{bmatrix}5 & 5\\5 & 5\end{bmatrix}$$

## Step 2: Compute eigenvectors $\mathbf{z}_i$
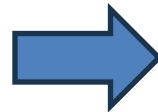
$$\det(\mathbf{S} - \lambda\mathbf{I}_2) = 0$$

$$\det\left(\begin{bmatrix}5 & 5\\5 & 5\end{bmatrix} - \lambda\begin{bmatrix}1 & 0\\0 & 1\end{bmatrix}\right) = 0$$

$$\det\left(\begin{bmatrix}5-\lambda & 5\\5 & 5-\lambda\end{bmatrix}\right) = 0$$

$$(5-\lambda)(5-\lambda) - 5\times 5 = 0$$

$$\lambda_1 = 10, \lambda_2 = 0$$

For $\lambda_1 = 10$, $(\mathbf{S} - \lambda\mathbf{I}_d)\mathbf{z}_1 = 0$

$$\begin{bmatrix}-5 & 5\\5 & -5\end{bmatrix}\begin{bmatrix}z_{1,1}\\z_{2,1}\end{bmatrix} = 0$$

$$z_{1,1} = z_{2,1}$$

Since $\mathbf{z}_1^T\mathbf{z}_1 = 1$, i.e., $z_{1,1}^2 + z_{2,1}^2 = 1$

$$z_{1,1} = z_{2,1} = \pm\frac{1}{\sqrt{2}}, \mathbf{z}_1 = \pm\begin{bmatrix}\frac{1}{\sqrt{2}}\\\frac{1}{\sqrt{2}}\end{bmatrix}, \text{var}[y_1] = \lambda_1 = 10$$

## Step 3: Compute Projection $y$ on $\mathbf{z}_1$

$$y_1 = \mathbf{z}_1^T\mathbf{x}_1 = \pm[\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}]\begin{bmatrix}1\\2\end{bmatrix} = \pm\frac{3}{\sqrt{2}}, y_2 = \mathbf{z}_1^T\mathbf{x}_2, y_3 = \mathbf{z}_1^T\mathbf{x}_3, y_4 = \mathbf{z}_1^T\mathbf{x}_4$$

# Practical Computation of PCA

- In reality, we compute the PCs via Singular Value Decomposition (SVD) on the centered data matrix.

$$\mathbf{S} = \frac{1}{n}\sum_{i=1}^{n}(\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

$$= \frac{1}{n}\mathbf{X}_{d\times n}\mathbf{X}_{d\times n}^T$$

$$= \frac{1}{n}\mathbf{U}_{d\times d}\mathbf{D}_{d\times n}\mathbf{V}_{n\times n}^T(\mathbf{U}_{d\times d}\mathbf{D}_{d\times n}\mathbf{V}_{n\times n}^T)^T$$

$$= \frac{1}{n}\mathbf{U}_{d\times d}\mathbf{D}_{d\times n}\mathbf{D}_{d\times n}^T\mathbf{U}_{d\times d}^T$$

$V \cdot D^T \cdot U^T$

**Centered data matrix:**
$$X_{d,n} = [(x_1 - \bar{x}), \ldots, (x_n - \bar{x})]$$

**SVD of $\mathbf{X}_{d,n}$:**
$$\mathbf{X}_{d\times n} = \mathbf{U}_{d\times d}\mathbf{D}_{d\times n}\mathbf{V}_{n\times n}^T$$

$$\mathbf{U}_{d\times d}\mathbf{U}_{d\times d}^T = \mathbf{U}_{d\times d}^T\mathbf{U}_{d\times d} = \mathbf{I}_d,$$
$$\mathbf{V}_{n\times n}\mathbf{V}_{n\times n}^T = \mathbf{V}_{n\times n}^T\mathbf{V}_{n\times n} = \mathbf{I}_n,$$
$$\mathbf{D}_{d\times n} = \begin{bmatrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_d & 0 & \cdots & 0 \end{bmatrix}$$

Then, the eigenvectors $\mathbf{z}_i$ of $\mathbf{S}$ are the columns of $\mathbf{U}$ and the eigenvalues $\lambda_i$ are the diagonal elements of $\mathbf{DD}^T/n$

Why?

# Practical Computation of PCA

Suppose, $\mathbf{U}_{d\times d} = [\mathbf{u}_1, \ldots, \mathbf{u}_d], \mathbf{D}_{d\times n} = \begin{bmatrix} \sigma_1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_d & 0 & \cdots & 0 \end{bmatrix}, \mathbf{DD}^T/n = \begin{bmatrix} \frac{\sigma_1^2}{n} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{\sigma_d^2}{n} \end{bmatrix}$

**Prove**: the eigenvectors $\mathbf{z}_i$ of $\mathbf{S}$ are the columns of $\mathbf{U}$ and the eigenvalues $\lambda_i$ are the diagonal elements of $\mathbf{DD}^T/n$

**Prove**: $\mathbf{S}\mathbf{u}_i = \dfrac{\sigma_i^2}{n}\mathbf{u}_i$

$\mathbf{S}\mathbf{u}_i = (\dfrac{1}{n}\mathbf{U}_{d\times d}\mathbf{D}_{d\times n}\mathbf{D}_{d\times n}^T\mathbf{U}_{d\times d}^T)\mathbf{u}_i$

$$= [\mathbf{u}_1, \ldots, \mathbf{u}_d]\begin{bmatrix} \frac{\sigma_1^2}{n} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{\sigma_d^2}{n} \end{bmatrix}\begin{bmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_d^T \end{bmatrix}\mathbf{u}_i$$

*i*-th entry

$$= [\mathbf{u}_1, \ldots, \mathbf{u}_d]\begin{bmatrix} \frac{\sigma_1^2}{n} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{\sigma_d^2}{n} \end{bmatrix}\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = [\mathbf{u}_1, \ldots, \mathbf{u}_d]\begin{bmatrix} 0 \\ \vdots \\ \frac{\sigma_i^2}{n} \\ \vdots \\ 0 \end{bmatrix} = \frac{\sigma_i^2}{n}\mathbf{u}_i$$

$U_{d,d}U_{d,d}^T = I_d$

$U_{d,d}$ is an orthonormal matrix

$u_i^T u_i = 1, \quad u_j^T u_i = 0$

28

# Practical Computation of PCA

- The new reconstructed sample using $p$ PCs in original space is:

$$\widehat{\mathbf{x}_i} = \overline{\mathbf{x}} + \mathbf{U}_{d \times p} \mathbf{U}_{d \times p}^T (\mathbf{x}_i - \overline{\mathbf{x}})$$

Step 1: project the original sample $\mathbf{x}_i$ onto the $p$ PCs

$$\mathbf{y} = \mathbf{U}_{d \times p}^T (\mathbf{x}_i - \overline{\mathbf{x}})$$

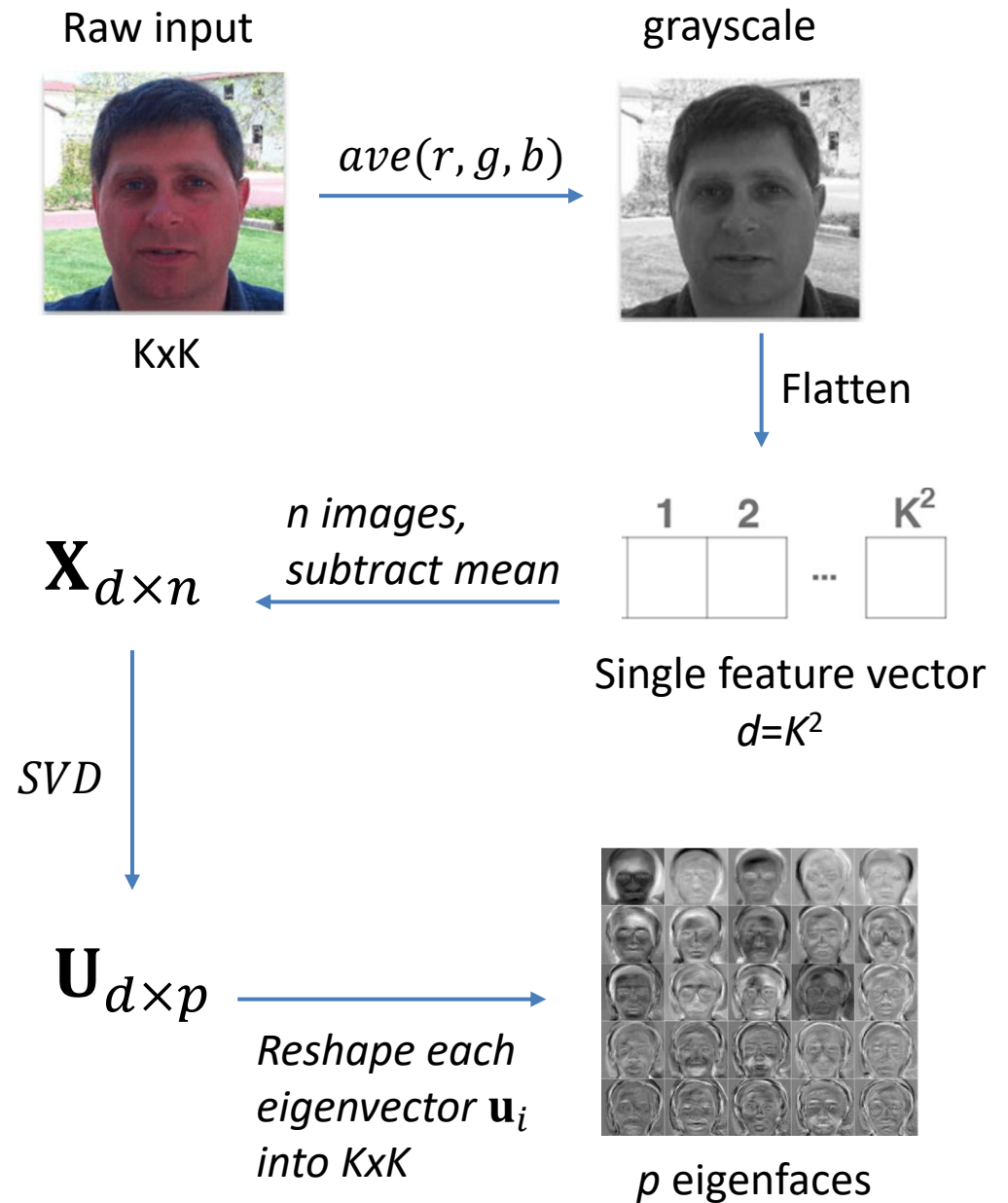Step 2: map $\mathbf{y}$ back to the original $d$-dimensional space

$$\widehat{\mathbf{x}_i^c} = \mathbf{U}_{d \times p} \mathbf{y} = \mathbf{U}_{d \times p} \mathbf{U}_{d \times p}^T (\mathbf{x}_i - \overline{\mathbf{x}})$$

$$\widehat{\mathbf{x}_i} = \widehat{\mathbf{x}_i^c} + \overline{\mathbf{x}} = \mathbf{U}_{d \times p} \mathbf{U}_{d \times p}^T (\mathbf{x}_i - \overline{\mathbf{x}}) + \overline{\mathbf{x}}$$

# Visualize PCs

Raw input



KxK

$$ave(r, g, b)$$

grayscale



Flatten



Single feature vector
$d = K^2$

$n$ images,
subtract mean

$$\mathbf{X}_{d \times n}$$

$SVD$

Eigenface images



$$\mathbf{U}_{d \times p}$$

Reshape each
eigenvector $\mathbf{u}_i$
into KxK

$p$ eigenfaces

# Reconstruction with PCs

Image Dimension:
100x100x3

Original | 250 Eigenfaces | 1000 Eigenfaces | 4000 Eigenfaces



$$\widehat{\mathbf{x}}_i = \overline{\mathbf{x}} + \mathbf{U}_{d \times p} \mathbf{U}_{d \times p}^T (\mathbf{x}_i - \overline{\mathbf{x}})$$

# PCA and Classification
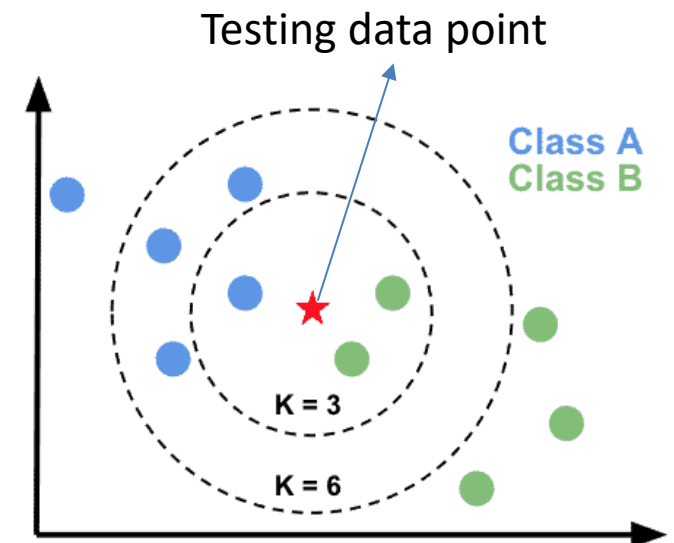
- Classification with PCA
  - Compute PCs using training dataset
  - Project both training and testing data into the obtained PCs space
  - For each testing data point, use K-Nearest-Neighbor (KNN) for classification
  - Issue: accuracy is sensitive to the number of PCs



How to determine the number of PCs?

Testing data point

Class A
Class B

K = 3

K = 6

To choose $p$ based on percentage of variation to retain, we can use the following criterion (smallest $p$):

$$\frac{\sum_1^p \lambda_i}{\sum_1^d \lambda_i} \geq Threshold\ (e.g., 0.95)$$

# PCA Remarks

- PCA
  - finds orthonormal basis for data
  - Sorts dimensions in order of "importance"
  - Discard low significance dimensions

- Uses:
  - Get compact description
  - Ignore noise
  - Improve classification (hopefully)

# PCA Limitations

- PCA is not always an optimal feature extraction for classification
  - PCA is based on the sample covariance, irrespective of class-membership
  - The projection axes chosen by PCA might not provide good discrimination power



PC direction

# PCA Limitations
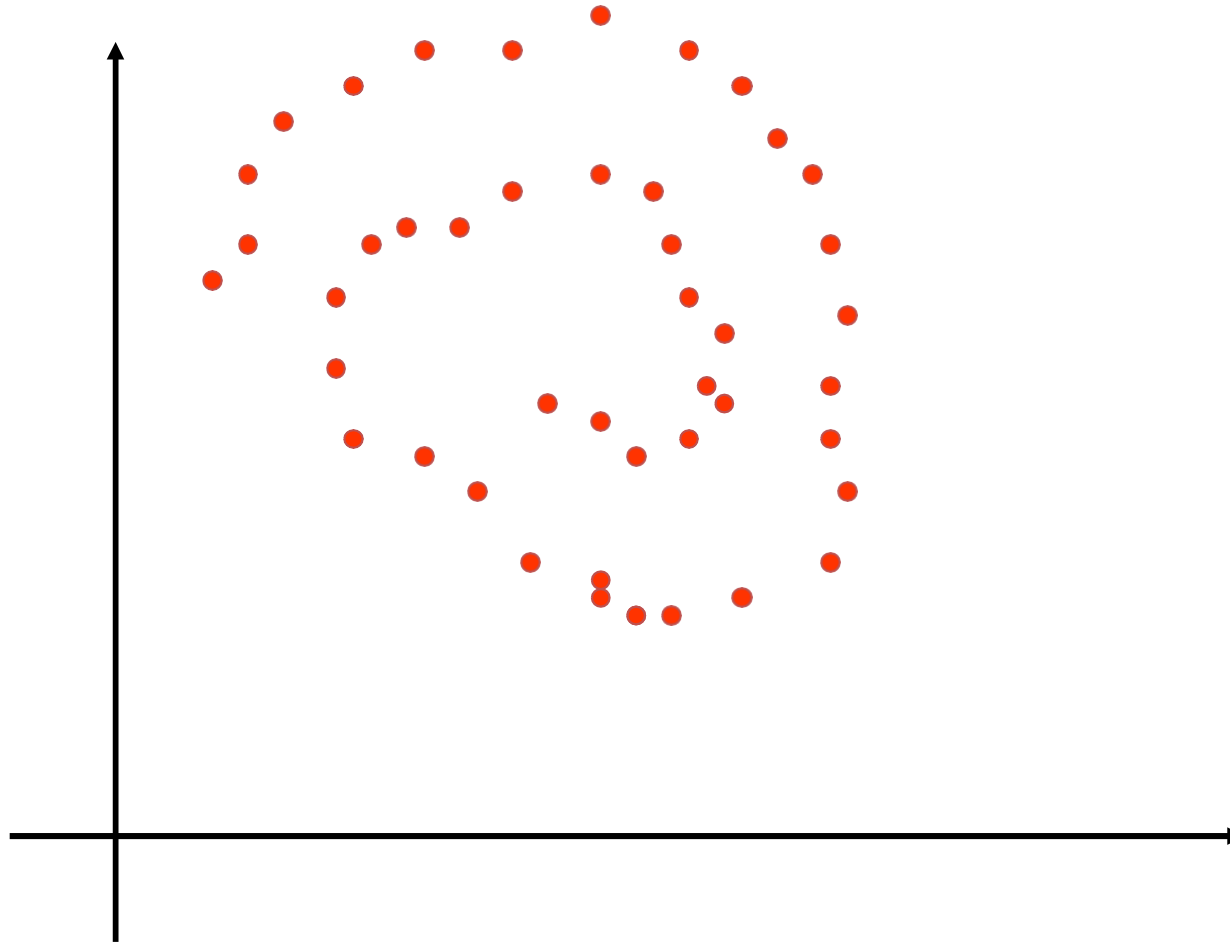
- PCA cannot capture NON-LINEAR structure



Note: Curvilinear Component Analysis can solve this case. Study this work if you are interested.

# PCA Limitations

- Lack of interpretability

  - Each PC projection is a linear combination of **all** input features

$$y_k = \mathbf{z}_k^T \mathbf{x}$$

  What does each PC actually represent?

  e.g., Face Recognition



  Not localized to a single feature such as "mouth", "eyebrows"

  e.g., Topic Modelling from Text

  A PC projection is a linear combination of occurrence of many words

$$PC = \begin{bmatrix} 0.4 \\ -0.3 \\ 0.2 \\ -0.1 \end{bmatrix}$$

  → "war"
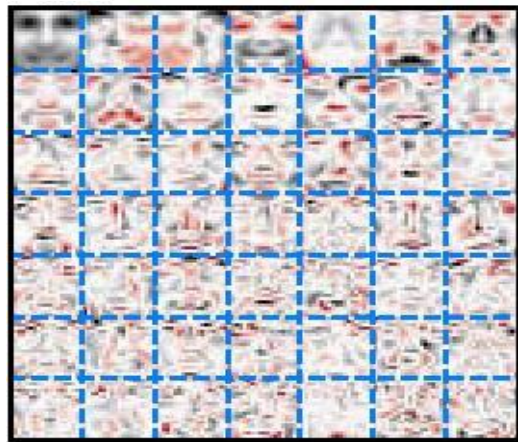  → "education"
  → "government"
  → "poem"

# PCA Limitations

- Lack of interpretability

$$\mathbf{x} = a_1\mathbf{z}_1 + \cdots + a_p\mathbf{z}_p$$

- Both PCs and coefficients can be positive or negative

What does "negative" contribution mean?

e.g., Face Recognition

Every vector can be expressed as the linear combination of basis vectors

$$\begin{bmatrix} 2 \\ 6 \end{bmatrix} = 2\begin{bmatrix} 1 \\ 0 \end{bmatrix} + 6\begin{bmatrix} 0 \\ 1 \end{bmatrix}$$
$$= 3\begin{bmatrix} 1 \\ 4 \end{bmatrix} + (-1)\begin{bmatrix} 1 \\ 6 \end{bmatrix}$$

Original

what does "subtracting an eigenface" mean?

PCs (Eigenfaces)   ×   Coefficients   =   Reconstructed

e.g., Topic Modelling from Text

$$PC = \begin{bmatrix} 0.4 \\ -0.3 \\ 0.2 \\ -0.1 \end{bmatrix}$$

→ "war"
→ "education"
→ "government"
→ "poem"

what does a **"negative occurrence" of a word** mean?

# Unsupervised Feature Learning

# Non-negative Matrix Factorization

(non-negative coefficients and basis vectors)

# NMF: Motivation

- Non-negative coefficients and basis vectors

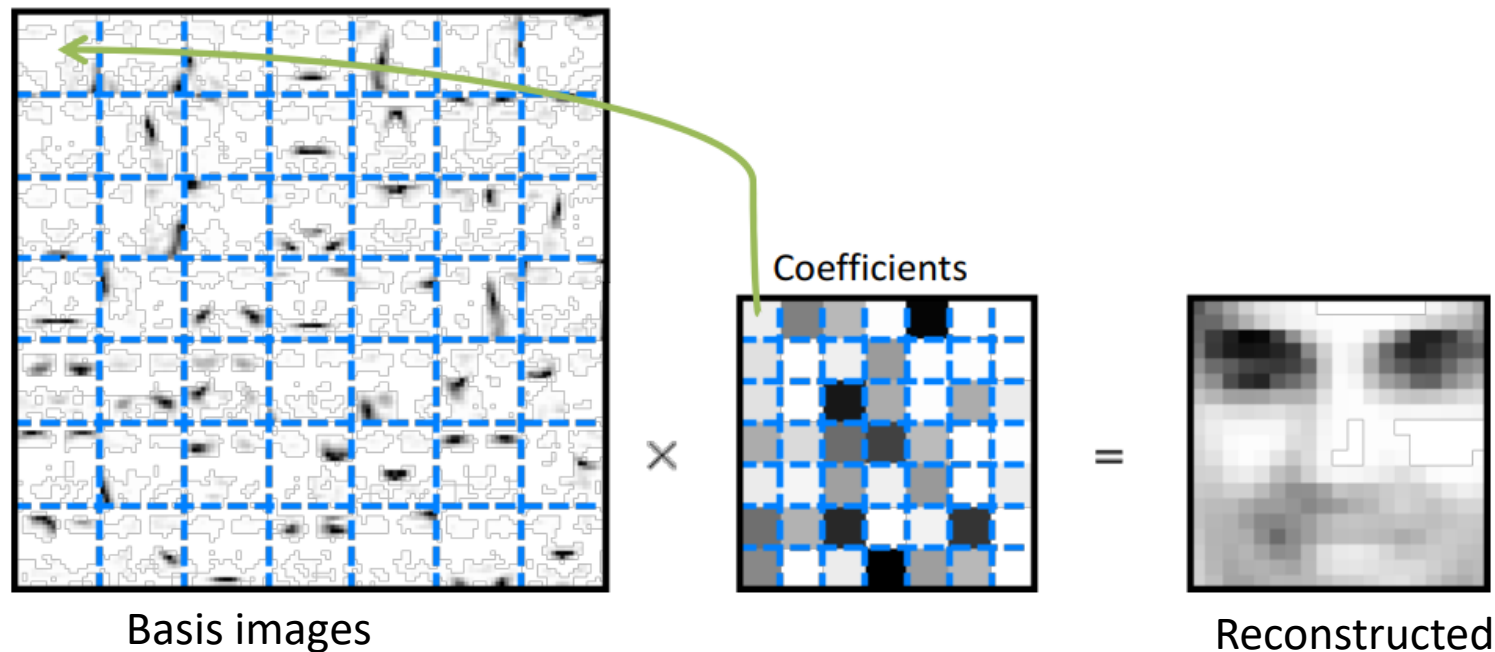  -More intuitive for non-negative datasets (e.g., images, word counts in documents, etc.)

  -Leads to basis vectors that represent parts



Basis images       ×       Coefficients       =       Reconstructed

# NMF: Mathematical Formulation

- Matrix factorization: $V \approx WH$    <span style="color:red">All entries in **V**, **W**, **H**$\geq 0$</span>

  - **V**: $(d \times n)$ Input data matrix, where $d$ is the number of input features, $n$ is the number of samples.   $\mathbf{V}_{d \times n} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$

  - **W**: $(d \times r)$ Basis matrix, where $r$ is the number of basis vectors.

$$\mathbf{W}_{d \times r} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_r]$$

  - **H**: $(r \times n)$ Coefficient matrix, each column of **H** is called encoding.

$$\mathbf{h}_1 \quad \mathbf{H}_{r \times n} = \begin{bmatrix} h_{1,1} & \cdots & h_{1,n} \\ \vdots & \ddots & \vdots \\ h_{r,1} & \cdots & h_{r,n} \end{bmatrix}$$

$$\mathbf{x}_i = h_{1,i}\mathbf{w}_1 + h_{2,i}\mathbf{w}_2 + \cdots + h_{r,i}\mathbf{w}_r = \mathbf{W}\mathbf{h}_i$$

**V**: Each column is a flattened vectorized image

**W**: Each column is a **non-negative** basis image

**H**: Each column is a **non-negative** coefficient vector



$\mathbf{x}_i$

$n=8$, $d=$KxK

Basis vector $\mathbf{w}_i$

# NMF: Objective Function

$$\min_{\mathbf{W},\mathbf{H}} \frac{1}{2} \boxed{\|\mathbf{V} - \mathbf{W}\mathbf{H}\|^2}, \text{ s.t. } \mathbf{W} \geq 0, \mathbf{H} \geq 0$$

$$\sum_{i,j} (\mathbf{V}_{i,j} - [\mathbf{W}\mathbf{H}]_{i,j})^2$$

**No closed form solution!**

## Gradient Descent!

**Recall** Goal: $\underset{\theta}{\text{argmin}} \, J(\theta)$

Initialize $\theta_0$ and learning rate $\eta$
**while** *true* **do**
    $\theta_{k+1} \leftarrow \theta_k - \eta \boxed{\nabla_\theta J(\theta_k)}$
    **if** *converge* **then**
        **return** $\theta_{k+1}$
    **end**
**end**

# NMF: Gradient Descent Solution

$$\min_{\mathbf{W}, \mathbf{H}} \frac{1}{2} \|\mathbf{V} - \mathbf{WH}\|^2, \text{ s.t. } \mathbf{W} \geq 0, \mathbf{H} \geq 0$$

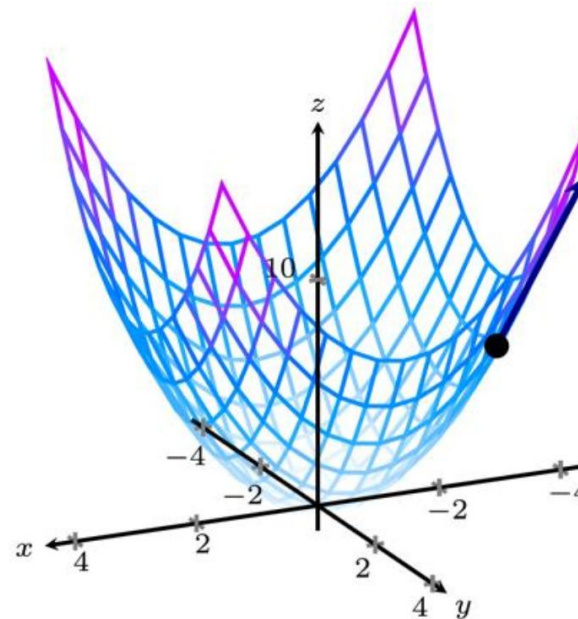No closed form solution!

**Gradient Descent!**

$$J(\mathbf{W}, \mathbf{H}) = \frac{1}{2}\|\mathbf{V} - \mathbf{WH}\|^2 = \frac{1}{2}\sum_{i,j}(\mathbf{V}_{i,j} - \sum_{l=1}^{r}\mathbf{W}_{i,l}\mathbf{H}_{l,j})^2$$

$$\nabla_{\mathbf{W}_{p,q}}J(\mathbf{W}, \mathbf{H}) = \sum_{i,j}(\mathbf{V}_{i,j} - \sum_{l=1}^{r}\mathbf{W}_{i,l}\mathbf{H}_{l,j})(-\frac{\partial}{\partial\mathbf{W}_{p,q}}\sum_{l=1}^{r}\mathbf{W}_{i,l}\mathbf{H}_{l,j})$$

Only when $i=p$, $l=q$, $\mathbf{W}_{i,l}$ depends on $\mathbf{W}_{p,q}$

$$\nabla_{\mathbf{x}}f = \frac{df(\mathbf{X})}{d\mathbf{X}}$$

$$= \begin{bmatrix} \frac{\partial f}{\partial x_{1,1}} & \cdots & \frac{\partial f}{\partial x_{1,K}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{d,1}} & \cdots & \frac{\partial f}{\partial x_{d,K}} \end{bmatrix}$$

$$= -\sum_{j}(\mathbf{V}_{p,j} - \sum_{l=1}^{r}\mathbf{W}_{p,l}\mathbf{H}_{l,j})\mathbf{H}_{q,j}$$

$$= \sum_{j}((\mathbf{WH})_{p,j} - \mathbf{V}_{p,j})\mathbf{H}_{q,j} = \sum_{j}(\mathbf{WH} - \mathbf{V})_{p,j}\mathbf{H}^{T}_{j,q}$$

$$= \left((\mathbf{WH} - \mathbf{V})\mathbf{H}^{T}\right)_{p,q}$$

$$\nabla_{\mathbf{W}}J(\mathbf{W}, \mathbf{H}) = (\mathbf{WH} - \mathbf{V})\mathbf{H}^{T}$$

Similarly, $\nabla_{\mathbf{H}}J(\mathbf{W}, \mathbf{H}) = \mathbf{W}^{T}(\mathbf{WH} - \mathbf{V})$

# NMF: Gradient Descent Solution

## Procedures:

Step 1. Initialize non-negative matrices **W**, **H**

Step 2. Iteratively update **W**, **H** through gradient descent until convergence

$$\mathbf{H}_{new} \leftarrow \mathbf{H} - \eta_{\mathbf{H}} \mathbf{W}^{\mathrm{T}} (\mathbf{WH} - \mathbf{V})$$

$$\mathbf{W}_{new} \leftarrow \mathbf{W} - \eta_{\mathbf{W}} (\mathbf{WH}_{new} - \mathbf{V}) \mathbf{H}_{new}^{\mathrm{T}}$$

How to ensure non-negativity?

### Method 1:

After each update, clip any negative values

$$\mathbf{W} \leftarrow \max(\mathbf{W}, \mathbf{0}), \mathbf{H} \leftarrow \max(\mathbf{H}, \mathbf{0})$$

### Method 2:

Set $\eta_{\mathbf{H}_{p,q}} = \dfrac{\mathbf{H}_{p,q}}{[\mathbf{W}^{\mathrm{T}}\mathbf{WH}]_{p,q}}, \eta_{\mathbf{W}_{p,q}} = \dfrac{\mathbf{W}_{p,q}}{[\mathbf{WH}_{new}\mathbf{H}_{new}^{\mathrm{T}}]_{p,q}},$

$$(\mathbf{H}_{new})_{p,q} \leftarrow \mathbf{H}_{p,q} \dfrac{[\mathbf{W}^{\mathrm{T}}\mathbf{V}]_{p,q}}{[\mathbf{W}^{\mathrm{T}}\mathbf{WH}]_{p,q}}, \quad (\mathbf{W}_{new})_{p,q} \leftarrow \mathbf{W}_{p,q} \dfrac{[\mathbf{VH}_{new}^{\mathrm{T}}]_{p,q}}{[\mathbf{WH}_{new}\mathbf{H}_{new}^{\mathrm{T}}]_{p,q}}.$$

**Multiplicative Update**

Proposed by D. Lee and H. Seung (NIPS 2000)

43

# NMF: Multiplicative Update

**Procedures:**

Step 1. Initialize non-negative matrices **W**, **H**

Step 2. Iteratively update **W**, **H** until convergence

$$(\mathbf{H}_{new})_{p,q} \leftarrow \mathbf{H}_{p,q} \frac{[\mathbf{W}^{\mathrm{T}}\mathbf{V}]_{p,q}}{[\mathbf{W}^{\mathrm{T}}\mathbf{W}\mathbf{H}]_{p,q}}$$

$$(\mathbf{W}_{new})_{p,q} \leftarrow \mathbf{W}_{p,q} \frac{[\mathbf{V}\mathbf{H}_{new}{}^{\mathrm{T}}]_{p,q}}{[\mathbf{W}\mathbf{H}_{new}\mathbf{H}_{new}{}^{\mathrm{T}}]_{p,q}}$$

Global Minimum?

# Examples: Compute W, H

Suppose $r = 2$,

**W** is initialized to $\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$,

**H** is initialized to $\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$

| Observations | Feature 1 | Feature 2 |
|---|---|---|
| 1 | 1 | 2 |
| 2 | 3 | 4 |
| 3 | 5 | 6 |
| 4 | 7 | 8 |

What are **H** and **W** after implementing 1 iteration of multiplicative update?
What is the corresponding reconstructed data matrix **V**?

# Examples: Compute W, H

Suppose *r* = 2,

**W** is initialized to $\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$,

**H** is initialized to $\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$

n=4, d=2

| Observations | Feature 1 | Feature 2 |
|---|---|---|
| 1 | 1 | 2 |
| 2 | 3 | 4 |
| 3 | 5 | 6 |
| 4 | 7 | 8 |

1st Iteration:

$$(\mathbf{H_{new}})_{p,q} \leftarrow \mathbf{H}_{p,q} \frac{[\mathbf{W}^T\mathbf{V}]_{p,q}}{[\mathbf{W}^T\mathbf{W}\mathbf{H}]_{p,q}}$$

$$\mathbf{H}_{new} \leftarrow \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \times \frac{\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}\begin{bmatrix} 1 & 3 & 5 & 7 \\ 2 & 4 & 6 & 8 \end{bmatrix}}{\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}} = \begin{bmatrix} 0.56 & 1.22 & 1.89 & 2.56 \\ 0.44 & 1.11 & 1.78 & 2.44 \end{bmatrix}$$

$$\mathbf{V}_{reconst} = \mathbf{W}_{new}\,\mathbf{H}_{new} = \begin{bmatrix} 1.33 & 3.17 & 5.01 & 6.85 \\ 1.68 & 3.84 & 5.99 & 8.15 \end{bmatrix}$$

$$(\mathbf{W_{new}})_{p,q} \leftarrow \mathbf{W}_{p,q} \frac{[\mathbf{V}\mathbf{H_{new}}^T]_{p,q}}{[\mathbf{W}\mathbf{H_{new}}\mathbf{H_{new}}^T]_{p,q}}$$
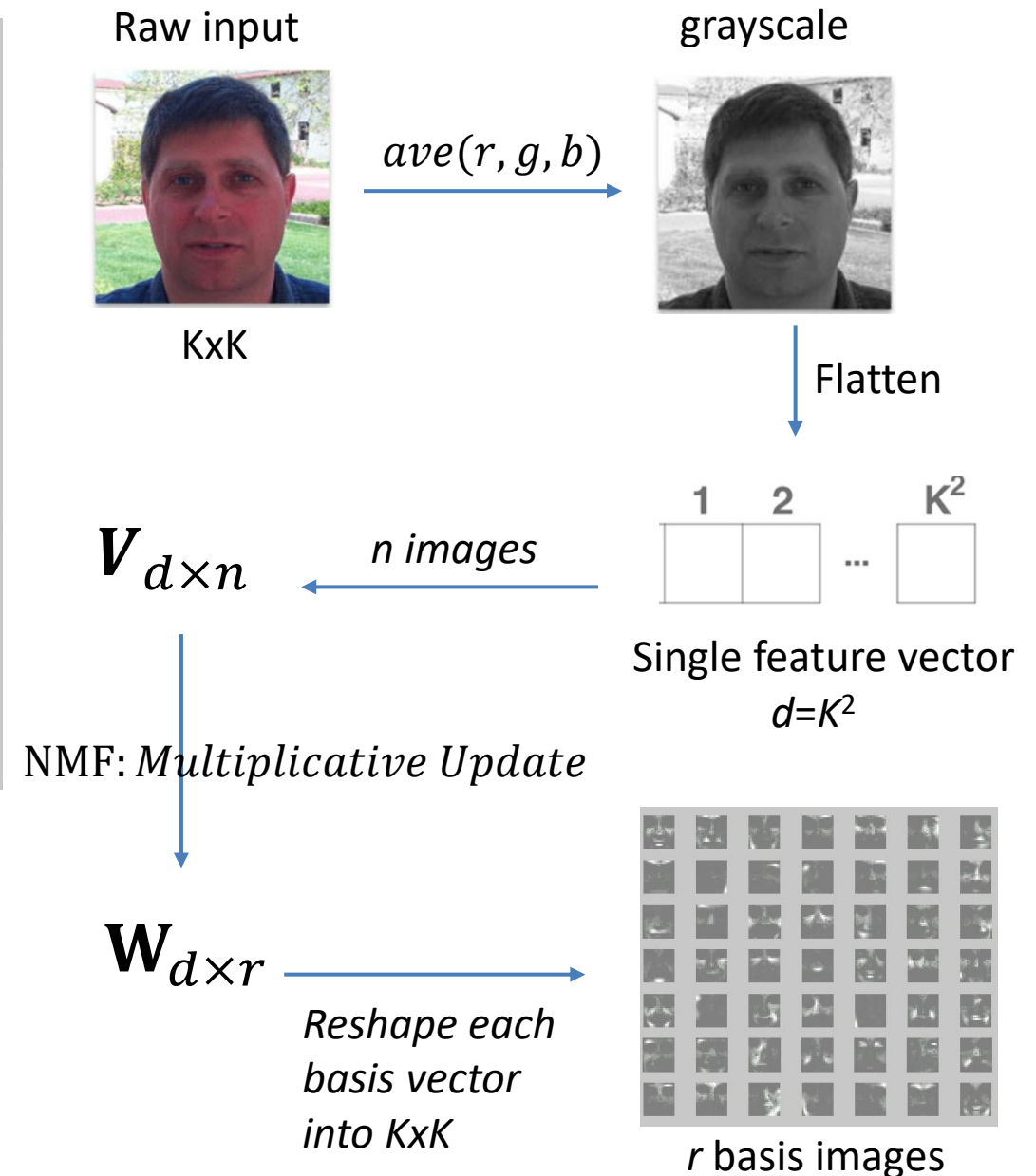
$$\mathbf{W}_{new} \leftarrow \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \times \frac{\begin{bmatrix} 1 & 3 & 5 & 7 \\ 2 & 4 & 6 & 8 \end{bmatrix}\begin{bmatrix} 0.56 & 0.44 \\ 1.22 & 1.11 \\ 1.89 & 1.78 \\ 2.56 & 2.44 \end{bmatrix}}{\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}\begin{bmatrix} 0.56 & 1.22 & 1.89 & 2.56 \\ 0.44 & 1.11 & 1.78 & 2.44 \end{bmatrix}\begin{bmatrix} 0.56 & 0.44 \\ 1.22 & 1.11 \\ 1.89 & 1.78 \\ 2.56 & 2.44 \end{bmatrix}} = \begin{bmatrix} 0.92 & 1.84 \\ 2.16 & 1.08 \end{bmatrix}$$

| Observations | Basis 1 | Basis 2 |
|---|---|---|
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |

# Visualize Basis Images

Raw input

grayscale

$ave(r, g, b)$

KxK

Flatten

$$V_{d \times n}$$

n images

| 1 | 2 | | $K^2$ |
|---|---|---|---|
| | | ... | |

Single feature vector
$d=K^2$

NMF: *Multiplicative Update*

NMF basis images

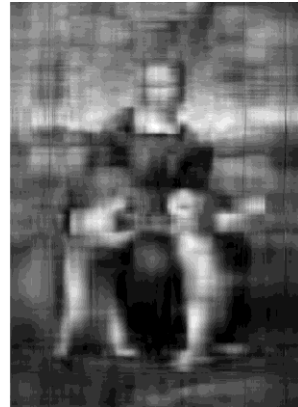$$W_{d \times r}$$

*Reshape each basis vector into KxK*

r basis images

# Reconstruction from W

Suppose we want to reconstruct the $i$-th image among the $n$ samples

$$\widehat{x}_i = \mathbf{WH}_{(:,i)}$$

Original

Image Dimension:
549x767

r=10    r=20    r=40    r=80

r=100    r=150    r=200    r=250

# Discussions

- For a new image, how to obtain the reconstruction coefficients?

- How to use NMF for classification?

- How to choose *r*

# Summary

| | PCA | NMF |
|---|---|---|
| **Representation** | Holistic | Part-based |
| **Basis Vector** | (PCs)<br>Globalized;<br>Orthogonal;<br>Sign ambiguity | (**W**)<br>Localized;<br>Nonnegative |
| **Features** | (PC projections)<br>Sign ambiguity | (**H**)<br>Nonnegative |
| **Optimization** | Global Optimum<br>(Eigenvalue Problem) | Local Optimum<br>(Multiplicative Update) |

# Papers to Read and Self-Study

- D. Lee and H. Seung. [Algorithms for Non-negative Matrix Factorization](#) NIPS (2000).

- ICA (Independent Component Analysis):
  [http://en.wikipedia.org/wiki/Independent_component_analysis](http://en.wikipedia.org/wiki/Independent_component_analysis)

- CCA (Canonical Correlation Analysis):
  [http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.101.6359&rep=rep1&type=pdf](http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.101.6359&rep=rep1&type=pdf)