



HUST

ĐẠI HỌC BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

ONE LOVE. ONE FUTURE.



ĐẠI HỌC
BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY
OF SCIENCE AND TECHNOLOGY

Machine Learning

IT3190E

Lecture: Introduction

ONE LOVE. ONE FUTURE.

Contents

- Lecture 1: Introduction to Machine Learning
- Lecture 2: Linear regression
- **Lecture 3+4: Clustering**
- Lecture 5: Decision tree and Random forest
- Lecture 6: Neural networks
- Lecture 7: Support vector machines
- Lecture 8: Performance evaluation
- Lecture 9: Probabilistic models
- Lecture 10: Ensemble learning
- Lecture 11: Reinforcement learning
- Lecture 12: Regularization
- Lecture 13: Discussion on some advanced topics



SOICT

HUST

ĐẠI HỌC BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

Dimensionality Reduction: Principal Component Analysis

Dam Quang Tuan

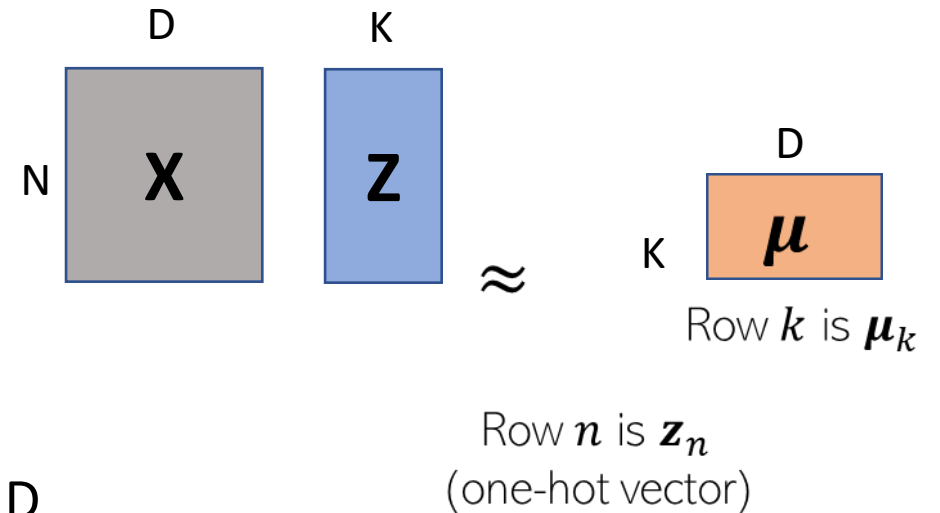
K-means loss function: recap

$\mathbf{z}_n = [z_{n1}, z_{n2}, \dots, z_{nK}]$
denotes a length K one-hot
encoding of \mathbf{x}_n

- Remember the matrix factorization view of the k-means loss function?

$$L(\mu, \mathbf{X}, \mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$$

$$= \underbrace{\|\mathbf{X} - \mathbf{Z}\mu\|_F^2}_{\text{matrix factorization view}}$$



- We approximated an $N \times D$ matrix with
 - An $N \times K$ matrix and a
 - $K \times D$ matrix
- This could be storage efficient if K is much smaller than D

Dimensionality Reduction

- A broad class of techniques
- Goal is to compress the original representation of the inputs
- Example: Approximate each input $\mathbf{x}_n \in \mathbb{R}^D$, $n = 1, 2, \dots, N$ as a linear combination of $K < \min\{D, N\}$ “basis” vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$, each also $\in \mathbb{R}^D$

Note: These “basis” vectors need not necessarily be linearly independent. But for some dim. red. techniques, e.g., classic principal component analysis (PCA), they are

$$\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{w}_k = \mathbf{W} \mathbf{z}_n$$

$\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ is $D \times K$

$\mathbf{z}_n = [z_{n1}, z_{n2}, \dots, z_{nK}]$ is $K \times 1$

- We have represented each $\mathbf{x}_n \in \mathbb{R}^D$ by a K -dim vector \mathbf{z}_n (a new feat. rep)
- To store N such inputs $\{\mathbf{x}_n\}_{n=1}^N$, we need to keep \mathbf{W} and $\{\mathbf{z}_n\}_{n=1}^N$
 - Originally we required $N \times D$ storage, now $N \times K + D \times K = (N + D) \times K$ storage
 - If $K \ll \min\{D, N\}$, this yields substantial storage saving, hence good compression

Can think of \mathbf{W} as a **linear** mapping that transforms low-dim \mathbf{z}_n to high-dim \mathbf{x}_n

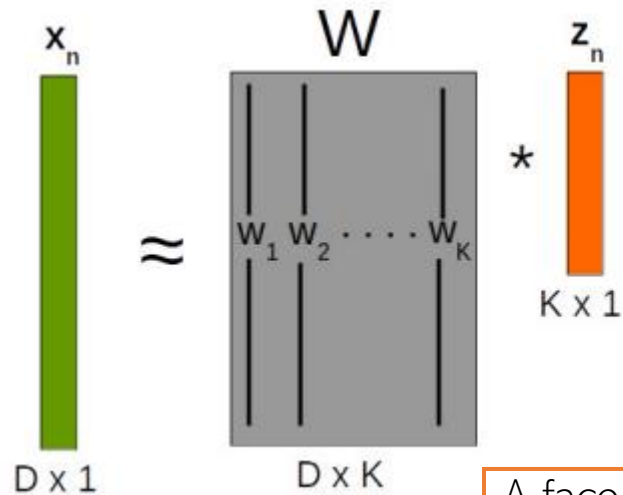
Some dim-red techniques assume a nonlinear mapping function f such that $\mathbf{x}_n = f(\mathbf{z}_n)$

For example, f can be modeled by a kernel or a deep neural net



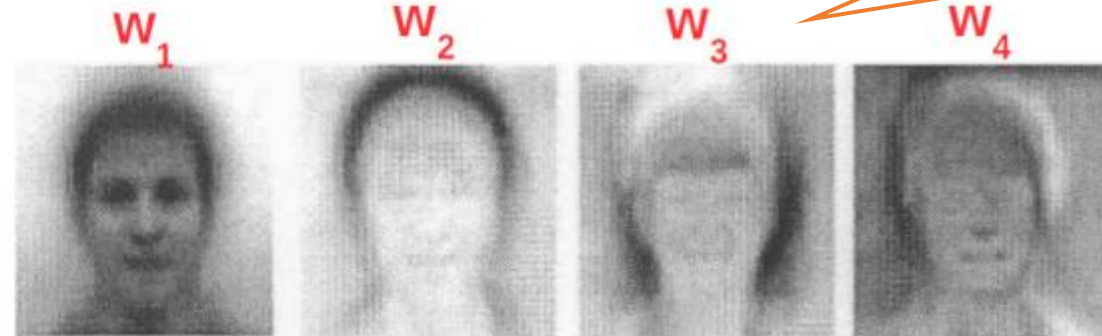
Dimensionality Reduction

■ Dim-red for face images

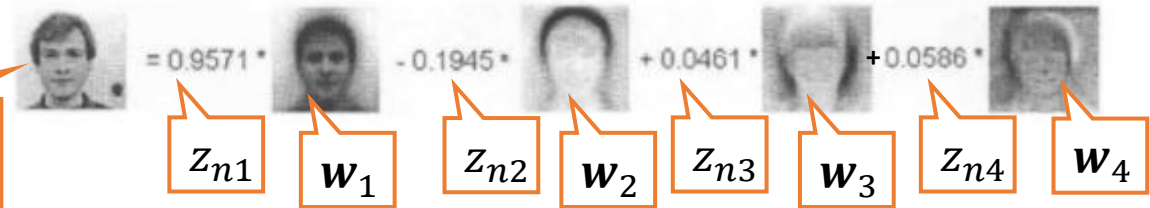


Each “basis” image is like a “template” that captures the common properties of face images in the dataset

$K=4$ “basis” face images



A face image $x_n \in \mathbb{R}^D$



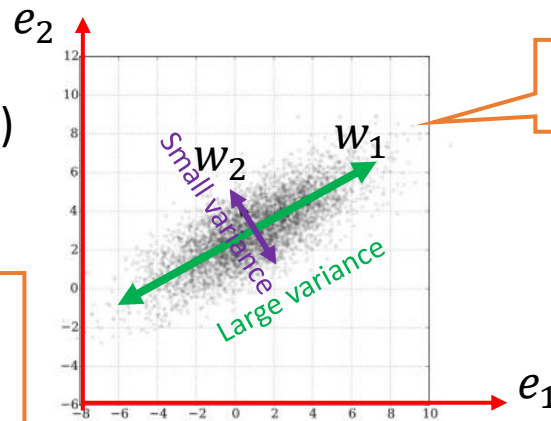
- In this example, $z_n \in \mathbb{R}^K$ ($K = 4$) is a low-dim feature rep. for $x_n \in \mathbb{R}^D$ Like 4 new features
- Essentially, each face image in the dataset now represented by just 4 real numbers ☺
- Different dim-red algos differ in terms of how the basis vectors are defined/learned
 - .. And in general, how the function f in the mapping $x_n = f(z_n)$ is defined

Principal Component Analysis (PCA)

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
 - Learning directions (co-ordinate axes) that capture maximum variance in data

e_1, e_2 : Standard co-ordinate axis ($\mathbf{x} = [x_1, x_2]$)

w_1, w_2 : New co-ordinate axis ($\mathbf{z} = [z_1, z_2]$)



PCA is essentially doing a change of axes in which we are representing the data

Each input will still have 2 co-ordinates, in the new co-ordinate system, equal to the distances measured from the new origin

To reduce dimension, can only keep the co-ordinates of those directions that have largest variances (e.g., in this example, if we want to reduce to one-dim, we can keep the co-ordinate z_1 of each point along w_1 and throw away z_2). We won't lose much information

- Learning projection directions that result in smallest reconstruction error

$$\operatorname{argmin}_{\mathbf{W}, \mathbf{Z}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{W}\mathbf{z}_n\|^2 = \operatorname{argmin}_{\mathbf{W}, \mathbf{Z}} \|\mathbf{X} - \mathbf{Z}\mathbf{W}\|^2$$

Subject to orthonormality constraints: $\mathbf{w}_i^T \mathbf{w}_j = 0$ for $i \neq j$ and $\|\mathbf{w}_i\|^2 = 1$

- PCA also assumes that the projection directions are orthonormal

Principal Component Analysis: the algorithm

- Center the data (subtract the mean $\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$ from each data point)
- Compute the $D \times D$ covariance matrix \mathbf{S} using the centered data matrix \mathbf{X} as

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad (\text{Assuming } \mathbf{X} \text{ is arranged as } N \times D)$$

- Do an eigendecomposition of the covariance matrix \mathbf{S} (many methods exist)
- Take top $K < D$ leading eigenvectors $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K\}$ with eigvalues $\{\lambda_1, \lambda_2, \dots, \lambda_K\}$
- The K -dimensional projection/embedding of each input is

$$\mathbf{z}_n \approx \mathbf{W}_K^T \mathbf{x}_n$$

$\mathbf{W}_K = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ is the “projection matrix” of size $D \times K$

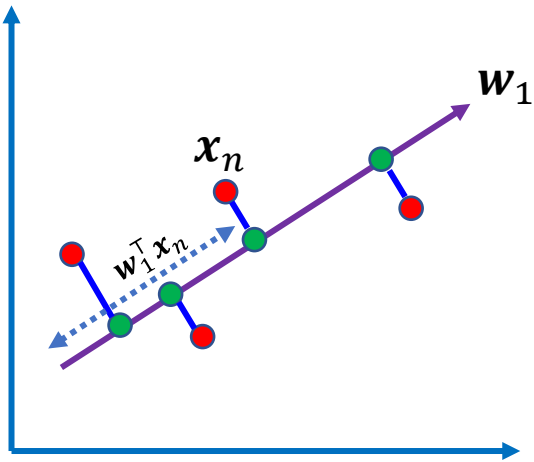
Note: Can decide how many eigvecs to use based on how much variance we want to capture (recall that each λ_k gives the variance in the k^{th} direction (and their sum is the total variance))



Understanding PCA: The variance perspective

Solving PCA by Finding Max. Variance Directions

- Consider projecting an input $\mathbf{x}_n \in \mathbb{R}^D$ along a direction $\mathbf{w}_1 \in \mathbb{R}^D$
- Projection/embedding of \mathbf{x}_n (red points below) will be $\mathbf{w}_1^\top \mathbf{x}_n$ (green pts below)



Mean of projections of all inputs:

$$\frac{1}{N} \sum_{n=1}^N \mathbf{w}_1^\top \mathbf{x}_n = \mathbf{w}_1^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right) = \mathbf{w}_1^\top \boldsymbol{\mu}$$

Variance of the projections:

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{w}_1^\top \mathbf{x}_n - \mathbf{w}_1^\top \boldsymbol{\mu})^2 = \frac{1}{N} \sum_{n=1}^N \{\mathbf{w}_1^\top (\mathbf{x}_n - \boldsymbol{\mu})\}^2 = \mathbf{w}_1^\top \mathbf{S} \mathbf{w}_1$$

\mathbf{S} is the $D \times D$ cov matrix of the data:

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top$$

- Want \mathbf{w}_1 such that variance $\mathbf{w}_1^\top \mathbf{S} \mathbf{w}_1$ is maximized

$$\operatorname{argmax}_{\mathbf{w}_1} \mathbf{w}_1^\top \mathbf{S} \mathbf{w}_1 \quad \text{s.t.} \quad \mathbf{w}_1^\top \mathbf{w}_1 = 1$$

Need this constraint otherwise the objective's max will be infinity

For already centered data, $\boldsymbol{\mu} = \mathbf{0}$ and $\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top$

Max. Variance Direction

Variance along the direction \mathbf{w}_1

- Our objective function was $\operatorname{argmax}_{\mathbf{w}_1} \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1$ s.t. $\mathbf{w}_1^T \mathbf{w}_1 = 1$
- Can construct a Lagrangian for this problem

$$\operatorname{argmax}_{\mathbf{w}_1} \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 + \lambda_1 (1 - \mathbf{w}_1^T \mathbf{w}_1)$$

- Taking derivative w.r.t. \mathbf{w}_1 and setting to zero gives $\mathbf{S} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1$
- Therefore \mathbf{w}_1 is an **eigenvector** of the cov matrix \mathbf{S} with eigenvalue λ_1
- Claim:** \mathbf{w}_1 is the eigenvector of \mathbf{S} with largest eigenvalue λ_1 . Note that

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1^T \mathbf{w}_1 = \lambda_1$$

- Thus variance $\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1$ will be max. if λ_1 is the largest eigenvalue (and \mathbf{w}_1 is the corresponding top eigenvector; also known as the first **Principal Component**)
- Other large variance directions can also be found likewise (with each being orthogonal to all others) using the eigendecomposition of cov matrix \mathbf{S} (this is PCA)

Note: Total variance of the data is equal to the sum of eigenvalues of \mathbf{S} , i.e., $\sum_{d=1}^D \lambda_d$

PCA would keep the top $K < D$ such directions of largest variances

Note: In general, \mathbf{S} will have D eigvecs



Understanding PCA: The reconstruction perspective

Alternate Basis and Reconstruction

- Representing a data point $\mathbf{x}_n = [x_{n1}, x_{n2}, \dots, x_{nD}]^T$ in the standard orthonormal basis $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_D\}$

$$\mathbf{x}_n = \sum_{d=1}^D x_{nd} \mathbf{e}_d$$

\mathbf{e}_d is a vector of all zeros except a single 1 at the d^{th} position. Also, $\mathbf{e}_d^T \mathbf{e}_{d'} = 0$ for $d \neq d'$

- Let's represent the same data point in a new orthonormal basis $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_D\}$

z_{nd} is the projection of \mathbf{x}_n along the direction \mathbf{w}_d since $z_{nd} = \mathbf{w}_d^T \mathbf{x}_n = \mathbf{x}_n^T \mathbf{w}_d$ (verify)

$$\mathbf{x}_n = \sum_{d=1}^D z_{nd} \mathbf{w}_d$$

$\mathbf{z}_n = [z_{n1}, z_{n2}, \dots, z_{nD}]^T$ denotes the co-ordinates of \mathbf{x}_n in the new basis

- Ignoring directions along which projection z_{nd} is small, we can approximate \mathbf{x}_n as

$$\mathbf{x}_n \approx \hat{\mathbf{x}}_n = \sum_{d=1}^K z_{nd} \mathbf{w}_d = \sum_{d=1}^K (\mathbf{x}_n^T \mathbf{w}_d) \mathbf{w}_d = \sum_{d=1}^K (\mathbf{w}_d \mathbf{w}_d^T) \mathbf{x}_n$$

Note that $\|\mathbf{x}_n - \sum_{d=1}^K (\mathbf{w}_d \mathbf{w}_d^T) \mathbf{x}_n\|^2$ is the **reconstruction error** on \mathbf{x}_n . Would like it to minimize w.r.t. $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$

- Now \mathbf{x}_n is represented by $K < D$ dim. rep. $\mathbf{z}_n = [z_{n1}, z_{n2}, \dots, z_{nK}]$ and (verify)

Also, $\mathbf{x}_n \approx \mathbf{W}_K \mathbf{z}_n$

$\mathbf{W}_K = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ is the "projection matrix" of size $D \times K$

Minimizing Reconstruction Error

- We plan to use only K directions $[\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ so would like them to be such that the total reconstruction error is minimized

$$\mathcal{L}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = \sum_{n=1}^N \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|^2 = \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{d=1}^K (\mathbf{w}_d \mathbf{w}_d^\top) \mathbf{x}_n \right\|^2 = C - \sum_{d=1}^K \mathbf{w}_d^\top \mathbf{S} \mathbf{w}_d \quad (\text{verify})$$

Constant; doesn't depend on the \mathbf{w}_d 's

Variance along \mathbf{w}_d

- Each optimal \mathbf{w}_d can be found by solving

$$\operatorname{argmin}_{\mathbf{w}_d} \mathcal{L}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = \operatorname{argmax}_{\mathbf{w}_d} \mathbf{w}_d^\top \mathbf{S} \mathbf{w}_d$$

- Thus minimizing the reconstruction error is equivalent to maximizing variance
- The K directions can be found by solving the eigendecomposition of \mathbf{S}

- Note: $\sum_{d=1}^K \mathbf{w}_d^\top \mathbf{S} \mathbf{w}_d = \operatorname{trace}(\mathbf{W}_K^\top \mathbf{S} \mathbf{W}_K)$

- Thus $\operatorname{argmax}_{\mathbf{W}_K} \operatorname{trace}(\mathbf{W}_K^\top \mathbf{S} \mathbf{W}_K)$ s.t. orthonormality on columns of \mathbf{W}_K is the same as solving the eigendec. of \mathbf{S} (recall that Spectral Clustering also required solving this)

Principal Component Analysis

- Center the data (subtract the mean $\boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$ from each data point)
- Compute the $D \times D$ covariance matrix \mathbf{S} using the centered data matrix \mathbf{X} as

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad (\text{Assuming } \mathbf{X} \text{ is arranged as } N \times D)$$

- Do an eigendecomposition of the covariance matrix \mathbf{S} (many methods exist)
- Take top $K < D$ leading eigenvectors $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K\}$ with eigvalues $\{\lambda_1, \lambda_2, \dots, \lambda_K\}$
- The K -dimensional projection/embedding of each input is

$$\mathbf{z}_n \approx \mathbf{W}_K^T \mathbf{x}_n$$

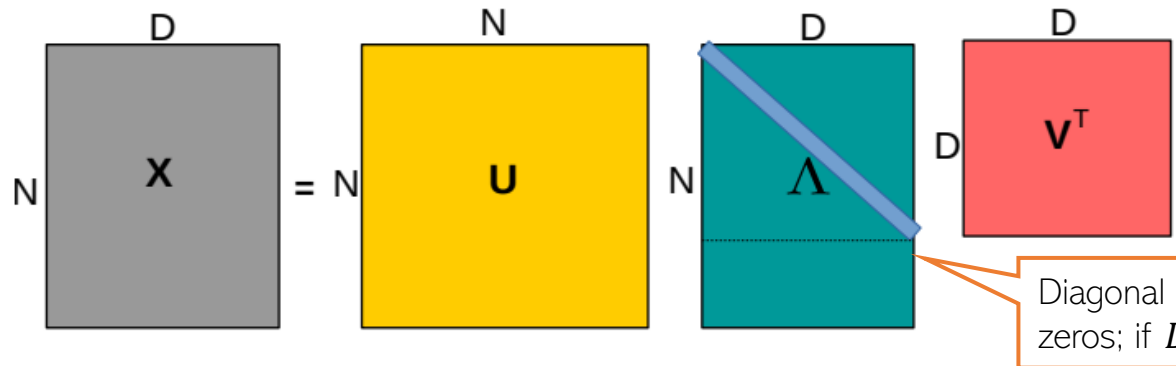
$\mathbf{W}_K = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$ is the
“projection matrix” of size $D \times K$

Note: Can decide how many eigvecs to use based on how much variance we want to capture (recall that each λ_k gives the variance in the k^{th} direction (and their sum is the total variance))



Singular Value Decomposition (SVD)

- Any matrix \mathbf{X} of size $N \times D$ can be represented as the following decomposition



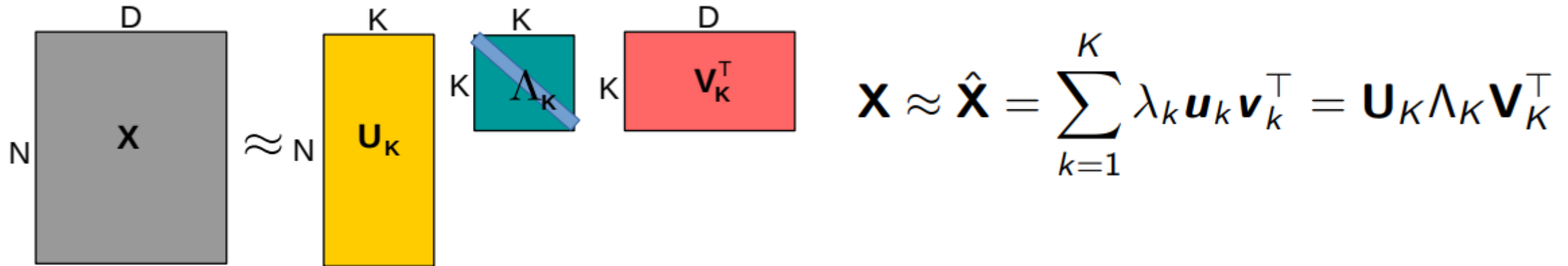
$$\mathbf{X} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T = \sum_{k=1}^{\min\{N,D\}} \lambda_k \mathbf{u}_k \mathbf{v}_k^T$$

Diagonal matrix. If $N > D$, last $D - N$ rows are all zeros; if $D > N$, last $D - N$ columns are all zeros

- $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N]$ is $N \times N$ matrix of **left singular vectors**, each $\mathbf{u}_n \in \mathbb{R}^N$
 - \mathbf{U} is also orthonormal
- $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D]$ is $D \times D$ matrix of **right singular vectors**, each $\mathbf{v}_d \in \mathbb{R}^D$
 - \mathbf{V} is also orthonormal
- $\mathbf{\Lambda}$ is $N \times D$ with only $\min(N, D)$ diagonal entries - **singular values**
- Note: If \mathbf{X} is symmetric then it is known as eigenvalue decomposition ($\mathbf{U} = \mathbf{V}$)

Low-Rank Approximation via SVD

- If we just use the top $K < \min\{N, D\}$ singular values, we get a rank- K SVD

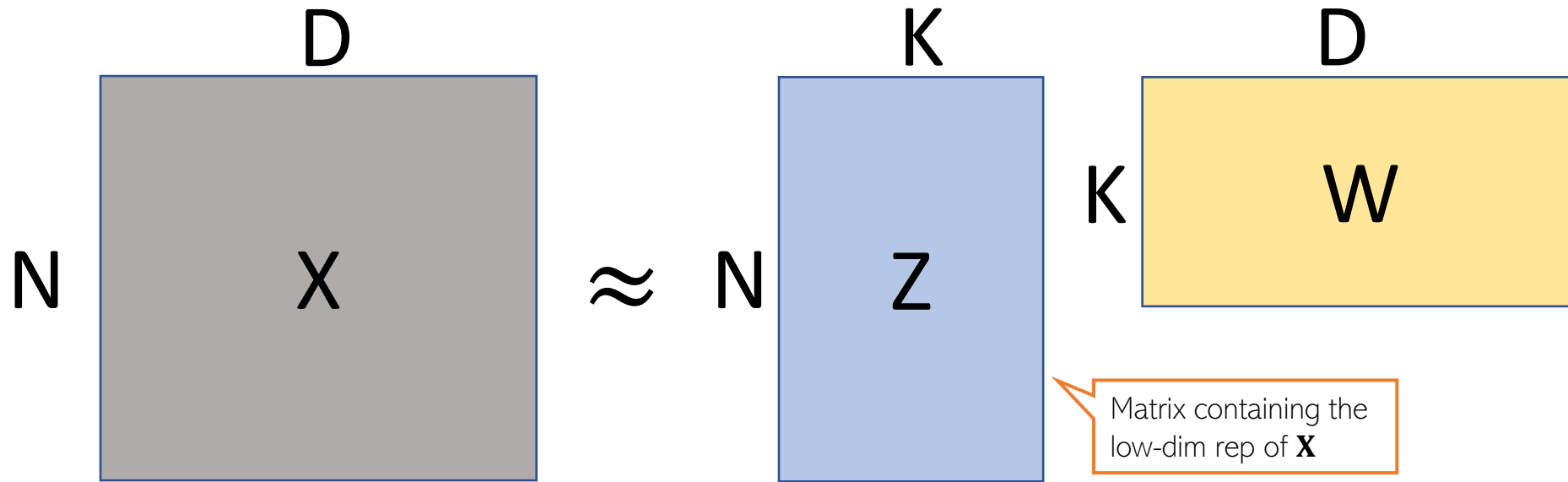


$$\mathbf{X} \approx \mathbf{U}_K \mathbf{\Lambda}_K \mathbf{V}_K^T \quad \mathbf{X} \approx \hat{\mathbf{X}} = \sum_{k=1}^K \lambda_k \mathbf{u}_k \mathbf{v}_k^T = \mathbf{U}_K \mathbf{\Lambda}_K \mathbf{V}_K^T$$

- Above SVD approx. can be shown to minimize the reconstruction error $\|\mathbf{X} - \hat{\mathbf{X}}\|$
 - Fact: SVD gives the best rank- K approximation of a matrix
- PCA is done by doing SVD on the covariance matrix \mathbf{S} (left and right singular vectors are the same and become eigenvectors, singular values become eigenvalues)

Dim-Red as Matrix Factorization

- If we don't care about the orthonormality constraints, then dim-red can also be achieved by solving a matrix factorization problem on the data matrix \mathbf{X}



$$\{\hat{\mathbf{Z}}, \hat{\mathbf{W}}\} = \operatorname{argmin}_{\mathbf{Z}, \mathbf{W}} \|\mathbf{X} - \mathbf{ZW}\|^2$$

If $K < \min\{D, N\}$, such a factorization gives a low-rank approximation of the data matrix \mathbf{X}

- Can solve such problems using ALT-OPT
- Can impose various constraints on \mathbf{Z} and \mathbf{W} (e.g., sparsity, non-negativity, etc)



References

CS771: Intro to Machine Learning (Fall 2021), Nisheeth Srivastava, IIT Kanpur



HUST

Thanks