# CZ4041/CE4041: Machine Learning

## Lesson 12: Dimensionality Reduction
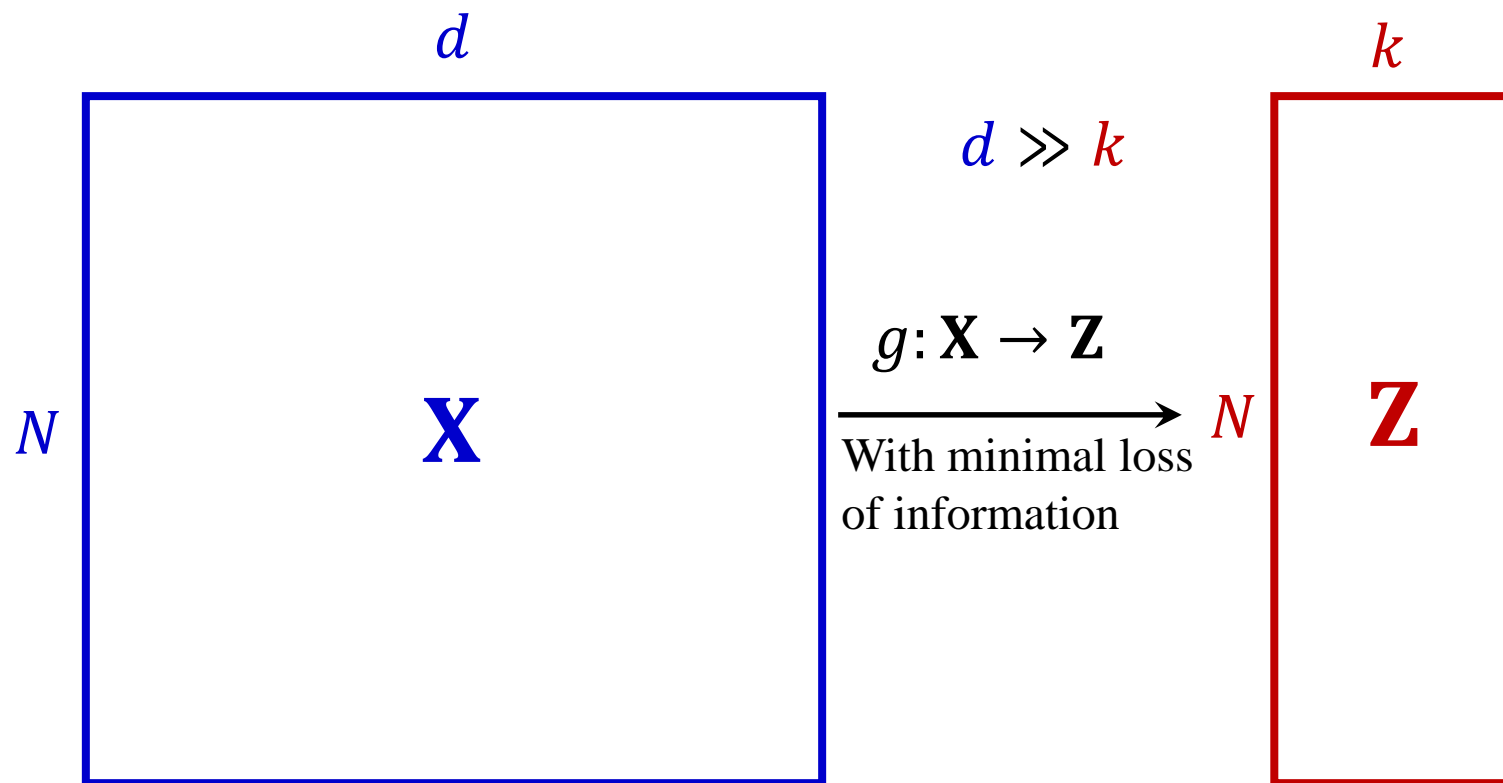
Kelly KE

School of Computer Science and Engineering,
NTU, Singapore

# High-level Idea

- To summarize observed high-dimensional data points with low-dimensional vectors

$d$

$k$

$d \gg k$

$g : \mathbf{X} \to \mathbf{Z}$

With minimal loss of information

$N$    **X**      $N$    **Z**

# **Why Dimensionality Reduction**

- Avoid curse of dimensionality
  - Distance-based methods, e.g., $K$-NN Classifiers, $K$-means
- Reduce amount of time and memory required by other machine learning algorithms
- Allow data to be more easily visualized
  - 2D or 3D
- Reduce noise, and thus improve the performance of the downstream machine learning tasks

# **Dimensionality Reduction Approaches**

- Feature Selection
  - To select a subset of $k$ features from the original $d$ features to represent each data instance
    - Brute-force approach
    - Greedy search

- Feature Extraction
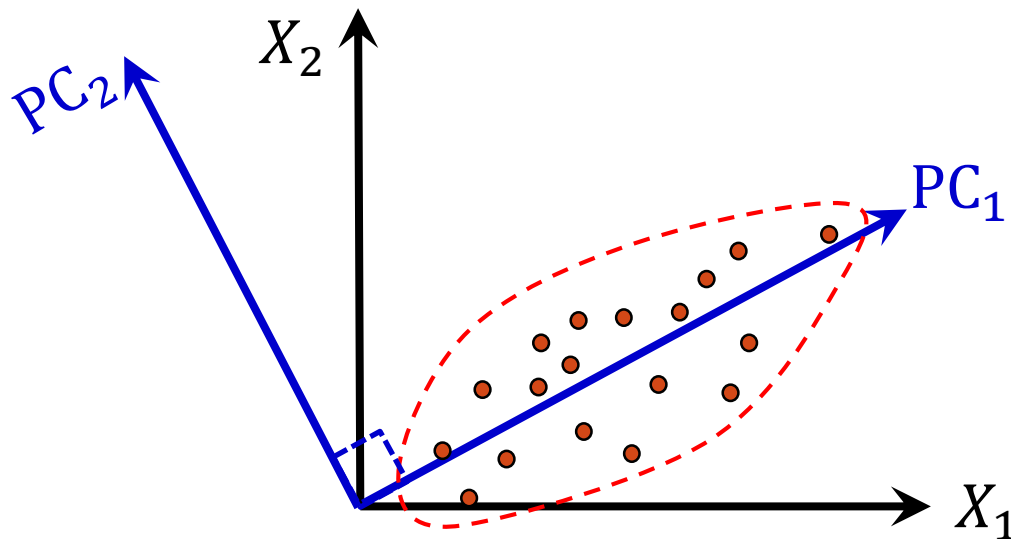  - To learn $k$ <u>new</u> features from the original $d$ features to represent each data instance
    - Linear combination of original features
      - Principal component analysis
    - Nonlinear combination of original features

# Principal Component Analysis

- One of the most widely-used (unsupervised) dimensionality reduction methods

- Takes a data matrix of $N$ data points by $d$ features, and summarizes it by principal components that are linear combinations of the original $d$ variables

- The first $k$ components display as much as possible of the variation among data instances

# PCA: Geometric Rationale

- Goal: to find a projection or rotation of the original $d$-dimensional coordinate system to capture the largest amount of variation in data
  - Ordered s.t. the 1st principal component has the highest variance, the 2nd component has the next highest variance, …, the $d$-th component has the lowest variance
  - Principal components are orthogonal to each other

# PCA: Algorithm

Input: $\mathcal{D} = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N\}$ a set of observed data

1. Centering the data points s.t. the mean is **0**

$$\widehat{\boldsymbol{\mu}} = \frac{1}{N}\sum_{i=1}^{N} \boldsymbol{x}_i \quad\longrightarrow\quad \boldsymbol{x}_i = \boldsymbol{x}_i - \widehat{\boldsymbol{\mu}}$$

2. Compute sample covariance matrix

$$\widetilde{\boldsymbol{\Sigma}} = \frac{1}{N-1}\sum_{i=1}^{N} \boldsymbol{x}_i \boldsymbol{x}_i{}^T$$

Each $\boldsymbol{u}_i$ is of $d$ dimensions

3. Compute eigenvectors of $\widetilde{\boldsymbol{\Sigma}}$, $\{\boldsymbol{u}_1, \boldsymbol{u}_2, \ldots, \boldsymbol{u}_d\}$, which are sorted based on their eigenvalues in non-increasing order, i.e., $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$

4. Select the first $k$ eigenvectors to construct principal components

# PCA: Algorithm (Illustration)



Centered data matrix $\quad$ Covariance matrix $\quad$ Matrix of eigenvectors

Projection matrix with top $k$ eigenvectors

# **Derivation of PCA**

- [The variance preservation view](#)
  - The first $k$ components display as much as possible of the variation among data instances
- The minimum reconstruction view
  - The first $k$ components convey maximum useful information of original data instances

Appendix (optional)

# Eigenvalues & Eigenvectors

- Given a $d$-by-$d$ square matrix $\mathbf{A}$, if there exists a non-zero $d$-dimensional vector $\boldsymbol{u}$, s.t.

$$\mathbf{A}\boldsymbol{u} = \lambda\boldsymbol{u} \qquad \text{scalar}$$

  then $\boldsymbol{u}$ is an eigenvector of $\mathbf{A}$, and $\lambda$ is called the corresponding eigenvalue

- Notes:
  - There are $d$ eigenvectors and eigenvalues
  - An eigenvalue can be positive, negative or zero
  - An eigenvector cannot be a zero vector
  - Eigenvectors are orthogonal to each other

# Properties of Eigenvalues

Given a square matrix $\mathbf{A}$ ($d$-by-$d$)

- $\mathbf{A}$ is invertible ($\mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$ or $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$) if all the eigenvalues of $\mathbf{A}$ are non-zero (positive or negative)

- If all the eigenvalues of $\mathbf{A}$ are non-negative, then $\mathbf{A}$ is a positive semi-definite matrix:

  $\text{For any non}-\text{zero vector } \boldsymbol{x} \in \mathbb{R}^{d \times 1}, \text{we have } \boldsymbol{x}^T \mathbf{A} \boldsymbol{x} \geq 0$

- If all the eigenvalues of $\mathbf{A}$ are positive, then $\mathbf{A}$ is a positive definite matrix:

  $\text{For any non}-\text{zero vector } \boldsymbol{x} \in \mathbb{R}^{d \times 1}, \text{we have } \boldsymbol{x}^T \mathbf{A} \boldsymbol{x} > 0$

# Properties of Eigenvalues (cont.)

- Recall: when inducing a closed form solution of regularized linear regression model, we mentioned that if a matrix **A** can be written as

$$\mathbf{A} = \mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}, \text{where } \mathbf{X} \in \mathbb{R}^{N \times d}, \mathbf{I} \in \mathbb{R}^{d \times d} \text{and } \lambda > 0$$

then **A** is always invertible:

$$\exists\, \mathbf{A}^{-1}, \text{s.t., } \mathbf{A}^{-1}\mathbf{A} = \mathbf{I} \text{ or } \mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$$

# Properties of Eigenvalues (cont.)

- We first prove $\mathbf{A}$ is positive definite
  - For any non-zero vector $\boldsymbol{x} \in \mathbb{R}^{d \times 1}$

$$\boldsymbol{x}^T \mathbf{A} \boldsymbol{x} = \boldsymbol{x}^T \big(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}\big) \boldsymbol{x}$$

$$= \boldsymbol{x}^T \big(\mathbf{X}^T \mathbf{X}\big) \boldsymbol{x} + \boldsymbol{x}^T (\lambda \mathbf{I}) \boldsymbol{x}$$

Denote $\boldsymbol{z} = \mathbf{X}\boldsymbol{x}$

$$= \boldsymbol{z}^T \boldsymbol{z} + \lambda \boldsymbol{x}^T \boldsymbol{x}$$

$$= \|\boldsymbol{z}\|_2^2 + \lambda \|\boldsymbol{x}\|_2^2$$

$\|\boldsymbol{z}\|_2^2 \geq 0$ and $\|\boldsymbol{z}\|_2^2 = 0$ if and only if $\boldsymbol{z} = \mathbf{0}$

$\|\boldsymbol{x}\|_2^2 > 0$ because $\boldsymbol{x} \neq \mathbf{0} \Rightarrow \lambda \|\boldsymbol{x}\|_2^2 > 0$ as long as $\lambda > 0$

$$\boldsymbol{x}^T \mathbf{A} \boldsymbol{x} > 0$$

# Properties of Eigenvalues (cont.)

- As **A** is positive definite, all of its eigenvalues are positive, i.e., non-zero

- Recall: **A** is invertible if all the eigenvalues of **A** are non-zero (either positive or negative)

- Therefore, if a matrix **A** can be written as

$$\mathbf{A} = \mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}, \text{ where } \mathbf{X} \in \mathbb{R}^{N \times d}, \mathbf{I} \in \mathbb{R}^{d \times d} \text{ and } \lambda > 0$$
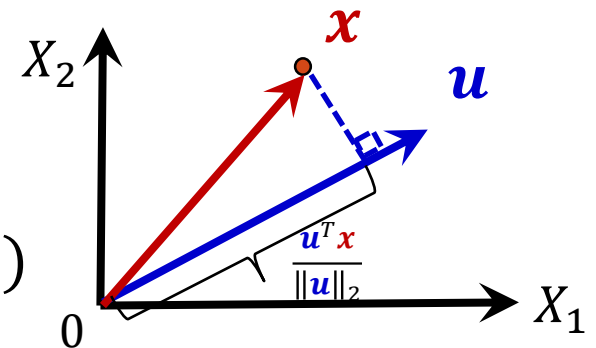
then **A** is invertible!

# PCA: Variance Preservation

- The first $k$ components display as much as possible of the variation among data instances

- Consider a projection of a data point $\boldsymbol{x}$ onto a vector going through the origin, represented by $\boldsymbol{u}$

- The projection of $\boldsymbol{x}$ onto $\boldsymbol{u}$ is

$$\frac{\boldsymbol{u}^T \boldsymbol{x}}{\|\boldsymbol{u}\|_2} = \frac{\|\boldsymbol{u}\|_2 \|\boldsymbol{x}\|_2 \cos(\theta)}{\|\boldsymbol{u}\|_2} = \|\boldsymbol{x}\|_2 \cos(\theta)$$

- For simplicity, consider $\boldsymbol{u}$ with unit length, i.e., $\|\boldsymbol{u}\|_2 = 1$

- The projected instances $\mathcal{D} = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N\}$ onto $\boldsymbol{u}$ are

$$\{\boldsymbol{u}^T \boldsymbol{x}_1, \boldsymbol{u}^T \boldsymbol{x}_2, \ldots, \boldsymbol{u}^T \boldsymbol{x}_N\}$$

# Variance Preservation (cont.)

- In PCA, data points are centered at the beginning

$$\frac{1}{N}\sum_{i=1}^{N} x_i = 0$$

- After projection onto $u$, the mean of data points is still 0

$$\frac{1}{N}\sum_{i=1}^{N} u^T x_i = u^T \frac{1}{N}\sum_{i=1}^{N} x_i = 0$$

- The variance of the data points projected onto $u$ is

$$\frac{1}{N-1}\sum_{i=1}^{N}(u^T x_i - 0)^2 = \frac{1}{N-1}\sum_{i=1}^{N}(u^T x_i)^2$$

Each row is a data instance

$$= u^T \widetilde{\Sigma} u \quad\rightarrow\quad \widetilde{\Sigma} = \frac{1}{N-1}\sum_{i=1}^{N} x_i x_i^T = \frac{1}{N-1} X^T X$$

# Variance Preservation (cont.)

- The goal of PCA (for simplicity, projected on 1 principal component only) is to find $\boldsymbol{u}$ that maximizes the variance, expecting to maximally preserve distinction among data

- The resultant optimization problem is

$$\max_{\boldsymbol{u}} \quad \boldsymbol{u}^T \widetilde{\boldsymbol{\Sigma}} \, \boldsymbol{u}$$

$$\text{s.t.} \quad \|\boldsymbol{u}\|_2^2 = 1$$

- It can be solved by forming the Lagrangian

$$\boldsymbol{u}^T \widetilde{\boldsymbol{\Sigma}} \boldsymbol{u} + \lambda \left(1 - \boldsymbol{u}^T \boldsymbol{u}\right)$$

- By setting the gradient w.r.t. $\boldsymbol{u}$ to zero, we have

$$2\widetilde{\boldsymbol{\Sigma}}\boldsymbol{u} - 2\lambda\boldsymbol{u} = \boldsymbol{0} \quad\longrightarrow\quad \boxed{\widetilde{\boldsymbol{\Sigma}}\boldsymbol{u} = \lambda\boldsymbol{u}}$$

The desired direction $\boldsymbol{u}$ is an eigenvector of $\widetilde{\boldsymbol{\Sigma}}$

$\widetilde{\boldsymbol{\Sigma}}$ has $d$ eigenvectors, which one?

# **Variance Preservation (cont.)**

- Recall that the variance of the projected dataset $\mathcal{D} = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N\}$ is $\boldsymbol{u}^T \widetilde{\boldsymbol{\Sigma}} \, \boldsymbol{u}$

- By substituting $\widetilde{\boldsymbol{\Sigma}} \boldsymbol{u} = \lambda \boldsymbol{u}$ into the above formula, the projected variance becomes

$$\boldsymbol{u}^T \widetilde{\boldsymbol{\Sigma}} \, \boldsymbol{u} = \boldsymbol{u}^T \lambda \boldsymbol{u} \; = \lambda \underbrace{\boldsymbol{u}^T \boldsymbol{u}}_{\|\boldsymbol{u}\|_2^2 \; (\|\boldsymbol{u}\|_2^2 = 1)} = \lambda$$

- To find a direction that maximizes the projected variance is to find the eigenvector $\boldsymbol{u}$ of $\widetilde{\boldsymbol{\Sigma}}$ with the largest eigenvalue

- Generalized to multiple components case: let $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$ be the eigenvalues of $\widetilde{\boldsymbol{\Sigma}}$, and $\boldsymbol{u}_1, \boldsymbol{u}_2, \ldots, \boldsymbol{u}_d$ be the corresponding eigenvectors, and choose the top $k$ eigenvectors as the principal components

18

# **Determine Value of $k$**

- Wrapper approaches
  - Dimensionality reduction is usually an intermediate step for some downstream tasks, such as classification, regression, clustering
  - Use cross-validation based on the performance of the final task to tune the value of $k$

# **Determine Value of $k$ (cont.)**

- Based on the percentage of variance preserved

$$p_{\text{var}} = \frac{\sum_{i=1}^{k} \lambda_i}{\sum_{i=1}^{d} \lambda_i} \times 100$$

  - All the $\lambda_i$'s are nonnegative
  - Predefine a value for the percentage of variance to determine the value of $k$

# Compute Eigenvalues and Eigenvectors

- How to compute eigenvalues and eigenvectors of $\widetilde{\mathbf{\Sigma}} = \frac{1}{N-1}\mathbf{X}^T\mathbf{X}$ ?

- In a general case, if a $d$-by-$d$ square matrix $\mathbf{A}$ can be written as

$$\mathbf{A} = \mathbf{X}^T\mathbf{X}, \text{ where } \mathbf{X} \in \mathbb{R}^{N \times d}$$

  then eigenvectors and eigenvalues of $\mathbf{A} = \mathbf{X}^T\mathbf{X}$ can be computed by performing Singular Value Decomposition (SVD) on $\mathbf{X}$

  - As $\mathbf{A}$ is positive semi-definite, all of its eigenvalues are non-negative.

# Orthogonal Vectors

- Two vectors $\boldsymbol{v}_1$ and $\boldsymbol{v}_2$ are said to be orthogonal if they are perpendicular to each other, i.e., the inner or dot product of two vectors is 0
  - $\boldsymbol{v}_1 \cdot \boldsymbol{v}_2 = 0$
- A set of vectors $\{\boldsymbol{v}_1, \dots, \boldsymbol{v}_d\}$ are mutually orthogonal if every pair of vectors are orthogonal
  - $\boldsymbol{v}_i \cdot \boldsymbol{v}_j = 0$, for any $i \neq j$

$$\boldsymbol{v}_1 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \quad \boldsymbol{v}_2 = \begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad \boldsymbol{v}_3 = \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix}$$

$$\boldsymbol{v}_1 \cdot \boldsymbol{v}_2 = \boldsymbol{v}_1 \cdot \boldsymbol{v}_3 = \boldsymbol{v}_2 \cdot \boldsymbol{v}_3 = 0$$

# Orthonormal Vectors

- A set of vectors $\{v_1, \ldots, v_d\}$ are mutually orthonormal if every pair of vectors are orthogonal, and the $L_2$ norm of each vector is 1

    - $\boldsymbol{v}_i \cdot \boldsymbol{v}_j = 0$, for any $i \neq j$
    - $\|\boldsymbol{v}_i\|_2 = \sqrt{\boldsymbol{v}_i \cdot \boldsymbol{v}_i} = 1$

- A set of orthogonal vectors $\{\boldsymbol{v}_1, \ldots, \boldsymbol{v}_d\}$ can be normalized to orthonormal via $\left\{\dfrac{\boldsymbol{v}_1}{\|\boldsymbol{v}_1\|_2}, \ldots, \dfrac{\boldsymbol{v}_d}{\|\boldsymbol{v}_d\|_2}\right\}$

$$\boldsymbol{v}_1 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \quad \|\boldsymbol{v}_1\|_2 = \sqrt{2} \qquad \boldsymbol{v}_2 = \begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad \|\boldsymbol{v}_2\|_2 = 2 \qquad \boldsymbol{v}_3 = \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix} \quad \|\boldsymbol{v}_3\|_2 = 2$$

$$\boldsymbol{v}_1' = \frac{1}{\sqrt{2}}\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \quad \|\boldsymbol{v}_1'\|_2 = 1 \qquad \boldsymbol{v}_2' = \frac{1}{2}\begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad \|\boldsymbol{v}_2'\|_2 = 1 \qquad \boldsymbol{v}_3' = \frac{1}{2}\begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix} \quad \|\boldsymbol{v}_3'\|_2 = 1$$

# **Orthonormal Vectors (cont.)**

- Given a matrix $\mathbf{V} = (\boldsymbol{v}_1, \ldots, \boldsymbol{v}_d)$, where $\boldsymbol{v}_i$ is an $N$-dimensional column vector, and $N \geq d$

- If the columns of $\mathbf{V}$ are mutually orthonormal, then we have

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}_d$$

$$\mathbf{I}_d = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$
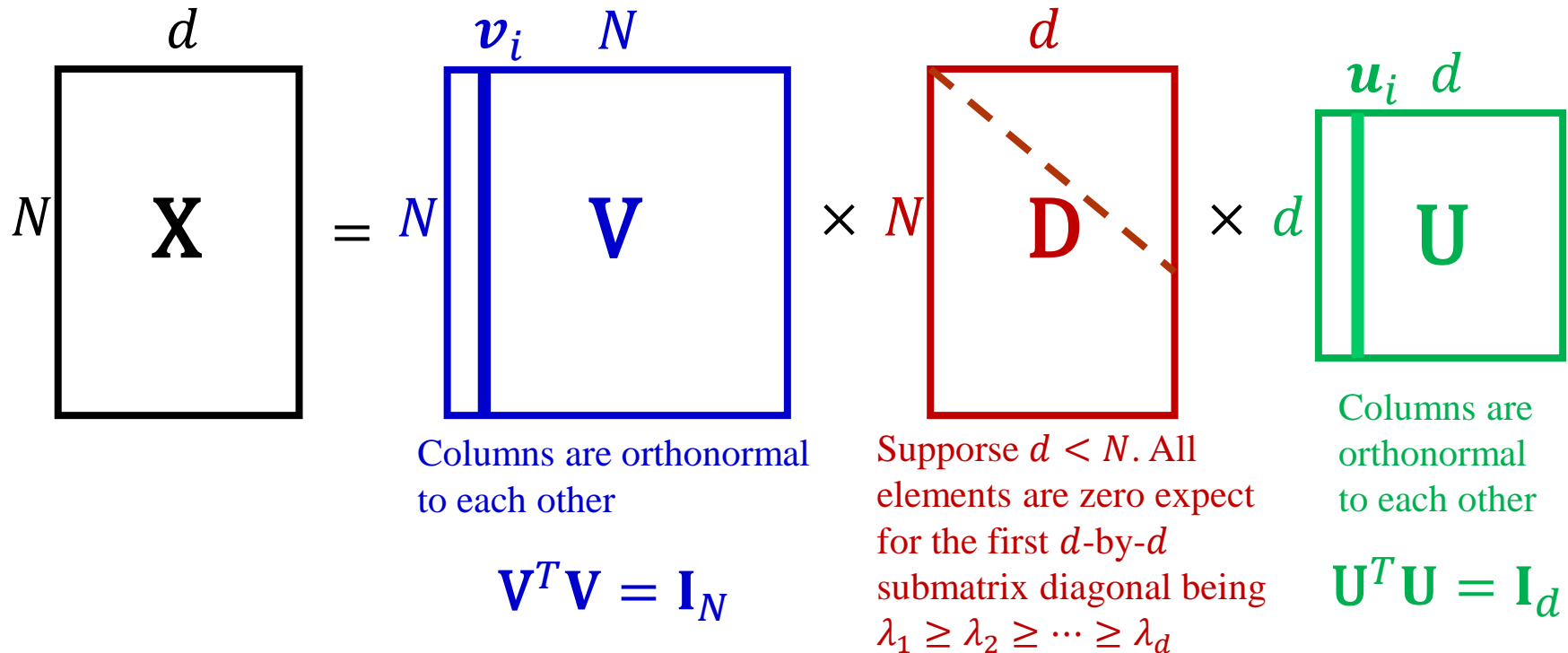
# Singular Value Decomposition (SVD)

- The SVD of $\mathbf{X}$ ($N$-by-$d$) has the following form

$$\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^{T}$$

$\mathbf{X}$ ($N \times d$) = $\mathbf{V}$ ($N \times N$, columns $\boldsymbol{v}_i$) $\times$ $\mathbf{D}$ ($N \times d$) $\times$ $\mathbf{U}$ ($d \times d$, columns $\boldsymbol{u}_i$)

Columns are orthonormal to each other

$$\mathbf{V}^T\mathbf{V} = \mathbf{I}_N$$

Supporse $d < N$. All elements are zero expect for the first $d$-by-$d$ submatrix diagonal being $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$

Columns are orthonormal to each other

$$\mathbf{U}^T\mathbf{U} = \mathbf{I}_d$$

# Obtain Eigenvectors via SVD

- Perform SVD on $\mathbf{X}$ to get $\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^T$
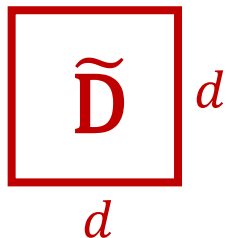
- Then $\mathbf{A}$ can be rewritten as

$$\mathbf{V}^T\mathbf{V} = \mathbf{I}_N$$

$$\mathbf{A} = \mathbf{X}^T\mathbf{X} = \left(\mathbf{V}\mathbf{D}\mathbf{U}^T\right)^T \mathbf{V}\mathbf{D}\mathbf{U}^T = \mathbf{U}\mathbf{D}^T\boxed{\mathbf{V}^T\mathbf{V}}\mathbf{D}\mathbf{U}^T$$
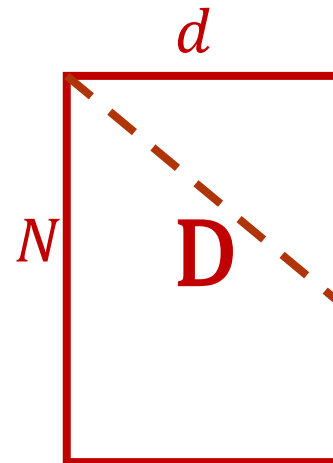
$$\mathbf{A} = \mathbf{U}\mathbf{D}^T\mathbf{D}\mathbf{U}^T$$

Denote $\widetilde{\mathbf{D}} = \mathbf{D}^T\mathbf{D}$

$$= \mathbf{U}\widetilde{\mathbf{D}}\mathbf{U}^T$$

$\widetilde{\mathbf{D}}$ $d$-by-$d$ diagonal matrix with diagonal elements $\lambda_1^2 \geq \lambda_2^2 \geq \cdots \geq \lambda_d^2 \geq 0$

$\mathbf{D}$ — $d$ by $N$

Supporse $d < N$. All elements are zero expect for the first $d$-by-$d$ submatrix diagonal being $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$

# Eigen Components via SVD (cont.)

$$\mathbf{A} = \mathbf{U}\widetilde{\mathbf{D}}\mathbf{U}^T$$

$\boxed{\widetilde{\mathbf{D}}}$ $d$   Diagonal matrix with diagonal elements $\lambda_1^2 \geq \lambda_2^2 \geq \cdots \geq \lambda_d^2 \geq 0$

$d$

$$\mathbf{U}^T\mathbf{U} = \mathbf{I}_d$$

$$\boxed{\mathbf{AU}} = \mathbf{U}\widetilde{\mathbf{D}}\,\cancel{\mathbf{U}^T\mathbf{U}} = \boxed{\mathbf{U}\widetilde{\mathbf{D}}}$$

$$(\mathbf{A} \times \boldsymbol{u}_1, \mathbf{A} \times \boldsymbol{u}_2, \dots, \mathbf{A} \times \boldsymbol{u}_d) \quad = \quad [\lambda_1^2 \times \boldsymbol{u}_1, \lambda_2^2 \times \boldsymbol{u}_2, \dots, \lambda_d^2 \times \boldsymbol{u}_d]$$

$$\mathbf{A}\boldsymbol{u}_i = \lambda_i^2 \boldsymbol{u}_i, i = 1, \dots, d$$

Each column $\boldsymbol{u}_i$ of $\mathbf{U}$ is an eigenvector of $\mathbf{A}$ with the eigenvalue $\lambda_i^2$

# Reference (Optional)

- For feature subset selection:
  - <u>An Introduction to Variable and Feature Selection</u>, Isabelle Guyon, Andre Elisseeff, in JMLR 2003

- For dimensionality reduction:
  - <u>Dimensionality Reduction: A Comparative Review,</u> L.J.P. van der Maaten and E. O. Postma and H. J. van den Herik, Technical Report, 2008
  - https://lvdmaaten.github.io/drtoolbox/

# Thank you!

# Derivation of PCA

- The variance preservation view
  - The first $k$ components display as much as possible of the variation among data instances

- The minimum reconstruction view
  - The first $k$ components convey maximum useful information of original data instances

# Minimum Reconstruction Error

- Given any <u>orthonormal</u> basis $\boldsymbol{v}_1, \boldsymbol{v}_2, \ldots, \boldsymbol{v}_d$, a data point $\boldsymbol{x}_i$ (has been centered) can be written as

$$\boldsymbol{x}_i = \sum_{j=1}^{d} \alpha_{ij} \boldsymbol{v}_j \qquad \alpha_{ij} = \boldsymbol{v}_j^T \boldsymbol{x}_i \qquad \boxed{\sum_{j=1}^{d} \boldsymbol{v}_j^T \boldsymbol{x}_i \boldsymbol{v}_j = \boldsymbol{x}_i \sum_{j=1}^{d} \boldsymbol{v}_j^T \boldsymbol{v}_j = \boldsymbol{x}_i}$$

- Consider the $k$-term approximation of $\boldsymbol{x}_i$:

$$\hat{\boldsymbol{x}}_i \approx \sum_{j=1}^{k} \alpha_{ij} \boldsymbol{v}_j$$

- The error of the approximate over all data points is

$$E = \frac{1}{N} \sum_{i=1}^{N} \|\hat{\boldsymbol{x}}_i - \boldsymbol{x}_i\|_2^2 = \frac{1}{N} \sum_{i=1}^{N} \left\| \sum_{j=k+1}^{d} \alpha_{ij} \boldsymbol{v}_j \right\|_2^2 = \frac{1}{N} \sum_{i=1}^{N} \sum_{j=k+1}^{d} \alpha_{ij}^2$$

# Minimum Reconstruction Error (cont.)

- The error of the approximate over all data points

$$E = \frac{1}{N}\sum_{i=1}^{N}\|\widehat{x}_i - x_i\|_2^2 = \frac{1}{N}\sum_{i=1}^{N}\sum_{j=k+1}^{d}\alpha_{ij}^2$$

$$= \frac{1}{N}\sum_{i=1}^{N}\sum_{j=k+1}^{d}v_j^T x_i x_i^T v_j \approx \sum_{j=k+1}^{d}v_j^T\widetilde{\Sigma}v_j$$

- Suppose $k = d - 1$, i.e., we aim to remove a single dimension, then resultant optimization problem is

$$\min_{v_d}\quad v_d^T\widetilde{\Sigma}\,v_d$$

$$\text{s.t.}\quad \|v_d\|_2^2 = 1$$

32

# Minimum Reconstruction Error (cont.)

- By setting the gradient of the Lagrangian w.r.t. $\boldsymbol{v}$ to zero, we have

$$2\widetilde{\boldsymbol{\Sigma}}\boldsymbol{v}_d - 2\lambda\boldsymbol{v}_d = \boldsymbol{0} \longrightarrow \boxed{\widetilde{\boldsymbol{\Sigma}}\boldsymbol{v}_d = \lambda\boldsymbol{v}_d}$$

  The desired direction $\boldsymbol{v}_d$ is an eigenvector of $\widetilde{\boldsymbol{\Sigma}}$

  $\widetilde{\boldsymbol{\Sigma}}$ has $d$ eigenvectors, which one?

- Our goal is to minimize the reconstruction error $\boldsymbol{v}_{\boldsymbol{d}}^T\widetilde{\boldsymbol{\Sigma}}\,\boldsymbol{v}_d$

$$\boldsymbol{v}_{\boldsymbol{d}}^T\widetilde{\boldsymbol{\Sigma}}\,\boldsymbol{v}_d = \boldsymbol{v}_{\boldsymbol{d}}^T\lambda\boldsymbol{v}_d = \lambda\boldsymbol{v}_{\boldsymbol{d}}^T\boldsymbol{v}_d = \lambda$$

- Therefore, $\boldsymbol{v}_d$ should be the eigenvector $\boldsymbol{u}_d$ of $\widetilde{\boldsymbol{\Sigma}}$ with the smallest eigenvalue because $\boldsymbol{u}_d^T\widetilde{\boldsymbol{\Sigma}}\boldsymbol{u}_d = \lambda_d$

- Similarly, the other dimensions to remove are subsequently the eigenvectors corresponding to the least eigenvalues