

CZ4041/SC4000: Machine Learning

Lesson 12: Dimensionality Reduction

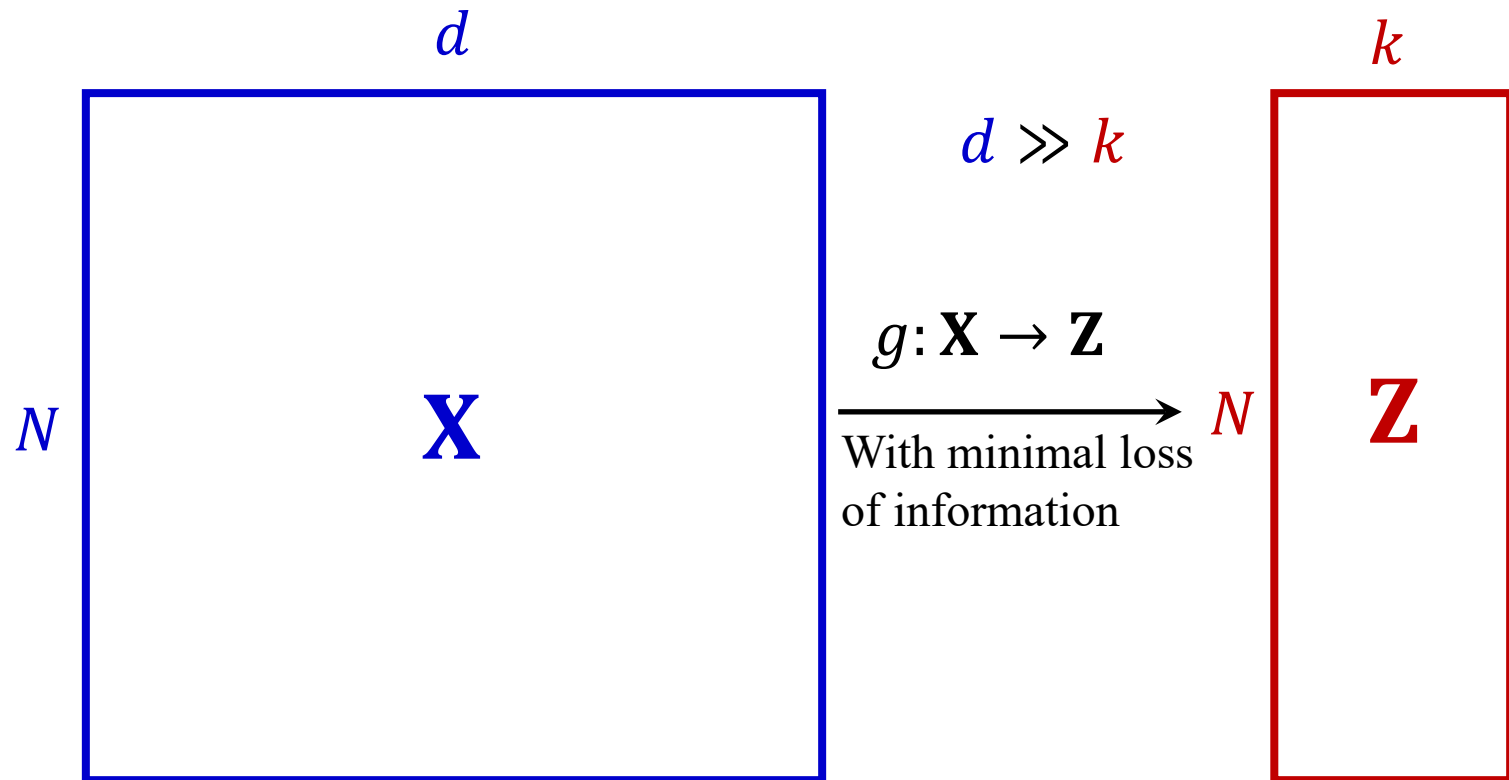
Li Boyang, Albert

School of Computer Science and Engineering,
NTU, Singapore

Acknowledgements: some content is adapted from the lecture notes of Xiaojin Zhu @University of Wisconsin–Madison. Slides are modified from the version prepared by Dr. Sinno Pan.

High-level Idea

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



Why Dimensionality Reduction

- To avoid curse of dimensionality
 - Decision boundary in SVM: $\mathbf{w} \cdot \mathbf{x} + b = 0$
 - Linear Regression: $f(x) = \mathbf{w} \cdot \mathbf{x}$
 - Perceptron: $f(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x})$
 - One parameter to learn for every input dimension
 - Difficult to accurately estimate the best parameters when the number of data points is small

Why Dimensionality Reduction

- To identify the features, or the transformations of features that capture the most important data characteristics
- All measurements contain error or noise.

Mathematically we may write

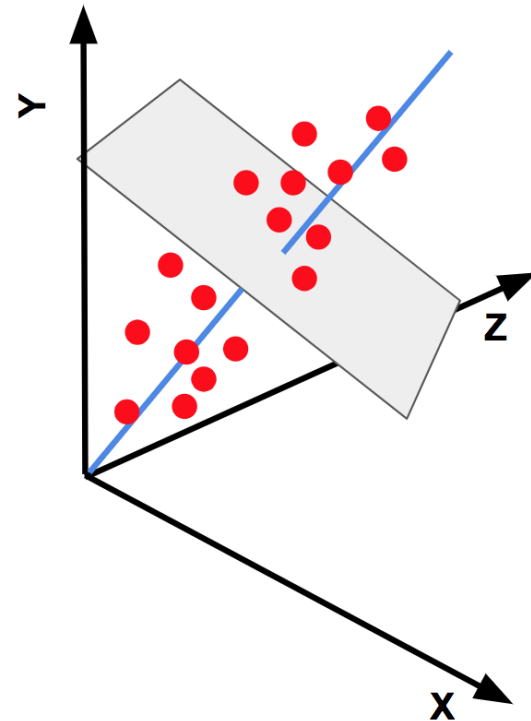
$$\mathbf{x} = \tilde{\mathbf{x}} + \mathbf{e}, \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma I)$$

Why Dimensionality Reduction

- All measurements contain error or noise

$$\mathbf{x} = \tilde{\mathbf{x}} + \mathbf{e}, \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

- Valid hypothesis: Data reside on a straight line, but the noise is isotropic (equal amount of variation) in all three dimensions.



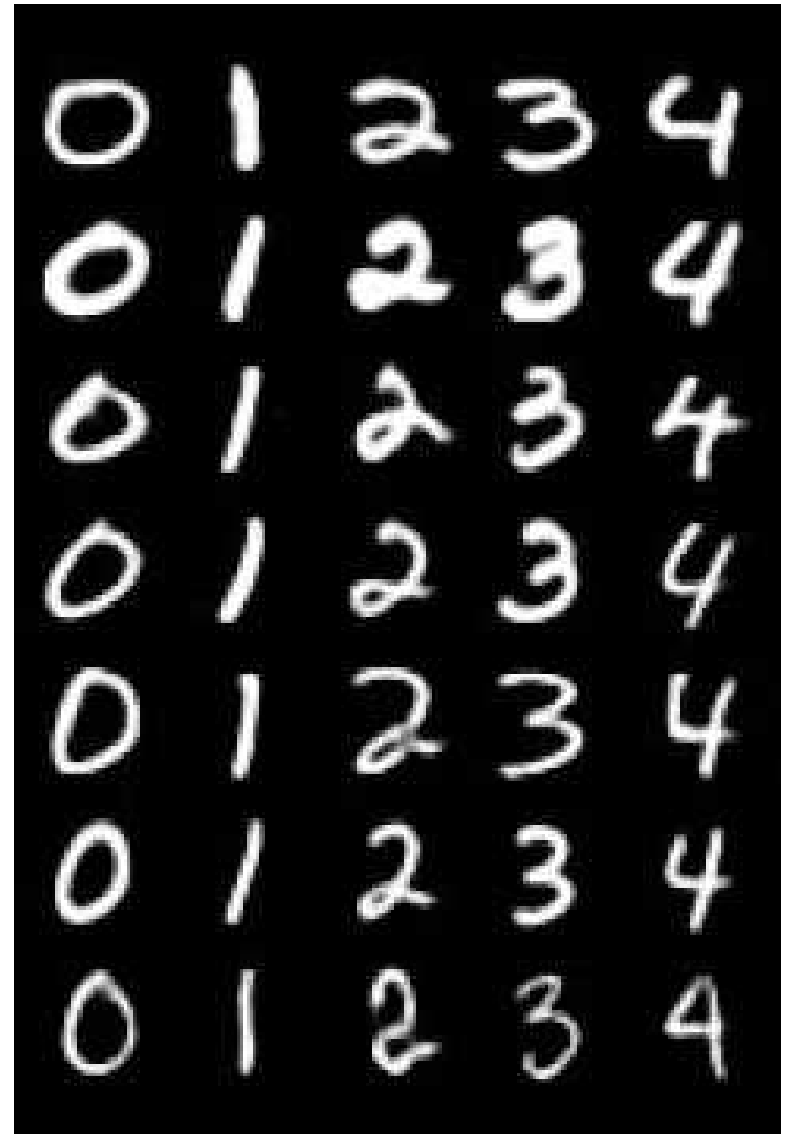
What is noise?

- These photos depict the same person
- Different lighting, make-up, facial hair, expressions, etc., create visual differences
- The data variation that we care about is far smaller than the pixel-level differences.



What is noise?

- The MNIST dataset
- Each column shows the same number.
- Some noise, which we do not have names for, change how they look.
- The data variation that we care about is far smaller than the pixel-level differences

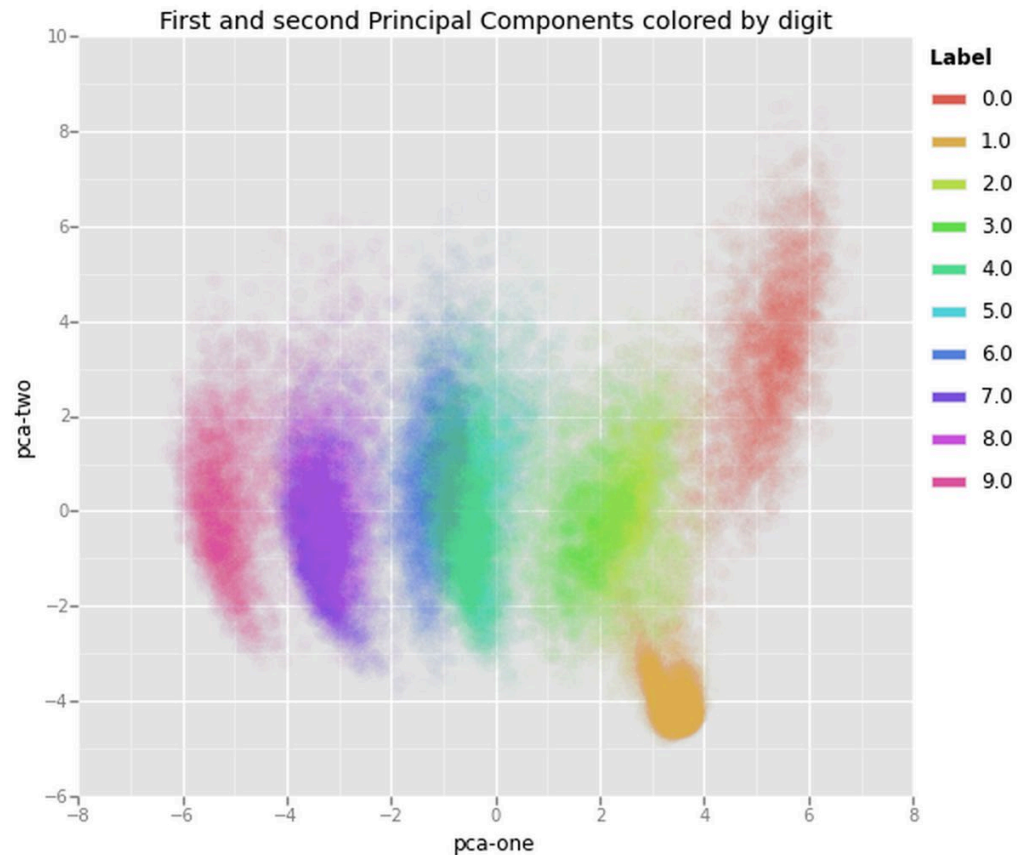


Why Dimensionality Reduction

- Thus, we should identify important variations in the data and discard the noise.
- Benefits
 - Reduces storage requirements
 - Allows visualization in 2D or 3D
 - Reduces noise and improve the performance of machine learning

A Case Study: MNIST

- We use PCA to reduce the number of dimensions to 2 and visualize the results.
- Some clustering structure



Dimensionality Reduction Approaches

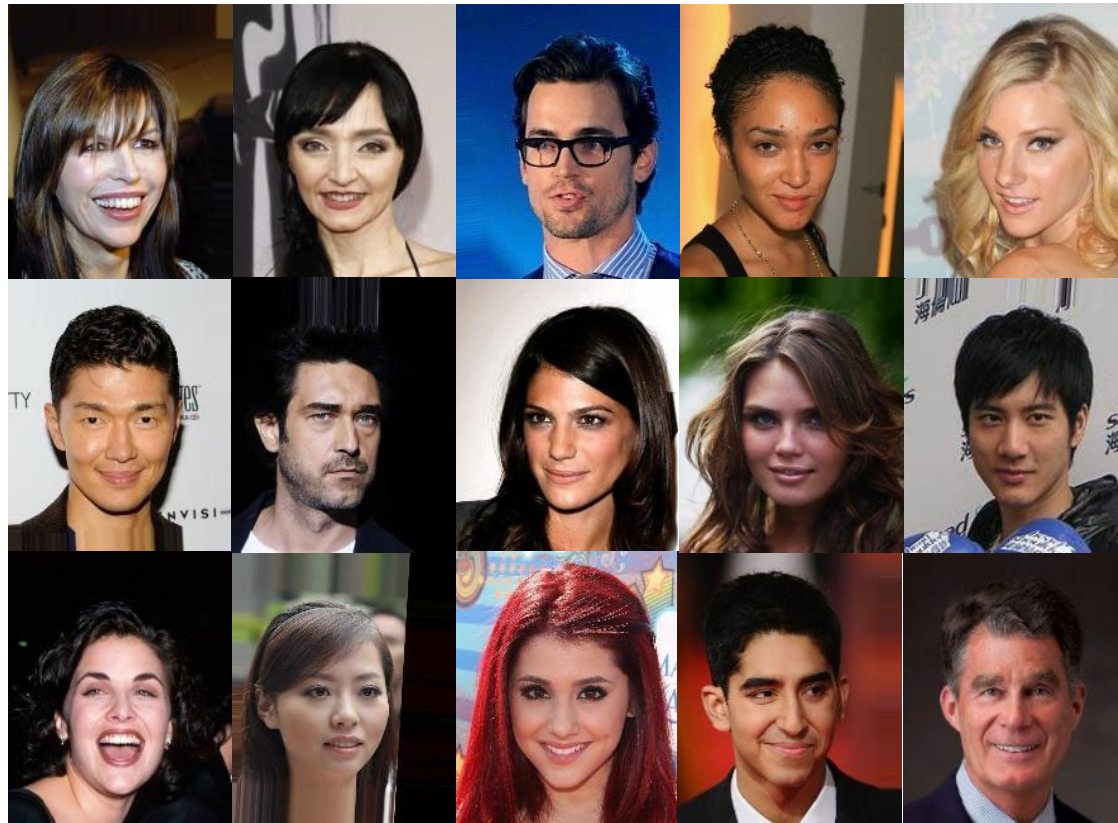
- Feature Selection
 - To **select** a subset of k features from the original d features to represent each data instance
 - Brute-force approach
 - Greedy search
- Feature Extraction
 - To **learn** k new features from the original d features to represent each data instance
 - Linear combination of original 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

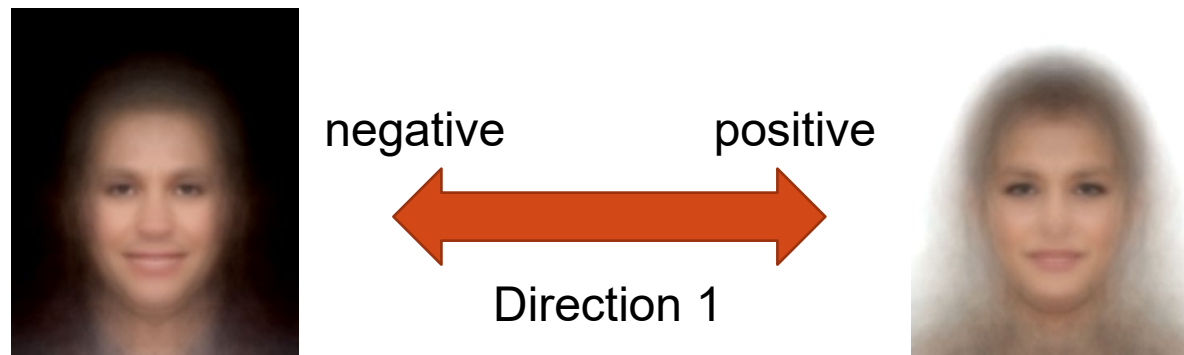
A Case Study: Eigenface

- PCA on images of human faces
- A technique initially created for face identification.

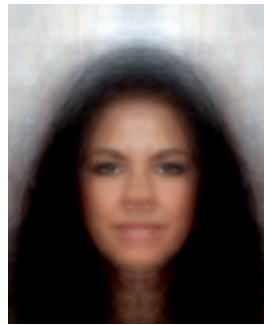


A Case Study: Eigenface

- We use PCA to identify 10 directions that capture most variations in the data.
- We visualize them by adding each direction to the average face, without which the faces are hard to interpret and look scary.
- The first direction seems to be mostly about lighting and background.



Long dark hair



negative

positive



Direction 2



Right facing

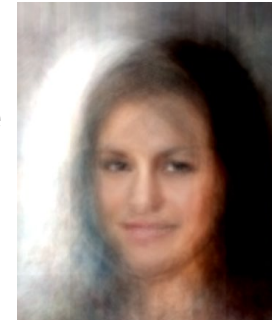


negative

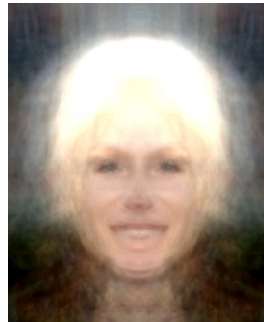
positive



Direction 3



Light-colored
hair or bald

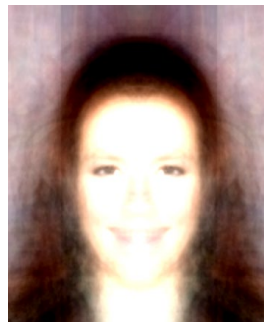


negative

positive



Direction 4



negative

positive



Direction 5



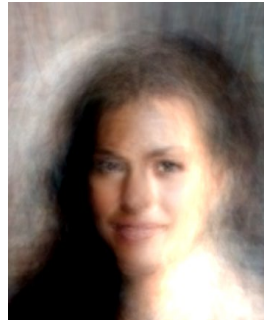
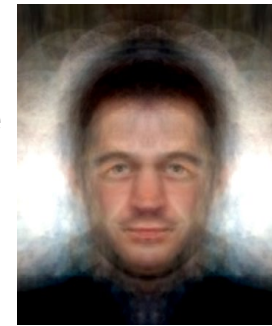


negative

positive



Direction 6

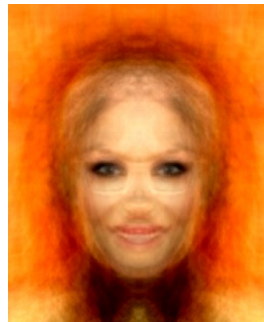
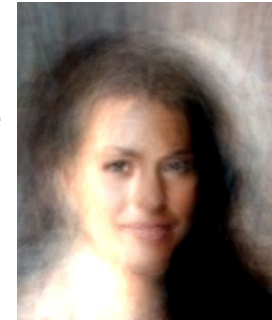


negative

positive



Direction 7

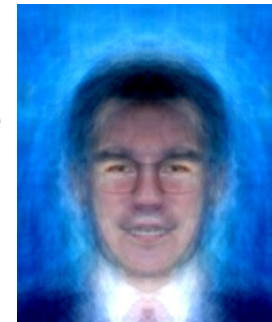


negative

positive



Direction 8



Background +
suit?

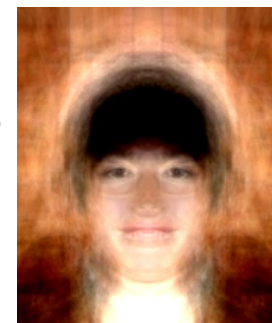


negative

positive



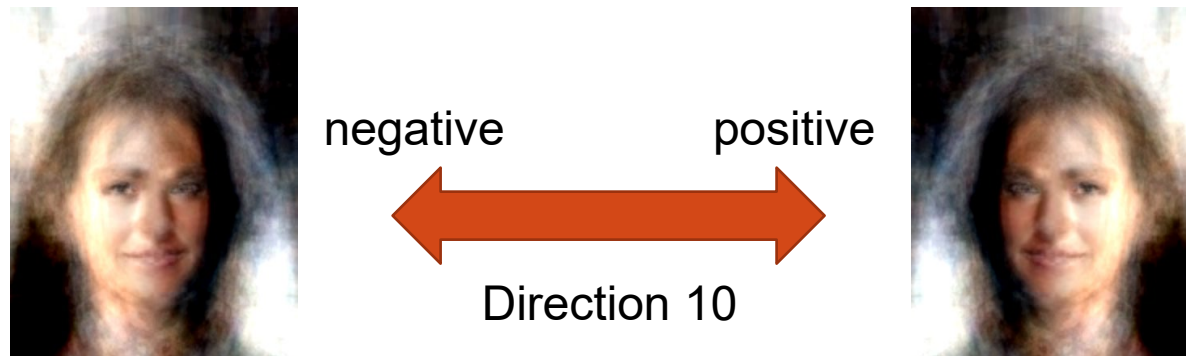
Direction 9



Light above
eyebrows or
below chin

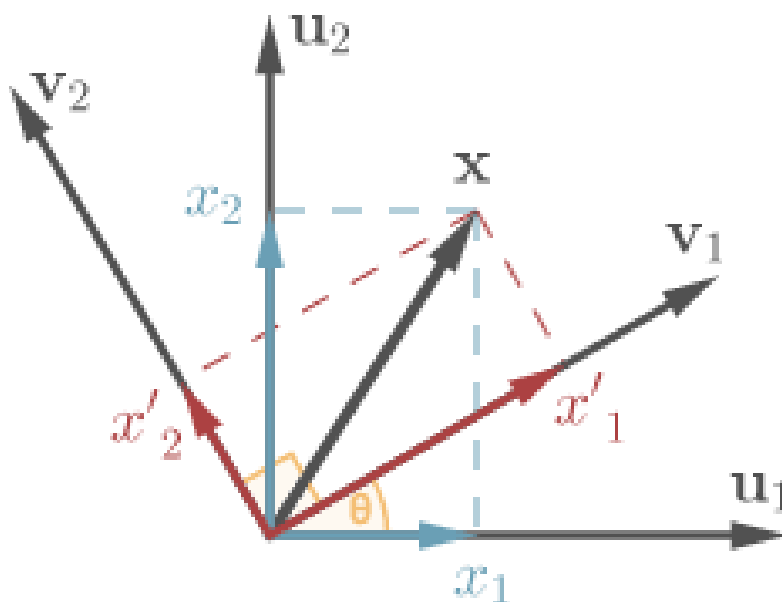
A Case Study: Eigenface

- PCA directions capture variations in data, which correspond to some semantic categories.
- But they are not always interpretable.
- PCA is linear. It would struggle to capture semantic categories that require non-linear variations, such as formal vs. casual dress.



Linear Algebra: Change of Basis

- Vector $\mathbf{x} = (x_1, x_2) = x_1\mathbf{u}_1 + x_2\mathbf{u}_2$ where $\mathbf{u}_1 = (1, 0)$ and $\mathbf{u}_2 = (0, 1)$ are the basis vectors.
- Can we change the basis vectors to \mathbf{v}_1 and \mathbf{v}_2 ?
- $\mathbf{x} = x'_1\mathbf{v}_1 + x'_2\mathbf{v}_2$. We just need to find x'_1 and x'_2 by projecting \mathbf{x} onto \mathbf{v}_1 and \mathbf{v}_2 .

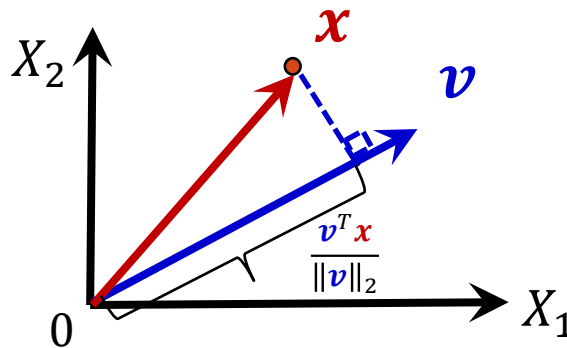


Linear Algebra: Change of Basis

- Consider a projection of a data point \mathbf{x} onto a vector \mathbf{v}
- The projection of \mathbf{x} onto \mathbf{v} is

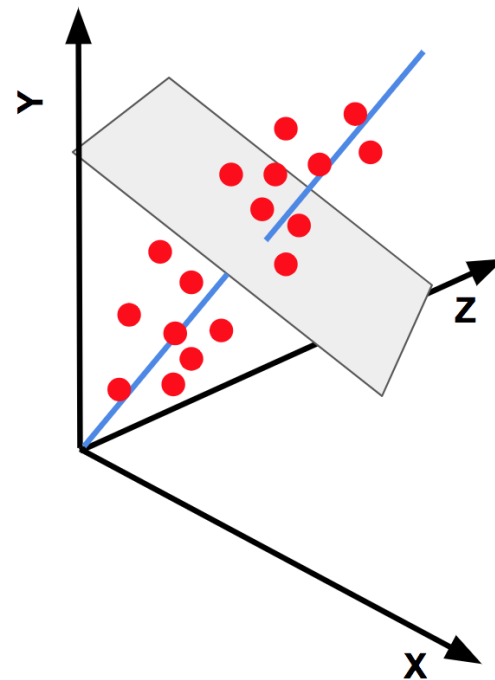
$$\frac{\mathbf{v}^T \mathbf{x}}{\|\mathbf{v}\|_2} = \frac{\|\mathbf{v}\|_2 \|\mathbf{x}\|_2 \cos(\theta)}{\|\mathbf{v}\|_2} = \|\mathbf{x}\|_2 \cos(\theta)$$

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



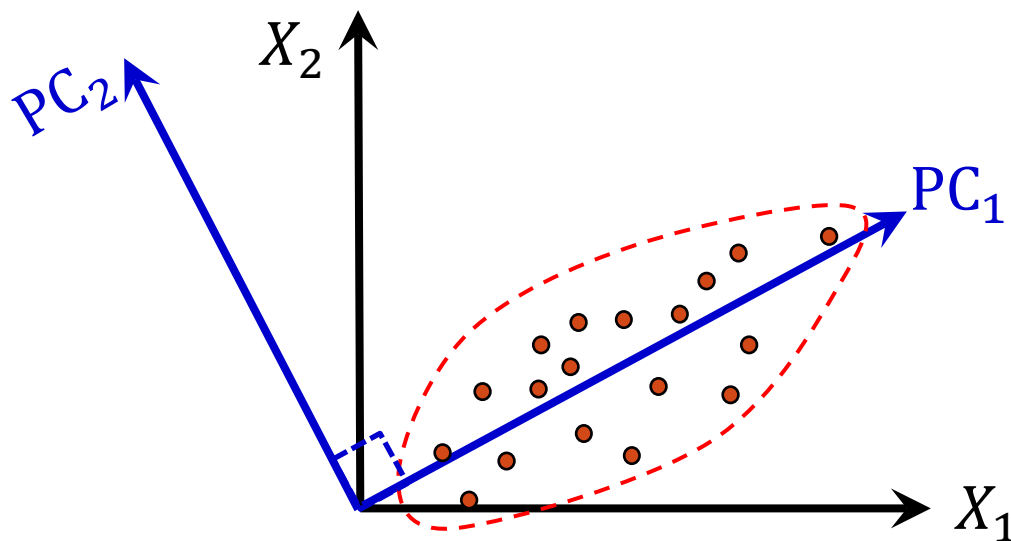
PCA: Geometric Rationale

- Data may only occupy a small subspace of the high-dimensional \mathbb{R}^d space.
- We want to pick a few ($< d$) basis vectors that data are projected onto.
- Which basis vectors would you pick?



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: $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ a set of observed data

1. Centering the data points s.t. the mean is $\mathbf{0}$

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

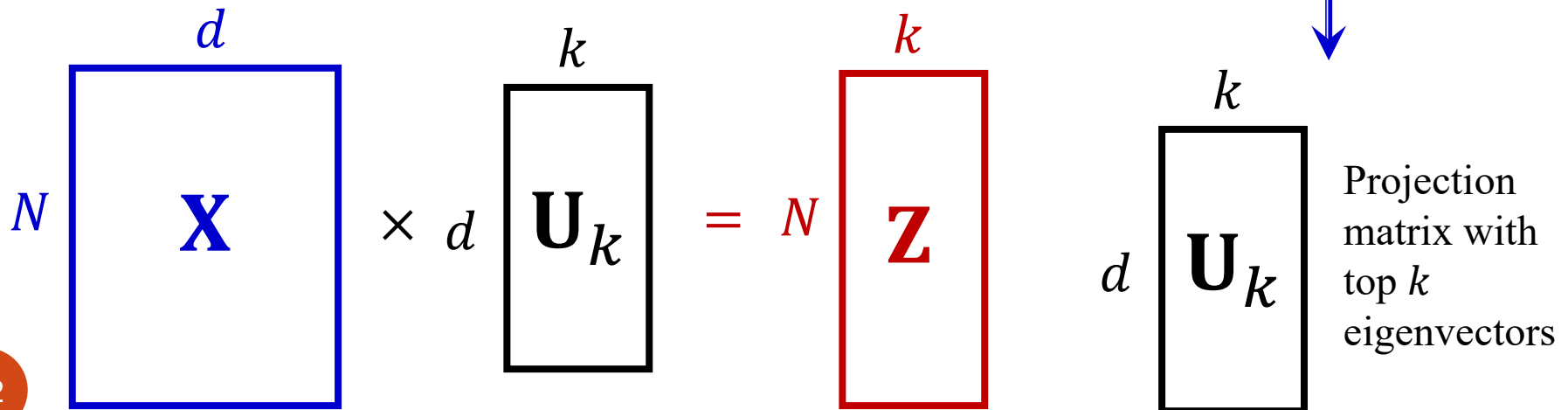
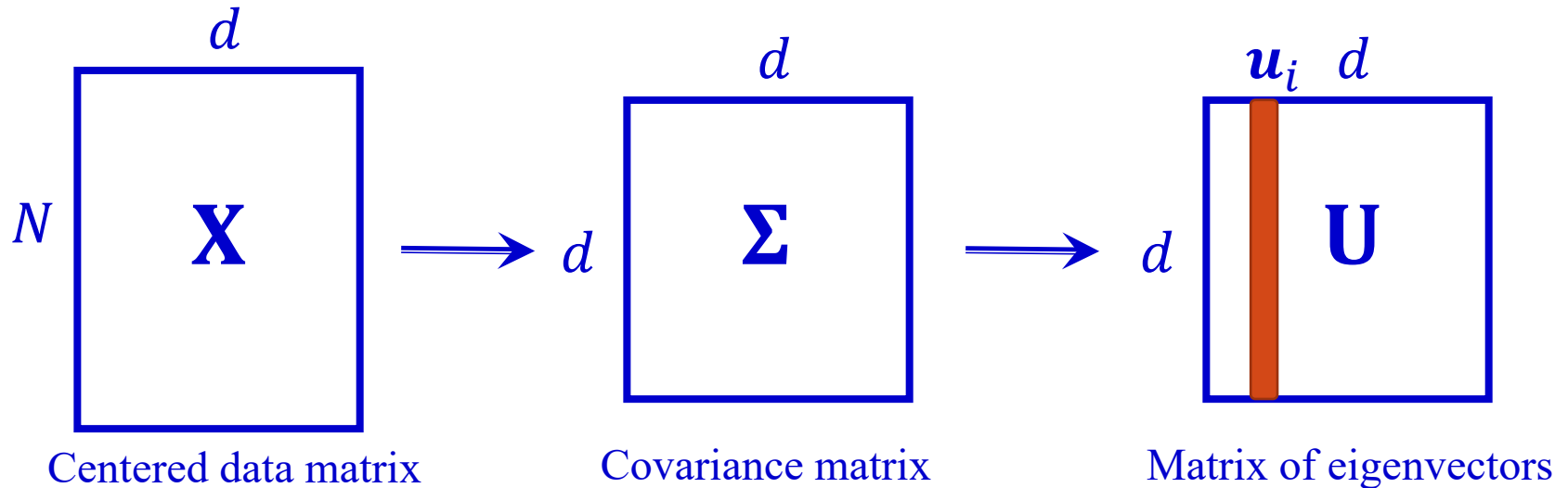
2. Compute sample covariance matrix

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

Each \mathbf{u}_i is of d dimensions

3. Compute eigenvectors of $\tilde{\boldsymbol{\Sigma}}$, $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d\}$, which are sorted based on their eigenvalues in non-increasing order, i.e., $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$
4. Select the first k eigenvectors to construct principal components

PCA: Algorithm (Illustration)



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 \mathbf{u} , s.t.

$$\mathbf{A}\mathbf{u} = \lambda\mathbf{u}$$

scalar



then \mathbf{u} is an eigenvector of \mathbf{A} , and λ 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
 - For symmetric matrices, 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:

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

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

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

Properties of Eigenvalues (cont.)

- Recall: when inducing a closed form solution of regularized linear regression model, we mentioned that if a matrix \mathbf{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 \mathbf{A} is always invertible:

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



Properties of Eigenvalues (cont.)

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

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

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

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

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

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

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

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

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

Properties of Eigenvalues (cont.)

- As \mathbf{A} is positive definite, all of its eigenvalues are positive, i.e., non-zero
- Recall: \mathbf{A} is invertible if all the eigenvalues of \mathbf{A} are non-zero (either positive or negative)
- Therefore, if a matrix \mathbf{A} can be written as

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

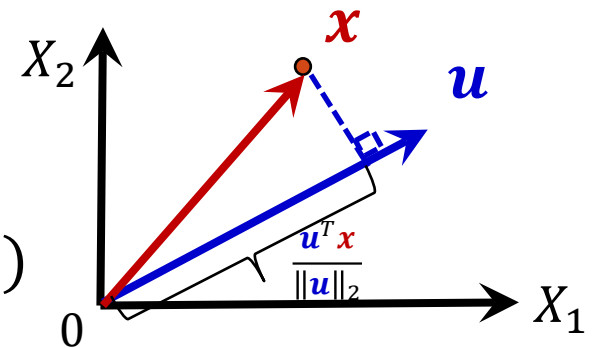
then \mathbf{A} is invertible!

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 \mathbf{x} onto a vector going through the origin, represented by \mathbf{u}

- The projection of \mathbf{x} onto \mathbf{u} is

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



- For simplicity, consider \mathbf{u} with unit length, i.e., $\|\mathbf{u}\|_2 = 1$
- The projected instances $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ onto \mathbf{u} are

$$\{\mathbf{u}^T \mathbf{x}_1, \mathbf{u}^T \mathbf{x}_2, \dots, \mathbf{u}^T \mathbf{x}_N\}$$

Variance Preservation (cont.)

- In PCA, data points are centered at the beginning

$$\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = 0$$

- After projection onto \mathbf{u} , the mean of data points is still 0

$$\frac{1}{N} \sum_{i=1}^N \mathbf{u}^T \mathbf{x}_i = \mathbf{u}^T \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = 0$$

- The variance of the data points projected onto \mathbf{u} is

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

$$= \mathbf{u}^T \tilde{\Sigma} \mathbf{u} \quad \tilde{\Sigma} = \frac{1}{N-1} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$$

Each row is a data instance

Variance Preservation (cont.)

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

$$\begin{aligned} \max_{\mathbf{u}} \quad & \mathbf{u}^T \tilde{\Sigma} \mathbf{u} \\ \text{s.t.} \quad & \|\mathbf{u}\|_2^2 = 1 \end{aligned}$$

- It can be solved by forming the Lagrangian

$$\mathbf{u}^T \tilde{\Sigma} \mathbf{u} + \lambda(1 - \mathbf{u}^T \mathbf{u})$$

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

$$2\tilde{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = \mathbf{0} \quad \longrightarrow \quad \boxed{\tilde{\Sigma}\mathbf{u} = \lambda\mathbf{u}} \quad \begin{array}{l} \text{The desired direction } \mathbf{u} \\ \text{is an eigenvector of } \tilde{\Sigma} \end{array}$$

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



Variance Preservation (cont.)

- Recall that the variance of the projected dataset $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ is $\mathbf{u}^T \tilde{\Sigma} \mathbf{u}$
- By substituting $\tilde{\Sigma} \mathbf{u} = \lambda \mathbf{u}$ into the above formula, the projected variance becomes

$$\mathbf{u}^T \tilde{\Sigma} \mathbf{u} = \mathbf{u}^T \lambda \mathbf{u} = \lambda \mathbf{u}^T \mathbf{u} = \lambda \|\mathbf{u}\|_2^2 \quad (\|\mathbf{u}\|_2^2 = 1)$$

- To find a direction that maximizes the projected variance is to find the eigenvector \mathbf{u} of $\tilde{\Sigma}$ with the largest eigenvalue
- Generalized to multiple components case: let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ be the eigenvalues of $\tilde{\Sigma}$, and $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d$ be the corresponding eigenvectors, and choose the top k eigenvectors as the principal components

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 λ_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 $\tilde{\Sigma} = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$?
- 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}

Orthogonal Vectors

- Two vectors \mathbf{v}_1 and \mathbf{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
 - $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$
- A set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_d\}$ are mutually orthogonal if every pair of vectors are orthogonal
 - $\mathbf{v}_i \cdot \mathbf{v}_j = 0$, for any $i \neq j$

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

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

Orthonormal Vectors

- A set of vectors $\{v_1, \dots, v_d\}$ are mutually orthonormal if every pair of vectors are orthogonal, and the L_2 norm of each vector is 1
 - $v_i \cdot v_j = 0$, for any $i \neq j$
 - $\|v_i\|_2 = \sqrt{v_i \cdot v_i} = 1$
- A set of orthogonal vectors $\{v_1, \dots, v_d\}$ can be normalized to orthonormal via $\left\{ \frac{v_1}{\|v_1\|_2}, \dots, \frac{v_d}{\|v_d\|_2} \right\}$

$$v_1 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \quad \|v_1\|_2 = \sqrt{2} \quad v_2 = \begin{pmatrix} 1 \\ \sqrt{2} \\ 1 \end{pmatrix} \quad \|v_2\|_2 = 2 \quad v_3 = \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix} \quad \|v_3\|_2 = 2$$

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

Orthonormal Vectors (cont.)

- Given a matrix $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_d)$, where \mathbf{v}_i is an N -dimensional column vector, and $N \geq d$
- If the columns of \mathbf{V} are mutually orthonormal, then we have

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

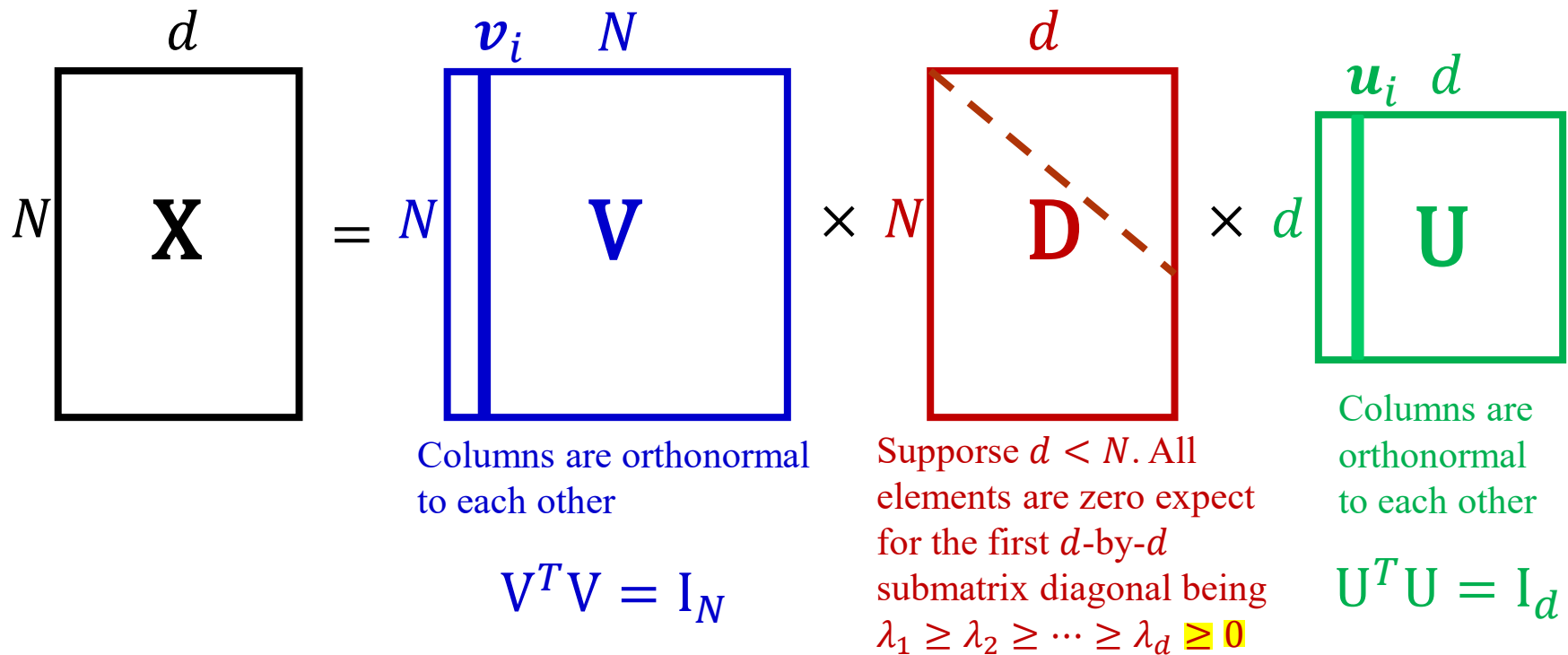
$$\mathbf{I}_d = \left(\begin{array}{cccc} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 1 \end{array} \right) \left. \vphantom{\begin{array}{cccc} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 1 \end{array}} \right\} d$$

d

Singular Value Decomposition (SVD)

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

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



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

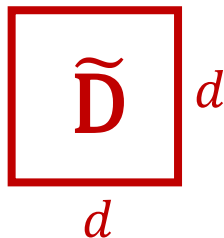
$$\mathbf{A} = \mathbf{X}^T \mathbf{X} = (\mathbf{V}\mathbf{D}\mathbf{U}^T)^T \mathbf{V}\mathbf{D}\mathbf{U}^T = \mathbf{U}\mathbf{D}^T \underbrace{\mathbf{V}^T \mathbf{V}}_{=\mathbf{I}_N} \mathbf{D}\mathbf{U}^T$$



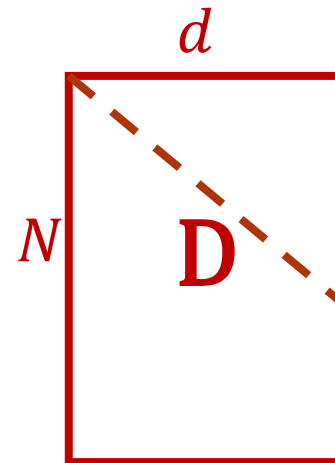
$$\mathbf{A} = \mathbf{U}\mathbf{D}^T \mathbf{D}\mathbf{U}^T$$

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

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



d -by- d diagonal matrix with diagonal elements $\lambda_1^2 \geq \lambda_2^2 \geq \dots \geq \lambda_d^2 \geq 0$



Suppose $d < N$. All elements are zero except for the first d -by- d submatrix diagonal being $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

Eigen Components via SVD (cont.)

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

$$\tilde{\mathbf{D}}$$

Diagonal matrix with diagonal elements $\lambda_1^2 \geq \lambda_2^2 \geq \dots \geq \lambda_d^2 \geq 0$

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

$$\boxed{\mathbf{A}\mathbf{U}} = \mathbf{U}\tilde{\mathbf{D}}\cancel{\mathbf{U}^T\mathbf{U}} = \boxed{\mathbf{U}\tilde{\mathbf{D}}}$$

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

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

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

Reference (Optional)

- For feature subset selection:
 - An Introduction to Variable and Feature Selection, Isabelle Guyon, Andre Elisseeff, in JMLR 2003
- For dimensionality reduction:
 - Dimensionality Reduction: A Comparative Review, 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 orthonormal basis $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$, a data point \mathbf{x}_i (has been centered) can be written as

$$\mathbf{x}_i = \sum_{j=1}^d \alpha_{ij} \mathbf{v}_j \quad \alpha_{ij} = \mathbf{v}_j^T \mathbf{x}_i$$

$$\sum_{j=1}^d \mathbf{v}_j^T \mathbf{x}_i \mathbf{v}_j = \mathbf{x}_i \sum_{j=1}^d \mathbf{v}_j^T \mathbf{v}_j = \mathbf{x}_i$$

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

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

- The error of the approximate over all data points is

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

Minimum Reconstruction Error (cont.)

- The error of the approximate over all data points

$$\begin{aligned}
 E &= \frac{1}{N} \sum_{i=1}^N \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|_2^2 = \frac{1}{N} \sum_{i=1}^N \sum_{j=k+1}^d \alpha_{ij}^2 \\
 &= \frac{1}{N} \sum_{i=1}^N \sum_{j=k+1}^d \mathbf{v}_j^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{v}_j \approx \sum_{j=k+1}^d \mathbf{v}_j^T \tilde{\Sigma} \mathbf{v}_j
 \end{aligned}$$

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

$$\begin{aligned}
 \min_{\mathbf{v}_d} \quad & \mathbf{v}_d^T \tilde{\Sigma} \mathbf{v}_d \\
 \text{s.t.} \quad & \|\mathbf{v}_d\|_2^2 = 1
 \end{aligned}$$

Minimum Reconstruction Error (cont.)

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

$$2\tilde{\Sigma}\mathbf{v}_d - 2\lambda\mathbf{v}_d = \mathbf{0} \longrightarrow \boxed{\tilde{\Sigma}\mathbf{v}_d = \lambda\mathbf{v}_d}$$

The desired direction \mathbf{v}_d is an eigenvector of $\tilde{\Sigma}$

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

- Our goal is to minimize the reconstruction error $\mathbf{v}_d^T \tilde{\Sigma} \mathbf{v}_d$

$$\mathbf{v}_d^T \tilde{\Sigma} \mathbf{v}_d = \mathbf{v}_d^T \lambda \mathbf{v}_d = \lambda \mathbf{v}_d^T \mathbf{v}_d = \lambda$$

- Therefore, \mathbf{v}_d should be the eigenvector \mathbf{u}_d of $\tilde{\Sigma}$ with the smallest eigenvalue because $\mathbf{u}_d^T \tilde{\Sigma} \mathbf{u}_d = \lambda_d$
- Similarly, the other dimensions to remove are subsequently the eigenvectors corresponding to the least eigenvalues