# CS 6316 Machine Learning

## Dimensionality Reduction

Yangfeng Ji

Information and Language Processing Lab
Department of Computer Science
University of Virginia

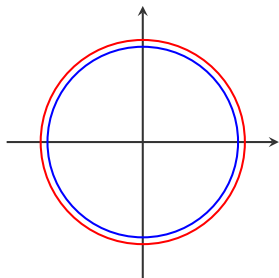## Overview

1. Reducing Dimensions

2. Principal Component Analysis

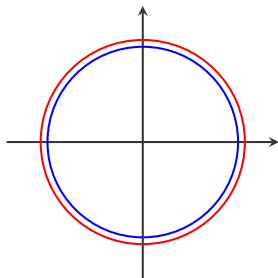3. A Different Viewpoint of PCA

# Reducing Dimensions

What is the volume difference between two $d$-dimensional balls with radii $r_1 = 1$ and $r_2 = 0.99$

What is the volume difference between two $d$-dimensional balls with radii $r_1 = 1$ and $r_2 = 0.99$



- $d = 2$: $\frac{1}{2}\pi(r_1^2 - r_2^2) \approx 0.03$
- $d = 3$: $\frac{4}{3}\pi(r_1^3 - r_2^3) \approx 0.12$
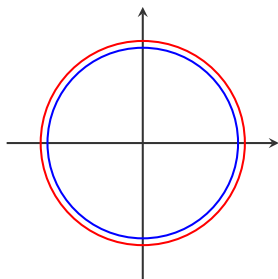
What is the volume difference between two $d$-dimensional balls with radii $r_1 = 1$ and $r_2 = 0.99$



- $d = 2$: $\frac{1}{2}\pi(r_1^2 - r_2^2) \approx 0.03$
- $d = 3$: $\frac{4}{3}\pi(r_1^3 - r_2^3) \approx 0.12$
- General form: $\frac{\pi^{d/2}}{\Gamma(\frac{d}{2}+1)}(r_1^d - r_2^d)$ with $r_2^d \to 0$ when $d \to \infty$
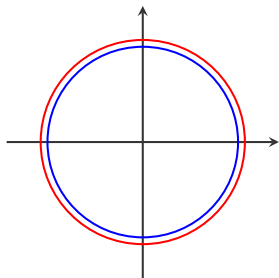  - E.g., $r_2^{500} = 0.00657$

# Curse of Dimensionality

What is the volume difference between two $d$-dimensional balls with radii $r_1 = 1$ and $r_2 = 0.99$



- $d = 2$: $\frac{1}{2}\pi(r_1^2 - r_2^2) \approx 0.03$
- $d = 3$: $\frac{4}{3}\pi(r_1^3 - r_2^3) \approx 0.12$
- General form: $\frac{\pi^{d/2}}{\Gamma(\frac{d}{2}+1)}(r_1^d - r_2^d)$ with $r_2^d \to 0$ when $d \to \infty$
  - E.g., $r_2^{500} = 0.00657$

Question: what will happen if we uniformly sample from a $d$-dimensional ball?

3

If we randomly sample 1K unit vectors from a $d$-dimensional space and calculate the the Euclidean distance between any two vectors, then the distance distribution looks like

If we randomly sample 1K unit vectors from a $d$-dimensional space and calculate the the Euclidean distance between any two vectors, then the distance distribution looks like



Figure: $d = 100$

# Curse of Dimensionality (II)

If we randomly sample 1K unit vectors from a $d$-dimensional space and calculate the the Euclidean distance between any two vectors, then the distance distribution looks like



Figure: $d = 500$

# Curse of Dimensionality (II)

If we randomly sample 1K unit vectors from a $d$-dimensional space and calculate the the Euclidean distance between any two vectors, then the distance distribution looks like
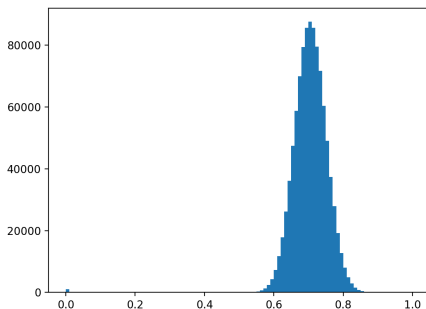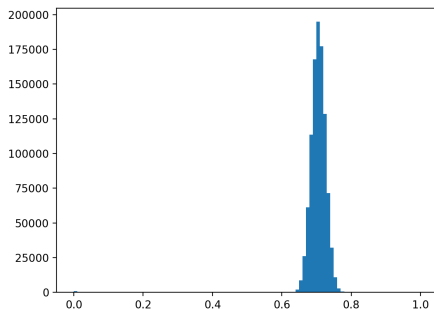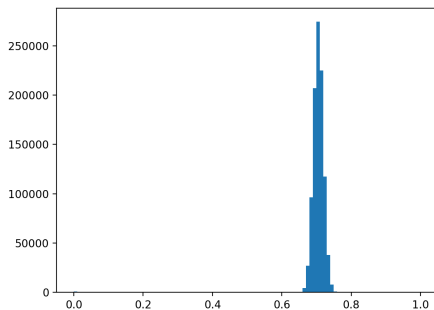


Figure: $d = 1000$

Dimensionality Reduction is the process of taking data in a high dimensional space and mapping it into a new space whose dimensionality is much smaller.

# Dimensionality Reduction

Dimensionality Reduction is the process of taking data in a high dimensional space and mapping it into a new space whose dimensionality is much smaller.
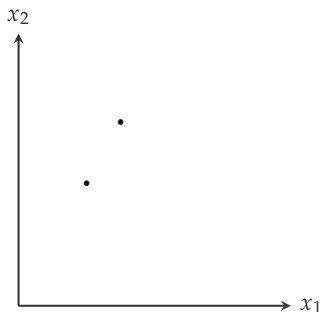
Mathematically, it means

$$f : x \rightarrow \tilde{x} \tag{1}$$

where $x \in \mathbb{R}^d$, $\tilde{x} \in \mathbb{R}^n$ with $n < d$

For the purpose of reducing dimensions, we can project $x = (x_1, x_2)$ into the direction along $x_1$ or $x_2$



Question: Given these two data examples, which direction we should pick? $x_1$ or $x_2$?
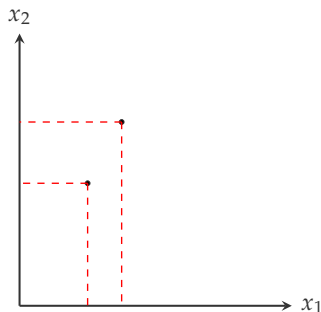
For the purpose of reducing dimensions, we can project $x = (x_1, x_2)$ into the direction along $x_1$ or $x_2$



Question: Given these two data examples, which direction we should pick? $x_1$ or $x_2$?

There is a better solution if we are allowed to rotate the coordinate

There is a better solution if we are allowed to rotate the coordinate



Pick $u_1$, then we preserve all the variance of the examples

# Reducing Dimensions: A toy example (III)

Consider a general case, where the examples do not lie on a perfect line



[Bishop, 2006, Section 12.1]

# Reducing Dimensions: A toy example (III)

Consider a general case, where the examples do not lie on a perfect line



We can follow the same idea by finding a direction that can preserve most of the variance of the examples

[Bishop, 2006, Section 12.1]

# Principal Component Analysis

Given a set of example $S = \{x_1, \ldots, x_m\}$

▶ Centering the data by removing the mean $\bar{x} = \frac{1}{m} \sum_{i=1}^{m} x_i$

$$x_i \leftarrow x_i - \bar{x} \quad \forall i \in [m] \tag{2}$$

Given a set of example $S = \{x_1, \ldots, x_m\}$

▶ Centering the data by removing the mean $\bar{x} = \frac{1}{m} \sum_{i=1}^{m} x_i$

$$x_i \leftarrow x_i - \bar{x} \quad \forall i \in [m] \tag{2}$$

▶ Assume the direction that we would like to project the data is $u$, then the objective function is the data variance

$$J(u) = \frac{1}{m} \sum_{i=1}^{m} (u^\top x_i)^2 \tag{3}$$

# Formulation

Given a set of example $S = \{x_1, \ldots, x_m\}$

▶ Centering the data by removing the mean $\bar{x} = \frac{1}{m} \sum_{i=1}^{m} x_i$

$$x_i \leftarrow x_i - \bar{x} \quad \forall i \in [m] \tag{2}$$

▶ Assume the direction that we would like to project the data is $u$, then the objective function is the data variance

$$J(u) = \frac{1}{m} \sum_{i=1}^{m} (u^\top x_i)^2 \tag{3}$$

▶ Maximize $J(u)$ is trivial, if there is no constriant on $u$. Therefore, we set $\|u\|_2^2 = u^\top u = 1$

# Covariance Matrix

The definition of $J(u)$ can be written as

$$
\begin{align}
J(u) &= \frac{1}{m} \sum_{i=1}^{m} (u^{\mathsf{T}} x_i)^2 \tag{4} \\
&= \frac{1}{m} \sum_{i=1}^{m} u^{\mathsf{T}} x_i u^{\mathsf{T}} x_i \tag{5} \\
&= \frac{1}{m} \sum_{i=1}^{m} u^{\mathsf{T}} x_i x_i^{\mathsf{T}} u \tag{6} \\
&= u^{\mathsf{T}} \left( \frac{1}{m} \sum_{i=1}^{m} x_i x_i^{\mathsf{T}} \right) u \tag{7} \\
&= u^{\mathsf{T}} \Sigma u \tag{8}
\end{align}
$$

where $\Sigma$ is the data covariance matrix

# Optimization

▶ The optimization of finding a single direction projection is

$$\max_{u} J(u) = u^\mathsf{T} \Sigma u \tag{9}$$

$$\text{s.t.} \quad u^\mathsf{T} u = 1 \tag{10}$$

▶ The optimization of finding a single direction projection is

$$\max_{u} J(u) = u^{\mathsf{T}} \Sigma u \qquad (9)$$

$$\text{s.t.} \qquad u^{\mathsf{T}} u = 1 \qquad (10)$$

▶ It can be converted to an unconstrained optimization problem with a Lagrange multiplier

$$\max_{u} \left\{ u^{\mathsf{T}} \Sigma u + \lambda (1 - u^{\mathsf{T}} u) \right\} \qquad (11)$$

# Optimization

▶ The optimization of finding a single direction projection is

$$\max_{u} J(u) = u^{\mathsf{T}} \Sigma u \tag{9}$$

$$\text{s.t.} \quad u^{\mathsf{T}} u = 1 \tag{10}$$

▶ It can be converted to an unconstrained optimization problem with a Lagrange multiplier

$$\max_{u} \left\{ u^{\mathsf{T}} \Sigma u + \lambda (1 - u^{\mathsf{T}} u) \right\} \tag{11}$$

▶ The optimal solution is given by

$$\Sigma u - \lambda u = 0 \tag{12}$$

$$\Sigma u = \lambda u \tag{13}$$

There are two observations from

$$\mathbf{\Sigma}\boldsymbol{u} = \lambda\boldsymbol{u} \tag{14}$$

- First, $\lambda$ is an eigenvalue of $\mathbf{\Sigma}$ and $\boldsymbol{u}$ is the corresponding eigenvector

There are two observations from

$$\Sigma u = \lambda u \tag{14}$$

- First, $\lambda$ is an eigenvalue of $\Sigma$ and $u$ is the corresponding eigenvector
- Second, multiplying $u^\mathsf{T}$ on both sides, we have

$$u^\mathsf{T} \Sigma u = \lambda \tag{15}$$

In order to maximize $J(u)$, $\lambda$ has to the largest eigenvalue $u$ is the corresponding eigen vector.

# Principal Component Analysis

- As $u$ indicates the first major direction that can preserve the data variance, it is called the first principal component

# Principal Component Analysis

▶ As $u$ indicates the first major direction that can preserve the data variance, it is called the first principal component

▶ In general, with eigen decomposition, we have

$$U^\mathsf{T} \Sigma U = \Lambda \tag{16}$$

  ▶ Eigenvalues $\Lambda = \mathrm{diag}(\lambda_1, \ldots, \lambda_d)$
  ▶ Eigenvectors $U = [u_1, \ldots, u_d]$

Assume in $\mathbf{\Lambda} = \text{diag}(\lambda_1, \ldots, \lambda_d)$,

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \tag{17}$$

## Principal Component Analysis (II)

Assume in $\mathbf{\Lambda} = \text{diag}(\lambda_1, \ldots, \lambda_d)$,

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \tag{17}$$

To reduce the dimensionality of $x$ from $d$ to $n$, with $n < d$

▶ Take the first $n$ eigenvectors in $U$ and form

$$\tilde{U} = [u_1, \ldots, u_n] \in \mathbb{R}^{d \times n} \tag{18}$$

Assume in $\mathbf{\Lambda} = \operatorname{diag}(\lambda_1, \ldots, \lambda_d)$,

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \tag{17}$$

To reduce the dimensionality of $x$ from $d$ to $n$, with $n < d$

▶ Take the first $n$ eigenvectors in $U$ and form

$$\tilde{U} = [u_1, \ldots, u_n] \in \mathbb{R}^{d \times n} \tag{18}$$

▶ Reduce the dimensionality of $x$ as

$$\tilde{x} = \tilde{U}^\top x \in \mathbb{R}^n \tag{19}$$

# Principal Component Analysis (II)

Assume in $\mathbf{\Lambda} = \operatorname{diag}(\lambda_1, \ldots, \lambda_d)$,

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \tag{17}$$

To reduce the dimensionality of $x$ from $d$ to $n$, with $n < d$

▶ Take the first $n$ eigenvectors in $U$ and form

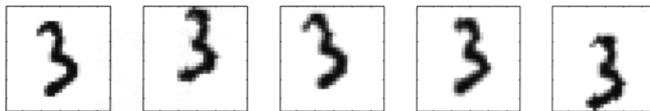$$\tilde{U} = [u_1, \ldots, u_n] \in \mathbb{R}^{d \times n} \tag{18}$$

▶ Reduce the dimensionality of $x$ as

$$\tilde{x} = \tilde{U}^\mathsf{T} x \in \mathbb{R}^n \tag{19}$$

▶ The value of $n$ can be determined by the following

$$\frac{\sum_{i=1}^{n} \lambda_i}{\sum_{i=1}^{d} \lambda_i} \approx 0.95 \tag{20}$$

# Applications: Image Processing

Reduce the dimensionality of an image dataset from $28 \times 28 = 784$ to $M$



(a) Original data
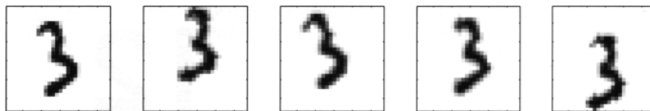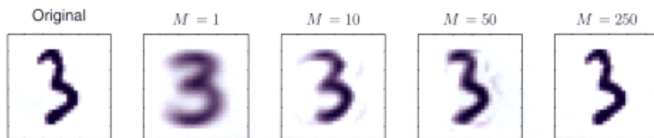
[Bishop, 2006, Section 12.1]

Reduce the dimensionality of an image dataset from $28 \times 28 = 784$ to $M$



(a) Original data



(b) With the first $M$ principal components

[Bishop, 2006, Section 12.1]

# A Different Viewpoint of PCA

Another way to formulate the objective function of PCA

$$\min_{W, U} \sum_{i=1}^{m} \|x_i - UWx_i\|_2^2 \qquad (21)$$

where

- $W \in \mathbb{R}^{n \times d}$: mapping $x_i$ from the original space to a lower-dimensional space $\mathbb{R}^n$
- $U \in \mathbb{R}^{d \times n}$: mapping back the original space $\mathbb{R}^d$

[Shalev-Shwartz and Ben-David, 2014, Chap 23]

Another way to formulate the objective function of PCA

$$\min_{W,U} \sum_{i=1}^{m} \|x_i - UWx_i\|_2^2 \qquad (21)$$

where

- $W \in \mathbb{R}^{n \times d}$: mapping $x_i$ from the original space to a lower-dimensional space $\mathbb{R}^n$
- $U \in \mathbb{R}^{d \times n}$: mapping back the original space $\mathbb{R}^d$
- Dimensionality reduction is performed as $\tilde{x} = Ux$, while $W$ make sure the reduction does not loss much information

[Shalev-Shwartz and Ben-David, 2014, Chap 23]

Consider the optimization problem

$$\min_{W,V} \sum_{i=1}^{m} \|x_i - UWx_i\|_2^2 \qquad (22)$$

- Let $W, U$ be a solution of equation 24
  [Shalev-Shwartz and Ben-David, 2014, Lemma 23.1]
    - the columns of $U$ are orthonormal
    - $W = U^\top$

# Optimization

Consider the optimization problem

$$\min_{W,V} \sum_{i=1}^{m} \|x_i - UWx_i\|_2^2 \qquad (22)$$

- Let $W$, $U$ be a solution of equation 24
  [Shalev-Shwartz and Ben-David, 2014, Lemma 23.1]
  - the columns of $U$ are orthonormal
  - $W = U^\top$
- The optimization problem can be simplified as

$$\min_{U^\top U = I} \sum_{i=1}^{m} \|x_i - UU^\top x_i\|_2^2 \qquad (23)$$

The solution will be the same.

If we extend the both mappings to be nonlinear, then the model becomes a simple encoder-decoder neural network model

$$\min_{W,V} \sum_{i=1}^{m} \|x_i - \tanh(U \cdot \tanh(Wx_i))\|_2^2 \qquad (24)$$

where

- $\tilde{x} = \tanh(Wx_i)$ is a simple encoder
- $x = \tanh(U\tilde{x})$ is a simple decoder
- No closed-form solutions of $W, U$, although the backpropagation algorithm still applies here

# Reference

Bishop, C. M. (2006).
*Pattern recognition and machine learning.*
Springer.

Shalev-Shwartz, S. and Ben-David, S. (2014).
*Understanding machine learning: From theory to algorithms.*
Cambridge university press.