

CS 4774 Machine Learning

Dimensionality Reduction

Yangfeng Ji

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

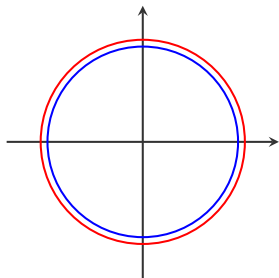


1. Reducing Dimensions
2. Principal Component Analysis
3. A Different Viewpoint of PCA

Reducing Dimensions

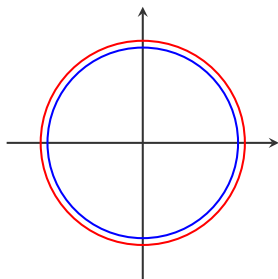
Curse of Dimensionality

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



Curse of Dimensionality

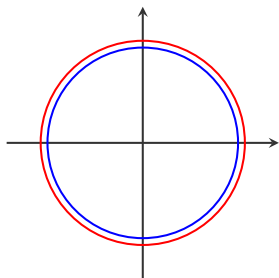
What is the volume difference between two d -dimensional balls with radii $r_1 = 1$ and $r_2 = 0.99$



- ▶ $d = 2$: $\frac{1}{2}\pi(r_1^2 - r_2^2) \approx 0.03$
- ▶ $d = 3$: $\frac{4}{3}\pi(r_1^3 - r_2^3) \approx 0.12$

Curse of Dimensionality

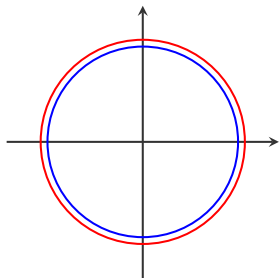
What is the volume difference between two d -dimensional balls with radii $r_1 = 1$ and $r_2 = 0.99$



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

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

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

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

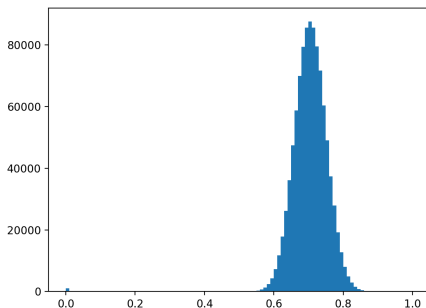


Figure: $d = 100$

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

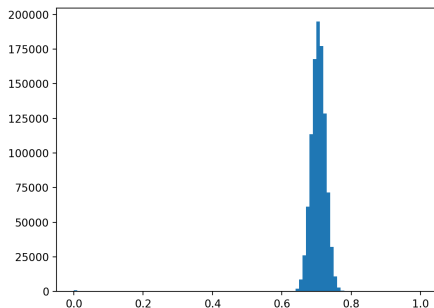


Figure: $d = 500$

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

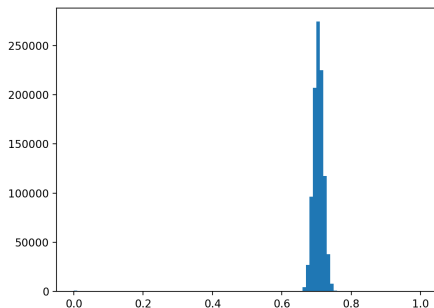


Figure: $d = 1000$

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.

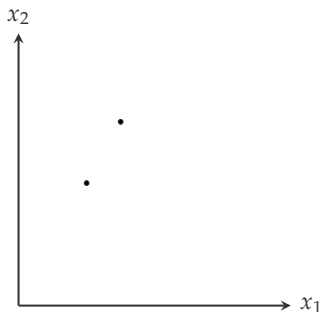
Mathematically, it means

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

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

Reducing Dimensions: A toy example

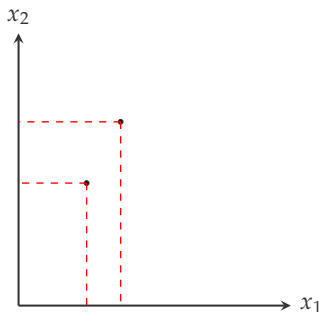
For the purpose of reducing dimensions, we can project $x = (x_1, x_2)$ into the direction along x_1 or x_2



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

Reducing Dimensions: A toy example

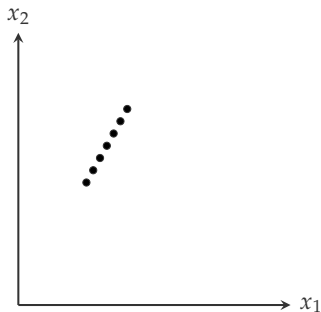
For the purpose of reducing dimensions, we can project $x = (x_1, x_2)$ into the direction along x_1 or x_2



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

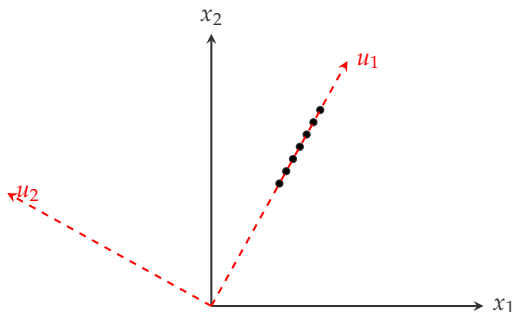
Reducing Dimensions: A toy example (II)

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



Reducing Dimensions: A toy example (II)

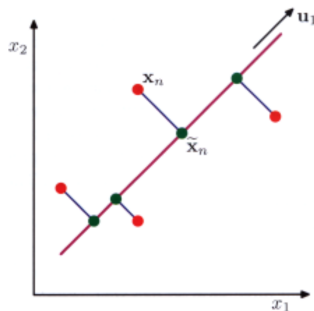
There is a better solution if we are allowed to rotate the coordinate



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

Reducing Dimensions: A toy example (III)

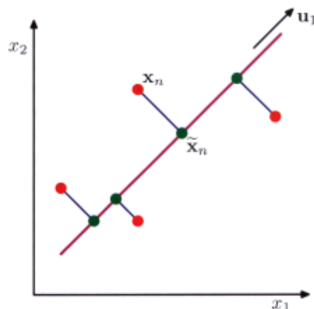
Consider a general case, where the examples do not lie on a perfect line



[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 $S = \{x_1, \dots, x_m\}$

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

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

Given a set of example $S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$

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

$$\mathbf{x}_i \leftarrow \mathbf{x}_i - \bar{\mathbf{x}} \quad \forall i \in [m] \quad (2)$$

- ▶ Assume the direction that we would like to project the data is \mathbf{u} , then the objective function is the data variance

$$J(\mathbf{u}) = \frac{1}{m} \sum_{i=1}^m (\mathbf{u}^\top \mathbf{x}_i)^2 \quad (3)$$

Formulation

Given a set of example $S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$

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

$$\mathbf{x}_i \leftarrow \mathbf{x}_i - \bar{\mathbf{x}} \quad \forall i \in [m] \quad (2)$$

- ▶ Assume the direction that we would like to project the data is \mathbf{u} , then the objective function is the data variance

$$J(\mathbf{u}) = \frac{1}{m} \sum_{i=1}^m (\mathbf{u}^\top \mathbf{x}_i)^2 \quad (3)$$

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

The definition of $J(\mathbf{u})$ can be written as

$$J(\mathbf{u}) = \frac{1}{m} \sum_{i=1}^m (\mathbf{u}^\top \mathbf{x}_i)^2 \quad (4)$$

$$= \frac{1}{m} \sum_{i=1}^m \mathbf{u}^\top \mathbf{x}_i \mathbf{u}^\top \mathbf{x}_i \quad (5)$$

$$= \frac{1}{m} \sum_{i=1}^m \mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u} \quad (6)$$

$$= \mathbf{u}^\top \left(\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{u} \quad (7)$$

$$= \mathbf{u}^\top \mathbf{\Sigma} \mathbf{u} \quad (8)$$

where $\mathbf{\Sigma}$ is the data covariance matrix

- ▶ The optimization of finding a single direction projection is

$$\max_u J(\mathbf{u}) = \mathbf{u}^\top \Sigma \mathbf{u} \quad (9)$$

$$\text{s.t.} \quad \mathbf{u}^\top \mathbf{u} = 1 \quad (10)$$

- ▶ The optimization of finding a single direction projection is

$$\max_u J(\mathbf{u}) = \mathbf{u}^T \Sigma \mathbf{u} \quad (9)$$

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

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

$$\max_u \{ \mathbf{u}^T \Sigma \mathbf{u} + \lambda(1 - \mathbf{u}^T \mathbf{u}) \} \quad (11)$$

- ▶ The optimization of finding a single direction projection is

$$\max_u J(\mathbf{u}) = \mathbf{u}^T \Sigma \mathbf{u} \quad (9)$$

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

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

$$\max_u \{ \mathbf{u}^T \Sigma \mathbf{u} + \lambda(1 - \mathbf{u}^T \mathbf{u}) \} \quad (11)$$

- ▶ The optimal solution is given by

$$\Sigma \mathbf{u} - \lambda \mathbf{u} = 0 \quad (12)$$

$$\Sigma \mathbf{u} = \lambda \mathbf{u} \quad (13)$$

Two Observations

There are two observations from

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

- ▶ First, λ is an eigenvalue of Σ and \mathbf{u} is the corresponding eigenvector

Two Observations

There are two observations from

$$\Sigma \mathbf{u} = \lambda \mathbf{u} \quad (14)$$

- ▶ First, λ is an eigenvalue of Σ and \mathbf{u} is the corresponding eigenvector
- ▶ Second, multiplying \mathbf{u}^\top on both sides, we have

$$\mathbf{u}^\top \Sigma \mathbf{u} = \lambda \quad (15)$$

In order to maximize $J(\mathbf{u})$, λ has to be the **largest** eigenvalue and \mathbf{u} is the corresponding eigenvector.

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 \mathbf{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

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

- ▶ Eigenvalues $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$
- ▶ Eigenvectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_d]$

Principal Component Analysis (II)

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

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \quad (17)$$

Principal Component Analysis (II)

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

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \quad (17)$$

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

- ▶ Take the first n eigenvectors in \mathbf{U} and form

$$\tilde{\mathbf{U}} = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{R}^{d \times n} \quad (18)$$

Principal Component Analysis (II)

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

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \quad (17)$$

To reduce the dimensionality of \mathbf{x} from d to n , with $n < d$

- ▶ Take the first n eigenvectors in \mathbf{U} and form

$$\tilde{\mathbf{U}} = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{R}^{d \times n} \quad (18)$$

- ▶ Reduce the dimensionality of \mathbf{x} as

$$\tilde{\mathbf{x}} = \tilde{\mathbf{U}}^T \mathbf{x} \in \mathbb{R}^n \quad (19)$$

Principal Component Analysis (II)

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

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \quad (17)$$

To reduce the dimensionality of \mathbf{x} from d to n , with $n < d$

- ▶ Take the first n eigenvectors in \mathbf{U} and form

$$\tilde{\mathbf{U}} = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{R}^{d \times n} \quad (18)$$

- ▶ Reduce the dimensionality of \mathbf{x} as

$$\tilde{\mathbf{x}} = \tilde{\mathbf{U}}^T \mathbf{x} \in \mathbb{R}^n \quad (19)$$

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

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

What if we want to reconstruct x (\mathbb{R}^d) from \tilde{x} (\mathbb{R}^n)?

- ▶ The answer is

$$x_{\text{pca}} = \tilde{U}\tilde{x} \in \mathbb{R}^d$$

What if we want to reconstruct x (\mathbb{R}^d) from \tilde{x} (\mathbb{R}^n)?

- ▶ The answer is

$$x_{\text{pca}} = \tilde{U}\tilde{x} \in \mathbb{R}^d$$

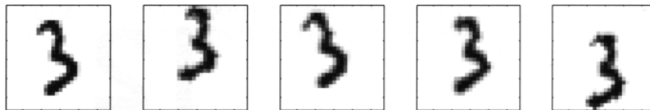
- ▶ Therefore, we have

$$x_{\text{pca}} = \tilde{U}\tilde{U}^T x$$

as a reasonable approximation of x

Applications: Image Processing

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

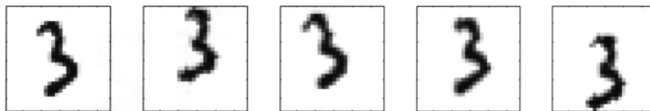


(a) Original data

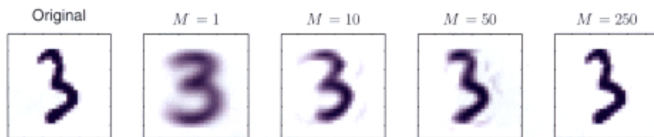
[Bishop, 2006, Section 12.1]

Applications: Image Processing

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

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

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

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

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

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

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

Consider the optimization problem

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

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

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

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

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

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



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.