

Big Data Analysis

(MA60306)

Bibhas Adhikari

Spring 2022-23, IIT Kharagpur

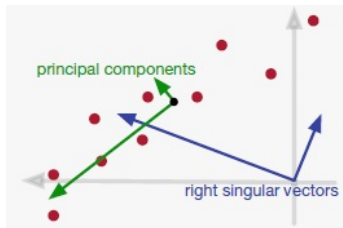
Lecture 11
February 2, 2023

Principal component analysis

Observation the best fitting subspace is a subspace!! i.e. it is a plane which passes through origin

Principal component analysis

Observation the best fitting subspace is a subspace!! i.e. it is a plane which passes through origin



Principal component analysis (PCA) - an extension of SVD when the desired **subspace** V does not pass through origin but it goes through the mean of all the data points! So use SVD after a preprocessing step, called centering to shift the data matrix to its mean at the origin!

PCA

Centering - adjusting the given data matrix $A \in \mathbb{R}^{n \times d}$ such that each column has mean value 0.

PCA

Centering - adjusting the given data matrix $A \in \mathbb{R}^{n \times d}$ such that each column has mean value 0.

→ Define $\bar{\mathbf{a}}_j = \frac{1}{n} \sum_{i=1}^n A_{ij}$. the average of each column

PCA

Centering - adjusting the given data matrix $A \in \mathbb{R}^{n \times d}$ such that each column has mean value 0.

→ Define $\bar{\mathbf{a}}_j = \frac{1}{n} \sum_{i=1}^n A_{ij}$. the average of each column

→ Define \tilde{A} with $\tilde{A}_{ij} = A_{ij} - \bar{\mathbf{a}}_j$, the ij -th entry of the centered matrix \tilde{A}

PCA

Centering - adjusting the given data matrix $A \in \mathbb{R}^{n \times d}$ such that each column has mean value 0.

→ Define $\bar{\mathbf{a}}_j = \frac{1}{n} \sum_{i=1}^n A_{ij}$. the average of each column

→ Define \tilde{A} with $\tilde{A}_{ij} = A_{ij} - \bar{\mathbf{a}}_j$, the ij -th entry of the centered matrix \tilde{A}

→ Another way: define the centering matrix $C_n = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T$, where $\mathbf{1}$ is the all-one vector. Then

$$\tilde{A} = C_n A$$

PCA

Centering - adjusting the given data matrix $A \in \mathbb{R}^{n \times d}$ such that each column has mean value 0.

→ Define $\bar{\mathbf{a}}_j = \frac{1}{n} \sum_{i=1}^n A_{ij}$. the average of each column

→ Define \tilde{A} with $\tilde{A}_{ij} = A_{ij} - \bar{\mathbf{a}}_j$, the ij -th entry of the centered matrix \tilde{A}

→ Another way: define the centering matrix $C_n = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T$, where $\mathbf{1}$ is the all-one vector. Then

$$\tilde{A} = C_n A$$

→ The matrix C_n is a projection matrix!! Where does it project?

PCA

Centering - adjusting the given data matrix $A \in \mathbb{R}^{n \times d}$ such that each column has mean value 0.

- Define $\bar{\mathbf{a}}_j = \frac{1}{n} \sum_{i=1}^n A_{ij}$. the average of each column
- Define \tilde{A} with $\tilde{A}_{ij} = A_{ij} - \bar{\mathbf{a}}_j$, the ij -th entry of the centered matrix \tilde{A}
- Another way: define the centering matrix $C_n = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T$, where $\mathbf{1}$ is the all-one vector. Then

$$\tilde{A} = C_n A$$

- The matrix C_n is a projection matrix!! **Where does it project?**
- Let SVD of $\tilde{A} = C_n A = U \Sigma V^T$. Then the singular values of \tilde{A} are called the *principal values*, and the k singular vectors corresponding to the k largest singular values are called *top-k principal directions/vectors*

PCA

Let

$$A = \begin{bmatrix} 1 & 5 \\ 2 & 3 \\ 3 & 10 \end{bmatrix}.$$

Then the center vector is $\bar{a} = [2, 6]$

The centered matrix is

$$\tilde{A} = \begin{bmatrix} -1 & -1 \\ 0 & -3 \\ 1 & 4 \end{bmatrix}$$

PCA

Another interpretation of PCA

→ We introduce a complete orthonormal set of d -dimensional vectors $\mathbf{v}_j, 1 \leq j \leq d$ that satisfy $\mathbf{v}_i^T \mathbf{v}_j = \delta_{i,j}$

PCA

Another interpretation of PCA

- We introduce a complete orthonormal set of d -dimensional vectors $\mathbf{v}_j, 1 \leq j \leq d$ that satisfy $\mathbf{v}_i^T \mathbf{v}_j = \delta_{i,j}$
- Then any data point \mathbf{x}_i can be written as

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

i.e. this corresponds to a rotation of the coordinate system to a new system defined by the \mathbf{v}_j , and the original d components $\{x_{i1}, \dots, x_{id}\}$ are replaced by an equivalent set $\{\alpha_{i1}, \dots, \alpha_{id}\}$

PCA

Another interpretation of PCA

- We introduce a complete orthonormal set of d -dimensional vectors $\mathbf{v}_j, 1 \leq j \leq d$ that satisfy $\mathbf{v}_i^T \mathbf{v}_j = \delta_{i,j}$
- Then any data point \mathbf{x}_i can be written as

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

i.e. this corresponds to a rotation of the coordinate system to a new system defined by the \mathbf{v}_j , and the original d components $\{x_{i1}, \dots, x_{id}\}$ are replaced by an equivalent set $\{\alpha_{i1}, \dots, \alpha_{id}\}$

- Obviously, $\alpha_{ij} = \mathbf{x}_i^T \mathbf{v}_j$ and hence

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

PCA

- Goal: to approximate the data points using a representation involving a restricted number $k < d$ of variables corresponding to a projection onto a lower-dimensional subspace

- Goal: to approximate the data points using a representation involving a restricted number $k < d$ of variables corresponding to a projection onto a lower-dimensional subspace
- The k -dimensional subspace can be represented WLOG by the first k vectors, and so we approximate each data point \mathbf{x}_i by

$$\tilde{\mathbf{x}}_i = \sum_{j=1}^k z_{ij} \mathbf{v}_j + \sum_{j=k+1}^d \mathbf{b}_j \mathbf{v}_j$$

where $\{z_{ij}\}$ depend on the particular data point, and $\{\mathbf{b}_j\}$ are constants that are the same for all data points

- We are free to choose the $\{\mathbf{v}_j\}$, $\{z_{ij}\}$, and $\{\mathbf{b}_j\}$ so as to minimize the distortion introduced by the reduction in dimensionality

PCA

→ The distortion measure that we consider is the squared distance between the original data point \mathbf{x}_i , and its approximation $\tilde{\mathbf{x}}_i$, averaged over the data set i.e. to minimize

$$J = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2$$

PCA

- The distortion measure that we consider is the squared distance between the original data point \mathbf{x}_i , and its approximation $\tilde{\mathbf{x}}_i$, averaged over the data set i.e. to minimize

$$J = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2$$

- First, consider this minimization wrt $\{z_{ij}\}$:

Homework Substituting $\tilde{\mathbf{x}}_i$, setting the derivative with respect to z_{ij} to zero, and making use of the orthonormality conditions, one can obtain

$$z_{ij} = \mathbf{x}_i^T \mathbf{v}_j, j = 1, \dots, d$$

PCA

- The distortion measure that we consider is the squared distance between the original data point \mathbf{x}_i , and its approximation $\tilde{\mathbf{x}}_i$, averaged over the data set i.e. to minimize

$$J = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2$$

- First, consider this minimization wrt $\{z_{ij}\}$:

Homework Substituting $\tilde{\mathbf{x}}_i$, setting the derivative with respect to z_{ij} to zero, and making use of the orthonormality conditions, one can obtain

$$z_{ij} = \mathbf{x}_i^T \mathbf{v}_j, j = 1, \dots, d$$

Homework Similarly, setting the derivative of J wrt \mathbf{b}_j to zero gives

$$\mathbf{b}_j = \bar{\mathbf{x}}^T \mathbf{v}_j, j = k + 1, \dots, d$$

where $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$

→ Then we have $x_i - \tilde{x}_i = \sum_{j=k+1}^d \{(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{v}_j\} \mathbf{v}_j$

→ Thus

$$J = \frac{1}{n} \sum_{i=1}^n \sum_{j=k+1}^d \left(\mathbf{x}_i^T \mathbf{v}_j - \bar{\mathbf{x}}^T \mathbf{v}_j \right)^2 = \sum_{j=k+1}^d \mathbf{v}_j^T S \mathbf{v}_j$$

where S is the covariance matrix defined by

$$S = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

→ Then we have $\mathbf{x}_i - \tilde{\mathbf{x}}_i = \sum_{j=k+1}^d \{(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{v}_j\} \mathbf{v}_j$

→ Thus

$$J = \frac{1}{n} \sum_{i=1}^n \sum_{j=k+1}^d \left(\mathbf{x}_i^T \mathbf{v}_j - \bar{\mathbf{x}}^T \mathbf{v}_j \right)^2 = \sum_{j=k+1}^d \mathbf{v}_j^T S \mathbf{v}_j$$

where S is the covariance matrix defined by

$$S = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Homework Then show that the general solution to the minimization for J for arbitrary d and $k < d$ is obtained by choosing the $\{\mathbf{v}_j\}$ as the eigenvectors of the the covariance matrix S ¹

¹Chapter 12, C. M. Bishop, Pattern Recognition and Machine Learning, Springer, 2009

Multidimensional scaling

Multidimensional scaling (MDS) - is a data analysis technique which translates distances and dissimilarities into a visual representation through a 'geometric' picture

Multidimensional scaling

Multidimensional scaling (MDS) - is a data analysis technique which translates distances and dissimilarities into a visual representation through a 'geometric' picture

- The input for an MDS algorithm usually is not an object data set but similarities of a set of objects
- Here we will use the term 'distance' in a generic sense, meaning it reflects a dissimilarity/similarity between pairs of the objects

Multidimensional scaling

Multidimensional scaling (MDS) - is a data analysis technique which translates distances and dissimilarities into a visual representation through a 'geometric' picture

- The input for an MDS algorithm usually is not an object data set but similarities of a set of objects
- Here we will use the term 'distance' in a generic sense, meaning it reflects a dissimilarity/similarity between pairs of the objects
- Suppose there are n objects in a set and the similarities between all pairs are measured. Then we can define an $n \times n$ distance matrix $D = [d_{ij}]$, where d_{ij} represents the similarity/distance between the objects i and j

Multidimensional scaling

Multidimensional scaling (MDS) - is a data analysis technique which translates distances and dissimilarities into a visual representation through a 'geometric' picture

- The input for an MDS algorithm usually is not an object data set but similarities of a set of objects
- Here we will use the term 'distance' in a generic sense, meaning it reflects a dissimilarity/similarity between pairs of the objects
- Suppose there are n objects in a set and the similarities between all pairs are measured. Then we can define an $n \times n$ distance matrix $D = [d_{ij}]$, where d_{ij} represents the similarity/distance between the objects i and j
- The objects are configured as virtual points in a low dimensional linear Euclidean space, and this point set is called *configuration* such that the Euclidean distances between the points have closest relation to the similarities

MDS

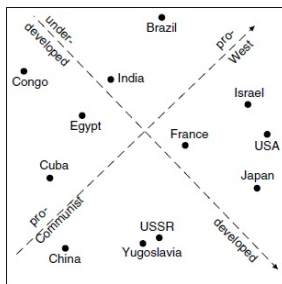
Example Similarity ratings for 12 nations (Wish, 1971)

Nation		1	2	3	4	5	6	7	8	9	10	11	12
Brazil	1	—											
Congo	2	4.83	—										
Cuba	3	5.28	4.56	—									
Egypt	4	3.44	5.00	5.17	—								
France	5	4.72	4.00	4.11	4.78	—							
India	6	4.50	4.83	4.00	5.83	3.44	—						
Israel	7	3.83	3.33	3.61	4.67	4.00	4.11	—					
Japan	8	3.50	3.39	2.94	3.83	4.22	4.50	4.83	—				
China	9	2.39	4.00	5.50	4.39	3.67	4.11	3.00	4.17	—			
USSR	10	3.06	3.39	5.44	4.39	5.06	4.50	4.17	4.61	5.72	—		
U.S.A.	11	5.39	2.39	3.17	3.33	5.94	4.28	5.94	6.06	2.56	5.00	—	
Yugoslavia	12	3.17	3.50	5.11	4.28	4.72	4.00	4.44	4.28	5.06	6.67	3.56	—

MDS

Example Similarity ratings for 12 nations (Wish, 1971)

Nation		1	2	3	4	5	6	7	8	9	10	11	12
Brazil	1	—											
Congo	2	4.83	—										
Cuba	3	5.28	4.56	—									
Egypt	4	3.44	5.00	5.17	—								
France	5	4.72	4.00	4.11	4.78	—							
India	6	4.50	4.83	4.00	5.83	3.44	—						
Israel	7	3.83	3.33	3.61	4.67	4.00	4.11	—					
Japan	8	3.50	3.39	2.94	3.83	4.22	4.50	4.83	—				
China	9	2.39	4.00	5.50	4.39	3.67	4.11	3.00	4.17	—			
USSR	10	3.06	3.39	5.44	4.39	5.06	4.50	4.17	4.61	5.72	—		
U.S.A.	11	5.39	2.39	3.17	3.33	5.94	4.28	5.94	6.06	2.56	5.00	—	
Yugoslavia	12	3.17	3.50	5.11	4.28	4.72	4.00	4.44	4.28	5.06	6.67	3.56	—



MDS

Example Distances between ten cities

	1	2	3	4	5	6	7	8	9	10
1	0	569	667	530	141	140	357	396	570	190
2	569	0	1212	1043	617	446	325	423	787	648
3	667	1212	0	201	596	768	923	882	714	714
4	530	1043	201	0	431	608	740	690	516	622
5	141	617	596	431	0	177	340	337	436	320
6	140	446	768	608	177	0	218	272	519	302
7	357	325	923	740	340	218	0	114	472	514
8	396	423	882	690	337	272	114	0	364	573
9	569	787	714	516	436	519	472	364	0	755
10	190	648	714	622	320	302	514	573	755	0

MDS

Example Distances between ten cities

	1	2	3	4	5	6	7	8	9	10
1	0	569	667	530	141	140	357	396	570	190
2	569	0	1212	1043	617	446	325	423	787	648
3	667	1212	0	201	596	768	923	882	714	714
4	530	1043	201	0	431	608	740	690	516	622
5	141	617	596	431	0	177	340	337	436	320
6	140	446	768	608	177	0	218	272	519	302
7	357	325	923	740	340	218	0	114	472	514
8	396	423	882	690	337	272	114	0	364	573
9	569	787	714	516	436	519	472	364	0	755
10	190	648	714	622	320	302	514	573	755	0

