

Principal Components Analysis: Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie Lindstrom

Pacific Lutheran University

2023

Principal Components Analysis

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

PCA is a mathematical tool used to reduce the dimension of a data set in the hope of obtaining useful and understandable takeaways.

Linear Combination

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Definition (Linear Combination)

A linear combination is any sum of vectors multiplied by constants, such as

$$\alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \alpha_3 \mathbf{u}_3 + \alpha_m \mathbf{u}_m = \sum_{j=1}^m \alpha_j \mathbf{u}_j \quad (1)$$

where \mathbf{u} are vectors and α are constants. [2]

Linear Transformation

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Definition (Linear Transformation)

A mapping f from \mathbb{R}^n to \mathbb{R}^m is **linear** if

$$f(ax + by) = af(\mathbf{x}) + bf(\mathbf{y})$$

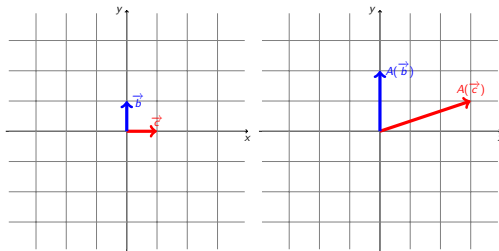
for all vectors \mathbf{x} and \mathbf{y} in \mathbb{R}^n and for all scalars a and b . [2]

Linear Transformation Example

Here, A is a square matrix representing some linear transformation.

$$A = \begin{bmatrix} 3 & 0 \\ 1 & 2 \end{bmatrix}$$

$$A(\vec{b}) = \begin{bmatrix} 3 & 0 \\ 1 & 2 \end{bmatrix} * \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} (3 * 0) + (0 * 1) \\ (1 * 0) + (2 * 1) \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \end{bmatrix}$$



A can transform the vectors \vec{b} and \vec{c} to what we see on the right. Notice \vec{b} stays on its span

Eigenvectors and Eigenvalues

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Definition

Let \mathbf{A} be any square matrix, real or complex. A number λ is an **eigenvalue** of \mathbf{A} if the equation

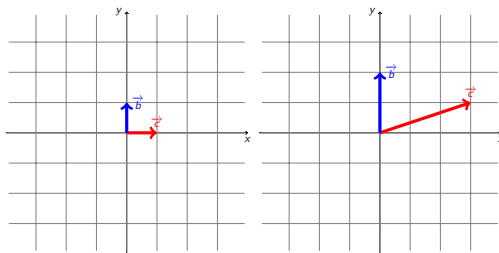
$$\mathbf{A}\vec{v} = \lambda\vec{v} \quad (2)$$

is true for some nonzero vector \vec{v} . The vector \vec{v} is an eigenvector associated with the eigenvalue λ . [2]

- left hand side of the equation: matrix vector multiplication
- right hand side: scalar vector multiplication

Visualizing Eigenvectors and Eigenvalues

$$A = \begin{bmatrix} 3 & 0 \\ 1 & 2 \end{bmatrix}$$



$$A\vec{v} = \lambda\vec{v} \quad (3)$$

Here, \vec{b} is an eigenvector because it remains on its span (the y-axis) when transformed by A . However, \vec{c} is not an eigenvector of A .

Factoring Matrices

Whenever we see a matrix written in the form

$$A = BC$$

we have a factorization of a matrix into the product of the two matrices B and C . [2]

Factorization can be useful to visualize and extract information from our matrix.

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

Eigen Decomposition

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

One form of factorization is the eigen-decomposition.

Recall Equation 3

$$A\vec{v} = \lambda\vec{v} \quad (4)$$

Take U to be the matrix where each column is an eigenvector of A . Take Λ to be a diagonal matrix that stores the eigenvalues of A . Therefore, we can rewrite Equation 3 as:

$$AU = \Lambda U$$

Thinking in terms of a change of basis, we can write this as

$$A = U\Lambda U^{-1}$$

where U^{-1} is the inverse matrix of U .

Eigen Decomposition Continued

Principal Components Analysis:
Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie Lindstrom

Linear Algebra Review

Singular Value Decomposition

Principal Components Analysis

Running PCA

References

$$A = U\Lambda U^{-1}$$

$$A = \begin{bmatrix} 3 & 0 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} * \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} * \begin{bmatrix} 1 & 0 \\ -1 & 1 \end{bmatrix}$$

Therefore, the eigen-decomposition is a way of rewriting matrix A in a form that clearly shows what the eigenvectors and eigenvalues are.

- Eigenvalues: 3, 2
- Eigenvectors: $[1,1]$, $[0,1]$

Singular Value Decomposition

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

- A generalization of the eigen-decomposition.
- Takes a matrix A and decomposes it into three factors
- Each factor gives useful information about the original matrix A
- Matrix A can be a non-square matrix (unlike the eigen-decomposition) [1].

Singular Value Decomposition

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Theorem

Any $m \times n$ matrix can be put into the factored form \mathbf{PDQ} , where \mathbf{P} and \mathbf{Q} are unitary and \mathbf{D} is diagonal. [2]

A complex matrix U is unitary if

$$UU^H = U^H U = I$$

Singular Value Decomposition

Given matrix A , the SVD is:

$$A = PDQ$$

P: A matrix of the left singular vectors of A

D: A diagonal matrix with the singular values on its diagonal

Q: A matrix of the right singular vectors of A

Written in matrix form:

$$A = PDQ$$

$$\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} * \begin{bmatrix} a & \cdot \\ \cdot & b \\ \cdot & \cdot \end{bmatrix} * \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$$

Singular Value Decomposition

Principal Components Analysis:
Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie Lindstrom

Linear Algebra Review

Singular Value Decomposition

Principal Components Analysis

Running PCA

References

Singular Vectors: The left and right SVs combined span all of matrix A

- Left Singular Vectors: Spans the column space of A . The normalized eigenvectors of AA^T
- Right Singular Vectors: Spans the row space of A . The normalized eigenvectors of $A^T A$

Singular values: The square roots of the eigenvalues of AA^T and $A^T A$

The largest singular vectors explain the most variance of matrix A , and the singular values associated with each singular vector tell us how much variance is explained.

Singular Value Decomposition

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Why are we referring to AA^T and $A^T A$ instead of just matrix A ?

- A is not necessarily square, and may not have positive eigenvalues that exist
- The matrices AA^T and $A^T A$ are positive semidefinite - which means they are symmetric, have positive eigenvalues, and eigenvectors are pairwise orthogonal
- Still useful to find the eigenvectors and eigenvalues of AA^T and $A^T A$

SVD in Relation to PCA

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

PCA is built on SVD. Eigenvectors, known as components, are retained based on the amount of variance that they explain (quantified by eigenvalues).

PCA Overview

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

The four goals of PCA [1]

- Extract only the most important information from the data set
- Compress/ reduce the size of the data set via keeping only the important information
- Simplify the description of the data set
- Analyze the structure of (and relationships between) the variables

PCA Overview

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

- PCA uses the vectors in the original data table and finds a set of new, orthogonal vectors, called **Principal Components**
- Original vectors can be expressed as linear combinations of the new principal components.
- The primary principal component captures the largest variance

Change of Basis and Dimension Reduction

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

The original basis of the dataset before PCA is done has one variable per axis = confusing

After PCA, the basis is usually the first principal components.

Therefore, we can create plots such as biplots which simply have a principal component on the x and y axis = less confusing

Applying PCA to Iraqi Migration Data

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

PCA was performed using R on a data set. This data set holds the results of a survey sourced from the International Organization for Migration. The survey was conducted in Iraq asking internally displaced persons about their access to certain needs, such as distance from clinics, access to clean water, etc.

- 206 variables
- 3,718 observations
- Collected April- June 2022

Data Cleaning

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

- Identify quantitative variables
- Remove commas from values with more than 3 digits
- Rename variables
- Create smaller dataframes with desired variables
- Addressed missing data... replaced missing with 0 for Residence good/ bad

Data Cleaning

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Noticed that the number of families and number of individuals were actually just scaled by 6...

Step 1: Calculate Correlation Matrix

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Correlation Matrix: In order to standardize data, the original matrix, A , is centered so that the mean of each column equals 0, and A will be normalized by dividing each variable by the norm.

- Therefore, values will range from 0 to $+/- 1$.
- The diagonals will be 1 and this is a symmetric matrix. [1]

The code used in R to do this is shown below:

```
1 #Create a correlation matrix from the original (cleaned) data set A
2 cor1<-cor(originalData)
3
4 #Print this correlation matrix
5 cor1
6
```

Step 2: Use SVD to find Eigenvectors and Eigenvalues of Correlation Matrix

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

SVD is used to factor the correlation matrix into a form that shows the eigenvectors and eigenvalues of the correlation matrix.

Step 3: Get First Principal Component from Largest Eigenvalue

Principal Components Analysis:
Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie Lindstrom

Linear Algebra Review

Singular Value Decomposition

Principal Components Analysis

Running PCA

References

The largest eigenvalue tells us the eigenvector associated with capturing the most variance. This eigenvector becomes the first Principal Component.

R Code

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Running PCA in R on the correlation matrix is shown below.

```
13  
14 #Run PCA on the correlation matrix  
15 data.pca<-princomp(cor_matrix)  
16  
17 #Print results to show each component  
18 data.pca  
19
```

Results

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Biplot: Shows the loadings on a graph to help visualize relationships between variables and interpret the principal components.

- **Loadings:** The length and the direction of the vectors. The loadings can tell us how much each variable contributes to a certain principal component

Initial Results

Principal Components Analysis:
Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie Lindstrom

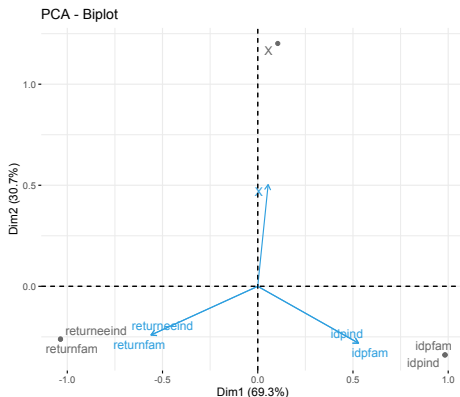
Linear Algebra Review

Singular Value Decomposition

Principal Components Analysis

Running PCA

References



- X-axis: PC1, y-axis: PC2
- vector direction: which PC that variable contributes more to
- vector length: how much the variable contributes to a PC

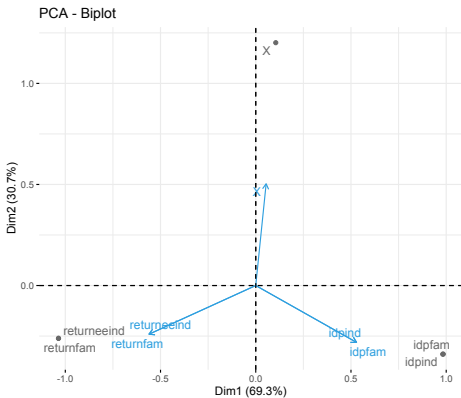
Initial Results

Principal Components Analysis: Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie
Lindstrom

Running PCA

References



Oops! Noticed that I accidentally included the index as a variable for my first run of PCA on a mini dataframe of 4 variables.

Initial Results

Principal Components Analysis:
Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie Lindstrom

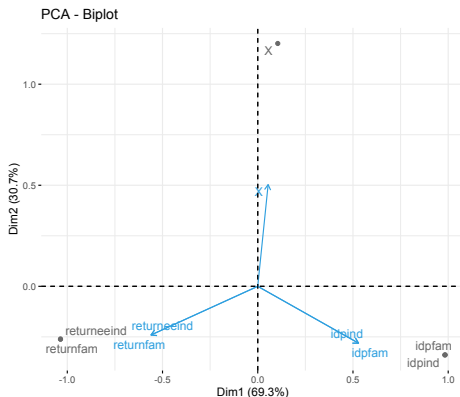
Linear Algebra Review

Singular Value Decomposition

Principal Components Analysis

Running PCA

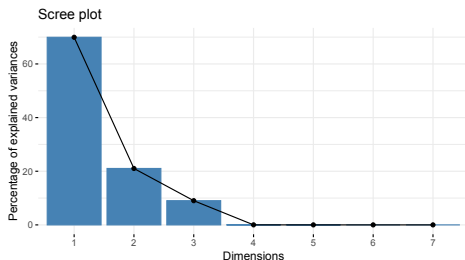
References



- the variable for the index, X , contributes almost entirely to PC2
- the number of individuals and families are slightly negatively correlated to PC2

Results

Scree plot



A scree plot graphs the explained variance versus the principal component number.

Biplot

Principal Components Analysis: Using Math to Extract Underlying Structures of Iraqi Migration Survey Data

Jackie Lindstrom

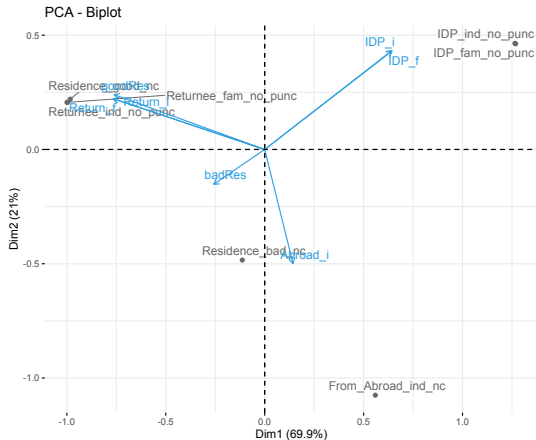
Linear Algebra Review

Singular Value Decomposition

Principal Components Analysis

Running PCA

References



Takeaways

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposition

Principal
Components
Analysis

Running PCA

References

- Variables for number of individuals and families are multiples of each other
- Residence good and residence bad load on opposite direction vectors
- PCA uses linear algebra to help us understand a dataset and make it less complicated
- Data cleaning is very involved

References

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

- [1] Herve Abdi and Lynne J. Williams. “Principal Component Analysis”. In: *WIREs Comp Stat* 2.4 (2010), pp. 433–459. DOI: 10.1002/wics.101.
- [2] Ward Cheney and David Kincaid. *Linear Algebra Theory and Applications*. Jones and Bartlett Publishers, 2009. ISBN: 9780763750206.

Acknowledgements

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

- International Organization for Migration and H. Tran: collecting and sharing their data with me
- Dr. Justice for patience, support, and tea
- PLU Math Department

Positive Semi-Definite Matrices

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

A **positive semi-definite matrix** P can be obtained as the product of a matrix X and its transpose X^T .

$$P = XX^T \quad (5)$$

The properties of a positive semi-definite matrix are:

- Symmetric (perhaps define at beginning?)
- Positive eigenvalues
- Pairwise orthogonal eigenvectors

Positive Semi-Definite Matrices

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

In order to connect this to eigenvalues and eigenvectors, we recall U , the matrix where each column is an eigenvector of A , and Λ , the matrix that stores the eigenvalues of A . We can write P as

$$P = U\Lambda U^{-1} \quad (6)$$

Finding Eigenvalues

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Therefore, finding the eigenvalues and eigenvectors for a matrix A involves finding values for \vec{v} and λ that satisfy this equation. We can rewrite Equation 3 as

$$A\vec{v} = (\lambda I)\vec{v} \quad (7)$$

where I is the identity matrix with 1's down the diagonal. Now we have matrix multiplication on both sides of our equation.

Using algebra, we can rearrange equation 4 to give us:

$$(A - \lambda I)\vec{v} = 0 \quad (8)$$

This notation is helpful in showing that the eigenvectors, \vec{v} , remain on their span for some linear transformation represented by matrix A .

Unitary Definition

Principal
Components
Analysis:
Using Math to
Extract
Underlying
Structures of
Iraqi
Migration
Survey Data

Jackie
Lindstrom

Linear Algebra
Review

Singular Value
Decomposi-
tion

Principal
Components
Analysis

Running PCA

References

Definition

A real matrix U is orthogonal if

$$UU^T = U^T U = I$$

A complex matrix U is unitary if

$$UU^H = U^H U = I$$

[2]