

Principal Component Analysis + Singular Value Decomposition

13 April 2023

Sumit Kumar Yadav

Department of Management Studies
Indian Institute of Technology, Roorkee

Recap and Today

- PCA theory
- PCA some more theory
- PCA Implementation
- SVD Theory + Implementation

PCA Motivation

- To reduce the number of dimensions in the data
- To visualize the data
- To avoid over fitting

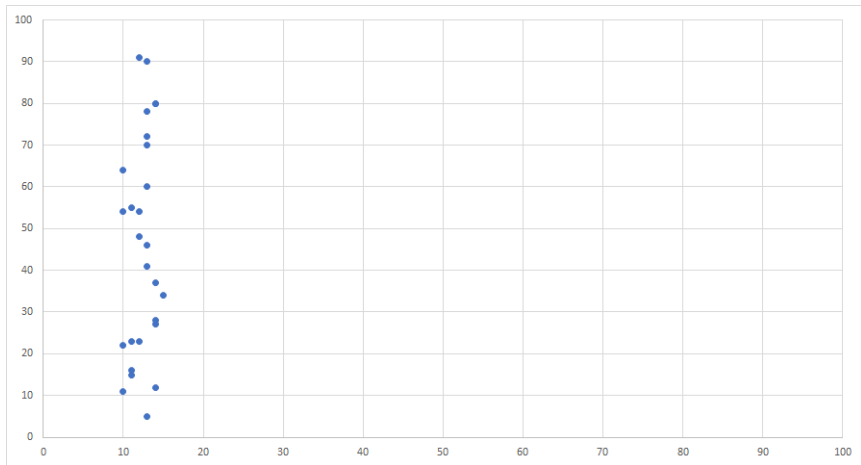
PCA Motivation

- To reduce the number of dimensions in the data
- To visualize the data
- To avoid over fitting

	Location	City	Society	Ambience	Airport
House 1	8	4	9	1	5
House 2	9	6	9	5	5
House 3	10	8	9	7	5
House 4	10	5	9	6	5
House 5	5	4	9	2	5
House 6	2	7	9	9	5
House 7	7	5	9	8	6
House 8	3	4	9	8	5
House 9	4	2	9	7	5
House 10	1	4	9	10	5

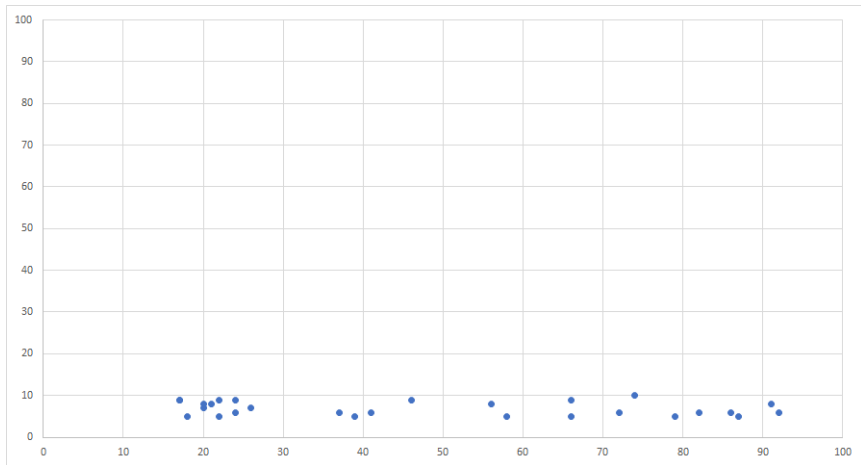
PCA Motivation

Which direction to skip?



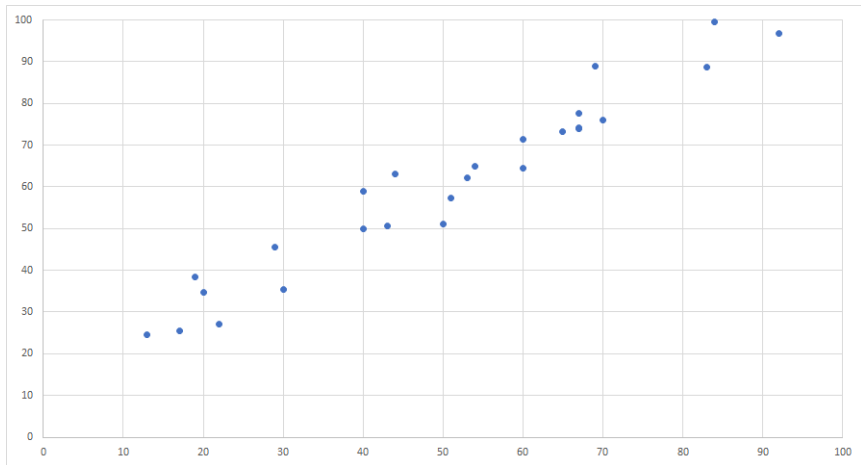
PCA Motivation

Which direction to skip?



PCA Motivation

Which direction to skip?



Looks like we are trying to capture as much variance as possible while reducing the number of dimensions.

What is a vector?

- a mathematical object that encodes a length and direction
- A vector is often represented as a 1-dimensional array of numbers, referred to as components and is displayed either in column form or row form
- Represented geometrically, vectors typically represent coordinates within a n -dimensional space

Vectors in R^n /Concept of Vector Space

- Vector Space - Closed under addition and multiplication
- Vectors on line $y = 2 * x + 1$ form a vector space??
- Addition/Subtraction (Graphical Representation)
- Multiplication by a scalar
- Dot Product (Inner Product)
- Length/Magnitude
- Angle between two vectors - $\cos \theta = \frac{\vec{x} \cdot \vec{y}}{||\vec{x}|| \cdot ||\vec{y}||}$
- Dot product of perpendicular vectors = ??

Linearly Independent and Dependent Vectors

- concept of zero vector

Consider m vectors in \mathbb{R}^n ,

If $\alpha_1 \vec{v}_1 + \alpha_2 \vec{v}_2 + \cdots + \alpha_n \vec{v}_n$ implies

$\alpha_1 = \alpha_2 = \cdots = \alpha_m = \vec{0}$, then

the vectors are said to be linearly independent.

In mathematics, a matrix (plural matrices) is a rectangular array or table of numbers, symbols, or expressions, arranged in rows and columns, which is used to represent a mathematical object or a property of such an object. - Wikipedia

- Matrix Size
- Representing any element of a matrix
- $(1,n)$ and $(m,1)$ matrices
- Square Matrix

Operations on Matrices

- Addition/Subtraction
- Multiplication of Matrix by a scalar
- Transpose of a Matrix
- Multiplication of two Matrices

Matrix as Linear Transformation

Types of Matrices

- Diagonal
- Identity
- Upper Triangular
- Symmetric
- Singular
- Mirror Matrix

- Inverse of a Matrix
- Trace of a Matrix
- Relationship of Eigenvalues with trace and Determinant
- Symmetric Matrix
- Orthogonal Matrix

Eigenvalues and Eigenvectors

Let M be a $n \times n$ matrix. A non-zero vector \vec{X} is said to be an eigenvector of M corresponding to eigenvalue λ if

Spectral Decomposition

A $n \times n$ matrix can be written in terms of its eigenvalues and eigenvectors as follows -

$$M = \lambda_1 \vec{v}_1 \vec{v}_1^T + \lambda_2 \vec{v}_2 \vec{v}_2^T + \dots + \lambda_n \vec{v}_n \vec{v}_n^T$$

Positive Semi-definite Matrix

$$x^T M x \geq 0$$

All the eigenvalues are greater than zero.

We have a Data Matrix(D), whose dimension is $n * f$.

$$\begin{bmatrix} d_{11} & d_{12} & \dots & d_{1f} \\ d_{21} & d_{22} & \dots & d_{2f} \\ \vdots & \vdots & \dots & \vdots \\ d_{(n-1)1} & d_{(n-1)2} & \dots & d_{(n-1)f} \\ d_{n1} & d_{n2} & \dots & d_{nf} \end{bmatrix}$$

We want to reduce the number of features to f_s .

We have a Data Matrix(D), whose dimension is $n * f$.

$$\begin{bmatrix} d_{11} & d_{12} & \dots & d_{1f} \\ d_{21} & d_{22} & \dots & d_{2f} \\ \vdots & \vdots & \dots & \vdots \\ d_{(n-1)1} & d_{(n-1)2} & \dots & d_{(n-1)f} \\ d_{n1} & d_{n2} & \dots & d_{nf} \end{bmatrix}$$

We want to reduce the number of features to f_s .

Let us solve a simpler problem first. Let us say we want to reduce the number of dimensions/features to just 1. Which is the best way to go about it?

We are looking for a unit vector $\vec{v} = (v_1, v_2, \dots, v_f)$ such that the variance of $D\vec{v}$ is as large as possible.

We are looking for a unit vector $\vec{v} = (v_1, v_2, \dots, v_f)$ such that the variance of $D\vec{v}$ is as large as possible.

Variance($D\vec{v}$) can be written as $\vec{v}^T \Sigma \vec{v}$, where Σ is the variance covariance matrix of D . (Try to prove this or atleast verify this)

We are looking for a unit vector $\vec{v} = (v_1, v_2, \dots, v_f)$ such that the variance of $D\vec{v}$ is as large as possible.

Variance($D\vec{v}$) can be written as $\vec{v}^T \Sigma \vec{v}$, where Σ is the variance covariance matrix of D . (Try to prove this or atleast verify this)

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1f} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2f} \\ \vdots & \vdots & \dots & \vdots \\ \sigma_{(f-1)1} & \sigma_{(f-1)2} & \dots & \sigma_{(f-1)f} \\ \sigma_{1f} & \sigma_{f2} & \dots & \sigma_f^2 \end{bmatrix}$$

PCA Optimization problem

Maximize - $\vec{v}^T \Sigma \vec{v}$

such that $\|\vec{v}\| = 1$

PCA Optimization problem

Maximize - $\vec{v}^T \Sigma \vec{v}$

such that $\|\vec{v}\| = 1$

Σ can be shown to be positive semi-definite. Thus, all the eigenvalues are greater than zero, and all eigenvectors are orthogonal.

PCA Optimization problem

Maximize - $\vec{v}^T \Sigma \vec{v}$

such that $\|\vec{v}\| = 1$

Σ can be shown to be positive semi-definite. Thus, all the eigenvalues are greater than zero, and all eigenvectors are orthogonal.

$$\Sigma = \lambda_1 \vec{v}_1 \vec{v}_1^T + \lambda_2 \vec{v}_2 \vec{v}_2^T + \dots + \lambda_n \vec{v}_n \vec{v}_n^T$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$$

PCA Optimization problem

Maximize - $\vec{v}^T(\lambda_1 \vec{v}_1 \vec{v}_1^T + \lambda_2 \vec{v}_2 \vec{v}_2^T + \dots + \lambda_p \vec{v}_p \vec{v}_p^T) \vec{v}$

such that $\|\vec{v}\| = 1$

PCA Optimization problem

Maximize - $\vec{v}^T(\lambda_1 \vec{v}_1 \vec{v}_1^T + \lambda_2 \vec{v}_2 \vec{v}_2^T + \dots + \lambda_n \vec{v}_n \vec{v}_n^T) \vec{v}$

such that $\|\vec{v}\| = 1$

All \vec{v}_i 's are unit vectors and orthogonal to each other, so the best choice for \vec{v} is ??

What about reducing to 2 dimensions?

Variance is no longer a scalar which can be compared across all possible 2 dimensions where data can be projected.

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & - \\ \sigma_{12} & \sigma_2^2 & \dots & - \\ \vdots & \vdots & \ddots & \vdots \\ - & - & \dots & \sigma_f^2 \end{bmatrix} \quad \lambda_1 \geq \lambda_2 \dots \geq \lambda_f$$

$\downarrow \quad \quad \downarrow \quad \quad \quad \downarrow$
 $v_1 \quad \quad \quad v_2 \quad \quad \quad v_f$

What about reducing to 2 dimensions?

Variance is no longer a scalar which can be compared across all possible 2 dimensions where data can be projected.

We have an additional constraint that the co-variance terms of the reduced dimensions should be zero.

$$D_{n \times f} \quad \bigg| \quad D_{\text{trans}} = D \begin{matrix} \uparrow E_{f \times f} \\ \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \vec{v}_3 & \dots & \vec{v}_f \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 1 & 1 & & 1 \end{bmatrix} \end{matrix}$$

$$D_{\text{cov trans}} = D_{\text{trans}} E^T$$

What about reducing to 2 dimensions?

Variance is no longer a scalar which can be compared across all possible 2 dimensions where data can be projected.

We have an additional constraint that the co-variance terms of the reduced dimensions should be zero.

Fortunately, this happens if we keep going in the same sequence of the principal components.

What about reducing to 2 dimensions?

Variance is no longer a scalar which can be compared across all possible 2 dimensions where data can be projected.

We have an additional constraint that the co-variance terms of the reduced dimensions should be zero.

Fortunately, this happens if we keep going in the same sequence of the principal components.

We can now try to maximize the sum of the variance terms in the 2 projected directions.

So, it looks like that the most logical thing to do is to pick up the eigenvector corresponding to the second largest eigenvector. This is what we do in PCA.

Reconstructing the data

One could start with the objective that we want to minimize the reconstruction error, and we would get the same result as what we have described above.

Basic Ideas - Any k -dimensional vector can be represented in terms of k -orthonormal vectors.

Reconstructing the data

One could start with the objective that we want to minimize the reconstruction error, and we would get the same result as what we have described above.

Basic Ideas - Any k -dimensional vector can be represented in terms of k -orthonormal vectors.

Eg - Consider the vector $(2,3)$ in 2-D space and two orthonormal vectors $(1,0)$ and $(0,1)$

$$(2,3) = ((2,3) \cdot (1,0))(1,0) + ((2,3) \cdot (0,1))(0,1)$$

Reconstructing the data

One could start with the objective that we want to minimize the reconstruction error, and we would get the same result as what we have described above.

Basic Ideas - Any k -dimensional vector can be represented in terms of k -orthonormal vectors.

Eg - Consider the vector $(2,3)$ in 2-D space and two orthonormal vectors $(1,0)$ and $(0,1)$

$$(2,3) = ((2,3) \cdot (1,0))(1,0) + ((2,3) \cdot (0,1))(0,1)$$

Too trivial, isn't it??

Reconstructing the data

Eg - Consider the vector $(2,3)$ in 2-D space and two orthonormal vectors $(0.6,0.8)$ and $(-0.8,0.6)$

$$(2,3) = ((2,3) \cdot (0.6,0.8))(0.6,0.8) + ((2,3) \cdot (-0.8,0.6))(-0.8,0.6)$$

Reconstructing the data

Eg - Consider the vector (2,3) in 2-D space and two orthonormal vectors (0.6,0.8) and (-0.8,0.6)

$$(2,3) = ((2,3) \cdot (0.6,0.8))(0.6,0.8) + ((2,3) \cdot (-0.8,0.6))(-0.8,0.6)$$

Thus, in general for a k-dimensional vector \vec{d} , and k-orthonormal vectors $\vec{v}_1, \vec{v}_2, \vec{v}_3, \dots, \vec{v}_k$

$$\vec{d} = (\vec{d} \cdot \vec{v}_1)\vec{v}_1 + (\vec{d} \cdot \vec{v}_2)\vec{v}_2 + (\vec{d} \cdot \vec{v}_3)\vec{v}_3 + \dots + (\vec{d} \cdot \vec{v}_k)\vec{v}_k$$

PCA reconstruction

$$\vec{v}_1, \vec{v}_2, \vec{v}_3, \dots, \vec{v}_f$$

In PCA, once we have found out the eigenvalues and eigenvectors and projected the data in the lower dimensional space(f_s), the way to reconstruct for a data point back to the original dimensions (f) is -

$$\vec{d}_1 = (\vec{d}_1 \cdot \vec{v}_1) \vec{v}_1 + (\vec{d}_2 \cdot \vec{v}_2) \vec{v}_2 + \dots + (\vec{d}_f \cdot \vec{v}_f) \vec{v}_f$$
$$\vec{d}_2 = (\vec{d}_2 \cdot \vec{v}_1) \vec{v}_1 + (\vec{d}_2 \cdot \vec{v}_2) \vec{v}_2 + \dots + (\vec{d}_2 \cdot \vec{v}_f) \vec{v}_f$$
$$\vdots$$

$$\vec{d}_n = (\vec{d}_n \cdot \vec{v}_1) \vec{v}_1 + (\vec{d}_n \cdot \vec{v}_2) \vec{v}_2 + \dots + (\vec{d}_n \cdot \vec{v}_f) \vec{v}_f$$

In PCA, once we have found out the eigenvalues and eigenvectors and projected the data in the lower dimensional space(f_s), the way to reconstruct for a data point back to the original dimensions (f) is -

$$\vec{d}_{inv} = (\vec{d} \cdot \vec{v}_1)\vec{v}_1 + (\vec{d} \cdot \vec{v}_2)\vec{v}_2 + (\vec{d} \cdot \vec{v}_3)\vec{v}_3 + \dots + (\vec{d} \cdot \vec{v}_{f_s})\vec{v}_{f_s}$$

In PCA, once we have found out the eigenvalues and eigenvectors and projected the data in the lower dimensional space(f_s), the way to reconstruct for a data point back to the original dimensions (f) is -

$$\vec{d}_{inv} = (\vec{d} \cdot \vec{v}_1)\vec{v}_1 + (\vec{d} \cdot \vec{v}_2)\vec{v}_2 + (\vec{d} \cdot \vec{v}_3)\vec{v}_3 + \dots + (\vec{d} \cdot \vec{v}_{f_s})\vec{v}_{f_s}$$

This is essentially what we get after the matrix manipulations we saw in the previous session.

Percentage of Variance Captured in PCA

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & - & - & \sigma_{1f} \\ \sigma_{12} & \sigma_2^2 & & & \vdots \\ \vdots & & \ddots & & \vdots \\ \sigma_{1f} & \sigma_{2f} & - & & \sigma_f^2 \end{bmatrix} \quad \sigma_1^2 + \sigma_2^2 + \dots + \sigma_f^2$$

$$\lambda_1 + \lambda_2 + \dots + \lambda_f = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_f^2$$

Singular Value Decomposition - A way to approximate a general matrix of dimension $m \times n$.

$$M = \lambda_1 v_1 v_1^T + \lambda_2 v_2 v_2^T + \dots + \lambda_{1000} v_{1k} v_{1k}^T$$

Singular Value Decomposition - A way to approximate a general matrix of dimension $m \times n$.

Can we use PCA to approximate the matrix of size $n \times n$?

$$(1000) \times 5 + 5 = 5005$$

$$1000 \times 1000 = 10 \text{ lakh}$$

Singular Value Decomposition - A way to approximate a general matrix of dimension $m \times n$.

Can we use PCA to approximate the matrix of size $n \times n$? Can we do it always??

Let A be a general matrix of size $m \times n$.
What can be said about AA^T and $A^T A$

Let A be a general matrix of size $m \times n$.
What can be said about AA^T and $A^T A$

Thank you for your attention