



Lecture 13

MFDS Team



BITS Pilani

Pilani | Dubai | Goa | Hyderabad



- ▶ In the previous lecture, we discussed dimensionality reduction using PCA
- ▶ In this lecture, we will see the use of linear algebra in practical implementation of PCA.
- ▶ We will also study the challenges encountered when PCA is used in problems of larger dimensions.
- ▶ Finally, we elaborate the key steps of PCA in practice.



- ▶ We derived the matrix \mathbf{B} used in generation of lower-dimensional representation \mathbf{z} and the compressed data $\tilde{\mathbf{x}}$.
- ▶ The data covariance matrix \mathbf{S} was used to derive \mathbf{B}
- ▶ Recall that linear relationship connecting the original data \mathbf{x} , its low-dimensional code \mathbf{z} and the compressed data $\tilde{\mathbf{x}}$:
 $\mathbf{z} = \mathbf{B}^T \mathbf{x}$, and $\tilde{\mathbf{x}} = \mathbf{B} \mathbf{z}$.

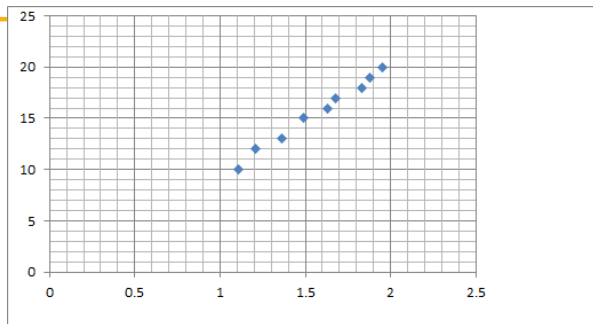
Consider the data given below:

x1	x2	x3	x4	x5	x6	x7	x8	x9	Mean
1.11	1.21	1.36	1.49	1.63	1.68	1.83	1.88	1.95	1.571
10	12	13	15	16	17	18	19	20	15.555

$$S = \frac{1}{9}XX^T = \begin{bmatrix} 0.079 & 0.888 \\ 0.888 & 10.024 \end{bmatrix}$$

The largest eigenvalue of S is $\lambda = 10.103$.

The corresponding eigenvector is $\begin{bmatrix} 0.088 \\ 0.996 \end{bmatrix}$.



The compressed or the reduced data is then given by $z = b_1^T X$

z1	z2	z3	z4	z5	z6	z7	z8	z9
-5.57	-3.57	-2.56	-0.56	0.45	1.45	2.46	3.46	4.46



- ▶ In the previous sections, we obtained the basis of the principal subspace as the eigenvectors that are associated with the largest eigenvalues of the data covariance matrix.

$$S = \frac{1}{N} \sum_{i=1}^N x_n x_n^T$$

- ▶ Equivalently we get

$$S = \frac{1}{N} X X^T$$

- ▶ Note that $X = [x_1, \dots, x_N] \in \mathbb{R}^{D \times N}$
- ▶ Note that X is a $D \times N$ matrix,

Eigenvector Computation and Low Rank Approximation



- ▶ To get the eigenvalues and the corresponding eigenvectors of S , we can follow two approaches
- ▶ We can perform an eigen-decomposition and compute the eigenvalues and eigenvectors of S directly.
- ▶ We can also use a singular value decomposition.
- ✓ ▶ Since S is symmetric and factorizes into XX^T , the eigenvalues of S are the squared singular values of X .
- ▶ Assume the SVD of X as $X = U\Sigma V^T$. Then

$$S = \frac{1}{N}XX^T = \frac{1}{N}U\Sigma\Sigma^T U^T$$

$$(U\Sigma V^T)(U\Sigma V^T)^T \\ U\Sigma \underbrace{V^T V}_I \Sigma^T U^T$$

dir of s.v \equiv First PCA comp.

max singular value \equiv max eigen val



col. of $U \equiv$ eigen val. of S

\Leftarrow eigen val. of $S =$ squared Sing. val. of X

Example of PCA Revisited



Consider the data given below:

x1	x2	x3	x4	x5	x6	x7	x8	x9	Mean
1.11	1.21	1.36	1.49	1.63	1.68	1.83	1.88	1.95	1.57
10	12	13	15	16	17	18	19	20	15.55

Consider $X = U\Sigma V^T$

$$U = \begin{bmatrix} 0.0883 & 0.9961 \\ 0.9961 & -0.0883 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 9.535 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0.084 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

- ✓▶ The columns of U are the eigenvectors of S .
- ▶ The eigenvalues λ_d of S are related to the singular values of X via

✓
$$\lambda_d = \frac{\sigma_d^2}{N}$$

$$S \approx \frac{1}{N} U \Sigma^2 U^T$$
$$S \approx U \begin{bmatrix} \frac{\sigma_1^2}{N} & & \\ & \frac{\sigma_2^2}{N} & \\ & & \ddots \end{bmatrix} U^T$$

- ▶ This relationship between the eigenvalues of S and the singular values of X provides the connection between the maximum variance view and the singular value decomposition.

- ▶ To maximize the variance of the projected data PCA chooses the columns of U to be the eigenvectors that are associated with the M largest eigenvalues of the data covariance matrix S
- ▶ The Eckart-Young theorem offers a direct way to estimate the low-dimensional representation.
- ▶ Consider the best rank-M approximation of X defined as \tilde{X}_M

$$\tilde{X}_M = \operatorname{argmin}_{\operatorname{rank}(A) \leq M} \|X - A\|_2$$

Handwritten notes:

- $\sum_{i=1}^n \sigma_i u_i v_i^T$ (with an arrow pointing to the $\|X - A\|_2$ term)
- $\sum_{i=1}^K \sigma_i u_i v_i^T$ (with an arrow pointing to the $\|X - A\|_2$ term)
- $\|X - A\|_2$ is labeled as "Spectral norm" (with an arrow pointing to the $\|X - A\|_2$ term)
- $\|X - A\|_2$ is labeled as "Frobenius norm" (with an arrow pointing to the $\|X - A\|_2$ term)



- ▶ The Eckart-Young theorem states that the best rank M approximation \tilde{X}_M is given by truncating the SVD at the top- M singular value.

$$\tilde{X}_M = U_M \Sigma_M V_M^T$$

- ▶ U_M is an orthogonal matrix
- ▶ V_M is an orthogonal matrix
- ▶ Σ_M has M largest singular values of X as diagonal entries



- ▶ Finding eigenvalues and eigenvectors is also important in other fundamental machine learning methods that require matrix decompositions
- ▶ In theory we can solve for the eigenvalues as roots of the characteristic polynomial.
- ▶ However for matrices larger than 4 by 4 this is not possible because we would need to find the roots of polynomial of degree 5 or higher.



- ▶ However the Abel-Ruffini theorem states that there exists no algebraic solution to this problem for polynomials of degree 5 or more.
- ▶ Therefore, in practice, solve for eigenvalues or singular values using iterative methods, which are implemented in all modern packages for linear algebra — *power method*
- ▶ In many applications we only require a few eigenvectors.



- ▶ It would be wasteful to compute the full decomposition, and then discard all eigenvectors with eigenvalues that are beyond the first few.
- ▶ It turns out that if we are interested in only the first few eigenvectors (with the largest eigenvalues), then iterative processes, which directly optimize these eigenvectors, are computationally more efficient than a full eigen-decomposition

↙ true in most Practical Cases



- ▶ In the extreme case of only needing the first eigenvector, a simple method called the power iteration is very efficient.
- ▶ Power iteration chooses a random vector x_0 that is not in the null space of S and follows the iteration for $k = 0, 1, \dots$

$$x_{k+1} = \frac{Sx_k}{\|Sx_k\|}$$

- ▶ This sequence of vectors converges to the eigenvector associated with the largest eigenvalue of S .



1. In order to do PCA, we need to compute the data covariance matrix.
2. In D dimensions, the data covariance matrix is a $D \times D$ matrix.
3. Computing the eigenvalues and eigenvectors of this matrix is computationally expensive as it scales cubically in D .
4. Therefore, PCA, as we discussed earlier, will be infeasible in very high dimensions.
5. In the following, we provide a solution to this problem for the case that we have substantially fewer data points than dimensions, i.e., $N \ll D$

- ▶ Assume we have a centered dataset x_1, \dots, x_N , where $x_i \in \mathbb{R}^D$.
- ▶ Then the data covariance matrix is given as $S = \frac{1}{N}XX^T$
- ▶ Consider the eigenvectors equation of S

$$Sb_m = \lambda_m b_m$$

- ▶ By substituting the definition of S

$$\frac{1}{N}XX^T b_m = \lambda_m b_m$$

↪ ①

\uparrow \nwarrow $N \times D$
 $D \times N$

i.e. $D \gg N$
(# features \gg # samples)

PCA Computation
is costly

- left*
▶ Multiply by X^T , we get
^

$$\frac{1}{N} X^T X X^T b_m = \lambda_m X^T b_m \quad - \bigcirc$$

- ▶ If $c_m = X^T b_m$, then

$$\frac{1}{N} X^T X c_m = \lambda_m c_m \quad - \bigcirc$$

eval
s' eval eval

- ▶ The nonzero eigenvalues of XX^T is same as the nonzero eigenvalues of $X^T X$.

λ_m from \bigcirc
 λ_m from \bigcirc

- Now that we have the eigenvectors of $\frac{1}{N}X^T X$,

$$\frac{1}{N}X^T X c_m = \lambda_m c_m$$

Handwritten notes:
③ even of $X^T X$ is c_m
→ use when $D \gg N$
Annotations: $N \times D$ (under X^T), $D \times N$ (under X)

- We need to derive the eigenvectors of XX^T , which we still need for PCA
- Multiply by X to get

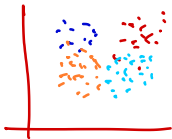
$$\frac{1}{N}XX^T X c_m = \lambda_m X c_m$$

Handwritten notes:
eigen vec. of XX^T is $X c_m$
Annotations: S (over X), S (over X^T)

- Here, we recover the data covariance matrix again. $\rightarrow \frac{1}{N} \sum x_i x_i^T$
- This means that we recover $X c_m$ as an eigenvector of S .

Different Views of PCA

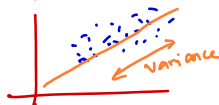
dim Reduction



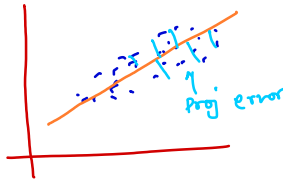
High dim. data
Projected in low dim

PCA

max. Variance



min Proj. error



Simple Auto
Encoder

HW

linear
layers





In the following, we will go through the individual steps of PCA using a running example.

- ▶ We are given a two dimensional dataset
- ▶ we want to use PCA to project it onto a one-dimensional subspace.
- ▶ The key steps are given below
 - ▶ Mean subtraction
 - ▶ Standardization
 - ▶ Eigen-decomposition of the covariance matrix
 - ▶ Projection



Mean subtraction

1. We start by centering the data by computing the mean of the dataset and subtracting it from every single data point.
- ✓ 2. This ensures that the dataset has mean 0.
3. Mean subtraction is not strictly necessary but reduces the risk of numerical problems

Key Steps of PCA in Practice



Standardization

	x_1	x_2	...
f_1	0	4	2 5
f_2	50	110	70 68

1. Divide the data points by the standard deviation σ_d of the dataset for every dimension d .
2. Now the data is unit free, and it has variance 1 along each axis, which is indicated by the standardization
3. This step completes the standardization of the data.

Figure: PCA

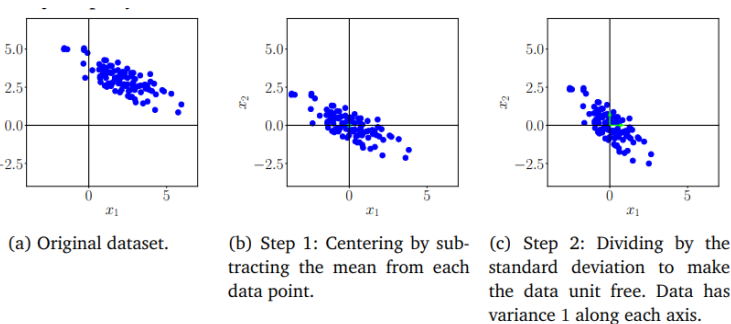


Figure 10.11 Steps of PCA. (a) Original dataset; (b) centering; (c) divide by standard deviation; (d) eigendecomposition; (e) projection; (f) mapping back to original data space.



Eigen-decomposition of the covariance matrix

1. Compute the data covariance matrix
2. Compute its eigenvalues and corresponding eigenvectors.
- ✓ 3. Since the covariance matrix is symmetric, the spectral theorem states that we can find an orthonormal basis of eigenvectors.
4. The eigenvectors are scaled by the magnitude of the corresponding eigenvalue.



Projection

- ✓ 1. We can project any data point $x_* \in \mathbb{R}^d$ onto the principal subspace:
- 2. To get this right, we need to standardize x_* using the mean and standard deviation of the training data in the d th dimension

$$x_*^{(d)} = \frac{x_*^{(d)} - \mu_d}{\sigma_d}, \quad d = 1, \dots, D$$

- ✓ 3. Here $x_*^{(d)}$ is the d th component of x_* .



1. We obtain the projection as

$$\tilde{x} = BB^T x_*$$

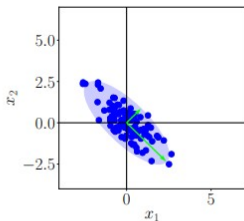
2. The coordinates are

$$z_* = B^T x_*$$

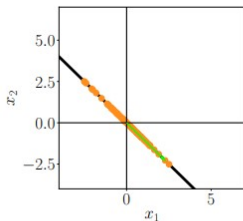
with respect to the basis of the principal subspace.

3. Here, B is the matrix that contains the eigenvectors that are associated with the largest eigenvalues of the data covariance matrix as columns.

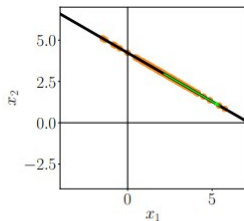
Key Steps of PCA in Practice



(d) Step 3: Compute eigenvalues and eigenvectors (arrows) of the data covariance matrix (ellipse).



(e) Step 4: Project data onto the principal subspace.



(f) Undo the standardization and move projected data back into the original data space from (a).



1. We derived PCA by maximizing the variance in the projected space
2. We took high-dimensional data $x \in \mathbb{R}^D$ and used a matrix B to find a lower-dimensional representation $z \in \mathbb{R}^M$
3. The columns of B are the eigenvectors of the data covariance matrix S that are associated with the largest eigenvalues.
4. Once we have a low-dimensional representation z , we can get a high-dimensional version of it as Bz .