



## Lecture 12

Math Foundations Team



**BITS Pilani**

Pilani | Dubai | Goa | Hyderabad

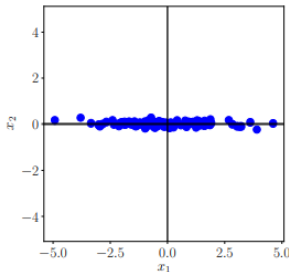


- ▶ We will look at principle components analysis and dimension reduction in this lecture.
- ▶ High-dimensional data is hard to visualize and interpret, can we project this data into lower dimensions while preserving the semantics of the data so as to draw the same conclusions as if we interpreted the higher dimensional data?
- ▶ Higher dimensional data is often overcomplete, in that there are redundant dimensions which can be explained by a combination of other dimensions.
- ▶ Dimensions in higher-dimensional data might be correlated, so the actual data may have an intrinsic lower-dimensional structure

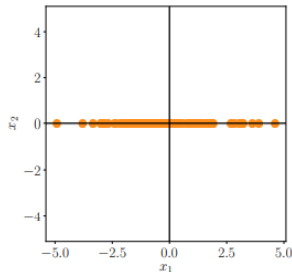
# Principle components analysis



- ▶ PCA is a technique for linear dimensionality reduction. It was first proposed by Pearson in 1900 and was independently rediscovered by Hotelling in 1933.



(a) Dataset with  $x_1$  and  $x_2$  coordinates.

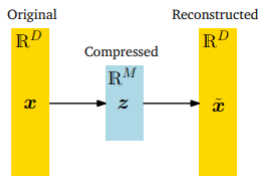


(b) Compressed dataset where only the  $x_1$  coordinate is relevant.



- ▶ **Given :** Let  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}, \mathbf{x}_n \in \mathbb{R}^D$  an independent, identically distributed dataset, with mean  $\mathbf{0}$ .
- ▶ Thus, the data covariance matrix  $\mathbf{S} = \frac{1}{N} \sum_{n=1}^{n=N} \mathbf{x}_n \mathbf{x}_n^T$ .
- ▶ **Aim :** To find projections  $\tilde{\mathbf{x}}_n \in U \subseteq \mathbb{R}^D$  of datapoints  $\mathbf{x}_n \in \mathbb{R}^D$  which are as similar as possible to the original datapoints but  $\dim(U) = M < D$ .
- ▶ We are looking for a lower-dimensional compressed representation  $\mathbf{z}_n$  of  $\mathbf{x}_n$  such that  $\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n$  where the projection matrix  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_m] \in \mathbb{R}^{D \times M}$ .
- ▶ The columns of  $\mathbf{B}$  are orthonormal which means  $\mathbf{b}_i^T \mathbf{b}_j = 0$  when  $i \neq j$  and  $\mathbf{b}_i^T \mathbf{b}_i = 1$ .

- ▶ The figure below shows how  $\mathbf{z}$  represents the lower-dimensional representation of the compressed data  $\tilde{\mathbf{x}}$  and plays the role of a bottleneck which controls the information flow between  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$ .



- ▶ There exists a linear relationship between 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  $\tilde{\mathbf{x}} = \mathbf{B} \mathbf{z}$  for a suitable matrix  $\mathbf{B}$ .



- ▶ We can interpret information content in the data as how "space-filling" it is and describe the information contained in the data by looking at the spread of the data.
- ▶ We can capture spread of the data using the concept of variance.
- ▶ PCA can then be viewed as a dimensionality reduction algorithm that maximizes the variance in the low-dimensional representation of the data to retain as much information as possible.
- ▶ Mathematically our aim is to find a matrix  $\mathbf{B}$  so that we can retain as much information as possible by projecting the data on the columns  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_M$  of the matrix.



- ▶ If  $\mu$  is the mean of the data. Centred data means that we work with data columns  $\mathbf{x} - \mu$ , rather than the original columns  $\mathbf{x}$ .
- ▶ Note that
$$\mathbb{V}_z(\mathbf{z}) = \mathbb{V}_x(\mathbf{B}^T(\mathbf{x} - \mu)) = \mathbb{V}_x(\mathbf{B}^T\mathbf{x} - \mathbf{B}^T\mu) = \mathbb{V}_x(\mathbf{B}^T\mathbf{x}).$$
- ▶ Thus by considering  $\mathbf{x} - \mu$ , the variance does not change. Therefore we assume that the data has a mean of  $\mathbf{0}$  for this lecture.
- ▶ Letting the mean be  $\mathbb{E}_x(\mathbf{x}) = \mathbf{0}$  means
$$\mathbb{E}_z(\mathbf{z}) = \mathbb{E}_x(\mathbf{B}^T\mathbf{x}) = \mathbf{B}^T\mathbb{E}_x(\mathbf{x}) = \mathbf{0}$$
- ▶ And the data covariance matrix  $\mathbf{S} = \frac{1}{N} \sum_{n=1}^{n=N} \mathbf{x}_n \mathbf{x}_n^T$ .



- ▶ We maximize the variance of the low-dimensional code by following a sequential approach.
- ▶ **Aim 1** : To maximize the variance  $V_1$  of the first coordinate  $z_{1n}$  of  $\mathbf{z} \in R^M$
- ▶ i.e to maximize  $V_1 = \mathbb{V}(z_1) = \frac{1}{N} \sum_{n=1}^N z_{1n}^2$  since the data  $\mathbf{x}$  is independent.
- ▶ Now  $z_{1n} = \mathbf{b}_1^T \mathbf{x}_n$ , and can be viewed as the orthogonal projection of  $\mathbf{x}_n$  onto the one-dimensional subspace spanned by  $\mathbf{b}_1$ .





$$\begin{aligned}\text{Then we have } V_1 &= \frac{1}{N} \sum_{n=1}^{n=N} (\mathbf{b}_1^T \mathbf{x}_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^{n=N} \mathbf{b}_1^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{b}_1 \\ &= \mathbf{b}_1^T \left( \sum_{n=1}^{n=N} \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1 \\ &= \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1.\end{aligned}$$

Arbitrarily increasing the magnitude of the vector  $\mathbf{b}_1$  will increase the variance - so we seek to maximize the variance subject to  $\|\mathbf{b}_1\| = 1$ .



- Therefore to find the direction  $\mathbf{b}_1$  that maximizes variance can be set up as a constrained optimization problem

$$\begin{aligned} \max \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 \text{ subject to} \\ \|\mathbf{b}_1\| = 1 \end{aligned}$$

- To solve this problem we set up the Lagrangian  $\mathcal{L}(\mathbf{x}, \lambda) = \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 + \lambda(1 - \mathbf{b}_1^T \mathbf{b}_1)$ .

- ▶ To solve the Lagrangian, let  $\frac{\partial \mathbb{L}}{\partial \lambda} = 0$  and  $\frac{\partial \mathbb{L}}{\partial \mathbf{b}_1} = \mathbf{0}$ .
- ▶ So, we get

$$\begin{aligned}\frac{\partial \mathbb{L}}{\partial \lambda} &= 1 - \mathbf{b}_1^T \mathbf{b}_1 = 0 \\ \frac{\partial \mathbb{L}}{\partial \mathbf{b}_1} &= 2\mathbf{b}_1^T \mathbf{S} - 2\lambda \mathbf{b}_1^T = 0\end{aligned}$$

- ▶ On simplification, we have  $\mathbf{S}\mathbf{b}_1 = \lambda \mathbf{b}_1$  and  $\mathbf{b}_1^T \mathbf{b}_1 = 1$ .
- ▶ Thus we find that the direction  $\mathbf{b}_1$  we seek is an eigenvector of the covariance matrix  $\mathbf{S}$  and  $\lambda$  is its corresponding eigenvalue.



- ▶ Putting the result of the previous slide into the objective function of the constrained optimization problem ie.  
 $\max \mathbf{b}_1^T \mathbf{S} \mathbf{b}_1$  we have  $\mathbf{b}_1^T \mathbf{S} \mathbf{b}_1 = \mathbf{b}_1^T \lambda \mathbf{b}_1 = \lambda$ .
- ▶ Our objective function boils to maximizing  $\lambda$  which means we are looking for the eigenvector of  $\mathbf{S}$  that corresponds to its largest eigenvalue.
- ▶ This is the first principal component.
- ▶ Let us now examine the inner workings of the Lagrangian method.

# Why does this method work?



- ▶ Suppose we have the following constrained optimization problem:

$$\begin{aligned} \max f(x, y) \text{ subject to} \\ g(x, y) = c \end{aligned}$$

- ▶ We note that at the optimal solution  $(x_0, y_0)$ , if we move a small distance  $(\delta x, \delta y)$ , we must continue to remain on the surface  $g(x, y) = c$ . This means  $g(x_0 + \delta x, y_0 + \delta y) = c = g(x, y)$ .
- ▶ But  $dg = g(x_0 + \delta x, y_0 + \delta y) - g(x, y) = \nabla g \cdot (\delta x, \delta y) = 0$ .
- ▶ Therefore  $\nabla g$  is orthogonal to the displacement vector  $(\delta x, \delta y)$ .

# Why does this method work?



- ▶ As we move along the displacement vector  $(\delta x, \delta y)$ , the value of the objective function also cannot change because otherwise we can get a better solution by moving along the displacement vector or its negative direction.
- ▶ Thus we have  $\nabla f \cdot (\delta x, \delta y) = 0$ , so  $\nabla f$  is orthogonal to the displacement vector  $(\delta x, \delta y)$ .
- ▶ Therefore,  $\nabla f$  and  $\nabla g$  are orthogonal to the displacement vector  $(\delta x, \delta y)$  and hence  $\nabla f = \lambda \nabla g$
- ▶ This leads us to  $\frac{\partial \mathbb{L}}{\partial \mathbf{b}_1} = 0$ .
- ▶  $\frac{\partial \mathbb{L}}{\partial \lambda} = 0$  merely enforces the constraint in the original constrained optimization problem.



- ▶ Assume that we have found the first  $m - 1$  principal components as the  $m - 1$  eigenvectors of  $S$  that are associated with the largest  $m - 1$  eigenvalues of  $S$ .
- ▶ Since  $S$  is symmetric we can use the spectral theorem to use the  $m - 1$  eigenvectors to construct an orthonormal basis of an  $m - 1$ -dimensional basis of  $\mathbb{R}^D$ .
- ▶ The  $m$ th principal component can be found by subtracting from the data the contribution of the first  $m - 1$  components  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{m-1}$ . Essentially we are trying to find principal components that compress the remainder of the information.



- ▶  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  contains the data points  $\mathbf{x}_k$  as column vectors.
- ▶ Then, new data matrix  $\hat{\mathbf{X}}$  is given by

$$\begin{aligned}\hat{\mathbf{X}} &= \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X} \\ &= \mathbf{X} - \mathbf{B}_{m-1} \mathbf{X}\end{aligned}$$

- ▶ Here  $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$  is a projection matrix that projects  $\mathbf{X}$  onto the subspace spanned by  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{m-1}$ .





- To find the  $m$ th principal component we maximize

$$\begin{aligned} V_m = \mathbb{V}[z_m] &= \frac{1}{N} \sum_{n=1}^{n=N} z_{mn}^2 \\ &= \frac{1}{N} \sum_{n=1}^{n=N} (\mathbf{b}_m^T \hat{\mathbf{x}}_n)^2 \\ &= \mathbf{b}_m^T \hat{\mathbf{S}} \mathbf{b}_m. \end{aligned}$$

- Here  $\hat{\mathbf{S}}$  is the data covariance matrix of the transformed data set represented by  $\hat{\mathbf{x}} = [\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N]$ .



- ▶ As before we set up a constrained optimization problem to find the first principal component, and establish that the optimal solution  $\mathbf{b}_m$  is the eigenvector of  $\hat{\mathbf{S}}$  that corresponds to the largest eigenvalue.
- ▶ We now establish that  $\mathbf{b}_m$  is an eigenvector of the original data matrix  $\mathbf{X}$ .
- ▶ More generally the sets of eigenvectors for  $\hat{\mathbf{S}}$  and  $\mathbf{S}$  are the same.

# Eigenvectors of $\mathbf{S}$ and $\hat{\mathbf{S}}$



- ▶ We now show that the eigenvectors of  $\mathbf{S}$  and  $\hat{\mathbf{S}}$  are the same.
- ▶ Let  $\mathbf{b}_i$  be an eigenvector of  $\mathbf{S}$ , i.e.  $\mathbf{S}\mathbf{b}_i = \lambda\mathbf{b}_i$ .
- ▶ Now we can write

$$\begin{aligned}\hat{\mathbf{S}}\mathbf{b}_i &= \frac{1}{N}(\mathbf{X} - \mathbf{B}_{m-1}\mathbf{X})(\mathbf{X} - \mathbf{B}_{m-1}\mathbf{X})^T \mathbf{b}_i \\ &= (\mathbf{S} - \mathbf{S}\mathbf{B}_{m-1}^T - \mathbf{B}_{m-1}\mathbf{S} + \mathbf{B}_{m-1}\mathbf{S}\mathbf{B}_{m-1}^T)\mathbf{b}_i \\ &= (\mathbf{S} - \mathbf{S}\mathbf{B}_{m-1} - \mathbf{B}_{m-1}\mathbf{S} + \mathbf{B}_{m-1}\mathbf{S}\mathbf{B}_{m-1})\mathbf{b}_i\end{aligned}$$

- ▶ Note that in the last line we have used the fact that  $\mathbf{B}_{m-1}$  is a projection matrix and is therefore symmetric.



- ▶ **Case 1:**  $i \geq m$ .
- ▶  $\mathbf{b}_i$  is an eigenvector not among the first  $m - 1$  components.
- ▶ Since  $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$  and  $\mathbf{b}_m$  is orthogonal to the  $\mathbf{b}_i, 1 \leq i \leq m - 1$ , we have  $\mathbf{B}_{m-1} \mathbf{b}_i = 0$ .
- ▶ Plugging this into the last equation on the previous slide, we see that  $\hat{\mathbf{S}} \mathbf{b}_i = (\mathbf{S} - \mathbf{B}_{m-1} \mathbf{S}) \mathbf{b}_i = \mathbf{S} \mathbf{b}_i = \lambda_i \mathbf{b}_i$ .
- ▶ Thus  $\mathbf{S} \mathbf{b}_m = \lambda_m \mathbf{b}_m$ .  $\lambda_m$  is the  $m$ th largest eigenvalue of  $\mathbf{S}$  and is also the largest eigenvalue of  $\hat{\mathbf{S}}$  because of the way the constrained optimization problem is set up.



- ▶ **Case 2:**  $i \leq m - 1$ .
- ▶ We have  $\mathbf{B}_{m-1}\mathbf{b}_i = \sum_{j=1}^{m-1} \mathbf{b}_j \mathbf{b}_j^T \mathbf{b}_i = \mathbf{b}_i$ .
- ▶ Plugging this into  $\hat{\mathbf{S}}\mathbf{b}_i = (\mathbf{S} - \mathbf{S}\mathbf{B}_{m-1} - \mathbf{B}_{m-1}\mathbf{S} + \mathbf{B}_{m-1}\mathbf{S}\mathbf{B}_{m-1})\mathbf{b}_i$ , we get  $\mathbf{S}\mathbf{b}_i = \mathbf{0}$ .
- ▶ Thus the vectors  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{m-1}$  are eigenvectors for  $\hat{\mathbf{S}}$  which are associated with the eigenvalue 0.

- ▶ Since  $V_m = \mathbf{b}_m^T \mathbf{S} \mathbf{b}_m = \lambda_m$ , we see that the variance of the data projected onto the  $m$ th principal component is  $\lambda_m$ .
- ▶ To find an  $M$ -dimensional subspace that retains as much information as possible, PCA tells us to choose the columns of the matrix  $\mathbf{B}$  as the  $M$  eigenvectors of the data covariance matrix  $\mathbf{S}$  that have the largest eigenvalues.



- ▶ We derived the PCA as an algorithm that maximizes the variance in the projected space to retain as much information as possible.
- ▶ Now we can also derive the PCA using a projection perspective to minimize the average reconstruction error. The original data is modeled as  $\mathbf{x}_n$  and the reconstruction is modeled as  $\tilde{\mathbf{x}}_n$ . We seek to minimize the distance between  $\mathbf{x}_n$  and  $\tilde{\mathbf{x}}_n$ .

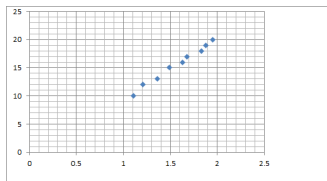
# An Example



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	1.57111111
10	12	13	15	16	17	18	19	20	15.555555

Data when plotted, we get





After subtracting mean from the data, then we get

$$\begin{aligned} S &= \frac{1}{9}XX' \\ &= \begin{pmatrix} 0.079498765 & 0.888271605 \\ 0.888271605 & 10.02469136 \end{pmatrix} \end{aligned}$$

The largest eigenvalue of  $S$  is 10.103 and the corresponding eigenvector with unit norm (first principal component) is  $b_1 = [0.088269, 0.996097]^T$ .

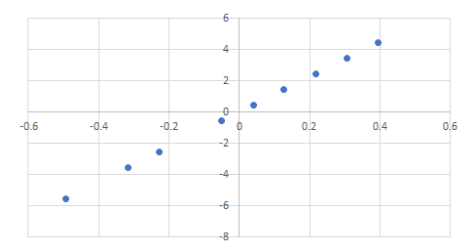
The compressed or the reduced data is then given by  $z = b_1^T X$

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

We can project the data onto the principal subspace by  $\tilde{X} = b_1 z$ . To obtain our projection in the original data space (i.e., before standardization), we need to undo the standardization.

$$\tilde{X} = \tilde{X} + \mu$$

# Projected data onto the Principal Subspace



# Projected data onto the Principal Subspace in the original data space

