

## Lecture 17: Dimensionality Reduction

Lecturer: Swaprava Nath

Scribe(s): SG33, SG34

**Disclaimer:** *These notes aggregate content from several texts and have not been subjected to the usual scrutiny deserved by formal publications. If you find errors, please bring to the notice of the Instructor.*

## 17.1 Dimensionality reduction

Dimensionality reduction involves condensing dataset information while preserving essential characteristics. The primary objective is to reduce algorithmic complexity and enhance model performance by removing redundant or irrelevant features without losing much information. Suppose we are given some input data, with dimension  $d$ , the problem we want to solve is whether we can reduce the dimension of the data to some  $m < d$ , without losing, or compromising a lot of the information.

We have two methods to achieve this:

- Unsupervised:  $D = \{x_i\} \implies$  Principal Component Analysis (PCA)
- Supervised:  $D = \{x_i, y_i\} \implies$  Linear Discriminant Analysis

## 17.2 Principle Component Analysis

Principal Component Analysis (PCA) is a fundamental technique that helps extract the most informative features and transform high-dimensional data into a more manageable low-dimensional form. It identifies a set of orthogonal axes called principal components (PCs) that capture the maximum variance upon orthogonal projection of the data onto these axes.

Say we have some data as shown in the figure. Let the vectors be  $u_1$  and  $u_2$  onto which projections are done.

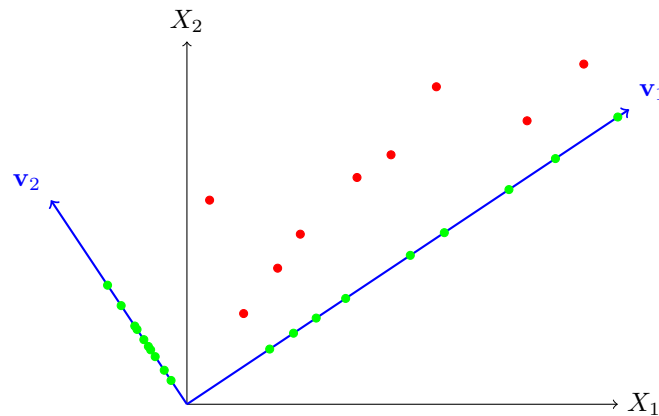


Figure 17.1: PCA illustration

We can say that  $v_1$  is a better basis vector than  $v_2$  since the variations among the data points in the original dataset are best depicted by the projections on  $v_1$  rather than  $v_2$ .

Consider a dataset  $D = \{x_1, x_2, \dots, x_n\}, x_i \in \mathbb{R}^d$

Our aim is to map dataset  $D$  onto  $m$  dimensions where  $m < n$ .

The question becomes which vectors to choose as our new basis vectors. To be able to capture the most important trends or patterns among the given data points, they have to be projected onto a vector that **maximises the variance of projected data**

### 17.2.1 Optimisation

Here, we're interested in getting a unit vector  $\mathbf{u}_1 \in \mathbb{R}^d$ , such that the variance of the projected data is maximum. And this vector will be our first PC. So here, for the sake of notation, we define the projection of  $\mathbf{x}_1$  on  $\mathbf{u}_1$  as :

$$\text{Proj}_{\mathbf{u}_1}(\mathbf{x}_1) = (\mathbf{x}_1^T \cdot \mathbf{u}_1) \mathbf{u}_1 = (\mathbf{u}_1^T \cdot \mathbf{x}_1) \mathbf{u}_1$$

where  $\mathbf{u}_1$  is a unit vector

Now, to calculate the variance, we get the mean of the projected data as:

$$\begin{aligned} \text{Mean of the projected data on } \mathbf{u}_1 &= \frac{1}{n} \sum_{i=1}^n \text{Proj}_{\mathbf{u}_1}(x_i) \\ &= \frac{1}{n} \sum_{i=1}^n (\mathbf{u}_1^T \cdot \bar{x}_i) \mathbf{u}_1 \\ &= (\mathbf{u}_1^T \cdot \bar{x}) \mathbf{u}_1 \end{aligned}$$

Now, we find the variance of the projected data. We know that each of the projected data points is  $(\mathbf{u}_1^T \cdot \mathbf{x}_i) \mathbf{u}_1$ . Since all projected points are in the direction of  $\mathbf{u}_1$ , and  $\mathbf{u}_1$  is a unit vector, we can therefore write variance in terms of distance of  $(\mathbf{u}_1^T \cdot \mathbf{x}_i) \mathbf{u}_1$  from the  $(\mathbf{u}_1^T \cdot \bar{x}) \mathbf{u}_1$  as :

$$\begin{aligned} \text{Variance of the projections} &= \frac{1}{n} \sum_{i=1}^n (u_1^T \cdot x_i - u_1^T \bar{x})^2 \\ &= \frac{1}{n} \sum_{i=1}^n u_1^T (x_i - \bar{x}) (x_i - \bar{x})^T u_1 \\ &= u_1^T S u_1 \quad , \text{ where } S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) (x_i - \bar{x})^T \end{aligned}$$

Here, we take a small digression to take a note of  $\mathbf{S}$ , we can write  $\mathbf{X} = [\mathbf{x}_1 - \bar{x} \quad \mathbf{x}_2 - \bar{x} \quad \dots \quad \mathbf{x}_n - \bar{x}]$ , where  $\mathbf{x}_i$ 's and  $\bar{x} \in \mathbb{R}^d$ . We have  $\mathbf{S}' = \mathbf{X}\mathbf{X}^T$ , and we get that

$$\mathbf{S}'(\mathbf{i}, \mathbf{j}) = \sum_{k=1}^n (x_{ki} - \bar{x}_i) (x_{kj} - \bar{x}_j)$$

Now, also we have,

$$\mathbf{S}(\mathbf{i}, \mathbf{j}) = \frac{1}{n} \sum_{k=1}^n (x_{ki} - \bar{x}_i) (x_{kj} - \bar{x}_j)$$

Therefore, we have  $\mathbf{S} = \frac{1}{n}\mathbf{X}\mathbf{X}^T$ , hence  $\mathbf{S}$  is covariance matrix, of the vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$

Now, we are back to the original discussion, that is, getting a unit vector  $\mathbf{u}_1$  such that the variance of the projected data is maximum. So far, we have reduced the optimisation to the following problem:

$$\begin{aligned} \max_{\mathbf{u}_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \\ \text{st } \mathbf{u}_1^T \mathbf{u}_1 = 1 \end{aligned}$$

Now the question comes: is it a convex optimization problem? Here the optimisation function  $f(\mathbf{u}_1) = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ , is convex because  $\nabla^2 f(\mathbf{u}_1) = 2\mathbf{S} = \frac{2}{n}\mathbf{X}\mathbf{X}^T$  is positive semi definite. The problem is not convex because the feasible region itself is not convex, and instead of minimizing it, we are maximizing it. So, it looks like we will find another way out.

### 17.2.1.1 Orthonormal Eigenvectors

Before we state the main theorem that we are going to use, we first introduce a few terminologies:

- For a matrix  $\mathbf{A}$ , we say a vector  $\mathbf{v}$  to be eigenvector with eigenvalue  $\lambda$ , if following holds:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

- If two vectors  $\mathbf{v}, \mathbf{u}$ , satisfies :

$$\mathbf{v}^T \cdot \mathbf{v} = 1, \mathbf{u}^T \cdot \mathbf{u} = 1 \text{ and } \mathbf{v}^T \cdot \mathbf{u} = \mathbf{u}^T \cdot \mathbf{v} = 0$$

then we say, the vectors  $\mathbf{v}, \mathbf{u}$  as orthonormal vectors.

**Theorem** ( *Spectral theorem for real matrices* )

A real symmetric matrix  $\mathbf{A}$  has a **orthonormal** basis of  $\mathbb{R}^d$ , formed by eigenvectors of  $\mathbf{A}$ .

Although we ain't going to use the actual version of the theorem, instead, we use the equivalent version, which states that *every real symmetry matrix  $\mathbf{A}$  has  $d$  orthonormal eigenvectors.*

### 17.2.1.2 Optimising the reduced problem

Here, we can observe that  $\mathbf{S}$  is a real symmetry matrix,

$$\mathbf{S}^T = \left( \frac{1}{n} \mathbf{X} \mathbf{X}^T \right)^T = \frac{1}{n} (\mathbf{X} \mathbf{X}^T)^T = \frac{1}{n} (\mathbf{X}^T)^T \mathbf{X}^T = \frac{1}{n} \mathbf{X} \mathbf{X}^T = \mathbf{S}$$

Now, using the theorem stated above, we can say that  $\exists \mathbf{v}_i, \lambda_i$  such that ,

$$\begin{aligned} \mathbf{S}\mathbf{v}_i = \lambda_i \mathbf{v}_i \quad \text{where } \begin{cases} \mathbf{v}_i^T \cdot \mathbf{v}_j = 0 & \text{if } i \neq j, i, j \in [1, 2, \dots, d] \\ \mathbf{v}_i^T \cdot \mathbf{v}_i = 1 & \text{if } i \in [1, 2, \dots, d] \end{cases} \\ \text{for } i = 1, 2, \dots, d \end{aligned}$$

Which can be equivalently rewritten as:

$$\begin{aligned} \mathbf{S} [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_d] &= [\lambda_1 \mathbf{v}_1 \ \lambda_2 \mathbf{v}_2 \ \cdots \ \lambda_d \mathbf{v}_d] \\ &= [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_d] \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_d \end{bmatrix} \end{aligned}$$

Denoting  $[\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_d]$  by  $\mathbf{V}$ , and by the condition that  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$  are orthonormal vectors we have that

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}$$

we have,

$$\mathbf{S} \mathbf{V} = \mathbf{V} \boldsymbol{\Sigma}$$

where,

$$\boldsymbol{\Sigma} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_d \end{bmatrix}$$

Therefore,

$$\mathbf{S} = \mathbf{V} \boldsymbol{\Sigma} \mathbf{V}^T$$

Now, putting this back into the reduced optimization problem, we have the following optimisation problem:

$$\begin{aligned} \max_{\mathbf{u}_1} \mathbf{u}_1^T \mathbf{V} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{u}_1 \\ \text{st } \mathbf{u}_1^T \mathbf{u}_1 = 1 \end{aligned}$$

Since the vectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$  form a basis of  $\mathbb{R}^d \exists \alpha_1, \alpha_2, \dots, \alpha_d$  such that  $\mathbf{u}_1 = \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \dots + \alpha_d \mathbf{v}_d$  which is equivalently,

$$\mathbf{u}_1 = \mathbf{V} \boldsymbol{\alpha}^T \text{ where } \boldsymbol{\alpha} = [\alpha_1 \ \alpha_2 \ \cdots \ \alpha_d]$$

$\therefore$  Our optimisation problem reduces as:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \boldsymbol{\alpha} (\mathbf{V}^T \mathbf{V}) \boldsymbol{\Sigma} (\mathbf{V}^T \mathbf{V}) \boldsymbol{\alpha}^T \quad \text{st } \boldsymbol{\alpha} (\mathbf{V}^T \mathbf{V}) \boldsymbol{\alpha}^T \\ \max_{\boldsymbol{\alpha}} \boldsymbol{\alpha} \boldsymbol{\Sigma} \boldsymbol{\alpha}^T \quad \text{st } \boldsymbol{\alpha} \boldsymbol{\alpha}^T \\ \max_{\boldsymbol{\alpha}} \sum_{i=1}^d \alpha_i^2 \lambda_i \quad \text{st } \sum_{i=1}^d \alpha_i^2 = 1 \end{aligned}$$

Suppose,  $\lambda_j = \max(\lambda_1, \lambda_2, \dots, \lambda_d)$ , then clearly the solution to our optimisation is  $\boldsymbol{\alpha}$  such that  $\alpha_i = \mathbb{1}(i = j) \forall i \in [1, 2, \dots, d]$ , we summarise the result for first **PC**  $\mathbf{u}_1$  as follow:

$\begin{aligned} \max_{\mathbf{u}_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 \text{ st } \mathbf{u}_1^T \cdot \mathbf{u}_1 = 1 \text{ and } \mathbf{S} = \mathbf{S}^T \\ \implies \mathbf{u}_1 = \mathbf{v}_j \text{ where } \lambda_j = \max(\lambda_1, \lambda_2, \dots, \lambda_d) \text{ and } \mathbf{S} \mathbf{v}_i = \lambda_i \mathbf{v}_i \forall i \in [1, 2, \dots, d] \end{aligned}$
--

### 17.2.2 Generalisation for reduction to $m$ dimensions

After we have the first PC  $\mathbf{u}_1$ , we can find the second PC  $\mathbf{u}_2$ , by solving the following optimisation problem:

$$\begin{aligned} \max_{\mathbf{u}_1} \mathbf{u}_2^T \mathbf{S} \mathbf{u}_2 \text{ st } \mathbf{u}_2^T \cdot \mathbf{u}_2 = 1 \text{ and } \mathbf{u}_1^T \cdot \mathbf{u}_2 = 0 \\ \implies \mathbf{u}_2 = \mathbf{v}_k \text{ st } \lambda_k = \max \{ \{ \lambda_1, \lambda_2, \dots, \lambda_d \} \setminus \{ \lambda_j \} \} \\ \text{where } \lambda_j = \max \{ \lambda_1, \lambda_2, \dots, \lambda_d \}, \mathbf{S} \mathbf{v}_i = \lambda_i \mathbf{v}_i \forall i \in [1, 2, \dots, d] \end{aligned}$$

In other words, if we order the eigenvectors according to their eigenvalues, then  $\mathbf{u}_1$  is the largest eigenvector, and  $\mathbf{u}_2$  is the second largest eigenvector of  $\mathbf{S}$ .

To generalize in the end, to reduce to  $m$  dimension, we just take the projection of the data points to the first  $m$  eigenvectors of  $\mathbf{S}$ , where eigenvectors are being ordered for the eigenvalues. In a nutshell, take the first  $m$  largest eigenvectors, say  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ , and then to reduce the dimension to  $m$ , take the projection of each of the data points  $x_i \in \mathbb{R}^d$ , in  $\text{Span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m\}$ .

### 17.2.3 PCA Algorithm

1. Compute mean of data points,  $\bar{x}$
2. Mean centre the data i.e., compute  $x_i - \bar{x}, \forall i$
3. Compute the covariance matrix,  $S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$
4. Do eigenvalue decomposition of  $S = V \Sigma V^T$
5. Pick  $m$  top eigenvectors (corresponding to largest eigenvalues):  $u_1, \dots, u_m \rightarrow$  Principal Components
6. Create a projection matrix,  $U = [u_1, \dots, u_m]$
7. For each  $x_i \in D$ ,  $U^T x_i$  will give the projection on  $m$ -dimensional space

Note that for  $U^T x$  is a  $m$ -dimensional vector.

$$U^T x = \begin{bmatrix} u_1^T x \\ u_2^T x \\ \vdots \\ u_m^T x \end{bmatrix}$$

### 17.3 Supervised Learning

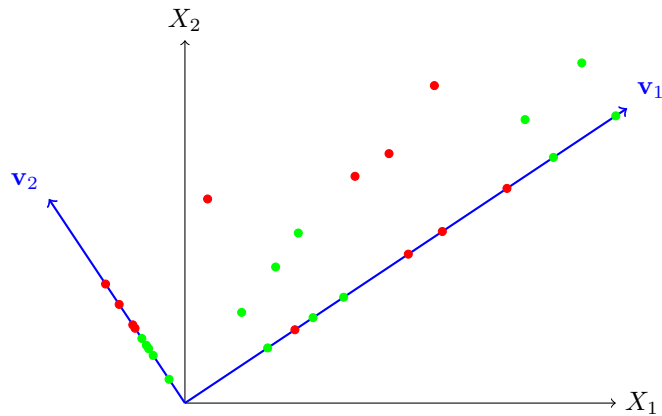


Figure 17.2: PCA's Blindness towards Class Discrimination

Utilizing the vector  $v_1$  from PCA for data projection results in the loss of class distinctions. Thus, employing vector  $v_2$  would be more advantageous in this context.

As we can see, PCA fails to preserve the intrinsic class characteristics. We need a projection that provides us with more resolution between classes.

We need a good measure to evaluate these projections.

$$J(u) = |\mu_1^T u - \mu_2^T u| \text{ where } \mu_i \text{ is the mean of class } i$$

$J(u)$  is one possible measure, but it alone is not good enough since it only considers the means of various classes.

#### Objective:

- Different classes are well separated.
- Data within the same class is not well separated.