

## Principal Component Analysis:

A very basic ML algorithm that can be used for dimensionality reduction. This algorithm performs a lossy compression on given data that can be used to reduce the amount of storage the data takes up. It can also be used to learn representation of data that is non-linear.

### Notations and definitions:

$n$ : Scalar denoting the dimension of original data.

$m$ : Scalar denoting number of examples.

$l$ : Scalar denoting the dimension after PCA is performed on data

$\mathbf{X}$ : The feature matrix of dimension  $m \times n$ . Every row of this matrix is a single vector  $\mathbf{x}_i$  denoting the  $i^{\text{th}}$  example that is  $n$ -dimension.

$\mathbf{C}$ : The code matrix of dimension  $m \times l$ . Every row of this matrix is a single vector  $\mathbf{c}_i$  denoting the  $i^{\text{th}}$  example that is of  $l$ -dimension.

$f(\mathbf{x})$ : The encoder function. This is the first function we need to derive to perform the dimensionality reduction of  $\mathbf{X}$  to get  $\mathbf{C}$ . Basically,  $\mathbf{c}_i = f(\mathbf{x}_i)$ .

$g(\mathbf{c})$ : The decoder function. This is the second function we need to derive to perform the dimensional expansion of  $\mathbf{C}$  to obtain  $\mathbf{X}$  back. Here, we will define a matrix  $\mathbf{D}$  to estimate the functions, listed next.  $\mathbf{x}_i = g(\mathbf{c}_i) = g(f(\mathbf{x}_i))$ .

$\mathbf{D}$ : A matrix of shape  $n \times l$  that provides the mapping between  $\mathbf{C}$  and  $\mathbf{X}$ , thus defining the decoding function as  $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$ . We restrict the columns of  $\mathbf{D}$  to stay orthogonal to each other. This will simplify our calculation because our assumption makes  $\mathbf{D}^T = \mathbf{D}^{-1}$ .

### Derivation of the encoder function, $f(\mathbf{x})$ :

To find the encoder function, we need to find the optimal code point for every  $\mathbf{x}_i$  because  $\mathbf{c}_i = f(\mathbf{x}_i)$ . We can minimize the Euclidean distance between the original data and the data computed after we decode the code point. This distance can be given using the L2 norm. Further, we will get rid of the  $i$  subscript:

$$\|\mathbf{x} - g(\mathbf{c})\|_2$$

Since this needs to be minimized, we will just be proceeding with the square of the L2 norm because minimizing the square will give us the same optimal point for  $\mathbf{c}$ :

$$(\|\mathbf{x} - g(\mathbf{c})\|_2)^2$$

$$= (\mathbf{x} - g(\mathbf{c}))^T (\mathbf{x} - g(\mathbf{c}))$$

$$((\|\mathbf{x}\|_2)^2 = \mathbf{x}^T \mathbf{x})$$

$$= \mathbf{x}^T \mathbf{x} - \mathbf{x}^T g(\mathbf{c}) - g(\mathbf{c})^T \mathbf{x} + g(\mathbf{c})^T g(\mathbf{c})$$

(distributive property)

$$= \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T g(\mathbf{c}) + g(\mathbf{c})^T g(\mathbf{c})$$

( $g(\mathbf{c})^T \mathbf{x} = (\mathbf{x}^T g(\mathbf{c}))^T$ ; transpose of a scalar)

$$= \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{D}^T \mathbf{D}\mathbf{c}$$

(substituting  $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$ )

$$= \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{c}$$

( $\mathbf{D}^T \mathbf{D} = \mathbf{I}$ ; read definition of  $\mathbf{D}$  above)

The next step is to differentiate this function w.r.t.  $\mathbf{c}$  and equating it to 0 to obtain the value of  $\mathbf{c}$  at minima:

$$\nabla_{\mathbf{c}} (\mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{D} \mathbf{c} + \mathbf{c}^T \mathbf{c}) = 0$$

$$\Rightarrow \nabla_{\mathbf{c}} \mathbf{x}^T \mathbf{x} - \nabla_{\mathbf{c}} (2\mathbf{x}^T \mathbf{D} \mathbf{c}) + \nabla_{\mathbf{c}} \mathbf{c}^T \mathbf{c} = 0$$

$$\Rightarrow 0 - 2\mathbf{D}^T \mathbf{x} + 2\mathbf{c} = 0$$

$$\Rightarrow \mathbf{c} = \mathbf{D}^T \mathbf{x}$$

For the matrices  $\mathbf{X}$  and  $\mathbf{C}$ ,

$$\mathbf{C} = \mathbf{X} \mathbf{D}$$

(rearranging for the feature matrix  $\mathbf{X}$ )

Or,

$$f(\mathbf{x}) = \mathbf{x} \mathbf{D}$$

Now we can easily encode the data and perform the dimensionality reduction. However, for both the encoder and the decoder functions, we need the numerical value of  $\mathbf{D}$ .

### Deriving the value of $\mathbf{D}$ :

Again, we will be minimizing the distance between the input and their reconstructions. However, since the value of  $\mathbf{D}$  is not trivial anymore, we shall minimize the norm across all values of  $\mathbf{X}$ . The square of the Frobenius norm can be used for this (basically an L2 norm including all elements of the matrix):

$$\|\mathbf{X} - g(\mathbf{C})\|_F \quad \text{s.t. } \mathbf{D}^T \mathbf{D} = \mathbf{I}$$

$$= (\|\mathbf{X} - g(f(\mathbf{X}))\|_F)^2$$

$$(\mathbf{C} = f(\mathbf{X}))$$

$$= (\|\mathbf{X} - g(\mathbf{X} \mathbf{D})\|_F)^2$$

$$(f(\mathbf{X}) = \mathbf{X} \mathbf{D})$$

$$= (\|\mathbf{X} - \mathbf{X} \mathbf{D} \mathbf{D}^T\|_F)^2$$

$$(g(\mathbf{D}^T \mathbf{X}) = \mathbf{X} \mathbf{D} \mathbf{D}^T)$$

$$= \text{Tr}((\mathbf{X} - \mathbf{X} \mathbf{D} \mathbf{D}^T)^T (\mathbf{X} - \mathbf{X} \mathbf{D} \mathbf{D}^T))$$

$$(\|\mathbf{X}\|_F = \sqrt{\text{Tr}(\mathbf{X}^T \mathbf{X})}; \text{Tr represents the Trace operator})$$

$$= \text{Tr}(\mathbf{X}^T \mathbf{X} - \mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{D}^T - \mathbf{D} \mathbf{D}^T \mathbf{X}^T \mathbf{X} + \mathbf{D} \mathbf{D}^T \mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{D}^T)$$

(Distributive property; maintaining shape along terms)

$$= \text{Tr}(\mathbf{X}^T \mathbf{X}) - \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{D}^T) - \text{Tr}(\mathbf{D} \mathbf{D}^T \mathbf{X}^T \mathbf{X}) + \text{Tr}(\mathbf{D} \mathbf{D}^T \mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{D}^T)$$

(Trace is distributive)

$$= \text{Tr}(\mathbf{X}^T \mathbf{X}) - 2\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{D}^T) + \text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{D}^T \mathbf{D} \mathbf{D}^T)$$

(Terms inside the Trace operator can be interchanged)

$$= \text{Tr}(\mathbf{X}^T \mathbf{X}) - 2\text{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{D}^T) + \text{Tr}(\mathbf{X}^T \mathbf{X})$$

$$(\mathbf{D}^T \mathbf{D} = \mathbf{I})$$

$$= 2\text{Tr}(\mathbf{X}^T \mathbf{X}) - 2\text{Tr}(\mathbf{D}^T \mathbf{X}^T \mathbf{X} \mathbf{D})$$

To minimize this function and find an optimal point for  $\mathbf{D}$ , we can ignore terms not involving  $\mathbf{D}$  and constants:

$$\min -\text{Tr}(\mathbf{D}^T \mathbf{X}^T \mathbf{X} \mathbf{D})$$

$$= \max \text{Tr}(\mathbf{D}^T \mathbf{X}^T \mathbf{X} \mathbf{D})$$

------(1)

---

To maximize this, we can exploit the eigen decomposition for matrices:

$$\mathbf{A} = \mathbf{Q}\mathbf{L}\mathbf{Q}^{-1}$$

Where  $\mathbf{Q}$  is an orthogonal matrix of with columns as vectors representing the different eigenvalues of  $\mathbf{A}$  and  $\mathbf{L}$  is a diagonal matrix where  $L_{i,i}$  is the eigenvalue corresponding to the  $i^{\text{th}}$  column eigenvector of  $\mathbf{Q}$ .

Rearranging,

$$\mathbf{L} = \mathbf{Q}^T \mathbf{A} \mathbf{Q} \quad (\mathbf{Q} \text{ is orthogonal})$$

---

This implies that in the equation (1) to be maximized above, when a column  $\mathbf{D}_i$  is the eigenvector of  $\mathbf{X}^T \mathbf{X}$ , the value that the equation takes is the eigenvalue of  $\mathbf{X}^T \mathbf{X}$  and to maximize it, we need the  $I$  maximum eigenvalues of  $\mathbf{X}^T \mathbf{X}$ .

Therefore, the required value of  $\mathbf{D}$  consists of  $I$  eigenvectors corresponding to the  $I$  maximum eigenvalues of  $\mathbf{X}^T \mathbf{X}$  as its columns. The eigenvectors of  $\mathbf{X}^T \mathbf{X}$  can be calculated easily by finding the right-singular eigenvectors of  $\mathbf{X}$  by performing Singular Value Decomposition on  $\mathbf{X}$ .

Once  $\mathbf{D}$  is calculated, the algorithm is as listed below.

### The algorithm

SVD( $\mathbf{X}$ ): Calculates the singular value decomposition of  $\mathbf{X}$  (returns left-singular eigenvectors, eigenvalues and right-singular eigenvectors)

- $\mathbf{U}, \mathbf{S}, \mathbf{V} := \text{SVD}(\mathbf{X})$

-Sort  $\mathbf{S}$  and correspondingly,  $\mathbf{V}$  (descending)

- $\mathbf{D}$  = matrix of first  $I$  columns of  $\mathbf{V}$

To encode:  $\mathbf{C} = \mathbf{X} \times \mathbf{D}$ ;    To decode:  $\mathbf{X} = \mathbf{C} \times \mathbf{D}^T$  (rearranging for the feature matrix  $\mathbf{X}$ )

### Minimum I:

As  $I$  decreases, the compression will be lossier, i.e., when decoded, the error will be higher. An obvious way of finding minimum  $I$  is by finding the value of  $I$  where the disparity between decoded value and original data exceeds a small error (say 1%). An easier way is finding the number  $I$  for which the sum of  $I$  maximum eigenvalues accounts for 99% of the sum over all eigenvalues (here,  $\mathbf{s}$  is the vector of eigenvalues obtained by SVD):

$$\frac{\sum_{i=1}^I S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

\*inspired by Deep Learning by Ian Goodfellow

\*scale the data while coding

\*example code in Python provided