

Dimension Reduction

Sunmook Choi felixchoi@korea.ac.kr

August 24, 2023

1 Principal Component Analysis

1.1 Prerequisites

Definition 1 (Orthogonal Complement). *Let V be a vector space and U be its subspace. The orthogonal complement of U is the set*

$$U^\perp = \{\mathbf{v} \in V : \langle \mathbf{u}, \mathbf{v} \rangle = 0, \forall \mathbf{u} \in U\}. \quad (1.1)$$

Then, the set U^\perp is a subspace of V .

Definition 2 (Direct Sum). *Let U_1, U_2 be subspaces of V . For each $\mathbf{v} \in V$, if there exist $\mathbf{u}_1 \in U_1$ and $\mathbf{u}_2 \in U_2$ uniquely such that $\mathbf{v} = \mathbf{u}_1 + \mathbf{u}_2$, then V is the direct sum of U_1 and U_2 , and we write $V = U_1 \oplus U_2$.*

Theorem 3. *If U is a subspace of \mathbb{R}^n , then the following hold.*

- $\dim(U) + \dim(U^\perp) = n$
- $\mathbb{R}^n = U \oplus U^\perp$ (orthogonal decomposition of \mathbb{R}^n)
- $(U^\perp)^\perp = U$

Definition 4 (Projection). *Let V be a vector space. If $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$ is an orthonormal basis for the subspace U of V , then the orthogonal projection of $\mathbf{v} \in V$ onto U is the vector*

$$\text{proj}_U \mathbf{v} = \langle \mathbf{v}, \mathbf{u}_1 \rangle \mathbf{u}_1 + \dots + \langle \mathbf{v}, \mathbf{u}_k \rangle \mathbf{u}_k. \quad (1.2)$$

Definition 5 (Orthogonal Decomposition). *Let V be a vector space with $\dim(V) = n$ and let U be its subspace. Let $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$ and $\{\mathbf{u}_{k+1}, \dots, \mathbf{u}_n\}$ be orthonormal bases for U and*

U^\perp , respectively. Then, the orthogonal decomposition of $\mathbf{v} \in V$ is

$$\mathbf{v} = \sum_{i=1}^k a_i \mathbf{u}_i + \sum_{i=k+1}^n b_i \mathbf{u}_i \in U \oplus U^\perp \quad (1.3)$$

where $a_i = \langle \mathbf{v}, \mathbf{u}_i \rangle$ for $i = 1, \dots, k$ and $b_i = \langle \mathbf{v}, \mathbf{u}_i \rangle$ for $i = k+1, \dots, n$. Here, $\sum_{i=1}^k a_i \mathbf{u}_i$ and $\sum_{i=k+1}^n b_i \mathbf{u}_i$ are the orthogonal projections onto U and U^\perp , respectively.

Definition 6 (Variance and Covariance). For a random variable X , the variance of X is defined to be

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2. \quad (1.4)$$

For two random variables S and Y , the covariance of X and Y is defined to be

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]. \quad (1.5)$$

For a random vector $X = (X_1, \dots, X_n)^\top$, the covariance matrix Σ is defined whose (i, j) -entry is defined to be

$$\Sigma_{ij} = \Sigma_{ji} = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]. \quad (1.6)$$

Definition 7 (Sample Covariance Matrix). Let $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ be a dataset, that is, let $\mathbf{x}_i \in \mathbb{R}^d$ be a sample data for $i = 1, \dots, n$. Let $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_d)$ be a sample mean vector of X such that $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$. Then the sample covariance matrix \mathcal{S} is defined to be

$$\mathcal{S} = \frac{1}{n-1} (X - \bar{\mathbf{x}})^\top (X - \bar{\mathbf{x}}) \in \mathbb{R}^{d \times d}. \quad (1.7)$$

If the population mean of X_j is known as μ_j , then the sample covariance is defined to be

$$\mathcal{S} = \frac{1}{n} (X - \boldsymbol{\mu})^\top (X - \boldsymbol{\mu}) \in \mathbb{R}^{d \times d} \quad (1.8)$$

where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)$.

Remark 8. The denominator $n-1$ in Eq. 7 is due to Bessel's correction, which makes the estimator unbiased. A sample covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$ is positive semi-definite, that is, $\mathbf{z}^\top \Sigma \mathbf{z} \geq 0$ for any $\mathbf{z} \in \mathbb{R}^d$.

1.2 Principal Components

For any sample covariance matrix $\mathcal{S} \in \mathbb{R}^{d \times d}$, it is positive semi-definite and real symmetric. Positive semi-definiteness implies that its eigenvalues are nonnegative, and real symmetry implies that it is orthogonally diagonalizable.

Let $\lambda_1, \dots, \lambda_d$ be eigenvalues of \mathcal{S} such that $\lambda_1 \geq \dots \geq \lambda_d \geq 0$, and let \mathbf{v}_i 's be orthonormal eigenvectors corresponding to λ_i 's for $i = 1, \dots, d$. Here, we call the eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ the top k principal components of X where $k \leq d$.

1.2.1 Problem Setup

Given a dataset $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ and a positive integer $k \leq d$, find the best linear projection \tilde{X} of X onto a lower dimensional subspace U of \mathbb{R}^d , $\dim(U) = k$. Here, the ‘best’ linear projection $\tilde{X} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n]^\top$ is the linear projection which minimizes

$$\frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{x}}_i - \mathbf{x}_i\|^2. \quad (1.9)$$

The compressed/encoded data can be expressed in two ways: $\mathbf{z}_i \in U$ and $\tilde{\mathbf{x}}_i \in \mathbb{R}^d$ which lie in different dimensional spaces. In this sense, $\tilde{\mathbf{x}}_i \in \mathbb{R}^d$ is called a reconstructed data. Note that PCA can be also considered finding \tilde{X} which has the smallest reconstruction loss. We additionally assume that the population mean of each dimension of data is zero. Then, the sample covariance \mathcal{S} of X will be $\mathcal{S} = \frac{1}{n} X^\top X = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i$.

1.2.2 Projection Perspective

Claim 1. *The vector $\tilde{\mathbf{x}}_i$ is the projection of \mathbf{x}_i onto the subspace spanned by top k principal components of X .*

Proof. Let $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ be a set of orthonormal vectors and U be the spanned subspace. The projection $\tilde{\mathbf{x}}_i$ of \mathbf{x}_i onto U is

$$\tilde{\mathbf{x}}_i = \sum_{j=1}^k (\mathbf{x}_i \mathbf{v}_j^\top) \mathbf{v}_j, \quad (1.10)$$

and we want to find

$$\mathcal{B}^* = \arg \min_{\mathcal{B}} \frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{x}}_i - \mathbf{x}_i\|^2. \quad (1.11)$$

By simple computation, we have the following.

$$\|\tilde{\mathbf{x}}_i - \mathbf{x}_i\|^2 = \|\tilde{\mathbf{x}}_i\|^2 - 2\tilde{\mathbf{x}}_i^\top \mathbf{x}_i + \|\mathbf{x}_i\|^2 \quad (1.12)$$

$$= \sum_{j=1}^k (\mathbf{x}_i \mathbf{v}_j^\top)^2 - 2 \sum_{j=1}^k (\mathbf{x}_i \mathbf{v}_j^\top) \mathbf{v}_j \mathbf{x}_i^\top + \|\mathbf{x}_i\|^2 \quad (\because \mathcal{B} \text{ is orthonormal})$$

$$= \|\mathbf{x}_i\|^2 - \sum_{j=1}^k (\mathbf{x}_i \mathbf{v}_j^\top)^2 \quad (1.13)$$

Therefore, it is equivalent to find

$$\mathcal{B}^* = \arg \max_{\mathcal{B}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k (\mathbf{x}_i \mathbf{v}_j^\top)^2. \quad (1.14)$$

Since $\mathbf{x}_i \mathbf{v}_j^\top$ is a scalar, we know that $\mathbf{x}_i \mathbf{v}_j^\top = \mathbf{v}_j \mathbf{x}_i^\top$. Then we have

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k (\mathbf{x}_i \mathbf{v}_j^\top)^2 = \frac{1}{n} \sum_{j=1}^k \sum_{i=1}^n (\mathbf{v}_j \mathbf{x}_i^\top) (\mathbf{x}_i \mathbf{v}_j^\top) = \sum_{j=1}^k \mathbf{v}_j \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i \right) \mathbf{v}_j^\top = \sum_{j=1}^k \mathbf{v}_j \mathcal{S} \mathbf{v}_j^\top. \quad (1.15)$$

Now, we are into an optimization problem of finding

$$\mathcal{B}^* = \arg \max_{\mathcal{B}} \sum_{j=1}^k \mathbf{v}_j \mathcal{S} \mathbf{v}_j^\top \quad (1.16)$$

where \mathcal{B} is orthonormal. In order to find a solution, we solve the following Lagrangian:

$$\mathcal{L}(\mathbf{v}_1, \dots, \mathbf{v}_k, \lambda_1, \dots, \lambda_k) = \sum_{j=1}^k \left[\mathbf{v}_j \mathcal{S} \mathbf{v}_j^\top - \lambda_j (\mathbf{v}_j \mathbf{v}_j^\top - 1) \right] \quad (1.17)$$

Recall that if $y = \mathbf{x} A \mathbf{x}^\top$ and $z = \mathbf{x} \mathbf{x}^\top$, then $\frac{dy}{d\mathbf{x}} = 2\mathbf{x} A$ and $\frac{dz}{d\mathbf{x}} = 2\mathbf{x}$ for $\mathbf{x} \in \mathbb{R}^d$ and $A \in \mathbb{R}^{d \times d}$.

$$\frac{\partial \mathcal{L}}{\partial \lambda_j} = \mathbf{v}_j \mathbf{v}_j^\top - 1 = 0 \quad \Rightarrow \quad \mathbf{v}_j \mathbf{v}_j^\top = \|\mathbf{v}_j\|^2 = 1 \quad (1.18)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}_j} = 2\mathbf{v}_j \mathcal{S} - 2\lambda_j \mathbf{v}_j = 0 \quad \Rightarrow \quad \mathbf{v}_j \mathcal{S} = \lambda_j \mathbf{v}_j \quad \text{and} \quad \mathbf{v}_j \mathcal{S} \mathbf{v}_j^\top = \lambda_j \quad (1.19)$$

Thus, \mathbf{v}_j is the eigenvector of \mathcal{S} corresponding to the j th largest eigenvalue, that is, it is the top j th principal component of X . \square

We can also prove the claim by using mathematical induction.

- Prove that it holds for $k = 1$.
- Assuming that it holds for $k - 1$, prove that it also holds for k . That is, if the top $k - 1$ principal components $\{\mathbf{v}_1, \dots, \mathbf{v}_{k-1}\}$ maximizes $\sum_{j=1}^{k-1} \mathbf{v}_j \mathcal{S} \mathbf{v}_j^\top$, then $\{\mathbf{v}_1, \dots, \mathbf{v}_{k-1}, \mathbf{v}_k\}$ also maximizes $\sum_{j=1}^k \mathbf{v}_j \mathcal{S} \mathbf{v}_j^\top$.
- Notice that the choice of k does not affect the resulting set of vectors \mathbf{v}_j .
- The low-dimension k only decides how many vectors of \mathbf{v}_j 's we use to make a linear projection.

We have concluded that PCA minimizes the reconstruction loss $\frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{x}}_i - \mathbf{x}_i\|^2$. Let's take a look at what the minimized loss will be. Let $\{\mathbf{v}_1, \dots, \mathbf{v}_d\}$ be the all principal components of X , that is, $\mathbf{x}_i = \sum_{j=1}^d (\mathbf{x}_i \mathbf{v}_j^\top) \mathbf{v}_j$. Then, for the projection $\tilde{\mathbf{x}}_i$ onto U , we have $\mathbf{x}_i - \tilde{\mathbf{x}}_i = \sum_{j=k+1}^d (\mathbf{x}_i \mathbf{v}_j^\top) \mathbf{v}_j \in U^\perp$ for all $i = 1, \dots, n$. Thus the reconstruction loss is

$$\frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{x}}_i - \mathbf{x}_i\|^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=k+1}^d (\mathbf{x}_i \mathbf{v}_j^\top)^2 = \sum_{j=k+1}^d \mathbf{v}_j \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i \right) \mathbf{v}_j^\top = \sum_{j=k+1}^d \lambda_j. \quad (1.20)$$

That is, the loss is the sum of the eigenvalues corresponding to the remaining principal components of X .

1.2.3 Maximum Variance Perspective

“Retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional code.” (Harold Hotelling, 1933)

Fix $k \leq d$. Let $\tilde{X} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n]^\top$ be the “best” linear projection of X in the projection perspective, that is, $\tilde{\mathbf{x}}_i = \sum_{j=1}^k (\mathbf{x}_i \mathbf{v}_j^\top) \mathbf{v}_j$ where \mathbf{v}_j is the top j th principal component of X .

Claim 2. *PCA finds the subspace that maximizes the variance of the projected data, that is, \tilde{X} can be obtained by a basis $\mathcal{B} = \{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ that maximizes the variance V_j of $\{\mathbf{x}_1 \mathbf{v}_j^\top, \dots, \mathbf{x}_n \mathbf{v}_j^\top\}$ with $V_1 \geq V_2 \geq \dots \geq V_k$.*

Proof. Assuming that $\mathbb{E}_{\mathbf{x}}[\mathbf{x} \mathbf{v}_j^\top] = 0$, we have $\text{Var}_{\mathbf{x}}[\mathbf{x} \mathbf{v}_j^\top] = \mathbb{E}_{\mathbf{x}}[(\mathbf{x} \mathbf{v}_j^\top)^2] - \mathbb{E}_{\mathbf{x}}[\mathbf{x} \mathbf{v}_j^\top]^2 = \mathbb{E}_{\mathbf{x}}[(\mathbf{x} \mathbf{v}_j^\top)^2]$.

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i \mathbf{v}_j^\top)^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{v}_j \mathbf{x}_i^\top) (\mathbf{x}_i \mathbf{v}_j^\top) = \frac{1}{n} \mathbf{v}_j \left(\sum_{i=1}^n \mathbf{x}_i^\top \mathbf{x}_i \right) \mathbf{v}_j^\top = \mathbf{v}_j \mathcal{S} \mathbf{v}_j^\top \quad (1.21)$$

The sample variance of X is equal to the objective in the projection perspective, which completes the proof. \square

The maximum variance that PCA can capture/explain is $\sum_{j=1}^k \lambda_j$, and the lost variance by data compression via PCA is $\sum_{j=k+1}^d \lambda_j$.

2 Autoencoder

Autoencoder (AE) is an encoder-decoder network to find a latent space of given dataset. It tries to find latent vectors by reducing the dimension of data and by restoring the original data from them. The natural bottleneck network forces the encoder network to compress information into latent vectors. Additionally, the encoder network should find features of the original data so that the decoder network can reconstruct the original ones from the latent vectors. AE has three important attributes: dimensionality reduction, data-specific, and lossy.

- (i) **Dimensionality Reduction:** The encoder compresses an input into a low dimensional latent vector and the decoder reconstructs the input from this latent vector as an output.
- (ii) **Data Specific:** To find well-trained encoder-decoder network which compresses data meaningfully, data should have strong correlations between input features, and/or data should be similar to what they have been trained on.
- (iii) **Lossy:** Since the dimension of latent vector is smaller than the original input, the reconstructed input is bound to have lower quality than the original input.

Optimizing the network can be considered as an unsupervised learning since it does not need any label for input. The training algorithm of AE is described as below.

Algorithm 1 Training Autoencoder

Input: Data batches $\mathcal{D} = \{\mathcal{B}_j\}_{j=1}^B$, $\mathcal{B}_j = \{x_i^{(j)}\}_i$, encoder Ψ , and decoder Φ

- 1: **for** each batch $j = 1, \dots, B$ **do**
 - 2: Initialize $L = 0$ $\triangleright L$ for accumulation of loss
 - 3: **for** each data $x_i^{(j)} \in \mathcal{B}_j$ **do**
 - 4: Compute $\tilde{x}_i^{(j)} = \Phi(\Psi(x_i^{(j)}))$ and the loss $\mathcal{L}(x_i^{(j)}, \tilde{x}_i^{(j)})$
 - 5: $L = L + \frac{1}{B}\mathcal{L}(x_i^{(j)}, \tilde{x}_i^{(j)})$
 - 6: **end for**
 - 7: Take a gradient descent step on $\nabla_{\theta} L$ $\triangleright \theta$: the set of parameters in Ψ and Φ
 - 8: **end for**
-

Dimensionality reduction methods can be used for data visualization. In PCA, it attempts to find a lower dimensional hyperplane which maximizes the variance of original data. On the other hand, since AE is defined as nonlinear neural networks, it is capable of learning a nonlinear manifold describing the data in a lower dimensionality.

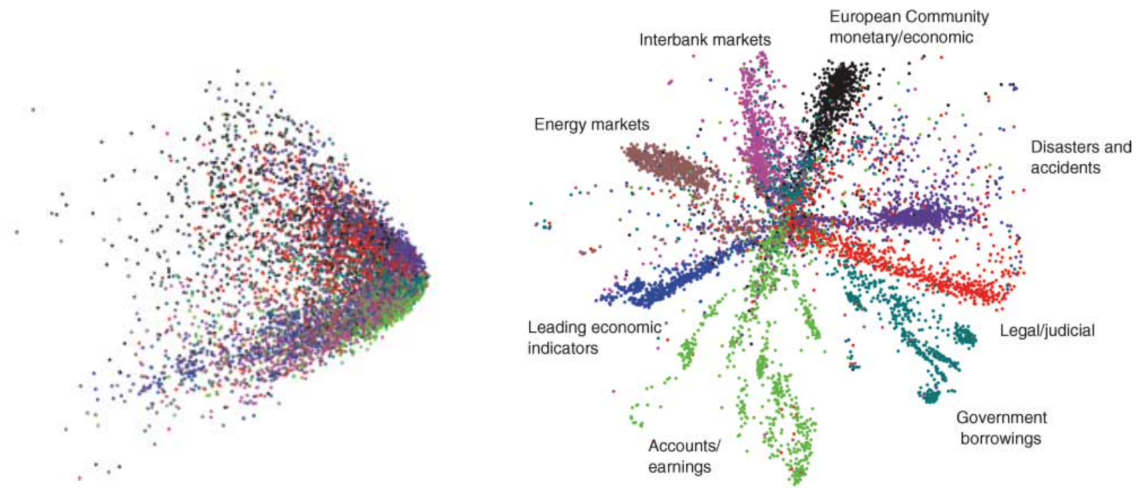


Figure 2.1: The left shows the two-dimensional codes produced by LSA (a well-known document retrieval method based on PCA), while the right shows the codes produced by an autoencoder.

Figure 2 shows that autoencoder clearly learns a nonlinear manifold, which can provide clear visualization.

2.1 Denoising Autoencoder

References

- [1] G.E. Hinton and R. R. Salakhutdinov, *Reducing the Dimensionality of Data with Neural Networks*, Science (2006).