

# Principal Component Analysis

This version: September 27, 2023

## 1 Dimension Reduction: Encoding and Decoding

Suppose that we have  $p$  variables. Denote our data as  $\mathbf{X}$ , which is an  $n \times p$  matrix, and each observation is a  $p \times 1$  vector. In dimension reduction, we aim to find a function  $f(\cdot)$

$$f : \mathbb{R}^p \rightarrow \mathbb{R}^q, \quad q < p,$$

that can “encode” an observation  $\mathbf{x}_i \in \mathbb{R}^p$  into a lower-dimensional space  $\mathbf{z}_i = f(\mathbf{x}_i) \in \mathbb{R}^q$ . Ideally, we would also like to have a function  $g(\cdot)$

$$g : \mathbb{R}^q \rightarrow \mathbb{R}^p,$$

that “decodes” the observation back to the original space  $\hat{\mathbf{x}}_i = g(\mathbf{z}_i)$ . Since the encoder maps observations into a lower-dimensional space, information loss in the process is inevitable. A good encoder/decoder pair should minimize the reconstruction error  $L$

$$L = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2^2.$$

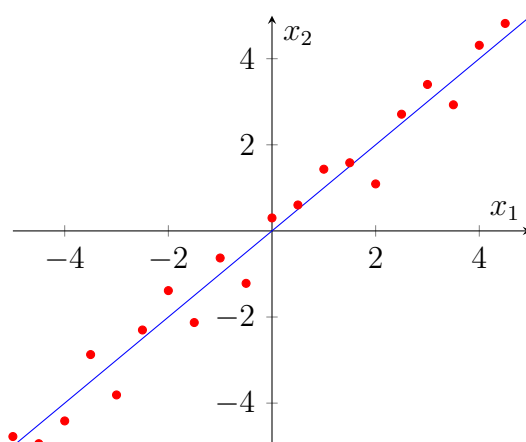
**Remark 1.** For a vector  $\mathbf{x} = (x_1, x_2, \dots, x_p) \in \mathbb{R}^p$ , the notation  $\|\mathbf{x}\|_l$  denotes its  $l$ -norm:

$$\begin{aligned} \|\mathbf{x}\|_l &= \left( \sum_{j=1}^p |x_j|^l \right)^{\frac{1}{l}} \\ &= (|x_1|^l + |x_2|^l + \dots + |x_p|^l)^{\frac{1}{l}}. \end{aligned}$$

For example, when  $l = 2$ , it is also known as the  $L_2$  norm, which corresponds to the Euclidean norm that we have been using since high school. Later in this course, we will also use the  $L_1$  norm, denoted as  $\|x\|_1$ . However, unless otherwise specified, we usually consider the  $L_2$  norm.

How do we find an encoder/ decoder? Consider the following scatter plot of our data  $p = 2$ :

Is our data two-dimensional or one-dimensional?



From the graph, we can see that most of our observations are close to the line 45-degree. So, why not project our data onto the line, which is an one dimensional space. Consider the following example.

**Example 1** (“Eyeball” PCA). Suppose  $p = 2$ . From the scatter plot above, it appears that we can reduce the dimension of our data from  $p = 2$  to  $q = 1$  with minimal information loss by projecting the data onto the 45-degree line. Specifically, let

$$\mathbf{A} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix},$$

and consider the function

$$f(\mathbf{x}) = \mathbf{A}^T \mathbf{x}$$

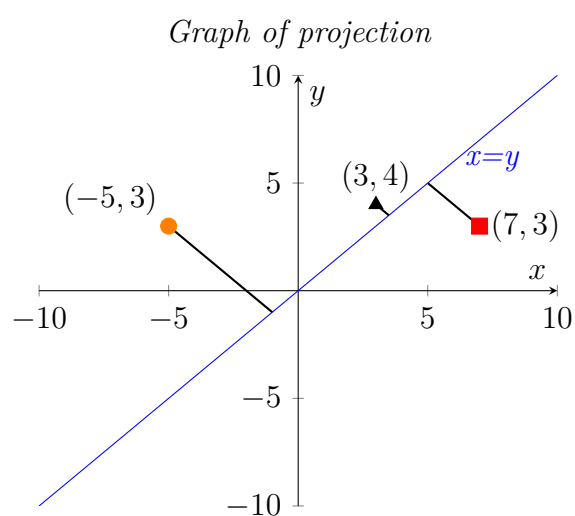
that maps from  $\mathbb{R}^2$  to  $\mathbb{R}^1$ , and

$$g(z) = \mathbf{A}z$$

that maps from  $\mathbb{R}^1$  to  $\mathbb{R}^2$ .

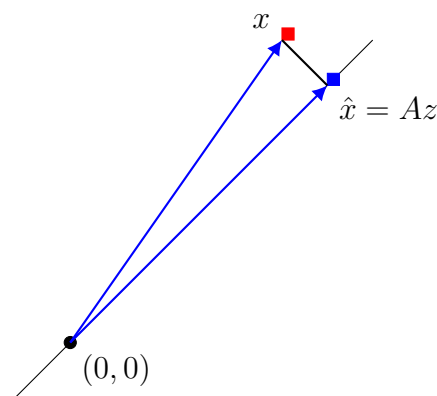
You can verify that the encoder  $f(\cdot)$  is the function that projects observations onto the 45-degree line, and the compressed observation is given by  $\hat{\mathbf{x}} = g(z) = g(f(\mathbf{x}))$ , with the reconstruction error given by  $\|\mathbf{x} - \hat{\mathbf{x}}\|$ . See the figures and table below for an illustration.

**Quiz 1.** Can you find  $z \in \mathbb{R}$ ,  $\hat{\mathbf{x}} \in \mathbb{R}^2$ , and the reconstruction error  $\|\mathbf{x} - \hat{\mathbf{x}}\|$  for  $\mathbf{x} = (3, 4)^T$ ?



| $x$       | $A^T x$               | $\hat{x} = Az$               | $\ x - \hat{x}\ $    |
|-----------|-----------------------|------------------------------|----------------------|
| $(3, 4)$  | $\frac{7}{\sqrt{2}}$  | $(\frac{7}{2}, \frac{7}{2})$ | $\frac{1}{\sqrt{2}}$ |
| $(7, 3)$  | $\frac{10}{\sqrt{2}}$ | $(5, 5)$                     | $2\sqrt{2}$          |
| $(-5, 3)$ | $\frac{-2}{\sqrt{2}}$ | $(-1, -1)$                   | $2\sqrt{5}$          |

Graph of  $\mathbf{x} = \hat{\mathbf{x}}, \hat{\mathbf{z}}$  and  $\hat{\mathbf{x}} - \hat{\mathbf{z}}$



Example 1 essentially demonstrates how *principal component analysis (PCA)* operates in practice. PCA seeks to exploit the correlation between variables to capture the predominant variation in the data. In the extreme case, when the correlation coefficient is close to 1, two variables are nearly linearly dependent (i.e.,  $y \approx ax + b$  for some  $a, b \in \mathbb{R}$ ), and they lie along a line, as illustrated in Example 1.

For  $p = 2$  or  $p = 3$ , we have the luxury of visualizing the data to detect whether there exists a lower-dimensional space (a line or a plane) that adequately represents our data. However, this visual inspection, or "eyeballing," becomes unfeasible for dimensions beyond  $p = 3$ . What is our alternative? To address cases where  $p > 3$ , we must generalize our approach by abstracting the ideas of what we applied in lower dimensions.

## 2 Beyond $p = 3$ : Deriving PCA from Scratch

As illustrated in Example 1, the core concept of PCA is to identify a lower-dimensional space onto which our observations can be projected. How can we formalize this concept for  $p > 3$ ?

**The key to generalization is abstraction.** For  $q = 1$ , we may characterize PCA as

“finding (1) a line to (2) project onto, with the aim of minimizing the (3) reconstruction error.” How might this statement be generalized for  $q > 1$ ? Essentially, three questions need to be addressed:

1. What is the entity onto which we are projecting when  $q > 1$ ?
2. How is projection defined for  $q > 1$ ?
3. How is reconstruction defined for  $q > 1$ ?

**What is the entity onto which we are projecting when  $q > 1$ ?** For  $q = 1$ , we are looking for a line; for  $q = 2$ , we seek a plane, and so forth. Since one vector is required to span a line, and two vectors to define a plane, in general, we need  $q$  vectors  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_q$  to construct a  $q$ -dimensional space. In PCA, our goal is to identify these vectors, termed as the “principal components”, which form the lower-dimensional space. Here,  $\mathbf{a}_1$  is the first principal component (PC),  $\mathbf{a}_2$  is the second PC, and so on.

| Project to | Find       | Need vectors (“principal components”)                       |
|------------|------------|---|
| $q = 1$    | a line     | $\mathbf{a}_1 \in \mathbb{R}^p$                             |
| $q = 2$    | a plane    | $\mathbf{a}_1, \mathbf{a}_2 \in \mathbb{R}^p$               |
| $q = 3$    | a 3D space | $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3 \in \mathbb{R}^p$ |

**Remark 2.** Note that if the original data is  $p$ -dimensional, we have at most  $p$  PC’s.

Given that the vectors  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_q \in \mathbb{R}^p$  fundamentally define a lower-dimensional subspace (of dimension  $q$ ) in  $\mathbb{R}^p$ , we impose the conditions:

$$\begin{aligned} \|\mathbf{a}_j\|_2 &= \sqrt{\mathbf{a}_j^T \mathbf{a}_j} = 1, & (\text{unit length}) \\ \mathbf{a}_j^T \mathbf{a}_k &= 0, & (\text{orthogonality for } j \neq k) \end{aligned}$$

The principal components  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_q$  possess unit length as they fundamentally represent directions. Ensuring their orthogonality, as we will see below, simplifies the definition of the projection onto the space they span.

**Quiz 2.** Which of the following vectors can possibly be the first two PC’s for a data with three variables ( $p = 3$ )?

$$\begin{array}{ll}
1. \mathbf{a}_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \mathbf{a}_2 = \begin{bmatrix} 3 \\ 2 \end{bmatrix} & 3. \mathbf{a}_1 = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}, \mathbf{a}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \\
2. \mathbf{a}_1 = \begin{bmatrix} 1 \\ \sqrt{2} \\ 0 \end{bmatrix}, \mathbf{a}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} & 4. \mathbf{a}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \mathbf{a}_2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}
\end{array}$$

Now, **how do we define project on  $q > 1$ ?**<sup>1</sup> It turns out that if we define

$$\mathbf{A}_{p \times q} = \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_q \end{bmatrix}$$

then for an observation  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T \in \mathbb{R}^p$ :

$$\mathbf{A}^T \mathbf{x}_i = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_q^T \end{bmatrix}_{q \times p} \cdot \mathbf{x}_{p \times 1} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{x}_i \\ \mathbf{a}_2^T \mathbf{x}_i \\ \vdots \\ \mathbf{a}_q^T \mathbf{x}_i \end{bmatrix} \in \mathbb{R}^q.$$

So, the  $j$ -th component in  $\mathbf{A}^T \mathbf{x}_i$  is simply the inner product of  $\mathbf{a}_j$  and  $\mathbf{x}_i$ . Since  $\mathbf{a}_j$  has unit length, then the projection of  $x_i$  on  $\mathbf{a}_j$  is simply

$$\mathbf{a}_j^T \mathbf{x}_i \cdot \mathbf{a}_j.$$

So, the matrix multiplication  $\mathbf{A}^T \mathbf{x}_i$  encodes  $\mathbf{x}_i$  in terms of its coordinates on  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_q$ . In matrix notation, the projection of  $\mathbf{x}_i$  onto the space spanned by  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_q$  is then

---

<sup>1</sup>Note that, in math, an understanding of the case for  $q = 2$  would automatically generalize to any finite-dimensional case. So it is fine to assume  $q = 2$  for the discussion.

given by

$$\begin{aligned}\hat{\mathbf{x}}_i &= \mathbf{a}_1^T \mathbf{x}_i \cdot \mathbf{a}_1 + \mathbf{a}_2^T \mathbf{x}_i \cdot \mathbf{a}_2 + \cdots + \mathbf{a}_q^T \mathbf{x}_i \cdot \mathbf{a}_q \\ &= \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_q \end{bmatrix} \begin{bmatrix} \mathbf{a}_1^T \mathbf{x}_i \\ \mathbf{a}_2^T \mathbf{x}_i \\ \vdots \\ \mathbf{a}_q^T \mathbf{x}_i \end{bmatrix} \\ &= \mathbf{A} \mathbf{A}^T \mathbf{x}_i.\end{aligned}$$

**Summary:** to reduce the dimension from  $p$  to  $q$ , PCA aims to find  $q$  vectors  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_q \in \mathbb{R}^p$  so that one can construct a  $q$ -dimensional space to project our data onto. Our goal is find the optimal vectors (the principal components) that the reconstruction error is smallest. That is, we want to solve the minimization problem:

$$\begin{aligned}\min_{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p \in \mathbb{R}^p} & \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{A} \mathbf{A}^T \mathbf{x}_i\|_2^2 \\ \text{s.t. } & \mathbf{a}_j^T \mathbf{a}_j = 1, \quad j = 1, 2, \dots, q \\ & \mathbf{a}_j^T \mathbf{a}_k = 0, \quad j \neq k\end{aligned}$$

where

$$\mathbf{A}_{p \times q} = \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_q \end{bmatrix}.$$

and  $\mathbf{A} \mathbf{A}^T \mathbf{x}_i \in \mathbb{R}^q$  is the projected (compressed) version of  $\mathbf{x}_i \in \mathbb{R}^p$ .

**Quiz 3.** Suppose that  $p = 3$ . What is matrix  $\mathbf{x}$  for the projection on the plane  $x = y$ ?

### 3 Solving PCA

So far, we have spent a lot of time defining PCA. But, how can we solve it, i.e., how do we find the PC's  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p$ ?

PCA attempts to summarize variation in the random vector  $\mathbf{X}$  with few principal components (PC). In below, we will define what are *principal components* and explain how

to derive them.

In fact, we can solve PCA iteratively, that is, we can first find the first PC,  $\mathbf{a}_1$ , then  $\mathbf{a}_2$ , ..., and  $\mathbf{a}_q$ . So, to find the first PC, we solve

$$\begin{aligned} \min_{\mathbf{a}_1 \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{a}_1 \mathbf{a}_1^T \mathbf{x}_i\|_2^2 \\ \text{s.t. } \mathbf{a}_1^T \mathbf{a}_1 = 1. \end{aligned}$$

Given  $\mathbf{a}_1$ , we can then solve  $\mathbf{a}_2$  by minimizing the remaining part  $\mathbf{x}_i - \mathbf{a}_1 \mathbf{a}_1^T \mathbf{x}_i$

$$\begin{aligned} \min_{\mathbf{a}_2 \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \|(\mathbf{x}_i - \mathbf{a}_1 \mathbf{a}_1^T \mathbf{x}_i) - \mathbf{a}_2 \mathbf{a}_2^T \mathbf{x}_i\|_2^2 \\ \text{s.t. } \mathbf{a}_2^T \mathbf{a}_2 = 1, \\ \mathbf{a}_2^T \mathbf{a}_1 = 0. \end{aligned}$$

While we motivated PCA by minimizing reconstruction errors, an alternative definition of PCA is by “maximizing variances”. For example, the first PC can be found by solving the following problem:<sup>2</sup>

$$\max_{\mathbf{a}_1 \in \mathbb{R}^p} \text{Var}(\mathbf{a}_1' \mathbf{X}) \quad \text{s.t. } \mathbf{a}_1' \mathbf{a}_1 = 1.$$

**Remark 3.** The restriction  $\mathbf{a}_i' \mathbf{a}_i = 1$  is necessary, otherwise  $\text{Var}(\mathbf{a}_i' \mathbf{X})$  can be made arbitrary large by multiplying the coefficient with a constant.

The rest of the principal components can be defined iteratively. The second PC,  $PC_2 = \mathbf{a}_2'$  is defined by the optimization problem:

$$\begin{aligned} \max_{\mathbf{a}_2 \in \mathbb{R}^p} \quad & \text{Var}(\mathbf{a}_2' \mathbf{X}) \\ \text{s.t. } \quad & \mathbf{a}_2^T \mathbf{a}_2 = 1, \\ & \mathbf{a}_1^T \mathbf{a}_2 = 0. \end{aligned}$$

Notice that we require the second PC has to be uncorrelated with the first PC. We can think of it as we are trying to have the second PC explain the variation that is not explained by the first PC.

---

<sup>2</sup>For a rigor argument, see our textbook (Murphy 2022).



Similarly, the  $j$ th PC,  $\mathbf{a}'_j$  is the solution to

$$\begin{aligned} \max_{\mathbf{a}_j \in \mathbb{R}^p} \quad & \text{Var}(\mathbf{a}'_j \mathbf{X}) \\ \text{s.t.} \quad & \mathbf{a}_j^T \mathbf{a}_j = 1, \\ & \mathbf{a}_j^T \mathbf{a}_{j'} = 0 \text{ for } j' = 1, 2, \dots, j-1. \end{aligned}$$

**Remark 4.** Recall that

$$\mathbf{a}'_j \mathbf{x} = \langle \mathbf{a}_j, \mathbf{x} \rangle = \begin{pmatrix} a_{j1} & a_{j2} & \dots & a_{jp} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} = a_{j1}x_1 + a_{j2}x_2 + \dots + a_{jp}x_p,$$

so  $\mathbf{a}'\mathbf{x}$  is a scalar.

### 3.1 Solve PCA when $\Sigma$ Diagonal

**Remark 5.** We can verify that

$$\text{Var}(\mathbf{a}^T \mathbf{X}) = \mathbf{a}^T \Sigma \mathbf{a}.$$

See the following example.

**Example 2.** Let  $\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$  and  $\mathbf{a} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ .

$$\begin{aligned} \text{Var}(\mathbf{a}^T \mathbf{X}) &= \text{Var} \left( \begin{pmatrix} 1 & 2 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \right) \\ &= \text{Var}(X_1 + 2X_2) \\ &= \text{Var}(X_1) + 4\text{Var}(X_2) \\ &= 1 + 4 \cdot 2 = 9 \end{aligned}$$

Alternatively, we can calculate  $\text{Var}(\mathbf{a}^T \mathbf{X})$  by

$$\mathbf{a}^T \Sigma \mathbf{a} = \begin{pmatrix} 1 & 2 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 & 2 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} = 9.$$

Indeed, we end up with the same result, as the previous remarks implied.

By remark 5, the optimization problem can be written as

$$\max_{\mathbf{a}_1 \in \mathbb{R}^p} \quad \mathbf{a}_1^T \Sigma \mathbf{a}_1 \quad \text{s.t.} \quad \mathbf{a}_1^T \mathbf{a}_1 = 1.$$

In his book “How to Solve It”, mathematician and Probabilist George Pólya famously said

“If you can’t solve a problem, then there is an easier problem you can solve:  
find it.”

Let’s follow his suggestion by assuming the covariance matrix

$$\Sigma = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix}_{p \times p},$$

where  $\lambda_1 > \lambda_2 > \dots > \lambda_p$ . That is, we assume  $\Sigma$  is now a diagonal matrix with positive, decreasing diagonal elements.

Now we are ready to solve the optimization problem for  $PC_1$ :

$$\max_{\mathbf{a}_1 \in \mathbb{R}^p} \quad \mathbf{a}_1^T \Sigma \mathbf{a}_1 \quad \text{s.t.} \quad \mathbf{a}_1^T \mathbf{a}_1 = 1.$$

Notice the objective function is now

$$\begin{aligned}\mathbf{a}_1^T \Sigma \mathbf{a}_1 &= \begin{pmatrix} a_{11} & \dots & a_{1p} \end{pmatrix} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} a_{11} \\ \vdots \\ a_{1p} \end{pmatrix} \\ &= \lambda_1 a_{11}^2 + \lambda_2 a_{12}^2 + \dots + \lambda_p a_{1p}^2.\end{aligned}$$

By redefining  $b_{1i} = a_{1i}^2$ , the maximization problem becomes

$$\max_{\mathbf{b}} \quad \lambda_1 b_{11} + \dots + \lambda_p b_{1p} \quad \text{s.t.} \quad b_{11} + \dots + b_{1p} = 1.$$

It is plain to see that  $b_{11} = 1, b_{12} = \dots = b_{1p} = 0$  attains the maximum, while subjecting to the constraint. Hence,  $a_{11} = \pm 1, a_{12} = \dots = a_{1p} = 0$  is the solution to the  $PC_1$  problem.<sup>3</sup>

The first PC is hence

$$\mathbf{a}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

which points to the direction with largest variance.

Next we solve the  $PC_2$  problem

$$\begin{aligned}\max_{\mathbf{a}_2 \in \mathbb{R}^p} \quad & \text{Var}(\mathbf{a}_2^T \mathbf{X}) \\ \text{s.t.} \quad & \mathbf{a}_2^T \mathbf{a}_2 = 1 \\ & \mathbf{a}_1^T \mathbf{a}_2 = 0\end{aligned}$$

---

<sup>3</sup>PCs are only uniquely defined up to reflection over the origin. It is not hard to see that if  $\mathbf{a}^*$  is a solution to the PCA problem, then  $-\mathbf{a}^*$  is also a solution.

The constraint  $\mathbf{a}_1^T \mathbf{a}_2 = 0$  implies  $a_{21} = 0$ . So  $\mathbf{a}_2 = \begin{pmatrix} 0 \\ a_{22} \\ \vdots \\ a_{2p} \end{pmatrix}$ , and the problem reduces to

$$\begin{aligned} \mathbf{a}_2^T \Sigma \mathbf{a}_2 &= \begin{pmatrix} 0 & a_{22} & \dots & a_{2p} \end{pmatrix} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} 0 \\ a_{22} \\ \vdots \\ a_{2p} \end{pmatrix} \\ &= \begin{pmatrix} a_{22} & \dots & a_{2p} \end{pmatrix} \begin{pmatrix} \lambda_2 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} a_{22} \\ \vdots \\ a_{2p} \end{pmatrix}. \end{aligned}$$

The problem is now the same as finding  $PC_1$ , with one less variable. Hence the solution to the second component is  $a_{22} = \pm 1$ ,  $a_{21} = a_{23} = \dots = a_{2p} = 0$ , and the second PC is  $\mathbf{a}_2 = (0, 1, 0, \dots, 0)^T$ . Similarly, we can see that  $j$ -th PC are just the unit vectors in  $\mathbb{R}^p$  in which whose  $j$ -th component is 1 and zero otherwise.

Conclusion: When  $X_j$ 's are uncorrelated, PCA is basically keeping  $X_j$ 's with largest variances.

### 3.2 Solve PCA when $\Sigma$ is not Diagonal

What if  $\Sigma$  is not diagonal? Luckily, we have the following theorem:

**Theorem 1** (Real Spectral Theorem). *If  $\Sigma$  is symmetric, then there exists a  $p \times p$  matrix  $\mathbf{P}$  such that*

$$\Sigma = \mathbf{P} \mathbf{D} \mathbf{P}^{-1},$$

where  $\mathbf{D}$  is a diagonal matrix, and

$$\mathbf{P}^T \mathbf{P} = \mathbf{I}_{p \times p},$$

i.e.,  $\mathbf{P}^{-1} = \mathbf{P}^T$ .

**Example 3.**

$$\begin{pmatrix} 34 & 12 \\ 12 & 41 \end{pmatrix} = \begin{pmatrix} \frac{3}{5} & \frac{-4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{pmatrix} \begin{pmatrix} 50 & 0 \\ 0 & 25 \end{pmatrix} \begin{pmatrix} \frac{3}{5} & \frac{4}{5} \\ \frac{-4}{5} & \frac{3}{5} \end{pmatrix}.$$

Since  $\Sigma$  is symmetric, it can now be written as  $\Sigma = \mathbf{P}\mathbf{D}\mathbf{P}^T$ , where the notation is followed by the Real Spectral theorem. Let  $\mathbf{b}_i = \mathbf{P}^T \mathbf{a}_i$ ,

$$\mathbf{a}_i^T \Sigma \mathbf{a}_i = \mathbf{a}_i^T \mathbf{P} \mathbf{D} \mathbf{P}^T \mathbf{a}_i = (\mathbf{P}^T \mathbf{a}_i)^T \mathbf{D} (\mathbf{P}^T \mathbf{a}_i) = \mathbf{b}_i^T \mathbf{D} \mathbf{b}_i.$$

$$\mathbf{a}_i^T \mathbf{a}_i = \mathbf{a}_i^T (\mathbf{P} \mathbf{P}^{-1}) \mathbf{a}_i = \mathbf{a}_i^T (\mathbf{P} \mathbf{P}^T) \mathbf{a}_i = (\mathbf{P}^T \mathbf{a}_i)^T (\mathbf{P}^T \mathbf{a}_i) = \mathbf{b}_i^T \mathbf{b}_i$$

Restate the maximization problem in terms of  $\mathbf{b}_i$  gives

$$\max_{\mathbf{b}_i \in \mathbb{R}^p} \mathbf{b}_i^T \mathbf{D} \mathbf{b}_i \quad \text{s.t.} \quad \mathbf{b}_i^T \mathbf{b}_i = 1,$$

and we know how to solve it since  $\mathbf{D}$  is diagonal.

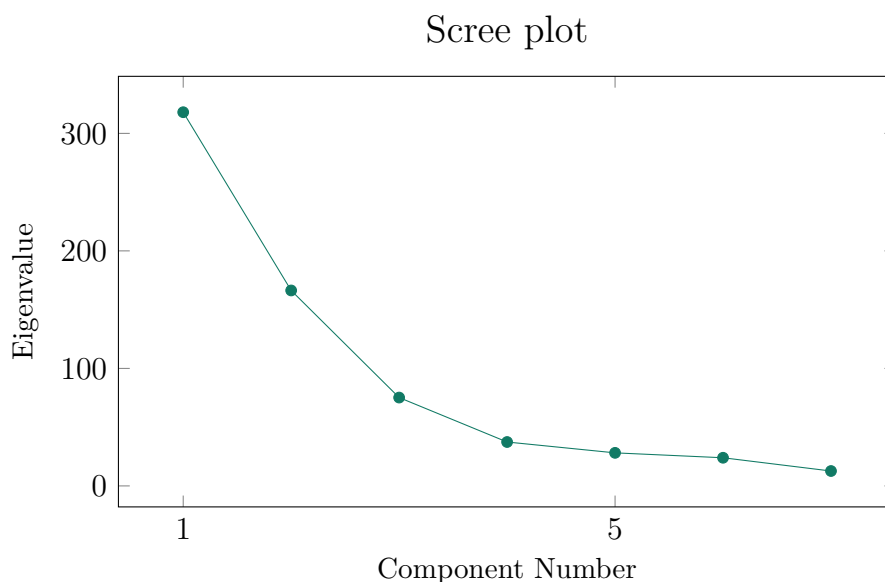
## 4 Applications of PCA

Here we list examples of applications of PCA.

1. Data compression
2. Summarize data and explorative analysis
3. Use PC in regression
4. Use PC in clustering
5. Factor analysis

## 5 Choosing $q$

One might ask how much variance is summarized in  $q$ , and how we choose such  $q$ . A scree plot might be helpful to answer the question.



The y-axis of a scree plot is the sorted eigenvalues of the covariance matrix, in a decreasing order, and the x-axis is the corresponding rank of the eigenvalue. As a rule of thumb,  $q$  is chosen at the elbow of the scree plot, so one might choose  $q = 2$  for the example above.

Suppose that the covariance matrix  $\Sigma = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_p \end{pmatrix}_{p \times p}$  is a diagonal matrix,

then we have

$$\text{Var}(PC_1) = \text{Var}(X_1) = \lambda_1$$

$$\text{Var}(PC_2) = \text{Var}(X_2) = \lambda_2$$

$$\vdots$$

$$\text{Var}(PC_p) = \text{Var}(X_p) = \lambda_p$$