

Principal Component Analysis

Jianxin Wu

LAMDA Group

National Key Lab for Novel Software Technology

Nanjing University, China

wujx2001@gmail.com

May 4, 2017

Contents

1	Introduction	2
1.1	Dimensionality and inherent dimensionality	2
1.2	Dimensionality reduction	4
1.3	PCA and the subspace methods	5
2	PCA to zero-dimensional subspace	5
2.1	The idea-formalization-optimization exercise	6
2.2	A simple optimization	6
2.3	A few notes	6
3	PCA to one-dimensional subspace	7
3.1	New formalization	7
3.2	Optimality condition and simplification	8
3.3	The eigen-decomposition connection	9
3.4	The solution	10
4	PCA to more dimensions	10
5	The complete PCA algorithm	11
6	Variance analysis	12
6.1	PCA from maximization of variance	13
6.2	A simpler derivation	14
6.3	How many dimensions do we need?	14

7	When to use or not to use PCA?	15
7.1	PCA for Gaussian data	15
7.2	PCA for non-Gaussian data	15
7.3	PCA for data with outlier	17
8	The whitening transform	18
9	Eigen-decomposition vs. SVD	18
	Exercises	18

1 Introduction

In this note, we will introduce the Principal Component Analysis (PCA) technique. The main purpose of this note is comprehensibility. We focus on the motivation and ideas behind PCA, rather than the mathematics. However, we will maintain correctness of the equations. The goal is that an average non-math-major student with undergraduate mathematical background (about linear algebra and some basic probability) will read through this note without any difficulty. And, we will start from the motivation.

1.1 Dimensionality and inherent dimensionality

Let us consider some data which are two dimensional, and we will use (x, y) to denote one such data point. Hence, the (natural) dimensionality of our data is 2. In different scenarios, we will encounter data that exhibit different properties, which we illustrate in Figure 1. We will discuss these examples one by one.

- The data set has 2 degrees of freedom. As illustrated by those points in Figure 1a, we need exactly two values (x, y) to specify any single data point in this case. Hence, we say that the data set has an *inherent dimensionality* of 2, which is the same as its natural dimensionality. In fact, the data points in this figure are generated in a way that make y independent of x , as shown in Table 1, if we consider x and y as random variables. We cannot find any relationship among the ten 2-dimensional data points in this figure.
- The world, however, usually generates data that are not independent. In other words, x and y are dependent in most cases (especially considering those data we need to deal with in machine learning or pattern recognition). In our examples, we expect that x and y are usually correlated with each other. Figure 1b exhibits a special type of correlation: linear correlation. As shown by the data generation command in Table 1, y is a linear function of x . Hence, in this figure, there is only 1 degree of freedom—once we know x , we can immediately determine the value of y ; and vice versa. We only need one value to *completely* specify a pair (x, y) ,

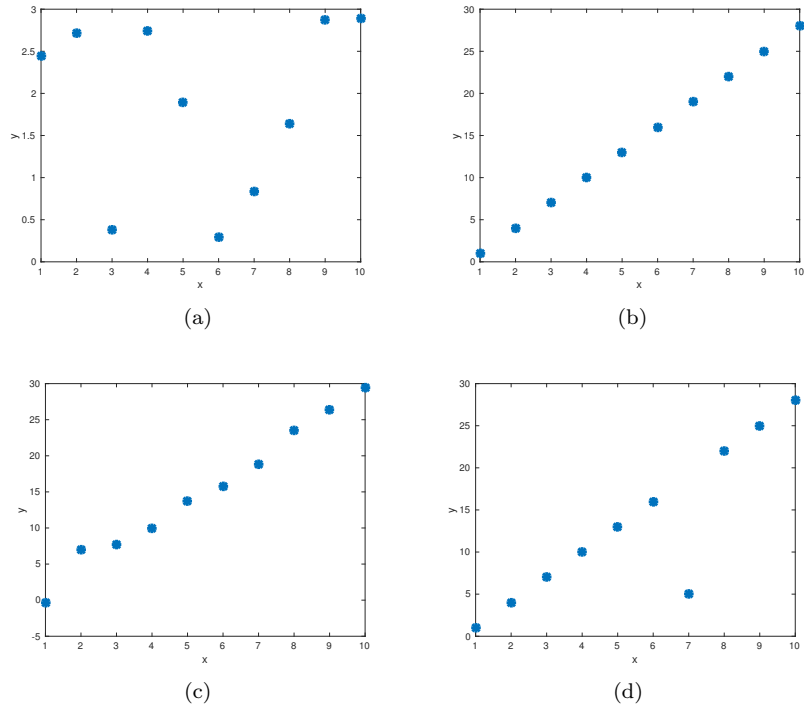


Figure 1: Illustration of various types of relationships within dimensions.

which could be x , y , or even a *linear combination of x and y* .¹ Hence, we say that the inherent dimensionality of the data in this figure is 1, which is obviously smaller than its natural dimensionality. In the figure, the 10 points are aligned in one line $y = 3x - 2$.

- The world, once again, is not benign as what is shown in Figure 1b. In many cases, there exist obvious correlations among dimensions of the original data. However, the correlation is rarely a perfect linear relationship. As shown in Figure 1c, many people will agree that x and y are roughly related by a linear relationship (i.e., the 10 points are roughly aligned in a line), but with noticeable deviations. As shown in Table 1, indeed between x and y there is the same linear relationship as that one in Figure 1b, but in addition affected by a Gaussian noise. The noise is usually not welcome and we want to get rid of them. Hence, it is still reasonable to say that the data in Figure 1c has an inherent dimensionality of 1, because the number of meaningful degree of freedom is still 1.

¹As we will soon see, a linear combination of the original dimensions is the central part of PCA.

Table 1: Matlab code to generate the data points in Figure 1.

Generating x	<code>x=1:10;</code>
Generating y for Fig. 1a	<code>y=rand(1,10)*3;</code>
Generating y for Fig. 1b	<code>y=3*x-2;</code>
Generating y for Fig. 1c	<code>y=3*x-2; y=y+randn(size(y));</code>
Generating y for Fig. 1d	<code>y=3*x-2; y(7) = 5;</code>

- The world could be even more hostile to us, as shown by the example in Figure 1d. If we remove the 7th data point, we get a perfect linear relationship. This point is significantly different than others, and we call it an *outlier*. It is difficult to say that data in Figure 1d has an inherent dimensionality of 1, given the existence of the outlier. In other words, outliers seem more difficult to handle than noise. However, if we are able to remove the outlier using some sophisticated techniques, the inherent dimensionality is still 1.

1.2 Dimensionality reduction

With these four examples, we may safely conclude that the inherent dimensionality in both Figure 1b and 1c is 1. In other words, we can use only 1 variable z to represent such data points. Finding a lower dimensional representation of the original (relatively) higher dimensional vector is called *dimensionality reduction*. In spite of the fact that fewer dimensions are used, we expect the dimensionality reduction process to keep the useful information in the original data.

The potential benefits of dimensionality reduction can be many-fold, with a few listed below.

- Lower resource requirements. A direct consequence of lower dimensionality is that less memory is needed to store the data, either in the main memory or hard disk. An equally obvious benefit is that fewer CPU cycles are requested.
- Removal of noise. As shown in Figure 1c, we can recover the linear relationship and reduce the dimensionality to 1. However, a benign side effect is that the noise may well be removed from the data in this process, which happens in PCA in many problems. Less noisy data usually lead to higher accuracy.
- Explanation and understanding. If the outcome of dimensionality reduction happen to coincide well with the underlying hidden factors that generate the original data, the new lower-dimensional representation will be helpful in explaining how the data is generated and in understanding its properties. When we get in touch with a new dataset, we can reduce it to two or three dimensional (by PCA or other methods). A visualization of

the dimensionality reduced data will give us some hints on the properties of that dataset.

1.3 PCA and the subspace methods

As far as we know, there is no precise definition for the inherent dimensionality of a data set or a data generation process / distribution in general. In the above examples, we use the degrees of freedom of the data generation process to describe their inherent dimensionality. Data generation, however, could be non-linear, which will be more complex than the linear relationship shown in Figure 1. However, since PCA only considers linear relationships, we will ignore non-linear relationships or non-linear dimensionality reduction methods in this note.

Now consider a vector $\mathbf{x} \in \mathbb{R}^D$. A linear relationship among its components can be expressed as a simple equation

$$\mathbf{x}^T \mathbf{w} + b = 0. \tag{1}$$

From basic linear algebra, we know that any \mathbf{x} that satisfies Equation 1 resides in a subspace of \mathbb{R}^D , whose cardinality is $D - 1$. If there are more linear constraints on \mathbf{x} , \mathbf{x} will reside in a subspace with still lower dimensionality.

Hence, the problem of linear dimensionality reduction can be seen as finding the linear constraints or finding the lower dimensional subspace of \mathbb{R}^D , which are also called the subspace methods. Subspace methods differ from each other in how they find the constraints or subspaces. They have different evaluation metrics (e.g., which subspace is considered the best?) or assumption (e.g., do we know the category label of \mathbf{x} ?).

PCA is possibly the simplest subspace method.

2 PCA to zero-dimensional subspace

Let us start from an extreme case. What if the lower dimensional subspace is only a single point? In other words, what if its dimensionality is 0?

Suppose we are given a set of instantiations of \mathbf{x} , $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, which form the training set for us to learn PCA parameters. Note that there is a slight abuse of notation here. The symbol \mathbf{x} may refer to a single data point. However, it may be used to refer to the underlying distribution (or random variable) that generates the training set.

If there is no noise and a zero-dimensional subspace exists to represent this set, the only possibility is that $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_N$. We can use \mathbf{x}_1 to represent any example in X without requiring any additional information. Storing \mathbf{x}_1 requires D dimensions. However, the average dimensions needed for every \mathbf{x}_i is only $\frac{D}{N}$. Because $\lim_{N \rightarrow \infty} \frac{D}{N} = 0$, it is reasonable to say that this is a zero-dimensional representation.

However, if noise exists, there will be $1 \leq i < j \leq N$ such that $\mathbf{x}_i \neq \mathbf{x}_j$. How shall we find the *best* zero-dimensional representation in the presence of noise?

We still need to find a vector \mathbf{m} that represents every element in X . The key issue is: how do we decide the optimality?

2.1 The idea-formalization-optimization exercise

The *idea* can be inspired by the noise-free case. If we assume that the noise scale is small, we want to find an \mathbf{m} , which is *close to all the elements in X* . This choice has two nice properties. First, it coincides well with our intuition. Second, when the noise scale is 0, “close to” can be changed to “equal to”, and degenerates to the noise free case.

The next step is to *formalize* this idea. It is natural to translate “close to” as “small distance”, and translate “distance to all elements in X is small” to “the sum of distances to all elements in X is small”. Hence, we can write our idea precisely in the mathematical language, as

$$\mathbf{m}^* = \arg \min_{\mathbf{m}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{m}\|^2, \quad (2)$$

where the right hand side is an optimization problem, and \mathbf{m}^* is the parameter that solves the optimization problem (which is the best \mathbf{m} we seek.)

The final step is: how can we get \mathbf{m}^* , or, how to optimize Equation 2?

2.2 A simple optimization

With some background in vector calculus, it is easy to solve Equation 2. We can denote

$$J = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{m}\|^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{m})^T (\mathbf{x}_i - \mathbf{m}), \quad (3)$$

and get

$$\frac{\partial J}{\partial \mathbf{m}} = \frac{2}{N} \sum_{i=1}^N (\mathbf{m} - \mathbf{x}_i), \quad (4)$$

where the partial derivative rules can be found in *the Matrix Cookbook*. Setting this term to 0 gives us the following optimality condition:

$$\mathbf{m} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \bar{\mathbf{x}}. \quad (5)$$

That is, the best zero-dimensional representation is the average of the data, which we denote as $\bar{\mathbf{x}}$.

2.3 A few notes

In spite of its simplicity, we have a few comments on the zero-dimensional reduction.

- When we encounter a new problem, the *idea-formalization-optimization* process can be very helpful. We first inspect the problem (maybe with some initial trials or visualization to understand the data property), which gives us some ideas to its solution. We then define proper notations and convert our ideas into a precise mathematical form, which often appears in an optimization problem format. The final step is to solve it, either by ourselves or use the abundant tools that are available.
- It is also worth-noting that Equation 2 is in fact slightly different from our translation. First, it is the squared distance that are summed, rather than the distance. Second, the term $\frac{1}{N}$ converts the sum into an average. The term $\frac{1}{N}$ will not change the optimization's solution. It is introduced following the tradition in the literature and sometimes it helps simplify the optimization or reduce numerical difficulties. The change from distance to squared distance, as shown in the optimization, makes the optimization much easier. In addition, a small squared distance implies a small distance. Hence, our idea is still valid. It is always good to *tune the mathematical translation* so long as the tuning still matches our idea and it makes the optimization easier to do.
- Although our idea starts from the assumption that the elements in X are the same data point subject to different noise, this assumption is not used at all in the formalization or optimization steps. Thus, for *any* data set X , we can *generalize* and say that its average is the best zero-dimensional representation (under the sum of squared distance evaluation metric).

3 PCA to one-dimensional subspace

Now we are ready to move one step further to the one-dimensional subspace, where we can use one additional value to represent each \mathbf{x}_i in addition $\bar{\mathbf{x}}$.

3.1 New formalization

Any element \mathbf{x} in a one-dimensional subspace can be represented as $\mathbf{x} = \mathbf{x}_0 + a\mathbf{w}$ for some a , \mathbf{x}_0 and \mathbf{w} , and vice versa, in which \mathbf{x}_0 and \mathbf{w} are determined by the subspace and a is determined by the element (recall your linear algebra!)

Now that we already have the zero-dimensional representation, we should set $\mathbf{x}_0 = \bar{\mathbf{x}}$. Hence, for any \mathbf{x}_i , the new one dimensional representation is a_i . Using this new representation, we can find an approximation of \mathbf{x}_i as $\mathbf{x}_i \approx \bar{\mathbf{x}} + a_i\mathbf{w}$, and the difference (or residue) $\mathbf{x}_i - (\bar{\mathbf{x}} + a_i\mathbf{w})$ is considered to be caused by noise, which we want to minimize.

The parameters we need to find are: a_i ($1 \leq i \leq N$) and \mathbf{w} . We denote $\mathbf{a} = (a_1, a_2, \dots, a_N)^T$, and define an objective J to minimize the average squared

distance:

$$J(\mathbf{w}, \mathbf{a}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - (\bar{\mathbf{x}} + a_i \mathbf{w})\|^2 \quad (6)$$

$$= \frac{1}{N} \sum_{i=1}^N \|a_i \mathbf{w} - (\mathbf{x}_i - \bar{\mathbf{x}})\|^2 \quad (7)$$

$$= \sum_{i=1}^N \frac{a_i^2 \|\mathbf{w}\|^2 + \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 - 2a_i \mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}})}{N}. \quad (8)$$

3.2 Optimality condition and simplification

Now we calculate the partial derivatives, and set them to 0 as

$$\frac{\partial J}{\partial a_i} = \frac{2}{N} (a_i \|\mathbf{w}\|^2 - \mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}})) = 0 \quad \forall i, \quad (9)$$

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{2}{N} \sum_{i=1}^N (a_i^2 \mathbf{w} - a_i (\mathbf{x}_i - \bar{\mathbf{x}})) = 0. \quad (10)$$

Equation 9 gives us the solution for a_i as

$$a_i = \frac{\mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}})}{\|\mathbf{w}\|^2} = \frac{(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}}{\|\mathbf{w}\|^2}. \quad (11)$$

Before proceeding to process Equation 10, a further examination of $a_i \mathbf{w}$ shows that

$$a_i \mathbf{w} = \frac{\mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}}) \mathbf{w}}{\|\mathbf{w}\|^2} = \frac{(c\mathbf{w})^T (\mathbf{x}_i - \bar{\mathbf{x}}) (c\mathbf{w})}{\|c\mathbf{w}\|^2}, \quad (12)$$

for any non-zero scalar value c . In other words, we have the freedom to specify that

$$\|\mathbf{w}\| = 1, \quad (13)$$

and this additional constraint will not change the optimization problem's solution!

We choose $\|\mathbf{w}\| = 1$ because this choice greatly simplifies our optimization problem. Now we have

$$a_i = \mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}}) = (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}. \quad (14)$$

Plugging it back to the optimization objective, we get a much simplified version

$$J(\mathbf{w}, \mathbf{a}) = \frac{1}{N} \sum_{i=1}^N [\|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 - a_i^2], \quad (15)$$

by noting that $a_i \mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}}) = a_i^2$ and $a_i^2 \|\mathbf{w}\|^2 = a_i^2$.

Hence, we know the optimal parameters are obtained via *maximizing*

$$\frac{1}{N} \sum_{i=1}^N a_i^2. \quad (16)$$

because $\|\mathbf{x}_i - \bar{\mathbf{x}}\|^2$ does not depend on either \mathbf{w} or \mathbf{a} .

One note we want to add is that various transformations can greatly simplify our optimization problem, and it is worthwhile to pay attention to such simplification opportunities. In fact, in this derivation, we could specify the constraint $\|\mathbf{w}\| = 1$ *before* finding the optimality conditions. It is easy to observe that

$$J(\mathbf{w}, \mathbf{a}) = J(c\mathbf{w}, \frac{1}{c}\mathbf{a}) \quad (17)$$

for any $c \neq 0$, and there is no constraint on \mathbf{w} or \mathbf{a} in the original formulation. Hence, if $(\mathbf{w}^*, \mathbf{a}^*)$ is an optimal solution that minimizes J , so will be $(c\mathbf{w}, \frac{1}{c}\mathbf{a})$ for any $c \neq 0$. That is, for an optimal solution $(\mathbf{w}^*, \mathbf{a}^*)$, $(\frac{1}{\|\mathbf{w}^*\|}\mathbf{w}^*, \|\mathbf{w}^*\|\mathbf{a}^*)$ will also be an optimal solution.

Obviously the norm of $\frac{1}{\|\mathbf{w}^*\|}\mathbf{w}^*$ is 1. Hence, we can specify $\|\mathbf{w}\| = 1$ without changing the optimization objective, but will greatly simplify the optimization from the very beginning. It is always beneficial to find such simplifications and transformations before we attempt to solve the optimization task.

3.3 The eigen-decomposition connection

Now we turn our attention to Equation 10, which tells us that

$$\frac{1}{N} \left(\sum_{i=1}^N a_i^2 \right) \mathbf{w} = \frac{1}{N} \sum_{i=1}^N a_i (\mathbf{x}_i - \bar{\mathbf{x}}). \quad (18)$$

And, plugging a_i into Equation 18, we can simplify its right hand side as

$$\frac{1}{N} \sum_{i=1}^N a_i (\mathbf{x}_i - \bar{\mathbf{x}}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) a_i \quad (19)$$

$$= \frac{\sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}}{N} \quad (20)$$

$$= \text{Cov}(\mathbf{x})\mathbf{w}, \quad (21)$$

where $\text{Cov}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$ is the covariance matrix of \mathbf{x} computed from the training set X .

Hence, Equation 18 now gives us

$$\text{Cov}(\mathbf{x})\mathbf{w} = \frac{\sum_{i=1}^N a_i^2}{N} \mathbf{w}, \quad (22)$$

which immediately reminds us of the eigen-decomposition equation—this equation tells us that the optimal \mathbf{w} must be an eigenvector of $\text{Cov}(\mathbf{x})$, and $\frac{\sum_{i=1}^N a_i^2}{N}$

is the corresponding eigenvalue! The constraint in Equation 13 also fits in this eigen-decomposition interpretation, because eigenvectors are also constrained to have unit ℓ_2 norm.

3.4 The solution

$\text{Cov}(\mathbf{x})$ has many eigenvectors and their corresponding eigenvalues. However, Equation 16 reminds us that we want to maximize $\frac{\sum_{i=1}^N a_i^2}{N}$, while Equation 22 tells us that $\frac{\sum_{i=1}^N a_i^2}{N}$ is the eigenvalue corresponding to \mathbf{w} . Hence, it is trivial to choose among all the eigenvectors—choose the one that corresponds to the largest eigenvalue!

Now we have everything to compute the one-dimensional reduction. From X we can compute the covariance matrix $\text{Cov}(\mathbf{x})$, and we set $\mathbf{w}^* = \boldsymbol{\xi}_1$, where $\boldsymbol{\xi}_1$ is the eigenvector of $\text{Cov}(\mathbf{x})$ which corresponds to the largest eigenvalue. The optimal new one-dimensional representation for \mathbf{x}_i is then $a_i^* = \boldsymbol{\xi}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})$.

Given the one-dimensional representation, the original input \mathbf{x} is approximated as

$$\mathbf{x} \approx \bar{\mathbf{x}} + (\boldsymbol{\xi}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})) \boldsymbol{\xi}_1. \quad (23)$$

Because $(\boldsymbol{\xi}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})) \boldsymbol{\xi}_1$ equals the projection of $\mathbf{x}_i - \bar{\mathbf{x}}$ onto $\boldsymbol{\xi}_1$,² we also call $\boldsymbol{\xi}_1$ the first *projection direction*, and call $\boldsymbol{\xi}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})$ the projected value of \mathbf{x}_i .

4 PCA to more dimensions

Now we generalize PCA to two or more dimensions, thanks to the spectral decomposition of the covariance matrix.

It is obvious that $\text{Cov}(\mathbf{x})$ is a real symmetric matrix, and furthermore it is a positive-semidefinite matrix. According to matrix analysis theory, $\text{Cov}(\mathbf{x})$ has D eigenvectors $\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_D$ whose elements are all real numbers. And, the eigenvalues corresponding to them are $\lambda_1, \lambda_2, \dots, \lambda_D$, which are all real numbers, and satisfying $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$. The spectral decomposition states that

$$\text{Cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T. \quad (24)$$

The eigenvectors of real symmetric matrices satisfy that for any $i \neq j$, $\boldsymbol{\xi}_i^T \boldsymbol{\xi}_j = 0$ and $\|\boldsymbol{\xi}_i\| = 1$ for $1 \leq i \leq D, 1 \leq j \leq D$. Hence, if we construct a $D \times D$ matrix E , whose i -th column is formed by $\boldsymbol{\xi}_i$, we have

$$EE^T = E^T E = I. \quad (25)$$

²The projection of \mathbf{x} onto \mathbf{y} is $\frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{y}\|^2} \mathbf{y}$, and $\|\boldsymbol{\xi}_1\| = 1$.

Then, we can show that

$$\mathbf{x} = \bar{\mathbf{x}} + (\mathbf{x} - \bar{\mathbf{x}}) \quad (26)$$

$$= \bar{\mathbf{x}} + EE^T(\mathbf{x} - \bar{\mathbf{x}}) \quad (27)$$

$$= \bar{\mathbf{x}} + (\boldsymbol{\xi}_1^T(\mathbf{x} - \bar{\mathbf{x}}))\boldsymbol{\xi}_1 + (\boldsymbol{\xi}_2^T(\mathbf{x} - \bar{\mathbf{x}}))\boldsymbol{\xi}_2 + \cdots + (\boldsymbol{\xi}_D^T(\mathbf{x} - \bar{\mathbf{x}}))\boldsymbol{\xi}_D, \quad (28)$$

for any $\mathbf{x} \in \mathbb{R}^D$, even if \mathbf{x} does not follow the same relationships inside the training set X (in other words, being an outlier).

Comparing Equation 23 with Equation 28, a guess of PCA with more dimensions naturally comes to us: $\boldsymbol{\xi}_i$ should be the i -th projection direction, and the coefficient is $\boldsymbol{\xi}_i^T(\mathbf{x} - \bar{\mathbf{x}})$.

This conjecture is correct, and is easy to prove following the procedure in Section 3. We will omit the details here, but leave that to the readers.

5 The complete PCA algorithm

The complete principal component analysis algorithm is described in Algorithm 1.

Algorithm 1 The PCA algorithm

- 1: **Input:** a D -dimensional training set $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ and the new (lower) dimensionality d (with $d \leq D$)
- 2: Compute the mean $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$
- 3: Compute the covariance matrix $\text{Cov}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}_i)(\mathbf{x}_i - \bar{\mathbf{x}}_i)^T$
- 4: Find the spectral decomposition of $\text{Cov}(\mathbf{x})$, obtaining the eigenvectors $\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_D$ and their corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_D$. Note that the eigenvalues are sorted, such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$
- 5: For any $\mathbf{x} \in \mathbb{R}^D$, its new lower dimensional representation is:

$$\mathbf{y} = \left(\boldsymbol{\xi}_1^T(\mathbf{x} - \bar{\mathbf{x}}), \boldsymbol{\xi}_2^T(\mathbf{x} - \bar{\mathbf{x}}), \dots, \boldsymbol{\xi}_d^T(\mathbf{x} - \bar{\mathbf{x}}) \right)^T \in \mathbb{R}^d, \quad (29)$$

and the original \mathbf{x} can be approximated as

$$\mathbf{x} \approx \bar{\mathbf{x}} + (\boldsymbol{\xi}_1^T(\mathbf{x} - \bar{\mathbf{x}}))\boldsymbol{\xi}_1 + (\boldsymbol{\xi}_2^T(\mathbf{x} - \bar{\mathbf{x}}))\boldsymbol{\xi}_2 + \cdots + (\boldsymbol{\xi}_d^T(\mathbf{x} - \bar{\mathbf{x}}))\boldsymbol{\xi}_d \quad (30)$$

Let E_d be the $D \times d$ matrix which consists of the first d columns in E , the new lower dimensional representation can be succinctly written as

$$\mathbf{y} = E_d^T(\mathbf{x} - \bar{\mathbf{x}}), \quad (31)$$

and the approximation is

$$\mathbf{x} \approx \bar{\mathbf{x}} + E_d E_d^T(\mathbf{x} - \bar{\mathbf{x}}). \quad (32)$$

Dimensions of the new representation are called the *principal components*, hence the name Principal Component Analysis (PCA). Sometimes a typo will happen, which spells PCA as Principle Component Analysis, but that is not correct.

6 Variance analysis

Note that we have used \mathbf{y} to denote the new lower dimensional representation in Algorithm 1. Let y_i be the i -th dimension of \mathbf{y} , we can compute its expectation:

$$\mathbb{E}[y_i] = \mathbb{E}[\boldsymbol{\xi}_i^T (\mathbf{x} - \bar{\mathbf{x}})] = \boldsymbol{\xi}_i^T \mathbb{E}[\mathbf{x} - \bar{\mathbf{x}}] = \boldsymbol{\xi}_i^T \mathbf{0} = 0, \quad (33)$$

where $\mathbf{0}$ is a vector whose elements are all zero.

We can further calculate its variance:

$$\text{Var}(y_i) = \mathbb{E}[y_i^2] - (\mathbb{E}[y_i])^2 \quad (34)$$

$$= \mathbb{E}[y_i^2] \quad (35)$$

$$= \mathbb{E}[\boldsymbol{\xi}_i^T (\mathbf{x} - \bar{\mathbf{x}}) \boldsymbol{\xi}_i^T (\mathbf{x} - \bar{\mathbf{x}})] \quad (36)$$

$$= \mathbb{E}[\boldsymbol{\xi}_i^T (\mathbf{x} - \bar{\mathbf{x}}) (\mathbf{x} - \bar{\mathbf{x}})^T \boldsymbol{\xi}_i] \quad (37)$$

$$= \boldsymbol{\xi}_i^T \text{Cov}(\mathbf{x}) \boldsymbol{\xi}_i \quad (38)$$

$$= \boldsymbol{\xi}_i^T \lambda_i \boldsymbol{\xi}_i \quad (39)$$

$$= \lambda_i. \quad (40)$$

Hence, y_i is zero-mean and its variance is λ_i .

Equations 22 tell us that for the first new dimension,

$$\frac{\sum_{i=1}^N a_i^2}{N} = \lambda_1, \quad (41)$$

And, Equation 15 tells us that

$$J(\boldsymbol{\xi}_1) = \frac{1}{N} \sum_{i=1}^N [\|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 - a_i^2] = \frac{\sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}{N} - \lambda_1. \quad (42)$$

That is, $\frac{\sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}{N}$ is the average squared distance for the zero dimensional representation, λ is the cost that is reduced by introducing the $\boldsymbol{\xi}_1$ projection direction as a new dimension. It is also easy to prove that

$$J(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_k) = \frac{\sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}{N} - \lambda_1 - \lambda_2 - \dots - \lambda_k \quad (43)$$

for $1 \leq k \leq D$ and

$$J(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_D) = 0. \quad (44)$$

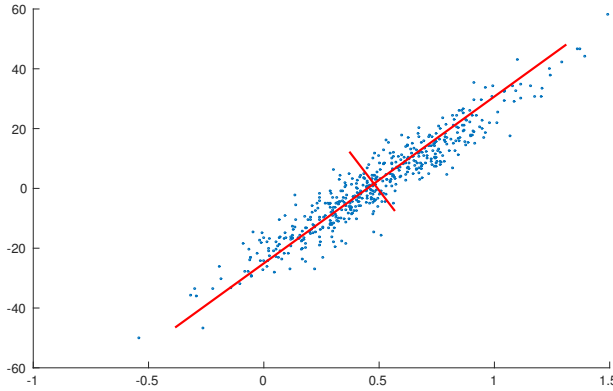


Figure 2: Variance of projected values.

Hence, every new dimension is responsible for reducing the reconstruction distance between any point \mathbf{x} and its approximation $\bar{\mathbf{x}} + \sum_{i=1}^k (\boldsymbol{\xi}_i^T (\mathbf{x} - \bar{\mathbf{x}})) \boldsymbol{\xi}_i$. And, we know that the i -th new dimension reduces the expected squared distance by λ_i , and that distance will reduce to 0 if all eigenvectors are used. From these observations, we get

- If all eigenvectors are used, PCA is simply a *rotation*, because $\mathbf{y} = E^T(\mathbf{x} - \bar{\mathbf{x}})$ for any \mathbf{x} , and E is an orthogonal matrix. We also know that in this case the norm is unchanged: $\|\mathbf{y}\| = \|\mathbf{x} - \bar{\mathbf{x}}\|$.
- We know that the larger an eigenvalue, the larger its associated eigenvector (projection direction) will reduce the approximation error.
- We also know that the eigenvalue λ_i is the expectation of the square of the i -th new dimension (cf. Equation 40). In general, we will expect the average scale of y_i (the i -th dimension in \mathbf{y}) to be larger than that of y_j if $i < j$.

6.1 PCA from maximization of variance

Based on the above observations, we can also interpret PCA as maximizing the variance of the projected values to a certain direction, as illustrated in Figure 2.

The two red lines show the two eigenvectors of the covariance matrix. It is not surprising that the long red line has the highest variance of its projected values. In other words, PCA can also be derived by maximizing the variance of the projected values to a projection direction, and the optimal solution of this formulation will be the same as what we obtain from the minimization of average squared distances. Hence, we can use the following terms interchangeably: variance (of projected values), eigenvalue, and reduction in approximation error.

6.2 A simpler derivation

Let us only work out the first projection direction that maximizes the variance of the projected values. Given any projection direction \mathbf{w} , the projected point of a data point \mathbf{x} onto \mathbf{w} will be (cf. footnote 2 in page 10)

$$\frac{\mathbf{x}^T \mathbf{w}}{\|\mathbf{w}\|^2} \mathbf{w} = \mathbf{x}^T \mathbf{w} \mathbf{w} . \quad (45)$$

In the above formula, we assume $\|\mathbf{w}\| = 1$, because as aforementioned, the norm of \mathbf{w} does not change our optimization. Hence, the projected value for \mathbf{x} is $\mathbf{x}^T \mathbf{w}$. The mean of all projected value is $\mathbb{E}[\mathbf{x}^T \mathbf{w}] = \bar{\mathbf{x}}^T \mathbf{w}$.

Next, we compute its variance, as

$$\text{Var}(\mathbf{x}^T \mathbf{w}) = \mathbb{E} [(\mathbf{x}^T \mathbf{w} - \bar{\mathbf{x}}^T \mathbf{w})^2] \quad (46)$$

$$= \mathbf{w}^T \mathbb{E} [(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T] \mathbf{w} \quad (47)$$

$$= \mathbf{w}^T \text{Cov}(\mathbf{x}) \mathbf{w} . \quad (48)$$

This is the second time we see this optimization problem:

$$\max_{\mathbf{w}} \mathbf{w}^T \text{Cov}(\mathbf{x}) \mathbf{w} \quad \text{subject to} \quad \|\mathbf{w}\| = 1 . \quad (49)$$

This is an sufficient evidence that the variance maximization perspective leads to exactly the same PCA solution as the one we just obtained by minimizing the approximation error. The variance perspective, however, is much easier to derive the PCA solution. The approximation perspective, although requiring a little more efforts, reveals more properties of PCA, such as the relationship between the eigenvalue, approximation error, and variance. Hence, we started from the approximation perspective in this note.

6.3 How many dimensions do we need?

The equivalence of these terms also gives us hints on how to choose d , the number of dimensions in the new representation.

If one eigenvalue is 0, then its corresponding eigenvector (projection direction) is not useful in keep information of the original data distribution at all. All data points that have the same characteristic as the training set will have a constant projected value for this eigenvector. Hence, this eigenvalue and its associated eigenvector can be safely removed.

When an eigenvalue is quite small, there is good reason to conjecture that: this projection direction does not contain useful information either. Rather, it could be there due to white (or other types of) noise. Hence, removing such directions are usually encouraged, and in many cases will increase the accuracy of our new representation in subsequent classification system (or other tasks).

We want to keep a reasonably large portion of the variance, such that the remaining eigenvalues / variance are likely caused by noise. Hence, the rule of thumb is often to choose to cut off if the accumulated eigenvalues has exceeded

90% of the sum of all eigenvalues. In other words, we choose d to be the first integer that satisfies

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_d}{\lambda_1 + \lambda_2 + \dots + \lambda_D} > 0.9. \quad (50)$$

Although 0.9 seems to be the widely used cutoff threshold, other values (such as 0.85 or 0.95) can also be used.

7 When to use or not to use PCA?

This question will end this note. And, before we touch this tough question, we start by a simpler one: How will PCA affect \mathbf{x} if it follows a normal distribution?

7.1 PCA for Gaussian data

Let us suppose $\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$. Usually we do not know the exact value of the mean $\boldsymbol{\mu}$ or the covariance matrix Σ . However, the maximum likelihood estimation of these terms can be obtained as $\bar{\mathbf{x}}$ and $\text{Cov}(\mathbf{x})$, respectively.

Let Λ denote the diagonal matrix consists of the eigenvalues of $\text{Cov}(\mathbf{x})$, i.e., $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_D)$. Following the properties of normal distributions,³ it is easy to verify that the new representation \mathbf{y} is also normally distributed, whose parameters are estimated as

$$\mathbf{y} \sim N(\mathbf{0}, \Lambda), \quad (51)$$

if all projection direction are used. That is, PCA performs a rotation, such that the axes of the normal distribution are parallel to the coordinate system's axes. A direct consequence of this result is that different components of \mathbf{y} are *independent* to each other.

If only the first d eigenvectors are used, then we can define a $d \times d$ matrix Λ_d such that $\Lambda_d = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$, and

$$\mathbf{y}_d \sim N(\mathbf{0}, \Lambda_d). \quad (52)$$

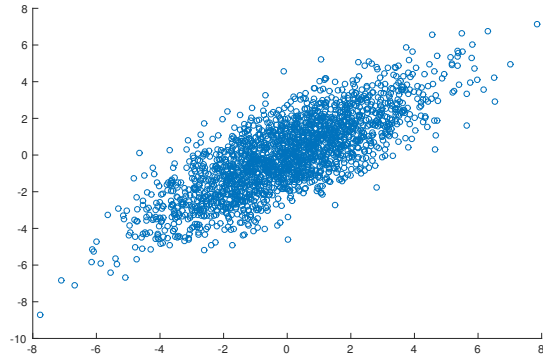
The projected dimensions are independent to each other too. Figure 3 shows an example, in which 3a contains 2000 normally distributed two-dimensional data, generated in Matlab by `x=randn(2000,2)*[2 1;1 2]`. After the PCA operation, Figure 3b shows that these data points are rotated back to follow an ellipsoidal normal distribution (i.e., whose covariance matrix is diagonal).

7.2 PCA for non-Gaussian data

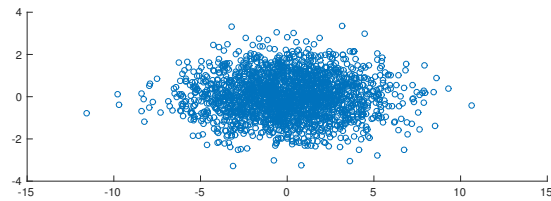
We will, however, expect many non-Gaussian data. Figure 4 shows the eigenvalue distribution of a non-Gaussian data.

We can observe an *exponential* trend in how the eigenvalue decreases. Because of this exponential decay trend, the first few eigenvalues may quickly

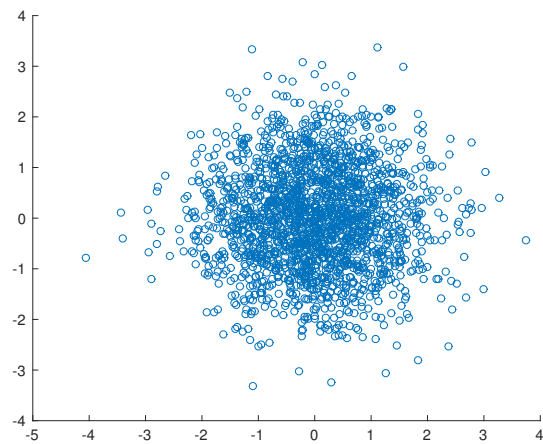
³Please refer to my note on this topic.



(a) Input data



(b) PCA



(c) whitening

Figure 3: PCA and the whitening transform applied to Gaussian data. **3a** is the 2-dimensional input data. After PCA, the data are rotated such that two major axes of the data are parallel to the coordinate system's axes in **3b** (i.e., the normal distribution becomes an ellipsoidal one). After the whitening transform, the data has the same scale in the two major axes in **3c** (i.e., the normal distribution becomes a spherical one). Note that the x - and y -axis in every subfigure has different scales.

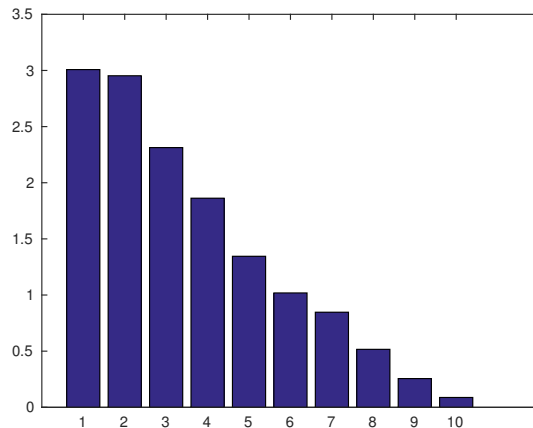


Figure 4: Eigenvalues shown in decreasing order.

Table 2: PCA outcome for data in Figure 1b, 1c, and 1d.

Data source	λ_1	λ_2	projection direction
Figure 1b	825	0	(3.00,1)
Figure 1c	770	0.75	(2.90,1)
Figure 1d	859	16.43	(3.43,1)

accumulate a high proportion of the total variance (or, sum of eigenvalues). Hence, when the eigenvalues show an exponential decay trend, the last few dimensions may be noise, and it is reasonable to apply PCA to such data. In Figure 4, the last 3 eigenvalues only account for 6% of the total variance. We can safely set $d = 7$ in this example (where $D = 10$).

If the data is not Gaussian, \mathbf{y} will have mean $\mathbf{0}$ and covariance matrix Λ (or Λ_d) after PCA. Thus, we know the dimensions in \mathbf{y} are not correlated. However, they are not necessarily independent.

7.3 PCA for data with outlier

Outlier can cause serious trouble for PCA. We compute PCA for the data used in Figure 1, and in Table 2, we listed the PCA results for the data in the last three subfigures of Figure 1.

When there is no noise, PCA successfully estimates the major projection direction as (3, 1), which fits Figure 1b. The white noise are applied to every data point in Figure 1c. However, it only slightly increases λ_2 from 0 to 0.75, and slightly changes the projection direction to (2.9, 1). But, one single outlier in the data of Figure 1d significantly changes the projection direction to (3.43, 1), and leads to a large λ_2 (16.43). Overall, PCA is not effective in the existence of

outliers.

8 The whitening transform

Sometimes we have reasons to require that the dimensions in \mathbf{y} have roughly the same scale. However, PCA ensures that $\mathbb{E}[y_1^2] \geq \mathbb{E}[y_2^2] \geq \dots \geq \mathbb{E}[y_D^2]$. The whitening transform is a simple variation of PCA, and can achieve this goal.

The whitening transform gets the new lower dimensional representation as

$$\mathbf{y} = (E_d \Lambda^{-1/2})^T (\mathbf{x} - \bar{\mathbf{x}}). \quad (53)$$

This equation differs from Equation 31 by an additional term $\Lambda^{-1/2}$, which guarantees that $\mathbb{E}[y_1^2] = \mathbb{E}[y_2^2] = \dots = \mathbb{E}[y_D^2]$ after the whitening transform. However, we have to remove any projection direction whose corresponding eigenvalue is 0 in the whitening transform.

As shown in Figure 3c, after the whitening transform the dataset follow a spherical normal distribution.

9 Eigen-decomposition vs. SVD

When either the number of data points N or the dimensionality D is large, especially when D is large, eigen-decomposition could be computationally very expensive. The Singular Value Decomposition (SVD) is usually used to replace eigen-decomposition in this scenario. The covariance matrix $\text{Cov}(\mathbf{x})$ do not need to be explicitly computed. Simply using the data matrix X , SVD can compute the (left and right) singular vectors and singular values. Depending on whether $N > D$ or $D > N$, the eigenvectors of $\text{Cov}(\mathbf{x})$ will match either the left or right singular vectors. And, the singular values, when squared, will match the eigenvalues.

Exercises

1. Let X be an $m \times n$ matrix with singular value decomposition $X = U\Sigma V^T$, where $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_{\min(m,n)})^T$ contains the singular values of X .
 - (a) What are the eigenvalues and eigenvectors of XX^T ?
 - (b) What are the eigenvalues and eigenvectors of $X^T X$?
 - (c) What is the relationship between the eigenvalues of XX^T and $X^T X$?
 - (d) What is the relationship between the singular values of X and the eigenvalues of XX^T ($X^T X$)?
 - (e) If $m = 2$ and $n = 100000$, how will you compute the eigenvalues of $X^T X$?

2. We study the effect of the average vector in PCA in this problem. Use the following Matlab code to generate a dataset with 5000 examples and compute their eigenvectors. If we forget to normalize the data by minus the average vector from every example, is there a relationship between the first eigenvector (i.e., the one associated with the largest eigenvalue) and the average vector? Observe these vectors while changing the value of `scale` in the set $\{1, 0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001\}$. What is the correct eigenvector (in which the average vector are removed from every example) if `scale` changes?

```
% set the random number seed to 0 for reproducibility
rng(0);
avg = [1 2 3 4 5 6 7 8 9 10];
scale = 0.001;
% generate 5000 examples, each 10 dim
data = randn(5000,10)+repmat(avg*scale,5000,1);

m = mean(data); % average
m1 = m / norm(m); % normalized avearge

% do PCA, but without centering
[~, S, V] = svd(data);
S = diag(S);
e1 = V(:,1); % first eigenvector
% do correct PCA with centering
newdata = data - repmat(m,5000,1);
[U, S, V] = svd(newdata);
S = diag(S);
new_e1 = V(:,1); % first eigenvector

% correlation between first eigenvector (new & old) and mean
corr1 = corrcoef(avg,e1)
corr2 = corrcoef(e1,new_e1)
```

3. Complete the following experiments using Matlab or GNU Octave. Write your own code to implement both PCA and whitening—you can use functions such as `eig` or `svd`, but do not use functions that directly finish the task for you (e.g., the `princomp` function).
- Generate 2000 examples using `x=randn(2000,2)*[2 1;1 2]`, in which the examples are two-dimensional. Use the `scatter` function to plot these 2000 examples.
 - Perform the PCA operation on these examples and keep all the 2 dimensions. Use the `scatter` function to plot the examples after PCA.
 - Perform the whitening operation on these examples and keep all the 2 dimensions. Use the `scatter` function to plot the examples after PCA.
 - Why PCA is a rotation of data if all dimensions are kept in the PCA operation? Why is this operation useful?
4. (Givens rotation) Givens rotations (which are named after James Wallace Givens, Jr., a mathematician of USA) are useful in setting certain elements

in a vector to 0. A Givens rotation involves two indexes i, j and an angle θ , which generates a matrix of the form:

$$G(i, j, \theta) = \begin{bmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \cdots & c & \cdots & s & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & -s & \cdots & c & \cdots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \quad (54)$$

in which $c = \cos \theta$ and $s = \sin \theta$. The diagonal entries in the matrix $G(i, j, \theta)$ are all 1, except $G(i, j, \theta)_{i,i} = G(i, j, \theta)_{j,j} = c$. Most of the off-diagonal entries are 0, except $G(i, j, \theta)_{i,j} = s$ and $G(i, j, \theta)_{j,i} = -s$. Let $G(i, j, \theta)$ be of size $m \times m$ and $\mathbf{x} \in \mathbb{R}^m$.

(a) What is the effect of left multiplying $G(i, j, \theta)^T$ to \mathbf{x} ? That is, what is the difference between \mathbf{x} and $\mathbf{y} = G(i, j, \theta)^T \mathbf{x}$?

(b) If we want to enforce $y_j = 0$, what is your choice of θ (or equivalently, c and s values)?

(c) Evaluation of trigonometric functions or their inverse functions are expensive. Without using trigonometric functions, how shall you determine the matrix $G(i, j, \theta)$? That is, how to determine the values c and s ?

(d) If $G(i, j, \theta)^T$ left multiplies a matrix A of size $m \times n$, how does it alter A ? How can we use a Givens rotation to change one entry of a matrix A to 0? What are the computational complexity of applying this transformation?

(e) (QR decomposition) Let A be a real matrix of size $m \times n$. Then, there exists an orthogonal real matrix Q (of size $m \times m$) and an upper triangular real matrix R (of size $m \times n$), such that $A = QR$.⁴ This decomposition for any real matrix A is called the *QR decomposition*. How shall you make use of Givens rotations to produce a QR decomposition?

5. (Jacobi's method) One method to calculate the principal components is Jacobi's approach, invented by and named after Carl Gustav Jacob Jacobi, a German mathematician. Let X be a real symmetric matrix of size $n \times n$.

(a) If G is an orthogonal $n \times n$ real matrix and $\mathbf{x} \in \mathbb{R}^n$, prove that $\|\mathbf{x}\| = \|G\mathbf{x}\|$ and $\|\mathbf{x}\| = \|G^T \mathbf{x}\|$ (i.e., a rotation will not change the length of a vector).

(b) The Frobenius norm of an $m \times n$ matrix A is defined as $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{i,j}^2} = \sqrt{\text{tr}(AA^T)}$. If G is an orthogonal $n \times n$ real matrix,

⁴ R is an upper triangular matrix if any entry below the main diagonal is 0, i.e., $R_{ij} = 0$ so long as $i > j$.

and X is any real matrix of size $n \times n$, prove that $\|X\|_F = \|G^T X G\|_F$.

(c) For a real symmetric matrix X , Jacobi defined a loss function for it, as the Frobenius norm of the off-diagonal elements in X , i.e., $\mathbf{off}(X) = \sqrt{\sum_{i=1}^n \sum_{j=1, j \neq i}^n x_{i,j}^2}$. The basic building block in Jacobi's method for eigen-decomposition is to find an orthogonal matrix J such that $\mathbf{off}(J^T X J) < \mathbf{off}(X)$. Explain why this basic step is useful in finding the eigenvectors and eigenvalues of X ?

(d) How do you decide an orthogonal matrix J such that the (i, j) -th and (j, i) -th entries in $J^T X J$ are both zeros ($i \neq j$)? (Hint: Let $J(i, j, \theta)$ be a Givens rotation matrix)

(e) The classical Jacobi method iterates between the following steps: 1) Find one off-diagonal entry in X that has the largest absolute value, i.e., $|x_{pq}| = \max_{i \neq j} |x_{ij}|$; 2) $X \leftarrow J(p, q, \theta)^T X J(p, q, \theta)$. This process converges if $\mathbf{off}(X)$ is smaller than a predefined threshold ϵ . Prove that one iteration will not increase $\mathbf{off}(X)$.

(f) Given the specific choice in the classical Jacobi method, prove that it always converges.