

Module 4

Problem Setting.

In this session we will introduce the setting of PCA and the high level idea.

"Assume we have a dataset x in \mathbb{R}^D consisting of N vectors"

$$X = \{x_1, \dots, x_N\} \quad x_i \in \mathbb{R}^D$$

Our objective is to find the low Dim representation of the data that is as similar to x as possible

Before we start, let's briefly review 3 important Concepts

First, every vector in \mathbb{R}^D can be represented as a linear combination of the basis vectors

Let's write it down: " x_n can be written

as the sum of i equal 1 to D , of $B(b_i)$ in times b_i , and we will assume that b_i

are an orthonormal basis of \mathbb{R}^D "

2

$$\textcircled{1} \quad X_n = \sum_{i=1}^D \beta_{in} b_i$$

If we assume that we use the dot product as
an inner product, and b_i to be an orthonormal
basis,

second: we can also write " β_{in} as X_n^T transpose
times b_i ", which means we can interpret
 β_{in} to be the orthogonal projection of X_n onto
one dimensional subspace (spanned) spanned by
the i th basis vector

$$\textcircled{2} \quad \beta_{in} = X_n^T b_i$$

Third: If we have an orthonormal basis b_1 to b_m ,
of \mathbb{R}^D and we define B to be the
matrix that consists of these orthonormal
basis vectors, $B(b_1, \dots, b_m)$, then the
projection of X onto the subspace

"We can write as \tilde{X} (tilde) is B times $B^{\text{transpose}}$ times X ", that means \tilde{X} (tilde) is the orthogonal projection of X onto the subspace, spanned by the M basis vectors

$$(3) \quad B = (b_1, \dots, b_m)$$

$$\tilde{X} = BB^T X$$

And $BB^T X$ is the coordinate of \tilde{X} w.r.t. the basis vectors collected in matrix B

This $(BB^T X)$ is also called the coordinates or CODE

Now let's have a look at PCA.

4
The key idea in PCA is to find a lower dim
representation \tilde{X}_n of X_n that can be expressed
using fewer basis vectors
(see p. 1)

Prerequisites

① we assume the data is centered, that
means the dataset has means 0,

② and we also assume that b_1 to b_D are
an orthonormal basis (ONB) of \mathbb{R}^D

Basically we can write any \tilde{X}_n in the
following way

" \tilde{X}_n can be written as a sum, $i=1$ to M ,
of β_i in times b_i plus
sum of $i=M+1$ to D of β_i in times b_i
still living in \mathbb{R}^D "

$$\textcircled{B} \quad \tilde{X}_n = \sum_{i=1}^M \beta_i n b_i + \sum_{i=M+1}^D \beta_i n b_i \in \mathbb{R}^D \quad \textcircled{A}$$

So we took our general way of writing any vector in \mathbb{R}^D which comes from property ①, and we split the sum in property ① into 2 sums.

One is living in an M dimensional subspace.

Other one is living in $D-M$ dimensional subspace which is an orthogonal component to the particular subspace.

In PCA we ignore the second term.

So we get rid of

~~$$\sum_{i=M+1}^D \beta_i n b_i$$~~

This part

And then we call the subspace that is spanned by the basis vectors b_1 to b_M , the principal subspace

Although \tilde{X}_n is still a D dimensional vector, it ~~thus~~ lives in a M Dim subspace of \mathbb{R}^D , and only M coordinates ~~to~~ β_{m1} to β_{mM} , ~~are~~ unnecessary in order to represent it.

So these ones β (A) and are the coordinates of the β (page 5)

The β (betas)'s are ~~are~~ also called the code or coordinates of \tilde{X}_n w.r.t to basis vectors b_1, \dots, b_M

And the setting now is as follows:

Assuming we have data X_1, \dots, X_N , we want to find parameters β_{in} and

orthonormal basis vectors b_i , such that the average squared reconstruction error is minimized

And we can write the average squared reconstruction error, as follows:

Can write " J " ^{← average squared reconstruction error} is going to be 1 over N times
sum n equal 1 to N || then we write X_n minus
 \hat{X}_n (fitted) squared"

$$J = \frac{1}{N} \sum_{n=1}^N \|X_n - \hat{X}_n\|^2$$

Let's have a look at an example:

we have data living in 2 Dim, and now we
want to find a good 1 Dim subspace such
that the squared or average squared reconstruction
error of the data, original data points in the
corresponding projection is minimized (see pic)

Here I am plotting the original data set with
the corresponding properties onto 1 Dim

Subspace, and cycling through a couple of options
of Subspace and you can see that some
of the properties are significantly more
informative than others.

And in PCA we go to find the best one:
(see PCs)

Our approach is to compute partial derivatives of J wrt the parameters: ⑨

The parameters are:

- B_{in}

- b_i 's

We set the partial derivatives of J wrt the parameters to zero, and solve for the optimal parameters.

But one ~~easy~~ observation we can already make. and that observation is that the parameters only enter this loss function (J) through \hat{X}_n

This means in order to get our partial derivatives, we need to apply the Chain rule, so

So we can write $\mathcal{L}J$ over \mathcal{L} either β_{in} or b_i (10)

Can be written as $\mathcal{L}J$ over $\mathcal{L}\tilde{X}_n$ times $\mathcal{L}\tilde{X}_n$ over \mathcal{L} either β_{in} or b_i

$$\frac{\mathcal{L}J}{\mathcal{L}\{\beta_{in}, b_i\}} = \frac{\mathcal{L}J}{\mathcal{L}\tilde{X}_n} \frac{\mathcal{L}\tilde{X}_n}{\mathcal{L}\{\beta_{in}, b_i\}}$$

And the first part we can already compute, and we get

" $\frac{\mathcal{L}J}{\mathcal{L}\tilde{X}_n}$ is minus 2 over N , times X_n minus \hat{X}_n transpose"

$$\frac{\mathcal{L}J}{\mathcal{L}\tilde{X}_n} = -\frac{2}{N} (X_n - \hat{X}_n)^T$$

and other derivatives we compute in next session.