

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-Dens 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 $B(B^{-1})_i$
times b_i , and we will assume that b_i

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

① $X_n = \sum_{i=1}^D \beta_{in} b_i$

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

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

② $\beta_{in} = X_n^T b_i$

Third; If we have an orthonormal basis b_1, \dots, b_m , of \mathbb{R}^D \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} (tilda) is B times B transpose times X ", that means \tilde{X} (tilda) is the orthogonal projection of X onto the subspace, spanned by the M basis vectors³

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

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

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.

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 next slide)

- (P2)
1. Centered data: $E[X] = 0$
 2. ONB b_1, \dots, b_D

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

Finally we can write any \tilde{x}_n in the following way

" x_n can be written as a sum, $i=1$ to M , of $\beta_{i1}(\beta)$ times b_i plus sum of $i=M+1$ to D of β_i times b_i still living in \mathbb{R}^D "

$$\textcircled{B} \quad \textcircled{A}$$

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

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

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

The part

And then we call the subspace that's spanned by the
basis vectors b_i for m , the principal subspace

Although \tilde{X}_n is still a D dimensional vector, it (we) lives in a M dim subspace of \mathbb{R}^D , and only m coordinates $\beta_{1n}, \dots, \beta_{mn}$ are necessary in order to represent it.

So the ones $\beta_{1n}, \dots, \beta_{mn}$ are the coordinates of the \tilde{X}_n (Betas).

The β 's 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 orthogonal basis vectors b_i , such that the average Squared reconstruction error is minimized

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

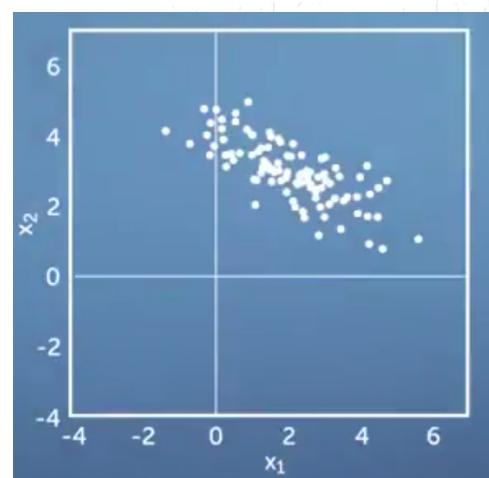
Can write " J " is going to be 1 over N times
average Squared reconstruction error $\quad \text{⑦}$

Sum n equal 1 to N || Then we write X_n minus
 \hat{X}_n (fitted) Squared"

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

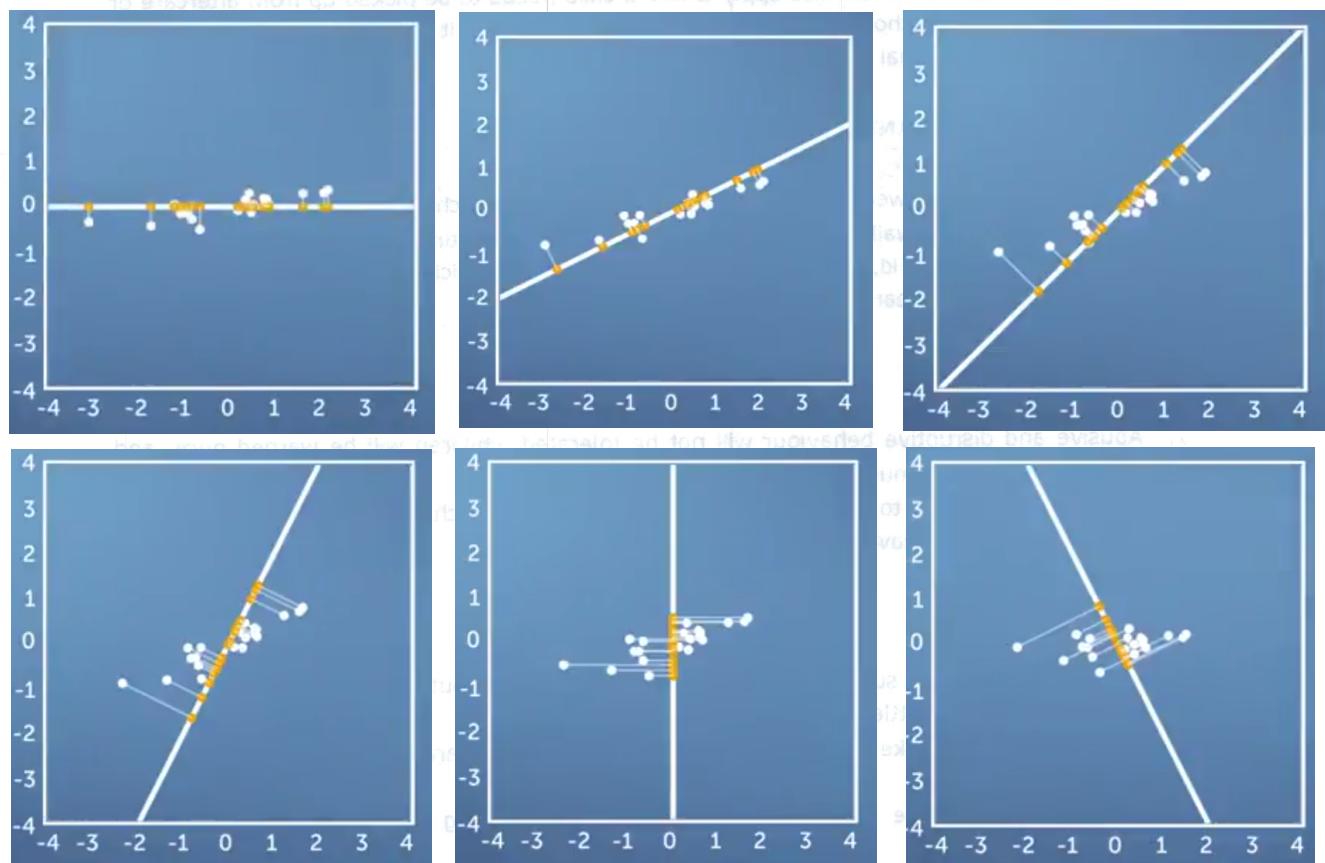
Let's have a look at an example:

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 projections onto 1D and

subspaces, and cycling through a couple of options of subspaces and you can see that some of the projections are significantly more informative than others.
And in PCA we going to find the best one!
(see PCs)



Our approach is to compute partial derivatives
of J w.r.t the parameters:

The parameters are :

- B_{in}

- b_i 's

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

But one observation we can already make.
One observation is that the parameters
and that observation x_n enter this loss function (J) through
only x_n .

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

So we can write " $\frac{\partial J}{\partial \theta}$ over ∂ either β_{in} or b_i " ⑩

Can be written as $\frac{\partial J}{\partial \theta}$ over \hat{x}_n times \hat{x}_n over ∂ either β_{in} or b_i

$$\frac{\partial J}{\partial \{\beta_{in}, b_i\}} = \frac{\frac{\partial J}{\partial \hat{x}_n}}{\hat{x}_n} \frac{\partial \hat{x}_n}{\partial \{\beta_{in}, b_i\}}$$

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

$$\frac{\partial J}{\partial \hat{x}_n} = \text{"minus 2 over } N \text{, times } x_n \text{ minus } \\ x_n \text{ transpose"}$$

$$\frac{\partial J}{\partial \hat{x}_n} = -\frac{2}{N} (x_n - \hat{x}_n)^T$$

and other derivatives we compute in next session.