

## Module 4

### Optimal basis vectors

In the last session we found that minimizing the Average Squared Residual Error is equivalent to minimizing the projection of variance of data when projected onto the subspace that we'll ignore in PCA.

In this session we will exploit this insight and determine an orthonormal basis of the  $n$  dimensional principal subspace using the results from earlier (see pre).

$$\begin{array}{l} \textcircled{A} \quad \tilde{x}_n = \sum_{j=1}^M \beta_{jn} b_j \\ \textcircled{B} \quad J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 \\ \textcircled{C} \quad \frac{\partial J}{\partial \tilde{x}_n} = -\frac{2}{N} (x_n - \tilde{x}_n)^T \\ \textcircled{D} \quad \beta_{jn} = x_n^T b_j, j=1..M \\ \textcircled{E} \quad x_n - \tilde{x}_n = \sum_{j=M+1}^P (b_j^T x_n) b_j \\ \textcircled{F} \quad J = \sum_{j=M+1}^P b_j^T S b_j \end{array}$$

we can write (at) our loss function as:

$$J = \sum_{j=M+1}^D b_j^T S b_j$$

Where  $S$  is the data Covariance Matrix

Minimizing this objective, requires us to find the orthonormal basis matrices the subspace that we will ignore

and when we have that basis, we take the orthogonal complement as the basis of the principal subspace

Remember that the orthogonal Complement of a subspace  $U$ , consisting of all vectors in the original vector space, that are orthogonal to every vector in  $U$ .

Let us start with an example to determine the  $b$  vectors.

And that in 2 dimensions, what we wish to find a 1 dim subspace such that the variance of data when projected onto that subspace is maximized.

So we looking at 2 basis vectors  $b_1$  and  $b_2$  in  $\mathbb{R}^2$ :  $b_1, b_2$

$b_1$  will be spanning the principal subspace,  
 $b_2$  is orthogonal complement, i.e. the subspace we will ignore.

We also have the constraint that  $b_1$  and  $b_2$  are orthonormal, i.e. " $b_i^T b_j = \delta_{ij}$ "

$$\underbrace{b_i^T b_j}_{= \delta_{ij}}$$

meaning that, this dot product is 1 if  $i=j$ , and zero otherwise

"So,  $\frac{\partial L}{\partial t}$  over  $\partial t$  is 1 minus  $b_2^T b_2$ , and this is 0, if and only if  $b_2^T b_2$  is 1"

$$\frac{\partial L}{\partial t} = 1 - b_2^T b_2 = 0 \Leftrightarrow b_2^T b_2 = 1$$

(So we recover our constraint)

So now let's look at the partial derivative

of  $L$  wrt  $b_2$

We get:

$$\frac{\partial L}{\partial b_2} = 2b_2^T S - 2\lambda b_2 = 0$$

"and zero, if and only if,  $S$  times  $b_2$  is  $\lambda$  times  $b_2$ "

$$\Leftrightarrow Sb_2 = \lambda b_2$$

Here we end up with an Eigenvalue problem

$b_2$  is an eigenvector of (rooted) data covariance matrix

and the Lagrange multipliers play the role of the corresponding Eigenvalue

If we now go back to the loss function, we can use this expression:

We can write " $f$  which was  $b_2$  transpose times  $b_2$ "

$$J = b_2^T S b_2$$

"But now we know that  $S$  times  $b_2$  can be written

as  $\lambda$  times  $b_2$ , so we get  $b_2$  transpose times  $b_2$  times  $\lambda$ "

$$= \underbrace{b_2^T b_2}_{\text{orthonormal}} \lambda$$

And because we have an orthonormal basis<sup>7</sup>  
"we end up with  $\lambda^k$ "

$$\rightarrow = \lambda$$

as a loss function.

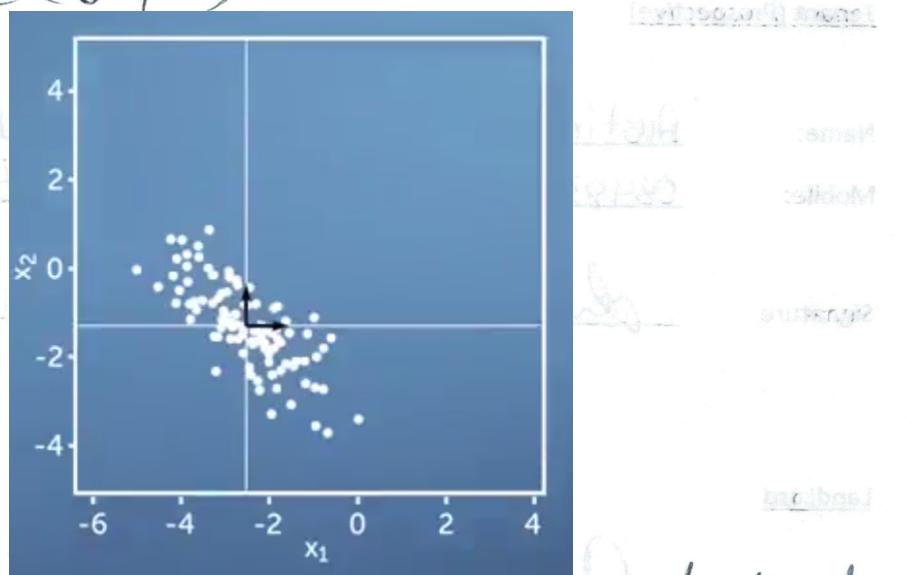
Therefore the average squared reconstruction error  
is minimized if  $\lambda$  is the smallest eigenvalue  
of the data covariance matrix

That means we need to choose  $b_2$  as  
the corresponding eigenvector and  
and that one will span the subspace that  
we will ignore.

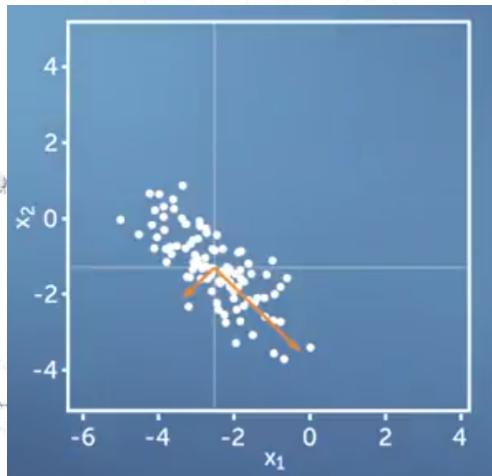
If  $b_1$  which spans the principal subspace,  
is then the eigen vector that belongs to  
the largest eigen value of the data  
covariance matrix

Kapitel und that the eigenvectors of the Covariance matrix are already orthogonal to each other because of the symmetry of the Covariance matrix.

If we look at a two Dim example, of the car data (see pic):



Then the best projection that we can get, that retains most of the information is the one that projects onto the subspace that's spanned by the Eigenvector of the data Covariance matrix which belongs to the largest eigenvalue, and which is indicated by this long arrow over here (see pic)



Let's go to the general case.

If we want to find the n dim. principal subspace of a D dimensional dataset; and we chose for the basis vectors "b<sub>j</sub>" where j equals M+1 to D"

$$b_j, \quad j = M+1, \dots, D$$

if we optimize these ones, we end up with the same kind of eigenvalue problem that we had earlier with the simple example

we end up with "S times b<sub>j</sub> equal A<sub>j</sub> times b<sub>j</sub> for j equal M+1 to D"

$$\sum b_j = \lambda_j b_j ; j = M+1, \dots, D$$

And the loss function is given by the sum of the corresponding eigenvalues.

We can write:

$$J = \sum_{j=M+1}^D \lambda_j$$

Now we can see that

also in the general case the average reconstruction error is minimized, if we choose the basis vectors that spans the ignored subspace to be eigen vectors of the data co-variance matrix belonging to the smallest eigenvalues.

This equivalently means, that the principal subspace is spanned by the eigen vectors belonging to the  $m$  largest eigenvalues of the data co-variance matrix.

The nicely aligns with properties of the covariance matrix.

- ① The eigenvectors of the covariance matrix are orthogonal to each other, because of symmetry.
- ② The eigenvector belonging to the largest eigenvalue, points ~~test~~ in the direction of data with largest variance.
- ③ and the variance in that direction is given by the corresponding eigenvalue.

Similarly, the eigenvector belonging to the second largest eigenvalue, points in the direction and largest variance of data and so on...

In this session we identified the orthonormal basis of the principal subspace as the eigenvectors of the data covariance matrix that are associated with the largest eigenvalues.

In next session we going to put all pieces together and run through the PCA algorithm in detail.