Module 4

Optimal basis vector

In the last session we found that minimizing the average squared reconstruction error is equivalent to minimizing the projection of variance of data when projected onto the subspace that will 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)

We can write (both) 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 that spans the subspace that we will ignore and when we have that basis, we take the orthogonal compliment as the basis of the principal subspace

Remember that the orthogonal compliment of a subspace $U$, consist 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

Can det Stats in 2 Dimensions, where we wish to find a 1 Dim subspace such that the variance of data when projected onto that subspace is maximized.

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

$b_1$ 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 "$b_i^T b_j$ is Delta $(\delta)_{ij}$"

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

meaning that, this dot product is 1 if $i=j$, and zero (0) otherwise

"So, $\frac{dL}{d\lambda}$ over $d\lambda$ is 1 minus $b_2$ transpose times $b_2$, and this is 0, if and only if $b_2^T$ transpose $b_2$ is 1"

$$\frac{\partial L}{\partial \lambda} = 1 - b_2^T b_2 = 0 \Longleftrightarrow b_2^T b_2 = 1$$

(So we recover our constraint)

So now let's have a look at the partial derivative of $L$ w.r.t $b_2$

we get:

Ⓐ from First term    Ⓑ from Second term    needs to be zero

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

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

$$\Longleftrightarrow S b_2 = \lambda b_2$$

Here we end up with an Eigenvalue problem

$b_2$ is an Eigenvector of (vector) data covariance matrix

and the Lagrange multiplier plays the role of the corresponding Eigenvector value

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

we can write "J which was $b_2$ transpose times $S$ times $b_2$"

$$J = b_2^T S b_2 \longrightarrow$$

"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 $\lambda$"

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

and because we have an orthonormal basis,

" we end up with $k \lambda$ "

$$ \rightarrow \quad = \lambda $$

as our 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 eigen~~value~~ vector, ~~and~~ and that one will span the subspace that we will ignore.

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

Keep in mind that the eigenvectors of the Covariance matrix are already orthogonal to each other because of the symmetry of the Co-variance matrix.

If we look at a two dim example, of the or 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 data Covariance matrix which belong to the largest eigenvalue, and that's indicated by these long arrow over here (see pic)

Lets go to the general case.

If we want to find the n Dim. principal subspace of a D Dimensional dataset, and we solve for the basis vectors "bj where j equals M+1 to D"

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

We do optimize these "ones", we end up with the same kind of eigenvalue problems that we had earlier with the simple example we end up with "S times bj equals $\lambda_j$ times bj for j equal M+1 to D"

$$S b_j = \lambda_j b_j \; ; \; j = M+1, \ldots, \Delta$$

And the loss function is given by the sum of the corresponding eigen values.

we can write:

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

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 that belong to the smallest eigenvalues.

This equivalently means, that the principal subspace is spanned by the eigen vectors belonging to the the M largest eigen values of the data co-variance matrix

This, 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 eigen vector belonging to the largest
eigenvalue, points 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 eigen value, points in the
direction and largest variance of data
and so on ....

In this session we identified the orthonomal [12] 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.