

Module 4

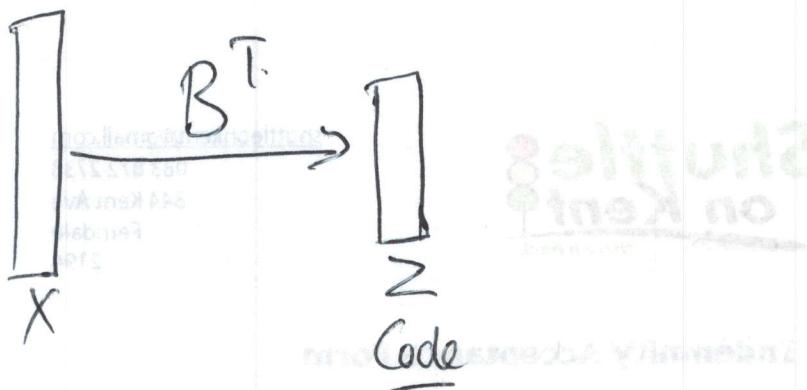
Other perspectives of PCA (optional)

We derived PCA from the perspective of minimizing the average squared reconstruction error.

However PCA can also be interpreted from different perspectives.

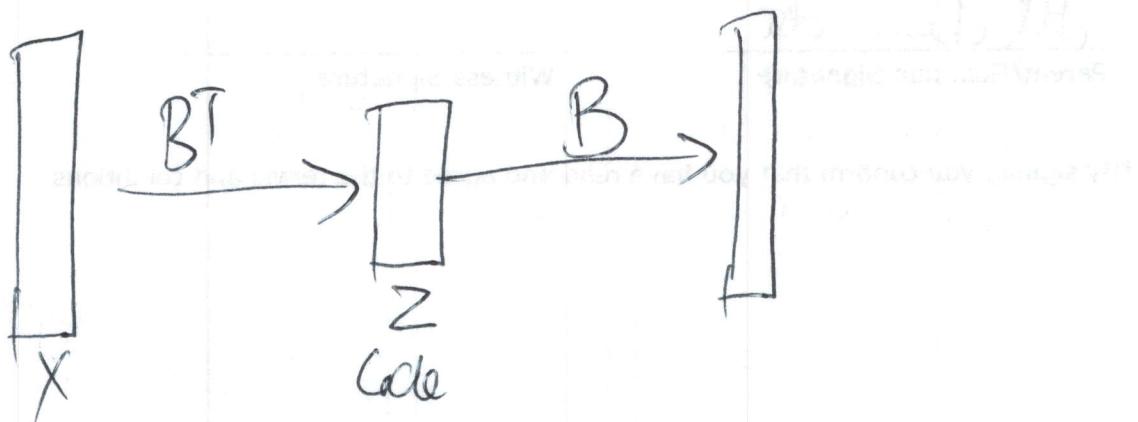
In the session we'll have a brief look at some of these interpretations. Let's start with a recap what we have done so far:

We took a high-Dim vector \mathbf{x} and we projected it onto a lower-Dim representation, \mathbf{z} , using the matrix B^T . The columns of this matrix B^T will be eigenvectors of the data covariance matrix that are associated with the largest eigenvalues.

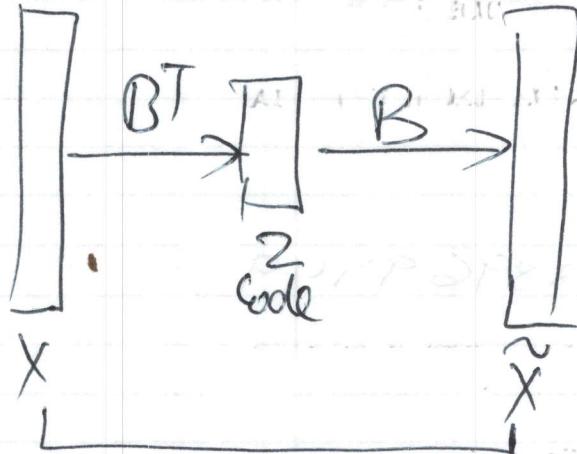


The 2 values are the coordinate of our datapoint w.r.t the basis vectors which span the principal subspace, and that is also called the code of our data point.

Once we have that lower dimensional representation Z , we can get a higher Dim of it, by multiplying B onto Z to get a higher dimensional version of Z of original data space.

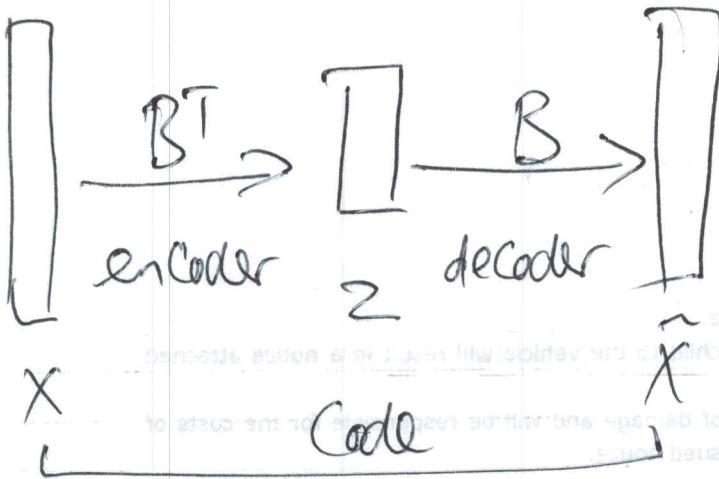


"we often set the PCA parameters, such that,
the reconstruction error between X and the
reconstruction \hat{X} is minimized"



We can also think of PCA as a linear autoencoder.
An autoencoder encodes a data point x , and
tries to decode it to something similar to the
same data point.

The mapping from the data to the code's
Called the encoder, and mapping
from the code to the original data space is
called the decoder.



DOMAINE TO PROPERTY

If the encoder and decoder are linear mappings, then we get the PCA Solution, when we minimize the Squared autoencoder loss.

If we replace the linear mapping of PCA, with a non-linear mapping, we get

an a Non-linear autoencoder.

A prominent example of this is a deep autoencoder with a linear functions of the encoder and decoder ~~is~~ are replaced with deep neural networks.

Another interpretation of PCA is related to information theory.

We can think of the code as a smaller compressed version of the original data point. When we reconstruct our original data using the code, we don't get the exact data point back, but a slightly distorted or noisy version of it.

This means that our compression is lossy.

Intuitively, we want to maximise the correlation between the original data and the lower dimensional code.

More formally this would be related to the mutual information.

We would then get the same structure to PCA we discussed earlier in this course, by

maximising the mutual information, a core concept in information theory.

6
we we derive PCA, using projections, we reformulated the average reconstruction error loss, as minimizing the variance of the data that's projected onto the orthogonal complement of the principal subspace.

Minimizing the variance is equivalent to maximizing the variance of the data when projected onto the principal subspace. If we think of variance as information contained in the data, this means that PCA can also be interpreted as a method that retains as much information as possible.

We can also look at PCA from the perspective of a latent variable model.

"we assume that we have an unknown lower dimensional data X , and we assume that we have a linear relationship

between Z and X'



Generally we can then write that $X \leftarrow Bz + \mu(m)$, and maybe some noise (ε) and assume that the noise is isotropic, with mean 0 and covariance matrix ($\sigma^2 I$)

times I)"

$$X = Bz + \mu + \varepsilon,$$

$$\varepsilon \sim N(0, \sigma^2 I)$$

"We also further assume that the distribution of this Z is a standard normal, so $P(Z)$ is gaussian with mean 0, and covariance the I matrix"

②

$$p(z) = N(0, I)$$

\downarrow

X

estimate
line no

We can now write down the likelihood of the model, "the likelihood $\rightarrow p(x \text{ given } z)$ and that's the Gaussian distribution in x , with mean $B_2 + \mu$ and Covariance matrix, sigma squared I ".

$$p(x|z) = N(x | B_2 + \mu, \sigma^2 I)$$

We can also compute the "marginal likelihood as $p(x)$ as the integral of $p(x \text{ given } z)$, times the distribution on z, dz "

- $p(x) = \int p(x|z)p(z)dz$

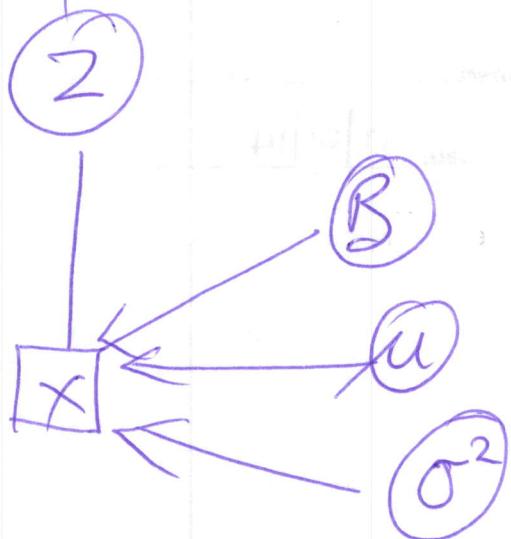
"and that turns out to be " gaussian distribution
of \mathbf{X} with mean μ and covariance matrix B times
 B^T plus sigma squared I "

$$= N(\mathbf{x} | \mu, BB^T + \sigma^2 I)$$

The parameters of this model are:

- μ
- B
- σ^2

and we can write them explicitly down in our
model up here



We can now determine the parameters of the model using maximum likelihood estimation.

We will find that μ is the mean of the data and B is a matrix that contains the eigenvectors that correspond to the largest eigenvalues.

To get the linear Prior Code of datapoint, we can apply Bayes theorem to invert the linear relationship between $\Sigma_{\text{and}} X$.

In particular we going to get, "P of Z, given X as P of X given Z, that's the likelihood (p_{xz}), times $P of Z$ divided by marginal likelihood $P of X$ ".

$$P(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

In this session, we looked at 5 different perspectives of PCA, that lead to different objectives.

- ① Minimizing the squared reconstruction error
- ② Minimizing the autoencoder loss
- ③ Maximizing the mutual information
- ④ Maximizing the variance of the projected data.
- ⑤ maximizing the likelihood in latent variable model

All five different perspectives give us the

SAME solution to the PCA problem

The strengths and weaknesses of individual perspective become more clear and when we consider properties of real data