

Notes from Understanding Variational Autoencoders

Saturday, August 27, 2022

7:07 PM

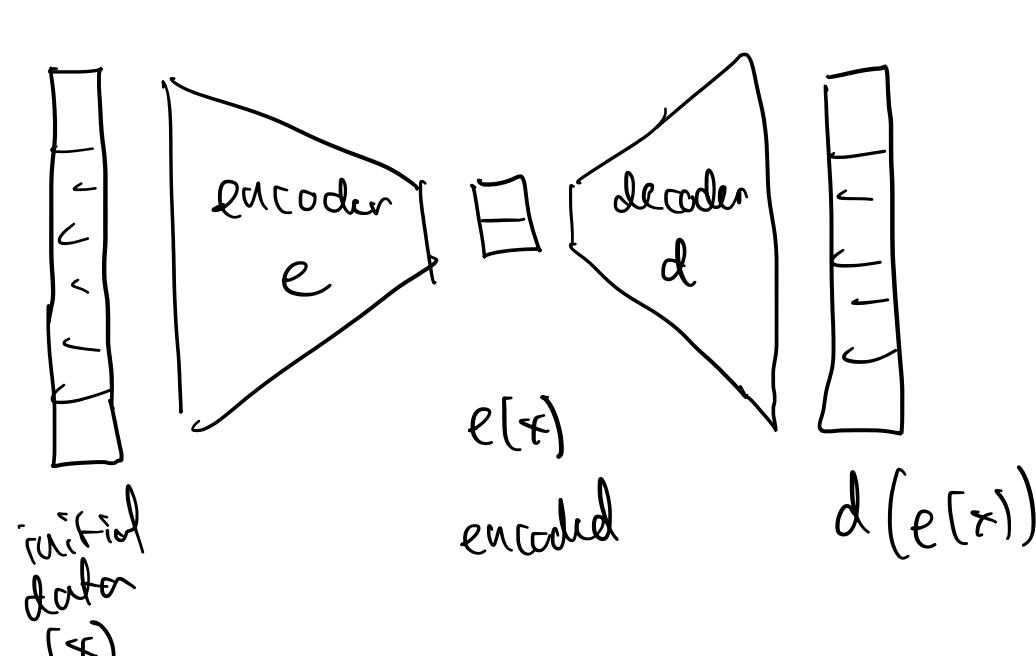
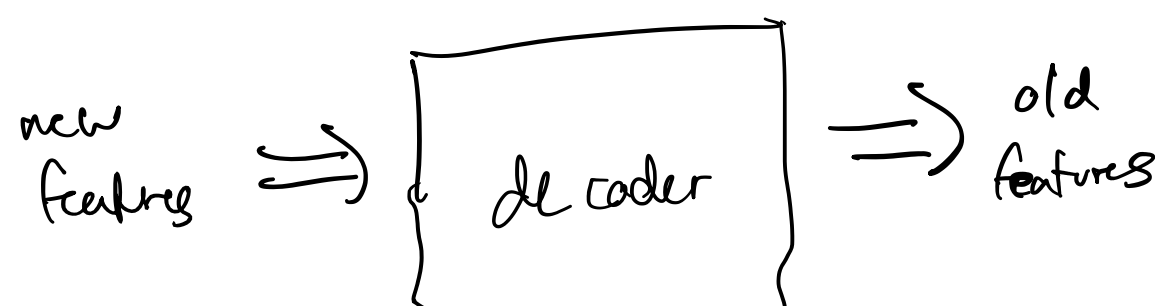
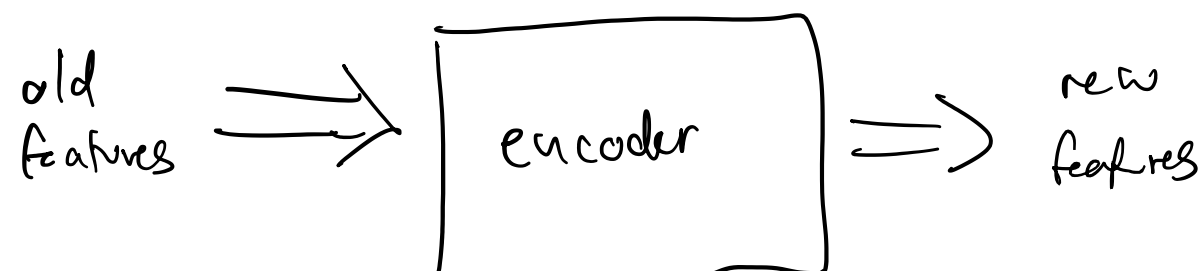
VAE - autoencoder whose encodings distribution is regularised during training to ensure latent space has good properties

notation

for R.V. z , $p(z)$ is distribution of R.V.

Dimensionality reduction, PCA, and autoencoders

- dimensionality reduction - process of reducing the number of features that describe some data (through selection or extraction)



dimensionality reduction problem

find the e and d that minimize reconstruction error

$$\{ (e^*, d^*) \mid \arg\min_{(e,d) \in E \times D} E(x, d(e(x)))$$

$$E(x, d(e(x))) \rightarrow \text{reconstruction error}$$

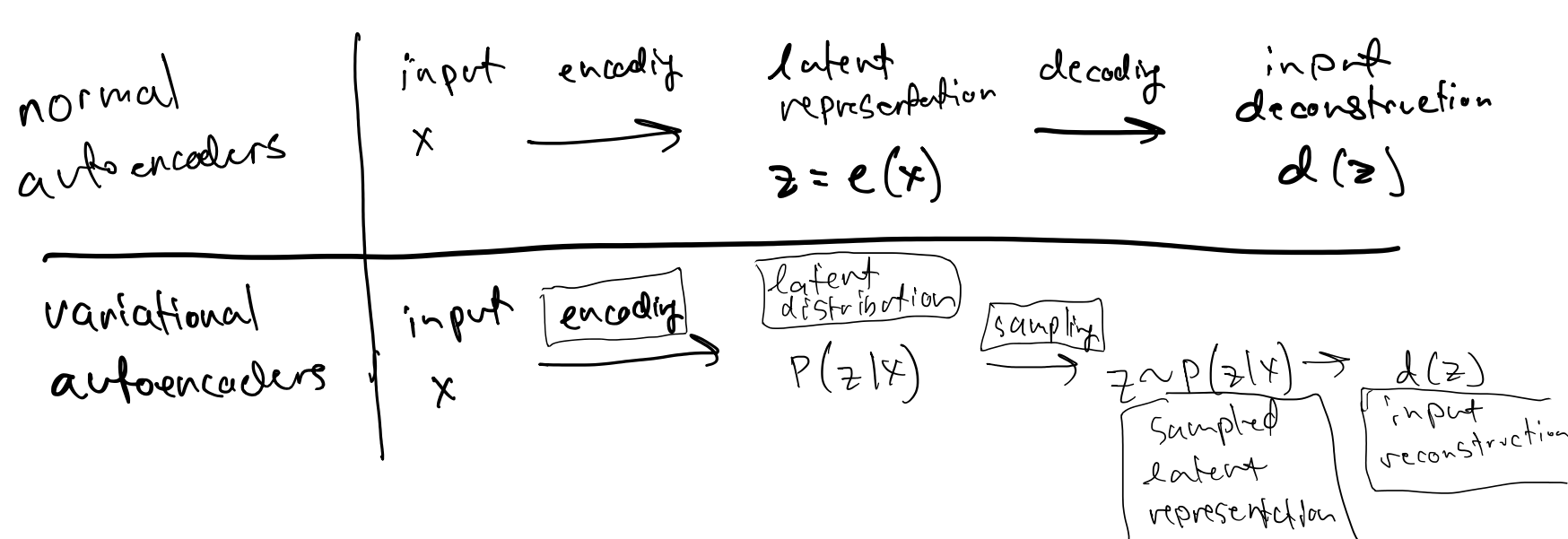
- issues w/ regular autoencoders
 - 1 - lack of interpretable latent space
 - 2 - need to keep major part of data structure information in reduced representations

Generation - sample latent space and pass it to the decoder

- difficult to guarantee w/ traditional autoencoders

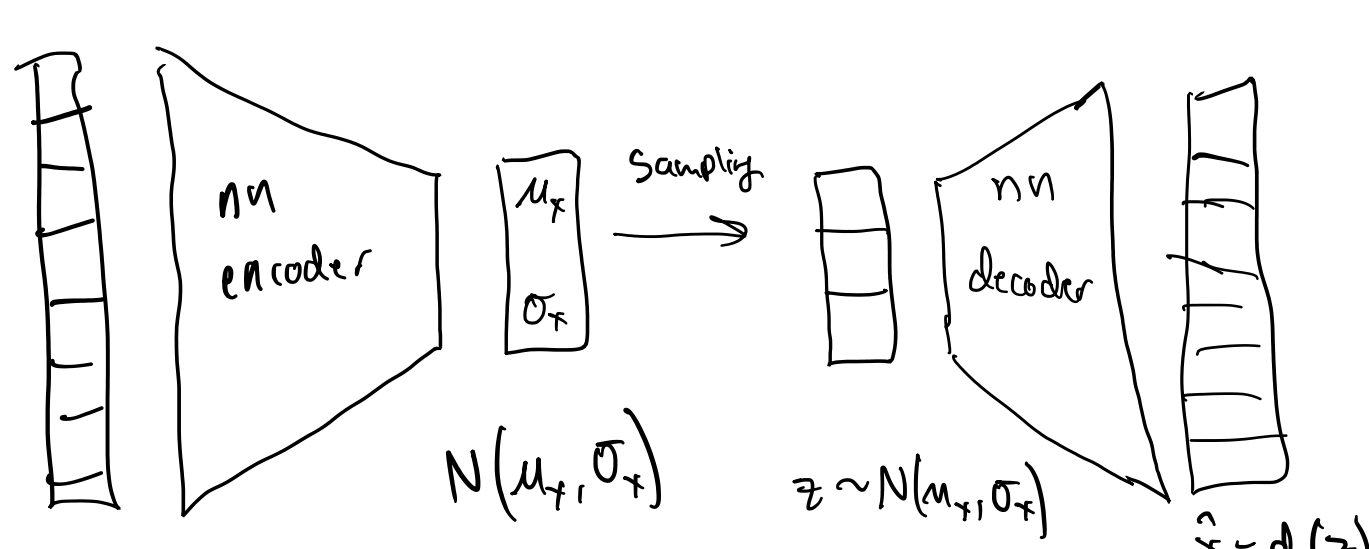
- no guarantee that model is actually learning meaningful representations

VAE: instead of encoding inputs as a single point, encode a distribution over the latent space



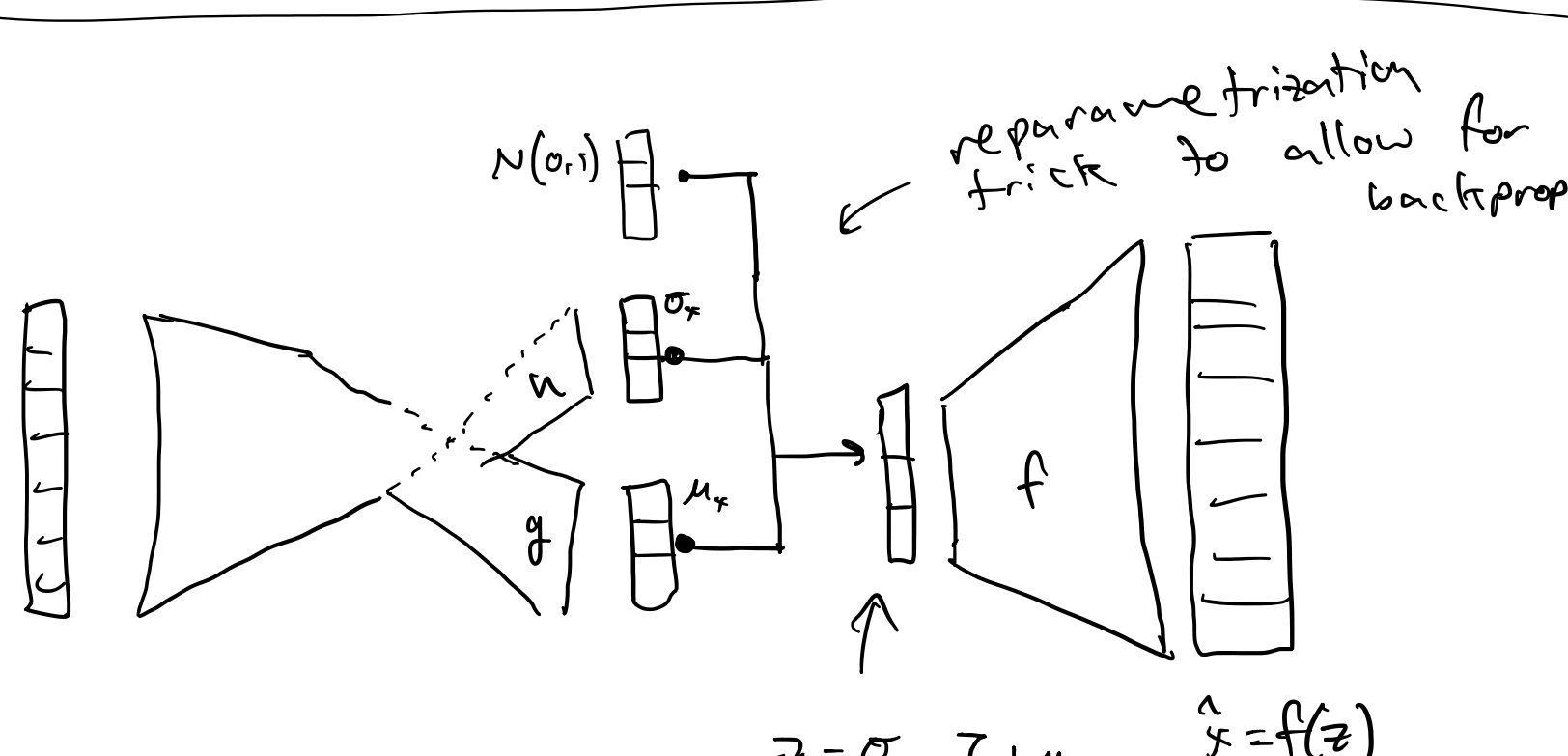
encoded distributions are chosen to be normal so encoder returns mean and covariance matrix describing distribution

VAE loss function



$$\text{loss} = \underbrace{\|x - \hat{x}\|^2}_{\text{reconstruction loss}} + \underbrace{\text{KL}[N(\mu_x, \sigma_x), N(0, 1)]}_{\text{regularisation loss}}$$

\uparrow standard normal distribution



$$\text{loss} = C \|x - \hat{x}\|^2 + \text{KL}[N(g(x), u(x)), N(0, 1)]$$