

$\underline{x} \in \mathbb{R}^d$  : Inputs

$\underline{z} = f_E(\underline{x}, \underline{\theta}_E) \in \mathbb{R}^c$  : Latent variable / representation /  
code  $\rightarrow$  hidden code of  $\underline{x}$

$$\begin{aligned}\hat{\underline{x}} &= f_D(\underline{z}, \underline{\theta}_D) = f_D(f_E(\underline{x}, \underline{\theta}_E), \underline{\theta}_D) \\ &= f(\underline{x}, \underline{\theta}) \in \mathbb{R}^d\end{aligned}$$

The output is a reconstruction of the input.

$$\underline{\theta}_E, \underline{\theta}_D, \underline{\theta} = \begin{bmatrix} \underline{\theta}_E \\ \underline{\theta}_D \end{bmatrix} \rightarrow \text{Parameters to be learnt}$$

$\underline{y} = \underline{x}$  : Self-reconstruction

$\Rightarrow$  Self-supervised learning

$\subset$   
subset

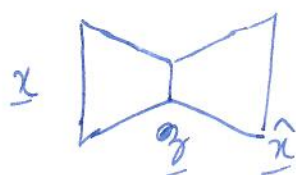
unsupervised learning

Supervised learning using unlabelled data (see Ch 10)

Encoder/Decoder:

- any DNN
- Symmetric structure

Always under complete AE  $\rightarrow c < d$  [eg image compression]



NN learns the most important  
bottleneck  $\rightarrow$  attributes

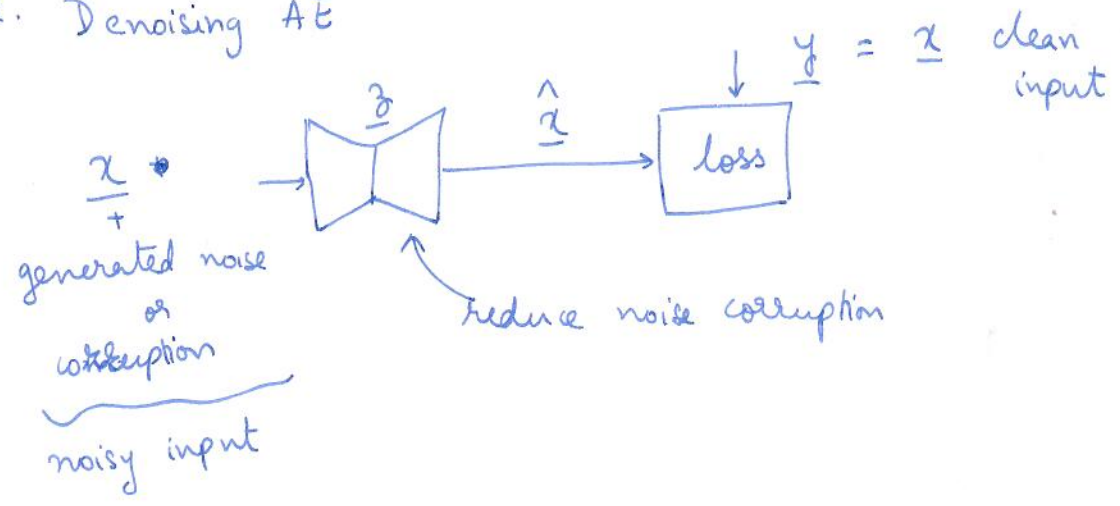
Overcomplete AE  $\rightarrow c > d \rightarrow$  never used

9-5

9-6

## Applications of Auto Encoders

1. dimension reduction (data compression)  $c < d$
2. Denoising AE



$\underline{z}$  : low-dimensional

→ keep only relevant information for clean input  $\underline{x}$

→ drop noise / corruptions information as it is not there in output

9-7

⋮

9-9

Performs good on known examples but not good on unknown domains which are similar to the others but a lot vast also.

Theory of VAE:

$p(\underline{x})$ : unknown distribution of  $\underline{x}$ ,  
only samples  $\underline{x}(1) \dots \underline{x}(N)$  available

We don't have  $\underline{y}$

$q(\underline{x}, \underline{\theta})$ : parametric approximation of  $p(\underline{x})$

$\underline{\theta}$ : parameter vector of VAE

Goal:  $\min_{\underline{\theta}} D_{KL}(p||q)$

As in ch 3.4:

$$\text{Cost f: } \mathcal{L}(\underline{\theta}) = D_{KL}(p||q) \\ = E_{\underline{x} \sim p} \left( \ln \frac{p(\underline{x})}{q(\underline{x}; \underline{\theta})} \right)$$

$$= \underbrace{E_{\underline{x} \sim p} \ln p(\underline{x})}_{\text{independent of } \underline{\theta}} - \underbrace{E_{\underline{x} \sim p} \ln q(\underline{x}; \underline{\theta})}_{\text{unknown}}$$

Replace  $p(\underline{x})$  by  $\hat{p}(\underline{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\underline{x} - \underline{x}(n))$

$$\mathcal{L}(\underline{\theta}) \approx \text{constant} - E_{\underline{x} \sim \hat{p}} \ln q(\underline{x}; \underline{\theta})$$

$$= \text{const} + \frac{1}{N} \sum_{n=1}^N \underbrace{-\ln q(\underline{x}(n); \underline{\theta})}_{\text{loss } \mathcal{L}(\underline{x}(n); \underline{\theta}) \text{ of VAE}}$$

to the derivation  
Diff ~~to~~ supervised learning in ch 3.4:

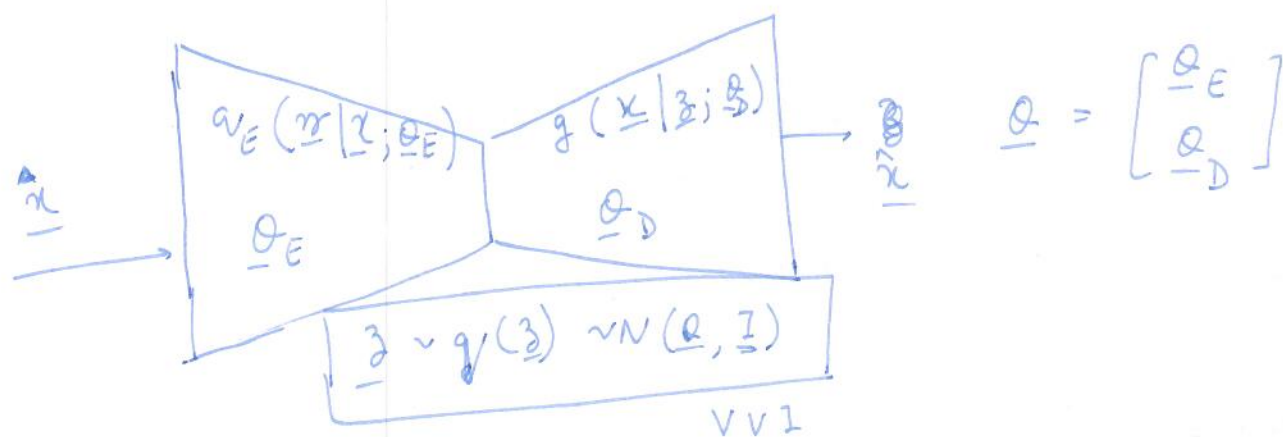
No  $\underline{y}$   $p(\underline{x}, \underline{y}) \rightarrow p(\underline{x})$   $q(\underline{y}|\underline{x}; \underline{\theta}) \rightarrow q(\underline{x}; \underline{\theta})$

6

How to model  $q(\underline{x}; \underline{\theta})$ ?

VAE:  $\underline{x}$  is generated by  $\underline{z}$  with a known distribution

$$\underline{z} \sim q(\underline{z}) \quad \text{eg } \underline{z} \sim N(\underline{0}, \underline{I})$$



$$\begin{aligned} \text{i.e. } q(\underline{x}; \underline{\theta}) &= \int p(\underline{x}, \underline{z}; \underline{\theta}) d\underline{z} \\ &= \int \underbrace{p(\underline{x} | \underline{z}; \underline{\theta})}_{\text{decoder}} \underbrace{q(\underline{z})}_{\text{known}} d\underline{z} \end{aligned}$$

Since the distribution of the decoder is unknown due to multiple neural networks, the integration is hard to calculate. So, we can't calculate  $\nabla L$

1<sup>st</sup> trick of sol<sup>n</sup>:

Replace  $-\ln q(\underline{x}; \underline{\theta})$  by its Variational upperbound.

Then we minimize the Variational upperbound.

9-13

14

15