

AUTOENCODERS

- Try to reconstruct input as its output
- Minimise reconstruction loss (diff b/w x, x')
- Encoder: $Z = g(w_e z + b_e)$
- Decoder: $x' = f(w_d z + b_d)$
- Data can be represented non-linearly (unlike PCA)
- No interpretability
- Self-supervised → does prediction (unlike unsupervised) data labels implicit in data (predict input from input)

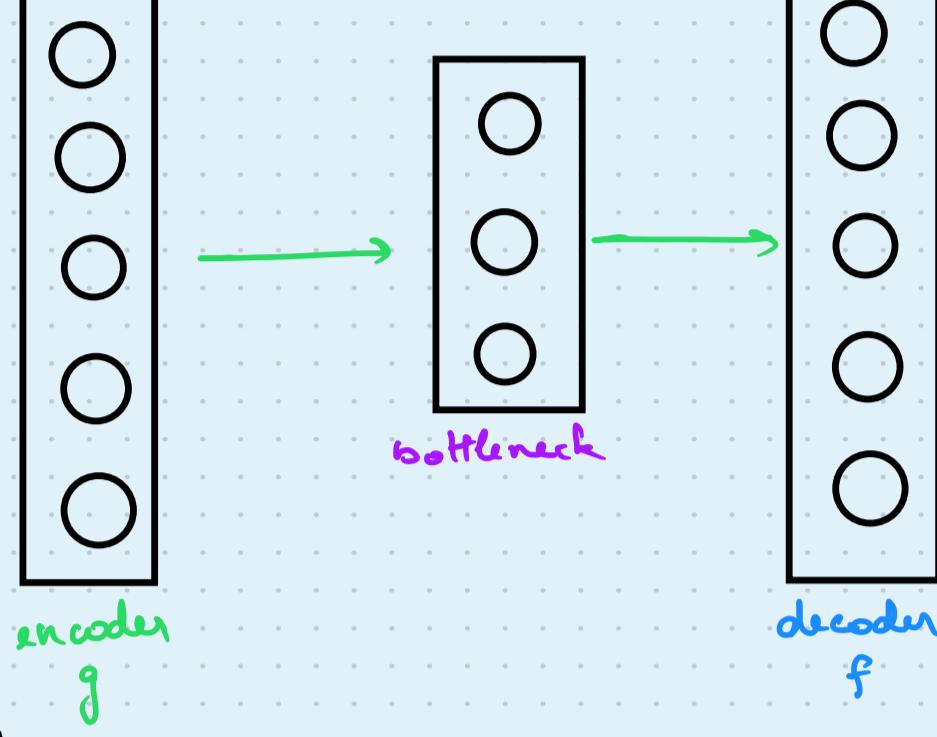
Limitations of basic AE

- Can memorise training data instead of generalising
- [from above] sensitive to input perturbations
- No structure in latent space

UNDERCOMPLETE AE

Bottleneck layer → ensures $\dim(\text{latent space}) \ll \dim(\text{input space})$

Thus prevents memorisation



loss function:

$$\ell(x, x') = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - x'^{(i)})^2$$

$$= \frac{1}{N} \sum_{i=1}^N [x^{(i)} - f(g(z^{(i)}))]^2$$

{ Basically MSE

SPARSE AE

- Only looks for "most important clues" → activate sparse subset of hidden neurons acc. to most relevant features

How to pick relevant features? → L1 penalisation

$$L_1 = \text{Reconstruction loss} + \lambda \|\mathbf{h}\|_1$$

no. of hidden activations
controls sparsity

Implementation key:

Small hidden layer → limits capacity, ensures latent space dimensionality is low

Sparcity constraints → further limit active neurons

Effect:

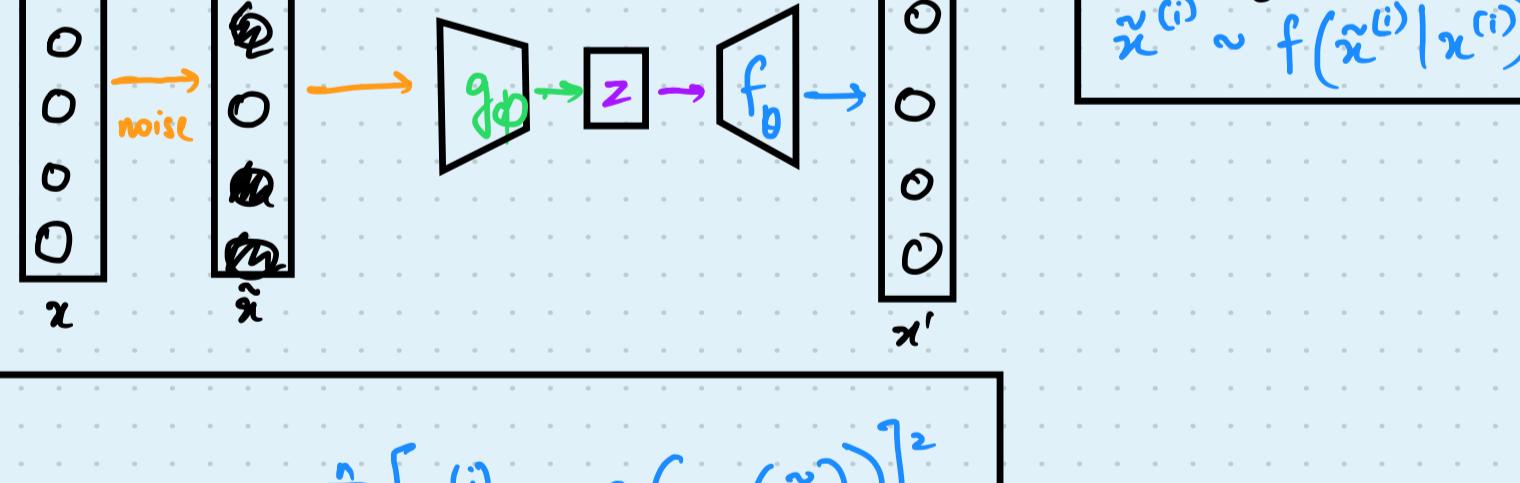
- Forces model to prioritise features essential for reconstruction
- Encourages disentangling → separating distinct aspects of data into different neurons

Trade-off:

- Higher compression → better feature extraction
- Higher reconstruction loss

DENOISING AE

- Able to reconstruct original input from corrupted input
- More robust model



loss:

$$\ell(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [x^{(i)} - f_\theta(g_\phi(\tilde{x}^{(i)}))]^2$$

clean input

encoder-decoder on noised input

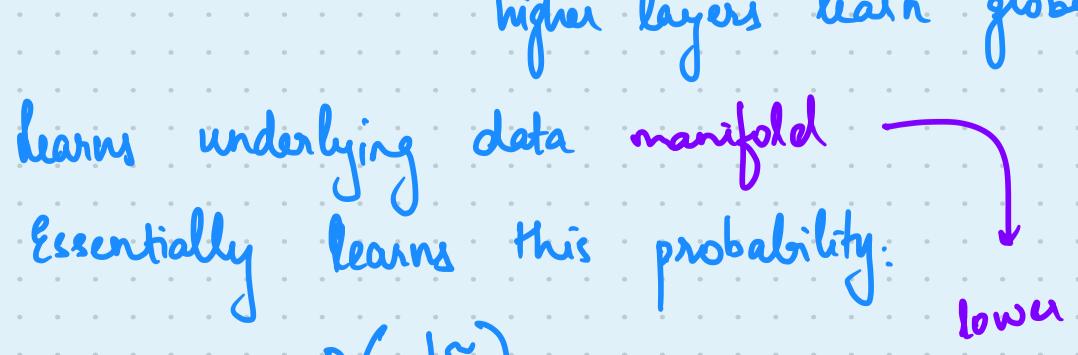
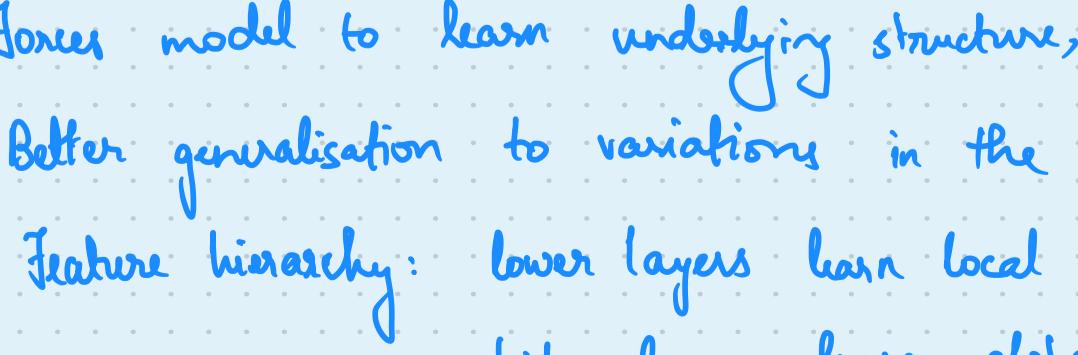
Types of noise

① Gaussian: Add random values from $N(\mu, \sigma)$

② Salt and pepper: Set random pixels to 0 or 1

③ Masking: Hide random parts

④ Dropout: Randomly set activations to 0



Why this works

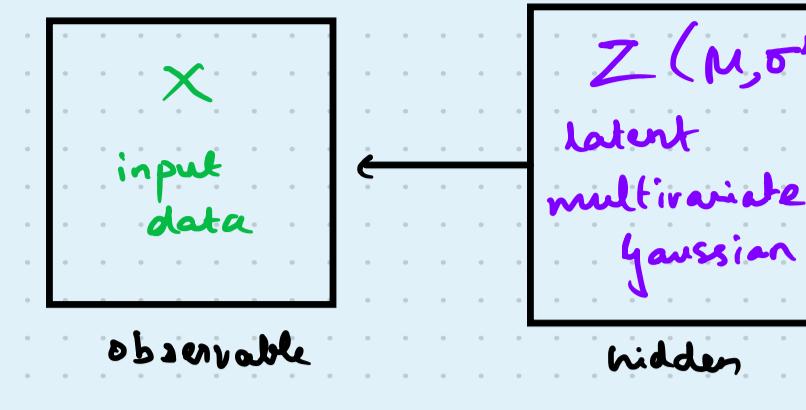
- Forces model to learn underlying structure, cannot memorise
- Better generalisation to variations in the test set
- Feature hierarchy: lower layers learn local features
- higher layers learn global features

- learns underlying data manifold
- Essentially learns this probability: $P(x|\tilde{x})$

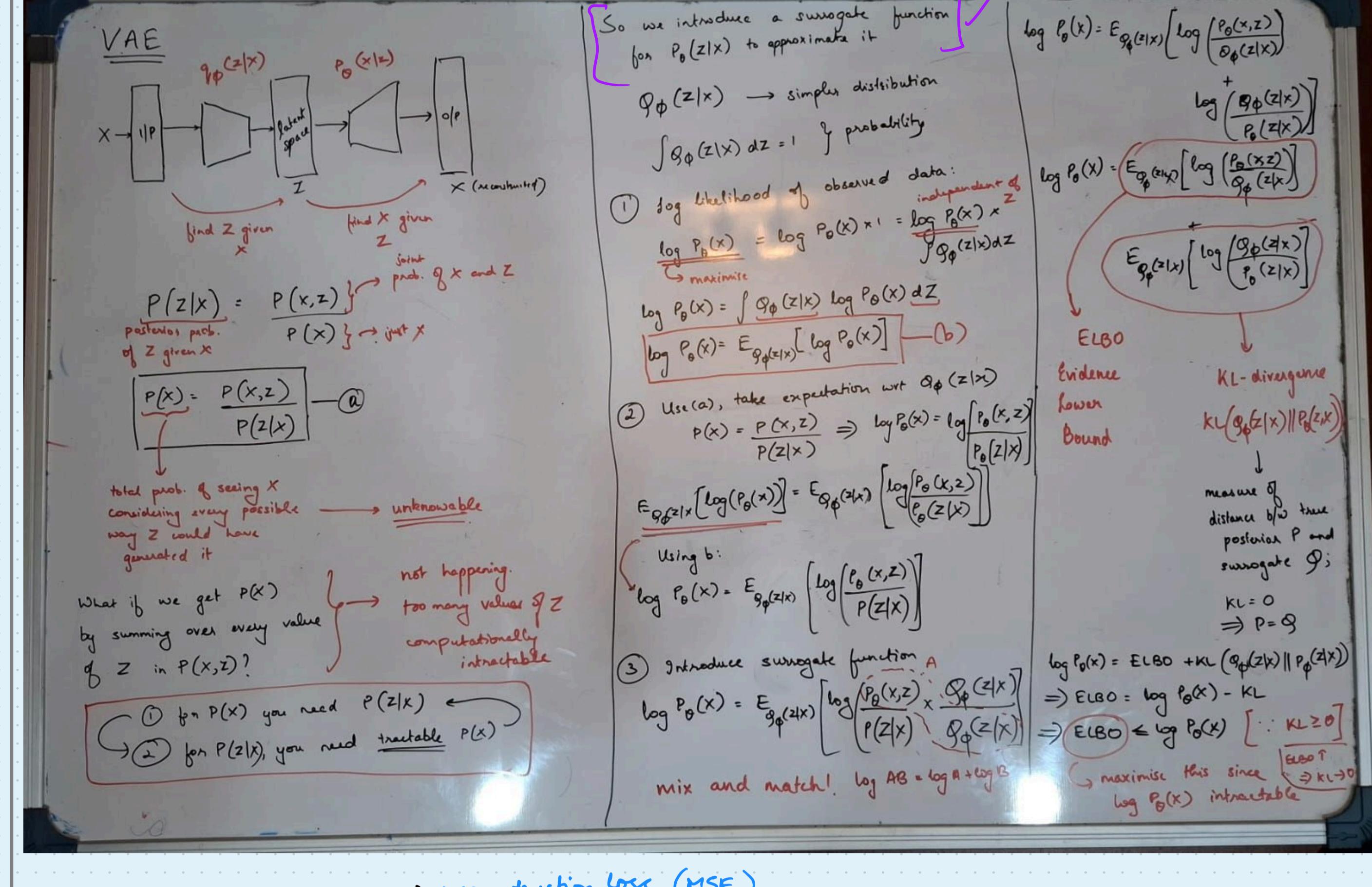
lower dimensional space
within entire possible input space
that has meaningful data

Variational Autoencoder (VAE)

- Generative model: learn a latent space distribution and sample from that instead of reconstructing
- Assume latent space to be some random variable Z
- Also assume input X is conditioned on Z .
- $Z \rightarrow$ modelled as multivariate Gaussian



Tweak Z according to assumptions to learn distribution of X



ELBO → reconstruction loss (MSE)
KL divergence b/w prior Z and learnt Z

Reparameterisation Trick

maximisation of ELBO → calculate gradient.

$$\text{ELBO} = E_{q_\phi(z|x)} \left[\log \left(\frac{P(x,z)}{q_\phi(z|x)} \right) \right]$$

Calculating gradient with Z → Z → stochastic, random values
very difficult → gradient shows very high variance

sampling is a non-differentiable operation

Solution: Separate the randomness from model parameters!

① Z supposed to be from $N(\mu, \sigma)$

where μ, σ learnt by encoder ($P(z|x)$)

② Sample an external source of noise for randomness:

$$E \sim N(0, 1)$$

std. normal

③ Calculate latent Z using a deterministic function.

$$Z = \mu + \sigma E$$