

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} \quad \left| \quad p(x|y) = \frac{p(y|x)p(x)}{p(y)} \right. \quad p(y|x)$$

Generative Models

(Recap of Autoencoder, Variational Autoencoder)

$$p(y|x) \propto p(x|y)p(y)$$

$$p(x|y) \propto \boxed{p(y|x)p(x)}$$

Discriminative Model

$$p(y|x) \leftarrow$$

Generative Model ?

$$p(x,y)$$

$$p(\underline{x}) \stackrel{\text{?}}{=} p(x_1, x_2, x_3, \dots, x_d)$$

X

① \underline{x} : Graph structured
(relational) data

RNN

Tirtharaj Dash

Dept. of CS & IS and APPCAIR
BITS Pilani, Goa Campus

December 02, 2021

different kinds of \underline{x}

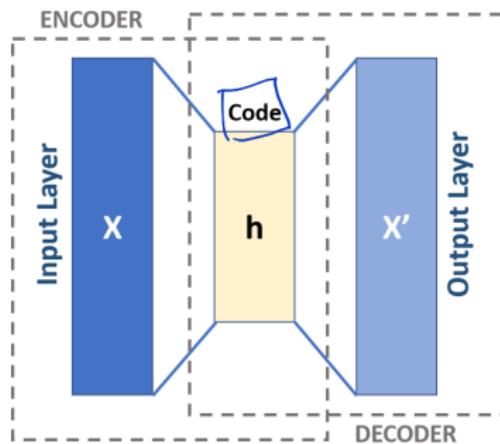
① \underline{x} : Tabular structure
MLP every row in
the table
rep. an instance

② \underline{x} : visual data
 CNN (spatial)

③ \underline{x} : Temporal data.
 RNN (sequences)

Autoencoder – Discriminative Model

- AE is a neural network that learns to copy its input to its output.
- It has an internal (hidden) layer that describes a **code** used to represent the input.



Autoencoder – Discriminative Model

- It is constituted by two main parts: an encoder that maps the input into the code, and a decoder that maps the code to a reconstruction of the input.
- However, simply copying input to output would just *duplicate* the signal (rather than generalising). Why?
- Instead, AE reconstructs the input approximately, preserving the most relevant aspects of the data (we can call this: some important latent aspects).

Autoencoder – Discriminative Model

- Let \mathbf{x} be an input example. The encoder and decoder do the following:

$$\underline{Enc : \mathcal{X} \mapsto \mathcal{H}}$$

$$\underline{Dec : \mathcal{H} \mapsto \mathcal{X}}$$

- Where, Enc and Dec are the functions obtained by minimising a reconstruction loss.

$$\underline{\mathcal{L}(\mathbf{x}, \hat{\mathbf{x}})}$$

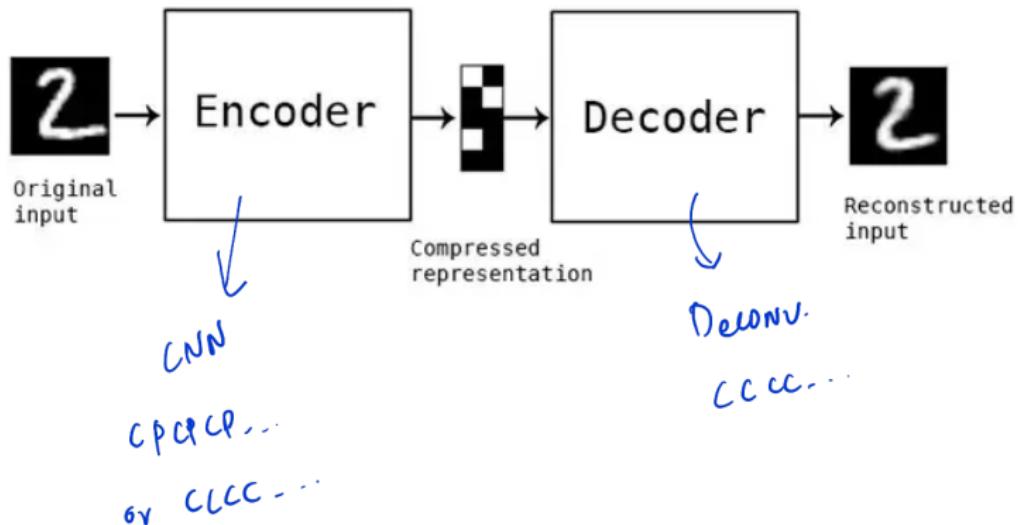
$$Enc, Dec = \arg \min_{Enc, Dec} \|\mathbf{x} - (\underline{Dec} \circ \underline{Enc})(\mathbf{x})\|^2$$

- In the simplest case, both encoder and decoder are single layered. That is:

$$\underline{\mathbf{h} = \sigma(\mathbf{w}\mathbf{x} + b)}$$

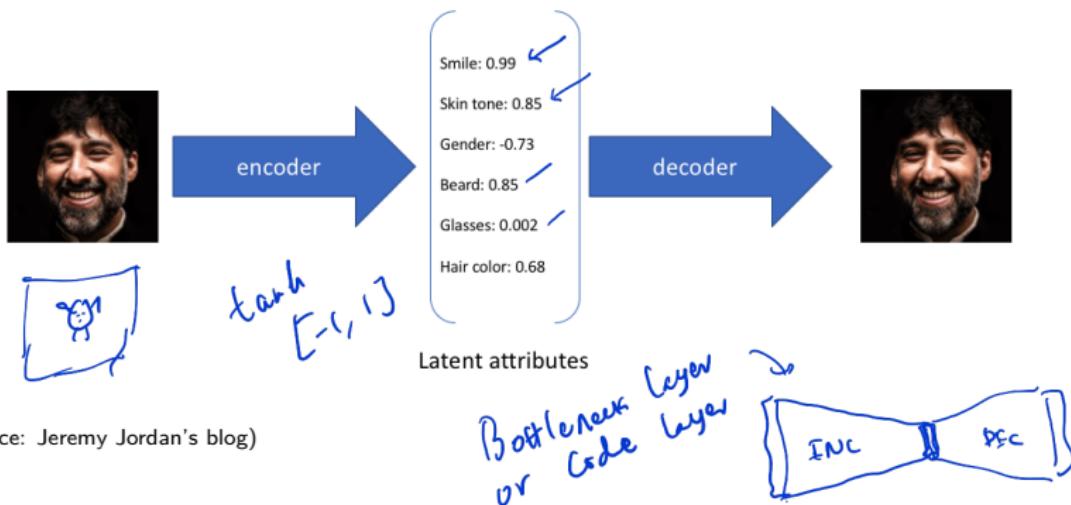
Autoencoder – Discriminative Model

- \mathbf{h} is referred to as code or latent variable or latent representation.



Autoencoder – Discriminative Model

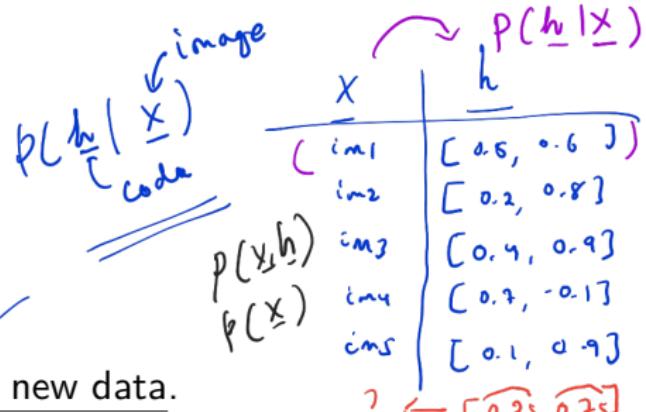
- Each hidden dimension represents some latent feature learned about the input.
- For example (the features mentioned are hypothetical for demonstration purpose only):



(Fig source: Jeremy Jordan's blog)

Autoencoder – Discriminative Model

- This is a discriminative model.
- AE can be used for:
 - compressing data
 - greedy layerwise pre-training
- AE cannot be used for: generating new data.
- For generating new data, the model needs to learn a joint distribution of some kind $p(\mathbf{x})$ or $p(\mathbf{x}, \mathbf{h})$.
- Or, a model can be considered as “generative” when the input latent variable has probability distribution associated with it – a kind of autoencoder that does this is ‘Variational AE’.

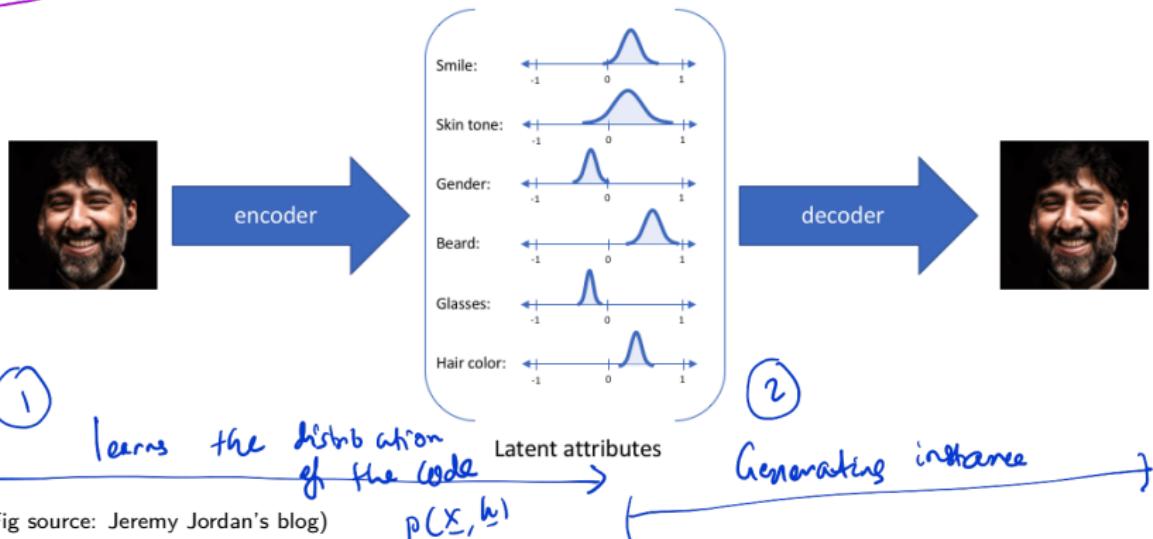


Given $\underline{\mathbf{x}}$, learns ~~a point estimate $\underline{\mathbf{h}}$~~ $\underline{\mathbf{h}}$ \leftarrow distribution.

VAE

"Variational" AE : Generative AE

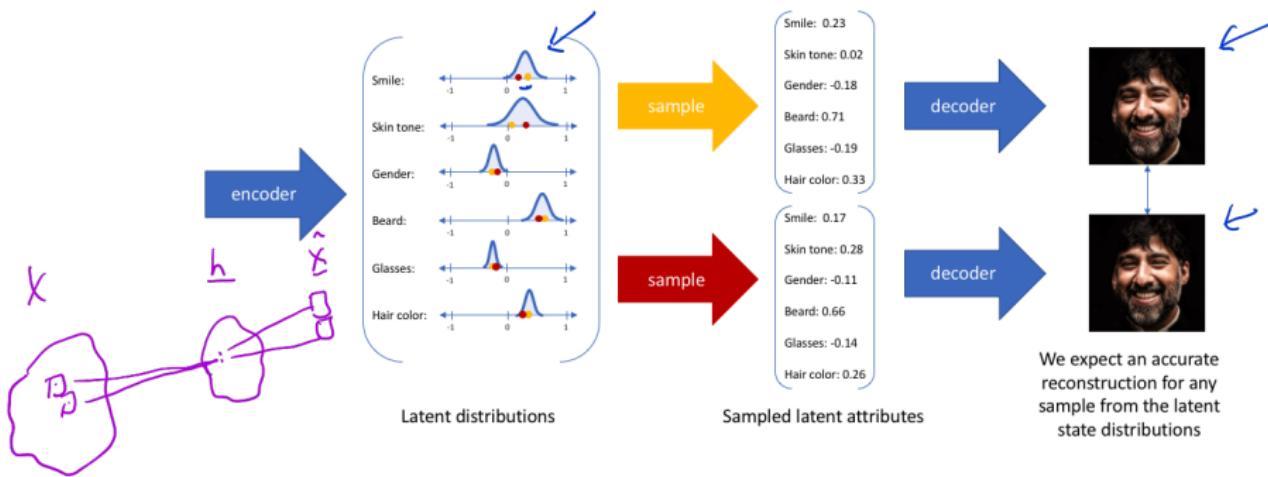
- VAE imposes a specific probabilistic structure on the hidden units.



- The AE network is sometimes called 'recognition model' whereas the decoder network is sometimes referred to as the 'generative model'.

VAE

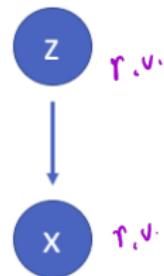
- The encoder model outputs a range of possible values (a statistical distro) from which, we can randomly sample and input to the decoder to re-construct the input. This enforces a continuous and smooth latent space representation.
- It is expected that values that are close enough in the latent space would result in similar reconstructions.



VAE

(Fig source: Jeremy Jordan's blog)

- Suppose there exists some hidden variable z which generates an observation x . As a Bayesian network, it looks like:



- The difficulty is that we know nothing of z , we can only see x . But, we can infer some characteristics of z :

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

Diagram annotations:

- A green oval encloses $p(z|x)$.
- A blue oval encloses $p(x)$.
- An arrow labeled "likelihood" points from $p(x|z)$ to $p(z|x)$.
- An arrow labeled "prior" points from $p(z)$ to $p(z|x)$.
- An arrow labeled "marginal" points from $p(x)$ to $p(x)$.
- Below the equation, $p(z|x) \propto p(x|z)p(z)$ is written.
- To the right of the equation, $\int p(x|z)p(z) dz$ is written under a horizontal line.

$x \rightarrow z$ →
unknown.

VAE

problem Computing the denominator $p(x)$ is hard:

Exact estimate for
 $p(z|x)$ is
not possible.

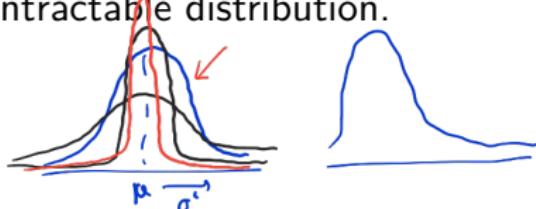
$$p(x) = \int_{\mathbb{R}^d} p(x|z)p(z)dz$$

solution Variational inference

- Let's approximate $p(z|x)$ by another distribution $q(z|x)$ such that it has a tractable distribution.
- If we can define the parameters of $q(z|x)$ such that it is very similar to $p(z|x)$, we can use it to perform approximate inference of the intractable distribution.

approximate if

$\phi \xrightarrow{(1)} q$
 $\downarrow (2)$
params?



VAE

$$p \xrightarrow{\text{distance}} q$$

$$D_{KL}(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$

- If we minimise KL divergence between $q(z|x)$ and $p(z|x)$, then these two distros will become similar to each other.

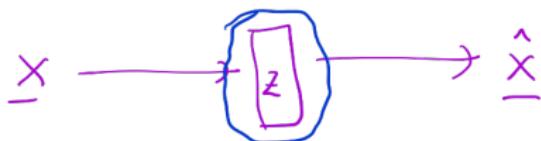
$$\sum q(z|x) \log \frac{q(z|x)}{p(z|x)} \rightarrow \min KL(q(z|x) \parallel p(z|x))$$

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

Which in turn can be solved by maximising the following (not showing the proof here):

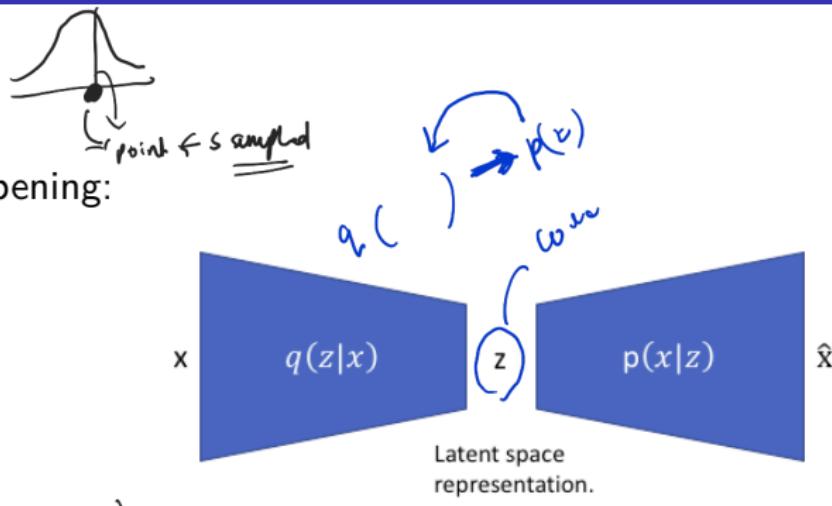
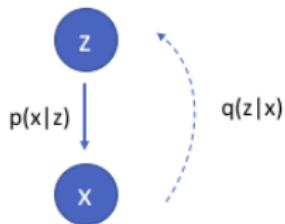
$$E_{q(z|x)} \log p(x|z) - KL(q(z|x) \parallel p(z))$$

The first term represents reconstruction likelihood. The second term enforces that q distro is similar to the true distro $p(z)$.

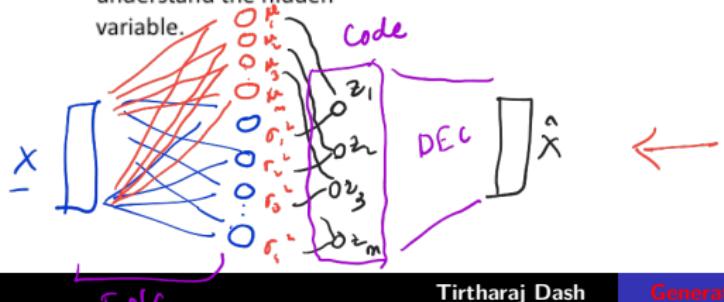


VAE

- This is what is happening:



We'd like to use our observations to understand the hidden variable.



Neural network mapping x to z .

Neural network mapping z to x .

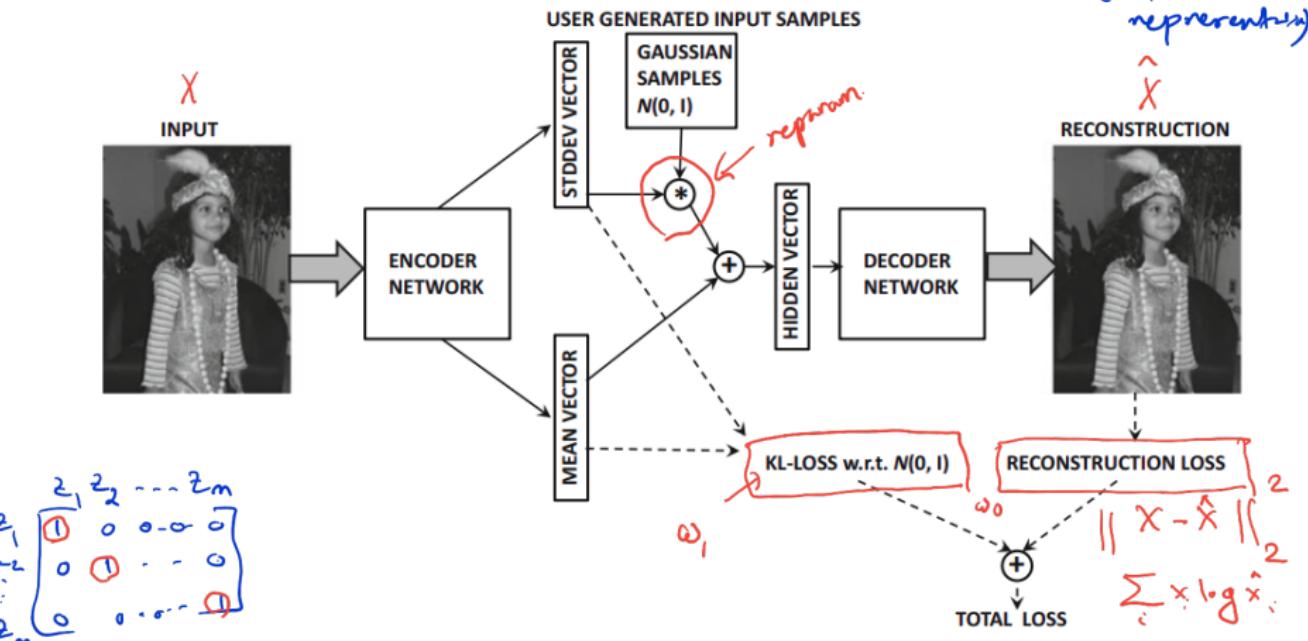
$$\begin{aligned}
 & p(z) \\
 & \downarrow \\
 & q(z|x) \quad \text{parameters} \\
 & = \\
 & q : \mathcal{N}(\mu, \sigma^2)
 \end{aligned}$$

VAE

- For VAE, we assume that the 'true' distribution $p(z)$ is Gaussian $\mathcal{N}(0, I)$ (I is the identity matrix).
- Therefore, VAE's loss function is this:

~~distr.~~ of z

(latent representation)



VAE



(Fig source: Charu Agarwal's NNDL book)

- The loss function can be written as the sum of the reconstruction loss, and sum of KL divergences over all dimensions of the latent space \mathbf{z} :

$$L(\mathbf{x}, \hat{\mathbf{x}}) + \sum_j KL(q_j(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

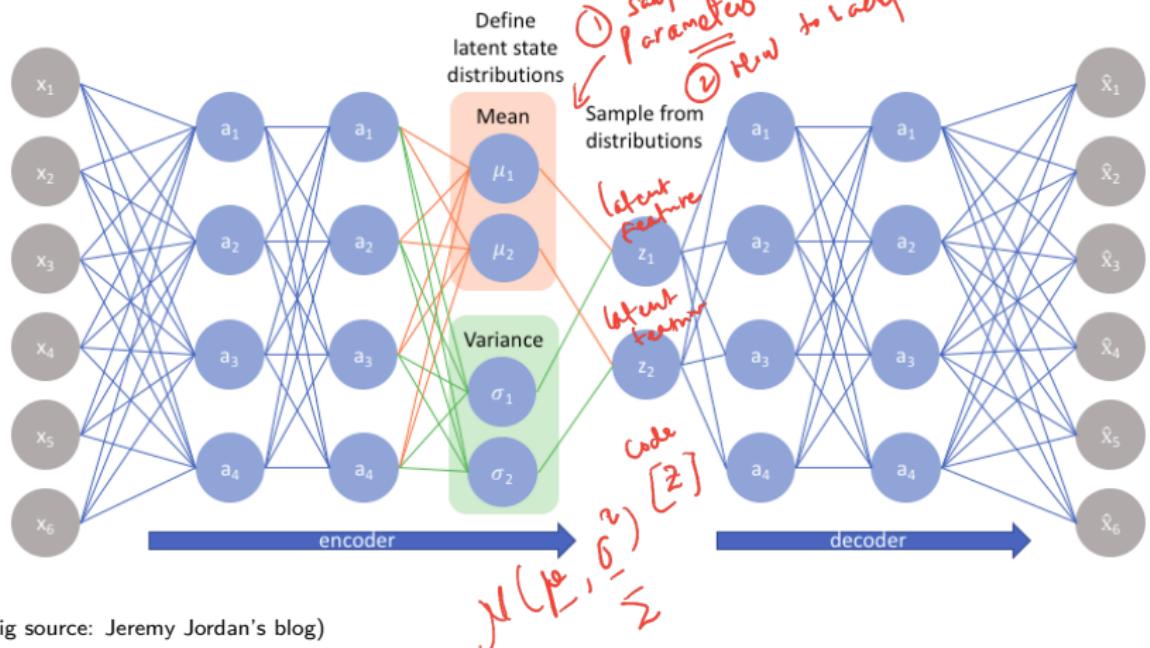
*Recon.
loss*

$$KL(q || p)$$

$$= -\sum q \log \frac{q}{p}$$

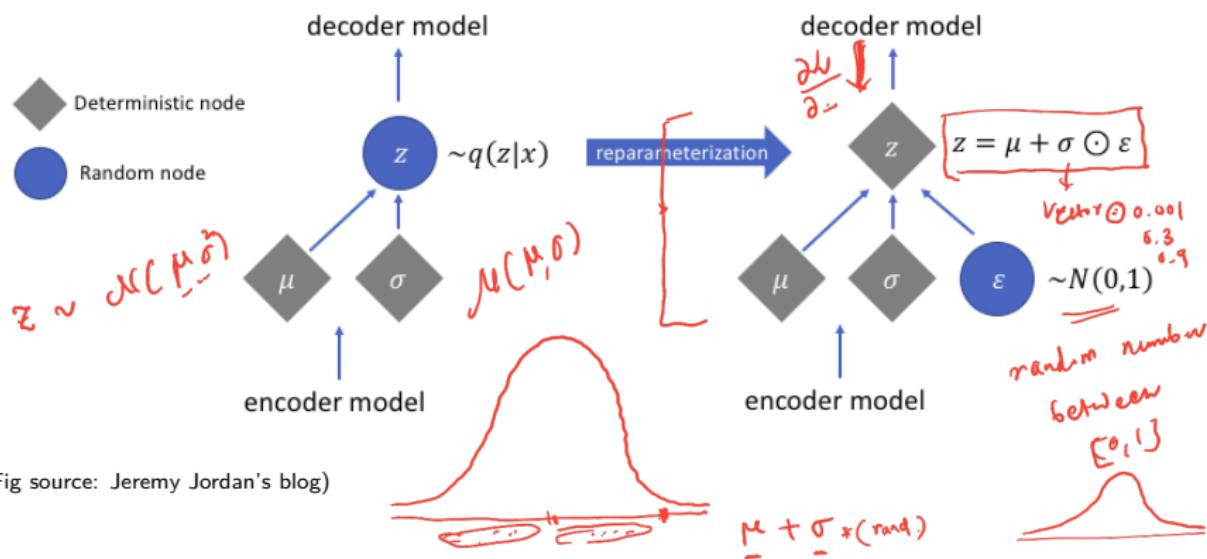
VAE

- So, if we see closely at these the NN structure is:



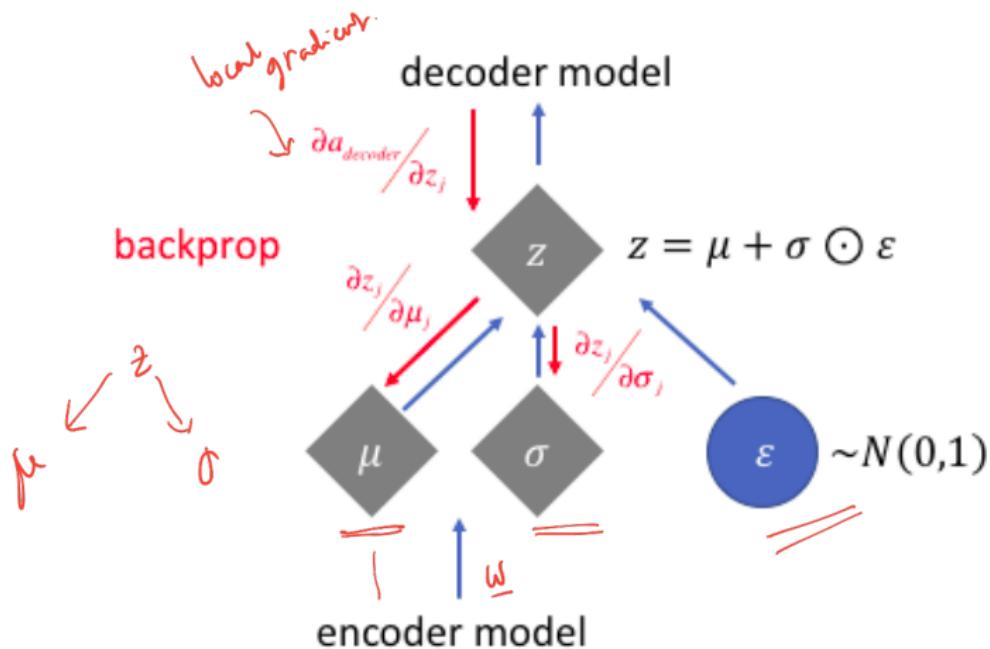
VAE

- Note z 's are sampled from a normal distro. But, backprop cannot handle sampling directly. We use a trick called 're-parameterisation trick' to ensure that backprop can go inside sampling.
- This is:



(Fig source: Jeremy Jordan's blog)

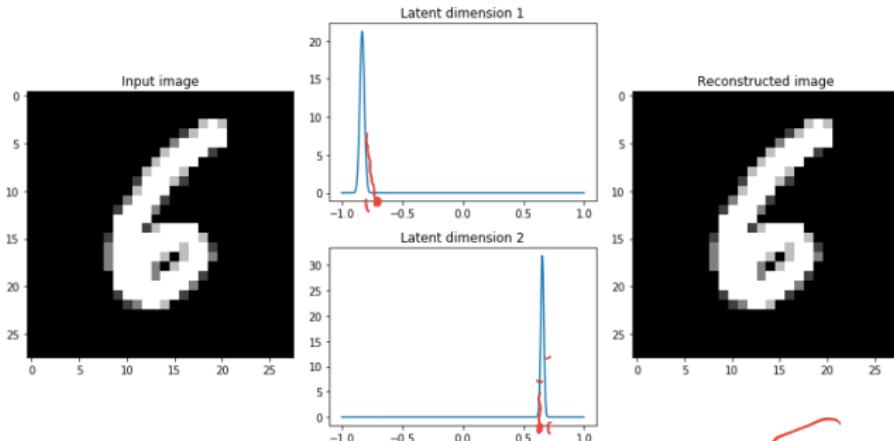
VAE



(Fig source: Jeremy Jordan's blog)

VAE

- A simple experimental visualisation with 2 latent feature dimensions:



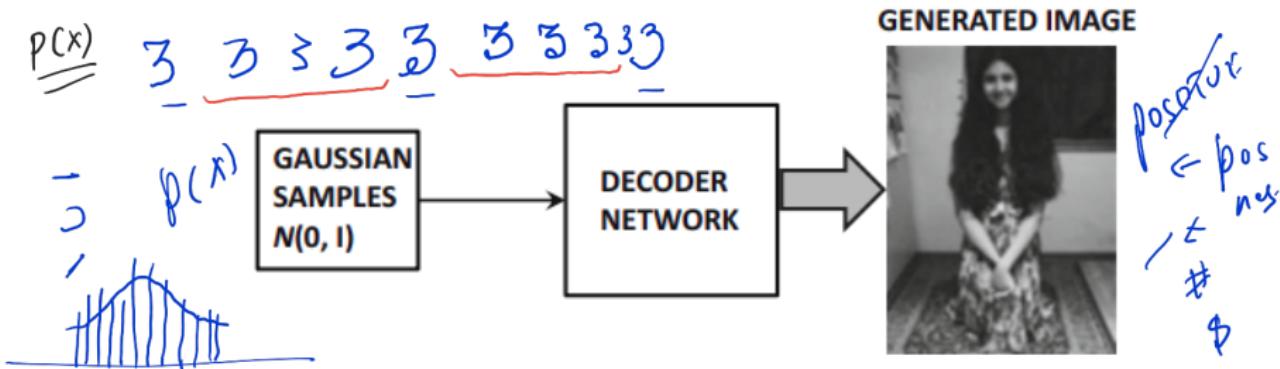
Loop ? • $z_1 = -0.86$

Strange
line ? • $z_2 = 0.1$



VAE – as generative model

- The decoder network can be used as a generative model capable of generating new data similar to what was observed during training.
- Specifically, we will be sampling data from the (prior) distribution $p(z)$ that was constructed during training.



(Fig source: Charu Agarwal's NNDL book)

- In practice:
 - Requires large amount of data
 - Requires large amount of training iterations

