

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

$$P(y|x)$$

## Generative Models

(Recap of Autoencoder, Variational Autoencoder)

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

$$P(x|y) \propto \frac{P(y|x)P(x)}{P(x,y)}$$

Tirtharaj Dash

Discriminative Model

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

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

Generative Model?

$$P(\underline{x}, y)$$

$$P(\underline{x}) \rightarrow P(x_1, x_2, x_3, \dots, x_d)$$

different kinds of  $\underline{x}$

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

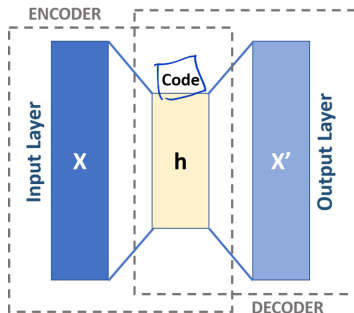
✓ ②  $\underline{x}$ : visual data (spatial)  
CNN

✓ ③  $\underline{x}$ : Temporal data (sequential)  
RNN, Trans

X  
④  $\underline{x}$ : Graph-structured (relational) data  
GNN

# 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:

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

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

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

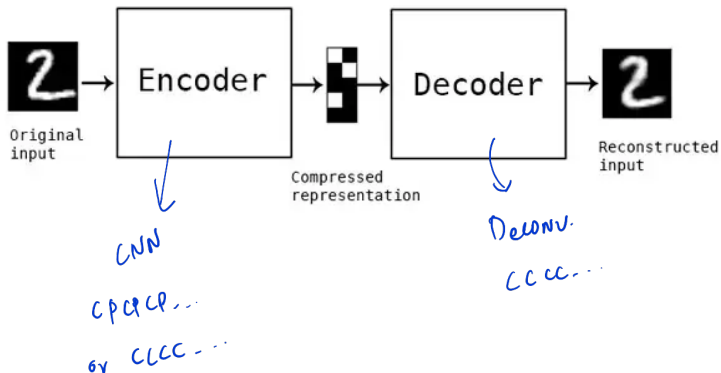
$$Enc, Dec = \arg \min_{Enc, Dec} \mathcal{L}(\underline{\mathbf{x}}, \underline{\hat{\mathbf{x}}}) ||\mathbf{x} - (\underline{Dec} \circ \underline{Enc})(\mathbf{x})||^2 \quad \leftarrow$$

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

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

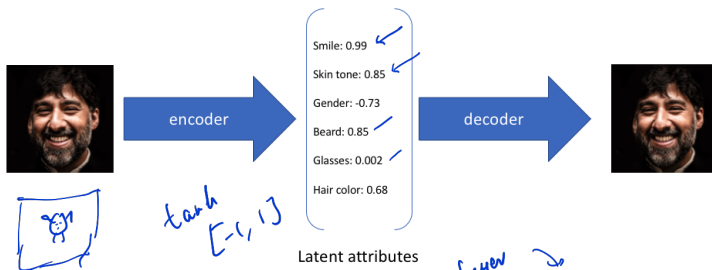
# Autoencoder – Discriminative Model

- $h$  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):



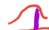

(Fig source: Jeremy Jordan's blog)

# 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’.

Handwritten notes:  $p(h|x)$  (code),  $p(x, h)$ ,  $p(x)$

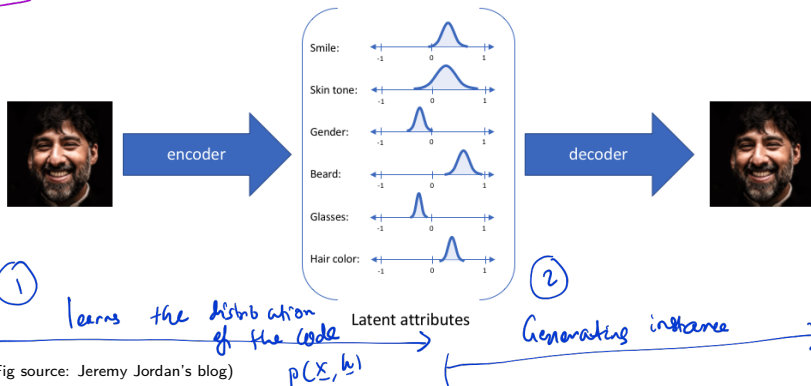
$x$	$h$	$p(h x)$
im1	[0.5, 0.6]	
im2	[0.2, 0.8]	
im3	[0.4, 0.9]	
im4	[0.7, -0.1]	
im5	[0.1, 0.9]	
?	[0.35, 0.75]	

Handwritten notes: ? ←  

AE  
 Given  $x$ , leaves ~~a point estimate~~  $h$  ← distribution.

# "Variational" AE : Generative AE

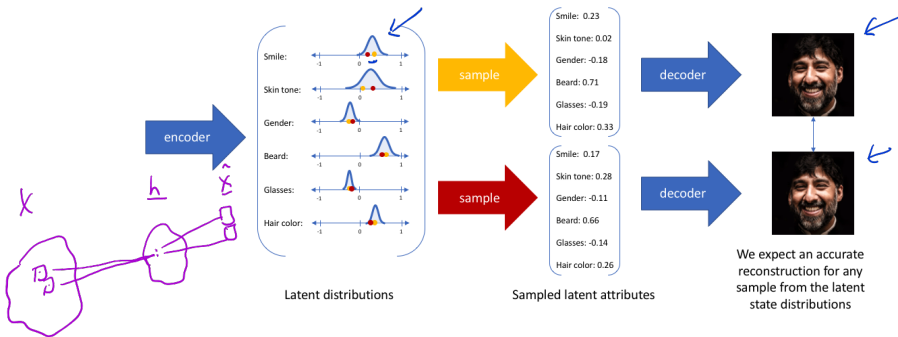
- VAE imposes a specific probabilistic structure on the hidden units.



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



- 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.



(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  $z$ :



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

$p(z|x)$  is circled in green.  $p(x|z)p(z)$  is labeled "likelihood".  $p(x)$  is circled in blue and labeled "prior".  
 The equation is also labeled "marginal" with an arrow pointing to the denominator.  
 Below the equation:  $p(z|x) \propto p(x|z)p(z)$   
 At the bottom right:  $\int p(x|z)p(z) dz$

**Problem** Computing the denominator  $p(\mathbf{x})$  is hard:

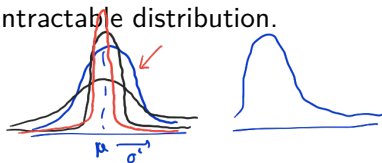
$$p(\mathbf{x}) = \int_{\mathbf{z} \in \mathbb{R}^d} \boxed{p(\mathbf{x}|\mathbf{z})} p(\mathbf{z}) d\mathbf{z}$$

Exact estimate for  $p(\mathbf{z}|\mathbf{x})$  is not possible.

↓  
approximate it

**Solution** Variational inference

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



$\phi \xrightarrow{(1)} q$   
↓  
 $\text{params?}$

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

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

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

$$\sum q(z|x) \log \frac{q(z|x)}{p(z|x)}$$

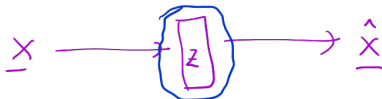
$$\min KL(q(z|x) \parallel p(z|x))$$

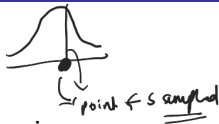
$$\frac{p(x|z) p(z)}{p(x)}$$

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

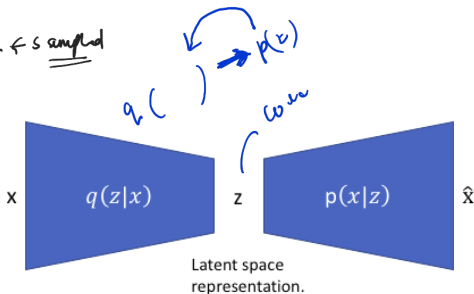
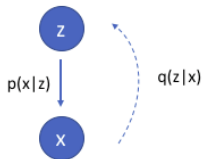
$$E_{q(z|x)} \log p(x|z) - KL(q(z|x) \parallel p(z))$$

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





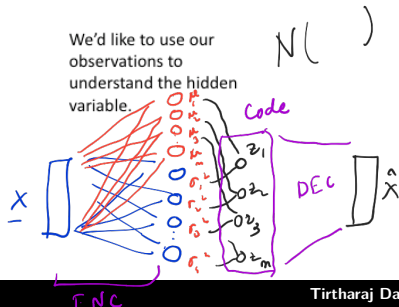
- This is what is happening:



Neural network  
mapping  $x$  to  $z$ .

Neural network  
mapping  $z$  to  $x$ .

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



$$p(z) \rightarrow q(z|x) \text{ — Parameters = } q: \mathcal{N}(\mu, \sigma^2)$$