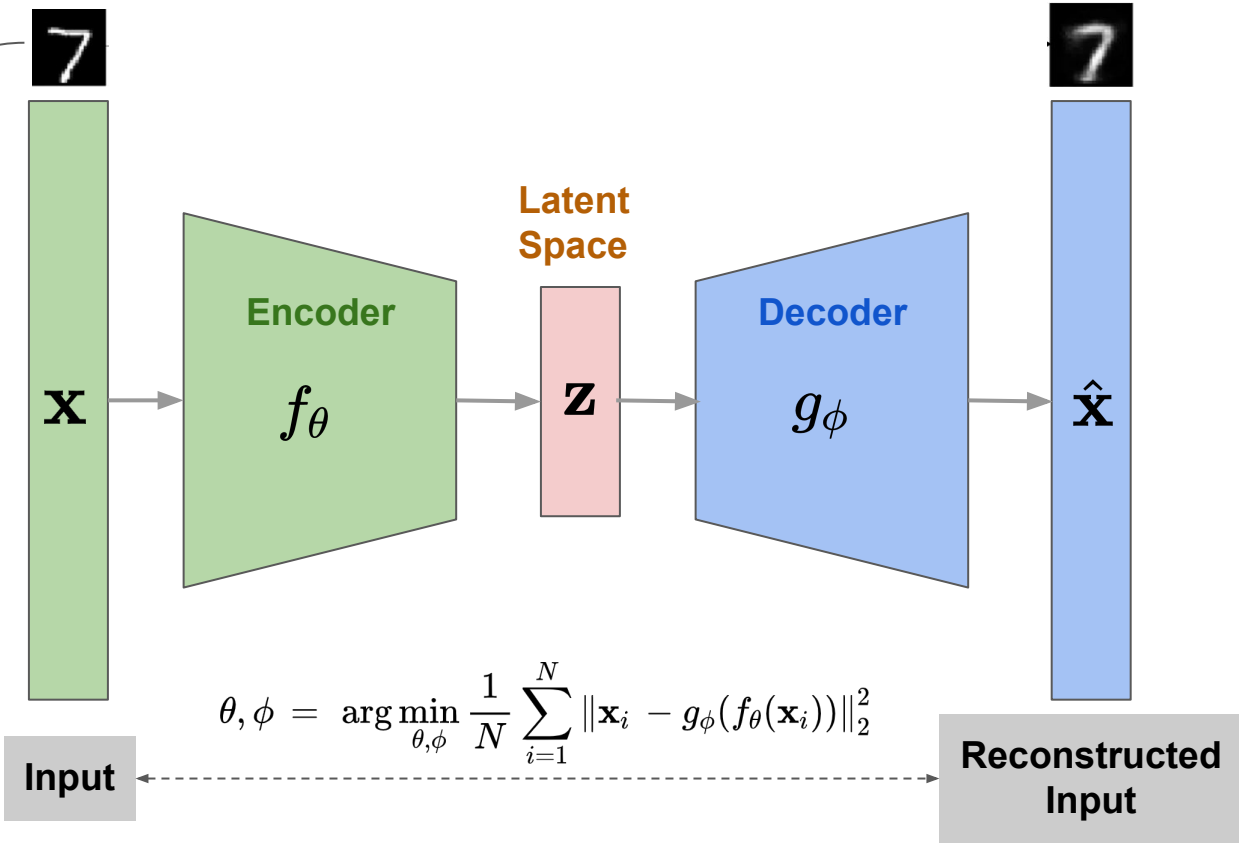


Variational Autoencoder (VAE)

University of Victoria - PHYS-555

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

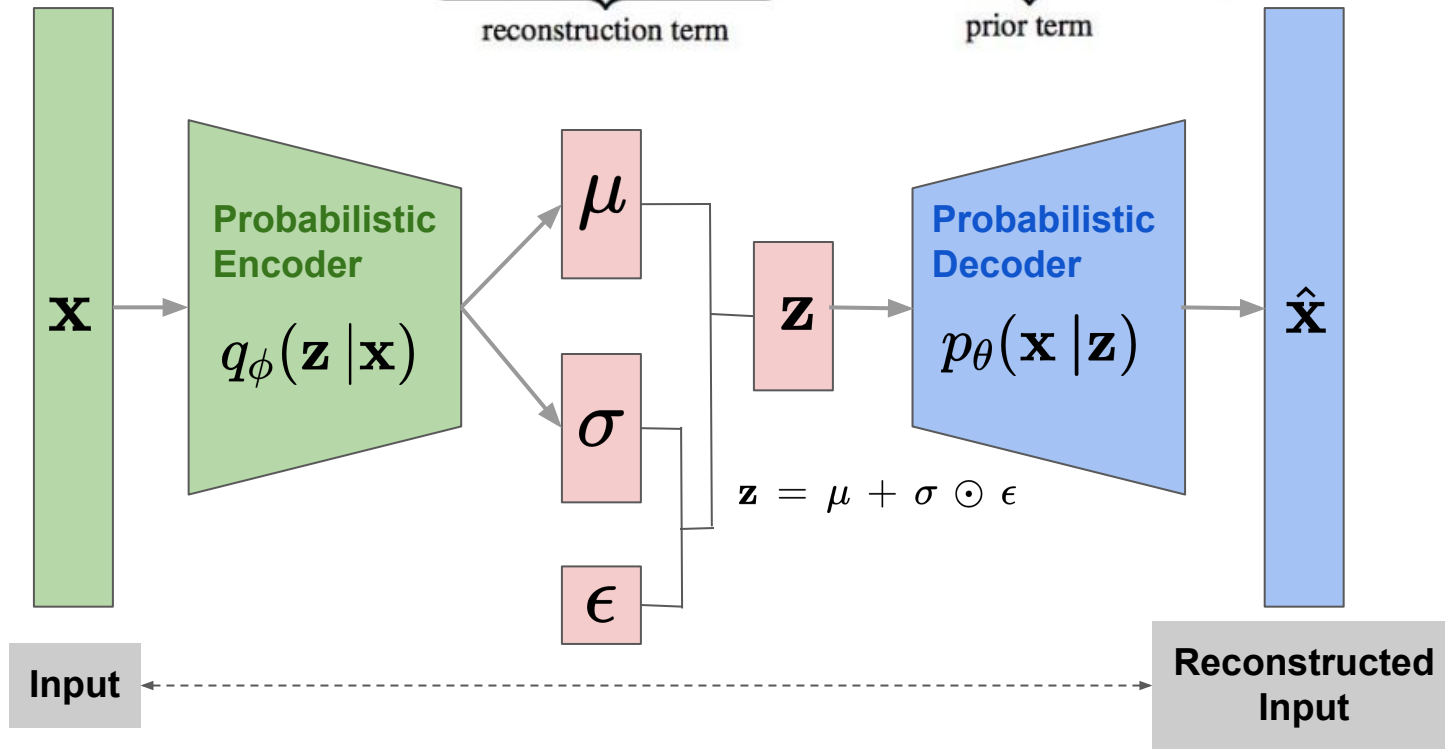


0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9

7

$$\mathcal{L}(x; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z})}_{\text{reconstruction term}} - \underbrace{D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{prior term}}$$

7



Why do we need a VAE?

VAE are generative models: sample a latent vector, decode to get a new sample.

Q: Why can't we use a (deterministic) autoencoder?

A: If we choose an arbitrary latent vector, we are not guaranteed to be close to any points in the training set and the reconstruction is garbage

Q: How can we avoid it?

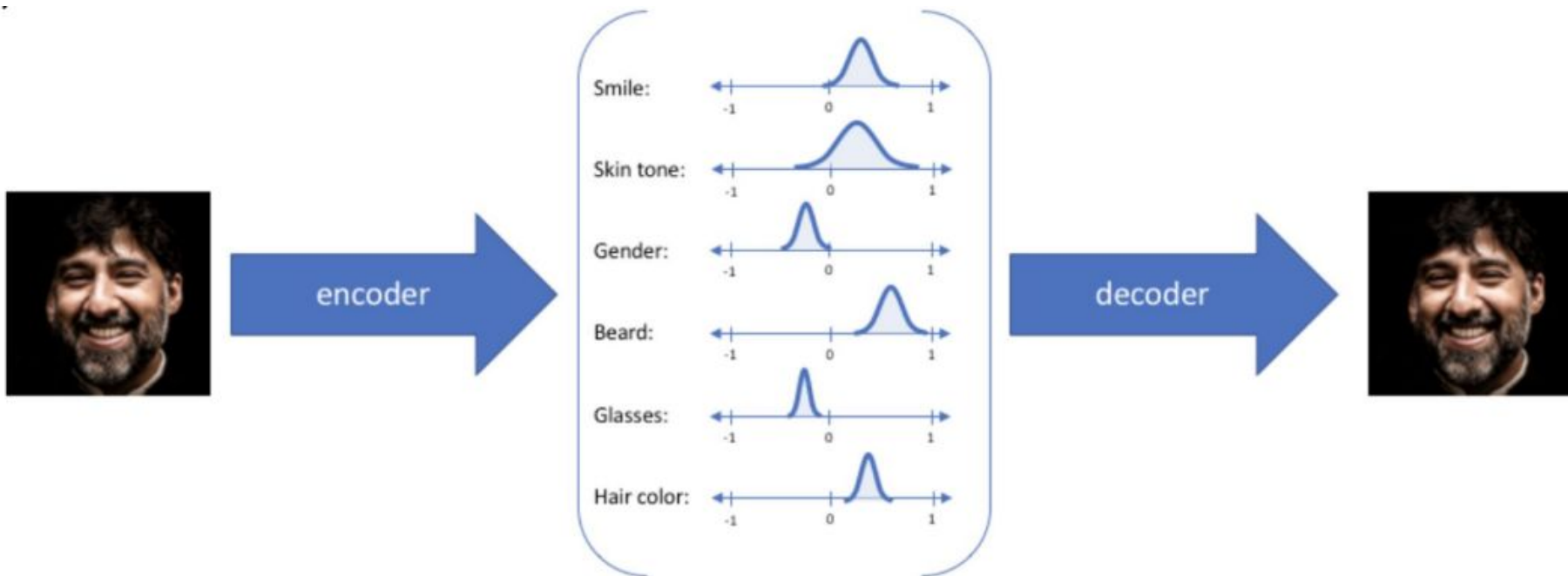
A: Have a more compact latent space.

Q: How can we do this:

A: 1) encode in balls rather than points 2) bring balls closer together

VAE: General Idea

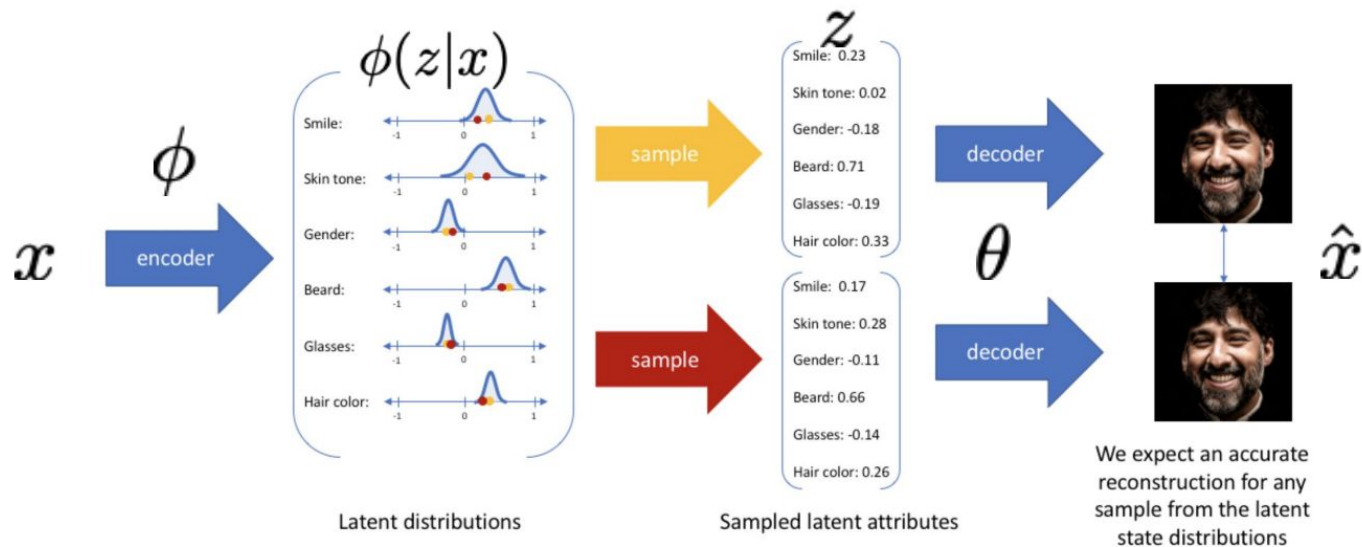
Map the complicated data distribution to a simpler distribution (encoder) we can sample from ([Kingma & Welling 2014](#)) to generate images (decoder)



Encode into Distributions

Q: Why encode into distributions rather than deterministic values?

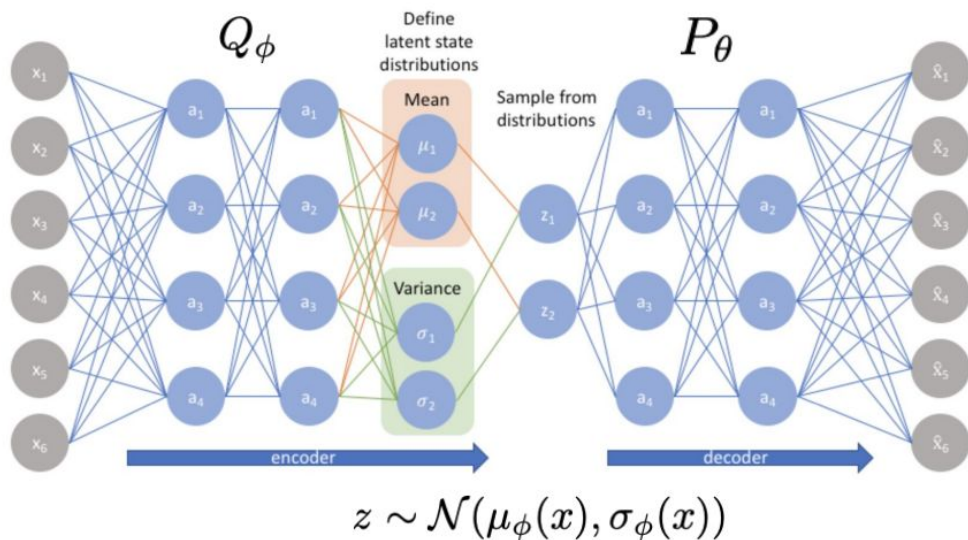
A: This creates balls in latent space and ensures that close points in latent space lead to the same reconstruction and we expect to give more “meaning” to the latent space.



Impose Structure

Q : How can I bring the balls together to compactify latent space ?

A : Make sure that $q(z|x)$ for different x 's are close together !



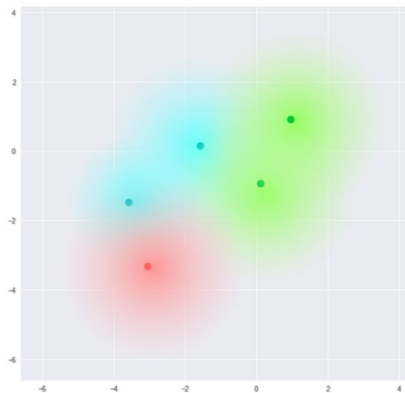
Impose Structure

Q : How do we keep the balls close together?

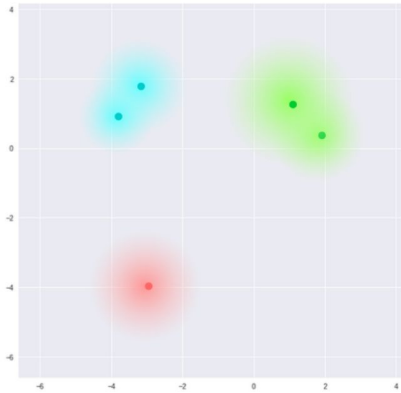
A : By adding springs the balls which pull them towards the center

Q : How?

A : KL divergence with $N(0,1)$ prior. $\mathcal{L}_{KL} = \mathbb{E}_{x \sim dataset} [D_{KL}\{Q_{\phi}(z|x)||p(z)\}], p \sim \mathcal{N}(0, 1)$

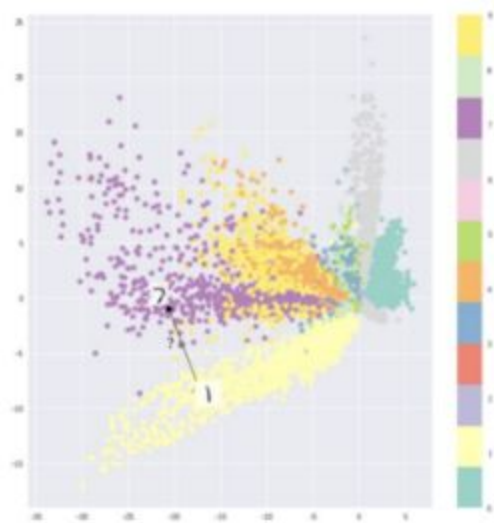


What we require

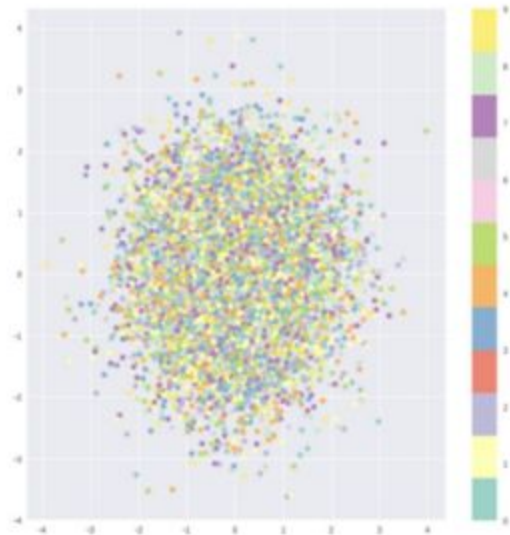


What we may inadvertently end up with

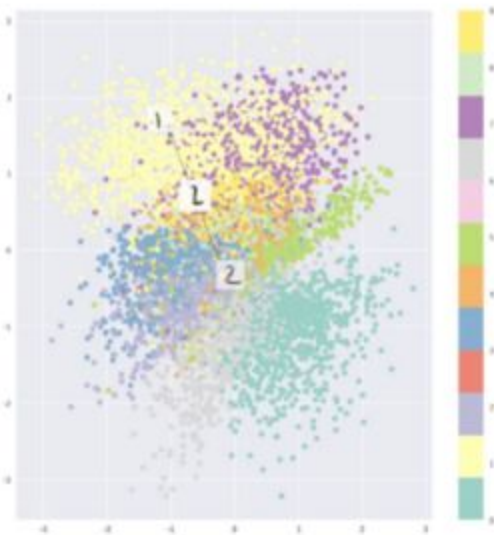
Only reconstruction loss



Only KL divergence



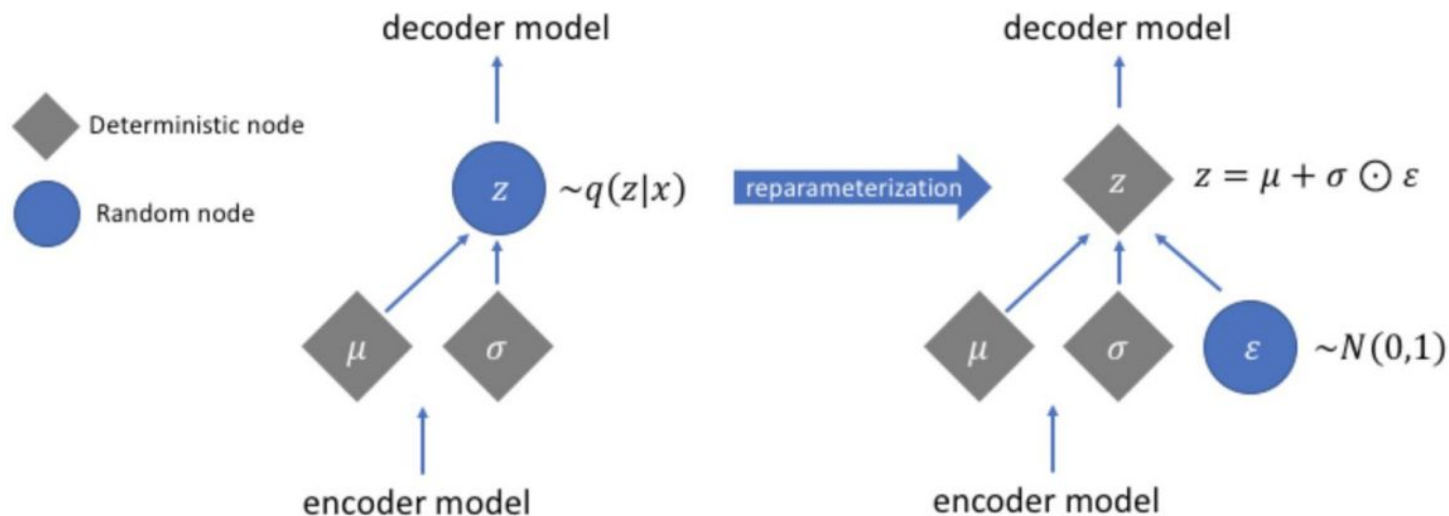
Combination



The Reparameterization Trick

Q: How can we backpropagate when one of the nodes is non-deterministic?

A: Put the random process outside the network



Variational inference

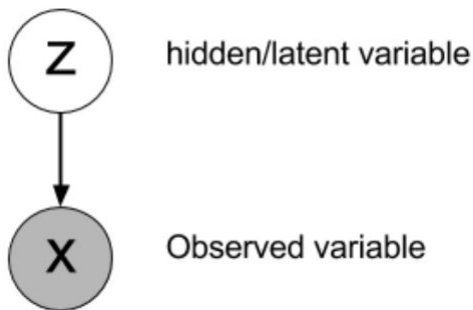
- Latent variable model:

$$p(X, Z) = p(Z)p(X|Z)$$

Diagram illustrating the decomposition of the joint probability $p(X, Z)$ into a prior $p(Z)$ and a likelihood $p(X|Z)$.

- $p(Z)$ is labeled "prior" (blue box).
- $p(X|Z)$ is labeled "likelihood of the observation X given Z" (red box).

- \mathbf{X} – observed data, \mathbf{Z} – latent variable



Variational inference

- Variational inference problem:

The diagram shows the equation $p(Z|X) = \frac{p(X, Z)}{p(X)}$. An orange bracket above the left side of the equation is connected to a blue box containing the text "posterior distribution". A green oval around the denominator $p(X)$ is connected to a green box containing the text "evidence".

$$p(Z|X) = \frac{p(X, Z)}{p(X)}$$

- Computing marginal distribution of the data:

$$p(X) = \int p(X, Z) dz$$

- Intractable in many cases

Approximate posterior

- Goal of variational inference: find approximate posterior distribution similar to the true posterior distribution
- KL divergence measures how two distributions diverge

- Discrete:

$$\text{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

- Continuous:

$$\text{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

- Minimize the KL divergence between the approximate posterior distribution $Q(Z)$ and the true posterior distribution $P(Z)$

Approximate posterior

- From a family of distributions τ find a distribution $q(Z)$ that minimizes KL divergence with the true posterior distribution $p(Z|X)$:

$$q^*(Z) = \operatorname{argmin}_{q(Z) \in \tau} \text{KL}(q(Z) || p(Z|X)) \quad [1]$$

- Can't optimize directly because of $p(Z|X)$

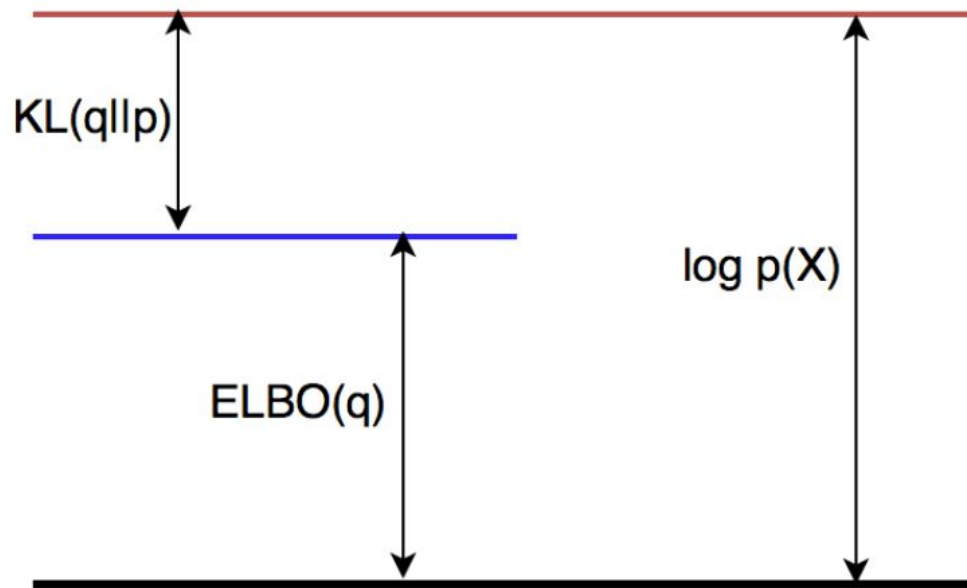
- Rewrite Eq. [1] as:

$$\begin{aligned}
 \text{KL}(q(Z)||p(Z|X)) &= \mathbb{E}_{q(Z)} \left[\log \frac{q(Z)}{p(Z|X)} \right] \\
 &= \mathbb{E}_{q(Z)} [\log q(Z)] - \mathbb{E}_{q(Z)} [\log p(Z|X)] \\
 &= \mathbb{E}_{q(Z)} [\log q(Z)] - \mathbb{E}_{q(Z)} \left[\log \frac{p(Z, X)}{p(X)} \right] \\
 &= \mathbb{E}_{q(Z)} [\log q(Z)] - \mathbb{E}_{q(Z)} [\log p(Z, X)] + \mathbb{E}_{q(Z)} [\log p(X)]
 \end{aligned} \tag{2}$$

- Rearrange:

$$\underbrace{\{\mathbb{E}_{q(Z)} [\log p(Z, X)] - \mathbb{E}_{q(Z)} [\log q(Z)]\}}_{\text{Evidence Lower Bound (ELBO)}} + \underbrace{\text{KL}(q(Z)||p(Z|X))}_{\text{log evidence}} = \log p(X) \tag{3}$$

Evidence Lower Bound (ELBO)



Evidence Lower Bound (ELBO)

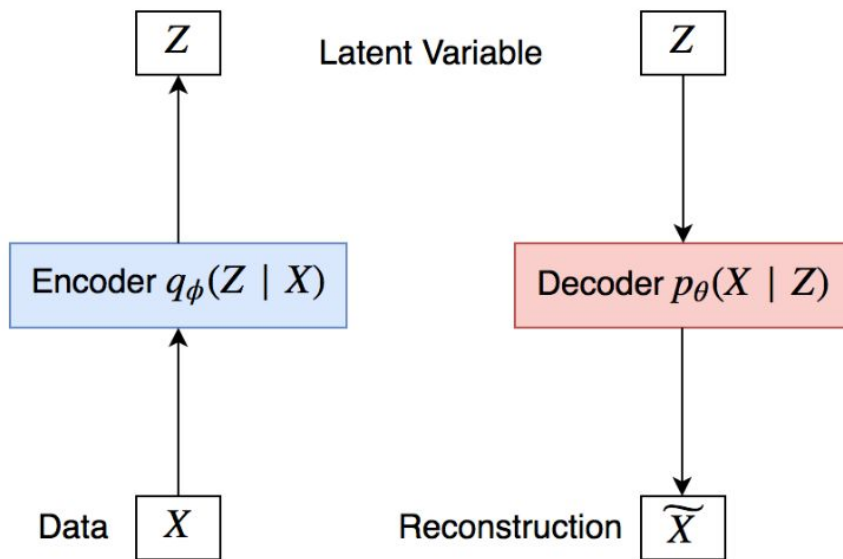
- ELBO can be re-written as:

$$\begin{aligned}\text{ELBO}(q) &= \mathbb{E}_{q(Z)} [\log p(Z, X)] - \mathbb{E}_{q(Z)} [\log q(Z)] \\ &= \mathbb{E}_{q(Z)} [\log p(X|Z)p(Z)] - \mathbb{E}_{q(Z)} [\log q(Z)] \\ &= \mathbb{E}_{q(Z)} [\log p(X|Z)] + \mathbb{E}_{q(Z)} [\log p(Z)] - \mathbb{E}_{q(Z)} [\log q(Z)] \\ &= \mathbb{E}_{q(Z)} [\log p(X|Z)] + \mathbb{E}_{q(Z)} \left[\log \frac{p(Z)}{q(Z)} \right] \\ &= \underbrace{\mathbb{E}_{q(Z)} [\log p(X|Z)]}_{\text{Expected log likelihood of the data}} - \underbrace{\text{KL}(q(Z)||p(Z))}_{\text{KL divergence between the approximate posterior and the prior}}\end{aligned}$$

[4]

- Maximize expected log likelihood of data and minimize KL divergence

Variational Autoencoder



Variational Autoencoder

- Replace approximate posterior $q(Z)$ with $q_{\phi}(Z|X)$
- Given a dataset: $\mathcal{D} = \{x^{(n)}\}_{n=1}^N$
- The likelihood of a data point (log evidence) and ELBO are related as follows:

$$\begin{aligned}\log p_{\theta}(x^{(n)}) &\geq \mathbb{E}_{z \sim q_{\phi}(z|x^{(n)})} \left[\log \left\{ \frac{p_{\theta}(x^{(n)}, z)}{q_{\phi}(z|x^{(n)})} \right\} \right] \\ &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(n)})} [\log p_{\theta}(x^{(n)}|z)] - \text{KL} (q_{\phi}(z|x^{(n)}) \| p(z))\end{aligned}\quad [5]$$

Variational Autoencoder

- Maximizing ELBO(q) is equivalent to minimizing $-\text{ELBO}(q)$
- Loss function for the VAE:

$$\begin{aligned} J^{(n)} &= -\mathbb{E}_{z \sim q_{\phi}(z|x^{(n)})} [\log p_{\theta}(x^{(n)}|z)] + \text{KL}(q_{\phi}(z|x^{(n)}) \| p(z)) \\ &= J_{\text{rec}}(\theta, \phi, x^{(n)}) + \text{KL}(q_{\phi}(z|x^{(n)}) \| p(z)) \end{aligned}$$

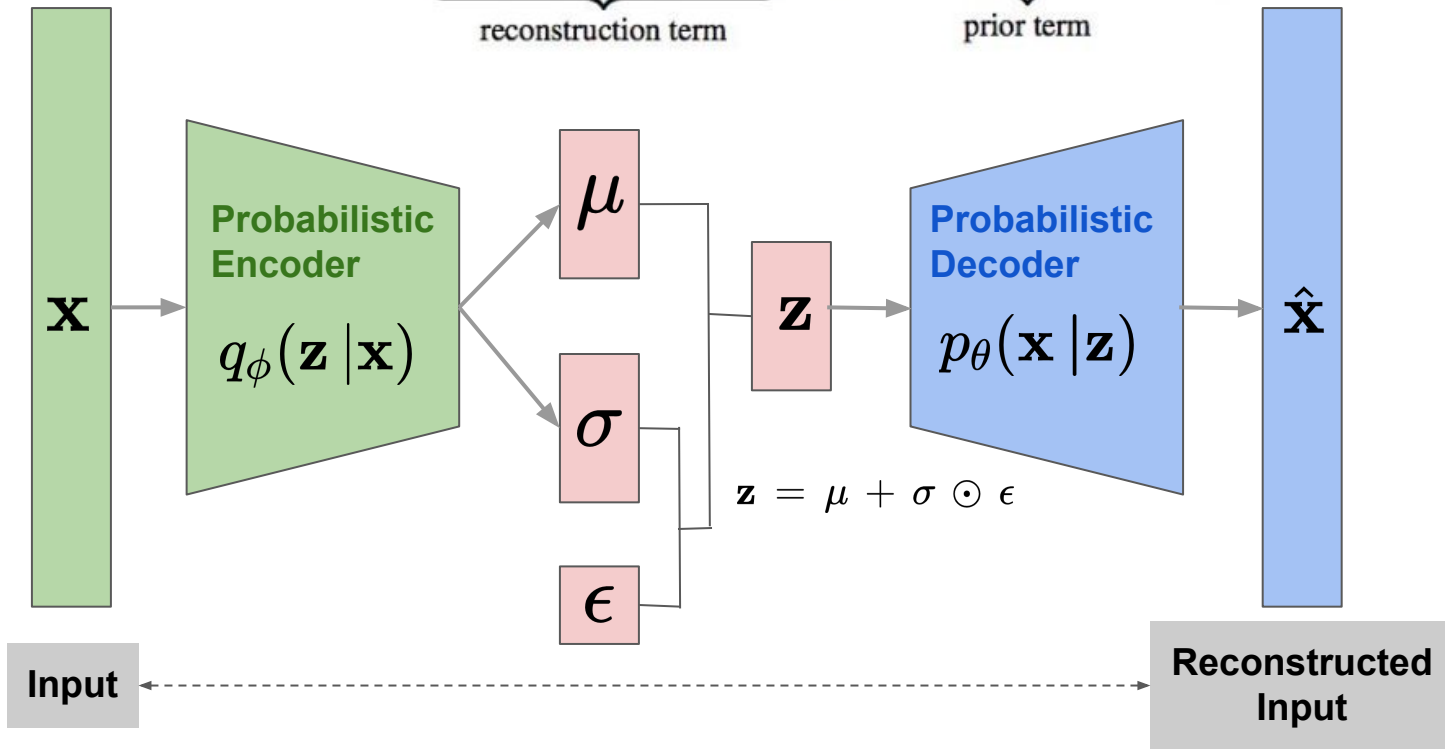
- Reconstruction loss: expected negative log likelihood of the data
- Prior is typically Gaussian $N(0, I)$

0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9

7

$$\mathcal{L}(x; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z})}_{\text{reconstruction term}} - \underbrace{D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{prior term}}$$

7



VAE Extension: β -VAE

- Simple modification to VAE framework

$$\mathcal{L}(x; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z})}_{\text{reconstruction term}} - \underbrace{\beta D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{prior term}}$$

- Encourages disentangled representation in latent code with higher beta at the expense of reconstruction quality

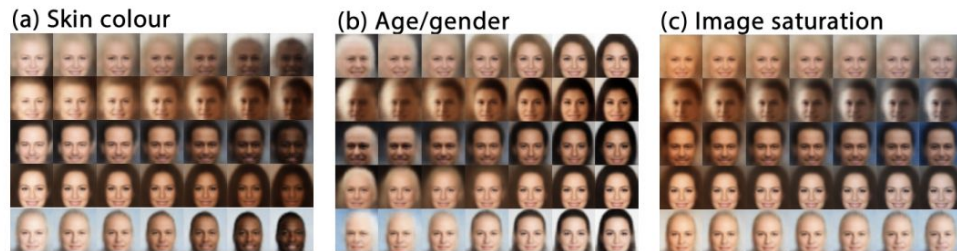


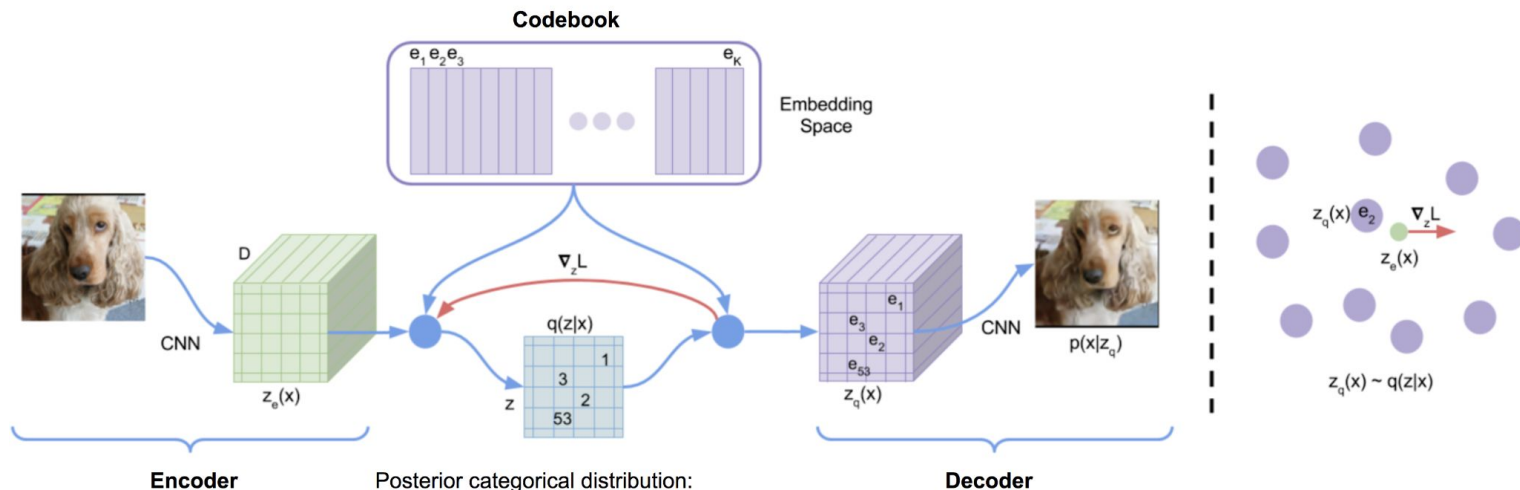
Figure 4: **Latent factors learnt by β -VAE on celebA:** traversal of individual latents demonstrates that β -VAE discovered in an unsupervised manner factors that encode skin colour, transition from an elderly male to younger female, and image saturation.

Vector Quantized VAE

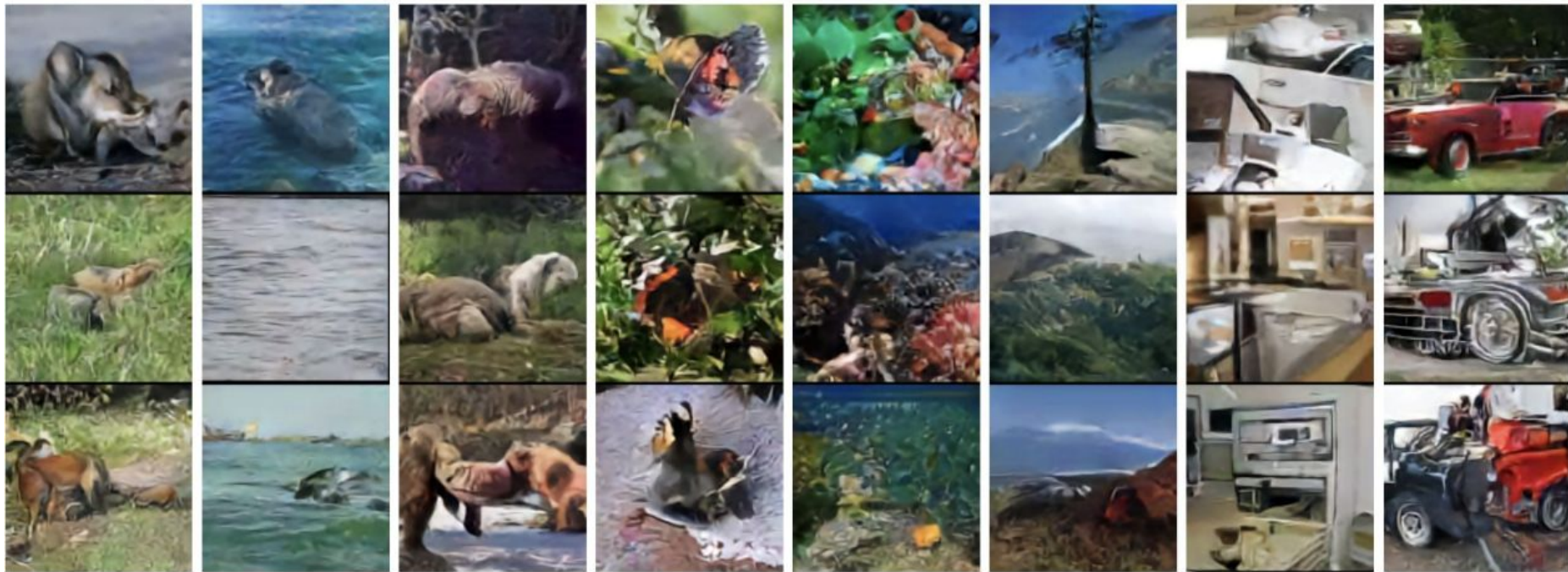
VQ-VAE

- Discrete latent code
- Autoregressive decoder
- Train prior over latent code can be added

(VQ-VAE-2)



$$q(\mathbf{z} = \mathbf{e}_k | \mathbf{x}) = \begin{cases} 1 & \text{if } k = \arg \min_i \|\mathbf{z}_e(\mathbf{x}) - \mathbf{e}_i\|_2 \\ 0 & \text{otherwise.} \end{cases}$$



VQ-VAE high quality samples

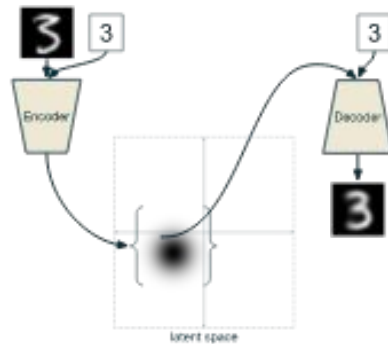
Conditional VAE

If we have labels, we can learn to generate data *given* a label to hopefully improve or guide our generation.

Guide the generator towards more familiar samples.

Implementation: concatenate labels to dense layers in the encoder (or latent code), do the same for the decoder.

Training time



Inference time

