

DUL: Variational Autoencoders

Naeemullah Khan
naeemullah.khan@kaust.edu.sa



جامعة الملك عبدالله
للتكنولوجيا
King Abdullah University of
Science and Technology

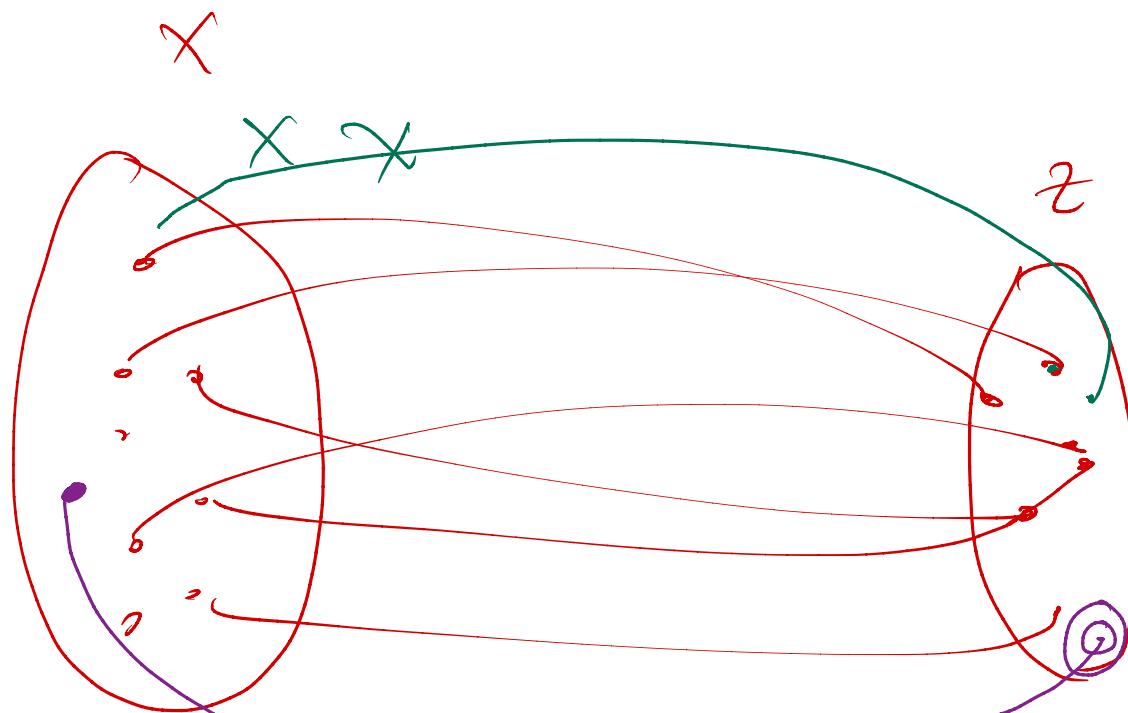
KAUST Academy
King Abdullah University of Science and Technology

May 27, 2025

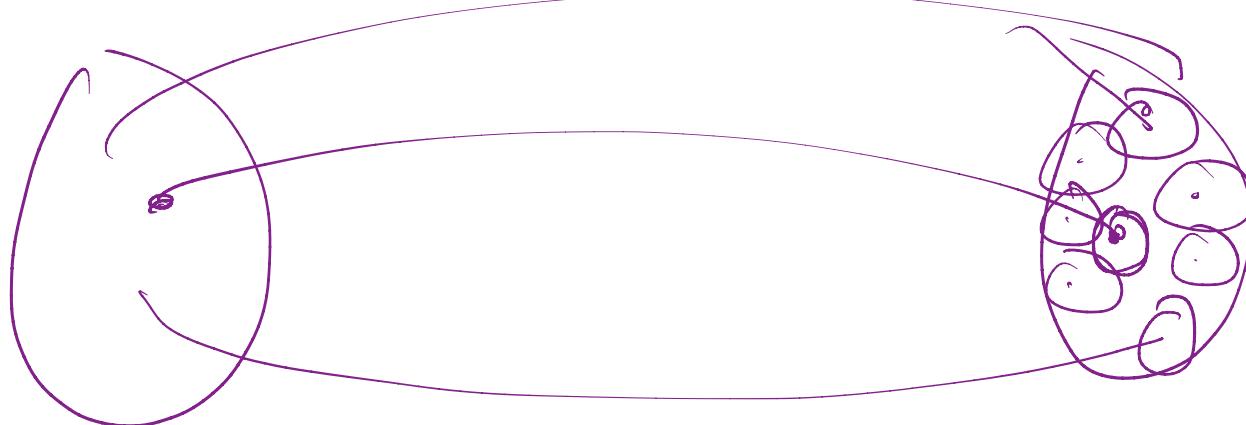
Table of Contents

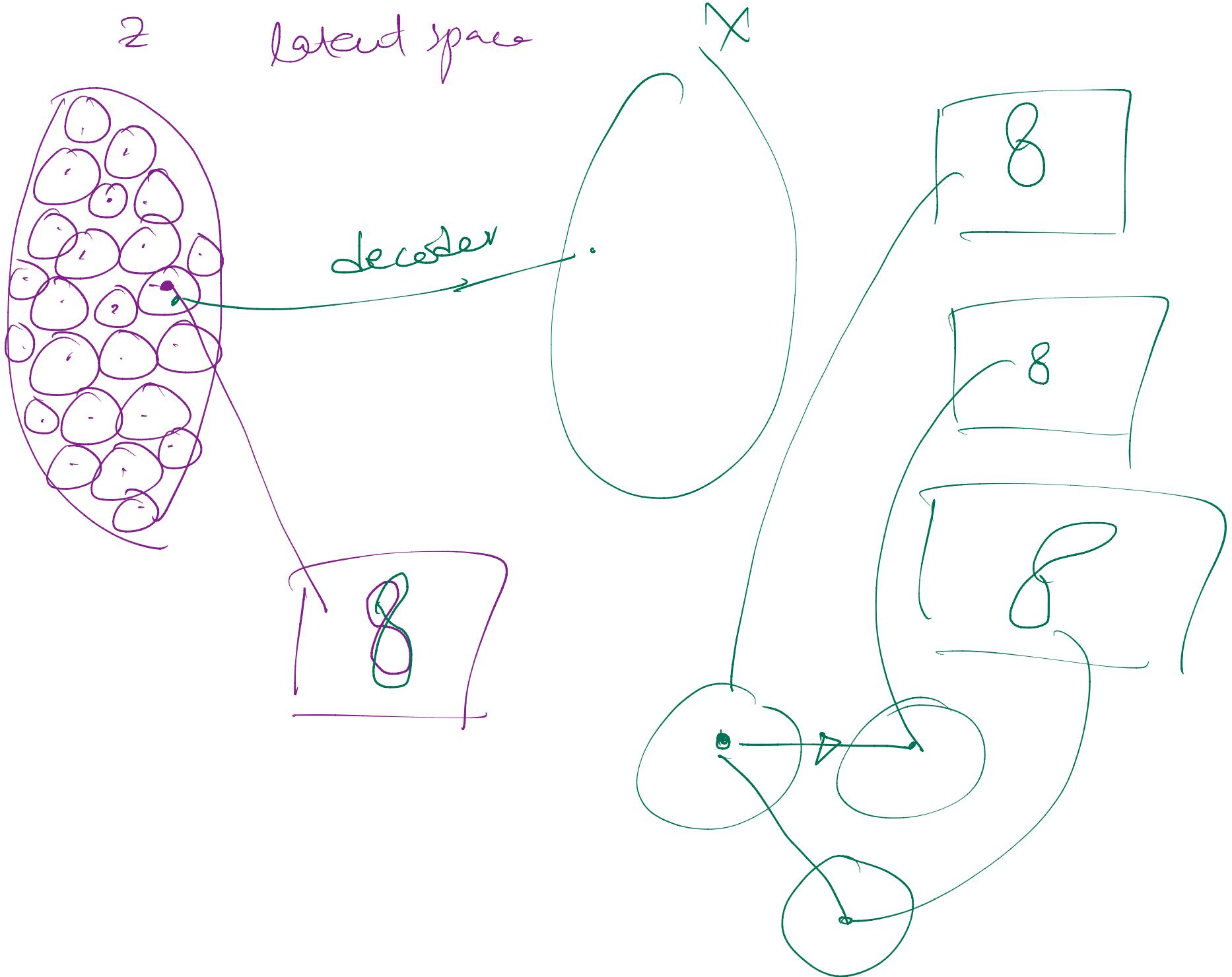
1. Motivation
2. Autoencoders Recap
3. Introduction to Variational Autoencoders
4. Latent Variable Models
5. Variational Inference
6. Evidence Lower Bound (ELBO)
7. Reparameterization Trick
8. Loss Function
9. Results
10. Variants and Extensions
11. Limitations and Challenges

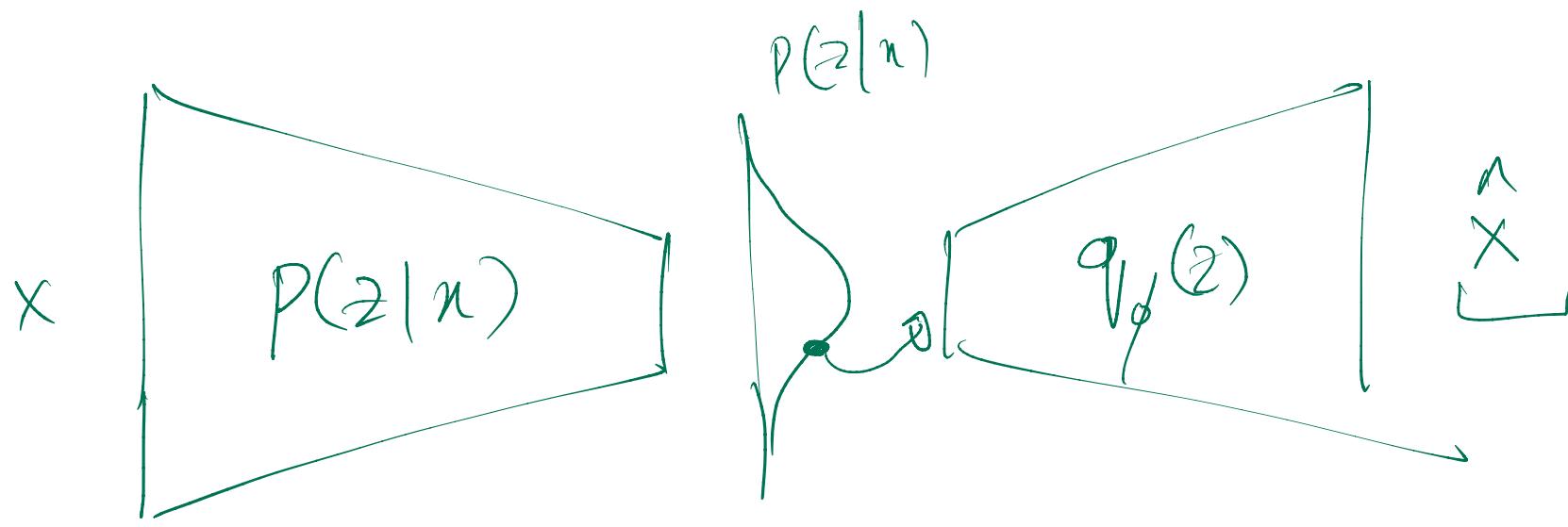
Auto Encoders



VAEs







$$\min \quad \|x - \hat{x}\|$$

$$P(z|n) = \underbrace{p(x|z) \cdot p(z)}_{\text{intractable}} \xrightarrow{\text{decoder}}$$

stuck!!

$$q(z|n) \leftarrow \mathcal{N}(\mu, \sigma^2 I)$$

$$KL(q \parallel p) = \int_z q(z|n) \log \frac{q(z|n)}{p(z|n)} dz$$

$$KL(q||p) = \int q(z|x) \log \frac{q(z|x)}{p(z|x)} dz$$

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

$$= \int q(z|u) \left\{ \log \frac{q(z|\alpha)}{p(\alpha|z)} + \log p(\alpha) \right\} dz$$

$$= \int q(z|x) \log \frac{q(z|x)}{p(x,z)} dz + \int q(z|x) \log p(x) dz$$

$\stackrel{\text{def}}{=} p(x|z) \cdot p(z)$

$$\log P(x) = \int_{-\infty}^x p(z|y) dz$$

$$KL(q \parallel p) = \int q(z|x) \log \frac{q(z|x)}{p(x,z)} dz + \log p(x)$$

$$\log p(x) = KL(q \parallel p) + \int q(z|x) \log \frac{p(x,z)}{q(z|x)} dz$$

$$\log p(x) \geq \int q(z|x) \log \frac{p(x,z)}{q(z|x)} dz$$

$$\max \int q_j(z|n) \log \frac{p(x|z)}{q_j(z|n)} dz$$

$$\max \int q_j(z|n) \log \frac{(p(x|z) p(\theta))}{q_j(z|n)} dz$$

$$\max \int q_j(z|n) \log p(x|z) dz + \int q_j(z|n) \log \frac{p(\theta)}{q_j(z|n)} dz$$

$$\min \int -q_j(z|n) \log p(x|z) + \int q_j(z|n) \log \left(\frac{q_j(z|n)}{R^z} \right) dz$$

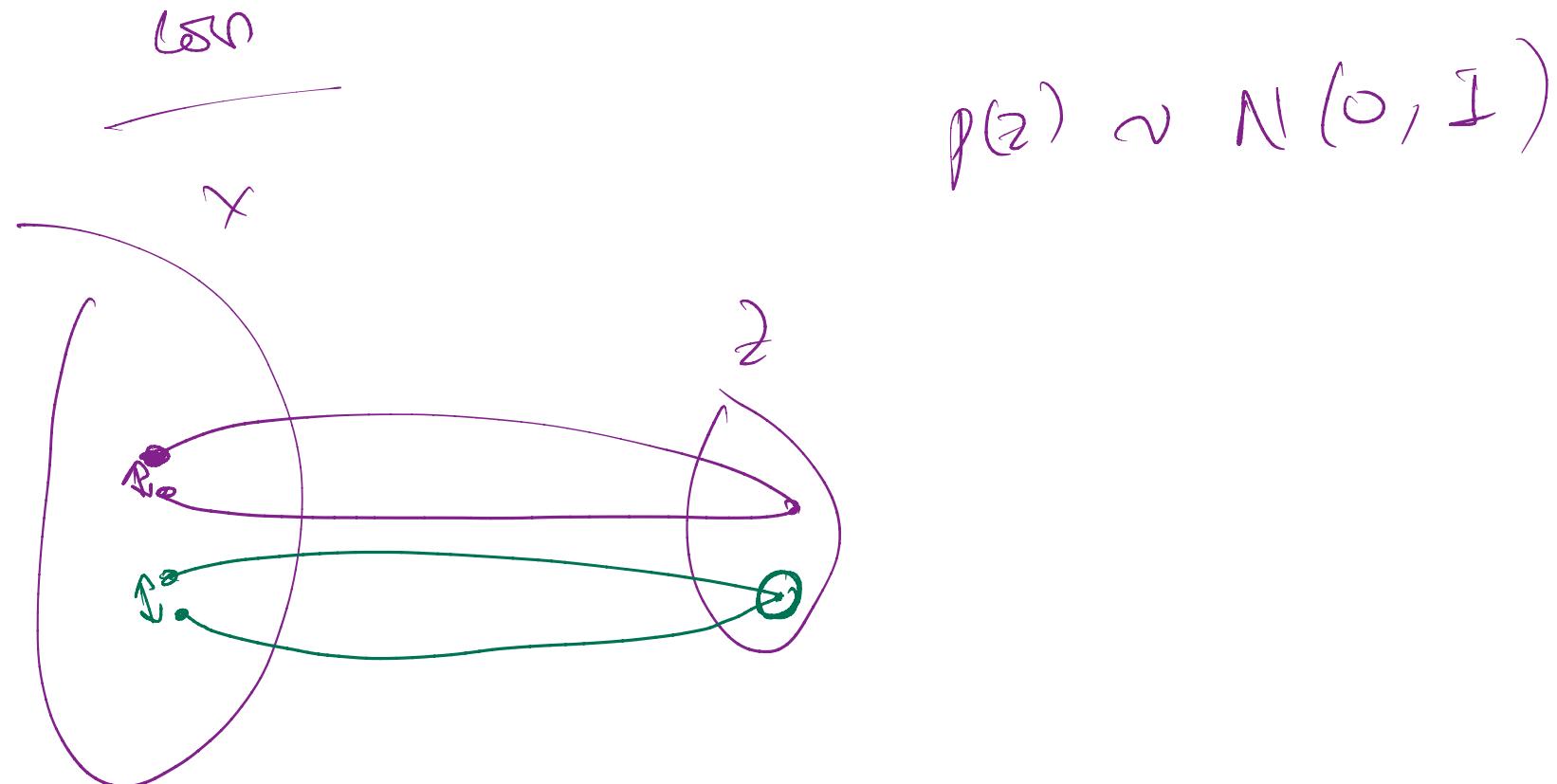
→

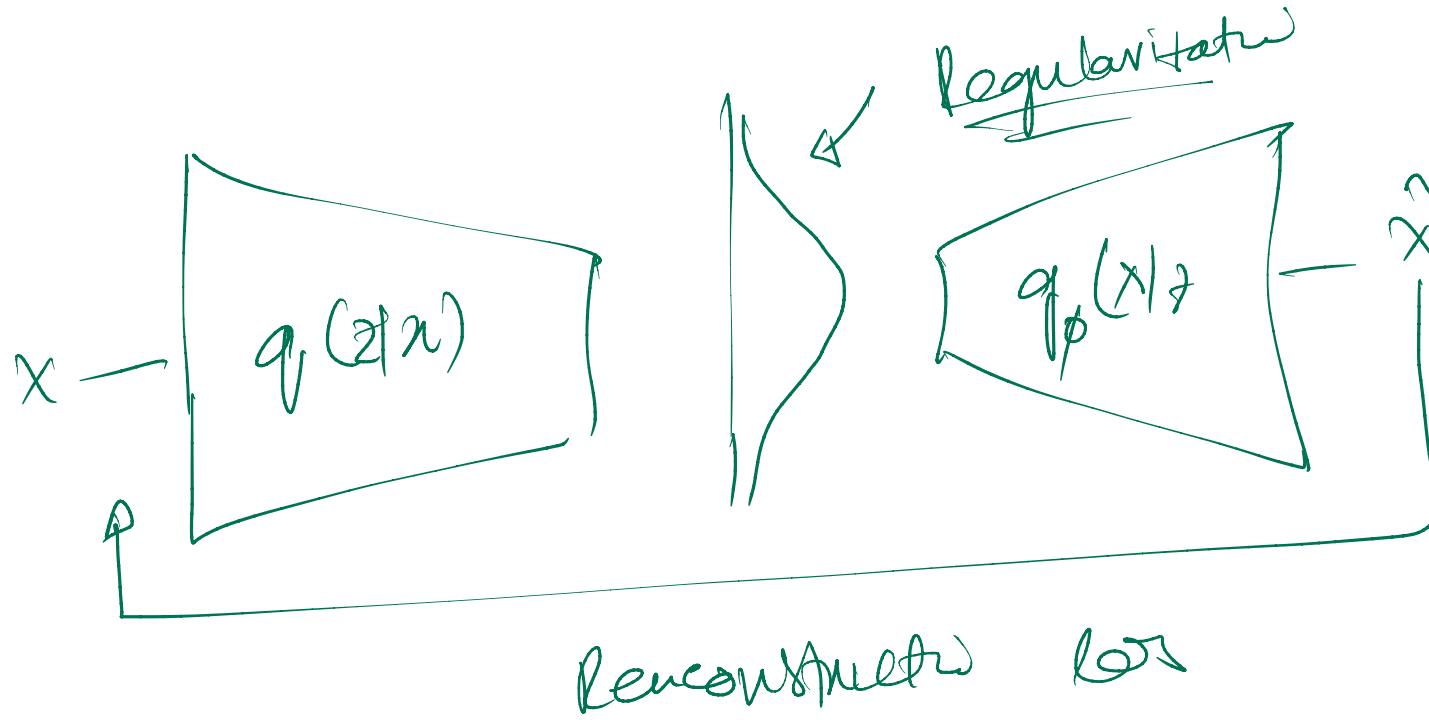
$$-(-\|\hat{x} - x\|_2^2) + KL(q_j(z|n) || p(\theta))$$

→

$$\min \|\hat{x} - x\|_2^2 + KL(q_j(z|n) || p(\theta))$$

$$\min \underbrace{\|x - \hat{x}\|_2^2}_{\text{Reconstruction}} + \underbrace{\text{KL}(q(z|x) \parallel p(z))}_{\text{Regularization}}$$



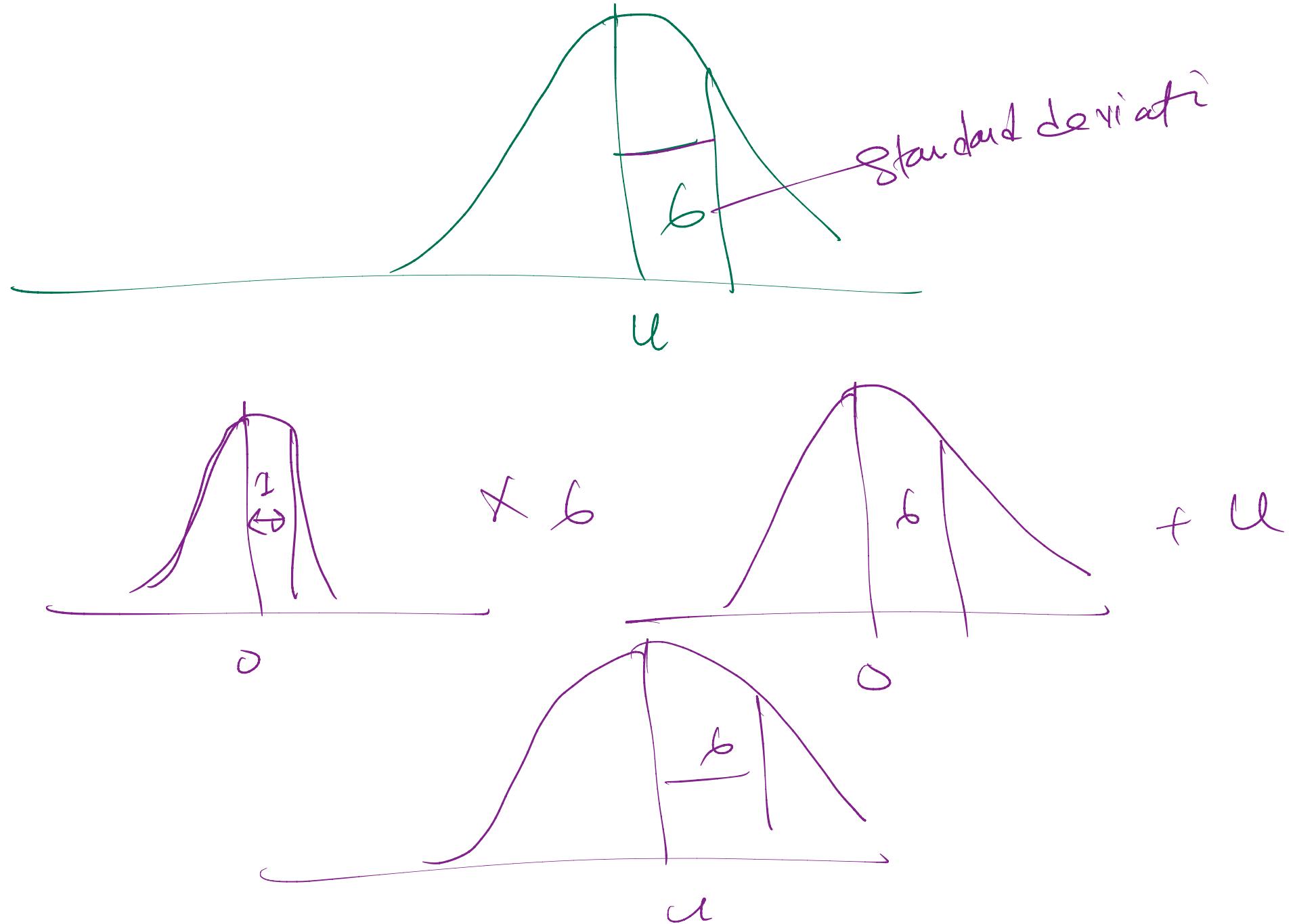


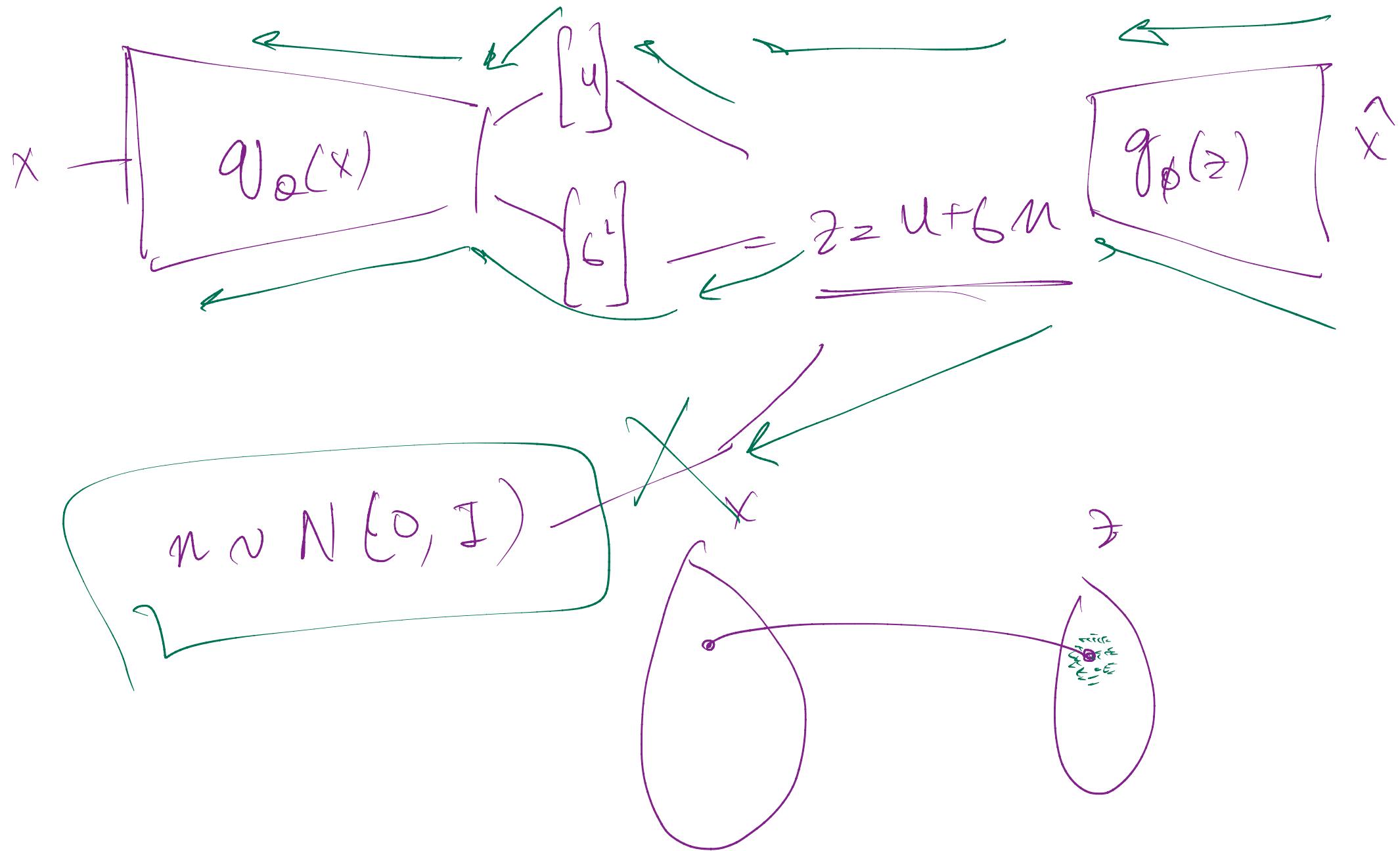
Reparametrisation trick

$$u \sim N(\mu, \delta^2 I)$$

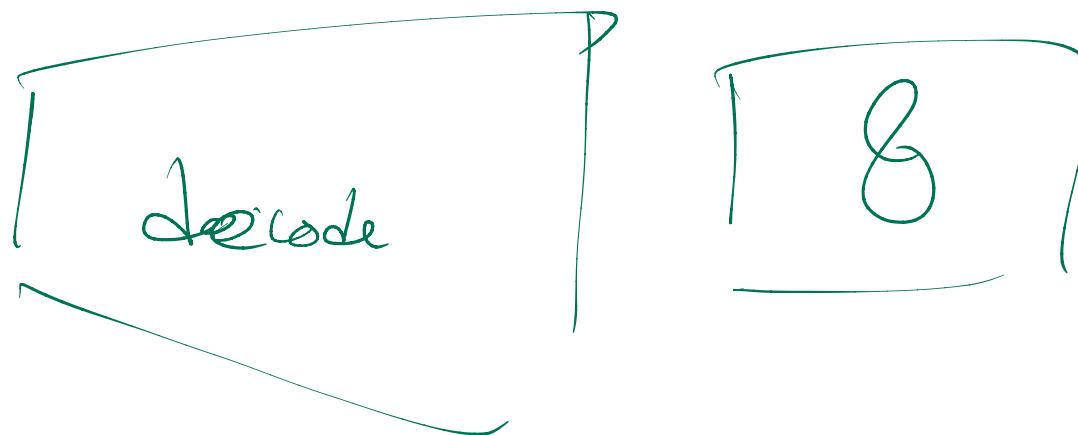
$$u = y + b n$$

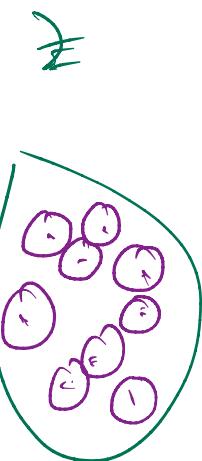
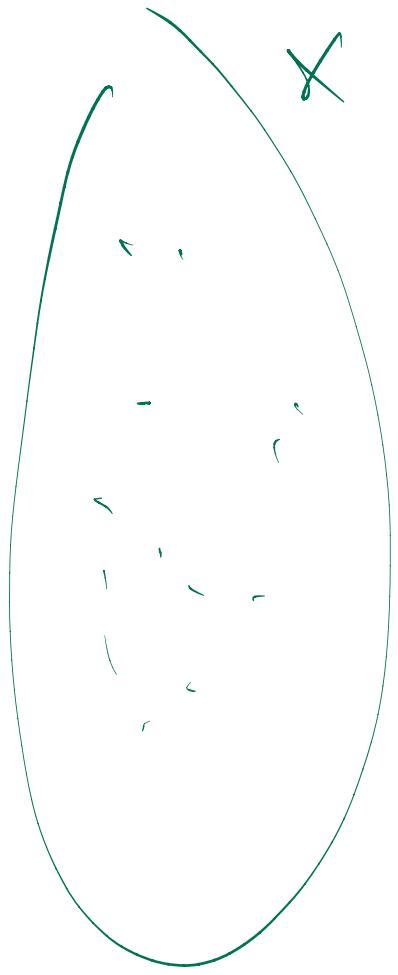
$$n \sim N(0, I)$$





$$z \sim N(0, I)$$





$$KL(q(z|n) \parallel p(z)) = \int q(z|n) \log \frac{q(z|n)}{p(z)} dz$$

$$\mathbb{E}_{z \sim q(z|n)} \left[\log \frac{q(z|n)}{p(z)} \right]$$

$$\mathbb{E}_{z \sim q} \left[\log \left(\frac{\exp^{-\frac{(z - \mu)^2}{2\sigma^2}}}{\sqrt{2\pi\sigma^2}} \right) \right]$$

$$q(z|x) \sim N(\mu, \sigma^2 I)$$

$$p(z) \sim N(0, I)$$

$$\mathbb{E}_x [f(x)] = \int p(x) f(x) dx$$

$$\frac{E}{2\pi g} \left[\log \left(\frac{1}{6} \cdot \exp \left(-\frac{(z-u)^2}{2\delta^2} + \frac{u^2}{2} \right) \right) \right]$$

$$\frac{E}{2\pi g} \left[\log \frac{1}{6} + -\frac{(z-u)^2}{2\delta^2} + \frac{u^2}{2} \right]$$

$$\log \frac{1}{6} + \frac{E}{2\pi g} \left[-\frac{z^2 + 2zu - u^2}{2\delta^2} + \frac{z^2 + u^2}{2\delta^2} \right]$$

$$\log \frac{1}{6} + \left[-\frac{(6^2 \cancel{+ u^2}) + 2\cancel{z^2} - \cancel{u^2} + (6^2 + u^2) \delta^2}{2\delta^2} \right]$$

$$\log\left(\frac{1}{6}\right) + \left[\frac{-6^2 + 6^4 + 6^{12}}{2 \cdot 6^2} \right]$$

$$\frac{2}{2} \log\left(\frac{1}{6}\right) + \left[\frac{-1 + 6^2 + 6^2}{2} \right]$$

$$\frac{1}{2} \left[\log \frac{1}{6^2} - 1 + 6^2 + 6^2 \right]$$

$$\frac{1}{2} \left[-\log 6^2 - 1 + 6^2 + 6^2 \right]$$

$$\|x - \hat{x}\|_2$$

$$T_0 \in \mathbb{C}^2$$

$$z \sim (0, 1)$$

$$z \sim \frac{1}{2\pi c} \exp \frac{(2-z)^2}{2c}$$

$$a \in \log S^{\perp}$$

$$E(z) = u$$

$$z \sim N(u, \sigma^2)$$

$$E(z^2)$$

$$\text{Var}(z) = E(z^2) - (E(z))^2$$

$$E(z^2) = \underbrace{\text{Var}(z)}_{\sigma^2} + (E(z))^2$$

$$\sigma^2 + u^2$$

Reconstruction

$$p(\hat{x}|x) \propto \exp - \|\hat{x} - x\|_2^2$$

$$\log p(\hat{x}|x) = -\|\hat{x} - x\|_2^2$$

$$N(0, \sigma^2 I)$$

$$\hat{x} = x + n$$

Why VAEs?

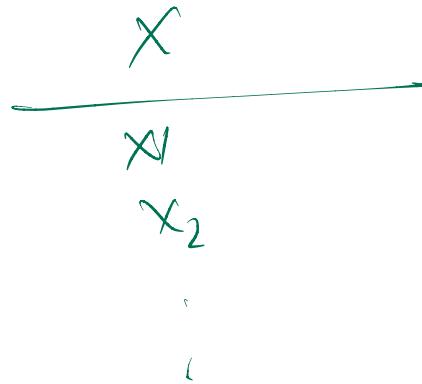
- ▶ Traditional autoencoders are deterministic and lack generative capabilities.
- ▶ Need for models that can generate new, diverse data samples.

Applications

- ▶ Image generation (e.g., generating new faces or handwritten digits).
- ▶ Data compression and denoising.
- ▶ Anomaly detection in various domains.

Auto Regressive Models

$$x \sim P(x)$$

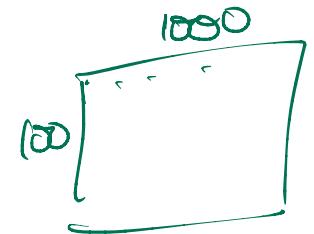


$$\underbrace{P_\theta(n)}$$

~~x_{N+1}~~

$$\arg \max_{\theta} P_\theta(x_1, \dots, x_N) \stackrel{\text{iid}}{=} \arg \max_Q \prod_{i=1}^N P_\theta(x_i)$$
$$\arg \max_{\theta} \sum_{i=1}^N \log P_\theta(x_i)$$

$$\text{arg max}_{\theta} \sum_{i=1}^n \log p_{\theta}(x_i)$$



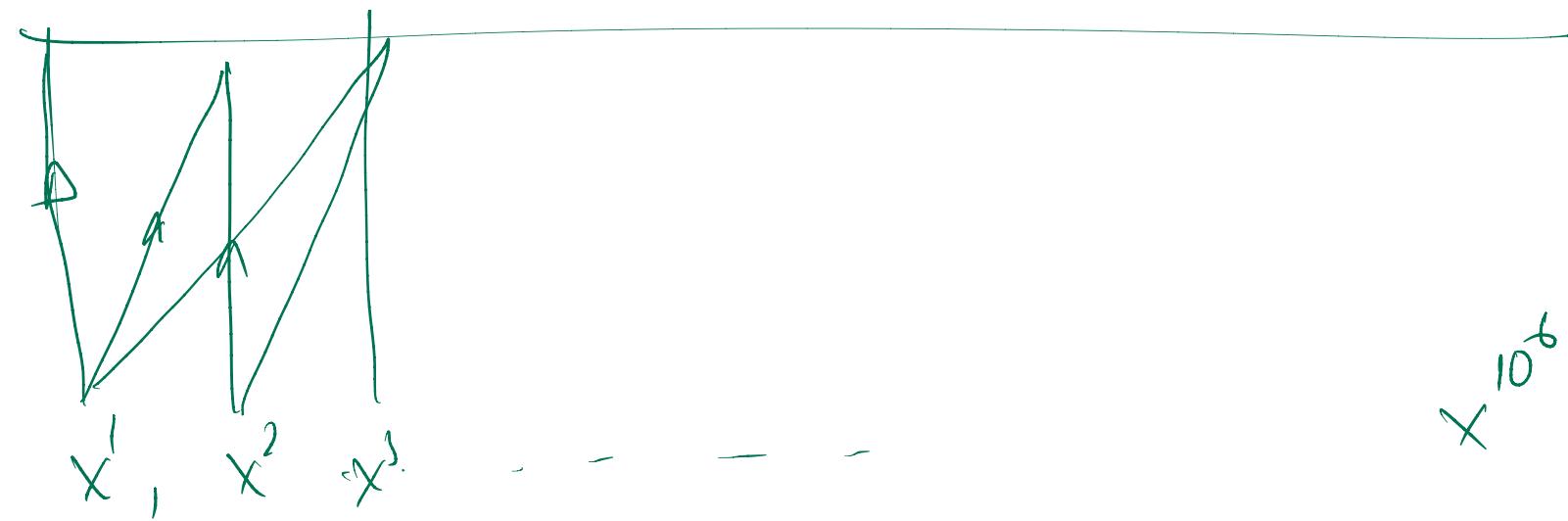
$$\text{arg max}_{\theta} \sum_{i=1}^N \log \prod_{j=1}^M p_{\theta}(x_i^j)$$

x_i^j

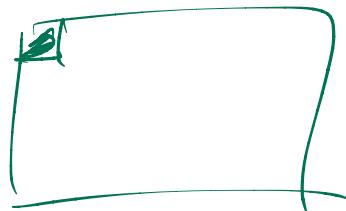
$$p(x^1, x^2, x^3, \dots, x^{10^6}) = p(x^1) p(x^2|x^1) p(x^3|x^1, x^2) p(x^7|x^1, x^2, x^3) \dots$$



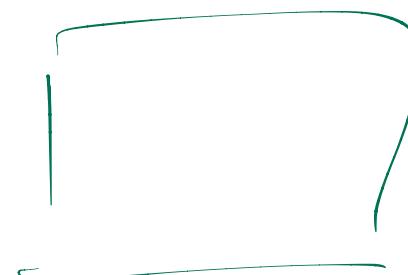
NN



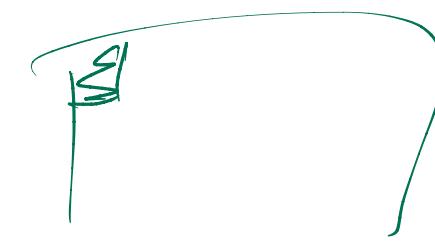
$$P(x_2 \mid \alpha) = 0.5 (\alpha_1')$$



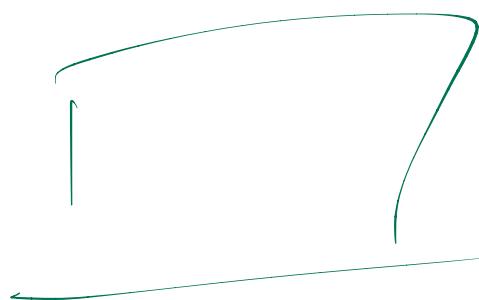
$$P(x^2 \mid x^1) = \alpha_2^2 + \alpha_2^1 x^1$$



$$P(x^3 \mid x^1, x^2) = \alpha_3^3 + \alpha_3^1 x^1 + \alpha_3^2 x^2$$



1 + 2 + 3 ... N



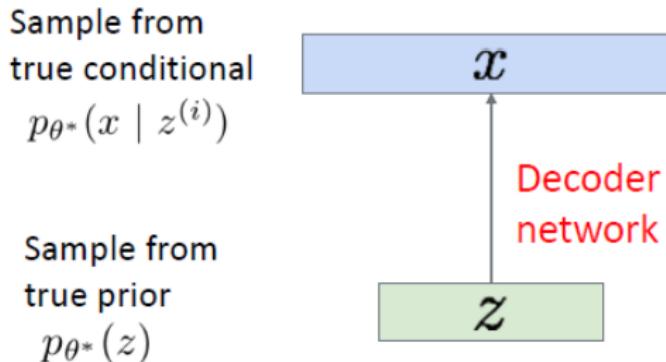
Autoencoders:

- ▶ **Architecture:**
 - **Encoder:** Compresses input x into a latent representation z .
 - **Decoder:** Reconstructs the input x from the latent representation z .
- ▶ **Objective:** Minimize reconstruction error $\|x - \hat{x}\|^2$.
- ▶ **Limitation:** Cannot generate new data; lacks a probabilistic foundation.

- ▶ Autoencoders do not work as generative models because the latent space Z is too discrete.
- ▶ **Solution:** Let us rectify that.
- ▶ **Variational Autoencoders:** Probabilistic spin on autoencoders will let us sample from the model to generate data.

Variational Autoencoders: Introduction

- ▶ We want a generative model, which given a prior z outputs a new sample from data.
- ▶ To make the latent space Z continuous, let's choose it to be Gaussian
- ▶ So now, we want to estimate the true parameters θ^* of this generative model.



- ▶ **Probabilistic Encoder:** Instead of encoding an input x to a fixed point z , the encoder learns a distribution over the latent variable z conditioned on the input:

$$q(z | x)$$

Typically, this is a multivariate Gaussian with parameters $\mu(x)$ and $\sigma(x)$ learned by a neural network.

- ▶ **Probabilistic Decoder:** Given a sample z from the latent distribution, the decoder reconstructs the input by modeling:

$$p(x | z)$$

This allows for generating new data by sampling from the latent space.

- ▶ **Latent Space:** The model learns a continuous latent space where similar data points are close together. This latent space captures the underlying structure or factors of variation in the data.

Objectives and Training: The training objective of a VAE is to maximize the **Evidence Lower Bound (ELBO)**:

$$\log p(x) \geq \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{\text{KL}}(q(z|x)\|p(z))$$

- ▶ The first term, $\mathbb{E}_{q(z|x)}[\log p(x|z)]$, encourages accurate reconstruction.
- ▶ The second term, $D_{\text{KL}}(q(z|x)\|p(z))$, regularizes the latent space by making the approximate posterior $q(z|x)$ close to the prior $p(z)$, usually $\mathcal{N}(0, I)$.

Why Use VAEs?:

- ▶ Enable generative modeling: generate new samples by sampling from the latent distribution.
- ▶ Learn smooth, structured, and interpretable latent representations.
- ▶ Provide a principled probabilistic framework for inference and generation.

Latent Variable Models:

- ▶ **Definition:** Models that assume the data is generated from some unobserved (latent) variables.
- ▶ **Goal:** Learn a mapping from observed data x to latent variables z and vice versa.
- ▶ **Generative Process:**
 - Sample latent variable z from a prior distribution $p(z)$.
 - Generate data x from the latent variable using a likelihood function $p(x|z)$.
- ▶ **Inference Problem:** Given observed data x , infer the posterior distribution $p(z|x)$.
- ▶ **Challenge:** Direct computation of posterior $p(z|x)$ is often intractable, leading to the need for approximations.

- ▶ Now, how to train this model?
- ▶ How about following the same strategy as in FVSBNs? Learn model parameters to maximize the likelihood of training data.

$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

- ▶ But there is a problem here. It is Intractible to compute $p(x|z)$ for every z !
- ▶ Intuitively, need to figure out which z corresponds to each x in the dataset, but such mapping is unknown.
- ▶ This also makes posterior density $p(z|x)$ intractable because it depends on $p_{\theta}(x)$

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

► Solution

- Let's approximate $p(z|x)$ with another distribution by another distribution $q(z)$
- If $q(z)$ is a tractable distribution e.g. Gaussian distribution
- **Approach:** We can adjust parameters of $q(z)$ and make it as close to $p(z|x)$, i.e. $q(z) \approx p(z|x)$.
- **Goal:** Minimize the Kullback-Leibler (KL) divergence $KL(q||p)$.

VAE: Variational Inference (cont.)

$$\begin{aligned} KL(q(z)||p(z|x)) &= - \sum q(z) \log \frac{p(z|x)}{q(z)} \\ &= - \sum q(z) \log \frac{\frac{p(x,z)}{p(x)}}{q(z)} \\ &= - \sum q(z) \log \left(\frac{p(x,z)}{q(z)} \frac{1}{p(x)} \right) \\ &= - \sum q(z) \left[\log \frac{p(x,z)}{q(z)} + \log \frac{1}{p(x)} \right] \\ &= - \sum q(z) \left[\log \frac{p(x,z)}{q(z)} - \log p(x) \right] \\ &= - \sum_z q(z) \log \frac{p(x,z)}{q(z)} + \log p(x) \sum_z q(z) \\ &= - \sum_z q(z) \log \frac{p(x,z)}{q(z)} + \log p(x) \quad \because \sum_z q(z) = 1 \end{aligned}$$

VAE: Variational Inference (cont.)



$$KL(q(z)||p(z|x)) = - \sum_z q(z) \log \frac{p(x, z)}{q(z)} + \log p(x)$$

► We can also write above equation as:

$$\log p(x) = KL(q(z)||p(z|x)) + \sum_z q(z) \log \frac{p(x, z)}{q(z)}$$

VAE: Variational Inference (cont.)

- ▶ Given x , $\log p(x)$ is a constant
- ▶ $KL(q(z)||p(z|x))$ is the quantity we wanted to minimize
- ▶ Assume $L = \sum_z q(z) \log \frac{p(x,z)}{q(z)}$, then

$$\text{constant} = KL + L$$

$$L \leq \log p(x) \quad \therefore kl \geq 0$$

- ▶ Instead of minimizing KL we can maximise L

What is ELBO?:

- ▶ Allows us to optimize the model.
- ▶ It provides a lower bound on the log-likelihood of the data, which we aim to maximize during training.
- ▶ The ELBO consists of two main components:
 - **Reconstruction term:** Measures how well the model can reconstruct the input data from the latent representation.
 - **Regularization term:** Encourages the learned latent distribution to be close to a prior distribution (usually Gaussian).

- ▶ Mathematically, the ELBO can be expressed as:

$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x)\|p(z))$$

where:

- $q_\phi(z|x)$ is the approximate posterior distribution (encoder).
 - $p_\theta(x|z)$ is the likelihood of the data given the latent variable (decoder).
 - D_{KL} is the Kullback-Leibler divergence between the approximate posterior and the prior distribution $p(z)$.
- ▶ The goal is to maximize the ELBO with respect to the model parameters θ and ϕ .
 - ▶ By maximizing the ELBO, we ensure that the model learns a meaningful latent representation while also being able to generate new samples.

Looking at Lower bound L

$$\begin{aligned}
 L &= \sum_z q(z) \log \frac{p(x, z)}{q(z)} \\
 &= \sum_z q(z) \log \frac{p(x|z)p(z)}{q(z)} \\
 &= \sum_z q(z) \left[\log p(x|z) + \log \frac{p(z)}{q(z)} \right] \\
 &= \underbrace{\sum_z q(z) \log p(x|z)}_{\text{Expectation } E_{q(z)}(\log p(x|z))} + \underbrace{\sum_z q(z) \log \frac{p(z)}{q(z)}}_{-KL(q(z)||p(z))}
 \end{aligned}$$

So,

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

VAE: Evidence Lower Bound (ELBO) (cont.)

- ▶ $E_{q(z)}(\log p(x|z))$ is conceptually reconstruction
- ▶ We can assume z to be Standard Normal Distribution

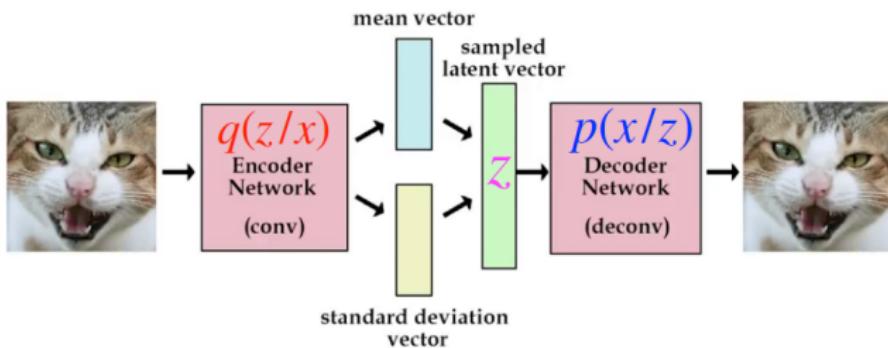
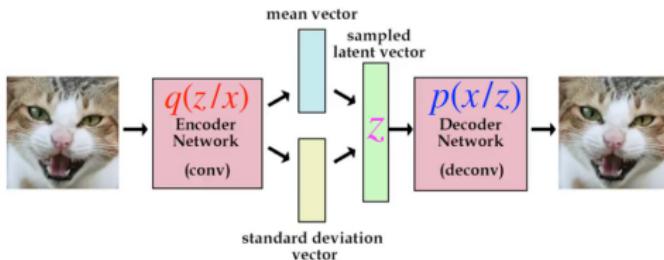


Figure 2: Variational Autoencoder Architecture

VAE: Evidence Lower Bound (ELBO) (cont.)



$$p(x|\hat{x}) = e^{-|x-\hat{x}|^2}$$

$$\log e^{-|x-\hat{x}|^2} = -|x - \hat{x}|^2$$

$$L = E_{q(z)}(-|x - \hat{x}|^2) - KL(q(z)||p(z))$$

$$\min |x - \hat{x}| + KL(q(z|x)||\mathcal{N}(0, 1))$$

$$\min |x - \hat{x}| - 0.5 * (1 + \log \sigma^2 - \sigma^2 - \mu^2)$$

For full derivation of KL Loss, read [here](#)

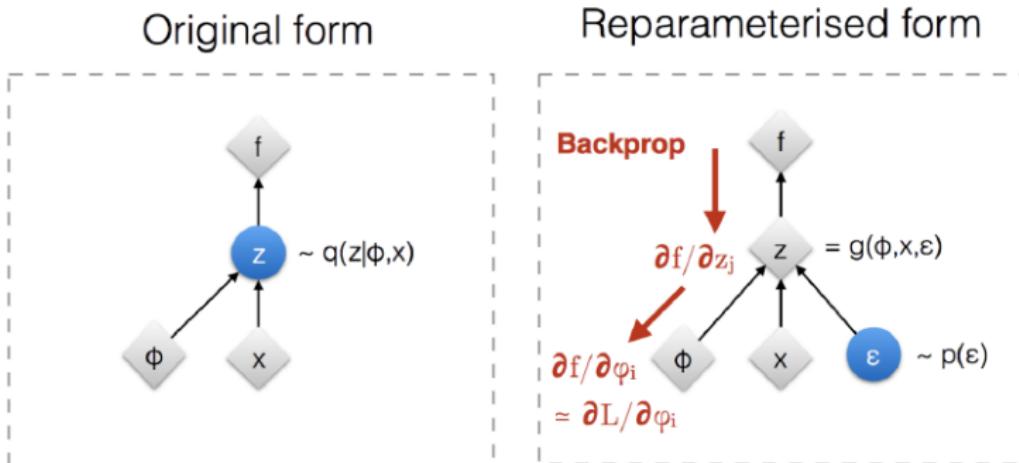
Problem: Cannot backpropagate through stochastic sampling.

Solution: Reparameterize z as:

$$z = \mu + \sigma \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

Benefit: Enables gradient-based optimization by making the sampling operation differentiable.

VAE: Reparameterization Trick (cont.)



◆ : Deterministic node
● : Random node

[Kingma, 2013]
[Bengio, 2013]
[Kingma and Welling 2014]
[Rezende et al 2014]

Figure 3: Reparameterization trick to make back propagation possible

VAE: Reparameterization Trick (cont.)

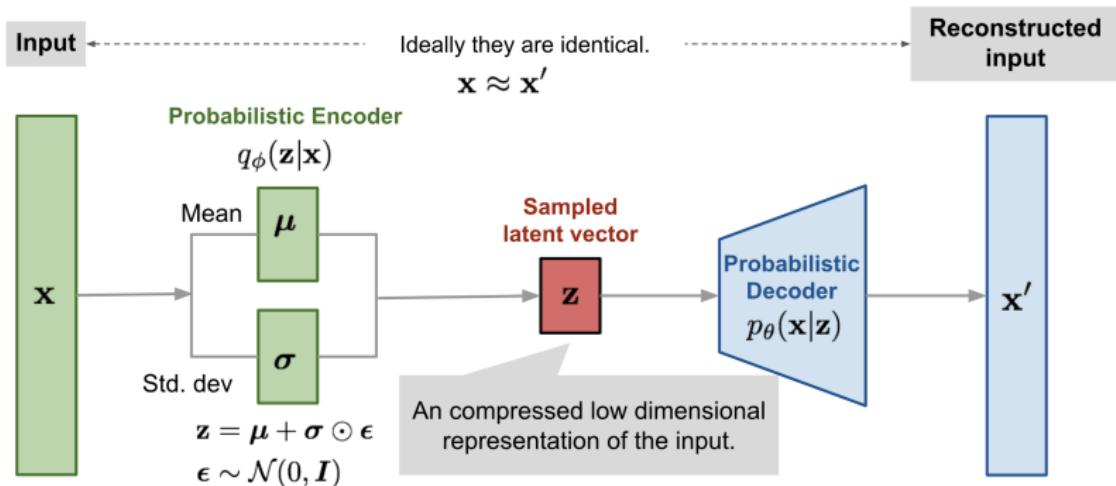


Figure 4: Variational Autoencoder with reparameterization trick

Total Loss:

$$L = \text{Reconstruction Loss} + \text{KL Divergence}$$

Reconstruction Loss:

- ▶ Measures how well \hat{x} matches x .
- ▶ Common choices: Mean Squared Error (MSE) or Binary Cross-Entropy.

KL Divergence:

- ▶ Measures how much $q(z | x)$ diverges from the prior $p(z)$.
- ▶ Encourages the latent space to follow a standard normal distribution.

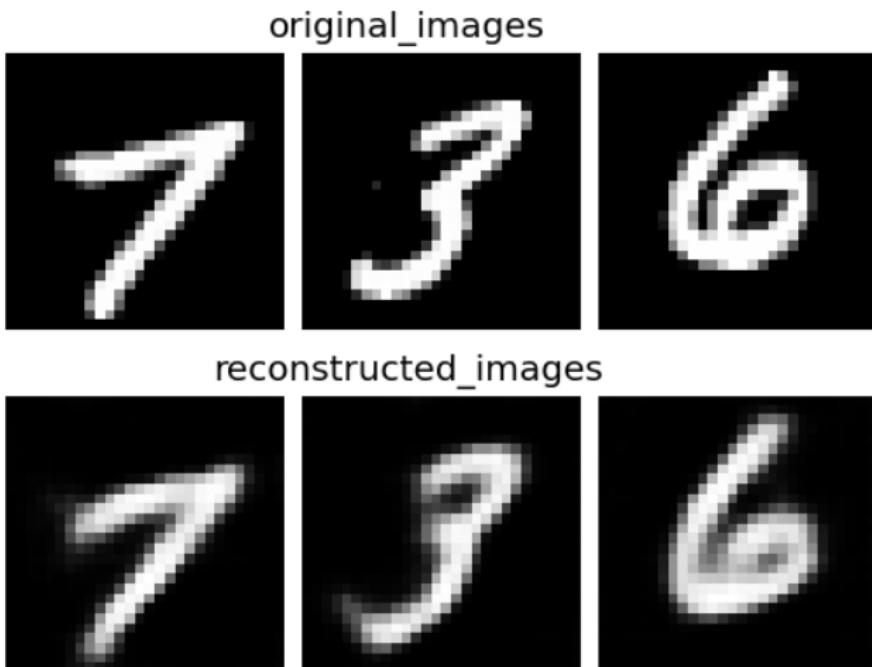


Figure 5: Image reconstruction with variational autoencoders on MNIST digits dataset

VAE: Results (cont.)

generated_images

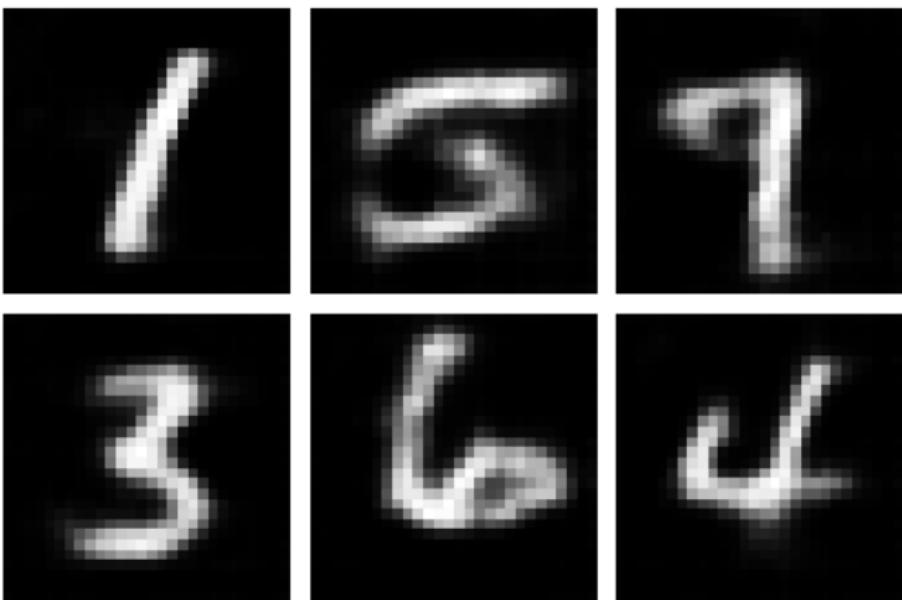


Figure 6: Image generation with variational autoencoders on MNIST digits dataset. Sample an encoding vector from $\mathcal{N}(0, 1)$ and passed it through decoder

VAE: Results (cont.)

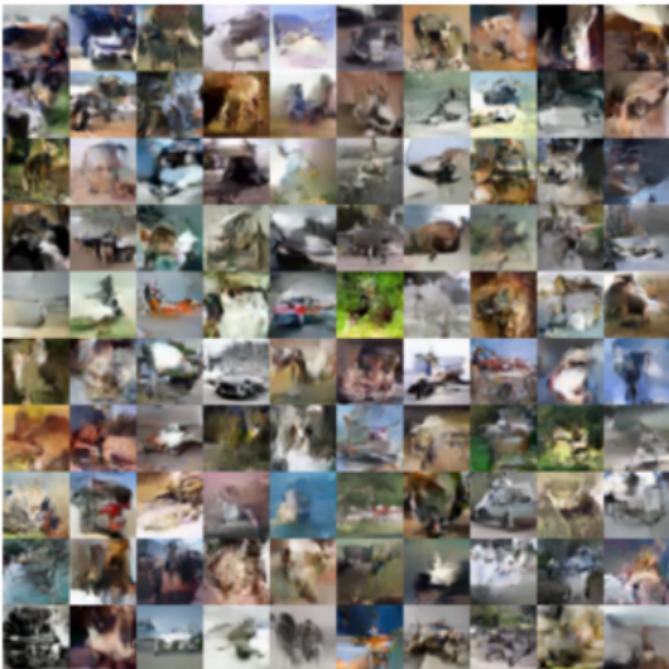


Figure 7: Image generation with variational autoencoders on CIFAR-10 32x32 dataset

β -VAE:

- ▶ Introduces a hyperparameter β to control the trade-off between reconstruction and regularization.
- ▶ Encourages disentangled representations.

Conditional VAE (CVAE):

- ▶ Incorporates additional information y (e.g., class labels) into the encoder and decoder.
- ▶ Enables generation of data conditioned on y .

Discrete VAE:

- ▶ Utilizes discrete latent variables.
- ▶ Techniques like Gumbel-Softmax are used for differentiable sampling.

- ▶ **Basic Idea:** Different neurons in latent space should be uncorrelated, i.e. they all try to learn something different about input data.
- ▶ **Implementation:**

$$\mathcal{L}(\theta, \phi; x, z, \beta) = E_{q_\phi(z|x)}(\log p_\theta(x|z)) - \beta KL(q_\phi(z|x)||p(z))$$

- ▶ Increasing the β is forcing variational autoencoder to encode the information in only few latent variables

Disentangled Variational Autoencoders (β -VAEs) (cont.)



Figure 8: Azimuthal rotation in β -VAEs and simple VAEs. β -VAEs produce more disentangled rotation, whereas some other features also change in simple VAEs.

Limitations:

- ▶ **Gaussian Assumption:** The assumption that the latent variables follow a Gaussian distribution may not hold for all datasets.
- ▶ **Over-smoothing:** The model may produce overly smooth reconstructions, losing fine details in the data.
- ▶ **Sensitivity to Hyperparameters:** The performance of VAEs can be sensitive to the choice of hyperparameters, such as the weight of the KL divergence term.
- ▶ **Computational Complexity:** Training VAEs can be computationally expensive, especially for large datasets or complex models.
- ▶ **Evaluation Metrics:** Evaluating the quality of generated samples can be subjective and challenging, as traditional metrics may not capture the nuances of the data.

Challenges:

- ▶ **Training Instability:** VAEs can be difficult to train, especially with complex datasets.
- ▶ **Mode Collapse:** The model may generate samples from only a subset of the latent space.
- ▶ **Balancing Reconstruction and Regularization:** Finding the right balance between reconstruction loss and KL divergence can be tricky.
- ▶ **Posterior Collapse:** In some cases, the model may ignore the latent variables, leading to poor representations.
- ▶ **Limited Expressiveness:** The Gaussian assumption for the latent space may not capture complex data distributions.
- ▶ **Disentanglement:** Achieving disentangled representations can be challenging, especially in high-dimensional spaces.

Future Directions:

- ▶ **Improved Training Techniques:** Developing better optimization methods to stabilize training.
- ▶ **Advanced Architectures:** Exploring more complex latent variable models, such as Normalizing Flows or Hierarchical VAEs.
- ▶ **Better Regularization Techniques:** Investigating alternative regularization methods to improve disentanglement and representation quality.
- ▶ **Hybrid Models:** Combining VAEs with other generative models (e.g., GANs) to leverage their strengths.
- ▶ **Application-Specific Variants:** Tailoring VAEs for specific applications, such as text or video generation.

- ▶ Add a probabilistic spin to Autoencoders to make them generative models
- ▶ Assume Z to be from Gaussian Distribution.
- ▶ But $p(z|x)$ is intractable.
- ▶ **Solution:** Approximate $p(z|x)$ with Gaussian distribution $q(z)$.
- ▶ To minimize the KL divergence between them maximize the Evidence lower bound

$$L = \sum_z q(z) \log \frac{p(x, z)}{q(z)}$$

- ▶ But image produced by variational autoencoders are blurry.

Reference Slides

- ▶ Fei-Fei Li "Generative Deep Learning" CS231
- ▶ Hao Dong "Deep Generative Models"
- ▶ Hung-Yi Lee "Machine Learning"
- ▶ Murtaza Taj "Deep Learning" CS437
- ▶ Aykut Erdem, COMP547: Deep Unsupervised Learning, Koc University

Credits

Dr. Prashant Aparajeya

Computer Vision Scientist — Director(AISimply Ltd)

p.aparajeya@aisimply.uk

This project benefited from external collaboration, and we acknowledge their contribution with gratitude.