

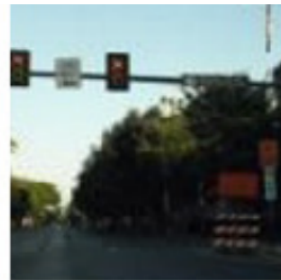
Deep Generative Models

Advanced Statistical Inference

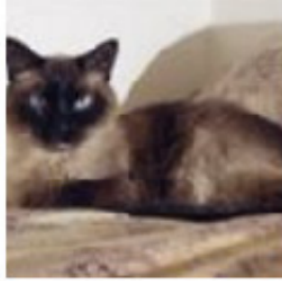
Simone Rossi

Generative models

Given training data from an unknown distribution $p_{\text{data}}(\boldsymbol{x})$, we want to learn a model $p_{\text{model}}(\boldsymbol{x})$ that approximates the true distribution.



Train from $\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})$.

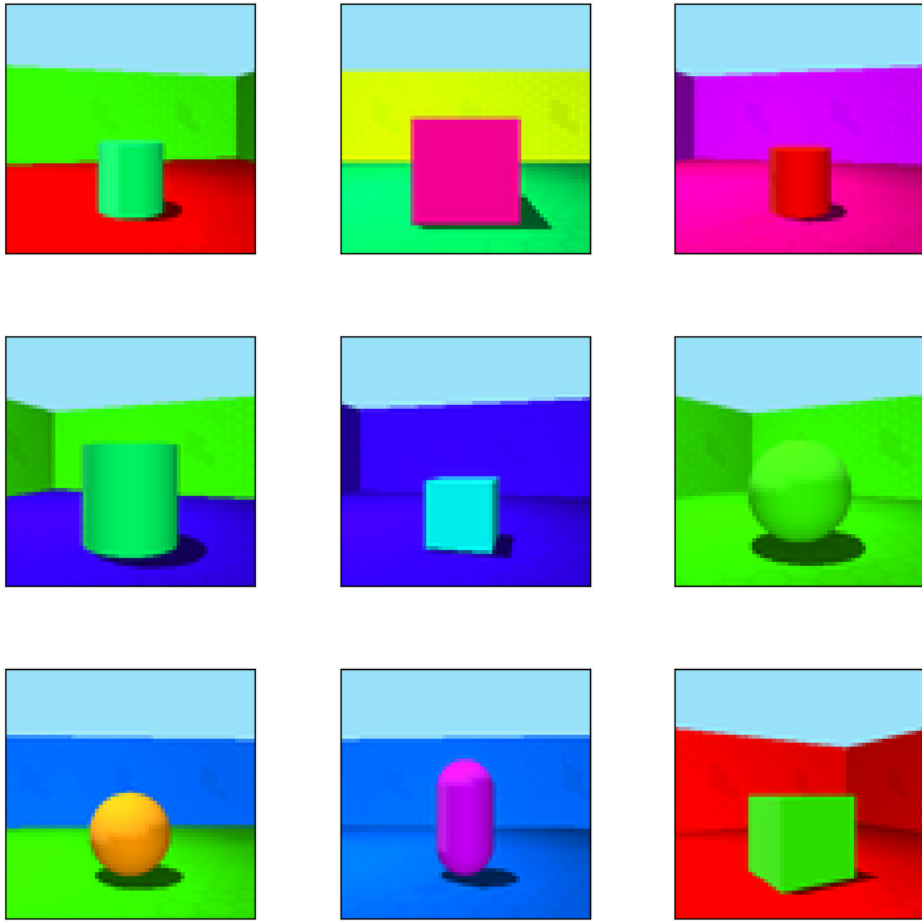


Generate from $\boldsymbol{x} \sim p_{\text{model}}(\boldsymbol{x})$.

Latent variable models

Our objective is to learn the data distribution $p(\boldsymbol{x})$, but we suppose that each data point \boldsymbol{x} is associated with a **latent variable** \boldsymbol{z} .

For example, image we want to learn the distribution of images of objects, like the one below:



Each image is a huge vector of pixels $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$, where H is the height, W is the width and C is the number of channels (in this case, $64 \times 64 \times 3$).

Each image can be described by a set of 6 latent variables: floor, wall and object color, shape, orientation and scale.

Latent variable models

Latent variables are unobserved variables, we cannot measure them directly, but we can infer them from the observed data.

In math terms, we model our data distribution as a **marginal** of the joint distribution of the observed data and the latent variables:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad \text{or} \quad p(\mathbf{x}) = \sum_i p(\mathbf{x}, \mathbf{z}_i)$$

where $p(\mathbf{x}, \mathbf{z})$ is the **joint distribution** of the observed data and the latent variables.

Linear latent variable models

Gaussian mixture model (GMM)

Conditional distribution of the data given the latent variable:

$$p(\mathbf{x} \mid \mathbf{z} = k) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ are the mean and covariance of the k -th Gaussian component.

Prior distribution of the latent variable:

$$p(\mathbf{z}) = \text{Categorical}(\mathbf{z} \mid \boldsymbol{\pi})$$

where π_k are the mixing coefficients, $\boldsymbol{\mu}_k$ is the mean of the k -th Gaussian component, and $\boldsymbol{\Sigma}_k$ is the covariance of the k -th Gaussian component.

Probabilistic PCA (PPCA)

Conditional distribution of the data given the latent variable:

$$p(\mathbf{x} \mid \mathbf{z}) = \mathcal{N}(\mathbf{x} \mid \mathbf{W}\mathbf{z} + \mathbf{b}, \sigma^2 \mathbf{I})$$

where \mathbf{W} is a linear transformation matrix, \mathbf{b} is a bias vector, and $\sigma^2 \mathbf{I}$ is the noise covariance.

Prior distribution of the latent variable:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$$

Variational autoencoders (VAEs)

Variational autoencoders (VAEs)

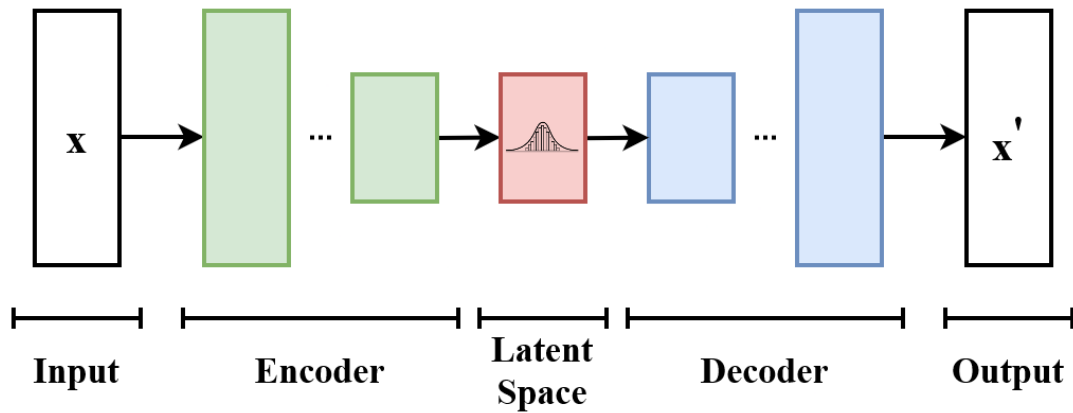
Idea:

1. We want to learn a generative model $p_{\text{model}}(\mathbf{x})$ that approximates the true data distribution $p_{\text{data}}(\mathbf{x})$.
2. We introduce a latent variable \mathbf{z} to explain the observed data \mathbf{x} .
3. We model the relationship between the data and the latent variable with a more complex function.
4. We use a variational inference approach to learn the model parameters.

Variational autoencoders (VAEs)

Disclaimer:

- Variational autoencoders (VAEs) are often misrepresented as a simple autoencoder with a probabilistic twist.
- This is **wrong** (sort of)! It hides the real nature of the encoder and decoder!



Neural networks as flexible likelihood models

Consider a maximum likelihood training objective for a generative model with latent variables:

$$p(\mathbf{x}) = \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z}) d\mathbf{z}$$

Objective: increase the flexibility of $p(\mathbf{x} | \mathbf{z})$ to approximate the true data distribution $p_{\text{data}}(\mathbf{x})$.

...

Solution: define a noisy observation model

$$p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{f}(\boldsymbol{\theta}, \mathbf{z}), \sigma^2 \mathbf{I})$$

where $\mathbf{f}(\boldsymbol{\theta}, \mathbf{z})$ is a function of the latent variable \mathbf{z} and the model parameters $\boldsymbol{\theta}$.

...

Note: $p(\mathbf{x} | \mathbf{z})$ is *not* a deterministic function, it is a probabilistic model that outputs a distribution over the data given the latent variable. If it were deterministic, then $p(\mathbf{x}) = 0$ almost everywhere.

...

Notation: we will use $p(\mathbf{x}; \boldsymbol{\theta})$ and $p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})$ to remember that these two distributions depend on the model parameters $\boldsymbol{\theta}$.

Neural networks as flexible likelihood models

$p(\mathbf{x}; \boldsymbol{\theta}) = \int p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})p(\mathbf{z}) d\mathbf{z}$ is well-defined

...

- Can we compute it?
 - **No**, $\mathbf{f}(\boldsymbol{\theta}, \mathbf{z})$ is very complex, so there is no hope of finding a closed-form solution.

...

- Can we maximize it w.r.t. $\boldsymbol{\theta}$?
 - **No**, we cannot compute the integral, so we cannot maximize the likelihood directly.

...

Solution: try to maximize a lower bound on the likelihood $\log p(\mathbf{x}; \boldsymbol{\theta})$

$$\log \mathbb{E}[X] \geq \mathbb{E}[\log X]$$

Variational inference to the rescue

We can use Jensen's inequality to obtain a lower bound on the log-likelihood.

Suppose we have a distribution $q(\mathbf{z})$ (we will see later how to choose it):

$$\begin{aligned} \log p(\mathbf{x}; \boldsymbol{\theta}) &= \log \int p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z} \\ &= \log \int q(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \\ &\geq \int q(\mathbf{z}) \log \left(\frac{p(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) \right) d\mathbf{z} \quad (\text{Jensen's inequality}) \\ &= \mathbb{E}_{q(\mathbf{z})} \frac{p(\mathbf{z})}{q(\mathbf{z})} + \mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) \end{aligned}$$

Variational inference to the rescue

Remember the **Variational Inference** lecture?

We can rewrite the lower bound as:

$$\log p(\mathbf{x}; \boldsymbol{\theta}) \geq \mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) - \text{KL}(q(\mathbf{z}) || p(\mathbf{z}))$$

...

1. The first term is the expected log-likelihood of the data given the latent variable. Since we assume a Gaussian model:

$$\begin{aligned} \log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) &= \log \mathcal{N}(\mathbf{x} | \mathbf{f}(\boldsymbol{\theta}, \mathbf{z}), \sigma^2 \mathbf{I}) \\ &= \log \left(\frac{1}{(2\pi\sigma^2)^{D/2}} \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{f}(\boldsymbol{\theta}, \mathbf{z})\|^2 \right) \right) \\ &= -\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{f}(\boldsymbol{\theta}, \mathbf{z})\|^2 + \text{const} \end{aligned}$$

...

Meaning that we want to minimize the expected reconstruction error between the data \mathbf{x} and the model output $\mathbf{f}(\boldsymbol{\theta}, \mathbf{z})$, when \mathbf{z} is sampled from $q(\mathbf{z})$.

Variational inference to the rescue

2. The second term is the KL divergence between the variational distribution $q(\mathbf{z})$ and the prior distribution $p(\mathbf{z})$:

$$\text{KL}(q(\mathbf{z})\|p(\mathbf{z})) = \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z}$$

This term acts as a regularizer, encouraging the variational distribution $q(\mathbf{z})$ to be close to the prior distribution $p(\mathbf{z})$.

Optimization objective

$$\log p(\mathbf{x}; \boldsymbol{\theta}) \geq \mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z})} \log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) - \text{KL}(q(\mathbf{z})\|p(\mathbf{z}))$$

We would like to maximize the lower bound on the log-likelihood, but we want the bound to be tight

...

Remember that the KL divergence is always non-negative, so we can rewrite the objective as:

$$\log p(\mathbf{x}; \boldsymbol{\theta}) - \mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta}) = \text{KL}(q(\mathbf{z})\|p(\mathbf{z} | \mathbf{x}))$$

This means that the closer the variational distribution $q(\mathbf{z})$ is to the true posterior distribution $p(\mathbf{z} | \mathbf{x})$, the tighter the bound will be.

Approximate posterior inference

- To fit q , we need to define a **variational distribution** $q(\mathbf{z})$ that approximates the true posterior distribution $p(\mathbf{z} | \mathbf{x})$.
- The variational distribution is often chosen to be Gaussian, i.e., $q(\mathbf{z}; \boldsymbol{\nu}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})$, where $\boldsymbol{\nu} = \{\boldsymbol{\mu}, \boldsymbol{\sigma}^2\}$ are the variational parameters.
- Gaussian is good because it is easy to sample from, it has a closed-form KL divergence, and it is possible to use the reparameterization trick to compute gradients of the ELBO w.r.t. the variational parameters $\boldsymbol{\nu}$.

...

! Important

The **reparameterization trick** allows us to sample from the variational distribution $q(\mathbf{z}; \boldsymbol{\nu})$ in a way that is differentiable w.r.t. the variational parameters $\boldsymbol{\nu}$.

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Training a VAE

1. We need to maximize the likelihood $\log p(\mathbf{x}; \boldsymbol{\theta})$. But we cannot do it directly.
2. We derive a lower bound on the likelihood, the **evidence lower bound (ELBO)** $\mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta})$, which we can maximize instead. But we need to define a variational distribution $q(\mathbf{z})$ that approximates the true posterior distribution $p(\mathbf{z} | \mathbf{x})$.
3. We parametrize the variational distribution $q(\mathbf{z}; \boldsymbol{\nu})$ with a Gaussian distribution, where $\boldsymbol{\nu} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$ are the variational parameters.

...

At the end, we are left with the following optimization problem:

$$\max_{\boldsymbol{\theta}, \boldsymbol{\nu}} \mathcal{L}_{\text{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\nu}) = \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\nu})} \log p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) - \text{KL}(q(\mathbf{z}; \boldsymbol{\nu}) \| p(\mathbf{z}))$$

Note:

- The ELBO is a function of both the likelihood parameters $\boldsymbol{\theta}$ and the variational parameters $\boldsymbol{\nu}$.
- The likelihood are **shared** across all data points, while the variational parameters are **per-sample**.

Training a VAE via joint optimization is hard

In classic VI (e.g., Bayesian logistic regression):

- ELBO depends **only** on $q(\boldsymbol{\theta})$
- Optimization is stable — just improve the approximation to the true posterior

In our case, for variational autoencoders:

- We optimize **both**:
 - $\boldsymbol{\theta}$: parameters of the generative model

- ν : variational parameters
- But these two are **coupled**:
 - ELBO depends on **both** sets of parameters
 - Changing θ shifts the optimal ν , changing ν affects how we learn θ .
 - No close form solution for ν given θ exists (no EM-like solution).

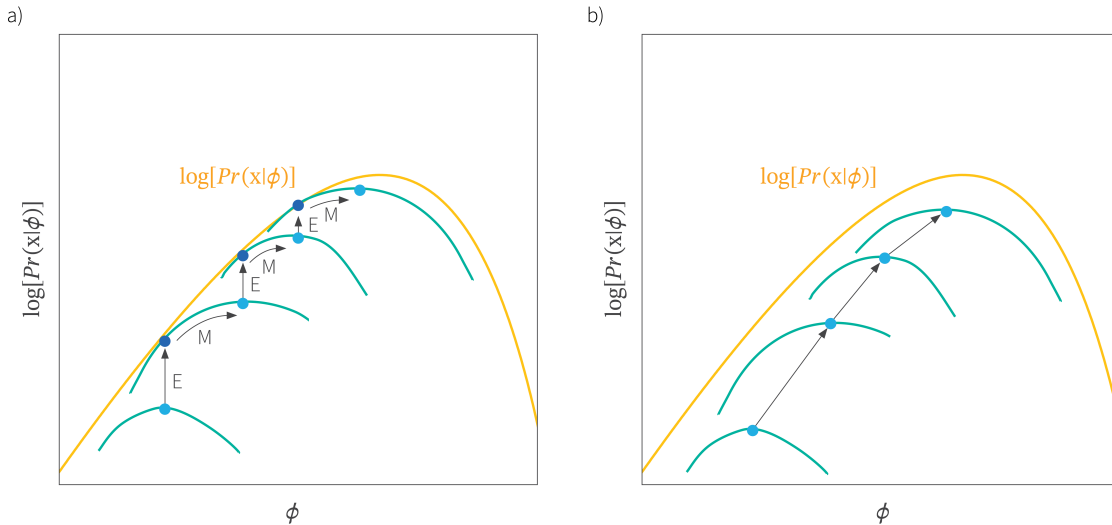
...

This may lead to unstable or inefficient training, especially when $f(\cdot)$ is complex (e.g., deep neural networks).

Expectation-Maximization vs Variational Inference

The joint optimization of θ and ν in VAEs is similar to the **Expectation-Maximization (EM)** algorithm, but with some key differences:

1. In EM, the E-step computes the exact posterior distribution $p(\mathbf{z} | \mathbf{x}; \theta)$, while in VAEs, we approximate it with a variational distribution $q(\mathbf{z}; \nu)$.
2. In EM, the M-step maximizes the exact expected log-likelihood, while in VAEs, we maximize the ELBO, which is a lower bound on the log-likelihood.

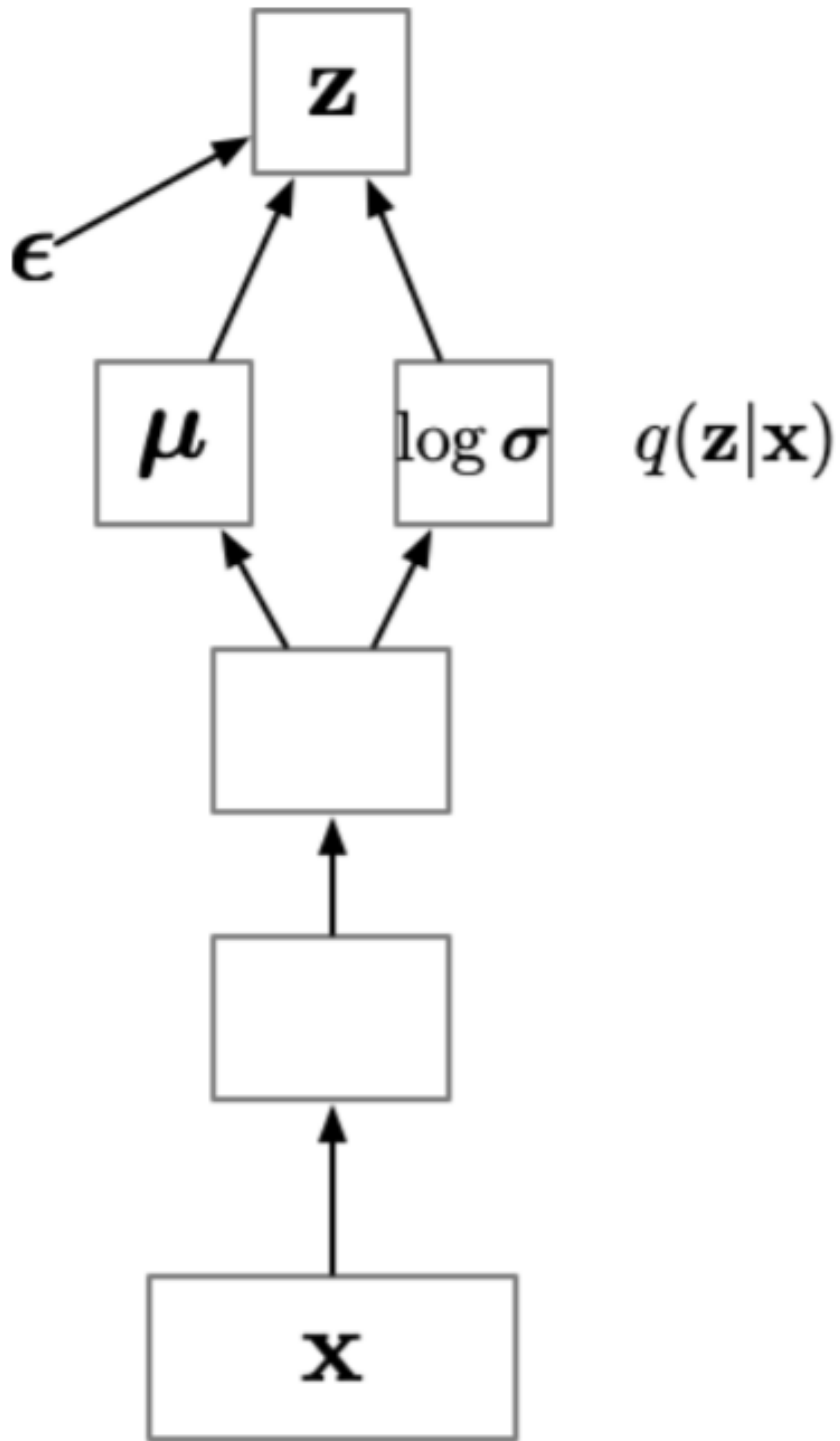


Amortization

- In the previous approach we learn a separate set of variational parameters ν for each data point x .
- This requires an expensive iterative optimization for each data point, which will take a long time to process large datasets.
- **Amortization:** instead of learning a separate set of variational parameters for each data point, we learn a function $g(\phi, x)$ that maps the data point x to the variational parameters ν .
- This function is often implemented as a neural network, which allows us to learn a complex mapping from the data to the variational parameters.

Amortization

- **Idea:** “amortize” the cost of computing the variational parameters ν by learning a function that predicts $\{\mu, \sigma^2\}$ as a function of the data x .
- The output of this function is generally μ and $\log \sigma$ (the log is needed to ensure that the variance is positive).
- The reparametrization trick is still used to sample from the variational distribution, the only difference is that the variational parameters are not free variables anymore, but are predicted by a neural network $g(\phi, x)$.
- Sometimes, this is noted as $q(z | x)$ to emphasize that the variational distribution depends on the data x , even though it’s not actually a conditional distribution.

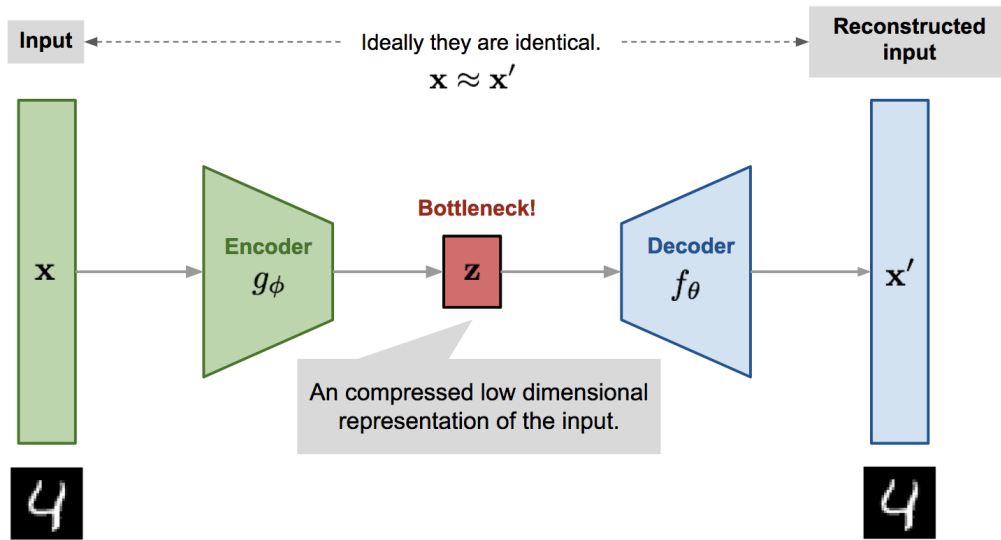


Putting it all together

- We have a neural network $g(\mathbf{x}; \phi)$ that predicts the variational parameters $\{\mu, \sigma^2\}$ from the data \mathbf{x} .
- We can sample from the variational distribution $q(\mathbf{z} | \mathbf{x}; \phi)$ using the reparameterization trick
- We compute the likelihood $\log p(\mathbf{x} | \mathbf{z}; \theta)$ using a neural network $f(\mathbf{z}; \theta)$ that maps the latent variable \mathbf{z} to the data space.
- We can now optimize the ELBO with respect to both the likelihood parameters θ and the amortization parameters ϕ :

$$\mathcal{L}_{\text{ELBO}}(\theta, \phi) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\phi)} \log p(\mathbf{x} | \mathbf{z}; \theta) - \text{KL}(q(\mathbf{z} | \mathbf{x}; \phi) || p(\mathbf{z}))$$

Putting it all together



Generating samples from the VAE

Follow the generation process:

1. Sample from the latent space $z \sim p(\mathbf{z})$, where $p(\mathbf{z})$ is the **prior** distribution.
2. Decode the latent variable z to generate a sample $x \sim p(\mathbf{x} | \mathbf{z})$ using the decoder network.

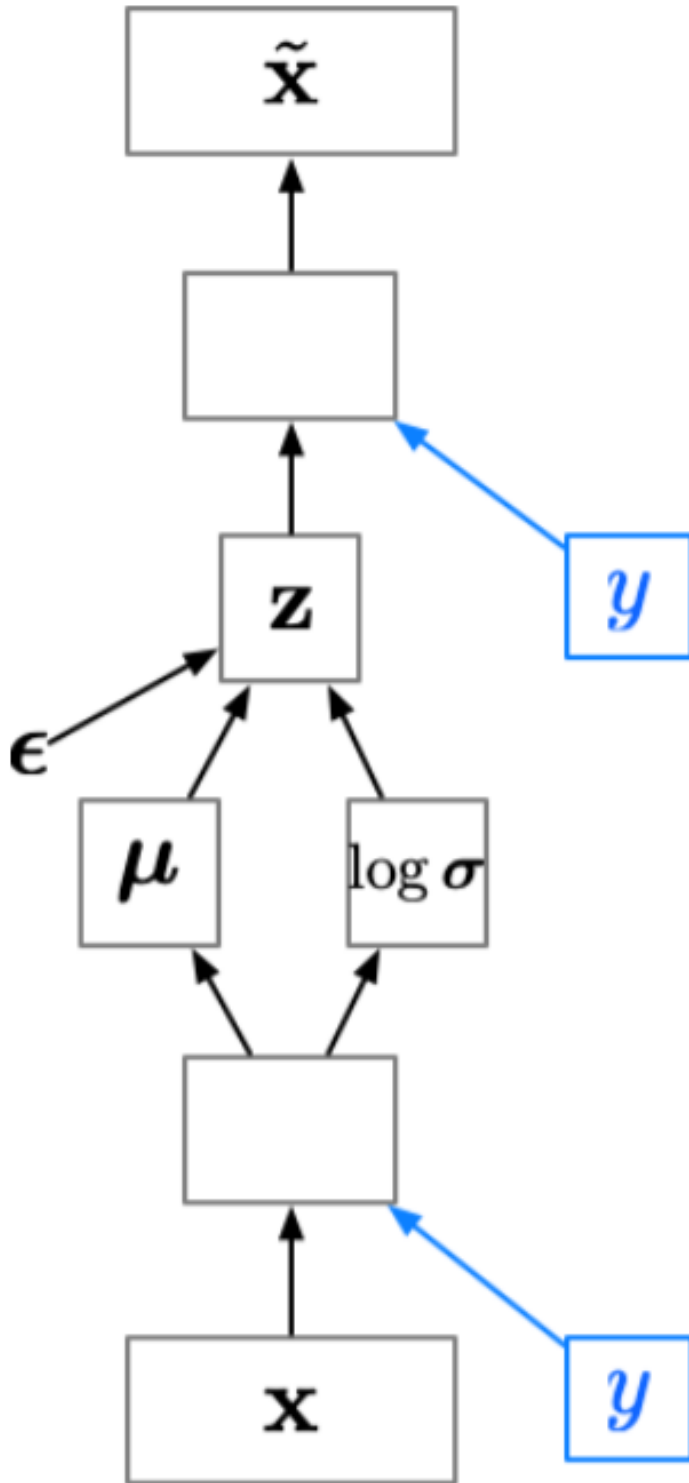
! Note

When using amortized inference, the encoder network is not used during generation, since we are not conditioning on any observed data

It does not make sense to compute $q(\mathbf{z} | \mathbf{x})$ when we don't have any \mathbf{x}

Conditional VAE (CVAE)

- Sometimes we want to generate samples conditioned on some additional information $p(\mathbf{x} | \mathbf{y})$, e.g., we want to generate images of a specific class.
- Since the latent code \mathbf{z} no longer has to model the image category, it can focus on modeling the stylistic features.
- If we're lucky, this lets us disentangle style and content. (Note: disentanglement is still a dark art.)
- **Conditional VAE** can be achieved by adding an additional input to the encoder and decoder networks, which is the additional information \mathbf{y} :



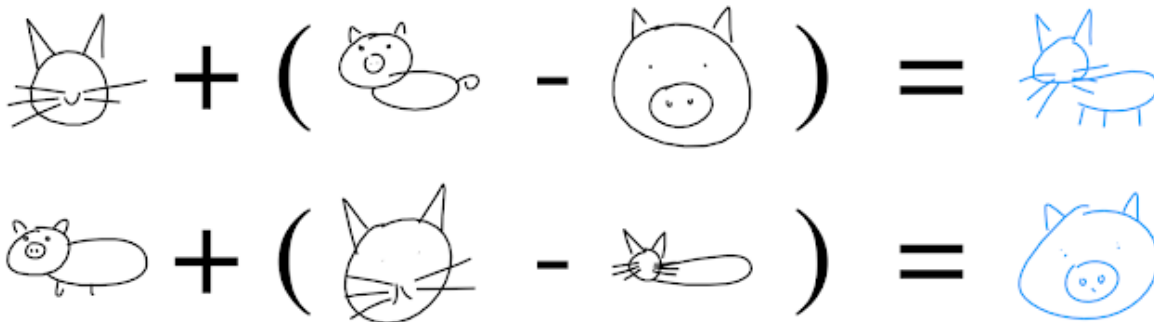
Latent space analysis

You can often get interesting results by interpolating between two vectors in the latent space:



...

You can also perform “arithmetic” in the latent space, e.g. by adding or subtracting vectors:



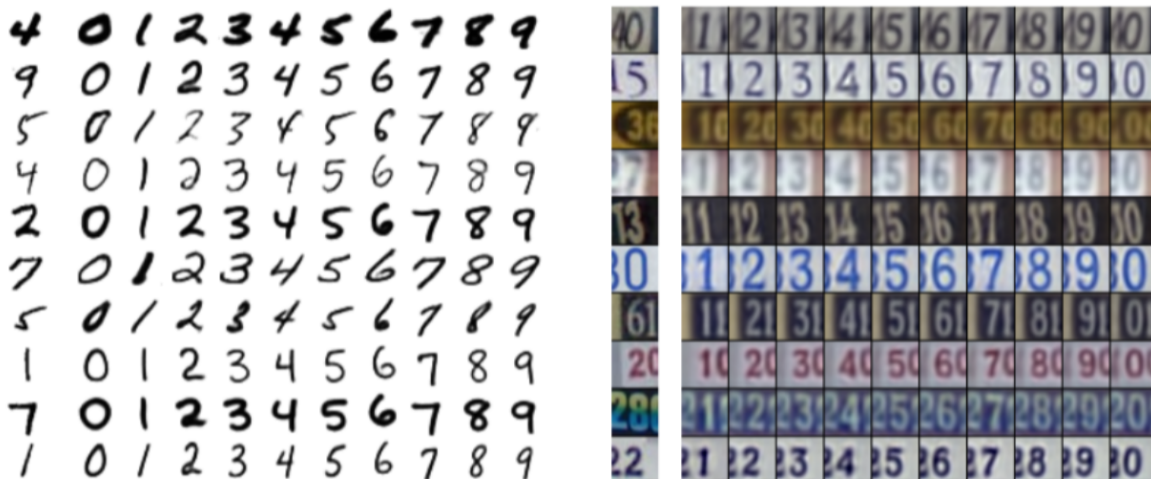
Latent space visualization

By varying two latent dimensions (i.e. dimensions of z) while holding y fixed, we can visualize the latent space.



Latent space interpolation

By varying the label y while holding z fixed, we can solve image analogies.



Problems with VAEs

- The samples may not be very good:
 - VAEs suffer from **mode collapse**, where the amortization network learns to predict the prior. This is due to the strong effect of the KL divergence term in the ELBO.
 - The spherical prior $p(z)$ is too restrictive, it can lead to a “blurry” output distribution (or very noisy samples, if we sample from the likelihood).

