

Deep Generative Models: Variational Auto Encoders

Fall Semester 2025

René Vidal

Director of the Center for Innovation in Data Engineering and Science (IDEAS)

Rachleff University Professor, University of Pennsylvania

Amazon Scholar & Chief Scientist at NORCE



The story up till now

- **Step 1:** We set out with our original goal of learning a model p_θ that gives maximum likelihood to our datapoints \mathbf{x}_i
- **Step 2:** We introduced latent variables \mathbf{z} such that $\mathbf{z} \sim p(\mathbf{z})$ and $\mathbf{x} | \mathbf{z} \sim p(\mathbf{x} | \mathbf{z})$, which gave us the marginalization $p(\mathbf{x}) = \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z}) d\mathbf{z}$
 - Step 2a: When we assumed $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$, and $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{z} + \mathbf{b}, \sigma^2 \mathbf{I})$, we could solve for $(\mathbf{W}, \mathbf{b}, \sigma^2)$ in closed form! This gave us PPCA.
- **Step 3:** We set up variational inference because sadly, not everything in life is Gaussian and linear. This gave us a new Evidence Lower Bound (ELBO) objective:
$$\max_{\theta} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i) = \max_{\theta} \max_{q(\cdot | \mathbf{x}_i), \forall i} \sum_{i=1}^N \int q(\mathbf{z} | \mathbf{x}_i) \log \frac{p_{\theta}(\mathbf{x}_i, \mathbf{z})}{q(\mathbf{z} | \mathbf{x}_i)} d\mathbf{z}$$
 - Step 3a: When the integral is easy to evaluate, we can alternate between optimizing w.r.t. θ with $q(\mathbf{x} | \mathbf{z})$ fixed and vice versa, leading to the Expectation Maximization (EM) algorithm.

The story continues with Variational Autoencoders

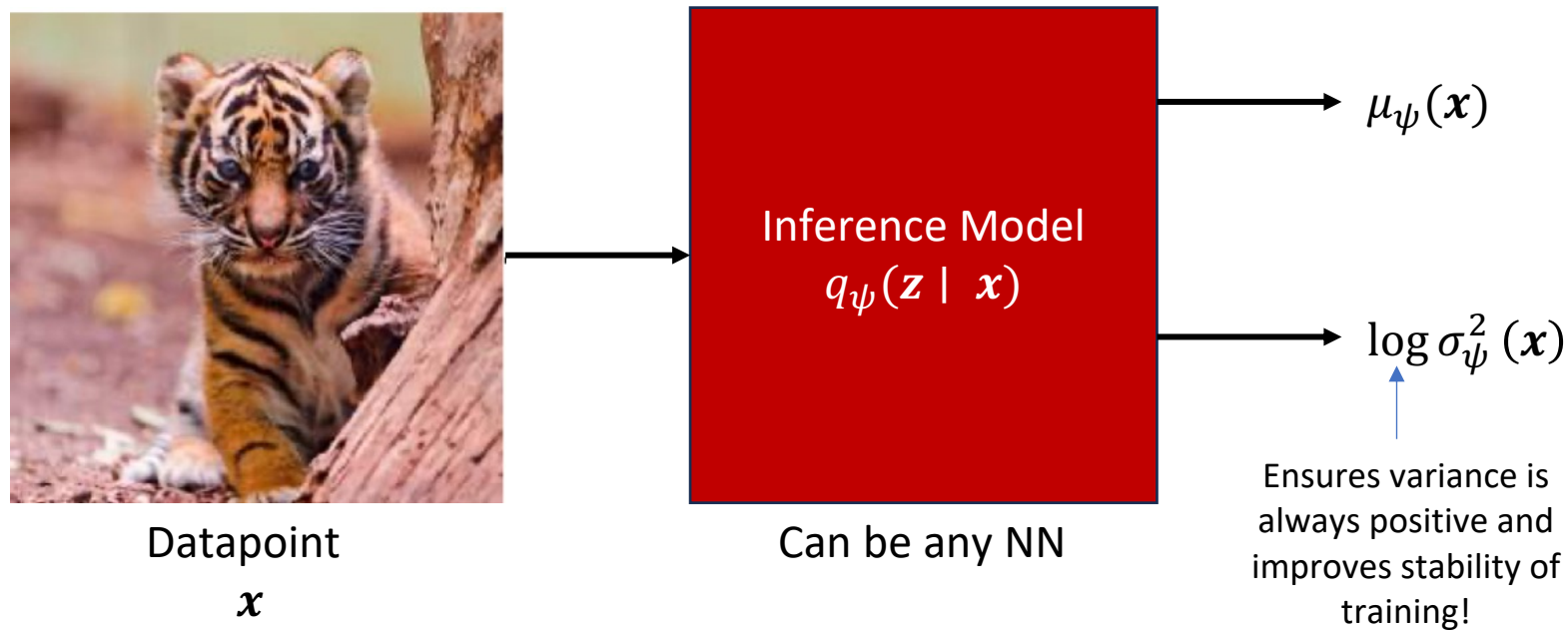
- Before introducing VAEs formally, let us decompose the ELBO further

$$\begin{aligned}\max_{\theta} \log p_{\theta}(\mathbf{x}_i) &= \max_{\theta} \max_q \sum_{i=1}^N \int q(\mathbf{z}|\mathbf{x}_i) \log \frac{p_{\theta}(\mathbf{x}_i, \mathbf{z})}{q(\mathbf{z}|\mathbf{x}_i)} d\mathbf{z} \\ &\geq \max_{\theta, \psi} \sum_i \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x}_i)} \log \frac{p_{\theta}(\mathbf{x}_i, \mathbf{z})}{q_{\psi}(\mathbf{z}|\mathbf{x}_i)} \quad \text{Evidence Lower Bound (ELBO)} \\ &= \max_{\theta, \psi} \sum_i \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x}_i)} \log p_{\theta}(\mathbf{x}_i|\mathbf{z}) \frac{p(\mathbf{z})}{q_{\psi}(\mathbf{z}|\mathbf{x}_i)} \\ &= \max_{\theta, \psi} \sum_i \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x}_i)} \log p_{\theta}(\mathbf{x}_i|\mathbf{z}) + \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x}_i)} \log \frac{p(\mathbf{z})}{q_{\psi}(\mathbf{z}|\mathbf{x}_i)} \\ &\quad \quad \quad \uparrow \\ &\quad \quad \quad -D_{KL}(q_{\psi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))\end{aligned}$$

Variational AutoEncoders (VAEs): Setup

- We have three models we need to define for VAE model

1. **Inference model $q_{\psi}(\mathbf{z} \mid \mathbf{x})$** : We will define as $q_{\psi}(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_{\psi}(\mathbf{x}), \sigma_{\psi}^2(\mathbf{x})\mathbf{I})$, i.e., a normal distribution with learned mean and covariance



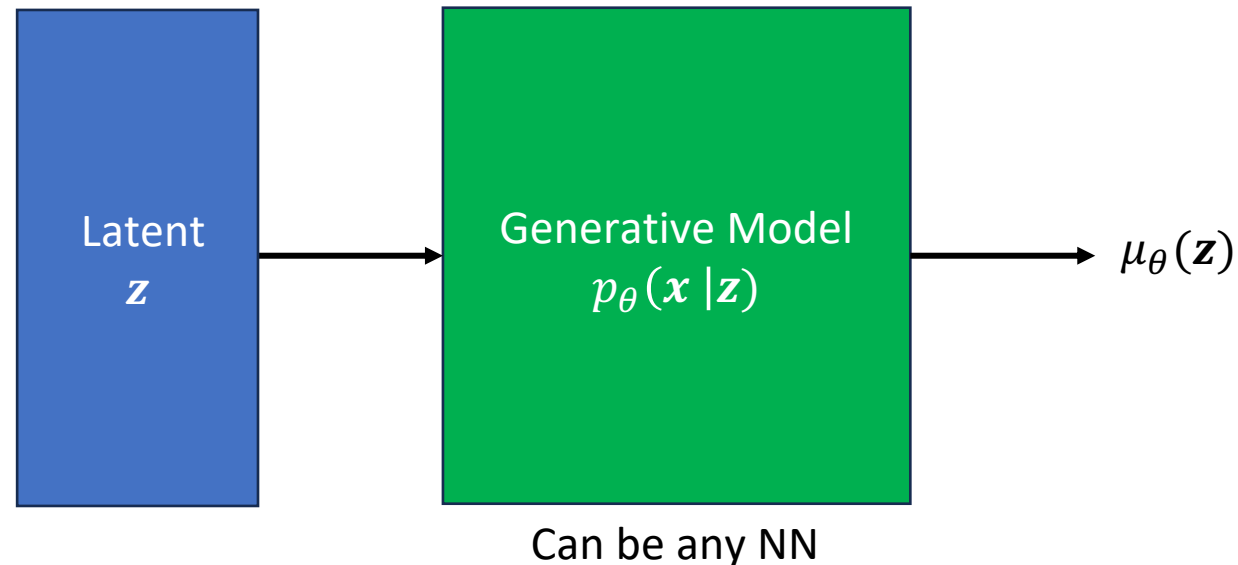
2. **Prior model $p(\mathbf{z})$** : We will define prior for latent variables as $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$

Variational AutoEncoders (VAEs): Setup

- We have three models we need to define for VAE model

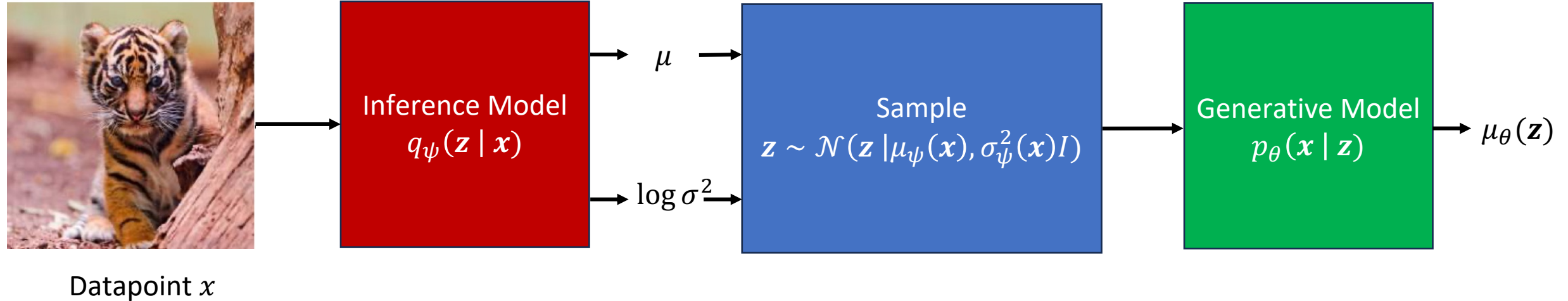
3. **Generative model $p_{\theta}(\mathbf{x} | \mathbf{z})$:** We will define as

- $p_{\theta}(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{z}; \mu_{\theta}(\mathbf{z}), \eta^2 \mathbf{I})$, i.e., a normal distribution with learned mean and variance
- $p_{\theta}(\mathbf{x} | \mathbf{z}) = \text{Cat}(\mathbf{z}; \pi_{\theta}(\mathbf{z}))$, i.e., a categorical distribution with learned class probabilities



- Note this can be defined in many different ways, yielding different models (such as a categorical distribution over 255 values of each pixel)

Variational Autoencoders: Training

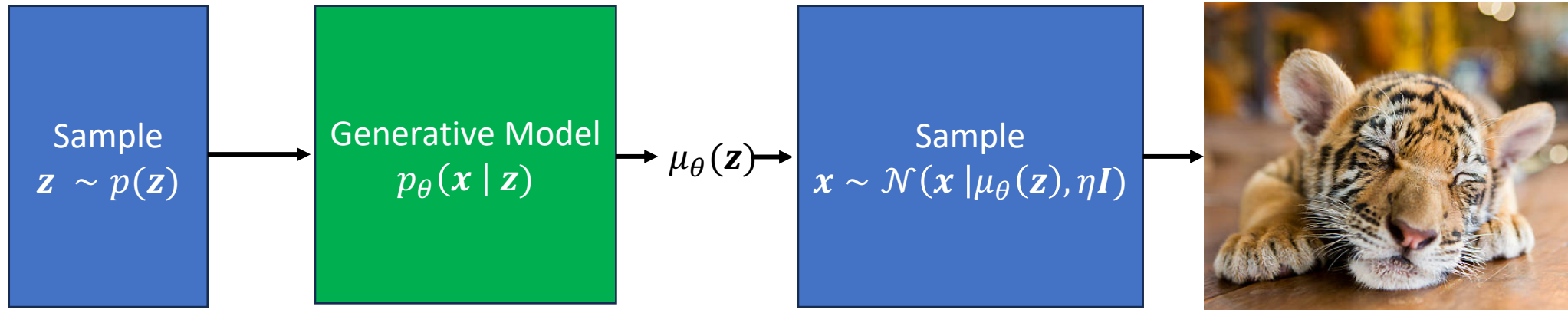


ELBO Objective

$$\mathbb{E}_{z \sim q_{\psi}(z|x)} [\log p_{\theta}(x | z) - KL(q_{\psi}(z | x) || p(z))]$$

Variational Autoencoders after Training

- Suppose we have learned VAE using the ELBO loss (details to follow).
- Then, as a generative model, we just sample $\mathbf{z} \sim p(\mathbf{z})$ and use the fixed generative model $p_{\theta}(\mathbf{x} | \mathbf{z})$



Computing the ELBO Loss

$$L_{\theta,\psi}(\mathbf{x}) := \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \log p_{\theta}(\mathbf{x} | \mathbf{z})}_{\text{Term 1 (Reconstruction Error)}} + \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \log \frac{p(\mathbf{z})}{q_{\psi}(\mathbf{z} | \mathbf{x})}}_{\text{Term 2 (KL Divergence)}}$$

- **Term 1 (Reconstruction Error)**

- Because $p_{\theta}(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mu_{\theta}(\mathbf{z}), \eta \mathbf{I})$, we have $\log p_{\theta}(\mathbf{x} | \mathbf{z}) = -\frac{1}{2\eta} \|\mathbf{x} - \mu_{\theta}(\mathbf{z})\|_2^2 + \text{const.}$
- We can approximate the expectation over $\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})$ by an average over $q_{\psi}(\mathbf{z} | \mathbf{x})$
- Recall that the expectation of a function $f(\mathbf{z})$ w.r.t. a random variable $\mathbf{z} \sim p$ can be approximated from i.i.d. samples $\mathbf{z}_j \sim p, j = 1, \dots, M$, via Monte Carlo averages

$$\mathbb{E}_{\mathbf{z} \sim p}[f(\mathbf{z})] = \int_{\mathbf{z}} f(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \approx \frac{1}{M} \sum_j f(\mathbf{z}_j)$$

- Applying the above formula to M i.i.d. samples $\mathbf{z}_j \sim q_{\psi}(\mathbf{z} | \mathbf{x})$ we get **reconstruction error**

$$\mathbb{E}_{\mathbf{z} \sim q_{\psi}}[\log p_{\theta}(\mathbf{x} | \mathbf{z})] \approx \frac{1}{M} \sum_j \log p_{\theta}(\mathbf{x} | \mathbf{z}_j) = \frac{1}{2\eta M} \sum_j \|\mathbf{x} - \mu_{\theta}(\mathbf{z}_j)\|_2^2$$

Computing the ELBO Loss

$$L_{\theta, \psi}(\mathbf{x}) := \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \log p_{\theta}(\mathbf{x} | \mathbf{z})}_{\text{Term 1}} + \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \log \frac{p(\mathbf{z})}{q_{\psi}(\mathbf{z} | \mathbf{x})}}_{\text{Term 2}}$$

- **Term 2 (Regularization to Prior)**

- Since $\mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \log \frac{p(\mathbf{z})}{q_{\psi}(\mathbf{z} | \mathbf{x})} = -KL(q_{\psi}(\mathbf{z} | \mathbf{x}) || p(\mathbf{z}))$, $q_{\psi}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mu_{\psi}(\mathbf{x}), \sigma_{\psi}^2(\mathbf{x})\mathbf{I})$, $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, the second term is a KL divergence between two d -dimensional Gaussians, which has a closed form solution

$$KL(\mathcal{N}(\mu_1, \sigma_1^2 \mathbf{I}) || \mathcal{N}(\mu_2, \sigma_2^2 \mathbf{I})) = \log \left(\frac{\sigma_2}{\sigma_1} \right) - \frac{d}{2} + \frac{d\sigma_1^2 + ||\mu_1 - \mu_2||_2^2}{2\sigma_2^2}$$

- Applying the above formula to our VAE model yields,

$$KL(q_{\psi}(\mathbf{z} | \mathbf{x}) || p(\mathbf{z})) = -\log(\sigma_{\psi}(\mathbf{x})) + \frac{d\sigma_{\psi}^2(\mathbf{x}) + ||\mu_{\psi}(\mathbf{x})||_2^2}{2} + \text{constant}$$

- Note this term does not require any sampling w.r.t. \mathbf{z} because the expectation is computed in closed form thanks for the KL formula for Gaussians.

Maximizing ELBO: How to Optimize?

- **ELBO objective:** We want to solve the following optimization problem:

$$\max_{\theta, \psi} \sum_i L_{\theta, \psi}(\mathbf{x}_i) = \sum_i \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x}_i)} \log \frac{p_{\theta}(\mathbf{x}_i, \mathbf{z})}{q_{\psi}(\mathbf{z} | \mathbf{x}_i)}$$

- **Simple idea:** Just alternate gradient ascent wrt θ, ψ on objective function

$$\theta_{k+1} = \theta_k + \alpha \frac{\delta L}{\delta \theta}(\theta_k, \psi_k)$$

$$\psi_{k+1} = \psi_k + \alpha \frac{\delta L}{\delta \psi}(\theta_k, \psi_k)$$

Stochastic Optimization of ELBO wrt θ

- **Issue:** Computing ∇_{θ} (ELBO) is not easy because of expectation w.r.t. latent z .
- **Solution:** Compute an unbiased estimator

$$\begin{aligned}\nabla_{\theta} L_{\theta, \psi}(\mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} [\log(p_{\theta}(\mathbf{x} | \mathbf{z}))] + \nabla_{\theta} \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{p(\mathbf{z})}{q_{\psi}(\mathbf{z} | \mathbf{x})} \right) \right] \\ &= \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x} | \mathbf{z})] \\ &= -\frac{1}{2\eta} \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \nabla_{\theta} \|\mathbf{x} - \mu_{\theta}(\mathbf{z})\|_2^2\end{aligned}$$

- We can take sample averages to compute an unbiased estimator
 - For each datapoint \mathbf{x} , compute $q_{\psi}(\mathbf{z} | \mathbf{x})$ through encoder, which gives $\mu_{\psi}(\mathbf{x}), \sigma_{\psi}^2(\mathbf{x})$
 - Sample $\mathbf{z}_j \sim \mathcal{N}(\mu_{\psi}(\mathbf{x}), \sigma_{\psi}^2(\mathbf{x})\mathbf{I})$, find $\mu_{\theta}(\mathbf{z}_j)$ through decoder
 - Estimate gradient from M samples: $\nabla_{\theta} L_{\theta, \psi}(\mathbf{x}) \approx -\frac{1}{2\eta M} \sum_j \nabla_{\theta} \|\mathbf{x} - \mu_{\theta}(\mathbf{z}_j)\|_2^2$

Stochastic Optimization of ELBO wrt ψ

- Here, we cannot just switch gradient and expectation because both are wrt to ψ

$$\nabla_{\psi} \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\psi}(\mathbf{z} | \mathbf{x})] \neq \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x})} \nabla_{\psi} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\psi}(\mathbf{z} | \mathbf{x})]$$

- Reparameterization trick**

- Because $q_{\psi}(\mathbf{z} | \mathbf{x}) = N(\mathbf{z}; \mu_{\psi}(\mathbf{x}), \sigma_{\psi}(\mathbf{x})I)$, we can rewrite samples $\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})$ as

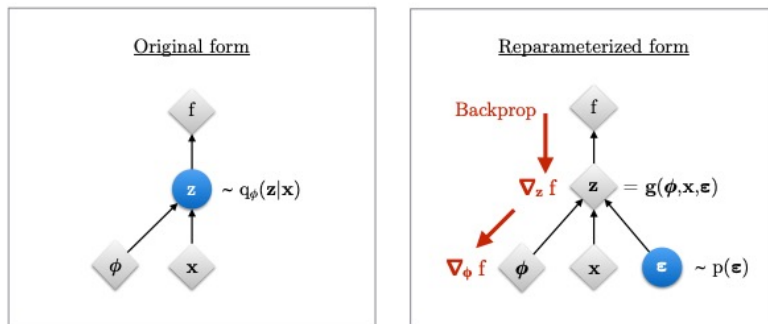
$$\mathbf{z}_{\psi} = g(\epsilon, \psi, \mathbf{x}) = \mu_{\psi}(\mathbf{x}) + \sigma_{\psi}(\mathbf{x})\epsilon \quad \text{for } \epsilon \sim N(\mathbf{0}, I)$$

- With this change of variables, we can rewrite the gradient of the loss as:

$$\nabla_{\psi} L_{\theta, \psi}(\mathbf{x}) = \nabla_{\psi} \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\psi}(\mathbf{z} | \mathbf{x})]$$

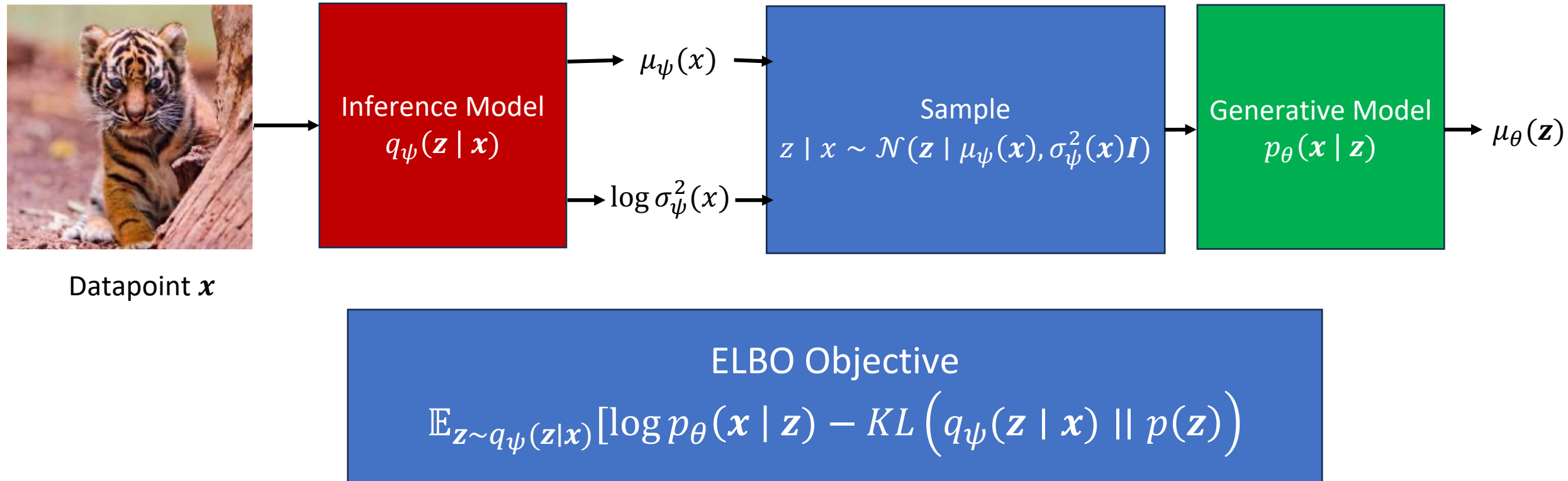
$$= \nabla_{\psi} \mathbb{E}_{\epsilon \sim N(\mathbf{0}, I)} [\log p_{\theta}(\mathbf{x}, \mathbf{z}_{\psi}) - \log q_{\psi}(\mathbf{z}_{\psi} | \mathbf{x})]$$

Gradient and expectation can now be switched! So as before, we compute an unbiased estimator by sampling many ϵ



Putting it all together

- Variational Autoencoder
 - We modelled inference and generative model as deep networks
 - We interpreted ELBO as an expected reconstruction error plus a KL-regularization to prior
 - Then, we rewrote the sampling in the latent space using the reparameterization trick
 - Finally, we derived stochastic gradient estimates to optimize the ELBO and learn a VAE



VAE's in Action

- $q_{\psi}(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\mathbf{z} \mid \mu_{\psi}(\mathbf{x}), \sigma_{\psi}^2(\mathbf{x}))$
- $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$
- $p_{\theta}(\mathbf{x} \mid \mathbf{z}) = \textit{Categorical}(\mathbf{x} \mid \pi_{\theta}(\mathbf{z}))$ Note this is different from the model considered up till now!
- Encoder network: $\mathbf{x} \in \mathbb{R}^D \rightarrow \text{Linear}(D, 256) \rightarrow \text{LeakyReLU} \rightarrow \text{Linear}(256, 2d) \rightarrow$
split into $\mu \in \mathbb{R}^d, \log \sigma^2 \in \mathbb{R}^d$
- Decoder network: $\mathbf{z} \in \mathbb{R}^d \rightarrow \text{Linear}(d, 256) \rightarrow \text{LeakyReLU} \rightarrow \text{Linear}(256, D) \rightarrow$
softmax

VAE's for Generation of MNIST Digits

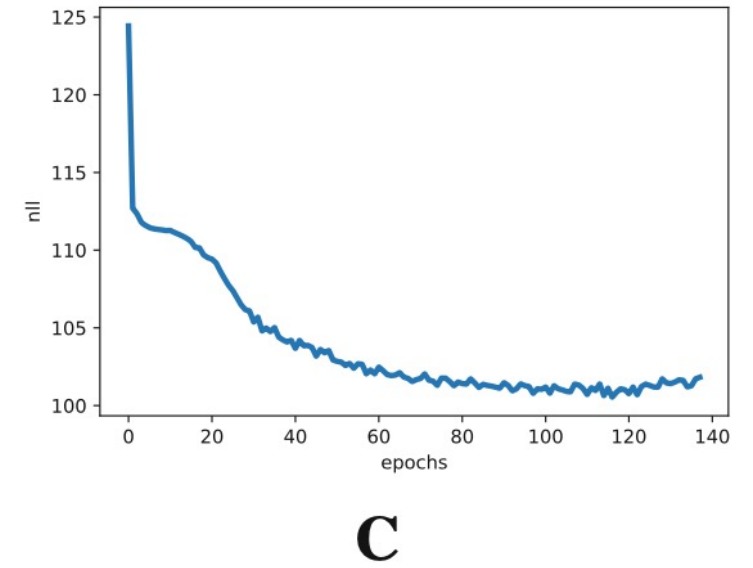


Fig. 4.4 An example of outcomes after the training: (a) Randomly selected real images. (b) Unconditional generations from the VAE. (c) The validation curve during training

VAEs: extensions

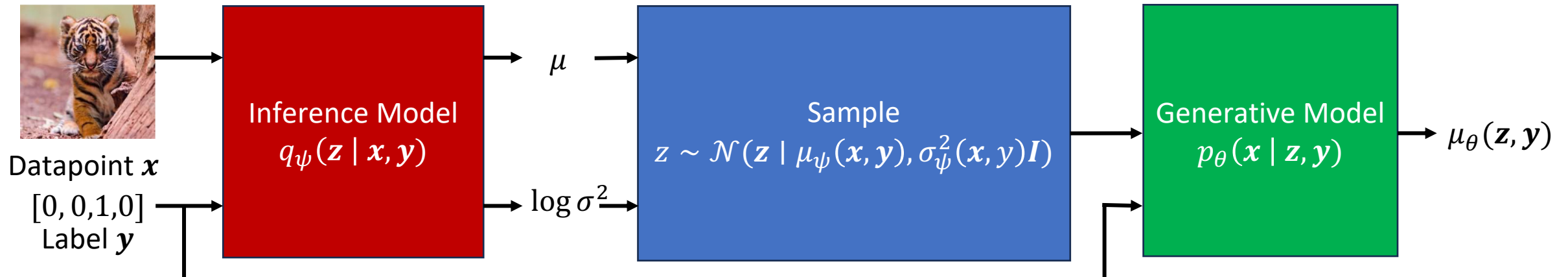
- So far, we have seen how to model $p(\mathbf{x})$ via VAEs.
- Many machine learning applications requires beyond learning $p(\mathbf{x})$:
 - Conditional generation, i.e., sampling from $p(\mathbf{x} \mid \mathbf{y})$ with a given class label \mathbf{y}
 - Classification, i.e., sampling from $p(\mathbf{y} \mid \mathbf{x})$ with a given sample \mathbf{x}
- Answer: We can extend the VAE paradigm to model $p(\mathbf{x}, \mathbf{y})$

Joint VAE

- Regular VAE: start from $\log p_\theta(x)$ and derive its ELBO $\int q(z | x) \log \frac{p_\theta(x, z)}{q(z | x)} dz$
- Here, since we want to model $p_\theta(x, y)$, let us derive ELBO for $\log p_\theta(x, y)$

$$\begin{aligned} & \max_{\theta} \mathbb{E}_{x, y \sim p_{\text{data}}} \log p_\theta(x, \mathbf{y}) \\ &= \max_{\theta} \max_{\psi} \mathbb{E}_{x, y \sim p_{\text{data}}} \left[\mathbb{E}_{q_\psi(z | x, \mathbf{y})} \log p_\theta(x | z, \mathbf{y}) - KL[q_\psi(z | x, \mathbf{y}) || p_\theta(z | \mathbf{y})] + \log p_\theta(\mathbf{y}) \right] \end{aligned}$$

- Details in the next slide. Spoiler alert: the ELBO derivation is almost the same as before
- Architecture:



Joint VAE: ELBO

- Let $q_\psi(z|x, \mathbf{y})$ be the variational distribution. Observe that

$$\begin{aligned}\log p_\theta(x, \mathbf{y}) &= \int q_\psi(\mathbf{z} | x, \mathbf{y}) \log p_\theta(x, \mathbf{y}) d\mathbf{z} = \int q_\psi(\mathbf{z} | x, \mathbf{y}) \log \frac{p_\theta(x, \mathbf{y}, \mathbf{z})}{p_\theta(\mathbf{z} | x, \mathbf{y})} d\mathbf{z} \\ &= \int q_\psi(\mathbf{z} | x, \mathbf{y}) \log \frac{p_\theta(x, \mathbf{y}, \mathbf{z})}{q_\psi(\mathbf{z} | x, \mathbf{y})} \frac{q_\psi(\mathbf{z} | x, \mathbf{y})}{p_\theta(\mathbf{z} | x, \mathbf{y})} d\mathbf{z} \\ &= \boxed{\int q_\psi(\mathbf{z} | x, \mathbf{y}) \log \frac{p_\theta(x, \mathbf{y}, \mathbf{z})}{q_\psi(\mathbf{z} | x, \mathbf{y})} d\mathbf{z}} + \boxed{\int q_\psi(\mathbf{z} | x, \mathbf{y}) \log \frac{q_\psi(\mathbf{z} | x, \mathbf{y})}{p_\theta(\mathbf{z} | x, \mathbf{y})} d\mathbf{z}} \\ &\quad \text{Evidence Lower Bound (ELBO)} \quad \text{KL}[q(\mathbf{z} | x, \mathbf{y}) || p_\theta(\mathbf{z} | x, \mathbf{y})]\end{aligned}$$

$$\begin{aligned}\text{Therefore, } \log p_\theta(x | \mathbf{y}) &= \max_{\psi} \mathbb{E}_{q_\psi(\mathbf{z} | x, \mathbf{y})} \log \frac{p_\theta(x, \mathbf{y}, \mathbf{z})}{q_\psi(\mathbf{z} | x, \mathbf{y})} \\ &= \max_{\psi} \mathbb{E}_{q_\psi(\mathbf{z} | x, \mathbf{y})} \left[\log \frac{p_\theta(\mathbf{z} | \mathbf{y})}{q_\psi(\mathbf{z} | x, \mathbf{y})} + \log p_\theta(x | \mathbf{y}, \mathbf{z}) + \log p_\theta(\mathbf{y}) \right] \\ &= \max_{\psi} \mathbb{E}_{q_\psi(\mathbf{z} | x, \mathbf{y})} \log p_\theta(x | \mathbf{y}, \mathbf{z}) - \text{KL}[q_\psi(\mathbf{z} | x, \mathbf{y}), p_\theta(\mathbf{z} | \mathbf{y})] + \log p_\theta(\mathbf{y})\end{aligned}$$

Joint VAE: Training objective

- The training objective is given by

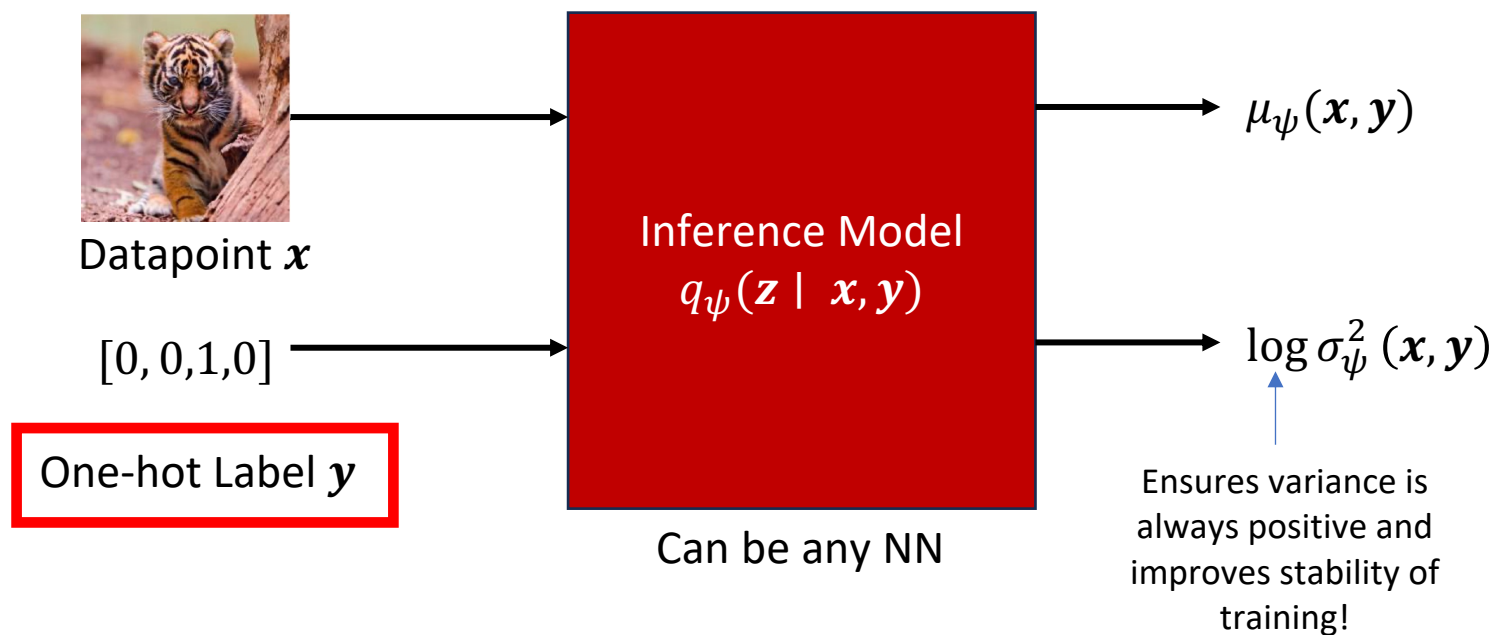
$$\begin{aligned} \max_{\theta} J_{\text{cond}} &:= \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \log p_{\theta}(\mathbf{x}, \mathbf{y}) \\ &= \max_{\theta} \max_{\psi} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \left[\mathbb{E}_{q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y})} \log p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y}) - KL[q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y}) || p_{\theta}(\mathbf{z} | \mathbf{y})] + \log p_{\theta}(\mathbf{y}) \right] \end{aligned}$$

- What are $q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y})$, $p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y})$, $p_{\theta}(\mathbf{z} | \mathbf{y})$?
 - They are similar to what we needed in regular VAEs, with an additional input \mathbf{y}
- What is $p_{\theta}(\mathbf{y})$?
 - This is a new term denoting the prior probability on class labels
- Answer: next two slides

Joint VAE: Setup

- We have three models we need to define for VAE model

1. **Inference model $q_{\psi}(\mathbf{z} \mid \mathbf{x}, \mathbf{y})$:** We will define as $q_{\psi}(\mathbf{z} \mid \mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{z}; \mu_{\psi}(\mathbf{x}, \mathbf{y}), \sigma_{\psi}^2(\mathbf{x}, \mathbf{y})\mathbf{I})$, i.e., a normal distribution with learned mean and covariance



2. **Latent prior:** $p(\mathbf{z} \mid \mathbf{y}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$

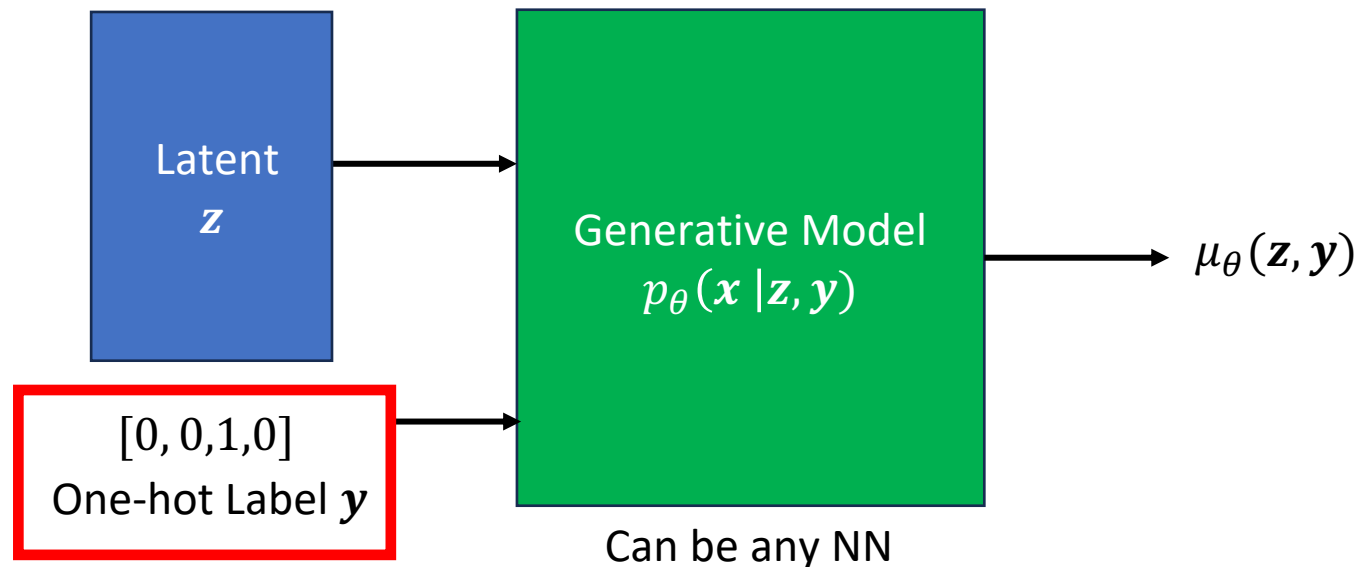
3. **Class prior:** $p_{\theta}(\mathbf{y}) = \text{Categorical}(\pi)$

Joint VAE: Setup

- We have three models we need to define for VAE model

4. **Generative model $p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y})$:** We will define as

- $p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y}) = \mathcal{N}(\mathbf{z}; \mu_{\theta}(\mathbf{z}, \mathbf{y}), \eta^2 \mathbf{I})$, i.e., a normal distribution with learned mean and variance
- $p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y}) = \text{Cat}(\mathbf{z}; \pi_{\theta}(\mathbf{z}, \mathbf{y}))$, i.e., a categorical distribution with learned class probabilities

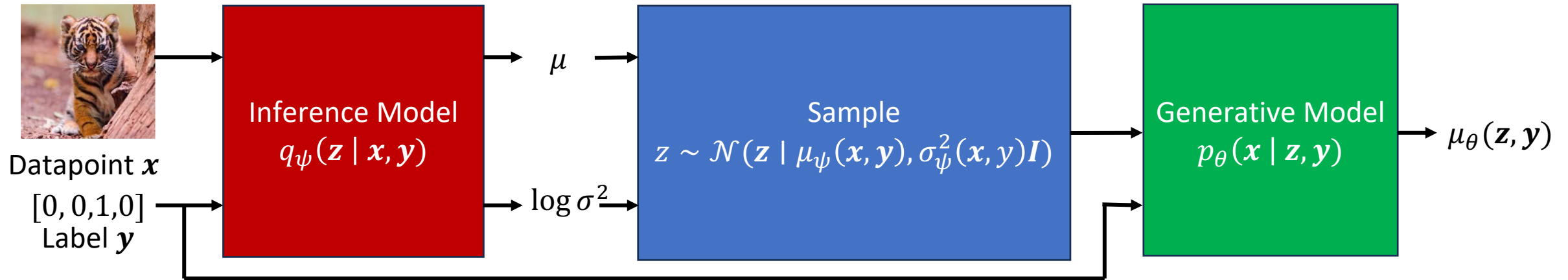


- Note this can be defined in many different ways, yielding different models (such as a categorical distribution over 255 values of each pixel)

Joint VAE: Training objective

- So far, the training objective is given by

$$\begin{aligned} \max_{\theta} J_{\text{cond}} &:= \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \log p_{\theta}(\mathbf{x}, \mathbf{y}) \\ &= \max_{\theta} \max_{\psi} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \left[\mathbb{E}_{q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y})} \log p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y}) - KL[q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y}) || p_{\theta}(\mathbf{z} | \mathbf{y})] + \log p_{\theta}(\mathbf{y}) \right] \end{aligned}$$



- This allows us to do conditional generation
 - Specify a class \mathbf{y} that you want to sample from
 - Sample from $p_{\theta}(\mathbf{z} | \mathbf{y})$, which is a Gaussian $\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$
 - Sample from the conditional generation $p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y})$

Joint VAE: Training objective

- So far, the training objective is given by

$$\begin{aligned} & \max_{\theta} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \log p_{\theta}(\mathbf{x}, \mathbf{y}) \\ &= \max_{\theta} \max_{\psi} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \left[\mathbb{E}_{q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y})} \log p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y}) - KL[q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y}) || p_{\theta}(\mathbf{z} | \mathbf{y})] + \log p_{\theta}(\mathbf{y}) \right] =: J_{\text{cond}} \end{aligned}$$

- Hold on... Classification, i.e., the distribution of \mathbf{y} given \mathbf{x} is not taken care of
- How can we also model classification?
 - Again, let us start with MLE principle: $\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \log q_{\psi}(\mathbf{x}, \mathbf{y})$
 - Define $q_{\psi}(\mathbf{x}, \mathbf{y}) = q_{\psi}(\mathbf{y} | \mathbf{x}) p_{\text{data}}(\mathbf{x})$, where $q_{\psi}(\mathbf{y} | \mathbf{x})$ is a classifier
 - $\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \log q_{\psi}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \log q_{\psi}(\mathbf{y} | \mathbf{x}) + \text{const. not depending on } \psi$
 - This gives the classification loss $J_{\text{cls}} := \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \log q_{\psi}(\mathbf{y} | \mathbf{x})$
- The final objective is given by $J_{\text{cond}} + \lambda J_{\text{cls}}$
 - They originate from the MLE principles, since it is the sum of two log-likelihoods

Joint VAE: summary

$$\max_{\theta} \max_{\psi} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} \left[\mathbb{E}_{q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y})} \log p_{\theta}(\mathbf{x} | \mathbf{z}, \mathbf{y}) - KL[q_{\psi}(\mathbf{z} | \mathbf{x}, \mathbf{y}) || p_{\theta}(\mathbf{z} | \mathbf{y})] + \log p_{\theta}(\mathbf{y}) + \lambda \log q_{\psi}(\mathbf{y} | \mathbf{x}) \right]$$

