

# Variational Inference and Variational Autoencoders

- variational inference
- variational autoencoders

# Unsupervised Learning

- *supervised* machine learning is about **learning functions from data**

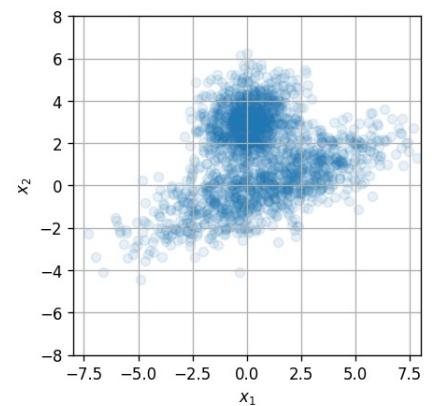
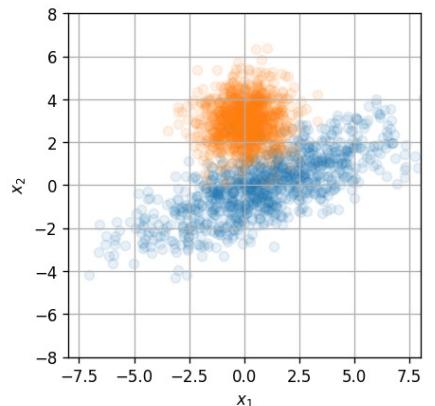
$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

- *unsupervised* machine learning is about finding interesting patterns in data

$$\mathcal{D} = \{(x_i)\}_{i=1}^n$$

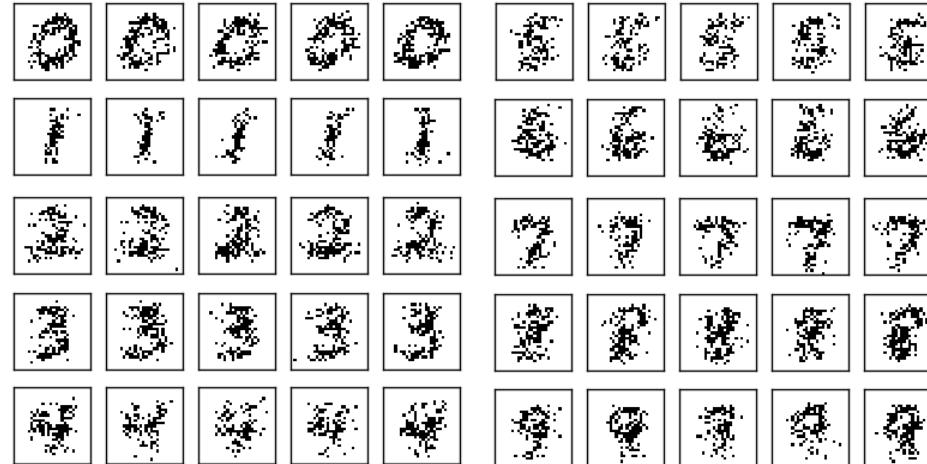
- *knowlege discovery*
- under-specified

- learn  $p(\mathbf{x})$  instead of  $p(\mathbf{x}, y)$  or  $p(y|\mathbf{x})$
- learn  $p_\theta(\mathbf{x})$



# Learning to Draw MNIST

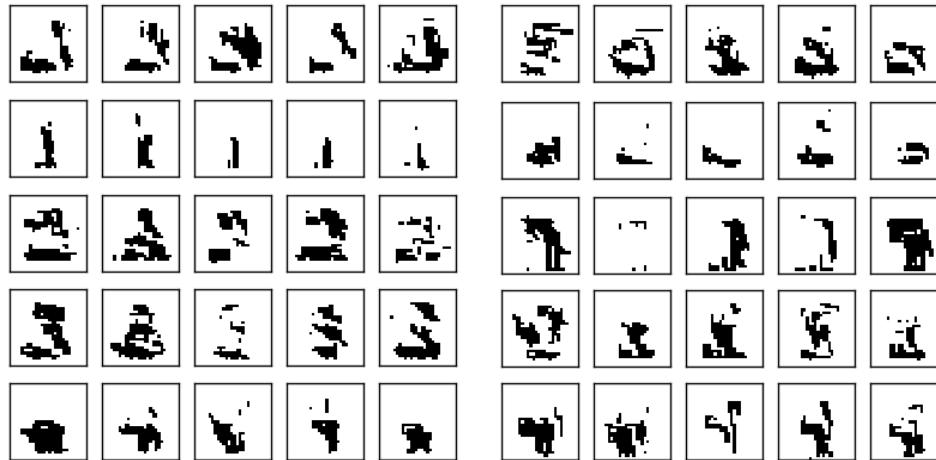
- Naive Bayes



- sample with 784 independent coin flips with different bias

# Learning to Draw MNIST

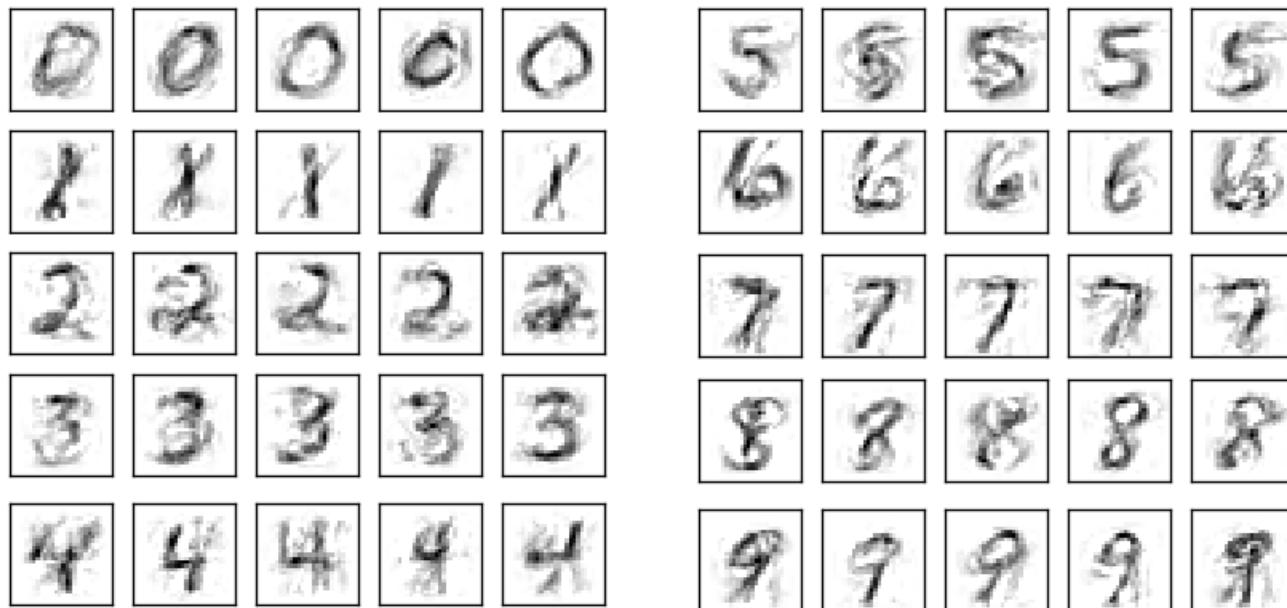
- Bayes Net / Belief Network



- ancestral sampling

# Learning to Draw MNIST

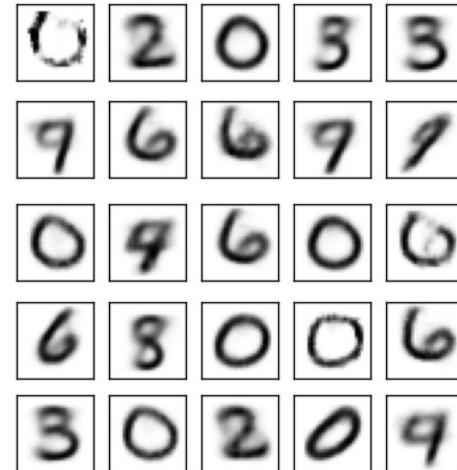
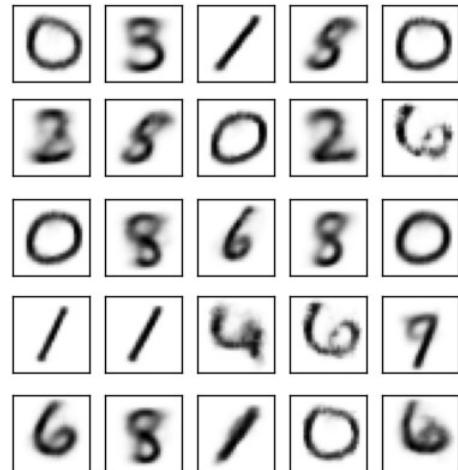
- multivariate normal



- sample  $\mathcal{N}(0, \mathbf{I})$ , compute  $\mathbf{Ax} + \boldsymbol{\mu}$

# Learning to Draw MNIST

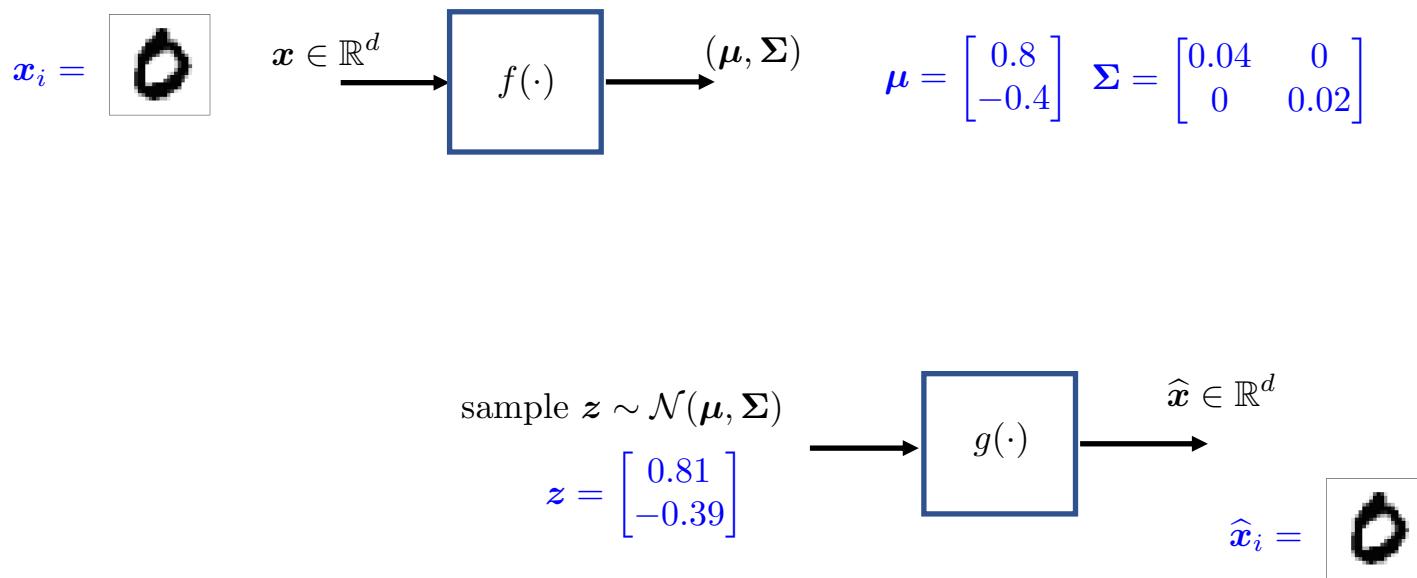
- variational autoencoder



- sample  $z \sim \mathcal{N}(0, I)$ , pass sample through function  $g(\cdot)$

# Variational Autoencoders (at test time)

- how can we sample from a VAE?



- how can we learn  $f(\cdot)$  and  $g(\cdot)$  from data?

- ensure that  $x \approx \hat{x}$
- try to force  $p(z) \sim \mathcal{N}(\mathbf{0}, I)$

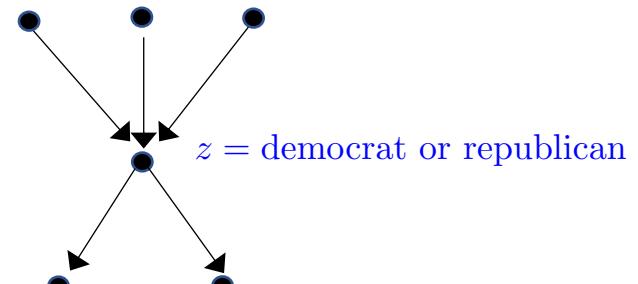
# Latent Variable Models

- goal: learn  $p(\mathbf{x})$  from  $\mathcal{D} = \{(\mathbf{x}_i)\}_{i=1}^n$

$$p(\mathbf{x}) = \int p(\mathbf{x}, z) dz$$

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

M/F, age, income, etc.

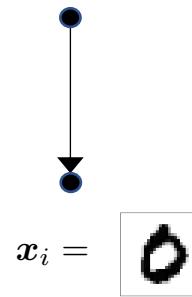


how they vote in two elections

- $z$  is a *latent* variable
- $z$  might be meaningful, or maybe just convenient
- MNIST:  $z$  = digit, pen width, slant, etc.
  - not principled/helpful to specify  $z$  manually
- if  $z$  is discrete, then *mixture* model

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

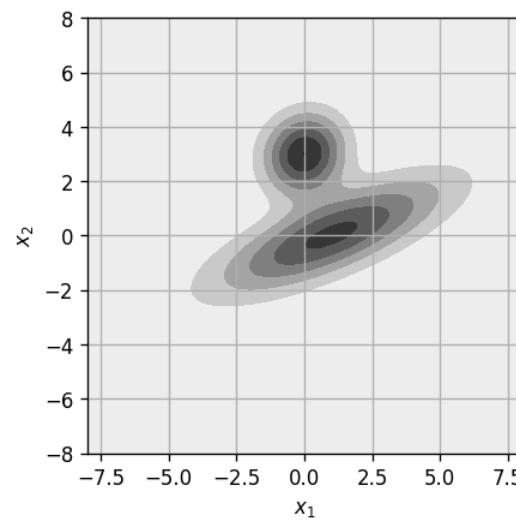
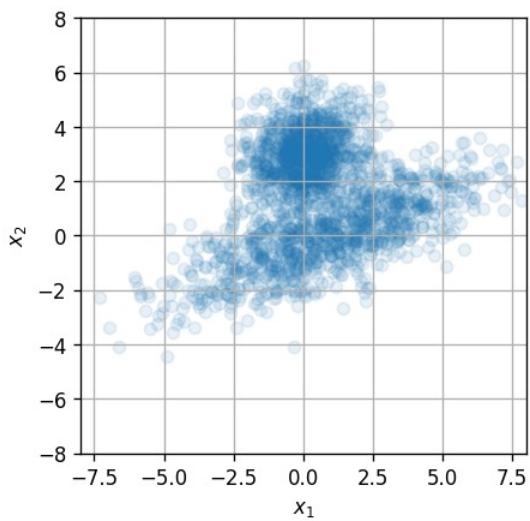
$z$  = intended digit, slant, width, etc.



$x_i =$



# Latent Variable Model Example - Gaussian Mixture Models (GMM)



$z$

$x$

$$p(\mathbf{x}) = 0.5 \mathcal{N} \left( \mathbf{x} \left| \begin{bmatrix} 0 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 & 0.1 \\ 0.1 & 1 \end{bmatrix} \right. \right) + 0.5 \mathcal{N} \left( \mathbf{x} \left| \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 8 & 3 \\ 3 & 2 \end{bmatrix} \right. \right)$$

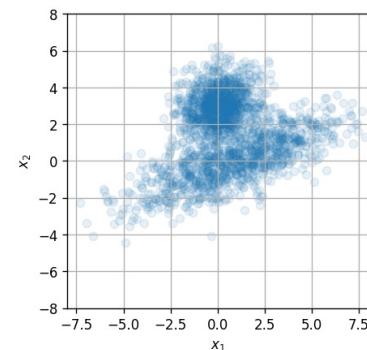
- Gaussian mixture model:

$$p(\mathbf{x}) = \sum_{k=1}^K p(z=k) \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

# Latent Variable Models

- goal: learn relationship between  $\mathbf{x}$  and  $\mathbf{z}$

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



- learning  $p(\mathbf{z}|\mathbf{x})$  without more assumptions is ill-posed

- need to specify the form of  $p(\mathbf{z})$
- need to specify the form of  $p(\mathbf{x}|\mathbf{z})$
- example: Gaussian mixtures and EM

$$p(\mathbf{z}|\mathbf{x}) = r_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \mu_k, \Sigma_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_i | \mu_{k'}, \Sigma_{k'})}$$

- starting point  $p(\mathbf{z}) \sim K$  categories  
and  $p(\mathbf{x}|\mathbf{z})$  Gaussian

# Variational Inference

- goal: learn  $p(\mathbf{z}|\mathbf{x})$

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



- $p(\mathbf{z}|\mathbf{x})$  is the posterior over the hidden variable  $\mathbf{z}$

- variational inference: estimate  $p(\mathbf{z}|\mathbf{x})$  using some other distribution  $q(\mathbf{z}|\mathbf{x})$

- restrict  $q(\mathbf{z}|\mathbf{x})$  to be from some tractable family, i.e,  $\mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\Sigma}(\mathbf{x}))$

- consider the KL divergence between  $q(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{z}|\mathbf{x})$ :

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

# Variational Inference

- consider the KL divergence between  $q(\mathbf{z}|\mathbf{x})$  and  $p(\mathbf{z}|\mathbf{x})$ :

$$\begin{aligned}
 0 \leq D(q(\cdot|x)||p(\cdot|x)) &= \int q(z|x) \log \frac{q(z|x)}{p(z|x)} dz, \quad \text{with equality if and only if } q(\cdot|x) = p(\cdot|x), \\
 &= \int q(z|x) \log \frac{q(z|x)p(x)}{p(x|z)p(z)} dz \\
 &= \int q(z|x) \log p(x) dz + \int q(z|x) \log \frac{1}{p(x|z)} dz + \int q(z|x) \log \frac{q(z|x)}{p(z)} dz \\
 &= \log p(x) - \mathbb{E}_{q(\cdot|x)}[\log p(x|Z)] - D(q(\cdot|x)||p(\cdot)).
 \end{aligned}$$



Therefore,

$$\log p(x) \geq E_{q(\cdot|x)}[\log p(x|Z)] + D(q(\cdot|x)||p(\cdot)).$$

Instead of choosing  $\mu(x)$  and  $\Sigma(x)$  to maximize  $\log p(x)$ , as in maximum likelihood, maximize the right-hand side and hope for the best!

# Variational Autoencoders

- how do we match data to a model? maximize the likelihood

$$\max \sum_i \log p(\mathbf{x}_i)$$



- we can try maximizing the lower bound

$$\log p(\mathbf{x}) \geq E_{z \sim q(z)} [\log p(\mathbf{x}|z)] - D(q(z)||p(z)) \quad (\text{called ELBO - evidence lower bound})$$

- how can we compute  $D(q(z)||p(z))$ ?

- we can choose  $p(z)$  to be  $\mathcal{N}(0, \mathbf{I})$  and  $q(z)$  to be  $\mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\Sigma}(\mathbf{x}))$

$$\begin{aligned} D(q(z)||p(z)) &= \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} \text{tr} (\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2} n \\ &= \frac{1}{2} \log \frac{1}{|\boldsymbol{\Sigma}_1(\mathbf{x})|} + \frac{1}{2} \text{tr} (\boldsymbol{\Sigma}_1(\mathbf{x})) + \boldsymbol{\mu}(\mathbf{x})^T \boldsymbol{\mu}(\mathbf{x}) - \frac{1}{2} n \end{aligned}$$

# Variational Autoencoders

- how can we compute/maximize  $E_{z \sim q(z)} [\log p(\mathbf{x}|\mathbf{z})]?$

- $g(\cdot)$  is a deterministic function

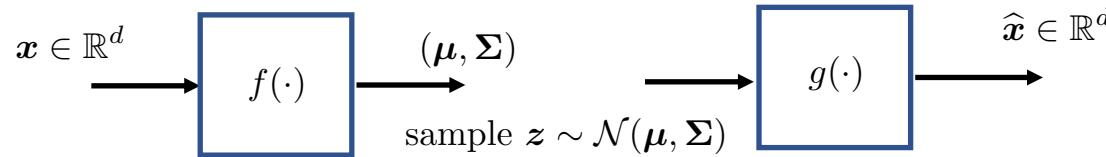
$$\log p(\mathbf{x}|\mathbf{z}) = \log p(\mathbf{x}|\hat{\mathbf{x}}) \quad \text{small caveat: if } g(\cdot) \text{ is an invertible mapping}$$



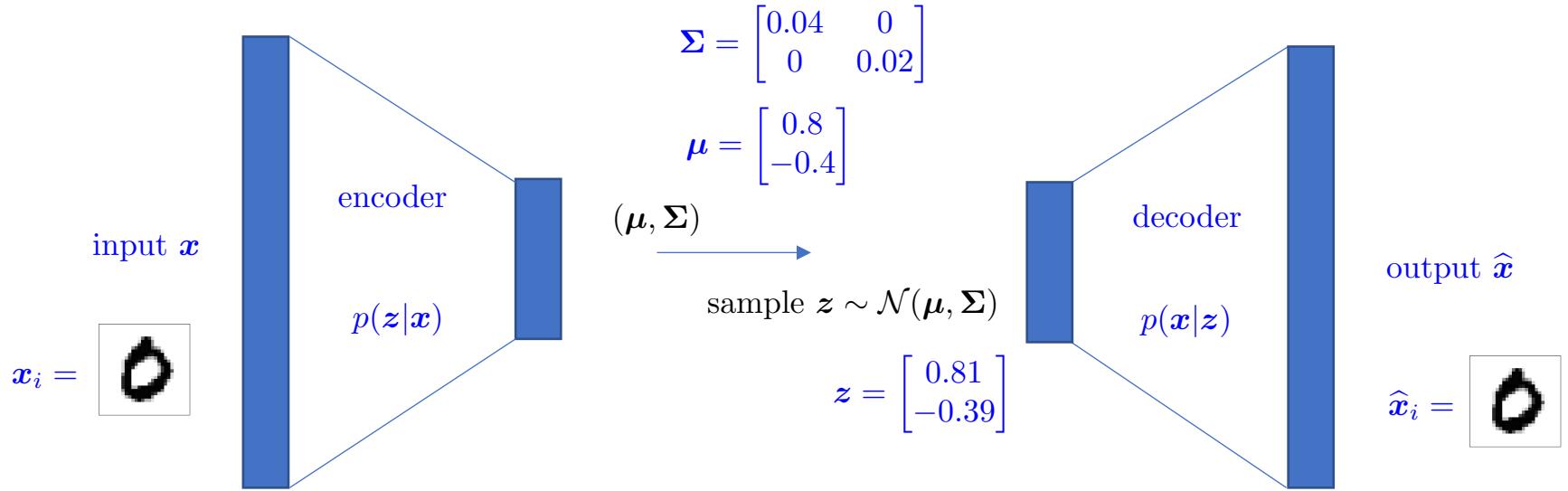
- big when  $\hat{\mathbf{x}} \approx \mathbf{x}$

- a squared error loss helps maximize  $p(\mathbf{x}|\mathbf{z})$

$$E_{z \sim q(z)} [\log p(\mathbf{x}|\mathbf{z})] \approx \frac{1}{n} \sum \log p(\mathbf{x}_i|\mathbf{z}_i) \approx \frac{1}{n} \sum \log p(\mathbf{x}_i|\hat{\mathbf{x}}_i)$$



# Variational Autoencoders



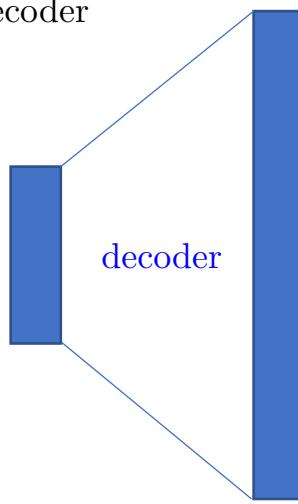
- ensure that  $\mathbf{x} \approx \hat{\mathbf{x}}$
- try to force  $p(\mathbf{z}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

# VAE's on MNIST

- how can we generate random MNIST images?

1. sample  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

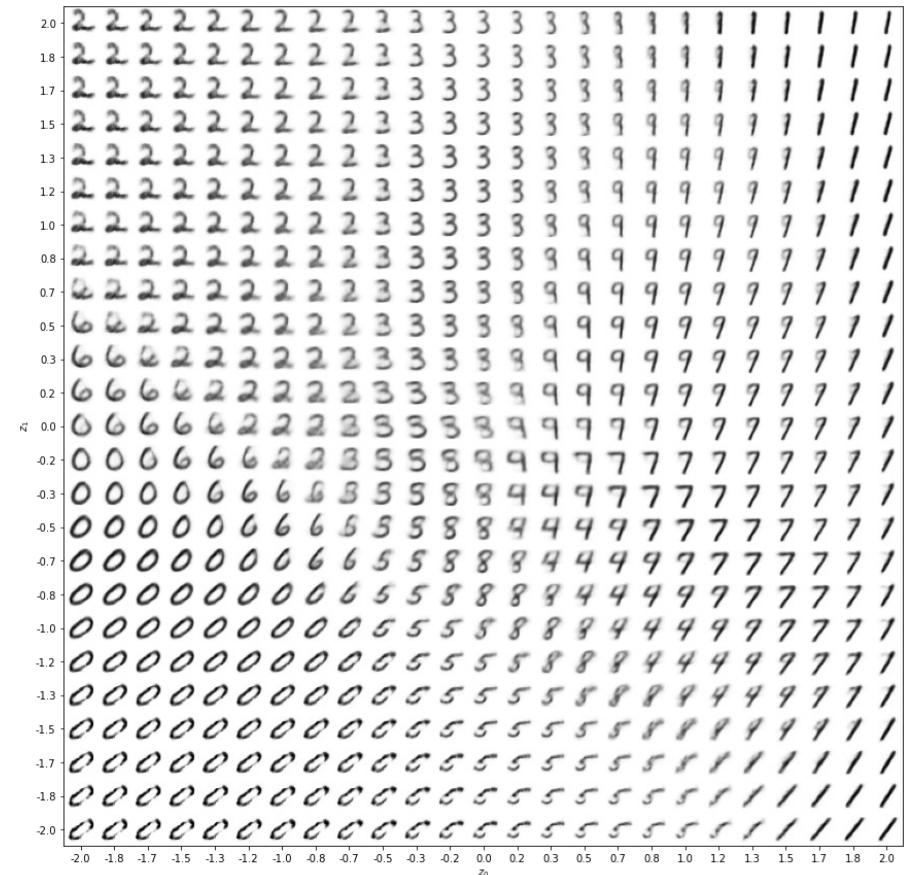
2. input sample into decoder side of neural network



- how can we creating this figure?

1. grid  $\mathbf{z} \in [-2, 2]^2$

2. input each grid point into decoder side of neural network



# VAE's on MNIST

- [Tutorial on Variational Autoencoders, Carl Doersch]
- [Variational Autoencoders, Ali Ghodsi]
- [Keras, <https://keras.io/examples/generative/vae/>]