

CS 4803 / 7643: Deep Learning

Topics:


- Variational Auto-Encoders (VAEs)
 - Reparameterization trick
- Generative Adversarial Networks (GANs)

Dhruv Batra
Georgia Tech

Administrativa

- Project submission instructions released
 - Due: 12/04, 11:55pm
 - Last deliverable in the class
 - Can't use late days
 - <https://piazza.com/class/jkujs03pgu75cd?cid=225>

Recap from last time



Variational Autoencoders (VAE)

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

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

$\sum p(x|z)p(z)$

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

Cannot optimize directly, derive and optimize lower bound on likelihood instead

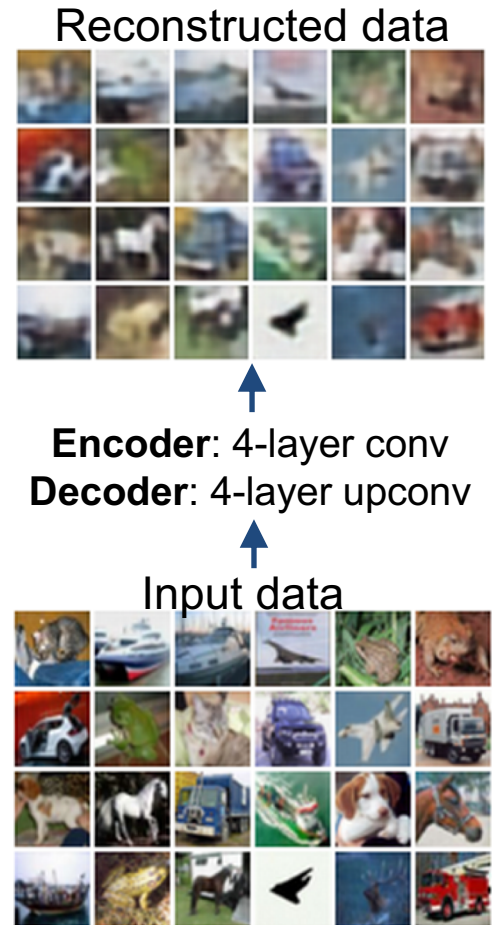
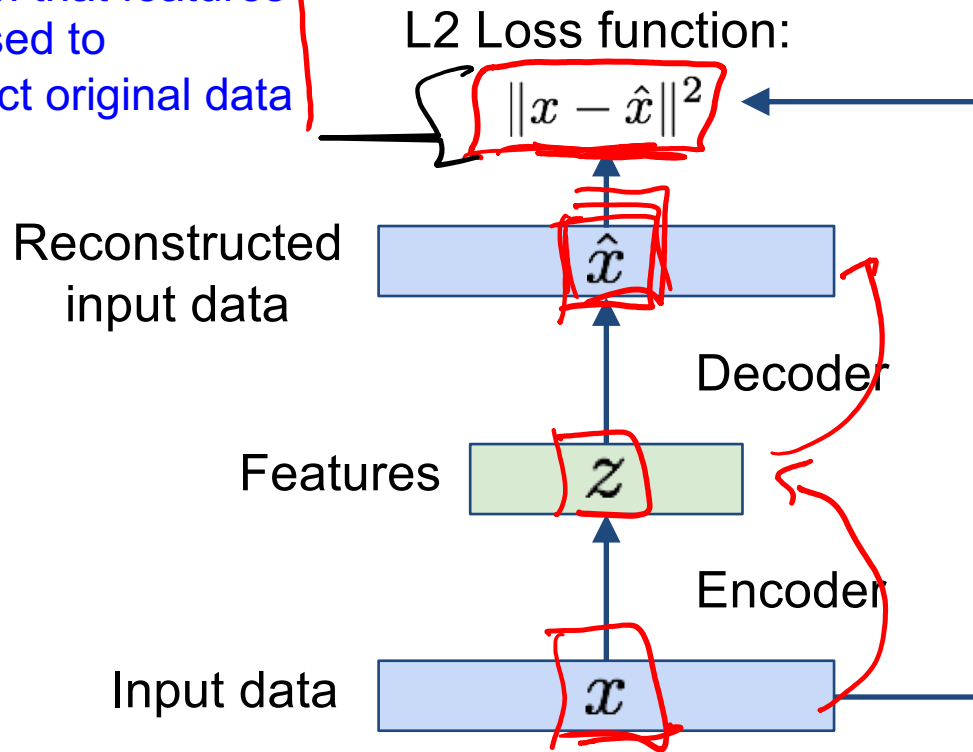
Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

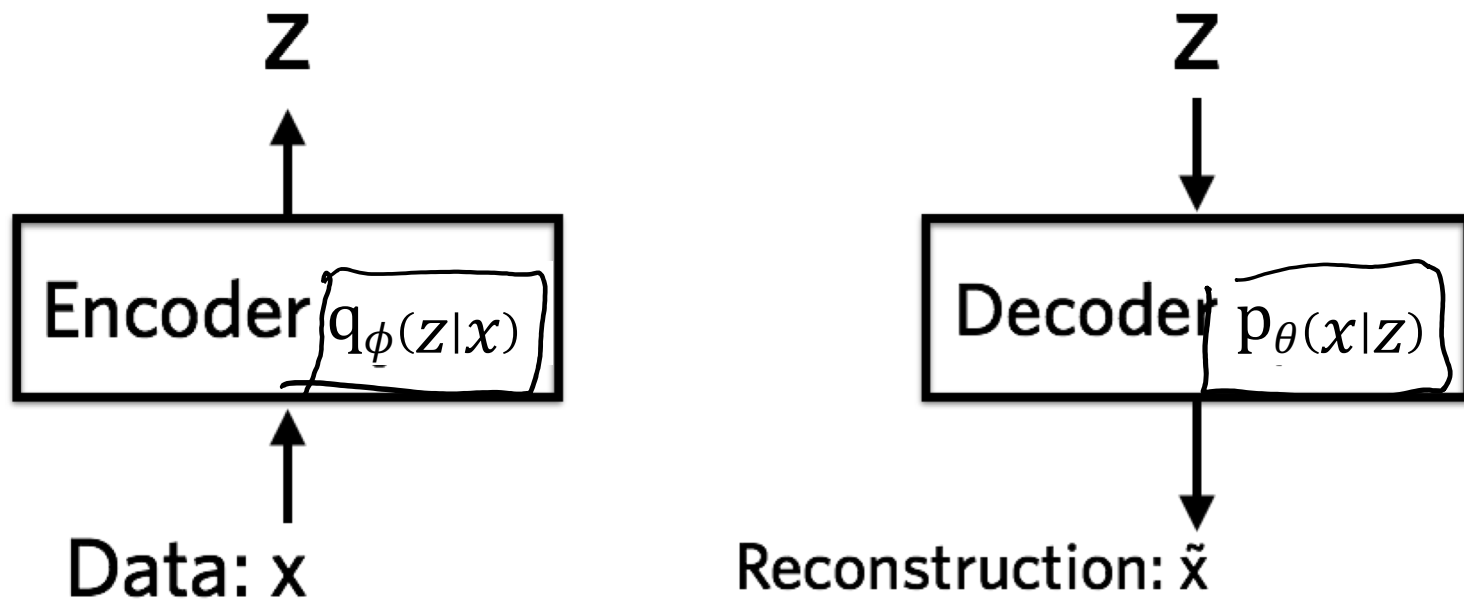
Autoencoders

Train such that features can be used to reconstruct original data



Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!



Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders

2. Variational Approximation

- Variational Lower Bound / ELBO

3. Amortized Inference Neural Networks

4. “Reparameterization” Trick

Key problem

$$\bullet \quad \boxed{P(z|x)} = \frac{P(z,x)}{P(x)} = \frac{P(x|z)P(z)}{\sum_z P(x|z)P(z)} \leftarrow$$

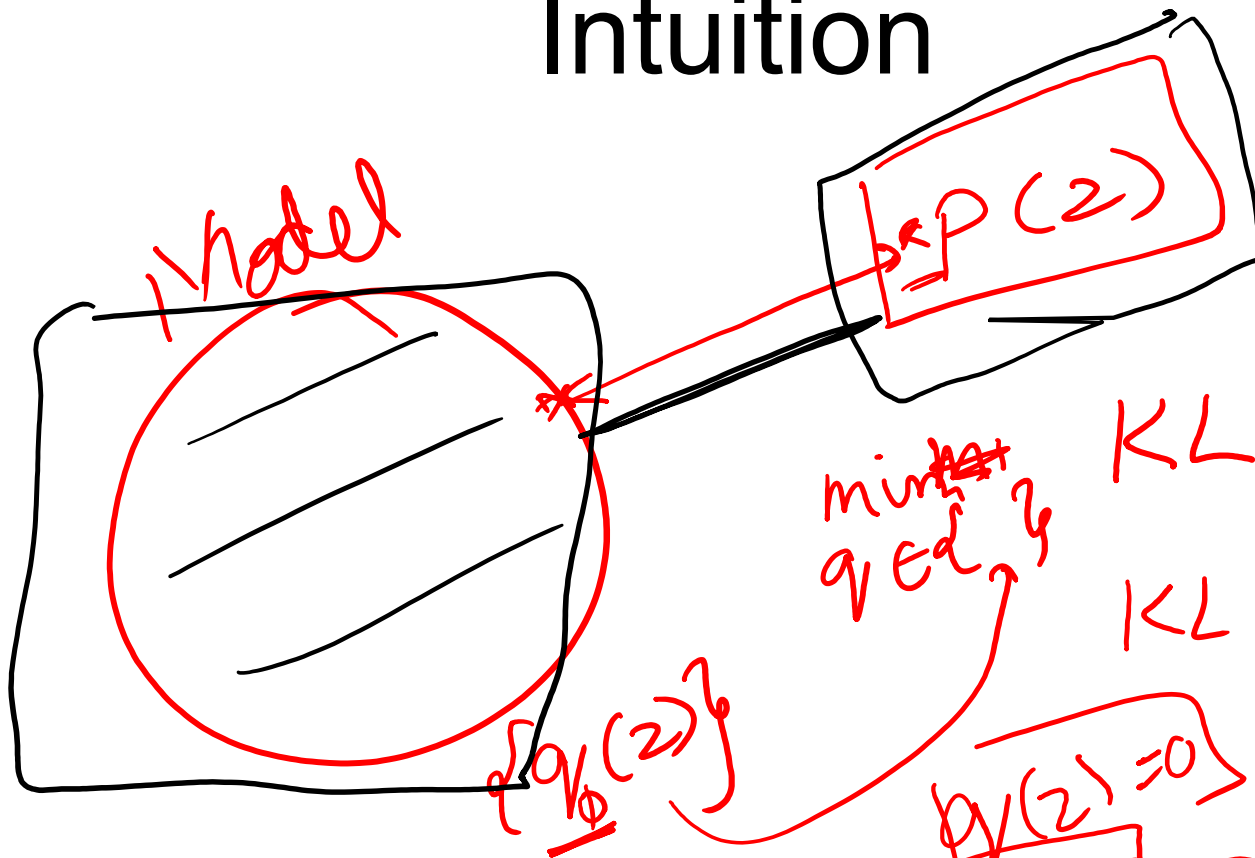
\downarrow

$q(z)$

What is Variational Inference?

- Key idea
 - Reality is complex
 - Can we approximate it with something “simple”?
 - Just need to make sure the simple thing is “close” to the complex thing.

Intuition



$$KL(P(z) \parallel q(z))$$

$$KL(q(z) \parallel P(z))$$

$$KL(P(z) \parallel q(z)) = \sum_z P(z) \log \frac{P(z)}{q(z)}$$

error

$$- \sum_{\text{sample}} \log q(z)$$

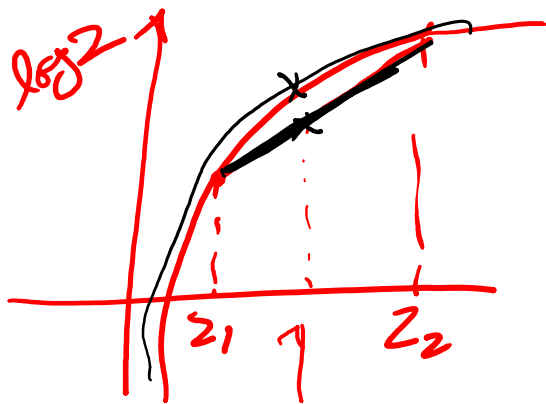
The general learning problem with missing data

- Marginal likelihood – \mathbf{x} is observed, \mathbf{z} is missing:

$$\begin{aligned} \boxed{ll(\theta : \mathcal{D})} &= \log \prod_{i=1}^N \underbrace{P(\bar{\mathbf{x}}_i | \theta)} \\ &= \sum_{i=1}^N \log P(\mathbf{x}_i | \theta) \quad \leftarrow P(\bar{\mathbf{x}}_i, \mathbf{z}) \\ &= \cancel{\sum_{i=1}^N} \log \sum_{\mathbf{z}} P(\mathbf{x}_i, \mathbf{z} | \theta) \\ &= \log \left[\sum_{\mathbf{z}} P(\bar{\mathbf{x}}_i | \theta) P(\mathbf{z} | \bar{\mathbf{x}}_i, \theta) \right] \\ &= \log \left[\underbrace{E_{P(\mathbf{z} | \bar{\mathbf{x}}_i, \theta)} \left[P(\bar{\mathbf{x}}_i | \theta) \right]} \right] \end{aligned}$$

Applying Jensen's inequality ←

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) g(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log g(\mathbf{z})$



$$\lambda z_1 + \lambda z_2$$

$$\lambda_1, \lambda_2 \geq 0$$

$$\lambda_1 + \lambda_2 = 1$$

$$f(\lambda_1 z_1 + \lambda_2 z_2) \geq \lambda_1 f(z_1) + \lambda_2 f(z_2)$$

$$f\left(\sum_{i=1}^K \lambda_i z_i\right) \geq \sum_{i=1}^K \lambda_i f(z_i)$$

$$f(E[z]) \geq E[f(z)]$$

$$f(E[g(z)]) \geq E[f(g(z))]$$

$z \rightarrow g(z)$

Applying Jensen's inequality

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) g(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log g(\mathbf{z})$

$$\rightarrow \log \sum_{\mathbf{z}} \underbrace{P(\vec{x}_i, \mathbf{z} | \theta)}_{Q_i(\mathbf{z})} \cdot \underbrace{Q_i(\mathbf{z})}_{Q_i(\mathbf{z})}$$

$$\geq \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\vec{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

$$F(\theta, Q_i)$$

"Free Energy"
Variational

VLB / ELBO ←

Applying Jensen's inequality

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) g(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log g(\mathbf{z})$

F

$$\underline{l(\theta : \mathcal{D})} = \sum_{i=1}^N \log \sum_{\mathbf{z}} Q_i(\mathbf{z}) \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

$P(\mathbf{z} | \theta) P(\bar{x}_i | z, \theta)$

$P(\bar{x}_i | \theta) P(\mathbf{z} | x_i, \theta)$

ELBO: Factorization #1

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q_i) = \sum_{i=1}^N \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

$\rightarrow P(\bar{x}_i | \theta) P(z | \bar{x}_i, \theta)$

$$\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \left(\frac{P(\bar{x}_i | \theta) P(z | \bar{x}_i, \theta)}{Q_i(\mathbf{z})} \right)$$

$$= \underbrace{\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log P(\bar{x}_i | \theta)}_1 + \left[\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(z | \bar{x}_i, \theta)}{Q_i(\mathbf{z})} \right]$$

$$F(\theta, Q_i) = \log P(\bar{x}_i | \theta) - \left[\text{KL} \left(\underbrace{Q_i(\mathbf{z})}_{\text{approx}} \parallel \underbrace{P(z | \bar{x}_i, \theta)}_{\text{target}} \right) \right]$$

$$\max_{\theta, Q_i} \boxed{F} = \underline{\text{ll}} - \underline{\text{KL}}$$

ELBO: Factorization #1

max
 θ, Q_i
 $F(\theta, Q_i)$

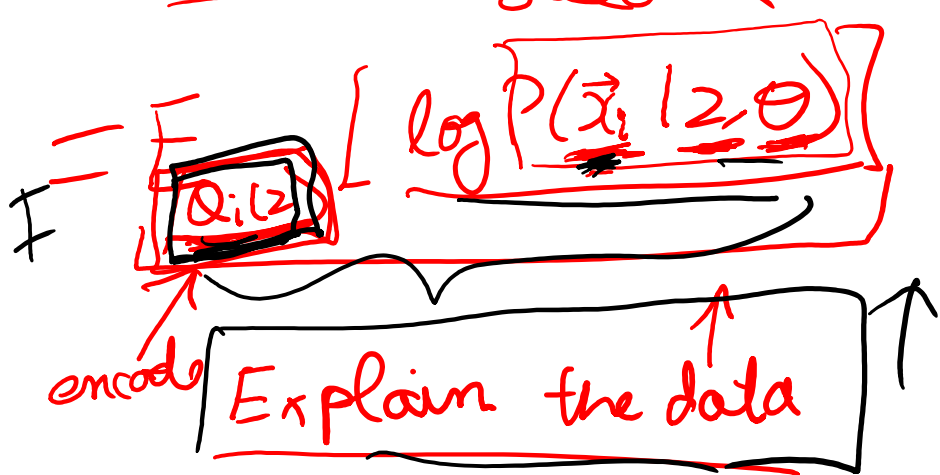
$$l(\theta : \mathcal{D}) \geq F(\theta, Q_i) = \sum_{i=1}^N \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

$\rightarrow P(\mathbf{z} | \theta) P(\mathbf{x}_i | \mathbf{z}, \theta)$

$$= \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{z} | \theta) P(\mathbf{x}_i | \mathbf{z}, \theta)}{Q_i(\mathbf{z})}$$

$$= \left[\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log P(\mathbf{x}_i | \mathbf{z}, \theta) \right] + \left[\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{z} | \theta)}{Q_i(\mathbf{z})} \right]$$

decoder



$$- \text{KL} [Q_i(\mathbf{z}) \parallel P(\mathbf{z} | \theta)]$$

Regularizer

Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders

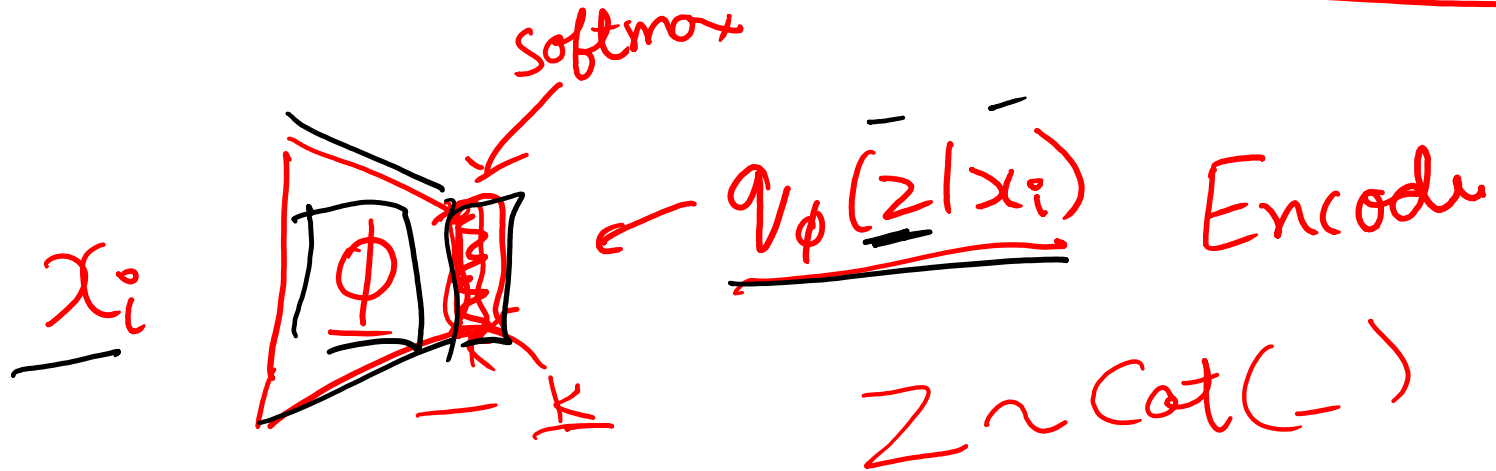
2. Variational Approximation

- Variational Lower Bound / ELBO

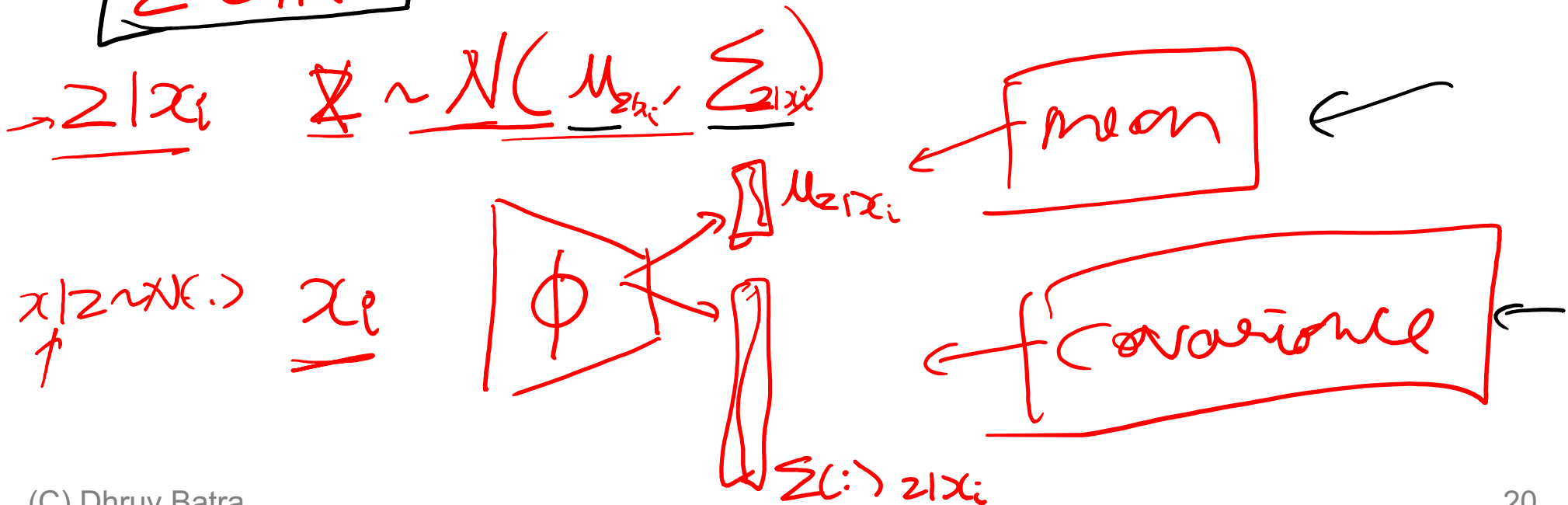
3. Amortized Inference Neural Networks

4. “Reparameterization” Trick

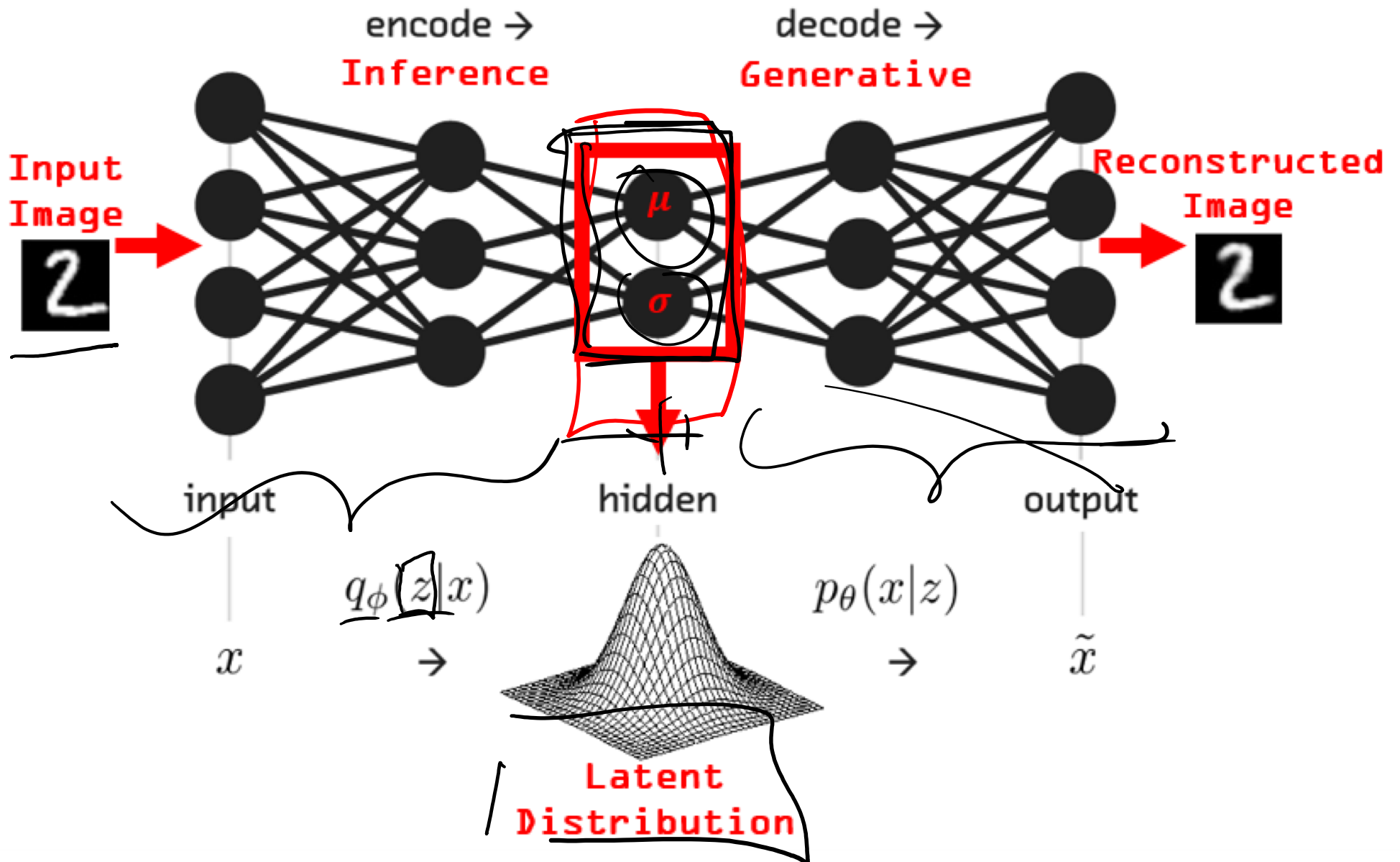
Amortized Inference Neural Networks



$$z \in \mathbb{R}^k$$

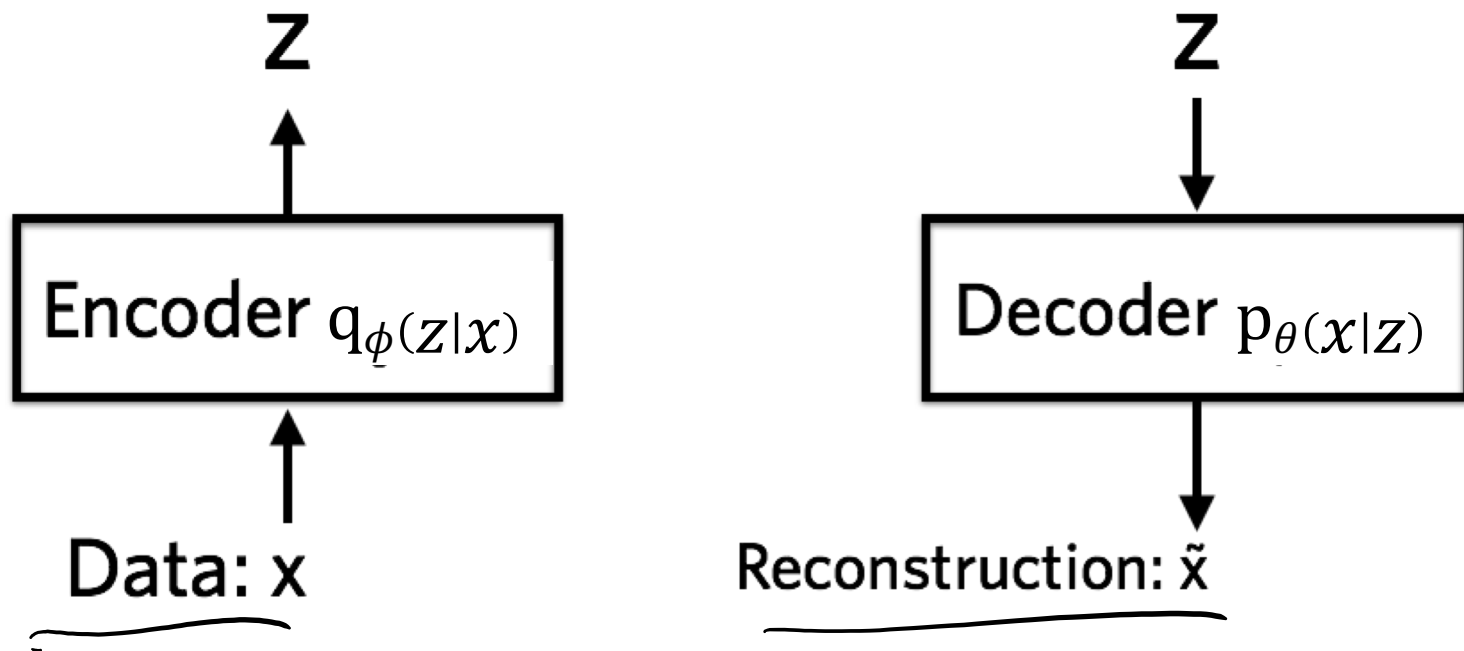


VAEs



Variational Autoencoders

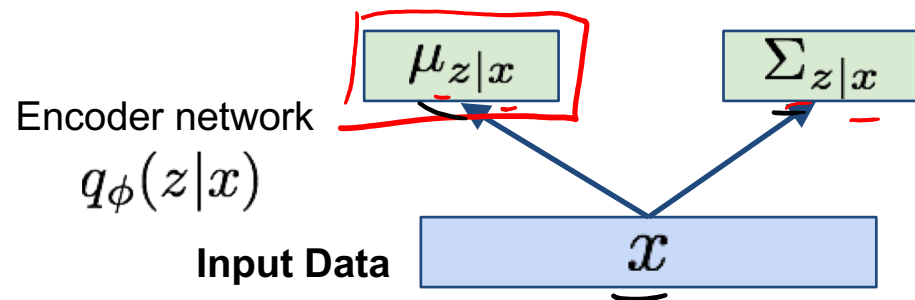
Probabilistic spin on autoencoders - will let us sample from the model to generate data!



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

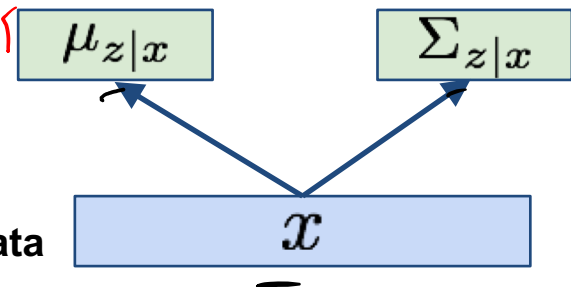
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

Make approximate posterior distribution close to prior

Encoder network

$$q_\phi(z|x)$$

Input Data

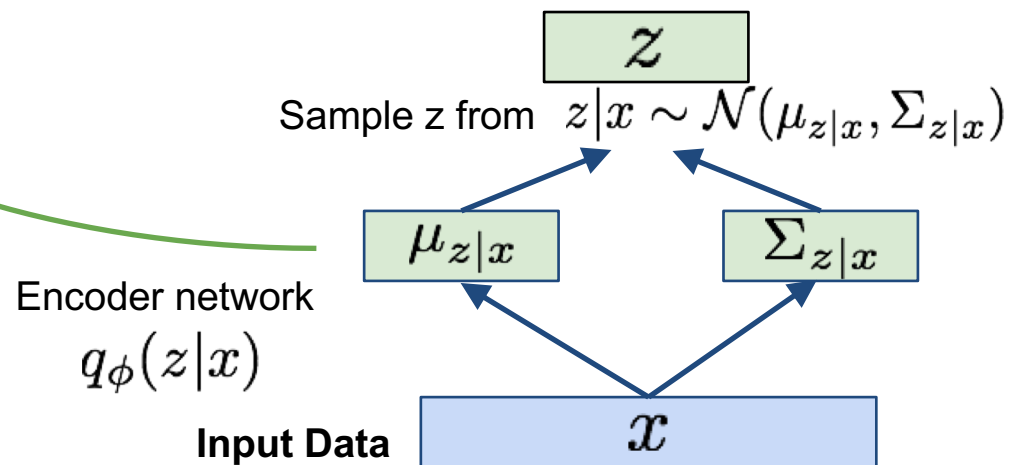


Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

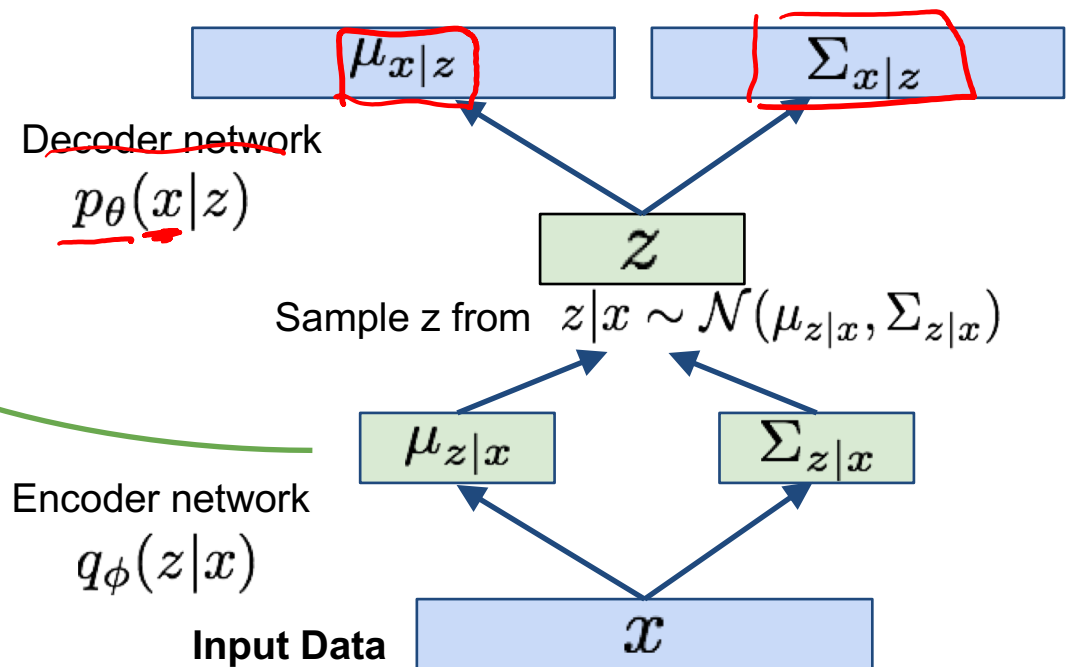


Variational Auto Encoders

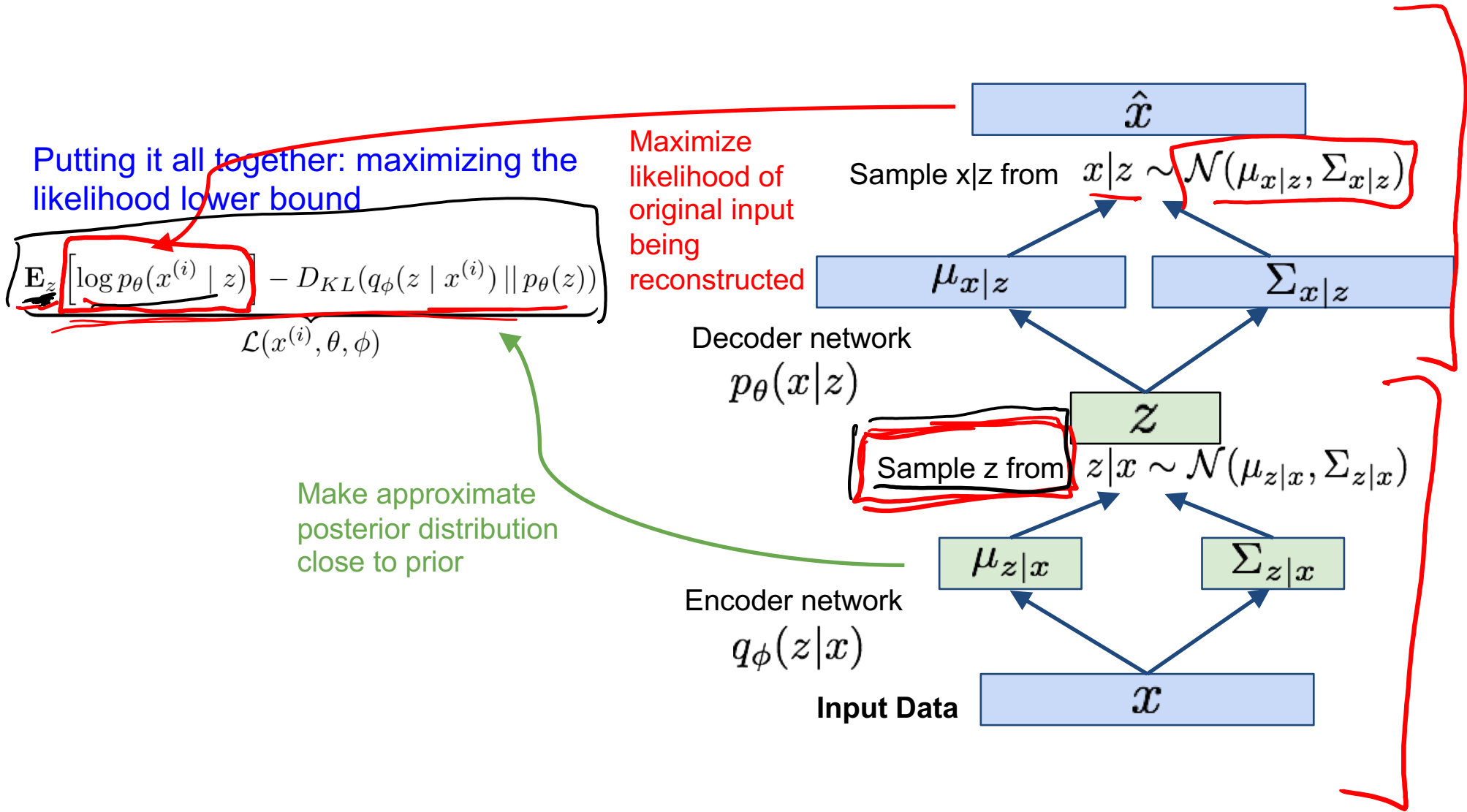
Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

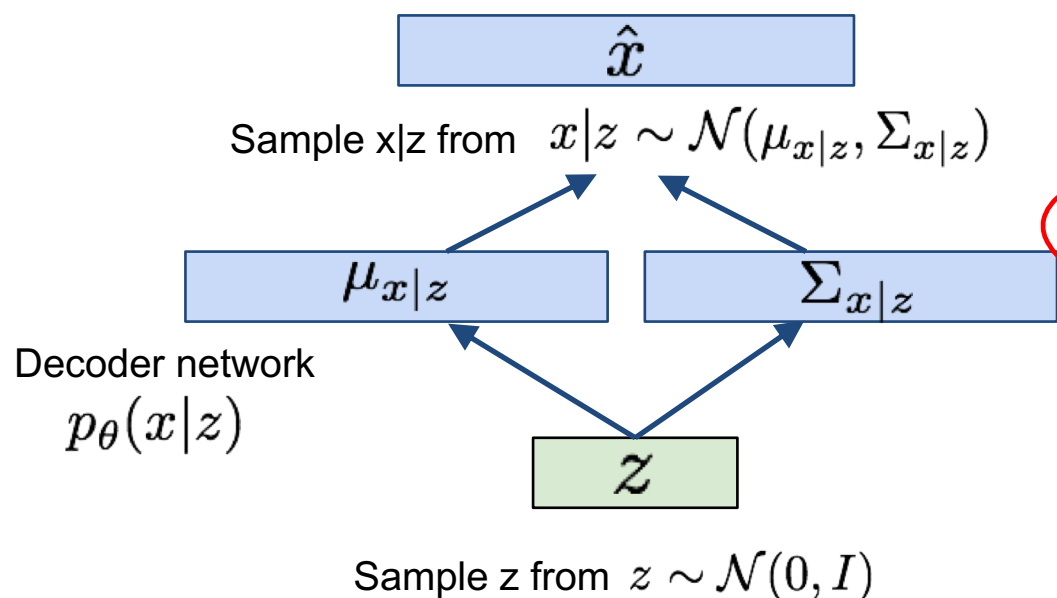


Variational Auto Encoders



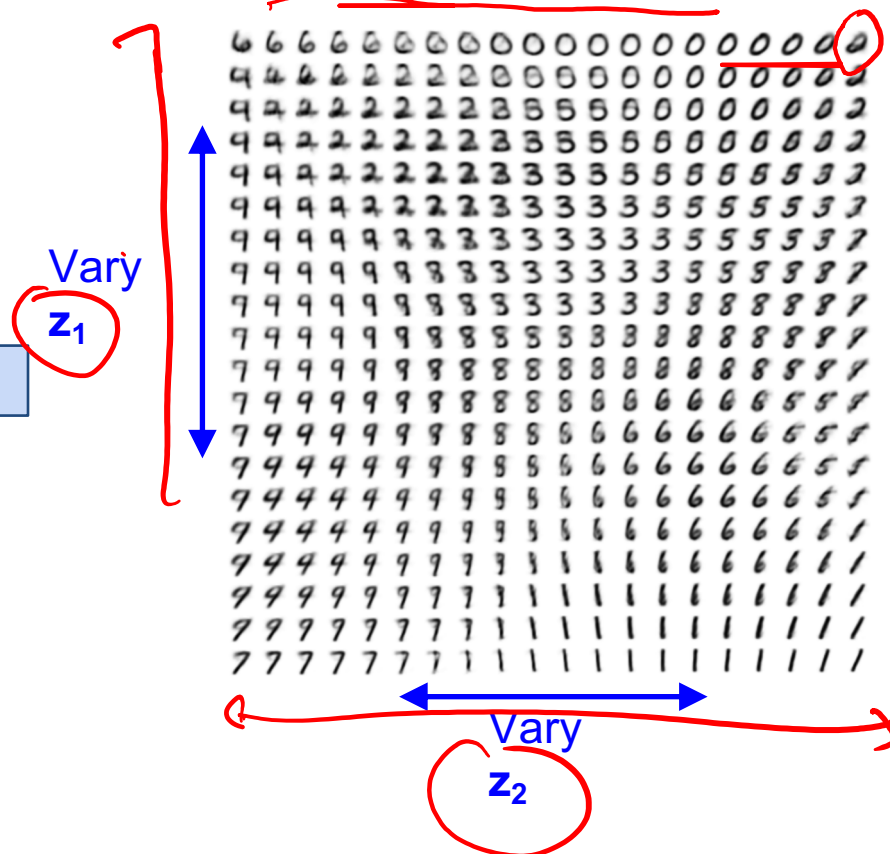
Variational Auto Encoders: Generating Data

Use decoder network. Now sample z from prior!



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

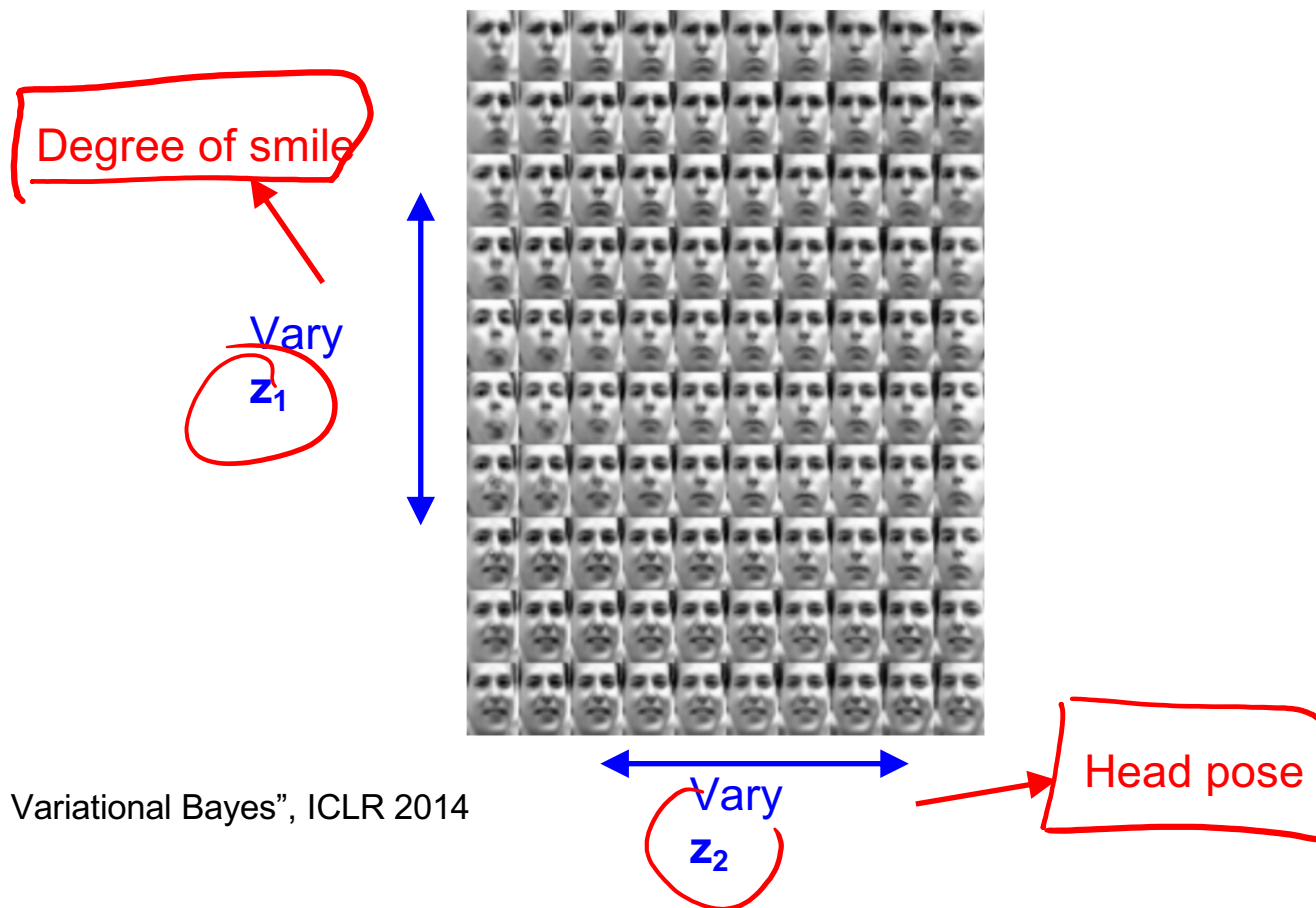
Data manifold for 2-d z



Variational Auto Encoders: Generating Data

Diagonal prior on \mathbf{z}
=> independent
latent variables

Different
dimensions of \mathbf{z}
encode
interpretable factors
of variation



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Plan for Today

- VAEs
 - Reparameterization trick
- Generative Adversarial Networks (GANs)

Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

Basic Problem

- Goal

$$\min_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [f(z)]$$

$$\theta^{(t)} \leftarrow \theta - \eta \nabla_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [f(z)]$$

Basic Problem

- Goal

$$\min_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [f(z)]$$

- Need to compute:

$$\begin{aligned} & \nabla_{\theta} \int f_{\theta}(z) p(z) dz \\ & \int \nabla_{\theta} f_{\theta}(z) \cdot p(z) dz \\ & \mathbb{E}_z [\nabla_{\theta} f_{\theta}(z)] \\ & \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} f_{\theta}(z_i) \end{aligned}$$

$$\nabla_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [f(z)]$$

$$\begin{aligned} & \nabla_{\theta} \int f(z) p_{\theta}(z) dz \\ & \int \nabla_{\theta} f(z) p_{\theta}(z) dz \\ & \int f(z) \nabla_{\theta} p_{\theta}(z) dz \end{aligned}$$

Does this happen in supervised learning?

$$\begin{aligned} \text{Goal } & \mathbb{E}_{x, y \sim P_{\text{data}}} \left[\min_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [l(y, \hat{y}(x, \theta))] \right] \\ &= \mathbb{E}_{x, y \sim P_{\text{data}}} \left[\nabla_{\theta} l(\dots, \theta) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[\nabla_{\theta} l(y_i, \hat{y}(x_i, \theta)) \right] \end{aligned}$$

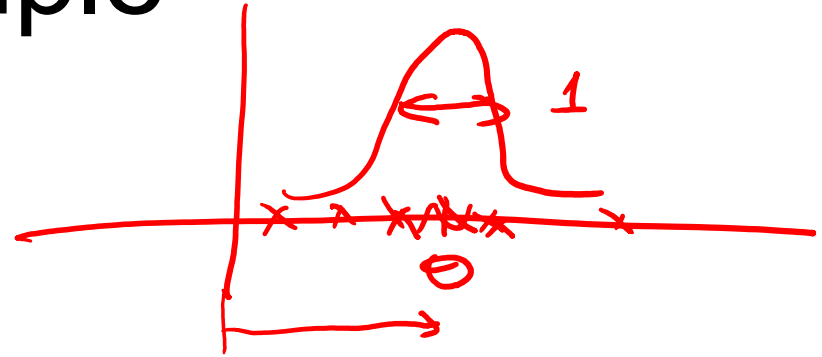
Example

$$z \sim N(\theta, 1)$$

$$f(z) = z^2$$

$$\min_{\theta} E_z[f(z)]$$

$$\min_{\theta} \int z^2 p(z) dz$$
$$\int z^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{(z-\theta)^2}{2}} dz$$



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

$$E[z^2] = \theta^2 + \text{Var}(z)$$

Example

Two Options

- ① • Score Function based Gradient Estimator
aka REINFORCE (and variants)

$$\underline{\nabla_{\theta} \mathbb{E}_z [f(z)]} = \underline{\mathbb{E}_z [f(z) \nabla_{\theta} \log p_{\theta}(z)]}$$

- ② • Path Derivative Gradient Estimator
aka “reparameterization trick”

$$\frac{\partial}{\partial \theta} \mathbb{E}_{z \sim p_{\theta}} [f(z)] = \frac{\partial}{\partial \theta} \mathbb{E}_{\epsilon} [f(g(\theta, \epsilon))] = \mathbb{E}_{\epsilon \sim p_{\epsilon}} \left[\frac{\partial f}{\partial g} \frac{\partial g}{\partial \theta} \right]$$

Option 1

- Score Function based Gradient Estimator
aka REINFORCE (and variants)

$$\nabla_{\theta} \mathbb{E}_z [f(z)] = \mathbb{E}_z [f(z) \nabla_{\theta} \log p_{\theta}(z)]$$

$$\nabla_{\theta} \int f(z) p_{\theta}(z) dz$$

$$= \int f(z) \left[\nabla_{\theta} p_{\theta}(z) \right] dz \cdot \frac{p_{\theta}(z)}{p_{\theta}(z)}$$

$$= \int f(z) \left[\nabla_{\theta} \log p_{\theta}(z) \right] \cdot p_{\theta}(z) dz$$

$$= \mathbb{E}_{z \sim p_{\theta}(z)} \left[f(z) \nabla_{\theta} \log p_{\theta}(z) \right] \approx \frac{1}{N} \sum$$

Example

$$P_{\theta}(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(z-\theta)^2}{2}}$$

$$\frac{\partial}{\partial \theta} \log P_{\theta}(z) = -\frac{(z-\theta)^2}{2} - \frac{1}{2} \log 2\pi$$
$$= \frac{-2(z-\theta) \cdot (-1)}{2} = (z-\theta)$$

$$\nabla_{\theta} = E \left[z^2 (z-\theta) \right]$$
$$\approx \frac{1}{N} \sum_{i=1}^N (z_i^2) (z_i - \theta)$$

Two Options

- Score Function based Gradient Estimator
aka REINFORCE (and variants)

$$\nabla_{\theta} \mathbb{E}_z [f(z)] = \mathbb{E}_z [f(z) \nabla_{\theta} \log p_{\theta}(z)]$$

- Path Derivative Gradient Estimator
aka "reparameterization trick"

$$\frac{\partial}{\partial \theta} \mathbb{E}_{z \sim p_{\theta}} [f(z)] = \frac{\partial}{\partial \theta} \mathbb{E}_{\epsilon} [f(g(\theta, \epsilon))] = \mathbb{E}_{\epsilon \sim p_{\epsilon}} \left[\begin{matrix} \frac{\partial f}{\partial g} & \frac{\partial g}{\partial \theta} \end{matrix} \right]$$

Option 2

- Path Derivative Gradient Estimator
aka “reparameterization trick”

$$\frac{\partial}{\partial \theta} \mathbb{E}_{z \sim p_{\theta}} [f(z)] = \frac{\partial}{\partial \theta} \mathbb{E}_{\epsilon} [f(g(\theta, \epsilon))] = \mathbb{E}_{\epsilon \sim p_{\epsilon}} \left[\frac{\partial f}{\partial g} \frac{\partial g}{\partial \theta} \right]$$

$\underline{z} \sim p_{\theta}(z)$

$\underline{z} = g(\underline{\theta}, \underline{\epsilon})$

$\underline{\epsilon} \sim U(-0.1, 1)$
 $\sim N(0, 1)$

$z \sim N(\mu, \sigma^2)$ $\epsilon \sim N(0, 1)$

$\rightarrow z = \underbrace{\mu}_{\text{constants}} + \sigma \underbrace{(\epsilon)}_{\text{simple RV}}$

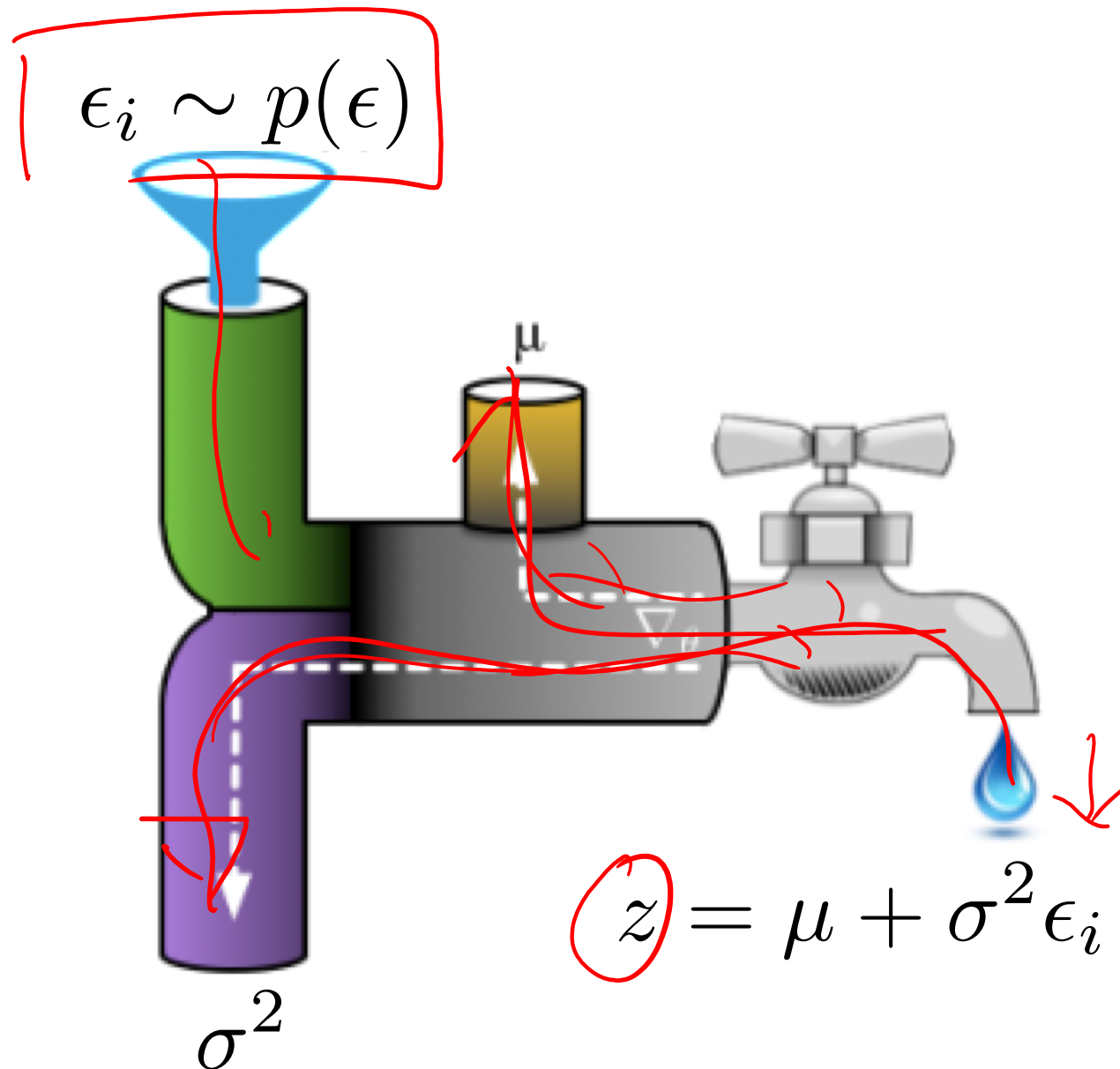
Option 2

- Path Derivative Gradient Estimator
aka “reparameterization trick”

$$\frac{\partial}{\partial \theta} \mathbb{E}_{z \sim p_{\theta}} [f(z)] = \frac{\partial}{\partial \theta} \mathbb{E}_{\epsilon \sim p_{\epsilon}} [f(g(\theta, \epsilon))] = \mathbb{E}_{\epsilon \sim p_{\epsilon}} \left[\frac{\partial f}{\partial g} \frac{\partial g}{\partial \theta} \right]$$

$$\begin{aligned} & \frac{\partial}{\partial \theta} \int f(g(\theta, \epsilon)) \cdot p(\epsilon) \cdot d\epsilon \quad \approx \frac{1}{N} \sum (\epsilon) \\ &= \int \frac{\partial}{\partial \theta} f(g(\theta, \epsilon)) \cdot p(\epsilon) \cdot d\epsilon \\ &= \int \frac{\partial f}{\partial g} \frac{\partial g}{\partial \theta} \cdot p(\epsilon) d\epsilon \end{aligned}$$

Reparameterization Intuition



Example

$$z \sim N(\theta, 1)$$

$$z = \theta + \varepsilon$$

$$f(z) = (\theta + \varepsilon)^2$$

$$\min_{\theta} \mathbb{E}_{\varepsilon \sim p(\varepsilon)} [(\theta + \varepsilon)^2]$$

$$= \mathbb{E}_{\varepsilon} [\nabla_{\theta} (\theta + \varepsilon)^2]$$

$$= \mathbb{E}_{\varepsilon} [2(\theta + \varepsilon)] = \underline{2\theta} + \mathbb{E}_{\varepsilon} [\varepsilon]$$

$$\approx \frac{1}{N} \sum_{i=1}^N 2(\theta + \varepsilon_i)$$

$$\varepsilon \sim N(0, 1)$$

$$\mathbb{E}_{\varepsilon} [\varepsilon] = 0$$

Two Options

- Score Function based Gradient Estimator
aka REINFORCE (and variants)

$$\nabla_{\theta} \mathbb{E}_z [f(z)] = \mathbb{E}_z [f(z) \nabla_{\theta} \log p_{\theta}(z)]$$

- Path Derivative Gradient Estimator
aka “reparameterization trick”

$$\frac{\partial}{\partial \theta} \mathbb{E}_{z \sim p_{\theta}} [f(z)] = \frac{\partial}{\partial \theta} \mathbb{E}_{\epsilon} [f(g(\theta, \epsilon))] = \mathbb{E}_{\epsilon \sim p_{\epsilon}} \left[\frac{\partial f}{\partial g} \frac{\partial g}{\partial \theta} \right]$$

Example

```
import numpy as np
```

```
N = 1000
```

```
theta = 2.0
```

```
x = np.random.randn(N) + theta
```

```
eps = np.random.randn(N)
```

```
grad1 = lambda x: np.sum(np.square(x)*(x-theta)) / x.size
```

```
grad2 = lambda eps: np.sum(2*(theta + eps)) / x.size
```

```
print grad1(x)
```

```
print grad2(eps)
```

```
4.46239612174
```

```
4.1840532024
```

20

Example

```
Ns = [10, 100, 1000, 10000, 100000]
reps = 100
```

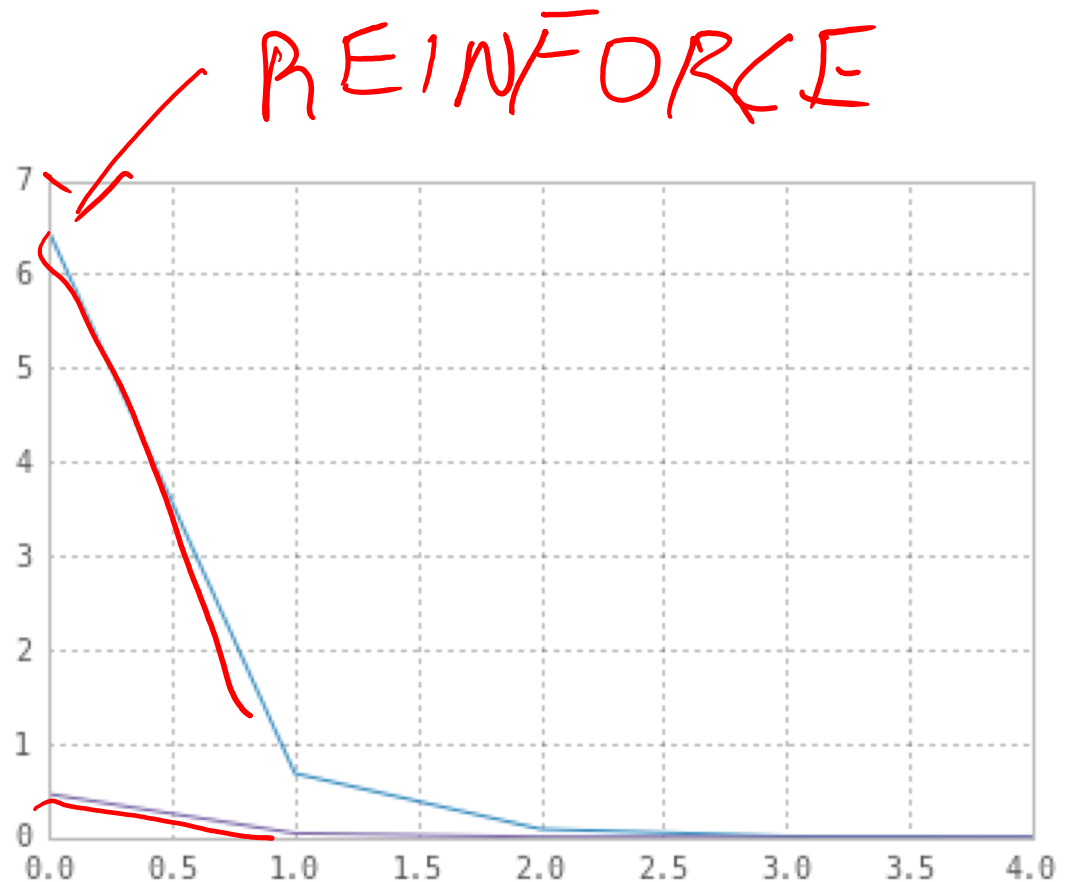
```
means1 = np.zeros(len(Ns))
vars1 = np.zeros(len(Ns))
means2 = np.zeros(len(Ns))
vars2 = np.zeros(len(Ns))
```

```
est1 = np.zeros(reps)
est2 = np.zeros(reps)
for i, N in enumerate(Ns):
    for r in range(reps):
        x = np.random.randn(N) + theta
        est1[r] = grad1(x)
        eps = np.random.randn(N)
        est2[r] = grad2(eps)
    means1[i] = np.mean(est1)
    means2[i] = np.mean(est2)
    vars1[i] = np.var(est1)
    vars2[i] = np.var(est2)
```

```
print means1
print means2
print
print vars1
print vars2
```

```
[ 3.8409546  3.97298803  4.03007634  3.98531095  3.99579423]
[ 3.97775271  4.00232825  3.99894536  4.00353734  3.99995899]
```

```
[ 6.45307927e+00  6.80227241e-01  8.69226368e-02  1.00489791e-02
 8.62396526e-04]
[ 4.59767676e-01  4.26567475e-02  3.33699503e-03  5.17148975e-04
 4.65338152e-05]
```



Variational Auto Encoders

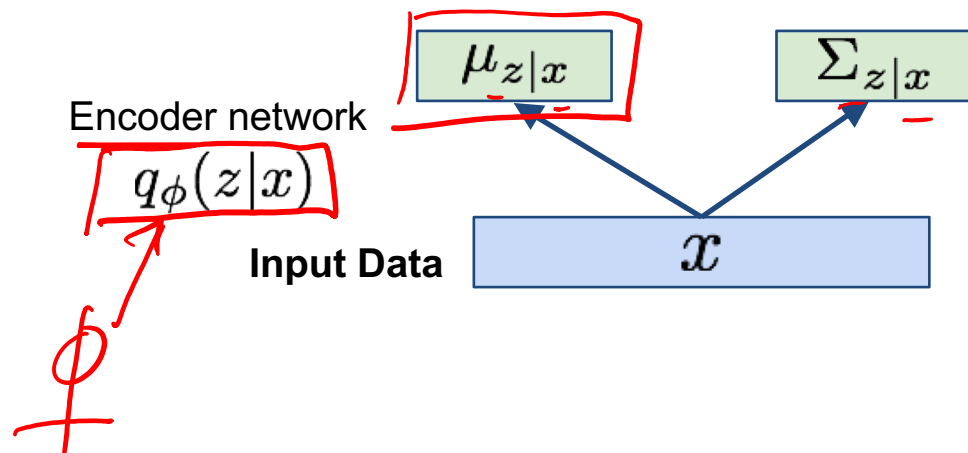
VAEs are a combination of the following ideas:

1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

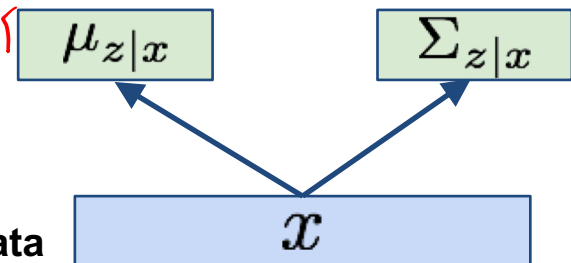
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\text{Make approximate posterior distribution close to prior}}$$

Make approximate posterior distribution close to prior

Encoder network

$$q_\phi(z|x)$$

Input Data



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

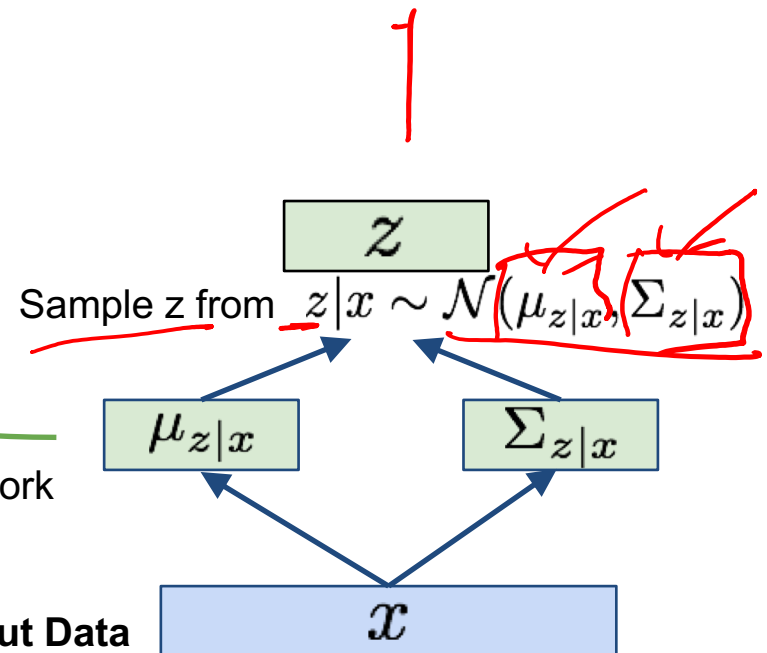
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Encoder network

$$q_\phi(z|x)$$

Input Data



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

