

CS 4803 / 7643: Deep Learning

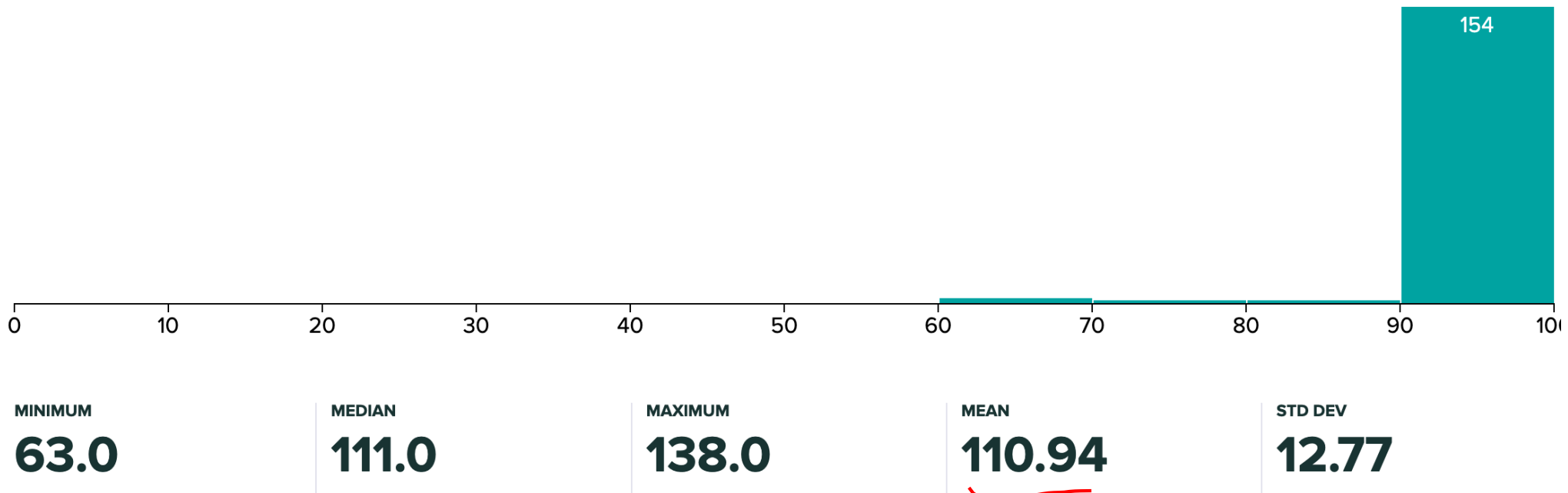
Topics:

- Variational Auto-Encoders (VAEs)
- Reparameterization trick

Dhruv Batra
Georgia Tech

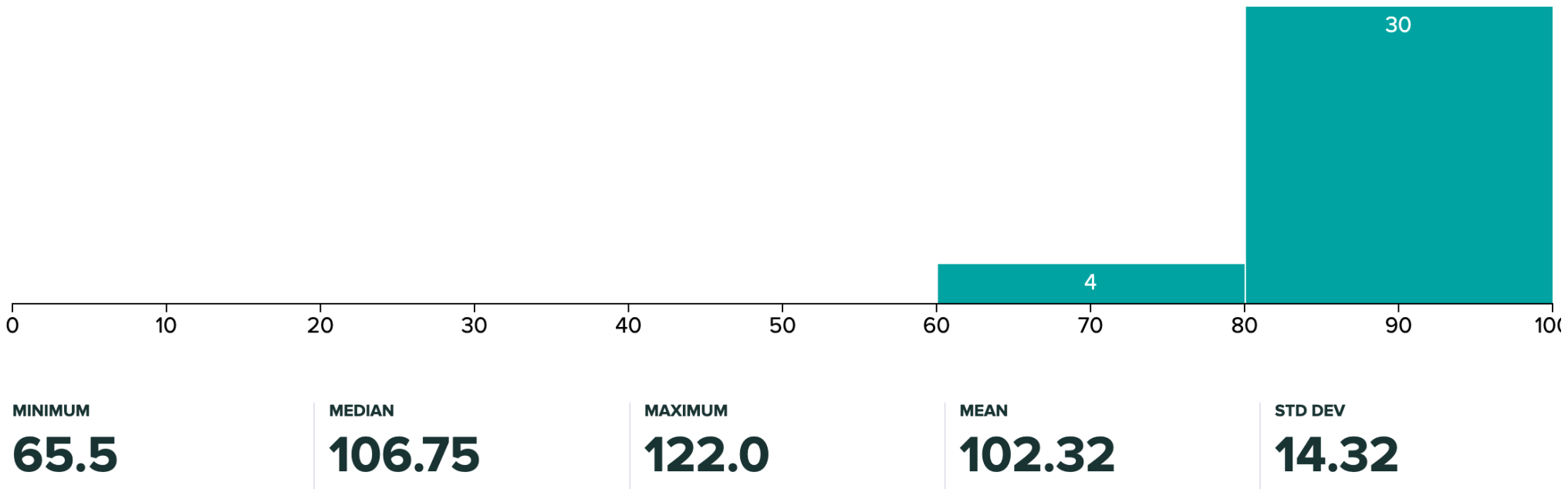
Administrativa

- HW4 Grades Released
 - Regrade requests close: 12/03, 11:55pm
 - Please check solutions first!
- Grade histogram: 7643
 - Max possible: 100 (regular credit) + 40 (extra credit)




Administrativa

- HW4 Grades Released
 - Regrade requests close: 12/03, 11:55pm
 - Please check solutions first!
- Grade histogram: 4803
 - Max possible: 100 (regular credit) + 40 (extra credit)



Recap from last time



Variational Autoencoders (VAE)

So far...

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

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

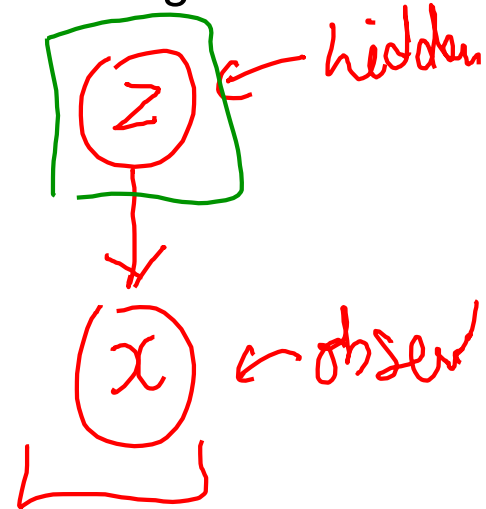
VAEs define intractable density function with latent z :

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

$$\sum_z p_{\theta}(z) p_{\theta}(x|z)$$

z continuous

z discrete



Variational Auto Encoders

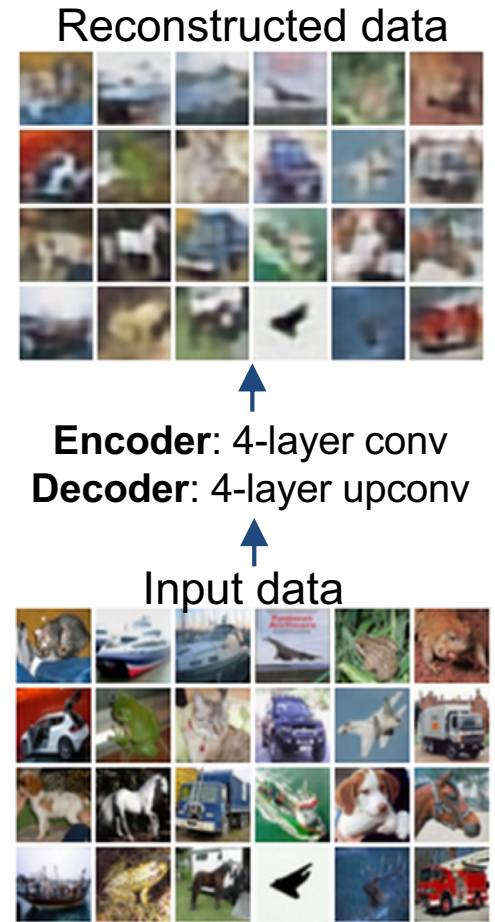
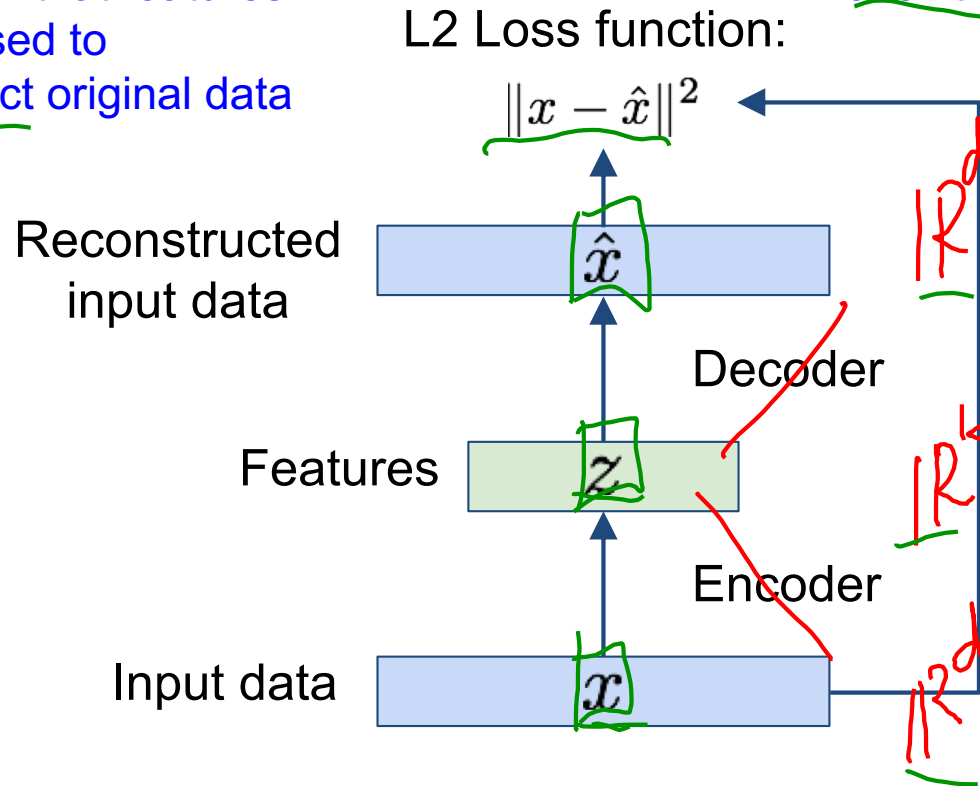
VAEs are a combination of the following ideas:

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

Autoencoders

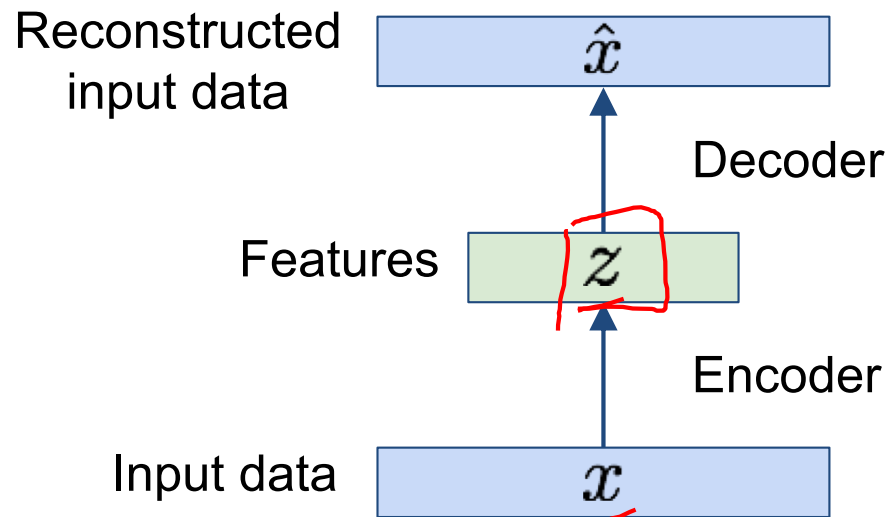
Train such that features can be used to reconstruct original data

Doesn't use labels!



Autoencoders

$$z = f_{\phi}(x)$$
$$\hat{x} = g_{\phi}(z)$$
$$p(z|x)$$
$$p(\hat{x}|z)$$

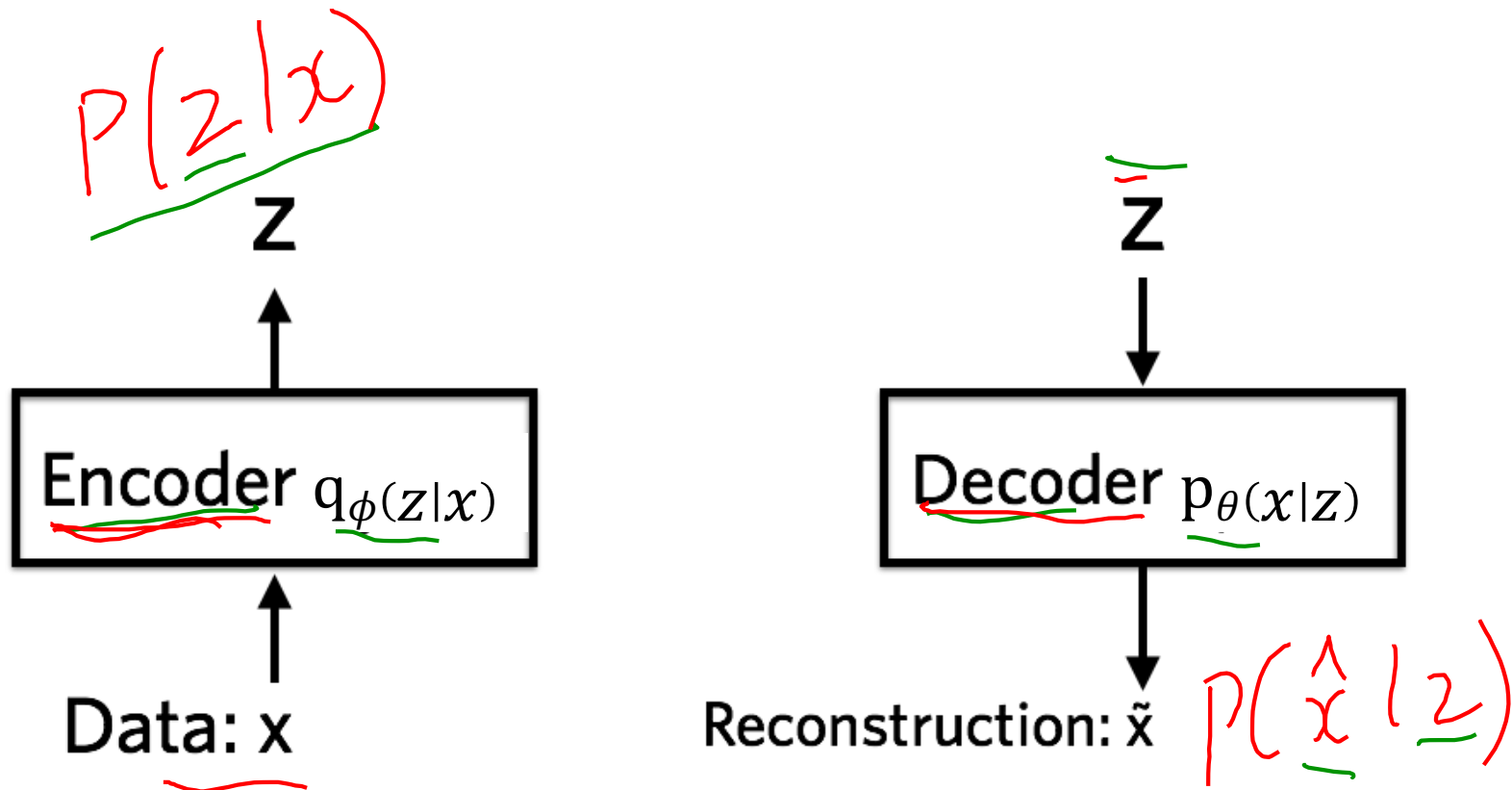


Autoencoders can reconstruct data, and can learn features to initialize a supervised model

Features capture factors of variation in training data. Can we generate new images from an autoencoder?

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

• $P(\vec{z}|\vec{x}) =$

$q_i(z)$

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

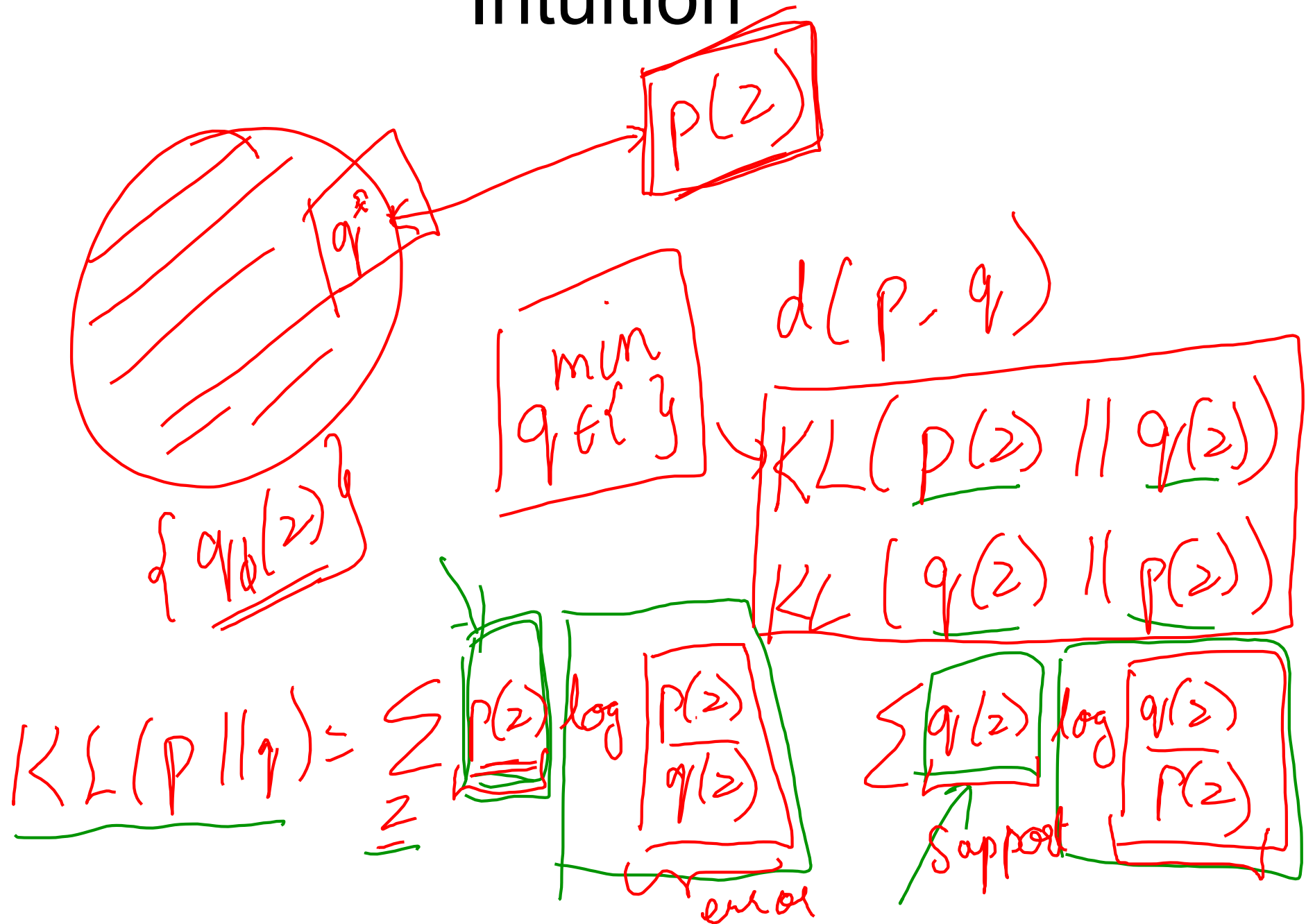
Hard



What is Variational Inference?

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

Intuition



The general learning problem with missing data

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

$$\begin{aligned}
 \underline{\underline{ll(\theta : \mathcal{D})}} &= \underline{\underline{\log}} \left[\prod_{i=1}^N P(\underline{\underline{\mathbf{x}_i}} \mid \theta) \right] \\
 &= \sum_{i=1}^N \log P(\underline{\underline{\mathbf{x}_i}} \mid \theta) \\
 &= \sum_{i=1}^N \log \left[\sum_{\mathbf{z}} P(\underline{\underline{\mathbf{x}_i}}, \mathbf{z} \mid \theta) \right]
 \end{aligned}$$

②
↓
①

$\mathcal{D} = \{ \vec{x}_1, \dots, \vec{x}_N \}$

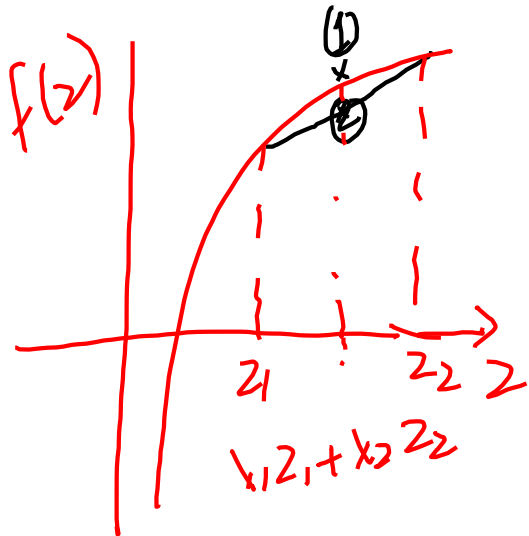
$P(\vec{x}, \mathbf{z})$

$P(x_i \mid \theta) P(\mathbf{z} \mid x_i, \theta)$

$\log \sum_{\mathbf{z}} P(\mathbf{z} \mid x_i, \theta) P(x_i \mid \theta)$
 $\approx \log \left[\mathbb{E} \left[P(\mathbf{z} \mid x_i, \theta) \right] \right] \left[P(x_i \mid \theta) \right]$

Jensen's inequality

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



$$\textcircled{1} \geq \textcircled{2}$$

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

$$f\left(\sum_{i=1}^{k(2)} \lambda_i z_i\right) \geq \sum_{i=1}^{k(2)} \lambda_i f(z_i) \quad \leftarrow z \rightarrow k$$

$$\begin{cases} \lambda_1, \lambda_2 \geq 0 \\ \lambda_1 + \lambda_2 = 1 \end{cases}$$

$$\equiv f(E[z]) \geq E_{P(z)}[f(z)] \quad \leftarrow z \rightarrow g(z)$$

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

Applying Jensen's inequality

(z)
↓
(x)

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

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

$$\ell(\theta) \geq F(\theta, Q_i)$$

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

"Free Energy" $F(\theta, Q_i)$

$$\max_{\theta} \ell \Rightarrow \max_{\theta, Q_i} F$$

Variational Lower Bound
Evidence Lower Bound (ELBO)

Evidence Lower Bound

- Define potential function $F(\theta, Q)$:

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

(VAES)

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

(GMMs)

$\rightarrow P(z | x_i, \theta) P(x_i | \theta)$

ELBO: Factorization #1 (GMMs)

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

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

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

$l(\theta) \Rightarrow F(\theta, Q_i) \Rightarrow l(\theta) = F(\theta, Q_i) + \text{KL}(\dots)$

ELBO: Factorization #2 (VAEs)

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

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

$$\stackrel{\text{(VAEs)}}{=} \left[\underbrace{F(Q_i(\mathbf{z}))}_{\text{Explains the data}} \left[\log P(\tilde{\mathbf{x}}_i | \mathbf{z}, \theta) \right] \right] = \left[\text{KL} \left(Q_i(\mathbf{z}) \parallel P(\mathbf{z} | \theta) \right) \right]$$

Explains the data

Regulariser
Be simple

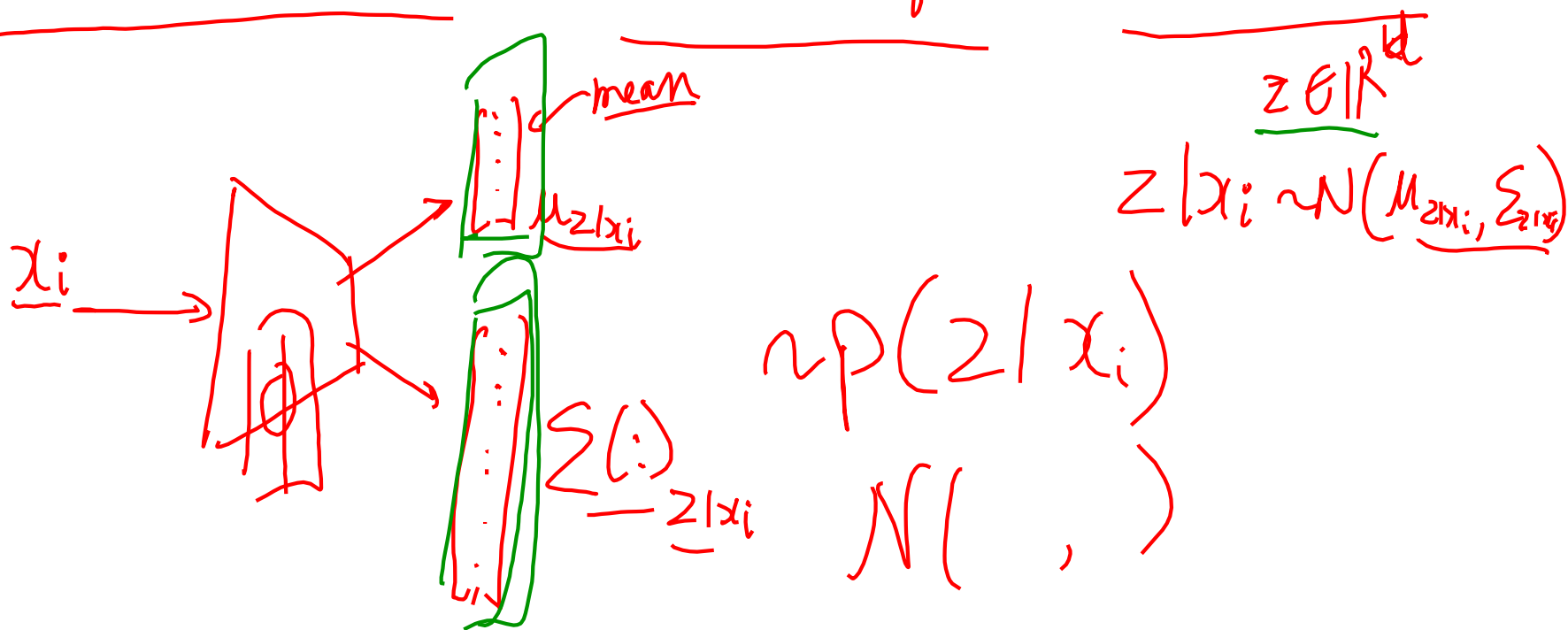
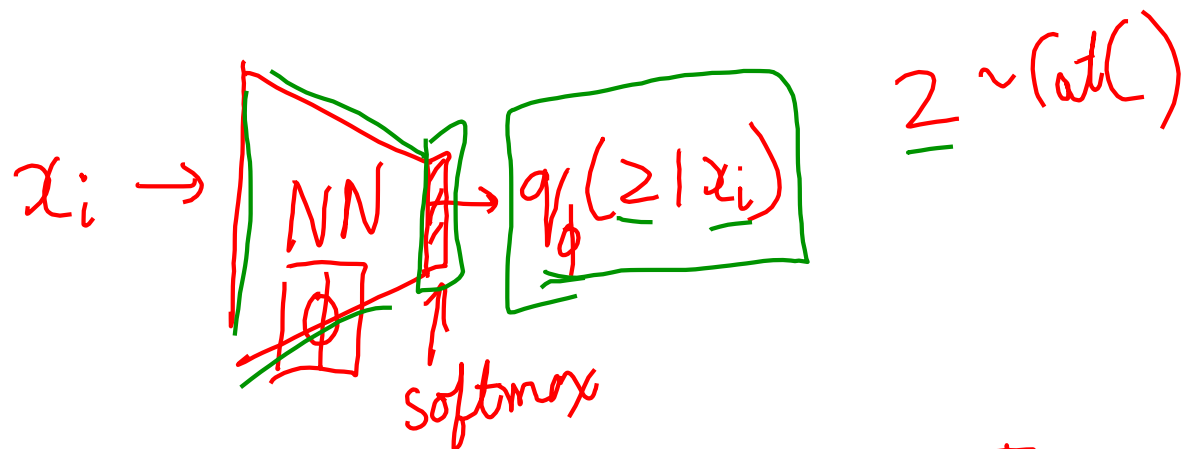
Variational Auto Encoders

VAEs are a combination of the following ideas:

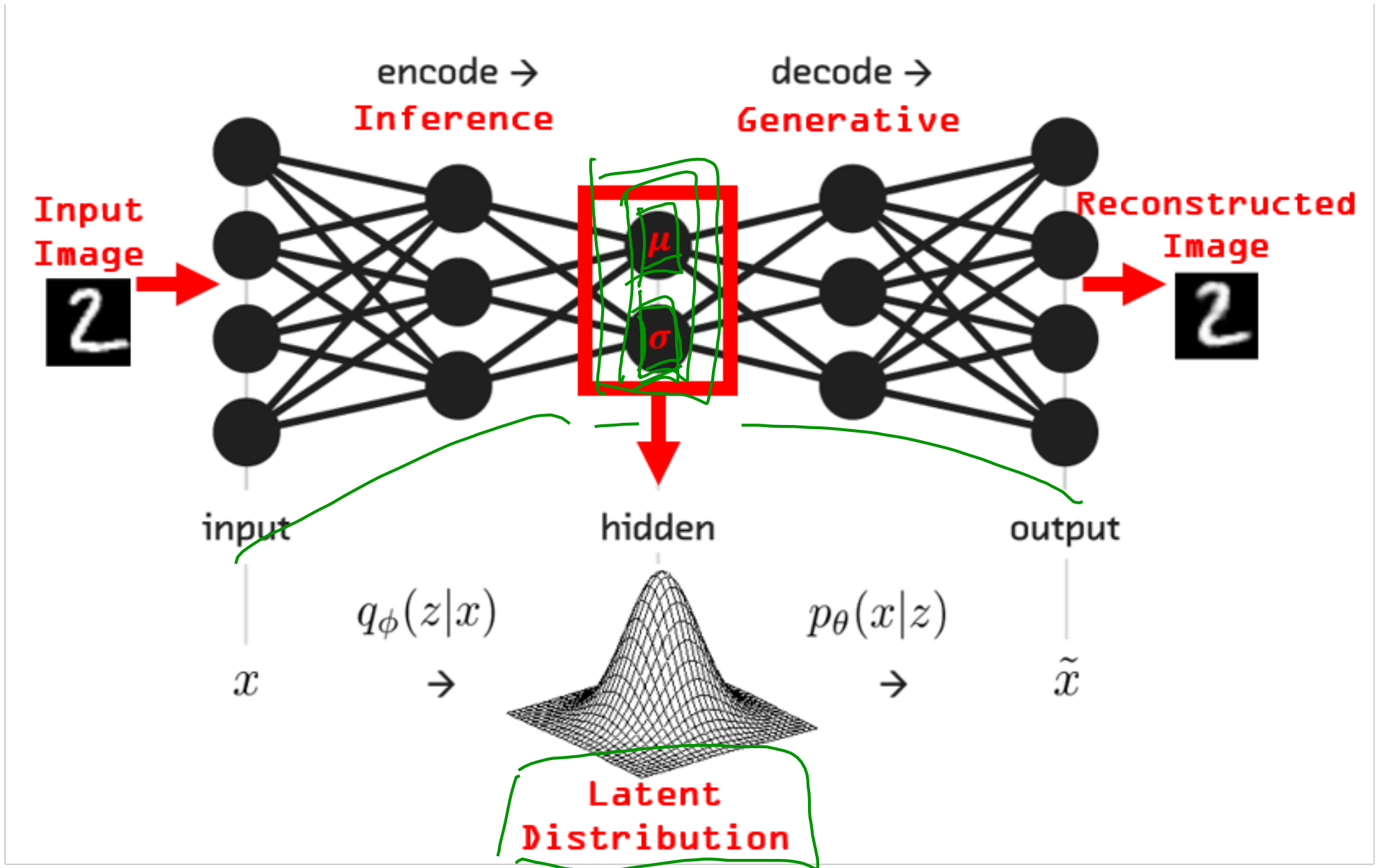
1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

Amortized Inference Neural Networks

$$Q_i(z) = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.3 \end{bmatrix}$$

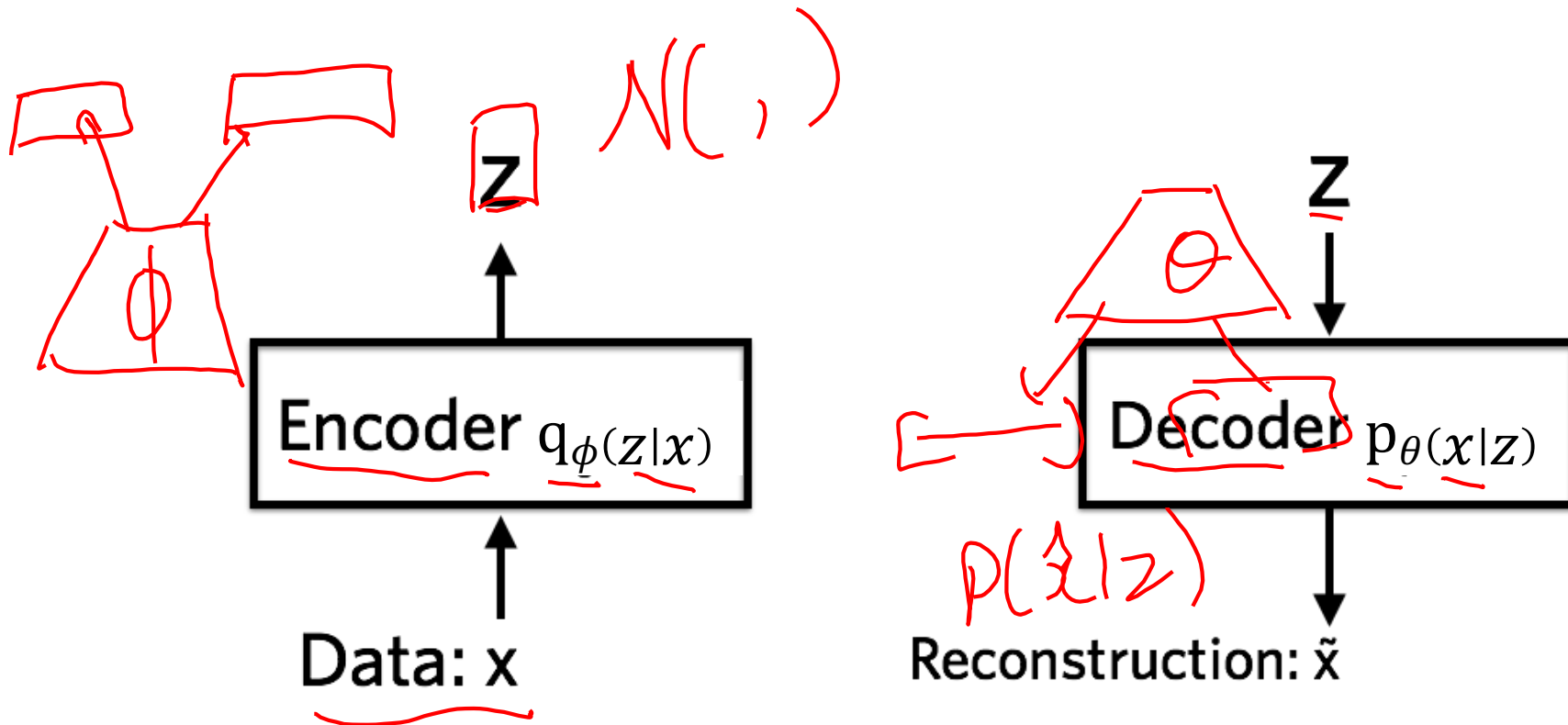


VAEs



Variational Autoencoders

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



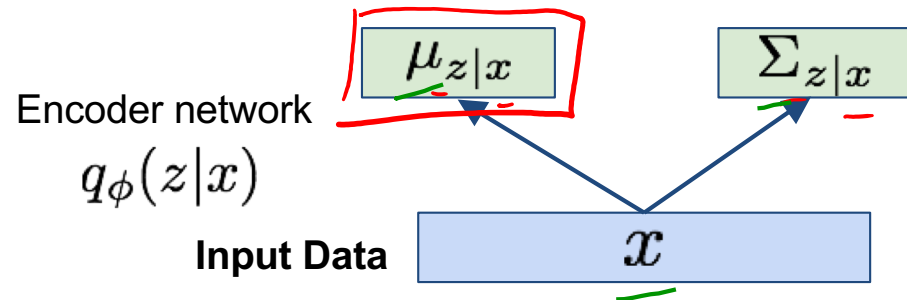
Variational Auto Encoders

$F(\theta, \phi)$

Putting it all together: maximizing the likelihood lower bound

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

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



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

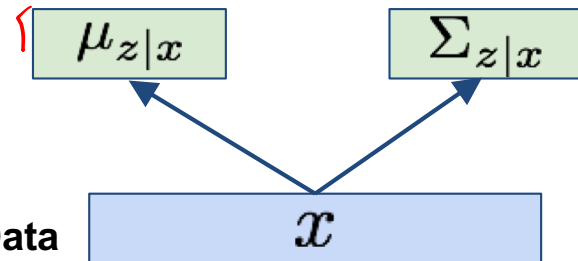
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\text{Make approximate posterior distribution close to prior}}$$

Make approximate posterior distribution close to prior

Encoder network

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

Input Data

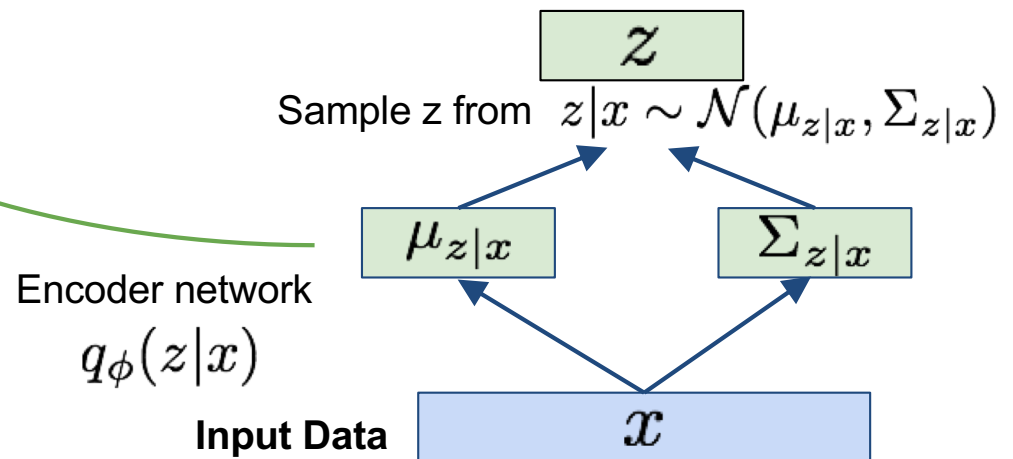


Variational Auto Encoders

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

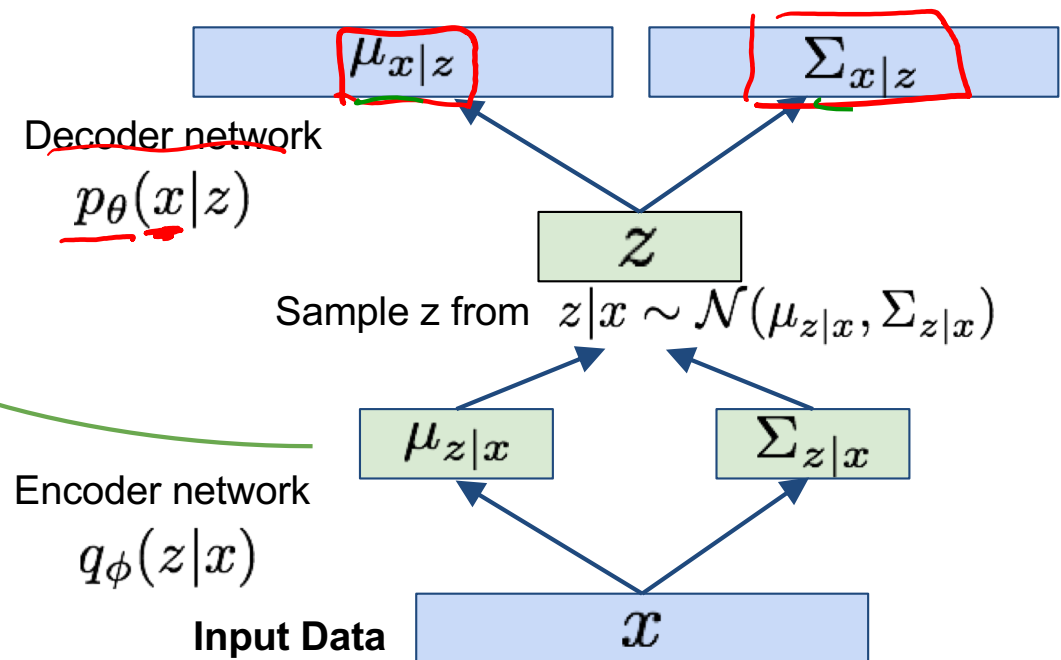


Variational Auto Encoders

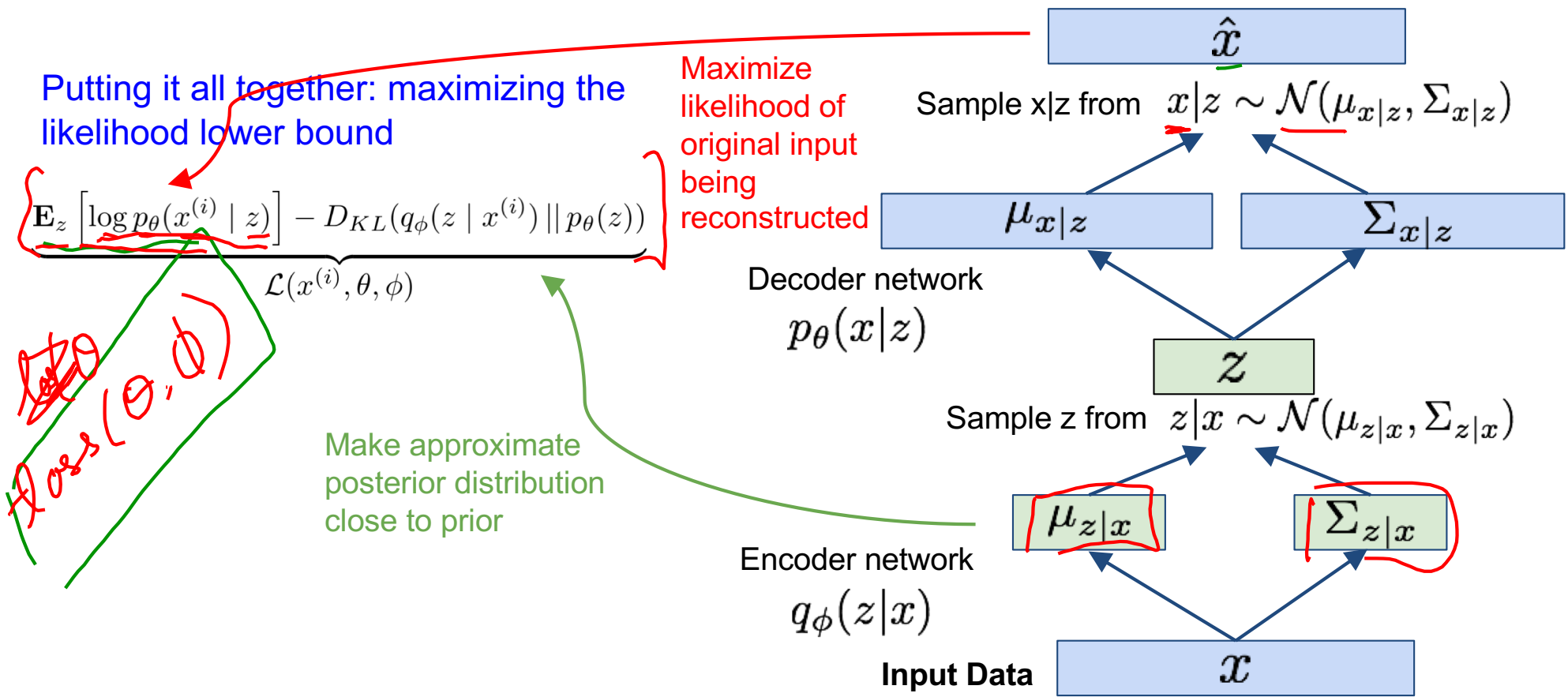
Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior

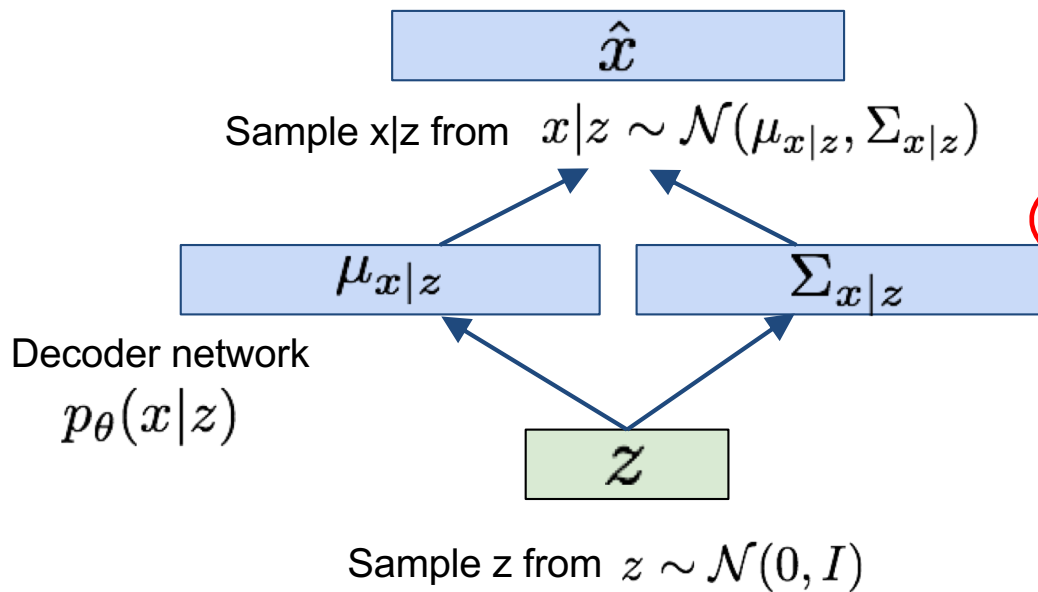


Variational Auto Encoders



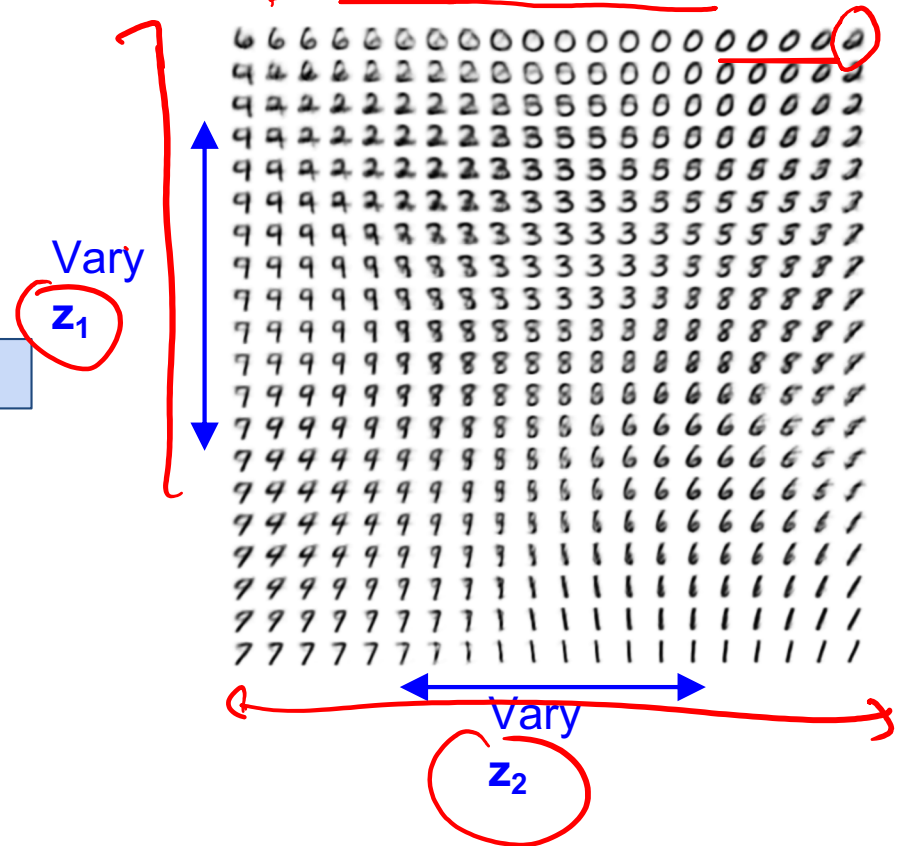
Variational Auto Encoders: Generating Data

Use decoder network. Now sample z from prior!



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

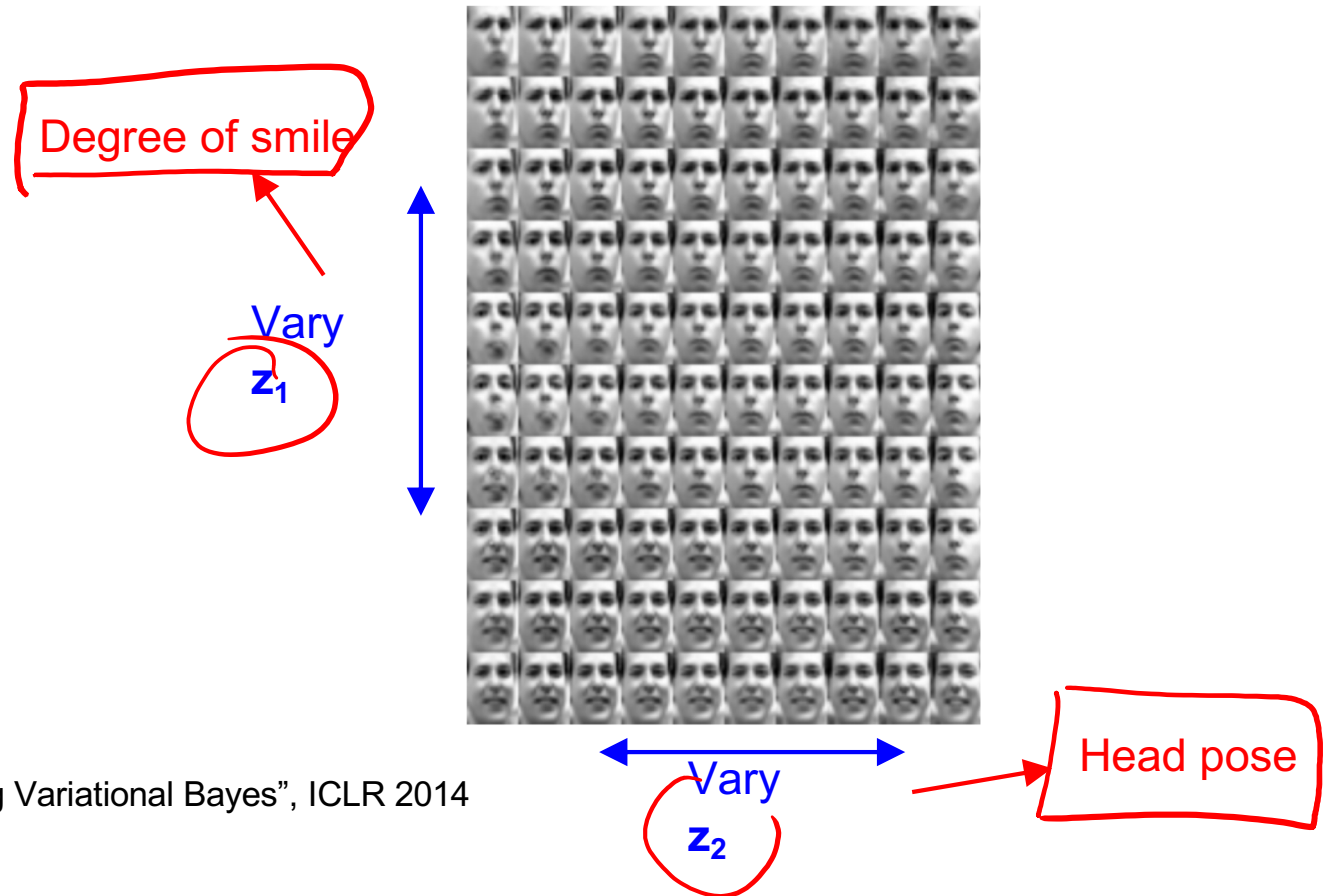
Data manifold for 2-d z



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

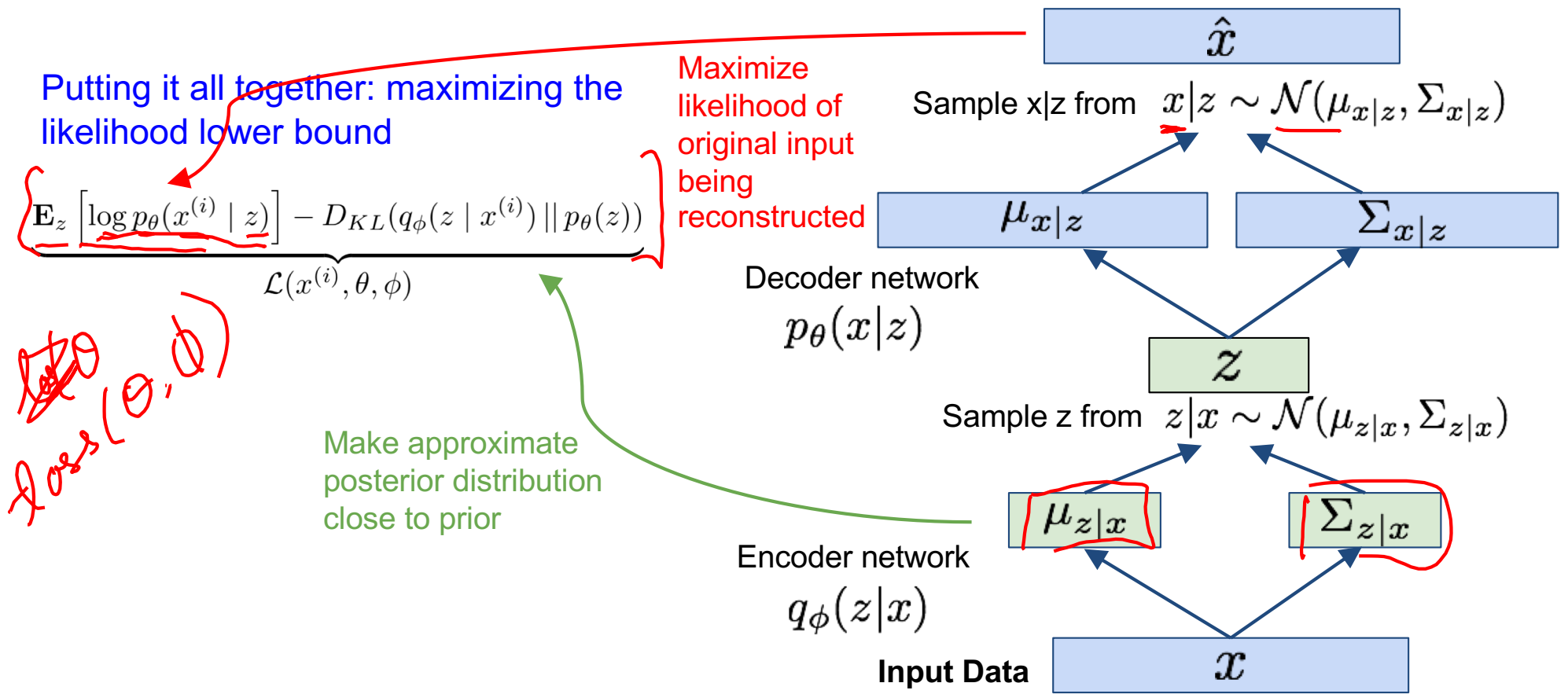
- VAEs
 - Reparameterization trick

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



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

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

max
 θ, ϕ

Basic Problem

$z \sim \text{cat}(\pi)$

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

Basic Problem

- Goal

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

Basic Problem

- Goal

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

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [f(z)]$$

- Need to compute:

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

$$f_{\theta}(z) \quad p(z)$$

$$\int \nabla_{\theta} f_{\theta}(z) p(z) dz$$

$$= \mathbb{E} [\nabla_{\theta} f_{\theta}(z)]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} f(z_i) \quad z_i \sim p(z)$$

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

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

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

Basic Problem

- Need to compute: $\nabla_{\theta} \mathbb{E}_{z \sim p_{\theta}(z)} [f(z)]$

Example

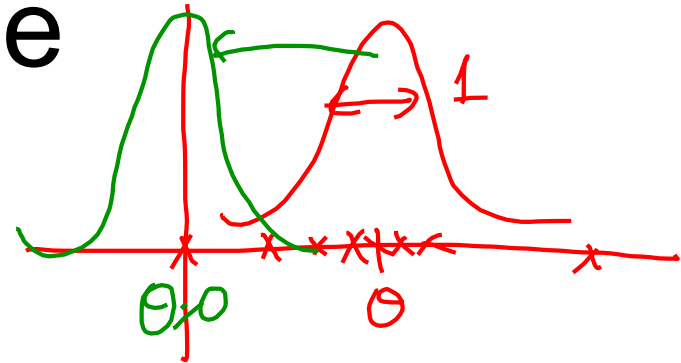
$$z \sim N(\theta, 1)$$

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

$$\min_{\theta} E_z[z^2]$$

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

"hard"



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

$$\min_{\theta} E[z^2] = \frac{\text{Var}(z)}{1} + \theta^2$$

"easy"

Does this happen in supervised learning?

- Goal

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

$$\min_{\theta} \nabla_{\theta} \mathbb{E}_{x, y \sim \boxed{P_{\text{data}}}}$$

$$\left[\begin{array}{c} \boxed{f_{\theta}(z)} \\ \downarrow \\ l(y^{\text{gt}}, \hat{y}(x, \theta)) \end{array} \right]$$

$$\nabla_{\theta} E \approx \frac{1}{N} \sum \nabla_{\theta} l(y_i, \hat{y}_i(x_i, \theta))$$

$\boxed{y_i, x_i \sim P_{\text{data}}}$

But what about other kinds of learning?

- Goal

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

VL

$$\max_{\theta, \phi}$$

$$\mathbb{E}_{q_{\phi}(z)} \left[\underbrace{\log p(\dots)}_{\theta, \phi} \right]$$

RL

$$\max_{\theta}$$

$$\mathbb{E}_{a_1 \dots a_T \sim \pi_{\theta}(a|s)} \left[\sum r_t(s_t, a_t) \right]$$

Two Options

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

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

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

- 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)

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

$$\begin{aligned} \nabla_{\theta} \int f(z) p_{\theta}(z) dz &= \int f(z) \nabla_{\theta} p_{\theta}(z) dz \cdot \frac{p_{\theta}(z)}{p_{\theta}(z)} \\ &= \int f(z) \nabla_{\theta} \log p_{\theta}(z) p_{\theta}(z) dz \\ &= \mathbb{E} [f(z) \nabla_{\theta} \log p_{\theta}(z)] \\ &\approx \frac{1}{N} \left(\begin{array}{c} \downarrow \\ \end{array} \right) \end{aligned}$$

Recall: Policy Gradients

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\mathcal{R}(\tau)] \\ &= \nabla_{\theta} \int \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau && \text{Expand expectation} \\ &= \int \nabla_{\theta} \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau && \text{Exchange integration and expectation} \\ &= \int \nabla_{\theta} \pi_{\theta}(\tau) \cdot \frac{\pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} \cdot \mathcal{R}(\tau) d\tau \\ &= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau && \nabla_{\theta} \log \pi(\tau) = \frac{\nabla_{\theta} \pi(\tau)}{\pi(\tau)} \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) \mathcal{R}(\tau)]\end{aligned}$$

Example

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

$$E_Z \left[f(Z) \nabla_{\theta} \log P_{\theta}(Z) \right]$$

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

$$\frac{\partial}{\partial \theta} \left[\log P_{\theta}(z) \right] = \left[\frac{-(z-\theta)^2}{2} - \frac{1}{2} \log 2\pi \right]$$

$$\downarrow$$
$$\frac{\partial}{\partial \theta} \left[\frac{-(z-\theta)^2}{2} \right]$$

$$= (z-\theta)$$

Gradient Estimator

$$E_Z \left[z^2 (z-\theta) \right]$$

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

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

Mental Break!

- VAE Demo
 - <https://www.siares.com/projects/variational-autoencoder>

Two Options

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

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

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

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

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

Option 2

- Path Derivative Gradient Estimator
aka “reparameterization trick”

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

$$\underline{z} \sim p_{\theta}(\underline{z}) \implies \underline{z} = g(\underline{\theta}, \underline{\epsilon}) \quad \epsilon \sim \text{“Standard” RV}$$

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

$$\underline{z} = \underline{\theta} \mu + \sigma \underline{\epsilon}$$

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

(μ, σ)

$$\begin{aligned} \epsilon &\sim N(0, 1) \\ &\sim U(0, 1) \\ &\sim \text{Bern}(0.5) \\ &\vdots \\ &\vdots \end{aligned}$$

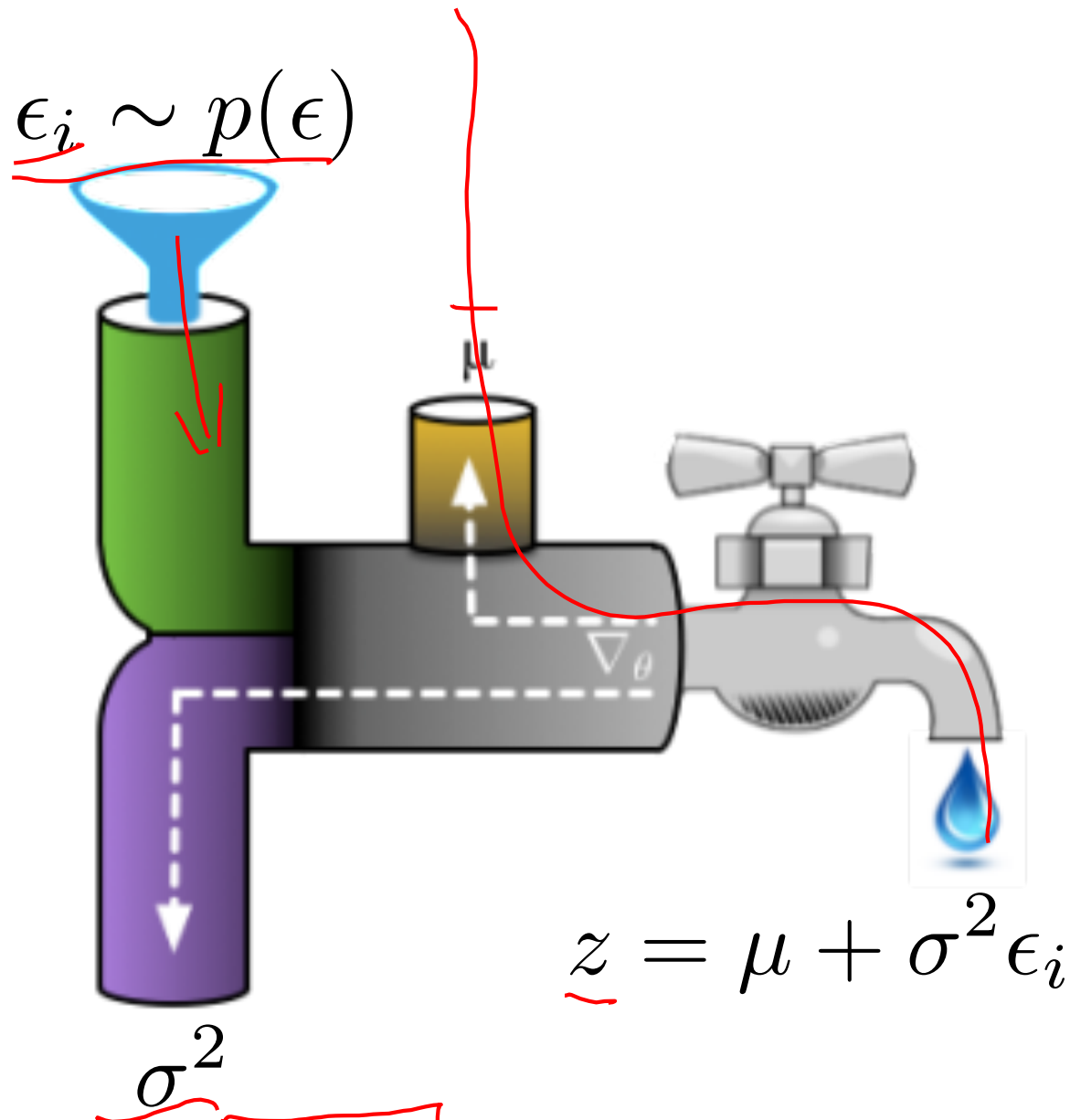
Option 2

- Path Derivative Gradient Estimator
aka “reparameterization trick”

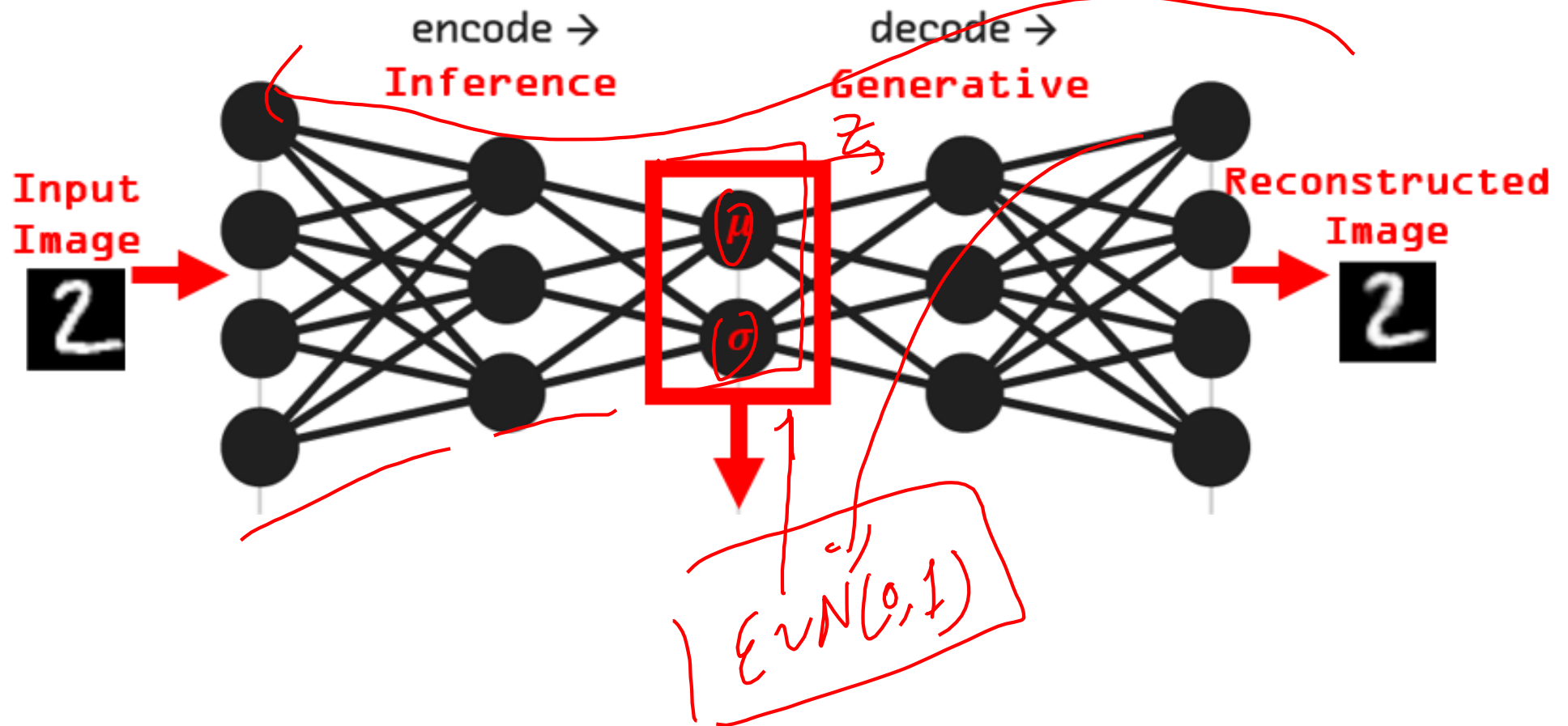
$$\frac{\partial}{\partial \theta} \mathbb{E}_{z \sim p_{\theta}} [f(z)] = \frac{\partial}{\partial \theta} \mathbb{E}_{\epsilon} [f(g(\theta, \epsilon))] = \mathbb{E}_{\epsilon \sim p_{\epsilon}} \left[\frac{\partial f}{\partial g} \frac{\partial g}{\partial \theta} \right]$$

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

Reparameterization Intuition



Reparameterization Intuition



Example

$$z \sim N(\theta, 1)$$
$$f(z) = z^2$$

$$\min_{\theta} E_z [z^2]$$

$$\frac{\partial}{\partial \theta} E_{\varepsilon} [(\theta + \varepsilon)^2]$$
$$= E_{\varepsilon} \left[\frac{\partial}{\partial \theta} (\theta + \varepsilon)^2 \right]$$

$$= E_{\varepsilon} [2(\theta + \varepsilon)]$$

$$z = g(\theta, \varepsilon)$$

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

$$z = \theta + \varepsilon$$

$$E [2\theta + 2\varepsilon]$$

$$\min_{\theta} \theta^2$$

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

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

Two Options

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

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

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

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

Example

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} L$$

```
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)
```

$\theta^{(t)} = 2$

$z_i^2 (z_i - \theta)$

$2(\theta + \epsilon_i)$

```
4.46239612174
4.1840532024
```

$\approx 2\theta$

Example

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

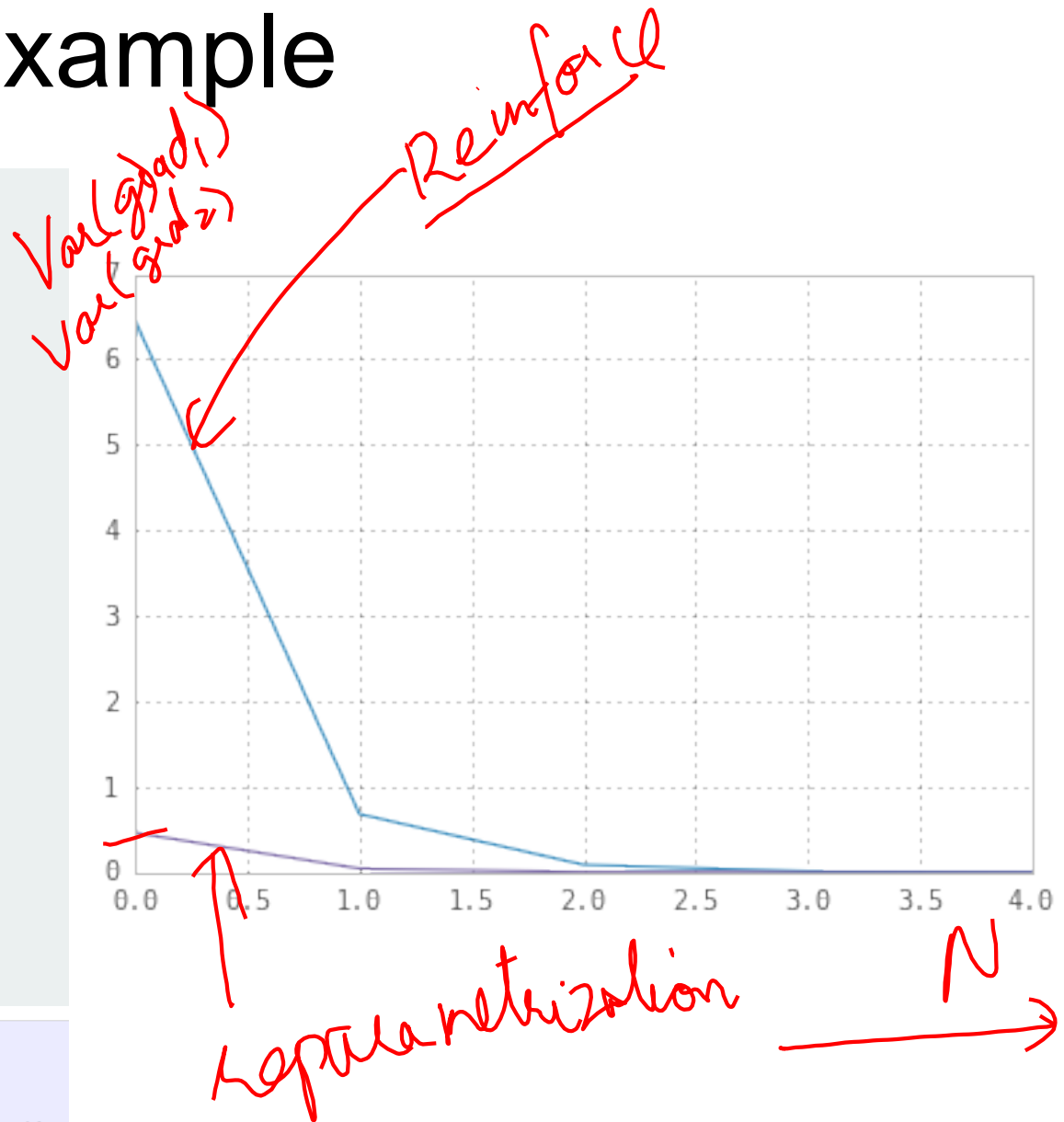
means1 = np.zeros(len(Ns))
vars1 = np.zeros(len(Ns))
means2 = np.zeros(len(Ns))
vars2 = np.zeros(len(Ns))

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:

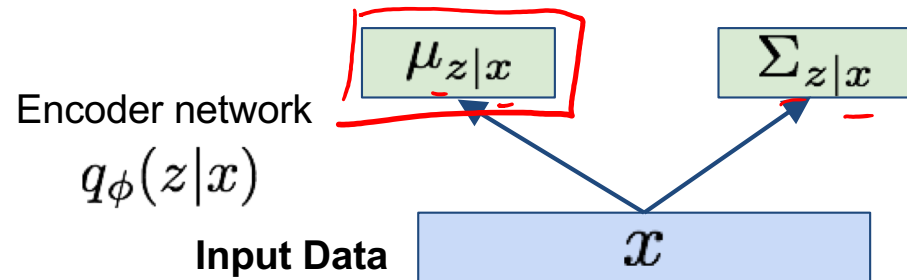
1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

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

Equal to!



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

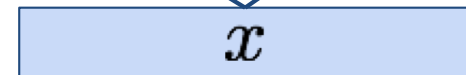
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\text{Make approximate posterior distribution close to prior}}$$

Make approximate posterior distribution close to prior

Encoder network

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

Input Data

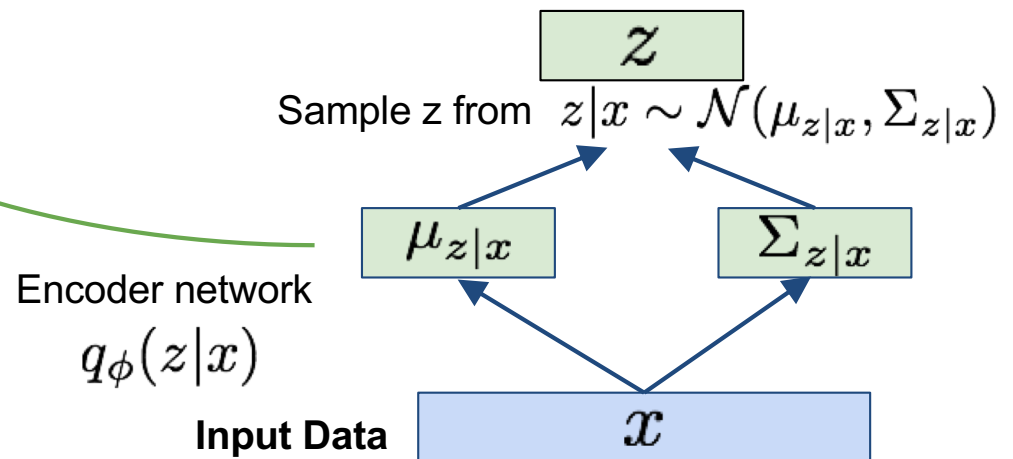


Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

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

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

