# CS 190I
# Deep Learning
# Variational Auto-Encoder

Lei Li (leili@cs)

UCSB

# Course Evaluation

- https://esci.id.ucsb.edu
- https://bit.ly/3FSqFs0
- Feedback is important and helpful for improving the course
- Encourage narrative comments
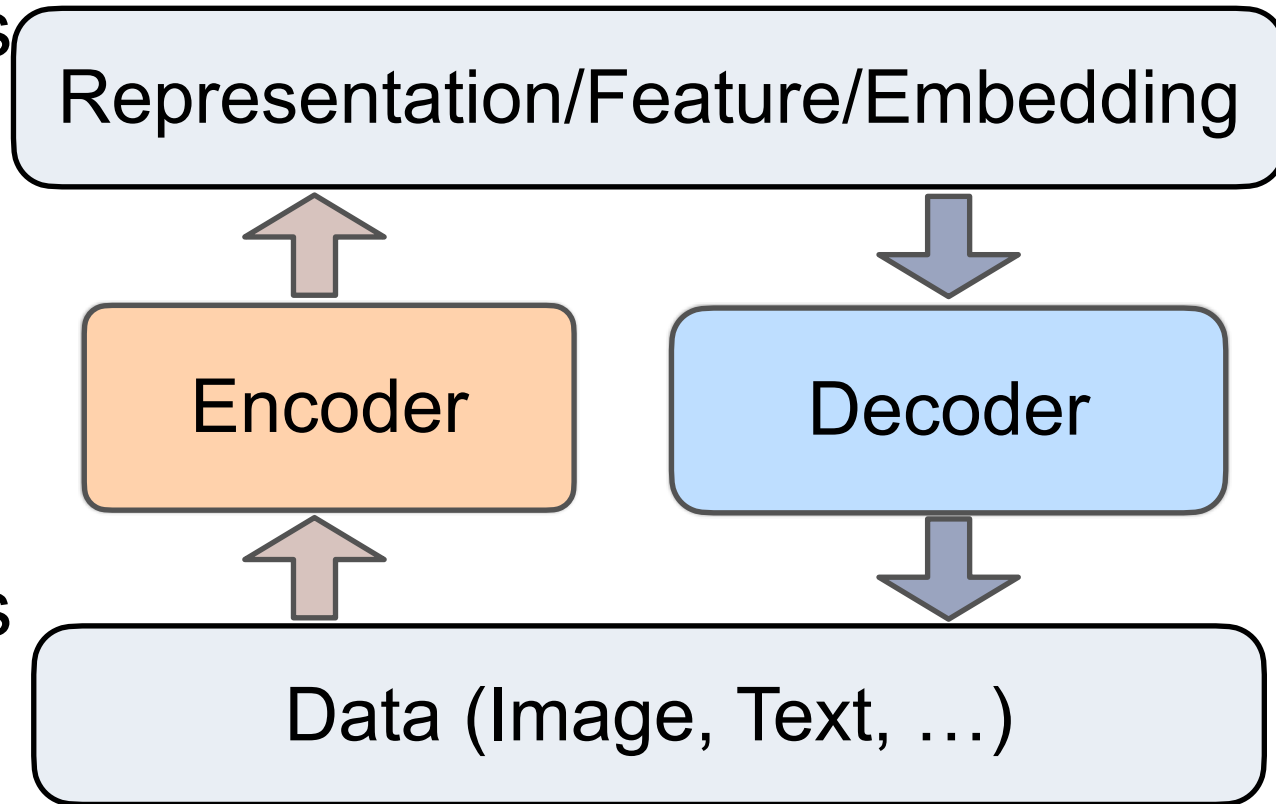- Bonus 5% to final exam, if response rate > 90% (48% today)

# Recap

- Graph neural network
  - message passed along graph edges
  - aggregate message/embedding by FFN
  - many variants: GCN, GAT, GraphSAGE
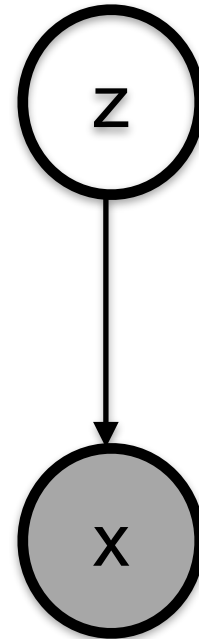
# Variational Auto-Encoder (VAE)

# VAE

- Hidden representations follow a prior distribution

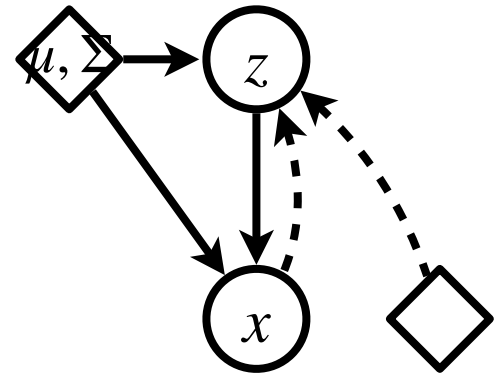- Encoder will produce a distribution of representations (posterior distribution)

Representation/Feature/Embedding

Encoder

Decoder

Data (Image, Text, …)

# Deep Latent Model

- z follows a prior distribution, e.g. Gaussian(0, I)

- p(x|z) is defined by a deep neural network $f(z; \theta)$

- To learn $\theta$, use $E_{(Z|X)}[\log p(X, Z; \theta)]$

# Graphical Model for VAE

- Assuming data X is generated from a latent variable $Z$

- Generation process
  - draw $Z \sim N(\mu, \Sigma)$
  - draw $X \mid Z \sim p(f(Z))$ , defined by a neural network f

- The goal is to maximize the data log-likelihood

$$\log p(X; \theta) = \log \int p(X \mid Z) p(Z) dZ$$

- Hard to optimize over $\theta$, if f(Z) is very complex such as a CNN, RNN, or Transformer.

# VAE

Objective: maximize the data loglikelihood

$$\max \ell(\theta) = \sum_n \log p(x_n; \theta)$$

$$= \sum_n \log \int p(x_n \mid z_n; \theta) p(z_n; \theta) dz_n$$

# VAE

$$\max \ell(\theta) = \sum_n \log p(x_n; \theta)$$

$$= \sum_n \log \int p(x_n \mid z_n; \theta) p(z_n; \theta) dz_n$$

- But $\log p(x; \theta)$ is intractable.

- For any distribution $q(z \mid x, \phi)$:

$q(z \mid x; \phi)$ is the posterior distribution from encoder!

$$\log p(x; \theta) \geq \mathrm{E}_{q(z \mid x; \phi)} \left[ \log \frac{p(x, z; \theta)}{q(z \mid x; \phi)} \right] = \mathrm{ELBO}$$

- Derivation via Jensen's inequality.
- Maximizing the ELBO instead of maximizing $\log p(x; \theta)$

# **Understanding ELBO**

$$\log p(X;\theta) \geq \mathrm{E}_q[\log \frac{p(X,Z;\theta)}{q(Z|X;\phi)}]$$

$$\max_{\theta} \max_{\phi} \mathrm{ELBO} = \sum_n \mathrm{E}_q\left[\log \frac{p(x_n|z_n;\theta)p_0(z_n)}{q(z_n|x_n;\phi)}\right]$$

$$=$$
$$=$$
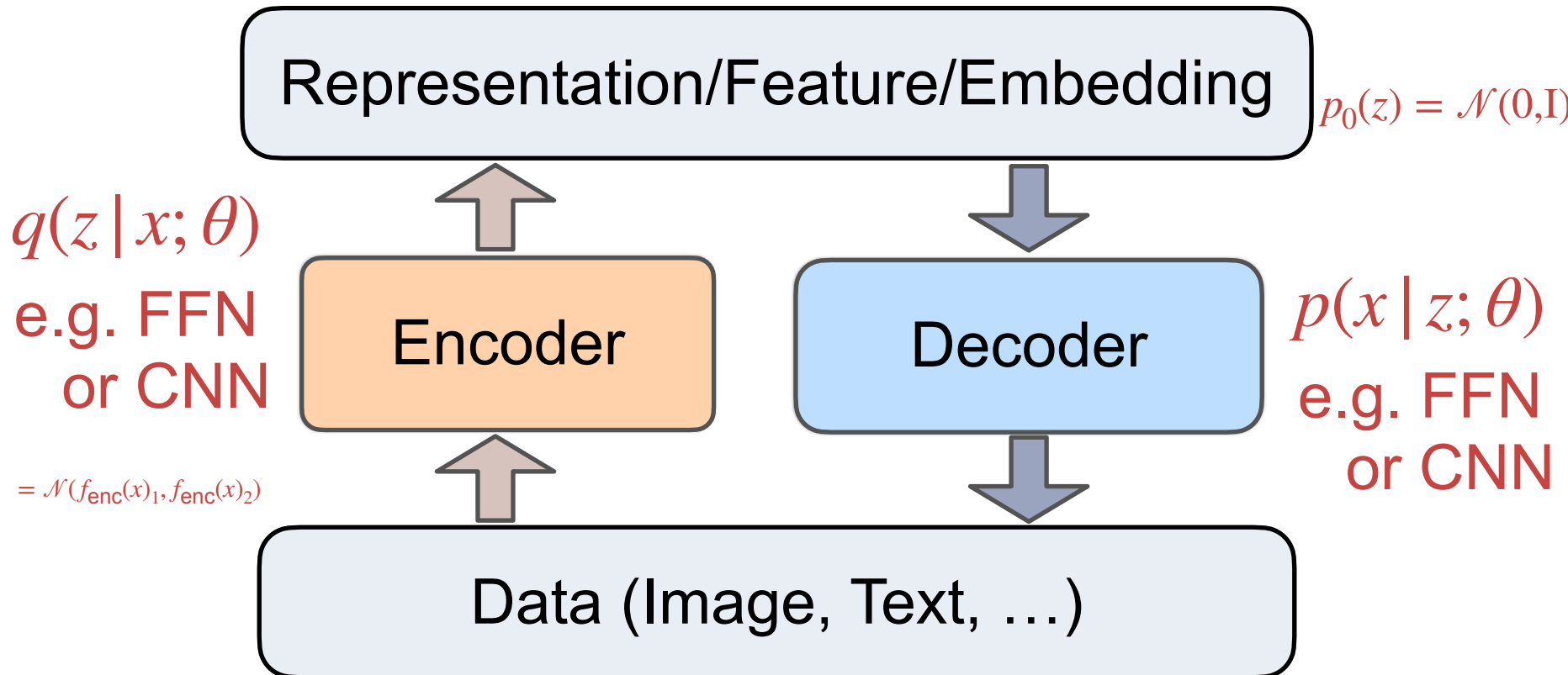
$$= \mathrm{E}_q\left[\log p(x_n|z_n;\theta)\right] - \mathrm{KL}\left(q(z_n|x_n;\phi)\|p_0(z_n)\right)$$

Reconstruction loss      Regularization

# VAE

Let $q(z|x;\phi)$ and $p(x|z;\theta)$ share the same parameter $\theta$

Representation/Feature/Embedding

$p_0(z) = \mathcal{N}(0, I)$

$q(z|x;\theta)$

e.g. FFN
or CNN

$= \mathcal{N}(f_{\text{enc}}(x)_1, f_{\text{enc}}(x)_2)$

Encoder

Decoder

$p(x|z;\theta)$

e.g. FFN
or CNN

Data (Image, Text, …)

# Training VAE

gradient descent(ascent for max)

$$\max_{\theta} \max_{\phi} \text{ELBO} = \sum_n \text{E}_{q(z_n|x_n;\theta)} \left[ \log \frac{p(x_n|z_n;\theta)p_0(z_n)}{q(z_n|x_n;\theta)} \right]$$

$$= \sum_n \text{E}_{q(z_n|x_n;\theta)} \left[ r(\theta, z_n, x_n) \right]$$

$$r(\theta, z_n, x_n) = \log \frac{p(x_n|z_n;\theta)p_0(z_n)}{q(z_n|x_n;\theta)}$$

Computing gradient:

$$\nabla_{\theta} \text{E}_{q(z_n|x_n;\theta)} \left[ r(\theta, z_n, x_n) \right]$$

# **Gradient of ELBO**

$$r(\theta, z_n, x_n) = \log \frac{p(x_n \mid z_n; \theta) p_0(z_n)}{q(z_n \mid x_n; \theta)}$$

Computing gradient:

$$\nabla_\theta \mathrm{E}_{q(z_n \mid x_n; \theta)} \left[ r(\theta, z_n, x_n) \right]$$

# Gradient of ELBO

$$r(\theta, z_n, x_n) = \log \frac{p(x_n \mid z_n; \theta) p_0(z_n)}{q(z_n \mid x_n; \theta)}$$

Computing gradient:

$$\nabla_\theta \mathrm{E}_{q(z_n|x_n;\theta)} \left[ r(\theta, z_n, x_n) \right] = \mathrm{E}_{q(z_n|x_n;\theta)} \left[ \nabla_\theta r(\theta, z_n, x_n) \right] + \int r(\theta, z_n, x_n) \nabla_\theta q(z_n \mid x_n; \theta) d_{z_n}$$

1. sample $z_n \sim q(z_n \mid x_n; \theta) = \mathcal{N}(f(x_n)_1, f(x_n)_2)$, then compute average of $\nabla_\theta r(\theta, z_n, x_n)$

# **Gradient of ELBO**

$$r(\theta, z_n, x_n) = \log \frac{p(x_n \mid z_n; \theta)p_0(z_n)}{q(z_n \mid x_n; \theta)}$$

Computing gradient:

$$\nabla_\theta \mathrm{E}_{q(z_n|x_n;\theta)}\left[r(\theta, z_n, x_n)\right] = \mathrm{E}_{q(z_n|x_n;\theta)}\left[\nabla_\theta r(\theta, z_n, x_n)\right] + \int r(\theta, z_n, x_n)\nabla_\theta q(z_n|x_n;\theta)d_{z_n}$$

2. rewrite as

$$\int r(\theta, z_n, x_n)\nabla_\theta q(z_n|x_n;\theta)d_{z_n} = \mathrm{E}_{q(z_n|x_n;\theta)}\left[r(\theta, z_n, x_n)\nabla_\theta \log q(z_n|x_n;\theta)\right]$$

then sample $z_n \sim q(z_n \mid x_n; \theta) = \mathcal{N}(f(x_n)_1, f(x_n)_2)$
compute average of $r(\theta, z_n, x_n)\nabla_\theta q(z_n \mid x_n; \theta)$

Problem — high variance

15

# **Reparameterization Trick**

$$q(z_n \mid x_n; \theta) = \mathcal{N}(f(x_n)_1, f(x_n)_2) = \mathcal{N}(\mu_\theta(x_n), \Sigma_\theta(x_n))$$

Treating $\epsilon \sim N(0,1)$, standard Gaussian distribution, then

$$\mathrm{E}_{q(z_n \mid x_n; \theta)} \left[ r(\theta, z_n, x_n) \right] = \mathrm{E}_{\epsilon \sim N(0,1)} \left[ r(\theta, z_n, x_n) \right]$$

where $z_n = \Sigma_\theta^{\frac{1}{2}}(x_n)\epsilon + \mu_\theta(x_n)$

Taking gradient does not depend on the distribution

# Reparameterization Trick

$$\nabla_\theta \mathrm{E}_{q(z_n|x_n;\theta)} \left[ \log \frac{p(x_n|z_n;\theta)p_0(z_n)}{q(z_n|x_n;\theta)} \right]$$

$$= \nabla_\theta \mathrm{E}_{q(z_n|x_n;\theta)} \left[ \log p(x_n|z_n;\theta) \right] - \mathrm{KL}\left( q(z_n|x_n;\phi) \| p_0(z_n) \right)$$

$$= \nabla_\theta \mathrm{E}_{\epsilon \sim N(0,1)} \left[ \log p(x_n|z_n;\theta) \right] - \mathrm{KL}\left( \mathcal{N}(\mu_\theta(x_n), \Sigma_\theta(x_n)) \| \mathcal{N}(0,1) \right)$$

$$= \nabla_\theta \mathrm{E}_{\epsilon \sim N(0,1)} \left[ \log p(x_n|z_n;\theta) \right] - \frac{1}{2} \left( \mu_\theta(x_n)^T \mu_\theta(x_n) + \mathrm{tr}(\Sigma_\theta)(x_n) - M - \mathrm{logDet}(\Sigma_\theta(x_n)) \right)$$

$$= \mathrm{E}_{\epsilon \sim N(0,1)} \left[ \nabla_\theta \log p(x_n|z_n;\theta) \right] - \nabla_\theta \frac{1}{2} \left( \mu_\theta(x_n)^T \mu_\theta(x_n) + \mathrm{tr}(\Sigma_\theta(x_n)) - M - \mathrm{logDet}(\Sigma_\theta(x_n)) \right)$$

where $z_n = \Sigma_\theta^{\frac{1}{2}}(x_n)\epsilon + \mu_\theta(x_n)$

# Compute Gradient using Reparameterization Trick

For each data point x_n, current parameter $\theta$

Step 1: sample $\epsilon \sim N(0,1)$

Step 2: using encoder forward to compute $\mu, \Sigma = f_{\mathrm{enc}}(x_n; \theta)$

Step 3: $z(\theta) = \Sigma^{\frac{1}{2}}\epsilon + \mu$

Step 4: using decoder forward to compute $p(x_n | z(\theta); \theta)$
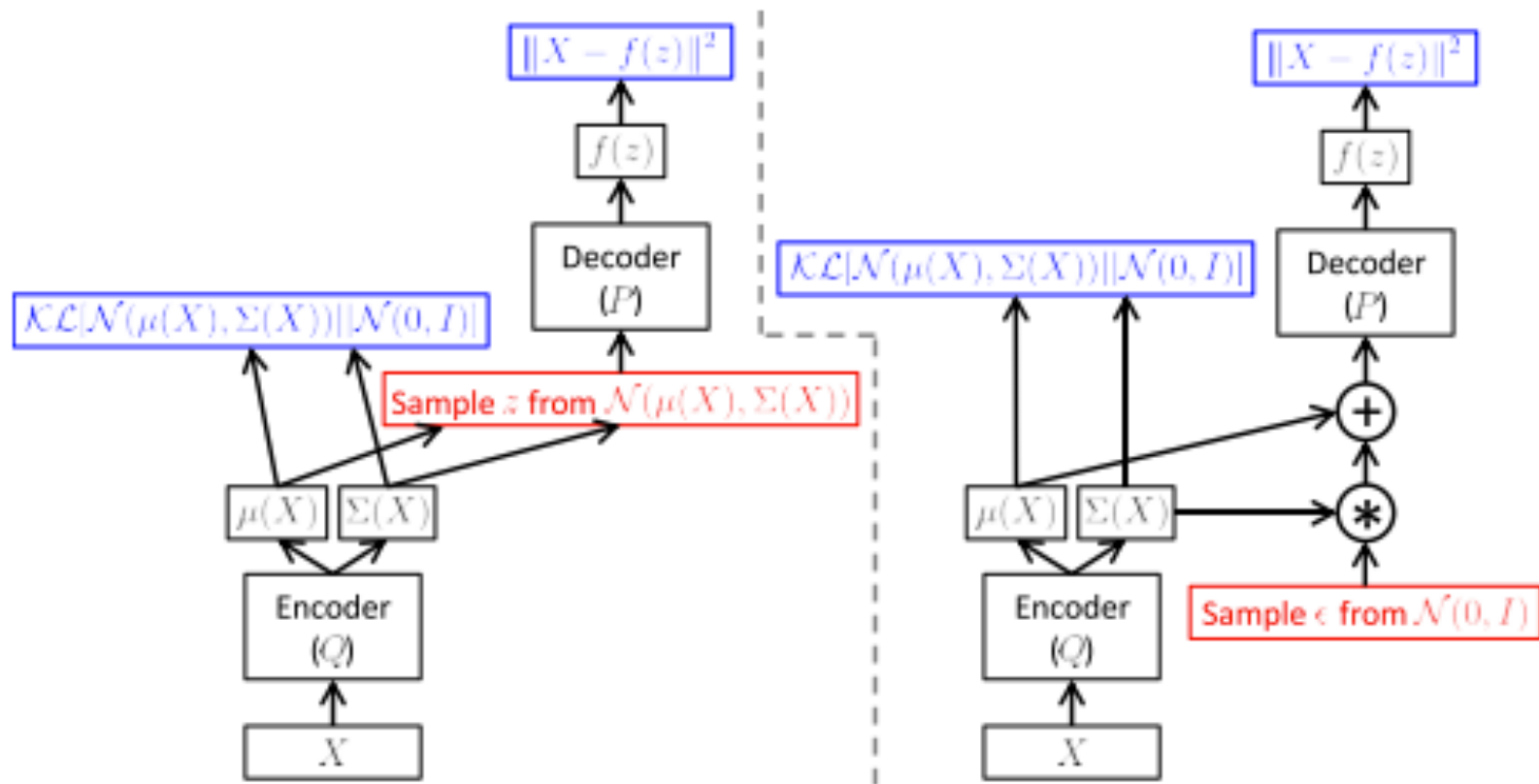
Step 5: define

$\mathrm{err} = \log p(x_n | z_n; \theta) - \beta \cdot \mathrm{KL}\left(q(z | x_n; \theta)\|p_0(z)\right)$ , then

using back-propagation to compute gradient for $\theta$

$$\frac{1}{2}\left(\mu_\theta(x_n)^T \mu_\theta(x_n) + \mathrm{tr}(\Sigma_\theta)(x_n) - M - \mathrm{logDet}(\Sigma_\theta(x_n))\right)$$

# Training VAE

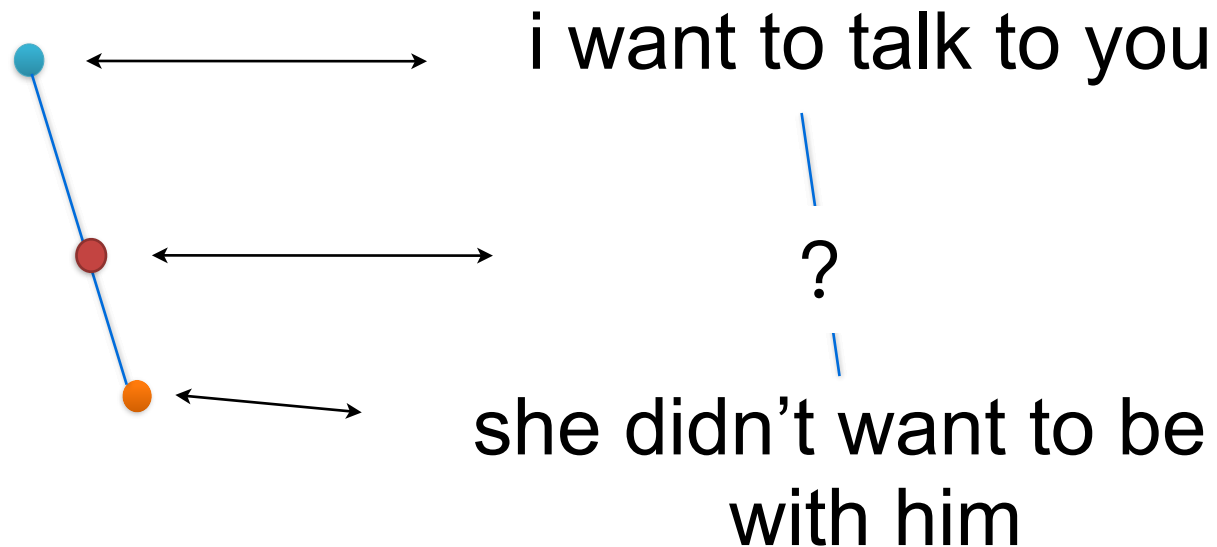- Reparameterization trick



*Tutorial on Variational Autoencoders*
(Doersch Carl, 2016)

# Sentence VAE

# Generating Sentence from Continuous vectors

- Key challenge: Interpolation in continuous space should yield reasonable sentences

i want to talk to you

?

she didn't want to be with him

# Conditional Sequence Generation

Given a latent variable z, a sequence of text tokens $x = (x_1, x_2, \ldots, x_t)$ can be generated with RNN (or LSTM, transformer), CRNN model:

$$p(x \mid z; \theta) = \prod_t p(x\_t \mid x_{<t}, z; \theta)$$

$$p(x_t \mid x_{<t}, z; \theta) = \text{softmax}(W \cdot h_t)$$

$$h_t = RNN(h_{t-1}, [x_{t-1}, \ z], \theta)$$

# VAE for Sentence Generation

Decoding:

$$z \sim N(0, I)$$

generate x from Transformer(z) or LSTM(z)

Encoding:
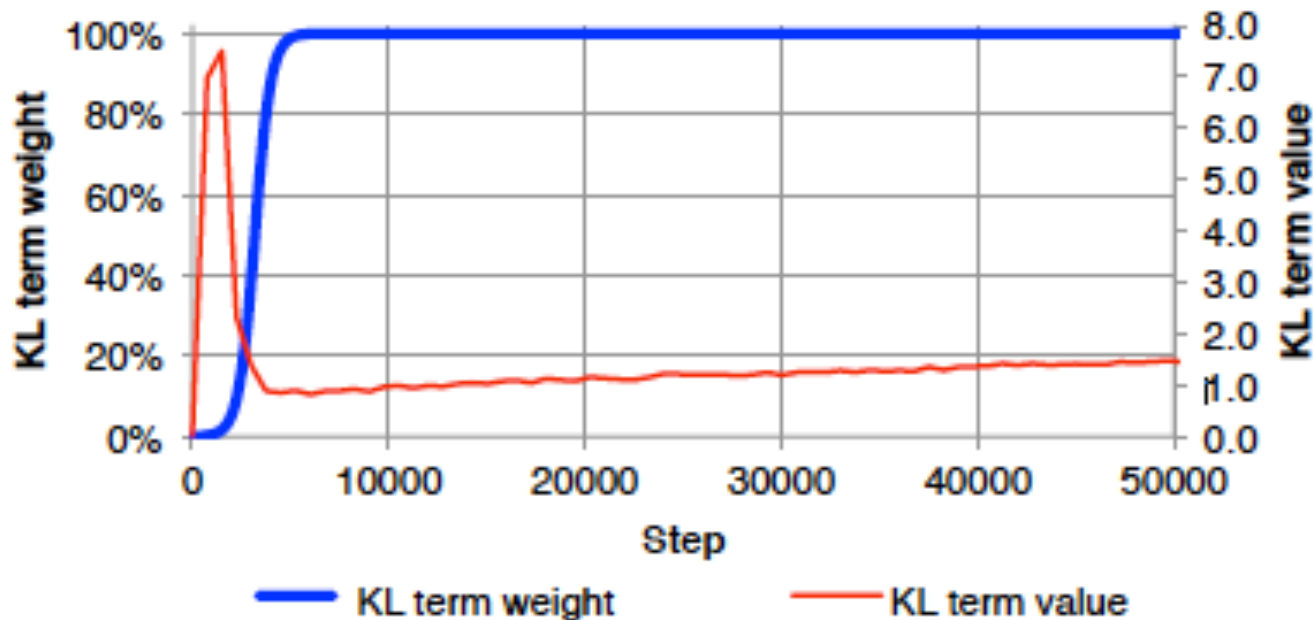
$$q(z \mid x) = N\left(\mu, \sigma^2\right)$$

$$\mu = W_1 \cdot h_t, \ \sigma^2 = \exp W_2 \cdot h_t$$

$$h_t = \text{Transformer}(x; \theta)$$



**Encoder**

Avg

Enco Enco Enco Enco

Enco Enco Enco Enco

0 1 2 3

Tok Tok Tok Tok

A cat sits on a mat

**Decoder**

A cat sits on a m

Softm Softm Softm

Deco Deco Deco

Deco Deco Deco

Deco Deco Deco

<BOS $y_1$ $y_2$

# Training VAE: Posterior Collapse

- KL term in ELBO collapses to zero and latent variable encodes little information.

- Solution: KL annealing & word dropout

# Examples on Sentence Interpolation

" i want to talk to you . "
"i want to be with you . "
"i do n't want to be with you . "
i do n't want to be with you .
she did n't want to be with him .

he was silent for a long moment .
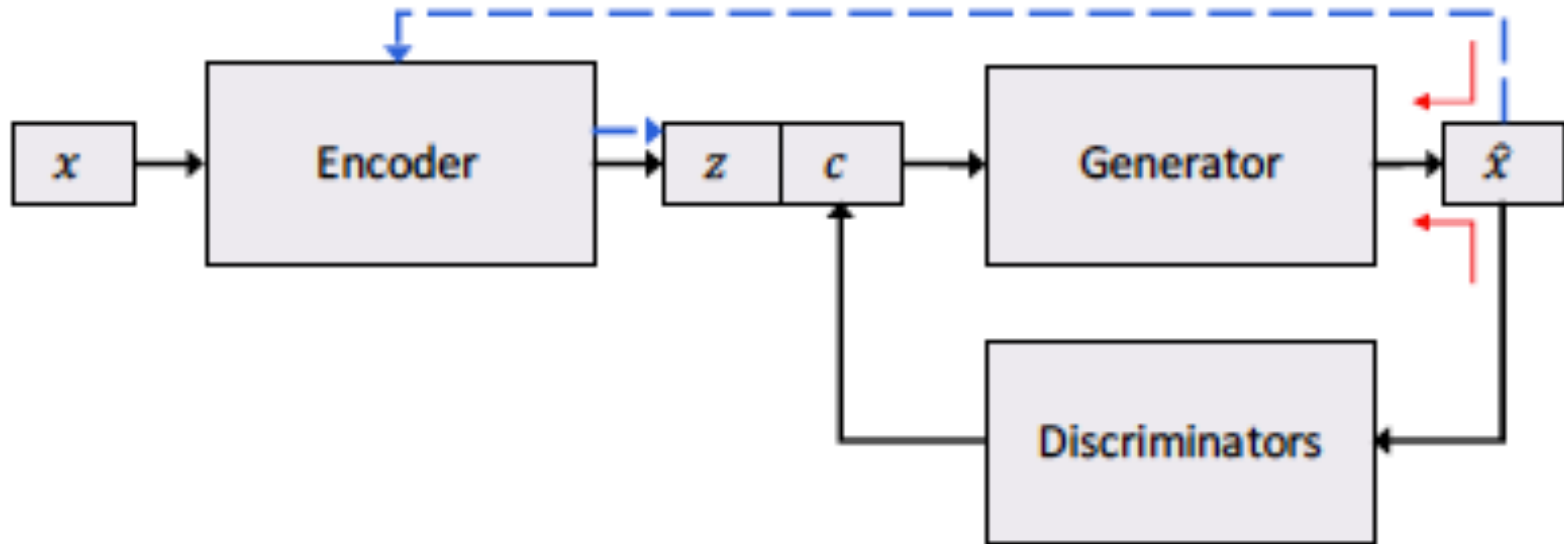he was silent for a moment .
it was quiet for a moment .
it was dark and cold .
there was a pause .
it was my turn .

# Variants

- Controllable sentence generation with both continuous and discrete labels



Toward Controlled Generation of Text, (Hu et. al. ICML 2017)

# Generating with Varying Semantic Label

the film is strictly routine !
the film is full of imagination .

after watching this movie , i felt that disappointed .
after seeing this film , i 'm a fan .

the acting is uniformly bad either .
the performances are uniformly good .

this is just awful .
this is pure genius .

the acting is bad .
the movie is so much fun .

none of this is very original .
highly recommended viewing for its courage , and ideas .

too bland
highly watchable

i can analyze this movie without more than three words .
i highly recommend this film to anyone who appreciates music .

Toward Controlled Generation of Text, (Hu et. al. ICML 2017)

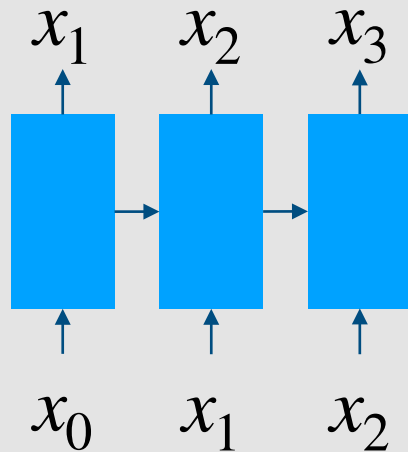# Deep Latent Variable Models for Text

- Interpretable Deep Latent Representation from Raw Text
  - Learning Exponential Family Mixture VAE [ICML 20]
- Disentangled Representation Learning for Text Generation
  - Data to Generation: VTM [ICLR 20b]
  - Learning syntax-semantic representation [ACL 19c]
- One model to acquire 4 language skills
  - Mirror Generative NMT [ICLR 20a]

# Learning Interpretable Latent Representation

Latent structure
dialog actions

**GENERATOR**

$x_1$  $x_2$  $x_3$

$x_0$  $x_1$  $x_2$

Sampling

"Remind me about the football game."
[action=remind]

"Will it be overcast tomorrow?"
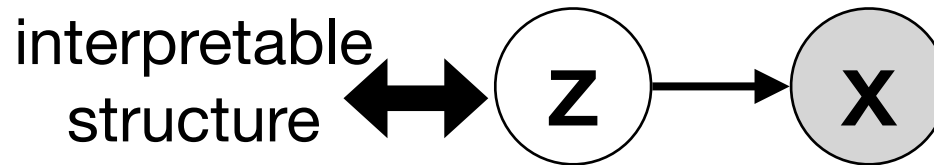[action=request]

......

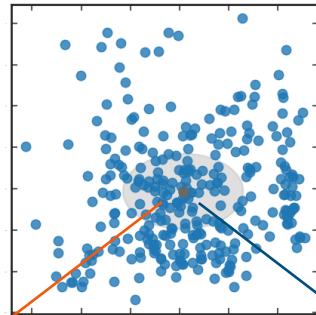Generate Sentences with interpretable factors

# How to Interpret Latent Variables in VAEs?

**Variational Auto-encoder (VAE)**

interpretable structure ⬌ ( **Z** ) → ( **X** )

(Kingma & Welling, 2013)

*z*:
**continuous** latent variables

difficult to interpret discrete factors

Will it be humid in New York today?

Remind me about my meeting.

# VAEs Introduce Latent Variables

**Variational Auto-encoder (VAE)**

interpretable structure ⬌ **C** → **X**

(Zhao et al, 2018b)

$c$: **discrete** latent variables

expressiveness is limited.

Remind me about my meeting.
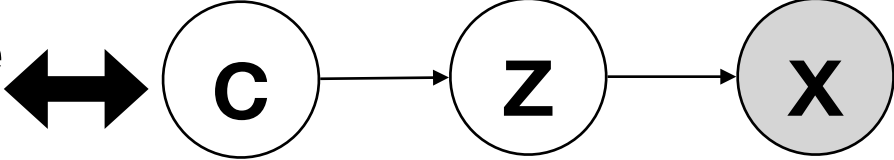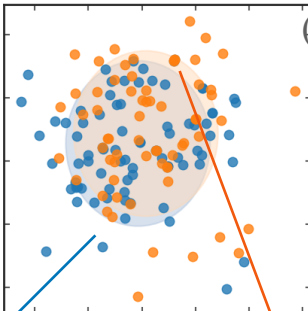
Remind me about the football game.

# Discrete Variables Could Enhance Interpretability - but one has to do it right!

## Gaussian Mixture Variational Auto-encoder (GM-VAE)

**interpretable structure** ⬛↔⬛ ( C ) → ( Z ) → ( X )

(Dilokthanakul et al., 2016; Jiang et al., 2017)

$c$: **discrete** component

$z$: **continuous** latent variable

Why?

How to fix it?

mode-collapse

Will it be overcast tomorrow?
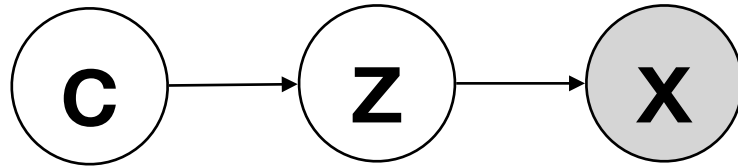
Remind me about the football game.

# Do it right for VAE w/ hierarchical priors - Dispersed Exponential-family Mixture VAE

Exponential-family Mixture VAE



adding dispersion term in training

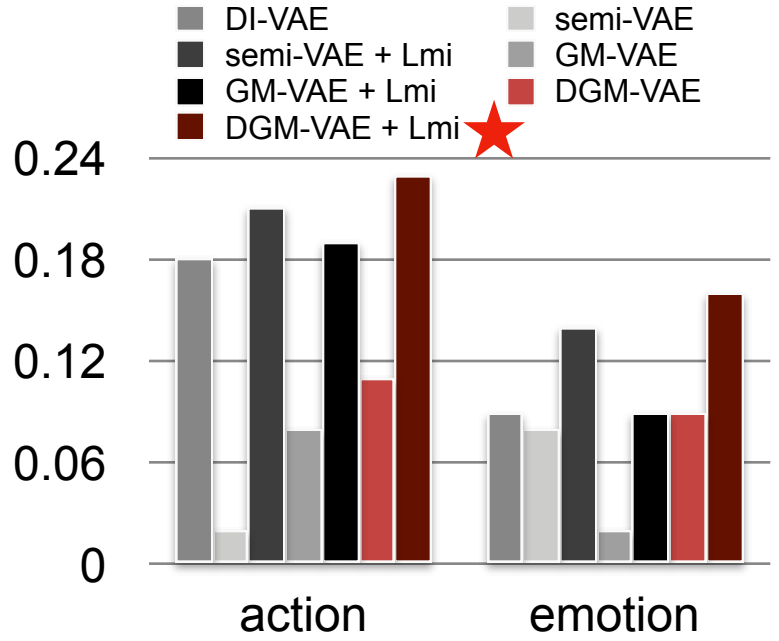## Dispersed EM-VAE

$$L(\theta; x) = \text{ELBO} + \beta \cdot \boxed{L_d,}$$

dispersion term
$$L_d = \mathbb{E}_{q_\phi(c|x)} A(\boldsymbol{\eta}_c) - A(\mathbb{E}_{q_\phi(c|x)} \boldsymbol{\eta}_c) \,.$$

33

DEM-VAE [W. Shi, H. Zhou, N. Miao, **Lei Li**, ICML 2020]
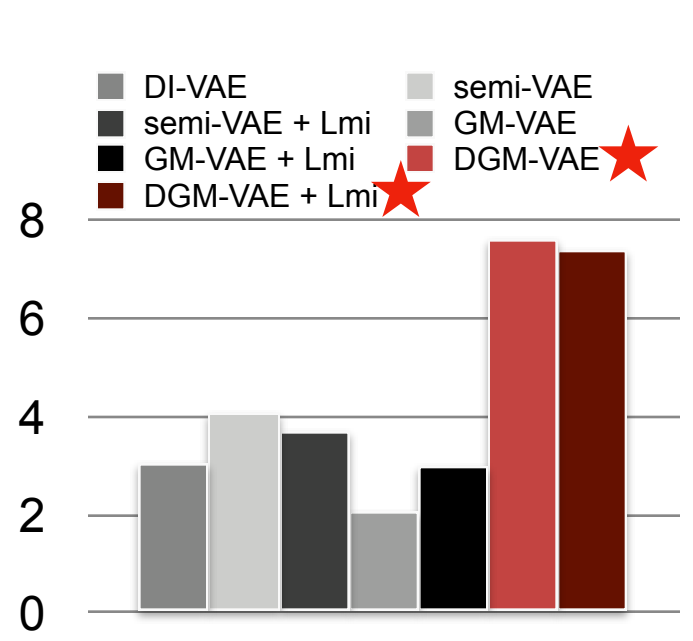
# Generation Quality and Interpretability

DGM-VAE obtains the best performance in interpretability and reconstruction

Homogeneity with golden label in DD



Best interpretability

BLEU of reconstruction in DD



Best reconstruction

DEM-VAE [W. Shi, H. Zhou, N. Miao, **Lei Li**, ICML 2020]

# Latent Variables Learned by DEM-VAE are Semantically Meaningful

Example actions and corresponding utterances (classified by $q_\phi(c|x)$)
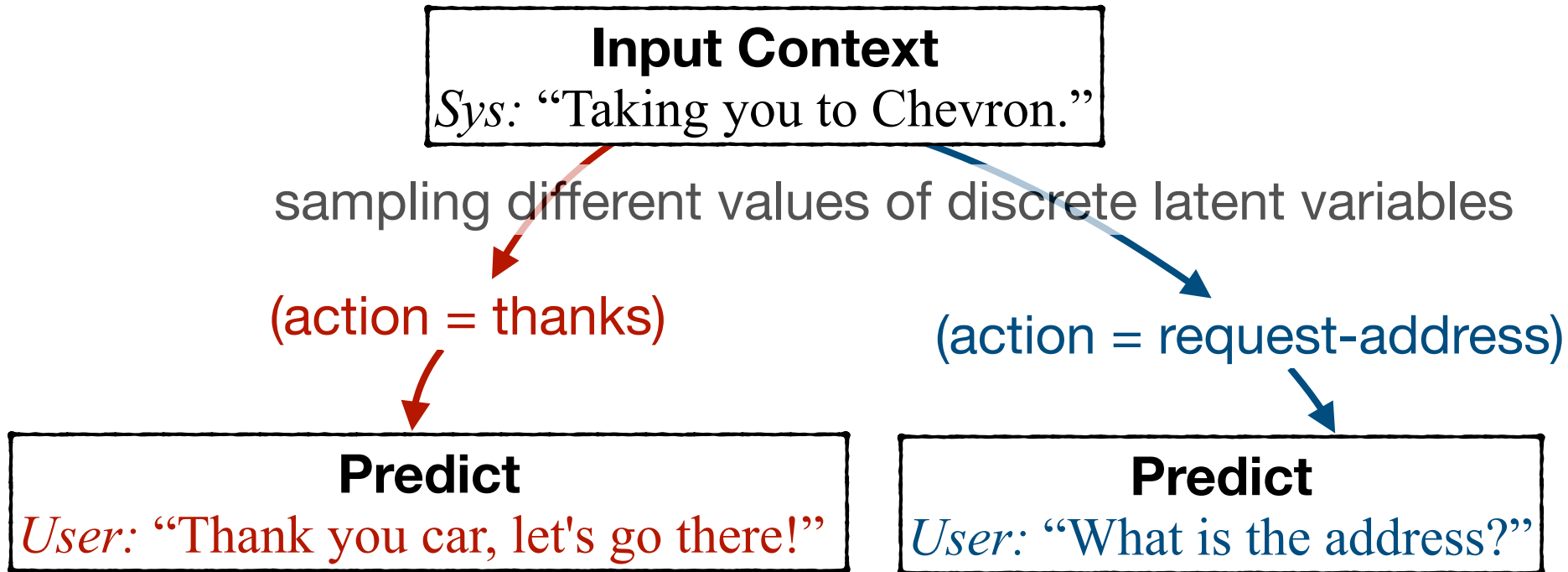
**Inferred action=Inform-route/address**
"There is a Safeway 4 miles away."
"There are no hospitals within 2 miles."
"There is Jing Jing and PF Changs."
…

**Inferred action =Request-weather**
"What is the weather today?"
"What is the weather like in the city?"
"What's the weather forecast in New York?"
…

Utterances of the same actions could be assigned with the same discrete latent variable $c$.

# Generate Sensible Dialog Response with DEM-VAE

Input Context
*Sys:* "Taking you to Chevron."

sampling different values of discrete latent variables

(action = thanks)

(action = request-address)

Predict
*User:* "Thank you car, let's go there!"

Predict
*User:* "What is the address?"

Responses with different actions are generated by sampling different values of discrete latent variables.

DEM-VAE [W. Shi, H. Zhou, N. Miao, **Lei Li**, ICML 2020]

# Summary

- Auto-Encoder: learning representation by reconstruction

- Variational Auto-Encoder: put prior on latent representation and use variational method to train

# Next Up

- Industrial talk on computer vision