CS769 Advanced NLP

Latent Variable Models

Junjie Hu



Slides adapted from Graham https://junjiehu.github.io/cs769-fall25/

Goal for Today

- Variational Auto-encoder (VAE)
 - 1. Maximize Evidence Lower Bound (ELBO)
 - 2. Reparamertizction trick
 - 3. Stabilizing Training
- Discrete Latent Variable Models (LVM)
 - 1. Enumerate
 - 2. Sampling
 - 3. Gumbel-softmax
- LVM Applications in NLP

Discriminative vs. Generative Models

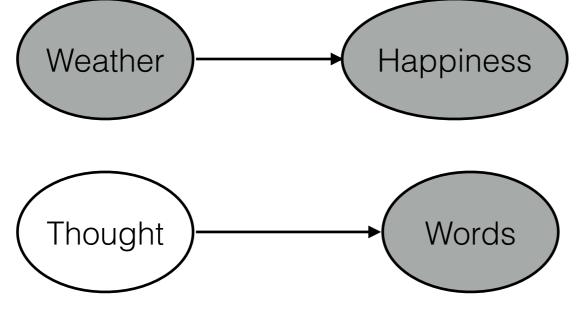
- **Discriminative model:** calculate the probability of output given input P(Y|X)
- Generative model: calculate the probability of a variable P(X), or multiple variables P(X,Y)
- Which of the following models are discriminative vs. generative?
 - Standard BiLSTM POS tagger
 - Language model

Types of Variables

- Random vs. deterministic variables:
 - Random variable (R.V.): is a math formalization of a object which depends on a random event. Namely, it defines a function that map an event to a numerical value. We usually use uppercase letter X, Y to denote R.V.. Once the event has been sampled, we get a real value of the R.V., we call this realization of a R.V., and usually use lowercase letter x, y to denote them.
 - Deterministic variable: a variable that does not depend on any random event.
 We typically use a, b, c to denote them.
- Example:
 - X denotes a R.V. of the weather, X can be realized by taking values from a set (domain) {"rainy", "sunny", "cloudy",}
 - P(X = "rainy") provides a probability of a realization of the R.V. (i.e., a particular event)

Types of Variables

- Observed vs. Latent:
 - Observed: something that we can observe and measure from our data, e.g.
 X or Y
 - Latent: a variable that we assume exists, but we aren't given the value, namely something we don't have data for
- In graphical models, observed variables are typically shaded nodes, and latent variables are typically unshaded nodes.



Marginal over latent variables is usually intractable (not computable in a reason time).

$$p(w) = \int_{T} p(T, w) dT$$

Latent Variable Models

- · A vector of latent variables ${\bf z}$ in a high-dim space ${\cal Z}$ which we can sample from some PDF $p({\bf z})$ over ${\cal Z}$
- For every observed x in our dataset, there is a z that causes the model to generate x
- A latent variable model (LVM) is a probability distribution over two sets of variables x, z:

$$p(\mathbf{x}|\mathbf{z};\theta)$$

Maximum Likelihood Framework

 Maximize the probability of each x in the training set under the generative process according to:

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z};\theta)p(\mathbf{z})d\mathbf{z}$$

Why LVM?

- Intuitively, latent variable z enables the model to first decide which property to generate before it assigns values to the output x
- Example:
 - Sample a LV z from the set [0,...9] before generating an image for the digit

 Sample a sentiment from {positive, negative} before generating a positive/negative movie review.

z="negative" -> x= "This is a bad movie."

What Types of Latent Variables?

- Latent continuous vector (e.g. variational autoencoder)
- Latent discrete vector (e.g. topic model)
- Latent structure (e.g. HMM or tree-structured model)

Variational Auto-encoders

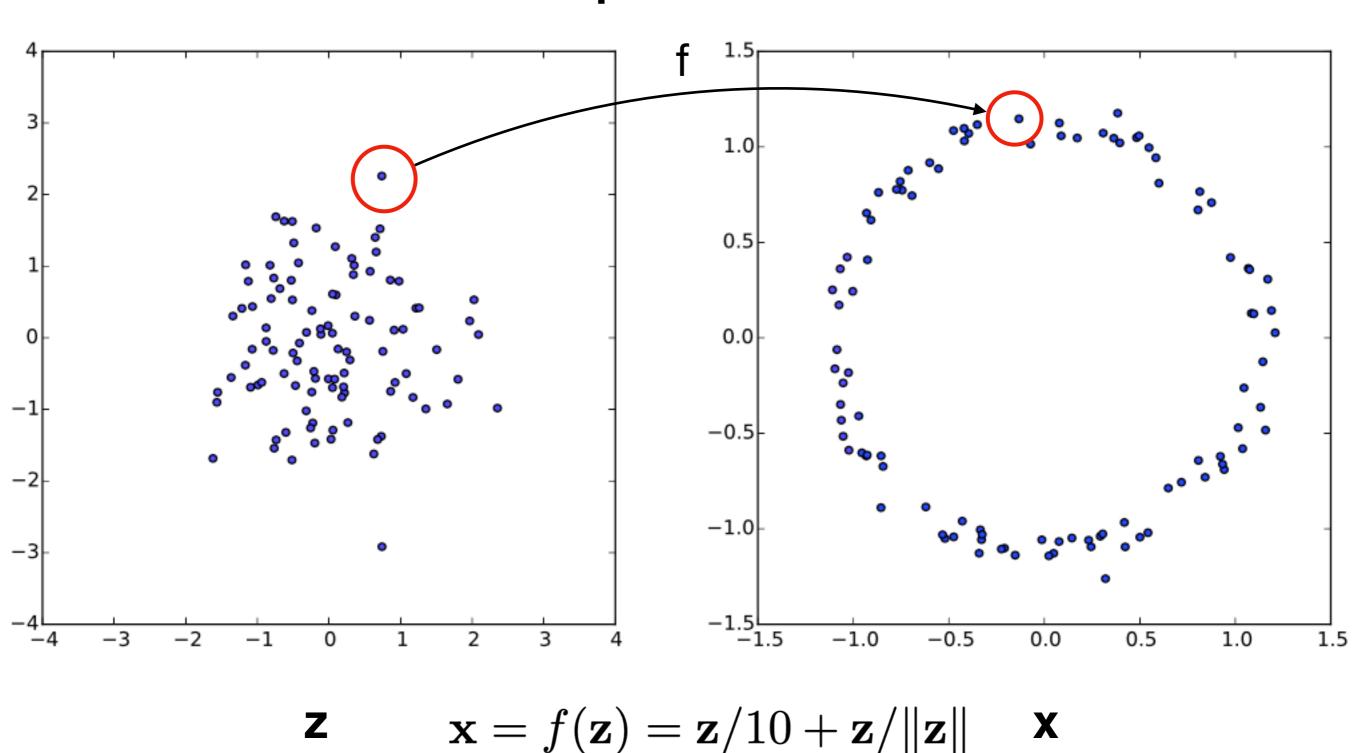
(Kingma and Welling 2014)

VAE as a Graphical Model

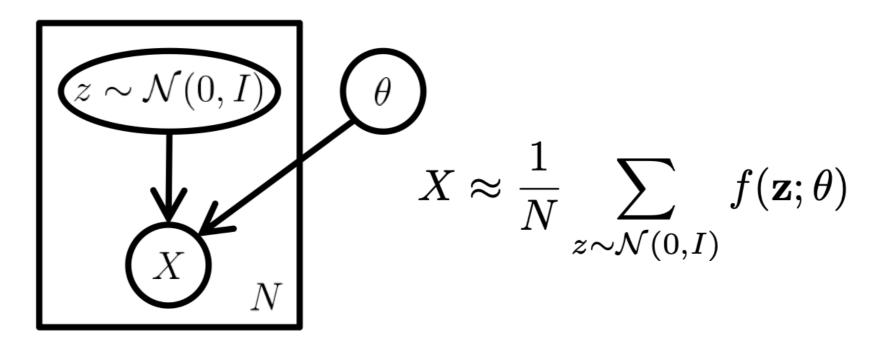
- We have a random variable X for our dataset
- We have a latent variable z sampled from a Gaussian
- We have a deterministic function $f(\mathbf{z}; \theta)$ that maps \mathbf{z} to the data space \mathcal{X} . If \mathbf{z} is a random variable (r.v.) in \mathcal{Z} , then $f(\mathbf{z}; \theta)$ is also a r.v. in \mathcal{X}
- If we repeat this sampling N times, we hope to approximate X by $f(\mathbf{z};\theta)$

$$X \approx \frac{1}{N} \sum_{z \sim \mathcal{N}(0,I)} f(\mathbf{z};\theta)$$

An Example (Goersch 2016)



A probabilistic perspective on Variational Auto-Encoder



- For each datapoint i:
 - Draw latent variables $z_i \sim p(z)$ (prior)
 - Draw data point $x_i \sim p_{\theta}(x|z)$
- Joint probability distribution over data and latent variables:

$$p(x,z) = p(z)p_{\theta}(x|z)$$

Similar to Naive Bayes, HMM: only the prior differs

What is Our Loss Function?

We would like to maximize the corpus log likelihood

$$\log P(\mathcal{X}) = \sum_{\boldsymbol{x} \in \mathcal{X}} \log P(\boldsymbol{x}; \theta)$$

For a single example, the marginal likelihood is

$$P(\boldsymbol{x}; \theta) = \int P(\boldsymbol{x} \mid \boldsymbol{z}; \theta) P(\boldsymbol{z}) d\boldsymbol{z}$$

We can approximate this by sampling zs then summing

$$P(\boldsymbol{x}; \theta) pprox \sum_{\boldsymbol{z} \in S(\boldsymbol{x})} P(\boldsymbol{x}|\boldsymbol{z}; \theta)$$
 where $S(\boldsymbol{x}) := \{\boldsymbol{z}'; \boldsymbol{z}' \sim P(\boldsymbol{z})\}$

Two tasks of interest:

- Learn the parameters θ of $p_{\theta}(x|z)$
- Inference over z with the *posterior* distribution: $p_{\theta}(z|x)$ given input x, what are its latent factors?

$$p_{\theta}(z|x) = \frac{p_{\theta}(x|z)p(z)}{p(x)}$$

$$p(x) = \int p(z)p_{\theta}(x|z)dz$$
 <- intractable

- However, as p(x) is intractable $p_{\theta}(z|x)$ is also intractable.
- Solution: Variational inference approximates the posterior $p_{\theta}(z|x)$ with a family of distributions $q_{\phi}(z|x)$

• Variational inference approximates the true posterior $p_{\theta}(z|x)$ with a family of distributions $q_{\phi}(z|x)$

minimize:
$$\mathrm{KL}(q_{\phi}(z|x)||p_{\theta}(z|x))$$

 One other target is to maximize the log-data likelihood which can be rewritten as:

$$\log p(x) = \text{ELBO} + \text{KL}(q_{\phi}(z|x)||p_{\theta}(z|x))$$

 Combining the above two, this is equivalent to maximize the Evidence Lower Bound (ELBO)

ELBO =
$$\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - \text{KL}(q_{\phi}(z|x)||p(z))$$

 $\text{KL}(q||p) \ge 0 => \log p(x) \ge \text{ELBO}$

• Variational inference approximates the true posterior $p_{\theta}(z|x)$ with a family of distributions $q_{\phi}(z|x)$

minimize:
$$\mathrm{KL}(q_{\phi}(z|x)||p_{\theta}(z|x))$$

 One other target is to maximize the log-data likelihood which can be rewritten as:

$$\log p(x) = \text{ELBO} + \text{KL}(q_{\phi}(z|x)||p_{\theta}(z|x))$$

Evidence Lower Bound (ELBO)

ELBO =
$$\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - \text{KL}(q_{\phi}(z|x)||p(z))$$

 $\text{KL}(q||p) \ge 0 => \log p(x) \ge \text{ELBO}$

• Variational inference approximates the true posterior $p_{\theta}(z|x)$ with a family of distributions $q_{\phi}(z|x)$

minimize: $\mathrm{KL}(q_{\phi}(z|x)||p_{\theta}(z|x))$

maximize: $\log p(x) = \text{ELBO} + \text{KL}(q_{\phi}(z|x)||p_{\theta}(z|x))$



to maximize ELBO





maximize: ELBO

Variational Auto-Encoders

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) \geq \text{ELBO}$$

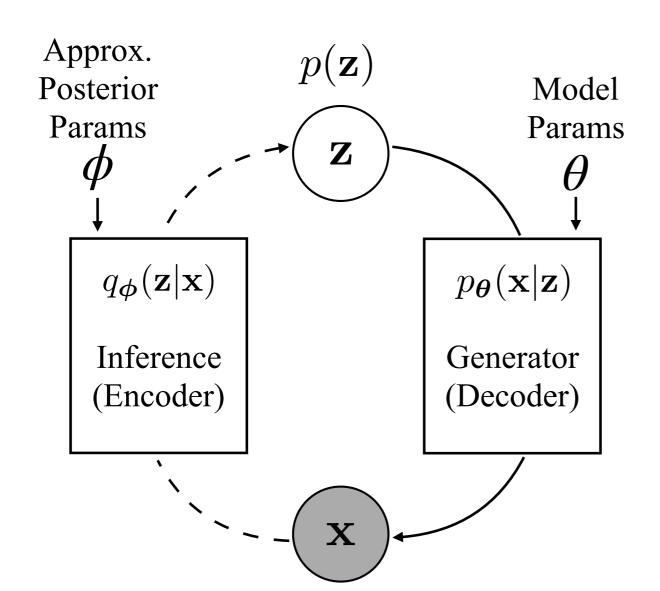
$$\underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{KL Regularizer}}$$

The inequality holds for any $q_{\phi}(z|x)$, but the lower bound is tight only if $p_{\theta}(z|x) = q_{\phi}(z|x)$

 $p_{\theta}(z|x)$ is intractable

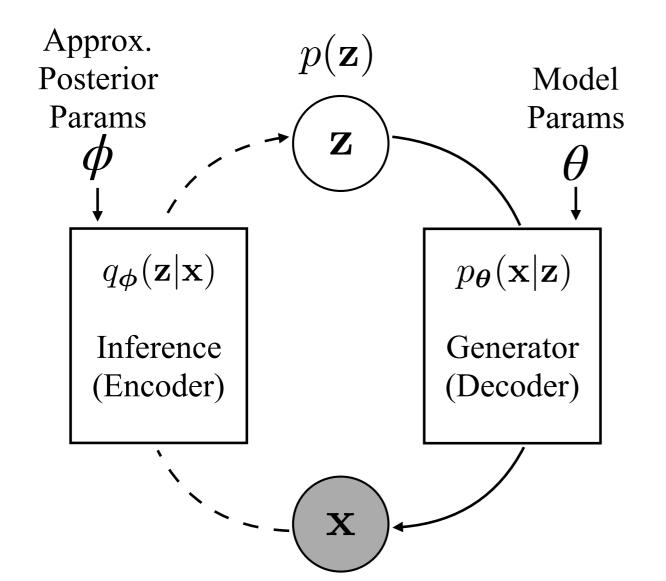
Variational Autoencoders

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) >= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{KL Regularizer}}$$



Variational Autoencoders

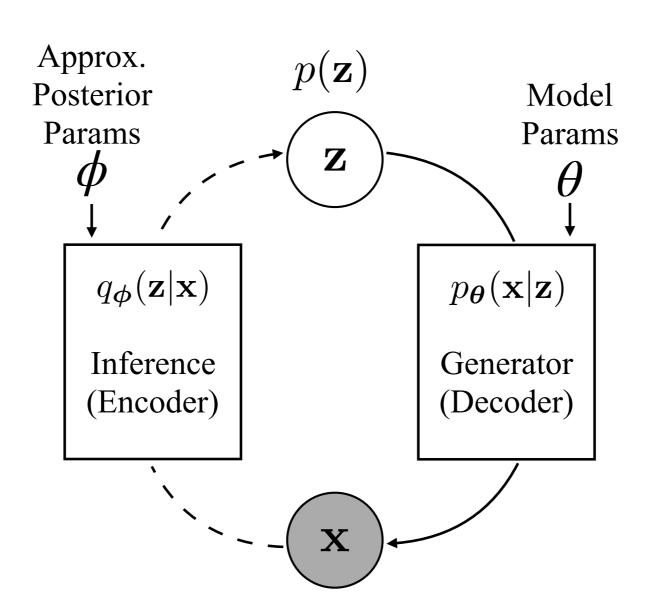
$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) >= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{KL Regularizer}}$$



Regularized Autoencoder

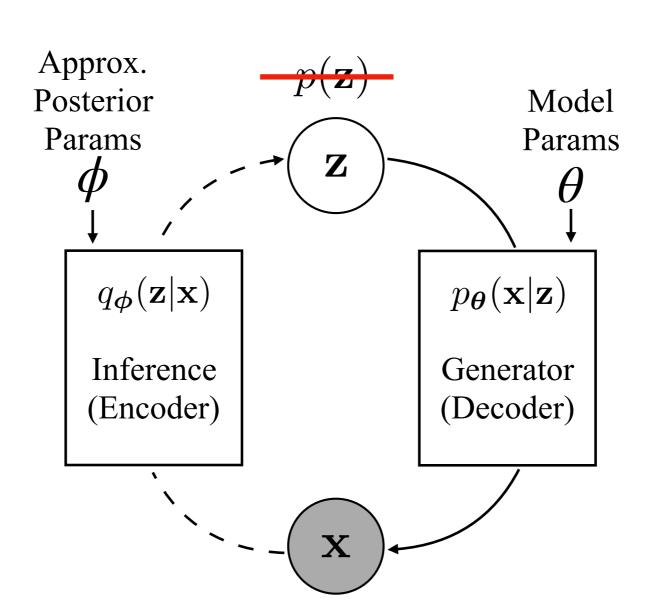
Why prior?

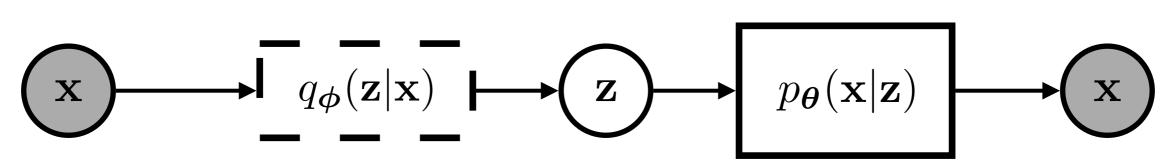
$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) >= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{KL Regularizer}}$$

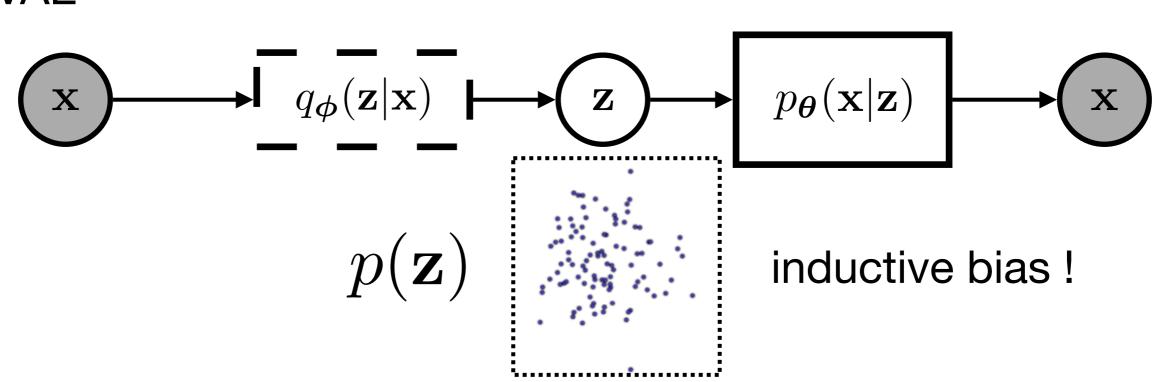


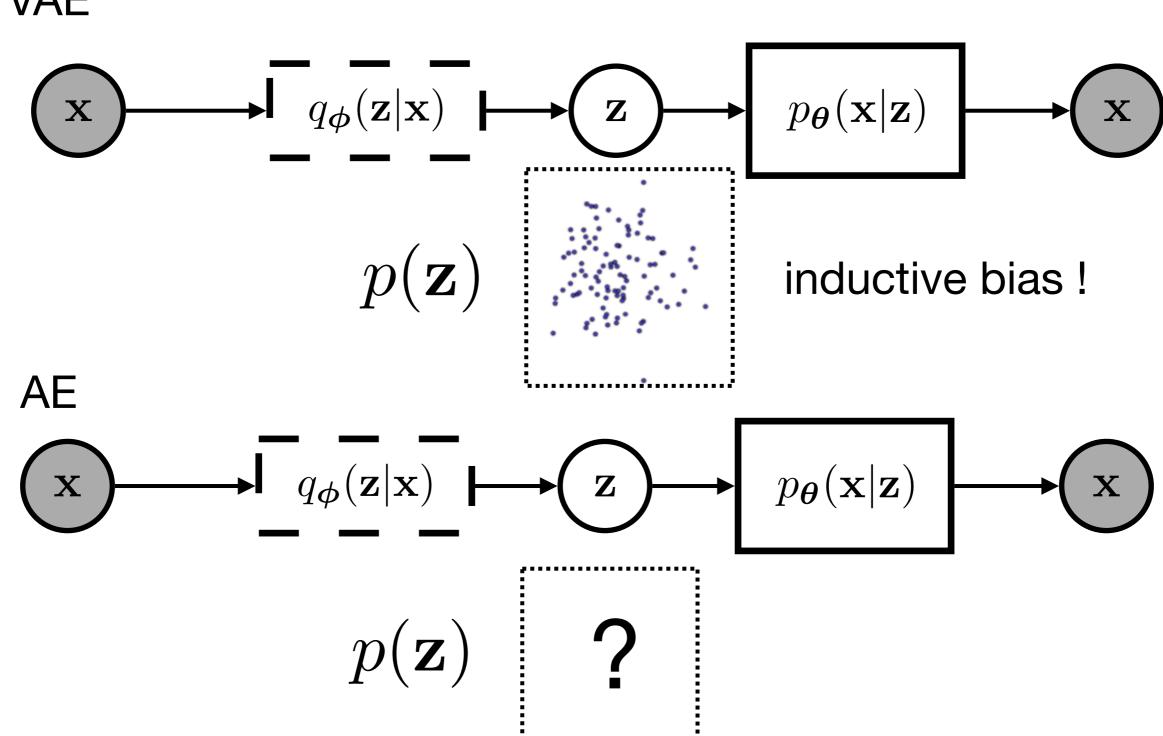
Why prior?

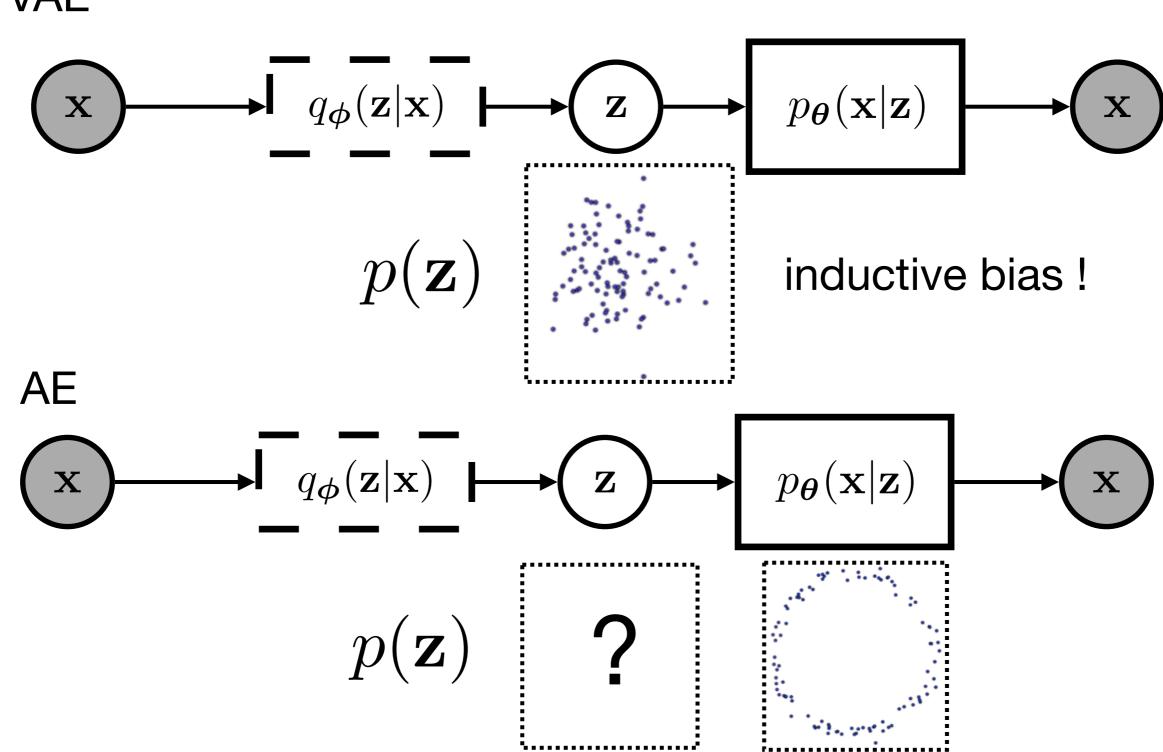
$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) >= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{KL Regularizer}}$$











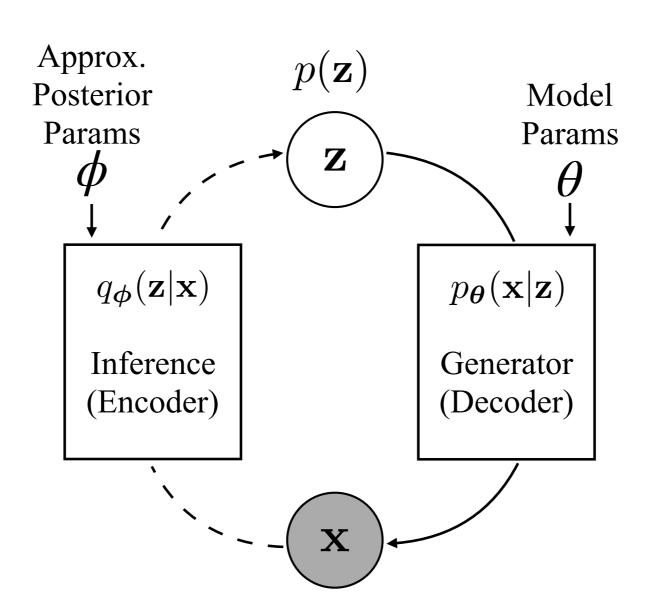
VAE $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ $q_{oldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$ $p(\mathbf{z})$ inductive bias! AE $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$ $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$

AE is not generative model:

- (1) Can't sample new data from AE
- (2) Can't compute the log likelihood of data x

Training of VAE

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) >= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Loss}} - \underbrace{D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{KL Regularizer}}$$



Problem! Sampling Breaks Backprop

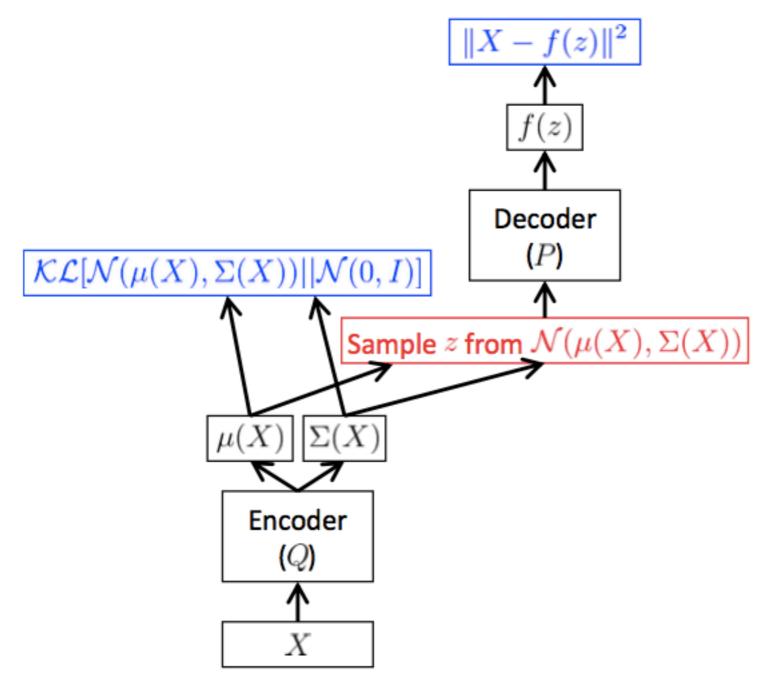
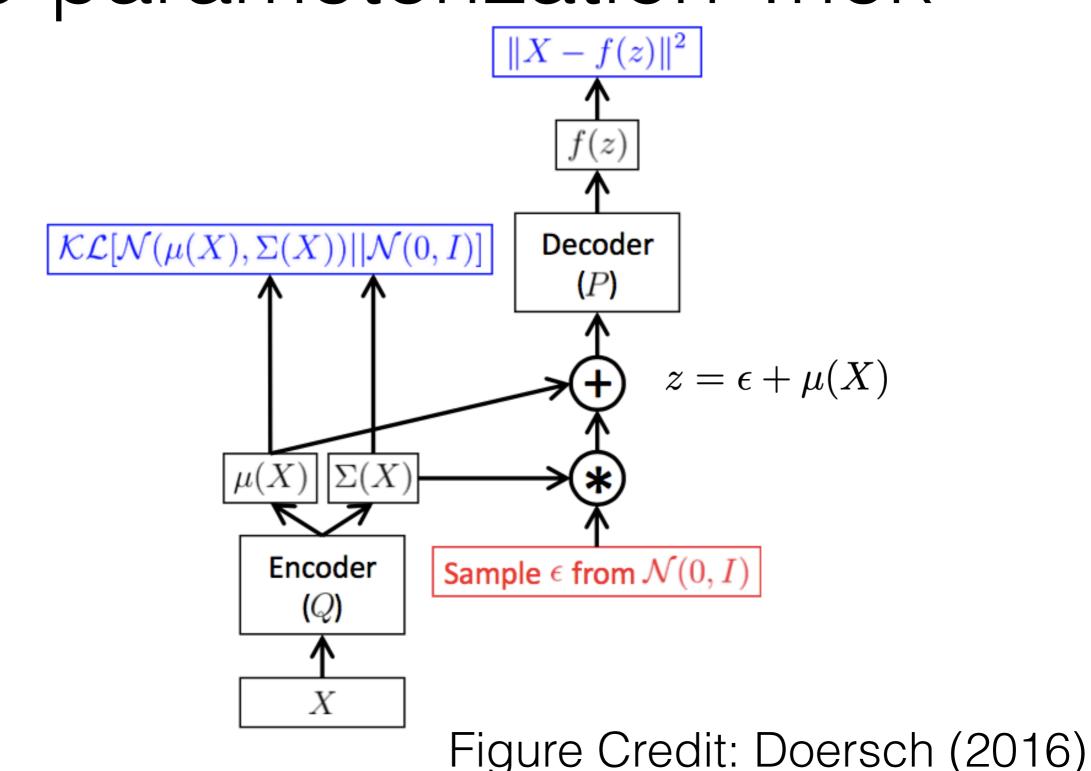


Figure Credit: Doersch (2016)

Solution: Re-parameterization Trick



An Example: Generating Sentences w/ Variational Autoencoders

Generating from Language Models

• Remember: using ancestral sampling, we can generate from a normal language model

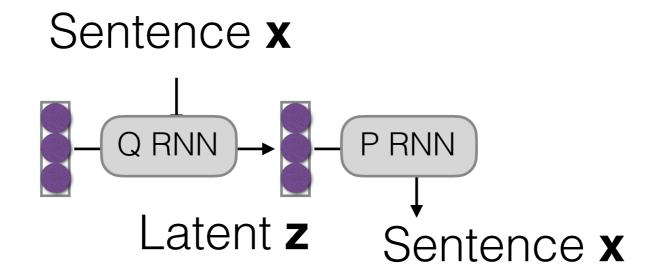
while
$$x_{j-1} != "": $x_j \sim P(x_j \mid x_1, ..., x_{j-1})$$$

We can also generate conditioned on something
 P(y|x) (e.g. translation, image captioning)

while
$$y_{j-1} != "": $y_j \sim P(y_j \mid X, y_1, ..., y_{j-1})$$$

Generating Sentences from a Continuous Space (Bowman et al. 2015)

- The VAE-based approach is conditional language model that conditions on a latent variable z
- Like an encoder-decoder, but latent representation is latent variable, input and output are identical



Motivation for Latent Variables

- Allows for a consistent latent space of sentences?
 - e.g. interpolation between two sentences

Standard encoder-decoder

i went to the store to buy some groceries.

i store to buy some groceries.

i were to buy any groceries.

horses are to buy any groceries.

horses are to buy any animal.

horses the favorite any animal.

horses the favorite favorite animal.

horses are my favorite animal.

VAE

```
"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.

 More robust to noise? VAE can be viewed as standard model + regularization.

Difficulties in Training

 Of the two components in the VAE objective, the KL divergence term is much easier to learn!

$$\mathbb{E}_{\boldsymbol{z} \sim Q(\boldsymbol{z} \mid \boldsymbol{x})}[\log P(\boldsymbol{x} \mid \boldsymbol{z})] - \mathcal{KL}[Q(\boldsymbol{z} \mid \boldsymbol{x}) || P(\boldsymbol{z})]$$

Requires good generative model

Just need to set the mean/variance of Q to be same as P

- Results in the model learning to rely solely on decoder and ignore latent variable (P(x|z) = P(x))
 - -> Posterior Collapse

Solution 1: KL Divergence Annealing

- Basic idea: Multiply KL term by a constant λ starting at zero, then gradually increase to 1
- Result: model can learn to use z before getting penalized

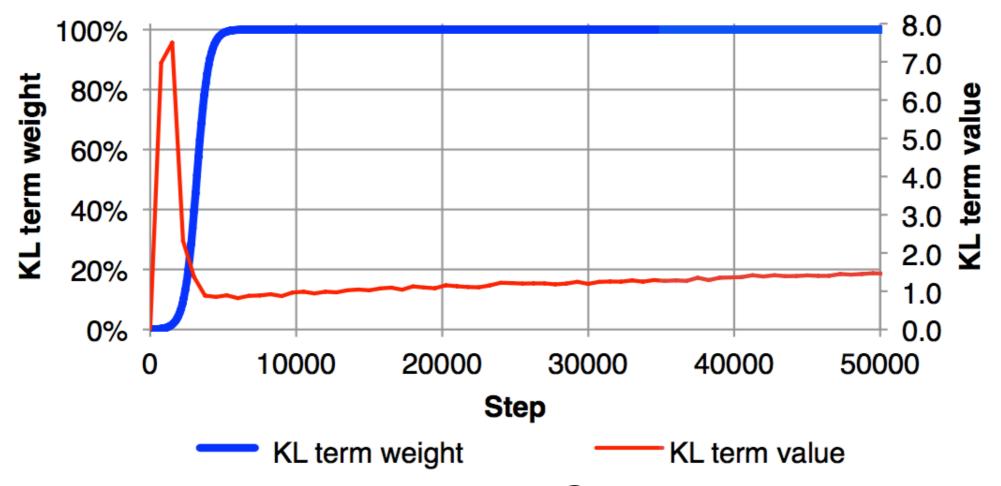


Figure Credit: Bowman et al. (2017)

Solution 2: Free bits / KL thresholding

 Free bits replaces the KL term in ELBO with a hinge loss that maximize each component of the original KL with a constant:

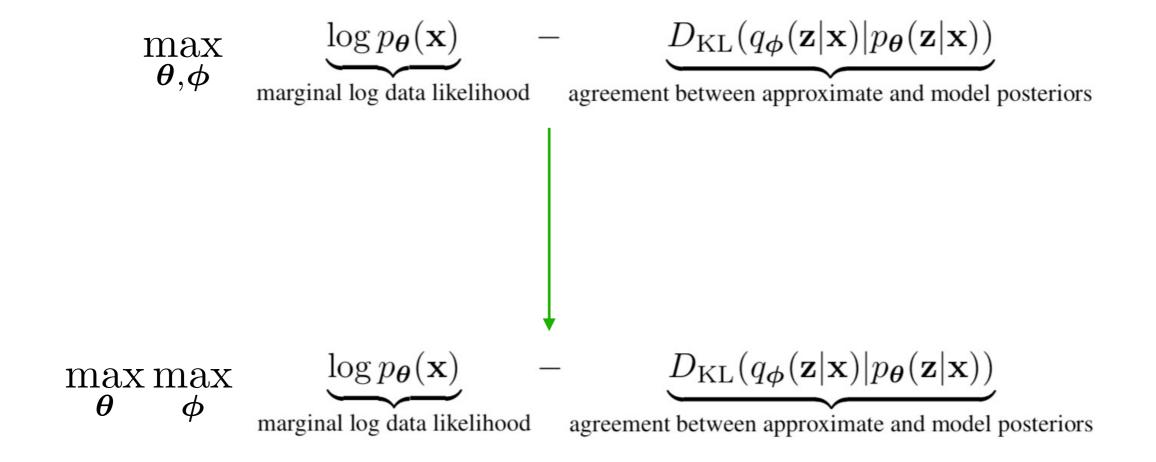
$$\sum_{i} \max[\lambda, D_{\mathrm{KL}}(q_{\phi}(z_i|x)||p(z_i))]$$

• λ :Target rate

Solution 3: Weaken the Decoder

- But theoretically still problematic: it can be shown that the optimal strategy is to ignore z when it is not necessary (Chen et al. 2017)
- Solution: weaken decoder P(x|z) so using z is essential
 - Use word dropout to occasionally skip inputting previous word in x (Bowman et al. 2015)
 - Use a convolutional decoder w/ limited context (Yang et al. 2017)

Solution 4: Aggressive Inference Network Learning



(He et al. 2019)

Handling Discrete Latent Variables

Discrete Latent Variables?

- Many variables are better treated as discrete
 - Part-of-speech of a word
 - Class of a question
 - Writer traits (left-handed or right-handed, etc.)
- How do we handle these?

Method 1: Enumeration

For discrete variables, our integral is a sum

$$P(\boldsymbol{x}; \theta) = \sum_{\boldsymbol{z}} P(\boldsymbol{x} \mid \boldsymbol{z}; \theta) P(\boldsymbol{z})$$

• If the number of possible configurations for **z** is small, we can just sum over all of them

Method 2: Sampling

- Randomly sample a subset of configurations of z and optimize with respect to this subset
- Various flavors:
 - Minimum risk training
 - Maximize ELBO loss
- Score function gradient estimator Policy Gradient Method
 - Unbiased estimator but high variance need to control variance

Method 3: Reparameterization

(Maddison et al. 2017, Jang et al. 2017)

Reparameterization also possible for discrete variables!

Original Categorical Sampling Method:

$$\hat{\boldsymbol{z}} = \text{cat-sample}(P(\boldsymbol{z} \mid \boldsymbol{x}))$$

Reparameterized Method

$$\hat{z} = \operatorname{argmax}(\log P(z \mid x) + \operatorname{Gumbel}(0,1))$$

where the Gumbel distribution is
 $\operatorname{Gumbel}(0,1) = -\log(-\log(\operatorname{Uniform}(0,1)))$

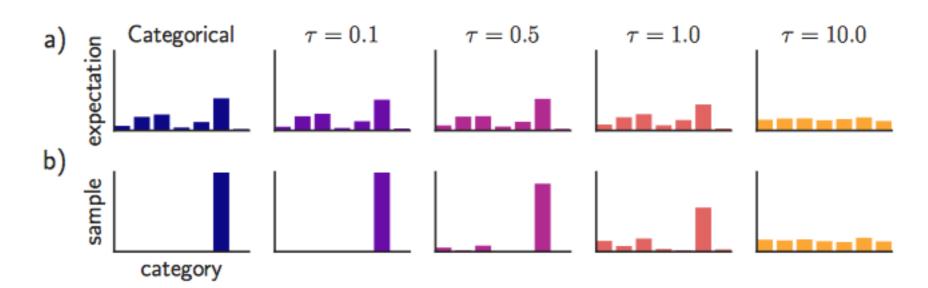
Backprop is still not possible, due to argmax

Gumbel-Softmax

- A way to soften the decision and allow for continuous gradients
- Instead of argmax, take softmax with temperature τ

$$\hat{\boldsymbol{z}} = \operatorname{softmax}((\log P(\boldsymbol{z} \mid \boldsymbol{x}) + \operatorname{Gumbel}(0,1))^{1/\tau})$$

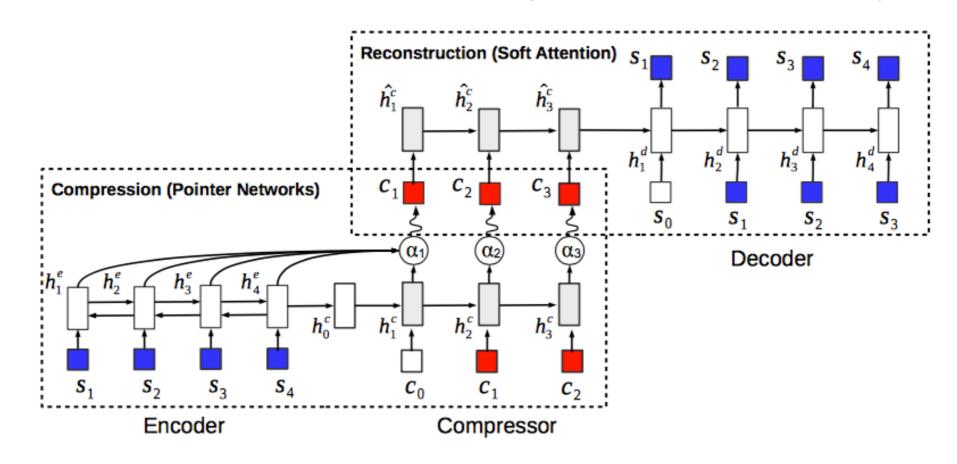
As τ approaches 0, will approach max



Application Examples in NLP

Symbol Sequence Latent Variables (Miao and Blunsom 2016)

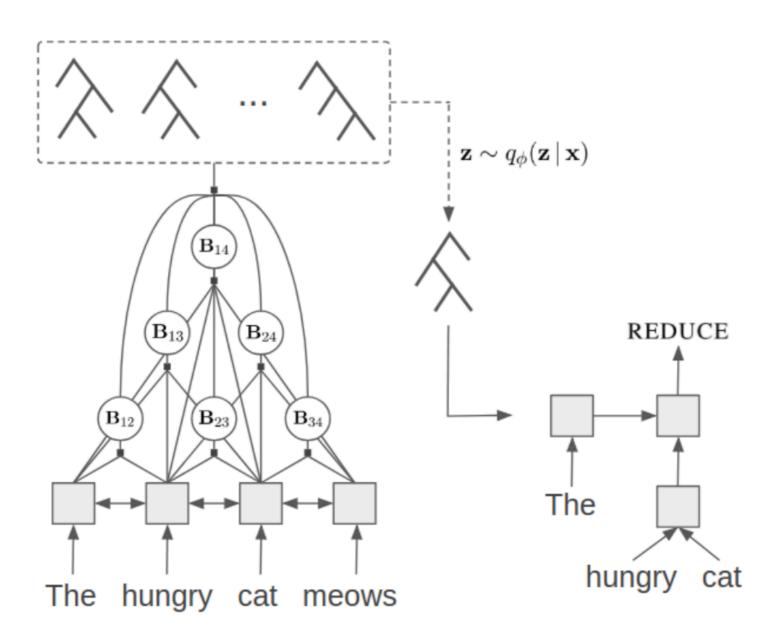
Encoder-decoder with a sequence of latent symbols



- Summarization in Miao and Blunsom (2016)
- Attempts to "discover" language (e.g. Havrylov and Titov 2017)
 - But things may not be so simple! (Kottur et al. 2017)

Unsupervised Recurrent Neural Network Grammars

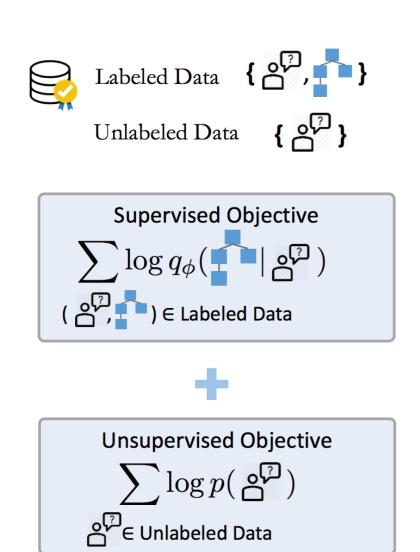
(Kim et al., 2019)

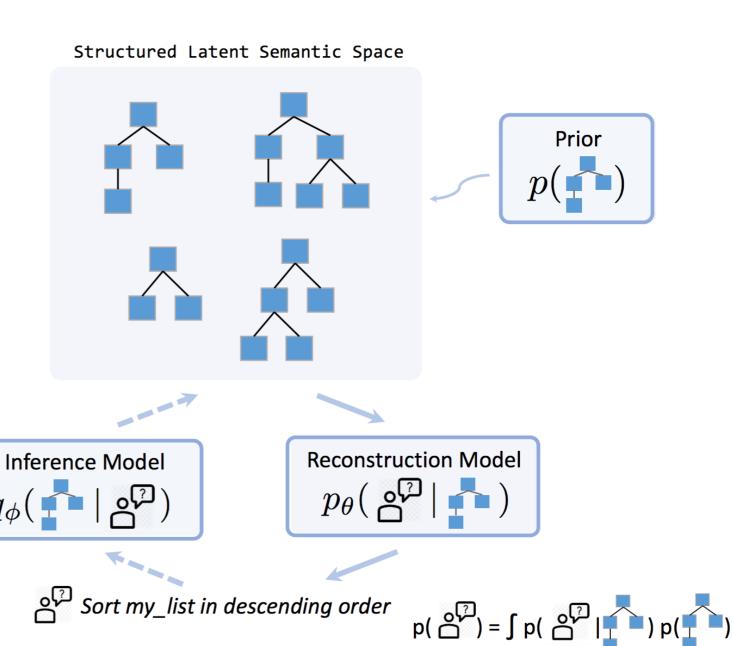


Inference Network $q_{\phi}(\mathbf{z} \mid \mathbf{x})$

Generative Model $p_{\theta}(\mathbf{x}, \mathbf{z})$

STRUCTVAE: Tree-structured Latent Variable Models for Semi-supervised Semantic Parsing (Yin et al. 2018)





Questions?