

# (Shallow and Deep) Generative Models

CS772A: Probabilistic Machine Learning

Piyush Rai

# Plan for today

- Latent Factor Models
- Latent Dirichlet Allocation (LDA)
- Deep generative models: Variational Autoencoders



# Factor Analysis and Probabilistic PCA

- Assumption: Latent variables  $\mathbf{z}_n \in \mathbb{R}^K$  typically assumed to have a Gaussian prior
  - If we want sparse latent variable, can use Laplace or spike-and-slab prior on  $\mathbf{z}_n$
  - More complex extensions of FA/PPCA use a mixture of Gaussians prior on  $\mathbf{z}_n$
- Assumption: Observations  $\mathbf{x}_n \in \mathbb{R}^D$  typically assumed to have a Gaussian likelihood
  - Other likelihood models (e.g., exp-family) can also be used if data not real-valued
- Relationship between  $\mathbf{z}_n$  and  $\mathbf{x}_n$  modeled by a noisy linear mapping

$$\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \epsilon_n = \sum_{k=1}^K \mathbf{w}_k z_{nk} + \epsilon_n$$

Zero-mean and diagonal or spherical Gaussian noise

Linear combination of the columns of  $\mathbf{W}$

$$p(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x}_n | \mathbf{z}_n) = \mathcal{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n, \Psi)$$

Diagonal for FA,  
spherical for PPCA

- Unknowns  $\mathbf{W}$ ,  $\mathbf{z}_n$ 's, and  $\Psi$  can be learned
  - EM, VI, MCMC



# Some Other Classical Models

## Gamma-Poisson latent factor model

Popular for modeling count-valued data (in text analysis, recommender systems, etc)

Non-negative priors often give a nice interpretability to such latent variable models (will see some more examples of such models shortly)

- Assumes  $K$ -dim non-negative latent variable  $\mathbf{z}_n$  and  $D$ -dim count-valued observations  $\mathbf{x}_n$
- An example: Each  $\mathbf{x}_n$  is the word-count vector representing a document

$$p(\mathbf{z}_n) = \prod_{k=1}^K \text{Gamma}(z_{nk} | a_k, b_k)$$

$$p(\mathbf{x}_n | \mathbf{z}_n) = \prod_{d=1}^D \text{Poisson}(x_{nd} | f(\mathbf{w}_d, \mathbf{z}_n))$$

This is the rate of the Poisson. It should be non-negative,  $\exp(\mathbf{w}_d^T \mathbf{z}_n)$ , or simply  $\mathbf{w}_d^T \mathbf{z}_n$  if  $\mathbf{w}_d$  is also non-negative (e.g., using a gamma/Dirichlet prior on it)

- This can be thought of as a probabilistic non-negative matrix factorization model

## Dirichlet-Multinomial/Multinoulli PCA

- Assumes  $K$ -dim non-negative latent variable  $\mathbf{z}_n$  and  $D$  categorical obs  $\mathbf{x}_n = \{x_{nd}\}_{d=1}^D$
- An example: Each  $\mathbf{x}_n$  is a document with  $D$  words in it (each word is a categorical value)

Also sums to 1

$$p(\mathbf{z}_n) = \text{Dirichlet}(\mathbf{z}_n | \boldsymbol{\alpha})$$

$$p(\mathbf{x}_n | \mathbf{z}_n) = \prod_{d=1}^D \text{Multinoulli}(x_{nd} | f(\mathbf{w}_d, \mathbf{z}_n))$$

This should give the probability vector of the multinoulli over  $x_{nd}$ . It should be non-negative and should sum to 1



# Latent Dirichlet Allocation (LDA) a.k.a. “Topic Model”



# Motivation: Multinomial Mixture Model for Text

- Assume  $D$  documents, and document  $d$  has  $N_d$  words in it
- We can represent doc  $d$  by a word count vector  $\mathbf{w}_d$
- Assuming a vocab of  $V$  unique words,  $\mathbf{w}_d$  is a  $V \times 1$  vector of counts
  - $w_{dv}$  = no of times word  $v$  appears in doc  $d$
- Let's model the docs by a mixture of  $K$  multinomial distributions, each  $V$ -dim
  - The  $k^{th}$  multinomial modeled by a  $V$ -dim prob vector  $\phi_k$  (sums to 1)
  - $\phi_k$  can be thought of as a "topic vector" (or just "topic"),  $\phi_{kv}$ : prob of word  $v$  in topic  $k$
- Generative model and plate diagram below

Each topic is a prob. distribution over word tokens

Each representing a "topic" ( $K$  topics)

Limitation: Each doc  $d$  belongs to a single cluster  $z_d$  and all words in a document assumed to be from the same topic. This is unrealistic/restrictive

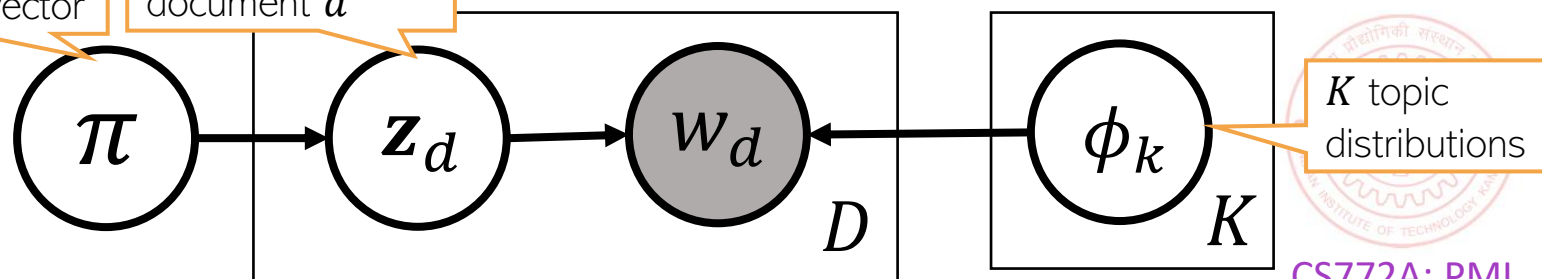
$$\mathbf{z}_d \sim \text{multinoulli}(\boldsymbol{\pi})$$

Topic Mixing proportion vector

Cluster/topic of document  $d$

$$\mathbf{w}_d \sim \text{multinomial}(\phi_{z_d}, N_d)$$

Counts will sum to  $N_d$



# Documents can be about multiple topics

## Seeking Life's Bare (Genetic) Necessities

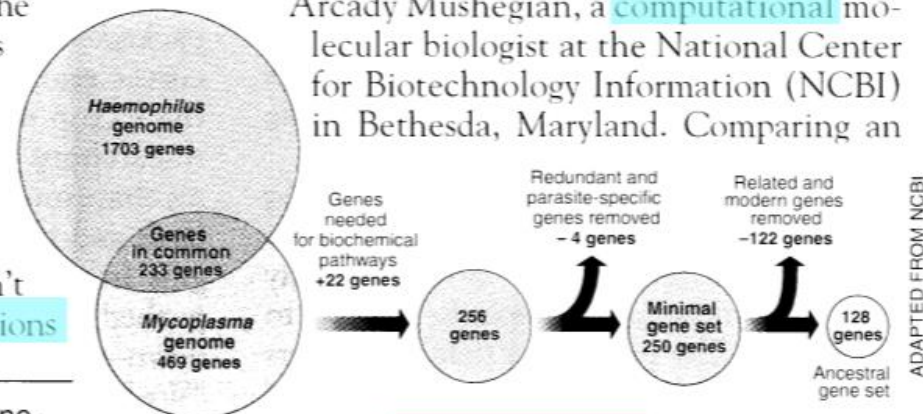
COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

“are not all that far apart,” especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers game**, particularly as more and more **genomes** are completely mapped and sequenced. “It may be a way of organizing any newly **sequenced genome**,” explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



**Stripping down.** **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

How do we find the word-topic associations in each document?

How do we use them to learn topics in the given text collection?

How do we learn low-dim document representations in terms of the topics they represent?





# A More Fine-Grained Mixture Model for Text

- Assume a corpus-level topic mixing proportions  $\alpha$  ( $K \times 1$  prob vector)
- Also assume doc-level topic mixing props  $\theta_d$  ( $K \times 1$  prob vector)
- Instead of assuming a single cluster  $\mathbf{z}_d$  for doc  $d$ , cluster each word in it
  - $\mathbf{z}_{d,n} \in \{1, 2, \dots, K\}$  denotes the cluster/topic of word  $w_{d,n} \in \{1, 2, \dots, V\}$

Each assumed a one-hot  $K \times 1$  vector

- Can obtain the “average” clustering for doc  $d$  using  $\theta_d$  or  $\bar{\mathbf{z}}_d = \frac{1}{N_d} \sum_{n=1}^{N_d} \mathbf{z}_{d,n}$

Locally-conjugate. Easy  
Gibbs sampling, VI, etc

Somewhat similar to  
Dir-Mult PCA model

- The generative model is as follows

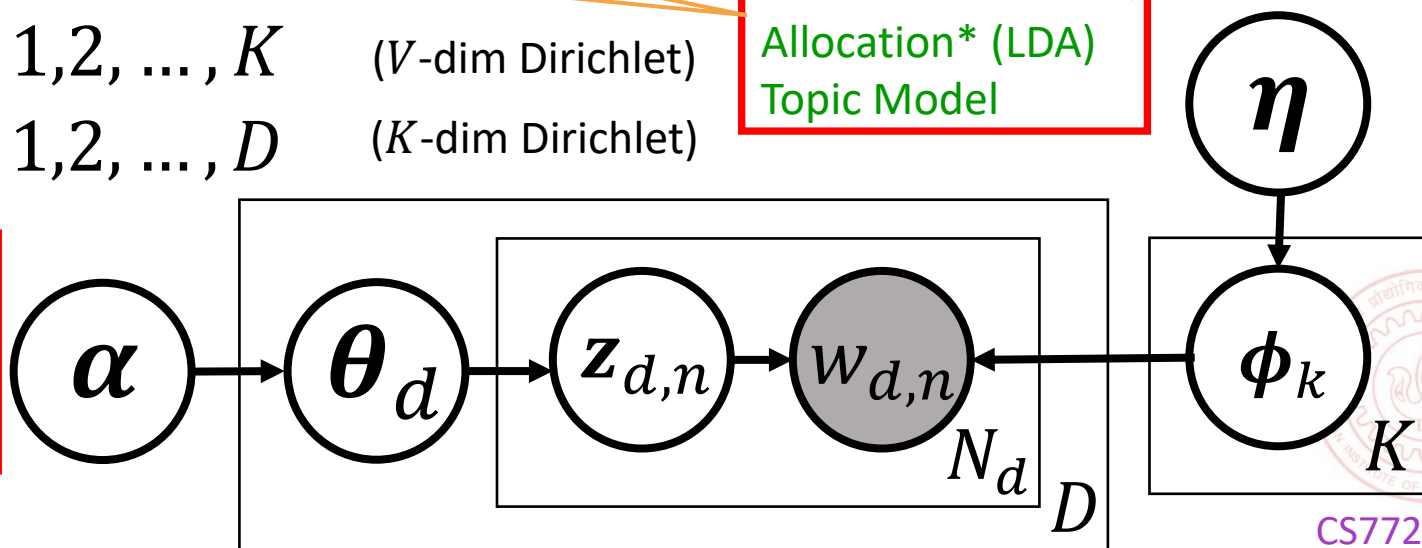
$$\phi_k \sim \text{Dirichlet}(\eta) \quad k = 1, 2, \dots, K \quad (V\text{-dim Dirichlet})$$

$$\theta_d \sim \text{Dirichlet}(\alpha) \quad d = 1, 2, \dots, D \quad (K\text{-dim Dirichlet})$$

Latent Dirichlet  
Allocation\* (LDA)  
Topic Model

$$\mathbf{z}_{d,n} \sim \text{multinoulli}(\theta_d)$$

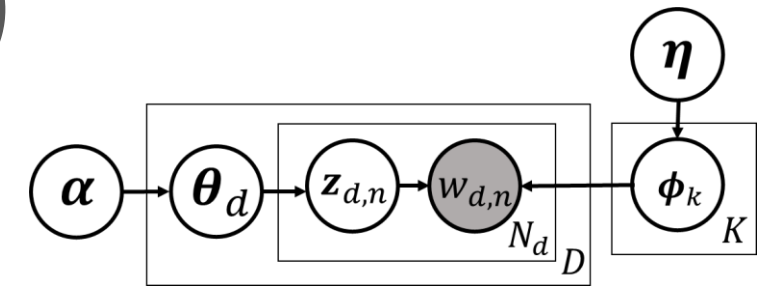
$$\mathbf{w}_{d,n} \sim \text{multinoulli}(\phi_{\mathbf{z}_{d,n}})$$





# Latent Dirichlet Allocation (LDA)

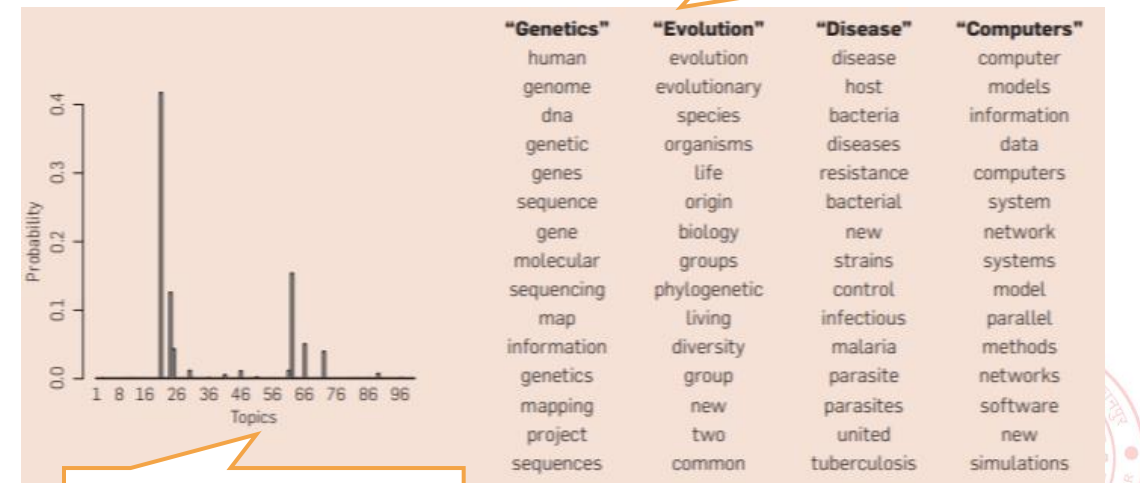
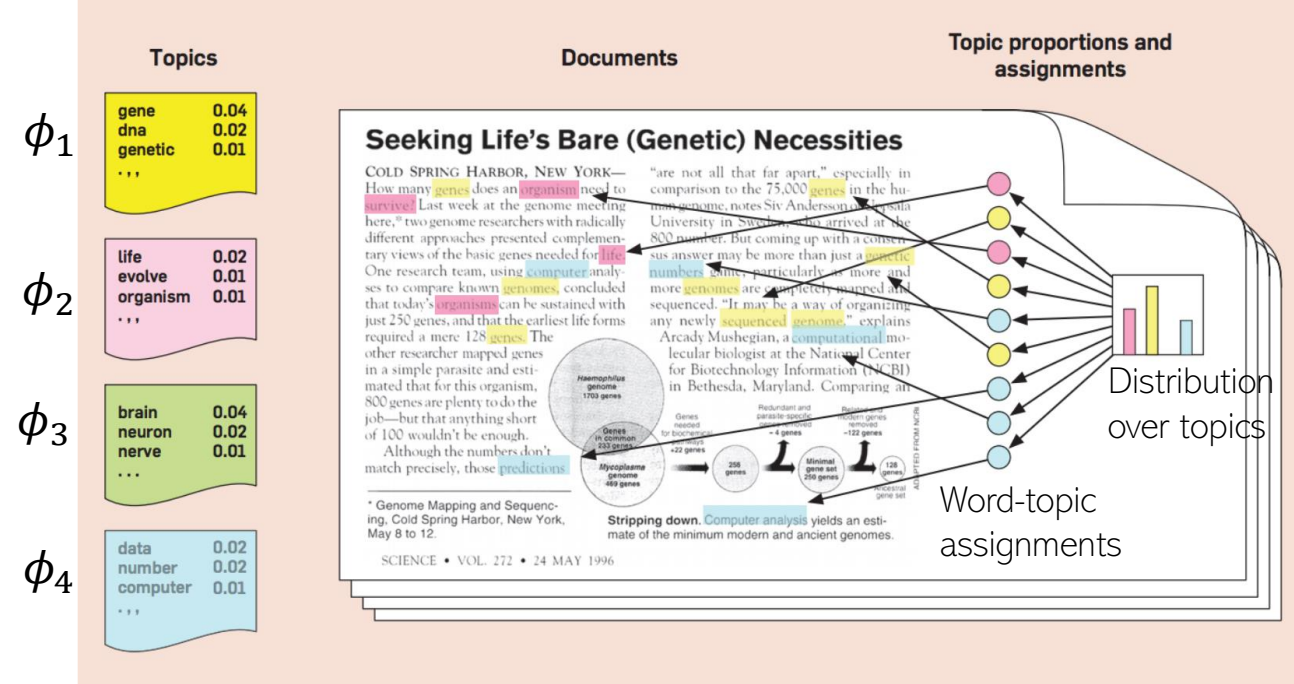
- A very widely used probabilistic model for text data
- Nice and easy insights into the text collection



- Each  $\phi_k = [\phi_{k1}, \dots, \phi_{kV}]$  can be interpreted as topic ( $\phi_{kv}$  = prob. of word  $v$  in topic  $k$ )
- $\theta_d = [\theta_{d1}, \dots, \theta_{dK}]$ : how much each topic is present in document  $d$  (topic distribution)
- $\bar{z}_d = \frac{1}{N_d} \sum_{n=1}^{N_d} z_{d,n}$  also has a similar interpretation as  $\theta_d$

A topic is a set of words that tend to co-occur together

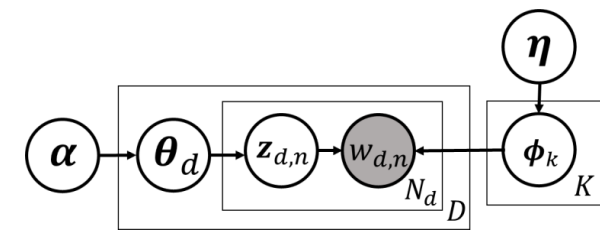
15 most frequent (most probable) words from four most prominent topics in this doc



Topic distribution for the document on left



# LDA: Inference and Evaluation



- LDA is locally conjugate. Many inference methods (VI, variational EM, Gibbs samp, etc)

$$p(\mathbf{Z}, \Theta, \Phi | \mathbf{W}, \alpha, \eta) = \frac{p(\mathbf{W} | \Phi, \mathbf{Z}) p(\mathbf{Z} | \Theta) p(\Phi | \eta) p(\Theta | \alpha)}{p(\mathbf{W} | \alpha, \eta)} \quad (\text{assuming hyperparams } \alpha, \eta \text{ are fixed})$$

- Can even collapse some variables and do collapsed Gibbs or collapsed VB
  - E.g., collapse  $\theta_d$  and  $\phi_k$  (if needed, these can be approximated using  $\mathbf{Z}$ )
- Many ways to evaluate how well LDA performs on some data
  - Extrinsic measures: Perform LDA and use its output for another task (e.g., classification)
  - Perplexity is another **intrinsic** measure to evaluate LDA-style models

Lower is better

Test set with  $M$  docs

Marginal likelihood of all words in the  $d^{\text{th}}$  test doc

$$\text{perplexity}(D_{\text{test}}) = \exp \left\{ - \frac{\sum_{d=1}^M \log p(\mathbf{w}_d)}{\sum_{d=1}^M N_d} \right\}$$



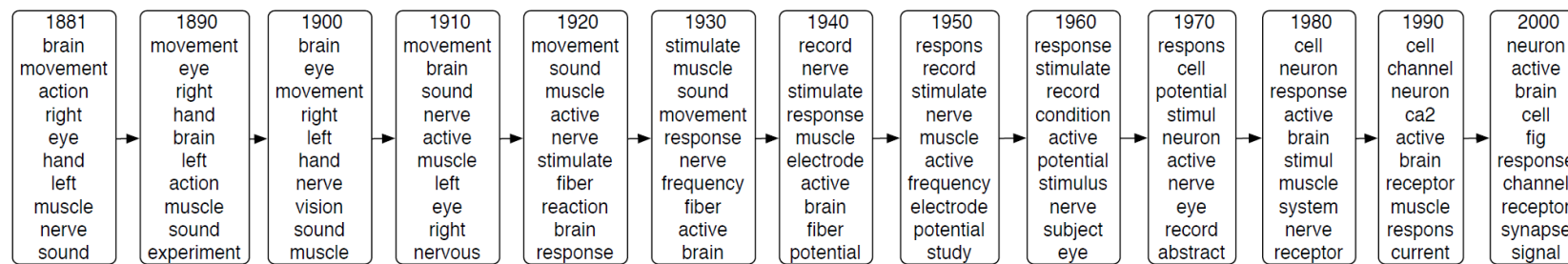
# LDA: Limitations and Extensions

- LDA assumes topics remain static over time (improvement: Dynamic Topic Model)

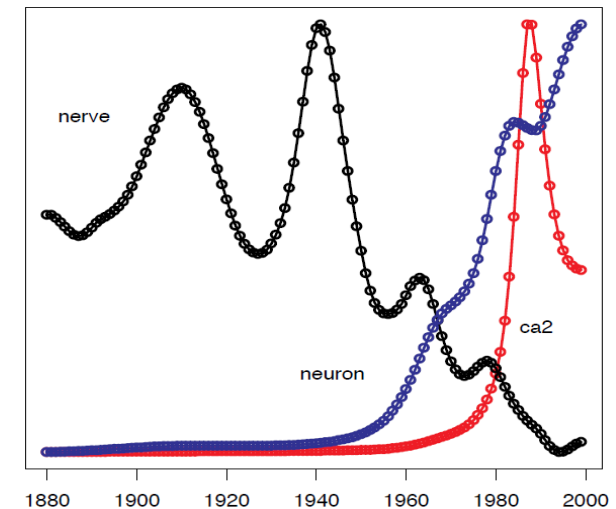
Assume a first-order Markov evolution for each topic w.r.t. time

$$w_k^t \sim \mathcal{N}(w_k^{t-1}, \sigma^2 I) \quad \phi_k^t = \mathcal{S}(w_k^t)$$

Simplex transformation (convert  $w_k^t$  into a probability vector)



Evolution of topic “Neuroscience”  
(learned from the journal Science)



- LDA assumes topics are uncorrelated (improvement: Corr-LDA)
  - Use a **logistic normal** distribution on  $\theta_d$  (cov matrix of log-normal makes component correlated)
- LDA ignores the sequential structure in the text (improvement: HMM-LDA)

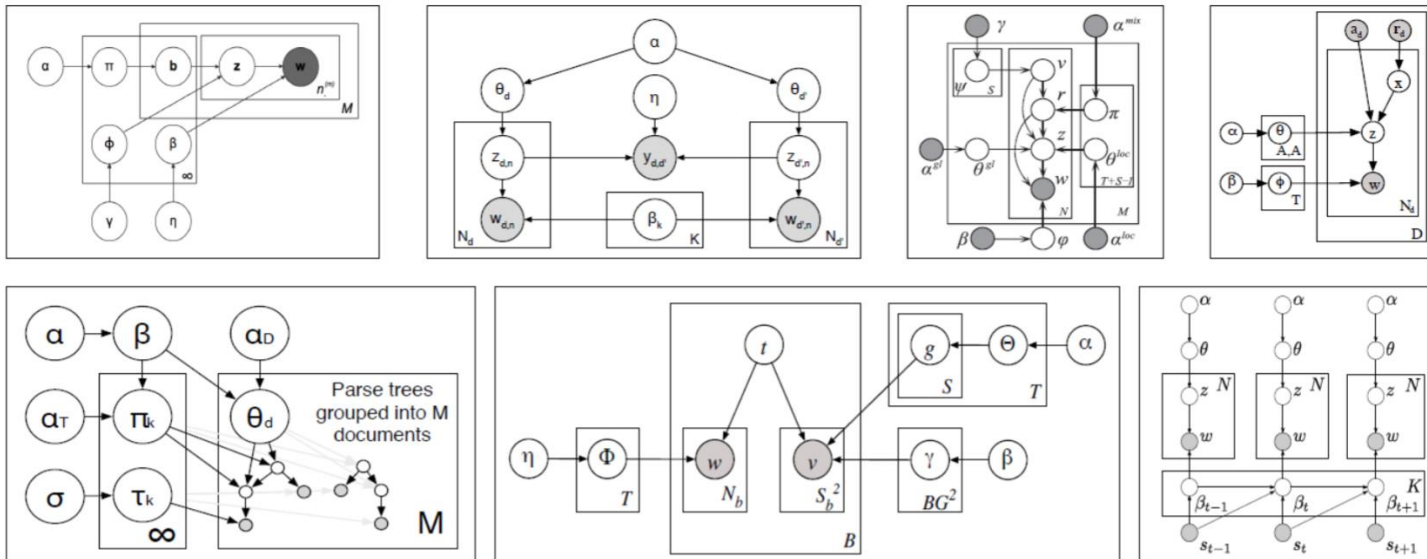


# LDA Extensions (Contd)

Also: “Neural” Topic Models are popular nowadays ( $z$  to  $x$  mapping and vice-versa modeled via deep nets)

- LDA for non-text data, e.g., images
  - Each image can be represented as a bag of “visual words” and LDA can be applied
- Supervised/Labeled LDA (when we have a label for each document)
- LDA for paired/multimodality data (e.g., images and text caption)
- LDA for graph-structured data instead of documents

Plate diagrams for some LDA extensions



LDA is also equivalent to doing a non-negative matrix fact. of the  $V \times D$  word-document matrix  $\mathbf{X}$  using a Poisson likelihood model\*

$$\mathbf{X} \sim \text{Poisson}(\Phi\Theta)$$

$\Phi$  ( $V \times K$ ) and  $\Theta$  ( $K \times D$ ) can be given any non-negative priors (Dirichlet/gamma)

This can be extended to “deep” matrix factorization\*\* (modeling  $\Theta$  using many layers)

\*Sec 4 and 5 of “Beta-Negative Binomial Process and Poisson Factor Analysis” (Zhou et al, 2012)

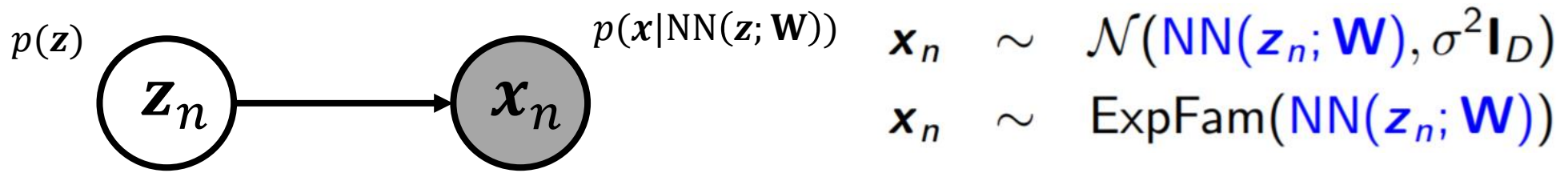
\*\* Poisson-gamma belief networks” (Zhou et al, 2015)





# Constructing Generative Models using Neural Nets<sup>13</sup>

- We can use a neural net to define the mapping from a  $K$ -dim  $\mathbf{z}_n$  to  $D$ -dim  $\mathbf{x}_n$



- If  $\mathbf{z}_n$  has a Gaussian prior, such models are called **deep linear Gaussian models** (DLGM)
- Since NN mapping can be very powerful, DLGM can generate very high-quality data
  - Take the trained network, generate a random  $\mathbf{z}$  from prior, pass it through the model to generate  $\mathbf{x}$

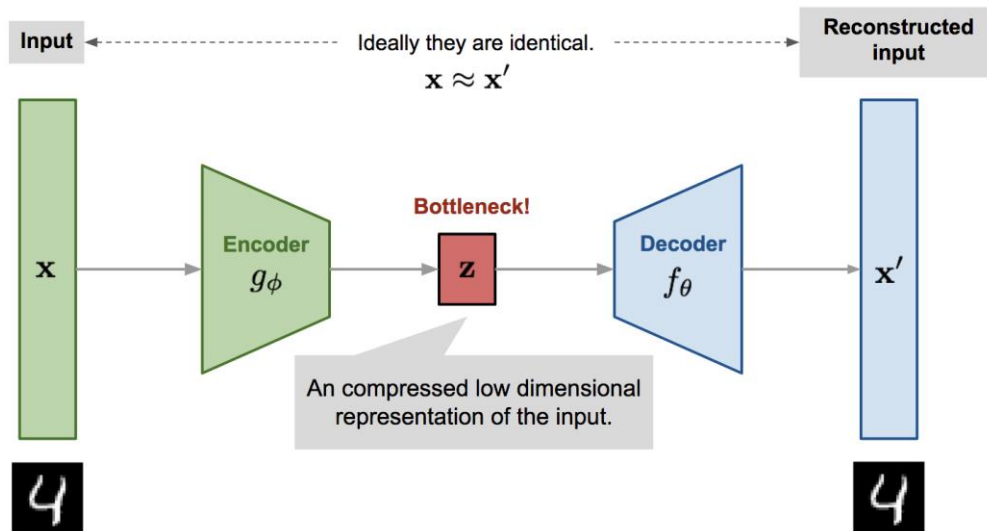


Some sample images generated by Vector Quantized Variational AutoEncoder (VQ-VAE), a state-of-the-art DLGM



# Variational Autoencoder (VAE)

- VAE\* is a probabilistic extension of autoencoders (AE)



$$L_{AE}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - f_{\theta}(g_{\phi}(\mathbf{x}^{(i)})))^2$$

- The basic difference is that VAE assumes a prior  $p(\mathbf{z})$  on the latent code  $\mathbf{z}$ 
  - This enables it to not just compress the data but also generate synthetic data
  - How: Sample  $\mathbf{z}$  from a prior and pass it through the decoder
- Thus VAE can learn good latent representation + generate novel synthetic data
- The name has “Variational” in it since it is learned using VI principles



# Variational Autoencoder (VAE)

- VAE has three main components
  - A prior  $p_{\theta}(\mathbf{z})$  over latent codes
  - A probabilistic decoder  $p_{\theta}(\mathbf{x}|\mathbf{z})$
  - A posterior or probabilistic encoder  $p_{\theta}(\mathbf{z}|\mathbf{x})$  approx. by an “inference network”  $q_{\phi}(\mathbf{z}|\mathbf{x})$

Here  $\theta$  collectively denotes all the parameters of the prior and likelihood

Using the idea of “Amortized Inference” (next slide)

- VAE is learned by maximizing the ELBO

ELBO for a single data point

$$\begin{aligned} \mathcal{L}(\theta, \phi | \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathbb{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z})) \end{aligned}$$

Here  $\phi$  collectively denotes all the parameters that define the inference network

Maximized to find the optimal  $\theta$  and  $\phi$

$q_{\phi}$  should reconstruct the data  $\mathbf{x}$  well from  $\mathbf{z}$  (high log-lik)

$q_{\phi}$  should also be simple (close to the prior)

- The **Reparametrization Trick** is commonly used to optimize the ELBO
- Posterior is inferred only over  $\mathbf{z}$ , and usually only point estimate on  $\theta$  and  $\phi$





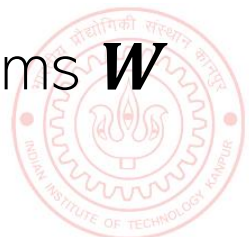
# Amortized Inference

- Latent variable models need to infer the posterior  $p(\mathbf{z}_n | \mathbf{x}_n)$  for each observation  $\mathbf{x}_n$
- This can be slow if we have lots of observations because
  1. We need to iterate over each  $p(\mathbf{z}_n | \mathbf{x}_n)$
  2. Learning the global parameters needs wait for step 1 to finish for all observations
- One way to address this is via Stochastic VI (already saw)
- Amortized inference is another appealing alternative (used in VAE and other LVMs too)

$$p(\mathbf{z}_n | \mathbf{x}_n) \approx q(\mathbf{z}_n | \phi_n) = q(\mathbf{z}_n | \text{NN}(\mathbf{x}_n; \mathbf{W}))$$

If  $q$  is Gaussian then the NN will output a mean and a variance

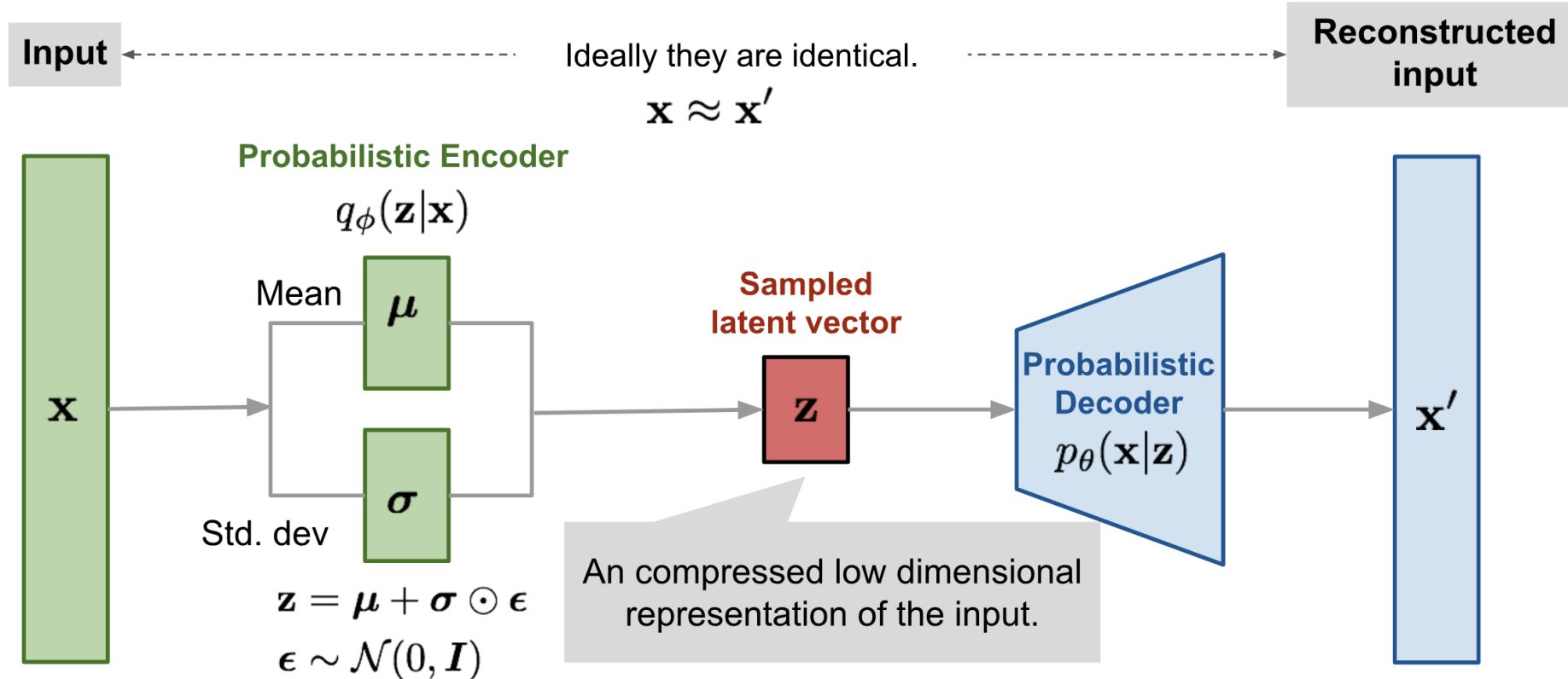
- Thus no need to learn  $\phi_n$ 's (one per data point) but just a single NN with params  $\mathbf{W}$ 
  - This will be our “encoder network” for learning  $\mathbf{z}_n$
  - Also very efficient to get  $p(\mathbf{z}_* | \mathbf{x}_*)$  for a new data point  $\mathbf{x}_*$



# Variational Autoencoder: The Complete Pipeline

- Both probabilistic encoder and decoder learned jointly by maximizing the ELBO

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \phi | \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z} | \mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z}))\end{aligned}$$



# VAE and Posterior Collapse

- VAEs may suffer from **posterior collapse**

Decoder is a neural net and can be arbitrarily powerful making this term very large

Consequently, KL will become close to zero collapsing posterior to the prior

$$\mathcal{L}(\theta, \phi | \mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z}))$$

- Thus, due to posterior collapse, reconstruction will still be good but the code  $\mathbf{z}$  may be garbage (not useful as a representation for  $\mathbf{x}$ )
- Several ways to prevent posterior collapse, e.g.,

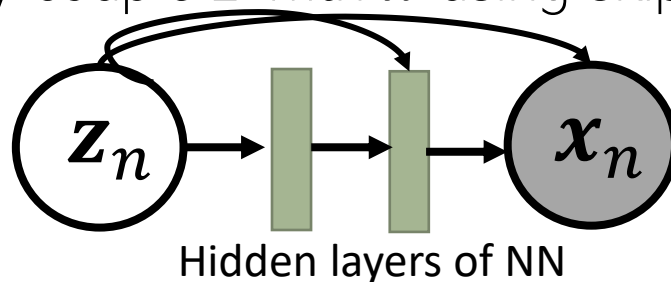
- Use KL annealing

A carefully tuned value between 0 and 1

$$\mathcal{L}(\theta, \phi | \mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \beta \text{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) || p_{\theta}(\mathbf{z}))$$

For example, keep the variance of  $q$  as fixed

- Avoid KL from becoming 0 using some  $q$  doesn't collapse to the prior
- More tightly couple  $\mathbf{z}$  with  $\mathbf{x}$  using skip-connections (Skip-VAE)



Besides these, MCMC (sometimes used for inference in VAE), or improved VI techniques can also help in preventing posterior collapse in VAEs

