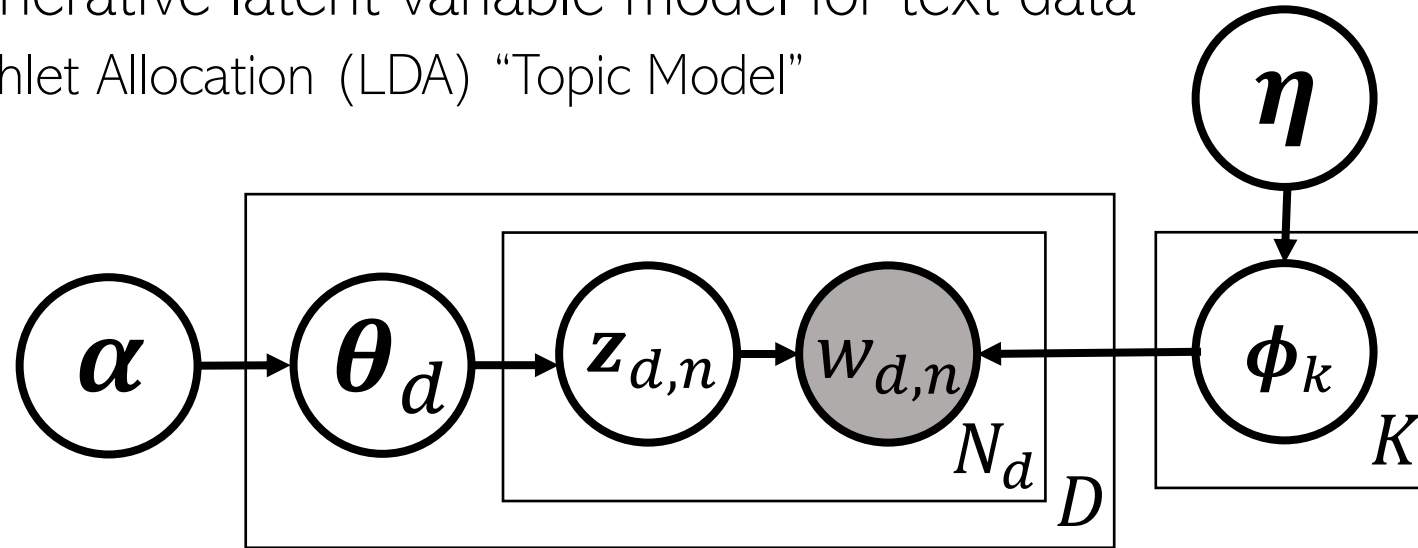# Assorted Topics (1)

CS772A: Probabilistic Machine Learning
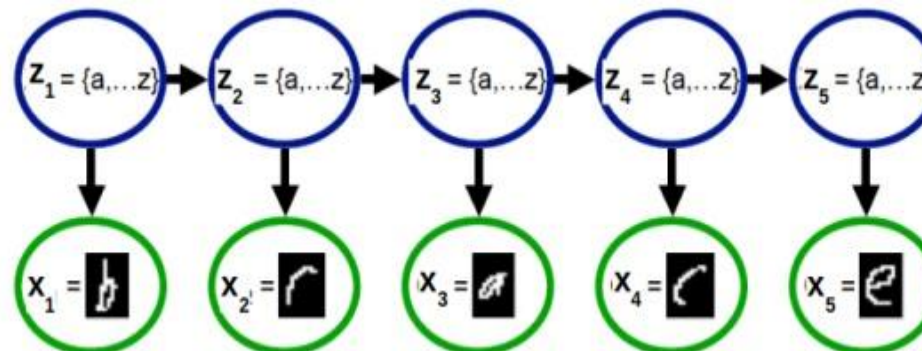
Piyush Rai

# Plan today..

- A classical generative latent variable model for text data
  - Latent Dirichlet Allocation (LDA) "Topic Model"



- Probabilistic models for sequential data
  - HMM, state-space models

# 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 $\boldsymbol{w_d}$

- Assuming a vocab of $V$ unique words, $\boldsymbol{w_d}$ is a $V \times 1$ vector of counts

  > Each topic is a prob. distribution over word tokens

  - $w_{dv} = $ no of times word $v$ appears in doc $d$

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

- 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

  > **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**
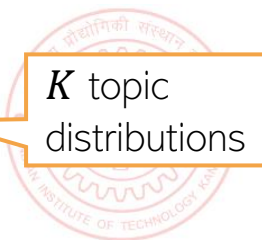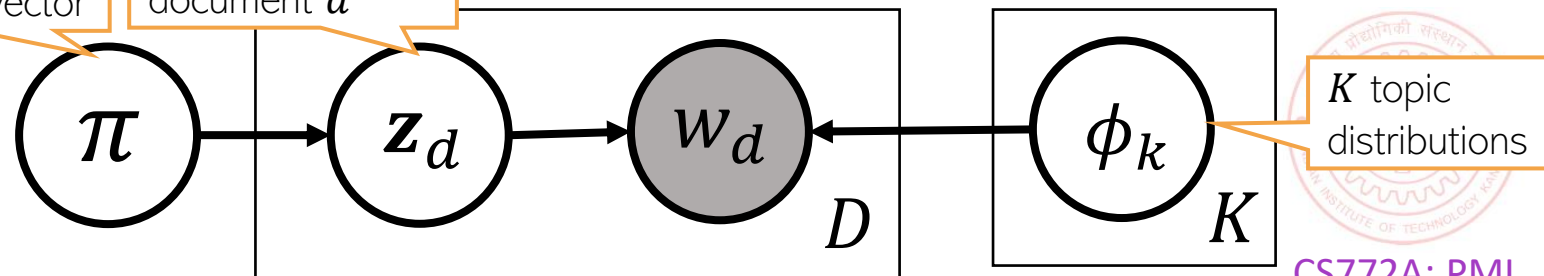
$$z_d \sim \text{multinoulli}(\pi)$$

$$w_d \sim \text{multinomial}(\phi_{z_d}, N_d)$$

> Topic Mixing proportion vector

> Cluster/topic of document $d$

> Counts will sum to $N_d$

> $K$ topic distributions

# Documents can be about multiple topics



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 <u>corpus-level</u> topic mixing proportions $\boldsymbol{\alpha}$ ($K \times 1$ prob vector)

- Also assume <u>doc-level</u> 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
  - $\mathbf{z}_{d,n} \in \{1, 2, \ldots, K\}$ denotes the cluster/topic of word $w_{d,n} \in \{1, 2, \ldots, V\}$

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

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

- The generative model is as follows

  > Locally-conjugate. Easy Gibbs sampling, VI, etc

  > Latent Dirichlet Allocation* (LDA) Topic Model

  > Somewhat similar to Dir-Mult PCA model

$$\phi_k \sim \text{Dirichlet}(\boldsymbol{\eta}) \quad k = 1, 2, \ldots, K \qquad (V\text{-dim Dirichlet})$$
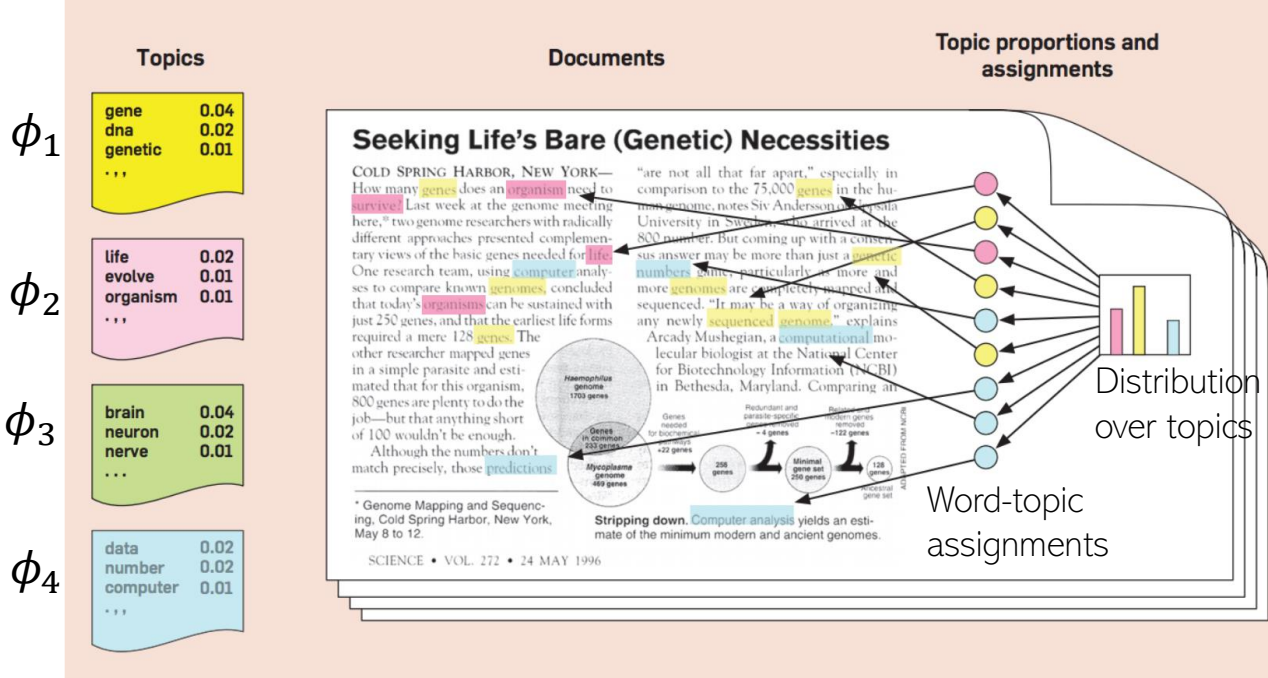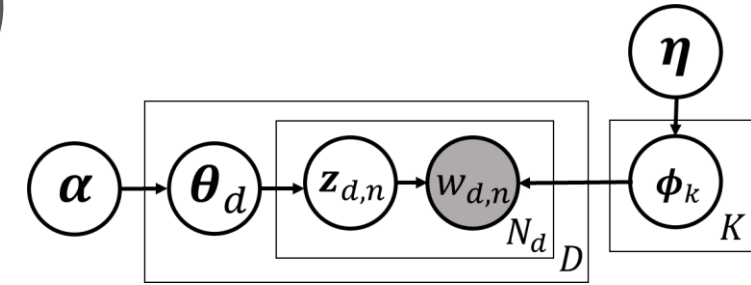
$$\theta_d \sim \text{Dirichlet}(\boldsymbol{\alpha}) \quad d = 1, 2, \ldots, D \qquad (K\text{-dim Dirichlet})$$

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

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

# Latent Dirichlet Allocation (LDA)

- A very widely used probabilistic model for text data
- Nice and easy insights into the text collection
    - Each $\phi_k = [\phi_{k1}, \ldots, \phi_{kV}]$ can be interpreted as topic ($\phi_{kv} =$ prob. of word $v$ in topic $k$)
    - $\theta_d = [\theta_{d1}, \ldots, \theta_{dK}]$: 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 $\theta_d$

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

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



Topics

$\phi_1$

| gene | 0.04 |
| dna | 0.02 |
| genetic | 0.01 |
| ... | |

$\phi_2$

| life | 0.02 |
| evolve | 0.01 |
| organism | 0.01 |
| ... | |

$\phi_3$

| brain | 0.04 |
| neuron | 0.02 |
| nerve | 0.01 |
| ... | |

$\phi_4$

| data | 0.02 |
| number | 0.02 |
| computer | 0.01 |
| ... | |

Documents

Topic proportions and assignments

Distribution over topics

Word-topic assignments

Topic distribution for the document on left

# LDA: Inference and Evaluation

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

$$p(\mathbf{Z}, \Theta, \Phi | \mathbf{W}, \alpha, \eta) = \frac{p(\mathbf{W}|\Phi, \mathbf{Z})p(\mathbf{Z}|\Theta)p(\Phi|\eta)p(\Theta|\alpha)}{p(\mathbf{W}|\alpha, \eta)}$$

(assuming hyperparams $\alpha, \eta$ 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

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

Test set with $M$ docs

Lower is better

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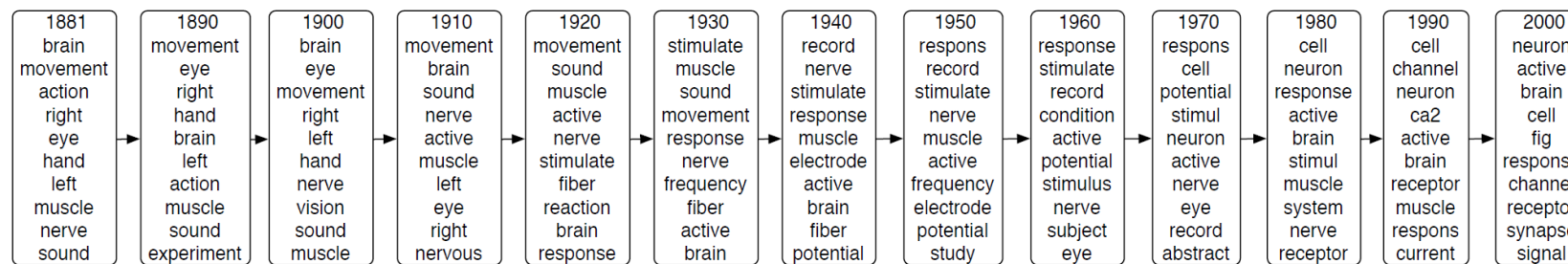
# LDA: Limitations and Extensions

- LDA assumes topics remain static over time (improvement: Dynamic Topic Model)
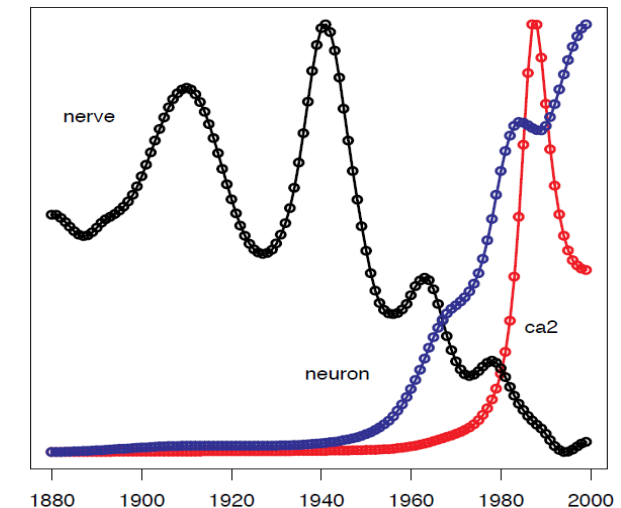
Assume a first-order Markov evolution for each topic w.r.t. time

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

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)
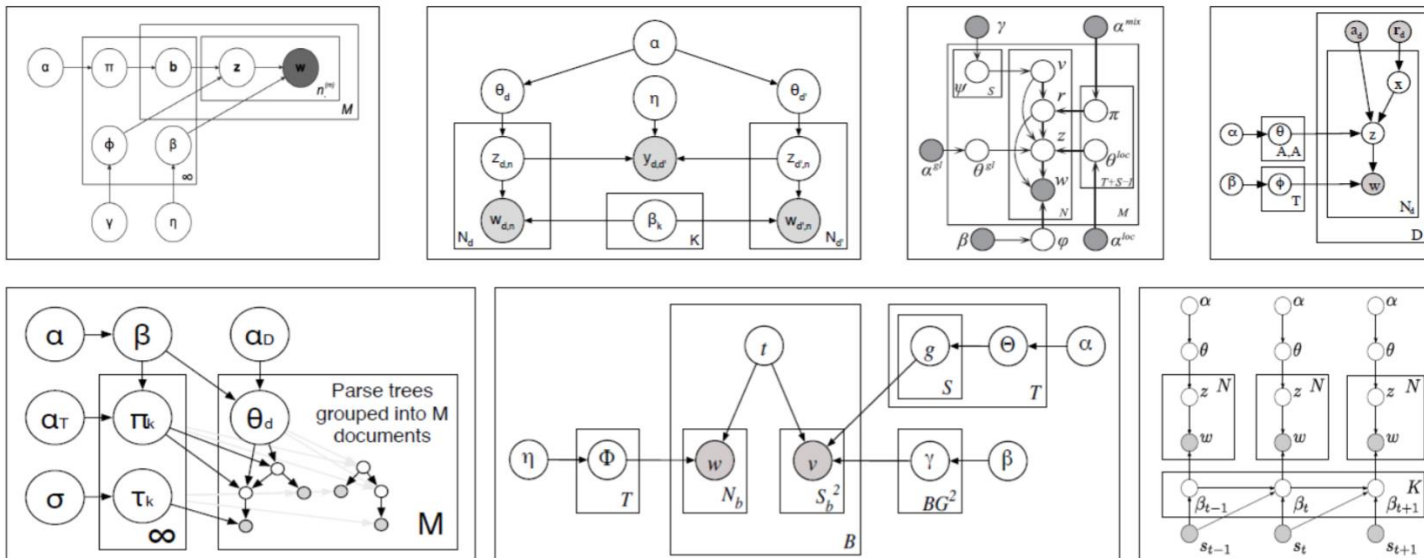
Fig courtesy: Dynamic Topic Models (Blei and Lafferty, 2006)

# LDA Extensions (Contd)

Also: "Neural" Topic Models are popular nowadays ($z$ to $x$ mapping and vice-versa modeled via deep nets). Also, some topic models use pre-computed word-embeddings rather than one-hot Representation of each word

- 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 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

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*

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

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

This can be extended to "deep" matrix factorization** (modeling $\boldsymbol{\Theta}$ 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)
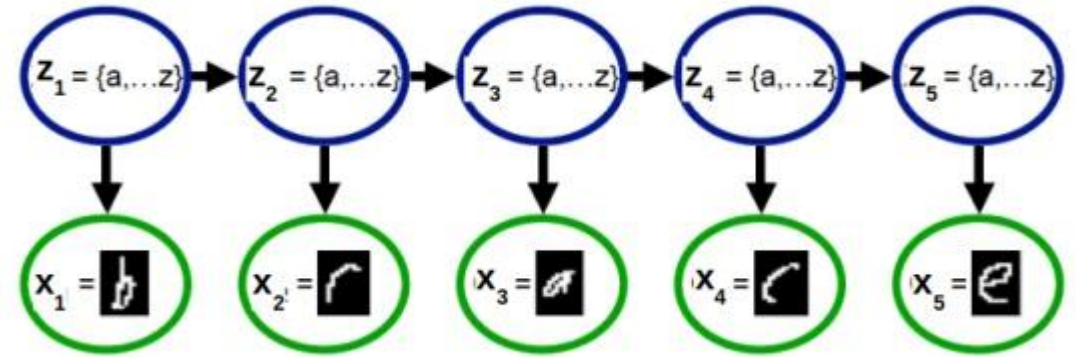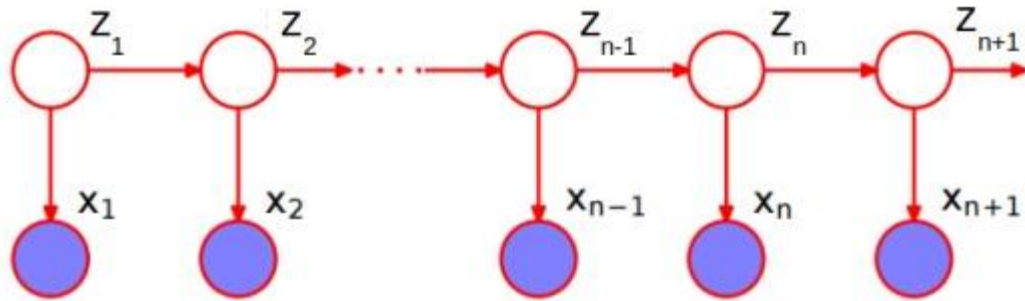
**Plate diagrams for some LDA extensions**

# Probabilistic Models for Sequential Data

# Latent Variable Models for Sequential Data

- Task: Given a sequence of observations, infer the latent state of each observation



Observation model $\longrightarrow$ $x_n|z_n \quad \sim \quad p(x_n|z_n) \qquad$ (i.i.d. draws of $x_n$ given $z_n$)

State-transition model $\longrightarrow$ $z_n|z_{n-1} \quad \sim \quad p(z_n|z_{n-1}) \qquad$ (first-order dependence b/w $z_n$'s)

- If $z_n$'s are discrete, we have a hidden Markov model (HMM) $\quad p(z_n|z_{n-1} = \ell) = \text{multinoulli}(\pi_\ell)$
- If $z_n$'s are real-valued, we have a state-space model (SSM) $\quad p(z_n|z_{n-1}) = \mathcal{N}(\mathbf{A}z_{n-1}, \mathbf{I}_K)$

# State-Space Models

- In the most general form, the state-transition and observation models of an SSM



Using 's' instead of 'z' to refer to states

Using 't' to denote the 'time-step'

HMM is similar to SSM except the state-transition model is a discrete distribution

$g_t, h_t$ can be linear or nonlinear functions

$$s_t | s_{t-1} = g_t(s_{t-1}) + \epsilon_t \quad \text{(must be a cont. dist. over } s_t\text{)}$$

$$x_t | s_t = h_t(s_t) + \delta_t \quad \text{(can be any dist. over } x_t\text{)}$$

- Assuming Gaussian noise in the state-transition and observation models

This is a Gaussian SSM

$$s_t | s_{t-1} \sim \mathcal{N}(s_t | g_t(s_{t-1}), Q_t)$$

$$x_t | s_t \sim \mathcal{N}(x_t | h_t(s_t), R_t)$$

If $g_t, h_t, Q_t, R_t$ are independent of $t$ then it is called a stationary model

$g_t, h_t, Q_t, R_t$ may be known or can be learned

# State-Space Models: A Simple Example

- Consider the linear Gaussian SSM

$$s_t | s_{t-1} = \mathbf{A}_t s_{t-1} + \epsilon_t$$
$$x_t | s_t = \mathbf{B}_t s_t + \delta_t$$

- Suppose $x_t \in \mathbb{R}^2$ denotes the (noisy) observed 2D location of an object
- Suppose $s_t \in \mathbb{R}^6$ denotes the "state" vector

$$s_t = [\text{pos1, vel1, accel1, pos2, vel2, accel2}]$$

- Here is an example SSM for this problem with pre-defined $\mathbf{A}_t$ and $\mathbf{B}_t$ matrices

$$\mathbf{A}_t$$

$$s_t = \begin{bmatrix} 1 & \Delta t & \frac{1}{2}(\Delta t)^2 & 0 & 0 & 0 \\ 0 & 1 & \Delta t & 0 & 0 & 0 \\ 0 & 0 & e^{-\alpha \Delta t} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & \Delta t & \frac{1}{2}(\Delta t)^2 \\ 0 & 0 & 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 0 & 0 & e^{-\alpha \Delta t} \end{bmatrix} s_{t-1} + \epsilon_t$$
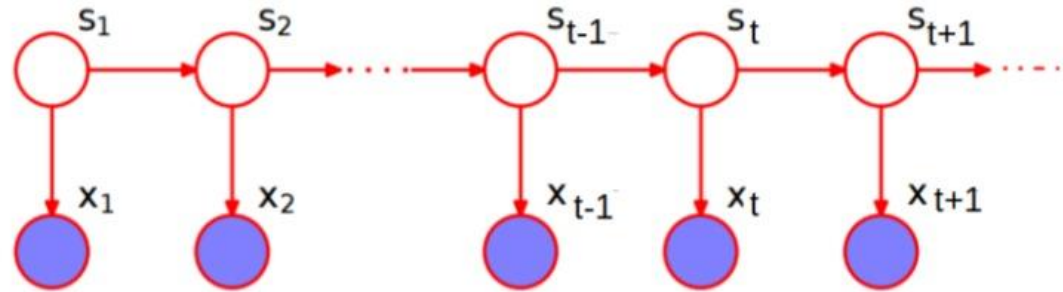
$$\mathbf{B}_t$$

$$x_t = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} s_t + \delta_t$$

# Typical Inference Task for Gaussian SSM

- One of the key tasks: Given sequence $x_1, x_2, \ldots, x_T$, infer latent $s_1, s_2, \ldots, s_T$



- Usually two ways of inferring the latent states
  - Infer $p(s_t|x_1, x_2, \ldots, x_t)$: Called the "filtering" problem

  A Gaussian

  Kalman Filtering is a popular algorithm for a linear Gaussian SSM

  Turns out to be another Gaussian

  $$p(s_t|x_1, x_2, \ldots, x_t) \propto \underbrace{p(x_t|s_t)}_{\mathcal{N}(x_t|\mathbf{B}s_t, \mathbf{R})} \int \underbrace{p(s_t|s_{t-1})}_{\mathcal{N}(s_t|\mathbf{A}s_{t-1}, \mathbf{Q})} p(s_{t-1}|x_1, x_2, \ldots, x_{t-1}) ds_{t-1}$$

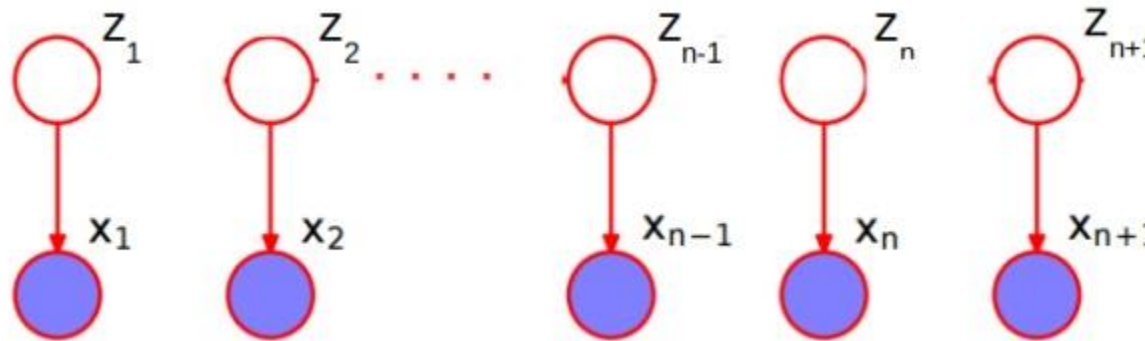  - Infer $p(s_t|x_1, x_2, \ldots, x_t, \ldots, x_T)$: Called the "smoothing" problem
- Some other tasks one can solve for using an SSM
  - Predicting future states $p(s_{t+h}|x_1, x_2, \ldots, x_t)$ for $h \geq 1$, given observations thus far
  - Predicting future observations $p(x_{t+h}|x_1, x_2, \ldots, x_t)$ for $h \geq 1$, given observations thus far

# A Special Case

- What if we have i.i.d. latent states, i.e.,. $p(z_n|z_{n-1}) = p(z_n)$?



- Discrete case (HMM) becomes a simple mixture model $p(z_n|z_{n-1} = \ell) = p(z_n) = \text{multinoulli}(\pi)$
- Real-valued case (SSM) becomes a PPCA model $p(z_n|z_{n-1}) = p(z_n) = \mathcal{N}(0, I_K)$ or $\mathcal{N}(\mu, \Psi)$
- Inference algos for HMM/SSM are thus very similar to that of mixture models/PPCA
  - Only main difference is how the latent variables $z_n$'s are inferred since they aren't i.i.d.
  - E.g., if using EM, only E step needs to change (Bishop Chap 13 has EM for HMM and SSM)