

# Probabilistic Models for Graphs, and Intro to Nonparametric Bayesian Modeling

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

March 25, 2019



# Modeling Graphs

- Often we wish to understand the underlying structure (e.g., communities/groups/topics) in a graph, predict links, classify nodes, visualize, etc.
- An example graph<sup>†</sup> (a 575,000 node citation network of papers; each node is a paper):

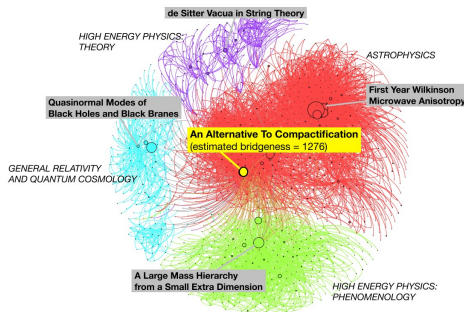


Fig. 1. The discovered community structure in a subgraph of the arXiv citation network (21). The figure shows the top four link communities that include citations to “An alternative to compactification” (22), an article that bridges several communities. We visualize the links between the articles and show some highly cited titles. Each community is labeled with its dominant subject area; nodes are sized by their bridgeness (39), an inferred measure of their impact on multiple communities. This is taken from an analysis of the full 575,000 node network.

<sup>†</sup> Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)

# Modeling Graphs

- Often we wish to understand the underlying structure (e.g., communities/groups/topics) in a graph, predict links, classify nodes, visualize, etc.
- An example graph<sup>†</sup> (a 575,000 node citation network of papers; each node is a paper):

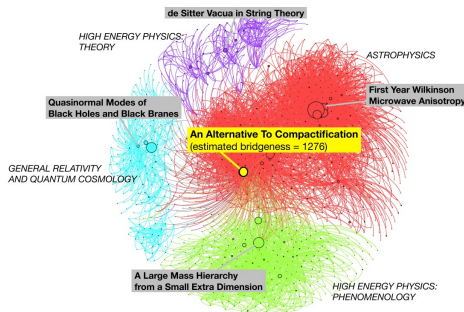


Fig. 1. The discovered community structure in a subgraph of the arXiv citation network (21). The figure shows the top four link communities that include citations to “An alternative to compactification” (22), an article that bridges several communities. We visualize the links between the articles and show some highly cited titles. Each community is labeled with its dominant subject area; nodes are sized by their bridgeness (39), an inferred measure of their impact on multiple communities. This is taken from an analysis of the full 575,000 node network.

- Statistical models of graphs can help us solve these problems

<sup>†</sup> Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)



# Probabilistic Models for Graphs

- Assume each entity/node  $n$  to have a latent representation vector (“embedding”)  $\mathbf{z}_n$  of size  $K$



# Probabilistic Models for Graphs

- Assume each entity/node  $n$  to have a latent representation vector (“embedding”)  $\mathbf{z}_n$  of size  $K$
- We can model each link/non-link  $A_{nm} \in \{0, 1\}$  via a probability model, e.g.

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \theta) = f(\mathbf{z}_n, \mathbf{z}_m, \theta)$$

where  $f$  is some function of  $\mathbf{z}_n, \mathbf{z}_m$  and params  $\theta$ , and returns a probability



# Probabilistic Models for Graphs

- Assume each entity/node  $n$  to have a latent representation vector (“embedding”)  $\mathbf{z}_n$  of size  $K$
- We can model each link/non-link  $A_{nm} \in \{0, 1\}$  via a probability model, e.g.

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \theta) = f(\mathbf{z}_n, \mathbf{z}_m, \theta)$$

where  $f$  is some function of  $\mathbf{z}_n, \mathbf{z}_m$  and params  $\theta$ , and returns a probability

- The overall probability of the observed graph

$$p(\mathbf{A} | \mathbf{Z}, \theta) = \prod_{n,m} p(A_{nm} | \mathbf{z}_n, \mathbf{z}_m, \theta)$$



# Probabilistic Models for Graphs

- Assume each entity/node  $n$  to have a latent representation vector (“embedding”)  $\mathbf{z}_n$  of size  $K$
- We can model each link/non-link  $A_{nm} \in \{0, 1\}$  via a probability model, e.g.

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \theta) = f(\mathbf{z}_n, \mathbf{z}_m, \theta)$$

where  $f$  is some function of  $\mathbf{z}_n, \mathbf{z}_m$  and params  $\theta$ , and returns a probability

- The overall probability of the observed graph

$$p(\mathbf{A} | \mathbf{Z}, \theta) = \prod_{n,m} p(A_{nm} | \mathbf{z}_n, \mathbf{z}_m, \theta)$$

- Various models differ in terms of what the embeddings  $\mathbf{z}_n$  look like, and what  $f$  is defined as



# Probabilistic Models for Graphs

- Assume each entity/node  $n$  to have a latent representation vector (“embedding”)  $\mathbf{z}_n$  of size  $K$
- We can model each link/non-link  $A_{nm} \in \{0, 1\}$  via a probability model, e.g.

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \theta) = f(\mathbf{z}_n, \mathbf{z}_m, \theta)$$

where  $f$  is some function of  $\mathbf{z}_n, \mathbf{z}_m$  and params  $\theta$ , and returns a probability

- The overall probability of the observed graph

$$p(\mathbf{A} | \mathbf{Z}, \theta) = \prod_{n,m} p(A_{nm} | \mathbf{z}_n, \mathbf{z}_m, \theta)$$

- Various models differ in terms of what the embeddings  $\mathbf{z}_n$  look like, and what  $f$  is defined as
- Some representative models





# Probabilistic Models for Graphs

- Assume each entity/node  $n$  to have a latent representation vector (“embedding”)  $\mathbf{z}_n$  of size  $K$
- We can model each link/non-link  $A_{nm} \in \{0, 1\}$  via a probability model, e.g.

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \theta) = f(\mathbf{z}_n, \mathbf{z}_m, \theta)$$

where  $f$  is some function of  $\mathbf{z}_n, \mathbf{z}_m$  and params  $\theta$ , and returns a probability

- The overall probability of the observed graph

$$p(\mathbf{A} | \mathbf{Z}, \theta) = \prod_{n,m} p(A_{nm} | \mathbf{z}_n, \mathbf{z}_m, \theta)$$

- Various models differ in terms of what the embeddings  $\mathbf{z}_n$  look like, and what  $f$  is defined as
- Some representative models
  - Latent Space Model
  - Stochastic Blockmodel
  - Mixed-Membership Blockmodel (MMSB)



# Latent Space Model for Graphs

- LSM<sup>†</sup> assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$

---

<sup>†</sup> Latent space approaches to social network analysis (Hoff et al, 2002)



# Latent Space Model for Graphs

- LSM<sup>†</sup> assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$
- The link probability is defined in terms of the Euclidean similarity between nodes in latent space

---

<sup>†</sup> Latent space approaches to social network analysis (Hoff et al, 2002)



# Latent Space Model for Graphs

- LSM<sup>†</sup> assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$
- The link probability is defined in terms of the Euclidean similarity between nodes in latent space
- One such possible model would be

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m) = \sigma(-||\mathbf{z}_n - \mathbf{z}_m||)$$

---

<sup>†</sup> Latent space approaches to social network analysis (Hoff et al, 2002)



# Latent Space Model for Graphs

- LSM<sup>†</sup> assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$
- The link probability is defined in terms of the Euclidean similarity between nodes in latent space
- One such possible model would be

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m) = \sigma(-||\mathbf{z}_n - \mathbf{z}_m||)$$

- A reasonable model for link prediction

---

<sup>†</sup> Latent space approaches to social network analysis (Hoff et al, 2002)



# Latent Space Model for Graphs

- LSM<sup>†</sup> assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$
- The link probability is defined in terms of the Euclidean similarity between nodes in latent space
- One such possible model would be

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m) = \sigma(-||\mathbf{z}_n - \mathbf{z}_m||)$$

- A reasonable model for link prediction
- However, the real-valued embeddings in LSM don't provide a good interpretability

---

<sup>†</sup> Latent space approaches to social network analysis (Hoff et al, 2002)



# Latent Space Model for Graphs

- LSM<sup>†</sup> assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$
- The link probability is defined in terms of the Euclidean similarity between nodes in latent space
- One such possible model would be

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m) = \sigma(-||\mathbf{z}_n - \mathbf{z}_m||)$$

- A reasonable model for link prediction
- However, the real-valued embeddings in LSM don't provide a good interpretability
  - Therefore not very ideal for discovering clusters etc.

---

<sup>†</sup> Latent space approaches to social network analysis (Hoff et al, 2002)



# Latent Space Model for Graphs

- LSM<sup>†</sup> assumes each node to have a real-valued embedding vector  $\mathbf{z}_n \in \mathbb{R}^K$
- The link probability is defined in terms of the Euclidean similarity between nodes in latent space
- One such possible model would be

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m) = \sigma(-||\mathbf{z}_n - \mathbf{z}_m||)$$

- A reasonable model for link prediction
- However, the real-valued embeddings in LSM don't provide a good interpretability
  - Therefore not very ideal for discovering clusters etc.
- Blockmodels and its variants has such properties as we will see next

---

<sup>†</sup> Latent space approaches to social network analysis (Hoff et al, 2002)





# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

- The full generative model looks like

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

- The full generative model looks like

- For  $n = 1, \dots, N$

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

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

- The full generative model looks like
  - For  $n = 1, \dots, N$   
$$\mathbf{z}_n \sim \text{multinoulli}(\boldsymbol{\pi})$$
  - For  $n = 1, \dots, N$

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

- The full generative model looks like

- For  $n = 1, \dots, N$

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

- For  $n = 1, \dots, N$

- For  $m = 1, \dots, n - 1$

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)





# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

- The full generative model looks like

- For  $n = 1, \dots, N$

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

- For  $n = 1, \dots, N$

- For  $m = 1, \dots, n - 1$

$$A_{nm} \sim \text{Bernoulli}(\eta_{\mathbf{z}_n, \mathbf{z}_m})$$

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Stochastic Blockmodel

- SBM<sup>†</sup> assumes each node belongs to a cluster/community or “block” (total  $K$  clusters)
- The node  $n$ 's cluster membership denoted by a one-hot vector  $\mathbf{z}_n$  of size  $K$
- Assume **probability of link** b/w a node with  $\mathbf{z}_n = k$  and another node with  $\mathbf{z}_m = \ell$  to be  $\eta_{k\ell}$

$$p(A_{nm} = 1 | \mathbf{z}_n, \mathbf{z}_m, \eta) = \eta_{\mathbf{z}_n, \mathbf{z}_m}$$

- The full generative model looks like

- For  $n = 1, \dots, N$

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

- For  $n = 1, \dots, N$

- For  $m = 1, \dots, n - 1$

$$A_{nm} \sim \text{Bernoulli}(\eta_{\mathbf{z}_n, \mathbf{z}_m})$$

- Note: In the fully Bayesian version,  $\pi$  and  $\eta$  can also be given priors

---

<sup>†</sup> Estimation and prediction for stochastic blockmodels for graphs with latent block structure (Snijders and Nowicki, 1997), Discovering latent classes in relational data (Kemp et al, 2004)



# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)



# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)



# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities
- For  $n = 1, \dots, N$

$$\pi_n \sim \text{Dirichlet}(\alpha, \dots, \alpha)$$

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)



# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities
- For  $n = 1, \dots, N$

$$\pi_n \sim \text{Dirichlet}(\alpha, \dots, \alpha)$$

- For  $n = 1, \dots, N$ 
  - For  $m = 1, \dots, n - 1$

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)



# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities
- For  $n = 1, \dots, N$

$$\pi_n \sim \text{Dirichlet}(\alpha, \dots, \alpha)$$

- For  $n = 1, \dots, N$ 
  - For  $m = 1, \dots, K$

$$z_{n \rightarrow m} \sim \text{multinoulli}(\pi_n)$$

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)



# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities
- For  $n = 1, \dots, N$

$$\pi_n \sim \text{Dirichlet}(\alpha, \dots, \alpha)$$

- For  $n = 1, \dots, N$

- For  $m = 1, \dots, n - 1$

$$z_{n \rightarrow m} \sim \text{multinoulli}(\pi_n)$$

$$z_{m \rightarrow n} \sim \text{multinoulli}(\pi_m)$$

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)





# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities
- For  $n = 1, \dots, N$

$$\pi_n \sim \text{Dirichlet}(\alpha, \dots, \alpha)$$

- For  $n = 1, \dots, N$

- For  $m = 1, \dots, n - 1$

$$\mathbf{z}_{n \rightarrow m} \sim \text{multinoulli}(\pi_n)$$

$$\mathbf{z}_{m \rightarrow n} \sim \text{multinoulli}(\pi_m)$$

$$A_{nm} \sim \text{Bernoulli}(\eta_{\mathbf{z}_{n \rightarrow m}, \mathbf{z}_{m \rightarrow n}})$$

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013)



# Mixed-Membership Stochastic Blockmodel

- Unlike SBM, the MMSB<sup>†</sup> assumes each node  $n$  to have a  $K \times 1$  probability vector  $\pi_n$
- $\pi_n$  denotes the probabilities of memberships of node  $n$  in each of the  $K$  communities
- For  $n = 1, \dots, N$

$$\pi_n \sim \text{Dirichlet}(\alpha, \dots, \alpha)$$

- For  $n = 1, \dots, N$

- For  $m = 1, \dots, n - 1$

$$\mathbf{z}_{n \rightarrow m} \sim \text{multinoulli}(\pi_n)$$

$$\mathbf{z}_{m \rightarrow n} \sim \text{multinoulli}(\pi_m)$$

$$A_{nm} \sim \text{Bernoulli}(\eta_{\mathbf{z}_{n \rightarrow m}, \mathbf{z}_{m \rightarrow n}})$$

- Unlike SBM in which node  $n$  has a unique one-hot  $\mathbf{z}_n$  vector, in MMSB, each node  $n$  has an **interaction-specific** cluster assignment

---

<sup>†</sup> Mixed-Membership Stochastic Blockmodel (Airoldi et al, 2008), Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013) ▶ ◀ ≡ ≡ ≡



# Modeling Graphs: Some Other Comments

- A lot of work on various extensions\* of LSM, SBM, MMSB, etc.
- A lot of work on scalable Bayesian inference# in these models (e.g., online MCMC/VI)

---

\* Nonparametric Bayesian modeling of complex networks: an introduction (Schmidt et al, 2013), # Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013), † Graph Convolutional Networks (Kipf and Welling, 2016), ‡ GraphRNN: Generating Realistic Graphs with Deep Auto-regressive Models (You et al, 2018)



# Modeling Graphs: Some Other Comments

- A lot of work on various extensions\* of LSM, SBM, MMSB, etc.
- A lot of work on scalable Bayesian inference# in these models (e.g., online MCMC/VI)
- Some of the recent trends in this area

---

\* Nonparametric Bayesian modeling of complex networks: an introduction (Schmidt et al, 2013), # Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013), † Graph Convolutional Networks (Kipf and Welling, 2016), ‡ GraphRNN: Generating Realistic Graphs with Deep Auto-regressive Models (You et al, 2018)



# Modeling Graphs: Some Other Comments

- A lot of work on various extensions\* of LSM, SBM, MMSB, etc.
- A lot of work on scalable Bayesian inference# in these models (e.g., online MCMC/VI)
- Some of the recent trends in this area
  - Combining classical prob. models of graphs with **graph neural nets** (e.g., Graph Convolutional Net†)

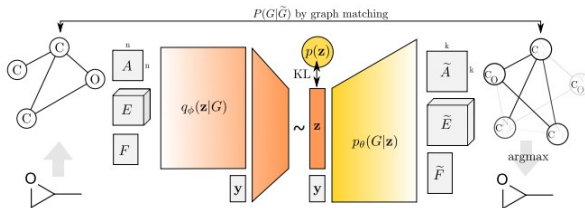
---

\* Nonparametric Bayesian modeling of complex networks: an introduction (Schmidt et al, 2013), # Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013), † Graph Convolutional Networks (Kipf and Welling, 2016), ‡ GraphRNN: Generating Realistic Graphs with Deep Auto-regressive Models (You et al, 2018)



# Modeling Graphs: Some Other Comments

- A lot of work on various extensions\* of LSM, SBM, MMSB, etc.
- A lot of work on scalable Bayesian inference# in these models (e.g., online MCMC/VI)
- Some of the recent trends in this area
  - Combining classical prob. models of graphs with **graph neural nets** (e.g., Graph Convolutional Net<sup>†</sup>)
  - Learning to generate graphs<sup>‡</sup> (just like image or text generation in deep learning)



\* Nonparametric Bayesian modeling of complex networks: an introduction (Schmidt et al, 2013), # Efficient discovery of overlapping communities in massive networks (Gopalan and Blei, 2013), <sup>†</sup> Graph Convolutional Networks (Kipf and Welling, 2016), <sup>‡</sup> GraphRNN: Generating Realistic Graphs with Deep Auto-regressive Models (You et al, 2018)

# Nonparametric Bayesian Modeling

(A way of learning the “right” model size/complexity)



# Motivating Problem: Mixture Models

- A mixture model can be used to cluster/partition the data into multiple groups
- Defined by  $K$  component distributions (e.g.,  $K$  Gaussians for a **Gaussian Mixture Model**)
- Every component distribution  $p(\mathbf{x}|\phi_k)$  has a mixing weight  $\pi_k \in (0, 1)$ , and  $\sum_{k=1}^K \pi_k = 1$





# Motivating Problem: Mixture Models

- A mixture model can be used to cluster/partition the data into multiple groups
- Defined by  $K$  component distributions (e.g.,  $K$  Gaussians for a **Gaussian Mixture Model**)
- Every component distribution  $p(\mathbf{x}|\phi_k)$  has a mixing weight  $\pi_k \in (0, 1)$ , and  $\sum_{k=1}^K \pi_k = 1$
- The distribution of any observation  $\mathbf{x}$

$$p(\mathbf{x}|\pi, \phi) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\phi_k)$$

where  $\pi = \{\pi_k\}_{k=1}^K$  and  $\phi = \{\phi_k\}_{k=1}^K$



# Motivating Problem: Mixture Models

- A mixture model can be used to cluster/partition the data into multiple groups
- Defined by  $K$  component distributions (e.g.,  $K$  Gaussians for a **Gaussian Mixture Model**)
- Every component distribution  $p(\mathbf{x}|\phi_k)$  has a mixing weight  $\pi_k \in (0, 1)$ , and  $\sum_{k=1}^K \pi_k = 1$
- The distribution of any observation  $\mathbf{x}$

$$p(\mathbf{x}|\pi, \phi) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\phi_k)$$

where  $\pi = \{\pi_k\}_{k=1}^K$  and  $\phi = \{\phi_k\}_{k=1}^K$

- Question: What's the “right” number of clusters?



# Motivating Problem: Mixture Models

- A mixture model can be used to cluster/partition the data into multiple groups
- Defined by  $K$  component distributions (e.g.,  $K$  Gaussians for a **Gaussian Mixture Model**)
- Every component distribution  $p(\mathbf{x}|\phi_k)$  has a mixing weight  $\pi_k \in (0, 1)$ , and  $\sum_{k=1}^K \pi_k = 1$
- The distribution of any observation  $\mathbf{x}$

$$p(\mathbf{x}|\pi, \phi) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\phi_k)$$

where  $\pi = \{\pi_k\}_{k=1}^K$  and  $\phi = \{\phi_k\}_{k=1}^K$

- Question: What's the “right” number of clusters? Can do Bayesian model comparison (try different  $K$ 's and compute marginal likelihood for each choice of  $K$ ).



# Motivating Problem: Mixture Models

- A mixture model can be used to cluster/partition the data into multiple groups
- Defined by  $K$  component distributions (e.g.,  $K$  Gaussians for a **Gaussian Mixture Model**)
- Every component distribution  $p(\mathbf{x}|\phi_k)$  has a mixing weight  $\pi_k \in (0, 1)$ , and  $\sum_{k=1}^K \pi_k = 1$
- The distribution of any observation  $\mathbf{x}$

$$p(\mathbf{x}|\pi, \phi) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\phi_k)$$

where  $\pi = \{\pi_k\}_{k=1}^K$  and  $\phi = \{\phi_k\}_{k=1}^K$

- Question: What's the “right” number of clusters? Can do Bayesian model comparison (try different  $K$ 's and compute marginal likelihood for each choice of  $K$ ). But that can be expensive.



# Motivating Problem: Mixture Models

- A mixture model can be used to cluster/partition the data into multiple groups
- Defined by  $K$  component distributions (e.g.,  $K$  Gaussians for a **Gaussian Mixture Model**)
- Every component distribution  $p(\mathbf{x}|\phi_k)$  has a mixing weight  $\pi_k \in (0, 1)$ , and  $\sum_{k=1}^K \pi_k = 1$
- The distribution of any observation  $\mathbf{x}$

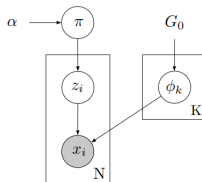
$$p(\mathbf{x}|\pi, \phi) = \sum_{k=1}^K \pi_k p(\mathbf{x}|\phi_k)$$

where  $\pi = \{\pi_k\}_{k=1}^K$  and  $\phi = \{\phi_k\}_{k=1}^K$

- Question: What's the “right” number of clusters? Can do Bayesian model comparison (try different  $K$ 's and compute marginal likelihood for each choice of  $K$ ). But that can be expensive.
- How about having a **single model** but allowing the number of clusters to “grow” with data?



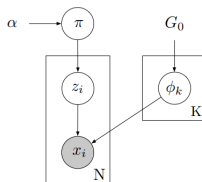
# Prelude: A Bayesian Mixture Model (with fixed $K$ )



- Assuming some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the generative model can be defined as follows



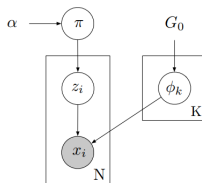
# Prelude: A Bayesian Mixture Model (with fixed $K$ )



- Assuming some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the generative model can be defined as follows
  - Draw mixture proportion vector  $\pi = [\pi_1, \dots, \pi_K]$  from the prior [Dirichlet\( \$\alpha/K, \dots, \alpha/K\$ \)](#)



# Prelude: A Bayesian Mixture Model (with fixed $K$ )

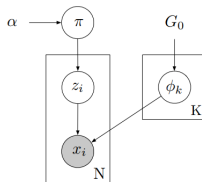


- Assuming some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the generative model can be defined as follows
  - Draw mixture proportion vector  $\pi = [\pi_1, \dots, \pi_K]$  from the prior [Dirichlet](#)( $\alpha/K, \dots, \alpha/K$ )
  - Draw parameters  $\{\phi_k\}_{k=1}^K$  of each mixture component i.i.d. from a prior “base distribution”  $G_0$  (note: choice of  $G_0$  depends on what the component distributions are; e.g., for Gaussians,  $G_0$  can be NIW)





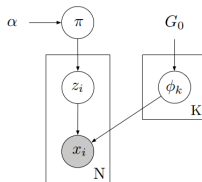
# Prelude: A Bayesian Mixture Model (with fixed $K$ )



- Assuming some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the generative model can be defined as follows
  - Draw mixture proportion vector  $\pi = [\pi_1, \dots, \pi_K]$  from the prior [Dirichlet](#)( $\alpha/K, \dots, \alpha/K$ )
  - Draw parameters  $\{\phi_k\}_{k=1}^K$  of each mixture component i.i.d. from a prior “base distribution”  $G_0$  (note: choice of  $G_0$  depends on what the component distributions are; e.g., for Gaussians,  $G_0$  can be NIW)
  - Draw the data: For each observation  $i = 1, \dots, N$



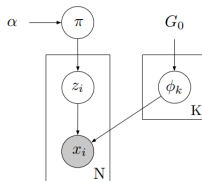
# Prelude: A Bayesian Mixture Model (with fixed $K$ )



- Assuming some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the generative model can be defined as follows
  - Draw mixture proportion vector  $\pi = [\pi_1, \dots, \pi_K]$  from the prior [Dirichlet](#)( $\alpha/K, \dots, \alpha/K$ )
  - Draw parameters  $\{\phi_k\}_{k=1}^K$  of each mixture component i.i.d. from a prior “base distribution”  $G_0$  (note: choice of  $G_0$  depends on what the component distributions are; e.g., for Gaussians,  $G_0$  can be NIW)
  - Draw the data: For each observation  $i = 1, \dots, N$ 
    - Draw a cluster id  $z_i \in \{1, \dots, K\}$  from [multinoulli](#)( $\pi$ )



# Prelude: A Bayesian Mixture Model (with fixed $K$ )



- Assuming some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the generative model can be defined as follows
  - Draw mixture proportion vector  $\pi = [\pi_1, \dots, \pi_K]$  from the prior [Dirichlet](#)( $\alpha/K, \dots, \alpha/K$ )
  - Draw parameters  $\{\phi_k\}_{k=1}^K$  of each mixture component i.i.d. from a prior “base distribution”  $G_0$  (note: choice of  $G_0$  depends on what the component distributions are; e.g., for Gaussians,  $G_0$  can be NIW)
  - Draw the data: For each observation  $i = 1, \dots, N$ 
    - Draw a cluster id  $z_i \in \{1, \dots, K\}$  from [multinoulli](#)( $\pi$ )
    - Suppose  $z_i = k$ . Draw  $\mathbf{x}_i$  from [p](#)( $\mathbf{x}|\phi_k$ )



# Bayesian Inference for Mixture Models: A Gibbs Sampler

- Assume a Gaussian Mixture Model (GMM) with  $\phi_k = (\mu_k, \Sigma_k)$ . Thus  $G_0$  can be NIW



# Bayesian Inference for Mixture Models: A Gibbs Sampler

- Assume a Gaussian Mixture Model (GMM) with  $\phi_k = (\mu_k, \Sigma_k)$ . Thus  $G_0$  can be NIW
- Due to conjugacy, we can easily derive a Gibbs sampler for this model



# Bayesian Inference for Mixture Models: A Gibbs Sampler

- Assume a Gaussian Mixture Model (GMM) with  $\phi_k = (\mu_k, \Sigma_k)$ . Thus  $G_0$  can be NIW
- Due to conjugacy, we can easily derive a Gibbs sampler for this model
- The basic Gibbs sampler is sketched below



# Bayesian Inference for Mixture Models: A Gibbs Sampler

- Assume a Gaussian Mixture Model (GMM) with  $\phi_k = (\mu_k, \Sigma_k)$ . Thus  $G_0$  can be NIW
- Due to conjugacy, we can easily derive a Gibbs sampler for this model
- The basic Gibbs sampler is sketched below
  - Randomly initialize  $Z, \pi, \phi$ . Repeat until we have enough samples



# Bayesian Inference for Mixture Models: A Gibbs Sampler

- Assume a Gaussian Mixture Model (GMM) with  $\phi_k = (\mu_k, \Sigma_k)$ . Thus  $G_0$  can be NIW
- Due to conjugacy, we can easily derive a Gibbs sampler for this model
- The basic Gibbs sampler is sketched below
  - Randomly initialize  $Z$ ,  $\pi$ ,  $\phi$ . Repeat until we have enough samples
    - Sample  $\pi$  as (due to Dirichlet-multinomial conjugacy)

$$\pi | \mathbf{Z} \sim \text{Dirichlet}(n_1 + \alpha/K, \dots, n_K + \alpha/K)$$

.. where  $n_k$  is the number of observations currently assigned to cluster  $k$





# Bayesian Inference for Mixture Models: A Gibbs Sampler

- Assume a Gaussian Mixture Model (GMM) with  $\phi_k = (\mu_k, \Sigma_k)$ . Thus  $G_0$  can be NIW
- Due to conjugacy, we can easily derive a Gibbs sampler for this model
- The basic Gibbs sampler is sketched below
  - Randomly initialize  $Z, \pi, \phi$ . Repeat until we have enough samples

- Sample  $\pi$  as (due to Dirichlet-multinomial conjugacy)

$$\pi | \mathbf{Z} \sim \text{Dirichlet}(n_1 + \alpha/K, \dots, n_K + \alpha/K)$$

.. where  $n_k$  is the number of observations currently assigned to cluster  $k$

- Sample each  $\mathbf{z}_i$  using

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \pi, \phi, \mathbf{X}) \propto p(\mathbf{z}_i | \pi) \times p(\mathbf{x}_i | \phi_k) = \pi_k \times p(\mathbf{x}_i | \phi_k)$$

.. (note: the above is equivalent to computing the above posterior probabilities for  $k = 1, \dots, K$  and drawing  $\mathbf{z}_i$  from a multinoulli with probability parameter vector given by these posterior probabilities)



# Bayesian Inference for Mixture Models: A Gibbs Sampler

- Assume a Gaussian Mixture Model (GMM) with  $\phi_k = (\mu_k, \Sigma_k)$ . Thus  $G_0$  can be NIW
- Due to conjugacy, we can easily derive a Gibbs sampler for this model
- The basic Gibbs sampler is sketched below
  - Randomly initialize  $Z, \pi, \phi$ . Repeat until we have enough samples

- Sample  $\pi$  as (due to Dirichlet-multinomial conjugacy)

$$\pi | \mathbf{Z} \sim \text{Dirichlet}(n_1 + \alpha/K, \dots, n_K + \alpha/K)$$

.. where  $n_k$  is the number of observations currently assigned to cluster  $k$

- Sample each  $\mathbf{z}_i$  using

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \pi, \phi, \mathbf{X}) \propto p(\mathbf{z}_i | \pi) \times p(\mathbf{x}_i | \phi_k) = \pi_k \times p(\mathbf{x}_i | \phi_k)$$

.. (note: the above is equivalent to computing the above posterior probabilities for  $k = 1, \dots, K$  and drawing  $\mathbf{z}_i$  from a multinoulli with probability parameter vector given by these posterior probabilities)

- Sample  $\phi_k = (\mu_k, \Sigma_k)$  from NIW posterior given  $\mathbf{Z}$  and  $\mathbf{X}$  (NIW posterior also has a closed form).

# A “Collapsed” Gibbs Sampler

- Let's integrate out  $\pi$  (due to Dirichlet-multinomial conjugacy) from prior  $p(\mathbf{z}_i|\pi)$  and eliminate  $\pi$



# A “Collapsed” Gibbs Sampler

- Let's integrate out  $\pi$  (due to Dirichlet-multinomial conjugacy) from prior  $p(\mathbf{z}_i|\pi)$  and eliminate  $\pi$

$$p(\mathbf{z}_i = k|\alpha, \mathbf{Z}_{-i}) = \int p(\mathbf{z}_i = k|\pi)p(\pi|\mathbf{Z}_{-i})d\pi = \int \pi_k p(\pi|\mathbf{Z}_{-i})d\pi = \mathbb{E}_{p(\pi|\mathbf{Z}_{-i})}[\pi_k] = \frac{n_k + \alpha/K}{\alpha + N - 1}$$

.. where  $n_k$  is the number of examples (other than  $\mathbf{x}_i$ ) assigned to cluster  $k$



# A “Collapsed” Gibbs Sampler

- Let's integrate out  $\pi$  (due to Dirichlet-multinomial conjugacy) from prior  $p(\mathbf{z}_i|\pi)$  and eliminate  $\pi$

$$p(\mathbf{z}_i = k|\alpha, \mathbf{Z}_{-i}) = \int p(\mathbf{z}_i = k|\pi)p(\pi|\mathbf{Z}_{-i})d\pi = \int \pi_k p(\pi|\mathbf{Z}_{-i})d\pi = \mathbb{E}_{p(\pi|\mathbf{Z}_{-i})}[\pi_k] = \frac{n_k + \alpha/K}{\alpha + N - 1}$$

.. where  $n_k$  is the number of examples (other than  $\mathbf{x}_i$ ) assigned to cluster  $k$

- The collapsed Gibbs sampler is sketched below



# A “Collapsed” Gibbs Sampler

- Let's integrate out  $\pi$  (due to Dirichlet-multinomial conjugacy) from prior  $p(\mathbf{z}_i|\pi)$  and eliminate  $\pi$

$$p(\mathbf{z}_i = k|\alpha, \mathbf{Z}_{-i}) = \int p(\mathbf{z}_i = k|\pi)p(\pi|\mathbf{Z}_{-i})d\pi = \int \pi_k p(\pi|\mathbf{Z}_{-i})d\pi = \mathbb{E}_{p(\pi|\mathbf{Z}_{-i})}[\pi_k] = \frac{n_k + \alpha/K}{\alpha + N - 1}$$

.. where  $n_k$  is the number of examples (other than  $\mathbf{x}_i$ ) assigned to cluster  $k$

- The collapsed Gibbs sampler is sketched below
  - Randomly initialize  $Z, \phi$ . Repeat until we have enough samples



# A “Collapsed” Gibbs Sampler

- Let's integrate out  $\pi$  (due to Dirichlet-multinomial conjugacy) from prior  $p(\mathbf{z}_i|\pi)$  and eliminate  $\pi$

$$p(\mathbf{z}_i = k|\alpha, \mathbf{Z}_{-i}) = \int p(\mathbf{z}_i = k|\pi)p(\pi|\mathbf{Z}_{-i})d\pi = \int \pi_k p(\pi|\mathbf{Z}_{-i})d\pi = \mathbb{E}_{p(\pi|\mathbf{Z}_{-i})}[\pi_k] = \frac{n_k + \alpha/K}{\alpha + N - 1}$$

.. where  $n_k$  is the number of examples (other than  $\mathbf{x}_i$ ) assigned to cluster  $k$

- The collapsed Gibbs sampler is sketched below
  - Randomly initialize  $\mathbf{Z}$ ,  $\phi$ . Repeat until we have enough samples
  - Sample each  $\mathbf{z}_i$  as

$$p(\mathbf{z}_i = k|\mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i|\phi_k)$$

.. (note: just like the uncollapsed case, this is equivalent to drawing  $\mathbf{z}_i$  from a multinoulli with probability parameter vector given by the above posterior probabilities evaluated for  $k = 1, \dots, K$ )



# A “Collapsed” Gibbs Sampler

- Let's integrate out  $\pi$  (due to Dirichlet-multinomial conjugacy) from prior  $p(\mathbf{z}_i|\pi)$  and eliminate  $\pi$

$$p(\mathbf{z}_i = k|\alpha, \mathbf{Z}_{-i}) = \int p(\mathbf{z}_i = k|\pi)p(\pi|\mathbf{Z}_{-i})d\pi = \int \pi_k p(\pi|\mathbf{Z}_{-i})d\pi = \mathbb{E}_{p(\pi|\mathbf{Z}_{-i})}[\pi_k] = \frac{n_k + \alpha/K}{\alpha + N - 1}$$

.. where  $n_k$  is the number of examples (other than  $\mathbf{x}_i$ ) assigned to cluster  $k$

- The collapsed Gibbs sampler is sketched below
  - Randomly initialize  $\mathbf{Z}$ ,  $\phi$ . Repeat until we have enough samples

- Sample each  $\mathbf{z}_i$  as

$$p(\mathbf{z}_i = k|\mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i|\phi_k)$$

.. (note: just like the uncollapsed case, this is equivalent to drawing  $\mathbf{z}_i$  from a multinoulli with probability parameter vector given by the above posterior probabilities evaluated for  $k = 1, \dots, K$ )

- Sample  $\phi_k = (\mu_k, \Sigma_k)$  from NIW posterior given  $\mathbf{Z}$  and  $\mathbf{X}$





Nonparametric Bayesian Mixture Model  
(.. which you get when you allow **unbounded**  $K$ )  
(.. i.e., you allow  $K \rightarrow \infty$ )



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- For the collapsed Gibbs sampler, we saw that  $p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i | \phi_k)$



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- For the collapsed Gibbs sampler, we saw that  $p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i | \phi_k)$
- As  $K \rightarrow \infty$ , the probability that  $\mathbf{x}_i$  will be assigned to an **existing cluster** (say  $k$ ) is given by

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto n_k \times p(\mathbf{x}_i | \phi_k)$$



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- For the collapsed Gibbs sampler, we saw that  $p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i | \phi_k)$
- As  $K \rightarrow \infty$ , the probability that  $\mathbf{x}_i$  will be assigned to an **existing cluster** (say  $k$ ) is given by

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto n_k \times p(\mathbf{x}_i | \phi_k)$$

- It's proportional to the no. of obs. already in cluster  $k$  (it's like a “rich gets richer” tendency)



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- For the collapsed Gibbs sampler, we saw that  $p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i | \phi_k)$
- As  $K \rightarrow \infty$ , the probability that  $\mathbf{x}_i$  will be assigned to an **existing cluster** (say  $k$ ) is given by

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto n_k \times p(\mathbf{x}_i | \phi_k)$$

- It's proportional to the no. of obs. already in cluster  $k$  (it's like a “rich gets richer” tendency)
- Now suppose there are a total of  $K_+$  occupied clusters (with number of data points  $\geq 1$ )



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- For the collapsed Gibbs sampler, we saw that  $p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i | \phi_k)$
- As  $K \rightarrow \infty$ , the probability that  $\mathbf{x}_i$  will be assigned to an **existing cluster** (say  $k$ ) is given by

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto n_k \times p(\mathbf{x}_i | \phi_k)$$

- It's proportional to the no. of obs. already in cluster  $k$  (it's like a “rich gets richer” tendency)
- Now suppose there are a total of  $K_+$  occupied clusters (with number of data points  $\geq 1$ )
- Number of empty clusters =  $K - K_+$ . Can think of all of these as a single unoccupied cluster



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- For the collapsed Gibbs sampler, we saw that  $p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto (n_k + \alpha/K)p(\mathbf{x}_i | \phi_k)$
- As  $K \rightarrow \infty$ , the probability that  $\mathbf{x}_i$  will be assigned to an **existing cluster** (say  $k$ ) is given by

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto n_k \times p(\mathbf{x}_i | \phi_k)$$

- It's proportional to the no. of obs. already in cluster  $k$  (it's like a “rich gets richer” tendency)
- Now suppose there are a total of  $K_+$  occupied clusters (with number of data points  $\geq 1$ )
- Number of empty clusters =  $K - K_+$ . Can think of all of these as a single unoccupied cluster
- The probability of  $\mathbf{x}_i$  being assigned to this new (so far empty) cluster

$$p(\mathbf{z}_i = k_{\text{new}} | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto \alpha \times p(\mathbf{x}_i | G_0) \quad (\text{WHY? Reason on the next slide!})$$

$$\text{where } p(\mathbf{x}_i | G_0) = \int p(\mathbf{x}_i | \phi_{\text{new}}) p(\phi_{\text{new}} | G_0) d\phi_{\text{new}}$$



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- The probability of  $\mathbf{x}_i$  being assigned to the new (so far empty) cluster

$$p(\mathbf{z}_i = k_{\text{new}} | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto \alpha \times p(\mathbf{x}_i | G_0)$$





# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- The probability of  $\mathbf{x}_i$  being assigned to the new (so far empty) cluster

$$p(\mathbf{z}_i = k_{\text{new}} | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto \alpha \times p(\mathbf{x}_i | G_0)$$

- The  $\alpha$  above is due to the fact that, for the conditional prior on  $\mathbf{z}_i$ , as  $K \rightarrow \infty$  we have

$$p(\mathbf{z}_i = k_{\text{new}} | \alpha, \mathbf{Z}_{-i}) = \frac{0 + (\alpha/K) \times (K - K_+)}{\alpha + N - 1} \rightarrow \frac{\alpha}{\alpha + N - 1}$$



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- The probability of  $\mathbf{x}_i$  being assigned to the new (so far empty) cluster

$$p(\mathbf{z}_i = k_{\text{new}} | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto \alpha \times p(\mathbf{x}_i | G_0)$$

- The  $\alpha$  above is due to the fact that, for the conditional prior on  $\mathbf{z}_i$ , as  $K \rightarrow \infty$  we have

$$p(\mathbf{z}_i = k_{\text{new}} | \alpha, \mathbf{Z}_{-i}) = \frac{0 + (\alpha/K) \times (K - K_+)}{\alpha + N - 1} \rightarrow \frac{\alpha}{\alpha + N - 1}$$

- Note: Another way -  $p(\mathbf{z}_i = k_{\text{new}} | \alpha, \mathbf{Z}_{-i}) = 1 - \sum_{k=1}^{K_+} \frac{n_k}{\alpha + N - 1} = \frac{\alpha}{\alpha + N - 1}$  (since  $\sum_{k=1}^{K_+} n_k = N - 1$ )



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- The probability of  $\mathbf{x}_i$  being assigned to the new (so far empty) cluster

$$p(\mathbf{z}_i = k_{\text{new}} | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto \alpha \times p(\mathbf{x}_i | G_0)$$

- The  $\alpha$  above is due to the fact that, for the conditional prior on  $\mathbf{z}_i$ , as  $K \rightarrow \infty$  we have

$$p(\mathbf{z}_i = k_{\text{new}} | \alpha, \mathbf{Z}_{-i}) = \frac{0 + (\alpha/K) \times (K - K_+)}{\alpha + N - 1} \rightarrow \frac{\alpha}{\alpha + N - 1}$$

- Note: Another way -  $p(\mathbf{z}_i = k_{\text{new}} | \alpha, \mathbf{Z}_{-i}) = 1 - \sum_{k=1}^{K_+} \frac{n_k}{\alpha + N - 1} = \frac{\alpha}{\alpha + N - 1}$  (since  $\sum_{k=1}^{K_+} n_k = N - 1$ )
- Also, instead of the likelihood, we use the **marginal likelihood**

$$p(\mathbf{x}_i | G_0) = \int p(\mathbf{x}_i | \phi_{\text{new}}) p(\phi_{\text{new}} | G_0) d\phi_{\text{new}}$$

.. because the new cluster hasn't yet been created and thus we don't have its  $\phi_{k_{\text{new}}}$



# Let's Have $K \rightarrow \infty$ (i.e., Unbounded)

- The probability of  $\mathbf{x}_i$  being assigned to the new (so far empty) cluster

$$p(\mathbf{z}_i = k_{\text{new}} | \mathbf{Z}_{-i}, \phi, \mathbf{X}) \propto \alpha \times p(\mathbf{x}_i | G_0)$$

- The  $\alpha$  above is due to the fact that, for the conditional prior on  $\mathbf{z}_i$ , as  $K \rightarrow \infty$  we have

$$p(\mathbf{z}_i = k_{\text{new}} | \alpha, \mathbf{Z}_{-i}) = \frac{0 + (\alpha/K) \times (K - K_+)}{\alpha + N - 1} \rightarrow \frac{\alpha}{\alpha + N - 1}$$

- Note: Another way -  $p(\mathbf{z}_i = k_{\text{new}} | \alpha, \mathbf{Z}_{-i}) = 1 - \sum_{k=1}^{K_+} \frac{n_k}{\alpha + N - 1} = \frac{\alpha}{\alpha + N - 1}$  (since  $\sum_{k=1}^{K_+} n_k = N - 1$ )
- Also, instead of the likelihood, we use the **marginal likelihood**

$$p(\mathbf{x}_i | G_0) = \int p(\mathbf{x}_i | \phi_{\text{new}}) p(\phi_{\text{new}} | G_0) d\phi_{\text{new}}$$

.. because the new cluster hasn't yet been created and thus we don't have its  $\phi_{k_{\text{new}}}$

- Note: Once the new cluster has been created (after a data point has been assigned to it), we also have to sample for  $\phi_{\text{new}}$  from its posterior.

# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )



# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$



# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$

# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$ 
  - For each observation  $i = 1, \dots, N$ , sample the cluster id  $z_i$





# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$ 
  - For each observation  $i = 1, \dots, N$ , sample the cluster id  $z_i$

$$p(z_i = k | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t)}, \mathbf{X}) \propto n_k^{(t-1)} \times p(\mathbf{x}_i | \phi_k^{(t-1)}) = \hat{\pi}_{ik} \quad (k = 1, \dots, K)$$



# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$ 
  - For each observation  $i = 1, \dots, N$ , sample the cluster id  $z_i$

$$p(z_i = k | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t)}, \mathbf{X}) \propto n_k^{(t-1)} \times p(\mathbf{x}_i | \phi_k^{(t-1)}) = \hat{\pi}_{ik} \quad (k = 1, \dots, K)$$

$$p(z_i = k_{\text{new}} | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t-1)}, \mathbf{X}) \propto \alpha^{(t-1)} \times p(\mathbf{x}_i | G_0) = \hat{\pi}_{ik_{\text{new}}}$$



# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$ 
  - For each observation  $i = 1, \dots, N$ , sample the cluster id  $z_i$

$$p(z_i = k | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t)}, \mathbf{X}) \propto n_k^{(t-1)} \times p(\mathbf{x}_i | \phi_k^{(t-1)}) = \hat{\pi}_{ik} \quad (k = 1, \dots, K)$$

$$p(z_i = k_{\text{new}} | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t-1)}, \mathbf{X}) \propto \alpha^{(t-1)} \times p(\mathbf{x}_i | G_0) = \hat{\pi}_{ik_{\text{new}}}$$

$$\mathbf{z}_i^{(t)} \sim \text{multinoulli}(\hat{\pi}_{i1}, \hat{\pi}_{i2}, \dots, \hat{\pi}_{ik_{\text{new}}})$$



# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$ 
  - For each observation  $i = 1, \dots, N$ , sample the cluster id  $z_i$

$$p(z_i = k | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t)}, \mathbf{X}) \propto n_k^{(t-1)} \times p(\mathbf{x}_i | \phi_k^{(t-1)}) = \hat{\pi}_{ik} \quad (k = 1, \dots, K)$$

$$p(z_i = k_{\text{new}} | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t-1)}, \mathbf{X}) \propto \alpha^{(t-1)} \times p(\mathbf{x}_i | G_0) = \hat{\pi}_{ik_{\text{new}}}$$

$$\mathbf{z}_i^{(t)} \sim \text{multinoulli}(\hat{\pi}_{i1}, \hat{\pi}_{i2}, \dots, \hat{\pi}_{ik_{\text{new}}})$$

$$\text{set } K = K + 1 \quad (\text{if } \mathbf{x}_i \text{ assigned to a new cluster})$$



# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$ 
  - For each observation  $i = 1, \dots, N$ , sample the cluster id  $z_i$

$$p(z_i = k | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t)}, \mathbf{X}) \propto n_k^{(t-1)} \times p(\mathbf{x}_i | \phi_k^{(t-1)}) = \hat{\pi}_{ik} \quad (k = 1, \dots, K)$$

$$p(z_i = k_{\text{new}} | \mathbf{Z}_{-i}^{(t-1)}, \phi^{(t-1)}, \mathbf{X}) \propto \alpha^{(t-1)} \times p(\mathbf{x}_i | G_0) = \hat{\pi}_{ik_{\text{new}}}$$

$$\mathbf{z}_i^{(t)} \sim \text{multinoulli}(\hat{\pi}_{i1}, \hat{\pi}_{i2}, \dots, \hat{\pi}_{ik_{\text{new}}})$$

$$\text{set } K = K + 1 \quad (\text{if } \mathbf{x}_i \text{ assigned to a new cluster})$$

- Sample the mixture component parameters  $\{\phi_k^{(t)}\}_{k=1}^K$  and  $\alpha^{(t)}$  from the respective CPs



# Gibbs Sampler for Nonparametric Bayesian Mixture Model

A brief sketch of a basic Gibbs sampler (samples  $\mathbf{Z}$  and  $\{\phi_k\}_{k=1}^K$ ) for this model with unbounded  $K$  (note: The mixing proportions  $\pi_k$ 's were collapsed from the prior  $p(\mathbf{z}_i|\pi)$ )

## Gibbs Sampler for NPBayes Mixture Model

- Set an initial  $K$ . Initialize  $\mathbf{Z}^{(0)}$  and  $\{\phi_k^{(0)}\}_{k=1}^K$
- For  $t = 1, \dots, T$ 
  - For each observation  $i = 1, \dots, N$ , sample the cluster id  $\mathbf{z}_i$

$$p(\mathbf{z}_i = k | \mathbf{Z}_{-i}^{(t-1)}, \phi_k^{(t)}, \mathbf{X}) \propto n_k^{(t-1)} \times p(\mathbf{x}_i | \phi_k^{(t-1)}) = \hat{\pi}_{ik} \quad (k = 1, \dots, K)$$

$$p(\mathbf{z}_i = k_{\text{new}} | \mathbf{Z}_{-i}^{(t-1)}, \phi_k^{(t-1)}, \mathbf{X}) \propto \alpha^{(t-1)} \times p(\mathbf{x}_i | G_0) = \hat{\pi}_{ik_{\text{new}}}$$

$$\mathbf{z}_i^{(t)} \sim \text{multinoulli}(\hat{\pi}_{i1}, \hat{\pi}_{i2}, \dots, \hat{\pi}_{ik_{\text{new}}})$$

$$\text{set } K = K + 1 \quad (\text{if } \mathbf{x}_i \text{ assigned to a new cluster})$$

- Sample the mixture component parameters  $\{\phi_k^{(t)}\}_{k=1}^K$  and  $\alpha^{(t)}$  from the respective CPs

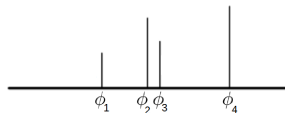
Note: "Markov Chain Sampling Methods for Dirichlet Process Mixture Models" (Neal, 2000) is an excellent reference for various MCMC sampling algorithms for nonparametric Bayesian mixture models (including collapsed versions that don't require sampling for  $\{\phi_k\}_{k=1}^K$ )

# Nonparametric Bayesian Mixture Models (A More Formal Perspective..)



# A Gentle Start: A Discrete Distribution

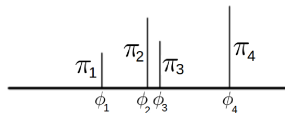
- Assume some space  $\Omega$  (e.g., the real line) and  $K$  “locations”  $\phi_1, \dots, \phi_K$  in that space





# A Gentle Start: A Discrete Distribution

- Assume some space  $\Omega$  (e.g., the real line) and  $K$  “locations”  $\phi_1, \dots, \phi_K$  in that space

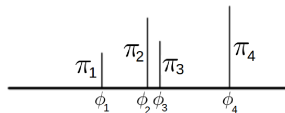


- Assume these locations have “weights”  $\pi_1, \dots, \pi_K$  where  $\pi_k \in (0, 1), \forall k$  and  $\sum_{k=1}^K \pi_k = 1$



# A Gentle Start: A Discrete Distribution

- Assume some space  $\Omega$  (e.g., the real line) and  $K$  “locations”  $\phi_1, \dots, \phi_K$  in that space

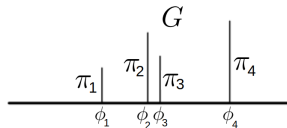


- Assume these locations have “weights”  $\pi_1, \dots, \pi_K$  where  $\pi_k \in (0, 1), \forall k$  and  $\sum_{k=1}^K \pi_k = 1$ 
  - Can think of  $\pi_k$  as how “popular” location  $\phi_k$  is



# A Gentle Start: A Discrete Distribution

- Assume some space  $\Omega$  (e.g., the real line) and  $K$  “locations”  $\phi_1, \dots, \phi_K$  in that space



- Assume these locations have “weights”  $\pi_1, \dots, \pi_K$  where  $\pi_k \in (0, 1), \forall k$  and  $\sum_{k=1}^K \pi_k = 1$ 
  - Can think of  $\pi_k$  as how “popular” location  $\phi_k$  is
- Then we can define a **discrete distribution**  $G$  as

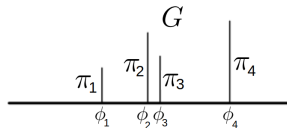
$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

where  $\delta_{\phi_k}$  is an “atom” or point-mass at location  $\phi_k$  ( $\delta_{\phi_k}(\phi) = 1$  iff  $\phi = \phi_k$ , and 0 otherwise)



# A Gentle Start: A Discrete Distribution

- Assume some space  $\Omega$  (e.g., the real line) and  $K$  “locations”  $\phi_1, \dots, \phi_K$  in that space



- Assume these locations have “weights”  $\pi_1, \dots, \pi_K$  where  $\pi_k \in (0, 1), \forall k$  and  $\sum_{k=1}^K \pi_k = 1$ 
  - Can think of  $\pi_k$  as how “popular” location  $\phi_k$  is
- Then we can define a **discrete distribution**  $G$  as

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

where  $\delta_{\phi_k}$  is an “atom” or point-mass at location  $\phi_k$  ( $\delta_{\phi_k}(\phi) = 1$  iff  $\phi = \phi_k$ , and 0 otherwise)

- Important: The support of this discrete distribution  $G$  is  $\{\phi_k\}_{k=1}^K$



# A Bayesian Construction of G

- Let's define appropriate priors on  $\{\phi_k\}_{k=1}^K$  and  $\{\pi_k\}_{k=1}^K$

$$(\pi_1, \dots, \pi_K) \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$



# A Bayesian Construction of G

- Let's define appropriate priors on  $\{\phi_k\}_{k=1}^K$  and  $\{\pi_k\}_{k=1}^K$

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \quad k = 1, \dots, K\end{aligned}$$



# A Bayesian Construction of $G$

- Let's define appropriate priors on  $\{\phi_k\}_{k=1}^K$  and  $\{\pi_k\}_{k=1}^K$

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \quad k = 1, \dots, K\end{aligned}$$

- $G_0$  (its support being  $\Omega$ ) is a “base distribution” (the choice depends on what  $\phi_k$ 's are)



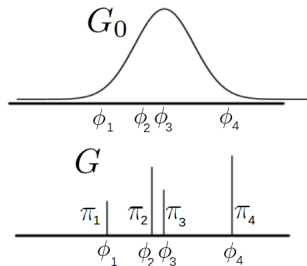
# A Bayesian Construction of G

- Let's define appropriate priors on  $\{\phi_k\}_{k=1}^K$  and  $\{\pi_k\}_{k=1}^K$

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \quad k = 1, \dots, K\end{aligned}$$

- $G_0$  (its support being  $\Omega$ ) is a “base distribution” (the choice depends on what  $\phi_k$ 's are)

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \quad k = 1, \dots, K \\ G &= \sum_{k=1}^K \pi_k \delta_{\phi_k}\end{aligned}$$





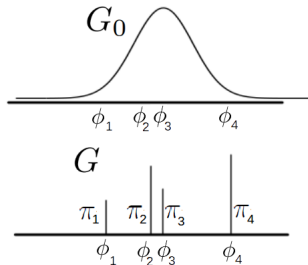
# A Bayesian Construction of $G$

- Let's define appropriate priors on  $\{\phi_k\}_{k=1}^K$  and  $\{\pi_k\}_{k=1}^K$

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \quad k = 1, \dots, K\end{aligned}$$

- $G_0$  (its support being  $\Omega$ ) is a “base distribution” (the choice depends on what  $\phi_k$ 's are)

$$\begin{aligned}(\pi_1, \dots, \pi_K) &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \quad k = 1, \dots, K \\ G &= \sum_{k=1}^K \pi_k \delta_{\phi_k}\end{aligned}$$

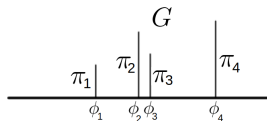


- Note that  $G$  is now a **random** distribution (or a random measure)



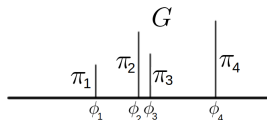
# Discrete Distributions Induce Clustering!

- Drawing values repeatedly from a discrete distribution leads to repetitions



# Discrete Distributions Induce Clustering!

- Drawing values repeatedly from a discrete distribution leads to repetitions

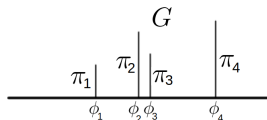


- E.g., Drawing 5 times from the above  $G$  is **guaranteed** to have **at least one repetition**



# Discrete Distributions Induce Clustering!

- Drawing values repeatedly from a discrete distribution leads to repetitions



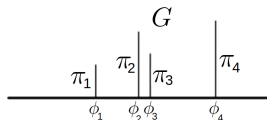
- E.g., Drawing 5 times from the above  $G$  is **guaranteed** to have **at least one repetition**
- Suppose we draw  $N > K$  “parameters”  $\theta_1, \dots, \theta_N$  i.i.d. from  $G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$

$$\theta_i \sim G \quad i = 1, \dots, N$$



# Discrete Distributions Induce Clustering!

- Drawing values repeatedly from a discrete distribution leads to repetitions



- E.g., Drawing 5 times from the above  $G$  is **guaranteed** to have **at least one repetition**
- Suppose we draw  $N > K$  “parameters”  $\theta_1, \dots, \theta_N$  i.i.d. from  $G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$

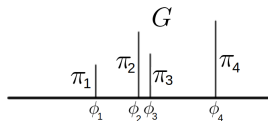
$$\theta_i \sim G \quad i = 1, \dots, N$$

.. then the collection  $(\theta_1, \dots, \theta_N)$  will have at most  $K$  unique parameters  $(\phi_1, \dots, \phi_K)$



# Discrete Distributions Induce Clustering!

- Drawing values repeatedly from a discrete distribution leads to repetitions



- E.g., Drawing 5 times from the above  $G$  is **guaranteed** to have **at least one repetition**
- Suppose we draw  $N > K$  “parameters”  $\theta_1, \dots, \theta_N$  i.i.d. from  $G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$

$$\theta_i \sim G \quad i = 1, \dots, N$$

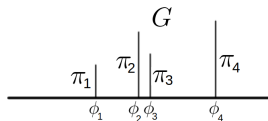
.. then the collection  $(\theta_1, \dots, \theta_N)$  will have at most  $K$  unique parameters  $(\phi_1, \dots, \phi_K)$

- Thus  $G$  induces a **clustering** of parameters  $\theta_i$ ’s, s.t.  $\theta_i$ ’s within any group are all **identical**



# Discrete Distributions Induce Clustering!

- Drawing values repeatedly from a discrete distribution leads to repetitions



- E.g., Drawing 5 times from the above  $G$  is **guaranteed** to have **at least one repetition**
- Suppose we draw  $N > K$  “parameters”  $\theta_1, \dots, \theta_N$  i.i.d. from  $G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$

$$\theta_i \sim G \quad i = 1, \dots, N$$

.. then the collection  $(\theta_1, \dots, \theta_N)$  will have at most  $K$  unique parameters  $(\phi_1, \dots, \phi_K)$

- Thus  $G$  induces a **clustering** of parameters  $\theta_i$ 's, s.t.  $\theta_i$ 's within any group are all **identical**
- Therefore  $G$  can be used in a mixture model where  $\theta_i$ 's define the params of mixture distributions

# Constructing a Mixture Model for Data

- Let us use  $G$  to construct a mixture model for some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$





# Constructing a Mixture Model for Data

- Let us use  $G$  to construct a mixture model for some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$
- How: Generate each observation  $\mathbf{x}_i, i = 1, \dots, N$ , assuming the following generative model

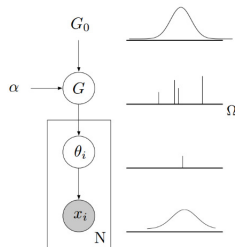
$$\pi \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$

$$\phi_k \sim G_0$$

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

$$\theta_i \sim G$$

$$x_i \sim p(x|\theta_i)$$



# Constructing a Mixture Model for Data

- Let us use  $G$  to construct a mixture model for some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$
- How: Generate each observation  $\mathbf{x}_i, i = 1, \dots, N$ , assuming the following generative model

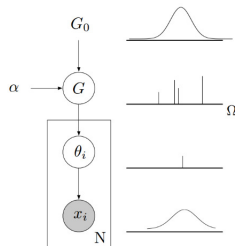
$$\pi \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$

$$\phi_k \sim G_0$$

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

$$\theta_i \sim G$$

$$x_i \sim p(x|\theta_i)$$



- Why drawing the parameters  $\{\theta_i\}_{i=1}^N$  from  $G$  clusters the observations  $\{\mathbf{x}_i\}_{i=1}^N$ ?



# Constructing a Mixture Model for Data

- Let us use  $G$  to construct a mixture model for some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$
- How: Generate each observation  $\mathbf{x}_i, i = 1, \dots, N$ , assuming the following generative model

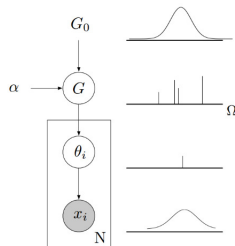
$$\pi \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$

$$\phi_k \sim G_0$$

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

$$\theta_i \sim G$$

$$x_i \sim p(x|\theta_i)$$



- Why drawing the parameters  $\{\theta_i\}_{i=1}^N$  from  $G$  clusters the observations  $\{\mathbf{x}_i\}_{i=1}^N$ ?
  - Reason: Since  $G$  is discrete,  $\theta_i$ 's generated by  $G$  won't be unique (only  $K$  unique values  $\{\phi_k\}_{k=1}^K$ )



# Constructing a Mixture Model for Data

- Let us use  $G$  to construct a mixture model for some observed data  $\mathbf{x}_1, \dots, \mathbf{x}_N$
- How: Generate each observation  $\mathbf{x}_i, i = 1, \dots, N$ , assuming the following generative model

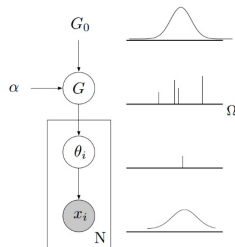
$$\pi \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$

$$\phi_k \sim G_0$$

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

$$\theta_i \sim G$$

$$x_i \sim p(x|\theta_i)$$



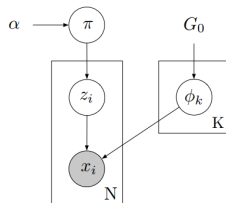
- Why drawing the parameters  $\{\theta_i\}_{i=1}^N$  from  $G$  clusters the observations  $\{\mathbf{x}_i\}_{i=1}^N$ ?
  - Reason: Since  $G$  is discrete,  $\theta_i$ 's generated by  $G$  won't be unique (only  $K$  unique values  $\{\phi_k\}_{k=1}^K$ )
  - Thus effectively the data is generated from not  $N$  separate distributions but only  $K < N$  unique distributions  $\{p(\mathbf{x}|\phi_k)\}_{k=1}^K$ , which naturally results in a clustering of data



# Two Equivalent Views

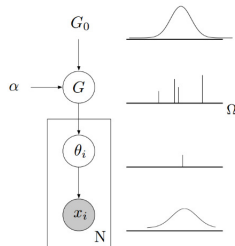
View-1 (the familiar one!): Clustering is explicitly described by the indicator  $z_i$

$$\begin{aligned}\pi &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \\ z_i &\sim \text{multinoulli}(\pi) \\ \mathbf{x}_i &\sim p(\mathbf{x}|\phi_k) \text{ if } z_i = k\end{aligned}$$



(Equivalent) View-2: Clustering is implicit (non-uniqueness of  $\theta_i$ 's denotes clustering of  $\mathbf{x}_i$ 's)

$$\begin{aligned}\pi &\sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \\ \phi_k &\sim G_0 \\ G &= \sum_{k=1}^K \pi_k \delta_{\phi_k} \\ \theta_i &\sim G \\ \mathbf{x}_i &\sim p(\mathbf{x}|\theta_i)\end{aligned}$$



# An Infinite Mixture Model

- Recall that our Bayesian mixture model construction for the finite  $K$  was

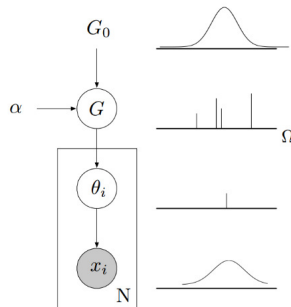
$$\pi \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$$

$$\phi_k \sim G_0$$

$$G = \sum_{k=1}^K \pi_k \delta_{\phi_k}$$

$$\theta_i \sim G$$

$$x_i \sim p(x|\theta_i)$$



- To get a mixture model with unbounded number of cluster, we need  $G$  of the form

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$

- How can we formally construct such a  $G$  that has potentially infinite mixture components?



# The Stick-Breaking Construction

- Sethuraman (1994) gave an “explicit” construction of  $G$  having infinite mixture components



# The Stick-Breaking Construction

- Sethuraman (1994) gave an “explicit” construction of  $G$  having infinite mixture components

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$





# The Stick-Breaking Construction

- Sethuraman (1994) gave an “explicit” construction of  $G$  having infinite mixture components

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$

- We basically need to generate  $\{\pi_k\}_{k=1}^{\infty}$  s.t.  $\pi_k \in (0, 1)$  and  $\sum_{k=1}^{\infty} \pi_k = 1$



# The Stick-Breaking Construction

- Sethuraman (1994) gave an “explicit” construction of  $G$  having infinite mixture components

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$

- We basically need to generate  $\{\pi_k\}_{k=1}^{\infty}$  s.t.  $\pi_k \in (0, 1)$  and  $\sum_{k=1}^{\infty} \pi_k = 1$
- Can be done using a stick-breaking construction for  $\{\pi_k\}_{k=1}^{\infty}$  as follows



# The Stick-Breaking Construction

- Sethuraman (1994) gave an “explicit” construction of  $G$  having infinite mixture components

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$$

- We basically need to generate  $\{\pi_k\}_{k=1}^{\infty}$  s.t.  $\pi_k \in (0, 1)$  and  $\sum_{k=1}^{\infty} \pi_k = 1$
- Can be done using a stick-breaking construction for  $\{\pi_k\}_{k=1}^{\infty}$  as follows

$$\beta_k \sim \text{Beta}(1, \alpha) \quad k = 1, \dots, \infty$$

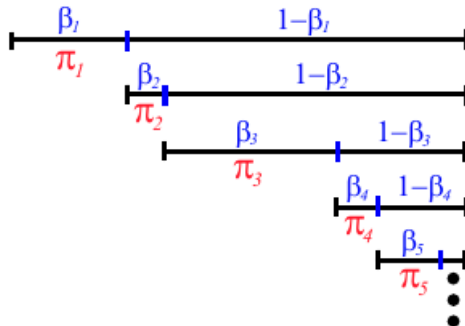
$$\pi_1 = \beta_1$$

$$\pi_k = \beta_k \prod_{\ell=1}^{k-1} (1 - \beta_{\ell}) \quad k = 2, \dots, \infty$$



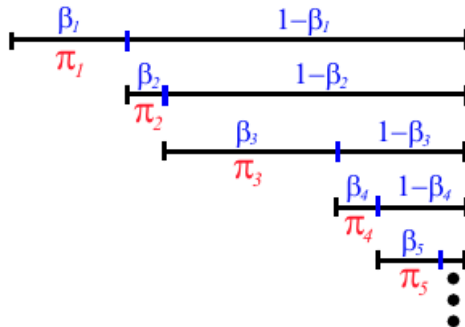
# The Stick-Breaking Construction

- Assume a stick of length 1 to begin with. Now recursively break it as follows:
  - Choose a random location  $\beta_k \in (0, 1)$  drawn from  $\text{Beta}(1, \alpha)$  at which to break the stick
  - Record  $\pi_k$  as " $\beta_k$  times the length of the remaining stick"



# The Stick-Breaking Construction

- Assume a stick of length 1 to begin with. Now recursively break it as follows:
  - Choose a random location  $\beta_k \in (0, 1)$  drawn from  $\text{Beta}(1, \alpha)$  at which to break the stick
  - Record  $\pi_k$  as " $\beta_k$  times the length of the remaining stick"



- Can show that  $\sum_{k=1}^K \pi_k = 1$  as  $K \rightarrow \infty$ . One easy way to verify this is by showing that  $1 - \sum_{k=1}^K \pi_k \rightarrow 0$  as  $K \rightarrow \infty$



# Next Class

- Some other (equivalent) ways of looking at nonparametric Bayesian mixture models
  - Dirichlet Process
  - Chinese Restaurant Process
  - Pólya-Urn Scheme
  - Hierarchical Dirichlet Process
- Some other examples of nonparametric Bayesian models

