

# Nonparametric Bayesian Modeling

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

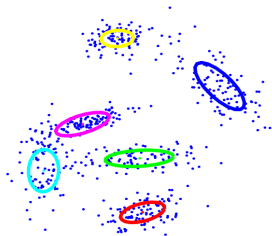
April 10, 2018

# Motivation: Mixture Models

- Suppose each observation is generated from a mixture model

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

$$\mathbf{x}_n \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}_n}, \boldsymbol{\Sigma}_{\mathbf{z}_n})$$



- How to learn  $K$ , i.e., the number of clusters for such a mixture model?
- Need to model/prior for  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$  that allows number of clusters to be unbounded

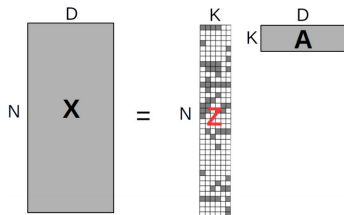
# Motivation: Latent Feature Models

- Suppose each observation is a subset sum of  $K$  latent features

$$z_{nk} \sim \text{Bernoulli}(\pi_k) \quad k = 1, \dots, K$$

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

- This can also be seen as a form of matrix factorization with one matrices being binary



- How to learn  $K$ , i.e., the number of latent features (i.e., number of columns/rows in  $\mathbf{Z}/\mathbf{A}$ )?
- Need to model/prior for  $\mathbf{Z}$  that allows number of columns in  $\mathbf{Z}$  to remain unbounded

# Motivation: Matrix Factorization

- Consider the following singular value decomposition model for an  $N \times M$  matrix  $\mathbf{X}$

$$\mathbf{X} = \sum_{k=1}^K \lambda_k \mathbf{u}_k \mathbf{v}_k^\top + \mathbf{E} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^\top + \mathbf{E}$$

where  $\mathbf{u}_k \in \mathbb{R}^N$  and  $\mathbf{v}_k \in \mathbb{R}^M$ ,  $\mathbf{\Lambda}$  is a  $K \times K$  diagonal matrix

- This is basically a weighted sum of rank-1 matrices
- How to learn  $K$ , i.e., the “rank” of the above factorization?
- Need a model/prior for  $\mathbf{\Lambda}$  that allows the size of  $\mathbf{\Lambda}$  to be unbounded

# Nonparametric Bayesian Modeling

- Enables constructing models that do not have an a priori fixed size
- Nonparametric does not mean no parameters
  - Instead, have a possibly infinite (unbounded) number of parameters
  - Note: We've already seen Gaussian Processes which is a nonparametric Bayesian model
- Usually constructed via one of the following ways
  - Take a finite model (e.g., a finite mixture model) and consider its “infinite limit”
  - Have a model that allows very large number of parameters but has a “shrinkage” effect, e.g.,

$$\mathbf{X} = \sum_{k=1}^K \lambda_k \mathbf{u}_k \mathbf{v}_k^\top + \mathbf{E} \quad \lambda_k \rightarrow 0 \quad \text{as} \quad k \rightarrow \infty$$

- We will look at some examples of both these approaches

# Being Nonparametric by taking Infinite Limit of Finite Models

# Finite Mixture Model

- Data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ , cluster assignments  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$ ,  $K$  clusters
- Denote the mixing proportion by a vector  $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$ ,  $\sum_{k=1}^K \pi_k = 1$

$$p(\boldsymbol{\pi}|\alpha) = \text{Dirichlet}\left(\frac{\alpha}{K}, \frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$$

$$p(\mathbf{z}_n|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_{nk}}$$

$$p(\mathbf{X}|\boldsymbol{\pi}) = \prod_{n=1}^N \sum_{k=1}^K \pi_k p(\mathbf{x}_n|\mathbf{z}_n = k)$$

- Integrating out  $\boldsymbol{\pi}$ , the marginal prior probability of cluster assignments  $\mathbf{Z}$

$$p(\mathbf{Z}|\alpha) = \int p(\mathbf{Z}|\boldsymbol{\pi})p(\boldsymbol{\pi}|\alpha)d\boldsymbol{\pi} = \frac{\Gamma(\alpha)}{\Gamma(N + \alpha)} \frac{\prod_{k=1}^K \Gamma(m_k + \frac{\alpha}{K})}{\Gamma(\frac{\alpha}{K})^K} \quad (\text{verify})$$

where  $m_k = \text{no. of points with } \mathbf{z}_n = k$

# Finite Mixture Model

- What is the prior  $p(\mathbf{z}_n | \mathbf{Z}_{-n}, \alpha)$ , i.e., the dist. of  $\mathbf{z}_n$  given cluster assignment  $\mathbf{Z}_{-n}$  of other points?

$$p(\mathbf{z}_n | \mathbf{Z}_{-n}, \alpha) = \frac{p(\mathbf{z}_n, \mathbf{Z}_{-n} | \alpha)}{p(\mathbf{Z}_{-n} | \alpha)} = \frac{p(\mathbf{Z} | \alpha)}{p(\mathbf{Z}_{-n} | \alpha)}$$

- Note that the above is a discrete distribution ( $\mathbf{z}_n$  takes one of  $K$  possible values)
- Using  $p(\mathbf{Z} | \alpha) = \frac{\Gamma(\alpha)}{\Gamma(N+\alpha)} \frac{\prod_{k=1}^K \Gamma(m_k + \frac{\alpha}{K})}{\Gamma(\frac{\alpha}{K})^K}$  (and applying the same to get  $p(\mathbf{Z}_{-n} | \alpha)$ ) we will have

$$\begin{aligned} p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) &= \frac{p(\mathbf{z}_n = j, \mathbf{Z}_{-n} | \alpha)}{p(\mathbf{Z}_{-n} | \alpha)} = \frac{\frac{\Gamma(\alpha)}{\Gamma(N+\alpha)} \frac{\Gamma(m_j + \frac{\alpha}{K}) \prod_{k \neq j} \Gamma(m_k + \frac{\alpha}{K})}{\Gamma(\frac{\alpha}{K})^K}}{\frac{\Gamma(\alpha)}{\Gamma(N-1+\alpha)} \frac{\Gamma(m_j - 1 + \frac{\alpha}{K}) \prod_{k \neq j} \Gamma(m_k + \frac{\alpha}{K})}{\Gamma(\frac{\alpha}{K})^K}} \\ &= \frac{m_{-n,j} + \frac{\alpha}{K}}{N - 1 + \alpha} \quad (m_{-n,j} = m_j - 1 \text{ denotes no. of other examples in cluster } j) \end{aligned}$$

- Thus **prior prob. of  $\mathbf{z}_n = j$  is prop. to  $m_{-n,j}$** , i.e., number of other examples assigned to cluster  $j$  (this is like a “rich gets richer” phenomenon; a popular cluster will attract more examples)



# Taking the Infinite Limit..

- We saw that  $p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) = \frac{m_{-n,j} + \frac{\alpha}{K}}{N-1+\alpha}$ . What if  $K \rightarrow \infty$  ?
- In that case, we will have  $p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) = \frac{m_{-n,j}}{N-1+\alpha}$
- Suppose only  $K_+$  clusters are currently occupied (that have at least one data point)
- Total probability of data point  $n$  going to any of these  $K_+$  clusters  $= \sum_{j=1}^{K_+} \frac{m_{-n,j}}{N-1+\alpha} = \frac{N-1}{N-1+\alpha}$
- Probability of data point  $n$  going to a **new (i.e., so far unoccupied) cluster**  $= \frac{\alpha}{N-1+\alpha}$
- Therefore in the limit of an unbounded number of clusters, we have

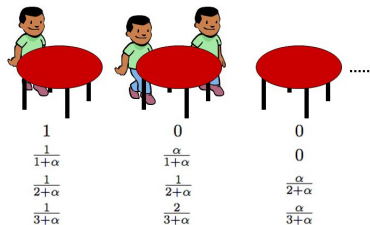
$$p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) = \begin{cases} \frac{m_{-n,j}}{N-1+\alpha} & \text{(prob. of going to } j = 1, \dots, K_+) \\ \frac{\alpha}{N-1+\alpha} & \text{(prob. of creating a new cluster } K_+ + 1) \end{cases}$$

- The above gives us a prior distribution for clustering with unbounded  $K$
- Note that the probability of starting a new cluster is proportional to Dirichlet hyperparam.  $\alpha$ 
  - The hyperparam  $\alpha$  can also be learned

# A metaphor for the same: Chinese Restaurant Process (CRP)

- Consider a restaurant with infinite number of tables (each table denotes a cluster)
- Customer 1 sits at a randomly chosen table (all tables are equivalent to begin with)
- Each subsequent customer  $n > 1$  sits using the following scheme
  - Sits at an already occupied table  $k$  with probability  $\frac{m_k}{n-1+\alpha}$
  - Sits at a new table with probability  $\frac{\alpha}{n-1+\alpha}$

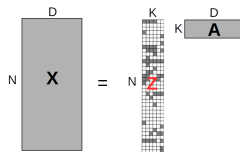
$$p(z_n = j | z_1, z_2, \dots, z_{n-1}, \alpha) = \begin{cases} \frac{m_{-n,j}}{n-1+\alpha} & (\text{for } j = 1, \dots, K_+) \\ \frac{\alpha}{n-1+\alpha} & (\text{create a new cluster } K_+ + 1) \end{cases}$$



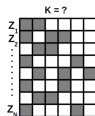
- The total number of occupied tables isn't fixed beforehand (connections to [Dirichlet Processes](#))

# Modeling Binary Matrices with Unbounded Number of Columns

- Recall the subset-sum problem



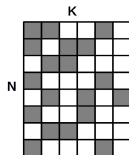
- So how do we model binary matrices for which number of columns  $K$  is *a priori* unknown?



- A nonparam. Bayesian model called “[Indian Buffet Process](#)” (IBP) defines a prior for such matrices
- Just like CRP, the IBP is a metaphor to describe the process that generates such matrices
- Note: In some models, the number of rows can be unknown (and number of columns fixed)
- The IBP model still applies (use IBP as the prior on the transpose of the matrix)

# Modeling Binary Matrices with Finite Many Columns

- Consider the generative process of an  $N \times K$  binary matrix  $Z$



- Rows denote the  $N$  examples, columns denote the  $K$  latent features
- Assume  $\pi_k \in (0, 1)$  to be probability of latent feature  $k$  being 1

$$z_{nk} \sim \text{Bernoulli}(\pi_k), \quad \pi_k \sim \text{Beta}(\alpha/K, 1)$$

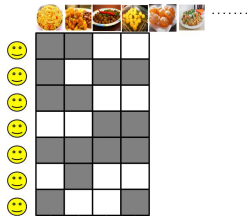
- Note: All  $z_{nk}$ 's are i.i.d. given  $\pi_k$
- For this model, the conditional probability of  $z_{nk} = 1$ , given other entries in column  $k$  of  $\mathbf{Z}$

$$p(z_{nk} = 1 | \mathbf{z}_{-n,k}) = \int p(z_{nk} = 1 | \pi_k) p(\pi_k | \mathbf{z}_{-n,k}) = \frac{m_{-n,k} + \frac{\alpha}{K}}{N + \frac{\alpha}{K}} \quad (\text{verify})$$

where  $m_{-n,k} = \sum_{i \neq n} z_{ik}$  denotes how many other entries in column  $k$  are equal to 1

# Another Metaphor: Indian Buffet Process

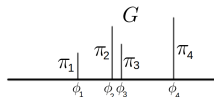
- For the finite  $K$  case, we saw that  $p(z_{nk} = 1 | \mathbf{z}_{-n,k}) = \frac{m_{-n,k} + \frac{\alpha}{K}}{N + \frac{\alpha}{K}}$
- As  $K \rightarrow \infty$ , we will have  $p(z_{nk} = 1 | \mathbf{z}_{-n,k}) = \frac{m_{-n,k}}{N}$  and  $p(z_{nk} = 0 | \mathbf{z}_{-n,k}) = \frac{N - m_{-n,k}}{N}$
- Note that this too exhibits a “rich-gets-richer” phenomenon (just like CRP)
- The Indian Buffet Process is a metaphor for this model. Assume a buffet with infinite dishes
  - Customer 1 selects  $\text{Poisson}(\alpha)$  dishes
  - The  $n$ -th customer selects:
    - Each already selected dish  $k$  with probability  $m_{-n,k}/n$   
( $m_k$  : how many previous customers before  $n$  selected dish  $k$ )
    - $\text{Poisson}(\alpha/n)$  new dishes (this can create new columns in  $\mathbf{Z}$ )
  - Note that as  $n$  grows, number of new dishes goes to zero (and the number of columns  $K$  converges to some finite number)
- The above can be used as a prior for  $\mathbf{Z}$ . Refer to (Griffiths and Ghahramani, 2011) for examples and other theoretical details of the model. Also has connections to [Beta Processes](#)



# Being Nonparametric using Models that have a Shrinkage Effect

# Nonparametric Mixture Models: Another Construction

- Consider a finite mixture model with  $K$  components with params  $(\mu_k, \Sigma_k)_{k=1}^K$



- Denoting  $(\mu_k, \Sigma_k) = \phi_k$ , this mixture model can be expressed as

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

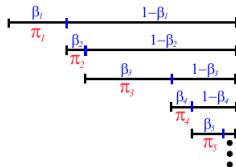
where  $\pi_k$  is the mixing proportion of component  $k$  and  $\sum_{k=1}^K \pi_k = 1$

- In a Bayesian setting, we usually assume  $[\pi_1, \dots, \pi_K] \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K)$
- We can make it a nonparametric model by constructing an infinite length probability vector

$$\pi_1, \pi_2, \pi_3, \dots, \quad \sum_{k=1}^{\infty} \pi_k = 1$$

- How to construct such an “infinite” length vector? Using an infinite dimensional Dirichlet?

# Nonparametric Mixture Models via Stick-Breaking Process



- Assume a mixture distribution  $G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$  and generate mixing proportions as

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

$$\pi_1 = \beta_1, \quad \pi_k = \beta_k \prod_{\ell=1}^{k-1} (1 - \beta_\ell) \quad k = 2, \dots, \infty$$

- Can show that  $\sum_{k=1}^{\infty} \pi_k = 1$
- Each “location”  $\phi_k$  can be drawn from an appropriate distribution  $\phi_k \sim G_0$
- Therefore, a finite mixture model can be made nonparametric by replacing the Dirichlet prior on  $\pi$  in a mixture model by a stick-breaking process prior as defined above



# Another Example: Multiplicative Gamma Process

- Consider the following probabilistic version of SVD

$$\mathbf{X} = \sum_{k=1}^K \lambda_k \mathbf{u}_k \mathbf{v}_k^\top + \mathbf{E}$$

- Consider the following prior on the “singular values”  $\lambda_k$

$$\lambda_k \sim \mathcal{N}(0, \tau_k^{-1})$$

$$\tau_k = \prod_{\ell=1}^k \delta_\ell$$

$$\delta_\ell \sim \text{Gamma}(\alpha, 1) \quad \text{where } \alpha > 1$$

- Note that as  $k$  becomes large,  $\tau_k$  gets larger and larger and  $\lambda_k$  shrinks to zero

# Some Comments

- Nonparametric Bayesian models have been widely used in several applications
  - Clustering, dim-red, regression/classification, time-series models such as HMM, and many others
- Nonparametric Bayesian models are not the only way to learn the right model size
- Marginal likelihood  $p(\mathcal{D}|\mathcal{M})$  can be used for model selection from a set of models  $\{\mathcal{M}_i\}_{i=1}^L$
- Other criteria such as Akaike or Bayesian Information Criteria are also commonly used
  - Usually defined as a sum of negative log-lik. and model size (models with smaller values preferred)

$$AIC = 2k - 2 \times \log\text{-lik}$$

$$BIC = k \log N - 2 \times \log\text{-lik}$$

where  $k$  denotes the number of parameters of the model,  $N$  denotes number of data points

- However, marginal likelihood, AIC/BIC, etc. **try multiple models** and then choose the best
- In contrast, NPBayes models learn a **single model** having an **unbounded complexity**
  - Also natural for streaming data where model selection is difficult/impractical to perform