

Introduction to Nonparametric Bayesian Modeling

CS698X: Topics in Probabilistic Modeling and Inference

Piyush Rai

Plan

- Need for nonparametric Bayesian modeling
- Some basic ideas
- Some examples of NPBayes modeling for
 - Mixture Models
 - Latent Feature Models and Matrix Factorization
- Some standard ways of constructing NPBayes models
 - Stick-breaking process, Dirichlet process
 - Some metaphors: Chinese Restaurant Process and Indian Buffet Process



Motivating Problem: Mixture Models

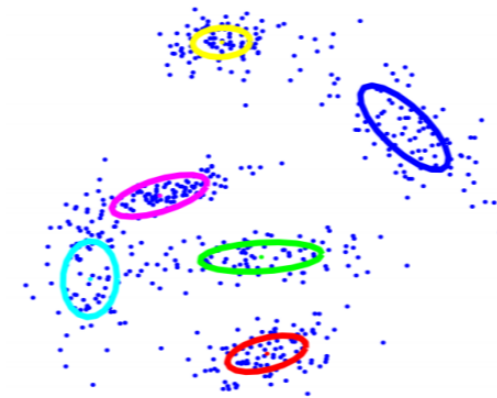
- Suppose each observation is generated from a K component mixture model

Cluster id of the n^{th} observation

K -dim probability vector of component mixing proportions

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

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



- How to learn K , i.e., the number of components (clusters) for such a mixture model?
- Can use marginal-likelihood based model selection but is expensive
 - Need to train the model several times for each possible value of K
- Also difficult if the data is streaming (hard to know beforehand how many clusters)
- How about a prior over $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N]$ (or $\boldsymbol{\pi}$) that allows learning the “right” K ?



Motivating Problem: Latent Feature Models

- Suppose each observation is a subset sum of K “basis vectors” (or “latent features”*)

Is k^{th} latent feature present in the n^{th} observation?

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

An example: Each text document in a collection being a subset sum of K “latent” themes or topics present in the collection

The n^{th} observation ($D \times 1$) expressed as a subset sum of the K latent features (each $D \times 1$), plus some noise

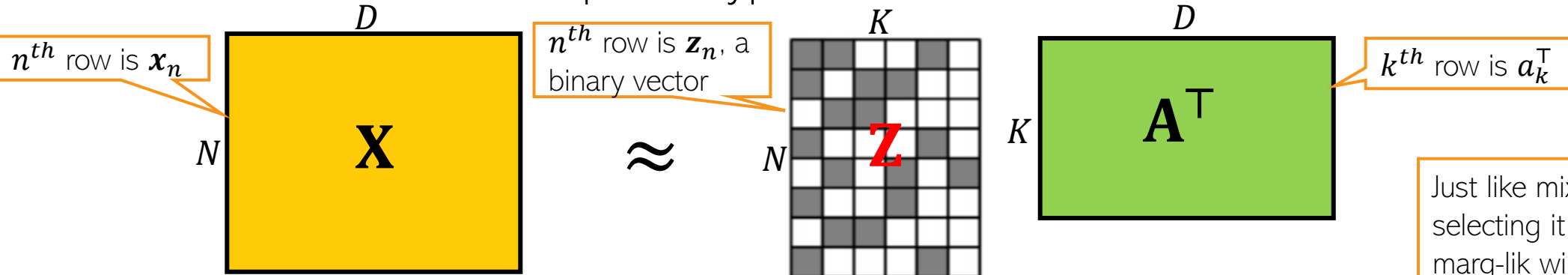
$$\mathbf{x}_n = \sum_{k=1}^K z_{nk} \mathbf{a}_k + \epsilon_n = \mathbf{A} \mathbf{z}_n + \epsilon_n$$

Noise (e.g., zero mean Gaussian)

The k^{th} latent feature ($D \times 1$)

A binary sparse matrix

- This can also be seen as special type of matrix factorization $\mathbf{X} = \mathbf{Z} \mathbf{A}^T + \mathbf{E}$



Just like mixture models, selecting it based on marg-lik will be expensive

- How about a prior over \mathbf{Z} (or \mathbf{A} or $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$) that allows learning the “right” K ?

Motivating Problem: SVD-style Matrix Factorization⁵

- Consider the following SVD-style decomposition for an $N \times M$ matrix \mathbf{X}

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

Rank 1 matrix

Zero mean Gaussian noise

- Each $\mathbf{u}_k \in \mathbb{R}^N$, $\mathbf{v}_k \in \mathbb{R}^M$, $\lambda_k \in \mathbb{R}$, and $\mathbf{\Lambda}$ is a $K \times K$ diag matrix with λ_k 's on diags
- This is basically a weighted sum of K rank-1 matrices
 - λ_k 's are the weights
 - λ_k 's are akin to the singular values in SVD
- How to learn K , i.e., the “rank” of the above factorization?
- How about a prior on $\mathbf{\Lambda}$, or \mathbf{U} or \mathbf{V} , that allows us to learn the “right” K ?



Nonparametric Bayesian Modeling

A vast area of research in ML and statistics. We will only be looking at a basic flavor of some approaches



- 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 params but has a “shrinkage” effect, e.g.,

And can potentially grow as we see more and more data (actual number will depend on the amount/properties of data)

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

- We will look at some examples of both these approaches



Being Nonparametric by Taking Infinite Limit of Finite Models



A Finite Mixture Model

- Data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, cluster assignments $\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N]$, K clusters
- Denote the mixing proportion by a vector $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$, $\sum_{k=1}^K \pi_k = 1$

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

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

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

a.k.a. “collapsing” a variable; one less variable to infer now

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

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

Number of points with $\mathbf{z}_n = k$



A Finite Mixture Model

- The prior distribution of \mathbf{z}_n given cluster assignment \mathbf{Z}_{-n} of other points?

A discrete distribution (multinoulli) since \mathbf{z}_n can take one of K possibilities

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

This “conditional” prior is needed when computing the posterior of \mathbf{z}_n since we have integrated out $\boldsymbol{\pi}$

- Using $p(\mathbf{Z} | \alpha) = \frac{\Gamma(\alpha)}{\Gamma(N + \alpha)} \prod_{k=1}^K \frac{\Gamma(m_k + \frac{\alpha}{K})}{\Gamma(\frac{\alpha}{K})^K}$ we have

Number of points in cluster j , not counting \mathbf{x}_n

$$p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) = \frac{p(\mathbf{z}_n = j, \mathbf{Z}_{-n} | \alpha)}{p(\mathbf{Z}_{-n} | \alpha)} = \frac{\frac{\Gamma(\alpha)}{\Gamma(N + \alpha)} \frac{\Gamma(m_j + \frac{\alpha}{K}) \prod_{k \neq j} \Gamma(m_k + \frac{\alpha}{K})}{\Gamma(\frac{\alpha}{K})^K}}{\frac{\Gamma(\alpha)}{\Gamma(N - 1 + \alpha)} \frac{\Gamma(m_j - 1 + \frac{\alpha}{K}) \prod_{k \neq j} \Gamma(m_k + \frac{\alpha}{K})}{\Gamma(\frac{\alpha}{K})^K}} = \frac{m_{-n,j} + \frac{\alpha}{K}}{N - 1 + \alpha}$$

- Note: Can also get this result using $p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) = \int p(\mathbf{z}_n = j | \boldsymbol{\pi}) p(\boldsymbol{\pi} | \mathbf{Z}_{-n}, \alpha) d\boldsymbol{\pi}$
- Thus prior prob. of $\mathbf{z}_n = j$ is proportional to how many other points are in cluster j
- Note that it also implies that mixture models have a **rich-gets-richer** property
 - Meaning: *a priori*, a cluster with more points is likely to attract more points



Taking the Infinite Limit..

- Since $p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) = \frac{m_{-n,j} + \frac{\alpha}{K}}{N-1+\alpha}$, as $K \rightarrow \infty$, $p(\mathbf{z}_n = j | \mathbf{Z}_{-n}, \alpha) = \frac{m_{-n,j}}{N-1+\alpha}$
- Suppose only K_+ clusters are currently occupied (i.e., have at least one data point)
- Total prob. of \mathbf{x}_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 \mathbf{x}_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 mixture models with unbounded K
 - Can combine it now with the suitable likelihood to infer the posterior* of \mathbf{Z}
- Note: Prob. of starting a new cluster is prop. to Dirichlet hyperparam α (can learn it)

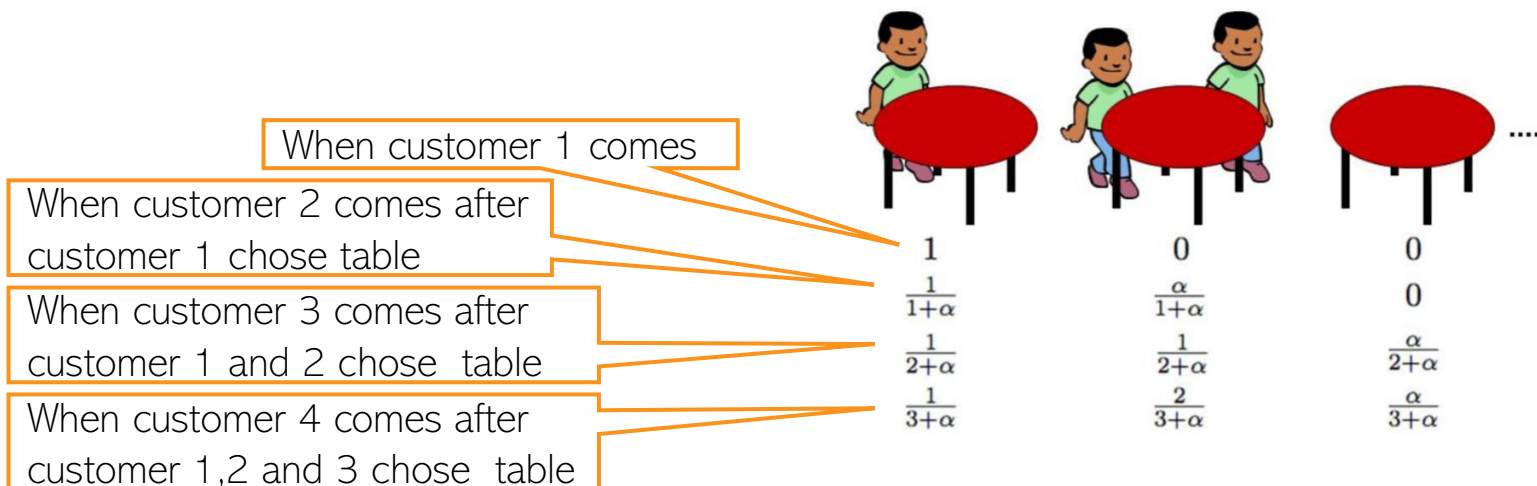
*Markov chain sampling methods for Dirichlet process mixture models, (Neal, 2000), Variational inference for Dirichlet process mixtures (Blei and Jordan, 2006)



A Metaphor: Chinese Restaurant Process (CRP)

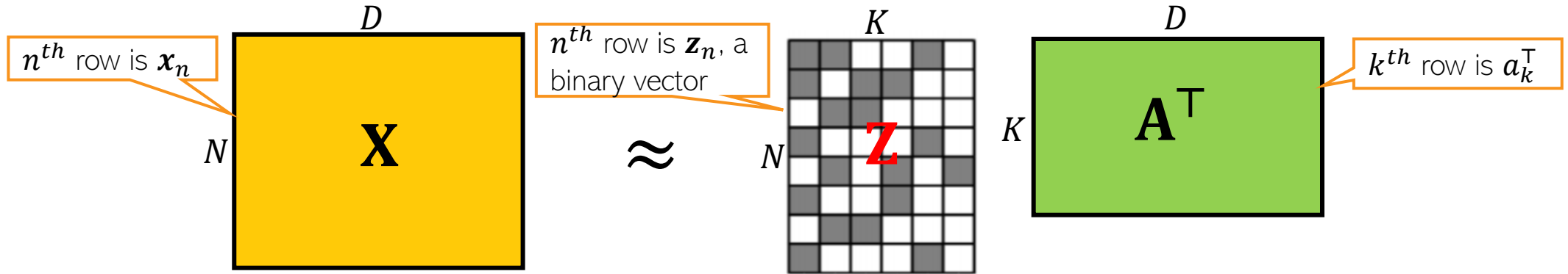
11

- Assume 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 $\frac{\alpha}{n-1+\alpha}$



Nonparametric Bayesian Latent Feature Model

- Recall the subset-sum problem: $\mathbf{x}_n = \mathbf{A}\mathbf{z}_n + \epsilon_n$



- To learn K , one option is to model \mathbf{Z} such that number of columns can be unbounded
- For the finite case of $N \times K$ matrix \mathbf{Z} , assume the following generative process

Prob. of an entry of column k being 1

Number of other entries in column k that are 1

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

Prob that this entry is 1 given values of other entries in column k

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

As $K \rightarrow \infty$,

$$p(z_{nk} = 1 | \mathbf{z}_{-n,k}) = \frac{m_{-n,k}}{N}$$

Proportional to how many other entries in column k are 1

Rich-gets-richer just like mixture models



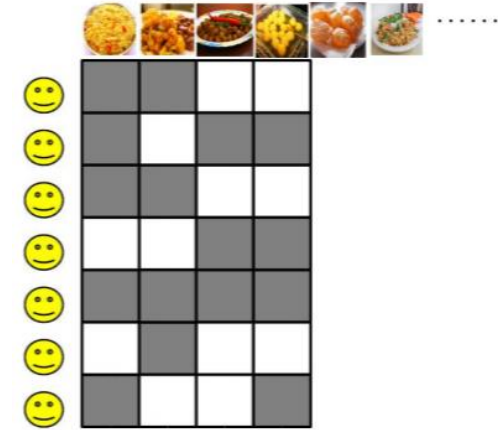
A Metaphor: Indian Buffet Process (IBP)

- IBP is a metaphor for latent feature model similar to CRP for mixture model

- Assume a buffet with infinite dishes

- Customer 1 selects **Poisson(α)** dishes
- Customer n makes selection as follows

- Select each already selected dish k with prob. $\frac{m_{-n,k}}{N}$
- Select **Poisson(α/n)** new dishes



Since new dishes are added as well

- Customer-dish assignment matrix is a binary matrix
 - Thus this “process” defines a prior over binary matrices without a pre-defined number of columns
- Note that as n grows, number of new dishes goes to zero (and the number of columns K converges to some finite number)
- Refer to (Griffiths and Ghahramani, 2011) for examples and other theoretical details. Also has connections to [Beta Process](#) (just like CRP has with Dirichlet Process)