



Probabilistic Graphical Models

01010001 Ω

Bayesian nonparametrics: Dirichlet Process

Eric Xing Lecture 23, April 13, 2020

Reading: see class homepage







How to pick number of cluster?









How to pick number of cluster?



Points on a 2D plane without color as attribute







T=1 Streaming Data T=2





- How to model this data?
- Mixture of Gaussians:

 $p(x_1, \dots, x_N | \pi, \{\mu_k\}, \{\Sigma_k\}) = \prod_{n=1}^{\infty} \sum_{k=1}^{K} \pi_k \mathcal{N}(x_k | \mu_k, \Sigma_k) \qquad \bigstar$

Parametric model: Fixed finite number of parameters.

×

×

×

×



×



×



- How to choose the mixing weights and mixture parameters?
- Bayesian choice: Put a prior on them and integrate out:

$$p(x_1, \dots, x_N)$$

= $\int \int \int \left(\prod_{n=1}^{\infty} \sum_{k=1}^{K} \pi_k \mathcal{N}(x_k | \mu_k, \Sigma_k) \right)$
 $p(\pi) p(\mu_{1:K}) p(\Sigma_{1:K}) d\pi d\mu_{1:K} d\Sigma_{1:K}$

- Where possible, use conjugate priors
 - Gaussian/inverse Wishart for mixture parameters
 - What to choose for mixture weights?





- The Dirichlet distribution is a distribution over the (K-1)-dimensional simplex.
- It is parametrized by a *K*-dimensional vector $(\alpha_1, \ldots, \alpha_K)$ such that $\alpha_k \ge 0, k = 1, \ldots, K$ and $\sum_k \alpha_k > 0$
- Its distribution is given by

$$\frac{\prod_{k=1}^{K} \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^{K} \alpha_k)} \prod_{k=1}^{K} \pi_k^{\alpha_k - 1}$$



Samples from the Dirichlet distribution

• If $\pi \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_K)$ then $\pi_k \ge 0$ for all k, and $\sum_{k=1}^K \pi_k = 1$.

• Expectation:
$$\mathbb{E}\left[(\pi_1, \dots, \pi_K)\right] = \frac{(\alpha_1, \dots, \alpha_K)}{\sum_k \alpha_k}$$





Conjugacy to the multinomial

• If $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$ and $x_n \stackrel{iid}{\sim} \theta$

$$p(\pi|x_1, \dots, x_n) \propto p(x_1, \dots, x_n|\pi)p(\pi)$$

$$= \left(\frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1}\right) \left(\frac{n!}{m_1! \dots m_K!} \pi_1^{m_1} \dots \pi_K^{m_K}\right)$$

$$\propto \frac{\prod_{k=1}^K \Gamma(\alpha_k + m_k)}{\Gamma(\sum_{k=1}^K \alpha_k + m_k)} \prod_{k=1}^K \pi_k^{\alpha_k + m_k - 1}$$

$$= \text{Dirichlet}(\pi|\alpha_1 + m_1, \dots, \alpha_K + m_K)$$



Distributions over distributions

- The Dirichlet distribution is a distribution over positive vectors that sum to one.
- We can associate each entry with a set of parameters
 - e.g. finite mixture model: each entry associated with a mean and covariance.
- □ In a Bayesian setting, we want these parameters to be *random*.
- We can combine the distribution over probability vectors with a distribution over parameters to get a distribution over distributions over parameters.



Example: finite mixture model

- Gaussian distribution: distribution over means.
 - Sample from a Gaussian is a real-valued number/vector.





Example: finite mixture model

- Gaussian distribution: distribution over means.
 - Sample from a Gaussian is a real-valued number/vector.
- Dirichlet distribution:
 - Sample from a Dirichlet distribution is a probability vector.





Example: finite mixture model

- Dirichlet prior
 - Each element of a Dirichlet-distributed vector is associated with a parameter value drawn from some distribution.
 - Sample from a Dirichlet prior is a probability distribution over parameters.





Properties of the Dirichlet distribution (collapsing)

Relationship to gamma distribution: If $\eta_k \sim \text{Gamma}(\alpha_k, 1)$,

$$\frac{(\eta_1,\ldots,\eta_K)}{\sum_k \eta_k} \sim \text{Dirichlet}(\alpha_1,\ldots,\alpha_K)$$

 \Box If $\eta_1 \sim \text{Gamma}(\alpha_1, 1)$ and $\eta_2 \sim \text{Gamma}(\alpha_2, 1)$ then

 $\eta_1 + \eta_2 \sim \text{Gamma}(\alpha_1 + \alpha_2, 1)$

• Therefore, if $(\pi_1 \ldots, \pi_K) \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_K)$ then

 $(\pi_1 + \pi_2, \pi_3, \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1 + \alpha_2, \alpha_3, \dots, \alpha_K)$



Properties of the Dirichlet distribution (splitting)

- □ The beta distribution is a Dirichlet distribution on the 1-simplex.
- Let $(\pi_1 \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$ and $\theta \sim \text{Beta}(\alpha_1 b, \alpha_1(1-b)), 0 < b < 1.$

• Then
$$(\pi_1\theta, \pi_1(1-\theta), \pi_2, \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1b_1, \alpha_1(1-b_1), \alpha_2, \dots, \alpha_K)$$

• More generally, if $\theta \sim \text{Dirichlet}(\alpha_1 b_1, \alpha_1 b_2, \dots, \alpha_1 b_N), \sum_i b_i = 1.$ then

 $(\pi_1\theta_1,\ldots,\pi_1\theta_N,\pi_2,\ldots,\pi_K) \sim \text{Dirichlet}(\alpha_1b_1,\ldots,\alpha_1b_N,\alpha_2,\ldots,\alpha_K)$



Properties of the Dirichlet distribution

• Renormalization:

If
$$(\pi_1 \dots, \pi_K) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$$

then $\frac{(\pi_2, \dots, \pi_K)}{\sum_{k=1}^K \pi_k} \sim ?$

$$\frac{(\pi_2, \dots, \pi_K)}{\sum_{k=1}^K \pi_k} \sim \text{Dirichlet}(\alpha_2, \dots, \alpha_K)$$



Choosing the number of clusters

Mixture of Gaussians – but how many components? What if we see more data – may find new components?



θ

 π

 X_i

G/IW





Parametric model:

- Assumes all data can be represented using a fixed, finite number of parameters.
 - □ Mixture of *K* Gaussians, polynomial regression.
- Nonparametric model:
- Number of parameters can grow with sample size.
- Number of parameters may be random.
 - Kernel density estimation.

Bayesian nonparametrics:

- Allow an *infinite* number of parameters *a priori*.
- A finite data set will only use a finite number of parameters.
- Other parameters are integrated out.



Bayesian nonparametric mixture models

- Make sure we always have more clusters than we need.
- Solution infinite clusters *a priori*!

$$p(x_n|\pi, \{\mu_k\}, \{\Sigma_k\}) = \sum_{k=1}^{\infty} \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)$$

- A finite data set will always use a finite but *random* number of clusters.
- How to choose the prior?
- We want something *like* a Dirichlet prior but with an infinite number of components.





Constructing an appropriate prior

- Start off with $\pi^{(2)} = (\pi_1^{(2)}, \pi_2^{(2)}) \sim \text{Dirichlet}\left(\frac{\alpha}{2}, \frac{\alpha}{2}\right)$
- Split each component according to the splitting rule:

$$\begin{split} \theta_1^{(2)}, \theta_2^{(2)} &\stackrel{iid}{\sim} \operatorname{Beta}\left(\frac{\alpha}{2} \cdot \frac{1}{2}, \frac{\alpha}{2} \cdot \frac{1}{2}\right) \\ \pi^{(4)} &= (\theta_1^{(2)} \pi_1^{(2)}, (1 - \theta_1^{(2)}) \pi_1^{(2)}, \theta_2^{(2)} \pi_2^{(2)}, (1 - \theta_2^{(2)}) \pi_2^{(2)} \\ &\sim \operatorname{Dirichlet}\left(\frac{\alpha}{4}, \frac{\alpha}{4}, \frac{\alpha}{4}, \frac{\alpha}{4}\right) \\ \end{split}$$
Repeat to get
$$\pi^{(K)} \sim \operatorname{Dirichlet}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$$

 \square As $K \to \infty$, we get a vector with infinitely many components



• Let *H* be a distribution on some space Ω – e.g. a Gaussian distribution on the real line.

• Let
$$\pi \sim \lim_{K \to \infty} \operatorname{Dirichlet}\left(\frac{\alpha}{K} \dots, \frac{\alpha}{K}\right)$$

• For
$$k = 1, \ldots, \infty$$
 let $\theta_k \sim H$.

• Then $G := \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k}$ is an infinite distribution over H.

• We write $G \sim \mathrm{DP}(\alpha, H)$



Samples from the Dirichlet process

- Samples from the Dirichlet process are *discrete*.
- We call the point masses in the resulting distribution, *atoms*.



• The base measure H determines the locations of the atoms.



Samples from the Dirichlet process

- The concentration parameter α determines the distribution over atom sizes.
- **□** Small values of α give *sparse* distributions.





Properties of the Dirichlet process

• For any partition A_1, \ldots, A_K of Ω (imagine different colors from a possibly infinite color library), the total mass assigned to each partition is distributed according to







• A Dirichlet process is the unique distribution over probability distributions on some space Ω , such that for any finite partition A_1, \dots, A_K of Ω ,

 $(P(A_1),\ldots,P(A_K)) \sim \text{Dirichlet}(\alpha H(A_1),\ldots,\alpha H(A_K)).$





Conjugacy of the Dirichlet process

• Let A_1, \ldots, A_K be a partition of Ω , and let H be a measure on Ω . Let $P(A_k)$ be the mass assigned by $G \sim DP(\alpha, H)$ to partition A_k . Then $(P(A_1), \ldots, P(A_K)) \sim Dirichlet(\alpha H(A_1), \ldots, \alpha H(A_K)).$

• If we see an observation in the \mathcal{J}^{h} segment, then $(P(A_1), \dots, P(A_j), \dots, P(A_K) | X_1 \in A_j)$ $\sim \text{Dirichlet}(\alpha H(A_1), \dots, \alpha H(A_j) + 1, \dots, \alpha H(A_K)).$

- This must be true for *all possible partitions of* Ω .
- This is only possible if the posterior of G, given an observation x, is given by $(\alpha H + \delta_m)$

$$G|X_1 = x \sim \mathrm{DP}\left(\alpha + 1, \frac{\alpha H + \delta_x}{\alpha + 1}\right)$$



- The Dirichlet process clusters observations.
- □ A new data point can either join an existing cluster, or start a new cluster.
- Question: What is the predictive distribution for a new data point?
- Assume H is a continuous distribution on Ω. This means for every point θ in Ω, $P_H(\theta) = 0$.
- First data point:
 - Start a new cluster.
 - Sample a parameter θ_1 for that cluster.



- We have now split our parameter space in two: the singleton θ_1 , and everything else.
- Let π_1 be the size of atom at θ_1 .
- The combined mass of all the other atoms is $\pi_* = 1 \pi_1$.
- $\Box \quad A \text{ priori}, \qquad (\pi_1, \pi_*) \sim \text{Dirichlet}(0, \alpha)$
- A posteriori, $(\pi_1, \pi_*)|X_1 = \theta_1 \sim \text{Dirichlet}(1, \alpha)$



• If we integrate out π_1 we get

$$P(X_2 = \theta_k | X_1 = \theta_1) = \int P(X_2 = \theta_k | (\pi_1, \pi_*)) P((\pi_1, \pi_* | X_1 = \theta_1) d\pi_1$$
$$= \int \pi_k \text{Dirichlet}((\pi_1, 1 - \pi_1) | 1, \alpha) d\pi_1$$
$$= \mathbb{E}_{\text{Dirichlet}(1, \alpha)} [\pi_k]$$
$$= \begin{cases} \frac{1}{1+\alpha} & \text{if } k = 1\\ \frac{\alpha}{1+\alpha} & \text{for new } k. \end{cases}$$



- Lets say we choose to start a new cluster, and sample a new parameter $\theta_2 \sim H$. Let π_2 be the size of the atom at θ_2 .
- A posteriori, $(\pi_1, \pi_2, \pi_*)|X_1 = \theta_1, X_2 = \theta_2 \sim \text{Dirichlet}(1, \alpha).$
- If we integrate out $\pi = (\pi_1, \pi_2, \pi_*)$, we get

$$P(X_3 = \theta_k | X_1 = \theta_1, X_2 = \theta_2)$$

= $\int P(X_3 = \theta_k | \pi) P(\pi | X_1 = \theta_1, X_2 = \theta_2) d\pi$
= $\mathbb{E}_{\text{Dirichlet}(1,1,\alpha)} [\pi_k]$
= $\begin{cases} \frac{1}{2+\alpha} & \text{if } k = 1 \\ \frac{1}{2+\alpha} & \text{if } k = 2 \\ \frac{\alpha}{2+\alpha} & \text{for new } k. \end{cases}$





□ In general, if m_k is the number of times we have seen $X_j = k$, and K is the total number of observed values,

$$P(X_{n+1} = \theta_k | X_1, \dots, X_n) = \int P(X_{n+1} = \theta_k | \pi) P(\pi | X_1, \dots, X_n) d\pi$$
$$= \mathbb{E}_{\text{Dirichlet}(m_1, \dots, m_K, \alpha)} [\pi_k]$$
$$= \begin{cases} \frac{m_k}{n+\alpha} & \text{if } k \leq K \\ \frac{\alpha}{n+\alpha} & \text{for new cluster.} \end{cases}$$

- We tend to see observations that we have seen before rich-get-richer property.
- We can always add new features *nonparametric*.







- Self-reinforcing property
- exchangeable partition of samples



- The resulting distribution over data points can be thought of using the following urn scheme.
- An urn initially contains a black ball of mass α.
- For n=1,2,... sample a ball from the urn with probability proportional to its mass.
- If the ball is black, choose a previously unseen color, record that color, and return the black ball plus a unit-mass ball of the new color to the urn.
- If the ball is not black, record it's color and return it, plus another unitmass ball of the same color, to the urn

[Blackwell and MacQueen, 1973]





- The distribution over partitions can be described in terms of the following restaurant metaphor:
- The first customer enters a restaurant, and picks a table.





- The distribution over partitions can be described in terms of the following restaurant metaphor:
- □ The first customer enters a restaurant, and picks a table.
- The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.





- The distribution over partitions can be described in terms of the following restaurant metaphor:
- □ The first customer enters a restaurant, and picks a table.
- The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.





- The distribution over partitions can be described in terms of the following restaurant metaphor:
- □ The first customer enters a restaurant, and picks a table.
- The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.





- The distribution over partitions can be described in terms of the following restaurant metaphor:
- □ The first customer enters a restaurant, and picks a table.
- The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.







- An interesting fact: the distribution over the clustering of the first N customers does not depend on the order in which they arrived.
- Homework: Prove to yourself that this is true.
- However, the customers are not independent they tend to sit at popular tables.
- We say that distributions like this are *exchangeable*.
- De Finetti's theorem: If a sequence of observations is exchangeable, there must exist a distribution given which they are iid.
- The customers in the CRP are iid given the underlying Dirichlet process by integrating out the DP, they become dependent.



The Stick-breaking Process



X

Stick breaking construction of DP

- We can represent samples from the Dirichlet process exactly.
- □ Imagine a stick of length 1, representing total probability.

■ For k=1,2,...

- Sample a beta(1, α) random variable b_k .
- Break off a fraction b_k of the stick. This is the k^{th} atom size
- Sample a random location for this atom.
- Recurse on the remaining stick.

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

$$\pi_k := b_k \prod_{j=1}^{k-1} (1 - b_k)$$

$$b_k \sim \text{Beta}(1, \alpha)$$
 [Sethuraman, 1994]

















$$G := \sum_{k=1}^{\infty} \pi_k \delta_{\theta_k} \sim \mathrm{DP}(\alpha, H)$$
$$\phi_n \sim G$$
$$x_n \sim f(\phi_n)$$







- We can integrate out *G* to get the CRP.
- Reminder: Observations in the CRP are exchangeable.
- Corollary: When sampling any data point, we can always rearrange the ordering so that it is the last data point.
- Let z_n be the cluster allocation of the *n*th data point.
- □ Let *K* be the total number of instantiated clusters.
- Then

$$p(z_n = k | x_n, z_{-n}, \phi_{1:K}) \propto \begin{cases} m_k f(x_n | \phi_k) & k \le K \\ \alpha \int_{\Omega} f(x_n | \phi) H(d\phi) & k = K+1 \end{cases}$$

 If we use a conjugate prior for the likelihood, we can often integrate out the cluster parameters



Problems with the collapsed sampler

- We are only updating one data point at a time.
- Imagine two "true" clusters are merged into a single cluster a single data point is unlikely to "break away".
- Getting to the true distribution involves going through low probability states → mixing can be slow.
- If the likelihood is not conjugate, integrating out parameter values for new features can be difficult.
- Neal [2000] offers a variety of algorithms.
- Alternative: Instantiate the latent measure.





- Topic models describe documents using a distribution over features.
- Each feature is a distribution over words
- Each document is represented as a collection of words (usually unordered – "bag of words" assumption).
- The words within a document are distributed according to a documentspecific mixture model
 - Each word in a document is associated with a feature.
- The features are shared between documents.
- The features learned tend to give high probability to semantically related words – "topics"



Latent Dirichlet allocation

- For each topic $k=1,\ldots,K$
 - Sample a distribution over words, $\beta \sim \text{Dir}(\eta_1, \dots, \eta_V)$
- For each document $m=1,\ldots,M$
 - Sample a distribution over topics, $\theta_m \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$
 - For each word $n=1,\ldots,N_m$
 - Sample a topic z_{mn} ~Discrete(θ_m)
 - Sample a word $w_{mk} \sim Discrete(\beta_z)$







Constructing a topic model with infinitely many topics

- □ LDA: Each distribution is associated with a distribution over *K* topics.
- Problem: How to choose the number of topics?
- Solution:
 - Infinitely many topics!
 - Replace the Dirichlet distribution over topics with a Dirichlet process!
- Problem: We want to make sure the topics are *shared* between documents



- □ In LDA, we have *M* independent samples from a Dirichlet distribution.
- □ The weights are different, but the topics are fixed to be the same.
- If we replace the Dirichlet distributions with Dirichlet processes, each atom of each Dirichlet process will pick a topic *independently* of the other topics.





- Because the base measure is *continuous*, we have zero probability of picking the same topic twice.
- If we want to pick the same topic twice, we need to use a *discrete* base measure.
- For example, if we chose the base measure to be $H = \sum_{k=1}^{\infty} \alpha_k \delta_{\beta_k}$, then we would have LDA again.
- We want there to be an infinite number of topics, so we want an *infinite, discrete* base measure.
- We want the location of the topics to be random, so we want an *infinite, discrete, random* base measure.



Hierarchical Dirichlet Process (Teh et al, 2006)

Solution: Sample the base measure from a Dirichlet process!

 $G_0 \sim \mathrm{DP}(\gamma, H)$ $G_m \sim \mathrm{DP}(\alpha, G_0)$





- Imagine a *franchise* of restaurants, serving an infinitely large, global menu.
- Each table in each restaurant orders a single dish.
- Let n_{rt} be the number of customers in restaurant *r* sitting at table *t*.
- Let m_{rd} be the number of tables in restaurant *r* serving dish *d*.
- Let m.d be the number of tables, across *all* restaurants, serving dish *d*.



- Customers enter the restaurants, and sit at tables according to the Chinese restaurant process
 - The first customer enters a restaurant, and picks a table.





- Customers enter the restaurants, and sit at tables according to the Chinese restaurant process
 - The first customer enters a restaurant, and picks a table.
 - The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.





- Customers enter the restaurants, and sit at tables according to the Chinese restaurant process
 - The first customer enters a restaurant, and picks a table.
 - The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.





- Customers enter the restaurants, and sit at tables according to the Chinese restaurant process
 - The first customer enters a restaurant, and picks a table.
 - The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.





- Customers enter the restaurants, and sit at tables according to the Chinese restaurant process
 - The first customer enters a restaurant, and picks a table.
 - The n^{th} customer enters the restaurant. He sits at an existing table with probability $m_k/(n-1+\alpha)$, where m_k is the number of people sat at table k. He starts a new table with probability $\alpha/(n-1+\alpha)$.





Each table in each restaurant picks a dish, with probability proportional to the number of times it has been served across *all* restaurants.





Each table in each restaurant picks a dish, with probability proportional to the number of times it has been served across *all* restaurants.





Each table in each restaurant picks a dish, with probability proportional to the number of times it has been served across *all* restaurants.





Each table in each restaurant picks a dish, with probability proportional to the number of times it has been served across *all* restaurants.





Recall: Graphical Model Representations of DP















Stick
$$(\alpha, \beta)$$
:
 $\pi'_{jk} \sim \operatorname{Beta}\left(\alpha\beta_{k}, \alpha\left(1 - \sum_{l=1}^{k}\beta_{l}\right)\right), \quad \pi_{jk} = \pi'_{jk}\prod_{l=1}^{k-1}\left(1 - \pi'_{jl}\right).$







Restaurants = documents; dishes = topics.

- Let H be a V-dimensional Dirichlet distribution, so a sample from H is a distribution over a vocabulary of V words.
- Sample a global distribution over topics,

$$G_0 := \sum_{k=1} \pi_k \delta_{\beta_k} \sim \mathrm{DP}(\alpha, H)$$

- For each document $m=1,\ldots,M$
 - Sample a distribution over topics, $G_m \sim DP(\gamma, G_0)$.
 - For each word $n=1,\ldots,N_m$
 - Sample a topic ϕ_{mn} ~Discrete (G₀).
 - Sample a word w_{mk} ~Discrete(ϕ_{mn}).







