

CS598JHM: Advanced NLP (Spring 2013)

<http://courses.engr.illinois.edu/cs598jhm/>

# Lecture 13:

# Dirichlet Processes

Julia Hockenmaier

[juliahmr@illinois.edu](mailto:juliahmr@illinois.edu)

3324 Siebel Center

Office hours: by appointment

# Finite mixture model

## Mixing proportions:

The prior probability of each component (assuming uniform  $\alpha$  )

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

## Mixture components:

The distribution over observations for each component

$$\theta_k^*|H \sim H \text{ (} H \text{ is typically a Dirichlet distribution)}$$

## Indicator variables:

Which component is observation  $i$  drawn from?

$$z_i|\pi \sim \text{Multinomial}(\pi)$$

## The observations:

The probability of observation  $i$  under component  $z_i$

$$x_i|z_i, \{ \theta_k^* \} \sim F(\theta_{z_i}^*) \text{ (} F \text{ is typically a categorical distribution)}$$

# Dirichlet Process $DP(\alpha, H)$

The **Dirichlet process**  $DP(\alpha, H)$  defines a distribution over distributions over a probability space  $\Theta$ .

Draws  $G \sim DP(\alpha, H)$  from this DP  
are **random distributions** over  $\Theta$

$DP(\alpha, H)$  has two parameters:

**Base distribution**  $H$ :

a distribution over the probability space  $\Theta$

**Concentration parameter**  $\alpha$ :

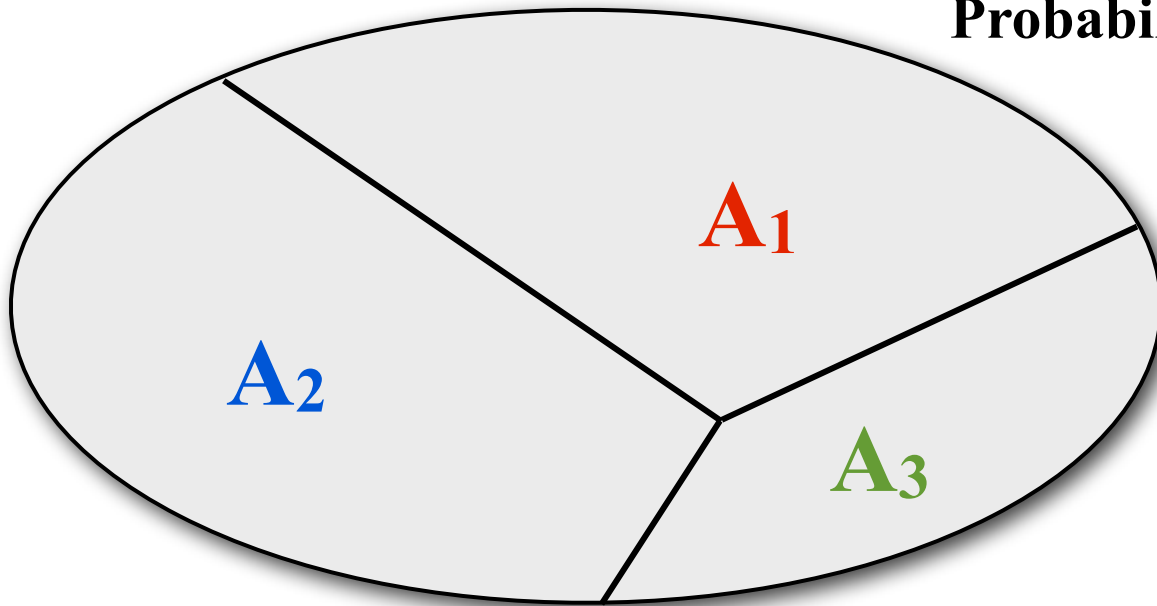
a positive real number

If  $G \sim DP(\alpha, H)$ , then for any finite measurable partition  $A_1 \dots A_r$  of  $\Theta$ :

$$(G(A_1), \dots, G(A_r)) \sim \text{Dirichlet}(\alpha H(A_1), \dots, \alpha H(A_r))$$

# The base distribution $H$

Probability space  $\Theta$



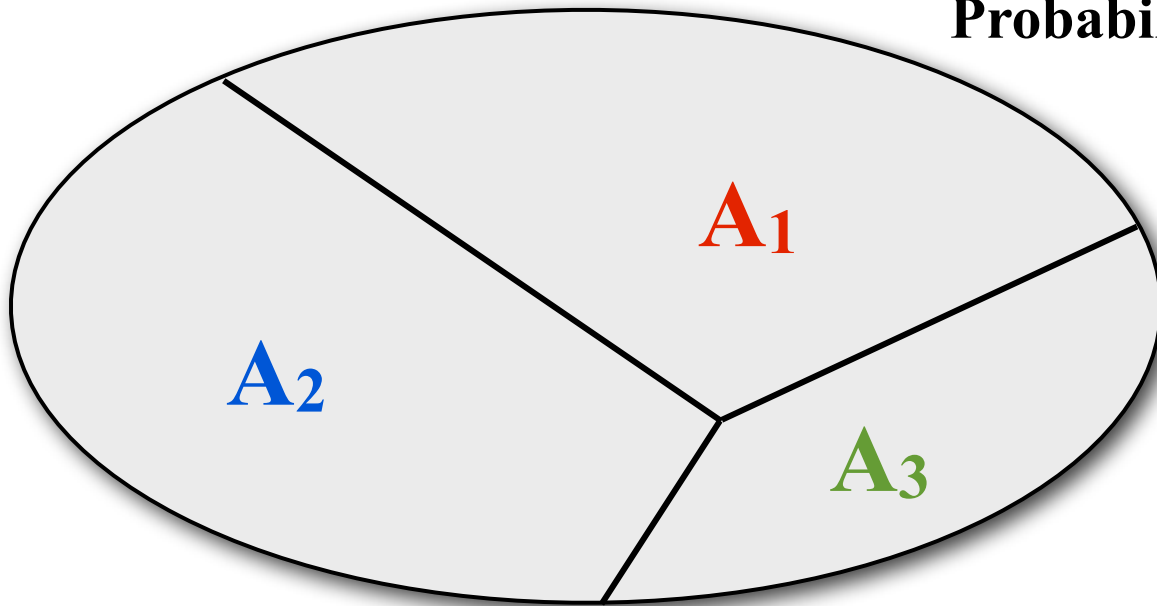
Since  $A_1, A_2, A_3$  partition  $\Theta$ , we can use the base distribution  $H$  to define a categorical distribution over  $A_1, A_2, A_3$ :

$$H(A_1) + H(A_2) + H(A_3) = 1$$

Note that we can use  $H$  to define a categorical distribution over *any* finite partition  $A_1 \dots A_r$  of  $\Theta$ , even if  $H$  is smooth

# Draws from the DP: $G \sim \text{DP}(\alpha, H)$

Probability space  $\Theta$



Every individual draw  $G$  from  $\text{DP}(\alpha, H)$  is also a distribution over  $\Theta$   
 $G$  also defines a categorical distribution over any partition of  $\Theta$

For *any* finite partition  $A_1 \dots A_r$  of  $\Theta$ , this categorical distribution is drawn from a Dirichlet prior defined by  $\alpha$  and  $H$ :

$$(G(A_1), G(A_2), G(A_3)) \sim \text{Dir}(\alpha H(A_1), \alpha H(A_2), \alpha H(A_3))$$

# The role of $H$ and $\alpha$

The base distribution  $H$  defines the **mean** (expectation) of  $G$ :

For any measurable set  $A \subseteq \Theta$ ,  $E[G(A)] = H(A)$

The concentration parameter  $\alpha$  is **inversely** related to the **variance** of  $G$  :

$$V[G(A)] = H(A)(1 - H(A))/(\alpha + 1)$$

$\alpha$  specifies how much mass is around the mean

The larger  $\alpha$ , the smaller the variance

$\alpha$  is also called the **strength parameter**: If we use  $DP(\alpha, H)$  as a prior,  $\alpha$  tells us how much we can deviate from the prior:

$$\text{As } \alpha \rightarrow \infty, G(A) \rightarrow H(A)$$

# The posterior of $G$ : $G|\theta_1, \dots, \theta_n$

Assume the distribution  $G$  is drawn from a DP:  $G \sim \text{DP}(\alpha, H)$

The **prior** of  $G$ :

$$(G(A_1), \dots, G(A_K)) \sim \text{Dirichlet}(\alpha H(A_1), \dots, \alpha H(A_K))$$

Given a sequence of observations  $\theta_1 \dots \theta_n$  from  $\Theta$

that are drawn from this  $G$ :  $\theta_i | G \sim G$

What is the **posterior of  $G$**  given the observed  $\theta_1 \dots \theta_n$  ?

For any finite partition  $A_1 \dots A_K$  of  $\Theta$ ,

define the number of observations in  $A_k$  :  $n_k = \#\{ i: \theta_i \in A_k \}$

The **posterior** of  $G$  given observations  $\theta_1 \dots \theta_n$

$$(G(A_1), \dots, G(A_K)) | \theta_1, \dots, \theta_n \sim \text{Dirichlet}(\alpha H(A_1) + \mathbf{n}_1, \dots, \alpha H(A_K) + \mathbf{n}_K)$$

# The posterior of $G$ : $G|\theta_1, \dots, \theta_n$

The observations  $\theta_1 \dots \theta_n$  define an **empirical distribution** over  $\Theta$ :

$$\frac{\sum_{i=1}^n \delta_{\theta_i}}{n}$$

← This is just a fancy way of saying  $P(A_k) = n_k/n$

The **posterior** of  $G$  given observations  $\theta_1 \dots \theta_n$

$$(G(A_1), \dots, G(A_K)) | \theta_1, \dots, \theta_n \sim \text{Dirichlet}(\alpha H(A_1) + n_1, \dots, \alpha H(A_K) + n_K)$$

The posterior is a DP with:

- **concentration parameter**  $\alpha + n$
- a **base distribution** that is a weighted average of  $H$  and the empirical distribution.

$$G | \theta_1, \dots, \theta_n \sim DP\left(\alpha + n, \frac{\alpha}{\alpha + n} H + \frac{n}{\alpha + n} \frac{\sum_{i=1}^n \delta_{\theta_i}}{n}\right)$$

The weight of the empirical distribution is proportional to the amount of data. The weight of  $H$  is proportional to  $\alpha$



# The Blackwell MacQueen urn

Assume each value in  $\Theta$  has a unique color.

$\theta_1 \dots \theta_n$  is a sequence of colored balls.

With probability  $\alpha / (\alpha + n)$ , the  $n+1$ th ball is drawn from  $H$

With probability  $n / (\alpha + n)$  the  $n+1$ th ball is drawn from an urn that contains all previously drawn balls.

Note that this implies that  $G$  is a discrete distribution, even if  $H$  is not.

# The clustering property of DPs

$\theta_1 \dots \theta_n$  induces a partition of the set  $1 \dots n$  into  $k$  unique values.

This means that the DP defines a distribution over such partitions.

The expected number of clusters  $k$  increases with  $\alpha$  but grows only logarithmically in  $n$ :

$$E[k | n] \approx \alpha \log(1 + n/\alpha)$$

# NLP 101: language modeling

**Task:** Given a stream of words  $w_1 \dots w_n$ , predict the next word  $w_{n+1}$  with a unigram model  $P(w)$

**Answer:**

If  $w_{n+1}$  is a word  $w$  we've seen before:

$$P(w_{n+1} = w) \propto \text{Freq}(w)$$

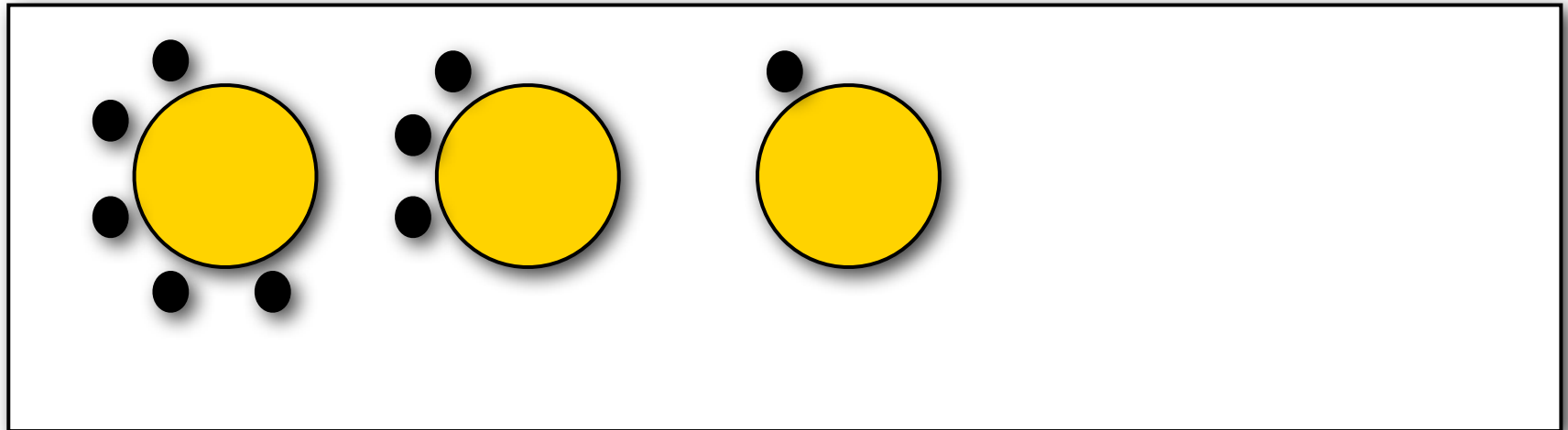
But what if  $w_{n+1}$  has never been seen before?

We need to reserve some mass for new events

$$P(w_{n+1} \text{ is a new word}) \propto \alpha$$

$$\begin{aligned} P(w_{n+1} = w) &= \text{Freq}(w)/(n+\alpha) && \text{if } \text{Freq}(w) > 0 \\ &= \alpha/(n+\alpha) && \text{if } \text{Freq}(w) = 0 \end{aligned}$$

# The Chinese restaurant process



The  $(i+1)$ th customer  $c_{i+1}$  sits:

- at an *existing* table  $t_k$  that already has  $n_k$  customers with probability  $n_k/(i+\alpha)$
- at *new* table with probability  $\alpha/(i+\alpha)$

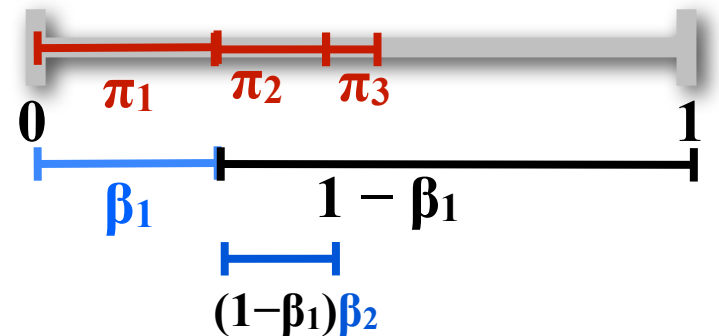
# The predictive distribution

$$\theta_{n+1} | \theta_1, \dots, \theta_n$$

The predictive distribution of  $\theta_{n+1}$  given a sequence of i.i.d. draws  $\theta_1, \dots, \theta_n \sim G$ , with  $G \sim DP(\alpha, H)$  and  $G$  marginalized out is given by the posterior base distribution given  $\theta_1, \dots, \theta_n$

$$\begin{aligned} P(\theta_{n+1} \in A) &= E[G(A) | \theta_1, \dots, \theta_n] \\ &= \frac{\alpha}{\alpha + n} H(A) + \frac{\sum_{i=1}^n \delta_{\theta_i}(A)}{\alpha + n} \end{aligned}$$

# The stick-breaking representation



$G \sim DP(\alpha, H)$  if:

- The component parameters are drawn from the base distribution:  $\theta_k^* \sim H$
- The weights of each cluster are defined by a stick-breaking process:

$$\beta_k \sim \text{Beta}(1, \alpha) \quad \pi_k \sim \beta_k \prod_{l=1 \dots k-1} (1 - \beta_l)$$

also written as  $\pi \sim \text{GEM}(\alpha)$  (Griffiths/Engen/McCloskey)

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

$\theta_k^*$

# Dirichlet Process Mixture Models

Each observation  $x_i$  is associated with a latent parameter  $\theta_i$

Each  $\theta_i$  is drawn i.i.d. from  $G$ ; each  $x_i$  is drawn from  $F(\theta_i)$

$$G|\alpha, H \sim DP(\alpha, H) \quad \theta_i|G \sim G \quad x_i|\theta_i \sim F(\theta_i)$$

Since  $G$  is discrete,  $\theta_i$  can be equal to  $\theta_j$

All  $x_i, x_j$  with  $\theta_i = \theta_j$  belong to the same mixture component

There are a countably infinite number of mixture components.

## Stick-breaking representation:

**Mixing proportions:**  $\pi|\alpha \sim GEM(\alpha)$

**Indicator variables:**  $z_i|\pi \sim Mult(\pi)$

**Component parameters:**  $\theta_k^* | H \sim H$

**Observations:**  $x_i|z_i, \{\theta_k^*\} \sim F(\theta_{z_i}^*)$

# Hierarchical Dirichlet Processes

Since both  $H$  and  $G$  are distributions over the same space  $\Theta$ , the base distribution of a DP can be a draw from another DP. This allows us to specify hierarchical Dirichlet Processes, where each group of data is generated by its own DP.

Assume a global measure  $G_0$  drawn from a DP:

$$G_0 \sim \text{DP}(\gamma, H)$$

For each group  $j$ , define another DP  $G_j$  with base measure  $G_0$ :

$$G_j \sim \text{DP}(\alpha_0, G_0)$$

(or  $G_j \sim \text{DP}(\alpha_j, G_0)$ , but it is common to assume all  $\alpha_j$  are the same)

$\alpha_0$  specifies the amount of variability around the prior  $G_0$

Since all groups share the same base  $G_0$ , all  $G_j$  use the same atoms (balls of the same colors)