

Tutorial on Bayesian Non-parametric methods

Erin Grant

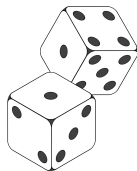
Department of Computer Science, University of Toronto

November 27th, 2015



Intuition for the Dirichlet Distribution (1)

Consider a six-sided die.



Intuition for the Dirichlet Distribution (1)

Consider a six-sided die.



The possible outcomes $\{1, 2, 3, 4, 5, 6\}$ of rolling the die:

- ▶ are discrete (countable);
- ▶ are disjoint;
- ▶ have probabilities:

$$P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = \frac{1}{6}.$$

Intuition for the Dirichlet Distribution (2)

How do we model manufacturing error that causes the die outcome probabilities $\{1, 2, 3, 4, 5, 6\}$ to be skewed?

Intuition for the Dirichlet Distribution (2)

How do we model manufacturing error that causes the die outcome probabilities $\{1, 2, 3, 4, 5, 6\}$ to be skewed?

We can use the **Dirichlet** distribution, so that the probabilities themselves are distributed:

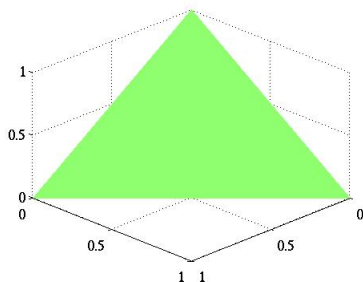
$$[P(1), P(2), P(3), P(4), P(5), P(6)] \sim \text{Dir}(\vec{\alpha})$$

where $\vec{\alpha} = (\alpha_1, \dots, \alpha_6)$ is a parameter vector of **pseudocounts** (prior expectations).

Math Background for the Dirichlet Distribution (1)

Probability simplex:

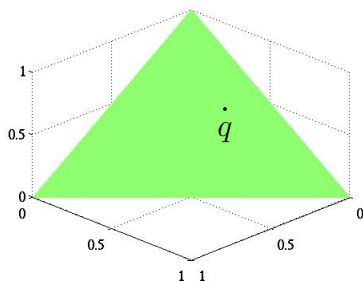
$$\{\vec{x} \in \mathbb{R}^k \mid x_1 + \cdots + x_k = 1, x_1, \dots, x_k \geq 0\}$$



The 2-dimensional probability simplex in \mathbb{R}^3 .

Math Background for the Dirichlet Distribution (2)

Each point q corresponds to a **probability mass function** (pmf) over k disjoint events.



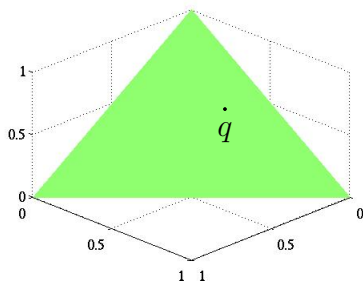
$$\text{E.g.: } q = [0.3, 0.1, 0.6]$$



pmf over three events with probabilities $\frac{3}{10}$, $\frac{1}{10}$, and $\frac{3}{5}$.

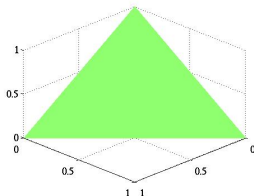
Math Background for the Dirichlet Distribution (3)

The **Dirichlet distribution** defines a probability for each point q in the simplex (i.e., it is a pmf over pmfs).

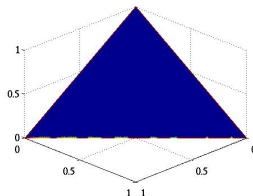


$$P(q) = \text{Dir}(\vec{\alpha}) \text{ for some } \vec{\alpha} = (\alpha_1, \alpha_2, \alpha_3).$$

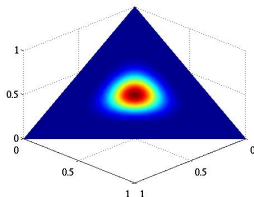
The Parameter Vector $\vec{\alpha}$



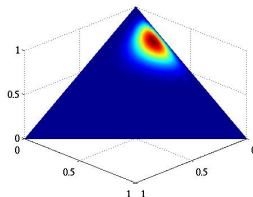
$$\alpha = [1, 1, 1]$$



$$\alpha = [.1, .1, .1]$$



$$\alpha = [10, 10, 10]$$



$$\alpha = [2, 5, 15]$$

Refresher: Bayesian Updating and Conjugacy (1)

We have a **likelihood** distribution $P(X \mid \theta)$, parametrised by some parameters θ , and a **prior** $P(\theta)$ over the parameters θ .

We want to infer a **posterior** distribution over the parameters θ after seeing some observations X_1, \dots, X_N .

We use Bayes' Rule:

$$P(\theta \mid X_1, \dots, X_N) = \frac{P(X_1, \dots, X_N \mid \theta) P(\theta)}{P(X_1, \dots, X_N)}$$

Refresher: Bayesian Updating and Conjugacy (1)

We have a **likelihood** distribution $P(X \mid \theta)$, parametrised by some parameters θ , and a **prior** $P(\theta)$ over the parameters θ .

We want to infer a **posterior** distribution over the parameters θ after seeing some observations X_1, \dots, X_N .

We use Bayes' Rule:

$$P(\theta \mid X_1, \dots, X_N) = \frac{P(X_1, \dots, X_N \mid \theta) P(\theta)}{P(X_1, \dots, X_N)}$$

But $P(X_1, \dots, X_N) = \int P(X_1, \dots, X_N \mid \theta') P(\theta') d\theta'$ is usually hard to compute.

Refresher: Bayesian Updating and Conjugacy (2)

But $P(X_1, \dots, X_N) = \int P(X_1, \dots, X_N \mid \theta') P(\theta') d\theta'$ is usually hard to compute.

Solution: Use a **conjugate** prior so that the posterior can be computed in closed form.

E.g., the Dirichlet is conjugate to the **Multinomial** distribution.

The Multinomial Distribution (1)

The Multinomial distribution gives the probability of N outcomes to be distributed over K categories:

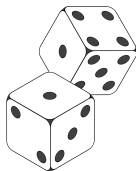
$$(X_1, \dots, X_N) \sim \text{Multi}(n, (q_1, \dots, q_K)),$$

where

- ▶ X_i is the number of times that the i th category occurred amongst the N events;
- ▶ $\vec{q} = (q_1, \dots, q_K)$ gives the probabilities for each of the K categories to occur.

The Multinomial Distribution (2)

Example: Roll a die 5 times. What are the outcomes?



Let $X_i \in \{1, 2, 3, 4, 5, 6\}$ be the outcome of the i th roll. Then

$$\vec{X} = (X_1, X_2, X_3, X_4, X_5) \sim \text{Multi}(n, \vec{q}),$$

with $n = 5$ and $\vec{q} = (\frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6})$.

Bayesian Updating for the Dirichlet-Multinomial

Use a Dirichlet prior over the probability vector \vec{Q} :

$$\vec{Q} \sim \text{Dir}(\vec{\alpha})$$



$$(\vec{X} \mid \vec{Q}) \sim \text{Multi}(n, \vec{Q})$$



$$(\vec{Q} \mid \vec{X} = \vec{x}) \sim \text{Dir}(\vec{\alpha} + \vec{x}).$$

Application: Dirichlet-Multinomial Mixture Model

$$\vec{Q} \mid \vec{\alpha} \sim \text{Dir}(\vec{\alpha})$$



component assignment
probability

$$Z_1, \dots, Z_N \mid \vec{Q} \sim \text{Multi}(N, \vec{Q})$$



component assignment
variables $Z_i \in \{1, \dots, k\}$

$$\theta_k \sim \mathcal{G}$$



prior over parameters of
component distribution

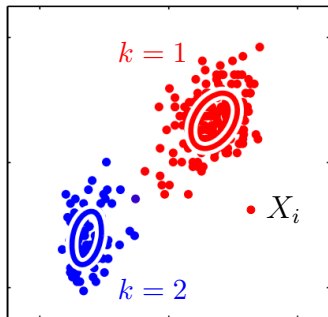
$$X_i \mid Z_i, \theta_{Z_i} \sim \mathcal{F}(\theta_{Z_i})$$



component likelihood
distribution

Application: Dirichlet-Multinomial Mixture Model

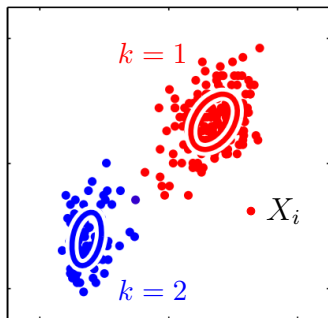
$$\vec{Q} \mid \vec{\alpha} \sim \text{Dir}(\vec{\alpha}) \quad \leftarrow \text{component assignment probability}$$



Determines $P(k)$.

Application: Dirichlet-Multinomial Mixture Model

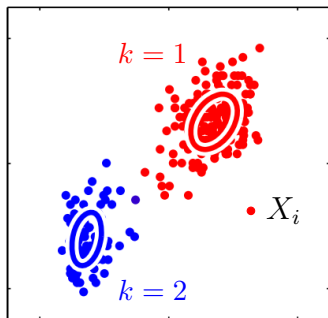
$$Z_1, \dots, Z_N \mid \vec{Q} \sim \text{Multi}(N, \vec{Q}) \quad \leftarrow \text{component assignment variables } Z_i \in \{1, \dots, k\}$$



Determines $P(Z_i = k \mid \vec{Q})$.

Application: Dirichlet-Multinomial Mixture Model

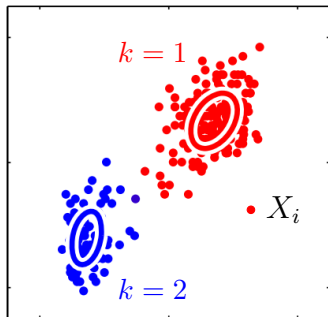
$\theta_k \sim \mathcal{G}$ ← prior over parameters of component distribution



Determines $P(\theta_k)$.

Application: Dirichlet-Multinomial Mixture Model

$$X_i \mid Z_i, \theta_{Z_i} \sim \mathcal{F}(\theta_{Z_i}) \quad \leftarrow \text{component likelihood distribution}$$



Determines $P(X_i \mid Z_i, \theta_{Z_i})$.

Application: Dirichlet-Multinomial Mixture Model

By conjugacy, the posterior is tractable to compute:

$$P(Z_1, \dots, Z_N, \theta_1, \dots, \theta_K \mid X_1, \dots, X_N)$$

Application: Dirichlet-Multinomial Mixture Model

By conjugacy, the posterior is tractable to compute:

$$P(Z_1, \dots, Z_N, \theta_1, \dots, \theta_K \mid X_1, \dots, X_N)$$

- ▶ Allows us to answer:
 - ▶ What cluster do the instances belong to? (Z_1, \dots, Z_N)
 - ▶ What are the properties of the clusters? $(\theta_1, \dots, \theta_K)$

Application: Dirichlet-Multinomial Mixture Model

By conjugacy, the posterior is tractable to compute:

$$P(Z_1, \dots, Z_N, \theta_1, \dots, \theta_K \mid X_1, \dots, X_N)$$

- ▶ Allows us to answer:
 - ▶ What cluster do the instances belong to? (Z_1, \dots, Z_N)
 - ▶ What are the properties of the clusters? $(\theta_1, \dots, \theta_K)$
- ▶ Applications:
 - ▶ X_i could be a document and Z_i the topic of X_i .

Background: Stochastic Processes (1)

A **stochastic process** is a collection of random variables indexed by some index set:

$$\{X_i\} \quad i \in \mathcal{I} \quad X_i \sim \mathcal{D}(\theta).$$

Background: Stochastic Processes (1)

A **stochastic process** is a collection of random variables indexed by some index set:

$$\{X_i\} \quad i \in \mathcal{I} \quad X_i \sim \mathcal{D}(\theta).$$

Any finite subset of these variables has a joint distribution; e.g.,

$$p(X_{j_1}, \dots, X_{j_n}) \sim \mathbb{D}(\theta), \quad (j_1, \dots, j_n) \subset \mathbb{N}.$$

Background: Stochastic Processes (2)

Why use a stochastic process instead of a set of random variables?

Background: Stochastic Processes (2)

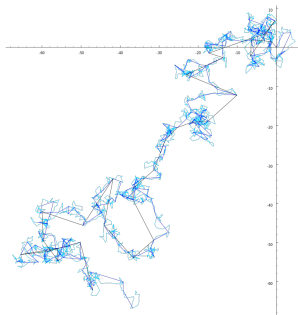
Why use a stochastic process instead of a set of random variables?

- ▶ Index set has unknown dimension
 - ▶ e.g., Topic modelling: words $w \in \{\text{corpus}\}$ are indexed by topics t , where the number of topics is not known but inferred from the data (Teh et al. [2006])

Background: Stochastic Processes (2)

Why use a stochastic process instead of a set of random variables?

- ▶ The index set is infinite-dimensional
 - ▶ e.g., Brownian motion: the random displacement $X(t)$ of a particle depend on a continuous time $t \in \mathbb{R}$



Intuition for the Dirichlet Process

We want to generalize the Dirichlet distribution to be a pmf over an infinite number of events:

$$\text{Dir}(\alpha_1, \dots, \alpha_K) \quad K \rightarrow \infty$$

It is a **non-parametric** model since the number of categories K is not fixed, and instead grows with the data.

Intuition for the Dirichlet Process

We want to generalize the Dirichlet distribution to be a pmf over an infinite number of events:

$$\text{Dir}(\alpha_1, \dots, \alpha_K) \quad K \rightarrow \infty$$

It is a **non-parametric** model since the number of categories K is not fixed, and instead grows with the data.

It becomes a **Dirichlet process** instead of a distribution.

Background: σ -Algebra (1)

A σ -Algebra over a set \mathcal{B} is a collection of subsets of \mathcal{B} that is *closed* under countably many of the following operations:

Complement: if $A \in \mathcal{B}$ then $A^C \in \mathcal{B}$;

Union: if $A_1, A_2, \dots \in \mathcal{B}$ then $\bigcup A_i \in \mathcal{B}$;

Intersection: if $A_1, A_2, \dots \in \mathcal{B}$ then $\bigcap A_i \in \mathcal{B}$.

Background: σ -Algebra (2)

Example: Let $\mathcal{B} = \{ \text{green}, \text{yellow}, \text{red} \}$ be the possible colours of a traffic light that you encounter while driving.

Then

$$\sigma = \left\{ \text{green}, \text{yellow}, \text{red}, \{ \text{green}, \text{yellow} \}, \{ \text{green}, \text{red} \}, \right. \\ \left. \{ \text{yellow}, \text{red} \}, \{ \text{green}, \text{yellow}, \text{red} \}, \emptyset \right\}$$

is a σ -algebra on \mathcal{B} .

Dirichlet Process: Introduction (1)

A Dirichlet process takes as its **index set** a σ -algebra over a space \mathcal{B} :

\forall sets $B \in \sigma(\mathcal{B})$, $\tilde{P}(B) \in [0, 1]$ is a random variable.

Dirichlet Process: Introduction (1)

A Dirichlet process takes as its **index set** a σ -algebra over a space \mathcal{B} :

\forall sets $B \in \sigma(\mathcal{B})$, $\tilde{P}(B) \in [0, 1]$ is a random variable.

Dirichlet Process: Introduction (2)

Marginalisation Property: For any finite partition (B_1, \dots, B_N) of the space \mathcal{B} , if $G \sim \text{DP}(\alpha, H)$, then

$$\begin{bmatrix} \tilde{P}(B_1) \\ \vdots \\ \tilde{P}(B_N) \end{bmatrix} \sim \text{Dir} \begin{pmatrix} \alpha H(B_1) \\ \vdots \\ \alpha H(B_N) \end{pmatrix}$$

Dirichlet Process: Introduction (2)

Marginalisation Property: For any finite partition (B_1, \dots, B_N) of the space \mathcal{B} , if $G \sim \text{DP}(\alpha, H)$, then

$$\begin{bmatrix} \tilde{P}(B_1) \\ \vdots \\ \tilde{P}(B_N) \end{bmatrix} \sim \text{Dir} \begin{pmatrix} \alpha H(B_1) \\ \vdots \\ \alpha H(B_N) \end{pmatrix}$$

Implication: Draws from the Dirichlet Process are random probability distributions.

Example of a “DP” Indexed by a Finite σ -algebra

$$\begin{bmatrix} \tilde{P}(\{\bullet, \bullet\}) \\ \tilde{P}(\bullet) \end{bmatrix} \sim \text{Dir} \begin{pmatrix} \alpha H(\{\bullet, \bullet\}) \\ \alpha H(\bullet) \end{pmatrix}$$

Example of a “DP” Indexed by a Finite σ -algebra

$$\begin{bmatrix} \tilde{P}(\{\text{green}, \text{yellow}, \text{red}\}) \\ \tilde{P}(\emptyset) \end{bmatrix} \sim \text{Dir} \begin{pmatrix} \alpha H(\{\text{green}, \text{yellow}, \text{red}\}) \\ \alpha H(\emptyset) \end{pmatrix}$$

Example of a “DP” Indexed by a Finite σ -algebra

$$\begin{bmatrix} \tilde{P}(\bullet) \\ \tilde{P}(\bullet) \\ \tilde{P}(\bullet) \end{bmatrix} \sim \text{Dir} \begin{pmatrix} \alpha H(\bullet) \\ \alpha H(\bullet) \\ \alpha H(\bullet) \end{pmatrix}$$

Dirichlet Process: Definition

$$[\tilde{P}(B_1), \dots, \tilde{P}(B_N)] \sim \text{Dir}(\alpha H(B_1), \dots, \alpha H(B_N))$$

H is the (non-random) **base distribution** over \mathcal{B} .

(Can be any probability distribution over \mathcal{B} .)

It determines the mean of the DP for any set B :

$$\mathbb{E}(\tilde{P}(B)) = H(B).$$

Dirichlet Process: Definition

$$[\tilde{P}(B_1), \dots, \tilde{P}(B_N)] \sim \text{Dir}(\alpha H(B_1), \dots, \alpha H(B_N))$$

α is a positive real **concentration parameter**.

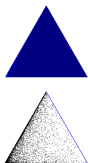
It determines the concentration of the DP about the mean:

$$\text{Var}(\tilde{P}(B)) = \frac{H(B)(1 - H(B))}{\alpha + 1}.$$

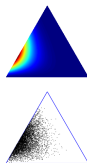
Example for a Finite σ -algebra

Suppose the base distribution is given by

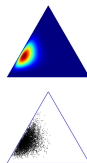
$$H(\text{green circle}) = 0.6 \quad H(\text{yellow circle}) = 0.1 \quad H(\text{red circle}) = 0.3.$$



$\alpha = 2$



$\alpha = 10$



$\alpha = 20$

Generalization to Infinite σ -algebra (1)

Suppose for a space \mathcal{B} , $\sigma(\mathcal{B})$ is (countably) infinite.

Generalization to Infinite σ -algebra (1)

Suppose for a space \mathcal{B} , $\sigma(\mathcal{B})$ is (countably) infinite.

Example: We want to create a generative unigram model of a text.

Generalization to Infinite σ -algebra (1)

Suppose for a space \mathcal{B} , $\sigma(\mathcal{B})$ is (countably) infinite.

Example: We want to create a generative unigram model of a text.

We want to assign non-zero probability to the next word in the document.

Generalization to Infinite σ -algebra (1)

Suppose for a space \mathcal{B} , $\sigma(\mathcal{B})$ is (countably) infinite.

Example: We want to create a generative unigram model of a text.

We want to assign non-zero probability to the next word in the document.

But how do we distribute the probability mass when the number of words (partitions) is unknown (and possibly infinite)?

Generalization to Infinite σ -algebra (2)

Solution: Model the document as a Dirichlet Process $DP(\alpha, H)$ (with H initially uniform).

Then a vocabulary of N symbols corresponds to a partition into N parts.

Generalization to Infinite σ -algebra (3)

Under this model, the probability of a symbol w is given by

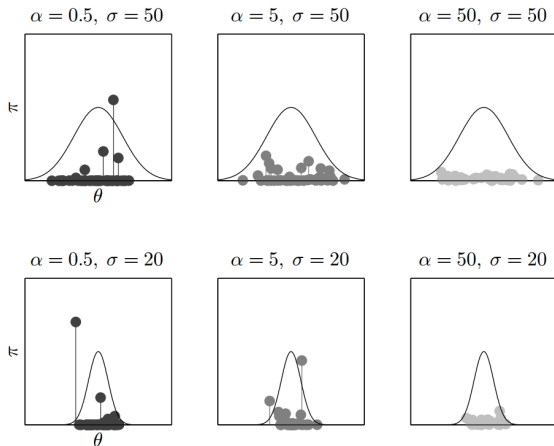
$$P(w) = \frac{F_w}{F + \alpha}$$

if w is seen, and

$$P(w) = \frac{\alpha}{F + \alpha}$$

if w is unseen.

Draws from a DP with a Gaussian H



Rows vary σ ; columns vary α .

Recall: Dirichlet-Multinomial Mixture Model (1)

$$\vec{Q} \mid \vec{\alpha} \sim$$

$$\text{Dir}(\vec{\alpha})$$



cluster assignment
probability

$$Z_1, \dots, Z_N \mid \vec{Q} \sim \text{Multi}(N, \vec{Q})$$



cluster assignment
variables

$$\theta_k \sim$$

$$\mathcal{G}$$



prior over parameters of
component distribution

$$X_i \mid Z_i, \theta_{Z_i} \sim$$

$$\mathcal{F}(\theta_{Z_i})$$



component distribution

Recall: Dirichlet-Multinomial Mixture Model (2)

$$(Q_1, \dots, Q_k) \mid \alpha_1, \dots, \alpha_k \sim \text{Dir}(\alpha_1, \dots, \alpha_k)$$

What if $K \rightarrow \infty$?

Dirichlet Process Mixture Model (1)

$Q \mid \alpha, H \sim \text{DP}(\alpha, H)$  random distribution over parameters

$\theta_i \mid Q \sim Q$  latent parameter for X_i

$X_i \mid \theta_i \sim \mathcal{F}(\theta_i)$  likelihood distribution

Dirichlet Process Mixture Model (1)

$Q \mid \alpha, H \sim \text{DP}(\alpha, H)$  random distribution over parameters

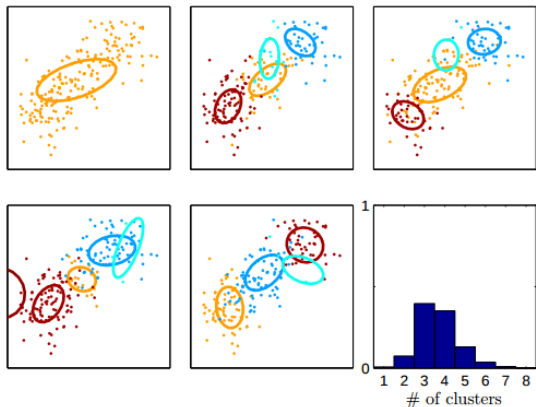
$\theta_i \mid Q \sim Q$  latent parameter for X_i

$X_i \mid \theta_i \sim \mathcal{F}(\theta_i)$  likelihood distribution

Q is discrete, so multiple θ_i can take on the same value (i.e., they **cluster**).

Dirichlet Process Mixture Model (2)

Posterior inference for the number of clusters can be done using Markov Chain Monte Carlo (MCMC):



Application: Dirichlet Processes in NLP

Language modelling: The number of words is unbounded.

Topic modelling: The number of topics is inferred from the data.

Discussion: Relevance to the Word Learning Model

In the present model, the meaning probability corresponds to the expected value of the posterior of a Dirichlet distribution:

$$p_t(f \mid w) = \frac{\text{assoc}_t(w, f) + \gamma}{\sum_{f'} \text{assoc}_t(w, f') + k \cdot \gamma}$$

However, we fix the number of features k ahead of time.

Discussion: Relevance to the Word Learning Model

In the present model, the meaning probability corresponds to the expected value of the posterior of a Dirichlet distribution:

$$p_t(f \mid w) = \frac{\text{assoc}_t(w, f) + \gamma}{\sum_{f'} \text{assoc}_t(w, f') + k \cdot \gamma}$$

However, we fix the number of features k ahead of time.

Can we effectively use a Dirichlet Process to allow the number of features to grow with the data?

References

- B. A. Frigyik, A. Kapila, and M. R. Gupta. Introduction to the Dirichlet Distribution and Related Processes. Technical report, 2010. URL <http://scholar.google.com/scholar?hl=en{%&}btnG=Search{%&}q=intitle:Introduction+to+the+Dirichlet+Distribution+and+Related+Processes{%#}0>.
- Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical dirichlet processes. *Journal of the american statistical association*, 101(476), 2006.