



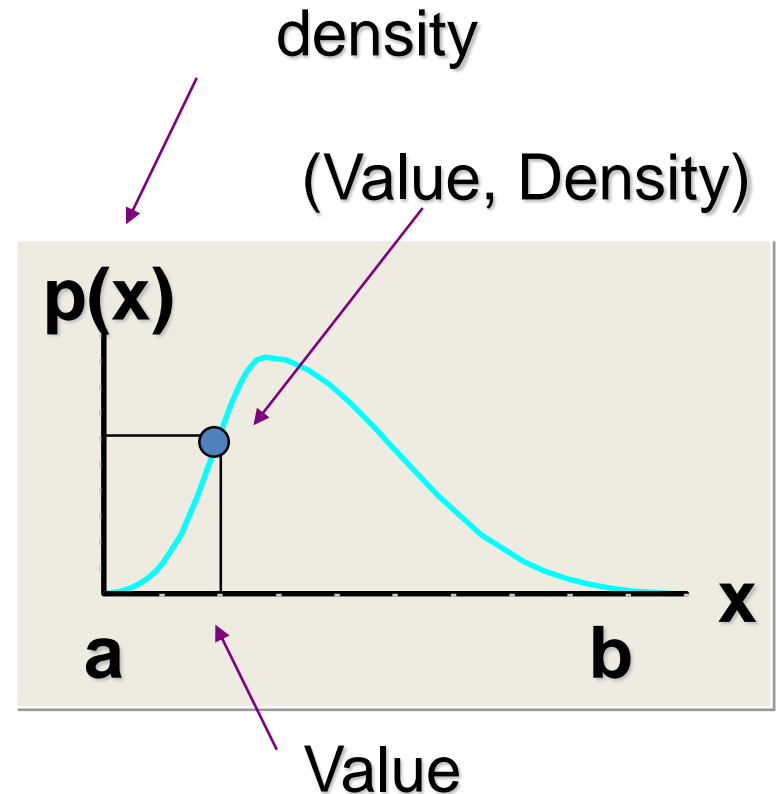
Bayesian ML : Posterior Distributions and Mixture Models

Suresh
Manandhar
Naamii, Nepal

Continuous Probability Density Function (pdf)

- Shows all values of x in the given interval $[a, b]$, the density $p(x)$
- $p(x)$ is a **probability density function (pdf)**
- Since probabilities need to sum to **1**, the corresponding condition is:

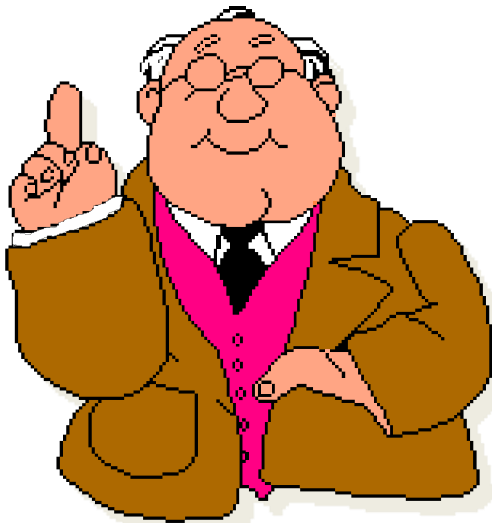
$$F(-\infty \leq X \leq \infty) \\ = \int_{-\infty}^{\infty} p(x) dx = 1$$



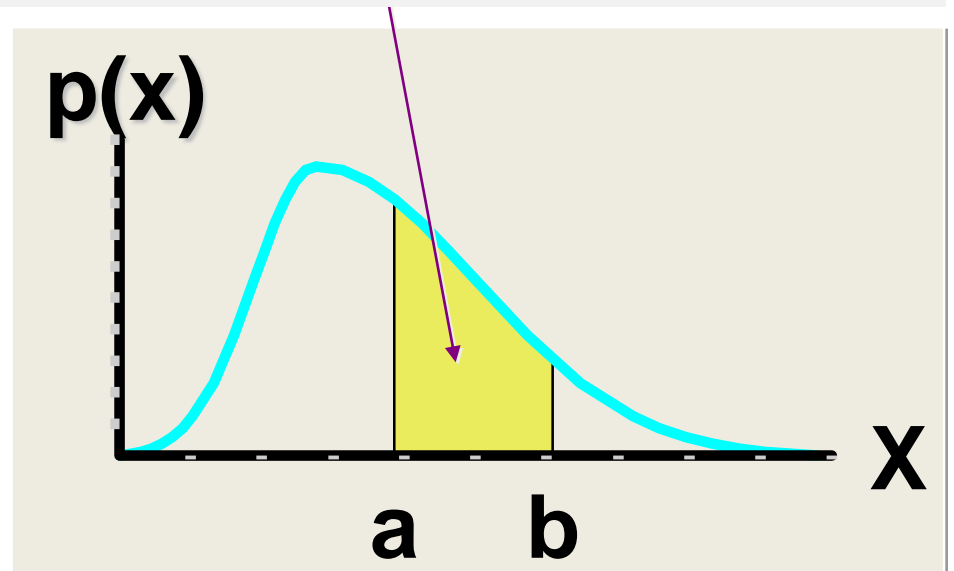
Cumulative density (cdf)

The probability that x lies in the interval $[a, b]$ is given by $F(a \leq X \leq b)$.

**Cumulative probability is
Area Under Curve!**



$$F(a \leq X \leq b) = \int_a^b p(x) dx$$



Some properties

$$F(-\infty \leq X \leq \infty) = \int_{-\infty}^{\infty} p(x) dx = 1$$

$$F(X = a) = F(a \leq X \leq a) = \int_a^a p(x) dx = 0$$

$$F(a) = F(-\infty \leq X \leq a) = \int_{-\infty}^a p(x) dx$$

- or, more generally:

$$F(x) = \int p(x) dx$$

- $F(x)$ is known as the **cumulative distribution function**. (**CDF**).
- The pdf and cdf are related by:

$$\frac{d}{dx} F(x) = p(x)$$

- The following provides an intuitive understanding:

$$F\left(a - \frac{\Delta}{2} \leq X \leq a + \frac{\Delta}{2}\right) = \int_{a - \frac{\Delta}{2}}^{a + \frac{\Delta}{2}} p(x) dx \simeq \Delta p(a)$$

Expectation and Variance

Weighted Average

$$E[X] = \int_{-\infty}^{\infty} x p(x) dx$$

Averaged Squared Distance From Mean

$$Var(X) = \int_{-\infty}^{\infty} (x - E[X])^2 p(x) dx$$

Bayes for Continuous random variables

Law of Total Probability

- For **discrete** random variables:

$$\begin{aligned} p(X = x) &= \sum_i p(X = x, Y = y_i) \\ &= \sum_i p(X = x | Y = y_i) p(Y = y_i) \end{aligned}$$

- For **continuous** random variables:

$$\begin{aligned} p(x) &= \int p(x, y) dy \\ &= \int p(x | y) p(y) dy \end{aligned}$$

Bayes for Continuous random variables

- From the definition of conditional probability:

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_y p(x|y)p(y)}$$

- Substituting for $p(x)$ we derive.

Bayes' theorem for continuous random variables:

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(x|y)p(y)}{\int p(x|y) p(y) dy}$$

Some continuous probability distributions

Distribution	Support/Range
■ Uniform distribution	(a, b)
■ Gaussian (or Normal) distribution	$(-\infty, +\infty)$
■ Exponential distribution	$(0, +\infty]$
■ Beta distribution	$[0, 1]$
■ Dirichlet distribution	$[0, 1]^n$
■ Multivariate Gaussian distribution	$(-\infty, +\infty)^n$

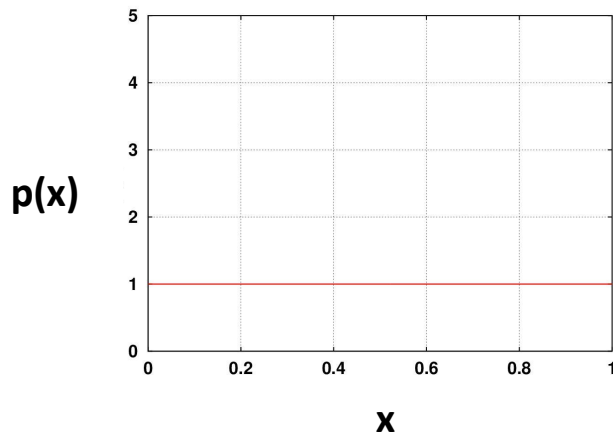
Uniform Distribution

- **Definition:** A random variable is uniformly distributed in the interval $(0, 1)$ if:

$$p(x) = \begin{cases} 1 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

- Hence:

$$\int_{-\infty}^{\infty} p(x) dx = \int_0^1 1 dx = x \Big|_0^1 = 1 - 0 = 1$$



Uniform Distribution

- **Definition:** A random variable is uniformly distributed in the interval $(0, 1)$ if:

$$p(x) = \begin{cases} 1 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

- Hence:

$$\int_{-\infty}^{\infty} p(x) dx = \int_0^1 1 dx = x \Big|_0^1 = 1 - 0 = 1$$

- Then for $a \leq X \leq b$:

$$F(0 < a \leq X \leq b < 1) = \int_a^b 1 dx = x \Big|_a^b = b - a$$

Uniform Distribution

- **Definition:** A random variable is uniformly distributed in the interval $(0, 1)$ if:

$$p(x) = \begin{cases} 1 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

- Hence:

$$\int_{-\infty}^{\infty} p(x) dx = \int_0^1 1 dx = x \Big|_0^1 = 1 - 0 = 1$$

- **Definition:** A random variable is uniformly distributed in the interval (a, b) if:

$$p(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$

- Hence:

$$\int_{-\infty}^{\infty} p(x) dx = \int_a^b \frac{1}{b-a} dx = \frac{x}{b-a} \Big|_a^b = \frac{b}{b-a} - \frac{a}{b-a} = 1$$

Example

- X is uniformly distributed over $(0, 10)$ calculate:
 $F(X < 3)$, $F(X > 3)$, $F(1 < X < 6)$

Example

- X is uniformly distributed over $(0, 10)$ calculate:

$$F(X < 3), \quad F(X > 3), \quad F(1 < X < 6)$$

- **Solution:**

$$p(x) = \begin{cases} \frac{1}{10} & \text{if } 0 < x < 10 \\ 0 & \text{otherwise} \end{cases}$$

$$F(X < 3) = \int_{-\infty}^3 p(x) dx = \frac{\int_0^3 dx}{10} = \frac{3}{10}$$

$$F(X > 3) = \int_3^{+\infty} p(x) dx = \frac{\int_3^{10} dx}{10} = \frac{7}{10}$$

$$F(1 < X < 6) = \int_1^6 p(x) dx = \frac{\int_1^6 dx}{10} = \frac{5}{10} = \frac{1}{2}$$

Normal Distribution

- The normal (or Gaussian) distribution gives the familiar bell shaped curve
- **Definition:** A random variable X is normally distributed if its probability density function is given by:

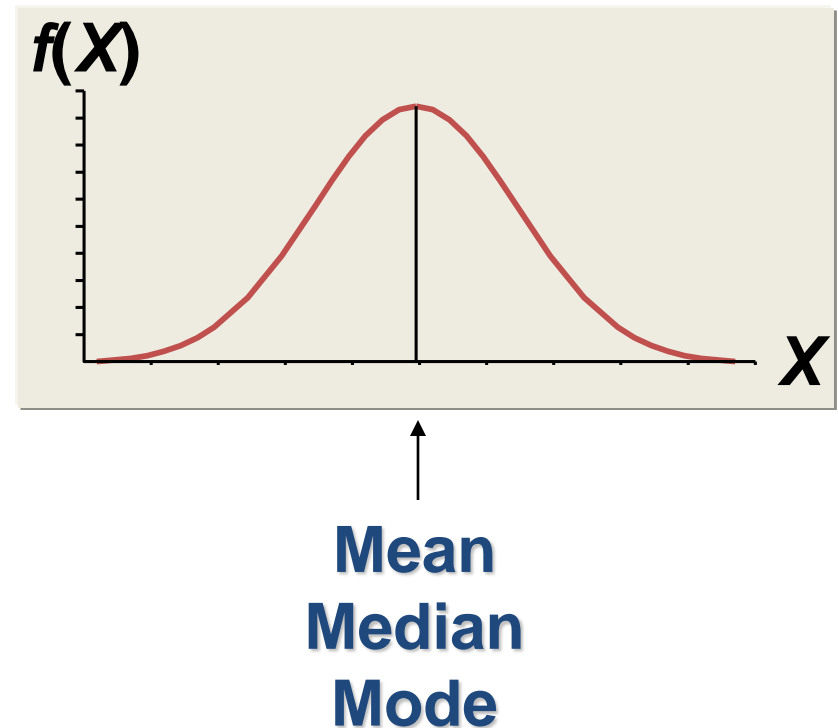
$$p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- The parameters μ and σ^2 are known as the mean and variance.
- It turns out that:

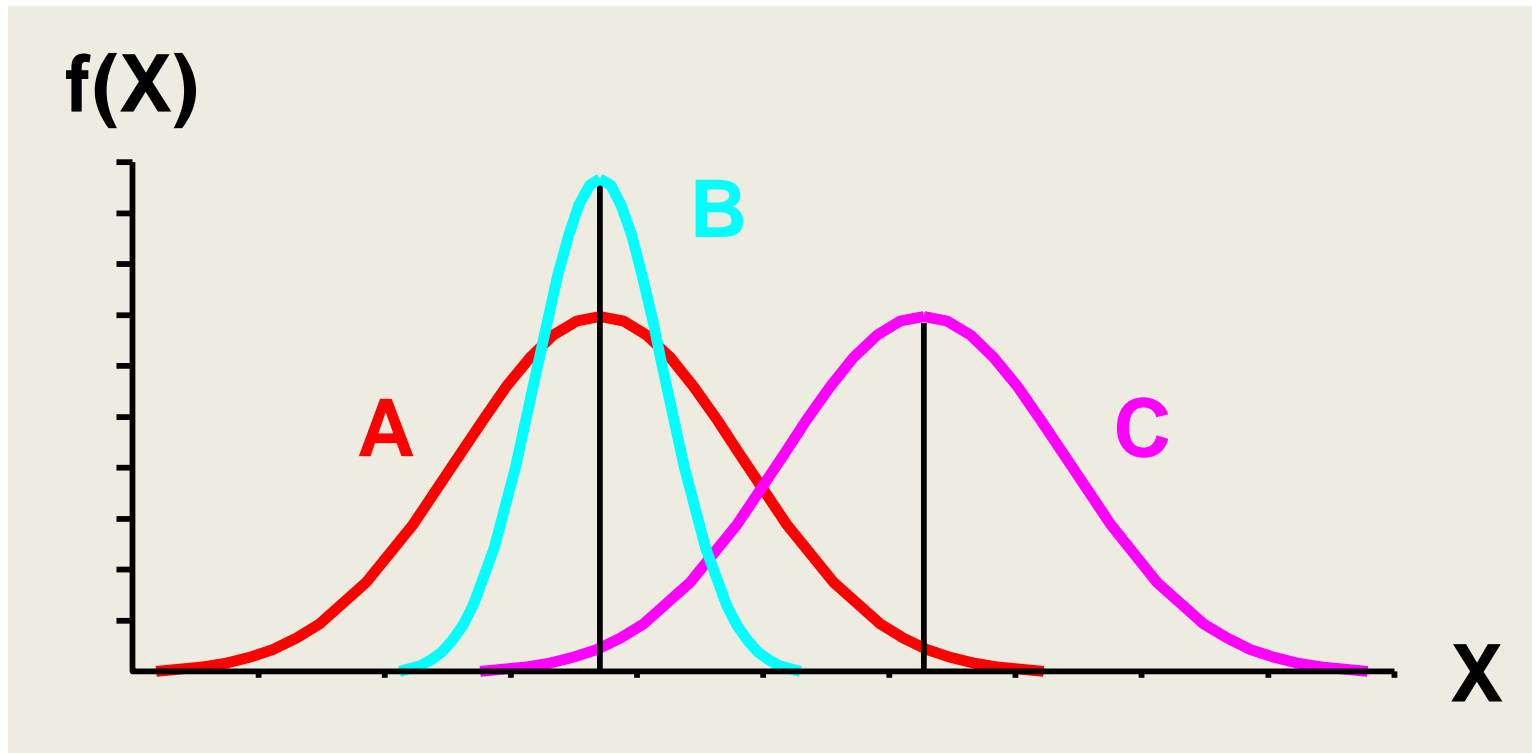
$$\begin{aligned} E[X] &= \mu \\ \text{Var}(X) &= \sigma^2 \end{aligned}$$

The Normal Distribution

- ☐ 1. 'Bell-Shaped' & Symmetrical
- ☐ 2. Mean, median, mode are equal
- ☐ 3. Random variable has infinite range



Effect of varying μ and σ^2

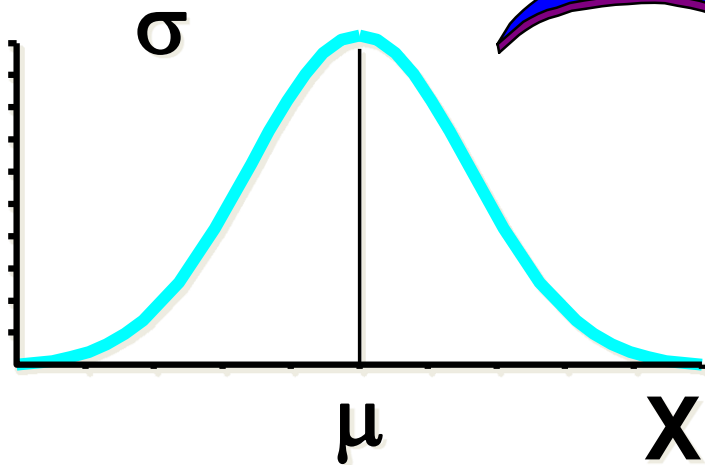


Standardize the Normal Distribution

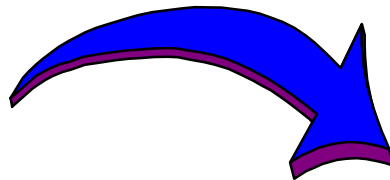
$$p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$p(X = x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

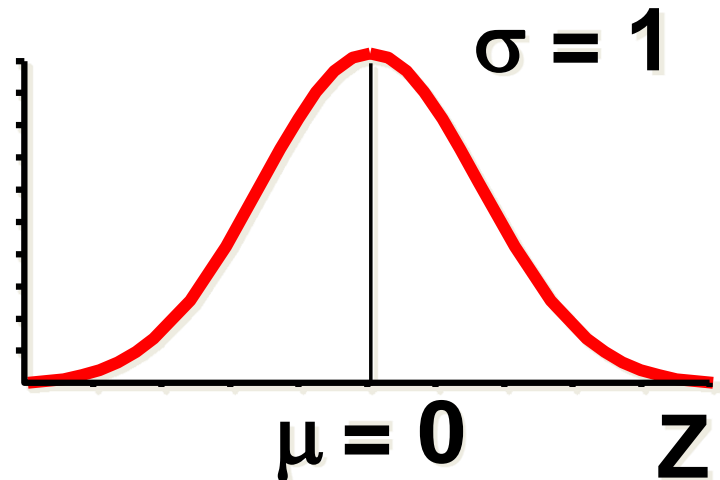
Normal
Distribution



$$Z = \frac{X - \mu}{\sigma}$$



Standardized
Normal Distribution



One table!

Some properties of Normal distribution

■ Central Limit Theorem

- *The sum of a set of independent random variables approaches the **Gaussian distribution** as the number of variables $\rightarrow \infty$, regardless of the distributions of the individual variables*
- (There are generalisations of the **CLT**.)
- **Example:** the sum of the face value of randomly drawn playing cards has an approximately Gaussian distribution
- Simplicity: specified by only two intuitive parameters
- Mathematically tractable
 - Many analyses turn out very simply with the Normal distribution

Beta Distribution

- Suppose that the coin factory makes coins that are not perfectly fair all the time
- So, most times the coins are fair but less often the coins are slightly unfair, and less-less often the coins are quite unfair
- So, there is a distribution over the bias of the coins
- **Question:** How to model this distribution?
- The output/sample from this distribution will be the **bias** of a coin i.e. a number between **0** and **1**

Gamma function

- The gamma function $\Gamma(N)$ generalises the factorial function to the reals such that:

$$\Gamma(N + 1) = N! \quad \text{for natural number } N$$

$$\Gamma(1) = 0! = 1$$

$$\Gamma\left(\frac{1}{2}\right) = \left(-\frac{1}{2}\right)! = \sqrt{\pi}$$

$$\Gamma(0) = (-1)! = \frac{\pi}{2}$$

- The gamma function can be used to *interpolate* for values for which the factorial is undefined.

Beta distribution

$$\mathbf{p}(\boldsymbol{\theta}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\boldsymbol{\theta}^{\boldsymbol{\alpha}-1}(\mathbf{1} - \boldsymbol{\theta})^{\boldsymbol{\beta}-1}}{B(\boldsymbol{\alpha}, \boldsymbol{\beta})} \sim \textit{Beta}(\boldsymbol{\alpha}, \boldsymbol{\beta})$$

- The $-\mathbf{1}$'s can be thought as mathematical convenience/convention
- To be able to show that such a thing exists, we need to show that:

$$\int_0^1 \mathbf{p}(\boldsymbol{\theta}|\boldsymbol{\alpha}, \boldsymbol{\beta}) d\boldsymbol{\theta} = \int_0^1 \frac{\boldsymbol{\theta}^{\boldsymbol{\alpha}-1}(\mathbf{1} - \boldsymbol{\theta})^{\boldsymbol{\beta}-1}}{B(\boldsymbol{\alpha}, \boldsymbol{\beta})} d\boldsymbol{\theta} = \mathbf{1}$$

- Equivalently, need to show that $B(\boldsymbol{\alpha}, \boldsymbol{\beta})$ is well defined:

$$B(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \int_0^1 \boldsymbol{\theta}^{\boldsymbol{\alpha}-1}(\mathbf{1} - \boldsymbol{\theta})^{\boldsymbol{\beta}-1} d\boldsymbol{\theta}$$

Beta function - derivation

Beta function: $B(\alpha, \beta) = \int_0^1 \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta$

- We will apply 'integration by parts', with $u = \theta^{\alpha-1}$, and, $dv = (1 - \theta)^{\beta-1} d\theta$ so the integral becomes

$$B(\alpha, \beta) = \int_0^1 u dv = uv - \int_0^1 v du$$

- $du = (\alpha - 1)\theta^{\alpha-2} d\theta$, and, $v = -\frac{1}{\beta}(1 - \theta)^{\beta}$. Thus:

$$\begin{aligned} &= \theta^{\alpha-1} \left(-\frac{1}{\beta} (1 - \theta)^{\beta} \right) \Big|_0^1 - \int_0^1 \left(-\frac{1}{\beta} (1 - \theta)^{\beta} \right) (\alpha - 1) \theta^{\alpha-2} d\theta \\ &= \frac{(\alpha - 1)}{\beta} \int_0^1 (1 - \theta)^{\beta} \theta^{\alpha-2} d\theta = \frac{(\alpha - 1)}{\beta} B(\alpha - 1, \beta + 1) \\ &= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} B(1, \beta + \alpha - 1) \\ &= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} \int_0^1 \theta^{1-1} (1 - \theta)^{\beta + \alpha - 2} d\theta \end{aligned}$$

$$\begin{aligned}
&= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} \int_0^1 (1 - \theta)^{\beta + \alpha - 2} d\theta \\
&= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} \left(-\frac{(1 - \theta)^{\beta + \alpha - 1}}{\beta + \alpha - 1} \right) \Big|_0^1 \\
&= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)(\beta + \alpha - 1)} \\
&= \frac{\Gamma(\alpha)}{\beta(\beta + 1) \dots (\beta + \alpha - 2)(\beta + \alpha - 1)} \\
&= \frac{\Gamma(\alpha)}{1 \dots (\beta - 2)(\beta - 1)\Gamma(\beta)} \\
&= \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}
\end{aligned}$$

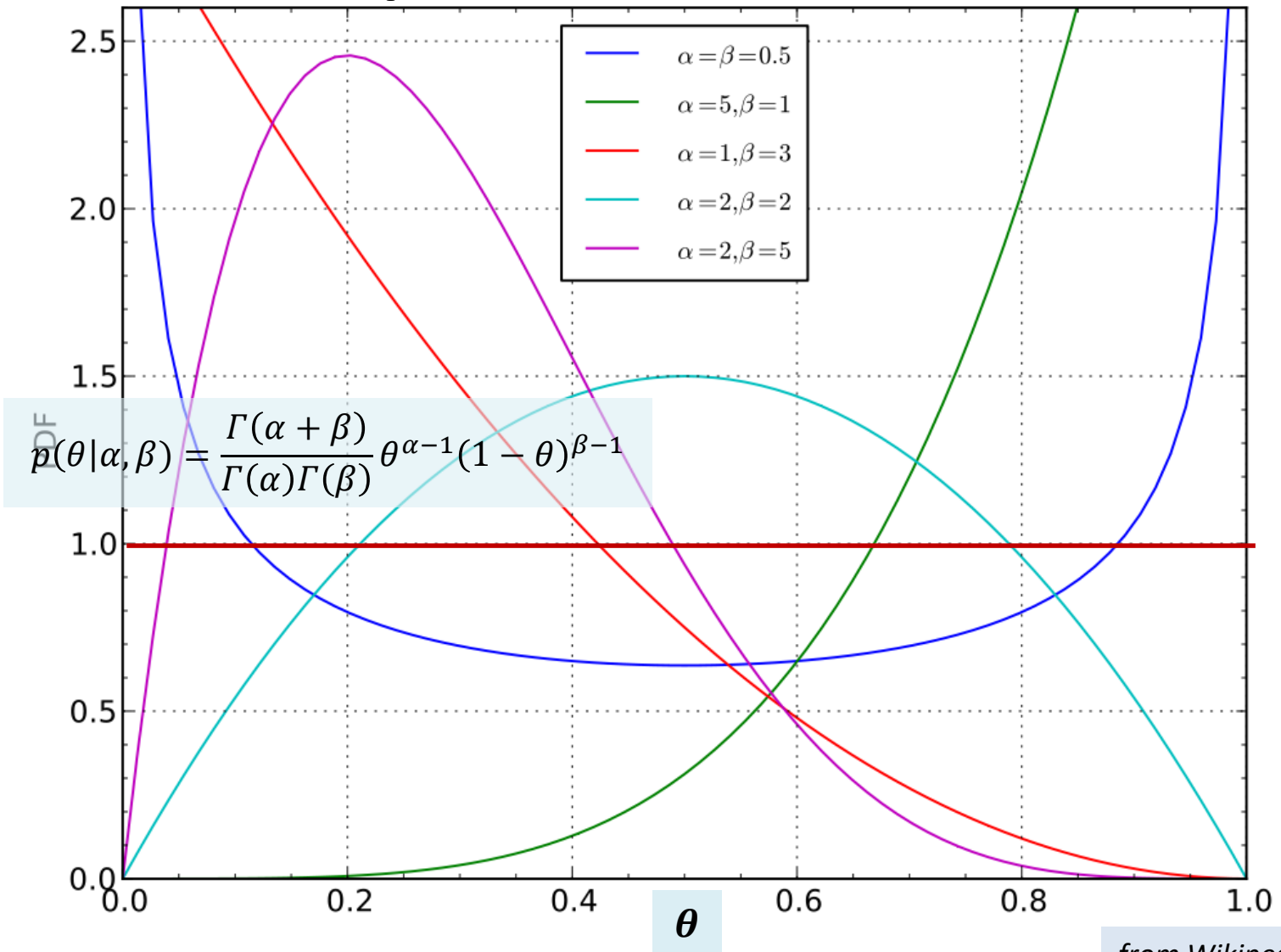
Strictly speaking this proof only applies to positive natural values of α, β

- Thus, $B(\alpha, \beta) = \int_0^1 \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$ hence:

$$p(\theta|\alpha, \beta) = \frac{\theta^{\alpha-1} (1 - \theta)^{\beta-1}}{B(\alpha, \beta)} = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

- For $\alpha = \beta = 1$, the distribution is **uniform** with $p(\theta|\alpha, \beta) = 1$

Shape of Beta distribution



Modelling the coin factory

- If the factory has a high quality control then choosing $\alpha = \beta = 100$ will give a very peaky pdf
- What would choosing a uniform distribution, $\alpha = \beta = 1$, give? Would that be a good factory?
- Similarly, what about values less than 1?

Modelling typical coins from coin factory

- Suppose, the factory inspector visits the coin factory and picks a coin at random from the factory
- Remember, the factory has beta parameters α, β
- The inspector flips the coin N times and sees

$$\mathbf{c} = (c_1, c_2)$$

i.e c_1 heads and c_2 tails with $N = c_1 + c_2$

- What is the probability of c_1 heads and c_2 tails?
- How do we compute this?
- Can we use Bayes?

N here equals $c_1 + c_2$

Modelling typical coins from coin factory

- What we want is: $p(c|\alpha, \beta)$?
- How do we derive this: (*use law of total probability*)

$$\begin{aligned} p(c|\alpha, \beta) &= \int p(c, \theta|\alpha, \beta) d\theta \\ &= \int \underbrace{p(c|\theta)}_{\substack{\text{Binomial} \\ \text{Likelihood}}} \underbrace{p(\theta|\alpha, \beta)}_{\substack{\text{Beta} \\ \text{Prior}}} d\theta \end{aligned}$$

- Integration can be solved analytically (i.e. by hand) if the two distributions have similar form.
- In this case, we say that that two distributions are **conjugate**.

Modelling coin tosses

- Once, we are happy with the choice of α, β for our factory we can plug this in

$$p(c|\alpha, \beta) = \int p(c, \theta|\alpha, \beta) d\theta = \int p(c|\theta) p(\theta|\alpha, \beta) d\theta$$

$$= \int \frac{N!}{c_1! c_2!} \theta^{c_1} (1 - \theta)^{c_2} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta$$

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int \theta^{c_1} (1 - \theta)^{c_2} \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta$$

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int \theta^{c_1+\alpha-1} (1 - \theta)^{c_2+\beta-1} d\theta$$

*conjugacy
helps*

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int \theta^{c_1+\alpha-1} (1 - \theta)^{c_2+\beta-1} d\theta$$

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)}{\Gamma(c_1 + c_2 + \alpha + \beta)}$$

*from definition of
Beta function*

N here equals $c_1 + c_2$

The Beta-Binomial Distribution

- What we have just derived is the beta-binomial distribution that gives the (averaged) probability of drawing c_1 heads and c_2 tails from a coin that has been drawn from a beta distribution with parameters α and β .
- The 'averaged' above means that the coin parameter θ has been **integrated out** (and hence no longer appears in the equations):

$$p(c|\alpha, \beta) = \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)}{\Gamma(c_1 + c_2 + \alpha + \beta)}$$

- Using the definition of the **beta function**:

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a + b)}$$

$$p(c|\alpha, \beta) = \frac{N!}{c_1! c_2!} \frac{B(c_1 + \alpha, c_2 + \beta)}{B(\alpha, \beta)}$$

- $\frac{N!}{c_1! c_2!} \neq \frac{1}{B(C_1+1, C_2+1)}$ since $\frac{1}{B(C_1+1, C_2+1)} = \frac{\Gamma(C_1+C_2+2)}{\Gamma(C_1+1)\Gamma(C_2+1)} = \frac{(N+1)!}{C_1! C_2!}$

Inference

Inference

- The inference problem can be viewed as determining the parameters of your model from observations
- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair?

Inference

- The inference problem can be viewed as determining the parameters of your model from observations
- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair? Or what is the probability that it is fair? How do you even ask such a question?

Inference

- The inference problem can be viewed as determining the parameters of your model from observations
- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair? Or what is the probability that it is fair? How do you even ask such a question?
- The key to Bayesian inference is the mechanism to integrate our prior beliefs into the modelling process and provide mathematically grounded answers to the above questions.

Inference

- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair?

$$p(\mathbf{c} = (\mathbf{c}_1 = 10, \mathbf{c}_2 = 10) | \theta) = \binom{20}{10} \theta^{10} (1 - \theta)^{10}$$

- What is $p(\theta | \mathbf{c} = (\mathbf{c}_1 = 10, \mathbf{c}_2 = 10))$?
- Using the definition of conditional probability:

$$p(\theta, \mathbf{c}) = p(\theta | \mathbf{c}) p(\mathbf{c}) = p(\mathbf{c} | \theta) p(\theta)$$

- Using Bayes':

$$p(\mathbf{c}) = \int_0^1 p(\theta, \mathbf{c}) d\theta = \int_0^1 p(\mathbf{c} | \theta) p(\theta) d\theta$$

- Substituting, we get:

$$p(\theta | \mathbf{c}) = \frac{p(\mathbf{c} | \theta) p(\theta)}{\int_0^1 p(\mathbf{c} | \theta) p(\theta) d\theta}$$

Inference

- Thus to find out $p(\theta | c = (c_1 = 10, c_2 = 10))$ we can use:

$$p(\theta | c) = \frac{p(c | \theta) p(\theta)}{\int_0^1 p(c | \theta) p(\theta) d\theta}$$

- However, we do not know what $p(\theta)$ is?
- Hmmm... What might this mean?

Inference

- Thus to find out $p(\theta | c = (c_1 = 10, c_2 = 10))$ we can use:

$$p(\theta | c) = \frac{p(c | \theta) p(\theta)}{\int_0^1 p(c | \theta) p(\theta) d\theta}$$

- However, we do not know what $p(\theta)$ is?
- Hmm... What might this mean?
- $p(\theta)$ is our prior belief about the coin with bias θ
- What does this mean?

Inference

$$p(\theta|c) = \frac{p(c|\theta)p(\theta)}{\int p(c|\theta)p(\theta)d\theta}$$

Posterior:

- $p(\theta|c)$ is the *posterior distribution*
- Gives the probability of the parameter given data

Prior:

- $p(\theta)$ is the *prior distribution*
- Gives the probability for different values of θ and quantifies our belief regarding θ

Likelihood:

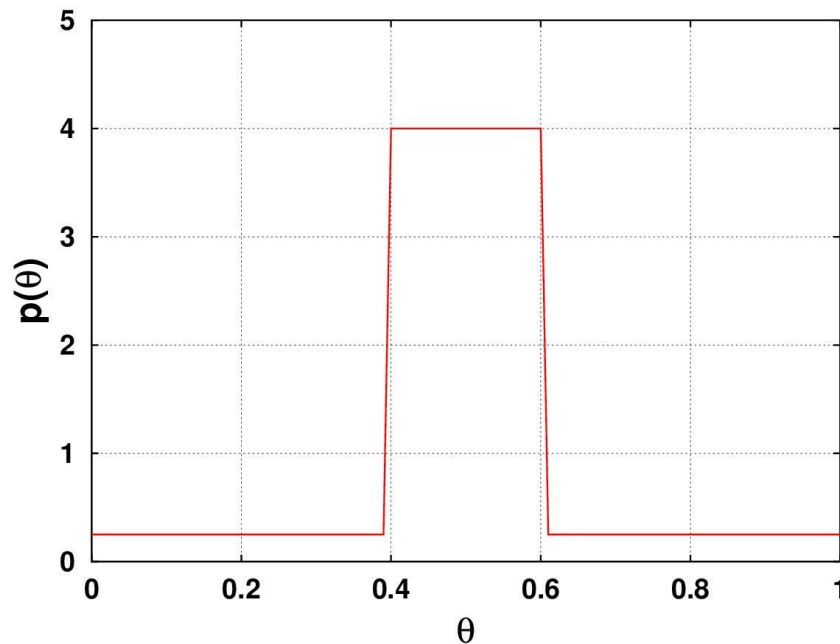
- $p(c|\theta)$ is the *likelihood* of the data given by the model
- Gives the probability of the data being generated by the model

Partition function/Normalising constant/Evidence:

- $\int_0^1 p(c|\theta)p(\theta)d\theta$ is the partition function/normalising constant/evidence.
- This is a constant needed to ensure that the probabilities sum to 1.

Inference

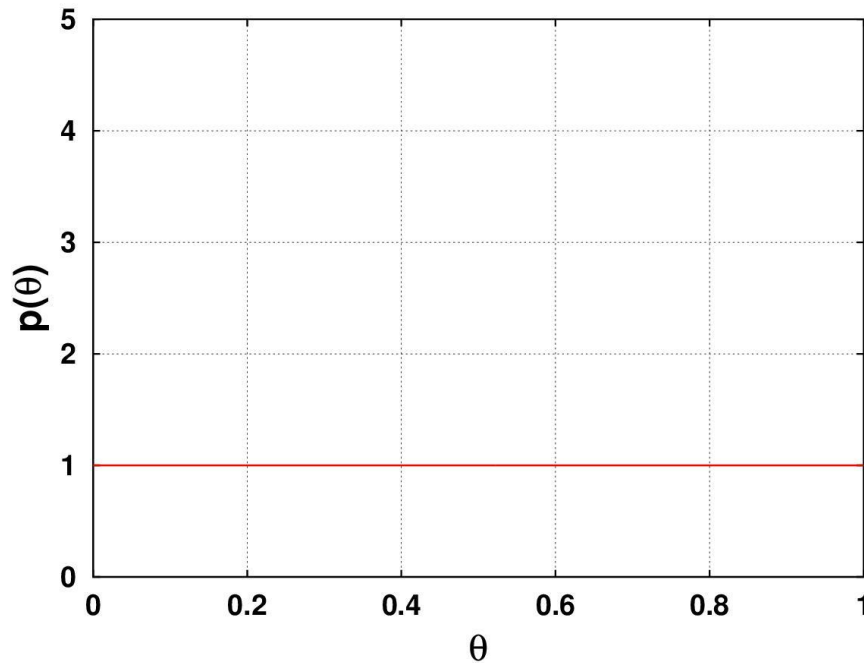
- $p(\theta)$ is our prior belief about the coin with bias θ
- Suppose, we say that, $p(\theta)$ is **piecewise uniform** with:
 - $F(0.4 \leq \theta \leq 0.6) = 0.8, F(0 \leq \theta < 0.4) = 0.1,$
 $F(0.6 < \theta \leq 1) = 0.1$



Inference

- Alternatively suppose that, $p(\theta)$ is **uniform** with:

$$p(\theta) = \begin{cases} 1 & \text{if } 0 < \theta < 1 \\ 0 & \text{otherwise} \end{cases}$$



Inference

- Alternatively still, we can assume that $p(\theta)$ is given by a **beta distribution** with parameters α, β

$$p(\theta|\alpha, \beta) = \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{B(\alpha, \beta)} = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1}(1-\theta)^{\beta-1}$$

- The uniform distribution can be recovered by $\alpha = 1, \beta = 1$

$$p(\theta|1, 1) = \frac{\theta^{1-1}(1-\theta)^{1-1}}{B(\alpha, \beta)} = \frac{\Gamma(1 + 1)}{\Gamma(1)\Gamma(1)} = 1$$

- The piecewise uniform distribution may also be approximated by suitable choices of α, β
- Choosing the **beta distribution** as a **prior** means that we can take advantage of **conjugacy**.

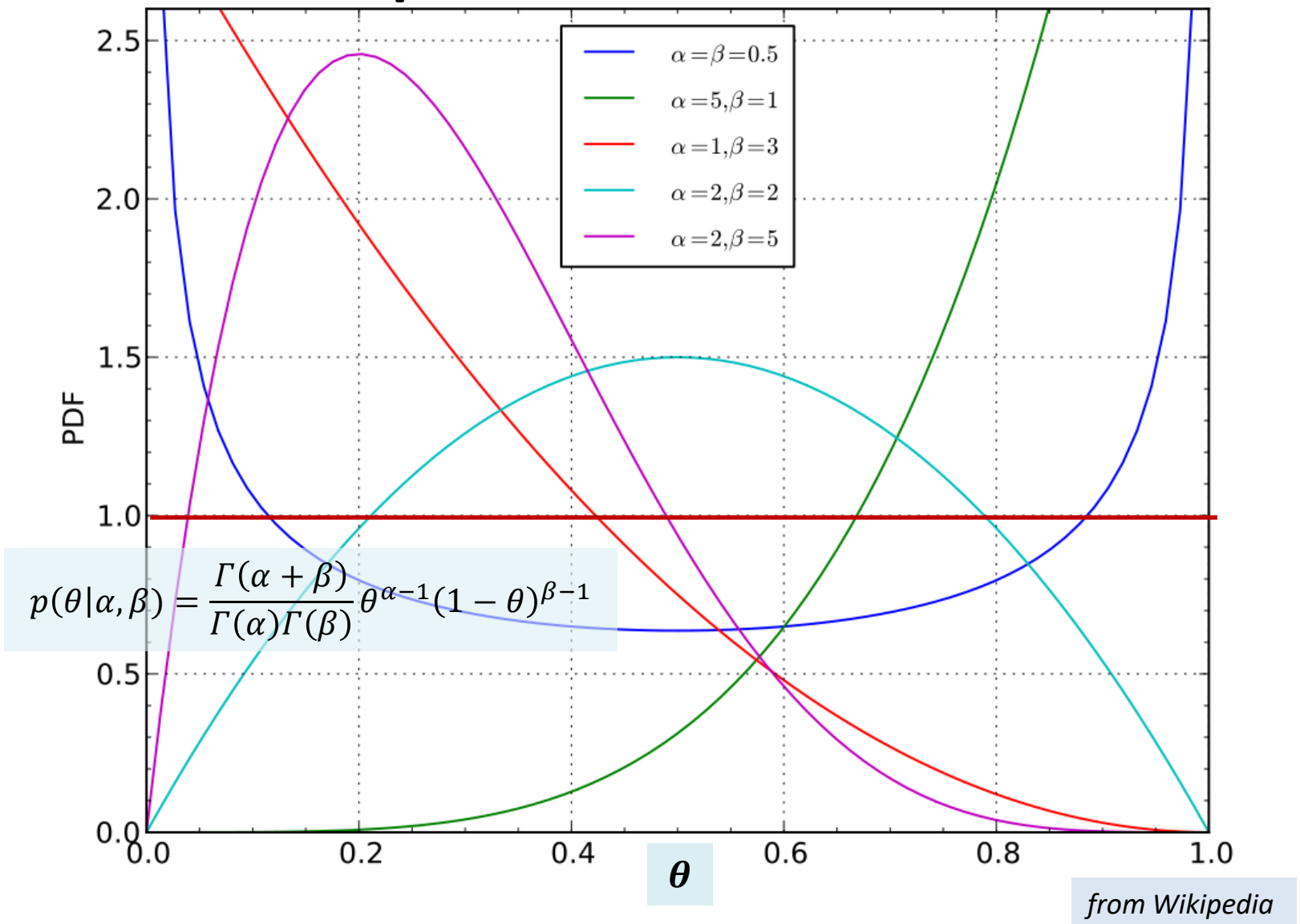
Inference: Binomial Posterior Distribution

- For the case when $p(\theta)$ is given by the **beta distribution** with parameters α, β :

$$\begin{aligned}
 p(\theta|c, \alpha, \beta) &= \frac{p(c|\theta, \alpha, \beta) p(\theta|\alpha, \beta)}{\int_0^1 p(c|\theta, \alpha, \beta) p(\theta|\alpha, \beta) d\theta} && \text{Remember } c = (c_1, c_2) \\
 &= \frac{\frac{N!}{c_1! c_2!} \theta^{c_1} (1 - \theta)^{c_2} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}}{\int_0^1 \frac{N!}{c_1! c_2!} \theta^{c_1} (1 - \theta)^{c_2} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta} \\
 &= \frac{\theta^{c_1 + \alpha - 1} (1 - \theta)^{c_2 + \beta - 1}}{\frac{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)}{\Gamma(c_1 + c_2 + \alpha + \beta)}} \\
 &= \frac{\Gamma(c_1 + c_2 + \alpha + \beta)}{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)} \theta^{c_1 + \alpha - 1} (1 - \theta)^{c_2 + \beta - 1} \\
 &= \frac{1}{B(c_1 + \alpha, c_2 + \beta)} \theta^{c_1 + \alpha - 1} (1 - \theta)^{c_2 + \beta - 1} \\
 &= p(\theta|\alpha + c_1, \beta + c_2) \sim \text{Beta}(\alpha + c_1, \beta + c_2)
 \end{aligned}$$

- Thus, posterior of binomial (under beta) is beta with **counts added into the parameters**

Shape of Beta distribution



Aside: Inference

- If θ is *piecewise uniform* with, $F(0.4 \leq \theta \leq 0.6) = 0.8$,
 $F(0 \leq \theta < 0.4) = 0.1$, $F(0.6 < \theta \leq 1) = 0.1$
- We can calculate, $p(\theta)$:

$$F(0.4 \leq \theta \leq 0.6) = 0.8 \text{ implies } p(\theta) = \frac{0.8}{0.2} = 4$$

$$F(0 \leq \theta < 0.4) = 0.1 \text{ implies } p(\theta) = \frac{0.1}{0.4} = \frac{1}{4}$$

$$F(0.6 < \theta \leq 1) = 0.1 \text{ implies } p(\theta) = \frac{0.1}{0.4} = \frac{1}{4}$$

Aside: Inference

$$\begin{aligned}
 p(\theta|c) &= \frac{p(c|\theta) p(\theta)}{\int_0^1 p(c|\theta) p(\theta) d\theta} \\
 &= \frac{p(c|\theta) p(\theta)}{\int_0^{0.4} p(c|\theta) \cdot \frac{1}{4} d\theta + \int_{0.4}^{0.6} p(c|\theta) \cdot 4 d\theta + \int_{0.6}^1 p(c|\theta) \cdot \frac{1}{4} d\theta} \\
 &= \frac{\theta^{c_1}(1-\theta)^{c_2} p(\theta)}{\frac{1}{4} IB(.4, c_1, c_2) + 4[IB(.6, c_1, c_2) - IB(.4, c_1, c_2)] + \frac{1}{4} [B(c_1, c_2) - IB(.6, c_1, c_2)]}
 \end{aligned}$$

- Where ***IB*** is the ***incomplete beta function*** defined by:

$$IB(a, c_1, c_2) = \int_0^a \theta^{c_1-1} (1-\theta)^{c_2-1} d\theta \quad 0 \leq a \leq 1$$

- Hence: ***IB***(1, c_1, c_2) = ***B***(c_1, c_2)
- Most math packages provide implementations of the incomplete beta function.
- Thus, posterior probability for any θ and c can be calculated.

Bayesian Inference

$$p(\theta|c) = \frac{p(c|\theta)p(\theta)}{\int p(c|\theta)p(\theta)d\theta}$$

- ***Being Bayesian*** typically means that we treat $p(\theta|c)$ as a distribution
- This means that we do not get a single value for the model parameter θ
- However, sometimes, we may want to find out what the '**best**' value for the model parameter θ would be.
- How do define what '**best**' means ?

MLE Inference

- This is the maximum likelihood estimate (where c is the data)

$$\theta_{MLE} = \arg \max_{\theta} p(c|\theta)$$

- Thus MLE is simply the maximum/**mode** of the model likelihood.
- We don't even need to define a prior, so (apparently) more “objective”.
- Typically deal with situations where there is only one mode, so can talk about **the** MLE.

MAP Inference

$$p(\theta|c) = \frac{p(c|\theta)p(\theta)}{\int p(c|\theta)p(\theta)d\theta}$$

MAP (Maximum a Posteriori):

- Assume that $p(\theta)$ is distributed as per **some given distribution**

$$\theta_{MAP} = \arg \max_{\theta} p(\theta|c) = \arg \max_{\theta} \frac{p(c|\theta)p(\theta)}{\int_0^1 p(c|\theta)p(\theta)d\theta}$$

- Since the denominator is a constant:

$$\theta_{MAP} = \arg \max_{\theta} p(c|\theta)p(\theta)$$

- Thus MAP estimate is the **mode** of the posterior distribution

Multinomial Distribution

- Probability of observing $\mathbf{c} = (c_1, \dots, c_k)$ heads in all possible ways out of $C = \sum_i c_i$ throws from a k-headed dice with probability of heads $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$ s.t. $\sum_i \theta_i = 1$

$$p(\text{heads} = \mathbf{c} | \boldsymbol{\theta}) = \frac{C!}{\prod_i c_i!} \prod_i \theta_i^{c_i}$$

- Points to note:
 - Generalises the binomial distribution to $k > 2$
 - Equivalent to the binomial for $k = 2$

The Dirichlet Distribution

- Dirichlet distribution generalises the Beta distribution to the $k - 1$ probability simplex

- The Beta distribution:

$$p(\theta|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \sim \text{Beta}(\alpha, \beta)$$

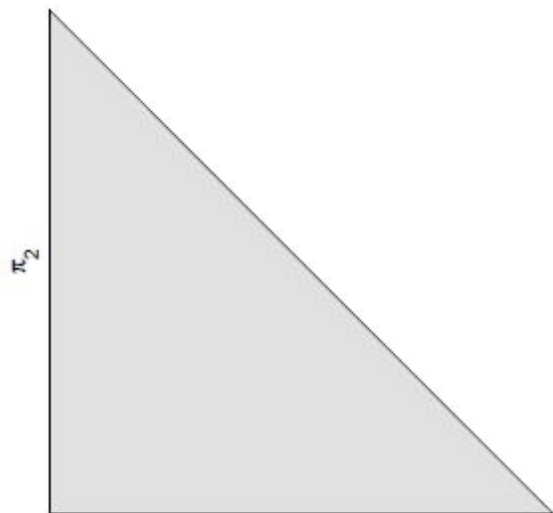
- The Dirichlet distribution:

$$p(\theta|\alpha) = \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i-1} \sim \text{Dir}(\alpha_1, \dots, \alpha_k)$$

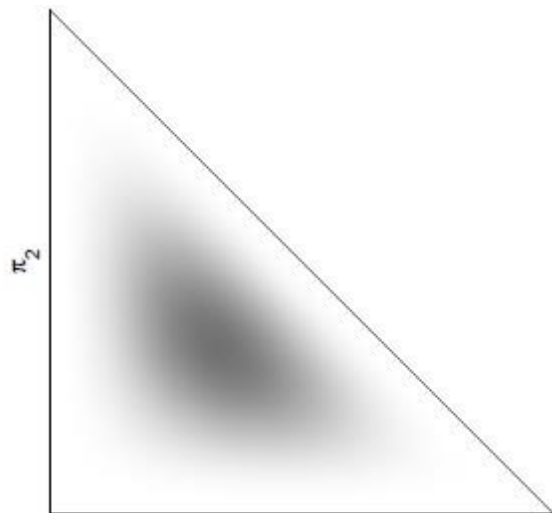
with:

$$\alpha = (\alpha_1, \dots, \alpha_k), \theta = (\theta_1, \dots, \theta_k), \sum_i \theta_i = 1, A = \sum_i \alpha_i$$

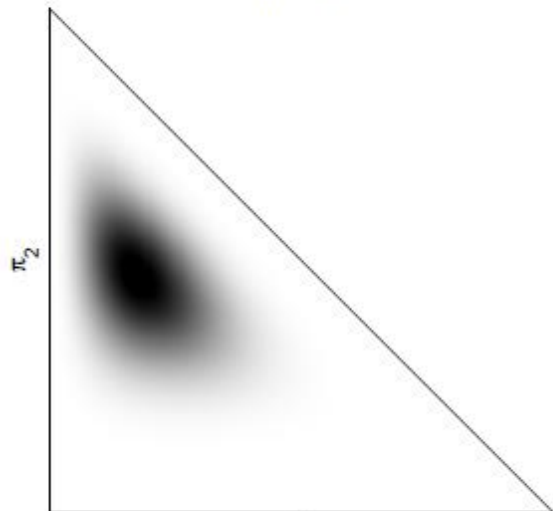
- So, the samples from the Dirichlet distribution can be used to model the bias in a k -sided dice.



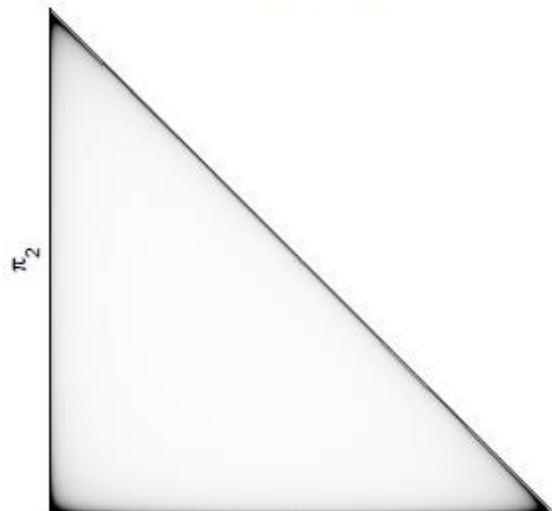
$$\pi \sim \text{Dir}(1, 1, 1)$$



$$\pi \sim \text{Dir}(4, 4, 4)$$



$$\pi \sim \text{Dir}(4, 9, 7)$$



$$\pi \sim \text{Dir}(0.2, 0.2, 0.2)$$

The Dirichlet-Multinomial

- Like for the Beta-Binomial distribution we can integrate out the Multinomial parameters
- Here, we are modelling the case where we need to predict the outcome $\mathbf{c} = (c_1, \dots, c_k)$ of a dice throw from a dice sampled from a factory with parameter $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_k)$

$$\begin{aligned}
 p(\mathbf{c}|\boldsymbol{\alpha}) &= \int p(\mathbf{c}, \boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta} = \int p(\mathbf{c}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta} \\
 &= \int \left[\frac{\mathbf{C}!}{\prod_i c_i!} \prod_i \theta^{c_i} \right] \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta^{\alpha_i-1} d\boldsymbol{\theta} \\
 &= \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \int \prod_i \theta^{c_i+\alpha_i-1} d\boldsymbol{\theta}
 \end{aligned}$$

- From earlier: $B(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$

- This generalises to: $\Delta(\boldsymbol{\alpha}) = \int_0^1 \prod_i \theta^{\alpha_i-1} d\boldsymbol{\theta} = \frac{\prod_i \Gamma(\alpha_i)}{\Gamma(A)}$

$$p(\mathbf{c}|\boldsymbol{\alpha}) = \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \int \prod_i \theta^{c_i+\alpha_i-1} d\boldsymbol{\theta} = \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{\prod_i \Gamma(c_i + \alpha_i)}{\Gamma(\mathbf{C} + A)}$$

Posterior Multinomial under Dirichlet prior

- Suppose we observe counts $\mathbf{c} = (c_1, \dots, c_k)$ from a dice sampled from our factory.
- We would like to predict the most likely parameters $\boldsymbol{\theta}$ for this dice.

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{c}, \boldsymbol{\alpha}) &= \frac{p(\mathbf{c}, \boldsymbol{\theta}|\boldsymbol{\alpha})}{p(\mathbf{c}|\boldsymbol{\alpha})} = \frac{p(\mathbf{c}, \boldsymbol{\theta}|\boldsymbol{\alpha})}{\int p(\mathbf{c}, \boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta}} = \frac{p(\mathbf{c}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\alpha})}{\int p(\mathbf{c}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta}} \\ &= \frac{\frac{\mathbf{C}!}{\prod_i c_i!} \prod_i \theta_i^{c_i} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i-1}}{\int \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{c_i+\alpha_i-1} d\boldsymbol{\theta}} = \frac{\prod_i \theta_i^{c_i+\alpha_i-1}}{\frac{\Gamma(\mathbf{C} + A)}{\prod_i \Gamma(c_i + \alpha_i)} \prod_i \theta_i^{c_i+\alpha_i-1}} \sim \text{Dir}(c_1 + \alpha_1, \dots, c_k + \alpha_k) \end{aligned}$$

- So, the shape of the posterior is exactly like that of the prior with counts added into the parameters

Posterior Dirichlet Categorical Conditional Distribution

- Using the Dirichlet Categorical we can compute the probability of observing class k having already observed counts $\mathbf{c} = (c_1, \dots, c_k)$.
- We will use an indicator variable $\mathbf{z} = k$ to indicate that the observed class is k :

$$\begin{aligned}
 p(\mathbf{z} = k | \mathbf{c}, \alpha) &= \frac{p(\mathbf{z} = k, \mathbf{c} | \alpha)}{p(\mathbf{c} | \alpha)} \\
 &= \frac{\frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{[\Gamma(c_k + \alpha_k + 1)] \prod_{i \neq k} \Gamma(c_i + \alpha_i)}{\Gamma(C + A)}}{\frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{\prod_i \Gamma(c_i + \alpha_i)}{\Gamma(C + A - 1)}} \\
 &= \frac{\Gamma(c_k + \alpha_k + 1) \Gamma(C + A - 1)}{\Gamma(c_k + \alpha_k) \Gamma(C + A)} \\
 &= \frac{c_k + \alpha_k}{C + A - 1}
 \end{aligned}$$

Not shown here –
The multinomial coefficients will cancel out. So, this distribution is the same for both Categorical and Multinomial

*Here: $C - 1 = \sum_i c_i$
So, the total number of items before adding the new \mathbf{z} is $C - 1$*

- This says that the probability of a new data point being assigned the class k is proportional to $c_k + \alpha_k$, having already observed \mathbf{c}
- Thus the Dirichlet exhibits the *rich-gets-richer* property

Finite Mixture Model

- Coin/Dice flips are generated from two (or more) coins
- So, the heads/tails/sides are possibly from different coins
- How can we model this?

Finite Mixture Model

- Coin/Dice flips are generated from two (or more) coins
- So, the heads/tails/sides are possibly from different coins
- How can we model this?
- In a finite mixture model we assume that the coin parameter θ is given by:

$$\theta = \pi_1 \theta_1 + \pi_2 \theta_2$$

where π_1, π_2 are the *mixture weights* or *proportions*.

- The mixture weights indicate how often each coin has been used. In general:

$$\theta = \sum_i \pi_i \theta_i$$

- If we assume, some prior for θ_i such as the Dirichlet then:

$$\theta = \sum_i \pi_i p(\theta_i | \alpha)$$

- This model will permit a single coin flip to be generated using multiple coins! So, it is overly general. We will restrict this in the following.

Latent variable models

- We assume that there is a hidden variable \mathbf{Z} or latent variable that says which mixture component/coin was used to generate the data.
- This makes sense for discrete data such as the coin example since a coin flip cannot be generated from more than one coin.
- Use of hidden variables considerably simplifies the problem formulation and inference
- Given $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$, we assume there is a corresponding latent variable $\mathbf{Z} = \{z_1, z_2, \dots, z_n\}$ that indicates the mixture component

$$p(x_i|\theta) = \sum_k p(x_i|\theta_k)p(z_i = k)$$
$$p(x_i|\theta) = \sum_k p(x_i|\theta_k)\pi_k$$

where π_k 's are the *mixture proportions* (or weights)

- Think of z_i 's like switches that indicate which coin was used.

Using Dirichlet priors

- To model π_k 's we can use Dirichlet priors

$$\begin{aligned} p(X, Z, \pi | \theta, \alpha) &= p(X, Z | \theta, \pi) p(\pi | \alpha) \\ &= p(X | Z, \theta) p(Z | \pi) p(\pi | \alpha) \end{aligned}$$

with $p(\pi | \alpha) \sim \text{Dir}\left(\frac{A}{K}, \dots, \frac{A}{K}\right)$ gives distribution over mixture weights/proportions

- Using the Dirichlet-Categorical distribution we can integrate out the mixture proportions π :

$$\begin{aligned} p(X, Z | \theta, \alpha) &= \int p(X, Z, \pi | \theta, \alpha) d\pi \\ &= p(X | Z, \theta) \int p(Z | \pi) p(\pi | \alpha) d\pi \\ &= p(X | Z, \theta) p(Z | \alpha) \\ &= p(Z | \alpha) \prod_i p(x_i | Z, \theta) \\ &= p(Z | \alpha) \prod_i p(x_i | \theta_{z_i}) \end{aligned}$$

$p(Z | \alpha)$ here would be the Dirichlet-Categorical distribution

Given Z the data points x_i 's are independent

Gibbs sampling (review)

- When we have multiple variables to sample e.g. $p(X_1, \dots, X_n)$
- Then the joint distribution can be estimated by successively estimating:
 - $p(X_1 | X_2, \dots, X_n)$
 - $p(X_i | X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$
- In the limit, this distribution converges to the joint $p(X_1, \dots, X_n)$
- More precisely:
$$p(X_i | X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n) = \frac{p(X_1, \dots, X_n)}{p(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)} \propto p(X_1, \dots, X_n)$$

Gibbs sampling for finite mixtures

- The mixture proportions are determined by \mathbf{Z} . In particular, the count c_k of the number of data points in class k is given by:

- $c_k = |\{z_i = k \mid z_i \in \mathbf{Z}\}|$

- We need to estimate $p(x_i, z_i \mid X_{-i}, Z_{-i}, \theta, \alpha)$:

$$p(x_i, z_i \mid X_{-i}, Z_{-i}, \theta, \alpha) p(X_{-i}, Z_{-i} \mid \theta, \alpha) = p(X, \mathbf{Z} \mid \theta, \alpha)$$

- Thus: $p(x_i, z_i \mid X_{-i}, Z_{-i}, \theta, \alpha) = \frac{p(X, \mathbf{Z} \mid \theta, \alpha)}{p(X_{-i}, Z_{-i} \mid \theta, \alpha)}$

$$= \frac{p(X \mid \mathbf{Z}, \theta)}{p(X_{-i} \mid Z_{-i}, \theta)} \frac{p(\mathbf{Z} \mid \alpha)}{p(Z_{-i} \mid \alpha)}$$

$$= p(x_i \mid z_i, \theta) \frac{p(\mathbf{Z} \mid \alpha)}{p(Z_{-i} \mid \alpha)}$$

$$= p(x_i \mid \theta_{z_i}) \frac{p(\mathbf{Z} \mid \alpha)}{p(Z_{-i} \mid \alpha)}$$

$$= p(x_i \mid \theta_{z_i}) p(z_i \mid Z_{-i}, \alpha)$$

Since:

$$\begin{aligned} p(z_i, Z_{-i} \mid \alpha) &= p(\mathbf{Z} \mid \alpha) \\ &= p(z_i \mid Z_{-i}, \alpha) p(Z_{-i} \mid \alpha) \end{aligned}$$

Posterior Dirichlet Categorical Conditional Distribution

- Using the Dirichlet Categorical we can compute the probability of observing class k having already observed counts $\mathbf{c} = (c_1, \dots, c_k)$:

$$\begin{aligned}
 p(\mathbf{z} = k | \mathbf{c}, \alpha) &= \frac{p(\mathbf{z} = k, \mathbf{c} | \alpha)}{p(\mathbf{c} | \alpha)} \\
 &= \frac{\frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{[\Gamma(c_k + \alpha_k + 1)] \prod_{i \neq k} \Gamma(c_i + \alpha_i)}{\Gamma(C + A)}}{\frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{\prod_i \Gamma(c_i + \alpha_i)}{\Gamma(C + A - 1)}} \\
 &= \frac{\Gamma(c_k + \alpha_k + 1) \Gamma(C + A - 1)}{\Gamma(c_k + \alpha_k) \Gamma(C + A)} \\
 &= \frac{c_k + \alpha_k}{C + A - 1}
 \end{aligned}$$

Not shown here –

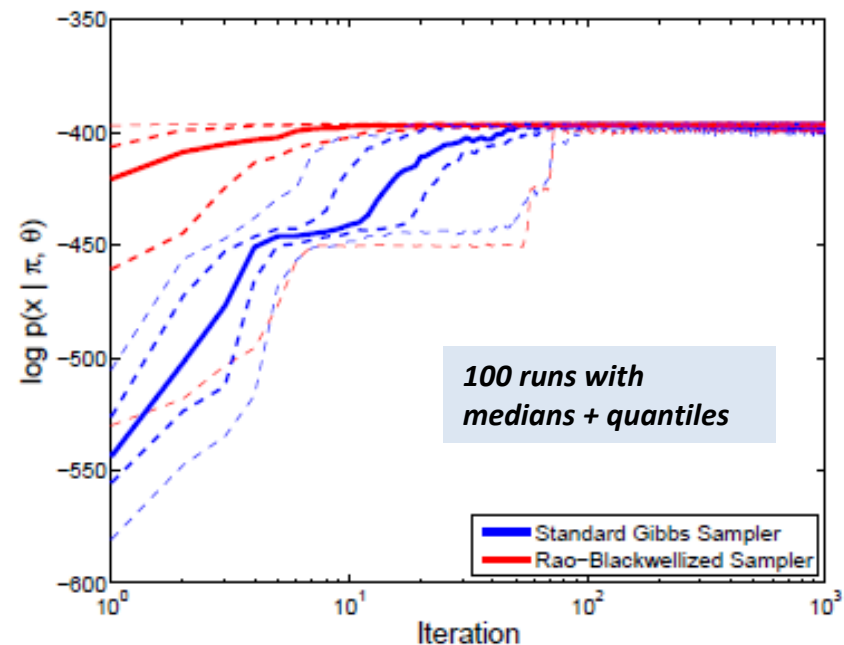
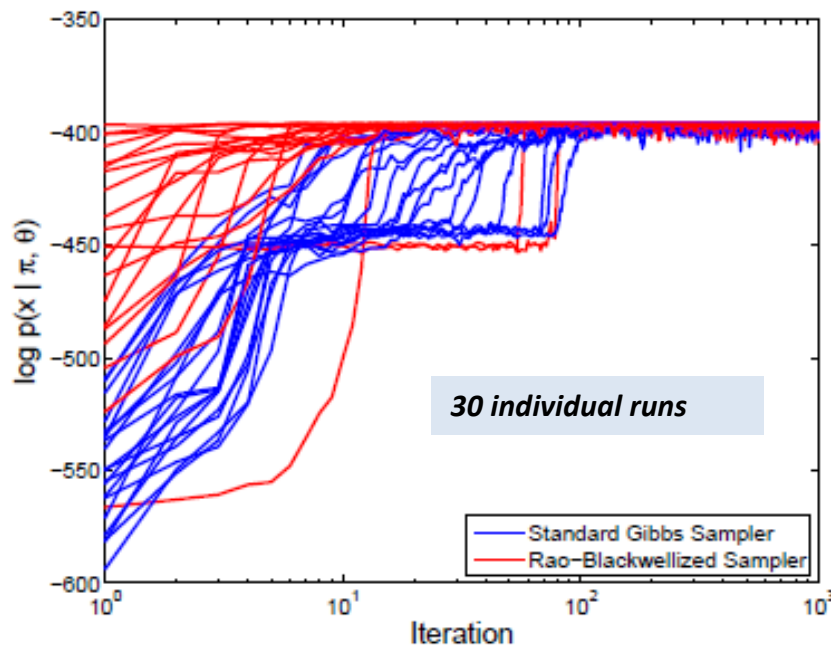
The multinomial coefficients will cancel out. So, this distribution is the same for both Categorical and Multinomial

Here: $C - 1 = \sum_i c_i$
So, the total number of items before adding the new \mathbf{z} is $C - 1$

- This says that the probability of a new data point being assigned the class k is proportional to $c_k + \alpha_k$, having already observed \mathbf{c}
- Thus the Dirichlet exhibits the *rich-gets-richer* property

Collapsed Gibbs Sampling

- Using the Dirichlet Multinomial to estimate the mixture proportions is a form of ***collapsed sampling*** or ***Rao-Blackwellized Gibbs sampling***
- Collapsed sampling results in Markov chains that converge much quicker and gets trapped in local optimal less often



From [Sudderth]

Finite Mixtures (θ known)

- In the simplest case, the (class density) parameters θ are known i.e. you know the bias of each coin
- So, the task is to determine the z_i 's i.e. to determine which coin was used for each x_i

Randomly assign $z_i = k$ with uniform probability

Repeat for each i :

Remove data point (x_i, z_i)

For each cluster k compute $p_k(x_i | \theta_k)$

Add data point (x_i, z_i) back by sampling $z_i = k$ using:

$$z_i = k \sim \frac{c_k + A/K}{N - 1 + A} p_k(x_i | \theta_k)$$

Mixture proportions can be estimated at the end

Finite Mixtures (θ unknown)

- In the general case, the (class density) parameters θ are unknown i.e. you do not know the bias of each coin
- So, the task is to determine both
 - the z_i 's i.e. to determine which coin was used for each x_i and
 - the class density parameters θ i.e. the bias for each coin

- In this case, the joint density we need to estimate:

$$\begin{aligned} p(X, Z, \theta | \alpha, \beta) &= p(X, Z | \theta, \alpha, \beta) p(\theta | \beta) \\ &= p(X | Z, \theta) p(Z | \alpha) p(\theta | \beta) \end{aligned}$$

- Following the derivation from page 34, gives:

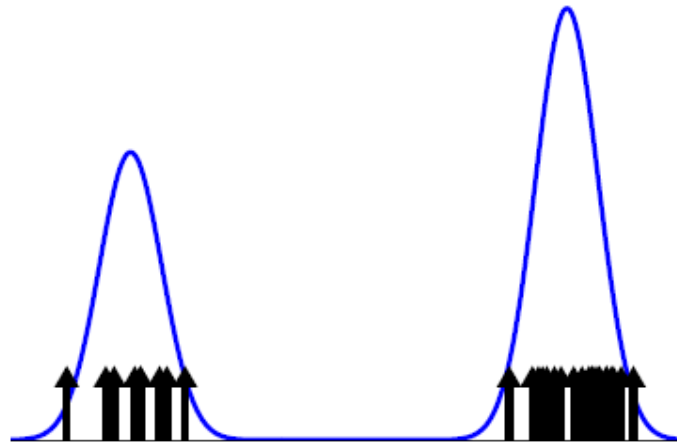
$$p(x_i, z_i | X_{-i}, Z_{-i}, \theta, \alpha) = p(x_i | \theta_{z_i}) p(z_i | Z_{-i}, \alpha)$$

- To sample θ we need to use the posterior distribution:

$$p(\theta_k | X, Z, \beta) = p(\theta_k | \{x_i | z_i = k \in Z\}, \beta_k)$$

- For the Multinomial with Dirichlet prior see page 27

Mixture of Univariate Gaussians



- Here the aim is to estimate the Gaussian parameters (μ , σ) and the mixture proportions using the finite mixture model

Mixture of Univariate Gaussians

- Randomly assign $\mathbf{z}_i = \mathbf{k}$ with uniform probability
- Repeat for each i :
 - Remove data point $(\mathbf{x}_i, \mathbf{z}_i)$
 - Estimate Gaussian parameters $\boldsymbol{\theta}_k = (\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k)$ for each cluster from current cluster assignments e.g. using MLE.

For MLE, we can ignore the prior $\boldsymbol{\beta}$

$$p(\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k | \mathbf{X}, \mathbf{Z}) = p(\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k | \{\mathbf{x}_i | \mathbf{z}_i = \mathbf{k} \in \mathbf{Z}\})$$

For each cluster \mathbf{k} compute $p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\sigma}_k)$

- Sample \mathbf{z}_i using:

$$\mathbf{z}_i = \mathbf{k} \sim \frac{c_k + A/K}{C - 1 + A} p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\sigma}_k)$$

- Mixture proportions can be estimated at the end

Mixture of Multinomials

- For each data point (x_i, z_i) randomly assign $z_i = k$ with uniform probability
- Repeat for each i :
 - Remove data point (x_i, z_i)
 - Sample Multinomial parameters θ_k for each cluster from current cluster assignments $X_k = \{x_i \in X | z_i = k\}$ using the posterior Multinomial distribution:

$$p(\theta|c, \beta) = \frac{\Gamma(C + B)}{\prod_i \Gamma(c_i + \beta_i)} \prod_i \theta_i^{c_i + \beta_i - 1}$$

- For each cluster k compute $p_k(x_i | \theta_k)$
- Add data point (x_i, z_i) back by sampling $z_i = k$ using:

$$z_i = k \sim \frac{c_k + \alpha_k}{C - 1 + A} p_k(x_i | \theta_k)$$

where $p_k(x_i | \theta_k)$ is the multinomial distribution for class k

- Mixture proportions can be estimated at the end

Collapsed Sampler for Mixture of Multinomials

- Like we did for the mixture proportions, the within-class multinomial parameters θ_k can be integrated out using the Dirichlet-Multinomial
- Thus instead of

$$p(x_i | \theta_{z_i})$$

we need to compute the factor:

$$\frac{p(X | Z, \beta)}{p(X_{-i} | Z_{-i}, \beta)}$$

- If we assume that for (x_i, z_i) the target destination class is k i.e. $z_i = k$:

$$\frac{p(X | Z, \beta)}{p(X_{-i} | Z_{-i}, \beta)} = \frac{\prod_k p(X_k | \beta)}{\prod_k p(X_{k,-i} | \beta)} = \frac{p(X_k | \beta)}{p(X_{k,-i} | \beta)}$$

- In the above, $X_k = \{x_i \in X | z_i = k\}$ (as discussed in previous slides)
- It is worth remembering that each x_i is a side of the dice
- Thus using the Dirichlet prior β in the Dirichlet-Multinomial distribution we can estimate the within-class likelihood $p(X_k | \beta)$

$$\begin{aligned}
\frac{p(X_k|\boldsymbol{\beta})}{p(X_{k,-i}|\boldsymbol{\beta})} &= \frac{\frac{\Gamma(B)}{\prod_j \Gamma(\beta_j)} \frac{[\Gamma(c_{x_i}^k + \beta_{x_i} + 1)] \prod_{j \neq x_i} \Gamma(c_j + \beta_j)}{\Gamma(C^j + B)}}{\frac{\Gamma(B)}{\prod_j \Gamma(\beta_j)} \frac{\prod_j \Gamma(c_j + \beta_j)}{\Gamma(C^j + B - 1)}} \\
&= \frac{\Gamma(c_{x_i}^k + \beta_{x_i} + 1)}{\Gamma(c_{x_i}^k + \beta_{x_i})} \frac{\Gamma(C^k + B - 1)}{\Gamma(C^k + B)} = \frac{c_{x_i}^k + \beta_{x_i}}{C^k + B - 1}
\end{aligned}$$

- j in the above, iterates over the total number of sides of the dice
- c_{x_i} is the count of the number of observations having the same side as x_i in class k

Fully collapsed sampler

- Randomly assign $\mathbf{z}_i = \mathbf{k}$ with uniform probability
- Repeat for each i :
 - Remove data point $(\mathbf{x}_i, \mathbf{z}_i)$
 - For each cluster \mathbf{k} compute $\frac{c_{\mathbf{x}_i}^{\mathbf{k}} + \beta_{\mathbf{x}_i}}{C^{\mathbf{k}} + B - 1}$
 - Add data point $(\mathbf{x}_i, \mathbf{z}_i)$ back by sampling $\mathbf{z}_i = \mathbf{k}$ using:
$$\mathbf{z}_i = \mathbf{k} \sim \frac{c_{\mathbf{k}} + A/K}{N - 1 + A} \left(\frac{c_{\mathbf{x}_i}^{\mathbf{k}} + \beta_{\mathbf{x}_i}}{C^{\mathbf{k}} + B - 1} \right)$$
- Mixture proportions can be estimated at the end

In the above $c_{\mathbf{x}_i}^{\mathbf{k}}$ is the count of the number of observations having the same side as \mathbf{x}_i in class \mathbf{k}

Topic models

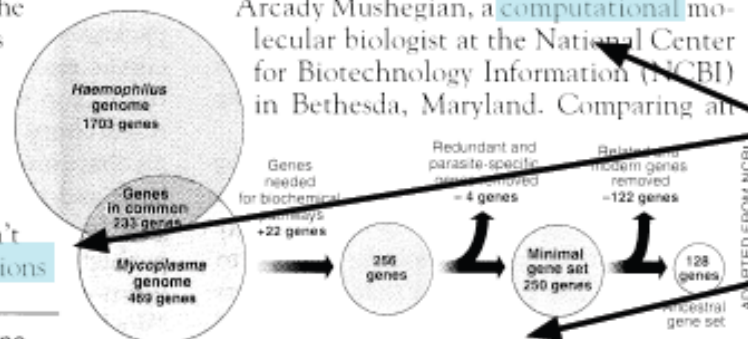
human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

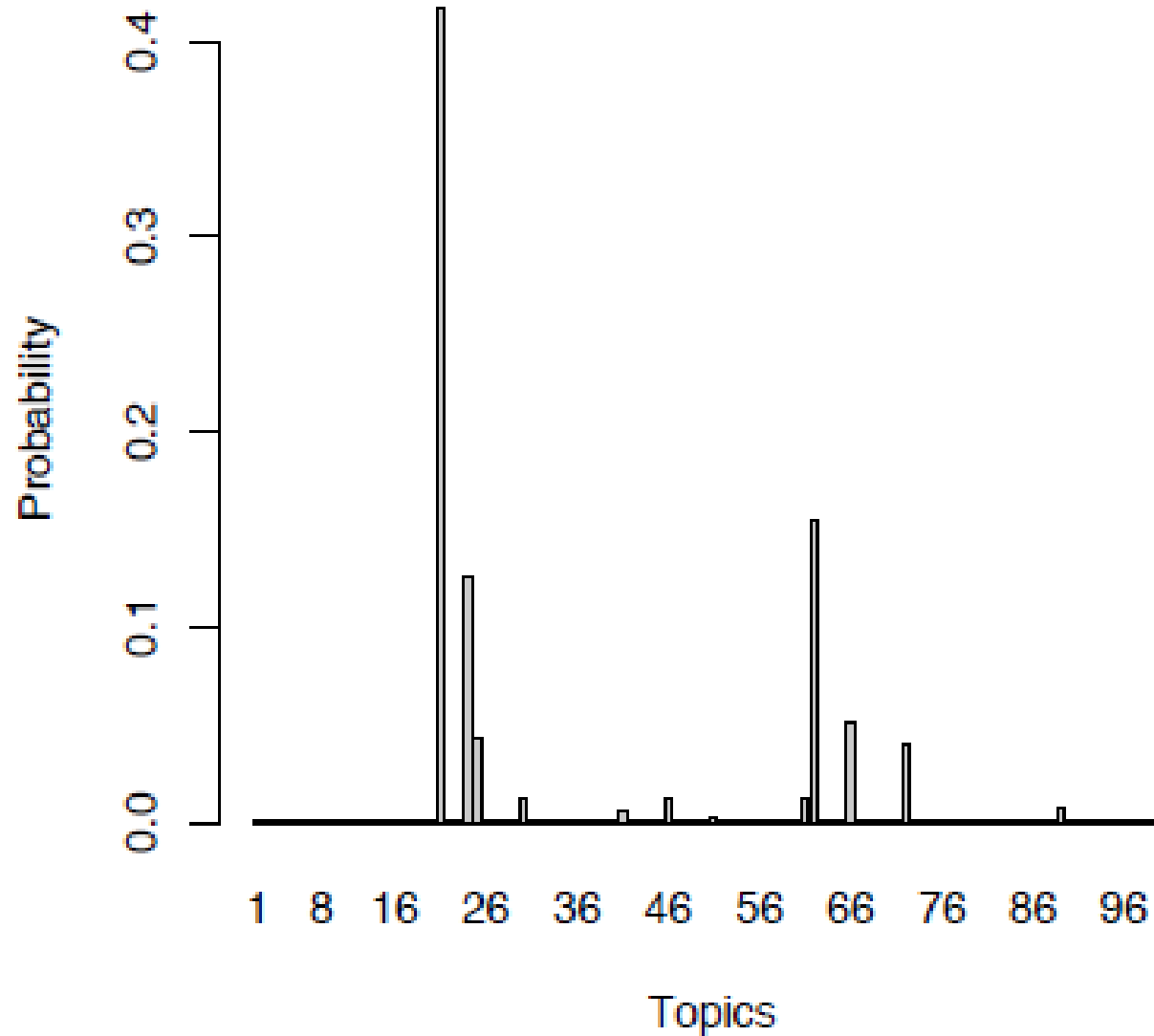
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

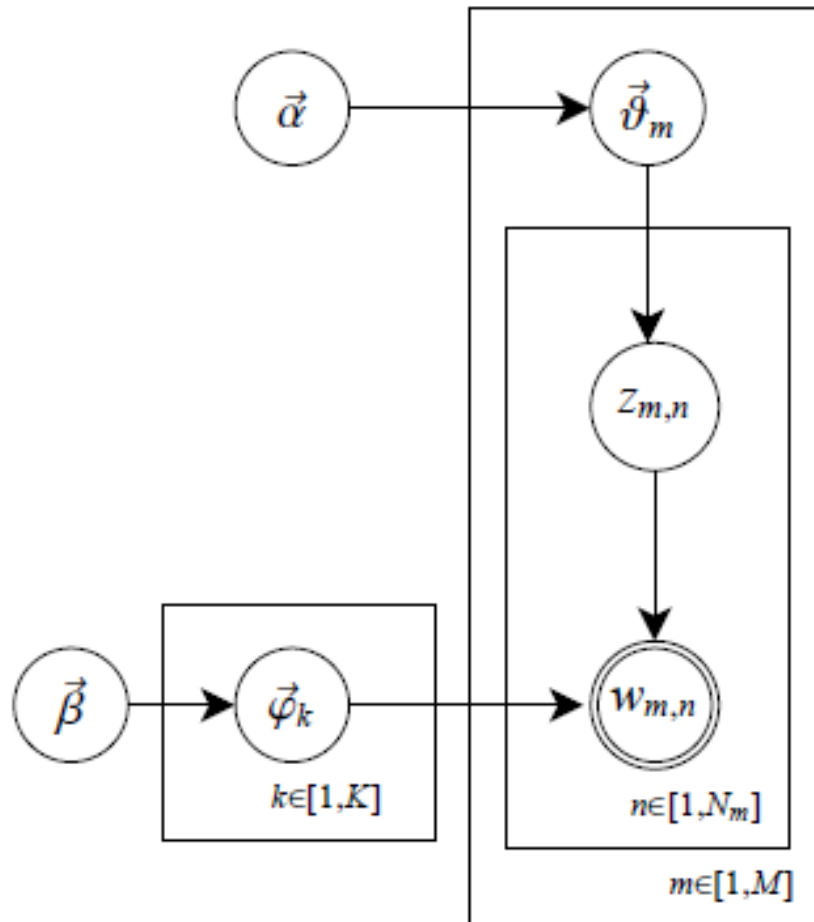
Document-topic distribution



Latent Dirichlet Allocation

- LDA is a hidden variable generative model for text documents
- The documents are modelled as bag-of-words or a Categorical distribution over words
- So, words are the sides of a dice and the number of times seen in a document is its count.
- However, each document is generated by a mixture of k Categorical distributions.
- But these distributions are shared across all documents
- So, all documents are generated by the same set of dices.
- The only difference between two documents are their mixture proportions

LDA plate diagram



- Use of Dirichlet Multinomials means that the multinomial parameters are integrated out.

Modeling the document-topic distribution

- Assume that we have K topics and D documents
- Let $\mathbf{c}_d = (c_{d,1}, \dots, c_{d,K})$ be a K -dimensional vector such that $c_{d,k}$ is the count of words with topic k in document d
- Let $C_d = \sum_k c_{d,k}$ is the count of words in document d
- Then, the *document-to-topic distribution* can be modelled using the Dirichlet Multinomial distribution:

$$p(\mathbf{c}_d | \boldsymbol{\alpha}) = \frac{\Gamma(A)}{\prod_k^K \Gamma(\alpha_k)} \frac{\prod_k^K \Gamma(c_{d,k} + \alpha_k)}{\Gamma(C_d + A)}$$

- Let $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_D)$ be a D -dimensional vector such that each $\mathbf{z}_i = (z_{i,1}, \dots, z_{i,|z_i|})$ is a vector of topic indicator variables for document i
- The topic assignments within each document are independent of the topic assignments in other documents.

$$p(\mathbf{Z} | \boldsymbol{\alpha}) = \prod_d^D \frac{\Gamma(A)}{\prod_k^K \Gamma(\alpha_k)} \frac{\prod_k^K \Gamma(c_{d,k} + \alpha_k)}{\Gamma(C_d + A)}$$

Modelling the topic-word distribution

- Let $\mathbf{n}_k = (n_{k,1}, \dots, n_{k,V})$ be a V -dimensional vector such that $n_{k,v}$ is the count of word \mathbf{v} in topic \mathbf{k} across the *whole* document collection
- Let $N_k = \sum_v n_{k,v}$ is the total number of words in topic \mathbf{k}
- Then we can use the Dirichlet Multinomial (again) to model the *topic-to-word* distribution:

$$p(\mathbf{n}_k | \mathbf{Z}, \boldsymbol{\beta}) = \frac{\Gamma(\mathbf{B})}{\prod_v \Gamma(\beta_v)} \frac{\prod_v \Gamma(n_{k,v} + \beta_v)}{\Gamma(N_k + \mathbf{B})}$$

- We need the \mathbf{Z} to collect all words belonging to a particular topic
- The topic-word distribution for one topic is independent of the same for another topic
- Let $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_D)$ be a D -dimensional vector such that each $\mathbf{w}_i = (\mathbf{w}_{i,1}, \dots, \mathbf{w}_{i,|z_i|})$ is a vector of words for document \mathbf{i} . Then:

$$p(\mathbf{W} | \mathbf{Z}, \boldsymbol{\beta}) = \prod_k \frac{\Gamma(\mathbf{B})}{\prod_v \Gamma(\beta_v)} \frac{\prod_v \Gamma(n_{k,v} + \beta_v)}{\Gamma(N_k + \mathbf{B})}$$

- Thus, the joint distribution becomes:

$$p(W, Z | \alpha, \beta) = p(W | Z, \beta) p(Z | \alpha)$$

$$= \left[\prod_k^K \frac{\Gamma(B)}{\prod_v^V \Gamma(\beta_v)} \frac{\prod_v^V \Gamma(n_{k,v} + \beta_v)}{\Gamma(N_k + B)} \right] \left[\prod_d^D \frac{\Gamma(A)}{\prod_k^K \Gamma(\alpha_k)} \frac{\prod_k^K \Gamma(c_{d,k} + \alpha_k)}{\Gamma(C_d + A)} \right]$$

- To do Gibbs sampling, we need to sample $z_i = k$ where $i = (m, n)$

$$p(w_i, z_i = k | W_{-i}, Z_{-i}, \alpha, \beta) = \frac{p(W, Z | \alpha, \beta)}{p(W_{-i}, Z_{-i} | \alpha, \beta)} = \frac{p(W | Z, z_i = k, \beta)}{p(W_{-i} | Z_{-i}, \beta)} \frac{p(Z, z_i = k | \alpha)}{p(Z_{-i} | \alpha)}$$

$$\frac{p(Z, z_i = k, i = (m, n) | \alpha)}{p(Z_{-i} | \alpha)} =$$

$$\frac{\left[\frac{\Gamma(A)}{\prod_j^K \Gamma(\alpha_j)} \frac{[\Gamma(c_{m,k} + \alpha_k + 1)] \prod_{j \neq k}^K \Gamma(c_{m,j} + \alpha_j)}{\Gamma(C_m + A)} \right] \prod_{d \neq m}^D \frac{\Gamma(A)}{\prod_j^K \Gamma(\alpha_j)} \frac{\prod_j^K \Gamma(c_{d,j} + \alpha_j)}{\Gamma(C_d + A)}}{\left[\frac{\Gamma(A)}{\prod_j^K \Gamma(\alpha_j)} \frac{[\Gamma(c_{m,k} + \alpha_k)] \prod_{j \neq k}^K \Gamma(c_{m,j} + \alpha_j)}{\Gamma(C_m + A - 1)} \right] \prod_{d \neq m}^D \frac{\Gamma(A)}{\prod_j^K \Gamma(\alpha_j)} \frac{\prod_j^K \Gamma(c_{d,j} + \alpha_j)}{\Gamma(C_d + A)}}$$

$$= \frac{\Gamma(c_{m,k} + \alpha_k + 1)}{\Gamma(C_m + A)} \frac{\Gamma(C_m + A - 1)}{\Gamma(c_{m,k} + \alpha_k)} = \frac{c_{m,k} + \alpha_k}{C_m + A - 1}$$

$$\begin{aligned}
& \frac{p(W, z_i = k, i = (m, n) | Z, \beta)}{p(W_{-i} | Z_{-i}, \beta)} = \\
& \frac{\left[\frac{\Gamma(B)}{\prod_v^V \Gamma(\beta_v)} \frac{[\Gamma(n_{k,i} + \beta_i + 1)] \prod_{v \neq i}^V \Gamma(n_{k,v} + \beta_v)}{\Gamma(N_k + B)} \right] \prod_{j \neq k}^K \frac{\Gamma(B)}{\prod_v^V \Gamma(\beta_v)} \frac{\prod_{v \neq i}^V \Gamma(n_{j,v} + \beta_v)}{\Gamma(N_j + B)}}{\left[\frac{\Gamma(B)}{\prod_v^V \Gamma(\beta_v)} \frac{[\Gamma(n_{k,i} + \beta_i)] \prod_{v \neq i}^V \Gamma(n_{k,v} + \beta_v)}{\Gamma(N_k + B - 1)} \right] \prod_{j \neq k}^K \frac{\Gamma(B)}{\prod_v^V \Gamma(\beta_v)} \frac{\prod_v^V \Gamma(n_{j,v} + \beta_v)}{\Gamma(N_j + B)}} \\
& = \frac{\Gamma(n_{k,i} + \beta_i + 1) \Gamma(N_k + B - 1)}{\Gamma(n_{k,i} + \beta_i) \Gamma(N_k + B)} \\
& = \frac{n_{k,i} + \beta_i}{N_k + B - 1}
\end{aligned}$$

▪ **Thus:**

$$p(w_i, z_i = k, i = (m, n) | W_{-i}, Z_{-i}, \alpha, \beta) = \frac{c_{m,k} + \alpha_k}{C_m + A - 1} \frac{n_{k,i} + \beta_i}{N_k + B - 1}$$

Note: Here $v \neq i$ means that the v^{th} word in the vocabulary is not the same as the i^{th} word in the document collection. The i^{th} word in the document collection being the same as the n^{th} word in the m^{th} document.

LDA algorithm sketch

- Randomly assign $\mathbf{z}_i = \mathbf{k}$ where $\mathbf{i} = (\mathbf{m}, \mathbf{n})$ with uniform probability. Remember \mathbf{i} is also a pointer to a word
- Update variables $\mathbf{c}_m, \mathbf{n}_k, \mathbf{C}_m, \mathbf{N}_k$ appropriately
- \mathbf{c} is $\mathbf{D} \times \mathbf{K}$ int matrix; \mathbf{n} is $\mathbf{K} \times \mathbf{V}$ int matrix
- \mathbf{C} is \mathbf{D} -dim vector of int; \mathbf{N} is \mathbf{K} -dim vector of int
- Repeat for each $\mathbf{i} = (\mathbf{m}, \mathbf{n})$:
 - Remove data point $(\mathbf{x}_i, \mathbf{z}_i)$
 - Sample $\mathbf{z}_i = \mathbf{k}$ using:

$$p(\mathbf{w}_i, \mathbf{z}_i = \mathbf{k} | \mathbf{W}, \mathbf{Z}_{-i}, \alpha, \beta) = \frac{\mathbf{c}_{m,k} + \alpha_k}{\mathbf{C}_m + \mathbf{A} - \mathbf{1}} \frac{\mathbf{n}_{k,i} + \beta_i}{\mathbf{N}_k + \mathbf{B} - \mathbf{1}}$$

- Update variables $\mathbf{c}_m, \mathbf{n}_k, \mathbf{C}_m, \mathbf{N}_k$ appropriately
- Mixture proportions can be estimated at the end

Infinite Mixture Models

Motivation

- Key problem with finite mixtures is that the number of mixture components (**K**) needs to be given in advance
- For large datasets/problems, it is often impossible to figure out what an optimal value for **K** should be
- The Dirichlet Process (**DP**) is an extension of the Dirichlet distribution to the infinite case

Dirichlet Categorical (revision)

- Like for the Beta-Binomial distribution we can integrate out the Categorical parameters
- Here, we are modelling the case where we need to predict the outcome $\mathbf{c} = (c_1, \dots, c_K)$ of a dice throw from a dice sampled from a factory with Dirichlet parameter $\alpha = (\alpha_1, \dots, \alpha_K)$

$$\begin{aligned}
 p(\mathbf{c}|\alpha) &= \int p(\mathbf{c}, \boldsymbol{\theta}|\alpha) d\boldsymbol{\theta} = \int p(\mathbf{c}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\alpha) d\boldsymbol{\theta} \\
 &= \int \left[\prod_i \theta_i^{c_i} \right] \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i-1} d\boldsymbol{\theta} \\
 &= \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \int \prod_i \theta_i^{c_i+\alpha_i-1} d\boldsymbol{\theta}
 \end{aligned}$$

- From earlier: $B(\alpha, \beta) = \int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$

- This generalises to: $\Delta(\alpha) = \int_0^1 \prod_i \theta_i^{\alpha_i-1} d\boldsymbol{\theta} = \frac{\prod_i \Gamma(\alpha_i)}{\Gamma(A)}$

$$p(\mathbf{c}|\alpha) = \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \int \prod_i \theta_i^{c_i+\alpha_i-1} d\boldsymbol{\theta} = \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{\prod_i \Gamma(c_i + \alpha_i)}{\Gamma(C + A)}$$

Posterior Conditional Distribution (revision)

- Having observed counts $\mathbf{c} = (c_1, c_2, \dots, c_K)$, the probability of observing class k is given by the posterior conditional distribution
- Using the Dirichlet Categorical we can compute the probability of observing class k having already observed \mathbf{c}

$$\begin{aligned}
 p(\mathbf{z} = k | \mathbf{c}, \boldsymbol{\alpha}) &= \frac{p(\mathbf{z} = k, \mathbf{c} | \boldsymbol{\alpha})}{p(\mathbf{c} | \boldsymbol{\alpha})} \\
 &= \frac{\frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{[\Gamma(c_k + \alpha_k + 1)] \prod_{i \neq k} \Gamma(c_i + \alpha_i)}{\Gamma(\mathbf{C} + A)}}{\frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{\prod_i \Gamma(c_i + \alpha_i)}{\Gamma(\mathbf{C} + A - 1)}} \\
 &= \frac{\Gamma(c_k + \alpha_k + 1) \Gamma(\mathbf{C} + A - 1)}{\Gamma(c_k + \alpha_k - 1) \Gamma(\mathbf{C} + A)} \\
 &= \frac{c_k + \alpha_k}{\mathbf{C} + A - 1}
 \end{aligned}$$

Here: $\mathbf{C} - 1 = \sum_i c_i$
 So, the total number of items before adding the new \mathbf{z} is $\mathbf{C} - 1$

- This says that the probability of a new data point being assigned the class k is proportional to $c_k + \alpha_k$, having already observed \mathbf{c}
- Thus the Dirichlet exhibits the *rich-gets-richer* property

The Chinese Restaurant Process

- Firstly, let's assume that we will use a symmetric Dirichlet prior i.e. assume that $\alpha = \left(\frac{A}{K}, \dots, \frac{A}{K}\right)$ i.e. each $\alpha_i = \frac{A}{K}$
- For the Dirichlet Categorical, the posterior conditional of observing class k having already observed counts c is:

$$p(z = k | c, \alpha) = \frac{c_k + \alpha_k}{C + A - 1}$$

- Now if we let $\alpha_k = \frac{A}{K}$ and re-parameterise (using A) then:

$$p(z = k | c, A) = \frac{c_k + \frac{A}{K}}{C + A - 1}$$

- Now we pose the question: ***What happens when $K \rightarrow \infty$?***

$$\lim_{K \rightarrow \infty} p(z = k | c, A) = \lim_{K \rightarrow \infty} \frac{c_k + \frac{A}{K}}{C + A - 1} = \frac{c_k}{C + A - 1}$$

- So:

$$p(\mathbf{z} = \mathbf{k} | \mathbf{c}, A) = \frac{c_k}{C + A - 1}$$

for a class \mathbf{k} that is already present i.e. $c_k > 0$

- Since, this is a distribution:

$$\sum_i^K p(\mathbf{z} = \mathbf{k} | \mathbf{c}, A) = 1$$

- But we find that:

$$\sum_i^K p(\mathbf{z} = \mathbf{k} | \mathbf{c}, A) = \sum_i^K \frac{c_i}{C + A - 1} = \frac{C - 1}{C + A - 1} < 1$$

$$1 - \frac{C - 1}{C + A - 1} = \frac{A}{C + A - 1}$$

- This remaining probability mass is for the remaining unseen classes.

- The **Chinese Restaurant Process (CRP)** is defined by the conditional distribution:

$$p(z = k | c, A) = \frac{c_k}{c + A - 1} \quad \text{if } k \text{ is an existing class}$$

$$p(z = K + 1 | c, A) = \frac{A}{c + A - 1} \quad \text{for a new class } K + 1$$

- More conventionally, it is written as:

$$p(z = k | c, \alpha) = \frac{c_k}{c + \alpha - 1} \quad \text{if } k \text{ is an existing class}$$

$$p(z = K + 1 | c, \alpha) = \frac{\alpha}{c + \alpha - 1} \quad \text{for new } K + 1$$

Here: α is no longer a vector but just a number.

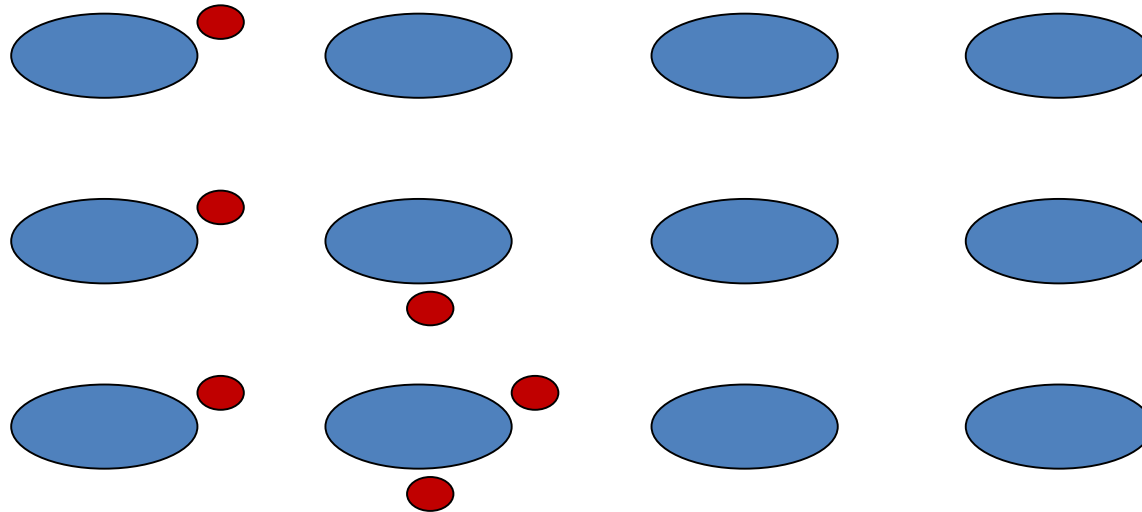
- So, the CRP is defined as:

$$p(\mathbf{z} = \mathbf{k} | \mathbf{c}, \alpha) = \mathbf{CRP}(\mathbf{z} = \mathbf{k} | \mathbf{c}, \alpha) = \frac{c_k}{c + \alpha - 1} \text{ for existing class } \mathbf{k}$$

$$p(\mathbf{z} = \mathbf{K} + \mathbf{1} | \mathbf{c}, \alpha) = \mathbf{CRP}(\mathbf{z} = \mathbf{k} + \mathbf{1} | \mathbf{c}, \alpha) = \frac{\alpha}{c + \alpha - 1} \text{ for new class}$$

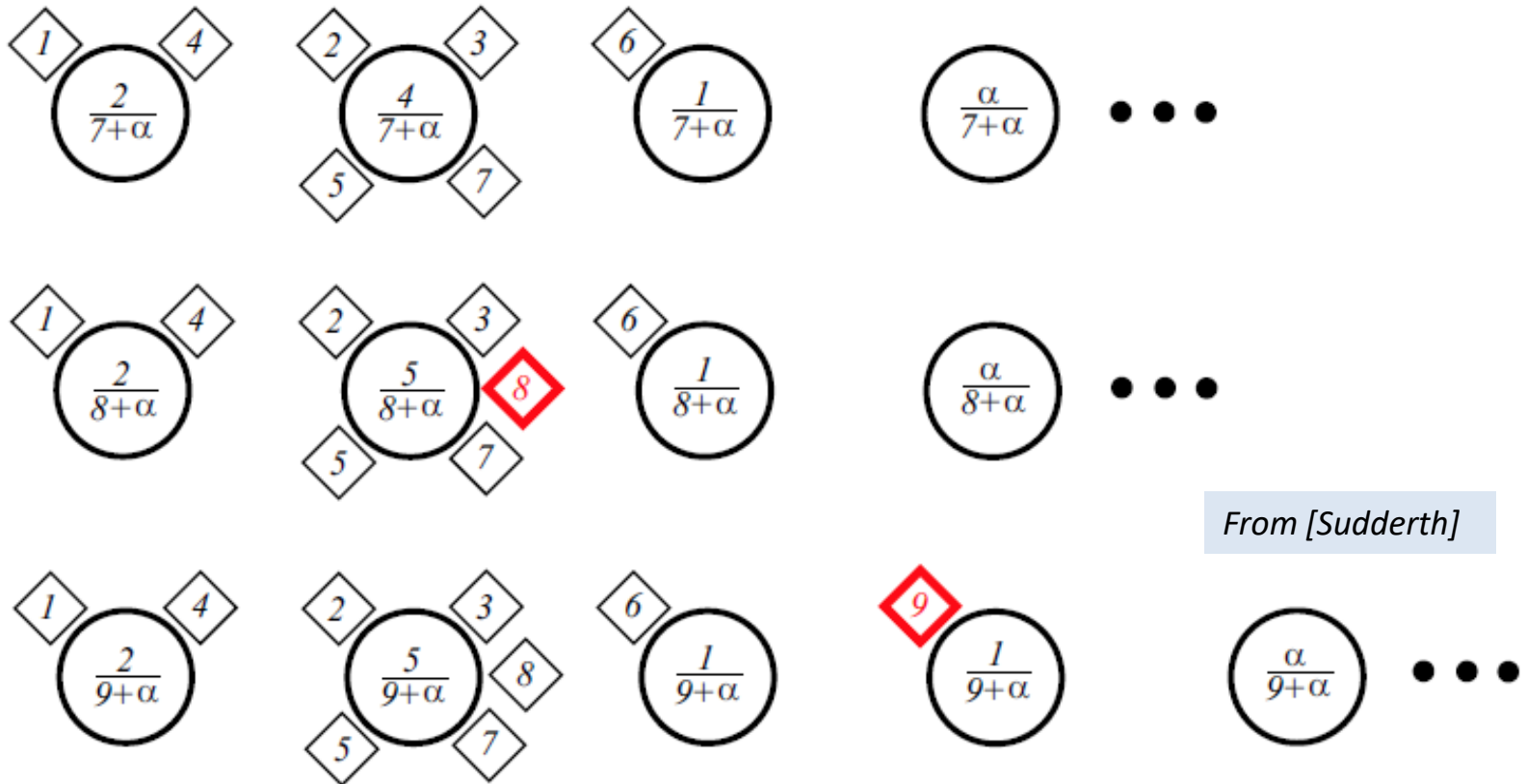
- α is the **concentration parameter** that controls how the mass is assigned to existing clusters
 - Smaller values of α concentrate the probability mass to **existing** clusters
 - Larger values of α concentrate the probability mass to **new** clusters
- The CRP is the infinite extension of the Dirichlet Categorical posterior conditional/predictive distribution.
 - Since, the mixture proportions $\boldsymbol{\pi}_k$ have been integrated out, the CRP is a distribution over counts.

The Chinese Restaurant



- The *first* customer occupies the first table with prob. $\frac{\alpha}{1+\alpha-1} = 1$
- The *second* customer either occupies the first table with prob. $\frac{1}{2+\alpha-1} = \frac{1}{1+\alpha}$ or can start a new table with prob. $\frac{\alpha}{2+\alpha-1} = \frac{\alpha}{1+\alpha}$
- The N^{th} customer occupies a table k with prob. $\frac{c_k}{N+\alpha-1}$ or can start a new table with prob. $\frac{\alpha}{N+\alpha-1}$

Another illustration of CRP



- Middle row : 8th customer joins the most popular table
- Last row : 9th customer starts a new table

The Polya urn scheme

An equivalent representation can be obtained by considering the following scheme

- An urn contains different coloured balls
- With probability $\frac{c-1}{c+\alpha-1}$:
 - Pick a ball from the urn (with uniform probability)
 - Put back the ball into the urn with an **extra** ball of the same colour
- With probability $\frac{\alpha}{c+\alpha-1}$:
 - Draw a new ball of different colour from somewhere else and add this into the urn

The Polya urn scheme is equivalent to the Chinese Restaurant metaphor.

The rich-gets-richer property is clearly evident.

Polya urn representation

- Equivalently, the CRP can be defined as:

$$p(\mathbf{z} = \mathbf{k} | \mathbf{Z}, \alpha) = CRP(\mathbf{z} = \mathbf{k} | \mathbf{c}, \alpha) = \sum_i \frac{I(\mathbf{z}_i = \mathbf{k})}{C + \alpha - 1} \quad \text{for}$$

existing class \mathbf{k}

$$p(\mathbf{z} = \mathbf{K} + \mathbf{1} | \mathbf{Z}, \alpha) = CRP(\mathbf{z} = \mathbf{K} + \mathbf{1} | \mathbf{c}, \alpha) = \frac{\alpha}{C + \alpha - 1}$$

for new class

Infinite Mixtures

- The infinite mixture model can be thought of as the extension of the finite mixture model to the infinite case

$$p(x_i|\theta) = \sum_{k=1}^{\infty} p(x_i|\theta_k)p(z_i = k)$$
$$p(x_i|\theta) = \sum_{k=1}^{\infty} p(x_i|\theta_k)\pi_k$$

- In practice, the ‘infinite’ here means that it is not fixed in advance.
- So, the number of clusters can grow or shrink.

Finite mixtures : Integrating out π_k (review)

- Using the Dirichlet-Categorical distribution we can integrate out the mixture proportions

$$\begin{aligned} p(X, Z | \theta, \alpha) &= \int p(X, Z, \pi | \theta, \alpha) d\pi \\ &= p(X | Z, \theta) \int p(Z | \pi) p(\pi | \alpha) d\pi \\ &= p(X | Z, \theta) p(Z | \alpha) \\ &= p(Z | \alpha) \prod_i p(x_i | Z, \theta) \\ &= p(Z | \alpha) \prod_i p(x_i | \theta_{z_i}) \end{aligned}$$

$$\begin{aligned} \text{So: } p(x_i, z_i | Z_{-i}, X_i, \theta, \alpha) &= \frac{p(X, Z | \theta, \alpha)}{p(X_{-i}, Z_{-i} | \theta, \alpha)} \\ &= \frac{p(X | Z, \theta)}{p(X_{-i} | Z_{-i}, \theta)} \frac{p(Z | \alpha)}{p(Z_{-i} | \alpha)} = p(x_i | \theta_{z_i}) p(z_i | Z_{-i}, \alpha) \end{aligned}$$

Extension to DP mixtures

- The above equations will remain the same for infinite mixtures.
- The only difference will be that to compute $p(\mathbf{z}_i | \mathbf{Z}_{-i}, \alpha)$ we will use CRP instead of the posterior conditional Dirichlet categorical
- The CRP is just a special case of the posterior conditional distribution of the Dirichlet categorical

- So:

$$p(x_i, z_i = k | \mathbf{Z}_{-i}, X_{-i}, \theta, \alpha) = p(x_i | \theta_{z_i}) p(z_i = k | \mathbf{Z}_{-i}, \alpha)$$

- Since, \mathbf{Z} , is just a count vector from the perspective of the Dirichlet Categorical

$$\begin{aligned} &= p(x_i | \theta_{z_i}) p(z_i = k | \mathbf{c}, \alpha) \\ &= p(x_i | \theta_{z_i}) \text{CRP}(z_i = k | \mathbf{c}, \alpha) \end{aligned}$$

- Thus:

$$p(x_i, z_i = k | \mathbf{Z}_{-i}, X_{-i}, \theta, \alpha) = p(x_i | \theta_{z_i}) \begin{cases} \frac{c_k}{C + \alpha - 1} & \text{if } k \text{ is an existing class} \\ \frac{\alpha}{C + \alpha - 1} & \text{if } k \text{ is a new class} \end{cases}$$

- If θ is **unknown** we sample θ_k using the posterior distribution:

$$p(\theta_k | X, Z, \beta) = p(\theta_k | \{x_i | z_i = k \in Z\}, \beta_k)$$

and then sample:

$$p(x_i, z_i = k | Z_{-i}, X_{-i}, \beta, \alpha) = p(x_i | \theta_k) \begin{cases} \frac{c_k}{C + \alpha - 1} & \text{if } k \text{ is an existing class} \\ \frac{\alpha}{C + \alpha - 1} & \text{if } k \text{ is a new class} \end{cases}$$

- **Alternatively**, we can **integrate out** θ_k if β_k is a *conjugate prior*
- This will give us a fully collapsed sampler

$$\begin{aligned} p(x_i | X_{-i}, z_i = k, \beta) &= p(x_i | \{x_j | z_j = k \in Z\}, \beta_k) \\ &= \int p(\theta_k | \{x_i | z_i = k \in Z\}, \beta_k) p(x_i | \theta_k) d\theta_k \end{aligned}$$

- In this case, $z_i = k$ can be sampled using:

$$p(x_i, z_i = k | Z_{-i}, X_{-i}, \beta, \alpha) =$$

$$p(x_i | X_{-i}, z_i = k) \begin{cases} \frac{c_k}{C + \alpha - 1} & \text{if } k \text{ is an existing class} \\ \frac{\alpha}{C + \alpha - 1} & \text{if } k \text{ is a new class} \end{cases}$$

Gibbs sampler for infinite mixtures

- Randomly assign $\mathbf{z}_i = \mathbf{k}$ for some max. K with uniform probability
- Repeat for each i :
 - Remove data point $(\mathbf{x}_i, \mathbf{z}_i)$
 - For each cluster \mathbf{k} compute $p_{\mathbf{k}}(\mathbf{x}_i | \boldsymbol{\theta}_{\mathbf{k}})$ if $\boldsymbol{\theta}_{\mathbf{k}}$ known or can be sampled; alternatively if conjugate, integrate out $\boldsymbol{\theta}_{\mathbf{k}}$ and compute $p_{\mathbf{k}}(\mathbf{x}_i | \{\mathbf{x}_i \in X | \mathbf{z}_i = \mathbf{k}\})$
 - Add data point $(\mathbf{x}_i, \mathbf{z}_i)$ back by sampling using:
$$\mathbf{z}_i = \mathbf{k} \sim \text{CRP}(\mathbf{z}_i = \mathbf{k} | \mathbf{Z}_{-i}, \alpha) p_{\mathbf{k}}(\mathbf{x}_i | \boldsymbol{\theta}_{\mathbf{k}})$$
- Mixture proportions can be estimated at the end

Discussion

- What we have derived is a DP mixture, i.e., an infinite mixture model
- But we have not actually explained what a Dirichlet process is, since we have only been using the posterior conditional distribution of the DP categorical
- This is like talking about the Dirichlet Categorical distribution while only using the equations for the posterior conditional of the Dirichlet categorical
- So what does the DP look like?
- And can we generate samples from the DP directly?
- A different question is: Is the DP exchangeable? Does it matter if we swap the order in which we sample data points in a different order?
- We will try to answer these.

Exchangeability

Definition: A joint distribution $p(x_1, \dots, x_n)$ is **exchangeable** if for any permutation $\pi(x_1, \dots, x_n) = (x_{\pi_1}, \dots, x_{\pi_n})$:

$$p(x_1, \dots, x_n) = p(x_{\pi_1}, \dots, x_{\pi_n})$$

- Clearly, if i.i.d. then:

$$p(x_1, \dots, x_n) = p(x_{\pi_1}, \dots, x_{\pi_n}) = p(x_1) \dots p(x_n)$$

hence exchangeable

- But does exchangeability imply i.i.d. ?

Exchangeability

Definition: A joint distribution $p(x_1, \dots, x_n)$ is **exchangeable** if for any permutation $\pi(x_1, \dots, x_n) = (x_{\pi_1}, \dots, x_{\pi_n})$:

$$p(x_1, \dots, x_n) = p(x_{\pi_1}, \dots, x_{\pi_n})$$

- Clearly, if i.i.d. then:

$$p(x_1, \dots, x_n) = p(x_{\pi_1}, \dots, x_{\pi_n}) = p(x_1) \dots p(x_n)$$

hence exchangeable

- But does exchangeability imply i.i.d. ?
- Answer: **NO**
- We will show that the Dirichlet process is exchangeable but not i.i.d.

- Lets assume that we are inserting 3 a's into the first cluster and 2 b's into the second cluster. Then:

$$\begin{aligned}
 & p(a, a, a, b, b | \alpha) \\
 &= p(a | \alpha) p(a | a, \alpha) p(a | a, a, \alpha) p(b | a, a, a, \alpha) p(b | a, a, a, b, \alpha) \\
 &= \frac{\alpha}{\alpha} \frac{1}{1 + \alpha} \frac{2}{2 + \alpha} \frac{\alpha}{3 + \alpha} \frac{1}{4 + \alpha}
 \end{aligned}$$

- Changing the order, we get:

$$\begin{aligned}
 & p(a, b, a, a, b | \alpha) \\
 &= p(a | \alpha) p(b | a, \alpha) p(a | a, b, \alpha) p(a | a, b, a, \alpha) p(b | a, b, a, a, \alpha) \\
 &= \frac{\alpha}{\alpha} \frac{\alpha}{1 + \alpha} \frac{1}{2 + \alpha} \frac{2}{3 + \alpha} \frac{1}{4 + \alpha}
 \end{aligned}$$

- Hence, the DP is exchangeable and clearly not i.i.d.
- Why is exchangeability important?

- Lets assume that we are inserting 3 a's into the first cluster and 2 b's into the second cluster. Then:

$$\begin{aligned}
 & p(a, a, a, b, b | \alpha) \\
 &= p(a | \alpha) p(a | a, \alpha) p(a | a, a, \alpha) p(b | a, a, a, \alpha) p(b | a, a, a, b, \alpha) \\
 &= \frac{\alpha}{\alpha} \frac{1}{1 + \alpha} \frac{2}{2 + \alpha} \frac{\alpha}{3 + \alpha} \frac{1}{4 + \alpha}
 \end{aligned}$$

- Changing the order, we get:

$$\begin{aligned}
 & p(a, b, a, a, b | \alpha) \\
 &= p(a | \alpha) p(b | a, \alpha) p(a | a, b, \alpha) p(a | a, b, a, \alpha) p(b | a, b, a, a, \alpha) \\
 &= \frac{\alpha}{\alpha} \frac{\alpha}{1 + \alpha} \frac{1}{2 + \alpha} \frac{2}{3 + \alpha} \frac{1}{4 + \alpha}
 \end{aligned}$$

- Hence, the DP is exchangeable and clearly not i.i.d.
- Why is exchangeability important?
- Without exchangeability sampling algorithms will give different solutions depending upon order.

The DP Categorical distribution

- Since the CRP is the posterior conditional distribution of the DP categorical, we can work out the full distribution using the same method as in the previous slide.
- Let $\mathbf{c} = (c_1, \dots, c_k)$ be a vector of counts, then, using the CRP:

$$\begin{aligned}
 p(\mathbf{c}|\alpha) &= \frac{\alpha^k}{(1-1+\alpha)(2-1+\alpha) \dots (\mathbf{C}-1+\alpha)} \prod_i (c_i - 1)! \\
 &= \frac{\alpha^k}{(1+(\alpha-1))(2+(\alpha-1)) \dots (\mathbf{C}+(\alpha-1))} \prod_i (c_i - 1)! \\
 &= \frac{(\alpha-1)!}{(\mathbf{C}+(\alpha-1))!} \alpha^k \prod_i (c_i - 1)! = \frac{(\alpha-1)!}{(\mathbf{C}+\alpha-1)!} \alpha^k \prod_i (c_i - 1)! \\
 &= \frac{\Gamma(\alpha)}{\Gamma(\mathbf{C}+\alpha)} \alpha^k \prod_i \Gamma(c_i)
 \end{aligned}$$

Contrast this with the Dirichlet Categorical:

$$p(\mathbf{c}|\alpha) = \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{\prod_i \Gamma(c_i + \alpha_i)}{\Gamma(\mathbf{C} + A)}$$

Dirichlet Process

- The DP is the infinite extension of the Dirichlet distribution
- So, the question is: ***How do we generate samples from the DP?***
- The Stick breaking process is a constructive method for generating samples from the DP.

Stick breaking

- Take a stick of unit length
- Break it into u_1 and $(1 - u_1)$
- The first portion is u_1 . Call this θ_1 i.e. let $\theta_1 = u_1$
- Break $(1 - u_1)$ into u_2 and $(1 - u_2)$
- The second portion is $u_2(1 - u_1) = \theta_2$
and remainder $(1 - u_2)(1 - u_1)$
-
- The k^{th} portion is

$$\theta_k = u_k \prod_{i=1}^{k-1} (1 - u_i) = u_k \left(1 - \sum_{i=1}^{k-1} \theta_i \right)$$

- with remainder

$$(1 - u_k) \prod_{i=1}^{k-1} (1 - u_i) = \prod_{i=1}^k (1 - u_i) = 1 - \sum_{i=1}^k \theta_i$$

Generating DP using Stick breaking

- Sample $u_k \sim \text{Beta}(1, \alpha)$ (simulates breaking stick of unit length)
- then:

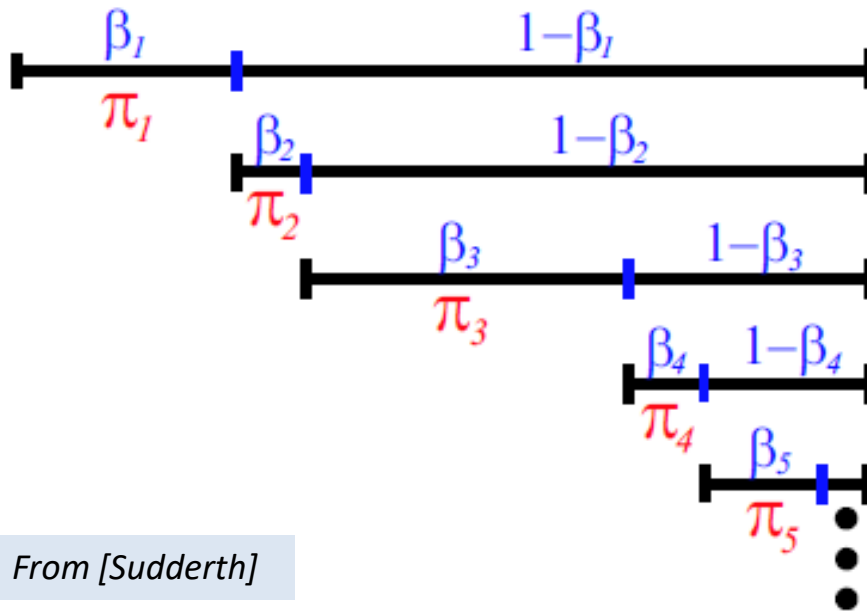
$$\theta_k = u_k \prod_{i=1}^{k-1} (1 - u_i) = u_k \left(1 - \sum_{i=1}^{k-1} \theta_i \right)$$

with remainder:

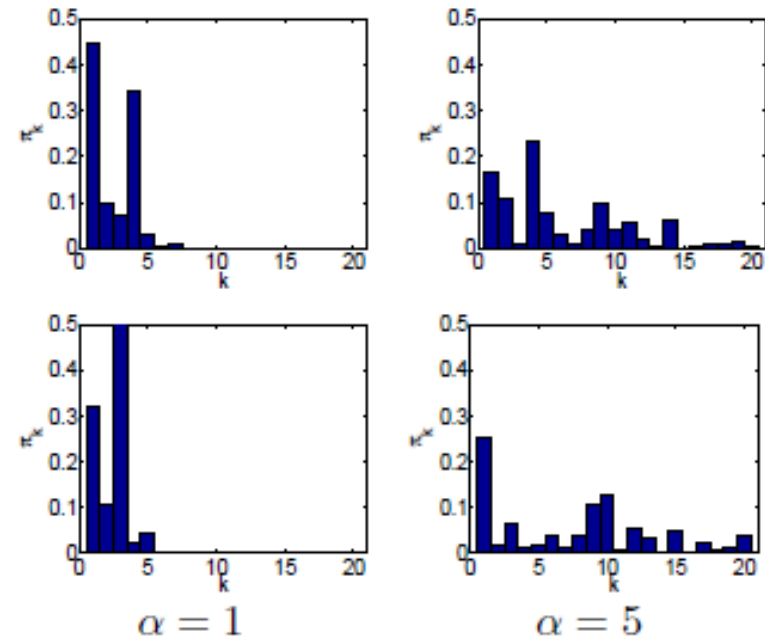
$$\prod_{i=1}^k (1 - u_i) = 1 - \sum_{i=1}^k \theta_i$$

- Since, the remainder is non-zero, we can go ad infinitum.
- Similar procedure for the Dirichlet distribution exists.

Example sample



From [Sudderth]



- Smaller values of α would result in very dense clusters
- Larger values of α would result in clusters that are more spread out

Hierarchical DP (HDP)

- The problem with standard mixture models is that each mixture component is independent of other mixture components.
- What if some information needs to be shared between mixture components?
- We already looked at an example of this when we used global codebooks for morphology learning.
- What we have so far is a scheme where the mixture proportions are generated by the CRP

- In a hierarchical CRP we assume that individual clusters have their own DPs but will borrow their cluster ids or the value of the z variable from their parent DPs.
- In a two stage DP, we will have one single parent DP.
- Lets assume that the parent DP has counts $\mathbf{p} = (p_1, p_2, \dots, p_L)$, with $\mathbf{P} - \mathbf{1} = \sum_i \mathbf{p}_i$
- The counts p_i indicates the number of different tables (i.e. mixture components) sharing the same cluster i :

$$\begin{aligned}
 p(x_i, z_i | \mathbf{Z}_{-i}, \mathbf{X}_{-i}, \mathbf{p}, \boldsymbol{\theta}, \alpha) &\propto p(x_i | \boldsymbol{\theta}_{z_i}) p(z_i | \mathbf{Z}_{-i}, \mathbf{p}, \alpha) \\
 &\propto p(x_i | \boldsymbol{\theta}_{z_i}) p(z_i | \mathbf{c}, \mathbf{p}, \alpha) \\
 &\propto p(x_i | \boldsymbol{\theta}_{z_i}) \text{hCRP}(z_i | \mathbf{c}, \mathbf{p}, \alpha)
 \end{aligned}$$

$$\text{hCRP}(z_i = k | \mathbf{c}, \mathbf{p}, \alpha) = \frac{c_k}{c + \alpha - 1} \quad \text{for existing class } k \text{ in } \mathbf{c}$$

$$\text{hCRP}(z_i = k | \mathbf{c}, \mathbf{p}, \alpha) = \frac{\alpha}{c + \alpha - 1} \text{CRP}(z_i = k | \mathbf{p}, \alpha)$$

- In a hierarchical CRP we assume that individual clusters have their own DPs but will borrow their cluster ids or the value of the z variable from their parent DPs.
- In a two stage DP, we will have one single parent DP.
- Lets assume that the parent DP has counts $\mathbf{p} = (p_1, p_2, \dots, p_L)$, with $P - 1 = \sum_i p_i$
- The counts p_i indicates the number of different tables (i.e. mixture components) sharing the same cluster i :

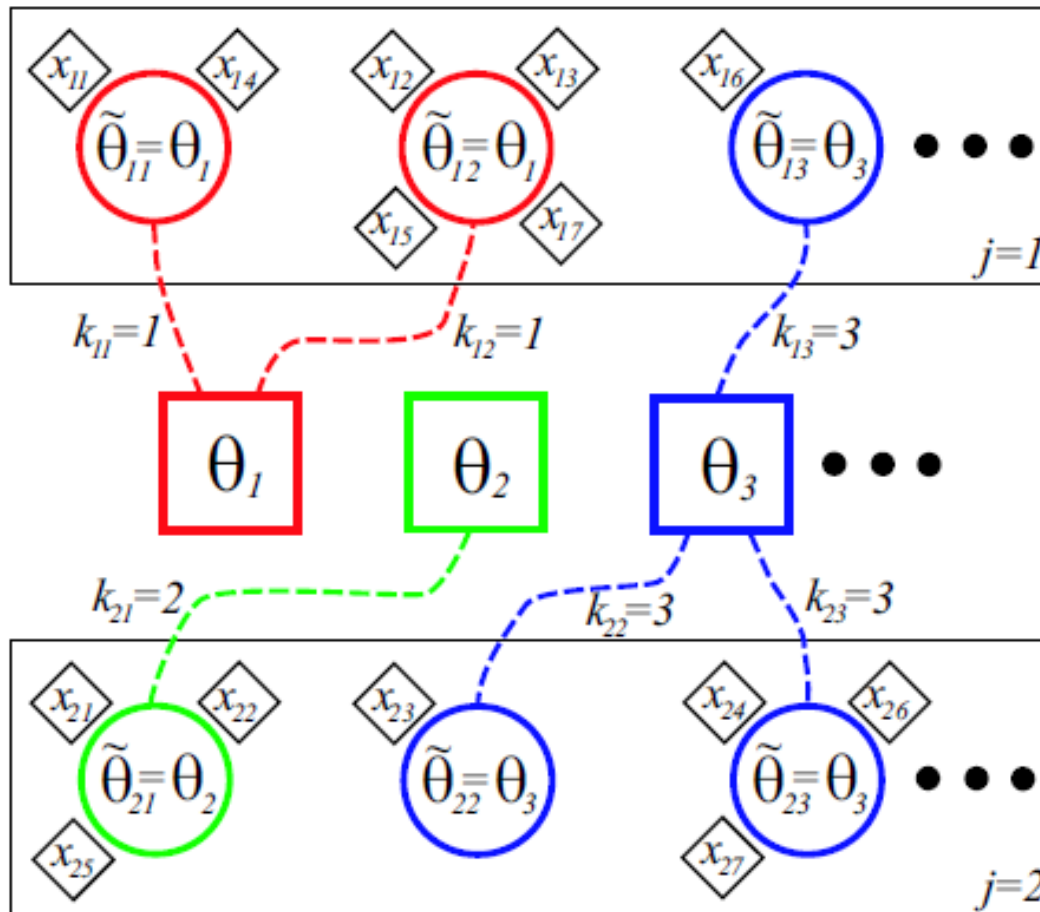
$$\begin{aligned}
 p(x_i, z_i | \mathbf{Z}_{-i}, \mathbf{X}_{-i}, \mathbf{p}, \boldsymbol{\theta}, \alpha) &\propto p(x_i | \boldsymbol{\theta}_{z_i}) p(z_i | \mathbf{Z}_{-i}, \mathbf{p}, \alpha) \\
 &\propto p(x_i | \boldsymbol{\theta}_{z_i}) p(z_i | \mathbf{c}, \mathbf{p}, \alpha) \\
 &\propto p(x_i | \boldsymbol{\theta}_{z_i}) \text{hCRP}(z_i | \mathbf{c}, \mathbf{p}, \alpha)
 \end{aligned}$$

$$\text{hCRP}(z_i = k | \mathbf{c}, \mathbf{p}, \alpha) = \frac{c_k}{C + \alpha - 1} \quad \text{for existing class } k \text{ in } \mathbf{c}$$

$$\text{hCRP}(z_i = k | \mathbf{c}, \mathbf{p}, \alpha) = \frac{\alpha}{C + \alpha - 1} \frac{p_k}{P + \beta - 1} \quad \text{for existing class } p_k \text{ in } \mathbf{p}$$

$$\text{hCRP}(z_i = L + 1 | \mathbf{c}, \mathbf{p}, \alpha) = \frac{\alpha}{C + \alpha - 1} \frac{\beta}{P + \beta - 1} \quad \text{for new class } L + 1 \text{ in } \mathbf{p}$$

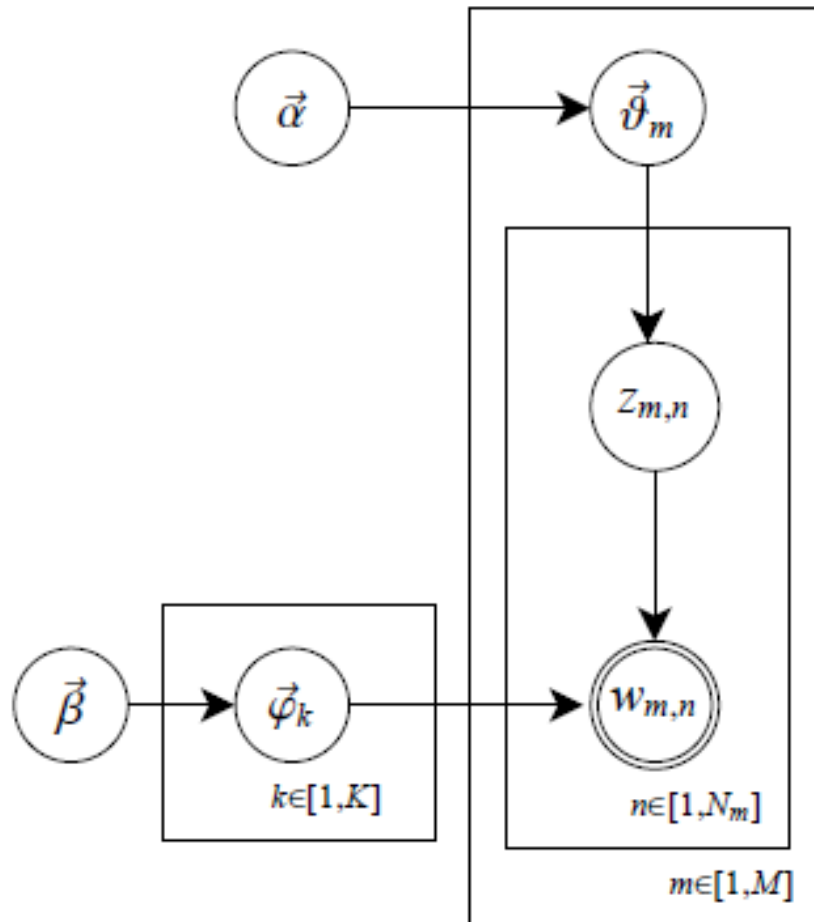
HDP illustration



From [Sudderth]

- Multiple tables in Local DPs may share the same table ids (i.e. dishes)

LDA plate diagram



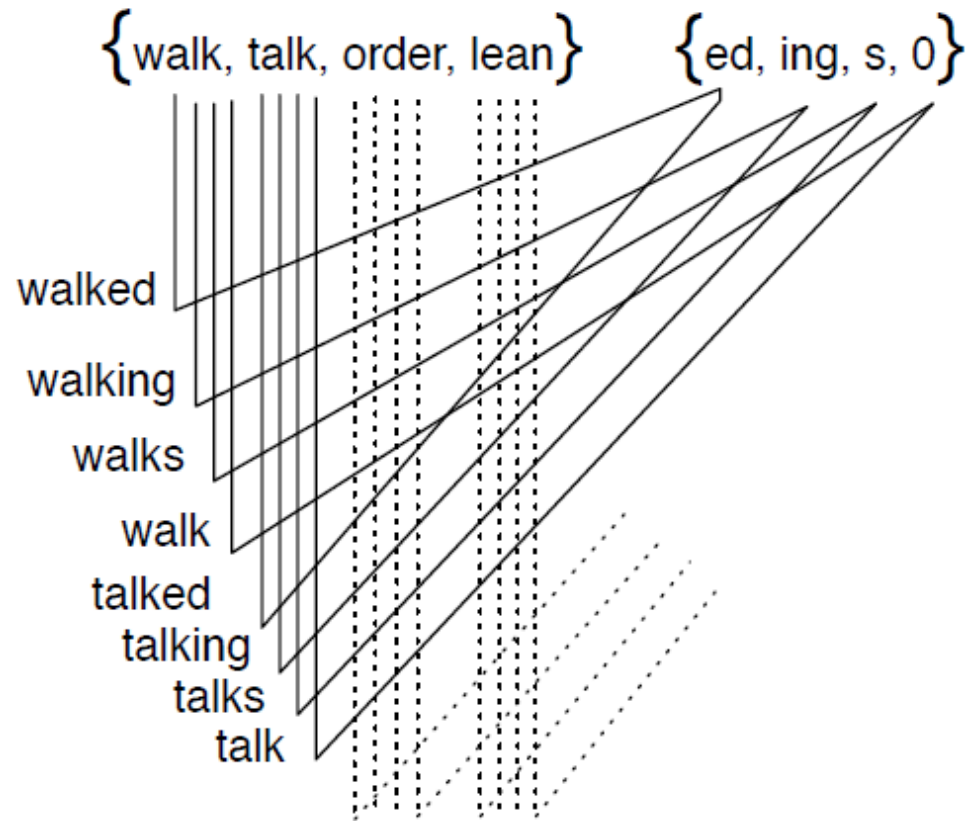
- Use of Dirichlet Multinomials means that the multinomial parameters are integrated out.

Non-parametric LDA

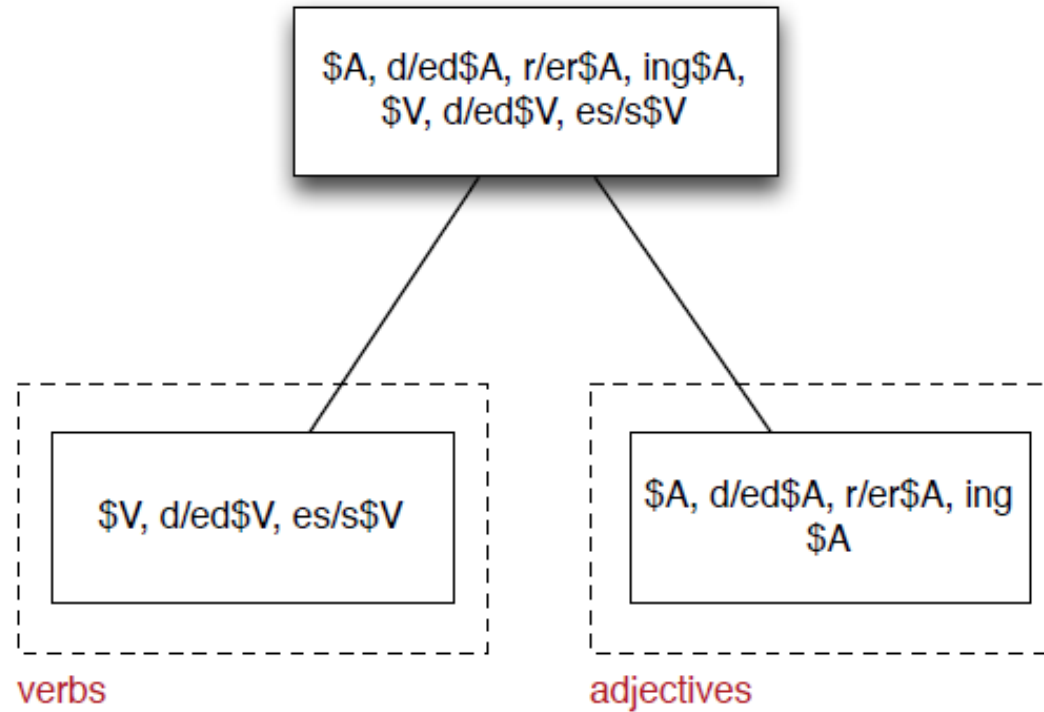
- For LDA where the number of topics need **not** be specified.
- We can employ hCRP scheme
- However, words in a LDA are generated jointly from topic distribution and topic indicators (i.e. the \mathbf{z} 's).
- So, we need to ensure that every time a new topic is created, not only should the topic be available for all documents, but topic-word distribution is also shared and updated

- A simple solution is to assume that the Topic-Word distribution is an infinite vector of Dirichlet
- Thus, $\beta = (\beta_1, \beta_2, \dots, \beta_\infty)$ where each $\beta_i = (\beta_{i_1}, \beta_{i_2}, \dots, \beta_{i_V})$ where V is the length of the vocabulary.
- Each of the β_i 's are independent from each other
- So, we let each document's Z vector to be generated by hCRP
- Each time a new topic is created within a document, hCRP ensures that the topic id is shared across all documents
- Now, since the β vector is infinite, for every (globally) new topic i there is a readily available parameter vector β_i that generates the words.
- β needs to be implemented e.g. as a hash table since, as topics die and are re-born not all topic ids will survive
- For example, there could be topics $\{1, 25, 47, 50\}$
- To manage this efficiently, keep a list of **dead topic ids** for reuse

Learning Hierarchical Morphology



Burcu Can and S. Manandhar. *Probabilistic Hierarchical Clustering of Morphological Paradigms*, EACL 2012.



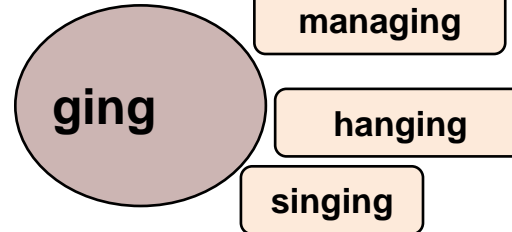
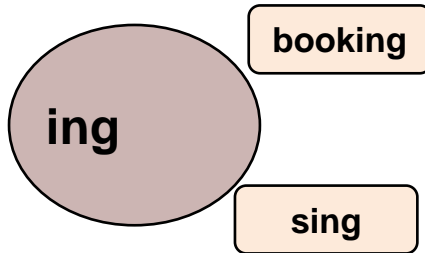
- Key idea is to learn morphological paradigms and arrange them in a hierarchy

Independent DPs for Stems/Suffixes

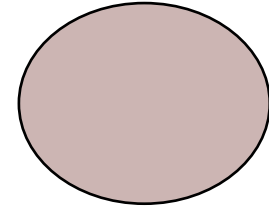
binging

new word

suffix restaurant



empty table

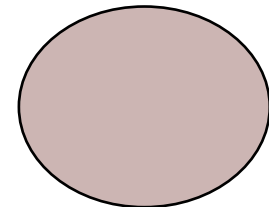
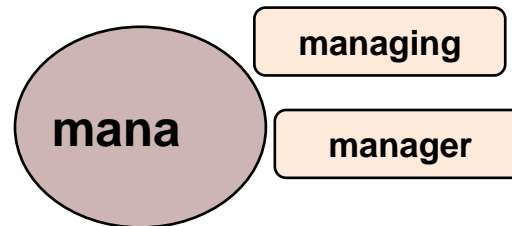
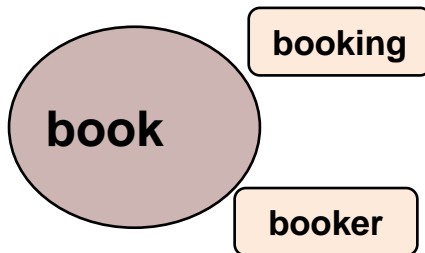


$$p(\text{ing}) = \frac{2}{(N-1+\beta)} * p(\text{ing})$$

$$p(\text{ging}) = \frac{3}{(N-1+\beta)} * p(\text{ging})$$

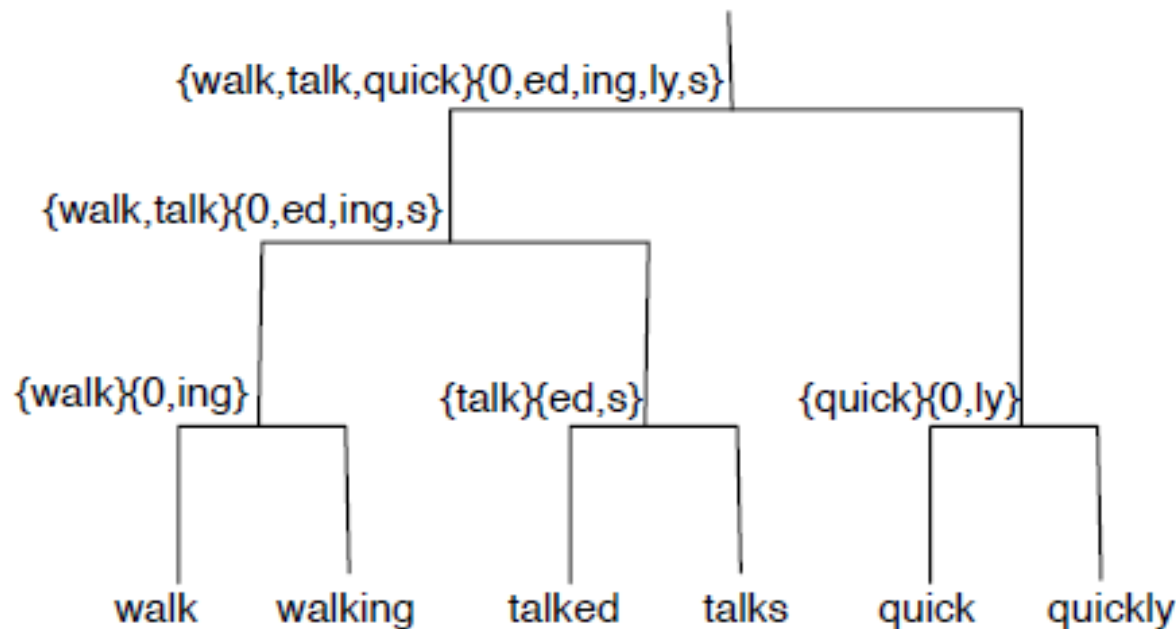
$$p(\text{ng}) = \frac{\beta}{(N-1+\beta)} * p(\text{ng})$$

stem restaurant



Data covered by a tree node

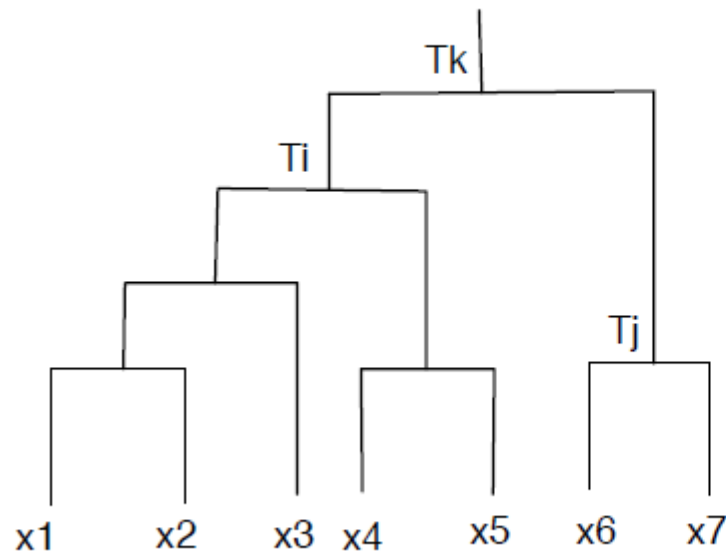
■ Dataset: $\mathbf{D} = \{w_1 = s_1 + m_1, \dots, w_n = s_n + m_n\}$



- Data points covered by a tree node are all words dominated by that node

Likelihood of a tree

■ Dataset: $\mathbf{D} = \{x_1, x_2, \dots, x_n\}$. T is the entire tree.



$$p(D_k | T_k) = p(D_k) p(D_i | T_i) p(D_j | T_j)$$

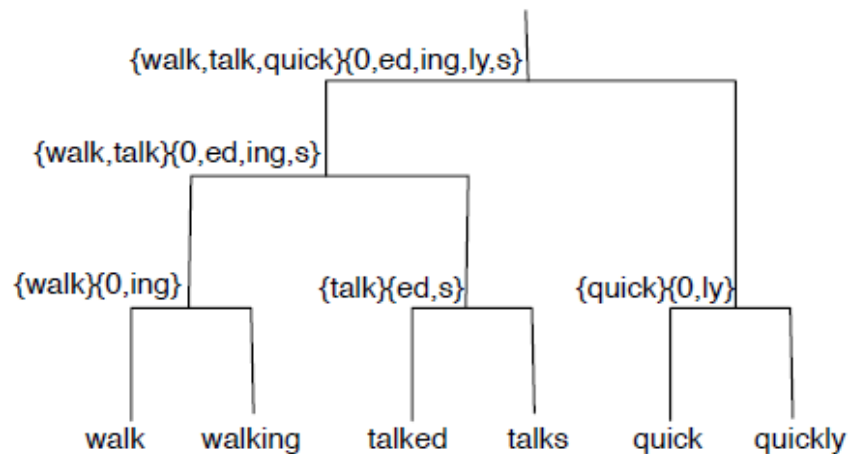
$$\mathbf{D}_i = \{x_1, x_2, x_3, x_4, x_5\}$$

$$\mathbf{D}_j = \{x_6, x_7\}$$

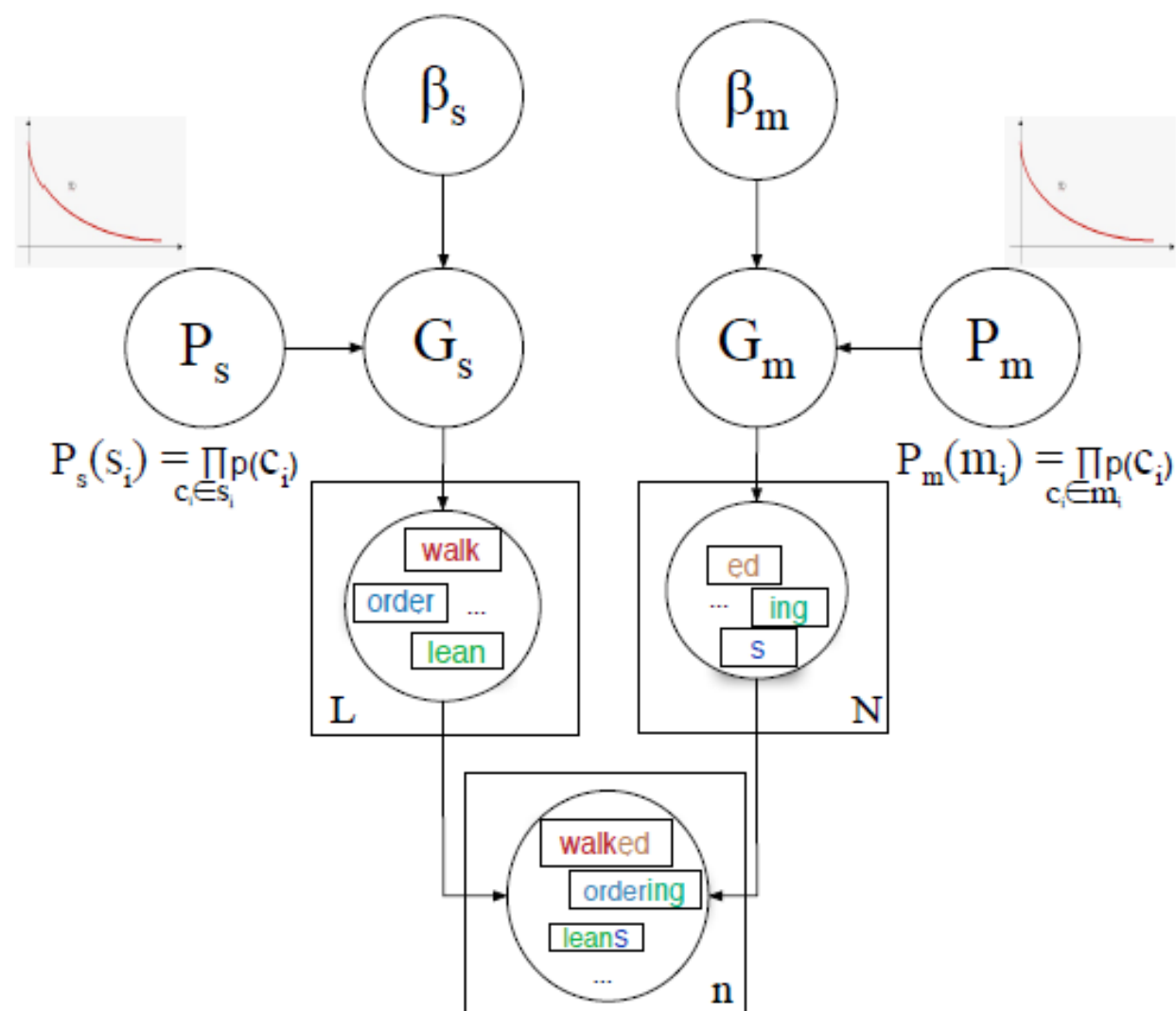
- Remember : Likelihood means $p(\mathbf{X}|\boldsymbol{\theta})$ i.e. the probability of data given the model parameters.
- Since, $\boldsymbol{\theta}$ is a tree node k here. $\boldsymbol{\theta}$ becomes T_k .
- And since, data covered by tree node is \mathbf{D}_k we write $p(\mathbf{D}_k | T_k)$

Local likelihood

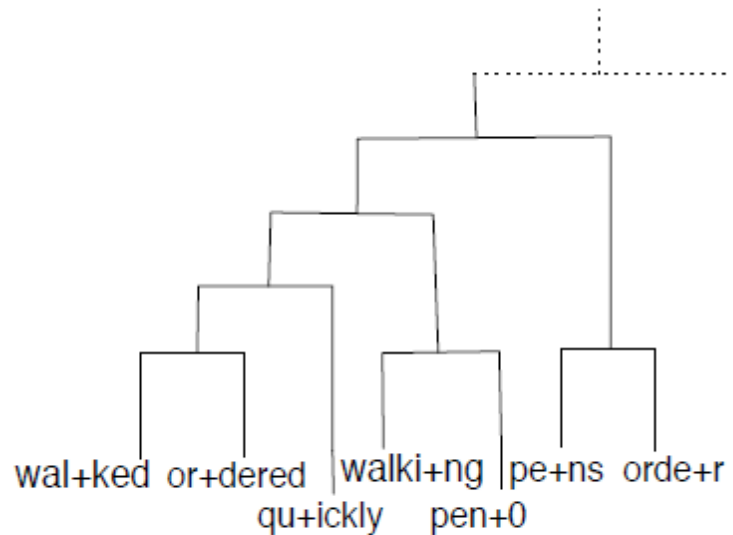
■ Dataset: $\mathbf{D} = \{w_1 = s_1 + m_1, \dots, w_n = s_n + m_n\}$



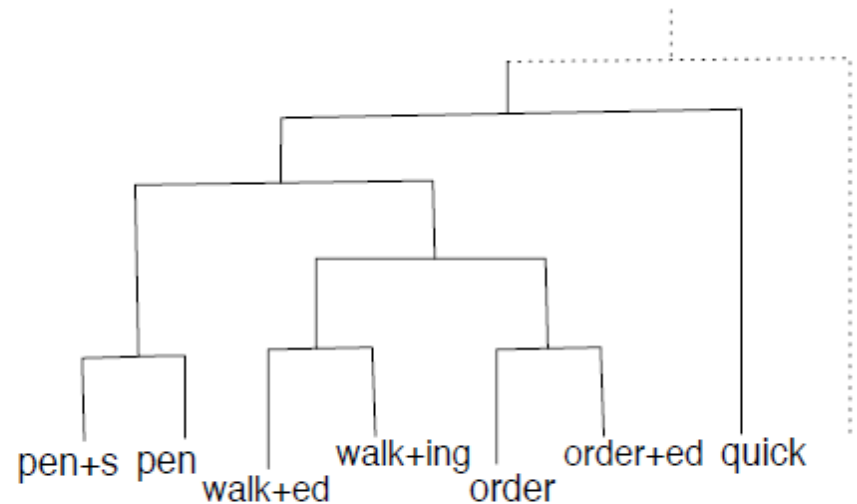
$$\begin{aligned} p(D_k) &= p(S_k)p(M_k) \\ &= p(s_1, s_2, \dots, s_n)p(m_1, m_2, \dots, m_n) \end{aligned}$$



Construct initial random tree



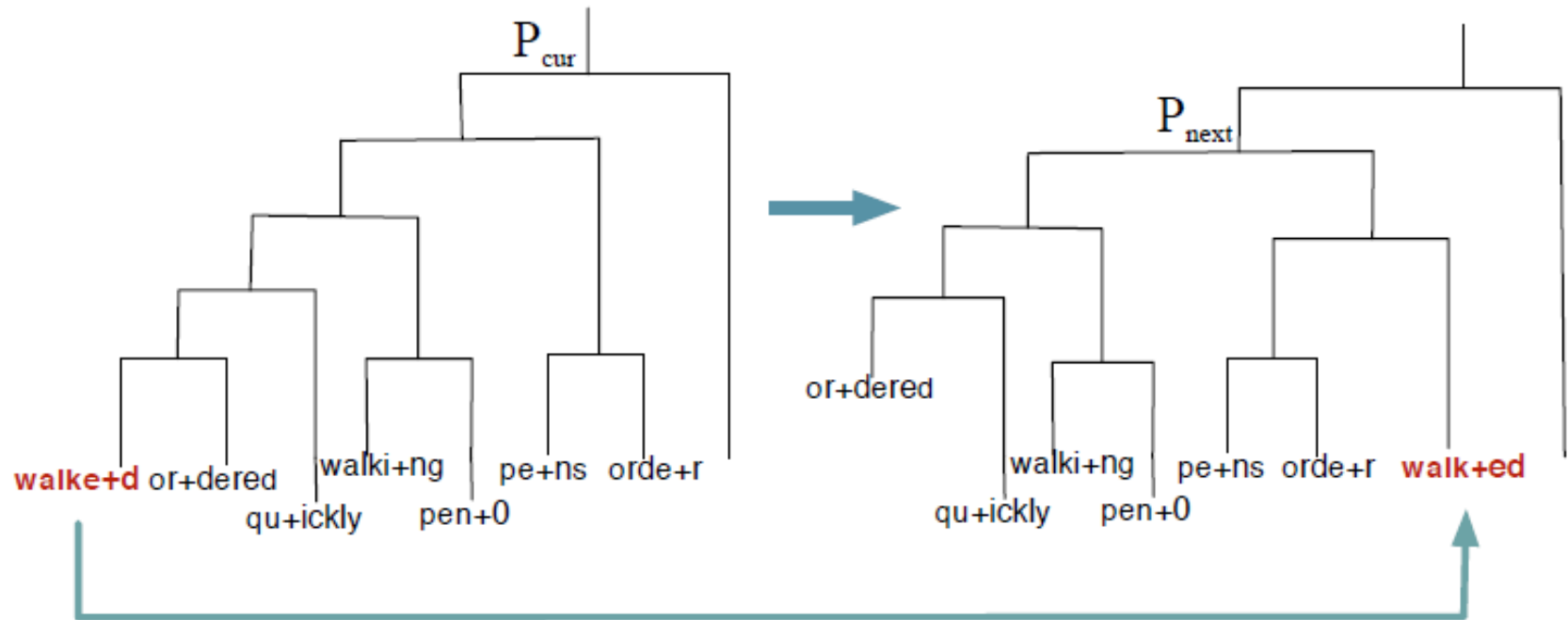
1) Construct the initial tree
(randomly generated)



2) Apply inference
(final tree)

- Generate a randomly created initial tree
- Move to the next tree configuration using sampling steps

Metropolis Hastings sampling



Accept if $p_{next}(D|T) > p_{cur}(D|T)$
 else accept with probability $P_{Acc} = \frac{p_{next}(D|T)}{p_{cur}(D|T)}$

- Gibbs sampling can be thought of as a special case of MH.
- In Gibbs, the 'if' part does not exist

Computing Likelihood

- The marginal likelihood of words in the node k is defined such that:

$$p(D_k) = p(s_1, s_2, \dots, s_n) p(m_1, m_2, \dots, m_n)$$

- The joint probability of stems $S = \{s_1, \dots, s_L\}$ is:

$$\begin{aligned} p(s_1, s_2, \dots, s_L) &= p(s_1) p(s_2 | s_1) \dots p(s_M | s_1, \dots, s_{M-1}) \\ &= \frac{\Gamma(\beta_s)}{\Gamma(L + \beta_s)} \beta_s^K \prod_{i=1}^K P_s(s_i) \prod_{i=1}^K (n_{s_i} - 1)! \end{aligned}$$

- The joint probability of suffixes $M = \{m_1, \dots, m_N\}$ is:

$$\begin{aligned} p(m_1, m_2, \dots, m_N) &= p(m_1) p(m_2 | m_1) \dots p(m_N | m_1, \dots, m_{N-1}) \\ &= \frac{\Gamma(\alpha)}{\Gamma(N + \alpha)} \alpha^T \prod_{i=1}^T P_m(m_i) \prod_{i=1}^T (n_{m_i} - 1)! \end{aligned}$$

T

- Since, MH based on using likelihoods and not conditional likelihoods as in Gibbs, we need to compute the full likelihood.
- So, this is just DP categorical distribution as worked out in slide 24

References

- **Tutorials**

Peter Orbanz. *Lecture Notes on Bayesian Nonparametrics*. Columbia University

Erik B. Sudderth. *Graphical Models for Visual Object Recognition and Tracking*. PhD Thesis, MIT, 2006.

Yee Whye Teh. *Introduction to Bayesian Nonparametrics*. Machine Learning Summer School (MLSS) 2011.

- **Papers on Dirichlet Processes**

Various papers by Radford Neal on Dirichlet Processes and sampling schemes

- **Hierarchical Morphology clustering based on the following paper**

Burcu Can and Suresh Manandhar. Probabilistic Hierarchical Clustering Of Morphological Paradigms. In *13th Conference of the European Chapter of the Association for computational Linguistics (EACL)*, 2012.