

Course Notes for EECS 9601: Advanced Probabilistic Machine Learning

John Paisley
Data Science Institute, Electrical Engineering
Columbia University

Spring 2022 (version date: February 1, 2022)

1	Poisson distribution and process, superposition and marking theorems	1
2	Completely random measures, Campbell's theorem, gamma process	11
3	Beta processes and the Poisson process	18
4	Beta processes and size-biased constructions	24
5	Dirichlet processes and a size-biased construction	30
6	Dirichlet process extensions, count processes	37
7	Exchangeability, Dirichlet processes and the Chinese restaurant process	44
8	Exchangeability, beta processes and the Indian buffet process	52
9	Another look at constructive definitions of the beta and Dirichlet process	56

Chapter 1

Poisson distribution and process, superposition and marking theorems

- The Poisson distribution is perhaps the fundamental discrete distribution and, along with the Gaussian distribution, one of the two fundamental distributions of probability.

Importance: Poisson → discrete r.v.'s
Gaussian → continuous r.v.'s

Definition: A random variable $X \in \{0, 1, 2, \dots\}$ is Poisson distributed with parameter $\lambda > 0$ if

$$P(X = n|\lambda) = \frac{\lambda^n}{n!}e^{-\lambda}, \quad (1.1)$$

denoted $X \sim \text{Pois}(\lambda)$.

Moments of Poisson

$$\mathbb{E}[X] = \sum_{n=1}^{\infty} nP(X = n|\lambda) = \sum_{n=1}^{\infty} \frac{\lambda^n}{(n-1)!}e^{-\lambda} = \lambda \underbrace{\sum_{n=1}^{\infty} \frac{\lambda^{n-1}}{(n-1)!}e^{-\lambda}}_{=1} \quad (1.2)$$

$$\begin{aligned} \mathbb{E}[X^2] &= \sum_{n=1}^{\infty} n^2 P(X = n|\lambda) = \lambda \sum_{n=1}^{\infty} \frac{n\lambda^{n-1}}{(n-1)!}e^{-\lambda} = \lambda \sum_{n=0}^{\infty} \frac{(n+1)\lambda^n}{n!}e^{-\lambda} \\ &= \lambda(\mathbb{E}[X] + 1) = \lambda^2 + \lambda \end{aligned} \quad (1.3)$$

$$\mathbb{V}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \lambda \quad (1.4)$$

Sums of Poisson r.v.'s (take 1)

- Sums of Poisson r.v.'s are also Poisson. Let $X_1 \sim \text{Pois}(\lambda_1)$ and $X_2 \sim \text{Pois}(\lambda_2)$. Then $X_1 + X_2 \sim \text{Pois}(\lambda_1 + \lambda_2)$.

Important interlude: Laplace transforms and sums of r.v.'s

- Laplace transforms give a very easy way to calculate the distribution of sums of r.v.'s (among other things).

Laplace transform

- Let $X \sim p(X)$ be a positive random variable and let $t > 0$. The Laplace transform of X is

$$\mathbb{E}[e^{-tX}] = \int_X e^{-tx} p(x) dx \quad (\text{sums when appropriate}) \quad (1.5)$$

Important property

- There is a one-to-one mapping between $p(x)$ and $\mathbb{E}[e^{-tX}]$. That is, if $p(x)$ and $q(x)$ are two distributions and $\mathbb{E}_p[e^{-tX}] = \mathbb{E}_q[e^{-tX}]$, then $p(x) = q(x)$ for all x . (p and q are the same distribution)

Sums of r.v.'s

- Let $X_1 \stackrel{ind}{\sim} p(x)$, $X_2 \stackrel{ind}{\sim} q(x)$ and $Y = X_1 + X_2$. What is the distribution of Y ?
- Approach: Take the Laplace transform of Y and see what happens.

$$\mathbb{E}e^{-tY} = \mathbb{E}e^{-t(X_1+X_2)} = \underbrace{\mathbb{E}[e^{-tX_1}e^{-tX_2}]}_{\text{by independence of } X_1 \text{ and } X_2} = \mathbb{E}[e^{-tX_1}]\mathbb{E}[e^{-tX_2}] \quad (1.6)$$

- So we can multiply the Laplace transforms of X_1 and X_2 and see if we recognize it.

Sums of Poisson r.v.'s (take 2)

- The Laplace transform of a Poisson random variable has a very important form that should be memorized.

$$\mathbb{E}e^{-tX} = \sum_{n=0}^{\infty} e^{-tn} \frac{\lambda^n}{n!} e^{-\lambda} = e^{-\lambda} \sum_{n=0}^{\infty} \frac{(\lambda e^{-t})^n}{n!} = e^{-\lambda} e^{\lambda e^{-t}} = e^{\lambda(e^{-t}-1)} \quad (1.7)$$

- Back to the problem: $X_1 \stackrel{ind}{\sim} \text{Pois}(\lambda_1)$, $X_2 \stackrel{ind}{\sim} \text{Pois}(\lambda_2)$, $Y = X_1 + X_2$.

$$\mathbb{E}e^{-tY} = \mathbb{E}[e^{-tX_1}]\mathbb{E}[e^{-tX_2}] = e^{\lambda_1(e^{-t}-1)} e^{\lambda_2(e^{-t}-1)} = e^{(\lambda_1+\lambda_2)(e^{-t}-1)} \quad (1.8)$$

We recognize that the last term is the Laplace transform of a $\text{Pois}(\lambda_1 + \lambda_2)$ random variable. We can therefore conclude that $Y \sim \text{Pois}(\lambda_1 + \lambda_2)$.

- Another way of saying this is that, if we draw $Y_1 \sim \text{Pois}(\lambda_1 + \lambda_2)$ and $X_1 \sim \text{Pois}(\lambda_1)$ and $X_2 \sim \text{Pois}(\lambda_2)$ and define $Y_2 = X_1 + X_2$. Then Y_1 is equal to Y_2 *in distribution*. (i.e., they may not be equal, but they have the same distribution. We write this as $Y_1 \stackrel{d}{=} Y_2$.)

- The idea extends to sums of more than two. Let $X_i \sim \text{Pois}(\lambda_i)$. Then $\sum_i X_i \sim \text{Pois}(\sum_i \lambda_i)$ since

$$\mathbb{E}e^{-t\sum_i X_i} = \prod_i \mathbb{E}e^{-tX_i} = e^{(\sum_i \lambda_i)(e^{-t}-1)}. \quad (1.9)$$

A conjugate prior for λ

- What if we have X_1, \dots, X_N that we believe to be generated by a $\text{Pois}(\lambda)$ distribution, but we don't know λ ?
- One answer: Put a prior distribution on λ , $p(\lambda)$, and calculate the posterior of λ using Bayes' rule.

Bayes' rule (review)

$$P(A, B) = P(A|B)P(B) = P(B|A)P(A) \tag{1.10}$$

$$\Downarrow$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \tag{1.11}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \tag{1.12}$$

Gamma prior

- Let $\lambda \sim \text{Gam}(a, b)$, where $p(\lambda|a, b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}$ is a gamma distribution. Then the posterior of λ is

$$p(\lambda|X_1, \dots, X_N) \propto p(X_1, \dots, X_N|\lambda)p(\lambda) = \left[\prod_{i=1}^N \frac{\lambda^{X_i}}{X_i!} e^{-\lambda} \right] \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}$$

$$\propto \lambda^{a+\sum_{i=1}^N X_i - 1} e^{-(b+N)\lambda} \tag{1.13}$$

$$\Downarrow$$

$$p(\lambda|X_1, \dots, X_N) = \text{Gam}(a + \sum_{i=1}^N X_i, b + N) \tag{1.14}$$

Note that $\mathbb{E}[\lambda|X_1, \dots, X_N] = \frac{a+\sum_{i=1}^N X_i}{b+N} \approx$ Empirical average of X_i
 (Makes sense because $\mathbb{E}[X|\lambda] = \lambda$)

$\mathbb{V}[\lambda|X_1, \dots, X_N] = \frac{a+\sum_{i=1}^N X_i}{(b+N)^2} \approx$ Empirical average/ N
 (Get more confident as we see more X_i)

- The gamma distribution is said to be the conjugate prior for the parameter of the Poisson distribution because the posterior is also gamma.

Poisson–Multinomial

- A sequence of Poisson r.v.'s is closely related to the multinomial distribution as follows:

Let $X_i \stackrel{ind}{\sim} \text{Pois}(\lambda_i)$ and let $Y = \sum_{i=1}^N X_i$.

Then what is the distribution of $\vec{X} = (X_1, \dots, X_N)$ given Y ?

We can use basic rules of probability. . .

$$\begin{aligned} P(X_1, \dots, X_N) &= P(X_1, \dots, X_N, Y = \sum_{i=1}^N X_i) \leftarrow (Y \text{ is a deterministic function of } X_{1:N}) \\ &= P(X_1, \dots, X_N | Y = \sum_{i=1}^N X_i) P(Y = \sum_{i=1}^N X_i) \end{aligned} \quad (1.15)$$

And so

$$P(X_1, \dots, X_N | Y = \sum_{i=1}^N X_i) = \frac{P(X_1, \dots, X_N)}{P(Y = \sum_{i=1}^N X_i)} = \frac{\prod_i P(X_i)}{P(Y = \sum_{i=1}^N X_i)} \quad (1.16)$$

- We know that $P(Y) = \text{Pois}(Y; \sum_{i=1}^N \lambda_i)$, so

$$\begin{aligned} P(X_1, \dots, X_N | Y = \sum_i X_i) &= \frac{\left[\prod_{i=1}^N \frac{\lambda_i^{X_i}}{X_i!} e^{-\lambda_i} \right]}{\left[\frac{(\sum_{i=1}^N \lambda_i)^{\sum_{i=1}^N X_i}}{(\sum_{i=1}^N X_i)!} e^{-\sum_{i=1}^N \lambda_i} \right]} \\ &= \frac{Y!}{X_1! \cdots X_N!} \prod_{i=1}^N \left(\frac{\lambda_i}{\sum_{j=1}^N \lambda_j} \right)^{X_i} \\ &\quad \downarrow \\ &\text{Mult}(Y; p_1, \dots, p_N), \quad p_i = \lambda_i / \sum_j \lambda_j \end{aligned} \quad (1.17)$$

- What is this saying?
 1. Given the sum of N independent Poisson r.v.'s, the individual values are distributed as a multinomial using the normalized parameters.
 2. We can sample X_1, \dots, X_N in two *equivalent* ways
 - a) Sample $X_i \sim \text{Pois}(\lambda_i)$ independently
 - b) First sample $Y \sim \text{Pois}(\sum_j \lambda_j)$, then $(X_1, \dots, X_N) \sim \text{Mult}\left(Y; \frac{\lambda_1}{\sum_j \lambda_j}, \dots, \frac{\lambda_N}{\sum_j \lambda_j}\right)$

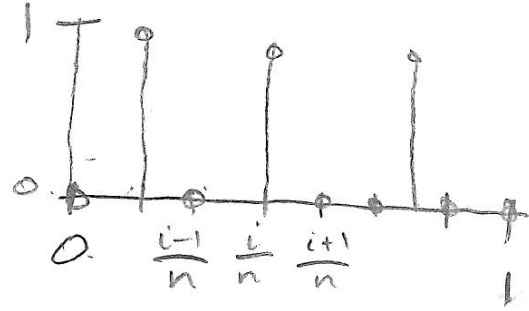
Poisson as a limiting case distribution

- The Poisson distribution arises as a limiting case of the sum over many binary events each having small probability.

Binomial distribution and Bernoulli process

- Imagine we have an array of random variables X_{nm} , where $X_{nm} \sim \text{Bern}\left(\frac{\lambda}{n}\right)$ for $m = 1, \dots, n$ and fixed $0 \leq \lambda \leq n$. Let $Y_n = \sum_{m=1}^n X_{nm}$ and $Y = \lim_{n \rightarrow \infty} Y_n$. Then $Y_n \sim \text{Bin}\left(n, \frac{\lambda}{n}\right)$ and $Y \sim \text{Pois}(\lambda)$.

Picture: Have n coins with bias λ/n evenly spaced between $[0, 1]$. Go down the line and flip each independently. In the limit $n \rightarrow \infty$, the total # of 1's is a $\text{Pois}(\lambda)$ r.v.



Proof:

$$\begin{aligned} \lim_{n \rightarrow \infty} P(Y_n = k | \lambda) &= \lim_{n \rightarrow \infty} \frac{n!}{(n-k)!k!} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} && (1.18) \\ &= \lim_{n \rightarrow \infty} \underbrace{\left[\frac{n(n-1) \cdots (n-k+1)}{n^k} \right]}_{\rightarrow 1} \underbrace{\left[\left(1 - \frac{\lambda}{n}\right)^{-k} \right]}_{\rightarrow 1} \underbrace{\left[\frac{\lambda^k}{k!} \left(1 - \frac{\lambda}{n}\right)^n \right]}_{\rightarrow \frac{\lambda^k}{k!} e^{-\lambda}} \\ &= \text{Pois}(k; \lambda) \end{aligned}$$

So $\lim_{n \rightarrow \infty} \text{Bin}\left(n, \frac{\lambda}{n}\right) = \text{Pois}(\lambda)$.

A more general statement

- Let λ_{nm} be an array of positive numbers such that

1. $\sum_{m=1}^n \lambda_{nm} = \lambda < \infty$
2. $\lambda_{nm} < 1$ and $\lim_{n \rightarrow \infty} \lambda_{nm} = 0$ for all m , i.e., $\lim_{n \rightarrow \infty} \max_m \lambda_{nm} = 0$

Let $X_{nm} \sim \text{Bern}(\lambda_{nm})$ for $m = 1, \dots, n$. Let $Y_n = \sum_{m=1}^n X_{nm}$ and $Y = \lim_{n \rightarrow \infty} Y_n$. Then $Y \sim \text{Pois}(\lambda)$.

Proof: (use Laplace transform)

1. $\mathbb{E} e^{-tY} = \lim_{n \rightarrow \infty} \mathbb{E} e^{-tY_n} = \lim_{n \rightarrow \infty} \mathbb{E} e^{-t \sum_{m=1}^n X_{nm}} = \lim_{n \rightarrow \infty} \prod_{m=1}^n \mathbb{E} e^{-tX_{nm}}$
2. $\mathbb{E} e^{-tX_{nm}} = \lambda_{nm} e^{-t \cdot 1} + (1 - \lambda_{nm}) e^{-t \cdot 0} = 1 - \lambda_{nm}(1 - e^{-t})$
3. So $\mathbb{E} e^{-tY} = \lim_{n \rightarrow \infty} \prod_{m=1}^n (1 - \lambda_{nm}(1 - e^{-t})) = \lim_{n \rightarrow \infty} e^{\sum_{m=1}^n \ln(1 - \lambda_{nm}(1 - e^{-t}))}$
4. $\ln(1 - \lambda_{nm}(1 - e^{-t})) = - \sum_{s=1}^{\infty} \frac{1}{s} \lambda_{nm}^s (1 - e^{-t})^s$ because $0 \leq 1 - \lambda_{nm}(1 - e^{-t}) < 1$
5. $\sum_{m=1}^n \ln(1 - \lambda_{nm}(1 - e^{-t})) = - \underbrace{\left(\sum_{m=1}^n \lambda_{nm} \right)}_{=\lambda} (1 - e^{-t}) - \underbrace{\sum_{s=2}^{\infty} \frac{1}{s} \left(\sum_{m=1}^n \lambda_{nm}^s \right) (1 - e^{-t})^s}_{\rightarrow 0 \text{ as } n \rightarrow \infty}$

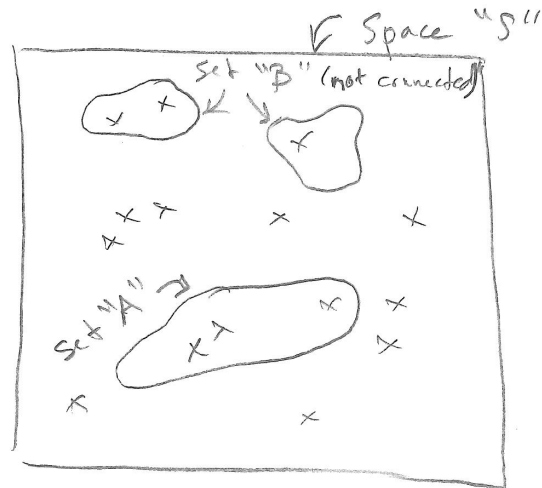
So $\mathbb{E} e^{-tY} = e^{-\lambda(1 - e^{-t})}$. Therefore $Y \sim \text{Pois}(\lambda)$.

Poisson process

- In many ways, the Poisson process is no more complicated than the previous discussion on the Poisson distribution. In fact, the Poisson process can be thought of as a “structured” Poisson distribution, which should hopefully be more clear below.

Intuitions and notations

- S : a space (think \mathbb{R}^d or part of \mathbb{R}^d)
- Π : a random countable subset of S (i.e., a random # of points and their locations)
- $A \subset S$: a subset of S
- $N(A)$: a counting measure. Counts how many points in Π fall in A (i.e., $N(A) = |\Pi \cap A|$)
- $\mu(\cdot)$: a measure on S
 - $\hookrightarrow \mu(A) \geq 0$ for $|A| > 0$
 - $\mu(\cdot)$ is non-atomic
 - $\hookrightarrow \mu(A) \rightarrow 0$ as $|A| \rightarrow \emptyset$



Think of μ as a scaled probability distribution that is continuous so that $\mu(\{x\}) = 0$ for all points $x \in S$.

Poisson processes

- A Poisson process Π is a countable subset of S such that
 1. For $A \subset S$, $N(A) \sim \text{Pois}(\mu(A))$
 2. For disjoint sets A_1, \dots, A_k , $N(A_1), \dots, N(A_k)$ are independent Poisson random variables.

$N(\cdot)$ is called a “Poisson random measure”. (See above for mapping from Π to $N(\cdot)$)

Some basic properties

- The most basic properties follow from the properties of a Poisson distribution.
 - a) $\mathbb{E}N(A) = \mu(A) \rightarrow$ therefore μ is sometimes referred to as a “mean measure”
 - b) If A_1, \dots, A_k are disjoint, then

$$N(\bigcup_{i=1}^k A_i) = \sum_{i=1}^k N(A_i) \sim \text{Pois}\left(\sum_{i=1}^k \mu(A_i)\right) = \text{Pois}\left(\mu\left(\bigcup_{i=1}^k A_i\right)\right) \quad (1.19)$$

Since this holds for $k \rightarrow \infty$ and $\mu(A_i) \searrow 0$, $N(\cdot)$ is “infinitely divisible”.

c) Let A_1, \dots, A_k be disjoint subsets of S . Then

$$P(N(A_1) = n_1, \dots, N(A_k) = n_k | N(\bigcup_{i=1}^k A_i) = n) = \frac{P(N(A_1)=n_1, \dots, N(A_k)=n_k)}{P(N(\bigcup_{i=1}^k A_i)=n)} \quad (1.20)$$

Notice from earlier that $N(A_i) \Leftrightarrow X_i$. Following the same exact calculations,

$$P(N(A_1) = n_1, \dots, N(A_k) = n_k | N(\bigcup_{i=1}^k A_i) = n) = \frac{n!}{n_1! \dots n_k!} \prod_{i=1}^k \left(\frac{\mu(A_i)}{\mu(\bigcup_{j=1}^k A_j)} \right)^{n_i} \quad (1.21)$$

Drawing from a Poisson process (break in the basic properties)

- Property (c) above gives a very simple way for drawing $\Pi \sim \text{PP}(\mu)$, though some thought is required to see why.
 1. Draw the total number of points $N(S) \sim \text{Pois}(\mu(S))$.
 2. For $i = 1, \dots, N(S)$ draw $X_i \stackrel{iid}{\sim} \mu/\mu(S)$. In other words, normalize μ to get a probability distribution on S .
 3. Define $\Pi = \{X_1, \dots, X_{N(S)}\}$.

Notice X_i is different from previously. Now it's a location, while before it was a count. Here, the Poisson distribution is determining the existence of an X_i , while the mean measure μ determines its value.

d) Final basic property (of these notes)

$$\mathbb{E} e^{-tN(A)} = e^{\mu(A)(e^{-t}-1)} \longrightarrow \text{an obvious result since } N(A) \sim \text{Pois}(\mu(A)) \quad (1.22)$$

Some more advanced properties

Superposition Theorem: Let Π_1, Π_2, \dots be a countable collection of independent Poisson processes with $\Pi_i \sim \text{PP}(\mu_i)$. Let $\Pi = \bigcup_{i=1}^{\infty} \Pi_i$. Then $\Pi \sim \text{PP}(\mu)$ with $\mu = \sum_{i=1}^{\infty} \mu_i$.

Proof: Remember from the definition of a Poisson process we have to show two things.

1. Let $N(A)$ be the PRM (Poisson random measure) associated with PP (Poisson process) Π and $N_i(A)$ with Π_i . Clearly $N(A) = \sum_{i=1}^{\infty} N_i(A)$, and since $N_i(A) \sim \text{Pois}(\mu_i(A))$ by definition, it follows that $N(A) \sim \text{Pois}(\sum_{i=1}^{\infty} \mu_i(A))$.
2. Let A_1, \dots, A_k be disjoint. Then $N(A_1), \dots, N(A_k)$ are independent because $N_i(A_j)$ are independent for all i and j .

Restriction Theorem: If we restrict Π to a subset of S , we still have a Poisson process. Let $S_1 \subset S$ and $\Pi_i = \Pi \cap S_1$. Then $\Pi_1 \sim \text{PP}(\mu_1)$, where $\mu_1(A) = \mu(S_1 \cap A)$. This can be thought of as setting $\mu = 0$ outside of S_1 , or just looking at the subspace S_1 and ignoring the rest of S .

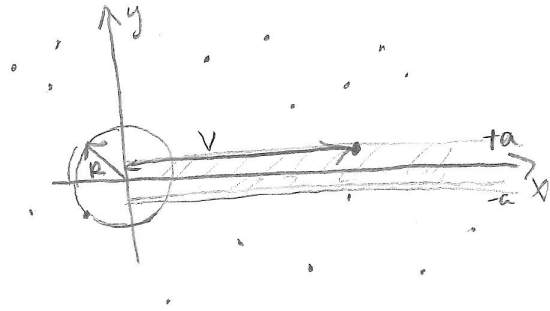
Mapping Theorem: This says that a one-to-one function $y = f(x)$ preserves the Poisson process. That is, if $\Pi_x \sim \text{PP}(\mu)$ and $\Pi_y = f(\Pi_x)$, then Π_y is also a Poisson process with the proper transformation made to μ . (See Kingman for details. We won't use this in this class.)

Example

Let N be a Poisson random measure on \mathbb{R}^2 with mean measure $\mu(A) = c|A|$.

$|A|$ is the area of A (“Lebesgue measure”)

Intuition : Think trees in a forest.



Question 1: What is the distance R from the origin to the nearest point?

Answer: We know that $R > r$ if $N(B_r) = 0$, where $B_r =$ ball of radius r . Since these are equivalent events, we know that

$$P(R > r) = P(N(B_r) = 0) = e^{-\mu(B_r)} = e^{-c\pi r^2}. \tag{1.23}$$

So $p(R = r) = 2c\pi r e^{-c\pi r^2} dr$, which is a Rayleigh distribution.

Question 2: Let each atom be the center of a disk of radius a . Take our line of sight as the x -axis. How far can we see?

Answer: The distance V is equivalent to the farthest we can extend a rectangle D_x with y -axis boundaries of $[-a, a]$. We know that $V > x$ if $N(D_x) = 0$. Therefore

$$P(V > x) = P(N(D_x) = 0) = e^{-\mu(D_x)} = e^{-2acx}. \tag{1.24}$$

So $p(V = x) = 2ace^{-2acx} dx$ which is a Gamma distribution.

Marked Poisson processes

- The other major theorem of these notes relates to “marking” the points of a Poisson process with a random variable.
- Let $\Pi \sim \text{PP}(\mu)$. For each $x \in \Pi$, associate a r.v. $y \sim p(y|x)$. We say that y has “marked” x . The results is also a Poisson process.

Theorem: Let μ be a measure on space S and $p(y|x)$ a probability distribution on space M . For each $x \in \Pi \sim \text{PP}(\mu)$ draw $y \sim p(y|x)$ and define $\Pi^* = \{(x_i, y_i)\}$. Then Π^* is a Poisson process on $S \times M$ with mean measure $\mu^* = \mu(dx)p(y|x)dy$.

Comment: If $N^*(C) = |\Pi^* \cap C|$ for $C \subset S \times M$, this says that $N^*(C) \sim \text{Pois}(\mu^*(C))$, where $\mu^*(C) = \int_C \mu(dx)p(y|x)dy$.

Proof: Need to show that $\mathbb{E}e^{-tN^*(C)} = \exp\left\{\int_C (e^{-t} - 1)\mu(dx)p(y|x)dy\right\}$

1. Note that $N^*(C) = \sum_{i=1}^{N(S)} \mathbf{1}\{(x_i, y_i) \in C\}$. $N(S)$ is PRM associated with $\Pi \sim \text{PP}(\mu)$.
2. Recall tower property: $\mathbb{E}f(A, B) = \mathbb{E}[\mathbb{E}[f(A, B)|B]]$.
3. Therefore, $\mathbb{E}e^{-tN^*(C)} = \mathbb{E}\left[\mathbb{E}\left[\exp\left\{-t \sum_{i=1}^{N(S)} \mathbf{1}\{(x_i, y_i) \in C\}\right\} \middle| \Pi\right]\right]$
4. Manipulating :

$$\begin{aligned} \mathbb{E}e^{-tN^*(C)} &= \mathbb{E}\left[\prod_{i=1}^{N(S)} \mathbb{E}[e^{-t\mathbf{1}\{(x_i, y_i) \in C\}} | \Pi]\right] \\ &= \mathbb{E}\left[\prod_{i=1}^{N(S)} \int_M \{e^{-t\mathbf{1}\{(x_i, y_i) \in C\}} + e^{-t\mathbf{1}\{(x_i, y_i) \notin C\}}\} p(y_i|x_i)dy_i\right] \end{aligned} \quad (1.25)$$

5. Continuing :

$$\begin{aligned} \mathbb{E}e^{-tN^*(C)} &= \mathbb{E}\left[\prod_{i=1}^{N(S)} \left[1 + \int_M (e^{-t} - 1)\mathbf{1}\{(x_i, y_i) \in C\}p(y_i|x_i)dy_i\right]\right] \\ &\quad \underbrace{\text{use } \mathbf{1}\{(x_i, y_i) \notin C\} = 1 - \mathbf{1}\{(x_i, y_i) \in C\}} \\ &= \mathbb{E}\left[\prod_{i=1}^n \left[1 + \int_S \int_M (e^{-t} - 1)\mathbf{1}\{(x, y) \in C\}p(y|x)dy \frac{\mu(dx)}{\mu(S)} \middle| N(S) = n\right]\right] \\ &\quad \underbrace{\text{Tower again using Poisson-multinomial representation}} \\ &= \mathbb{E}\left[\left(\left(1 + \int_C (e^{-t} - 1)p(y|x)dy \frac{\mu(dx)}{\mu(S)}\right)^{N(S)}\right)\right] \\ &\quad \underbrace{N(S) \sim \text{Pois}(\mu(S))} \end{aligned} \quad (1.26)$$

6. Recall that if $n \sim \text{Pois}(\lambda)$, then

$$\mathbb{E}z^n = \sum_{n=0}^{\infty} z^n \frac{\lambda^n}{n!} e^{-\lambda} = e^{-\lambda} \sum_{n=0}^{\infty} \frac{(z\lambda)^n}{n!} = e^{\lambda(z-1)}.$$

Therefore, this last expectation shows that

$$\begin{aligned} \mathbb{E}e^{-tN^*(C)} &= \exp\left\{\mu(S) \int_C (e^{-t} - 1)p(y|x)dy \frac{\mu(dx)}{\mu(S)}\right\} \\ &= \exp\left\{\int_C (e^{-t} - 1)\mu(dx)p(y|x)dy\right\}, \end{aligned} \quad (1.27)$$

thus $N^*(C) \sim \text{Pois}(\mu^*(C))$.

Example 1: Coloring

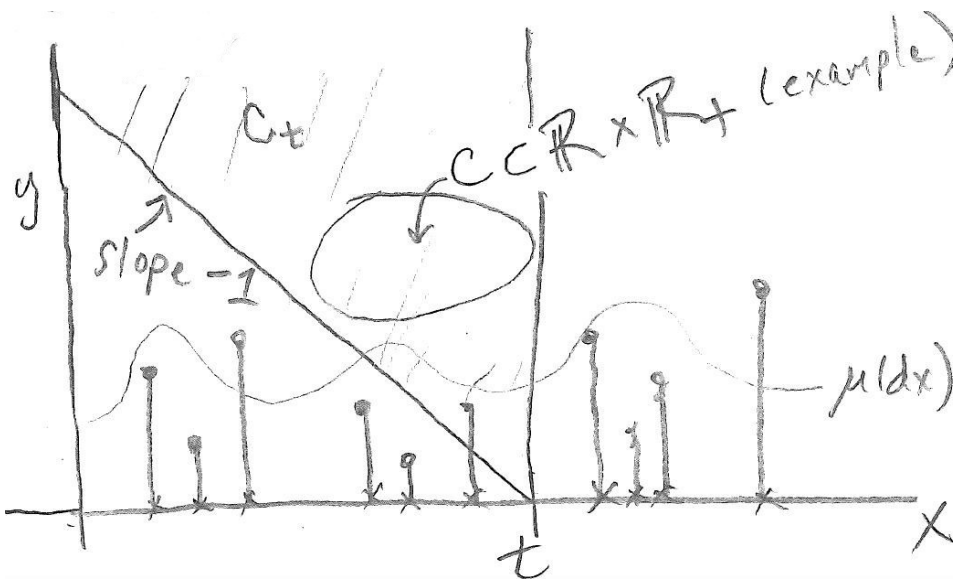
- Let $\Pi \sim \text{PP}(\mu)$ and let an $x \in \Pi$ be randomly colored from among K colors. Denote the color by y with $P(y = i) = p_i$. Then $\Pi^* = \{(x_i, y_i)\}$ is a PP on $S \times \{1, \dots, K\}$ with mean measure $\mu^*(dx \cdot \{y\}) = \mu(dx) \prod_{i=1}^K p_i^{\mathbb{1}(y=i)}$. If we want to restrict Π^* to the i th color, called Π_i^* , then we know that $\Pi_i^* \sim \text{PP}(p_i \mu)$. We can also restrict to two colors, etc. What if $p(y = i|x) = p_i(x)$? Then $\Pi_i^* \sim \text{PP}(\mu(dx)p_i(x))$.

Example 2: Using the extended space

- Imagine we have a 24-hour store and customers arrive according to a PP with mean measure μ on \mathbb{R} (time). In this case, let $\mu(\mathbb{R}) = \infty$, but $\mu([a, b]) < \infty$ for finite a, b (which gives the expected number of arrivals between times a and b). Imagine a customer arriving at time x stays for duration $y \sim p(y|x)$. At time t , what can we say about the customers in the store?
 - $\Pi \sim \text{PP}(\mu)$ and $\Pi^* = \{(x_i, y_i)\}$ for $x_i \in \Pi$ is $\text{PP}(\mu(dx)p(y|x)dy)$ because it's a marked Poisson process.
 - We can construct the marked Poisson process like below. The counting measure $N^*(C) \sim \text{Pois}(\mu^*(C))$, $\mu^* = \mu(dx)p(y|x)dy$.
 - The points in C_t (below) are those that arrive before time t and are still there at time t . It follows that

$$N^*(C_t) \sim \text{Pois} \left(\int_0^t \mu(dx) \int_{t-x}^{\infty} p(y|x)dy \right).$$

$N^*(C_t)$ is the number of customers in the store at time t .



Chapter 2

Completely random measures, Campbell's theorem, gamma process

Poisson process review

- Recall the definition of a Poisson random measure.

PRM definition: Let S be a space and μ a non-atomic (i.e., diffuse, continuous) measure on it (think a positive function). A random measure N on S is a PRM with mean measure μ if

- a) For every subset $A \subset S$, $N(A) \sim \text{Pois}(\mu(A))$.
- b) For disjoint sets A_1, \dots, A_k , $N(A_1), \dots, N(A_k)$ are independent r.v.'s

Poisson process definition: Let $X_1, \dots, X_{N(S)}$ be the $N(S)$ points in N (a random number) each having measure equal to one. Then the *point process* $\Pi = \{X_1, \dots, X_{N(S)}\}$ is a Poisson process, denoted $\Pi \sim \text{PP}(\mu)$.

Recall that to draw this (when $\mu(S) < \infty$) we can

- a) Draw $N(S) \sim \text{Pois}(\mu(S))$
- b) For $i = 1, \dots, N(S)$, draw $X_i \stackrel{iid}{\sim} \mu/\mu(S) \leftarrow$ normalize μ to get a probability measure.

Functions of Poisson processes

- Often models will take the form of a function of an underlying Poisson process: $\sum_{x \in \Pi} f(x)$.

Example (see Sec. 5.3 of Kingman): Imagine that star locations are distributed as $\Pi \sim \text{PP}(\mu)$. They're marked independently with a mass $m \sim p(m)$, giving a marked PP $\Pi^* \sim \text{PP}(\mu \times p dm)$. The gravitational field at, e.g., the origin is

$$\sum_{(x,m) \in \Pi^*} f((x,m)) = \sum_{(x,m) \in \Pi^*} \frac{Gm_x}{\|x\|_2^3} x \quad (G \text{ is a constant from physics})$$

- We can analyze these sorts of problems using PP techniques
- We will be more interested in the context of “completely random measures” in this class.

Functions: The finite case

- Let $\Pi \sim \text{PP}(\mu)$ and $|\Pi| < \infty$ (with probability 1). Let $f(x)$ be a positive function. Let $\mathcal{M} = \sum_{x \in \Pi} f(x)$. We calculate its Laplace transform (for $t < 0$):

$$\mathbb{E}e^{t\mathcal{M}} = \mathbb{E}e^{t\sum_{x \in \Pi} f(x)} = \mathbb{E} \left[\prod_{i=1}^{|\Pi|} e^{tf(x_i)} \right] \leftarrow \text{recall two things (below)} \quad (2.1)$$

Recall:

1. $|\Pi| = N(S) \leftarrow$ Poisson random measure for Π
2. Tower property $\mathbb{E}g(x, y) = \mathbb{E}[\mathbb{E}[g(x, y)|y]]$.

So:

$$\mathbb{E} \left[\prod_{i=1}^{N(S)} e^{tf(x_i)} \right] = \mathbb{E} \left[\mathbb{E} \left[\prod_{i=1}^{N(S)} e^{tf(x_i)} \mid N(S) \right] \right] = \mathbb{E} \left[\mathbb{E} \left[e^{tf(x)} \right]^{N(S)} \right]. \quad (2.2)$$

Since $N(S) \sim \text{Pois}(\mu(S))$, we use the last term to conclude

$$\begin{aligned} \mathbb{E} \left[\prod_{i=1}^{N(S)} e^{tf(x_i)} \right] &= \exp\{\mu(S)(\mathbb{E}[e^{tf(x)}] - 1)\} \\ &= \exp \int_S \mu(dx) (e^{tf(x)} - 1) \\ \left(\text{since } \mathbb{E}e^{tf(x)} = \int_S \frac{\mu(dx)}{\mu(S)} e^{tf(x)} \right) &\nearrow \quad \uparrow \\ &\text{recall that } \mathbb{E}[(e^t)^{N(A)}] = \exp \int_A \mu(dx) (e^t - 1). \\ &\quad (f(x) = 1 \text{ in this case}) \end{aligned}$$

And so for functions $\mathcal{M} = \sum_{x \in \Pi} f(x)$ of Poisson processes Π with an almost sure finite number of points

$$\mathbb{E}e^{t\mathcal{M}} = \exp \int_S \mu(dx) (e^{tf(x)} - 1). \quad (2.3)$$

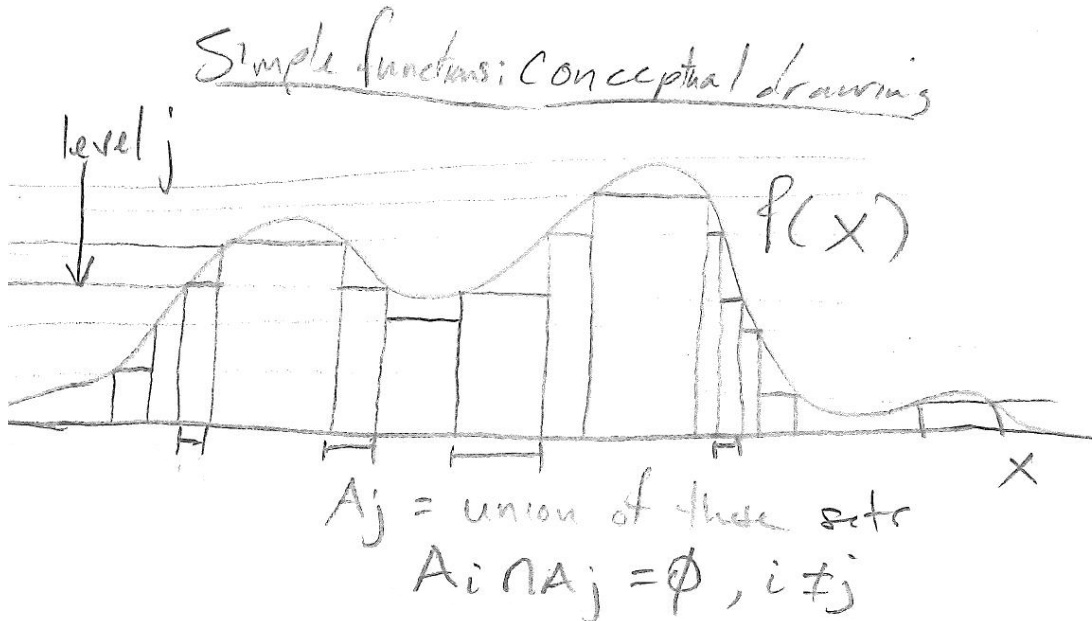
The infinite case

- In the case where $\mu(S) = \infty, N(S) = \infty$ with probability one. We therefore use a different proof technique that gives the same result.
- Approximate $f(x)$ with simple functions: $f_k = \sum_{i=1}^{k2^k} a_i \mathbb{1}_{A_i}$ using $k2^k$ equally-spaced levels in the interval $[0, k]$; $A_i = \{x : f(x) \in [\frac{i-1}{2^k}, \frac{i}{2^k})\}$ and $a_i = \frac{i-1}{2^k}$. In the limit $k \rightarrow \infty$, the width of those levels $\frac{1}{2^k} \searrow 0$, and so $f_k \rightarrow f$.
- For example, if $f(x) = x$, then $f_k(x) = \sum_{i=1}^{k2^k} \frac{i-1}{2^k} \mathbb{1}(x \in [\frac{i-1}{2^k}, \frac{i}{2^k}))$, so $f_k(x) \nearrow x$ as $k \rightarrow \infty$ for finite x .
- Important notation change: $\mathcal{M} = \sum_{x \in \Pi} f(x) \Leftrightarrow \int_S N(dx) f(x)$.
- Approximate f with f_k . Then with the notation change, $\mathcal{M}_k = \sum_{x \in \Pi} f_k(x) \Leftrightarrow \sum_{i=1}^{k2^k} a_i N(A_i)$.
- The Laplace functional is

$$\begin{aligned} \mathbb{E}e^{t\mathcal{M}} &= \lim_{k \rightarrow \infty} \mathbb{E}e^{t\mathcal{M}_k} = \lim_{k \rightarrow \infty} \mathbb{E} \prod_{i=1}^{k2^k} e^{ta_i N(A_i)} && \leftarrow N(A_i) \text{ and } N(A_j) \text{ are independent} \\ &= \lim_{k \rightarrow \infty} \exp \left\{ \sum_{i=1}^{k2^k} \mu(A_i) (e^{ta_i} - 1) \right\} && \leftarrow N(A_i) \sim \text{Pois}(\mu(A_i)) \\ &= \exp \int_S \mu(dx) (e^{tf(x)} - 1) && \leftarrow \text{integral as limit of infinitesimal sums} \end{aligned} \quad (2.4)$$

Mean and variance of \mathcal{M} : Using ideas from moment generating functions, it follows that

$$\begin{aligned} \mathbb{E}(\mathcal{M}) &= \underbrace{\int_S \mu(dx) f(x)}_{= \frac{d}{dt} \mathbb{E}e^{t\mathcal{M}}|_{t=0}}, & \mathcal{V}(\mathcal{M}) &= \underbrace{\int_S \mu(dx) f(x)^2}_{= \frac{d^2}{dt^2} \mathbb{E}e^{t\mathcal{M}}|_{t=0} - \mathbb{E}(\mathcal{M})^2} \end{aligned} \quad (2.5)$$



Finiteness of $\int_S N(dx)f(x)$

- The next obvious question when $\mu(S) = \infty$ (and thus $N(S) = \infty$) is if $\int_S N(dx)f(x) < \infty$. Campbell's theorem gives the necessary and sufficient conditions for this to be true.
- **Campbell's Theorem:** Let $\Pi \sim \text{PP}(\mu)$ and N the PRM. Let $f(x)$ be a non-negative function on S . Then with probability one,

$$\mathcal{M} = \int_S N(dx)f(x) \begin{cases} < \infty & \text{if } \int_S \min(f(x), 1)\mu(dx) < \infty \\ = \infty & \text{otherwise} \end{cases} \quad (2.6)$$

Proof: For $u > 0$, $e^{-u\mathcal{M}} = e^{-u\mathcal{M}}\mathbf{1}(\mathcal{M} < \infty) \nearrow \mathbf{1}(\mathcal{M} < \infty)$ as $u \searrow 0$. By dominated convergence, as $u \searrow 0$

$$\mathbb{E}e^{-u\mathcal{M}} = \mathbb{E}[e^{-u\mathcal{M}}\mathbf{1}(\mathcal{M} < \infty)] \nearrow \mathbb{E}[\mathbf{1}(\mathcal{M} < \infty)] = P(\mathcal{M} < \infty) \quad (2.7)$$

In our case: $P(\mathcal{M} < \infty) = \lim_{u \searrow 0} \exp \int_S \mu(dx)(e^{-uf(x)} - 1)$.

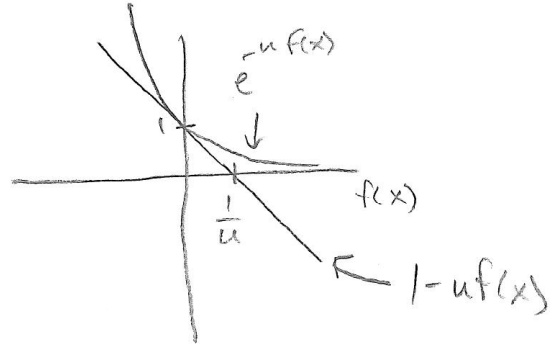
Sufficiency

1. For $P(\mathcal{M} < \infty) = 1$ as in the theorem, we need $\lim_{u \searrow 0} \int_S \mu(dx)(e^{-uf(x)} - 1) = 0$.

2. For $0 < u < 1$ we have

$$1 - f(x) < 1 - uf(x) < e^{-uf(x)} \quad \rightarrow$$

(Figure: $e^{-uf(x)}$ is convex in $f(x)$ and $1 - uf(x)$ is a 1st order Taylor expansion of $e^{-uf(x)}$ at $f(x) = 0$)



3. Therefore $1 - e^{-uf(x)} < f(x)$. Also, $1 - e^{-uf(x)} < 1$ trivially.

4. So:

$$0 \leq \int_S \mu(dx)(1 - e^{-uf(x)}) \leq \int_S \mu(dx) \min(f(x), 1) \quad (2.8)$$

5. If $\int_S \mu(dx) \min(f(x), 1) < \infty$, then by dominated convergence

$$\lim_{u \searrow 0} \int_S \mu(dx)(1 - e^{-uf(x)}) = \int_S \mu(dx) \left(1 - \exp \left\{ \lim_{u \searrow 0} -uf(x) \right\} \right) = 0 \quad (2.9)$$

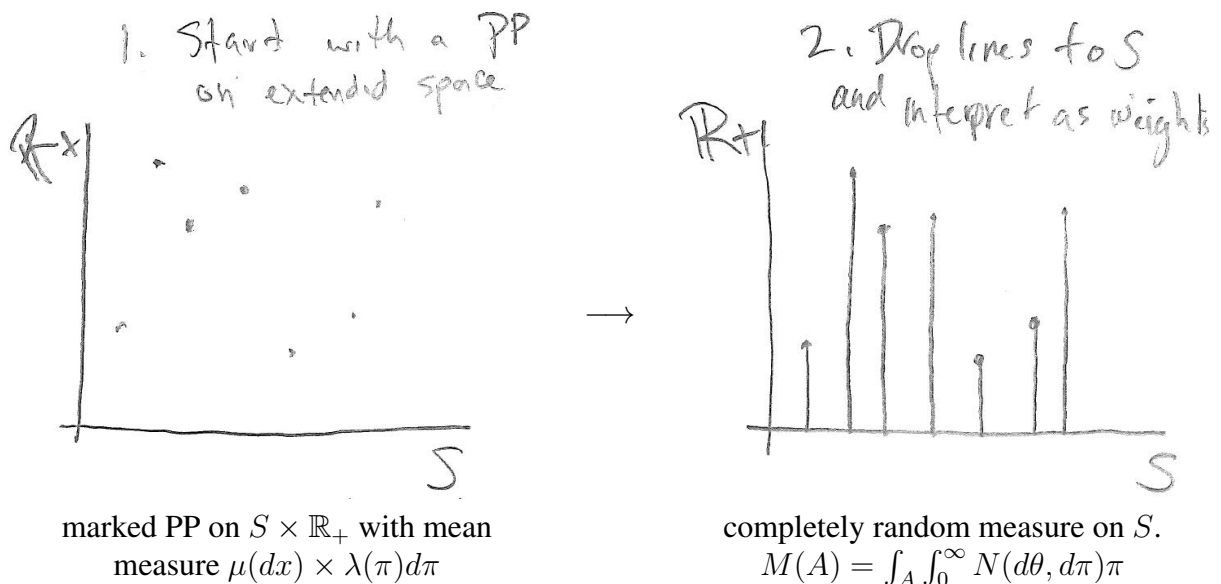
- This proves sufficiency. For necessity we can show that (see, e.g., Cinlar VI.2.13)

$$\int_S \min(f(x), 1)\mu(dx) = \infty \implies \int_S \mu(dx)(1 - e^{-uf(x)}) = \infty \implies \mathbb{E}e^{-u\mathcal{M}} = 0, \forall u \quad (2.10)$$

Completely random measures (CRM)

- **Definition of measure:** The set function μ is a measure on the space S if
 1. $\mu(\emptyset) = 0$
 2. $\mu(A) \geq 0$ for all $A \subset S$
 3. $\mu(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \mu(A_i)$ when $A_i \cap A_j = \emptyset, i \neq j$
- **Definition of completely random measure:** The set function M is a completely random measure on the space S if it satisfies #1 to #3 above and
 1. $M(A)$ is a random variable
 2. $M(A_1), \dots, M(A_k)$ are independent for disjoint sets A_i
- **Example:** Let N be the counting measure associated with $\Pi \sim \text{PP}(\mu)$. It's a CRM.
- We will be interested in the following situation: Let $\Pi \sim \text{PP}(\mu)$ and mark each $\theta \in \Pi$ with a r.v. $\pi \sim \lambda(\pi), \pi \in \mathbb{R}_+$. Then $\Pi^* = \{(\theta, \pi)\}$ is a PP on $S \times \mathbb{R}_+$ with mean measure $\mu(d\theta)\lambda(\pi)d\pi$
- If $N(d\theta, d\pi)$ is the counting measure for Π^* , then $N(C) \sim \text{Pois}(\int_C \mu(d\theta)\lambda(\pi)d\pi)$.
- For $A \subset S$, let $M(A) = \int_A \int_0^\infty N(d\theta, d\pi)\pi$. Then M is a CRM on S .
- M is a special case of sums of functions of Poisson processes with $f(\theta, \pi) = \pi$. Therefore we know that

$$\mathbb{E}e^{tM(A)} = \exp \int_A \int_0^\infty (e^{t\pi} - 1)\mu(d\theta)\lambda(\pi)d\pi. \tag{2.11}$$
- This works both ways: If we define M and show it has this Laplace transform, then we *know* there is a marked Poisson process "underneath" it with mean measure equal to $\mu(d\theta)\lambda(\pi)d\pi$.



Gamma processes

- **Definition:** Let μ be a non-atomic measure on S . Then G is a gamma process if for all $A \subset S$, $G(A) \sim \text{Gam}(\mu(A), c)$ and $G(A_1), \dots, G(A_k)$ are independent for disjoint A_1, \dots, A_k . We write $G \sim \text{GaP}(\mu, c)$. ($c > 0$)
- Before trying to intuitively understand G , let's calculate its Laplace transform. For $t < 0$,

$$\mathbb{E}e^{tG(A)} = \int_0^\infty \frac{c^{\mu(A)}}{\Gamma(\mu(A))} G(A)^{\mu(A)-1} e^{-G(A)(c-t)} dG(A) = \left(\frac{c}{c-t}\right)^{\mu(A)} \quad (2.12)$$

- Manipulate this term as follows (and watch the magic!)

$$\begin{aligned} \left(\frac{c}{c-t}\right)^{\mu(A)} &= \exp\left\{-\mu(A) \ln \frac{c-t}{c}\right\} \\ &= \exp\left\{-\mu(A) \int_c^{c-t} \frac{1}{s} ds\right\} \\ &= \exp\left\{-\mu(A) \int_c^{c-t} ds \int_0^\infty e^{-\pi s} d\pi\right\} \\ &= \exp\left\{-\mu(A) \int_0^\infty d\pi \int_c^{c-t} e^{-\pi s} ds\right\} \quad (\text{switched integrals}) \\ &= \exp\left\{\mu(A) \int_0^\infty (e^{t\pi} - 1) \pi^{-1} e^{-c\pi} d\pi\right\} \end{aligned} \quad (2.13)$$

Therefore, G has an underlying Poisson random measure on $S \times \mathbb{R}_+$

$$G(A) = \int_A \int_0^\infty N(d\theta, d\pi) \pi, \quad N(d\theta, d\pi) \sim \text{Pois}(\mu(d\theta) \pi^{-1} e^{-c\pi} d\pi) \quad (2.14)$$

- The mean measure of N is $\mu(d\theta) \pi^{-1} e^{-c\pi} d\pi$ on $S \times \mathbb{R}_+$. We can use this to answer questions about $G \sim \text{GaP}(\mu, c)$ using the Poisson process perspective.

1. How many total atoms? $\int_S \int_0^\infty \mu(d\theta) \pi^{-1} e^{-c\pi} d\pi = \infty \Rightarrow$ infinite # w.p. 1
(Tells us that there are an infinite number of points in any subset $A \subset S$ that have nonzero mass according to G)
2. How many atoms $\geq \epsilon > 0$? $\int_S \int_\epsilon^\infty \mu(d\theta) \pi^{-1} e^{-c\pi} d\pi < \infty \Rightarrow$ finite # w.p. 1
(w.r.t. #1, further tells us only a finite number have mass greater than ϵ)
3. Campbell's theorem: $f(\theta, \pi) = \pi \rightarrow \int_S \int_0^\infty \min(\pi, 1) \mu(d\theta) \pi^{-1} e^{-c\pi} d\pi < \infty$, therefore

$$G(A) = \int_A \int_0^\infty N(d\theta, d\pi) \pi < \infty \text{ w.p. } 1$$

(Tells us if we summed up the infinite number of nonzero masses in any set A , we would get a finite number even though we have an infinite number of nonzero things to add)

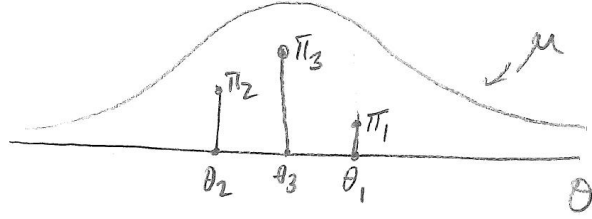
- **Aside:** We already knew #3 by definition of G , but this isn't always the order in which CRMs are defined. Imagine starting the definition with a mean measure on $S \times \mathbb{R}_+$.
- #3 shouldn't feel mysterious at all. Consider $\sum_{n=1}^\infty \frac{1}{n^2}$. It's finite, but for each n , $\frac{1}{n^2} > 0$ and there are an infinite number of settings for n . "Infinite jump processes" such as the gamma process replace deterministic sequences like $(1, 1/2^2, 1/3^2, \dots)$ with something random.

Gamma process as a limiting case

- Is there a more intuitive way to understand the gamma process?
- **Definition:** Let μ be a non-atomic measure on S and $\mu(S) < \infty$. Let $\pi_i \stackrel{iid}{\sim} \text{Gam}(\frac{\mu(S)}{K}, c)$ and $\theta_i \stackrel{iid}{\sim} \mu/\mu(S)$ for $i = 1, \dots, K$. If $G_k = \sum_{i=1}^K \pi_i \delta_{\theta_i}$, then $\lim_{K \rightarrow \infty} G_K = G \sim \text{GaP}(\mu, c)$.

Picture: In the limit $K \rightarrow \infty$, we have more and more atoms with smaller and smaller weights

$$G_K(A) = \sum_{i=1}^K \pi_i \delta_{\theta_i}(A) = \sum_{i=1}^K \pi_i \mathbb{1}(\theta_i \in A)$$



Proof: Use the Laplace transform

$$\begin{aligned} \mathbb{E}e^{tG_K(A)} &= \mathbb{E}e^{t \sum_{i=1}^K \pi_i \mathbb{1}(\theta_i \in A)} = \mathbb{E} \prod_{i=1}^K e^{t\pi_i \mathbb{1}(\theta_i \in A)} = \mathbb{E}[e^{t\pi \mathbb{1}(\theta \in A)}]^K \\ &= \mathbb{E} [e^{t\pi} \mathbb{1}(\theta \in A) + \mathbb{1}(\theta \notin A)]^K \\ &= [\mathbb{E}[e^{t\pi}]P(\theta \in A) + P(\theta \notin A)]^K \\ &= \left[\left(\frac{c}{c-t}\right)^{\frac{\mu(S)}{K}} \frac{\mu(A)}{\mu(S)} + 1 - \frac{\mu(A)}{\mu(S)} \right]^K \\ &= \left[1 + \frac{\mu(A)}{\mu(S)} \left(\left(\frac{c}{c-t}\right)^{\frac{\mu(S)}{K}} - 1 \right) \right]^K \\ &= \left[1 + \frac{\mu(A)}{\mu(S)} \left(\sum_{n=1}^{\infty} \frac{(\ln \frac{c}{c-t})^n}{n!} \left(\frac{\mu(S)}{K}\right)^n \right) \right]^K \leftarrow \text{(exponential power series)} \\ &= \left[1 + \frac{\mu(A)}{\mu(S)} \left(\frac{\mu(S)}{K} \ln \frac{c}{c-t} + O(1/K^2) \right) \right]^K \end{aligned} \tag{2.15}$$

- In the limit $K \rightarrow \infty$, this last equation converges to

$$\exp \left\{ \frac{\mu(A)}{\mu(S)} \mu(S) \ln \frac{c}{c-t} \right\} = \left(\frac{c}{c-t} \right)^{\mu(A)}$$

(recall that $\lim_{K \rightarrow \infty} (1 + \frac{a}{K} + O(K^{-2}))^K = e^a$)

- This is the Laplace transform of a $\text{Gam}(\mu(A), c)$ random variable.
- Therefore, $G_K(A) \rightarrow G(A) \sim \text{Gam}(\mu(A), c)$, which we've already defined and analyzed as a gamma process.

Chapter 3

Beta processes and the Poisson process

A sparse coding latent factor model

- We have a $d \times n$ matrix Y . We want to factorize it as follows:

$$d \begin{bmatrix} Y \end{bmatrix} \approx \begin{bmatrix} \Theta \end{bmatrix} \times \left(\begin{bmatrix} W \end{bmatrix} \circ \begin{bmatrix} Z \end{bmatrix} \right)$$

where

$$\left. \begin{array}{l} \theta_i \sim p(\theta) \quad i = 1, \dots, K \\ w_j \sim p(w) \quad j = 1, \dots, n \\ z_j \in \{0, 1\}^K \quad j = 1, \dots, n \end{array} \right\} \begin{array}{l} \text{“sparse coding” because each vector } Y_j \\ \text{only possesses the columns of } \Theta \text{ indicated} \\ \text{by } z_j \text{ (want } \sum_i z_{ji} \ll K \text{)} \end{array}$$

- Example: Y could be
 - a) gene data of n (or d) people,
 - b) patches extracted from an image for denoising (called “dictionary learning”)
- We want to define a “Bayesian nonparametric” prior for this problem. By this we mean that
 1. The prior can allow $K \rightarrow \infty$ and remain well defined
 2. As $K \rightarrow \infty$, the “effective rank” is finite (and relatively small)
 3. The model somehow learns this rank from the data during inference (not discussed)

A “beta sieves” prior

- Let $\theta_i \sim \mu/\mu(S)$ and w_j be drawn as above. Continue this generative model by letting

$$z_{ji} \stackrel{iid}{\sim} \text{Bern}(\pi_i), j = 1, \dots, n \quad (3.1)$$

$$\pi_i \sim \text{Beta}\left(\alpha \frac{\gamma}{K}, \alpha \left(1 - \frac{\gamma}{K}\right)\right), i = 1, \dots, K \quad (3.2)$$

- The set (θ_i, π_i) are paired. π_i gives the probability an observation picks θ_i . Notice that we expect $\pi_i \rightarrow 0$ as $K \rightarrow \infty$.
- Construct a completely random measure $H_K = \sum_{i=1}^K \pi_i \delta_{\theta_i}$.
- We want to analyze what happens when $K \rightarrow \infty$. We’ll see that it converges to a *beta process*.

Asymptotic analysis of beta sieves

- We have that $H_K = \sum_{i=1}^K \pi_i \delta_{\theta_i}$, $\pi_i \stackrel{iid}{\sim} \text{Beta}(\alpha\gamma/K, \alpha(1 - \gamma/K))$, $\theta_i \stackrel{iid}{\sim} \mu/\mu(S)$, where $\gamma = \mu(S) < \infty$. We want to understand $\lim_{K \rightarrow \infty} H_K$.
- Look at the Laplace transform of $H_K(A)$. Let $H(A) = \lim_{K \rightarrow \infty} H_K(A)$. Then

$$\mathbb{E}e^{tH(A)} = \lim_{K \rightarrow \infty} \mathbb{E}e^{tH_K(A)} = \underbrace{\lim_{K \rightarrow \infty} \mathbb{E}e^{t \sum_{i=1}^K \pi_i \mathbb{1}(\theta_i \in A)}}_{\text{sum} \rightarrow \text{product and use i.i.d. fact}} = \lim_{K \rightarrow \infty} \mathbb{E}[e^{t\pi \mathbb{1}(\theta \in A)}]^K \quad (3.3)$$

- Focus on $\mathbb{E}e^{t\pi \mathbb{1}(\theta \in A)}$ for a particular K -level approximation. We have the following (long) sequence of equalities:

$$\begin{aligned} \mathbb{E}e^{t\pi \mathbb{1}(\theta \in A)} &= \mathbb{E}[e^{t\pi \mathbb{1}(\theta \in A)} + \mathbb{1}(\theta \notin A)] = P(\theta \in A)\mathbb{E}e^{t\pi} + P(\theta \notin A) \\ &= 1 + \frac{\mu(A)}{\mu(S)} \left(\mathbb{E}e^{t\pi} - 1 \right) \leftarrow \mathbb{E}e^{t\pi} = 1 + \sum_{s=1}^{\infty} \frac{t^s}{s!} \prod_{r=0}^{s-1} \frac{\frac{\alpha\gamma}{K} + r}{\alpha + r} \\ &= 1 + \frac{\mu(A)}{\mu(S)} \sum_{s=1}^{\infty} \frac{t^s}{s!} \prod_{r=0}^{s-1} \frac{\frac{\alpha\gamma}{K} + r}{\alpha + r} \leftarrow \text{plugging in } \uparrow \\ &= 1 + \frac{\mu(A)}{K} \sum_{s=1}^{\infty} \frac{t^s}{s!} \prod_{r=1}^{s-1} \frac{r}{\alpha + r} + O\left(\frac{1}{K^2}\right) \leftarrow \text{separate out } r=0 \\ &= 1 + \frac{\mu(A)}{K} \sum_{s=1}^{\infty} \frac{t^s}{s!} \frac{\alpha\Gamma(\alpha)\Gamma(s)}{\Gamma(\alpha+s)} + O\left(\frac{1}{K^2}\right) \leftarrow \text{since } \Gamma(\alpha+1) = \alpha\Gamma(\alpha) \\ &= 1 + \frac{\mu(A)}{K} \sum_{s=1}^{\infty} \frac{t^s}{s!} \int_0^1 \alpha\pi^{s-1}(1-\pi)^{\alpha-1} d\pi + O\left(\frac{1}{K^2}\right) \leftarrow \text{normalizer of Beta}(s, \alpha) \\ &= 1 + \frac{\mu(A)}{K} \int_0^1 \sum_{s=1}^{\infty} \left(\frac{t\pi}{s!}\right) \alpha\pi^{-1}(1-\pi)^{\alpha-1} d\pi + O\left(\frac{1}{K^2}\right) \\ &= 1 + \frac{\mu(A)}{K} \int_0^1 (e^{t\pi} - 1) \alpha\pi^{-1}(1-\pi)^{\alpha-1} d\pi + O\left(\frac{1}{K^2}\right) \end{aligned}$$

Again using $\lim_{K \rightarrow \infty} (1 + a/K + O(K^{-2}))^K = e^a$, it follows that

$$\lim_{K \rightarrow \infty} \mathbb{E}[e^{t\pi \mathbf{1}(\theta \in A)}]^K = \exp \left\{ \mu(A) \int_0^1 (e^{t\pi} - 1) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi \right\}. \quad (3.4)$$

• We therefore know that

1. H is a completely random measure.
2. It has an associated underlying Poisson random measure $N(d\theta, d\pi)$ on $S \times [0, 1]$ with mean measure $\mu(d\theta) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi$.
3. We can write $H(A)$ as $\int_A \int_0^1 N(d\theta, d\pi) \pi$.

Beta process (as a CRM)

- **Definition:** Let $N(d\theta, d\pi)$ be a Poisson random measure on $S \times [0, 1]$ with mean measure $\mu(d\theta) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi$, where μ is a non-atomic measure. Define the CRM $H(A)$ as $\int_A \int_0^1 N(d\theta, d\pi) \pi$. Then H is called a beta process, $H \sim \text{BP}(\alpha, \mu)$ and

$$\mathbb{E}e^{tH(A)} = \exp \left\{ \mu(A) \int_0^1 (e^{t\pi} - 1) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi \right\}.$$

- (We just saw how we can think of H as the limit of a finite collection of random variables. This time we're just starting from the definition, which we could proceed to analyze regardless of the beta sieves discussion above.)

- Properties of H : Since H has a Poisson process representation, we can use the mean measure to calculate its properties (and therefore the asymptotic properties of the beta sieves approximation).

- Finiteness: Using Campbell's theorem, $H(A)$ is finite with probability one, since

$$\int_A \int_0^1 \underbrace{\min(\pi, 1)}_{=\pi} \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi \mu(d\theta) = \mu(A) < \infty \quad (\text{by assumption about } \mu) \quad (3.5)$$

- Infinite jump process: $H(A)$ is constructed from an infinite number of jumps, almost all infinitesimally small, since

$$N(A \times (0, 1]) \sim \text{Pois} \left(\mu(A) \int_0^1 \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi \right) = \text{Pois}(\infty) \quad (3.6)$$

- Finite number of "big" jumps: There are only a finite number of jumps greater than any $\epsilon > 0$ since

$$N(A \times [\epsilon, 1]) \sim \text{Pois} \left(\underbrace{\mu(A) \int_{\epsilon}^1 \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi}_{< \infty \text{ for } \epsilon > 0} \right) \quad (3.7)$$

As $\epsilon \rightarrow 0$, the value in the Poisson goes to infinity, so the infinite jump process arises in this limit. Since the integral over the magnitudes is finite, this infinite number of atoms is being introduced in a "controlled" way as a function of ϵ (i.e., not "too quickly")

- **Reminder and intuitions:** All of these properties are over instantiations of a beta process, and so all statements are made with probability one.
 - It’s not absurd to talk about beta processes that don’t have an infinite number of jumps, or integrate to something infinite (“not absurd” in the way that it is absurd to talk about a negative value drawn from a beta distribution).
 - The support of the beta process includes these events, but they have probability zero, so any $H \sim \text{BP}(\alpha, \mu)$ is guaranteed to have the properties discussed above.
 - It’s easy to think of H as one random variable, but as the beta sieves approximation shows, H is really a collection of an infinite number of random variables.
 - The statements we are making about H above aren’t like asking whether a beta random variable is greater than 0.5. They are larger scale statements about properties of this infinite collection of random variables as a whole.

- Another definition of the beta process links it to the beta distribution and our finite approximation:

- **Definition II:** Let μ be a non-atomic measure on S . For all infinitesimal sets $d\theta \in S$, let

$$H(d\theta) \sim \text{Beta}\{\alpha\mu(d\theta), \alpha(1 - \mu(d\theta))\},$$

then $H \sim \text{BP}(\alpha, \mu)$.

- We aren’t going to prove this, but the proof is actually very similar to the beta sieves proof.
- Note the difference from the gamma process, where $G(A) \sim \text{Gam}(\mu(A), c)$ for any $A \subset S$. The beta distribution only comes in the infinitesimal limit. That is

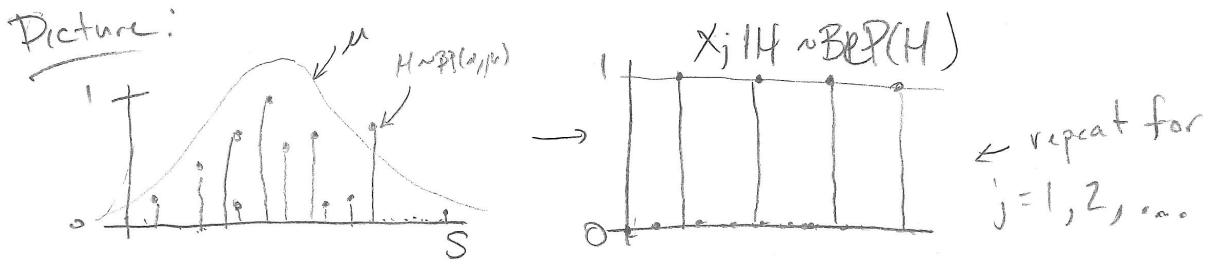
$$H(A) \not\sim \text{Beta}\{\alpha\mu(A), \alpha(1 - \mu(A))\},$$

when $\mu(A) > 0$. Therefore, we can only write beta distributions on things that equal zero with probability one. . . Compare this with the limit of the beta sieves prior.

- **Observation:** While $\mu(\{\theta\}) = \mu(d\theta) = 0$, $\int_A \mu(\{\theta\})d\theta = 0$ but $\int_A \mu(d\theta) = \mu(A) > 0$.
- This is a major difference between a measure and a function: μ is a measure, not a function. It also seems to me a good example of why these additional concepts and notations are necessary, e.g., why we can’t just combine things like $\mu(A) = \int_A p(\theta)d\theta$ into one single notation, but instead talk about the “measure” μ and it’s associated density $p(\theta)$ such that $\mu(d\theta) = p(\theta)d\theta$.
- This leads to discussions involving the Radon-Nikodym theorem, etc. etc.
- The level of our discussion stops at an appreciation for why these types of theorems exist and are necessary (as overly-obsessive as they may feel the first time they’re encountered), but we won’t re-derive them.

Bernoulli process

- The Bernoulli process is constructed from the infinite limit of the “z” sequence in the process $z_{ji} \sim \text{Bern}(\pi_i), \pi_i \stackrel{iid}{\sim} \text{Beta}(\alpha\gamma/K, \alpha(1 - \gamma/K)), i = 1, \dots, K$. The random measure $X_j^{(K)} = \sum_{i=1}^K z_{ji} \delta_{\theta_i}$ converges to a “Bernoulli process” as $K \rightarrow \infty$.
- Definition: Let $H \sim \text{BP}(\alpha, \mu)$. For each atom of H (the θ for which $H(\{\theta\}) > 0$), let $X_j(\{\theta\})|H \stackrel{iid}{\sim} \text{Bern}(H(\{\theta\}))$. Then X_j is a Bernoulli process, denoted $X_j|H \sim \text{BeP}(H)$.
- Observation: We know from the Poisson process that H has an infinite number of locations θ where $H(\{\theta\}) > 0$. Therefore, X is infinite as well.



Some questions about X

1. How many 1's in X_j ?
2. For $X_1, \dots, X_n|H \sim \text{BeP}(H)$, how many total locations are there with at least one X_j equaling one there (marginally speaking, with H integrated out)?

i.e., what is $|\{\theta : \sum_{j=1}^n X_j(\{\theta\}) > 0\}|$

- The Poisson process representation of the BP makes calculating this relatively easy. We start by observing that the X_j are marking the atoms, and so we have a marked Poisson process (or “doubly marked” since we can view (θ, π) as a marked PP as well).

Beta process marked by a Bernoulli process

- Definition: Let $\Pi \sim \text{PP}(\mu(d\theta)\alpha\pi^{-1}(1-\pi)^{\alpha-1}d\pi)$ on $S \times [0, 1]$ be a Poisson process underlying a beta process. For each $(\theta, \pi) \in \Pi$ draw a binary vector $z \in \{0, 1\}^n$ where $z_i|\pi \stackrel{iid}{\sim} \text{Bern}(\pi)$ for $i = 1, \dots, n$. Denote the distribution on z as $Q(z|\pi)$. Then $\Pi^* = \{(\theta, \pi, z)\}$ is a marked Poisson process with mean measure $\mu(d\theta)\alpha\pi^{-1}(1-\pi)^{\alpha-1}d\pi Q(z|\pi)$.
- There is therefore a Poisson process underlying the joint distribution of the hierarchical process

$$H \sim \text{BP}(\alpha, \mu), \quad X_i|H \stackrel{iid}{\sim} \text{BeP}(H), \quad i = 1, \dots, n.$$

- We next answer the two questions about X asked above, starting with the second one.

Question: What is $K_n^+ = \left| \{ \theta : \sum_{j=1}^n X_j(\{\theta\}) > 0 \} \right|$?

Answer: The transition distribution $Q(z|\pi)$ gives the probability of a vector z at a *particular* location $(\theta, \pi) \in \Pi$ (notice Q doesn't depend on θ).

All we care about is whether $z \in C = \{0, 1\}^n \setminus \vec{0}$ (i.e., has a 1 in it)

We make the following observations:

– The probability $Q(C|\pi) = P(z \in C|\pi) = 1 - P(z \notin C|\pi) = 1 - (1 - \pi)^n$.

– If we restrict the marked PP to C , we get the distribution on the value of K_n^+ :

$$K_n^+ = N(S, [0, 1], C) \sim \text{Pois} \left(\int_S \int_0^1 \mu(d\theta) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi \underbrace{Q(C|\pi)}_{= 1 - (1-\pi)^n} \right). \quad (3.8)$$

– It's worth stopping to remember that N is a *counting* measure, and think about what exactly $N(S, [0, 1], C)$ is counting.

* $N(S, [0, 1], C)$ is asking for the number of times event C happens (an event related to z), not caring about what the corresponding θ or π are (hence the S and $[0, 1]$).

* i.e., it's counting the thing we're asking for, K_n^+ .

– We can show that $1 - (1 - \pi)^n = \sum_{i=0}^{n-1} \pi(1 - \pi)^i \leftarrow$ geometric series

– It follows that

$$\begin{aligned} \int_S \int_0^1 \mu(d\theta) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi (1 - (1 - \pi)^n) &= \mu(S) \sum_{i=0}^{n-1} \alpha \int_0^1 (1 - \pi)^{\alpha+i-1} d\pi \\ &= \mu(S) \sum_{i=0}^{n-1} \frac{\alpha \Gamma(1) \Gamma(\alpha + i)}{\Gamma(\alpha + i + 1)} \\ &= \sum_{i=0}^{n-1} \frac{\alpha \mu(S)}{\alpha + i} \end{aligned} \quad (3.9)$$

• Therefore $K_n^+ \sim \text{Pois} \left(\sum_{i=0}^{n-1} \frac{\alpha \mu(S)}{\alpha + i} \right)$.

• Notice that as $n \rightarrow \infty$, $K_n^+ \rightarrow \infty$ with probability one, and that $\mathbb{E}K_n^+ \approx \alpha \mu(S) \ln n$.

• Also notice that we get the answer to the first question for free. Since X_j are i.i.d., we can treat each one marginally as if it were the first one.

• If $n = 1$, $X(S) \sim \text{Pois}(\mu(S))$. That is, the number of ones in each Bernoulli process is $\text{Pois}(\mu(S))$ -distributed.

Chapter 4

Beta processes and size-biased constructions

The beta process

- **Definition (review):** Let $\alpha > 0$ and μ be a finite non-atomic measure on S . Let $C \in S \times [0, 1]$ and N be a Poisson random measure with $N(C) \sim \text{Pois}(\int_C \mu(d\theta) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi)$. For $A \subset S$ define $H(A) = \int_A \int_0^1 N(d\theta, d\pi) \pi$. Then H is a beta process, denoted $H \sim \text{BP}(\alpha, \mu)$.

Intuitive picture (review)

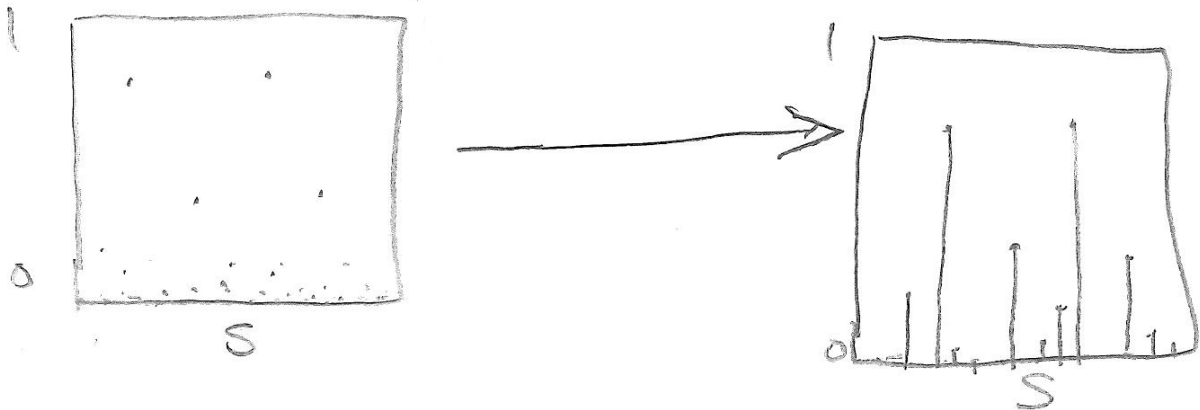


Figure 4.1 (left) Poisson process (right) CRM constructed from Poisson process. If $(d\theta, d\pi)$ is a point in the PP, $N(d\theta, d\pi) = 1$ and $N(d\theta, d\pi)\pi = \pi$. $H(A)$ is adding up π 's in the set $A \times [0, 1]$.

Drawing from this prior

- In general, we know that if $\Pi \sim \text{PP}(\mu)$, we can draw $N(S) \sim \text{Pois}(\mu(S))$ and $X_1, \dots, X_{N(S)} \stackrel{iid}{\sim} \mu/\mu(S)$ and construct Π from the X_i 's.
- Similarly, we have the reverse property that if $N \sim \text{Pois}(\gamma)$ and $X_1, \dots, X_N \stackrel{iid}{\sim} p(X)$, then the set $\Pi = \{X_1, \dots, X_N\} \sim \text{PP}(\gamma p(X) dX)$. (This inverse property will be useful later.)

- Since $\int_S \int_0^1 \mu(d\theta) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi = \infty$, this approach obviously won't work for drawing $H \sim \text{BP}(\alpha, \mu)$.
- The method of partitioning $[0, 1]$ and drawing $N(S \times (a, b])$ followed by

$$\theta_i^{(a)} \sim \mu/\mu(S), \quad \pi_i^{(a)} \sim \frac{\alpha \pi^{-1} (1 - \pi)^{\alpha-1}}{\int_a^b \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi} \mathbf{1}(a < \pi_i^{(a)} \leq b) \quad (4.1)$$

is possible (using the Restriction theorem from Lecture 1, independence of PP's on disjoint sets, and the first bullet of this section), but not as useful for Bayesian models.

- The goal is to find *size-biased* representations for H that are more straightforward. (i.e., that involve sampling from standard distributions, which will hopefully make inference easier)

Size-biased representation I (a “restricted beta process”)

- **Definition:** Let $\alpha = 1$ and μ be a non-atomic measure on S with $\mu(S) = \gamma < \infty$. Generate the following independent set of random variables

$$V_i \stackrel{iid}{\sim} \text{Beta}(\gamma, 1), \quad \theta_i \stackrel{iid}{\sim} \mu/\mu(S), \quad i = 1, 2, \dots \quad (4.2)$$

Let $H = \sum_{i=1}^{\infty} \left(\prod_{j=1}^i V_j \right) \delta_{\theta_i}$. Then $H \sim \text{BP}(1, \mu)$.

Proof: The proof uses the limiting case of the following finite approximation

- Let $\pi_i \sim \text{Beta}(\frac{\gamma}{K}, 1)$, $\theta_i \sim \mu/\mu(S)$ for $i = 1, \dots, K$. Let $H_K = \sum_{i=1}^K \pi_i \delta_{\theta_i}$. Then $\lim_{K \rightarrow \infty} H_K \sim \text{BP}(1, \mu)$. The proof is similar to the one last lecture.

Question 1: As $K \rightarrow \infty$, what is $\pi_{(1)} = \max\{\pi_1, \dots, \pi_K\}$?

Answer: Look at the CDF's. We want the function $P(\pi_{(1)} < V_1)$ for a $V_1 \in [0, 1]$. Because the π_i are independent,

$$P(\pi_{(1)} < V_1) = \lim_{K \rightarrow \infty} P(\pi_1 < V_1, \dots, \pi_K < V_1) = \lim_{K \rightarrow \infty} \prod_{i=1}^K P(\pi_i < V_1) \quad (4.3)$$

- $P(\pi_i < V_1) = \int_0^{V_1} \frac{\gamma}{K} \pi_i^{\frac{\gamma}{K}-1} d\pi_i = V_1^{\frac{\gamma}{K}}$
- $\lim_{K \rightarrow \infty} \prod_{i=1}^K P(\pi_i < V_1) = V_1^\gamma$
- Therefore, $\pi_{(1)} = V_1, \quad V_1 \sim \text{Beta}(\gamma, 1)$

Question 2: What is the second largest, denoted $\pi_{(2)} = \lim_{K \rightarrow \infty} \max\{\pi_1, \dots, \pi_K\} \setminus \{\pi_{(1)}\}$?

Answer: This is a little more complicated, but answering how to get $\pi_{(2)}$ shows how to get the remaining $\pi_{(i)}$.

$$\begin{aligned}
 P(\pi_{(2)} < t | \pi_{(1)} = V_1) &= \prod_{\pi_i \neq \pi_{(1)}} P(\pi_i < t | \pi_i < V_1) \leftarrow \text{condition is each } \pi_i < \pi_{(1)} = V_1 \\
 &= \prod_{\pi_i \neq \pi_{(1)}} \frac{P(\pi_i < t, \pi_i < V_1)}{P(\pi_i < V_1)} \\
 &= \prod_{\pi_i \neq \pi_{(1)}} \frac{P(\pi_i < t)}{P(\pi_i < V_1)} \leftarrow \text{since } t < V_1, \text{ first event contains second} \\
 &= \lim_{K \rightarrow \infty} \prod_{\pi_i \neq \pi_{(1)}} \frac{\int_0^t \frac{\gamma}{K} \pi_i^{\frac{\gamma}{K}-1} d\pi_i}{\int_0^{V_1} \frac{\gamma}{K} \pi_i^{\frac{\gamma}{K}-1} d\pi_i} \\
 &= \lim_{K \rightarrow \infty} \left[\left(\frac{t}{V_1} \right)^{\frac{\gamma}{K}} \right]^{K-1} = \left(\frac{t}{V_1} \right)^\gamma \tag{4.4}
 \end{aligned}$$

– So the density $p(\pi_{(2)} | \pi_{(1)} = V_1) = V_1^{-1} \gamma \left(\frac{\pi_{(2)}}{V_1} \right)^{\gamma-1}$. $\pi_{(2)}$ has support $[0, V_1]$.

– Change of variables: $V_2 := \pi_{(2)}/V_1 \rightarrow \pi_{(2)} = V_1 V_2, d\pi_{(2)} = V_1 dV_2$.

– Plugging in, $p(V_2 | \pi_{(1)} = V_1) = V_1^{-1} \gamma V_2^{\gamma-1} \cdot \underbrace{V_1}_{\text{Jacobian}} = \gamma V_2^{\gamma-1} = \text{Beta}(\gamma, 1)$

– The above calculation has shown two things:

1. V_2 is independent of V_1 (this is an instance of a “neutral-to-the-right process”)
2. $V_2 \sim \text{Beta}(\gamma, 1)$

- Since $\pi_{(2)} | \{\pi_{(1)} = V_1\} = V_1 V_2$ and V_1, V_2 are independent, we can get the value of $\pi_{(2)}$ using previously drawn V_1 and then drawing V_2 from $\text{Beta}(\gamma, 1)$ distributions.
- The same exact reasoning follows for $\pi_{(3)}, \pi_{(4)}, \dots$
- For example, for $\pi_{(3)}$, we have $P(\pi_{(3)} < t | \pi_{(2)} = V_1 V_2, \pi_{(1)} = V_1) = P(\pi_{(3)} < t | \pi_{(2)} = V_1 V_2)$ because conditioning on $\pi_{(2)} = V_1 V_2$ restricts $\pi_{(3)}$ to also satisfy condition of $\pi_{(1)}$.
- In other words, if we force $\pi_{(3)} < \pi_{(2)}$ by conditioning, we get the additional requirement $\pi_{(3)} < \pi_{(1)}$ for free, so we can condition on the $\pi_{(i)}$ immediately before.
- Think of $V_1 V_2$ as a single non-random (i.e., already known) value by the time we get to $\pi_{(3)}$. We can exactly follow the above sequence after making the correct substitutions and re-indexing.

Size-biased representation II

Definition: Let $\alpha > 0$ and μ be a non-atomic measure on S with $\mu(S) < \infty$. Generate the following random variables:

$$\begin{aligned} C_i &\sim \text{Pois}\left(\frac{\alpha\mu(S)}{\alpha+i-1}\right), \quad i = 1, 2, \dots \\ \pi_{ij} &\sim \text{Beta}(1, \alpha+i-1), \quad j = 1, \dots, C_i \\ \theta_{ij} &\sim \mu/\mu(S), \quad j = 1, \dots, C_i \end{aligned} \quad (4.5)$$

Define $H = \sum_{i=1}^{\infty} \sum_{j=1}^{C_i} \pi_{ij} \delta_{\theta_{ij}}$. Then $H \sim \text{BP}(\alpha, \mu)$.

Proof:

- We can use Poisson processes to prove this. This is a good example of how easy a proof can become when we recognize a hidden Poisson process and calculate its mean measure.
 - Let $H_i = \sum_{j=1}^{C_i} \pi_{ij} \delta_{\theta_{ij}}$. Then the set $\Pi_i = \{(\theta_{ij}, \pi_{ij})\}$ is a Poisson process because it contains a Poisson-distributed number of i.i.d. random variables.
 - As a result, the mean measure of Π_i is

$$\underbrace{\frac{\alpha\mu(S)}{\alpha+i-1}}_{\text{Poisson \# part}} \times \underbrace{(\alpha+i-1)(1-\pi)^{\alpha+i-2} d\pi}_{\text{distribution on } \pi} \times \underbrace{\mu(d\theta)/\mu(S)}_{\text{distribution on } \theta} \quad (4.6)$$

We can simplify this to $\alpha\mu(d\theta)(1-\pi)^{\alpha+i-2} d\pi$. We can justify this with the marking theorem (π marks θ), or just thinking about the joint distribution of (θ, π) .

- $H = \sum_{i=1}^{\infty} H_i$ by definition. Equivalently $\Pi = \bigcup_{i=1}^{\infty} \Pi_i$.
- By the superposition theorem, we know that Π is a Poisson process with mean measure equal to the sum of the mean measures of each Π_i .
- We can calculate this directly:

$$\sum_{i=1}^{\infty} \alpha\mu(d\theta)(1-\pi)^{\alpha+i-2} d\pi = \alpha\mu(d\theta)(1-\pi)^{\alpha-2} \underbrace{\sum_{i=1}^{\infty} (1-\pi)^i d\pi}_{= \frac{1-\pi}{\pi}} \quad (4.7)$$

- Therefore, we've shown that Π is a Poisson process with mean measure

$$\alpha\pi^{-1}(1-\pi)^{\alpha-1} d\pi\mu(d\theta).$$

- In other words, this second size-biased construction is the CRM constructed from integrating a PRM with this mean measure against the function $f(\theta, \pi) = \pi$ along the π dimension. This is the definition of a beta process.

Size-biased representation III

- **Definition:** Let $\alpha > 0$ and μ be a non-atomic measure on S with $\mu(S) < \infty$. The following is a constructive definition of $H \sim \text{BP}(\alpha, \mu)$.

$$C_i \sim \text{Pois}(\mu(S)), \quad V_{ij}^{(\ell)} \sim \text{Beta}(1, \alpha), \quad \phi_{ij} \sim \mu/\mu(S) \quad (4.8)$$

$$H = \sum_{i=1}^{\infty} \sum_{j=1}^{C_i} V_{ij}^{(i)} \prod_{\ell=1}^{i-1} (1 - V_{ij}^{(\ell)}) \delta_{\theta_{ij}} \quad (4.9)$$

- Like the last construction, the weights are decreasing in expectation as a function of i .

$$H = \sum_{j=1}^{C_1} V_{1j} \delta_{\theta_{1j}} + \sum_{j=1}^{C_2} V_{2j}^{(2)} (1 - V_{2j}^{(1)}) \delta_{\theta_{2j}} + \sum_{j=1}^{C_3} V_{3j}^{(3)} (1 - V_{3j}^{(2)}) (1 - V_{3j}^{(1)}) \delta_{\theta_{3j}} + \dots \quad (4.10)$$

- The structure is also very similar. We have a Poisson-distributed number of atoms in each group and they're marked with an independent random variable.
- Therefore, we can write $H = \sum_{i=1}^{\infty} H_i$, where each H_i has a corresponding Poisson process Π_i with mean measure $\mu(S) \times (\mu(d\theta)/\mu(S)) \times \lambda_i(\pi) d\pi$, where $\lambda_i(\pi)$ is the distribution of

$$\pi = V_i \prod_{j=1}^{i-1} (1 - V_j), \quad V_i \sim \text{Beta}(1, \alpha).$$

- By the superposition theorem, H has an underlying Poisson process Π with mean measure equal to the sum of each H_i 's mean measures: $\mu(d\theta) \sum_{i=1}^{\infty} \lambda_i(\pi) d\pi$.
- Therefore, all that remains is to calculate this sum (which is a little complicated).

Proof:

- We focus on $\lambda_i(\pi)$, which is the distribution on $\pi = f(V_1, \dots, V_i)$, where $f(V_1, \dots, V_i) = V_i \prod_{j=1}^{i-1} (1 - V_j)$, $V_j \sim \text{Beta}(1, \alpha)$.
- **Lemma:** Let $T \sim \text{Gam}(i - 1, \alpha)$. Then $e^{-T} \stackrel{d}{=} \prod_{j=1}^{i-1} (1 - V_j)$.
- **Proof:** Define $\xi_j = -\ln(1 - V_j)$. We can show by a change of variables that $\xi_j \sim \text{Exp}(\alpha)$. The function $-\ln \prod_{j=1}^{i-1} (1 - V_j) = \sum_{j=1}^{i-1} \xi_j$. Since the V_j are independent, the ξ_j are independent. We know that sums of i.i.d. exponential r.v.'s are gamma distributed, so $T = \sum_{j=1}^{i-1} \xi_j$ is distributed as $\text{Gam}(i - 1, \alpha)$. That is, $-\ln \prod_{j=1}^{i-1} (1 - V_j) \stackrel{d}{=} T \sim \text{Gam}(i - 1, \alpha)$ and the result follows because the same function of two equally distributed r.v.'s is also equally distributed.
- We split the proof into two cases, $i = 1$ and $i > 1$.
- **Case $i = 1$:** $V_{1j} \sim \text{Beta}(1, \alpha)$, therefore $\pi_{1j} = V_{1j} \sim \lambda_1(\pi) d\pi = \alpha(1 - \pi)^{\alpha-1} d\pi$.

- Case $i > 1$: $V_{ij} \sim \text{Beta}(1, \alpha)$, $T_{ij} \sim \text{Gam}(i - 1, \alpha)$, $\pi_{ij} = V_{ij}e^{-T_{ij}}$. We need to find the density of π_{ij} . Let $W_{ij} = e^{-T_{ij}}$. Then changing variables,

$$p_{W_i}(w|\alpha) = \frac{\alpha^{i-1}}{(i-2)!} w^{\alpha-1} (-\ln w)^{i-2}. \quad (4.11)$$

↑

plug $T_{ij} = -\ln W_{ij}$ into gamma distribution and multiply by Jacobian

- Therefore $\pi_{ij} = V_{ij}W_{ij}$ and using the product distribution formula

$$\begin{aligned} \lambda_i(\pi|\alpha) &= \int_{\pi}^1 w^{-1} p_V(\pi/w|\alpha) p_{W_i}(w|\alpha) dw \\ &= \frac{\alpha^i}{(i-2)!} \int_{\pi}^1 w^{\alpha-2} (-\ln w)^{i-2} (1 - \pi/w)^{\alpha-1} dw \\ &= \frac{\alpha^i}{(i-2)!} \int_{\pi}^1 w^{-1} (-\ln w)^{i-2} (w - \pi)^{\alpha-1} dw \end{aligned} \quad (4.12)$$

- This integral doesn't have a closed form solution. However, recall that we only need to calculate $\mu(d\theta) \sum_{i=1}^{\infty} \lambda_i(\pi) d\pi$ to find the mean measure of the underlying Poisson process.

$$\mu(d\theta) \sum_{i=1}^{\infty} \lambda_i(\pi) d\pi = \mu(d\theta) \lambda_1(\pi) d\pi + \mu(d\theta) \sum_{i=2}^{\infty} \lambda_i(\pi) d\pi \quad (4.13)$$

$$\begin{aligned} \mu(d\theta) \sum_{i=2}^{\infty} \lambda_i(\pi) d\pi &= \mu(d\theta) \sum_{i=2}^{\infty} d\pi \frac{\alpha^i}{(i-2)!} \int_{\pi}^1 w^{-1} (-\ln w)^{i-2} (w - \pi)^{\alpha-1} dw \\ &= \mu(d\theta) d\pi \alpha^2 \int_{\pi}^1 w^{-1} (w - \pi)^{\alpha-1} dw \underbrace{\sum_{i=2}^{\infty} \frac{(-\alpha \ln w)^{i-2}}{(i-2)!}}_{= e^{-\alpha \ln w} = w^{-\alpha}} \\ &= \mu(d\theta) d\pi \alpha^2 \underbrace{\int_{\pi}^1 w^{-(\alpha+1)} (w - \pi)^{\alpha-1} dw}_{= \left. \frac{(w-\pi)^{\alpha}}{\alpha \pi w^{\alpha}} \right|_{\pi}^1} \\ &= \mu(d\theta) \frac{\alpha(1-\pi)^{\alpha}}{\pi} d\pi \end{aligned} \quad (4.14)$$

- Adding $\mu(d\theta) \lambda_1(\pi) d\pi$ from Case 1 with this last value,

$$\mu(d\theta) \sum_{i=1}^{\infty} \lambda_i(\pi) d\pi = \mu(d\theta) \alpha \pi^{-1} (1 - \pi)^{\alpha-1} d\pi. \quad (4.15)$$

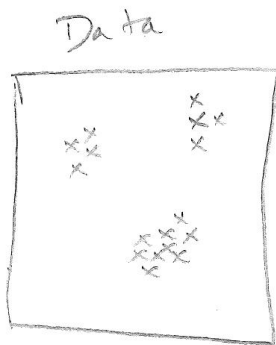
- Therefore, the construction corresponds to a Poisson process with mean measure equal to that of a beta process. It's therefore a beta process.

Chapter 5

Dirichlet processes and a size-biased construction

- We saw how beta processes can be useful as a Bayesian nonparametric prior for latent factor (matrix factorization) models.
- We'll next discuss BNP priors for mixture models.

Quick review



- 2-dimensional data generated from Gaussian with unknown mean and known variance.
- There are a small set of possible means and an observations picks one of them using a probability distribution.
- Let $G = \sum_{i=1}^K \pi_k \delta_{\theta_i}$ be the mixture distribution on mean parameters – θ_i : i th mean, π_i : probability of it

- For the n th observation,
 1. $c_n \sim \text{Disc}(\pi)$ picks mean index
 2. $x_n \sim N(\theta_{c_n}, \Sigma)$ generates observation

Priors on G

- Let μ be a non-atomic *probability* measure on the parameter space.
- Since π is a K -dimensional probability vector, a natural prior is Dirichlet.

Dirichlet distribution: A distribution on probability vectors

- Definition: Let $\alpha_1, \dots, \alpha_K$ be K positive numbers. The Dirichlet distribution density function is defined as

$$\text{Dir}(\pi | \alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_{i=1}^K \pi_i^{\alpha_i - 1} \quad (5.1)$$

- Goals: The goals are very similar to the beta process.
 1. We want $K \rightarrow \infty$
 2. We want the parameters $\alpha_1, \dots, \alpha_K$ to be such that, as $K \rightarrow \infty$, things are well-defined.
 3. It would be nice to link this to the Poisson process somehow.

Dirichlet random vectors and gamma random variables

- Theorem: Let $Z_i \sim \text{Gam}(\alpha_i, b)$ for $i = 1, \dots, K$. Define $\pi_i = Z_i / \sum_{j=1}^K Z_j$. Then

$$(\pi_1, \dots, \pi_K) \sim \text{Dir}(\alpha_1, \dots, \alpha_K). \quad (5.2)$$

Furthermore, π and $Y = \sum_{j=1}^K Z_j$ are independent random variables.

- Proof: This is just a change of variables.

$$- p(Z_1, \dots, Z_K) = \prod_{i=1}^K p(Z_i) = \prod_{i=1}^K \frac{b^{\alpha_i}}{\Gamma(\alpha_i)} Z_i^{\alpha_i-1} e^{-bZ_i}$$

$$- (Z_1, \dots, Z_K) := f(Y, \pi) = (Y\pi_1, \dots, Y\pi_{K-1}, Y(1 - \sum_{i=1}^{K-1} \pi_i))$$

$$- p_{Y,\pi}(Y, \pi) = P_Z(f(Y, \pi)) \cdot |J(f)| \quad \leftarrow J(\cdot) = \text{Jacobian}$$

$$- J(f) = \begin{bmatrix} \frac{\partial f_1}{\partial \pi_1} & \dots & \frac{\partial f_1}{\partial \pi_{K-1}} & \frac{\partial f_1}{\partial Y} \\ & \ddots & & \\ \frac{\partial f_K}{\partial \pi_1} & \dots & \frac{\partial f_K}{\partial \pi_{K-1}} & \frac{\partial f_K}{\partial Y} \end{bmatrix} = \begin{bmatrix} Y & 0 & \dots & \pi_1 \\ 0 & Y & 0 & \pi_2 \\ & & \ddots & \vdots \\ -Y & -Y & \dots & 1 - \sum_{i=1}^{K-1} \pi_i \end{bmatrix}$$

$$- \text{One can show that } |J(f)| = Y^{K-1}$$

- Therefore

$$\begin{aligned} p_Z(f(Y, \pi)) |J(f)| &= \prod_{i=1}^K \frac{b^{\alpha_i}}{\Gamma(\alpha_i)} (Y\pi_i)^{\alpha_i-1} e^{-bY\pi_i} Y^{K-1}, \quad (\pi_K := 1 - \sum_{i=1}^{K-1} \pi_i) \\ &= \underbrace{\left[\frac{b^{\sum_i \alpha_i}}{\Gamma(\sum_i \alpha_i)} Y^{\sum_i \alpha_i-1} e^{-bY} \right]}_{\text{Gam}(\sum_i \alpha_i, b)} \underbrace{\left[\frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_{i=1}^K \pi_i^{\alpha_i-1} \right]}_{\text{Dir}(\alpha_1, \dots, \alpha_K)} \end{aligned} \quad (5.3)$$

- We've shown that:

1. A Dirichlet distributed probability vector is a normalized sequence of independent gamma random variables with a constant scale parameter.
2. The sum of these gamma random variables is independent of the normalization because their joint distribution can be written as a product of two distributions.
3. This works in reverse: If we want to draw an independent sequence of gamma r.v.'s, we can draw a Dirichlet vector and scale it by an independent gamma random variable with first parameter equal to the sum of the Dirichlet parameters (and second parameter set to whatever we want).

Dirichlet process

- **Definition:** Let $\alpha > 0$ and μ a non-atomic *probability* measure on S . For all partitions of S , A_1, \dots, A_k , where $A_i \cap A_j = \emptyset$ for $i \neq j$ and $\cup_{i=1}^k A_i = S$, define the random measure G on S such that

$$(G(A_1), \dots, G(A_k)) \sim \text{Dir}(\alpha\mu(A_1), \dots, \alpha\mu(A_k)). \tag{5.4}$$

Then G is a Dirichlet process, denoted $G \sim \text{DP}(\alpha\mu)$.

Dirichlet processes via the gamma process

- Pick a partition of S , A_1, \dots, A_k . We can represent $G \sim \text{DP}(\alpha\mu)$ as the normalization of gamma distributed random variables,

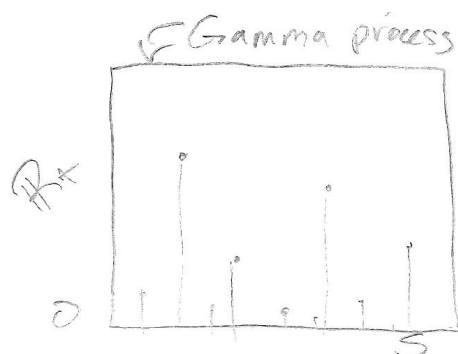
$$(G(A_1), \dots, G(A_k)) = \left(\frac{G'(A_1)}{G'(S)}, \dots, \frac{G'(A_k)}{G'(S)} \right), \tag{5.5}$$

$$G'(A_i) \sim \text{Gam}(\alpha\mu(A_i), b), \quad G'(S) = G'(\cup_{i=1}^k A_i) = \sum_{i=1}^k G'(A_i) \tag{5.6}$$

- Looking at the definition and how G' is defined, we realize that $G' \sim \text{GaP}(\alpha\mu, b)$. Therefore, a Dirichlet process is simply a normalized gamma process.
- Note that $G'(S) \sim \text{Gam}(\alpha, b)$. So $G'(S) < \infty$ with probability one and so the normalization G is well-defined.

Gamma processes and the Poisson process

- Recall that the gamma process is constructed from a Poisson process.
- Gamma process: Let N be a Poisson random measure on $S \times \mathbb{R}_+$ with mean measure $\alpha\mu(d\theta)z^{-1}e^{-bz}dz$. Define $G'(A_i) = \int_{A_i} \int_0^\infty N(d\theta, dz)z$. Then $G' \sim \text{GaP}(\alpha\mu, b)$.



- Since the DP is a rescaled GaP, this shares the same properties (from Campbells theorem)
 - For example, it’s an infinite jump process.
 - However, the DP is not a CRM like the GaP since $G(A_i)$ and $G(A_j)$ are not independent for disjoint sets A_i and A_j . This should be clear since G has to integrate to 1.

Dirichlet process as limit of finite approximation

- This is very similar to the previous discussion on limits of finite approximations to the gamma and beta process.
- Definition: Let $\alpha > 0$ and μ a non-atomic probability measure on S . Let

$$G_K = \sum_{i=1}^K \pi_i \delta_{\theta_i}, \quad \pi \sim \text{Dir}(\alpha/K, \dots, \alpha/K), \quad \theta_i \stackrel{iid}{\sim} \mu \quad (5.7)$$

Then $\lim_{K \rightarrow \infty} G_K = G \sim \text{DP}(\alpha\mu)$.

- Rough proof: We can equivalently write

$$G_K = \sum_{i=1}^K \left(\frac{Z_i}{\sum_{j=1}^K Z_j} \right) \delta_{\theta_i}, \quad Z_i \sim \text{Gam}(\alpha/K, b), \quad \theta_i \sim \mu \quad (5.8)$$

- If $G'_K = \sum_{i=1}^K Z_i \delta_{\theta_i}$, we've already proven that $G'_K \rightarrow G' \sim \text{GaP}(\alpha\mu, b)$. G_K is thus the limit of the normalization of G'_K . Since $\lim_{K \rightarrow \infty} G'_K(S)$ is finite almost surely, we can take the limit of the numerator and denominator of the gamma representation of G_K separately. The numerator converges to a gamma process and the denominator its normalization. Therefore, G_K converges to a Dirichlet process.

Some comments

- This infinite limit of the finite approximation results in an infinite vector, but the original definition was of a K dimensional vector, so is a Dirichlet process infinite or finite dimensional? Actually, the finite vector of the definition is constructed from an infinite process:

$$G(A_j) = \lim_{K \rightarrow \infty} G_K(A_j) = \lim_{K \rightarrow \infty} \sum_{i=1}^K \pi_i \delta_{\theta_i}(A_j). \quad (5.9)$$

- Since the partition A_1, \dots, A_k of S is of a continuous space we have to be able to let $K \rightarrow \infty$, so there has to be an infinite-dimensional process underneath G .
- The Dirichlet process gives us a way of defining priors on infinite discrete probability distributions on this continuous space S .
- As an intuitive example, if S is a space corresponding to the mean of a Gaussian, the Dirichlet process gives us a way to assign a probability to every possible value of this mean.
- Of course, by thinking of the DP in terms of the gamma and Poisson processes, we know that an infinite number of means will have probability zero, and infinite number will also have non-zero probability, but only a small handful of points in the space will have substantial probability. The number and locations of these atoms are random and learned during inference.
- Therefore, as with the beta process, size-biased representations of $G \sim \text{DP}(\alpha\mu)$ are needed.

A “stick-breaking” construction of $G \sim \text{DP}(\alpha\mu)$

- **Definition:** Let $\alpha > 0$ and μ be a non-atomic probability measure on S . Let

$$V_i \sim \text{Beta}(1, \alpha), \quad \theta_i \sim \mu \tag{5.10}$$

independently for $i = 1, 2, \dots$. Define

$$G = \sum_{i=1}^{\infty} V_i \prod_{j=1}^{i-1} (1 - V_j) \delta_{\theta_i}. \tag{5.11}$$

Then $G \sim \text{DP}(\alpha\mu)$.

- **Intuitive picture:** We start with a unit length stick and break off proportions.

$i=1$: $0 \mid \pi_1 \mid 1$ $G = V_1 \delta_{\theta_1} +$ $(1 - V_2)(1 - V_1)$ is what's left after the first two breaks. We take proportion V_3 of that for θ_3 and leave $(1 - V_3)(1 - V_2)(1 - V_1)$

$i=2$: $0 \mid \pi_1 \mid \pi_2 \mid 1$ $V_2(1 - V_1) \delta_{\theta_2} +$ ✓

$i=3$: $0 \mid \pi_1 \mid \pi_2 \mid \pi_3 \mid 1$ $V_3(1 - V_2)(1 - V_1) \delta_{\theta_3} + \dots$

Getting back to finite Dirichlets

- Recall from the definition that $(G(A_1), \dots, G(A_K)) \sim \text{Dir}(\alpha\mu(A_1), \dots, \alpha\mu(A_K))$ for all partitions A_1, \dots, A_K of S .
- Using the stick-breaking construction, we need to show that the vector formed by

$$G(A_k) = \sum_{i=1}^{\infty} V_i \prod_{j=1}^{i-1} (1 - V_j) \delta_{\theta_i}(A_k) \tag{5.12}$$

for $k = 1, \dots, K$ is distributed as $\text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$, where $\mu_k = \mu(A_k)$.

- Since $P(\theta_i \in A_k) = \mu(A_k)$, $\delta_{\theta_i}(A_k)$ can be equivalently represented by a K -dimensional vector $e_{Y_i} = (0, \dots, 1, \dots, 0)$, with the 1 in the position Y_i and $Y_i \sim \text{Disc}(\mu_1, \dots, \mu_K)$ and the rest 0.
- Letting $\pi_i = G(A_i)$, we therefore need to show that if

$$\pi = \sum_{i=1}^{\infty} V_i \prod_{j=1}^{i-1} (1 - V_j) e_{Y_i}, \quad V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad Y_i \stackrel{iid}{\sim} \text{Disc}(\mu_1, \dots, \mu_K) \tag{5.13}$$

Then $\pi \sim \text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$.

- Lemma: Let $\pi \sim \text{Dir}(a_1 + b_1, \dots, a_K + b_K)$. We can equivalently represent this as

$$\begin{aligned}\pi &= VY + (1 - V)W, \quad V \sim \text{Beta}(\sum_k a_k, \sum_k b_k), \\ Y &\sim \text{Dir}(a_1, \dots, a_K), \quad W \sim \text{Dir}(b_1, \dots, b_K)\end{aligned}\quad (5.14)$$

Proof: Use the normalized gamma representation: $\pi_i = Z_i / \sum_j Z_j$, $Z_i \sim \text{Gam}(a_i + b_i, c)$.

- We can use the equivalence

$$Z_i^Y \sim \text{Gam}(a_i, c), \quad Z_i^W \sim \text{Gam}(b_i, c) \quad \Leftrightarrow \quad Z_i^Y + Z_i^W \sim \text{Gam}(a_i + b_i, c) \quad (5.15)$$

- Splitting into two random variables this way we have the following normalized gamma representation for π

$$\begin{aligned}\pi &= \left(\frac{Z_1^Y + Z_1^W}{\sum_i Z_i^Y + Z_i^W}, \dots, \frac{Z_K^Y + Z_K^W}{\sum_i Z_i^Y + Z_i^W} \right) \\ &= \underbrace{\left(\frac{\sum_i Z_i^Y}{\sum_i Z_i^Y + Z_i^W} \right)}_{V \sim \text{Beta}(\sum_i a_i, \sum_i b_i)} \underbrace{\left(\frac{Z_1^Y}{\sum_i Z_i^Y}, \dots, \frac{Z_K^Y}{\sum_i Z_i^Y} \right)}_{Y \sim \text{Dir}(a_1, \dots, a_K)} \\ &\quad + \underbrace{\left(\frac{\sum_i Z_i^W}{\sum_i Z_i^Y + Z_i^W} \right)}_{1 - V} \underbrace{\left(\frac{Z_1^W}{\sum_i Z_i^W}, \dots, \frac{Z_K^W}{\sum_i Z_i^W} \right)}_{W \sim \text{Dir}(b_1, \dots, b_K)}\end{aligned}\quad (5.16)$$

- From the previous proof about normalized gamma r.v.'s, we know that the sums are independent from the normalized values. So V , Y , and W are all independent.

- Proof of stick-breaking construction:

- Start with $\pi \sim \text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$. Also, we use $\alpha_i \equiv \alpha\mu_i$ in parts below.

- Step 1:

$$\begin{aligned}\frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_{i=1}^K \pi_i^{\alpha_i - 1} &= \left(\sum_{j=1}^K \pi_j \right) \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_{i=1}^K \pi_i^{\alpha_i - 1} \\ &= \sum_{j=1}^K \frac{\alpha\mu_j}{\alpha\mu_j} \frac{\Gamma(\alpha)}{\prod_i \Gamma(\alpha_i)} \prod_{i=1}^K \pi_i^{\alpha_i + e_j(i) - 1} \\ &= \sum_{j=1}^K \mu_j \underbrace{\frac{\Gamma(1 + \alpha)}{\Gamma(1 + \alpha_j) \prod_{i \neq j} \Gamma(\alpha_i)} \prod_{i=1}^K \pi_i^{\alpha_i + e_j(i) - 1}}_{= \text{Dir}(\alpha\mu + e_j)}\end{aligned}\quad (5.17)$$

- Therefore, a hierarchical representation of $\text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$ is

$$Y \sim \text{Discrete}(\mu_1, \dots, \mu_K), \quad \pi \sim \text{Dir}(\alpha\mu + e_Y).\quad (5.18)$$

– Step 2:

From the lemma we have that $\pi \sim \text{Dir}(\alpha\mu + e_Y)$ can be expanded into the equivalent hierarchical representation $\pi = VY' + (1 - V)\pi'$, where

$$V \sim \text{Beta}\left(\underbrace{\sum_i e_Y(i)}_{=1}, \underbrace{\sum_i \alpha\mu_i}_{=\alpha}\right), \quad \underbrace{Y' \sim \text{Dir}(e_Y)}_{= e_Y \text{ with probability 1}}, \quad \pi' \sim \text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K) \quad (5.19)$$

– Combining Steps 1& 2:

We will use these steps to recursively break down a Dirichlet distributed random vector an infinite number of times. If

$$\pi = Ve_Y + (1 - V)\pi', \quad (5.20)$$

$$V \sim \text{Beta}(1, \alpha), \quad Y \sim \text{Disc}(\mu_1, \dots, \mu_K), \quad \pi' \sim \text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K),$$

Then from steps 1 & 2, $\pi \sim \text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$.

– Notice that there are $\text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$ r.v.'s on both sides. We “broke down” the one on the left. We can continue by “breaking down” the one on the right:

$$\pi = V_1e_{Y_1} + (1 - V_1)(V_2e_{Y_2} + (1 - V_2)\pi'') \quad (5.21)$$

$$V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad Y_i \stackrel{iid}{\sim} \text{Disc}(\mu_1, \dots, \mu_K), \quad \pi'' \sim \text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K),$$

π is still distributed as $\text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$.

– Continue this an infinite number of times:

$$\pi = \sum_{i=1}^{\infty} V_i \prod_{j=1}^{i-1} (1 - V_j) e_{Y_i}, \quad V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad Y_i \stackrel{iid}{\sim} \text{Disc}(\mu_1, \dots, \mu_K). \quad (5.22)$$

Still, following each time the right-hand Dirichlet is expanded we get a $\text{Dir}(\alpha\mu_1, \dots, \alpha\mu_K)$ random variable. Since $\lim_{T \rightarrow \infty} \prod_{j=1}^T (1 - V_j) = 0$, the term pre-multiplying this RHS Dirichlet vector equals zero and the limit above results, which completes the proof.

• Corollary:

If G is drawn from $\text{DP}(\alpha\mu)$ using the stick-breaking construction and $\beta \sim \text{Gam}(\alpha, b)$ independently, then $\beta G \sim \text{GaP}(\alpha\mu, b)$. Writing this out,

$$\beta G = \sum_{i=1}^{\infty} \beta \left(V_i \prod_{j=1}^{i-1} (1 - V_j) \right) \delta_{\theta_i} \quad (5.23)$$

- We therefore get a method for drawing a gamma process almost for free. Notice that α appears in both the DP and gamma distribution on β . These parameters must be the same value for βG to be a gamma process.

Chapter 6

Dirichlet process extensions, count processes

Gamma process to Dirichlet process

- Gamma process: Let $\alpha > 0$ and μ a non-atomic probability measure on S . Let $N(d\theta, dw)$ be a Poisson random measure on $S \times \mathbb{R}_+$ with mean measure $\alpha\mu(d\theta)we^{-cw}dw$, $c > 0$. For $A \subset S$, let $G'(A) = \int_A \int_0^\infty N(d\theta, dw)w$. Then G' is a gamma process, $G' \sim \text{GaP}(\alpha\mu, c)$, and $G'(A) \sim \text{Gam}(\alpha\mu(A), c)$.
- Normalizing a gamma process: Let's take G' and normalize it. That is, define $G(d\theta) = \frac{G'(d\theta)}{G'(S)}$. ($G'(S) \sim \text{Gam}(\alpha, c)$, so it's finite w.p. 1). Then G is called a Dirichlet process, written $G \sim \text{DP}(\alpha\mu)$.

Why? Take S and partition it into K disjoint regions, i.e., (A_1, \dots, A_K) , $A_i \cap A_j = \emptyset$, $i \neq j$, $\cup_i A_i = S$. Construct the vector

$$(G(A_1), \dots, G(A_K)) = \left(\frac{G'(A_1)}{G'(S)}, \dots, \frac{G'(A_K)}{G'(S)} \right). \quad (6.1)$$

Since each $G'(A_i) \sim \text{Gam}(\alpha\mu(A_i), c)$, and $G'(S) = \sum_{i=1}^K G'(A_i)$, it follows that

$$(G(A_1), \dots, G(A_K)) \sim \text{Dir}(\alpha\mu(A_1), \dots, \alpha\mu(A_K)). \quad (6.2)$$

This is the *definition* of a Dirichlet process.

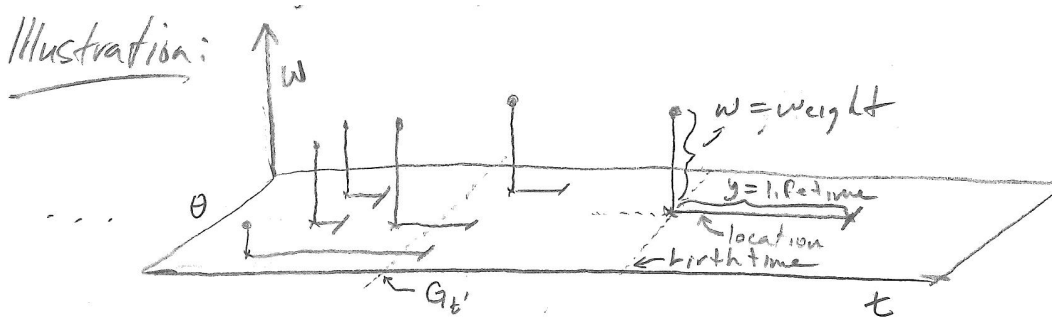
- The Dirichlet process has many extensions to suit the structure of different problems.
- We'll look at four, two that are related to the underlying normalized gamma process, and two from the perspective of the stick-breaking construction.
- The purpose is to illustrate how the basic framework of Dirichlet process mixture modeling can be easily built into more complicated models that address problems not perfectly suited to the basic construction.
- Goal is to make it clear how to continue these lines of thinking to form new models.

Example 1: Spatially and temporally normalized gamma processes

- Imagine we wanted a temporally evolving Dirichlet process. Clusters (i.e., atoms, θ) may arise and die out at different times (or exist in geographical regions)

Time-evolving model: Let $N(d\theta, dw, dt)$ be a Poisson random measure on $S \times \mathbb{R}_+ \times \mathbb{R}$ with mean measure $\alpha\mu(d\theta)w^{-1}e^{-cw}dw dt$. Let $G'(d\theta, dt) = \int_0^\infty N(d\theta, dw, dt)w$. Then G' is a gamma process with added “time” dimension t .

- (There’s nothing new from what we’ve studied: Let $\theta^* = (\theta, t)$ and $\alpha\mu(d\theta^*) = \alpha\mu(d\theta)dt$.)
- For each atom (θ, t) with $G'(d\theta, dt) > 0$, add a marking $y_t(\theta) \stackrel{ind}{\sim} \text{Exp}(\lambda)$.
- We can think of $y_t(\theta)$ as the lifetime of parameter θ born at time t .
- By the marking theorem, $N^*(d\theta, dw, dt, dy) \sim \text{Pois}(\alpha\mu(d\theta)w^{-1}e^{-cw}dw dt \lambda e^{-\lambda y} dy)$

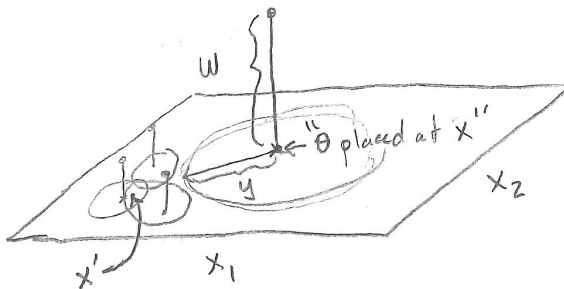


- At time t' , construct the Dirichlet process $G_{t'}$ by normalizing over all atoms “alive” at time t' . (Therefore, ignore atoms already dead or yet to be born.)

Spatial model: Instead of giving each atom θ a time-stamp and “lifetime,” we might want to give it a location and “region of influence”.

- Replace $t \in \mathbb{R}$ with $x \in \mathbb{R}^2$ (e.g., latitude-longitude). Replace dt with dx .
- Instead of $y_t(\theta) \sim \text{Exp}(\lambda) = \text{lifetime}$, $y_x(\theta) \sim \text{Exp}(\lambda) = \text{radius of ball at } x$.

Illustration:



- G'_x is the DP at location x' .
- It is formed by normalizing over all atoms θ for which

$$x' \in \text{ball of radius } y_x(\theta) \text{ at } x$$

Example 2: Another time-evolving formulation

- We can think of other formulations. Here's one where time is discrete. (We will build up to this with the following two properties).
- Even though the DP doesn't have an underlying PRM, the fact that it's constructed from a PRM means we can still benefit from its properties.

Superposition and the Dirichlet process: Let $G'_1 \sim \text{GaP}(\alpha_1\mu_1, c)$ and $G'_2 \sim \text{GaP}(\alpha_2\mu_2, c)$. Then $G'_{1+2} = G'_1 + G'_2 \sim \text{GaP}(\alpha_1\mu_1 + \alpha_2\mu_2, c)$. Therefore,

$$G_{1+2} = \frac{G'_{1+2}}{G'_{1+2}(S)} \sim \text{DP}(\alpha_1\mu_1 + \alpha_2\mu_2). \quad (6.3)$$

We can equivalently write

$$G_{1+2} = \underbrace{\frac{G'_1(S)}{G'_{1+2}(S)}}_{\text{Beta}(\alpha_1, \alpha_2)} \times \underbrace{\frac{G'_1}{G'_1(S)}}_{\text{DP}(\alpha_1\mu_1)} + \frac{G'_2(S)}{G'_{1+2}(S)} \times \underbrace{\frac{G'_2}{G'_2(S)}}_{\text{DP}(\alpha_2\mu_2)} \sim \text{DP}(\alpha_1\mu_1 + \alpha_2\mu_2) \quad (6.4)$$

From the lemma last week, these two DP's and the beta r.v. are all independent.

- Therefore,

$$G = \pi G_1 + (1 - \pi)G_2, \quad \pi \sim \text{Beta}(\alpha_1, \alpha_2), \quad G_1 \sim \text{DP}(\alpha_1\mu_1), \quad G_2 \sim \text{DP}(\alpha_2\mu_2) \quad (6.5)$$

is equal in distribution to $G \sim \text{DP}(\alpha_1\mu_1 + \alpha_2\mu_2)$.

Thinning of gamma processes (a special case of the marking theorem)

- We know that we can construct $G' \sim \text{GaP}(\alpha\mu, c)$ from the Poisson random measure $N(d\theta, dw) \sim \text{Pois}(\alpha\mu(d\theta)w^{-1}e^{-cw}dw)$. Mark each point (θ, w) in N with a binary variable $z \sim \text{Bern}(p)$. Then

$$N(d\theta, dw, z) \sim \text{Pois}(p^z(1-p)^{1-z}\alpha\mu(d\theta)w^{-1}e^{-cw}dw). \quad (6.6)$$

- If we view $z = 1$ as "survival" and $z = 0$ as "death," then if we only care about the atoms that survive, we have

$$N_1(d\theta, dw) = N(d\theta, dw, z = 1) \sim \text{Pois}(p\alpha\mu(d\theta)w^{-1}e^{-cw}dw). \quad (6.7)$$

- This is called "thinning." We see that $p \in (0, 1)$ down-weights the mean measure, so we only expect to see a fraction p of what we saw before.

- Still, a normalized thinned gamma process is a Dirichlet process

$$\dot{G}' \sim \text{GaP}(p\alpha\mu, c), \quad \dot{G} = \frac{\dot{G}'}{\dot{G}'(S)} \sim \text{DP}(p\alpha\mu). \quad (6.8)$$

- What happens if we thin twice? We're marking with $z \in \{0, 1\}^2$ and restricting to $z = [1, 1]$,

$$\ddot{G}' \sim \text{GaP}(p^2\alpha\mu, c) \quad \rightarrow \quad \ddot{G} \sim \text{DP}(p^2\alpha\mu). \quad (6.9)$$

- Back to the example, we again want a time-evolving Dirichlet process where new atoms are born and old atoms die out.
- We can easily achieve this by introducing new gamma processes and thinning old ones.

A dynamic Dirichlet process:

- At time t :
1. Draw $G_t^* \sim \text{GaP}(\alpha_t\mu_t, c)$.
 2. Construct $G'_t = G_t^* + \dot{G}'_{t-1}$, where \dot{G}'_{t-1} is the gamma process at time $t - 1$ thinned with parameter p .
 3. Normalize $G_t = G'_t / G'_t(S)$.

- Why is G_t still a Dirichlet process? Just look at G'_t :
 - Let $G'_{t-1} \sim \text{GaP}(\hat{\alpha}_{t-1}\hat{\mu}_{t-1}, c)$.
 - Then $\dot{G}'_{t-1} \sim \text{GaP}(p\hat{\alpha}_{t-1}\hat{\mu}_{t-1}, c)$ and $G'_t \sim \text{GaP}(\alpha_t\mu_t + p\hat{\alpha}_{t-1}\hat{\mu}_{t-1}, c)$.
 - So $G_t \sim \text{DP}(\alpha_t\mu_t + p\hat{\alpha}_{t-1}\hat{\mu}_{t-1})$.

- By induction,

$$G_t \sim \text{DP}(\alpha_t\mu_t + p\alpha_{t-1}\mu_{t-1} + p^2\alpha_{t-2}\mu_{t-2} + \cdots + p^{t-1}\alpha_1\mu_1). \quad (6.10)$$

- If we consider the special case where $\alpha_t\mu_t = \alpha\mu$ for all t , we can simplify this Dirichlet process

$$G_t \sim \text{DP}\left(\frac{1-p^t}{1-p}\alpha\mu\right). \quad (6.11)$$

In the limit $t \rightarrow \infty$, this has the steady state

$$G_\infty \sim \text{DP}\left(\frac{1}{1-p}\alpha\mu\right). \quad (6.12)$$

- Stick-breaking construction (review for the next process)

We saw that if $\alpha > 0$ and μ is any probability measure, atomic or non-atomic or mixed, then we can draw $G \sim \text{DP}(\alpha\mu)$ as follows:

$$V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad \theta_i \stackrel{iid}{\sim} \mu, \quad G = \sum_{i=1}^{\infty} V_i \prod_{j=1}^{i-1} (1 - V_j) \delta_{\theta_i} \quad (6.13)$$

- It's often the case where we have grouped data. For example, groups of documents where each document is a set of words.
- We might want to model each group (indexed by d) as a mixture G_d . Then, for observation n in group d , $\theta_n^{(d)} \sim G_d$, $x_n^{(d)} \sim p(x|\theta_n^{(d)})$.
- We might think that each group shares the same set of highly probable atoms, but has different distributions on them.
- The result is called a mixed-membership model.

Mixed-membership models and the hierarchical Dirichlet process (HDP)

- As the stick-breaking construction makes clear, when μ is non-atomic simply drawing each $G_d \stackrel{iid}{\sim} \text{DP}(\alpha\mu)$ won't work because it places all probability mass on a disjoint set of atoms.
- The HDP fixes this by “discretizing the base distribution.”

$$G_d | G_0 \stackrel{iid}{\sim} \text{DP}(\beta G_0), \quad G_0 \sim \text{DP}(\alpha\mu). \quad (6.14)$$

- Since G_0 is discrete, G_d has probability on the same subset of atoms. This is very obvious by writing the process with the stick-breaking construction:

$$G_d = \sum_{i=1}^{\infty} \pi_i^{(d)} \delta_{\theta_i}, \quad (\pi_1^{(d)}, \pi_2^{(d)}, \dots) \sim \text{Dir}(\alpha p_1, \alpha p_2, \dots) \quad (6.15)$$

$$p_i = V_i \prod_{j=1}^{i-1} (1 - V_j), \quad V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad \theta_i \stackrel{iid}{\sim} \mu.$$

- Nested Dirichlet processes

The stick-breaking construction is totally general: μ can be any distribution.

What if $\mu \rightarrow \text{DP}(\alpha\mu)$? That is, we define the base distribution to be a Dirichlet process.

$$G \sim \sum_{i=1}^{\infty} V_i \prod_{j=1}^{i-1} (1 - V_j) \delta_{G_i}, \quad V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad G_i \stackrel{iid}{\sim} \text{DP}(\alpha\mu). \quad (6.16)$$

(We write G_i to link to the DP, but we could have written $\theta_i \stackrel{iid}{\sim} \text{DP}(\alpha\mu)$ since that's what we've been using.)

- We now have a mixture model of mixture models. For example:

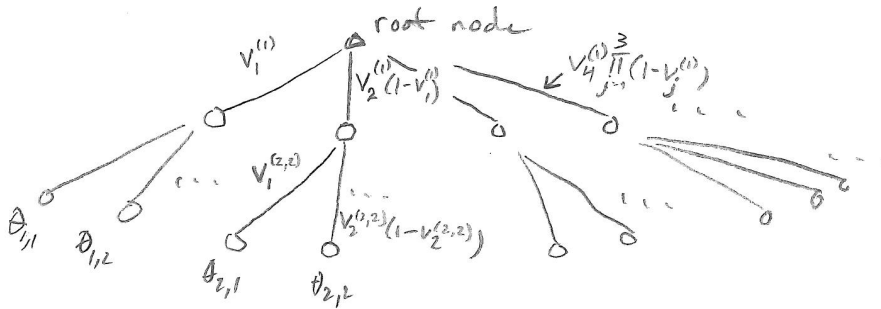
1. A group selects $G^{(d)} \sim G$ (picks mixture G_i according to probability $V_i \prod_{j<i} (1 - V_j)$).

2. Generates all of its data using this mixture. For the n th observation in group d , $\theta_n^{(d)} \sim G^{(d)}$, $X_n^{(d)} \sim p(X|\theta_n^{(d)})$.

- In this case we have all-or-nothing sharing. Two groups either share the atoms and the distribution on them, or they share nothing.

Nested Dirichlet process trees

- We can nest this further. Why not let μ in the nDP be a Dirichlet process also? Then we would have a three level tree.



- We can then pick paths down this tree to a leaf node where we get an atom.

Count Processes

- We briefly introduce count processes. With the Dirichlet process, we often have the generative structure

$$G \sim \text{DP}(\alpha\mu), \quad \theta_j^* | G \stackrel{iid}{\sim} G, \quad G = \sum_{i=1}^{\infty} \pi_i \delta_{\theta_i}, \quad j = 1, \dots, N \quad (6.17)$$

- What can we say about the count process $n(\theta) = \sum_{j=1}^N \mathbb{1}(\theta_j^* = \theta)$?

- Recall the following equivalent processes:

$$G' \sim \text{GaP}(\alpha\mu, c) \quad (6.18) \quad \text{and} \quad G' \sim \text{GaP}(\alpha\mu, c) \quad (6.20)$$

$$n(\theta) | G' \sim \text{Pois}(G'(\theta)) \quad (6.19) \quad n(S) \sim \text{Pois}(G'(S)) \quad (6.21)$$

$$\theta_{1:n(S)}^* \sim G' / G'(S) \quad (6.22)$$

- We can therefore analyze this using the underlying marked Poisson process. However, notice that we have to let the data size be random and Poisson distributed.

Marking theorem: Let $G' \sim \text{GaP}(\alpha\mu, c)$ and mark each (θ, w) for which $G'(\theta) = w > 0$ with the random variable $n|w \sim \text{Pois}(w)$. Then (θ, w, n) is a marked Poisson process with mean measure $\alpha\mu(d\theta)w^{-1}e^{-cw}dw \frac{w^n}{n!}e^{-w}$.

- We can restrict this to n by integrating over θ and w .

Theorem: The number of atoms having k counts is

$$\#_k \sim \text{Pois} \left(\int_S \int_0^\infty \alpha \mu(d\theta) w^{-1} e^{-cw} dw \frac{w^k}{k!} e^{-w} \right) = \text{Pois} \left(\frac{\alpha}{k} \left(\frac{1}{1+c} \right)^k \right) \quad (6.23)$$

Theorem: The total number of uniquely observed atoms is also Poisson

$$\#_{\text{unique}} = \sum_{k=1}^{\infty} \#_k \sim \text{Pois} \left(\sum_{k=1}^{\infty} \frac{\alpha}{k} \left(\frac{1}{1+c} \right)^k \right) = \text{Pois}(\alpha \ln(1+c^{-1})) \quad (6.24)$$

Theorem: The total number of counts is $n(S)|G' \sim \text{Pois}(G'(S))$, $G' \sim \text{GaP}(\alpha\mu, c)$. So $\mathbb{E}[n(S)] = \frac{\alpha}{c}$ ($= \sum_{k=1}^{\infty} k \mathbb{E}\#_k$)

Final statement: Let $c = \frac{\alpha}{N}$. If we expect a dataset of size N to be drawn from $G \sim \text{DP}(\alpha\mu)$, we expect that dataset to use $\alpha \ln(\alpha + N) - \alpha \ln \alpha$ unique atoms from G .

A quick final count process

- Instead of gamma process \rightarrow Poisson counts, we could have beta process \rightarrow negative binomial counts.
- Let $H = \sum_{i=1}^{\infty} \pi_i \delta_{\theta_i} \sim \text{BP}(\alpha, \mu)$.
- Let $n(\theta) = \text{negBin}(r, H(\theta))$, where the negative binomial random variable counts how many “successes” there are, with $P(\text{success}) = H(\theta)$ until there are r “failures” with $P(\text{failure}) = 1 - H(\theta)$.
- This is another count process that can be analyzed using the underlying Poisson process.

Chapter 7

Exchangeability, Dirichlet processes and the Chinese restaurant process

DP's, finite approximations and mixture models

- DP: We saw how, if $\alpha > 0$ and μ is a probability measure on S , for every finite partition (A_1, \dots, A_k) of S , the random measure

$$(G(A_1), \dots, G(A_k)) \sim \text{Dir}(\alpha\mu(A_1), \dots, \alpha\mu(A_k))$$

defines a Dirichlet process.

- Finite approximation: We also saw how we can approximate $G \sim \text{DP}(\alpha\mu)$ with a finite Dirichlet distribution,

$$G_K = \sum_{i=1}^K \pi_i \delta_{\theta_i}, \quad \pi_i \sim \text{Dir}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right), \quad \theta \stackrel{iid}{\sim} \mu.$$

- Mixture models: Finally, the most common setting for these priors is in mixture models, where we have the added layers

$$\theta_j^* | G \sim G, \quad X_j | \theta_j^* \sim p(X | \theta_j^*), \quad j = 1, \dots, n. \quad (Pr(\theta_j^* = \theta_i | G) = \pi_i)$$

- The values of $\theta_1^*, \dots, \theta_n^*$ induce a clustering of the data.
- If $\theta_j^* = \theta_{j'}^*$ for $j \neq j'$ then X_j and $X_{j'}$ are “clustered” together since they come from the same distribution.
- We’ve thus far focused on G . We now focus on the clustering of X_1, \dots, X_n induced by G .

Polya’s Urn model (finite Dirichlet)

- To simplify things, we work in the finite setting and replace the parameter θ_i with its index i . We let the indicator variables c_1, \dots, c_n represent $\theta_1^*, \dots, \theta_n^*$ such that $\theta_j^* = \theta_{c_j}$.

- Polya's Urn is the following process for generating c_1, \dots, c_n :

1. For the first indicator, $c_1 \sim \sum_{i=1}^K \frac{1}{K} \delta_i$

2. For the n th indicator, $c_n | c_1, \dots, c_{n-1} \sim \sum_{j=1}^{n-1} \frac{1}{\alpha+n-1} \delta_{c_j} + \frac{\alpha}{\alpha+n-1} \sum_{i=1}^K \frac{1}{K} \delta_i$

- In words, we start with an urn having $\frac{\alpha}{K}$ balls of color i for each of K colors. We randomly pick a ball, put it back in the urn and put another ball of the same color in the urn.

- Another way to write #2 above is to define $n_i^{(n-1)} = \sum_{j=1}^{n-1} \mathbb{1}(c_j = i)$. Then

$$c_n | c_1, \dots, c_{n-1} \sim \sum_{i=1}^K \frac{\frac{\alpha}{K} + n_i^{(n-1)}}{\alpha + n - 1} \delta_i.$$

To put it most simply, we're just sampling the next color from the empirical distribution of the urn at step n .

- What can we say about $p(c_1 = i_1, \dots, c_n = i_n)$ (write as $p(c_1, \dots, c_n)$) under this prior?

1. By the chain rule of probability, $p(c_1, \dots, c_n) = \prod_{j=1}^n p(c_j | c_1, \dots, c_{j-1})$.

2. $p(c_j = i | c_1, \dots, c_{j-1}) = \frac{\frac{\alpha}{K} + n_i^{(j-1)}}{\alpha + j - 1}$

3. Therefore,

$$p(c_{1:n}) = p(c_1)p(c_2|c_1)p(c_3|c_1, c_2) \cdots = \prod_{j=1}^n \frac{\frac{\alpha}{K} + n_{c_j}^{(j-1)}}{\alpha + j - 1} \quad (7.1)$$

- A few things to notice about $p(c_1, \dots, c_n)$

1. The denominator is simply $\prod_{j=1}^n (\alpha + j - 1)$

2. $n_{c_j}^{(j-1)}$ is incrementing by one. That is, after $c_{1:n}$ we have the counts $(n_1^{(n)}, \dots, n_K^{(n)})$. For each $n_i^{(n)}$ the numerator will contain $\prod_{s=1}^{n_i^{(n)}} (\frac{\alpha}{K} + s - 1)$.

3. Therefore,

$$p(c_1, \dots, c_n) = \frac{\prod_{i=1}^K \prod_{s=1}^{n_i^{(n)}} (\frac{\alpha}{K} + s - 1)}{\prod_{j=1}^n (\alpha + j - 1)} \quad (7.2)$$

- **Key:** The key thing to notice is that this does not depend on the order of c_1, \dots, c_n . That is, if we permuted c_1, \dots, c_n such that $c_j = i_{\rho(j)}$, where $\rho(\cdot)$ is a permutation of $(1, \dots, n)$, then

$$p(c_1 = i_1, \dots, c_n = i_n) = p(c_1 = i_{\rho(1)}, \dots, c_n = i_{\rho(n)}).$$

- The sequence c_1, \dots, c_n is said to be "exchangeable" in this case.

Exchangeability and independent and identically distributed (iid) sequences

- Independent sequences are exchangeable

$$p(c_1, \dots, c_n) = \prod_{i=1}^n p(c_i) = \prod_{i=1}^n p(c_{\rho(i)}) = p(c_{\rho(1)}, \dots, c_{\rho(n)}). \tag{7.3}$$

- Exchangeable sequences aren't necessarily independent (exchangeability is “weaker”). Think of the urn. c_j is clearly not independent of c_1, \dots, c_{j-1} .

Exchangeability and de Finetti

- **de Finetti's theorem:** A sequence is exchangeable if and only if there is a parameter π with distribution $p(\pi)$ for which the sequence is independent and identically distributed given π .
- In other words, for our problem there is a probability vector π such that $p(c_{1:n}|\pi) = \prod_{j=1}^n p(c_j|\pi)$.
- The problem is to find $p(\pi)$

$$\begin{aligned} p(c_1, \dots, c_n) &= \int p(c_1, \dots, c_n|\pi)p(\pi)d\pi \\ &= \int \prod_{j=1}^n p(c_j|\pi)p(\pi)d\pi \\ &= \int \prod_{j=1}^n \pi_{c_j} p(\pi)d\pi \\ &= \int \prod_{i=1}^K \pi_i^{n_i^{(n)}} p(\pi)d\pi \\ &\quad \downarrow \qquad \qquad \qquad \downarrow \\ \frac{\prod_{i=1}^K \prod_{s=1}^{n_i^{(n)}} (\frac{\alpha}{K} + s - 1)}{\prod_{j=1}^n (\alpha + j - 1)} &= \mathbb{E}_{p(\pi)} \left[\prod_{i=1}^K \pi_i^{n_i^{(n)}} \right] \end{aligned} \tag{7.4}$$

- Above, the first equality is always true. The second one is by de Finetti's theorem since c_1, \dots, c_n is exchangeable. (We won't proven this theorem, we'll just use it.) The following results. In the last equality, the left hand side was previously shown and the right hand side is what the second to last line is equivalently written as.
- By de Finetti and exchangeability of c_1, \dots, c_n , we therefore arrive at an expression for the moments of π according to the still unknown distribution $p(\pi)$.

- Because the moments of a distribution are unique to that distribution (like the Laplace transform), $p(\pi)$ has to be $\text{Dir}(\frac{\alpha}{K}, \dots, \frac{\alpha}{K})$, since plugging this in for $p(\pi)$ we get

$$\begin{aligned}
 \mathbb{E}_{p(\pi)} \left[\prod_{i=1}^K \pi_i^{n_i^{(n)}} \right] &= \int \prod_{i=1}^K \pi_i^{n_i} \frac{\Gamma(\alpha)}{\Gamma(\frac{\alpha}{K})^K} \prod_{i=1}^K \pi_i^{\frac{\alpha}{K}-1} d\pi \\
 &= \frac{\Gamma(\alpha) \prod_{i=1}^K \Gamma(\frac{\alpha}{K} + n_i)}{\Gamma(\frac{\alpha}{K})^K \Gamma(\alpha + n)} \int \underbrace{\frac{\Gamma(\alpha + n)}{\prod_{i=1}^K \Gamma(\frac{\alpha}{K} + n_i)} \prod_{i=1}^K \pi_i^{n_i + \frac{\alpha}{K} - 1}}_{= \text{Dir}(\frac{\alpha}{K} + n_1, \dots, \frac{\alpha}{K} + n_k)} d\pi \\
 &= \frac{\Gamma(\alpha) \prod_{i=1}^K \Gamma(\frac{\alpha}{K}) \prod_{s=1}^{n_i} (\frac{\alpha}{K} + s - 1)}{\Gamma(\frac{\alpha}{K})^K \Gamma(\alpha) \prod_{j=1}^n (\alpha + j - 1)} \\
 &= \frac{\prod_{i=1}^K \prod_{s=1}^{n_i^{(n)}} (\frac{\alpha}{K} + s - 1)}{\prod_{j=1}^n (\alpha + j - 1)} \tag{7.5}
 \end{aligned}$$

- This holds for all n and (n_1, \dots, n_k) . Since a distribution is defined by its moments, the result follows.
- Notice that we didn't *need* de Finetti since we could just hypothesize the existence of a π for which $p(c_{1:n}|\pi) = \prod_i p(c_i|\pi)$ and try to find it. It's more useful when the distribution is more "non-standard," or to prove that a π doesn't exist.
- Final statement: As $n \rightarrow \infty$, the distribution $\sum_{i=1}^K \frac{n_i^{(n)} + \frac{\alpha}{K}}{\alpha + n} \delta_i \rightarrow \sum_{i=1}^K \pi_i^* \delta_i$.
 - This is because there exists a π for which c_1, \dots, c_n are iid, and so by the law of large numbers the point π^* exists and $\pi^* = \pi$.
 - Since $\pi \sim \text{Dir}(\frac{\alpha}{K}, \dots, \frac{\alpha}{K})$, it follows that the empirical distribution converges to a *random* vector that is distributed as $\text{Dir}(\frac{\alpha}{K}, \dots, \frac{\alpha}{K})$.

The infinite limit (Chinese restaurant process)

- Let's go back to the original notation:

$$\theta_j^* | G_K \sim G_K, \quad G_K = \sum_{i=1}^K \pi_i \delta_{\theta_i}, \quad \pi \sim \text{Dir}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right), \quad \theta_i \stackrel{iid}{\sim} \mu.$$

- Following the exact same ideas (only changing notation). The urn process is

$$\theta_n^* | \theta_1^*, \dots, \theta_{n-1}^* \sim \sum_{i=1}^K \frac{\frac{\alpha}{K} + n_i^{(n-1)}}{\alpha + n - 1} \delta_{\theta_i}, \quad \theta_i \stackrel{iid}{\sim} \mu.$$

- We've proven that $\lim_{K \rightarrow \infty} G_K = G \sim \text{DP}(\alpha\mu)$. We now take the limit of the corresponding urn process.

- Re-indexing: At observation n , re-index the atoms so that $n_j^{(n-1)} > 0$ for $j = 1, \dots, K_{n-1}^+$ and $n_j^{(n-1)} = 0$ for $j > K_{n-1}^+$. ($K_{n-1}^+ = \#$ unique values in $\theta_{1:n-1}^*$) Then

$$\theta_n^* | \theta_1^*, \dots, \theta_{n-1}^* \sim \sum_{i=1}^{K_{n-1}^+} \frac{\frac{\alpha}{K} + n_i^{(n-1)}}{\alpha + n - 1} \delta_{\theta_i} + \frac{\alpha}{\alpha + n - 1} \sum_{i=1+K_{n-1}^+}^K \frac{1}{K} \delta_{\theta_i}. \quad (7.6)$$

- Obviously for $n \gg K$, $K_{n-1}^+ = K$ very probably, and just the left term remains. However, we're interested in $K \rightarrow \infty$ before we let n grow. In this case

1. $\frac{\frac{\alpha}{K} + n_i^{(n-1)}}{\alpha + n - 1} \rightarrow \frac{n_i^{(n-1)}}{\alpha + n - 1}$
2. $\sum_{i=1+K_{n-1}^+}^K \frac{1}{K} \delta_{\theta_i} \rightarrow \mu.$

- For #2, if you sample K times from a distribution and create a uniform measure on those samples, then in the infinite limit you get the original distribution back. Removing $K_{n-1}^+ < \infty$ of those atoms doesn't change this (we won't prove this).

The Chinese restaurant process

- Let $\alpha > 0$ and μ a probability measure on S . Sample the sequence $\theta_1^*, \dots, \theta_n^*, \theta^* \in S$ as follows:

1. Set $\theta_1^* \sim \mu$
2. Sample $\theta_n^* | \theta_1^*, \dots, \theta_{n-1}^* \sim \sum_{j=1}^{n-1} \frac{1}{\alpha + n - 1} \delta_{\theta_j^*} + \frac{\alpha}{\alpha + n - 1} \mu$

Then the sequence $\theta_1^*, \dots, \theta_n^*$ is a Chinese restaurant process.

- Equivalently define $n_i^{(n-1)} = \sum_{j=1}^{n-1} \mathbb{1}(\theta_j^* = \theta_i)$. Then

$$\theta_n^* | \theta_1^*, \dots, \theta_{n-1}^* \sim \sum_{i=1}^{K_{n-1}^+} \frac{n_i^{(n-1)}}{\alpha + n - 1} \delta_{\theta_i} + \frac{\alpha}{\alpha + n - 1} \mu.$$

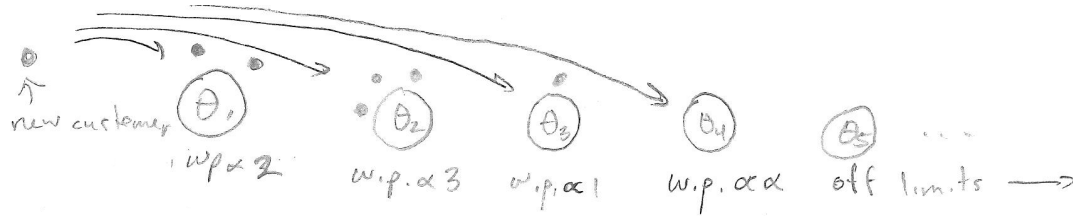
- As $n \rightarrow \infty$, $\frac{\alpha}{\alpha + n - 1} \rightarrow 0$ and

$$\sum_{i=1}^{K_{n-1}^+} \frac{n_i^{(n-1)}}{\alpha + n - 1} \delta_{\theta_i} \rightarrow G = \sum_{i=1}^{\infty} \pi_i \delta_{\theta_i} \sim \text{DP}(\alpha\mu). \quad (7.7)$$

- Notice with the limits that K first went to infinity and then n went to infinity. The resulting empirical distribution is a Dirichlet process because for finite K the de Finetti mixing measure is Dirichlet and the infinite limit of this finite Dirichlet is a Dirichlet process (as we've shown).

Chinese restaurant analogy

- An infinite sequence of tables each have a dish (parameter) placed on it that is drawn iid from μ . The n th customer sits at an occupied table with probability proportional to the number of customers seated there, or selects the first unoccupied table with probability proportional to α .



- “Dishes” are parameters for distributions, a “customer” is a data point that uses its dish to create its value.

Some properties of the CRP

- Cluster growth: What does the number of unique clusters K_n^+ look like as a function of n ?

$$K_n^+ = \sum_{j=1}^n \mathbb{1}(\theta_j^* \neq \theta_{\ell < j}^*) \leftarrow \text{this event occurs when we select a new table}$$

$$\mathbb{E}[K_n^+] = \sum_{j=1}^n \mathbb{E}[\mathbb{1}(\theta_j^* \neq \theta_{\ell < j}^*)] = \sum_{j=1}^n P(\theta_j^* \neq \theta_{\ell < j}^*) = \sum_{j=1}^n \frac{\alpha}{\alpha + j - 1} \approx \alpha \ln n$$

Where does this come from? Each θ_j^* can pick a new table. θ_j^* does so with probability $\frac{\alpha}{\alpha + j - 1}$.

- Cluster sizes: We saw last lecture that if we let the number of observations n be random, where

$$n|y \sim \text{Pois}(y), \quad y \sim \text{gam}(\alpha, \alpha/\hat{n}),$$

then we can analyze cluster size and number using the Poisson process.

- $\mathbb{E}[n] = \mathbb{E}[\mathbb{E}[n|y]] = \mathbb{E}[y] = \hat{n} \leftarrow$ expected number of observations
- $K_n^+ \sim \text{Pois}(\alpha \ln(\alpha + n) - \alpha \ln \alpha) \leftarrow$ total # clusters
- Therefore $\mathbb{E}[K_n^+] = \alpha \ln((\alpha + n)/\alpha)$ (compare with above where n is not random)
- $\sum_{i=1}^{K_n^+} \mathbb{1}(n_i^{(n)} = k) \sim \text{Pois}\left(\frac{\alpha}{k} \left(\frac{n}{\alpha + n}\right)^k\right) \leftarrow$ number of clusters with k observations

- It’s important to remember that n is random here. So in #2 and #3 above, we *first* generate n and then sample this many times from the CRP. For example, $\mathbb{E}[K_n^+]$ is slightly different depending on whether n is random or not.

Inference for the CRP

- **Generative process:** $X_n|\theta_n^* \sim p(X|\theta_n^*)$, $\theta_n^*|\theta_{1:n-1}^* \sim \sum_{i=1}^{n-1} \frac{1}{\alpha+n-1} \delta_{\theta_i^*} + \frac{\alpha}{\alpha+n-1} \mu$. $\alpha > 0$ is “concentration” parameter and we assume μ is a non-atomic probability measure.
- **Posterior inference:** Given the data X_1, \dots, X_N and parameters α and μ , the goal of inference is to perform the inverse problem of finding $\theta_1^*, \dots, \theta_N^*$. This gives the unique parameters $\theta_{1:K_N} = \text{unique}(\theta_{1:N}^*)$ and the partition of the data into clusters.
- Using Bayes rule doesn’t get us far (recall that $p(B|A) = \frac{p(A|B)p(B)}{p(A)}$).

$$p(\theta_1, \theta_2, \dots, \theta_{1:N}^* | X_{1:N}) = \left[\prod_{j=1}^N p(X_j | \theta_j^*) \right] p(\theta_{1:N}^*) \prod_{i=1}^{\infty} p(\theta_i) / \text{intractable normalizer} \quad (7.8)$$

- **Gibbs sampling:** We can’t calculate the posterior analytically, but we can sample from it: Iterate between sampling the atoms given the assignments and then sampling the assignments given the atoms.

Sampling the atoms $\theta_1, \theta_2, \dots$

- This is the easier of the two. For the K_N unique clusters in $\theta_1^*, \dots, \theta_N^*$ at iteration t , we need to sample $\theta_1, \dots, \theta_{K_N}$.

Sample θ_i : Use Bayes rule,

$$p(\theta_i | \theta_{-i}, \theta_{1:N}^*, X_{1:N}) \propto \underbrace{\left[\prod_{j:\theta_j^*=\theta_i} p(X_j | \theta_i) \right]}_{\text{likelihood}} \times \underbrace{p(\theta_i)}_{\text{prior } (\mu)} \quad (7.9)$$

- In words, the posterior of θ_i depends only on the data assigned to the i th cluster according to $\theta_1^*, \dots, \theta_N^*$.
- We simply select this subset of data and calculate the posterior of θ_i on this subset. When μ and $p(X|\theta)$ are conjugate, this is easy.

Sampling θ_j^* (seating assignment for X_j)

- Use exchangeability of $\theta_1^*, \dots, \theta_N^*$ to treat X_j as if it were the last observation,

$$p(\theta_j^* | X_{1:N}, \Theta, \theta_{-j}^*) \propto p(X_j | \theta_j^*, \Theta) p(\theta_j^* | \theta_{-j}^*) \leftarrow \text{also conditions on “future” } \theta_n^* \quad (7.10)$$

- Below is the sampling algorithm followed by the mathematical derivation

$$\text{set } \theta_j^* = \begin{cases} \theta_i & \text{w.p. } \propto p(X_j | \theta_i) \sum_{n \neq j} \mathbf{1}(\theta_n^* = \theta_i), \quad \theta_i \in \text{unique}\{\theta_{-j}^*\} \\ \theta_{\text{new}} \sim p(\theta | X_j) & \text{w.p. } \propto \alpha \int p(X_j | \theta) p(\theta) d\theta \end{cases} \quad (7.11)$$

- The first line should be straightforward from Bayes rule. The second line is trickier because we have to account for the infinitely remaining parameters. We'll discuss the second line next.
- First, return to the finite model and then take the limit (and assume the appropriate re-indexing).
- Define: $n_i^{-j} = \#\{\theta_n^* : \theta_n^* = \theta_i, n \neq j\}$, $K_{-j} = \#\text{unique}\{\theta_{-j}^*\}$.
- Then the prior on θ_j^* is

$$\theta_j^* | \theta_{-j}^* \sim \sum_{i=1}^{K_{-j}} \frac{n_i^{-j}}{\alpha + n - 1} \delta_{\theta_i} + \frac{\alpha}{\alpha + n - 1} \sum_{i=1}^K \frac{1}{K} \delta_{\theta_i} \quad (7.12)$$

- The term $\sum_{i=1}^K \frac{1}{K} \delta_{\theta_i}$ overlaps with the K_{-j} atoms in the first term, but we observe that $K_{-j}/K \rightarrow 0$ as $K \rightarrow \infty$.
- First: What's the probability a new atom is used in the infinite limit ($K \rightarrow \infty$)?

$$p(\theta_j^* = \theta_{new} | X_j, \theta_{-j}^*) \propto \lim_{K \rightarrow \infty} \alpha \sum_{i=1}^K \frac{1}{K} p(X_j | \theta_i) \quad (7.13)$$

Since $\theta_i \stackrel{iid}{\sim} \mu$,

$$\lim_{K \rightarrow \infty} \sum_{i=1}^K \frac{1}{K} p(X_j | \theta_i) = \mathbb{E}_\mu[p(X_j | \theta)] = \int p(X_j | \theta) \mu(d\theta). \quad (7.14)$$

Technically, this is the probability that an atom is selected from the second part of (7.12) above. We'll see why this atom is therefore "new" next.

- Second: Why is $\theta_{new} \sim p(\theta | X_j)$? (And why is it new to begin with?)
- Given that $\theta_j^* = \theta_{new}$, we need to find the index i so that $\theta_{new} = \theta_i$ from the second half of (7.12).

$$\begin{aligned} p(\theta_{new} = \theta_i | X_j, \theta_j^* = \theta_{new}) &\propto p(X_j | \theta_{new} = \theta_i) p(\theta_{new} = \theta_i | \theta_j^* = \theta_{new}) \\ &\propto \lim_{K \rightarrow \infty} p(X_j | \theta_i) \frac{1}{K} \Rightarrow p(X_j | \theta) \mu(d\theta) \end{aligned} \quad (7.15)$$

So $p(\theta_{new} | X_j) \propto p(X_j | \theta) \mu(d\theta)$.

Therefore, given that the atom associated with X_j is selected from the second half of (7.12), the probability it coincides with an atom in the first half equals zero (and so it's "new" with probability one). Also, the atom itself is distributed according to the posterior given X_j .

Chapter 8

Exchangeability, beta processes and the Indian buffet process

Marginalizing (integrating out) stochastic processes

- We saw how the Dirichlet process gives a discrete distribution on model parameters in a clustering setting. When the Dirichlet process is integrated out, the cluster assignments form a Chinese restaurant process:

$$\underbrace{p(\theta_1^*, \dots, \theta_N^*)}_{\text{Chinese restaurant process}} = \int \underbrace{\prod_{n=1}^N p(\theta_n^* | G)}_{\substack{\text{i.i.d. from discrete dist.} \\ \text{DP}}} \underbrace{p(G)}_{\text{DP}} dG \quad (8.1)$$

- There is a direct parallel between the beta-Bernoulli process and the “Indian buffet process”:

$$\underbrace{p(Z_1, \dots, Z_N)}_{\text{Indian buffet process}} = \int \underbrace{\prod_{n=1}^N p(Z_n | H)}_{\substack{\text{Bernoulli process} \\ \text{BP}}} \underbrace{p(H)}_{\text{BP}} dH \quad (8.2)$$

- As with the DP→CRP transition, the BP→IBP transition can be understood from the limiting case of the finite BP model.

Beta process (finite approximation)

- Let $\alpha > 0$, $\gamma > 0$ and μ a non-atomic probability measure. Define

$$H_K = \sum_{i=1}^K \pi_i \delta_{\theta_i}, \quad \pi_i \sim \text{Beta}(\alpha \frac{\gamma}{K}, \alpha(1 - \frac{\gamma}{K})), \quad \theta_i \sim \mu. \quad (8.3)$$

Then $\lim_{K \rightarrow \infty} H_K = H \sim \text{BP}(\alpha, \gamma\mu)$. (See Chapter 3 for proof.)

Bernoulli process using H_K

- Given H_K , we can draw the Bernoulli process $Z_n^K | H_K \sim BeP(H_K)$ as follows:

$$Z_n^K = \sum_{i=1}^K b_{in} \delta_{\theta_i}, \quad b_{in} \sim \text{Bernoulli}(\pi_i). \quad (8.4)$$

Notice that b_{in} should also be marked with K , which we ignore. Again we are particularly interested in $\lim_{K \rightarrow \infty} (Z_1^K, \dots, Z_N^K)$.

- To derive the IBP, we first consider

$$\lim_{K \rightarrow \infty} p(Z_{1:N}^K) = \lim_{K \rightarrow \infty} \int \prod_{n=1}^N p(Z_n^K | H_K) p(H_K) dH_K. \quad (8.5)$$

- We can think of Z_1^K, \dots, Z_N^K in terms of a binary matrix, $B_K = [b_{in}]$, where
 - each row corresponds to an atom, θ_i and $b_{in} \stackrel{iid}{\sim} \text{Bern}(\pi_i)$ for row i
 - each column corresponds to a Bernoulli process, Z_n^K for column n

Important: The rows of B are *independent* processes.

- Consider the process $b_{in} | \pi_i \stackrel{iid}{\sim} \text{Bern}(\pi_i)$, $\pi_i \sim \text{Beta}(\alpha \frac{\gamma}{K}, \alpha(1 - \frac{\gamma}{K}))$. The marginal process b_{i1}, \dots, b_{iN} follows an urn model with two colors.

Polya's urn (two-color special case)

- Start with an urn having $\alpha\gamma/K$ balls of color 1 and $\alpha(1 - \gamma/K)$ balls of color 2.
- Pick a ball at random, pit it back and put a second one of the same color

Mathematically:

- $b_{i1} \sim \frac{\gamma}{K} \delta_1 + (1 - \frac{\gamma}{K}) \delta_0$
- $b_{i,N+1} | b_{i1}, \dots, b_{iN} \sim \frac{\frac{\alpha\gamma}{K} + n_i^{(N)}}{\alpha + N} \delta_1 + \frac{\alpha(1 - \frac{\gamma}{K}) + N - n_i^{(N)}}{\alpha + N} \delta_0$

where $n_i^{(N)} = \sum_{j=1}^N b_{ij}$. Recall from exchangeability and de Finetti that

$$\lim_{K \rightarrow \infty} \frac{n_i^{(N)}}{N} \longrightarrow \pi_i \sim \text{Beta}(\alpha \frac{\gamma}{K}, \alpha(1 - \frac{\gamma}{K})) \quad (8.6)$$

- Last week we proved this in the context of the finite symmetric Dirichlet, $\pi \sim \text{Dir}(\frac{\alpha}{K}, \dots, \frac{\alpha}{K})$. The beta distribution is the two-dimensional special case of the Dirichlet and the proof can be applied to any parameterization besides symmetric.

- In the Dirichlet→CRP limiting case, K corresponds to the number of colors in the urn and $K \rightarrow \infty$ with the starting number of each color $\frac{\alpha}{K} \rightarrow 0$.
- In this case, there are always only two colors. However, the *number of urns* equals K , so the number of urn processes is going to infinity as the first parameter of each beta goes to zero.

An intuitive derivation of the IBP: Work with the urn representation of B_K . Again, b_{in} is entry (i, n) of B_K and the generative process for B_K is

$$b_{i,N+1}|b_{i1}, \dots, b_{iN} \sim \frac{\frac{\alpha\gamma}{K} + n_i^{(N)}}{\alpha + N} \delta_1 + \frac{\alpha(1 - \frac{\gamma}{K}) + N - n_i^{(N)}}{\alpha + N} \delta_0 \quad (8.7)$$

where $n_i^{(N)} = \sum_{j=1}^N b_{ij}$. Each row of B_K is associated with a $\theta_i \sim_{iid} \mu$ so we can reconstruct $Z_n^K = \sum_{i=1}^K b_{in} \delta_{\theta_i}$ using what we have.

- Let's break down $\lim_{K \rightarrow \infty} B_K$ into two cases.

Case $n = 1$: We ask how many ones are in the first column of B ?

$$\lim_{K \rightarrow \infty} \sum_{i=1}^K b_{i1} \sim \lim_{K \rightarrow \infty} \text{Bin}(K, \gamma/K) = \text{Pois}(\gamma) \quad (8.8)$$

So Z_1 has $\text{Pois}(\gamma)$ ones. Since the θ associated with these ones are i.i.d., we can “ex post facto” draw them i.i.d. from μ and re-index.

Case $n > 1$: For the remaining Z_n , we break this into two subcases.

- Subcase $n_i^{(n-1)} > 0$: $b_{in}|b_{i1}, \dots, b_{i,n-1} \sim \frac{n_i^{(n-1)}}{\alpha + n - 1} \delta_1 + \frac{\alpha + n - 1 - n_i^{(n-1)}}{\alpha + n - 1} \delta_0$
- Subcase $n_i^{(n-1)} = 0$: $b_{in}|b_{i1}, \dots, b_{i,n-1} \sim \lim_{K \rightarrow \infty} \frac{\alpha \frac{\gamma}{K}}{\alpha + n - 1} \delta_1 + \frac{\alpha(1 - \frac{\gamma}{K}) + n - 1}{\alpha + n - 1} \delta_0$
- For each i , $\left(\frac{\alpha\gamma}{\alpha+n-1}\right) \frac{1}{K} \delta_1 \rightarrow 0 \delta_1$, but there are also an infinite number of these indexes i for which $n_i^{(n-1)} = 0$. Is there a limiting argument we can again make to just ask how many ones there are total for these indexes with $n_i^{(n-1)} = 0$?
- Let $K_n = \#\{i : n_i^{(n-1)} > 0\}$, which is finite almost surely. Then

$$\lim_{K \rightarrow \infty} \sum_{i=1}^K b_{in} \mathbb{1}(n_i^{(n-1)} = 0) \sim \lim_{K \rightarrow \infty} \text{Bin}\left(K - K_n, \left(\frac{\alpha\gamma}{\alpha+n-1}\right) \frac{1}{K}\right) = \text{Pois}\left(\frac{\alpha\gamma}{\alpha+n-1}\right) \quad (8.9)$$

- So there are $\text{Pois}\left(\frac{\alpha\gamma}{\alpha+n-1}\right)$ new locations for which $b_{in} = 1$. Again, since the atoms are i.i.d. regardless of the index, we can simply draw them i.i.d. from μ and re-index.

Putting it all together: The Indian buffet process

- For $n = 1$: Draw $C_1 \sim \text{Pois}(\gamma)$, $\theta_1, \dots, \theta_{C_1} \stackrel{iid}{\sim} \mu$ and set $Z_1 = \sum_{i=1}^{C_1} \delta_{\theta_i}$.
- For $n > 1$: Let $K_{n-1} = \sum_{j=1}^{n-1} C_j$. For $i = 1, \dots, K_{n-1}$, draw

$$b_{in} | b_{i,1:n-1} \sim \frac{n_i^{(n-1)}}{\alpha + n - 1} \delta_1 + \frac{\alpha + n - 1 - n_i^{(n-1)}}{\alpha + n - 1} \delta_0. \quad (8.10)$$

Then draw $C_n \sim \text{Pois}\left(\frac{\alpha\gamma}{\alpha+n-1}\right)$ and $\theta_{K_{n-1}+1}, \dots, \theta_{K_{n-1}+C_n} \stackrel{iid}{\sim} \mu$ and set

$$Z_n = \sum_{i=1}^{K_{n-1}} b_{in} \delta_{\theta_i} + \sum_{i'=K_{n-1}+1}^{K_{n-1}+C_n} \delta_{\theta_{i'}}. \quad (8.11)$$

- By exchangeability and de Finetti, $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N Z_i \rightarrow H \sim \text{BP}(\alpha, \gamma\mu)$.

The IBP story: Start with α customers not eating anything.

1. A customer walks into an Indian buffet with an infinite number of dishes and samples $\text{Pois}(\gamma)$ of them.
 2. The n th customer arrives and samples from the previously sampled dishes with probability proportional to the number of previous customers who sampled it, and then samples $\text{Pois}\left(\frac{\alpha\gamma}{\alpha+n-1}\right)$ new dishes.
- In modeling scenarios, each dish corresponds to a factor (e.g., a one-dimensional subspace) and a customer samples a subset of factors.
 - Clearly, after n customers there are $\sum_{i=1}^n C_i \sim \text{Pois}\left(\sum_{i=1}^n \frac{\alpha\gamma}{\alpha+i-1}\right)$ dishes that have been sampled. (See Chapter 3 for another derivation of this quantity.)

Chapter 9

Another look at constructive definitions of the beta and Dirichlet process

- It isn't always obvious how equivalent representations for stochastic processes are arrived at, such as constructions for the Dirichlet and beta process. Often the proof of correctness requires the statement of equivalence as a starting point. We'll next look at alternative methods for deriving the constructions we've looked at that don't require the construction as a starting point.

BP constructions and the IBP

- From Chapter 4: Let $\alpha, \gamma > 0$ and μ a non-atomic probability measure. Let

$$C_i \sim \text{Pois}\left(\frac{\alpha\gamma}{\alpha + i - 1}\right), \quad \pi_{ij} \sim \text{Beta}(1, \alpha + i - 1), \quad \theta_{ij} \sim \mu \quad (9.1)$$

and define $H = \sum_{i=1}^{\infty} \sum_{j=1}^{C_i} \pi_{ij} \delta_{\theta_{ij}}$. Then $H \sim \text{BP}(\alpha, \gamma\mu)$.

- We proved this using Poisson process theory. Now we'll show how this construction can be arrived at in the first place.
- Imagine an urn with β_1 balls of one color and β_2 of another. The distribution on the first draw from this urn is

$$\frac{\beta_1}{\beta_1 + \beta_2} \delta_1 + \frac{\beta_2}{\beta_1 + \beta_2} \delta_0.$$

Let (b_1, b_2, \dots) be sequence generated from an urn process with this initial configuration. Then as we have seen, in the limit $N \rightarrow \infty$, $\frac{1}{N} \sum_{i=1}^N b_i \rightarrow \pi \sim \text{Beta}(\beta_1, \beta_2)$.

- Key: We can skip the whole urn procedure if we're only interested in π . That is, if we draw from the urn once, look at it and see it's color one, then the urn distribution is

$$\frac{\beta_1 + 1}{1 + \beta_1 + \beta_2} \delta_1 + \frac{\beta_2}{1 + \beta_1 + \beta_2} \delta_0$$

We can ask: what will happen if we continue after this first draw? What is the difference between this "posterior" configuration and another urn where this is defined to be the initial setting? It

shouldn't be hard to convince yourself that, in this case, as $N \rightarrow \infty$,

$$\frac{1}{N} \sum_{i=1}^N b_i | \{b_1 = 1\} \rightarrow \pi \sim \text{Beta}(\beta_1 + 1, \beta_2).$$

- This is because the sequence $(b_1 = 1, b_2, b_3, \dots)$ is equivalent to an urn process (b_2, b_3, \dots) where the initial configuration is $\beta_1 + 1$ of color 1 and β_2 of color 0.
- Furthermore, we know from de Finetti and exchangeability that if Z_1, \dots, Z_N are from an IBP, then $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N Z_i \rightarrow H \sim \text{BP}(\alpha, \gamma\mu)$. Right now we're only interested in this limit.

Derivation of the construction via the IBP

- For each Z_n , there are $C_n \sim \text{Pois}\left(\frac{\alpha\gamma}{\alpha+n-1}\right)$ new “urns” introduced with the initial configuration $\delta_1 + (\alpha + n - 1)\delta_0$, where δ_1 indicates the corresponding atom was used (it doesn't matter what that atom turned out to be). The initial distribution for the next draw is the normalization of this. From this point on, each new urns can be treated as independent processes. Notice that this is the case for every instantiation of the IBP.
- With the IBP, we continue this urn process. Instead, we can ask the limiting distribution of the urn immediately after instantiating it. We know from above that the new urns created at step n will converge to random variables drawn independently from a $\text{Beta}(1, \alpha + n - 1)$ distribution. We can draw this probability directly.
- The results is the construction written above.

Dirichlet stick-breaking construction and the CRP

- What about stick-breaking for the Dirichlet process? It turns out that knowledge about the urn representation of the DP provided by the Chinese restaurant process suggests this in a way very similar to how the IBP suggested the above construction of the beta process.

- First, what do we know about the CRP? We know that if $\theta_1^* \sim \mu$ and

$$\theta_{N+1}^* | \theta_{1:N}^* \sim G_N \equiv \frac{n_1^N}{\alpha + N} \delta_{\theta_1} + \frac{n_2^N}{\alpha + N} \delta_{\theta_2} + \dots + \frac{n_{K_N}^N}{\alpha + N} \delta_{\theta_{K_N}} + \frac{\alpha}{\alpha + N} \mu$$

then $\lim_{N \rightarrow \infty} G_N = G \sim \text{DP}(\alpha\mu)$. Here we've defined θ_i to be the i th unique atom generated by this process, n_i^N to be the number of atoms in $\theta_{1:N}^*$ equal to θ_i after N observations, and K_N to be the number of unique atoms contained in $\theta_{1:N}^*$.

- Now consider the first atom $\theta_1^* \sim \mu$. It might sound odd to say, but this is a probability one event. Therefore, we know that we will always have an urn that looks like $\alpha\mu + \delta_{\theta_1}$ where $\theta_1 \sim \mu$. The only thing that's random is what “color” the first observation θ_1 has, not that we start with an urn containing one “ball” equal to color θ_1 .

- Now ask: If I ran out this process and only cared about $\lim_{N \rightarrow \infty} \frac{n_1^N}{\alpha + N} \delta_{\theta_1}$ what do I know? It turns out that we can definitively say that $\lim_{N \rightarrow \infty} \frac{n_1^N}{\alpha + N} = V_1 \sim \text{Beta}(1, \alpha)$. Therefore, we are in a situation where

$$G_N \rightarrow V_1 \delta_{\theta_1} + (1 - V_1) G', \quad V_1 \sim \text{Beta}(1, \alpha), \quad \theta_1 \sim \mu$$

Similarly, as just stated, we know that $\lim_{N \rightarrow \infty} G_N = G \sim \text{DP}(\alpha \mu)$, so $V_1 \delta_{\theta_1} + (1 - V_1) G'$ must also be a Dirichlet process. The only question now is, what does G' equal to?

- It turns out that G' is also the limit of a CRP *that is independent of V_1 and θ_1* . I will leave it as a statement and not try to rigorously prove it further (since again this chapter is more intended to show the intuition of deriving rather than proving a construction).
- However, we're in the same recursive setting as when we first derived the stick-breaking construction. If we know that G' is the limit of an independent CRP, then we know that it will start with an urn that looks like $\alpha \mu + \delta_{\theta_2}$, where $\theta_2 \sim \mu$ is the atom generated the second time we choose to draw from μ , and the argument repeats.

- After the second nesting, we have

$$G_N \rightarrow V_1 \delta_{\theta_1} + (1 - V_1) V_2 \delta_{\theta_2} + (1 - V_1)(1 - V_2) G'', \quad V_1, V_2 \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad \theta_1, \theta_2 \stackrel{iid}{\sim} \mu$$

and in the limit:

$$G = \sum_{i=1}^{\infty} V_i \prod_{j=1}^{i-1} (1 - V_j) \delta_{\theta_i}, \quad V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha), \quad \theta_i \stackrel{iid}{\sim} \mu$$

Constructing beta processes by thinking about Dirichlet processes

- We next discuss a method for deriving the other stick-breaking construction of the beta process. This approach requires knowledge about the stick-breaking construction of the Dirichlet distribution (with the beta as a special case).
- Let μ be a non-atomic measure on S with $\gamma = \mu(S) < \infty$ and $\alpha > 0$. We saw that the following is a constructive definition of the beta process $H \sim \text{BP}(\alpha, \mu)$,

$$H = \sum_{i=1}^{\infty} \sum_{j=1}^{C_i} V_{ij}^{(i)} \prod_{l=1}^{i-1} (1 - V_{ij}^{(l)}) \delta_{\theta_{ij}}, \quad (9.2)$$

$$C_i \stackrel{iid}{\sim} \text{Pois}(\gamma), \quad V_{ij}^{(l)} \stackrel{ind}{\sim} \text{Beta}(1, \alpha), \quad \theta_{ij} \stackrel{iid}{\sim} \mu/\gamma.$$

- This construction sequentially incorporates into H a Poisson-distributed number of atoms drawn i.i.d. from μ/γ , with each group in this sequence indexed by i . The atoms receive weights as follows: Using an atom-specific stick-breaking construction, an atom in group i throws away the first $i - 1$ breaks of its stick and keeps the i th break as its weight.

Derivation via the finite approximation

- We prove the construction by constructing finite arrays of random variables and considering their limit. Remember that the finite approximation is

$$H_K = \sum_{k=1}^K \pi_k \delta_{\theta_k}, \quad \pi_k \sim \text{Beta}\left(\alpha \frac{\gamma}{K}, \alpha \left(1 - \frac{\gamma}{K}\right)\right), \quad \theta_k \sim \mu/\gamma$$

- We represent each beta-distributed random variable by its stick-breaking construction. Recall that using this construction, we can draw $\pi \sim \text{Beta}(a, b)$ as follows:
 1. Draw an infinite sequence of random variables (V_1, V_2, \dots) i.i.d. from $\text{Beta}(1, a + b)$ and a second sequence (Y_1, Y_2, \dots) i.i.d. $\text{Bern}\left(\frac{a}{a+b}\right)$.
 2. Construct $\pi_R = \sum_{i=1}^R V_i \prod_{j=1}^{i-1} (1 - V_j) \mathbb{1}(Y_i = 1)$.
 3. Then as $R \rightarrow \infty$, $\pi_R \rightarrow \pi \sim \text{Beta}(a, b)$.
- Using a truncated stick-breaking construction of the beta random variables in H_K , we can state the following: Draw two $K \times R$ arrays of independent random variables, $V_{ki} \sim \text{Beta}(1, \alpha)$ and $Y_{ki} \sim \text{Bern}(\gamma/K)$, and draw $\theta_k \sim \mu/\gamma$ for $k = 1, \dots, K$ and $i = 1, \dots, R$. Let

$$H_K^{(R)} = \sum_{k=1}^K \sum_{i=1}^R V_{ki} \prod_{j=1}^{i-1} (1 - V_{kj}) \mathbb{1}(Y_{ki} = 1) \delta_{\theta_k}.$$

Then $H_K^{(R)}$ converges in distribution to $H \sim \text{BP}(\alpha, \mu)$ by letting $K \rightarrow \infty$ and $R \rightarrow \infty$.

- This leads to the following thought process:
 - We know that in the limit as $K \rightarrow \infty$ and $R \rightarrow \infty$, $H_K^{(R)} \rightarrow H \sim \text{BP}(\alpha, \mu)$.
 - We first note that column sums of Y are marginally distributed as $\text{Bin}(K, \gamma/K)$, and are independent. The i th column sum value gives the number of atoms that receive probability mass at step i , with $Y_{ki} = 1$ indicating the k th indexed atom is one of them.
 - Let the set $\mathcal{I}_i^K = \{k : Y_{ki} = 1\}$ be the index set of these atoms at finite approximation level K . This set is constructed by selecting $C_i^K \sim \text{Bin}(K, \gamma/K)$ values from $\{1, \dots, K\}$ uniformly without replacement. In the limit $K \rightarrow \infty$, $C_i^K \rightarrow C_i$ with $C_i \sim \text{Pois}(\gamma)$.
 - Given $k \in \mathcal{I}_i^K$, we know that π_k has weight added to it from the i th break of its own stick-breaking construction. As a matter of accounting, we are interested in other values j for which $Y_{kj} = 1$, particularly when $j < i$.
- We next show that in the limit $K \rightarrow \infty$, the index values in the set $\mathcal{I}_i := \mathcal{I}_i^\infty$ are always unique from those in previous sets ($j < i$), meaning for a given column $L \leq R$, we see new index values with probability equal to one. We are therefore always adding probability mass to new atoms. This is significant because we can therefore create new sticks and break them on the fly, while letting the uniquely created atom for that weight be a proxy for the index that would have been chosen.

- Let E be the event that there exists a value $k \in \mathcal{I}_i \cap \mathcal{I}_j$ for some $i \neq j$ and $i, j \leq L \leq R$. Then $P_L(E) = \lim_{K \rightarrow \infty} P_L(\bigcup_{j < i \leq L} \mathcal{I}_i^K \cap \mathcal{I}_j^K \neq \emptyset | \mu)$. We can bound the probability of P_L as follows:

$$\begin{aligned}
 P_L(\bigcup_{j < i \leq L} \mathcal{I}_i^K \cap \mathcal{I}_j^K \neq \emptyset | \mu) &\leq \sum_{j < i \leq L} P_L(\mathcal{I}_i^K \cap \mathcal{I}_j^K \neq \emptyset | \mu) \\
 &\leq \sum_{j < i \leq L} \sum_{k=1}^K P_L(Y_{ki} Y_{kj} = 1 | \mu) \\
 &\leq \frac{L(L-1)}{2} \frac{\gamma^2}{K}.
 \end{aligned} \tag{9.3}$$

Therefore, for any finite integer $L \leq R$, in the limit $K \rightarrow \infty$ the atoms $\theta_k, k \in \mathcal{I}_L$, are different from all previously observed atoms with probability one since μ is a diffuse measure. Since this doesn't depend on R , we can just let $R \rightarrow \infty$ next.