- Clustering - How to choose $K$?

  - Ad hoc approaches:   data-driven, but little theory

  - Methods based on maximizing objective funs (M-estimators)
    e.g. clustering via $K$-means, spectral clustering
         has frequentist guarantees but hard to be data-driven.

  - Parametric methods such as AIC, BIC requires
    hidden assumptions

- Model-based approach:
  $$P(X|\phi) = \sum_{k=1}^{K} \pi_k \, p(x|\phi_k)$$

  Let $G = \sum_{k=1}^{K} \pi_k \delta_{\phi_k}$       mixing measure.

  Model:  $\phi \,|\, G \sim G$
          $x \,|\, \phi \sim p(x|\phi)$

- $G$ Random. Need a prior dist on $G$.
  want $K$ unbounded $(K = \infty) \implies$ Stick-breaking prior.
  This is a dist over $G$ of form:
  $$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k} \qquad \begin{cases} \pi & \text{Random} \\ \phi & \text{Random} \\ \sum_{1}^{\infty} \pi_k = 1 \end{cases}$$
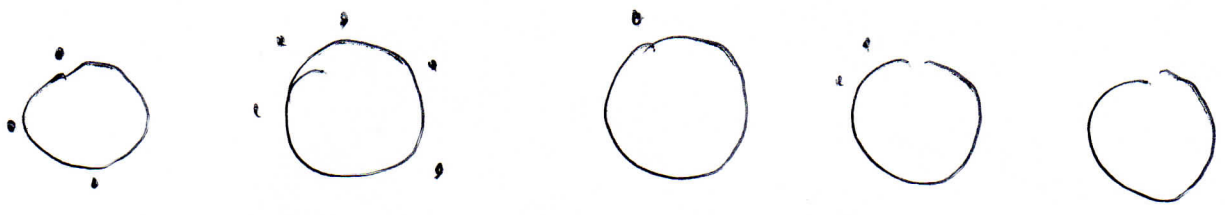
- <u>Chinese Restaurant Process</u> (CRP)
  A Random process in which Customers sit down in a Chinese Restaurant
  with infinite number of tables.

CRP: - first customer sits at the first table

- $m$-th subsequent customer sits at a table drawn from the following dist:

$\mathbb{P}$ ( sits at previously occupied table $i$ ) $\propto n_i$

$\mathbb{P}$ ( sits in a new (unoccupied) table ) $\propto d_0$

$n_i$ : # customers sitting at table $i$ ( among previous $m-1$ customers )



## CRP and Clustering.

CRP induces a √ partition of customers into disjoint clusters (subsets)

dist on and # of clusters (tables).

Let $\pi_k$ be the proportion of customers sitting at table $k$.

$\phi_k \overset{iid}{\sim} G_0$ (because $\pi_k$'s, $\phi_k$'s are).

Then we obtained a [Random] measure: $G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$

what can we say about $G$?

## Polya urn model

CRP is a particular example of a general class of prob. models known as Polya urn.

Let $\theta_1, \dots, \theta_n \overset{iid}{\sim} G$. there are duplicates among the $\theta_i$'s $\Rightarrow$ the Polya urn model is as follows:

$$\theta_1 \sim G_0$$

$$\theta_i \mid \theta_1, \dots, \theta_{i-1} \sim d_0 G_0 + \sum_{j=1}^{i-1} \delta_{\theta_j}$$

$$P(\theta_1, \theta_2) = P(\theta_1)\, P(\theta_2 | \theta_1)$$

$$\propto P(\theta_1) \left( \alpha_0 \, G_0(d\theta_2) + \sum_j \delta_{\theta_1}(\theta_2) \right)$$

$$\propto \alpha_0 \, G_0(d\theta_1)\, G_0(d\theta_2) + G_0(d\theta_1)\, \mathbb{1}(\theta_1 = \theta_2).$$

$$\Rightarrow P(\theta_1, \theta_2) = P(\theta_2, \theta_1).$$

$$\mathbb{P}(\theta_1, \theta_2, \theta_3) = \mathbb{P}(\theta_1, \theta_2)\, \mathbb{P}(\theta_3 | \theta_1, \theta_2)$$
$$= \mathbb{P}(\theta_2, \theta_1)\, \mathbb{P}(\theta_3 | \theta_2, \theta_1) = \mathbb{P}(\theta_2, \theta_1, \theta_3)$$

in general $\theta_1, \theta_2 \ldots, \theta_n \ldots$ are infinitely exchangeable!

By De Finetti's thm $\exists$ a random mixing dist $G$ s.t

$$\theta_1, \ldots \theta_n | G \overset{iid}{\sim} G.$$

So that $\quad P(\theta_1 \ldots \theta_n) = \int \prod_{i=1}^{n} P(\theta_i | G)\, d\Pi(G)$

what is $G$?

* **Stick- breaking process :** $\boxed{\{\pi_k\}_{k \geq 1}}$

Let $\beta_k \overset{iid}{\sim} Beta(1, \alpha_0) \qquad k = 1, 2, \ldots$

Define $\quad \pi_1 = \beta_1$
$$\pi_k = \beta_k \prod_{\ell=1}^{k-1}(1 - \beta_\ell) \qquad k = 2, 3, \ldots$$

Easy to check that $\quad \sum_{k=1}^{\infty} \pi_k = 1.$ $\qquad$ ( Breaking a stick to small pieces using beta proportions )

$$1 - \sum_{k=1}^{K} \pi_k = 1 - \beta_1 - \beta_2(1-\beta_1) - \beta_3(1-\beta_1)(1-\beta_2)$$
$$= (1-\beta_1)(1-\beta_2) - \beta_3(1-\beta_1)(1-\beta_2)$$
$$= \cdots$$
$$= (1-\beta_1) \cdots (1-\beta_K) \to 0 \text{ as } k \to \infty$$

Thm: Let $G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}$ where $\phi_k \overset{iid}{\sim} G_0$ then this is the dist $G$ that defines Polya's urn

in fact the Random measure $G$ is called a Dirichlet Process,
where does Dirichlet come in? Hint: Beta dist was used

Recall Dirichlet Dist (finite) & conjugacy (Dir.-Multinomial),

$$P(\pi | \alpha) = \frac{\Gamma(\Sigma \alpha_i)}{\prod \Gamma(\alpha_i)} \pi_1^{\alpha_1 - 1} \dots \pi_K^{\alpha_K - 1} \qquad (K < \infty)$$

$\pi | \alpha \sim Dir(\alpha)$

Let $Y_1, \dots, Y_n | \pi \overset{iid}{\sim} Mult(\pi) \qquad Y_i \in \{1, \dots, k\}$

Then $P(Y_1, \dots, Y_n) = \int \prod_{i=1}^{n} P(Y_i | \pi) \, dP(\pi | \alpha)$

$$\propto \int \prod_{k=1}^{K} \pi_k^{n_k} \prod_{k=1}^{K} \pi_k^{\alpha_k - 1} \, d\pi = \propto \int \prod_{k=1}^{K} \pi_k^{n_k + \alpha_k - 1} \, d\pi$$

$$= \frac{\Gamma(\alpha_1 + \dots + \alpha_K)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \cdot \frac{\Gamma(\alpha_1 + n_1) \dots \Gamma(\alpha_K + n_K)}{\Gamma(\alpha_1 + \dots + \alpha_K + n)}$$

$$P(Y_n | Y_1 \dots Y_{n-1}) = \frac{\Gamma(\alpha_1 + n_1) \dots \Gamma(\alpha_K + n_K)}{\Gamma(\alpha_1 + \dots + \alpha_K + n)} \cdot \frac{\Gamma(\alpha_1 + \dots + \alpha_K + n - 1)}{\Gamma(\alpha_1 + \tilde{n}_1) \dots \Gamma(\alpha_K + \tilde{n}_K)}$$

$$= \frac{\alpha_k + n_k}{\Sigma \alpha_k + n} \qquad \text{if } Y_n = k.$$

$$P(Y_n | Y_1 \dots Y_{n-1}) \propto \alpha_k + n_k \qquad \text{if } Y_n = k$$

Suppose that we can somehow let $K \to \infty$, yet $\alpha_k = \frac{\alpha}{K}$
so that $\Sigma \alpha_k = \alpha$ fixed, then

$$\theta_n | \theta_1 \dots \theta_{n-1} \propto n_k \qquad \text{if } Y_n = k.$$

there is a remaining probability that $Y_n \neq Y_1 \dots Y_{n-1}$

that probability $\longrightarrow \frac{\alpha}{\alpha + n}$ .

How is this related to ~~DP~~ ? } • As $K \to \infty$, $\pi \longrightarrow$ ~~Dirichlet~~ Dirichlet process

stick-breaking?

• Reorder $(\pi_1, \dots, \pi_K)$ in decreasing order
then $(\pi_1, \dots, \pi_K) \longrightarrow$ stick breaking process

- Now we're ready to state original def of Dirichler process.

- <u>Stochastic process</u> vs. RV's.

Elementary prob. thy : RV's are function $X: \Omega \longrightarrow \mathbb{R}$ real-valued

$\Omega$ : set of events

$$\mathbb{P}(X \leq t) = \mathbb{P}\left(\{X(\omega) < t\}\right).$$

A stochastic process can be viewed as a dist/measure on space of more complex objects such as functions or measures. therefore , Random function or random measures.

- <u>Gaussian process</u> as a dist on $\{X: T \longrightarrow \mathbb{R}\}$

is defined as a collection of RV's $\{X(t) : \Omega \to \mathbb{R}\}$ $t \in T$.

s.t. every finite collection of $\{X(t)\}_{t \in S}$ is a multivariate Gaussian

and that they are consistent with each other thru marginalization

- <u>Kolmogorov's thm</u> says that then there exists a meaningful dist on the space of function $\{X: T \longrightarrow \mathbb{R}\}$ which is called a Gaussian process..

For Gaussian, consistency is obviously satisfied : $(X_1, X_2) \sim N$ then $X_1$ is also $N$.

- <u>A probability measure</u> (dist) $G$ is a function that map from $\Omega$ to $[0,1]$
a probability space

$A \in \Omega$   $G(A) \equiv$ probability $\omega \in A$ under $G$.

i.e. if $X \sim G$ then $\mathbb{P}(X \in A) = G(A)$.

- To define a Random measure $G$ means the collection $\{G(A)\}_{A \subset \Omega}$ is a collection RV's.

**Def :** ⓐ Let $(\Omega, \mathcal{F})$ a prob. space.

if $G$ is a Random measure on $(\Omega, \mathcal{F})$ s.t

for any partition $(A_1, A_2, \ldots, A_k)$ of $\Omega$, $k \in \mathbb{N}$

$$(G(A_1), G(A_2), \ldots, G(A_k)) \sim \text{Dir}(\alpha G_0(A_1), \ldots, \alpha G_0(A_k)).$$

then $G$ is said to be dist according to a Dirichler proan

$G_0$ is called a base measure (mean measure)

$\alpha$ concentration parameter.

**Now** such DP exists due to Consistency property of finite dim

Dir dist : if $(Y_1, \ldots, Y_k) \sim \text{Dir}(\alpha_1, \ldots, \alpha_k)$, then $(Y_1 + Y_1, \ldots Y_{r_\ell} + \ldots Y_k)$
$$\sim \text{Dir}(\alpha_{1} + \alpha_{r_1}, \ldots, \alpha_{r_\ell} + \alpha_k)$$

**Fact :** if $G \sim DP(\alpha, G_0)$

then $G(A) \sim \text{Beta}(\alpha G_0(A), \alpha(1 - G_0(A)))$

$$\Rightarrow \begin{cases} \mathbb{E} \, G(A) = G_0(A) \\ \text{var} \, G(A) = \dfrac{\alpha G_0(A)(\alpha G_0(A)+1)}{\alpha^2(\alpha+1)} - G_0(A)^2 = \dfrac{G_0(A)(1 - G_0(A))}{\alpha+1} \end{cases}$$

$\alpha \quad G_0$
$$\searrow \swarrow$$
$G$
$$\downarrow$$
$\theta$

**Question :** if $\begin{cases} G \sim DP(\alpha, G_0) \\ \theta \mid G \sim G \end{cases}$
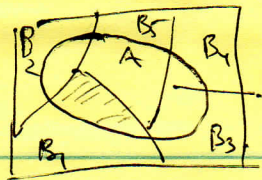
what is the posterior of $[G \mid \theta]$ ?

$\forall A \in \mathcal{F} \quad P(Y \in A \mid G(A)) = G(A)$ a.s.

$\Rightarrow P(X \in A) = \mathbb{E} \, G(A) = G_0(A)$

Take a measurable partition $(G(B_1) \ldots G(B_k)) \sim \text{Dir}(\alpha G_0(B_1) \ldots \alpha G_0(B_k))$

**Thm** $[G(B_1) \ldots G(B_k) \mid x] \sim \text{Dir}(\alpha G_0(B_1) + \delta_{B_1}(x) \ldots )$

$$P\left(X \in A, \ G(B_1) \overset{\leq y_1}{\ldots}, \ G(B_k) \overset{\leq y_k}{}\right)$$

$B_i' = A \cap B_i$

$$= \sum_{k=1}^{K} P\left(X \in B_i', \ G(B_1) \overset{\leq y_1}{\ldots}, \ G(B_k) \overset{\leq y_k}{}\right)$$

$$= \sum_{k=1}^{K} P(X \in B_i') \ P\left(G(B_1) \leq y_1, \ldots, \ G(B_k) \leq y_k \ \middle| \ X \in B_i'\right)$$

$$= \sum_{k=1}^{K} G_0(B_i') \ P_{ir}\left(G(B_i') + G(B_i'') \leq y_1, \ldots \ \middle| \ X \in B_i'\right)$$

Note that $\left(G(B_1'), G(B_1''), \ldots, G(B_k'), G(B_k')\right) \sim Dir\left(\alpha G_0 \cdots \right)$

By Dir-Mult Conjugacy:

$$\left[\left(G(B_1'), G(B_1''), \ldots\right) \ \middle| \ X \in B_i'\right]$$

$$\sim Dir\left(\ldots \alpha G_0(B_i') + \delta_{B_i'} \cdots \right).$$

$$= \sum_{k=1}^{K} G_0(B_i') \ P\left(G(B_1) \leq y_1, \ldots, \ G(B_k) \leq y_k\right) \ \middle| \ \alpha G_0(B_1), \ldots, \alpha G(B_k) \atop + \delta_{B_1} \quad + \delta_{B_k}$$

So $\quad P\left(G(B_1) \leq y_1, \ldots, G(B_k) \leq y_k \ \middle| \ X \in A\right)$

$$= \sum_{k=1}^{R} \frac{G_0(B_i')}{G_0(A)} \overset{\leftarrow A \cap B_i}{} \ P\left(G(B_1) \leq y_1, \ldots, G(B_k) \leq y_k\right) \ \middle| \ \alpha G_0(B_1) + \delta_{B_1}) \atop \alpha G_0(B_k) + \delta_{B_k})$$

if $A \subset B_i$ then $\quad P\left(G(B_1) \leq y_1, \ldots G(B_k) \leq y_k \ \middle| \ X \in A\right)$

$$= \quad P\left(G(B_1) \leq y_1, \ldots G(B_k) \leq y_k \ \middle| \ \alpha G_0(B_1) + \delta_{B_0} \cdots \right)$$

Since $\quad \left(G(B_1), \ldots, G(B_k)\right) \Big|_{X \in A} \sim Dir\left(\alpha G_0(B_1) + \delta_{B_2} \cdots \right)$

Sethuraman

Let $\beta_k \sim$ Beta$(1, \alpha_0)$ $\qquad \theta_k \stackrel{iid}{\sim} G_0$

$\pi_1 = \beta_1$

$\pi_k = \beta_k(1-\beta_1)\cdots(1-\beta_{k-1}) = \beta_k(1-\pi_1-\cdots-\pi_{k-1})$.

then let $G = \sum \pi_k \delta_{\theta_k}$

**Thm** $G \sim DP(\alpha, G_0)$.

$$G = \beta_1 \delta_{\theta_1} + (1-\beta_1)\left[\beta_2 \delta_{\theta_2} + (1-\beta_2)\beta_3 \delta_{\theta_3} + \cdots \right]$$

$$= \beta_1 \delta_{\theta_1} + (1-\beta_1) \; G'$$

$\qquad\qquad\qquad\qquad\qquad G' \perp\!\!\!\perp \beta_1, \theta_1$.

So $\quad G \stackrel{dist}{=} \beta_1 \delta_{\theta_1} + (1-\beta_1) G$.

Let $(B_1, \ldots, B_k)$ be a partition

$$P = (G(B_1), \ldots, G(B_k))$$

then

$$P = \beta_1 \left(\delta_{\theta_1}(B_1), \ldots, \delta_{\theta_1}(B_k)\right) + (1-\beta_1) P \qquad (\sharp)$$

Can easily check that

$$Dir\left(\alpha G_0(B_1), \ldots, \alpha G_0(B_k)\right) \stackrel{dist}{=} \beta\left(\left(\delta_{\theta_1}(B_1), \ldots, \delta_{\theta_1}(B_k)\right) + \right.$$

$$\left. (1-\beta_1) Dir\left(\alpha G_0(B_1), \ldots, \alpha G_0(B_k)\right)\right.$$

if $\theta_1 \in B_i$ then $\quad \beta_1(1, 0, \ldots, 0) + (1-\beta_1) Dir(\alpha G_0(B_1) \ldots )$

$$\stackrel{dist}{=} Dir(\alpha G_0 + \delta_{e_i}).$$

$$\sum_{j=1}^{k} P(\theta_j \in B_j) \; Dir(\alpha G_0 + \delta_{e_j}) = Dir(\alpha G_0).$$

next Lemma $N, U, W$ RV's, $W \in [-1, 1]$

$U, V$ taking value in some vector space

and
$$\begin{cases} V \overset{dist}{=} U + WV \\ P(|W| = 1) \neq 1 \end{cases} \quad \text{then} \quad V \text{ is unique.}$$

if $V \overset{dist}{\neq} V'$

$$V_{n+1} := U_n + W_n V_n$$
$$V'_{n+1} = U_n + W_n V'_n$$

$$\Rightarrow |V_{n+1} - V'_{n+1}| = |W_n| |V_n - V'_n| \longrightarrow 0 \quad \text{wp 1}$$

This then proves automatically that $\textcircled{D}P$ distributed Random measures are discrete wp 1.

## Alternative proof of discreteness

$$G \sim DP(\alpha, G_0)$$
$$X | G \sim G$$

then $G | X \sim DP(\alpha G_0 + \delta_X)$.

$$G | X_1, \dots X_n \sim DP\left(\alpha G_0 + \sum_{i=1}^{n} \delta_{X_i}\right)$$

$$DP\left(\alpha + n, \frac{\alpha}{\alpha+n} G_0 + \sum_{i=1}^{n} \frac{1}{\alpha+n} \delta_{X_i}\right).$$

as $n \to \infty$:
$$\begin{cases} \alpha + n \quad \text{Concentration} \to \infty \\ \text{base measure} \overset{W.}{\rightrightarrows} \frac{1}{n} \sum_{i=1}^{n} \delta_{X_i} \end{cases}$$