

# **ML Lecture #6**

## Bayesian nonparametrics

Julyan Arbel

# Outline

## Bayesian nonparametrics

- Pitman-Yor process, Gibbs-type process
- Hierarchical Dirichlet process ]
- Indian buffet process
- Practical: Dirichlet process mixture models in Pyro



# Outline

## Bayesian nonparametrics

- Pitman-Yor process, Gibbs-type process
- Hierarchical Dirichlet process
- Indian buffet process
- Practical: Dirichlet mixture models in Pyro

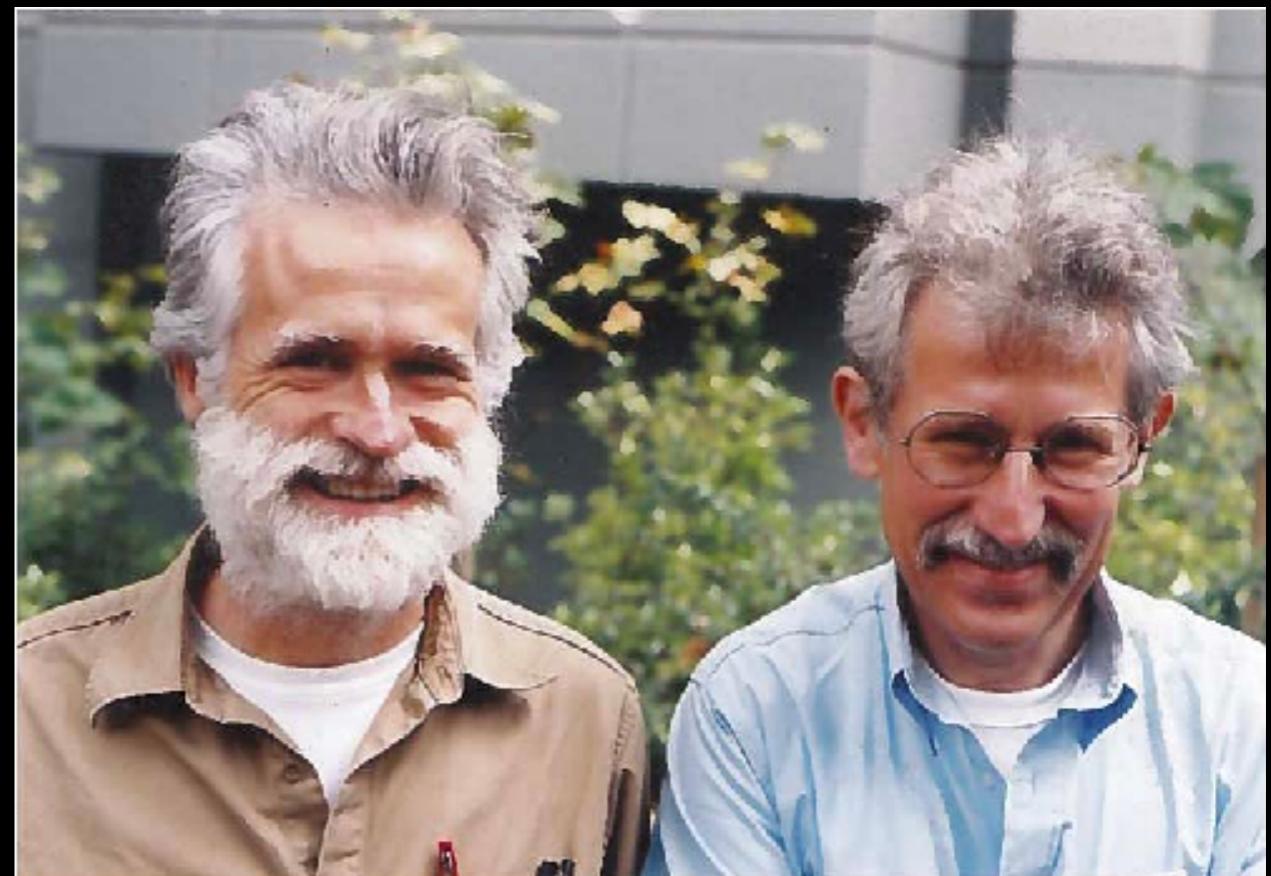
## Bayesian deep learning

- Maximum a posteriori = Regularized maximum likelihood
- Laplace approximation (MacKay, 1992, Neur. Comp.)
- Variational inference (Hinton and van Kamp, 1993, Barber & Bishop, 1998, NIPS)
- Monte Carlo dropout (Gal & Ghahramani, 2016, ICML)
- Practical: Bayesian neural networks in Pyro

# Bayesian nonparametrics

## Pitman-Yor process

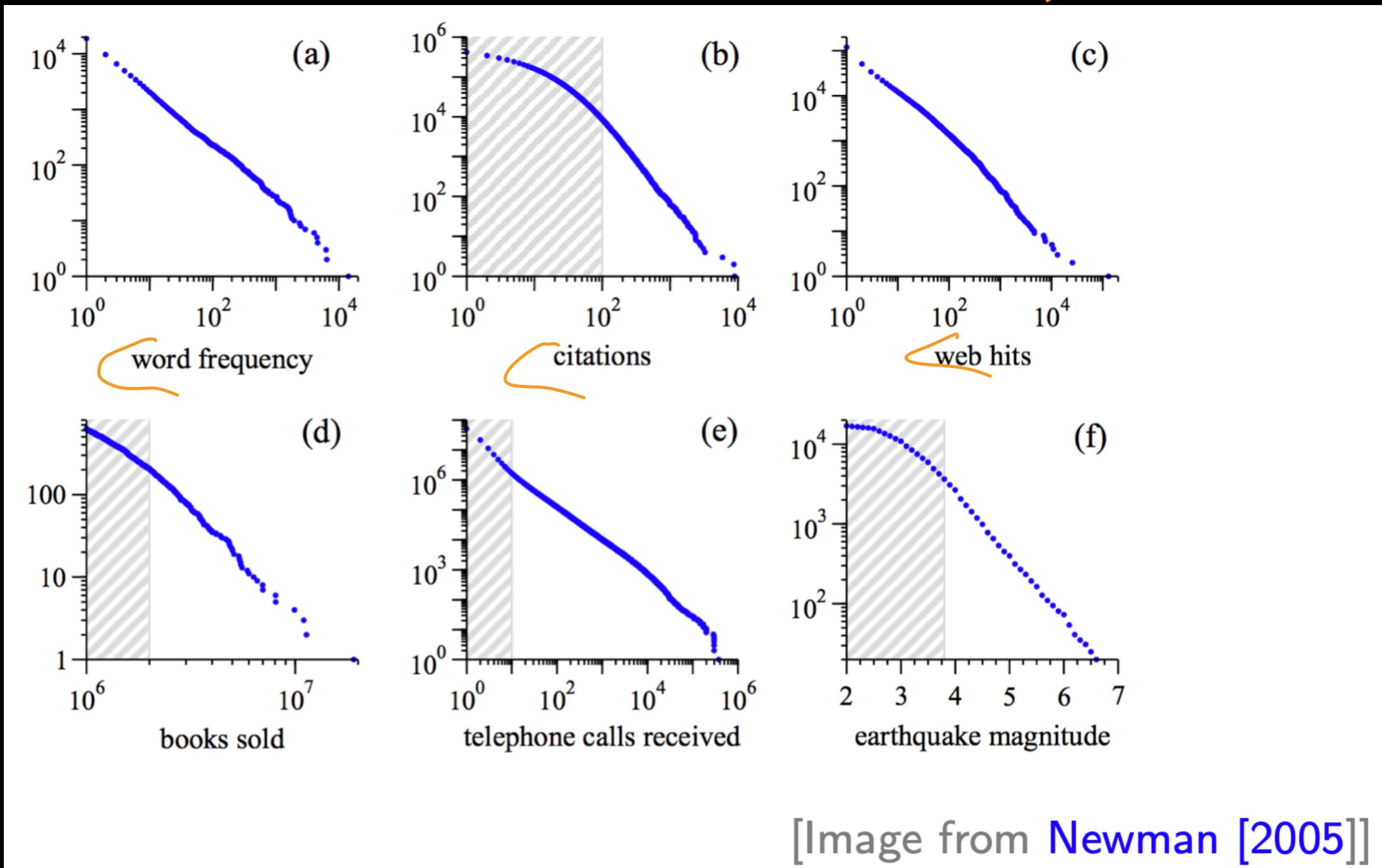
- Dirichlet process
  - Ferguson, 1973
- Pitman-Yor process
  - Perman, Pitman & Yor, 1992
  - Pitman & Yor, 1997



# Bayesian nonparametrics

## Pitman-Yor process

- Depart from log #clusters to power-law #clusters



## Dirichlet process

predictive distri:  $X_{n+1} | X_{1:n} \sim \frac{\alpha}{\alpha+n} P_0 + \frac{1}{\alpha+n} \sum_{i=1}^n \delta_{X_i}$

# clusters:  $K_n \approx \alpha \log n$

## Pitman-Yor process

$$X_{n+1} | X_{1:n} \sim \frac{\alpha + \sigma K_n}{\alpha + n} P_0 + \frac{1}{\alpha + n} \sum_{j=1}^{K_n} (n_j - \sigma) \delta_{x_j^*}$$

$x_j^*$ ,  $j = 1 \dots K_n$ , nb of distinct obs.

$$\sigma \in [0, 1]$$

$n_j$ : # of obs. = to  $x_j^*$

discount parameter

# clusters:  $K_n \approx S_{\alpha, \sigma} n^\sigma$

$S_{\alpha, \sigma}$ : diversity of PY

$\rightarrow$  Ex 2:  $E K_n$ ,  $\text{Var } K_n$ .

# Bayesian nonparametrics

## Pitman-Yor process: diversity

### Power Law and $\sigma$ -diversity

For  $\sigma > 0$  we have the almost sure convergence

$$n^{-\sigma} K_n \rightarrow S_{\sigma, \alpha},$$

where  $S_{\sigma, \alpha}$  is called  $\sigma$ -diversity of the PY,  
whose density is a polynomially tilted  
Mittag–Leffler density (ML):

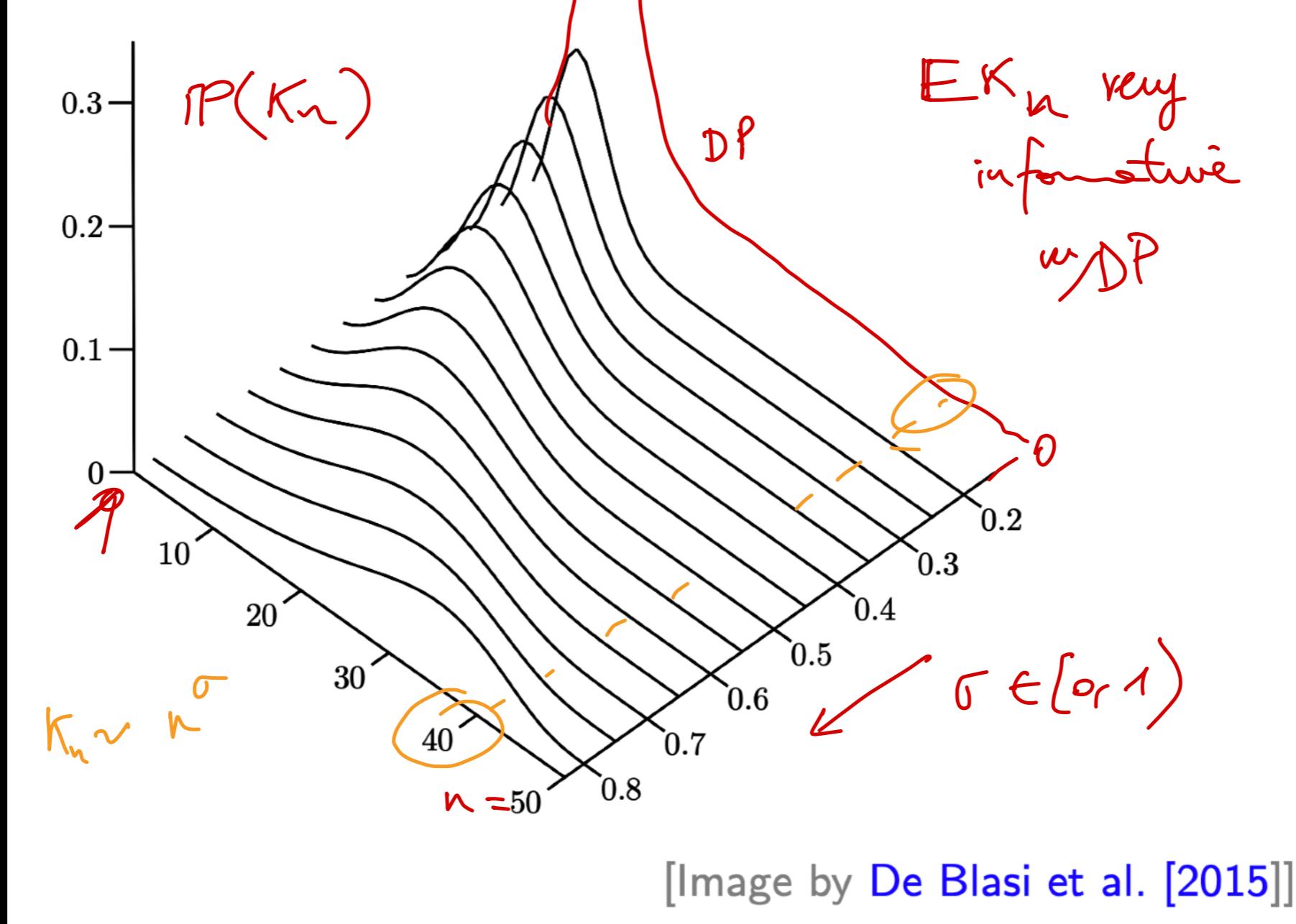
$$g_{\sigma, \alpha}(x) \propto x^{\alpha/\sigma} g_\alpha(x),$$

and  $g_\alpha$  is ML density.



# Bayesian nonparametrics

## Pitman-Yor process: number of clusters



$$\underline{\text{Notation}} \quad P \sim PY(P_0, \alpha, \sigma)$$

CRP: Chinese Restaurant process = clustering distribution  
aka EPPF: Exchangeable partition proba. function

$m_1, \dots, m_k$       k for  $K_n$

$$CRP : p(m_1, \dots, m_k)$$

Cust. seated at new table w.p.:

$$\frac{\alpha + \sigma K_n}{\alpha + m} \quad \begin{matrix} \text{discovery} \\ \text{probability} \end{matrix}$$

\_\_\_\_\_ occupied table w.p.:

$$\frac{(n_j - \sigma)}{\alpha + m} \quad \dots$$

$$\left\{ p(m_1, \dots, m_k) = \sigma^K \frac{\left(1 + \frac{\alpha}{\sigma}\right)_{k-1}}{\left(1 + \alpha\right)_{m-1}} \prod_{j=1}^k (1 - \sigma)^{m_j - 1} \right.$$

$$\sigma = 0 \rightarrow DP \ CRP$$

Ewens Sampling Formula

$$p(m_1, \dots, m_m)$$

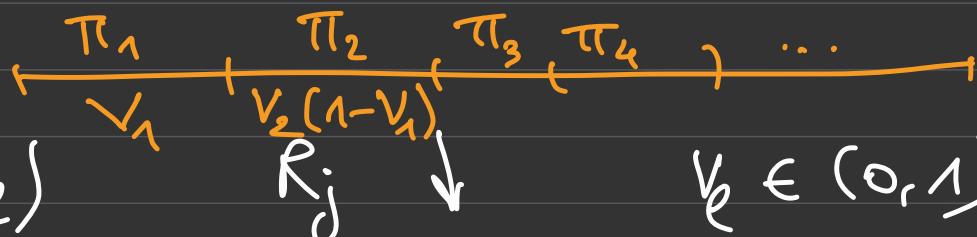
Stick-breaking representation for  $P_Y$ ,  $P_0$  is non atomic

$$\theta_j \stackrel{iid}{\sim} P_0, \quad V_j \stackrel{ind}{\sim} \text{Beta}(1-\sigma, \alpha + j\sigma) \quad (\text{DP w.r.t } \sigma = 0)$$

$$\pi_1 = v_1, \quad j > 1, \quad \left[ \pi_j = v_j \prod_{l=1}^{j-1} (1-v_l) \right] \quad \mathbb{E} V_j = \frac{1-\sigma}{1+\alpha+(j-1)\sigma}$$

Proposition:  $\rho := \sum_{j=1}^{\infty} \pi_j \delta_{\theta_j} \sim P_Y(P_0, \alpha, \sigma)$

$\rightarrow$  discrete random probability measure.  $\mathbb{E} R_j = (1-\dots)^j$



left-over  $R_j = \prod_{l=1}^j (1-v_l) \quad v_l \in (0, 1)$

$$\begin{aligned} \mathbb{E} R_i &= \mathbb{E} \prod_{j=1}^i (1-v_j) = \prod \mathbb{E} (1-v_j) = \prod \frac{\alpha + \sigma j}{1 + \alpha + (j-1)\sigma} \\ &= \prod \left( 1 - \frac{1-\sigma}{1 + \alpha + (j-1)\sigma} \right) \quad \log(1+x) \leq x \end{aligned}$$

$$\begin{aligned} \log \mathbb{E} R_i &= \sum \log \left( \frac{1-\sigma}{1 + \alpha + (j-1)\sigma} \right) \Rightarrow \mathbb{E} R_i \rightarrow 0 \\ &\leq - \sum \frac{1-\sigma}{1 + \alpha + (j-1)\sigma} \rightarrow -\infty \quad R_i \rightarrow 0 \text{ a.s.} \end{aligned}$$

$A, \beta$  measurable subsets of  $\mathbb{Q}$ . {Remember for DP : used  $P(A) \sim \text{Beta}(\dots, \dots)$ }

$$\begin{aligned} \mathbb{E}[P(A)] &= \mathbb{E}\left[\sum_{i=1}^{\infty} \pi_i \delta_{\theta_i}(A)\right] \\ &= \sum \mathbb{E}[\pi_i] \mathbb{E}[\delta_{\theta_i}(A)] = \sum \underbrace{\mathbb{E}[\pi_i]}_{\sim} P_o(A) \quad \text{because } \theta_i \sim P_o \\ &= P_o(A) \sum \mathbb{E}[\pi_i] = P_o(A) \underbrace{\mathbb{E}\left[\sum \pi_i\right]}_1 = P_o(A). \end{aligned}$$

$P_o$  serves as a base measure.

$$\text{Cov}(P(A), P(B)) = \underbrace{\mathbb{E}[P(A)P(B)]}_{(*)} - \underbrace{\mathbb{E}[P(A)]}_{P_o(A)} \underbrace{\mathbb{E}[P(B)]}_{P_o(B)}$$

$$X_1, X_2 \mid P \stackrel{\text{iid}}{\sim} P \quad (*) = P(X_1 \in A, X_2 \in B)$$

$$= \underbrace{P(X_1 \in A)}_{\mathbb{E}[P(A)] = P_o(A)} \underbrace{P(X_2 \in B \mid X_1 \in A)}_{(***)} \quad \begin{matrix} \alpha = 1 \\ n = 1 \end{matrix}$$

$$(****) = \frac{\alpha + \sigma}{\alpha + 1} P_o(B) + \frac{1 - \sigma}{\alpha + 1} P_o(B \mid A) \quad \begin{matrix} \text{predict. } X_2 \mid X_1 \sim \frac{\alpha + \sigma}{\alpha + 1} P_o(B) + \frac{1 - \sigma}{\alpha + 1} X_1 \\ \frac{P_o(A \cap B)}{P_o(A)} \end{matrix}$$

$$\text{Cov}(P(A), P(B)) = -\frac{1 - \sigma}{\alpha + 1} [P_o(A)P_o(B) - P_o(A \cap B)]$$

$$\text{Var}(P(A)) = \frac{1 - \sigma}{\alpha + 1} P_o(A)(1 - P_o(A))$$

## Posterior characterization for PY

Assume  $\{X_1, \dots, X_n | P \text{ iid } P\}$

$$\{P \sim PY(\rho_0, \alpha, \sigma) \quad PY\}$$

$$P | X_1, \dots, X_n \stackrel{\text{"d" }}{=} \underbrace{(1-q_n)}_{\text{dist. obs.}} \underbrace{P_n}_{\text{PY}} + q_n \underbrace{\sum_{j=1}^{k_n} p_j^* \delta_{x_j^*}}_{\text{dist. obs.}}$$

$$\left\{ \begin{array}{l} \square q_n \sim \text{Beta}(n - \sigma k_n, \alpha + \sigma k_n) \in (0, 1) \\ \square P_n \sim PY(\rho_0, \sigma, \underline{\alpha + \sigma k_n}) \\ \square (p_1^*, p_2^*, \dots, p_{k_n}^*) \sim Dir(m_1 - \sigma, \dots, m_{k_n} - \sigma) \end{array} \right.$$

# Bayesian nonparametrics

## Gibbs-type process

- Pitman, 2003

Define via EPPF (CRP) :

$$P(m_1, \dots, m_k) = \underbrace{V_{n,k}}_{n \geq 1} \prod_{j=1}^k (1-\sigma)^{m_j}$$

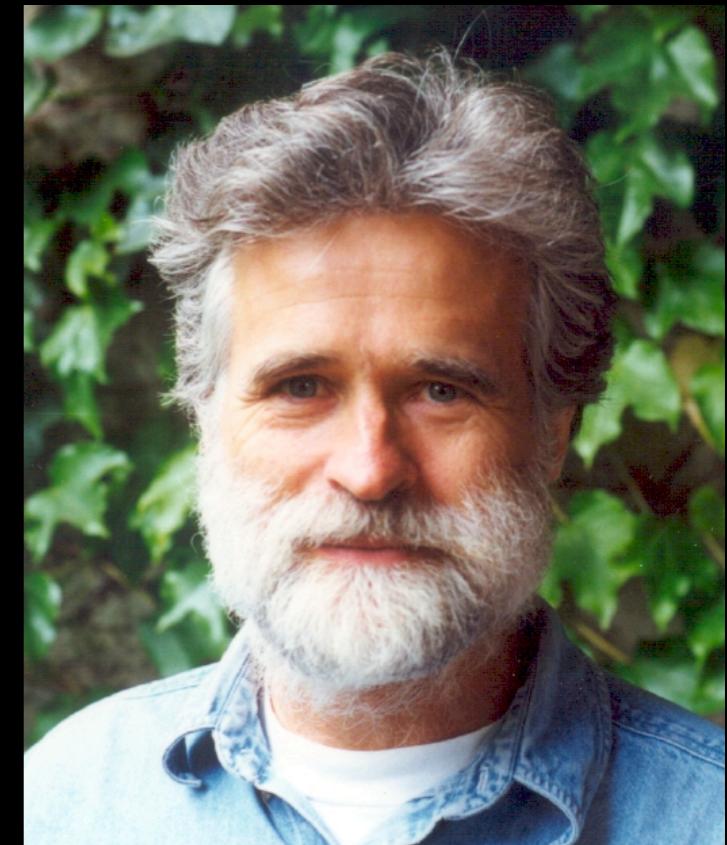
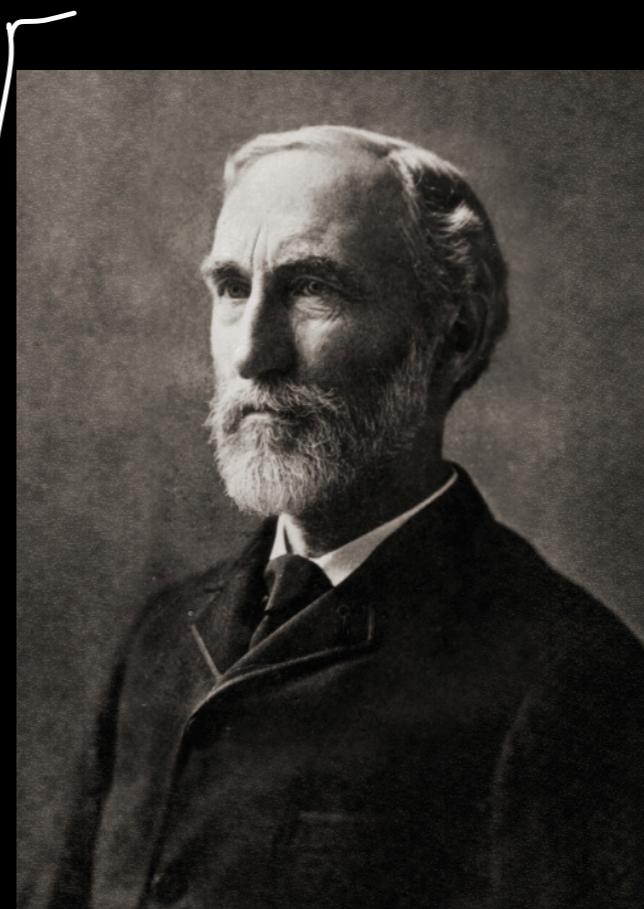
$$\begin{cases} n \geq 1 \\ 1 \leq k \leq n \\ \sigma \in (0, 1) \end{cases}$$

$$\begin{cases} V_{1,1} = 1 \\ V_{m,k} = (m - \sigma k) V_{m+1,k} + V_{m+1,k+1} \end{cases} \iff$$

Recover PY :  $V_{m,k} = \sigma^k \frac{(1 + \frac{\alpha}{\sigma})_{k-1}}{(1 + \alpha)_{n-1}}$

DP :  $\sigma \ll 0$

Predictive :  $X_{n+1} | X_{1:n} \sim \left[ \frac{V_{m+1,k+1}}{V_{m,k}} \right] P_0 + \frac{V_{m+1,k}}{V_{m,k}} \sum_{j=1}^{k_n} (n_j - \sigma) \delta_{x_j^*}$



Jim Pitman

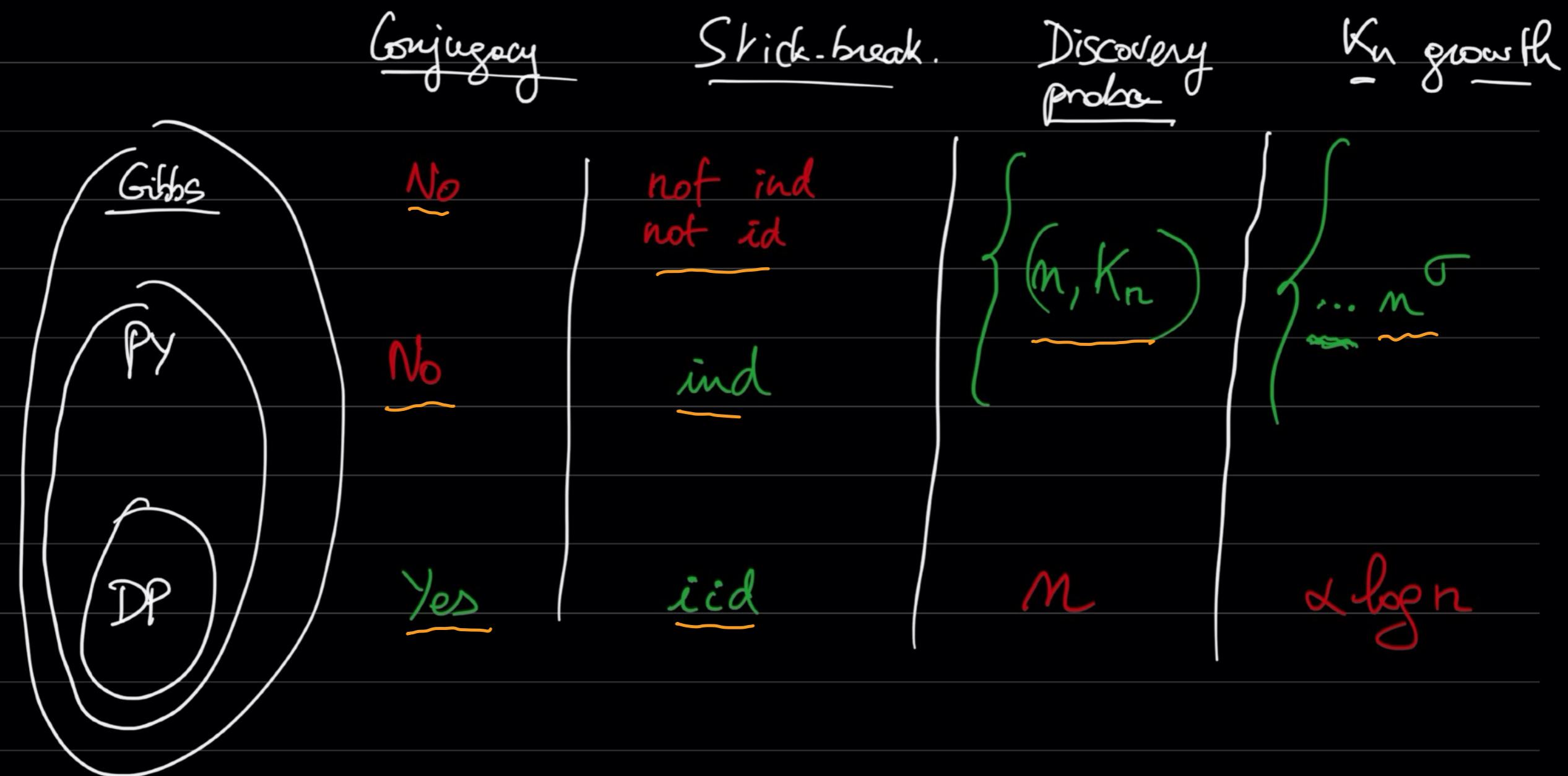
↑  
discovery proba.  
depends  $(n, K_n)$   
↑  
old quantity  
 $(m, k_n)$

Stick-breaking :  $v_j | v_{j-1}, \dots, v_1 \sim g(v_j | v_{j-1}, \dots, v_1)$

One class of Gibbs-type : normalized generalized gamma  
 $NG \neq PY$

# Bayesian nonparametrics

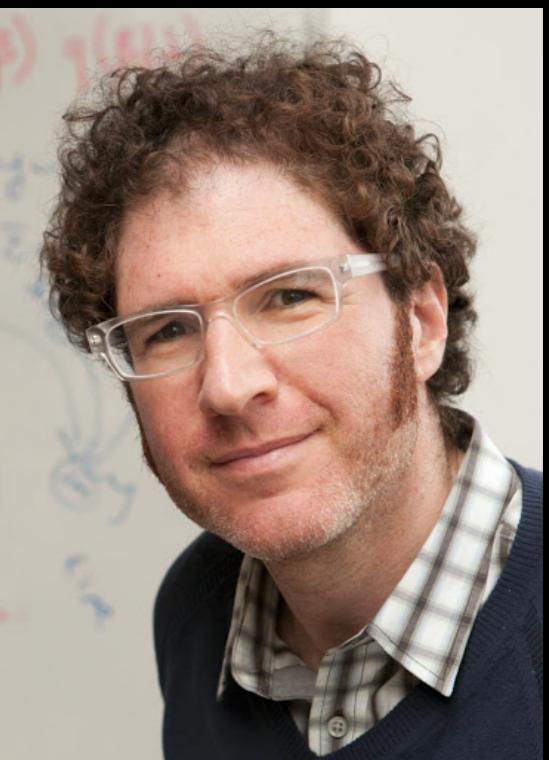
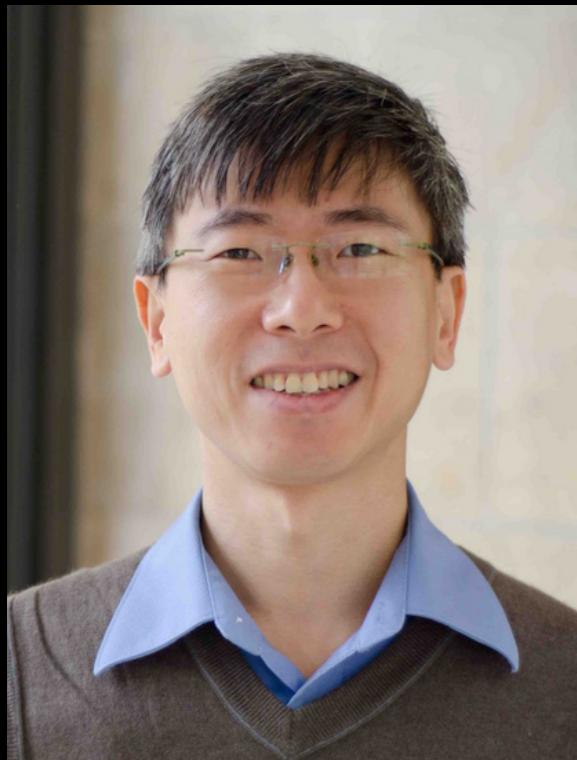
Dirichlet  $\subset$  Pitman-Yor  $\subset$  Gibbs-type



# Bayesian nonparametrics

## Hierarchical Dirichlet process

- Extension of Latent Dirichlet Allocation (Blei, Ng, and Jordan, 2003, JMLR)
- to infinite dimensional parameter space (Teh, Jordan, Beal, and Blei, JASA, 2006)



# Bayesian nonparametrics

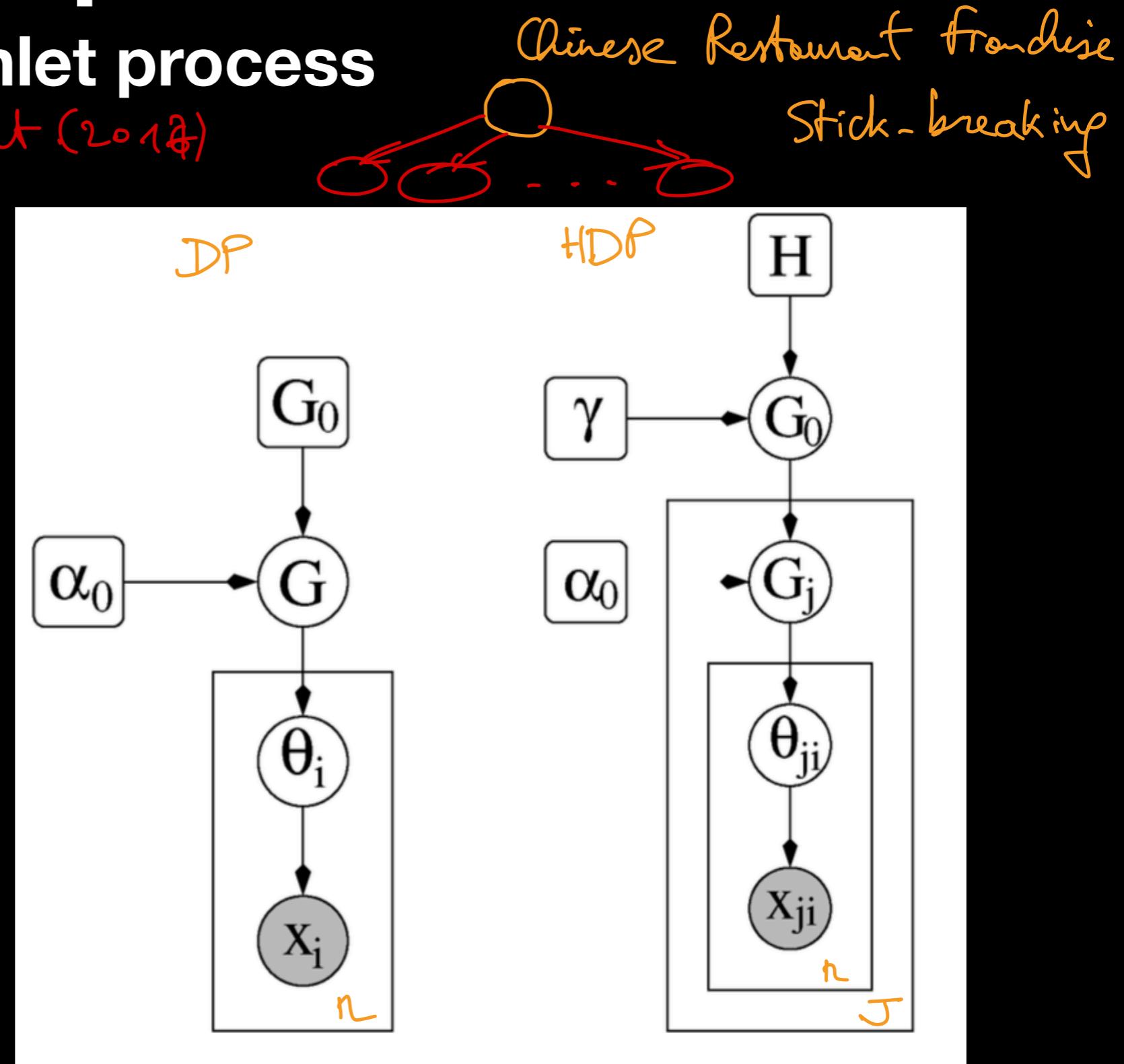
## Hierarchical Dirichlet process

Ghosal & van der Vaart (2012)

$$\text{DP: } \begin{cases} G | \alpha_0, \theta_0 \sim \text{DP}(\alpha_0, G_0) \\ \theta_i | G \stackrel{\text{iid}}{\sim} G \\ x_i | \theta_i \stackrel{\text{ll}}{\sim} F(x_i | \theta_i) \end{cases}$$

$$\text{HDP: } \begin{cases} G_0 | \gamma, H \sim \boxed{\text{DP}}(\gamma, H) \\ \cancel{G_j | \alpha_0, G_0 \stackrel{\text{iid}}{\sim} \text{DP}(\alpha_0, G_0)} \\ \theta_{ji} | G_j \sim G_j, j \dots \\ x_{ji} | \theta_{ji} \stackrel{\text{ll}}{\sim} F(x_{ji} | \theta_{ji}) \end{cases}$$

Something new here:  
base measure  $G_0$ : atomic.



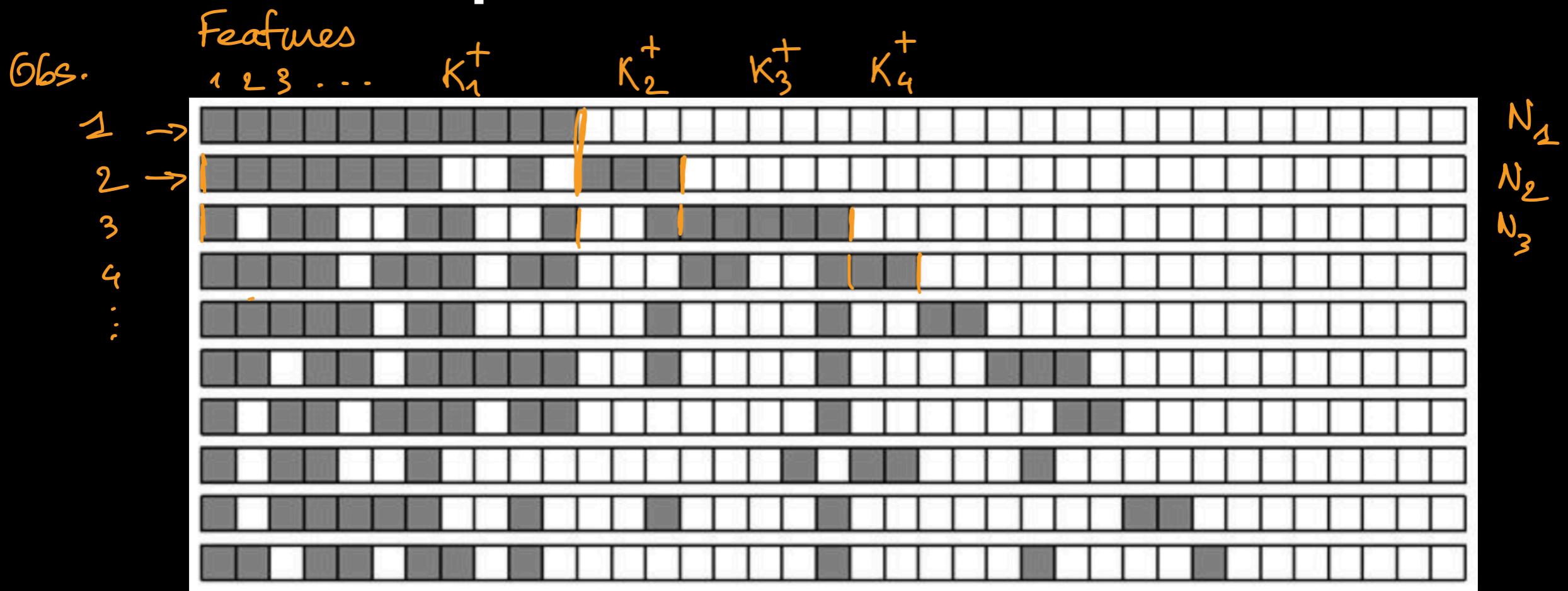
# Bayesian nonparametrics

## Indian buffet process

- Extension of Dirichlet process (Ferguson, 2003, AOS) to multiple classes/allocation models
  - Ghahramani & Griffiths, 2005.

# Bayesian nonparametrics

## Indian buffet process



Poisson lemmas :

if additivity  $N_1 \sim \text{Pois}(\lambda_1)$   
 $N_2 \sim \text{Pois}(\lambda_2)$   $\Downarrow \Rightarrow N_1 + N_2 \sim \text{Pois}(\lambda_1 + \lambda_2)$

if thinning :  $N \sim \text{Pois}(\lambda)$   
 $x_1, \dots, x_N \stackrel{iid}{\sim} \text{Ber}(\rho) \Rightarrow \sum_{i=1}^N x_i \sim \text{Pois}(\lambda\rho)$

Generative model for  $\text{IBP}(\gamma)$ ,  $\gamma > 0$

Customer 1: choose  $N_1 \sim \text{Pois}(\gamma)$

Customer 2: - choose some dishes of C1 w.p  $\frac{1}{2}$   
- ——— new dishes  $K_2^+ \sim \text{Pois}\left(\frac{\gamma}{2}\right)$

$$\Rightarrow N_2 \sim \text{Pois}\left(\frac{\gamma}{2}\right) + \text{Pois}\left(\frac{\gamma}{2}\right) = \text{Pois}(\gamma)$$

:

Customer  $i$ : denote  $n_1, \dots, n_K$  all sampled dishes

choose previous dishes w.p  $\frac{1}{i}$ , ie dish  $j$  w.p  $\frac{n_j}{i}$

new dishes  $\text{Pois}\left(\frac{\gamma}{i}\right)$

$$N_i \sim \underbrace{\text{Pois}\left(\frac{i-1}{i} \gamma\right)}_{\text{old}} + \text{Pois}\left(\frac{\gamma}{i}\right) = \text{Pois}(\gamma). \xrightarrow{\text{Pois}(\gamma \log n)}$$

Total nb of dishes:  $K_n = K_1^+ + \dots + K_n^+ = \underbrace{\text{Pois}(\gamma) + \text{Pois}\left(\frac{\gamma}{2}\right) + \dots + \text{Pois}\left(\frac{\gamma}{n}\right)}_{\text{Pois}(\gamma(1 + \frac{1}{2} + \dots + \frac{1}{n}))} \xrightarrow{n \rightarrow \infty} \infty$

# Bayesian nonparametrics

## Practical: Pyro, pyro.ai

Pyro is a universal probabilistic programming language (PPL) written in Python and supported by PyTorch on the backend. Pyro enables flexible and expressive **deep** probabilistic modeling, unifying the best of modern deep learning and Bayesian modeling. It was designed with these key principles:

**Universal**: Pyro can represent any computable probability distribution.

**Scalable**: Pyro scales to large data sets with little overhead.

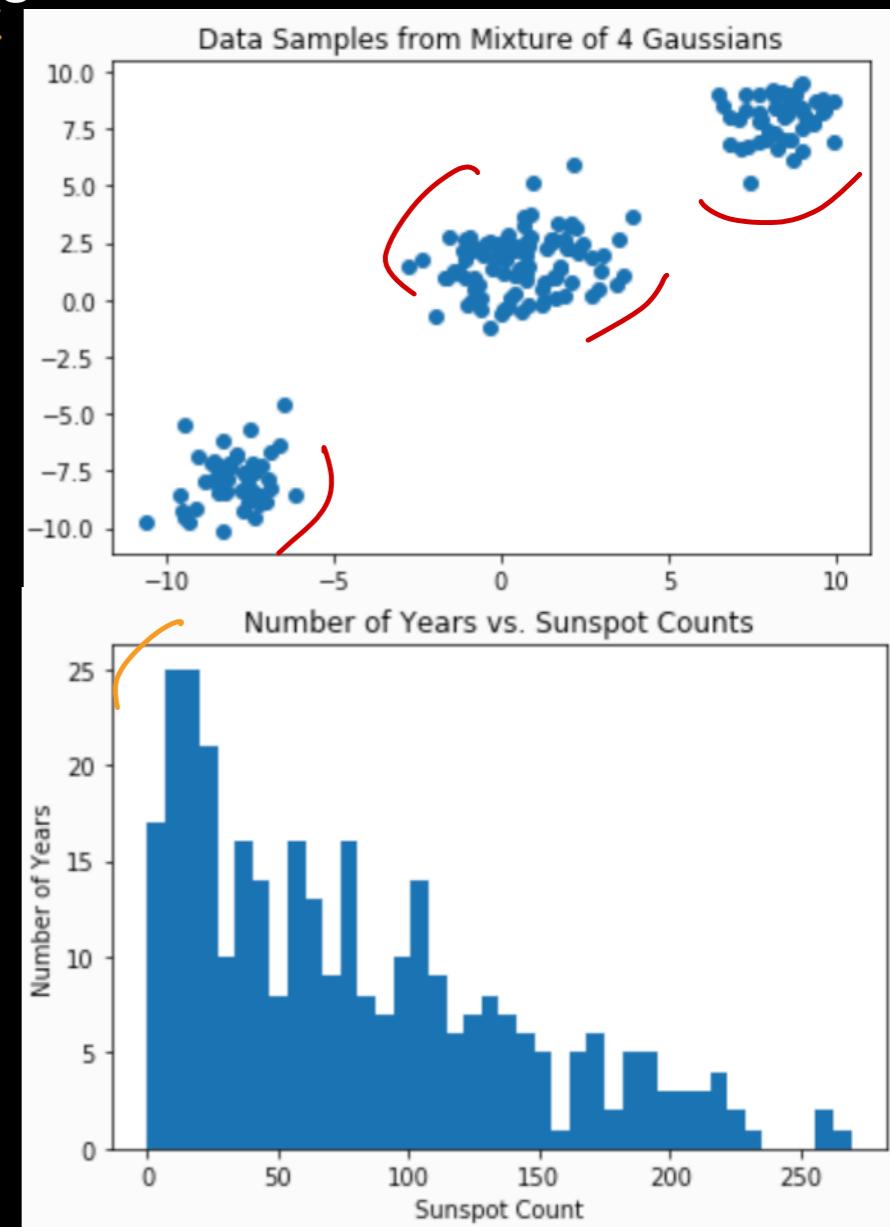
**Minimal**: Pyro is implemented with a small core of powerful, composable abstractions.

**Flexible**: Pyro aims for automation when you want it, control when you need it.

# Bayesian nonparametrics

## Practical: Dirichlet process mixture models in Pyro

- [https://pyro.ai/examples/dirichlet\\_process\\_mixture.html](https://pyro.ai/examples/dirichlet_process_mixture.html)
- Implements Variational inference based on stick-breaking
  - On Synthetic Mixture of Gaussians ✕
  - On Long Term Solar Observations ✕
- Objectives:
  - Install Pyro & PyTorch (or work on Colab) ||
  - Run the code
  - Write your own code for Pitman-Yor mixture models



# Bayesian deep learning

## Bayesian neural networks

- Maximum a posteriori = Regularized maximum likelihood
- Laplace approximation (MacKay, 1992, Neur. Comp.)
- Variational inference (Hinton and van Kamp, 1993, Barber & Bishop, 1998, NIPS)
- Monte Carlo dropout (Gal & Ghahramani, 2016, ICML)
- Implementation: Pyro & PyTorch

# Bayesian deep learning

## Pyro example on Bayesian neural networks

- Demonstrates how to use NUTS to do inference on a simple (small) Bayesian neural network with two hidden layers.

