

BML Lecture 4 / Bayesian nonparametrics

Julian ARBEL

What comes to your mind when you hear about
Bayesian nonparametrics ?

Infinite-dim spaces
Dirichlet process
Random process
Clustering

Nonparametric means : ④ infinite dimensional
or dimension that grows with n

Bayes : $P(\theta | x) \propto P(\theta) \cdot P(x | \theta)$

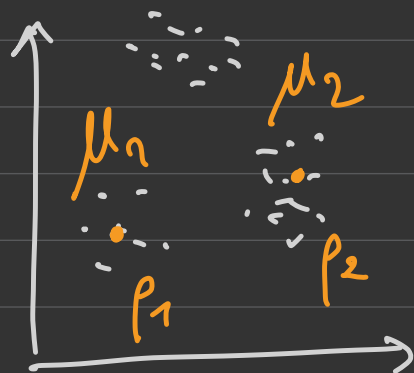
De Finetti theorem

infinite exchangeable data: $p(X_1, \dots, X_k) = p(X_{\sigma(1)}, \dots, X_{\sigma(k)})$
 $\forall k \geq 1, \forall \sigma \in S^k$

$$\Leftrightarrow p(X_1, \dots, X_k) = \int \prod_{i=1}^k p(X_i | \theta) \underline{P}(d\theta)$$

conditional iid. ④

Mixture Models, generative model



$n = 1 \dots N$

$K = 2$ clusters

$$\begin{cases} z_n \stackrel{\text{iid}}{\sim} \text{Categorical}(p_1, p_2, \dots, p_K) \\ x_n | z_n \stackrel{\text{iid}}{\sim} N(\mu_{z_n}, \Sigma) \end{cases} \quad \begin{matrix} p_1 + p_2 = 1 \\ p_1 > 0, p_2 > 0 \end{matrix}$$

Parameters $(\mu_1, \mu_2, \dots, \mu_K)$ (p_1, p_2, \dots, p_K)

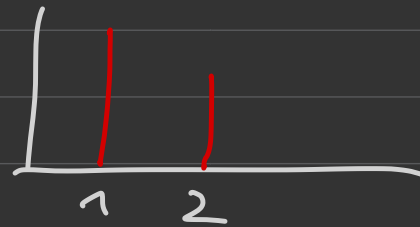
$$\mu_k \sim N(\mu_0, \Sigma_0)$$

$$p_1 \sim \text{Beta}(a_1, a_2) \\ p_2 = 1 - p_1$$

Now, $(p_1, \dots, p_K) \sim \text{Dirichlet}(a_1, \dots, a_K)$

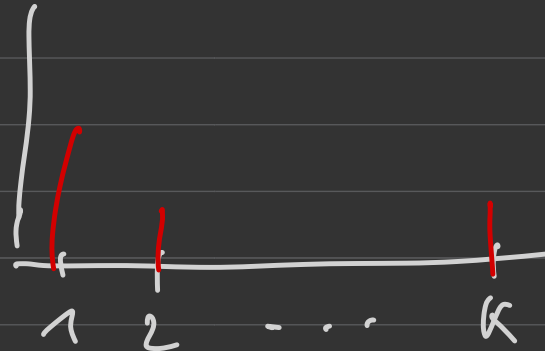
N obs. K clusters. What to do if $K > N$?

Beta



(p_1, p_2)

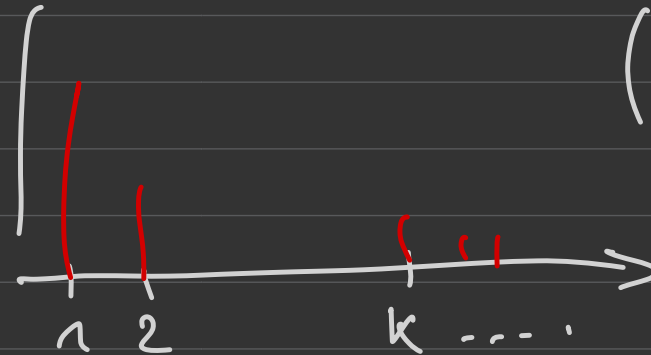
Dirichlet



(p_1, \dots, p_K)

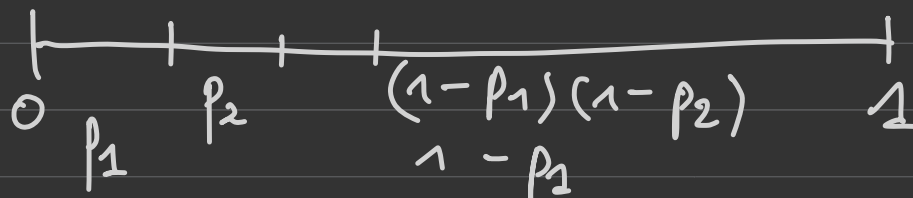
Stick-breaking

$K = \infty$



$(p_1, p_2, \dots) \in \text{Simplex}(\infty)$

$$\sum p_i = 1, p_i \geq 0$$



$(p_1, p_2, \dots) \sim \text{GEM}.$

Chapt 1 Introduction

Chapter 2 Dirichlet process

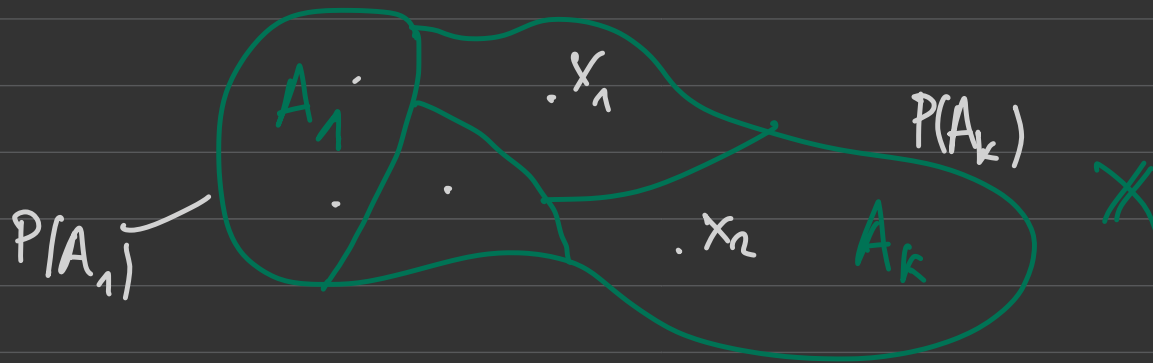
1. Definition Ferguson, 1973

Space X , P is a Dirichlet process (DP) on X

if $\exists \alpha > 0$ concentration
 P_0 (fixed) proba measure: **base measure**

$\forall k, \forall$ partition (A_1, \dots, A_k) of X

$$(P(A_1), \dots, P(A_k)) \sim \text{Dir}(\alpha P_0(A_1), \dots, \alpha P_0(A_k)).$$



$$P \sim \text{DP}(\alpha, P_0)$$

L_s is a proba dist. on X

$$x_1, \dots, x_n | P \sim P.$$

$\in X$

Moments

Let $A \subset X$. If $P \sim DP(\alpha, P_0)$, then

$$E[P(A)] = \frac{\alpha P_0(A)}{\alpha} = P_0(A)$$

$$(P(A), P(A^c)) \sim \text{Dir}(\underbrace{\alpha P_0(A)}_{a}, \underbrace{\alpha P_0(A^c)}_{b})$$

$$P(A) \sim \text{Beta}(a, b)$$

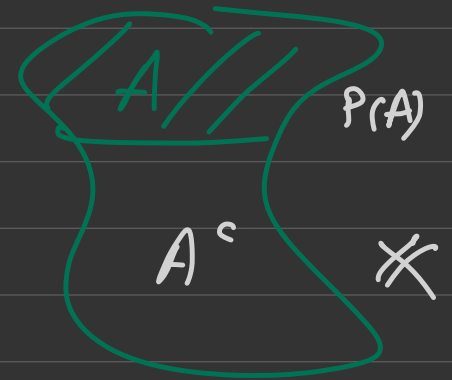
$$E \text{Beta}(a, b) = \frac{a}{a+b}$$

$$V[P(A)] = \frac{P_0(A)(1-P_0(A))}{1+\alpha} \text{ — contraction.}$$

$$V \text{ — } = \frac{ab}{(a+b)^2(a+b+1)}$$

$$\text{Cov}(P(A), P(B)) = \overset{\text{EXO}}{\dots} = \frac{P_0(A \cap B) - P_0(A)P_0(B)}{1+\alpha}$$

$$\text{ex. } A \cap B = \emptyset \Rightarrow \text{Cov} \leq 0$$



$$E[P(A)] = P_0(A) \quad \text{can be written} \quad \int P(A) dDP(P) = P_0(A)$$

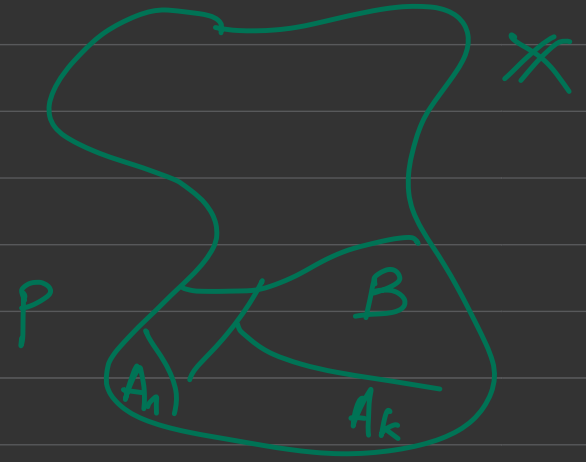
$$\begin{cases} P \sim DP(\alpha, P_0) \\ X|P \sim P \end{cases} \Rightarrow X \sim P_0$$

Self-similarity of the DP

$$P_B(A) = P(A \cap B)$$

$$P_B(A) = P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Theorem: if $P \sim DP(\alpha, P_0)$ on X , then $P_B \sim DP(\alpha P_{0|B})$.



$$(P(A_1), \dots, P(A_k), P(B^c)) \sim \text{Dir}_{k+1}$$

$$\left(\frac{P(A_1)}{P(B)}, \dots, \frac{P(A_k)}{P(B)} \right) \sim \text{Dir}_k \left(\underbrace{\alpha P_0(A_1)}_{\alpha P_{0|B}(A_1)}, \dots, \underbrace{\alpha P_0(A_k)}_{\alpha P_{0|B}(A_k)} \right)$$

Posterior conjugacy : X_1, \dots, X_n sample from DP :

$$\begin{cases} P \sim \text{DP}(\alpha, P_0) \\ X_1, \dots, X_n | P \sim P \end{cases} \quad (*) \quad \text{DP}(\alpha, P_0) \leftrightarrow \text{DP}(\alpha P_0)$$

$G = \alpha P_0$
 $\alpha = G(X)$ $P_0 = \frac{G}{\alpha}$

Th (Ferguson) : the posterior in (*) is

$$P(X_1, \dots, X_n) \sim \text{DP}\left(\alpha P_0 + \sum_{i=1}^n \delta_{X_i}\right)$$

Predictive distribution :

$$P(X_{n+1} | X_1, \dots, X_n) = \frac{\alpha}{\alpha + n} P_0 + \frac{1}{\alpha + n} \sum_{i=1}^n \delta_{X_i} = P_n$$

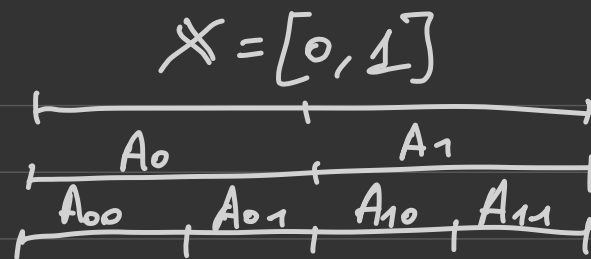
Updated parameter $G = \alpha P_0 + \sum_{i=1}^n \delta_{X_i}$ Dirac masses at X_i

— concentration $\alpha_n = G(X) = \alpha + n$

— base measure $P_n = \frac{G}{\alpha + n}$

proof: tail-free property

$$P \sim DP(\alpha, P_0) \text{ on } \mathbb{X}$$



then $\{P(A_0), P(A_1)\} \perp\!\!\!\perp \{P(A_{00}), P(A_{01}), P(A_{10}), P(A_{11})\}$.

Consider partition (A_1, \dots, A_k) of \mathbb{X} .

$$N_j = \#(i : X_i \in A_j)$$

tail-free property $\Rightarrow \underbrace{(P(A_1), \dots, P(A_k))}_{\text{prior}} | X_1 \dots X_n \stackrel{d}{=} \underbrace{(P(A_1), \dots, P(A_k))}_{\text{model}} | N_1, \dots, N_k$

prior $(P(A_1), \dots, P(A_k)) \sim \text{Dir}_k(\alpha P_0(A_1), \dots, \alpha P_0(A_k))$

model $(N_1, \dots, N_k) \sim \text{Multinomial}(n, (P(A_1), \dots, P(A_k)))$

posterior $(P(A_1), \dots, P(A_k)) | (N_1, \dots, N_k) \sim \text{Dir}_k(\alpha P_0(A_1) + N_1, \dots, \alpha P_0(A_k) + N_k)$

$$\Rightarrow P | N_1, \dots, N_k = P | X_1, \dots, X_n \sim DP\left(\alpha P_0 + \sum_{i=1}^n \delta_{X_i}\right)$$

Predictive

Fundamentals of BNP inference

$$P(X_{n+1} | X_1, \dots, X_n) = \underbrace{\frac{\alpha}{\alpha+n} P_0}_{\text{prior}} + \underbrace{\frac{n}{\alpha+n} \frac{1}{n} \sum_{i=1}^n \delta_{X_i}}_{\text{data}}$$

Ghosal
van der Vaart

Chapter 4
on DP.

Marginalization of DP : Polya-Urn
Blackwell-McQueen.

$$\begin{cases} P \sim DP \\ X_1 | P \sim P \end{cases} \Rightarrow X_1 \sim \underline{P_0}$$

$$\begin{cases} P \sim DP \\ X_1, X_2 | P \sim P \end{cases} \Rightarrow X_1 \sim P_0, \quad \underline{X_2 | X_1} \sim \frac{\alpha}{\alpha+1} P_0 + \frac{1}{\alpha+1} \delta_{X_1}$$

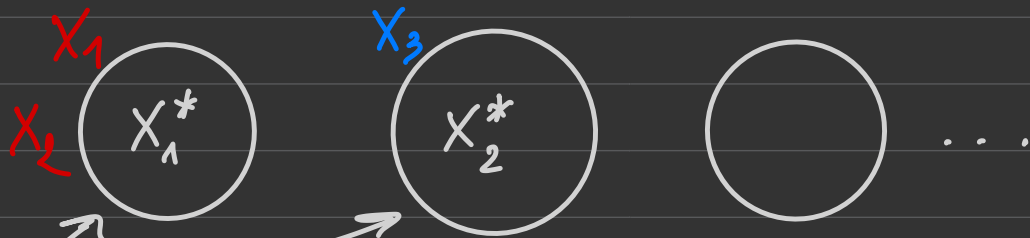
⋮

$$X_{n+1} | X_1, \dots, X_n \sim \frac{\alpha}{\alpha+n} P_0 + \frac{1}{\alpha+n} \sum_{i=1}^n \delta_{X_i}$$

α uncolored balls



Chinese Restaurant process : clustering distribution



rule : each new customer gets seated to tables
with probability to $\frac{\# \text{ of seated customers}}{\alpha + \# \text{ of seated customers}}$
 α for new table.

X_j^* unique observations (tables)

n_j : # customers on table X_j^*

$K = \# \text{ (populated) tables}$

$$P(X_{n+1} / X_1 \dots X_n) = \frac{\alpha}{\alpha + n} P_0 + \frac{n}{\alpha + n} \frac{1}{n} \sum_{j=1}^K n_j \delta_{X_j^*}$$

$$1, \frac{1}{\alpha + 1}$$

Theorem CRP $\longleftrightarrow x_1, \dots, x_n \mid P \sim DP$

$$\left\{ \begin{aligned} P(m_1, \dots, m_k) &= \frac{\alpha^k}{(\alpha)_n} \prod_{j=1}^k (m_j - 1)! = (*) \\ &= \alpha(\alpha+1) \dots (\alpha+n-1) \text{ ascending factorial} \end{aligned} \right.$$



$$(*) = \frac{x_1}{\alpha} \frac{x_2}{\alpha+1} \frac{x_3}{\alpha+2} \dots \frac{x_{n_1}}{\alpha+n_1-1} \frac{x_{n_1+1}}{\alpha+m_1} \frac{1}{\alpha+m_1+1} \dots \frac{x_{n_1+n_2}}{m_2-1} \dots$$

$$(*) = \frac{\alpha^k}{(\alpha)_n} \prod_{j=1}^k (m_j - 1)! .$$

Ewens sampling formula $p(m_1, \dots, m_\ell)$

m_ℓ : # tables with ℓ customers, $i = 1, \dots, m$

$$\sum_{\ell=1}^m m_\ell = k$$

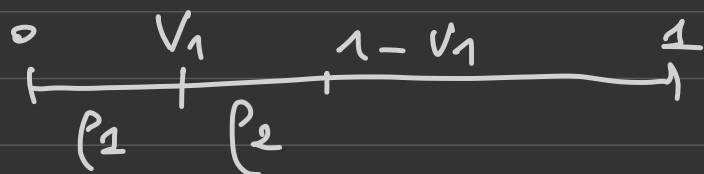
$$\sum_{\ell=1}^m \ell m_\ell = n$$

$$p(m_1, \dots, m_e) = \frac{n!}{\alpha(n)} \alpha^k \frac{1}{\prod_{l=1}^m l^{m_l} m_l!}$$

Stick-breaking for DP:

Let $V_i \stackrel{iid}{\sim} \text{Beta}(1, \alpha)$, $\theta_i \stackrel{iid}{\sim} P_0$.

Let $p_1 = V_1$, $p_i = V_i \prod_{l=1}^{i-1} (1 - V_l)$



if $P = \sum_{i=1}^{\infty} p_i \delta_{\theta_i}$

then $P \sim DP(\alpha, P_0)$.