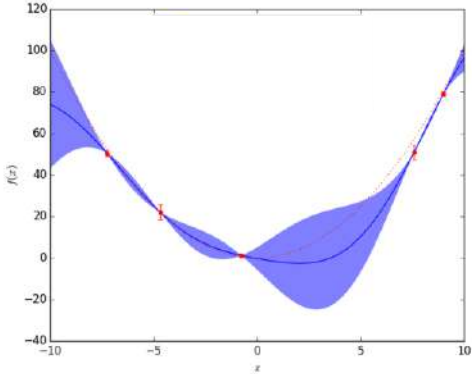




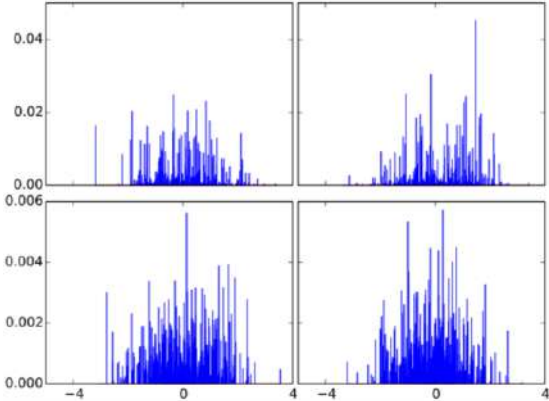
Non-parametric Models

Gaussian and Dirichlet Processes

Jose Gallego



Gaussian Processes



Dirichlet Processes

Gaussian Processes

Properties of Gaussians

- Normalization

$$\int p(x) dx = 1 \quad \frac{1}{(2\pi)^{D/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

- Marginalization

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_x & \Sigma_{xy} \\ \Sigma_{xy}^T & \Sigma_y \end{bmatrix}\right) \rightarrow x \sim \mathcal{N}(\mu_x, \Sigma_x)$$

- Addition

$$x \sim \mathcal{N}(\mu_x, \Sigma_x) \quad y \sim \mathcal{N}(\mu_y, \Sigma_y) \quad x+y \sim \mathcal{N}(\mu_x+\mu_y, \Sigma_x+\Sigma_y+2\Sigma_{xy})$$

- Conditioning

$$y|x \sim \mathcal{N}(\dots)$$

↳ characteristic function

- Product of densities

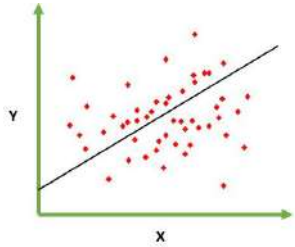
$$X \sim \mathcal{N}(0,1) \rightarrow X \cdot X \not\sim \mathcal{N} \rightarrow \sim \chi^2(1)$$

$$\phi(x+y) = \underline{\phi(x) \cdot \phi(y)}$$

$$\mathcal{N}(x|\mu_1, \Sigma_1) \cdot \mathcal{N}(x|\mu_2, \Sigma_2) = \mathcal{N}\left(\Sigma_3^{-1} \Sigma_1^{-1} \mu_1 + \Sigma_3^{-1} \Sigma_2^{-1} \mu_2, (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}\right)$$

↳ Σ_3

Linear Regression



$$y = f(x) = w^T x + \sigma \varepsilon$$

\uparrow
 $\mathcal{N}(0, \mathbb{I})$

$$p(y_i | x_i; w) = \mathcal{N}(w^T x_i, \sigma^2 \mathbb{I})$$

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$\begin{array}{l} \text{MLE} \\ \text{MAP} \end{array} \quad \underbrace{p(\mathcal{D} | w)} = \prod_i p(y_i | x_i; w) p(x_i) \quad \left| \quad \begin{array}{l} p(y_1, \dots, y_n | x_1, \dots, x_n, w) \\ = \prod_i p(y_i | x_i, w) \end{array} \right.$$

$$p(w | \mathcal{D}) = \log \downarrow + \frac{\log p(w)}{\hookrightarrow \|w\|^2}$$

$$\downarrow$$

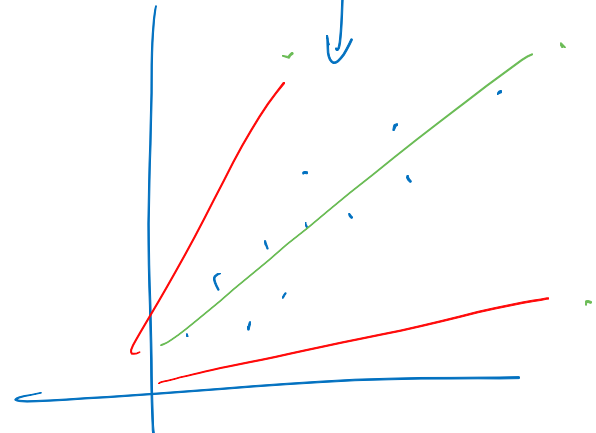
$$\underline{y = wx}$$

Bayesian Viewpoint (1)

$$\begin{aligned}
 p(y^* | x^k, \mathcal{D}) &= \int_w p(y, w | x, \mathcal{D}) dw = \int_w p(y | x, w) \cdot p(w | \mathcal{D}) dw \\
 &= \int_w \underbrace{p(y | x, w)}_G \cdot \underbrace{\frac{p(\Phi | w) \cdot p(w)}{p(\mathcal{D})}}_G dw = \mathcal{N}(\dots)
 \end{aligned}$$

Gaussian

"Use all of them"



Bayesian Viewpoint (2)

GP Assumption

$$P\left(\begin{bmatrix} y_1 \\ \vdots \\ y_n \\ y^* \end{bmatrix} \mid [x_1, \dots, x_n], x^*\right) \sim \text{some } \mathcal{N}(\mu, \Sigma).$$

$$P(y^* \mid \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, [x_1, \dots, x_n], x^*) \sim \mathcal{N}(\cdot, \cdot)$$

$$P(y^* \mid y_1, \dots, y_n, x_1, \dots, x_n, x^*) \sim \mathcal{N}\left(K_*^T K^{-1} y, K_{xx} - K_x^T K^{-1} K_x\right)$$

The role of Σ

$$\begin{bmatrix} 1 & 0.7 & 0 & \vdots \\ 0.7 & 1 & 0 & \vdots \\ 0 & 0 & 1 & 0 \\ \vdots & \vdots & 0 & \ddots \end{bmatrix}$$

me
negh
Imp

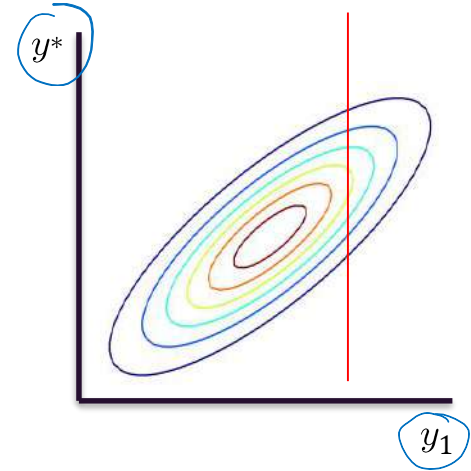
$$\sigma_{12} = \rho_{12} \sigma_1 \cdot \sigma_2$$

$\hookrightarrow [-1, 1]$

Positive Definite Kernel

$$\underline{k(x, x')} \geq 0 \quad \forall x, x'$$

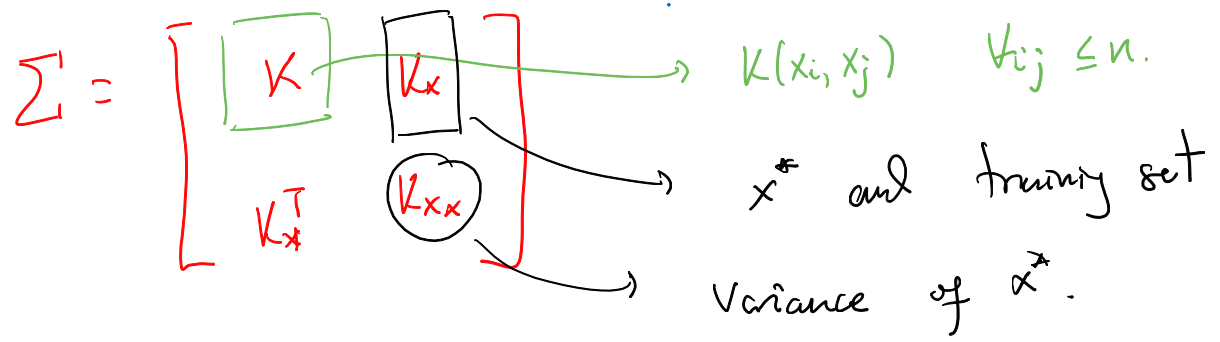
$$\forall \ell \quad \forall x_1, \dots, x_\ell \quad \underline{[K]_{ij}} = \underline{k(x_i, x_j)}$$



\hookrightarrow positive definite

Isn't this just kernel regression?

$$P(y^* | y_1, \dots, y_n, x_1, \dots, x_n) \sim \mathcal{N} \left(\underline{\underline{K_*^T K^{-1} y}}, \underline{\underline{K_{xx} - K_x^T K^{-1} K_x}} \right)$$



$$\underline{\underline{w^* = X^T (X^T X)^{-1} y}}$$

Demo

Demo

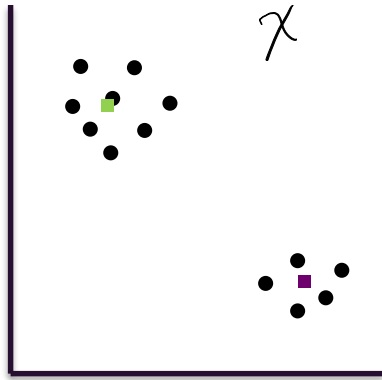
Problem: f is a function \sim infinite dimensional vector! But, the multivariate Gaussian distributions is for finite dimensional random vectors.

Definition: A GP is a (potentially infinite) collection of rvs such that the joint distribution of every finite subset of rvs is multivariate Gaussian.

$$y^* \mid y_1, \dots, y_n, x_1, \dots, x_n, x^* \sim \mathcal{N}(\mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{y}, k_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*)$$

Dirichlet Processes

Generative Model



$$z_n \stackrel{\text{iid}}{\sim} \text{Categorical}(\rho_1, \rho_2)$$

$$x_n | z_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_{z_n}, \Sigma)$$

$$\mu_k \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_0, \Sigma_0)$$

$$\rho_1 \sim \text{Beta}(a_1, a_2)$$

$$\rho_2 = 1 - \rho_1$$

Beta Distribution

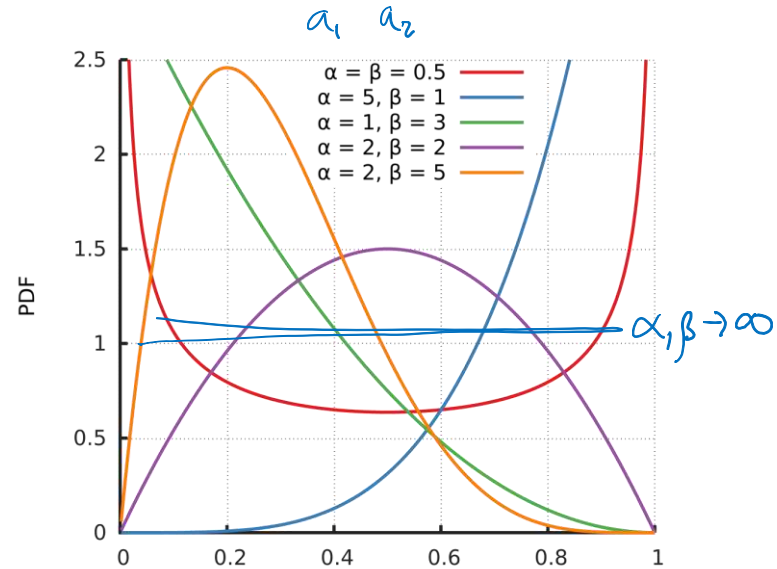
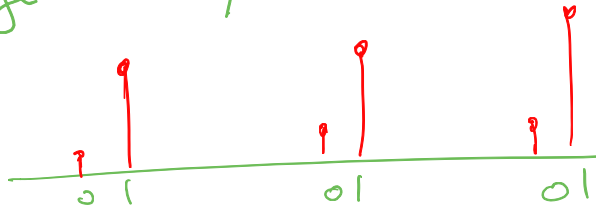
$$\text{Beta}(p_i | a_1, a_2) = \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1) \cdot \Gamma(a_2)} p_i^{a_1 - 1} (1 - p_i)^{a_2 - 1}$$

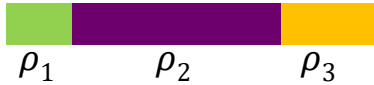
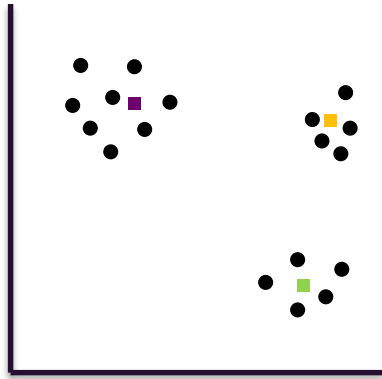
$$p_i \in (0, 1) \quad a_1, a_2 > 0$$

$$\alpha = \beta \rightarrow 0$$

$$[p_i, p_z] = [\text{Beta}, 1 - p_i]$$

mass @ 0 large $\Rightarrow p_i$ small





Same story

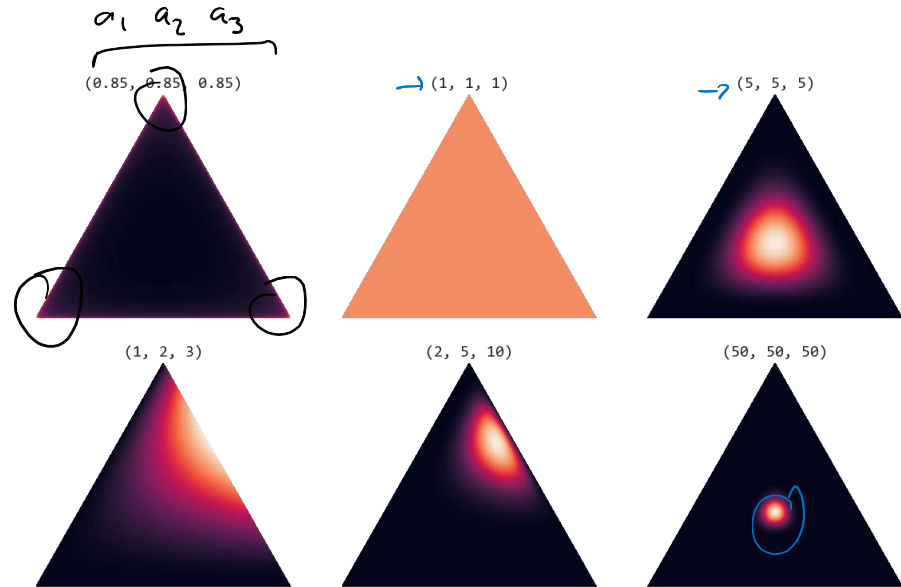
↳ Extension to Dirichlet distribution

$$\text{Cat}(\rho_1, \rho_2, \dots, \rho_k)$$

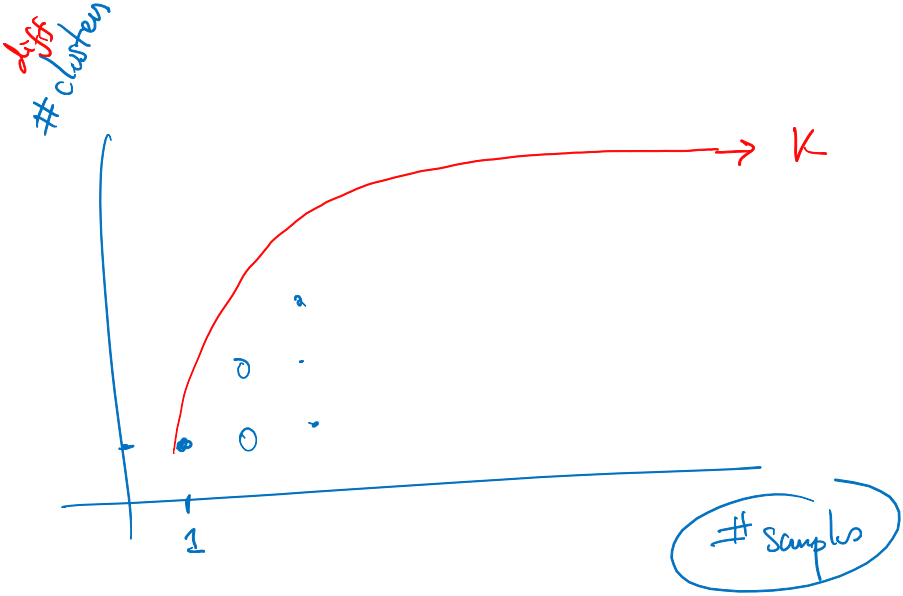
$$p \sim \text{Dir}(a_1, \dots, a_k)$$

Dirichlet Distribution

$$P(\beta_1, \dots, \beta_k \mid a_1, \dots, a_k)$$
$$= \frac{\Gamma(\sum_i a_i)}{\prod_i \Gamma(a_i)} \prod_i \beta_i^{a_i - 1}$$



What if $K > N$?



$\#$ components
vs
 $\#$ observed clusters

So, how do we choose K ?

We don't $K = \infty$

~~$\frac{1}{K}$~~

β_1, β_2, \dots

$$V_1 \sim \text{Beta}(\underline{a_1}, \underline{b_1})$$

$$1 \cdot V_1 = \beta_1$$

$$V_2 \sim \text{Beta}(\underline{a_2}, \underline{b_2})$$

$$V_2 (1 - p_1) = \beta_2$$

$$\underline{\beta_k} = V_k \left[\prod_{i=1}^{k-1} (1 - p_i) \right]$$

$\{(a_i, b_i)\}_{i=1}^{\infty}$

Infinitely many parameters!

Stick Breaking Process

$$V_i \sim \text{Beta}(1, \alpha)$$

$$\text{GEM}(\alpha)$$

$$\beta_1 = V_1$$

$$\beta_2 = V_2 (1 - \beta_1)$$

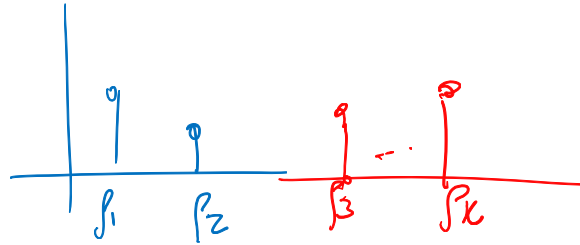
$$\vdots$$

$$1$$

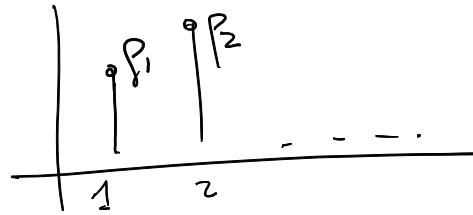
$$\beta_1, \beta_2, \dots \sim \text{GEM}(\alpha)$$

Distributions over Distributions over Distributions over Distributions...

Beta \rightarrow Distro (β_1, β_2)
 Dir \rightarrow Distro $(\beta_1, \beta_2, \dots, \beta_K)$



GEM \rightarrow Distro (\mathbb{N})
 $1, 2, 3, \dots$



$\beta = (\beta_1, \beta_2, \dots) \sim \text{GEM}(\alpha)$
 $\phi_k \stackrel{i.i.d.}{\sim} H \rightarrow$ distro over (Φ)

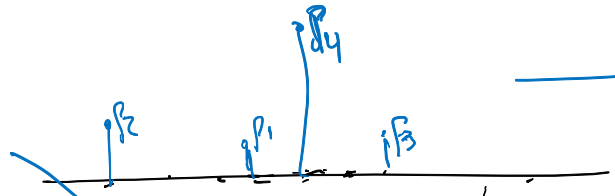
$$\Sigma = \sum_{i=1}^{\infty} \beta_i \int \phi_i$$



$$\Phi = \mathbb{R}$$

$$H \sim \mathcal{N}$$

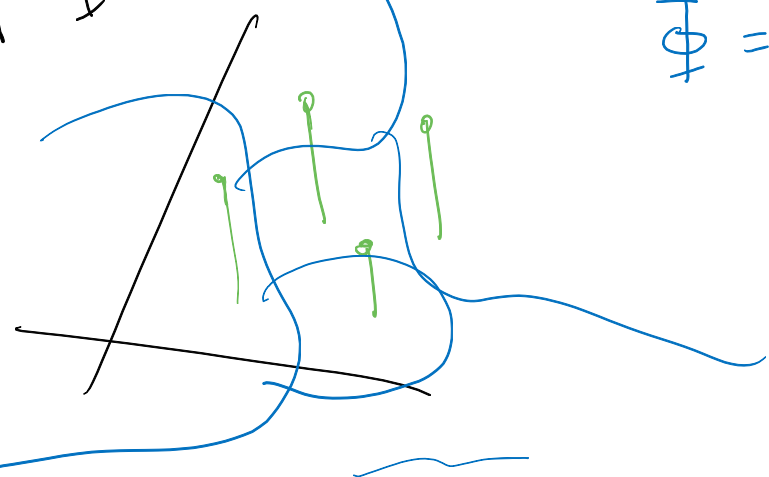
$$\mathcal{N} \quad \Phi = \mathbb{R}^2$$



$$\sum_i \underbrace{p_i} \underbrace{\delta_{\phi_i}} = G$$

$$\Phi = \bigsqcup_{i=1}^{\infty} A_i \leftarrow \text{(measurable) Partition of } \Phi$$

$$G(A_i) \in [0, 1]$$

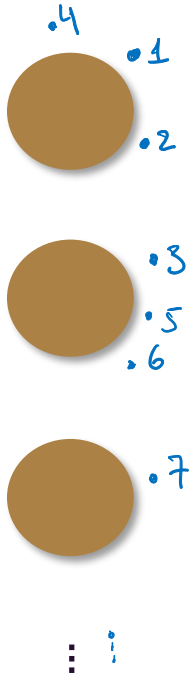


Marginal Cluster Assignments

$$[\beta_1, \beta_2, \beta_3] \sim \text{Dir} \begin{pmatrix} a_1 & a_2 & a_3 \\ 1 & 1 & 1 \end{pmatrix}$$
$$\beta_1 \sim \text{Beta} \left(1, \underbrace{2}_{\substack{a_2 \\ \sum_{i \neq 1} a_i}} \right)$$

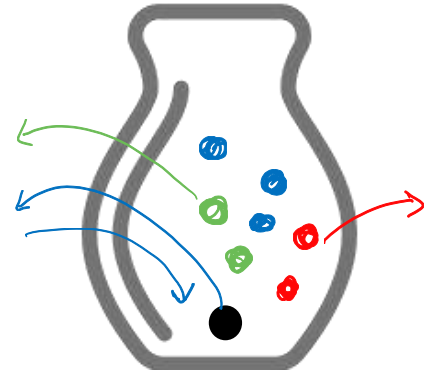
← Uniform over Δ_3 but not over β_1

Equivalent Constructions



Chinese Restaurant
Process

"Preferential
Attachment"



Problem: ρ_1, ρ_2, \dots is an *infinite dimensional* probability vector! But, the Dirichlet distributions is defined for finite dimensional spaces.

Definition: Let H be a distribution over Φ and $\alpha > 0$. Then we say that G is a Dirichlet process with base distribution H and concentration parameter α if for all finite (measurable) partition A_1, \dots, A_r ,

$$\underline{(G(A_1), \dots, G(A_r))} \sim \text{Dir}(\underline{\alpha H(A_1)}, \dots, \underline{\alpha H(A_r)}) \quad \leftarrow \text{Condition on finite objects}$$

Compare to finite number of vectors in G.P. \approx Finite partition of Φ