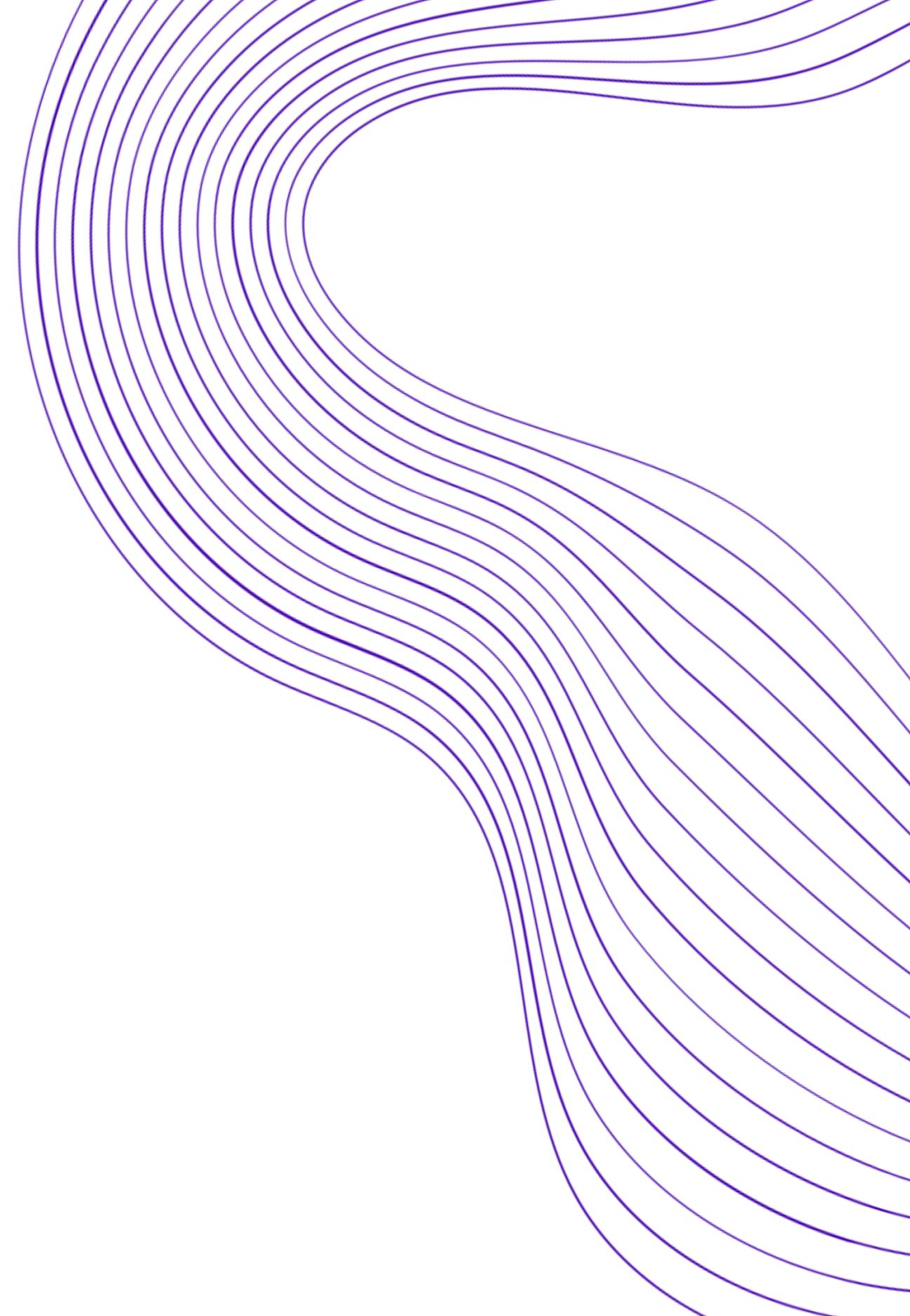


Bayesian

Non-Parametrics

JOSE GALLEGO-POSADA

IFT6269 | December 4, 2020



Contents

- Stochastic Processes
- Gaussian Processes
- Dirichlet Processes

Caveats

- Ideas over implementation/training details
- Hyperparameter selection
- Basic introduction - See suggested sources

3

◆ David MacKay

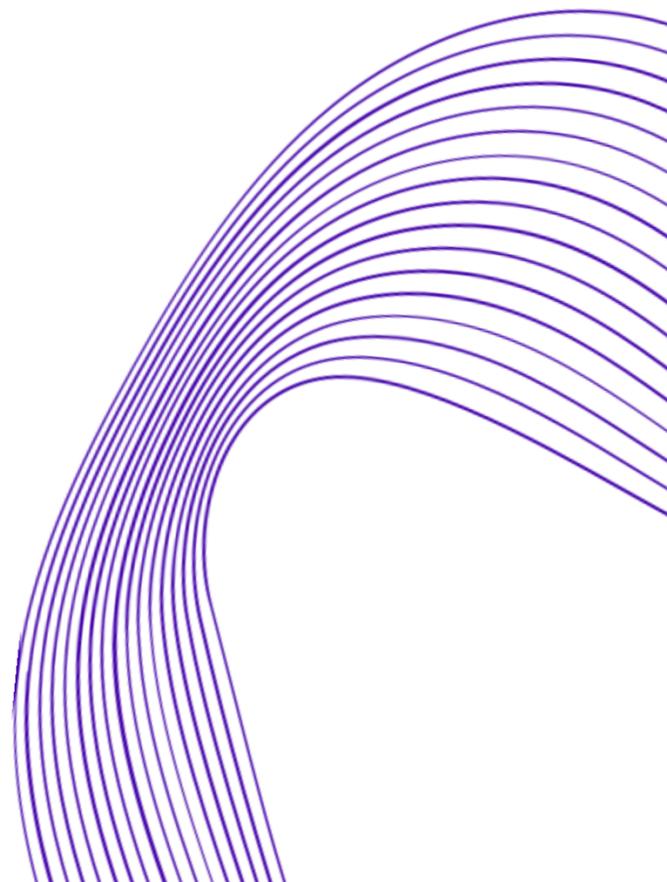
◆ Yee Whye Teh

◆ Kilian Weinberger

◆ Tamara Broderick

◆ Michael Jordan

Inspiration



Stochastic Process

A stochastic process is defined as a collection of random variables defined on a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$, where Ω is a sample space, \mathcal{F} is a sigma-algebra, and \mathbb{P} is a probability measure; and the random variables, indexed by some set T , all take values in the same space S , measurable with respect to some σ -algebra.

$$\{X(t) \mid t \in T\}$$

◆ Single random variable

$$T$$

$$\{1\}$$

$$\{X(t) \mid t \in T\}$$

$$\{X_1\}$$

◆ IID random variables

$$\{1, 2, \dots, n\}$$

$$\{X_1, \dots, X_n\}$$

$$X_j \stackrel{d}{=} X_k \text{ indep.}$$

◆ Deterministic function

$$f: T \rightarrow S ; T \checkmark$$

$$\{X(t) \triangleq f(t) \mid t \in T\}$$

◆ {Wiener, Poisson, etc.} Process

$$T = \mathbb{R}^+$$

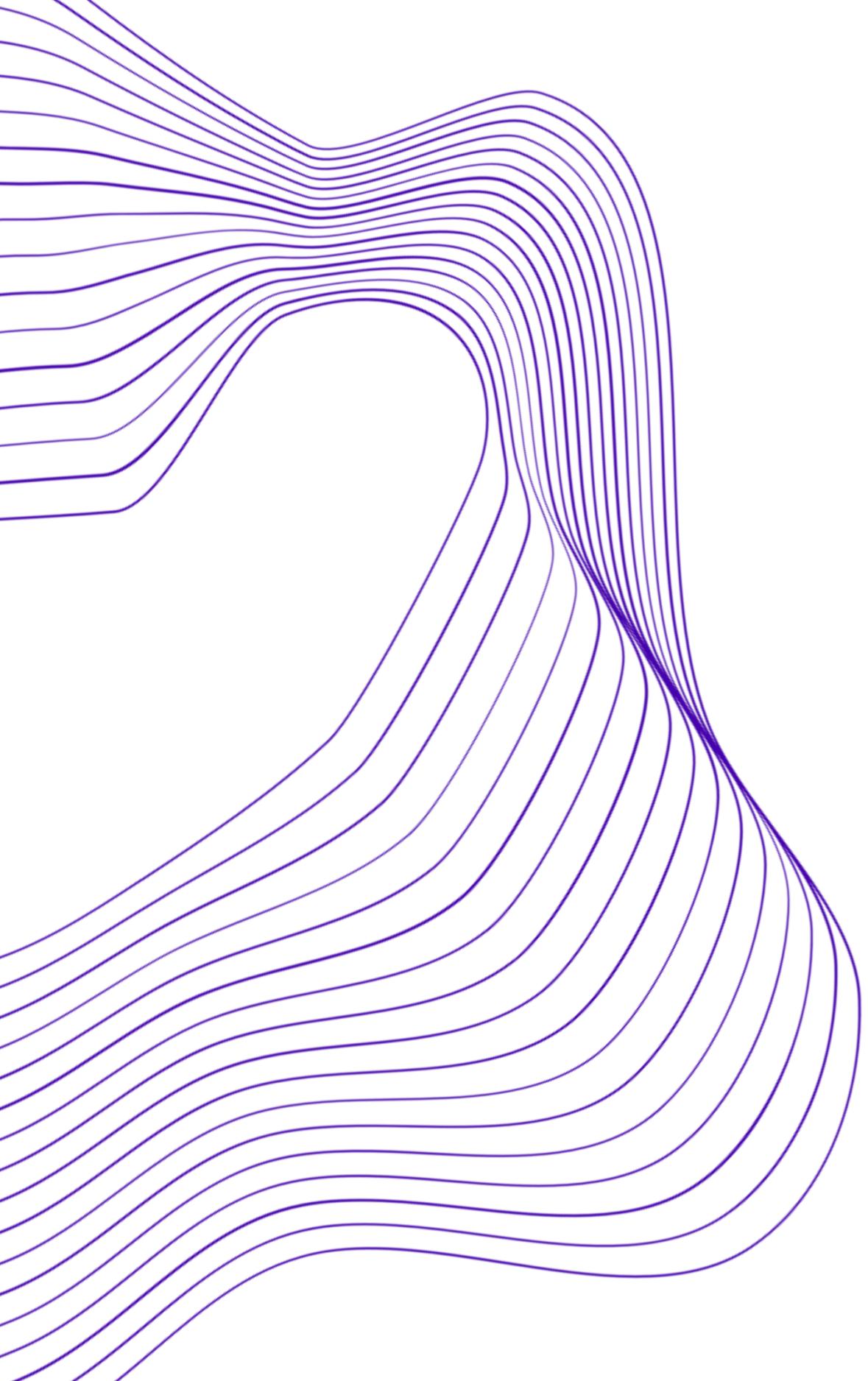
$$X(t) \sim \omega(0, t) ; \text{ a.s. cont.}$$

◆ Gaussian Process

◆ Dirichlet Process

Examples of Stochastic Process

Gaussian Processes

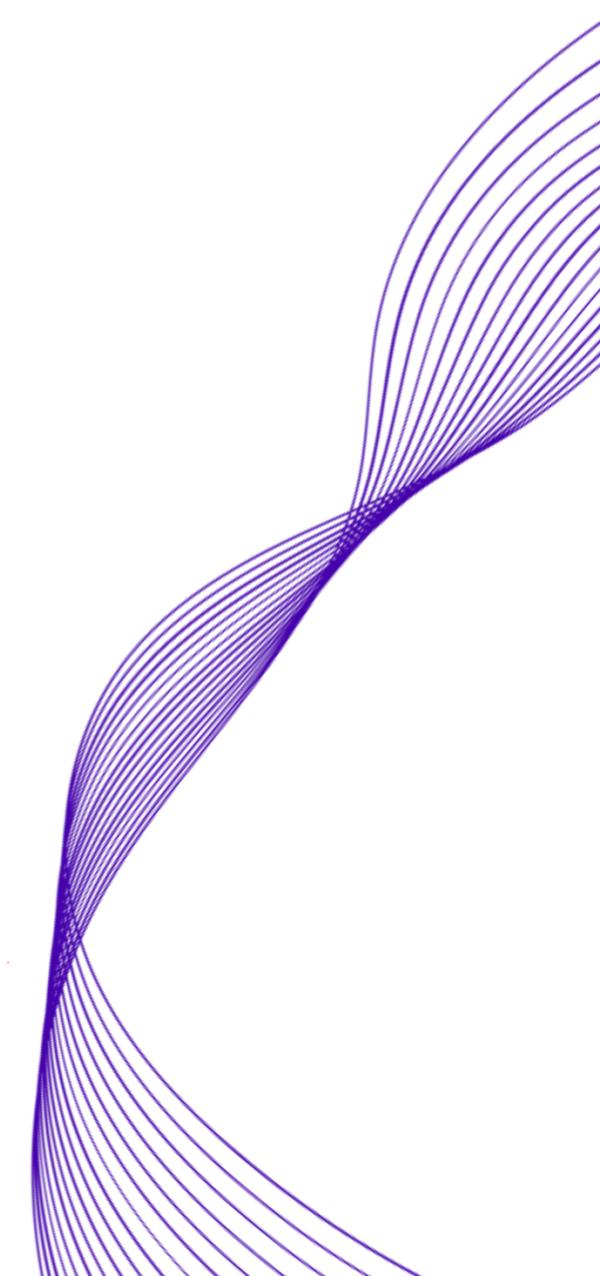


Inference and Prediction

$$\mathbb{P}(y(\cdot) | \mathbf{t}, \mathbf{X}) = \frac{\mathbb{P}(\mathbf{t} | y(\cdot), \mathbf{X}) \mathbb{P}(y(\cdot))}{\mathbb{P}(\mathbf{t} | \mathbf{X})}$$

Lik. *Prior*
Evidence

$$\mathbb{P}(t^* | \mathbf{x}^*, \mathbf{t}, \mathbf{X}) = \int \mathbb{P}(t^* | y(\cdot), \mathbf{x}^*) \mathbb{P}(y(\cdot) | \mathbf{t}, \mathbf{X}) dy$$



Gaussian Properties

◈ (Tractable) Normalization $\int \tilde{p}(x) dx = (2\pi)^{d/2} |\Sigma|^{\frac{1}{2}}$

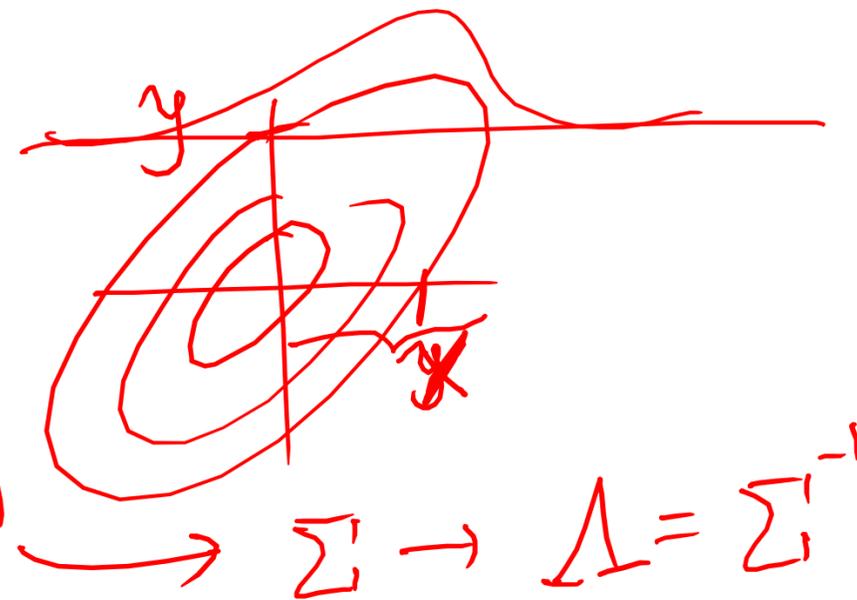
$\begin{bmatrix} X \\ Y \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix}, \begin{bmatrix} \Sigma_X & \Sigma_{XY} \\ \Sigma_{XY}^\top & \Sigma_Y \end{bmatrix} \right)$

◈ Marginalization

$X \sim \mathcal{N}(\mu_X, \Sigma_X)$

◈ Conditioning

$X|Y = \underline{y} \sim \mathcal{N}(\mu_X + \Sigma_{XY} \Sigma_Y^{-1} (y - \mu_Y), \Lambda_X^{-1})$



◈ Addition

$X+Y \sim \mathcal{N}(\mu_X + \mu_Y, \Sigma_X + \Sigma_Y + 2\Sigma_{XY})$

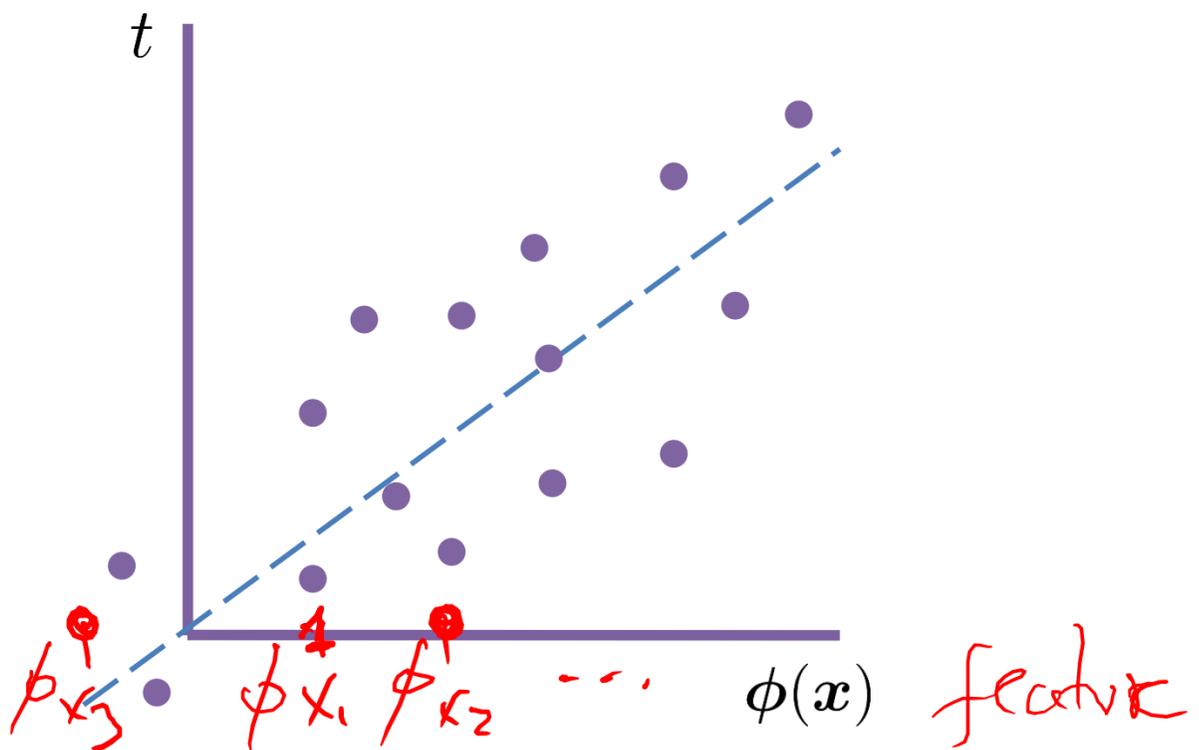
◈ Product of densities $\mathcal{N}(x | \mu_1, \Sigma_1) \cdot \mathcal{N}(x | \mu_2, \Sigma_2) \neq \mathcal{N}(x | \mu_3, \Sigma_3)$

Not prod of v.i.v.s $X \sim \mathcal{N}(0,1) \rightarrow X \cdot X = X^2 \sim \chi^2_{(1)}$

Linear Regression

Starting point:

- Locations $\{\mathbf{x}_i | i \in [1, N]\}$
- Basis functions $\{\phi_h(\mathbf{x})\}_{h=1}^H$
- Linear model $y(\phi(\mathbf{x}), \mathbf{w}) = \mathbf{w}^\top \phi(\mathbf{x})$
- Prior $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I})$
- Observation noise $t|y \sim \mathcal{N}(y, \sigma_\nu^2)$



This leads to $\mathbf{y}_N \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \Phi_N \Phi_N^\top)$

$\mathbf{t}_N \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \Phi_N \Phi_N^\top + \sigma_\nu^2 \mathbf{I})$

$(\Phi_N \Phi_N^\top)_{ij} = (\text{feat})^\top (\text{feat})$
 (PSD) similarity matrix
 similarity matrix $C(\mathbf{x}_i, \mathbf{x}_j)$
 $C(\mathbf{x}_i, \mathbf{x}_j)$...

y_1
 \vdots
 y_N

Gaussian Process

The probability distribution of a function $y(\mathbf{x})$ is a **Gaussian process** if for any **finite selection** of points $(\mathbf{x}_1, \dots, \mathbf{x}_N)$, the vector $[y(\mathbf{x}_1), \dots, y(\mathbf{x}_N)]$ follows a **Gaussian distribution**.

$$\bigcup_{N=1}^{\infty} \mathcal{X}^N$$

Index set T

$$\mathbb{R}^d$$

Value space S

Gaussians

RVs ←

(function values)

Habemus Datam!

... we finally observe $t_N = \{t_i | i \in [1, N]\}$ and would like to infer t_{N+1} for a new test point \mathbf{x}_{N+1} .

$$\mathbf{C}_{N+1} = \begin{bmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^\top & \kappa \end{bmatrix}$$

$$k_j = C(x_{N+1}, x_j)$$

$$\kappa_{N+1}$$

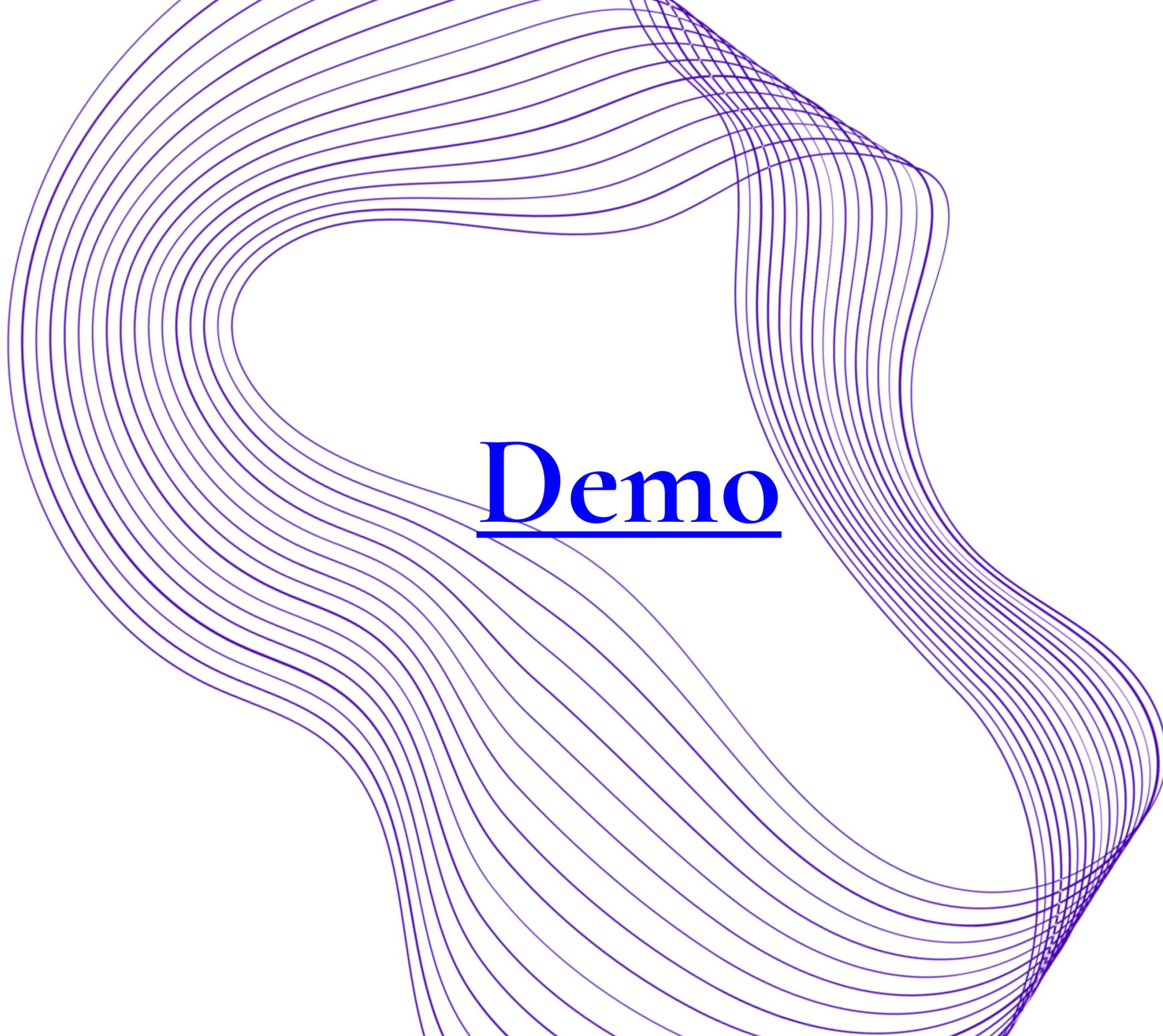
$$\mathbb{P}(t_{N+1} | \mathbf{t}_N) = \frac{\mathbb{P}(t_{N+1}, \mathbf{t}_N)}{\mathbb{P}(\mathbf{t}_N)} \propto \exp \left[-\frac{1}{2} \begin{bmatrix} \mathbf{t}_N & t_{N+1} \end{bmatrix} \mathbf{C}_{N+1}^{-1} \begin{bmatrix} \mathbf{t}_N \\ t_{N+1} \end{bmatrix} \right]$$

\mathbb{P}_N dependent

$$t^* | t_1, \dots, t_N, \mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_i^* \sim \mathcal{N} \left(\mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{t}_N, \kappa - \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{k} \right)$$

$$\mathbf{k}^\top \mathbf{C}_N^{-1}$$

Demo



$$P(t^* = 1 | \dots) = \sigma(\text{Regression})$$

GP Summary

D. Mackay / C. Bishop ^{Logits} \Rightarrow Rasmussen '06

Problem: y is a function aka “infinite dimensional vector”. But the multivariate Gaussian distribution is defined for **finite** dimensional vectors.

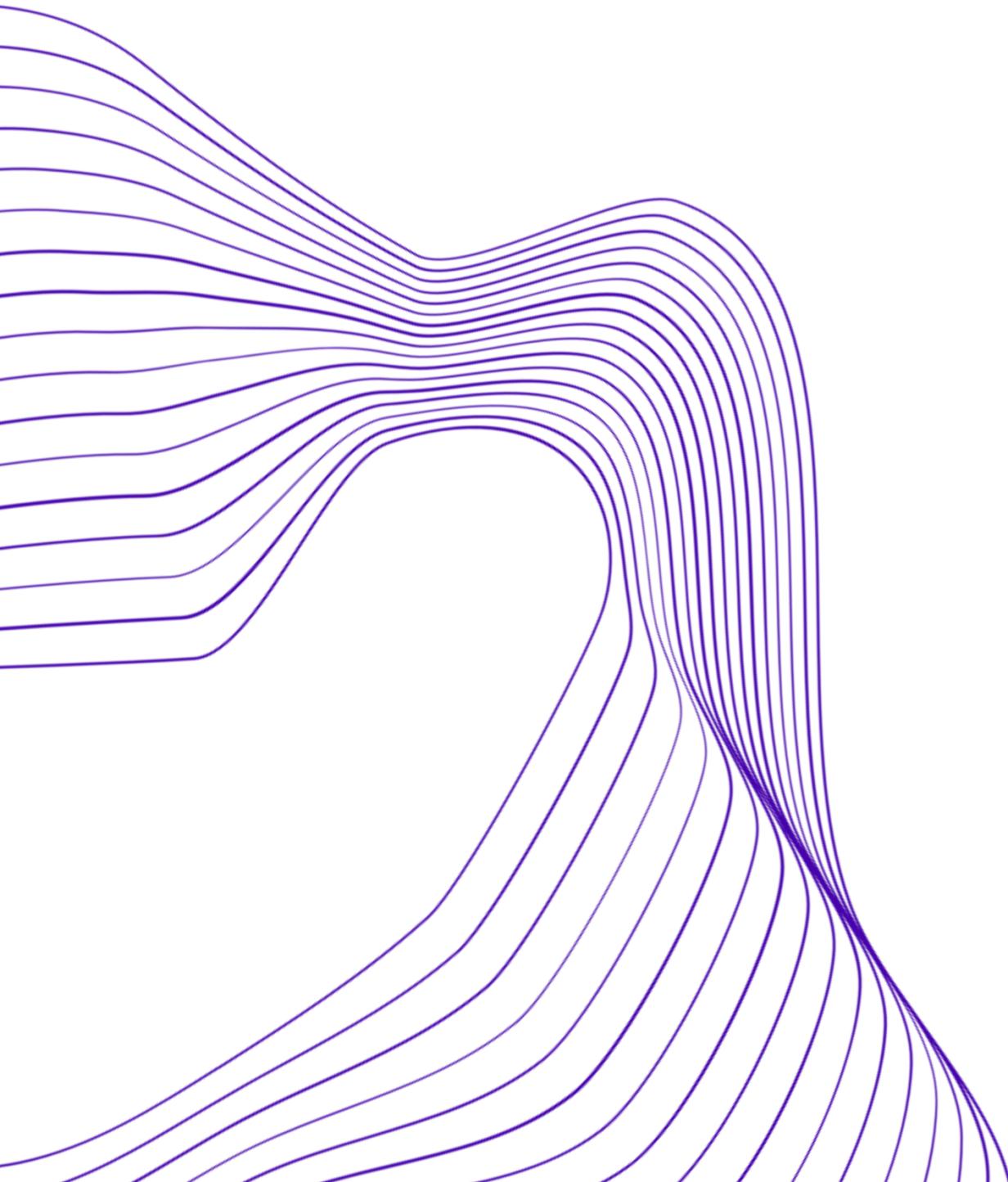
$y_N \sim$

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

$$t^* | t_1, \dots, t_N, x_1, \dots, x_n, x^* \sim \mathcal{N}(\mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{t}_N, \kappa - \mathbf{k}^\top \mathbf{C}_N^{-1} \mathbf{k})$$

$$[y(x_1) \dots y(x_n)] \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_N)$$

Dirichlet Processes



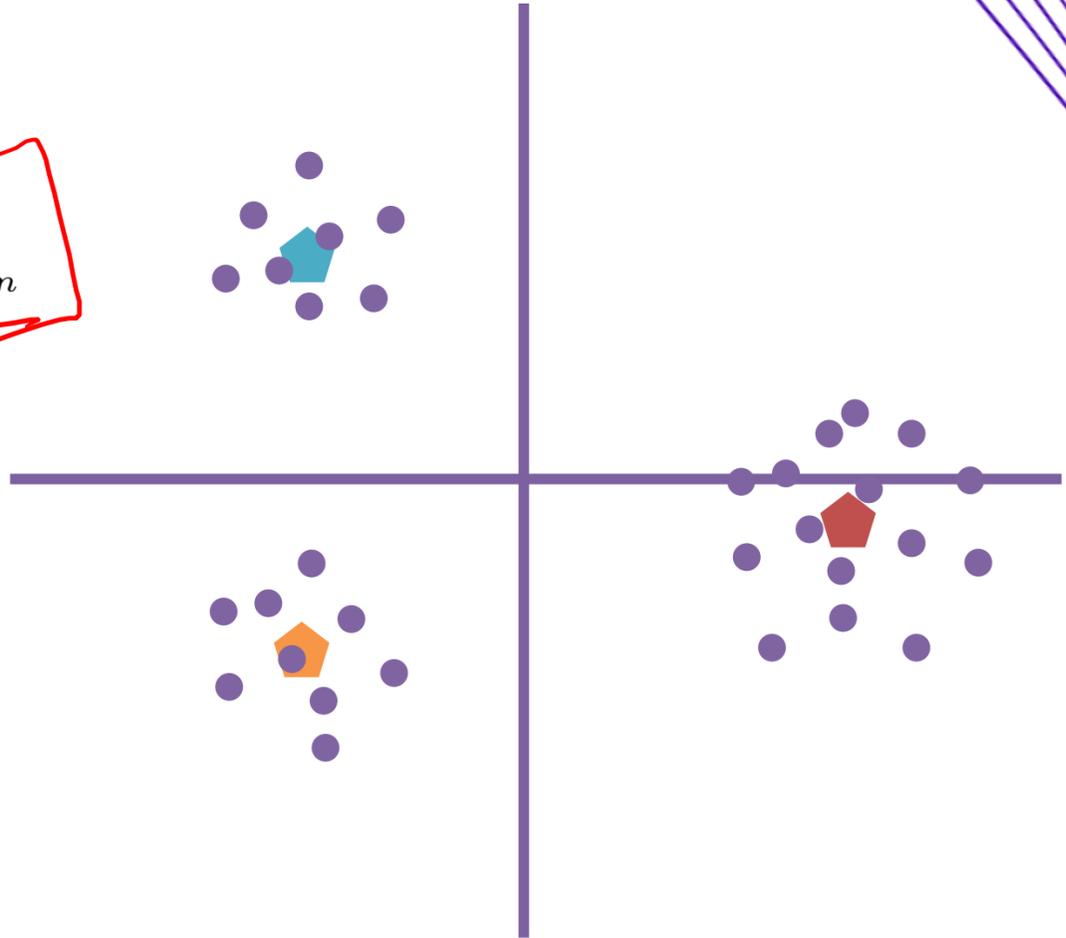
Model:

- Assignments $z_n | \boldsymbol{\rho} \sim \text{Categorical}(\rho_1, \dots, \rho_K)$
- Class conditionals $\mathbf{x}_n | z_n \sim \mathcal{N}(\boldsymbol{\mu}_{z_n}, \boldsymbol{\Sigma}) = F_{z_n}$

$\rho_i \geq 0 \quad \sum_i \rho_i = 1$

Prior assumptions:

- Number of classes K
- Class proportions $\boldsymbol{\rho} \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$
- Class parameters $\boldsymbol{\mu}_k \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) = H$



Generative Model

Dirichlet Distribution

over Distributions

$$\mathbb{P}(\rho_1, \dots, \rho_K | \alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \rho_k^{\alpha_k - 1}$$

$\rho_i \sim \text{Gamma}(\alpha_i, \theta)$

$$\vec{\rho} = \frac{\vec{\rho}}{\|\vec{\rho}\|_1}$$

Depends on all entries

$(n+1)!$ $(n+1) n!$

$$\Gamma(z + 1) = z\Gamma(z)$$

$\mathbb{P}(\rho_1, \rho_2 | \alpha_1, \alpha_2)$

$$= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)}$$

$\rho_1^{\alpha_1 - 1} \rho_2^{\alpha_2 - 1}$

Collapsibility

$$\rho_1, \dots, \rho_i + \rho_j, \dots, \rho_K \sim \text{Dir}(\alpha_1, \dots, \alpha_i + \alpha_j, \dots, \alpha_K)$$

$$\alpha_0 = \sum_i \alpha_i$$

$$\alpha_0 = 4.5$$

$$\alpha_0 = 15$$

$$(p_i, \sum_{j \neq i} p_j) \sim \text{Dir}(\alpha_i, \sum_{j \neq i} \alpha_j)$$

$$\alpha_1, \alpha_2, \alpha_3 = (1.5, 1.5, 1.5)$$

$$\alpha_1, \alpha_2, \alpha_3 = (5, 5, 5)$$

$$p_i \sim \text{Beta}(\alpha_i, \alpha_0 - \alpha_i)$$

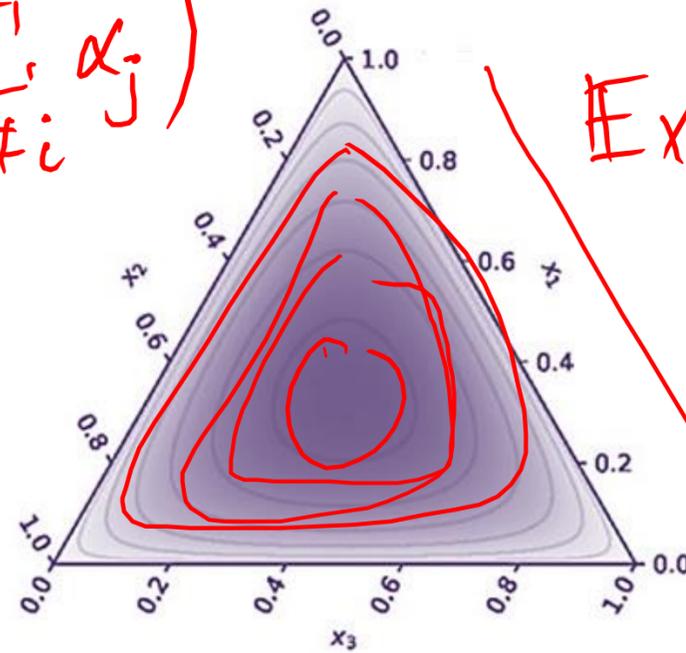
$$W \sim \text{Beta}(a, b) \Rightarrow \mathbb{E}W = \frac{a}{a+b}$$

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

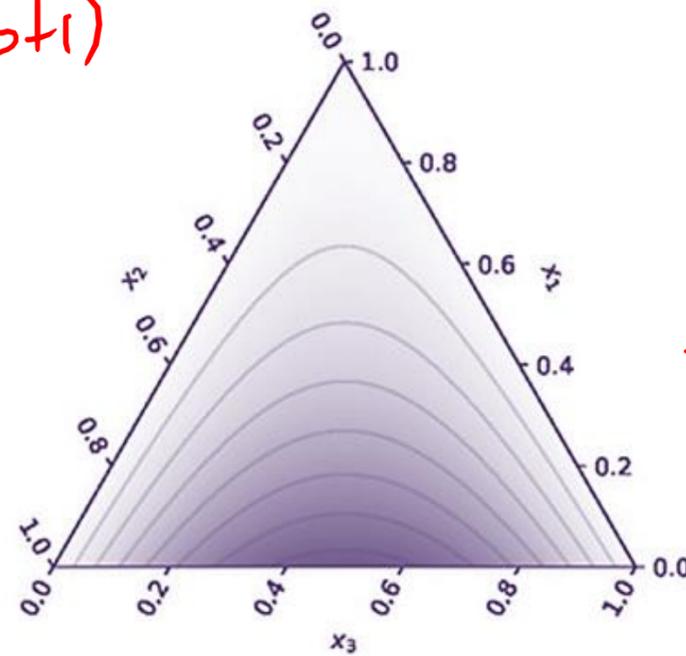
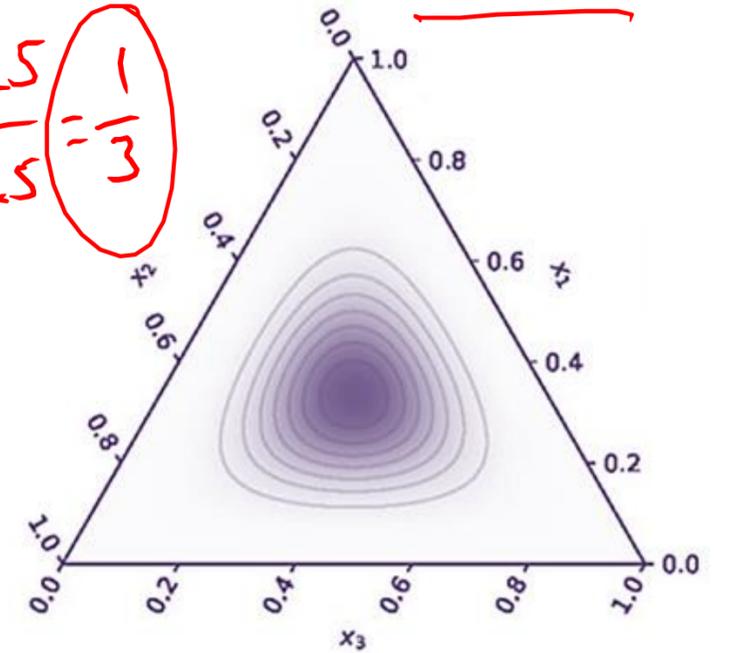
$$\mathbb{E}p_i = \frac{\alpha_i}{\alpha_i + \alpha_0 - \alpha_i} = \frac{\alpha_i}{\alpha_0}$$

$$\mathbb{V}p_i = \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}$$

$$\text{Cov}(p_i, p_j) = \frac{-\alpha_i \alpha_j}{\alpha_0^2(\alpha_0 + 1)}$$

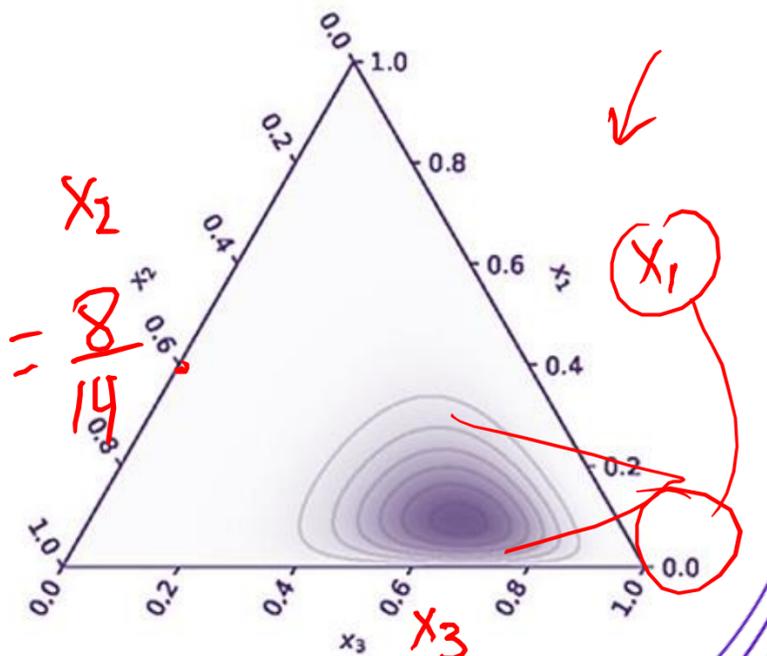


$$\mathbb{E}x_1 = \frac{1.5}{4.5} = \frac{1}{3}$$



$$\alpha_1, \alpha_2, \alpha_3 = (1, 2, 2)$$

$$\mathbb{E}x_3 = \frac{8}{14}$$



$$\alpha_1, \alpha_2, \alpha_3 = (2, 4, 8)$$

$$p_1 \sim \text{Beta}(2, 12)$$

$$\mathbb{E}p_1 = \frac{2}{14}$$

$$p_3 \sim \text{Beta}(8, 6)$$

Dirichlet Simulation

Stick-breaking

$$\rho_1 \sim \text{Beta}(\alpha_1, \sum_{k=1}^K \alpha_k - \alpha_1) \perp (\rho_2, \dots, \rho_K) \sim \text{Dir}(\alpha_2, \dots, \alpha_K)$$

Handwritten notes: ρ_1 is circled in red with ν_1 above it and 0.4 to its left. The Dirichlet distribution is boxed in purple with $1 - \rho_1$ circled in red and α_1 below it.

$$\nu_2 \sim \text{Beta}(\alpha_2, \sum_{k=2}^K \alpha_k - \alpha_2)$$

Handwritten notes: ν_2 is circled in red with 0.7 to its left. Below it, $\rho_2 = (1 - \nu_1)\nu_2 = 42\%$ is written in red, with 0.6 and 0.7 underlined.

$$\vdots$$

$$\nu_l \sim \text{Beta}(\alpha_l, \sum_{k=l}^K \alpha_k - \alpha_l)$$

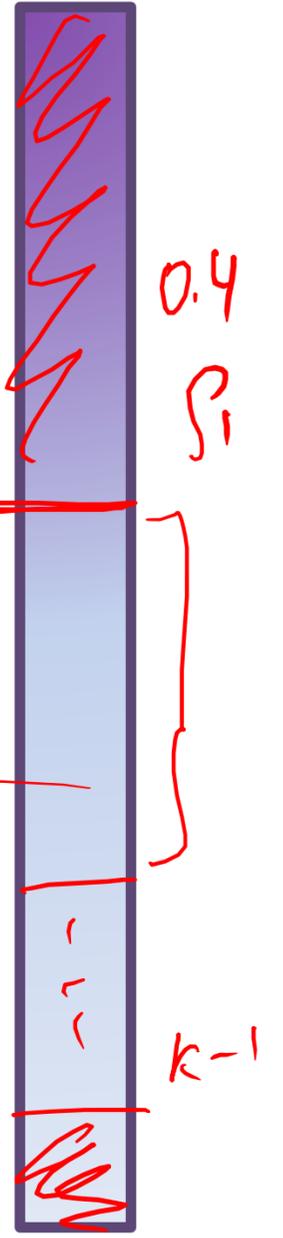
$$\rho_l = \left(\prod_{k=1}^{l-1} (1 - \nu_k) \right) \nu_l$$

Handwritten notes: ν_l is circled in red. 0.42 is written in red next to ρ_l .

$$\vdots$$

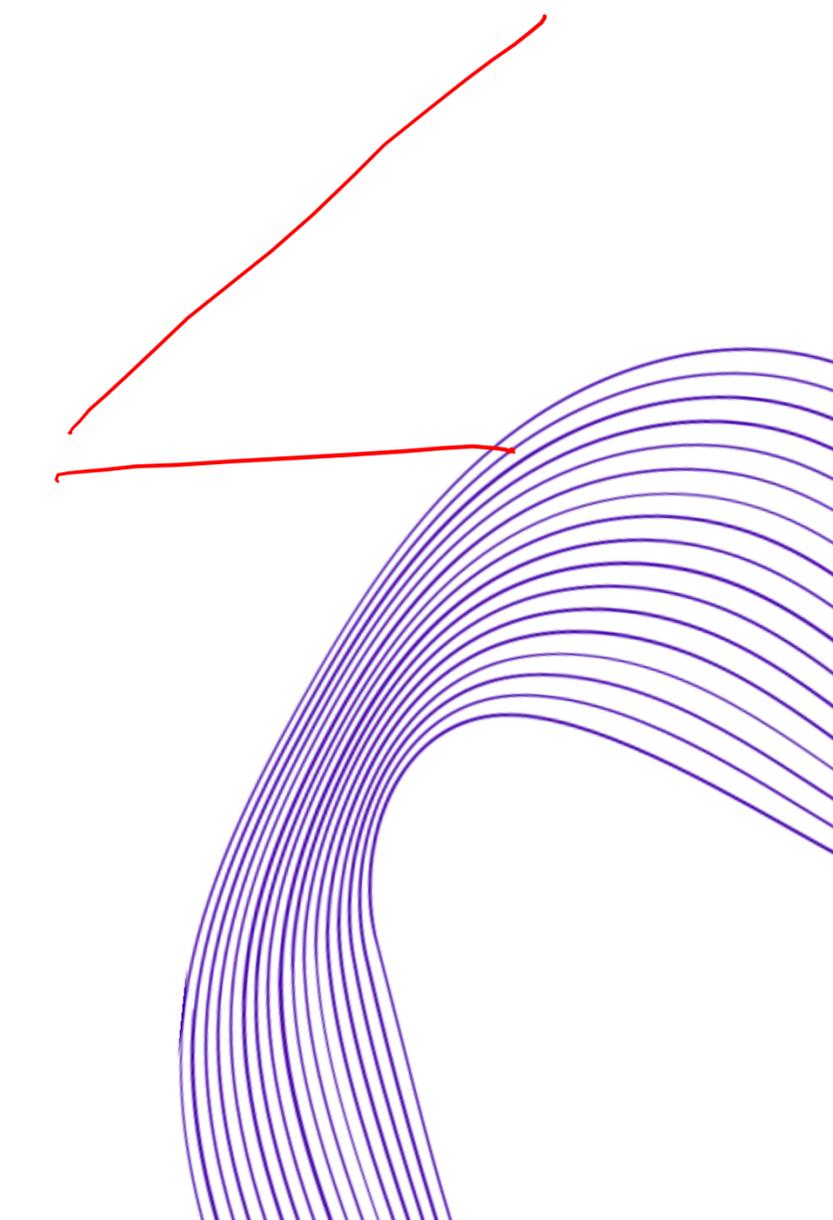
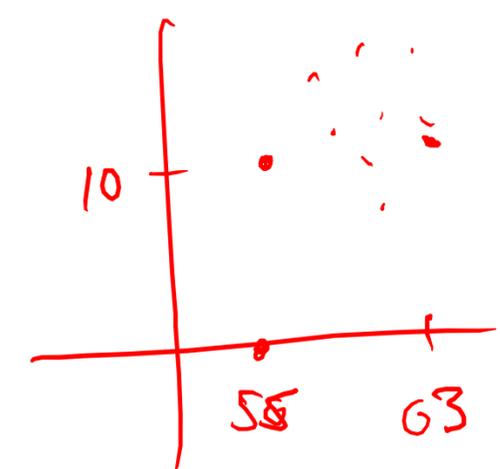
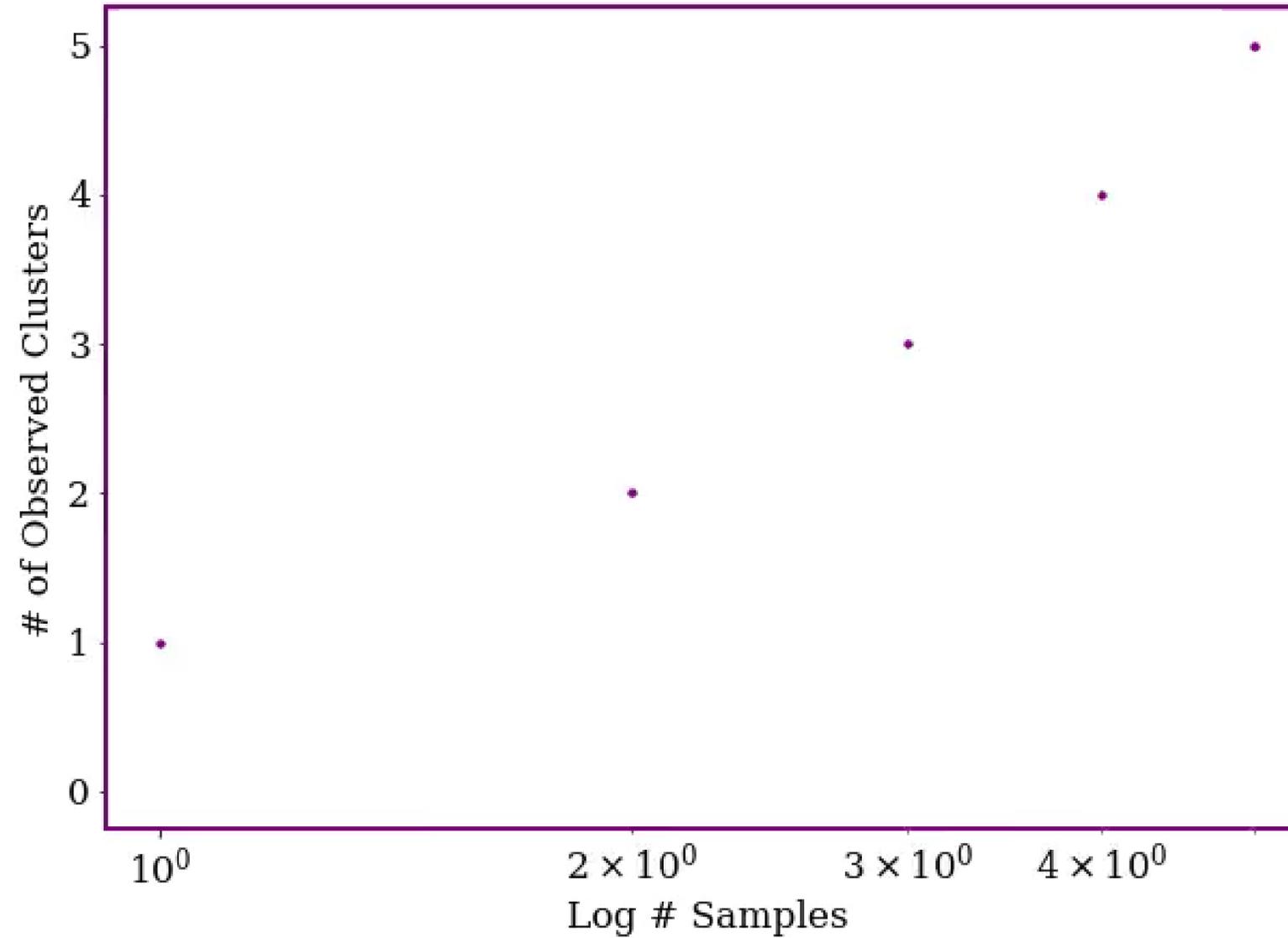
$$\rho_K = 1 - \sum_{k=1}^{K-1} \rho_k$$

Handwritten notes: ρ_K is circled in red. A red arrow points from this equation to a vertical bar on the right. The bar is divided into segments labeled 0.4 , 0.42 , and $k-1$.



fix $p \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_{10})$ 100 comp.

$Z_n \sim \text{Cat}(p)$



observed

10^4

Components vs Clusters

How do we choose K?

We don't! Consider an infinite mixture model.

$$\nu_1 \sim \text{Beta}(\alpha_1, \beta_1)$$

$$\nu_2 \sim \text{Beta}(\alpha_2, \beta_2)$$

⋮

$$\nu_l \sim \text{Beta}(\alpha_l, \beta_l)$$

⋮

$$\rho_1 = \nu_1$$

$$\rho_2 = (1 - \nu_1) \nu_2$$

⋮

$$\rho_l = \left(\prod_{k=1}^{l-1} (1 - \nu_k) \right) \nu_l$$

⋮

$$E \nu_l = \frac{1}{1 + \alpha}$$

$$\alpha \ll 1 \Rightarrow 1$$

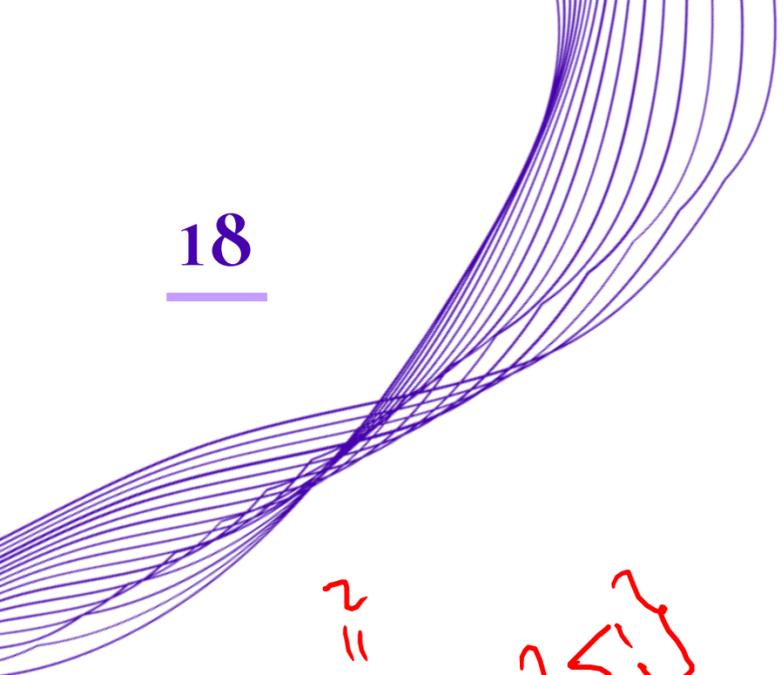
$$\alpha \gg 1 \Rightarrow 0$$

Griffiths, Engen & McCloskey

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

$$\rho_1, \rho_2, \dots \sim \text{GEM}(\alpha)$$

Conditions on $\{\alpha_l, \beta_l\}$ for proper normalization?



$$\rho_1, \rho_2, \dots \sim \text{GEM}(\alpha)$$

$$\phi_1, \phi_2, \dots \stackrel{\text{iid}}{\sim} H$$

e.g. $\mathcal{N}(\mu_0, \Sigma_0)$

Two sources of randomness

Distrn on S

Distro on S

$$\mathbb{E} \left[\sum_k \rho_k \mathbb{E} \int_S \delta_{\phi_k}(A) \right]$$

$$G = \sum_{k=1}^{\infty} \rho_k \delta_{\phi_k}$$

$$\mathbb{E} \left[\sum_k \rho_k H(A) \right]$$

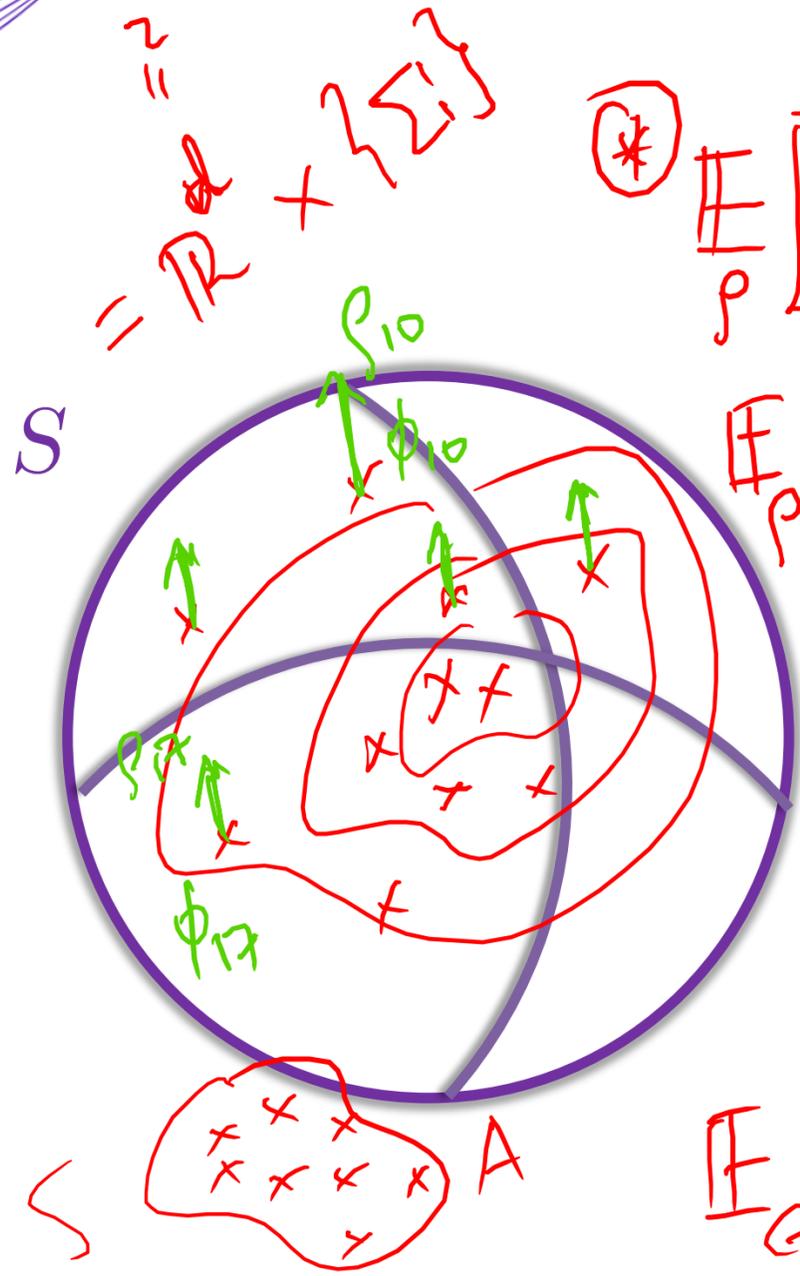
$H(A)$ random

Realizations of z_n are ϕ_k 's $\leftarrow (\mu_k, \Sigma)$

$$z_n | G \sim G$$

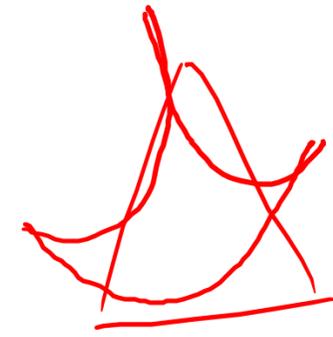
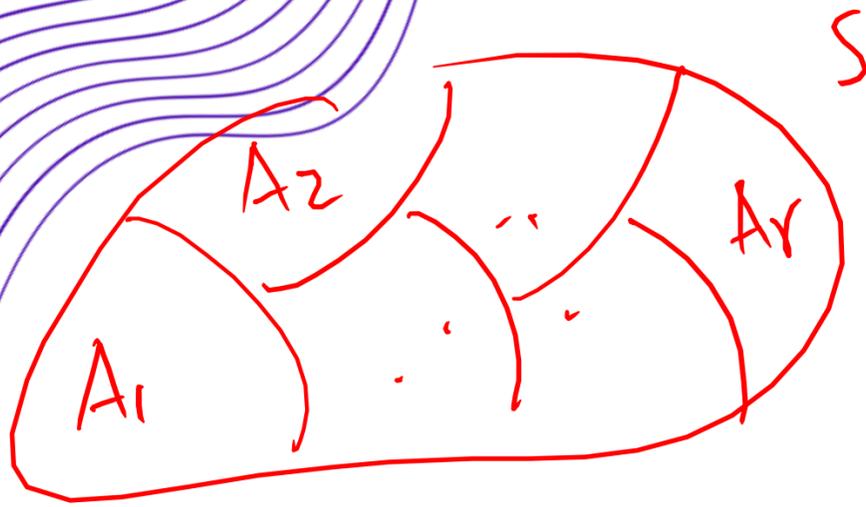
$$x_n | z_n \sim F_{z_n}$$

$\leftarrow \mathcal{N}(\mu_{z_n}, \Sigma)$



$$\mathbb{E}_G [G(A)] = \mathbb{E}_G \left[\sum_k \rho_k \delta_{\phi_k}(A) \right]$$

Random Measures



Dirichlet Process

Given a measurable space (S, \mathcal{M}) , a base probability distribution H and $\alpha > 0$, the probability distribution of a measure G is a **Dirichlet process** if for any finite partition $\{A_1, \dots, A_r\}$ of S , the vector $[G(A_1), \dots, G(A_r)] \sim \text{Dir}(\alpha H(A_1), \dots, \alpha H(A_r))$

Concentration parameter

\mathcal{M}

Index set T

Dirichlet

RVs

Ex Behavior

$\text{DP}(\alpha, H)$

$\alpha \rightarrow 0 ?$

$\alpha \rightarrow \infty ?$

DP Summary

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

Definition: A DP is a distribution over probability measures such that for **every finite partition**, the probabilities of the partition elements follow a joint Dirichlet distribution.

$$[G(A_1), \dots, G(A_r)] \sim \text{Dir}(\alpha H(A_1), \dots, \alpha H(A_r))$$

\uparrow
 X_1

\uparrow
 X_Y

Questions?

