

today: Bayesian non-parametrics  
 GP  
 DP

### Bayesian non-parametrics

non-parametric model  $\rightarrow$  infinite # of parameters  
 (or growing with # of data points)

- e.g. • KNN classifier  $\rightarrow$  boundary complexity grows with # of datapoints
- kernel density estimation  $\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n K(x, x_i)$

Bayesian non-parametric  $\rightarrow$  need prior over  $\infty$ -dim parameter  
 $\rightarrow$  define dist. on  $\infty$ -dim vector (stochastic process)

### stochastic process

collection of random variables indexed by a (potentially infinite) index set  $T$

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

examples:

$$T = \{1, \dots, n\} \quad X(t) = X_t \rightarrow \text{random vector } (X_1, \dots, X_n)$$

but also  $T = \mathbb{N}$   $\rightarrow$   $\infty$ -sequence  $(X_1, X_2, X_3, \dots, )$

or  $T = \mathbb{R}$   $\rightarrow$  "random function"

Gaussian process  $\rightarrow$  random functions

Dirichlet process  $\rightarrow$  random measures/distributions

### Gaussian process:

Motivation from Bayesian linear regression

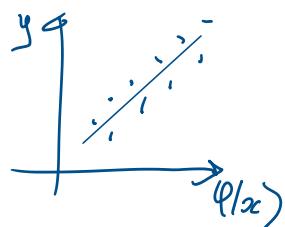
consider fixed location  $x_1, \dots, x_n$  [conditional  $y|x$ ]

model of  $Y|X$ :  $y = w^T \varphi(x) + \sigma_y \varepsilon$

thus  $Y|X, w \sim N(w^T \varphi(x), \sigma_y^2)$   
 (i.i.d.  $N(0, 1)$ )

Bayesian: prior on  $w \sim N(0, \sigma_w^{-2} I)$   $y_i \triangleq y|x_i$

$$\mathbb{E}[\mathbb{E}[y(x)|w]] = \mathbb{E}[w^T \varphi(x)] = \mathbb{E}[y]^T \varphi(x) = 0$$



$$\begin{aligned}\mathbb{E}[\mathbb{E}[y(x) | w]] &= \mathbb{E}[w^T \varphi(x)] = \mathbb{E}[\varphi(x)^T w] = 0 \\ \mathbb{E}[y_i y_j] &= \mathbb{E}[(w^T \varphi(x_i) + \sigma_y \varepsilon)(w^T \varphi(x_j) + \sigma_y \varepsilon)] \\ &= \mathbb{E}[w^T \varphi(x_i) w^T \varphi(x_j) + \underbrace{\sigma_y^2 \varepsilon^2}_{0}] + 0 \\ &= \mathbb{E}[\text{tr}(\varphi(x_i)^T w w^T \varphi(x_j))] + \sigma_y^2 \\ &= \text{tr}(\varphi(x_i)^T \mathbb{E}[w w^T] \varphi(x_j)) + \sigma_y^2 \\ &= \sigma_w^2 \underbrace{\varphi(x_i)^T \varphi(x_j)}_{K(x_i, x_j)} + \sigma_y^2 \\ &\quad K(x_i, x_j) \rightarrow \text{similarly}\end{aligned}$$

so generally, marginal on  $y$ 's [function values] a priori:

$$y_{1:n} \sim N(0, \sigma_w^2 \underbrace{\Phi_N \Phi_N^T}_{K \text{ (kernel matrix)}} + \sigma_y^2 I_n)$$

$$\Phi_N = \begin{pmatrix} -\varphi(x_1) \\ \vdots \\ -\varphi(x_n) \end{pmatrix}$$

Gaussian process: generalization of Gaussian to  $n$ -dimension parameterized by  $\mu(x)$   $\Sigma(x, x')$   
 prior mean covariance (kernel)

stochastic process  $y(x)$  where for any  $x_1, \dots, x_n$  (and  $n$ )  
 marginal:  $(y(x_1), \dots, y(x_n))$  follows a Gaussian with  
 mean  $\begin{pmatrix} \mu(x_1) \\ \vdots \\ \mu(x_n) \end{pmatrix}$  and covariance  $\Sigma$  s.t.  $\Sigma_{ij} = k(x_i, x_j)$

special case of GP: Bayesian linear regression

$$\text{use } \Sigma(x, x') = \sigma_w^2 \varphi(x)^T \varphi(x') + \sigma_y^2$$

but more generally, square exponential kernel

$$\Sigma(x, x') = \frac{C}{2} \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$$

(note 1-dim. e)

uncertainty length scale

15h46

so Bayesian inference: suppose observed  $y_1, \dots, y_n$  (for  $x_1, \dots, x_n$ )

simple? condition in Gaussian model?

$$k_{ij} \triangleq \Sigma(x_i, x_j)$$

simple? condition in Gaussian model?

$$Y_1, Y_2, \dots, Y_n, Y(x) \sim N(0, (C_n + k^T C_n^{-1} k))$$

$k_i \triangleq \mathbb{E}(x_i | x)$   
 to have output at  $x$   
 $\Sigma(x, x)$  covars with output at  $x$

$$\text{get } Y(x) \mid Y_1, \dots, Y_n \sim N\left(0 + k^T C_n^{-1} \vec{y}_{in}, K - k^T C_n^{-1} k\right)$$

that's it

note that only need  
to compute once

## ④ applications to Bayesian classification

- GPC use  $p(y=1|x) = \sigma(f(x))$
- hyperparameter selection  $\rightarrow$  can maximize the Marginal Likelihood  
 (for GP regression  $\rightarrow$  closed form formula)  
 $\rightarrow$  model selection

demos: <https://distill.pub/2019/visual-exploration-gaussian-processes/>  
<https://www.tml.fi/gp/>

## Dirichlet process:

\* used to model  $N$ -matrix model in Bayesian model

Bayesian matrix model (finite)

$$z \sim \text{Mult}(\pi) \quad \pi \in \Delta_K$$

$$x|z \sim p(x|z) \quad [\text{e.g. } N(x|\mu_z, \Sigma_z)]$$

Bayesian  $\rightarrow$  put prior on  $\pi$   $[\text{Dirichlet}(\alpha_1, \dots, \alpha_K)]$

and  $z$   $[\text{say } G_z \stackrel{\text{iid}}{\sim} G_0]$

would like  $K \rightarrow \infty$  can do  $G_z \setminus \pi$  together using DP

## Dirichlet process

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

random prob.  
measure

$G_0$ : dist. on  $\Theta$

$\alpha$ : concentration parameter

stochastic process  
induced by measurable sets  
of  $\Theta$

for every partition of  $\Theta$  in  $A_1, \dots, A_n$

$$\text{then } (G(A_1), G(A_2), \dots, G(A_n)) \sim \text{Dir}(\alpha G_0(A_1), \dots, \alpha G_0(A_n))$$

## stick breaking construction

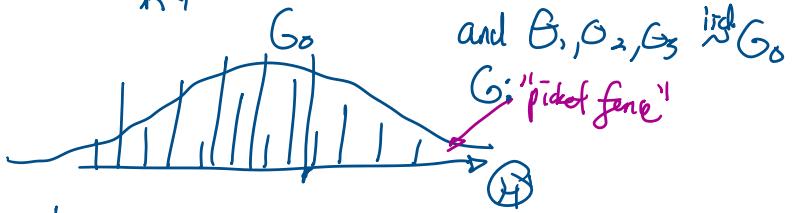
turns out that

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{g_k}$$

$\pi_k$ : stick length  
 $g_k$ : atom

where  $\pi = (\pi_1, \pi_2, \dots) \sim \text{GEM}(\alpha)$   
 and  $\pi_1, \pi_2, \dots \stackrel{\text{iid}}{\sim} G_\alpha$

turns out that  $\cup = \bigcup_{k=1}^K \pi_k G_0$  where  $\pi = (\pi_1, \pi_2, \dots) \sim \text{DEM}(\alpha)$



stick-breaking construction

$$\pi_1 \sim \text{Beta}_{\alpha}(1, \alpha)$$

$$\pi_2 = (1 - \pi_1) \nu_2$$

