

Bayesian Non-Parametrics

Plan:

- Bayesian Lin. Regression
 - Gaussian Process
 - Dirichlet Process / DP mixtures
 - de Finetti's Representation Theorem (dFRT).
-

Linear Regression:

$$y(n, \omega) = \omega_0 + \omega_1 n_1 + \dots + \omega_D n_D$$

$$y(n, \omega) = \omega_0 + \sum_{j=1}^{M-1} \omega_j \phi_j(x)$$

↳ basis
fn.

$$= \sum_{j=0}^{M-1} \omega_j \phi_j(x)$$

where $\phi_0(x) = 1$

$$= \omega^T \phi(x)$$

Gaussian basis fn:

$$\phi_j(x) = \exp \left\{ - \frac{(x - \mu_j)^2}{2s^2} \right\}$$

$$\text{MLE: } t = y(a, \omega) + \epsilon \quad \begin{matrix} \nearrow \text{zero mean} \\ \searrow \beta \text{ precision} \end{matrix}$$

$$\begin{aligned} P(t|x, \omega, \beta) &= N(t | y(a, \omega), \beta^{-1}) \\ &= \prod_{n=1}^N N(t_n | \omega^\top \phi(x_n), \beta^{-1}) \\ \omega_{ML} &= (\phi^\top \phi)^{-1} \phi^\top t. \end{aligned}$$

Bayesian Linear Regression:

β is known.

ω - Parameter to estimate

$$P(\omega) = N(\omega | m_0, S_0)$$

Posterior = Prior \times likelihood.

$$P(\omega | t) = N(\omega | m_n, S_n)$$

} ex.

$$m_N = \frac{S_N^{-1} (S_0^{-1} m_0 + \beta \phi^T t)}{S_N^{-1} + \beta \phi^T \phi}$$

$$S_N^{-1} = S_0^{-1} + \beta \phi^T \phi.$$

$$\omega_{MAP} = m_N$$

$$S_0 = \alpha^{-1} I \text{ with } \alpha \rightarrow 0 \quad m_N \rightarrow \omega_{ML}$$

$$\text{assume: } P(\omega | \alpha) = N(\omega | 0, \alpha^{-1} I)$$

$$m_N = \beta S_N \phi^T t$$

$$S_N^{-1} = \alpha I + \beta \phi^T \phi.$$

$$\ln P(\omega | t) = -\frac{\beta}{2} \sum_{n=1}^N \left\{ t_n - \omega^T \phi(x_n) \right\}^2$$

$$- \frac{\alpha}{2} \omega^T \omega + \text{const.}$$

$$P(t|E, \alpha, \beta) = \int P(t|\omega, \beta) P(\omega|E, \alpha, \beta) d\omega$$

$$P(E|x, t, \alpha, \beta) = N(t | m_n^T \phi(x), \sigma_n^2(x))$$

$$\sigma_n^2(x) = \frac{1}{\beta} + \phi(x)^T S_n \phi(x)$$

Kernel view

$$y(x, m_n) = m_n^T \phi(x)$$

$$= \beta \phi(x)^T S_n \phi^T E$$

$$= \sum_{n=1}^N \beta \phi(x)^T S_n \phi(x_n) t_n$$

$$y(x, m_n) = \sum_{n=1}^N k(x, x_n) t_n$$

where $k(x, x') = \beta \phi(x)^T S_n \phi(x')$

\hookrightarrow equivalent kernel /
Smoothing matrix.

$$\begin{aligned}\text{Cov}(y(x), y(x')) &= \text{Cov}[\phi(x)^T w, w^T \phi(x')] \\ &= \phi(x)^T S_N \phi(x') \\ &= \beta^{-1} k(x, x')\end{aligned}$$

Gaussian Processes:

functional view:

$$y(n) = w^T \phi(x)$$

$$p(w) = \mathcal{N}(w | 0, \alpha^{-1} I)$$

Every $w \sim p(w) \rightarrow$ fn. $y(n)$

distribution over fn. $y(n)$

$y(x)$

$y(n_1), \dots, y(n_n)$

$$x_1 \dots x_n \rightarrow (y(x_1), \dots, y(x_n))$$

$$y = \phi w$$

$$\mathbb{E}[y] = \phi \mathbb{E}[w] = 0$$

$$\text{Cov}[y] = \mathbb{E}[yy^T]$$

$$\begin{aligned} &= \phi \mathbb{E}[ww^T]\phi^T \\ &= \frac{1}{\alpha} \phi \phi^T = K \end{aligned}$$

↗ Gram matrix.

$$k_{nm} = k(x_n, x_m) = \frac{1}{\alpha} \phi(x_n)^T \phi(x_m)$$

$$k(x, x') = \exp(-\theta |x - x'|)$$

How to use this GP for regression?

$$t_n = y_n + \epsilon_n$$

$$P(t_n | y_n) = N(t_n | y_n, \beta^{-1})$$

$$P(t | y) = N(t | y, \beta^{-1} I_n)$$

$$P(y) = N(y | 0, k)$$

$$P(E) = \int P(E | y) P(y) dy$$

$$= N(t | 0, c) \quad \xrightarrow{\text{ex.}}$$



$$C(x_n, x_m) = k(x_n, x_m) + \beta^{-1} \delta_{nm}$$

$$P(t_{n+1}) = N(t_{n+1} | 0, C_{n+1})$$

$$C_{n+1} = \begin{pmatrix} C_n & k \\ k^T & c \end{pmatrix} \quad \xleftarrow{\text{ex.}}$$

where $k \rightarrow k(x_n, x_{n+1})$
 $c \rightarrow k(x_{n+1}, x_{n+1})$

$$P(t_{n+1} | t) = \mathcal{N}(m, \sigma^2)$$

$$m = k^T C_N^{-1} t$$

$$\sigma^2 = c - k^T C_N^{-1} k.$$

↓ exp.

One restriction:

Covariance matrix, \mathbf{k} should be
Positive definite

A-P represent

basis for library

invert $N \times N$ matrix.

$O(N^3)$

invert $M \times M$ matrix.

$O(M^3)$

one step

$M \ll N$

every step $O(N^2)$

$O(M^2)$

Main advantage of GP:

you can use inf. dimensional
base fn. \rightarrow kernel.

Parametric Model	Non-Parametric Model.
finite set of Param. θ	- infinite dimensional θ .
Given θ , Prediction is indep. of D .	$\theta \Rightarrow$ function.
$P(x \theta, D) = P(x \theta)$	\rightarrow amt of info. θ can capture grows as D grows.
- even if data is unbounded, Complexity of model is bounded.	
- Not flexible.	

Example :

1) fm approximation

Poly. regression Vs. C.P.

2) Classification

Log. regression Vs. GP classif.

3) Clustering.

mixture model, kMeans Vs.

DP mixtures .

4) time series .

HMM vs. infinite HMM .

de Finetti's Representation Theorem :

df-RT

Coin tossing:

I. I. D.

Independence:

$$P(x_{n+1} \mid x_1, \dots, x_n) = P(x_{n+1})$$

only with I.D., we can learn.

Exchangeability:-

$$\{0, 1, 0, 0, 1, 0\}$$

order does not matter.

$$P(F, F, F, G, G, G) >$$

$$P(G, G, G, F, F, F)$$

order matters.

exchangeable \rightarrow weaker than
indep

Defn: A set of R.V. $\{X_n\}$ is said to be exchangeable if given the joint density $p(x_1, \dots, x_n)$ we have

$$p(x_1, \dots, x_n) = P(x_{\sigma(1)}, \dots, x_{\sigma(n)})$$

σ - perm. of $(1 \dots n)$

IID \rightarrow exchangeable.

exchangeable $\not\rightarrow$ IID.

exchange \rightarrow ID.

~~d~~ def dFRT for Bernoulli R.V

Let $\{x_i\}_{i=1}^{\infty}$ be a seq. of
finitely exchange random variable.
i.e. $\forall n > 0$, each finite
Subsequence of $\{x_i\}_{i=1}^n$ is exchanged.

Then f a random variable θ
and a distrib. fn $F(\theta)$ such that

$$P\left(\lim_{n \rightarrow \infty} \frac{\bar{x}_i}{n} = \theta\right) = 1$$

with $\theta \sim F(\theta)$

and

$$P(x_1, \dots, x_n) = \int_0^1 \prod_i \theta^{x_i} (1-\theta)^{1-x_i} dF(\theta)$$