

Gaussian processes

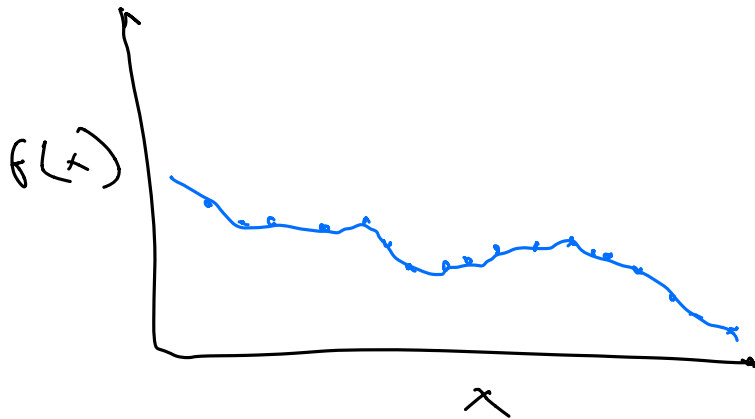
Idea: specify a prior $p(f)$

$$f = [f_1, f_2, f_3, \dots]$$

$$y_i = f(x_i) + \epsilon_i$$
$$\epsilon_i \sim N(0, \sigma^2)$$

GP:

any finite subset
of points $\{f_i\}$
are jointly Gaussian



Priors on functions

$$f_i = \omega^T \phi(x_i)$$

$$\omega \sim \mathcal{N}(0, \Sigma)$$

$$x_i \in \mathbb{R}^D$$

$$f_i \in \mathbb{R}$$

$$\phi: \mathbb{R}^D \rightarrow \mathbb{R}^H$$

$$\omega \in \mathbb{R}^H$$

for inputs x_1, \dots, x_N as matrices

$$f = \Phi \omega$$

$N \times 1$ $N \times H$ $H \times 1$

Q: what is $p(f)$?

$$E[f] = \Phi E[\omega] = 0$$

$$\begin{aligned} \text{Cov}(f) &= E[ff^T] - E[f]E[f]^T \\ &= \Phi E[\omega\omega^T] \Phi^T \\ &= \Phi \Sigma \Phi^T \equiv K \end{aligned}$$

Kernels
What is this $K = \Phi \Sigma \Phi^T = \text{Cov}(f)$?

$$K_{ij} = \underbrace{\text{Cov}(f(x_i), f(x_j))}_{\substack{\text{output} \\ \text{cov.}}} = \underbrace{\phi(x_i)^T \Sigma \phi(x_j)}_{k(x_i, x_j)}$$
$$f \sim N(0, K)$$

Kernel „trick“

Any $k(x, x')$ can be used as long as K is P.S.D.

e.g. $k(x, x') = e^{\frac{-\|x - x'\|_2^2}{2\ell^2}}$ } squared exponential kernel

$$\hookrightarrow \forall c, cKc \geq 0$$

Conditioning on data

Suppose we observe some x, f examples,
and we would like to predict at x^* .

$$\begin{bmatrix} f \\ f^* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \overset{C}{K(x, x)} & \overset{A^T}{K(x, x^*)} \\ \overset{A}{K(x^*, x)} & \overset{D}{K(x^*, x^*)} \end{bmatrix} \right)$$

$$f^* | x^*, f, x \sim \mathcal{N} \left(AC^{-1}f, D - AC^{-1}A^T \right) \quad \left. \vphantom{f^* | x^*, f, x} \right\} \text{Gaussian properties}$$

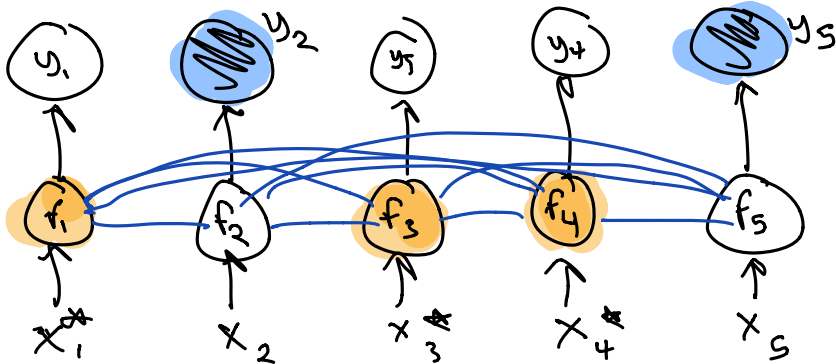
or y

Noisy data

Include a likelihood

$p(y_i | f_i)$, e.g. $y_i \sim \mathcal{N}(f_i, \sigma^2)$

$$\begin{bmatrix} y \\ f^* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \underbrace{K(x, x) + \sigma^2 I}_C & K(x, x^*) \\ K(x^*, x) & K(x^*, x^*) \end{bmatrix} \right)$$



In feedforward networks (Neal, 1996)

$$\begin{aligned} f(x_i) &= a + W^T g(b + V^T x_i) \\ &= a + \sum_{h=1}^H \omega_h \left[g(b + v_h^T x_i) \right] \end{aligned}$$

Annotations for the first equation:
- $a \in \mathbb{R}$ (scalar)
- $W \in \mathbb{R}^{H \times 1}$ (nonlinearity)
- $V \in \mathbb{R}^{D \times H}$ (nonlinearity)
- $x_i \in \mathbb{R}^D$ (nonlinearity)

Annotations for the second equation:
- ω_h (scalar)
- $v_h \in \mathbb{R}^D$ (scalar)

$$a \sim N(0, \sigma_a^2)$$

$$\omega_h \sim N(0, \sigma_\omega^2)$$

$$b \sim p(b)$$

$$v_h \stackrel{\text{iid}}{\sim} p(v)$$

Compute moments:

$$E[f(x)] = E[a] + \sum_h E[\omega_h] E[g(\dots)] = 0$$

$$\begin{aligned} E[f(x)f(x')] &= \sigma_a^2 + \sum_h \sigma_\omega^2 E[g(b + v_h^T x)g(b + v_h^T x')] \\ &= \sigma_a^2 + H \sigma_\omega^2 E[g(b + v^T x)g(b + v^T x')], v \stackrel{\text{iid}}{\sim} p(v) \end{aligned}$$

Annotations for the first equation:
- $E[f(x)f(x')] = \text{Cov}(x, x')$

Now let $p(w) = N(w | 0, \frac{s^2}{H})$.
sometimes analytic \rightarrow "arc-cosine kernel"

$$\text{Cov}(x, x') = \sigma_a^2 + s^2 E \left[g(b + v^T x) g(b + v^T x') \right]$$

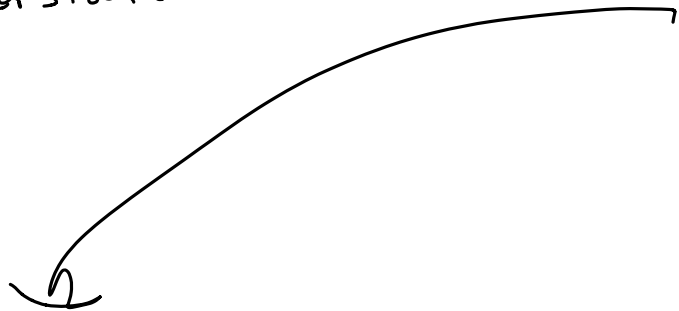
Let $g(\cdot)$ be bounded (e.g. sigmoid , \tanh) $\xrightarrow{p(b)p(v)}$ prior

CLT: as $H \rightarrow \infty$, this becomes Gaussian
(sample avg. iid. r.v.)

For small H : probably not a GP
large H : probably close

More layers?

Recursion:



$$\left[\begin{array}{l} z^{l-1} \sim GP \\ h^{l-1} = g(z^{l-1}) \\ z^l = b^l + \sum_j \omega_j^l h_j^{l-1} \end{array} \right]$$

$\omega_j^l \sim \mathcal{N}(0, \frac{\sigma_\omega^2}{H^l I})$

$$\Rightarrow E[z^l(x) z^l(x')^T] = \sigma_b^2 + \sigma_\omega^2 \int_{p(z^{l-1})} [g(z^{l-1}(x)) g(z^{l-1}(x'))^T]$$

Needs work to show it is actually a GP!

Simplified result :

Let $f_{\theta}(x)$ be a network, prior $p(\theta)$

Define :

$$m(x) = E_{p(\theta)} [f_{\theta}(x)] \quad (\text{often zero})$$

$$k(x, x') = E_{p(\theta)} [(f_{\theta}(x) - m(x)) (f_{\theta}(x') - m(x'))]$$

Now, use the GP $(m(\cdot), k(\cdot, \cdot))$!

Is this good?

- Deep learning "works" because of hierarchy of features.
- GPs w/ "complicated" covariance functions don't perform as well
- May be contradictory as much as useful!