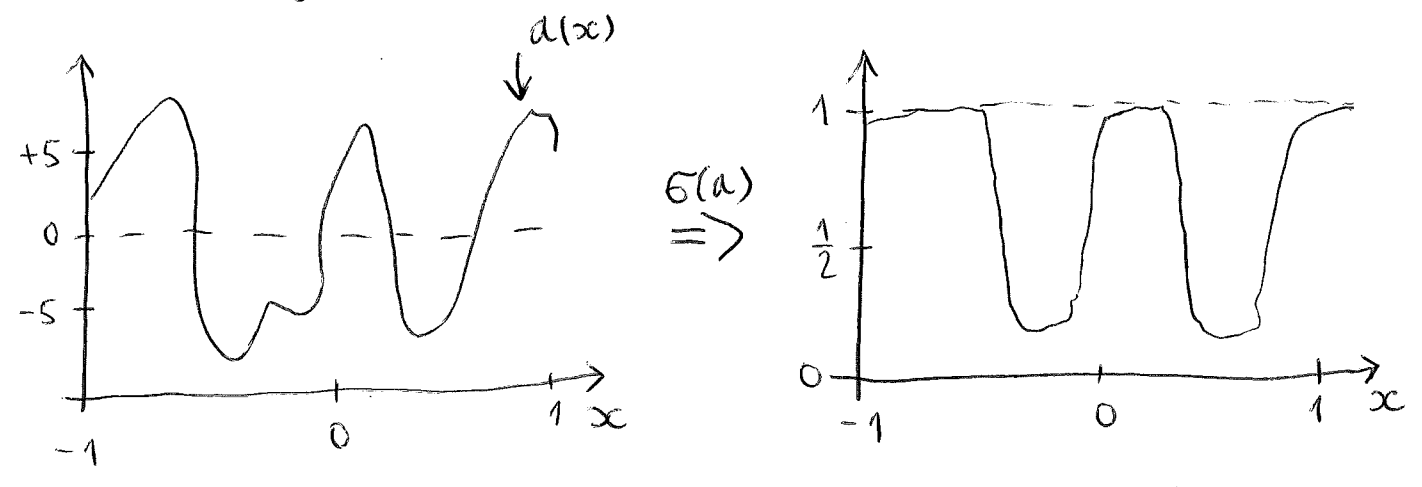


Gaussian processes for classification

Consider $k=2$, s.t. $t \in \{0, 1\}$.

Idea: define a gaussian process over $d(\vec{x})$ & then use $y = \sigma(a)$ to obtain a non-gaussian stochastic process over functions $y(\vec{x})$, s.t. $y(\vec{x}) \in (0, 1)$, $\forall x$.

For example, in the 1D case



Recall that $p(t|a) = \sigma(a)^t (1 - \sigma(a))^{1-t}$.

Training dataset: $\vec{x}_1 \dots \vec{x}_N \Rightarrow \vec{t} = t_1 \dots t_N$

also, $\vec{x}_{N+1} \Rightarrow t_{N+1}$

We need predictive distr'n $p(t_{N+1} | \vec{t})$

Define $\vec{a}_{N+1} = \overbrace{d(\vec{x}_1) \dots d(\vec{x}_{N+1})}$, then

$$p(\vec{a}_{N+1}) \stackrel{\text{gaussian process}}{=} \mathcal{N}(\vec{a}_{N+1} | \vec{0}, C_{N+1}), \text{ where } C_{N+1, nm} = C_{N+1}(\vec{x}_n, \vec{x}_m) = \underbrace{k(\vec{x}_n, \vec{x}_m)}_{\text{pos. semidefinite kernel f'n}} + \underbrace{\sum_{\text{small}} \delta_{nm}}_{\text{pos. const}}$$

In general, the kernel f'n $k(\vec{x}, \vec{x}')$ depends on hyperparameters $\vec{\theta}$.

Since $K=2$, we ^{can} focus on $p(t_{N+1}=1 | \vec{t}_N)$
 since $p(t_{N+1}=0 | \vec{t}_N) = 1 - p(t_{N+1}=1 | \vec{t}_N)$.

Then

$$p(t_{N+1}=1 | \vec{t}_N) = \int da_{N+1} \overbrace{p(t_{N+1}=1 | a_{N+1})}^{\sigma(a_{N+1})} \times \underbrace{p(a_{N+1} | \vec{t}_N)}_{\text{non-gaussian, need an approx'n}}$$

predictive distr'n

Use Jayplace:

$$p(a_{N+1} | \vec{t}_N) = \int d\vec{a}_N p(a_{N+1}, \vec{a}_N | \vec{t}_N) =$$

$$\stackrel{\uparrow}{=} \int d\vec{a}_N \frac{p(\vec{t}_N | a_{N+1}, \vec{a}_N) p(a_{N+1}, \vec{a}_N)}{p(\vec{t}_N)} =$$

Bayes theorem

$$= \frac{1}{p(\vec{t}_N)} \int d\vec{a}_N \underbrace{p(\vec{t}_N | \vec{a}_N)}_{\text{does not depend on } a_{N+1}} p(a_{N+1} | \vec{a}_N) p(\vec{a}_N) \stackrel{\downarrow \text{Bayes again}}{=} \frac{p(\vec{t}_N | \vec{a}_N) p(\vec{a}_N)}{p(\vec{t}_N)} = p(\vec{a}_N | \vec{t}_N)$$

$$= \int d\vec{a}_N \underbrace{p(a_{N+1} | \vec{a}_N)}_{\text{gaussian}} \underbrace{p(\vec{a}_N | \vec{t}_N)}_{\text{non-gaussian}}$$

Just as with regression,

$$p(a_{N+1} | \vec{a}_N) = \mathcal{N}(a_{N+1} | \vec{k}^T \mathbf{C}_N^{-1} \vec{a}_N, c - \vec{k}^T \mathbf{C}_N^{-1} \vec{k})$$

Now, use Laplace on $p(\vec{a}_N | \vec{t}_N)$:

go back to $p(\vec{a}_N | \vec{t}_N) = \frac{p(\vec{t}_N | \vec{a}_N) p(\vec{a}_N)}{p(\vec{t}_N)}$, where

$p(\vec{a}_N) = \mathcal{N}(\vec{a}_N | \vec{0}, C_N)$ and

$$p(\vec{t}_N | \vec{a}_N) = \prod_{n=1}^N \beta(a_n)^{t_n} (1 - \beta(a_n))^{1-t_n} = \prod_{n=1}^N e^{a_n t_n} \beta(-a_n)$$

$$\underbrace{\beta(a)^t (1 - \beta(a))^{1-t}}_{\frac{1}{(1+e^{-a})^t}} = \frac{e^{a(t-1)}}{1+e^{-a}} = \frac{e^{-a}}{1+e^{-a}} e^{at}$$

$$\left(\frac{e^{-a}}{1+e^{-a}} \right)^{1-t} \quad \beta(-a) = \frac{1}{1+e^a}$$

$$\text{Then } \log p(\vec{a}_N | \vec{t}_N) = \log p(\vec{a}_N) + \log p(\vec{t}_N | \vec{a}_N) - \log p(\vec{t}_N)$$

$$= -\frac{1}{2} \vec{a}_N^T C_N^{-1} \vec{a}_N - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |C_N| + \vec{t}_N^T \vec{a}_N - \sum_{n=1}^N \log(1+e^{a_n}) - \log p(\vec{t}_N)$$

" $\Psi(\vec{a}_N)$ "

$$\nabla_{\vec{a}} \log p(\vec{a}_N | \vec{t}_N) = \nabla_{\vec{a}} \Psi(\vec{a}_N) = \vec{t}_N - C_N^{-1} \vec{a}_N - \vec{\beta}_N$$

$$\frac{\partial}{\partial a_i} \sum_{n=1}^N \log(1+e^{a_n}) = \sum_n \frac{1}{1+e^{a_n}} e^{a_n} \delta_{ni} = \underbrace{\frac{e^{a_i}}{1+e^{a_i}}}_{\beta(a_i)}$$

Here, $\vec{\beta}_N = \underline{\beta(a_1) \dots \beta(a_N)}$

$$\nabla_{\vec{a}} \Psi(\vec{a}_N) = 0 \Rightarrow \underbrace{\vec{\beta}_N + C_N^{-1} \vec{a}_N}_{\text{nonlinear in } a_i} = \vec{t}_N$$

Need to find the solution by e.g.

Newton-Raphson method:

$$\nabla_{\vec{a}} \nabla_{\vec{a}} \Psi(\vec{a}_N) = -C_N^{-1} - W_N$$

$$\frac{\partial}{\partial a_j} \sigma(a_i) = \sigma(a_i) (1 - \sigma(a_i)) \delta_{ij}$$

$$W_N = \begin{pmatrix} \sigma(a_1)(1-\sigma(a_1)) & & 0 \\ & \ddots & \\ 0 & & \sigma(a_N)(1-\sigma(a_N)) \end{pmatrix}$$

↑
pos. def.

C_N pos. def. $\Rightarrow C_N^{-1}$ pos. def., yielding

Hessian $A = -\nabla_{\vec{a}} \nabla_{\vec{a}} \Psi = W_N + C_N^{-1}$ pos. def.

That means that $\log p(\vec{a}_N | \vec{t}_N)$ is concave everywhere \rightarrow there is a single global maximum, so NR should work very well.

Specifically,

$$\vec{a}_N^{\text{new}} = \vec{a}_N^{\text{old}} + A^{-1} \underbrace{\nabla_{\vec{a}} \log p(\vec{a}_N | \vec{t}_N)}_{\nabla_{\vec{a}} \Psi(\vec{a}_N)} =$$

$$= \vec{a}_N^{\text{old}} + \underbrace{(W_N + C_N^{-1})^{-1}}_{\text{when } C_N C_N^{-1} (W_N + C_N^{-1})^{-1} = C_N (\mathbb{I} + W_N C_N)^{-1}} [\vec{t}_N - \vec{b}_N - C_N^{-1} \vec{a}_N^{\text{old}}] \ominus$$

$$\ominus C_N (\mathbb{I} + W_N C_N)^{-1} [\vec{t}_N - \vec{b}_N + (\mathbb{I} + W_N C_N) C_N^{-1} \vec{a}_N^{\text{old}} - C_N^{-1} \vec{a}_N^{\text{old}}] \ominus$$

$$\textcircled{E} C_N (\mathbb{I} + W_N C_N)^{-1} [\vec{t}_N - \vec{b}_N + W_N \vec{a}_N^{\text{old}}].$$

Iterate to convergence \Rightarrow find \vec{a}_N^* s.t.

$$\nabla_{\vec{a}} \Psi(\vec{a}_N) \Big|_{\vec{a}_N^*} = \vec{0} \Rightarrow \underbrace{\vec{a}_N^* = C_N (\vec{t}_N - \vec{b}_N)}_{\text{non-linear eq'n in } \vec{a}_N^*}$$

Now, $H = -\nabla_{\vec{a}} \nabla_{\vec{a}} \Psi(\vec{a}_N) \Big|_{\vec{a}_N^*} = W_N \Big|_{\vec{a}_N^*} + C_N^{-1}$

Finally, Laplace:

$$p(\vec{a}_N | \vec{t}_N) \Rightarrow q(\vec{a}_N) = \mathcal{N}(\vec{a}_N | \vec{a}_N^*, H^{-1}).$$

$$\text{So, } p(a_{N+1} | \vec{t}_N) = \int d\vec{a}_N \mathcal{N}(a_{N+1} | \vec{k}^T C_N^{-1} \vec{a}_N, c - \vec{k}^T C_N^{-1} \vec{k}) \times \mathcal{N}(\vec{a}_N | \vec{a}_N^*, H^{-1})$$

$$(2.115): p(\vec{y}) = \int d\vec{x} p(\vec{y} | \vec{x}) p(\vec{x}) =$$

$$= \int d\vec{x} \mathcal{N}(\vec{y} | A\vec{x} + \vec{b}, L^{-1}) \mathcal{N}(\vec{x} | \vec{\mu}, \Lambda^{-1}) =$$

$$= \mathcal{N}(\vec{y} | A\vec{\mu} + \vec{b}, L^{-1} + A\Lambda^{-1}A^T)$$

$$\begin{cases} \vec{y} \rightarrow a_{N+1} \\ \vec{x} \rightarrow \vec{a}_N \end{cases} \Rightarrow \begin{cases} A\vec{x} + \vec{b} \rightarrow \vec{k}^T C_N^{-1} \vec{a}_N \left[\vec{b} \rightarrow 0 \right], L^{-1} \rightarrow c - \vec{k}^T C_N^{-1} \vec{k} \\ \vec{\mu} = \vec{a}_N^*, \Lambda^{-1} = H^{-1} \end{cases}$$

gives vector of length N

$$p(a_{N+1} | \vec{t}_N) = \mathcal{N}(a_{N+1} | \vec{k}^T C_N^{-1} \vec{a}_N^*, c - \vec{k}^T C_N^{-1} \vec{k} + \vec{k}^T C_N^{-1} H^{-1} C_N^{-1} \vec{k})$$

$(C_N^{-1})^T = C_N^{-1} \rightarrow C_N^{-1} \vec{k}$ III

$$\bar{k}^T C_N^{-1} \bar{a}_N^* = \bar{k}^T (\bar{t}_N - \bar{b}_N),$$

$$C_N (\bar{t}_N - \bar{b}_N)$$

$$C - \bar{k}^T [C_N^{-1} - C_N^{-1} (W_N + C_N^{-1})^{-1} C_N^{-1}] \bar{k} =$$

~~$$C - \bar{k}^T [C_N^{-1} - C_N^{-1} (W_N + C_N^{-1})^{-1} C_N^{-1}] \bar{k} =$$~~

$$= C - \bar{k}^T [C_N^{-1} - C_N^{-1} (W_N C_N + \mathbb{I})^{-1}] \bar{k} =$$

$$= C - \bar{k}^T [(C_N^{-1} \cdot (W_N C_N + \mathbb{I}) - C_N^{-1}) (W_N C_N + \mathbb{I})^{-1}] \bar{k} =$$

$$= C - \bar{k}^T [W_N (W_N C_N + \mathbb{I})^{-1}] \bar{k} = C - \bar{k}^T (C_N + W_N^{-1})^{-1} \bar{k}$$

$$\square \mathcal{N}(a_{N+1} | \underbrace{\bar{k}^T (\bar{t}_N - \bar{b}_N)}_{\mu}, \underbrace{C - \bar{k}^T (C_N + W_N^{-1})^{-1} \bar{k}}_{\sigma^2})$$

Finally,

$$p(t_{N+1}=1 | \bar{t}_N) = \int da_{N+1} \delta(a_{N+1}) \mathcal{N}(a_{N+1} | \mu, \sigma^2) \approx$$

$$\approx \delta(\kappa(\sigma^2) \mu), \text{ where}$$

$$\kappa(\sigma^2) = \frac{1}{\sqrt{1 + \frac{\kappa \sigma^2}{8}}}$$

(4.153)

————— 0 —————

What about hyperprms $\vec{\theta}$?

Maximize $p(\vec{t}_N | \vec{\theta}) \Rightarrow \vec{\theta}_{ML}$

Use $p(\vec{t}_N | \vec{\theta}) = \int d\vec{a}_N p(\vec{t}_N | \vec{a}_N) p(\vec{a}_N | \vec{\theta})$,

apply Laplace: $\int d\vec{z} f(\vec{z}) \approx f(\vec{z}_0) \frac{(2\pi)^{M/2}}{|A|^{1/2}}$

\swarrow M dim

$$\log p(\vec{t}_N | \vec{\theta}) \approx \underbrace{\log p(\vec{t}_N | \vec{a}_N^*) + \log p(\vec{a}_N^* | \vec{\theta})}_{\Psi(\vec{a}_N^*)} + \text{depends on } \vec{\theta}$$

$$+ \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |W_N + C_N^{-1}| \text{ depends on } \vec{\theta}$$

To maximize this, we need to find

$$\frac{\partial \log p(\vec{t}_N | \vec{\theta})}{\partial \theta_j} \text{ in terms of } \frac{\partial C_N}{\partial \theta_j} \text{ \& } \frac{\partial \vec{a}_N^*}{\partial \theta_j}$$

Note that $\begin{cases} C_N = C_N(\vec{\theta}), \\ \vec{a}_N^* = \vec{a}_N(\vec{\theta}). \end{cases}$

$$\text{Recall that } \begin{cases} p(\vec{t}_N | \vec{a}_N^*) = \prod_{n=1}^N e^{a_n^* t_n} \sigma(-a_n^*), \\ p(\vec{a}_N^* | \vec{\theta}) = \mathcal{N}(\vec{a}_N^* | \vec{0}, C_N) \end{cases}$$

Finally, note that

$$\frac{\partial \vec{a}_N^*}{\partial \theta_j} = \frac{\partial C_N}{\partial \theta_j} (\vec{t}_N - \vec{b}_N) - C_N \underbrace{\frac{\partial \vec{b}_N}{\partial \theta_j}}_{\text{vector with elements}}$$

$$\frac{\partial b(a_n^*)}{\partial \theta_j} = \underbrace{b(a_n^*) (1 - b(a_n^*))}_{W_N, nn} \frac{\partial a_n^*}{\partial \theta_j}$$

$$\frac{\partial \vec{b}_N}{\partial \theta_j} = W_N \frac{\partial \vec{a}_N^*}{\partial \theta_j}$$

$$\frac{\partial \vec{a}_N^*}{\partial \theta_j} = \frac{\partial C_N}{\partial \theta_j} (\vec{t}_N - \vec{b}_N), \text{ or}$$

$$(\mathbb{I} + C_N W_N)$$

$$\frac{\partial \vec{a}_N^*}{\partial \theta_j} = (\mathbb{I} + \underbrace{W_N C_N}_{C_N W_N})^{-1} \frac{\partial C_N}{\partial \theta_j} (\vec{t}_N - \vec{b}_N)$$

so $\frac{\partial \vec{a}_N^*}{\partial \theta_j}$ can be expressed through $\frac{\partial C_N}{\partial \theta_j}$

Differentiating C_N wrt θ_j is relatively straightforward as $c(\vec{x}_n, \vec{x}_m) = k(\vec{x}_n, \vec{x}_m) + \delta_{nm}$ depends on $\vec{\theta}$

with $\frac{\partial \log p(\vec{t}_N | \vec{\theta})}{\partial \theta_j}$ available, can use e.g. conjugate gradient to maximize $\log p(\vec{t}_N | \vec{\theta})$ & find $\vec{\theta}_{ML}$.