

## Gaussian processes

The problem:

$N$  data points  $\mathcal{X}_N, t_N = \left\{ \mathbf{x}_{(n)}, t_n \right\}_{n=1}^N$

↑  
vector of  
vectors

can be vectors

$t$  = targets

assume that a function  $y(\mathbf{x})$  underlies the observed data.

What is  $t_{N+1}$  at point  $\mathbf{x}_{(N+1)}$ ?

### Parametric approach

Express  $y(\mathbf{x})$  in terms of  $y(\mathbf{x}; w)$ , model with parameters  $w$

E.g. using a set of basis functions

$\{\phi_h(\mathbf{x})\}_{h=1}^H$  we have

$$y(\mathbf{x}; w) = \sum_{h=1}^H w_h \phi_h(\mathbf{x})$$

[Note: These can be non-linear functions such as

$$\phi_h(\mathbf{x}) = \exp \left[ -\frac{(\mathbf{x} - c_h)^2}{2r^2} \right]$$

(aka "radial basis functions" in this context ... Gaussians)

Still, since  $y(\mathbf{x}; w)$  depends linearly on  $w$ , this is a [linear model]

[W25] 2.

We have seen this problem before

$$p(w | t_N, \mathbf{x}_N) = \frac{p(t_N | w, \mathbf{x}_N) p(w)}{p(t_N | \mathbf{x}_N)}$$

here expressed slightly differently

$t_N$  = target values given  $\mathbf{x}_N$  positions  
do not include 'I' in these notes.

Having performed inference to obtain the posterior, a prediction can be made

$$p(t_{N+1} | \overset{\text{as}}{t_N}, \mathbf{x}_{N+1}) = d^H w p(t_{N+1} | w, \mathbf{x}_{N+1}) p(w | t_N, \mathbf{x}_N)$$

[Note : e.g. from R samples  $w^{(r)}$  from  $p(w | t_N, \mathbf{x}_N)$ ]

$$p(t_{N+1} | t_N, \mathbf{x}_{N+1}) \approx \frac{1}{R} \sum_{r=1}^R p(t_{N+1} | w^{(r)}, \mathbf{x}_{N+1})$$

[Note that the final result does not make explicit reference to our representation (parametric one) of the unknown function  $y(x)$ .]

Note on notation:

W2b / 3.

From now on:  $X_N$  = vector of  $N$  vectors

$t_N$  = vector of  $N$  scalars

$X^{(N+1)}$  = vector #  $N+1$

$t^{(N+1)}$  = target #  $N+1$

We will see that for models with fixed basis functions and Gaussian prior distributions (w. zero mean) on the parameters, the joint probability of all the observed data given the model,  $p(t_N | X_N)$  is a multivariate Gaussian with mean zero and with a covariance function determined by the basis functions; this implies that the conditional distribution  $p(t^{(N+1)} | t_N, X_{N+1})$  is also a Gaussian distribution whose mean depends linearly on  $t_N$ .

Consider a linear model

[W2b] 4

$$R_{nh} \equiv \phi_h(x^{(n)}) , \text{ i.e. } R \text{ is } N \times H$$

$\uparrow$   
number  
of points  
 $\{x^{(n)}\}_{n=1}^N$

$\uparrow$   
# of  
basis  
functions

$y_n$  is the vector of values  $y(x; w)$

$$y^{(n)} = \sum_h R_{nh} w_h \quad [y_n = R w]$$

Assume a Gaussian prior distribution for  $w$

$$p(w) = \mathcal{N}(w; 0, \Sigma_w^{-2} I)$$

since  $y$  is a linear function of  $w$ ,

it is also Gaussian distributed with mean zero. Its covariance matrix

$$Q = \langle yy^T \rangle = \langle R w w^T R^T \rangle = \Sigma_w^{-2} R R^T$$

$$\Rightarrow p(y) = \mathcal{N}(y; 0, \Sigma_w^{-2} R R^T)$$

This will be true for any selected points  $X_N$ ; which is the defining property of a Gaussian process

"The probability distribution of a function  $y(x)$  is a Gaussian process if for any selection of points:  $x_1, \dots, x_N$  the density  $p(y(x_1), \dots, y(x_N))$  is a Gaussian"

What about the target  
values?

W25 / 5

If  $t(n)$  is assumed to differ  
by additive Gaussian noise  
of variance  $\sigma_v^2$  from  $y(n)$ ,  
then  $t$  also has a Gaussian  
prior distribution

$$p(t) = \mathcal{N}(t; 0, Q + \sigma_v^2 I)$$

$$\equiv C = \sigma_w^2 R R^T + \sigma_v^2 I$$

What does  $Q$  and  $C$  look like?

$$Q_{nn'} = \sigma_w^2 \sum_h \phi_h(x^{(n)}) \phi_h(x^{(n')})$$

and

$$C_{nn'} = Q_{nn'} + \delta_{nn'} \sigma_v^2$$

i.e. the correlation between  
target values  $t^{(n)}$  and  $t^{(n')}$   
is determined by the points  
 $x^{(n)}, x^{(n')}$  and the behaviour  
of  $\phi$ .

In fact, we don't  
need the basis functions  
and parameters any more.

1x125 | 6

Instead of a prior distribution  
on functions  $p(y(x))$  in  
terms of basis functions and  
priors on parameters, we can  
summarize the prior simply  
by a covariance function

$$C(x, x')$$

generating a covariance matrix

$$Q_{nn'} = C(x^{(n)}, x^{(n')})$$

for any set of points  $\{x^{(n)}\}_{n=1}^N$

$Q$  must be non-negative-definite

$\Rightarrow$  constraint on valid cov. functions,

The cov. matrix for the target values

$$C_{nn'} = C(x^{(n)}, x^{(n')}) + \sigma^2 \delta_{nn'}$$

$$\Rightarrow p(t) = N(t; \theta, C) = \frac{1}{Z} e^{-\frac{1}{2} t^T C^{-1} t}$$

## Covariance functions

Given the constraint on the previous page we can construct a wide range of covariance functions. We can introduce hyperparameters  $\Theta$  so that

$$C_{mn} = C(x^{(m)}, x^{(n)}; \Theta) + \delta_{mn} N(x^{(n)}; \Theta)$$

↑  
noise model

Stationary cov. functions are translationally invariant

$$C(x, x'; \Theta) = D(x - x'; \Theta)$$

Also known as a kernel

A popular choice (known as radial basis functions; RBF, in this context — but really a Gaussian)

$$C(x, x'; \Theta) = \Theta_1 \exp\left[-\frac{1}{2} \sum_{i=1}^I \frac{(x_i - x'_i)^2}{r_i^2}\right] + \Theta_2$$

with hyperparameters  $\Theta = (\Theta_1, \Theta_2, \{r_i\})$

Determine  $\Theta$  by

• optimization (common)  
or marginalize over them.

↑  
or just one "r"