# Recap from last time

- · Genssian processes for regression
- · Conceptually: work with probability distr. on function space
- o Start:  $p(f_*, \overline{f} \mid X, \overline{X}_*) \leftarrow \frac{\text{Joint prior for all}}{\text{velevant function values}}$ (training + test)
- Goal:  $p(f_*|\bar{f},X,\bar{x})$  = Posterior pried for  $f_* = f(\bar{x}_*)$  given the now known training data
- a Start from Gaussian prior -> end with Gaussian posterior

· Know Mx, 6,2 analytically:

$$\int_{A_*} = w(\bar{x}_*) + k(\bar{x}_*, \chi) \sum_{k=1}^{\infty} (\bar{f} - w(\chi))$$

$$G_* = k(\bar{x}, \bar{x}_*) - k(\bar{x}_*, \chi) \sum_{k=1}^{\infty} k(\chi, \bar{x}_*)$$
our prediction (

- o  $M(\cdot)$  and  $k(\cdot,\cdot)$  are the mean and covariance functions we chose to define the mean vector and covariance matrix of over prior  $p(f_*, f \mid X, X_*)$   $(\Sigma = k(X,X))$
- · Note: So far we have assumed fully specified un(.) and (el.,.)

## Choosing and optimising covariance function (kernel)

- · By choosing kernel we define our prior on function space.
- o Encodes our expectations about the true function, but without specifying and assumed functional form for the true  $f(\bar{x})$
- · Kernels encode expectations about e.g.
  - smoothness
  - typical length scales
  - periodicity
  - trend (increase /decrease)
  - Symmetry

etc.

o This is the main modelling step in GP regression!

Choice of mean function  $m(\bar{x})$  usually much less important. Often set to  $m(\bar{x}) = 0$  or  $m(\bar{x}) = \text{constant}$ , But been in mind that predictions are "pulled" forwards  $m(\bar{x})$  in regions of  $\bar{x}$  space for away from known points

- o Need a fully specified prior -> fully specified kernel o Common approach :
  - 1) (onstruct a kernel with some free params. (hyperparameters)
  - 2) fit these hyperparameters in a wax. likelihood fit using training data (the training step in GPR)

· Note that this is a bit "un-Bayesian".

Proper Bayesian approach: Assign priors for the hyperpanameters
and marginalise over them

D: hyperparameters

GP posterior = 
$$p(f_* | \bar{f}, X, \bar{x}_*) = \int p(f_*, \bar{\theta} | \bar{f}, X, \bar{x}_*) d\bar{\theta}$$
  
=  $\int p(f_* | \bar{\theta}, \bar{f}, X, \bar{x}_*) p(\bar{\theta} | \bar{f}, \bar{X}) d\bar{\theta}$   
 $\propto \int p(f_* | \bar{\theta}, \bar{f}, X, \bar{x}_*) p(\bar{f} | \bar{\theta}, X) p(\bar{\theta}) d\bar{\theta}$ 

- · But this integration is usually very expensive
- o So rommon approach is to instead use a point est for  $\overline{\theta}$ , womely the  $\overline{\theta}$  value that maximises the likelihood

$$L(\bar{e}) = P(\bar{f} | \bar{e}, X) = N(m(X), \Sigma(\bar{e}))$$

- o Effectively like "peeking" at the data and setting a delta function prior  $p(\bar{b}) = \delta(\bar{b} \bar{\theta}_{ML})$ . ("Empirical Bayes")
- · Usually maximise the log-likelihood

$$\ln L(\bar{b}) = -\frac{1}{2} (\bar{f} - w(x))^T \sum_{(\bar{b})}^{-1} (\bar{f} - w(x)) - \frac{1}{2} \ln |\Sigma(\bar{b})| - \frac{N}{2} \ln (2\pi)$$

- o Note: For every  $\bar{\theta}$  we try we need a new matrix inverse  $\Sigma(\bar{\theta})$  and determinant  $|\Sigma(\bar{\theta})|$ 
  - Standard techniques for watrix inversion: O(43)
  - Can use Cholesky decomp. technique instead : O(42)

Still very slow when u grows large [See algo. 2.1 in Rusmussen and williams

Some kernels

Browse scikit-learn documentation to see Icenuel examples.

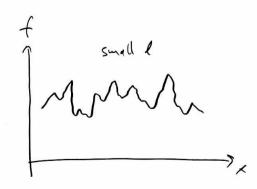
1) The squared-exponential kernel

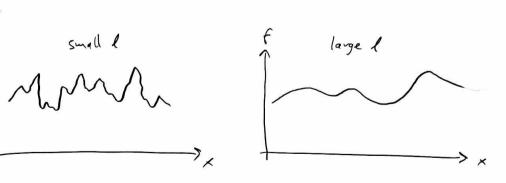
$$\langle (\bar{x}, \bar{x}') = G_f^2 e^{-\frac{1}{2}(\bar{x} - \bar{x}')^2}$$

Hyperpars.:
$$\bar{\theta} = \left\{ \mathcal{C}_{f}^{2}, 1 \right\}$$

As the Euclidean dist. between the input points X, X' increases, the covariance between  $f(\bar{x})$  and  $f(\bar{x}')$  decreases exponentially.

- o Universal Iceruel (can approx any cont. function given enough data)
- · Prefers very smooth functions
- · 60° : Scale factor, sets average dist away from the mean function
- · I : Length scale, sets how quickly correlation between f-values drop as x-dist increase. Sets the typical kength scale for "wiggles" in the function





o Example of stationary kernel, i.e. only depends on distance between x and x1, not the values of x and x' in an assolute sense. (The kernel acts the same across all of x space). A non-stationary learnel depends on the specific locations of x and x' in x space.

# 2) The Mattern Kernel

$$k(\bar{x},\bar{x}') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left[ \sqrt{2\nu'} \frac{|\bar{x}-\bar{x}'|}{\ell} \right]^{\nu} k_{\nu} \left[ \sqrt{2\nu'} \frac{|\bar{x}-\bar{x}'|}{\ell} \right]$$

· Common choices: 
$$\nu = \frac{1}{2}, \frac{3}{2}, \frac{5}{2}$$

#### o Other common laevuels

Browse the Kernel Cookbook for examples.

- Linear kernel
- Periodic Levuel
- Rational quadratic learnel (effectively infinite sum of squared-exp. kenels by different length scales)
- Noise kemel
- o Kernels can be summed and multiplied to construct new kernels

$$k(\bar{x},\bar{x}') = k_1(\bar{x},\bar{x}') k_2(\bar{x},\bar{x}')$$
: "AND operator" (Both k, and k, and

$$k(\bar{x},\bar{x}') = k_1(\bar{x},\bar{x}') + k_2(\bar{x},\bar{x}')$$
: "OR operator"

o (on use different kernel components for different input components, e.g. with  $\overline{x} = [x_1, x_2]$ 

$$k(\bar{x},\bar{x}') = e^{-\frac{1}{2}(\underbrace{x_1-x_1'})^2} e^{-\frac{1}{2}(\underbrace{x_2-x_2'})^2}$$

to allow differen lengthscales in different directions in x space. (But end up with more hyperpars, to determine.)

#### Noisy data

· So far we've focused on the case of voise-free data,
i.e. training data were values of the true f(x) directly.

$$\Rightarrow$$
 Got posterior  $p(f_* | \bar{f}, X, \bar{x}_*)$ 

- In this case  $p(f_*|\bar{f},X,\bar{x}_*) \to S(f_*-f_i)$  when  $\bar{x}_* \to \bar{x}_i$ , i.e. posterior collapses to deltafunction when  $\bar{x}_*$  is a a known point
- o Reasonable in theory (given assumption of training points with no uncent.), but in practice numerically problematic
- o Need to allow for uncertainty in data, either because this is actually the rase, or just for numerical stability

$$Y_i \equiv Y(\bar{x}_i) = f(\bar{x}_i) + \mathcal{E}_i$$

( Noise,  $E \sim \mathcal{N}(0, \sigma_E^2)$ 

( noise level

Note: We say "noise", but there doesn't have to be randomness involved.  $N(0, \sigma_{\epsilon}^2)$  rould just express our degree of certainty in the training values

· Now we must distinguish between

1) 
$$P(f_*|\overline{Y},X,\overline{X}_*)$$
: Degree of belief in underlying true function value  $f_*$  at  $\overline{X}_*$ 

o If we only add the noise variance to the training set part of the covariance matrix, we get case 7:

$$\Rightarrow P(f^* | \underline{\lambda}, \underline{\lambda}, \underline{x}^*) = V(\lambda^*, \underline{e}_s)$$

where
$$M_* = M(\overline{x}_*) + k(\overline{x}_*, X) \left[ \Sigma + G_{\varepsilon}^2 \overline{I} \right] (\overline{y} - M(X))$$

$$G_{\varepsilon}^2 = k(\overline{x}_*, \overline{x}_*) - k(\overline{x}_*, X) \left[ \Sigma + G_{\varepsilon}^2 \overline{I} \right]^2 k(X, \overline{x}_*)$$

o If we add noise variance to both 
$$\Sigma$$
 and  $k(\bar{x}_*,\bar{x}_*)$  we get case  $Z$ :

$$\Rightarrow P(y_* | \overline{y}, X, \overline{x}_*) = N(\mu_*, 6_*^2)$$

where

$$M_* = w(\bar{x}_*) + k(\bar{x}_*/X) \left[ \Sigma + G_{\varepsilon}^2 \bar{I} \right] \left( \bar{y} - w(X) \right)$$

$$G_{\kappa}^2 = k(\bar{x}_*, \bar{x}_*) + G_{\varepsilon}^2 - k(\bar{x}_*X) \left[ \Sigma + G_{\varepsilon}^2 \bar{I} \right] k(X, \bar{x}_*)$$

$$\int_{\text{Note}}^{\sqrt{2}} |x|^{-1} dx$$

- Common approach: Try to learn the typical noise level from data by treating  $6\pi^2$  as a hyperparameter, i.e. add a kernel term of the form  $k(\overline{x},\overline{x}') = \delta_{\overline{x}\overline{x}'} \delta_{\varepsilon}^2$ 
  - o Pitfall: In hyperparameter space, there will they often be a huge region towards large of where the model can obtain OK fit to data simply by "assuming" that all variation in the data is due to noise.

[ See documentation of sciket-leary for example ]

# Some final points

A fully connected, infinitely wide neural network with independent and identically distributed (i.i.d.)

network parameters a Gaussian process !

[Resically a conseq. of the central Limit Theorem.]

· Underest of uncertainty

· Regularisation of covariance matrix

Even in problems where there is zero data noise, the GP results can be better by introducing a small "mugget" / noise term along diagonal in rover matrix.

Makes it easier to compute inverse E- without bus of numerical precision.

- Keep an eye on the condition number  $K(\Xi) = \frac{\lambda_{max}}{\lambda_{min}}$ When  $K \to \infty$ ,  $\Xi \to \text{singular (not invertible)}$ 

- Typically gets worse with many (and very close) training points.

# How to deal with the large-data limit?

- o Basic problem: competation of 5" too expensive when N gets very large (>0(10")]
- o Many approaches, most are based on splitting up data / making & more sparse. Some examples:
  - Distributed Gaussian Processes, arxiv: 1502.02843
  - Generalited Robest Bayesian Committee Machine, arxiv: 1806.00720
  - Sparse Gaussian Processes (replace large training set with a smaller set of "inducing points")
  - Deep Kernel Leaving, arxiv: 1511.02222 (use large dataset to train a neural network which can act as keemel function for a GP)

+ much more !