

# Large Linear Gaussian Processes: Fast Multi-output GP Model Learning

## Vladimir Feinberg Li-Fang Cheng Barbara Engelhardt Kai Li

# Why multi-output GPs?

- 1) Create confidence intervals
- 2) Assign likelihoods
- 3) Generative sampling
- 4) Non-parametric modeling
- 5) Learn correlations between several outputs

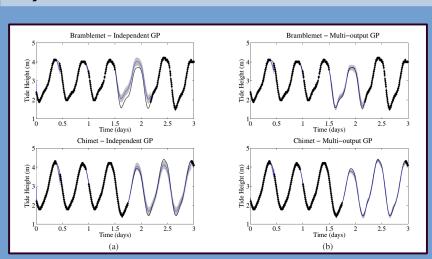


Figure 3 from [Osborne 2008

sensors Bramblemet and

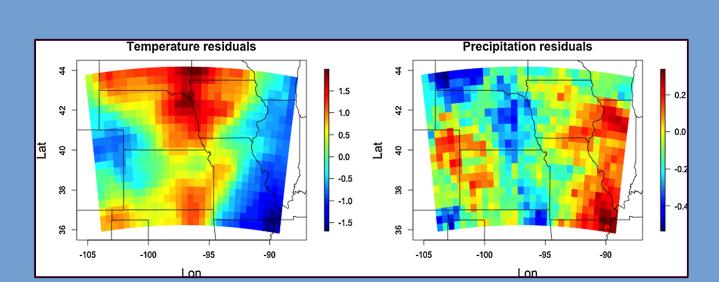
Chimet on the south coast of

England. The solid black line represents the actual held-out

data, the dots represent the

training data. The blue line is the expectation from the GP model and the grey shaded region is the 2 standard

deviation error bar



### Multi-output GP Approaches

#### LMC Kernel

*n* input-output pairs

D outputs total

\* What is  $\kappa_2^{1/2}$  in

dataset) below.

practice? Typically

O(n), see (financial

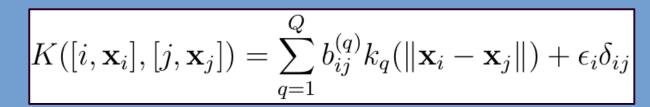
Q stationary kernels

Order-D matrices B<sub>g</sub> of average rank R determine *linear cross-correlation* 

between all pairs of our D outputs

Use approximation of m < n rank

An observation is an input [d, x], both output tag d (from 1 to D), and continuous, low dimensional x along with a corresponding real observation y.



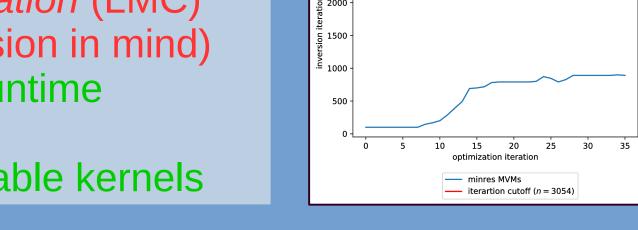
sensors on 4 points on the coast of England. Theoretical Complexity Comparison ( $\nabla$ L) Note there are ~QRD parameters in an LMC model. Log factors hidden.

Up-front Cost	Per-parameter
n <sup>3</sup>	n²
QRm <sup>3</sup>	nm
(n + min(QR, D <sup>2</sup> ) m) $\kappa_2^{1/2}$	n
	n³ QRm³

## Goals and Scope

Only apply to *linearly coregionalization* (LMC) Focus on time series (keep extension in mind) Asymptotically improve learning runtime Realize guarantees in practice

Apply to a wide range of differentiable kernels



#### How?

Use technique of Structured Kernel Interpolation from [Wilson et al 2015]

- Interpolate true input to a grid of m points for each  $D^2$  output-output pair

- Use stationarity in stationary subkernels  $k_{\alpha}(||x_{i+1} - x_{i+1}||) = k_{\alpha}(||x_i + \Delta - x_i - \Delta||) = k(||x_i - x_i||)$ Redundancy above gives rise to Toepliz matrix K<sub>g</sub> on the grid.

Describe cross-covariances on interpolating grid U with Kronecker product! Structured matrices → Fast Matrix Vector Products → Fast Inversions → Fast ∇L Use MINRES and [Gibbs and Mackay 1996].

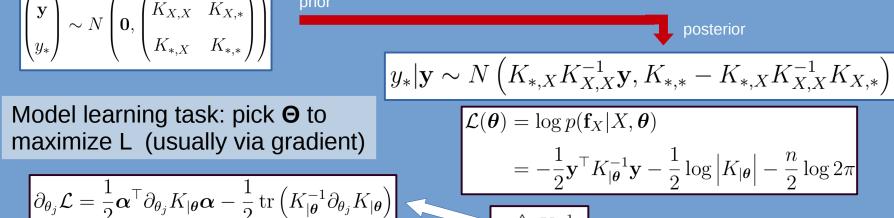
## GP Background

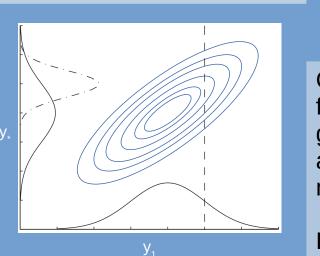
input data points X induces a covariance matrix  $K_{xx}$  whose ij-th entry is  $k(x_i, x_i)$ . Then the regressor is a random variable determined by a normal with that covariance.  $|\mathbf{y} \sim N(\mathbf{0}, K_{X,X})|$ Use marginalization to predict distribution at query input  $x_*$ .



GP Definition [Williams and Rasmussen 1996]

Model parameters  $\Theta$  parameterize the kernel k. The kernel k evaluated at

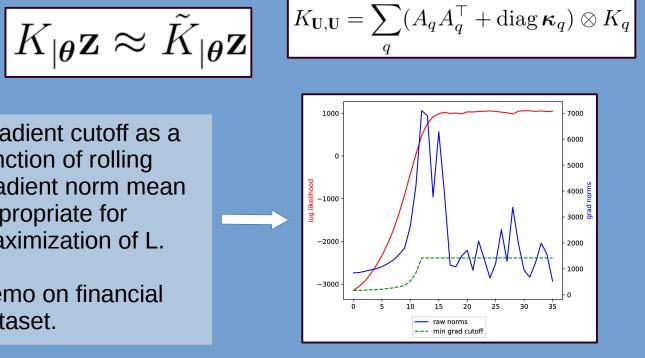




Demo of GP property from [Osborne 2008]

Gradient cutoff as a function of rolling gradient norm mean appropriate for maximization of L.

Demo on financial dataset.



the set of inputs, tagged by the output that they correspond to, and limited for mplicity to the one-dimensional case. Then  $\mathbf{y} \subset \mathbb{R}$  is the corresponding set of utputs. AdaDeltaUpdate is the update function for AdaDelta—we abstract way auxiliary AdaDelta variables.

Weather dataset D = 4, n = 16K

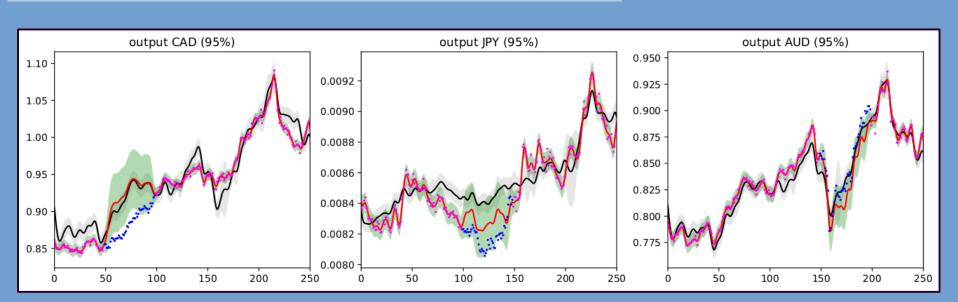
Impute missing periods from temperature

1: $\mathbf{procedure} \ \mathrm{LLGP}(\alpha, \mathbf{X}, \mathbf{y})$				
2: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}_0$ $\triangleright$ Initialization				
$g_{\max} = -\infty$				
4: repeat				
5: Construct a linear operator $\tilde{K}_{ \theta}$ from $\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}$ . $\triangleright$ MVM operator				
6: $g \leftarrow \nabla_{\theta} \mathcal{L}_{\tilde{K}_{ \theta}}$ $\triangleright$ Gradients from an operator, Algorithm 2				
7: $\boldsymbol{\theta} \leftarrow \text{AdaDeltaUpdate}(\boldsymbol{\theta}, g)$				
8: $g_{\text{max}} = \max(\ g\ , \ g_{\text{max}}\ )$				
9: $\mathbf{until} \ g\  \le \alpha g_{\max}$ $\triangleright \text{Cutoff}$				
10: return $ heta$				
11: end procedure				

#### Results

Financial dataset D = 13, n = 3K, m = 238 = n / DIndices for foreign exchange and commodities in 2007 Predict missing points (blue) from given ones (magenta) LLGP (mean in red, confidence in green) COGP (mean in black, confidence in grey)

METRIC	LLGP	COGP
SECONDS	64 (8)	296 (2)
SMSE	0.21 (0.01)	0.26 (0.03)
NLPD	-3.62 (0.07)	14.52 (3.10)



METRIC	$\begin{array}{c} \text{LLGP} \\ m = 500 \end{array}$	$\begin{array}{c} \text{LLGP} \\ m = 1000 \end{array}$	COGP
SECONDS SMSE	<b>60</b> ( <b>14</b> ) 0.09 (0.01)	259 (62) 0.09 (0.01)	1380 (12) <b>0.08</b> ( <b>0.00</b> )
NLPD	2.14 (0.58)	$1.54 \ (0.03)$	98.48 (1.30)

#### Conclusions

Beat state-of-the art in time and accuracy by using a scalable approximation (only quadratic growth in upfront cost per gradient)

Tested on a variety of kernels for accuracy (not shown)

Future work: rigorous iteration cutoff, extend to multidimensional input.

```
Algorithm 2 Given a linear MVM operator for a kernel, we compute the gradient of
ts data log likelihood. N_t is a fixed parameter for our stochastic trace approximation.
The output of this procedure is \nabla \mathcal{L}. MINRES(K, \mathbf{z}) computes K^{-1}\mathbf{z}.
  : procedure Gradient(K, \mathbf{y})
2: R \leftarrow \{\mathbf{r}_i\}_{i=1}^{N_t}, sampling \mathbf{r} \sim \text{Unif}\{\pm 1\}.
3: for z in \{y\} \cup R, in parallel do
           Store the result of MINRES(K, \mathbf{z}) as K^{-1}\mathbf{z}.
       end for
        Let \alpha = K^{-1}\mathbf{y}; computed before.
                                                                                                  \triangleright Compute \partial_{\theta_i} \mathcal{L}
             Create the operator L = \partial_{\theta_i} K_{|\theta}
             \{\mathbf{s}_i\}_{i=1}^{N_t} \leftarrow \{L(\mathbf{r}_i)\}_{i=1}^{N_t}
             r \leftarrow \boldsymbol{\alpha} \cdot L(\boldsymbol{\alpha})
.6: end procedure
```