

## INTRODUCTION

- Gaussian Processes (GP) are useful to solve non-linear regression problems, but they are limited to unimodal Gaussian output distributions, stationary functions and i.i.d noise scenarios.

### Objective

- Build a probabilistic model for general non-linear regression problems to deal with:
  - arbitrary output distribution (including multimodality)
  - non-stationary functions
  - heteroscedastic noise

### Paper Contributions

- A general model for non-linear regression: Infinite Mixture of Global GPs (IMoGGP). Novel interpretation as a single-p Dependent Dirichlet Process.
- An easy-to-implement MCMC sampling algorithm.
- Comparative results against Infinite Mixture of Experts (IMoE).

## MODEL

We want to estimate  $y \in \mathbb{R}$  given an input  $\mathbf{x} \in \mathbb{R}^D$  and a database  $\mathcal{D}_n = \{\mathbf{x}_i, y_i\}_{i=1}^n$ , that is

$$p(y|\mathbf{x}, \mathcal{D}_n). \quad (1)$$

Our model is based on the stick-breaking construction of a Dirichlet Process (DP):

$$\boldsymbol{\pi}|\alpha \sim \text{GEM}(\alpha) \quad (2)$$

$$z_i|\boldsymbol{\pi} \sim \text{Multinomial}(\boldsymbol{\pi}) \quad (3)$$

$$\theta_m|H \sim H \quad (4)$$

$$y_i|z_i, \{\theta_m\} \sim F(\theta_{z_i}), \quad (5)$$

where GEM stands for the stick-breaking prior by Griffiths, Engen and McCloskey,  $\alpha$  is the concentration parameter of the DP,  $z_i$  indicates the cluster assignment,  $\theta_m$  designates the cluster parameters,  $H$  is a base measure, and  $F(\cdot)$  is the likelihood function, typically Gaussian.

In the regression setting, each  $y_i$  is associated with an input  $\mathbf{x}_i$  and we can directly modify (5) as

$$y_i|z_i, \{\theta_m\}, \mathbf{x}_i \sim F(\theta_{z_i}(\mathbf{x}_i)), \quad (6)$$

$$\theta_m|H, \phi_m \sim H_{\phi_m}, \quad (7)$$

where we assume that  $F(\theta_{z_i}(\mathbf{x}_i))$  is Gaussian-distributed with mean  $\mu_{z_i}(\mathbf{x}_i)$  and variance  $\sigma_{z_i}^2(\mathbf{x}_i)$ , and  $H_{\phi_m}$  is a Gaussian process prior with hyperparameters  $\phi_m$ .

- Now each cluster parameter  $\theta_m$  corresponds to a latent function over the input space.
- We can interpret the model as a Single-p Dependent Dirichlet Process whose atoms are GP functions.

## PROPERTIES

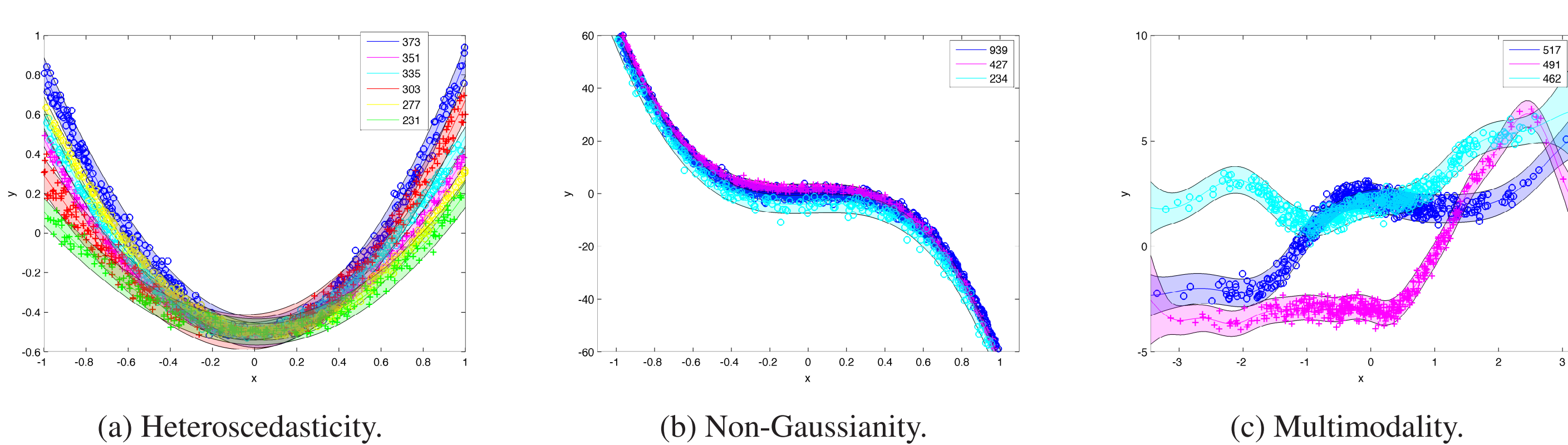


Figure 2: **Properties that can be captured by the IMoGGP model:** (a) non-stationary, heteroscedastic noise; (b) non-Gaussian likelihoods, specifically a Student's t with Gamma distributed noise; and, (c) multimodal predictive distributions.

## RESULTS

		(a)	(b)	(c)	(d)	(e)	(f)
PLLH	sGP	-0.0217	-3.4920	-3.3030	-0.5855	-1.6373	0.2033
	IMoE	0.7017	-2.1248	-2.1604	1.9452	-1.6308	<b>0.9943</b>
	IMoGGP	<b>0.9008</b>	<b>-2.1237</b>	<b>-1.2575</b>	<b>2.3587</b>	<b>-1.5723</b>	0.9846
MSE	sGP	0.0288	<b>4.8115</b>	4.2815	93.4640	0.7877	86.6815
	IMoE	0.0331	4.8394	5.2263	93.4640	0.7780	82.8929
	IMoGGP	<b>0.0287</b>	4.8500	<b>4.2703</b>	<b>43.6710</b>	<b>0.7754</b>	<b>82.4264</b>

Table 1: **Comparison of the single GP (sGP), the Infinite Mixture of Experts (IMoE) and the proposed Infinite Mixture of Global Gaussian Processes (IMoGGP).** The three first columns correspond to synthetic toy examples showed above: (a) Heteroscedasticity, (b) Non-Gaussianity, (c) Multimodality. The last three columns correspond to real databases available online: (d) Concrete, (e) Marathon, (f) RSSI.

## CONCEPTUAL COMPARISON OF METHODS

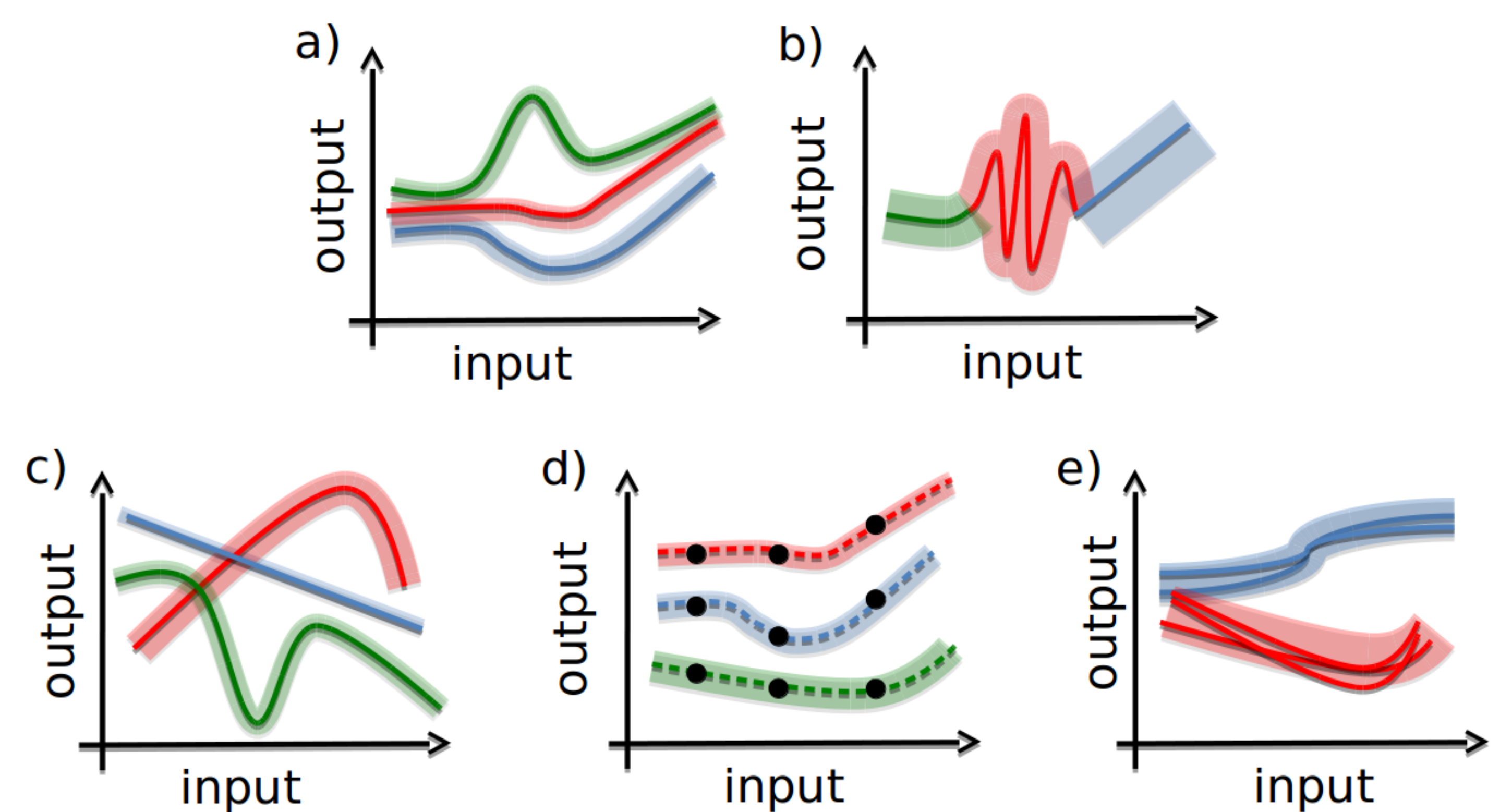


Figure 1: **Conceptual comparison of different approaches.** Sketch comparing a) Infinite Mixture of Global GPs (proposed approach), b) Infinite Mixture of Experts, c) Overlapping GPs for multi-tracking, d) Spatial Dirichlet Process, and e) time series clustering. Each color represents a different GP.

## INFERENCE FOR THE IMoGGP

**Algorithm 1** For each Gibbs sampling iteration:

- 1: Sample extended vector of mixture proportions:

$$\boldsymbol{\pi}|\mathbf{z}, \alpha \sim \text{Dirichlet}(n_1, \dots, n_K, \underbrace{\alpha/T \dots \alpha/T}_{T \text{ times}}) \quad (8)$$

- 2: Sample latent functions, i.e., cluster parameters  $\theta_m, m = 1, \dots, M^+$ :

$$p(\theta_m|\boldsymbol{\pi}, \mathbf{y}, \mathbf{X}, \mathbf{z}) \propto p(\theta_m|H_{\phi_m})p(\mathbf{y}|\mathbf{X}, \mathbf{z}, \theta_m) \quad (9)$$

- 3: Sample cluster assignments:

$$p(z_i|\boldsymbol{\pi}, y_i, \mathbf{x}_i, \mathbf{z}_{-i}, \{\theta_m\}) \propto p(z_i|\boldsymbol{\pi}) p(y_i|\mathbf{x}_i, \mathbf{z}, \{\theta_m\}) \quad (10)$$

- 4: Sample hyperparameters  $\phi_m, m = 1, \dots, M^+$ :

$$p(\phi_m|\boldsymbol{\pi}, \mathbf{y}, \mathbf{X}, \mathbf{z}) \propto p(\phi_m|H)p(\mathbf{y}|\mathbf{X}, \mathbf{z}, \phi_m) \quad (11)$$

- 5: Sample concentration parameter  $\alpha$  for the mixture model

- This algorithm is simple and computationally efficient, as it divides data into multiple GPs (we have smaller matrices to invert).

## FUTURE WORK

- Study sensibility to hyperparameters.
- Relax constant weights assumption (similar to Kernel-based Stick breaking process).
- Extend to higher dimensional problems (selection mechanism of input-dimension).

## FUNDING

This project has received funding from the European Union's Seventh Framework Programme for research, technological development and demonstration under grant agreement no. 316861.