



# Adaptive Multiple Importance Sampling for Gaussian Processes: Finding Difference Makers in Personality Impressions

Xiaoyu Xiong<sup>1</sup>, Maurizio Filippone<sup>2</sup>, Alessandro Vinciarelli<sup>1</sup>

1-The University of Glasgow email: x.xiong.1@research.gla.ac.uk, Alessandro.Vinciarelli@glasgow.ac.uk  
2- EURECOM, Sophia Antipolis, France email: maurizio.filippone@eurecom.fr



## Motivation

- ✓ Traditional parametric approaches for Social Signal Processing (SSP):  
falls short in taking into account uncertainty due to **model misspecification** or **overfitting**
- ✓ Current solutions: Bayesian treatment with nonparametric models like Gaussian Processes (GPs)
- ✓ GPs requires accurately characterizing posterior over covariance parameters:  
this is normally done by means of Markov chain Monte Carlo (MCMC) methods
- ✓ However, MCMC for learning GP parameters are inefficient due to its rejection of expensive proposals
- ✓ In this work, we propose an alternative inference framework for GPs based on Adaptive Multiple Importance Sampling (AMIS).
- ✓ The experimental results suggests AMIS is competitive with MCMC for GP models and suitable for SSP

## Bayesian Gaussian processes

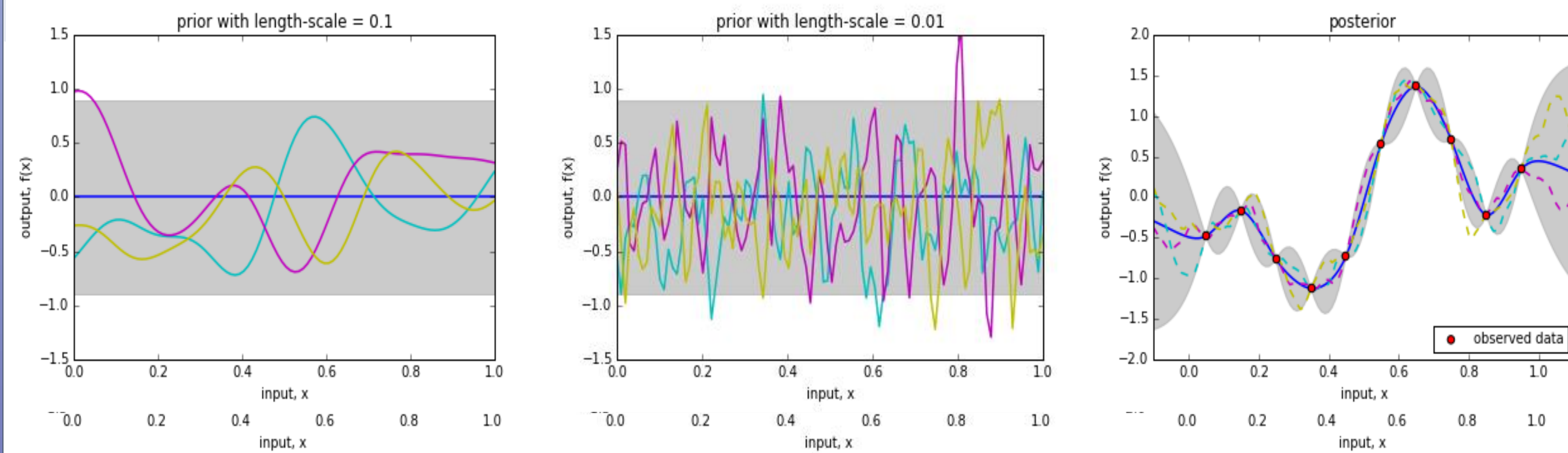
Gaussian process assumption  $p(f(X)|\theta) \sim \mathcal{N}(0, K(\theta))$

- ✓ Gaussian likelihood  $p(y|\theta) \sim p(y|f)p(f|\theta) \sim \mathcal{N}(0, K(\theta) + \lambda I)$  where  $\lambda$  is the variance of the Gaussian noise on  $y$ ,  $I$  is the identity matrix

Posterior  $\pi(\theta) := p(\theta|y, X) \propto p(y|\theta)p(\theta)$

- ✓ Non-Gaussian likelihood  $p(y|\theta) = \int p(y|f)p(f|\theta)df$ ,  $\tilde{p}(y|\theta) \approx \frac{1}{N_{imp}} \sum_{i=1}^{N_{imp}} \frac{p(y|f_i)p(f_i|\theta)}{q(f_i|\theta, y)}$  (estimated)

Posterior  $\pi(\theta) := p(\theta|y, X) \propto \tilde{p}(y|\theta)p(\theta)$



The entries of covariance K is determined by a kernel function

RBF kernel:  $k(x_i, x_j) = \sigma e^{-\frac{1}{2\tau^2} \|x_i - x_j\|^2}$ , where

$\sigma$  is the marginal variance of the function values at input locations  $x$ ,  $\tau$  is the length-scale, controlling the smoothness of functions

ARD kernel:  $k(x_i, x_j) = \sigma e^{-\sum_{r=1}^d \frac{1}{\tau_r^2} [x_{i(r)} - x_{j(r)}]^2}$ ,  $\tau_r$  is the length-scale of each feature of  $x$

**Inference in GPs is about learning the kernel parameters, e.g.,  $\theta = (\sigma, \tau)$  (RBF) or  $\theta = (\sigma, \tau_r)$  (ARD)**

- ✓ Computing expectation using MCMC

$$E_{\pi(\theta)} h(\theta) = \int h(\theta) \pi(\theta) d\theta = \frac{1}{N} \sum_{i=0}^N h(\theta^i)$$

- ✓ Computing expectation using Importance Sampling (IS)

$$E_{\pi(\theta)} h(\theta) = \int h(\theta) \frac{\pi(\theta)}{q(\theta)} q(\theta) d\theta = \frac{1}{N} \sum_{i=0}^N h(\theta^i) w_i, w_i = \frac{\pi(\theta^i)}{q(\theta^i)}$$

(Importance weight)

## MCMC for learning $\theta$

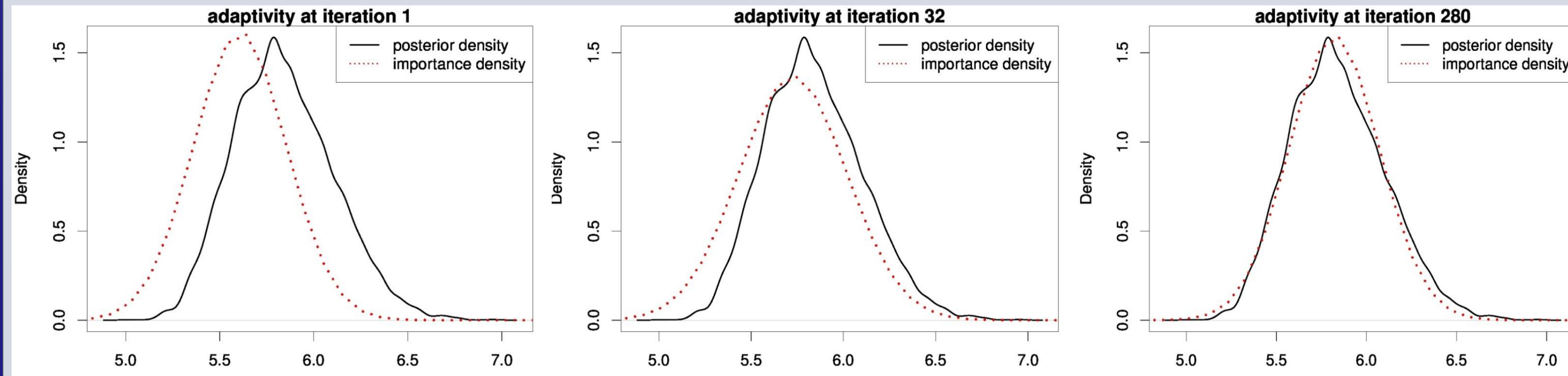
Hastings ratio:

$$\frac{p(y|\theta')p(\theta')}{p(y|\theta)p(\theta)} \text{ (MH for Gaussian likelihood)} \rightarrow \frac{\tilde{p}(y|\theta')p(\theta')}{\tilde{p}(y|\theta)p(\theta)} \text{ (Pseudo-Marginal MH for non-Gaussian likelihood)}$$

## Motivation for alternative sampling methods:

- ✓ In GPs with Gaussian likelihood, computing the marginal likelihood and its gradient with respect to  $\theta$  is expensive and standard MCMC algorithms reject proposals leading to a waste of computations
- ✓ In GPs with non-Gaussian likelihood, PM-MH may further cause inefficiencies due to large overestimation of the marginal likelihood

## Adaptivity of AMIS



- ✓ AMIS adaptively constructs an approximate posterior over  $\theta$  used to build an increasingly more accurate importance estimator:

$$\frac{1}{\sum_{t=0}^{T-1} \sum_{i=1}^{N_t} w_i^t} \sum_{t=0}^{T-1} \sum_{i=1}^{N_t} w_i^t h(\theta_i^t) \quad \text{where}$$

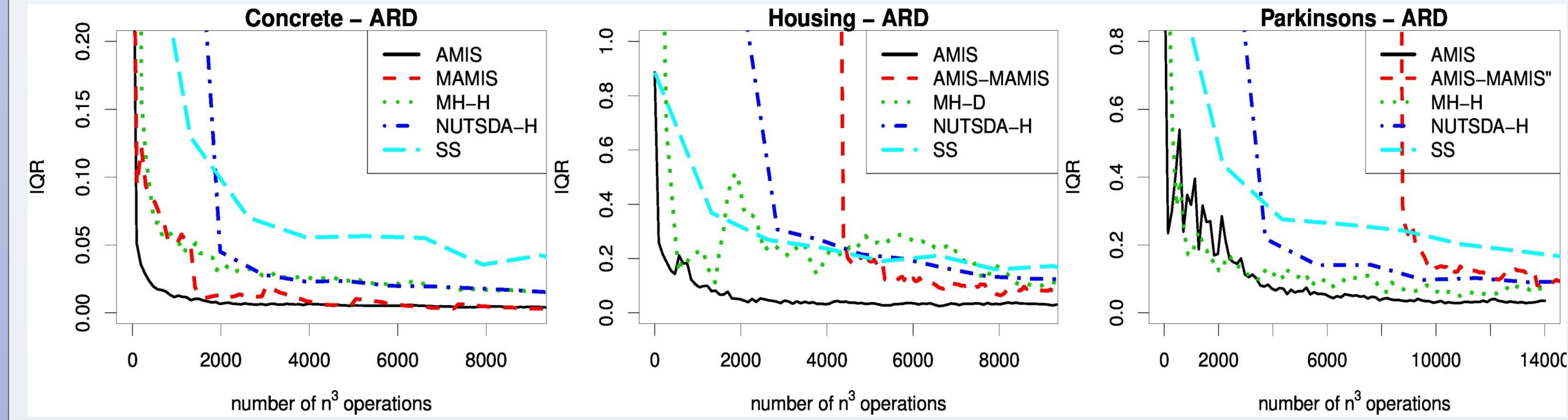
$$w_i^t = \frac{g(\theta^i)}{\sum_{t=0}^{T-1} \frac{1}{N_t} \sum_{i=1}^{N_t} N_t q_t(\theta_i^t; \gamma_t)}, g(\theta) = p(y|\theta)p(\theta) \quad \text{weight of AMIS for Gaussian likelihood}$$

$$w_i^t = \frac{\tilde{g}(\theta^i)}{\sum_{t=0}^{T-1} \frac{1}{N_t} \sum_{i=1}^{N_t} N_t q_t(\theta_i^t; \gamma_t)}, \tilde{g}(\theta) = \tilde{p}(y|\theta)p(\theta) \quad \text{weight of PM-AMIS for non-Gaussian likelihood}$$

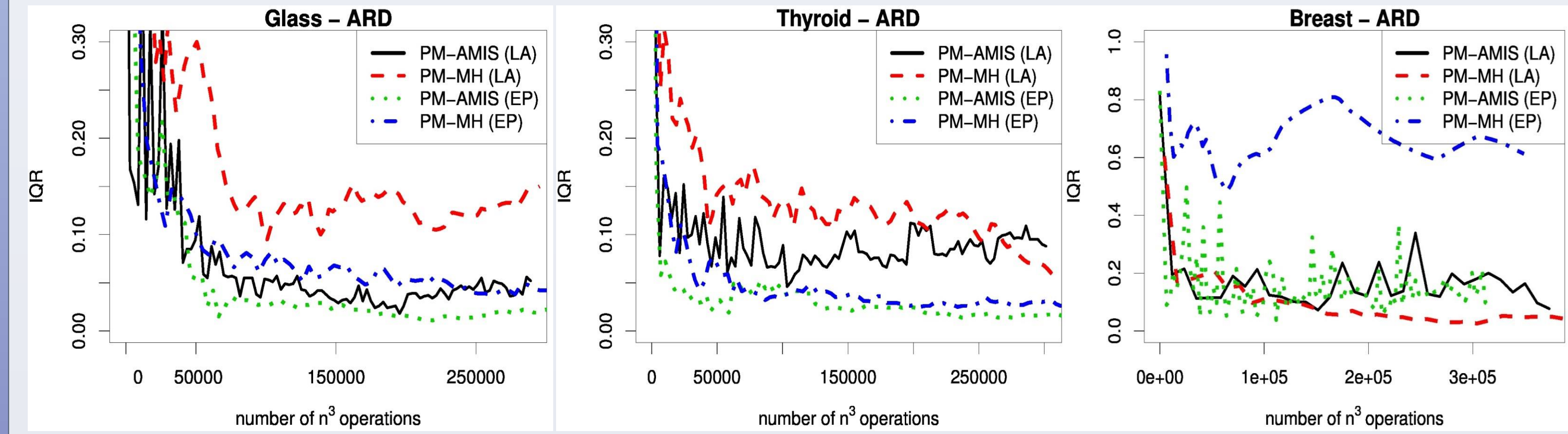
## Experiments

- ✓ IQR: Interquartile Range of the expectation of the norm of the parameters
- ✓ Convergence analysis: IQR against computational costs (number of  $n^3$  operations)

### GP regression - AMIS vs MH, Best of HMC family, slice sampling (SS)



### GP classification - PM-AMIS vs PM-MH



## Application of PM-AMIS for GPs: Personality inference from Flickr images

- ✓ **Data:** "fave", i.e., pictures that Flickr users have tagged as favourite, with personality traits attributed by Asian and UK assessors
- ✓ **Task:** predicting whether a Flickr user is perceived to be above median with respect to the Big-Five traits – Openness (Ope), Conscientiousness (Con), Extraversion (Ext), Agreeableness (Agr), Neuroticism (Neu)
- ✓ **Approaches:** Support Vector Machine (SVM) and PM-AMIS for GPs with G-ARD kernel:

$$k(x_i, x_j) = \sigma e^{-\sum_{r=1}^{N_g} \frac{1}{N_r \tau_r^2} \{\sum_{s \in G_r} [x_{i(s)} - x_{j(s)}]^2\}}$$

- $\tau_r$  is the length-scale parameter for group  $r$ ,  $N_r$  is the number of features in group  $r$ ,  $N_g$  is the number of groups,  $G_r$  is the set of indexes of the features that belong to group  $r$

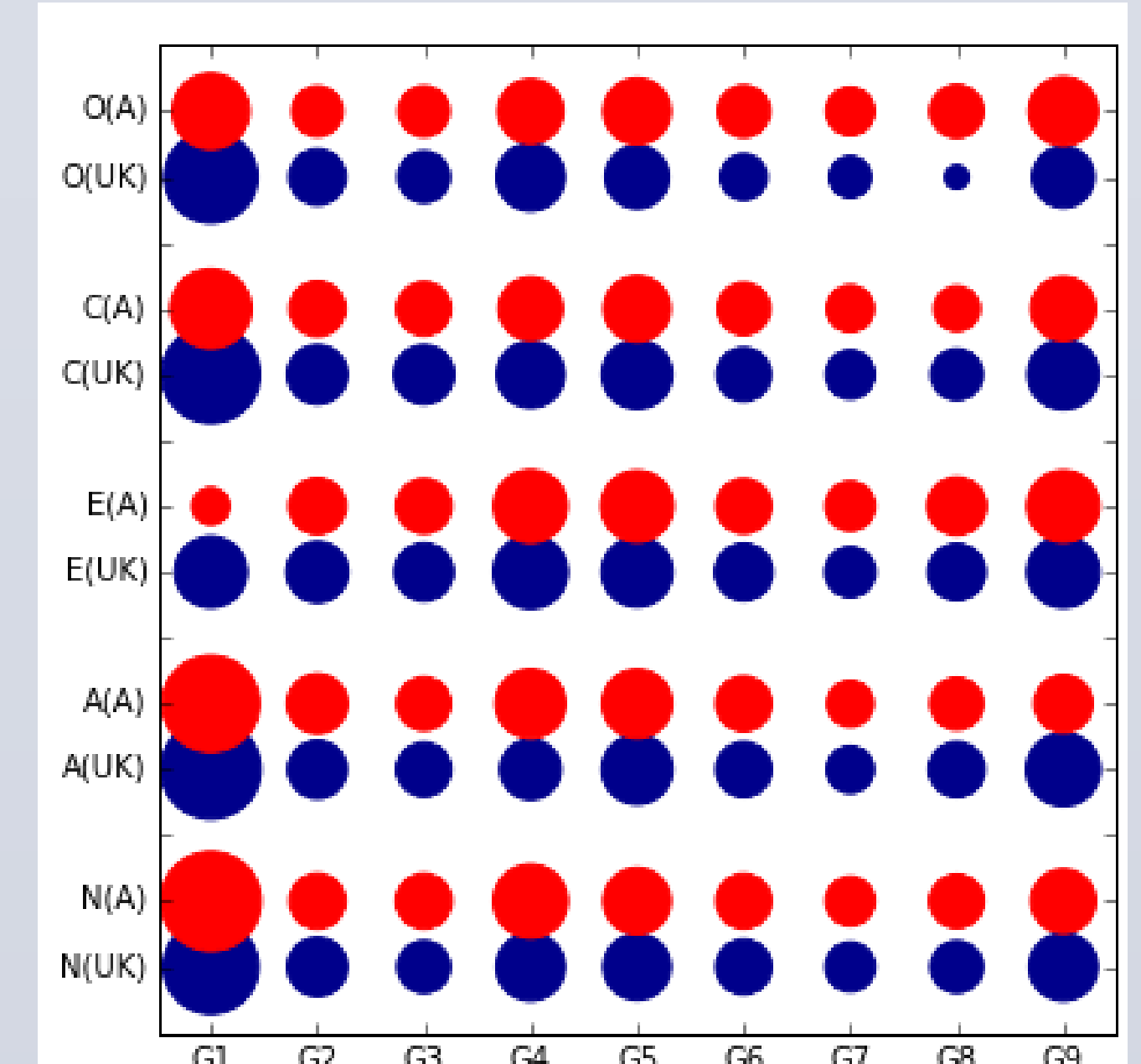
**Table 1: Groups of features (number in brackets denotes the number of features for each group)**

| G1           | G2                      | G3                        | G4                         | G5                 | G6                          | G7                      | G8           | G9                           |
|--------------|-------------------------|---------------------------|----------------------------|--------------------|-----------------------------|-------------------------|--------------|------------------------------|
| Faces<br>(1) | Colour<br>Properties(3) | Colour<br>Distribution(1) | Homogeneous<br>Regions (4) | Composition<br>(2) | Texture<br>Wavelets<br>(12) | GIST<br>filters<br>(24) | GLCM<br>(24) | Texture<br>Statistics<br>(4) |

- ✓ **Results:**

- PM-AMIS for GPs achieves comparable accuracies with SVM (see Table 2)
- G-ARD is able to identify the groups of features (G1, G5, G9, G4) that mostly influence personality impression
- Weights difference (Figure 1) across Asian and UK personality assessors suggests cultural difference on personality perception

**Figure 1: Coefficients of the G-ARD for the five traits (O, C, E, A, N) and two cultures Asian (A) and UK**



**Table 2: Prediction Accuracy**

|                        | Ope | Con | Ext | Agr | Neu |
|------------------------|-----|-----|-----|-----|-----|
| PM-AMIS for GPs (UK)   | 65% | 58% | 71% | 73% | 79% |
| SVM(UK)                | 59% | 62% | 71% | 74% | 77% |
| PM-AMIS for GPs (Asia) | 68% | 52% | 74% | 68% | 69% |
| SVM(Asia)              | 68% | 47% | 68% | 69% | 70% |

## Conclusions

- ✓ AMIS achieves faster convergence than MCMC for GPs while being easy to tune and implement and facilitating massive parallelization
- ✓ AMIS for GPs offers an efficient probabilistic framework for SSP

## References

- [1] Xiong, X., Šmídl, V. and Filippone, M. Adaptive Multiple Importance Sampling for Gaussian Processes, eprint arXiv:1508.01050v2, 2016.
- [2] Segalín, C., Perina, A., Cristani, M. and Vinciarelli, A. The pictures we like are our image: continuous mapping of favourite pictures into self-assessed and attributed personality traits. IEEE Transactions on Affective Computing (to appear), 2016.
- [3] Filippone, M., Girolami, M. Pseudo-marginal Bayesian inference for Gaussian processes. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2014.
- [4] Cornuet, J.-M., Marin, J.-M., Mira, A., Robert, C. P. Adaptive multiple importance sampling. Scandinavian Journal of Statistics 39, 798-812, 2012.