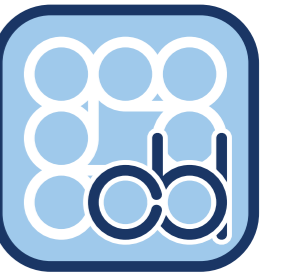




# Bayesian Active Learning for Gaussian Process Classification

Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, Máté Lengyel

Computational and Biological Learning Lab, Department of Engineering, University of Cambridge



Computational and Biological Learning  
University of Cambridge

## Introduction to Active Learning

- AL concerns designing learners that choose their training data.
- Applications include: sensor placement, information extraction, speech recognition, cognitive science, collaborative filtering, quantum tomography etc.
- Referred to as ‘Optimal Experimental Design’ in statistics.
- We revisit the Information Theoretic approach.

## Bayesian Information Theoretic AL

Latent parameters  $\theta \in \Theta$  govern dependence of  $\mathbf{y} \in \mathcal{Y}$  on input  $\mathbf{x} \in \mathcal{X}$  (discriminative model). Observe data,  $\mathcal{D}$ , Bayes rule yields the posterior distribution over parameters  $p(\theta|\mathcal{D})$ . Select  $\mathbf{x}_i$  (myopically) to minimize the posterior entropy:

$$\mathbf{x}_{\text{new}} = \arg \max_{\mathbf{x}} H[\theta|\mathcal{D}] - \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \mathcal{D})} [H[\theta|\mathbf{y}, \mathbf{x}, \mathcal{D}]]$$

Problems:

- Parameter space is often high dimensional, for GPs, it is infinite dimensional.
- Posterior updates required for all input/output combinations ( $\mathcal{O}(N_x N_y)$ ).

## Solution: Rearrange to Dataspace

$$\begin{aligned} H[\theta|\mathcal{D}] - \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \mathcal{D})} [H[\theta|\mathbf{y}, \mathbf{x}, \mathcal{D}]] \\ &= I[\mathbf{y}, \theta|\mathbf{x}, \mathcal{D}] \\ &= H[\mathbf{y}|\mathbf{x}, \mathcal{D}] - \mathbb{E}_{\theta \sim p(\theta|\mathcal{D})} [H[\mathbf{y}|\mathbf{x}, \theta]] \end{aligned}$$

- Output space is often low dimensional and  $\mathcal{O}(1)$  posterior updates required.
- We call this Bayesian Active Learning by Disagreement (BALD).
- Aside: equivalent to the Jensen-Shannon divergence.

## Review of Gaussian Processes

GPs provide a prior over functions  $f: \mathcal{X} \rightarrow \mathbb{R}$ :

$$f \sim \text{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

For regression/classification define likelihood functions respectively:

$$\mathbf{y}|\mathbf{x}, f \sim \mathcal{N}(f(\mathbf{x}), \sigma^2), \quad \mathbf{y}|\mathbf{x}, f \sim \text{Bernoulli}(\Phi(f(\mathbf{x}))) \quad \Phi(z) = \int_{-\infty}^z \mathcal{N}(0, 1) dz$$

For classification, posterior is intractable, make a Gaussian approximation (Expectation Propagation (EP), the Laplace approximation, Variational methods).

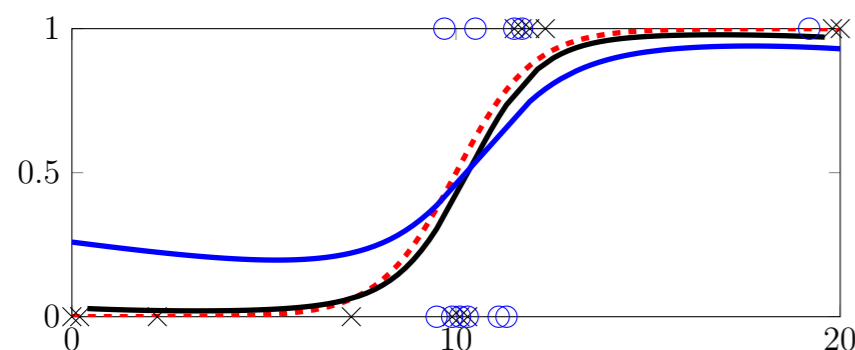


Figure 1: Toy active GPC problem. True generating function is (---). 15 actively selected samples are drawn using both BALD (x) and Maximum Entropy Sampling (o). The predictive distributions from BALD and MES are (—) and (—) respectively.

## BALD for GPC

Two terms need to be computed:

$$\begin{aligned} H[\mathbf{y}|\mathbf{x}, \mathcal{D}] &\approx h \left( \int \Phi(f_x) \mathcal{N}(f_x|\mu_{\mathbf{x}, \mathcal{D}}, \sigma_{\mathbf{x}, \mathcal{D}}^2) df_x \right) \\ &= h \left( \Phi \left( \frac{\mu_{\mathbf{x}, \mathcal{D}}}{\sqrt{\sigma_{\mathbf{x}, \mathcal{D}}^2 + 1}} \right) \right) \\ \mathbb{E}_{f \sim p(f|\mathcal{D})} [H[\mathbf{y}|f]] &\approx \int h(\Phi(f_x)) \mathcal{N}(f_x|\mu_{\mathbf{x}, \mathcal{D}}, \sigma_{\mathbf{x}, \mathcal{D}}^2) df_x \\ &\stackrel{2}{\approx} \int \exp \left( -\frac{f_x^2}{\pi \ln 2} \right) \mathcal{N}(f_x|\mu_{\mathbf{x}, \mathcal{D}}, \sigma_{\mathbf{x}, \mathcal{D}}^2) df_x \\ &= \frac{C}{\sqrt{\sigma_{\mathbf{x}, \mathcal{D}}^2 + C^2}} \exp \left( -\frac{\mu_{\mathbf{x}, \mathcal{D}}^2}{2(\sigma_{\mathbf{x}, \mathcal{D}}^2 + C^2)} \right) \end{aligned}$$

where:

$$h(p) = -p \log p - (1-p) \log(1-p), \quad C = \sqrt{\frac{\pi \ln 2}{2}}$$

- $\stackrel{1}{\approx}$  is a Gaussian approximation to intractable posterior.
- $\stackrel{2}{\approx}$  is a squared exponential approximation to  $h(\Phi(f_x))$  (binary entropy of Normal cdf).
- The objective function is smooth and differentiable.

## Summary

- Apply an approximate inference algorithm to get  $\mu_{\mathbf{x}, \mathcal{D}}$  and  $\sigma_{\mathbf{x}, \mathcal{D}}$  for each point of interest  $\mathbf{x}$ .
- Select  $\mathbf{x}$  that maximises:

$$h \left( \Phi \left( \frac{\mu_{\mathbf{x}, \mathcal{D}}}{\sqrt{\sigma_{\mathbf{x}, \mathcal{D}}^2 + 1}} \right) \right) - \frac{C}{\sqrt{\sigma_{\mathbf{x}, \mathcal{D}}^2 + C^2}} \exp \left( -\frac{\mu_{\mathbf{x}, \mathcal{D}}^2}{2(\sigma_{\mathbf{x}, \mathcal{D}}^2 + C^2)} \right) \quad (1)$$

## Related Algorithms

The following algorithms are closely related, often approximating the BALD objective:

- Uncertainty Sampling [Lewis and Gale, 1994] / Maximum Entropy Sampling [Sebastiani and Wynn, 2000].
- The Informative Vector Machine [Lawrence and Herbrich, 2001].
- Query by Committee [Freund et al., 1997].
- SVM-based active learning [Tong and Koller, 2001].

## Extensions

Further work that we have performed:

- Comparison to decision theoretic algorithms.
- Hyperparameter learning ( $\theta^+, \theta^-$  = params of interest, nuisance parameters):  $H[\mathbb{E}_{p(\theta^+, \theta^-|\mathcal{D})}[\mathbf{y}|\mathbf{x}, \theta^+, \theta^-]] - \mathbb{E}_{p(\theta^+|\mathcal{D})} [H[\mathbb{E}_{p(\theta^-|\theta^+, \mathcal{D})}[\mathbf{y}|\mathbf{x}, \theta^+, \theta^-]]]$
- Multiclass: combine criteria for  $K$  one-versus-all classifiers.
- Preference Learning: extend GP methods of [Chu and Ghahramani, 2005].

## Results

Experiments run on *pool-based* active learning. Test set accuracy plotted is against number of queries.

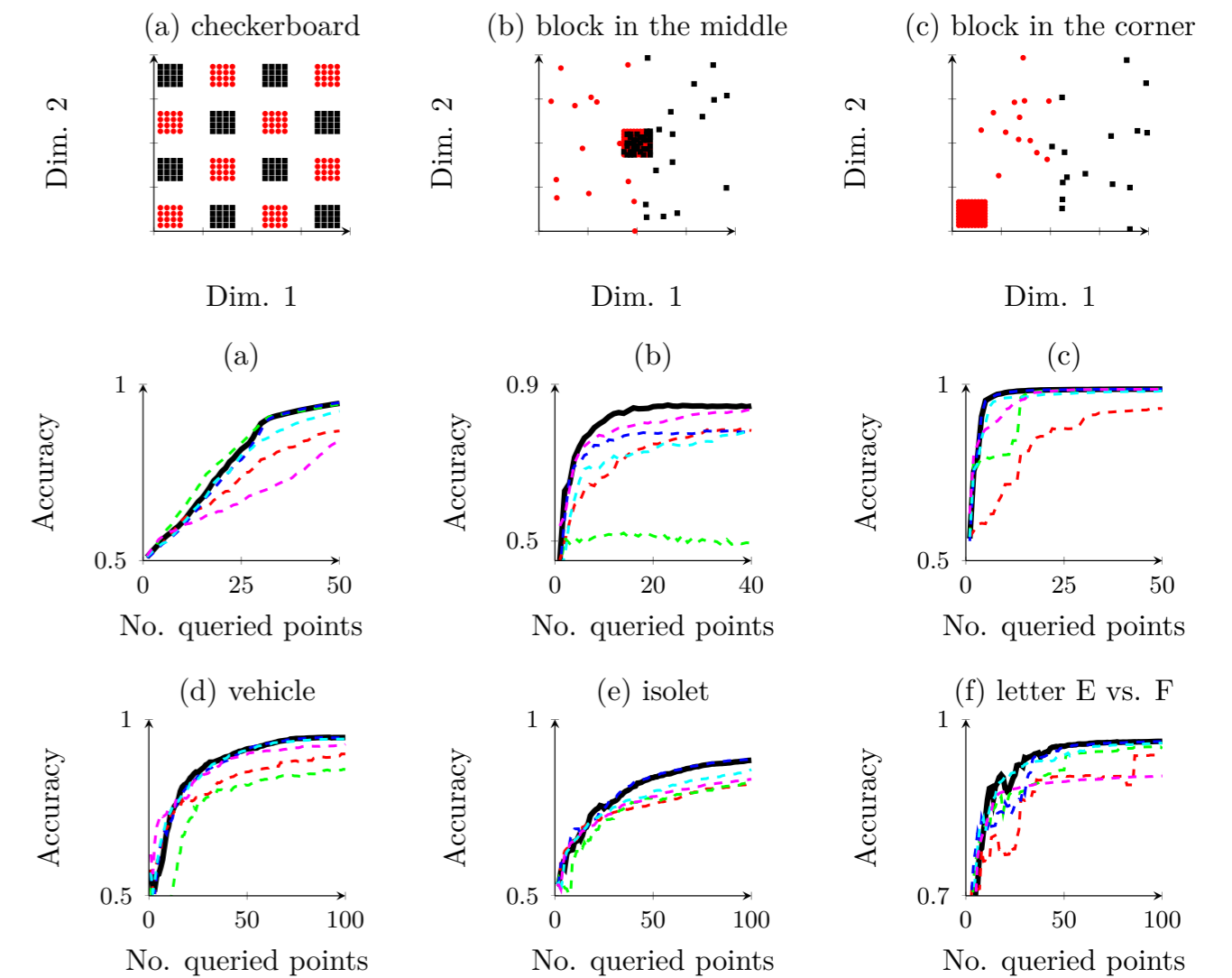


Figure 2: Top: Three 2D artificial datasets designed to test the algorithms in pathological scenarios. Middle: Results for corresponding artificial datasets using BALD (—), random query (---), MES (---), IVM (---), QBC (---), active SVM (---). Bottom: Results on three real-world datasets.

## Acknowledgements

This work is made possible by our sponsors: Google, Europe (Neil Houlsby) and Trinity College, Cambridge (Ferenc Huszár) and the Wellcome Trust (Máté Lengyel).

## References

- [Chu and Ghahramani, 2005] Chu W, Ghahramani Z. 2005. Preference learning with gaussian processes. In: Proceedings of the 22nd international conference on Machine learning. ACM, pp 137–144.
- [Freund et al., 1997] Freund Y, Seung H, Shamir E, Tishby N. 1997. Selective sampling using the query by committee algorithm. Machine Learning 28:133–168.
- [Lawrence and Herbrich, 2001] Lawrence N, Herbrich R. 2001. A sparse Bayesian compression scheme—the informative vector machine. In: NIPS 2001 workshop on kernel methods. Citeseer.
- [Lewis and Gale, 1994] Lewis D, Gale W. 1994. A sequential algorithm for training text classifiers. In: Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval. Springer-Verlag New York, Inc., pp 3–12.
- [Sebastiani and Wynn, 2000] Sebastiani P, Wynn H. 2000. Maximum entropy sampling and optimal Bayesian experimental design. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 62:145–157.
- [Tong and Koller, 2001] Tong S, Koller D. 2001. Support vector machine active learning with applications to text classification. Journal of Machine Learning Research 2:45–66.