# **Exploiting Structure in Bayesian Optimization using Multi-Fidelity Probabilistic Models**

Boris Bogaerts, Rudi Penne



#### Overview

- Application
- Machinery
- Theory
- Examples

#### Overview

- Application
- Machinery
- Theory
- Examples

## Practical problem



## **Preliminary conclusion**

- Submodular orienteering problem (coverage function)
- Multilinear extention can be estimated as an integral of the pointwise product of two functions
- Problem is factorized in a function dependent on the constraints, and a function dependent on the problem :  $\overline{F}_x = \int x(v)w(v)dv$

12-12-2019

• Functions are complex

#### Overview

- Application
- Machinery
- Theory
- Examples

## General optimization idea

We are interested in following problem

 $\int g(x).h(x)dx$ 

But we do not have access to these possibly complex functions, so we will learn them

$$P(f_g^*|x, y_g, x^*) \text{ and } P(f_h^*|x, y_h, x^*)$$

We also learn the dependence on the parameters which we want to optimize  $P(f_g^*|[x, \theta_g], y_g, [x, \theta_g]^*)$  and  $P(f_h^*|[x, \theta_h], y_h, [x, \theta_h]^*)$ 

## General optimization idea

The final optimization problem becomes

$$\operatorname{argmax}_{\theta_g,\theta_h} E\left[\int f_g^*(x)f_h^*(x)dx\right]$$

Our optimization strategy:

- Learn both distributions in advance (Gaussian Processes)
- In the optimization loop use optimized quadrature points to evaluate integral.
- Effectively sample the parameter space using Bayesian Optimisation

12-12-2019

#### **Gaussian Process**



Gaussian distribution over RKHS

Analytic exact expression for posterior after point measurements

#### Bayesian quadrature

Integration without close form expression



## Bayesian quadrature

#### Integration without close form expression





#### Bayesian quadrature **Bayesian quadrature** $f_2(x) = \sin(x+1) + 1$ Integration without close form expressic 2.5 1.5 **Trapezoidal rule** -0.5 b<sub>2</sub> a<sub>2</sub> $\int f(x)dx \approx \mathbf{y} \ \mathbf{K}^{-1} \ \left| \ K(X,x)dx \right|^{-1}$ У 2 $[f(x_1) \quad \cdots \quad f(x_n)] \begin{bmatrix} w_1 \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$ $\int f(x)dx \approx \sum_{i=1}^{n} w_i \big( f(x_{i-1}) + f(x_i) \big)$

#### **Product of Gaussian Processes**

Let  $h \sim GP(m_h(x), C_h(x, x'))$  and  $g \sim GP(m_g, C_g(x, x'))$  then we are interested the distribution over f.g given by its mean:

$$m_{gh}(x) = E[h(x)]E[g(x)] = m_h(x).m_g(x)$$

And covariance

$$C_{gh}(x,x') = C(f(x),g(x),f(x')g(x'))$$

The covariance of products of Gaussian random variables is analytically available (Bohrnstedt et al [7]) and for GP's given by:  $C_{ah}(x,x') = m_h(x)C_a(x,x')m_h(x') + m_g(x)C_h(x,x')m_g(x') + C_h(x,x')C_g(x,x')$ 

We choose to conveniently neglect other moments

#### Product quadrature

Further derivations are tedious!!! But three types of integrals need to be computed:

$$\int \mathcal{K}_{h}(X_{h_{i}}, x) \mathcal{K}_{g}(X_{g_{j}}, x) dx$$
$$\int \int \mathcal{K}_{h}(x, x') \mathcal{K}_{g}(x, x') dx' dx$$
$$\int \mathcal{K}_{h}(x, X_{h_{i}}) \int \mathcal{K}_{h}(X_{h_{j}}, x') \mathcal{K}_{g}(x, x') dx' dx$$

Computational cost Mean : O(MN)Variance :  $O(M^2N)$ 

12-12-2019

Where  $X_{h_i}$  and  $X_{h_i}$  are training points

A non-exhaustive list of possible kernel combinations is given by Briol et al [3]



#### Our choice of kernel

Product of one dimensional piecewise polynomials (1st order) with compact support

 $\mathcal{K}(x, x') = \prod_{i=1}^{D} \max(\theta_{1,i}(1 - \theta_{2,i}|x - x'|), 0)$ 

- Brownian motion prior (nonlinearities)
- Compact support
- Sparse



12-12-2019



12-12-2019





12-12-2019

## Note on difference with BQ

• Bayesian quadrature

Choose next quadrature point to minimize variance on the estimated integral

• Our setting

Remove n quadrature points that minimize the variance increase on the estimated integral



We use the fully independent training conditional approximation by Snelson and Ghahramani (2006) (**FITC**, using classification of Quiñonero-Candela et al(2005))

12-12-2019

We use the fully independent training conditional approximation by Snelson and Ghahramani (2006) (**FITC**, using classification of Quiñonero-Candela et al(2005))



We use the fully independent training conditional approximation by Snelson and Ghahramani (2006) (**FITC**, using classification of Quiñonero-Candela et al(2005))



We use the fully independent training conditional approximation by Snelson an Ghahramani (2006) (**FITC**, using classification of Quiñonero-Candela et al(2005))



We use the fully independent training conditional approximation by Snelson and Ghahramani (2006) (**FITC**, using classification of Quiñonero-Candela et al(2005))



- Positions of pseudo input locations u are optimized to maximize the marginal likelihood (details in the original work by Snelson an Ghahramani (2006))
- Inference is possible by marginalizing out pseudo inputs u
- Exact test conditional (Quiñonero-candela, Rasmussen, & Herbrich, 2005)

 $p(f_*|\mathbf{u}) = \mathcal{N}(K_{*u}K_{uu}^{-1}\mathbf{u}, K_{**} - K_{*u}K_{uu}^{-1}K_{u*})$ 



## Example: varying number of pseudo inputs



12-12-2019

## Example: varying number of pseudo inputs



12-12-2019

## Optimization

The goal is not to evaluate every point. We consider following optimization problem:

$$\underset{\Delta \in \mathbb{R}^n}{\operatorname{argmax}} E\left[\int P(f_g^* | x, y_g, x^*) P(f_h^* | x + \Delta, y_h, x^*) dx\right]$$

Strategy (Bayesian optimization):

- Evaluate functions at points  $\Delta^*$
- Put Gaussian process prior on all  $\Delta \in \mathbb{R}^n$
- Select new location with highest expected improvement/...

12-12-2019

#### Overview

- Application
- Machinery
- <u>Theory</u>
- Examples



$$\underset{\Delta \in \mathbb{R}^{n}}{\operatorname{argmax}} E\left[ \int P(f_{g}^{*} | x, y_{g}, x^{*}) \cdot P(f_{h}^{*} | x + \Delta, y_{h}, x^{*}) dx \right]$$

$$g(a, b)$$

Learning g is what we can do







#### Overview

- Application
- Machinery
- Theory
- <u>Examples</u>







12-12-2019













## The optimization puzzle: Approximate puzzle piece







20 pseudo-inputs

30 pseudo-inputs

40 pseudo-inputs

12-12-2019









12-12-2019



#### **Observation: Cost function becomes simpler/nicer !!!**



Puzzle	Piece
100 pseudo-inputs	30 pseudo-inputs







Puzzle	Piece
50 pseudo-inputs	15 pseudo-inputs





12-12-2019



## Interesting partial results



#### Free VR demo if you visit Antwerp

\*GP compressed (from 31 000 points) to only 50 pseudo-input points

12-12-2019

## **Preliminary conclusion**

- Application
- Machinery
- <u>Theory</u>

#### References

Briol, F.-X., Oates, C. J., Girolami, M. A., Osborne, M. A., & Sejdinovic, D. (2015). Probabilistic Integration. CoRR, abs/1512.00933.

Halcon, B. (n.d.). Robot-enabled LDV basketball modal data collection trial run. Retrieved from https://www.youtube.com/watch?v=xx2JI7fGf\_E%0A

Karimi, M., Lucic, M., Hassani, H., & Krause, A. (2017). Stochastic submodular maximization: The case of coverage functions. In Advances in Neural Information Processing Systems (pp. 6856–6866).

Quiñonero-candela, J., Rasmussen, C. E., & Herbrich, R. (2005). A unifying view of sparse approximate Gaussian process regression. Journal of Machine Learning Research.

https://doi.org/10.1163/016918609X12529286896877

Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian Processes using Pseudo-inputs. Advances in Neural Information Processing Systems 18. https://doi.org/10.1.1.60.2209

Zhang, H., & Vorobeychik, Y. (2016). Submodular Optimization with Routing Constraints. Proceedings of the 30th Conference on Artificial Intelligence (AAAI 2016).



12-12-2019

