
GLASSES: Relieving The Myopia Of Bayesian Optimisation

Anonymous Author(s)

Affiliation

Address

email

Abstract

We present GLASSES: Global optimisation with Look-Ahead through Stochastic Simulation and Expected-loss Search. The vast majority of global optimisation approaches in use are myopic, in only considering the impact of the next function value; the remaining approaches are able to consider only a handful of future evaluations. Our novel algorithm, GLASSES, permits the consideration of dozens of evaluations into the future. We show that the far-horizon planning thus enabled leads to substantive performance gains in empirical tests.

1 Introduction

Almost all global optimisation techniques are myopic, in considering no more than a single step into the future. We define the multi-step lookahead problem as the global optimisation of an function by considering the significance of the next function evaluation on function evaluations (steps) further into the future. It is clear that a solution to the problem would offer performance gains. For example, consider the case in which we have a budget of two evaluations with which to optimise a function $f(x)$ over the domain $\mathcal{X} = [0, 1] \subset \mathbb{R}$. If we are strictly myopic, our first evaluation will likely be at $x = 1/2$, and our second then at only one of $x = 1/4$ and $x = 3/4$. This myopic strategy thereby ignores half of the domain \mathcal{X} , regardless of the second choice. If we adopt a two-step lookahead approach, we will select function evaluations that will be more evenly distributed across the domain by the time the budget is exhausted, and will hence be more informative about f and its optimum.

There is a limited literature on the multi-step lookahead problem. [7] perform multi-step lookahead by optimising future evaluation locations, and sampling over future function values. This approach scales poorly with the number of future evaluations considered, and the authors present results for no more than two-step lookahead. [6] reframe the multi-step lookahead problem as a partially observed Markov decision process, and adopt a Monte Carlo tree search approach in solving it. Again, the scaling of the approach permits the authors to consider no more than six steps into the future.

There is a clear link between the multi-step lookahead problem and that considered in the literature as *batch* Bayesian optimisation. The two problems are distinct but related: the multi-step lookahead problem requires the challenging marginalisation over unknown future evaluation *locations*, in addition to the unknown future evaluation *values* also marginalised by batch approaches. Similarly to the state-of-the-art in multi-step lookahead, the batch literature provides only poor scaling with the number of evaluations. [5] present results for no more than six simultaneous function evaluations.[2, 1] use the surrogate model for f to generate ‘fake’ observations and avoid the marginalization step. This produce a large accumulation of errors that does not allow the use these techniques for the collection of large batches.

We propose an algorithm, GLASSES, that provides scaling superior to existing alternatives.

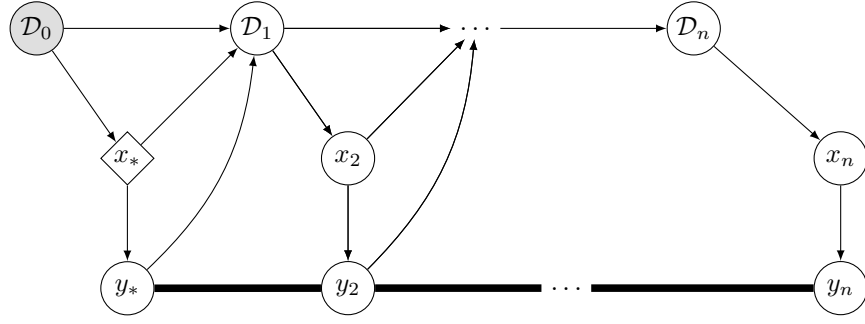


Figure 1: A Bayesian network describing the n -step lookahead problem. The shaded node (\mathcal{D}_0) is known, and the diamond node (x_*) is the current decision variable. All y nodes are correlated with one another under the GP model.

2 Conditional DPPs for step ahead locations

To approximate $\Lambda_n(\mathbf{x}_*)$ we compute $E_{p(\mathbf{y})}[\min(\mathbf{y}, \eta)]$ where $\mathbf{y} = \{y_1, \dots, y_n\}$ with $p(\mathbf{y}) \sim \mathcal{N}(\mathbf{y}; \mu, \Sigma)$ is the multivariate random vector that we obtain after evaluation the predictive distribution of the GP at locations $\mathbf{x}_*, \mathbf{x}_2, \dots, \mathbf{x}_n$. Point \mathbf{x}_* is fixed (is the point where we evaluate Λ_n) and we consider $\mathbf{x}_2, \dots, \mathbf{x}_n$ to be a sample of the conditional k-DPP (for $k = n$) on \mathbf{x}_* with kernel L (in principle, the one from the GP).

Let \mathbf{L} be the kernel matrix corresponding to the evaluation of L on a finite set Ω of potential points (pre-uniformly sampled in the domain of interest, for instance). The distribution obtained by conditioning on having observed $\mathbf{x}_* \in \Omega$ can be obtained as follows. Let $B \subset \Omega$ a non intersecting set with \mathbf{x}_* . We have that

$$p_L(\mathbf{x}_* \cup B | \mathbf{x}_* \subseteq Z) = \frac{p_L(Z = \mathbf{x}_* \cup B)}{p(\mathbf{x}_* \subseteq Z)} = \frac{\det(\mathbf{L}_{\mathbf{x}_* \cup B})}{\det(\mathbf{L} - \mathbf{I}_{\mathbf{x}_*})} \quad (1)$$

where $\mathbf{I}_{\mathbf{x}_*}$ is the matrix with ones in the diagonal entries indexed by elements of $\Omega - \mathbf{x}_*$ and zeros elsewhere. This conditional distribution is again a DPP over subsets of $\Omega - \mathbf{x}_*$ [3] with kernel

$$\mathbf{L}^{\mathbf{x}_*} = ((\mathbf{L} + \mathbf{I}_{\mathbf{x}_*})^{-1})_{\bar{\mathbf{x}}_*}^{-1}.$$

$[\cdot]_{\bar{\mathbf{x}}_*}$ represents the restriction of the matrix to all rows and columns not indexed by \mathbf{x}_* . The previous inverses exist if and only if the probability of \mathbf{x}_* appearing is nonzero, as is the case in our context. A second marginalization is later needed to generate samples of size k . Figure 2 shows an example of samples from a k-DPP and a conditional k-DPP with the same kernel.

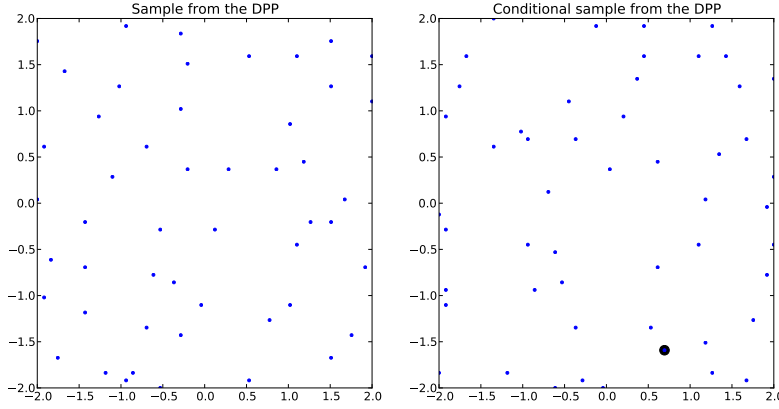


Figure 2: Left: sample from a k-DPP (k=50) for a SE kernel with length-scale 0.5. Right: sample from a k-DPP (k=50) conditional to x_1 (black dot) being in the selected set.

3 Computation of the Expected loss

Our goal is to compute $E_{p(\mathbf{y})}[\min(\mathbf{y}, \eta)]$ for $\mathbf{y} = \{y_1, \dots, y_n\}$ and $p_0(\mathbf{y}) \sim \mathcal{N}(\mathbf{y}; \mu, \Sigma)$. This approximates the n-steps ahead expected loss $\Lambda_n(\mathbf{x}_*)$. The goal of this section is to write $E_{p(\mathbf{y})}[\min(\mathbf{y}, \eta)]$ in a way that it is suitable to be computed by Expectation Propagation. Next proposition will do the work (I think).

Proposition 1

$$E_{p(\mathbf{y})}[\min(\mathbf{y}, \eta)] = \sum_{j=1}^n \int_{\mathbb{R}^n} y_j \prod_{i=1}^n t_{j,i}(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} + \eta \int_{\mathbb{R}^n} \prod_{i=1}^n h_i(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} \quad (2)$$

where $h_i(\mathbf{y}) = \mathbb{I}\{y_i > \eta\}$ and

$$t_{j,i}(\mathbf{y}) = \begin{cases} \mathbb{I}\{y_j \leq \eta\} & \text{if } i=j \\ \mathbb{I}\{0 \leq y_i - y_j\} & \text{otherwise.} \end{cases}$$

Proof 1 Denote by

$$\begin{aligned} E_{p(\mathbf{y})}[\min(\mathbf{y}, \eta)] &= \int_{\mathbb{R}^n} \min(\mathbf{y}, \eta) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} \\ &= \int_{\mathbb{R}^n - (\eta, \infty)^n} \min(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} + \int_{(\eta, \infty)^n} \eta \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} \end{aligned}$$

The first term can be written as follows:

$$\int_{\mathbb{R}^n - (\eta, \infty)^n} \min(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} = \sum_{j=1}^n \int_{P_j} y_j \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} \quad (3)$$

where $P_j := \{\mathbf{y} \in \mathbb{R}^n - (\eta, \infty)^n : y_j \leq y_i, \forall i \neq j\}$. We can do this because the regions P_j are disjoint and it holds that $\cup_{j=1}^n P_j = \mathbb{R}^n - (\eta, \infty)^n$. Also, note that the $\min(\mathbf{y})$ can be replaced within the integrals since within each P_j it holds that $\min(\mathbf{y}) = y_j$. Rewriting the integral in terms of indicator functions we have that

$$\sum_{j=1}^n \int_{P_j} y_j \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} = \sum_{j=1}^n \int_{\mathbb{R}^n} y_j \prod_{i=1}^n t_{j,i}(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} \quad (4)$$

Algorithm 1 Decission process of the GLASSES algorithm.

Input: dataset $\mathcal{D}_0 = \{(\mathbf{x}_0, y_0)\}$, number of remaining evaluations (n), representer points (r) and DPP replicates (s).

Fit a GP with kernel k to \mathcal{D}_0 .

Select $\mathbf{x}_{1*}, \dots, \mathbf{x}_{r*}$ representer points of the loss.

for $j = 1$ **to** r **do**

 Take s samples from a conditional n-DPP of kernel k given \mathbf{x}_{j*} .

 Approximate the expected loss at \mathbf{x}_j^* for the s samples computing $E[\min(\mathbf{y}, \eta)]$.

 Average the expected loss for the s samples and obtain $\tilde{\Lambda}_n(\mathbf{x}_j^*)$.

end for

Approximate $\Lambda_n(\mathbf{x}_*)$ fitting a GP₂ to $\{(\mathbf{x}_{j*}, \tilde{\Lambda}_n(\mathbf{x}_{j*}))\}_{j=1}^r$ with posterior mean μ_2 .

Returns: New location at $\arg \min_{\mathbf{x} \in \mathcal{X}} \{\mu_2(\mathbf{x})\}$.

where $t_{j,i}(y) = \mathbb{I}\{y_i \leq \eta\}$ if $j = i$ and $t_{j,i}(y) = \mathbb{I}\{y_j \leq y_i\}$ otherwise.

The second term can be written as

$$\int_{(\eta, \infty)^n} \eta \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} = \eta \int_{\mathbb{R}^n} \prod_{i=1}^n h_i(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} \quad (5)$$

where $h_i(\mathbf{y}) = \mathbb{I}\{y_i > \eta\}$. Merge (4) and (5) to conclude the proof.

All the elements in (2) can be rewritten in a way that can be computed using EP but the work in [4].

- The second term is a Gaussian probability on unbounded polyhedron in which the limits are aligned with the axis.
- The first term requires some more processing but it is still computable under the assumptions in [4]. Let \mathbf{w}_j the j th canonical vector. Then we have that

$$\int_{\mathbb{R}^n} y_j \prod_{i=1}^n t_{j,i}(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} = \mathbf{w}_j^T \int_{\mathbb{R}^n} \mathbf{y} \prod_{i=1}^n t_{j,i}(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y} \quad (6)$$

$$= \mathbf{w}_j^T E[\mathbf{y}] z_j \quad (7)$$

where the expectation is calculated over the normalized distribution over P_j , the one EP approximates with $q(\mathbf{y})$, and for z_j being the normalizing constant

$$z_j = \int_{\mathbb{R}^n} \prod_{i=1}^n t_{j,i}(\mathbf{y}) \mathcal{N}(\mathbf{y}; \mu, \Sigma) d\mathbf{y}$$

Because EP does moments matching, both the normalizing constant and the expectation are available.

References

- [1] Javad Azimi, Ali Jalali, and Xiaoli Fern. Dynamic batch Bayesian optimization. *CoRR*, abs/1110.3347, 2011.
- [2] Javad Azimi, Ali Jalali, and Xiaoli Zhang Fern. Hybrid batch Bayesian optimization. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- [3] Alexei Borodin and Eric M. Rains. Eynard–Mehta theorem, Schur process, and their Pfaffian analogs. *Journal of Statistical Physics*, 121(3-4):291–317, Nov 2005.
- [4] John P. Cunningham, Philipp Hennig, and Simon Lacoste-Julien. Gaussian probabilities and expectation propagation. *arXiv:1111.6832 [stat]*, Nov 2011. arXiv: 1111.6832.
- [5] David Ginsbourger, Rodolphe Le Riche, and Laurent Carraro. A multi-points criterion for deterministic parallel global optimization based on Gaussian processes. *HAL: hal-00260579*, 2009.

216 [6] Roman Marchant, Fabio Ramos, and Scott Sanner. Sequential Bayesian optimisation for spatial-
217 temporal monitoring. In *Proceedings of the International Conference on Uncertainty in Artificial*
218 *Intelligence*, 2014.

219 [7] Michael A. Osborne, Roman Garnett, and Stephen J. Roberts. Gaussian processes for global
220 optimization. In *3rd international conference on learning and intelligent optimization (LION3)*,
221 pages 1–15, 2009.

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269