Some random things I learned writing the BayesOpt book

Roman Garnett

BAYESIAN OPTIMIZATION

ROMAN GARNETT



bayesoptbook.com

(Semi-) Joking advice: Don't write a book...

Book timeline... (4 authors, January 2013)

Bayesian Optimization book



Nando de Freitas <nando@cs.ubc.ca> to me, Michael, Frank, Nando 💌

OK guys. I think it's time for us to do this seriously.

Expected Improvement with Noise

Step 1: build a model of (noisy) observations (*x*, *y*)

- latent function model, p(f)
- observation model, $p(y | x, \phi)$, $\phi = f(x)$ e.g., Gaussian noise

e.g., GP



Step 2: choose a utility function u(D), $D = \{(x, y)\}$



Step 3: give up on the optimal policy



Step 4: derive a policy via one-step lookahead (greedily maximize one-step expected gain in utility $D \rightarrow D'$)

$$\alpha(x;\mathcal{D}) = \mathbb{E}\big[u(\mathcal{D}') \mid x,\mathcal{D}\big] - u(\mathcal{D})$$



(...nothing to see here...) $u(\mathcal{D}_{\tau})$

Step 4: derive a policy via one-step lookahead

$$\alpha(x;\mathcal{D}) = \mathbb{E}\left[u(\mathcal{D}') \mid x,\mathcal{D}\right] - u(\mathcal{D})$$

(wrt noisy observation *y*! consequence: in general, penalizes high noise)

Prevalent in BayesOpt!

Utility

simple reward

global simple reward

information gain

Policy

expected improvement

knowledge gradient

mutual information (aka entropy search)

Noiseless expected improvement

utility (best seen value): $u(D) = \phi^* = \max \mathbf{f}$ marginal gain: $\max(\phi - \phi^*, 0)$ ϕ^*

expected utility easy to compute, has nice properties, etc.

Noiseless expected improvement



expected utility easy to compute, has nice properties, etc.

$$\alpha_{\rm EI}(x;\mathcal{D}) = (\mu - \phi^*) \Phi\left(\frac{\mu - \phi^*}{\sigma}\right) + \sigma \phi\left(\frac{\mu - \phi^*}{\sigma}\right)$$

Noiseless expected improvement

expected utility easy to compute, has nice properties, etc.



very tempting to start here and try to "fix" this!

$$\alpha_{\rm EI}(x;\mathcal{D}) = (\mu - \phi^*) \Phi\left(\frac{\mu - \phi^*}{\sigma}\right) + \sigma \phi\left(\frac{\mu - \phi^*}{\sigma}\right)$$

"Fixing" the expected utility

- plug-in estimators: use noiseless EI with "guess" of max f
- expectation of El

with respect to **f**

(Letham, et al. 2019)



Let's start with utility!

idea: consider gathering data to support a recommendation after optimization

action space: visited locations **x** utility: risk-neutral

optimal recommendation:

maximum of posterior mean on $\mathbf{x} = u(D)$



The noisy setting: Utility

maximum of posterior mean on $\mathbf{x} = u(D)$

- compatible with noiseless El!
- compatible with knowledge gradient! (just a different action space)



The difficulty

maximum of posterior mean can be anywhere!

local reasoning of just f(x), y not enough!



The fix (Frazier, et al. 2009)

posterior mean update is linear in observed value



The fix (Frazier, et al. 2009)

can compute piecewise linear update to max in $O(n^2 \log n)$



The fix (Frazier, et al. 2009)

sums of standard normal CDFs, PDFs as before



The result

- handles hetereoskedastic noise automatically / correctly
- handles correlations in / global nature of posterior mean
- noiseless El special case
- closed form



Alternative approaches



Why?

- ignores correlations in posterior mean update
- assumption of exact observations in expectation does not match true observation model
- (but honestly this is all fine for highish SNR)



Marginalizing Hyperparameters in Policy

Marginalizing hyperparameters







Standard approach

Let utility $u(D; \theta)$ depend on θ and integrate the hyperprameter-conditional acquisition function against the hyperparameter posterior

$$\int \alpha(x;\mathcal{D},\boldsymbol{\theta}) \, p(\boldsymbol{\theta} \mid \mathcal{D}) \, \mathrm{d}\boldsymbol{\theta}$$

Standard approach

Let utility $u(D; \theta)$ depend on θ and integrate the hyperprameter-conditional acquisition function against the hyperparameter posterior

$$\int \alpha(x; \mathcal{D}, \theta) p(\theta \mid \mathcal{D}) d\theta$$

blind to uncertainty in

 $\theta!$

Standard approach

Let utility $u(D; \theta)$ depend on θ and integrate the hyperprameter-conditional acquisition function against the hyperparameter posterior

$$\int \alpha(x; \mathcal{D}, \theta) p(\theta \mid \mathcal{D}) d\theta$$

blind to uncertainty in

 $\theta!$

Alternative approach

Define utility with respect to marginal model from the beginning!

E.g., for EI or KG, use θ marginal posterior mean (for a terminal recommendation we'd be marginalizing θ , right?)

$$\int \mu_{\mathcal{D}}(\boldsymbol{x};\boldsymbol{\theta}) \, \boldsymbol{p}(\boldsymbol{\theta} \mid \mathcal{D}) \, \mathrm{d}\boldsymbol{\theta}$$

Example

- function is f(x) = x or f(x) = -x
- knowledge gradient
- for standard approach, acquisition function is flat! (maximum of θ-conditional posterior mean always equal)
- for alternative approach, get sensible answers (prefer sampling on boundary)



History of BayesOpt

Who first proposed the following policies?

probability of improvement? expected improvement? upper confidence bound? knowledge gradient?

What I thought...

probability of improvement? expected improvement? upper confidence bound? knowledge gradient?

Harold Kushner, 1964 Jonas Mockus, 1972 Cox and John, 1998 Frazier, et al., 2009

I was wrong!

probability of improvement? expected improvement? upper confidence bound? knowledge gradient?

Harold Kushner, 1964

Jonas Mockus, 1972

Cox and John, 1998

Frazier, et al. 2009

Okay we can agree on this right? (1964)

H. J. KUSHNER RIAS, Inc., Baltimore, Md.

A New Method of Locating the Maximum Point of an Arbitrary Multipeak Curve in the Presence of Noise'

A versatile and practical method of searching a parameter space is presented. Theoretical and experimental results illustrate the usefulness of the method for such problems as the experimental optimization of the performance of a system with a very general multipeak performance function when the only available information is noise-distributed samples of the function. At present, its usefulness is restricted to optimization with respect to one system parameter. The observations are taken sequentially; but, as opposed to the gradient method, the observation may be located anywhere on the parameter interval. A sequence of estimates of the location of the curve maximum is generated. The location of the next observation may be interpreted as the location of the most likely competitor (with the current best estimate) for the location of the curve maximum. A Brownian motion stochastic process is selected as a model for the unknown function, and the observations are interpreted with respect to the model. The model gives the results a simple intuitive interpretation and allows the use of simple but efficient sampling procedures. The resulting process possesses some powerful convergence properties in the presence of noise; it is nonparametric and, despite its generality, is efficient in the use of observations. The approach seems quite promising as a solution to many of the problems of experimental system optimization.

Surprise twist! (Kushner, 1962)

A Versatile Stochastic Model of a Function of Unknown and Time Varying Form

HAROLD J. KUSHNER

Massachusetts Institute of Technology, Lincoln Laboratories, Lexington 73, Massachusetts

Submitted by Lotfi Zadeh

Properties of a random walk model of an unknown function are studied. The model is suitable for use in the following (among others) problem. Given a system with a performance function of unknown, time varying, and possibly multipeak form (with respect to a single system parameter), and given that the only information available are noise perturbed samples of the function at selected parameter settings, then determine the successive parameter settings such that the sum of the values of the observations is maximum. An attempt to avoid the optimal search problem through the use of several intuitively reasonable heuristics is presented.

Objective Model (Kushner, 1962)

X(†)

VAR X (1)

- Wiener process prior
- additive Gaussian noise



Policy desiderata (Kushner, 1962)

• sample densely

1. As N (the total number of observations) tends to infinity, every region of greater than zero size is sampled at least once.

2. For large N, the initial observations will tend to be information gathering (or play the long shot) and be taken near the point of maximum curve variance.

3. The final observations are taken at points where the expected "pay off" (in whatever sense the observations pay off) will be maximum.

Policy desiderata (Kushner, 1962)

- sample densely
- explore more at the beginning of search

1. As N (the total number of observations) tends to infinity, every region of greater than zero size is sampled at least once.

2. For large N, the initial observations will tend to be information gathering (or play the long shot) and be taken near the point of maximum curve variance.

3. The final observations are taken at points where the expected "pay off" (in whatever sense the observations pay off) will be maximum.

Policy desiderata (Kushner, 1962)

- sample densely
- explore more at the beginning of search
- exploit more at end of search

1. As N (the total number of observations) tends to infinity, every region of greater than zero size is sampled at least once.

2. For large N, the initial observations will tend to be information gathering (or play the long shot) and be taken near the point of maximum curve variance.

3. The final observations are taken at points where the expected "pay off" (in whatever sense the observations pay off) will be maximum.

Policies (Kushner, 1962)

Policy B: probability of improvement (will see again)

B. Sample at the t point (\hat{t}) at which $(\epsilon = \epsilon(N, n)$ is a positive sequence)

$$P(X_t \ge \bar{X}^* + \epsilon) = 1 - \Phi\left(\frac{\bar{X}_t + \epsilon}{\sqrt{\operatorname{Var} X_t}}\right)$$
(3.2)

is maximum.

Policies (Kushner, 1962)

Policy A: upper confidence bound!

A. The location of every observation is selected on the basis of a balance between properties 2 and 3. The simplest such balance is a linear weighing. We select the point at which

$$\sqrt{\operatorname{Var} X_t} + f(N, n) \left(\bar{X}_t - \bar{X}^* \right) \tag{3.1}$$

is maximum.

As far as I can tell...

upper confidence bound? probability of improvement? expected improvement? knowledge gradient?

Harold Kushner, 1962 Harold Kushner, 1964 1962

Further Development (Kushner, 1964)

H. J. KUSHNER RIAS, Inc., Baltimore, Md.

- same model
- probability of improvement
- (what happened to UCB?)

A New Method of Locating the Maximum Point of an Arbitrary Multipeak Curve in the Presence of Noise'

A versatile and practical method of searching a parameter space is presented. Theoretical and experimental results illustrate the usefulness of the method for such problems as the experimental optimization of the performance of a system with a very general multipeak performance function when the only available information is noise-distributed samples of the function. At present, its usefulness is restricted to optimization with respect to one system parameter. The observations are taken sequentially; but, as opposed to the gradient method, the observation may be located anywhere on the parameter interval. A sequence of estimates of the location of the curve maximum is generated. The location of the next observation may be interpreted as the location of the most likely competitor (with the current best estimate) for the location of the curve maximum. A Brownian motion stochastic process is selected as a model for the unknown function, and the observations are interpreted with respect to the model. The model gives the results a simple intuitive interpretation and allows the use of simple but efficient sampling procedures. The resulting process possesses some powerful convergence properties in the presence of noise; it is nonparametric and, despite its generality, is efficient in the use of observations. The approach seems quite promising as a solution to many of the problems of experimental system optimization.

Very thoughtful! (Kushner, 1964)

- very practical
- computational notes
- careful scheduling of improvement thresholds



Fig. 4 Experimental results with no observation noise; locations of observations

Very thoughtful! (Kushner, 1964)

- very practical
- computational notes
- scheduling of improvement thresholds
- handling noise





Aside: Gauss-Markov models



Okay but this is El right? (Mockus 1972)

ation

BAYESIAN METHODS OF SEARCH FOR AN EXTREMUM

I. B. Motskus

Avtomatika i Vychislitel'naya Tekhnika, Vol. 6, No. 3, pp. 53-62, 1972 UDC 62-50:519.83

The problem of finding those methods of saking an astronum that minimize the mathematical expetation of the loss functions for the sacred meth and possible downstions from it is formulated. Avaraging is carried out for a given a priori distribution. The solution is given in the form of a system of recursion relations for the case in which either the number of chapavations or the cost of a single observation is fixed. We also consider the case where the algortike has bounded "memory."

Many methods of searching for an extremum are known. However, their areas of spring tion are not sufficiently clear. Therefore, it is of interest to formulate mathematically and solve, the problem of combining methods that are optimal in some sense. This solution are the necessary for a mathematical theory of extremum search as well as for the solution are those practical problems in which the search costs play a visal role. In particular, the inequality of the second secon

Considerable sttention has been given to proving the convergence of various extreme search methods. More important, however, is the problem of finding methods which would permit optimal coordination of losses in the search for an extremum and costs due to the servetions" is fixed, this is equivalent to [inding methods of lasse strong.

Methods minimizing the maximum error for any functions of a given class may be terms minax methods. Their obvious disadvantage is emphasis on the lassf favorable conditions which are very soldom encountered in practice. This draws attention to the Bayesian going mility criterion for the case whare we are required to find a search method having the least mean error, where we average relative to an a priori distribution on the class of functions under consideration.

Methods minimizing the mean error in the above sense will, for herevity, be called Bayedian extremum-mearch methods. However, this term will be realmed even for the mote general case where we are required to minimize the mathematical expectation of the costs due to losses in searching and to deviation from the desired molution. We shall seak the minimum point. For convenience of presentation we first consider the comparatively simple case where the number of observations is fixed.

1. THE BAYESIAN EXTREMUM SEARCH METHOD WITH A FIXED NUMBER OF OBSERVATIONS

To describe Bayesian search methods it is necessary to specify an a priori probability distribution on the subsets of the class of functions for whose minimization we it and to use the method. In addition, we must determine the class of search methods from which the best is to be chosen, as well as the loss function, which determines the cost is case we find some point other than the true minimum when using some given search method.

To define Bayesian extremum search methods mathematically it is necessary to specify the following:

1. A set G of elementary events, i.e., the set of functions / to be minimized, with domain of definition $3{<}0{}^{\rm ev}$ taking on values in R, i.e., mapping X into R.

2. A family of finite-dimensional distribution functions $\mathbb{P}_{\chi(1)},\ldots,\chi_{|k|}(y_1,\ldots,y_k)=1,2,\ldots,$ satisfying the compatibility conditions [1]. Let us assume that the family fait distribution functions is defined so that for any fixed spect (t = 1,2,\ldots,k) we have the

$PII: f\{x(1)\} < y_1, ..., f(x(k)) < y_k, fund \} = F_{x(y_1,...,y_k)}(y_1, ..., y_k), k=1, 2, ...$

(1)

(2)

gers P is the a priori probability that the function values at points $x(1), \ldots, x(k)$ turn at to be less than the corresponding numbers y_1, \ldots, y_k .

If H is the set of all functions mapping X into R, then according to a theorem of ginegorov [1] Eq. (1) makes mones for any family of finite-dimensional distribution mations satisfying the compatibility conditions. Consequently, for any family of comciller finite-dimensional distribution functions there exists a random function whose priori probability distribution P satisfies condition (1). We shall minimize realizations of this random function.

If we wish to restrict the class of functions *f* to functions possessing certain proprise, then for (1) to make sense we must impose on the family of finite-dimensional matribution functions certain conditions basides those of compatibility. The conditions due which the probability measure P may be concentrated on a set of continuous functions for n = 1) are indicated in [1]. The question of the properties of realizations of ranbuse functions is considered in greater detail in [2].

Let $n_{\rm g}$ denote the value of f' at some fixed point x: $n_{\rm g}=V_{\rm g}(f)=f'({\rm x})$ where $V_{\rm g}(f)$ is a functional that maps 0 into 0. In accordance with (1), subsets of the form((0)<s.(=0)) are parameters (Consequently, $n_{\rm g}$ is a random variable (1).

3. A set δ of extremum search methods. Let us introduce some definitions. Let $\langle t \rangle$ denote the two-component vector

x(t) = (x(t), f(x(t))), t = 1, 2, ..., T, x(t) = X.

Let s_k denote a vector whose components are vectors, $s_k = \{c(1), \ldots, s(t)\}$. The set of rules of s_k will be denoted by \mathbb{S}_k . The vector s(t) contains the information obtained as a result of the t-th observation, while s_k contains the information obtained from all observations from the first to the t-th. It is clear that the extramum search must reflect the connections between the conductions of the previous observations. Therefore, we define the search method as a sequence of s_k of the previous observations. Therefore, we define the search method as a sequence of it is contained as $\{t_k = 1, 2, \ldots, 3\}$ that map \mathbb{E}_k onto X and determine the dependence of the oper observations s_k (t = $1, 2, \ldots, 3$) that map \mathbb{E}_k onto X and determine the dependence of the containets g(t = 1) on s_k :

att+13+d,120, 7+1, 2,..., Tr

The function δ_0 determines the dependence of $\kappa(1)$ on the a priori distribution \mathcal{T}_{*} for provity, the each other back with a system variable dimution \mathcal{T}_{*} , $\kappa_{*}(\alpha_{*}, \beta_{*}(\alpha_{*})) \in \mathcal{T}_{*}$ and \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable distribution of \mathcal{T}_{*} and \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} distribution \mathcal{T}_{*} and \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} distribution \mathcal{T}_{*} are a substantiable of the distribution \mathcal{T}_{*} distress distribution $\mathcal{$

According to the accepted definition of the search method 6 the functions $\delta_{\underline{k}}$ (t = 0,1,...,T) are measurable. Consequently, the functional $a_{\underline{k}}(f)$ is also measurable and the votor $x_{\underline{k}}$ is random. Let $\eta_{\underline{k}}$ denote the value of f at $x_{\underline{k}}$. Then $\eta_{\underline{k}}$ is a functional, which, for any fixed measurable and be dependence of the number $\eta_{\underline{k}}$ of the function with

$\pi_{i_0} = V_{\phi}(l) = H(x_{\phi}) = f(x_{\phi}(l)).$

The requirement of measurability of the functions δ_{\pm} septed β_{\pm} onto X is mainly of formal signification of is introduced in order to be able to speak of the mathematical expection of results detailed by the MADG4 of $(\delta_{\pm},\delta_{\pm},\ldots,\delta_{\pm})$ for each order of the E 11 for the definition of a measurable force

Nope, knowledge gradient! lol (Mockus 1972)

One of the simplifications for the solution of the equations (2) is "one-stage" method [1] [3] when at each stage it is assumed that the following observation is the last one. In such a case the sequence of observations is defined by the equations

$$E \{ u(z_n, f(x_{n+1}), x_{n+1}) | Z_n \} = \min_{x \in A} E\{ u(z_n, f(x), x) | Z_n \}$$
where

$$U(z_{n+1}) = \min_{x \in A} E\{f(x)|z_{n+1}\}, n=0,...,N.$$

The one-stage Bayesian method converges to the minimum of any continuous function under the conditions of theorem 1.



Maximum of posterior mean occurs at observation location...

As far as I can tell...

upper confidence bound? probability of improvement? expected improvement? knowledge gradient?

Harold Kushner, 1962 Harold Kushner, 1964 1962

Jonas Mockus, 1972

What about EI? (Šaltyanis, 1971)

ONE METHOD OF MULTIEXTREMUM OPTIMIZATION

V. R. Shaltyanis

Avtomatika i Vychislitel'naya Tekhnika, Vol. 5, No. 3, pp. 33-38, 1971

UDC 62-505

A nonlocal optimization method is proposed which utilizes all the information on the results of tests. The assumptions made lead to an algorithm which is optimum on average for one optimization step. Results of experimental investigations of the algorithm are given.

2. Choice of the loss function. Henceforth we will consider search for the minimum value of the target function, our assumption being that the treatment of the maximization problem will be similar. The smallest known value of the target function will be denoted by $w_p = \min_{j=\overline{1,p}} \omega_j$. The effectiveness of the effective ness of the (p + 1)-th trial will be measured by the difference $\Delta w_{p+1} = w_p - w_{p+1}$, while the average effectiveness will be measured by the mathematical expectation $M[\Delta w_{p+1}]$.

Expected Improvement (Šaltyanis, 1971)

- OU process prior on objective function
- experiments in up to 32 dimensions!
- very familiar comparison to random search...



Fig. 3. The quantity w as a function of p the number of tests p: 1) Monte Carlo method; 2) proposed method.

As far as I can tell...

upper confidence bound? probability of improvement? expected improvement? knowledge gradient?

Harold Kushner, 1962 Harold Kushner, 1964 1962 Šaltyanis, 1971 Jonas Mockus, 1972

Thank you!