

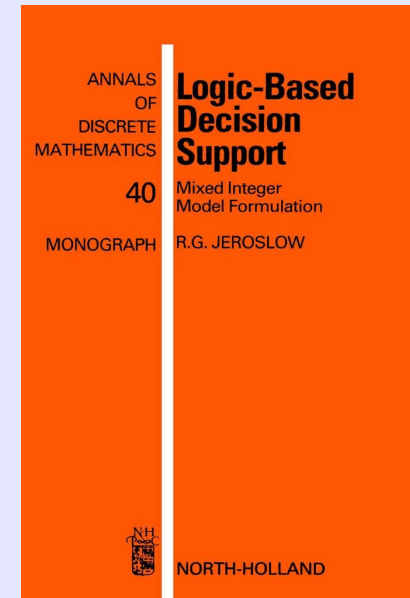
A Beautiful Paper

fabio.schoen@unifi.it
<http://gol.dinfo.unifi.it>

February 22, 2019

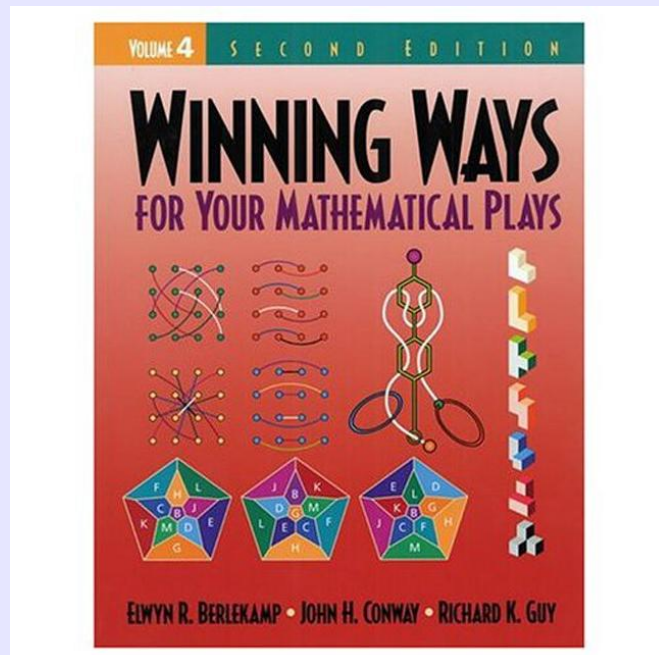
1

First attempts to find a topic: early passions



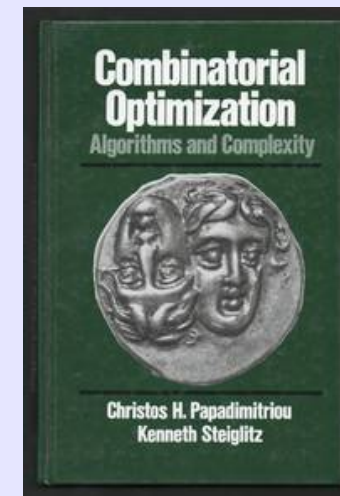
2

First attempts to find a topic: early passions

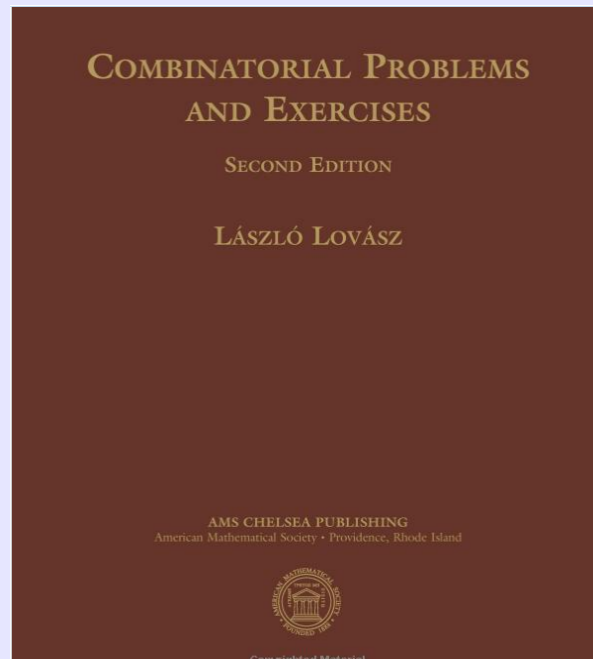


3

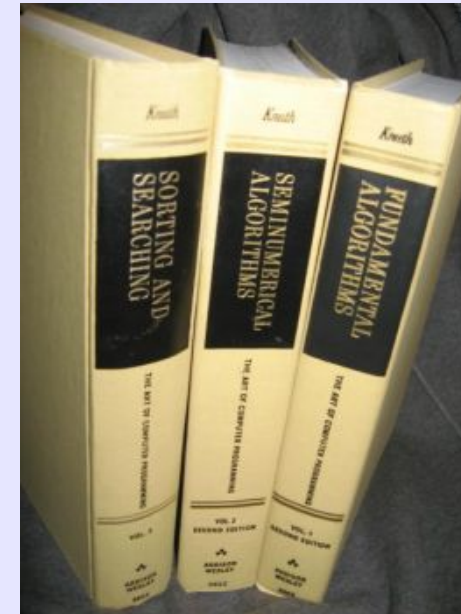
First attempts to find a topic: early passions



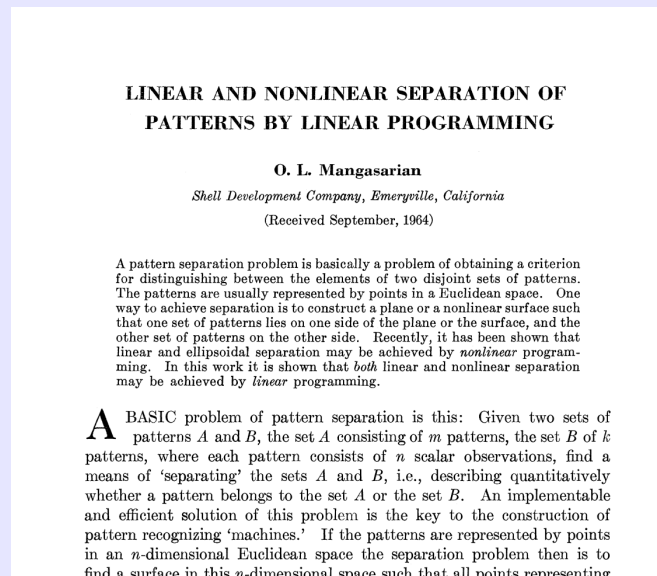
4



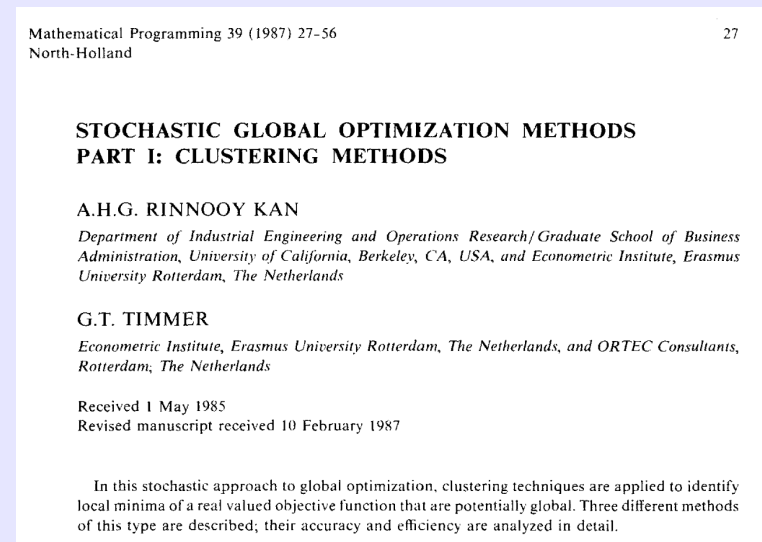
5



6



7



8



9

- Panos Pardalos
 - Reiner Horst
 - Chris Floudas
 - Manuel Bomze
 - Marco Locatelli
 - many many many others: Zelda, Gerardo, Eligius, Tibor, ...
- and ...

10

Don Jones



who...? Donald Jones, General Motors

Document title	Authors	Year	Source	Cited by
Efficient Global Optimization of Expensive Black-Box Functions	Jones, D.R., Schonlau, M., Welch, W.J.	1998	Journal of Global Optimization 13(4), pp. 455-492	2631
View abstract trova@unifi View at Publisher Related documents				
Lipschitzian optimization without the Lipschitz constant	Jones, D.R., Perttunen, C.D., Stuckman, B.E.	1993	Journal of Optimization Theory and Applications 79(1), pp. 157-181	1070
View abstract trova@unifi View at Publisher Related documents				
A Taxonomy of Global Optimization Methods Based on Response Surfaces	Jones, D.R.	2001	Journal of Global Optimization 21(4), pp. 345-383	1011
View abstract trova@unifi View at Publisher Related documents				

11



Journal of Global Optimization 13: 455-492, 1998.
© 1998 Kluwer Academic Publishers. Printed in the Netherlands.

455

Efficient Global Optimization of Expensive Black-Box Functions

DONALD R. JONES¹, MATTHIAS SCHONLAU^{2,*} and WILLIAM J. WELCH^{3,**}

¹Operations Research Department, General Motors R&D Operations, Warren, MI, USA; ²National Institute of Statistical Sciences, Research Triangle Park, NC, USA; ³Department of Statistics and Actuarial Science and The Institute for Improvement in Quality and Productivity, University of Waterloo, Waterloo, Ontario, Canada

(Accepted in final form 30 June 1998)

Abstract. In many engineering optimization problems, the number of function evaluations is severely limited by time or cost. These problems pose a special challenge to the field of global optimization, since existing methods often require more function evaluations than can be comfortably afforded. One way to address this challenge is to fit response surfaces to data collected by evaluating the objective and constraint functions at a few points. These surfaces can then be used for visualization, tradeoff analysis, and optimization. In this paper, we introduce the reader to a response surface methodology that is especially good at modeling the nonlinear, multimodal functions that often occur in engineering. We then show how these approximating functions can be used to construct an efficient global optimization algorithm with a credible stopping rule. The key to using response

12

Why did I choose this paper? For many reasons:

- it is very well and very clearly written
- it received an impressive number of citations (not only from machine learning)
- it is a fascinating idea (although not his own and although subject to many improvements) - behind its complications the very essence of GO is contained: automatically mixing exploration and exploitation
- I did not mention Don Jones in my GO book (I just cited the origins of the method (see next slide)) and I took this opportunity to remedy...

13

The problem

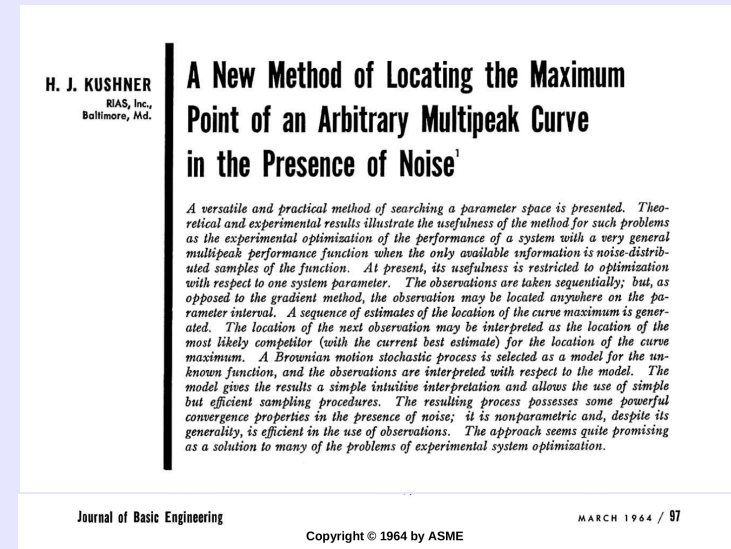
$$\min_{x \in S \subset \mathbb{R}^n} f(x)$$

where

- we are interested in (an approximation of) the global optimum
- S is a “simple” feasible set (e.g., the unit box)
- evaluating f at any feasible x is extremely expensive

15

Everything started and has its foundations in:



and generalized by Jonas Mockus.

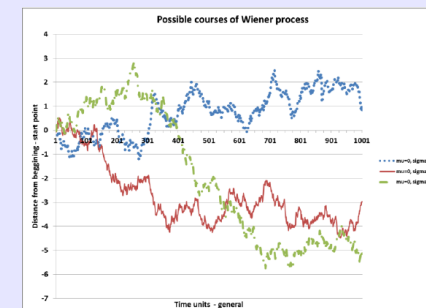
14

The model

Main idea: consider the objective function f as a realization of a stochastic process

Most frequently assumed model: a multidimensional Wiener process. In \mathbb{R}^1 a stochastic process $W()$ is Wiener if:

- W has independent increments: for any t and $u \geq 0$
 $W_{t+u} - W_t$ is independent of W_s for all $s < t$
- $W_{t+u} - W_t \sim \mathcal{N}(0, \sigma^2 u)$
- W is almost surely continuous



16

It is assumed that

$$f(x) = \mu + \varepsilon(x)$$

where:

- μ is an unknown constant
- $\varepsilon(x)$ is Gaussian with zero mean and variance σ^2
- $\text{Corr}(\varepsilon(x^i), \varepsilon(x^j)) = \exp(-d(x^i, x^j))$
- $d(x, y)$ is a generalized weighted Euclidean distance:

$$d(x, y) = \sum_{k=1}^n \theta_k |x_k^i - x_k^j|^{p_k}$$

where θ, p are parameters

17

Assume N observations are available $\{x^i, y^i\}_{i=1}^N$ with $y^i = f(x^i)$.

Let R be the Correlation matrix for these observations:

$$R_{ij} = \text{Corr}(x^i, x^j).$$

Then the likelihood turns out to be:

$$\mathcal{L}(\mu, \sigma^2, \theta, p) = \frac{1}{(2\pi\sigma^2)^{N/2} |R|^{1/2}} \exp\left(-\frac{(y - \mathbb{1}\mu)^T R^{-1} (y - \mathbb{1}\mu)}{2\sigma^2}\right)$$

and the maximum likelihood estimate of μ and σ^2 is:

$$\hat{\mu} = \mathbb{1}^T R^{-1} y / \mathbb{1}^T R^{-1} \mathbb{1}$$

$$\hat{\sigma}^2 = (y - \mathbb{1}\hat{\mu})^T R^{-1} (y - \mathbb{1}\hat{\mu}) / N$$

while $\hat{\theta}, \hat{p}$ can be found maximizing

$$\mathcal{L}(\hat{\mu}, \hat{\sigma}^2, \theta, p)$$

18

Prediction

Given x , where f has not been evaluated yet, the model can be used to predict $f(x)$.

Let r the correlation vector between x and the observations x^i .

Then the best unbiased estimator for $f(x)$ is

$$\hat{y}(x) = \hat{\mu} + r^T R^{-1} (y - \mathbb{1}\hat{\mu})$$

$$= \hat{\mu} + c^T r(x)$$

where

$$c = R^{-1} (y - \mathbb{1}\hat{\mu})$$

$$r_i(x) = \text{Corr}(x, x^i)$$

while the variance of this estimate will be

$$s^2(x) = \hat{\sigma}^2 (\mathbb{1} - r^T R^{-1} r)$$

19

Prediction: Branin function

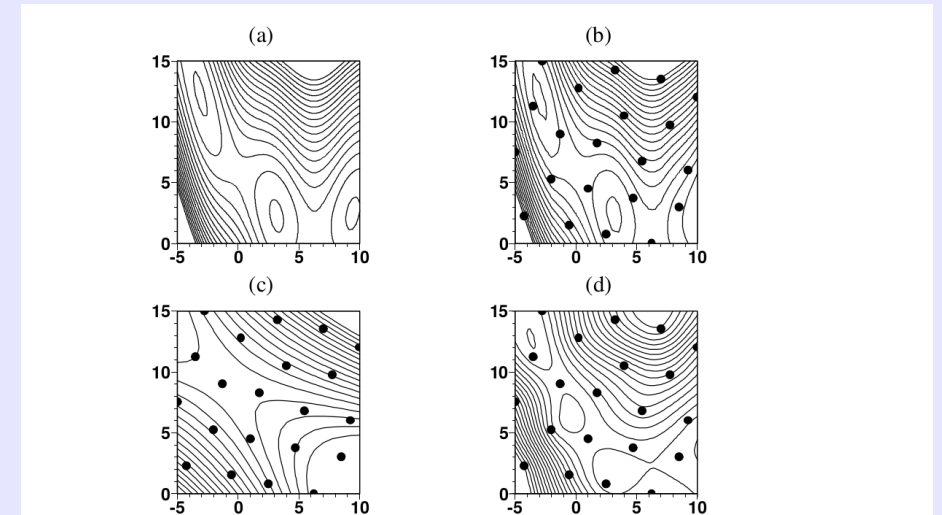


Figure 3. (a) Contours of the Branin test function; (b) contours of a DACE response surface based on the 21 sampled points shown as dots; (c) a quadratic surface fit to the 21 points; (d) a thin-plate spline fit to the 21 points.

20

First idea: evaluate f at the global minimizer of the predictor.
 Defects: greedy, stalling

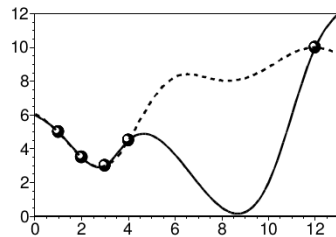


Figure 8. The solid line represents an objective function that has been sampled at the five points shown as dots. The dotted line is a DACE predictor fit to these points.

21

We should take variance into account

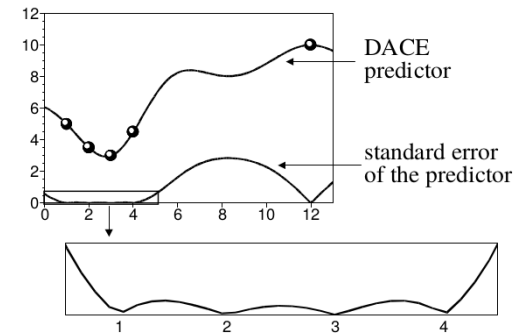
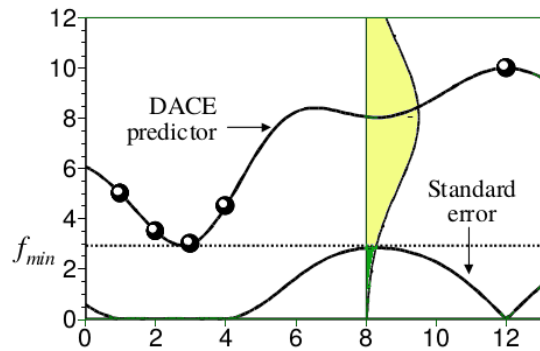


Figure 9. The DACE predictor and its standard error for a simple five-point data set.

Good choices should take into account good observations ([exploitation](#)) and unexplored regions ([exploration](#))

22

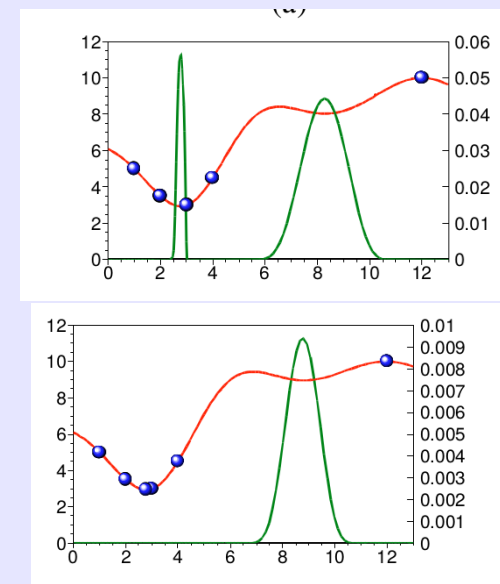


$$E[I(x)] = E[\max\{0, f_{\min} - f(x)\}]$$

$$= (f_{\min} - \hat{y}(x))\Phi((f_{\min} - \hat{y}(x))/s) + s\phi((f_{\min} - \hat{y}(x))/s)$$

with Φ, ϕ : standard normal density and CDF.

23



24

A Beautiful Paper

`fabio.schoen@unifi.it`
`http://gol.dinfo.unifi.it`

February 22, 2019

Grazie Marco!