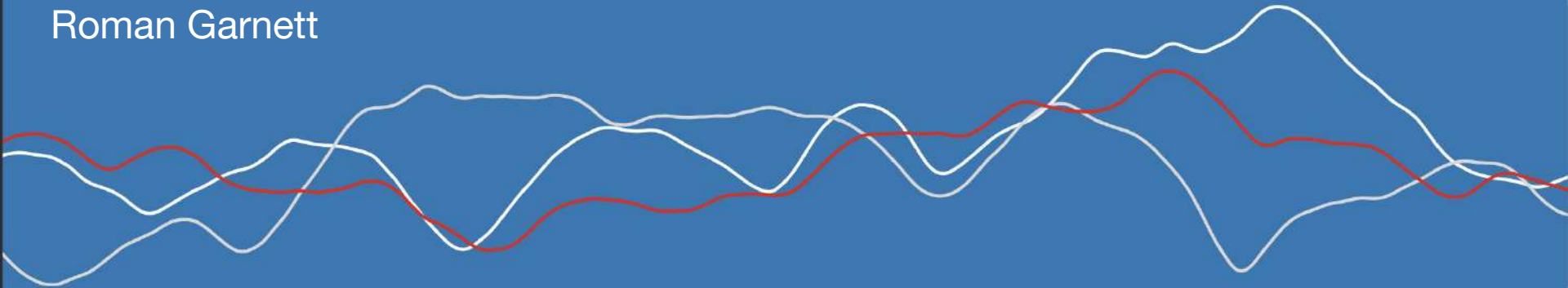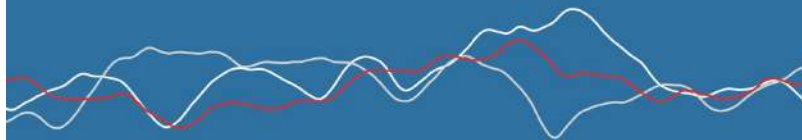# Some random things I learned writing the BayesOpt book

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>
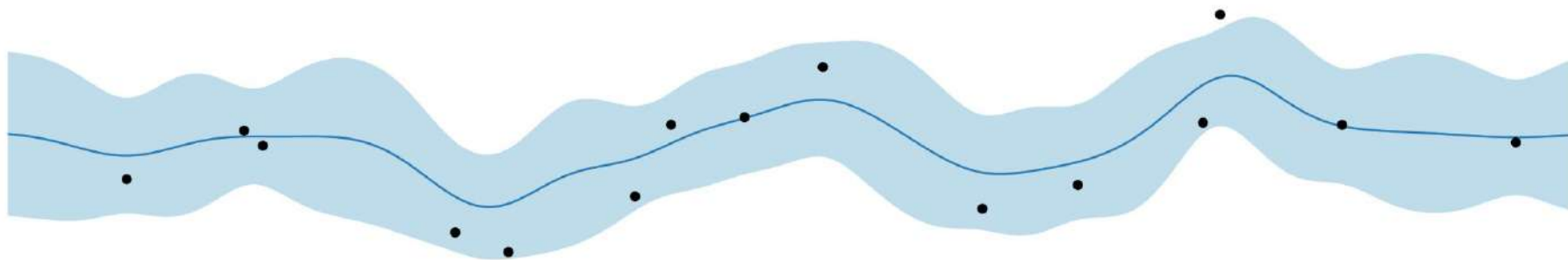to me, Michael, Frank, Nando

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

# Expected Improvement with Noise
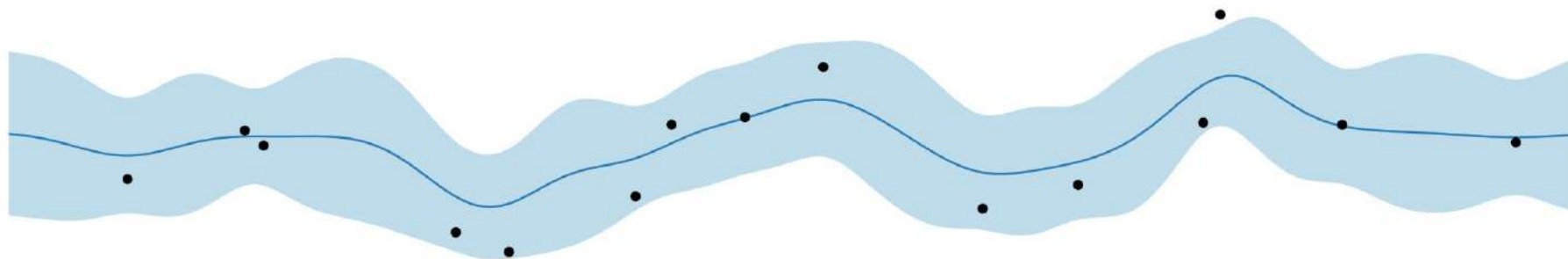
# Simple Recipe for Bayesian Experimental Design

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

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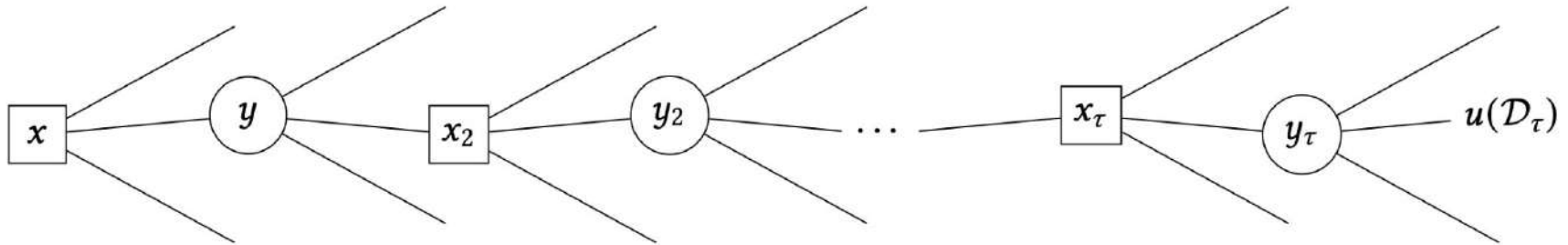- observation model, $p(y \mid x, \phi)$, $\phi = f(x)$    e.g., Gaussian noise

# Simple Recipe for Bayesian Experimental Design

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

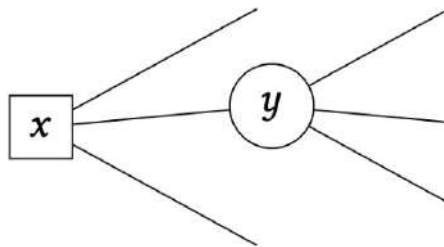# Simple Recipe for Bayesian Experimental Design

Step 3: give up on the optimal policy

# Simple Recipe for Bayesian Experimental Design

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...)

$u(\mathcal{D}_\tau)$

# Simple Recipe for Bayesian Experimental Design

Step 4: derive a policy via one-step lookahead

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

(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):

$$u(D) = \phi^* = \max \mathbf{f}$$



$p(\phi \mid x, \mathcal{D})$

$\max(\phi - \phi^*, 0)$

$\phi$

$\phi^*$

marginal gain:

$$\max(\phi - \phi^*, 0)$$

expected utility easy to compute, has nice properties, etc.
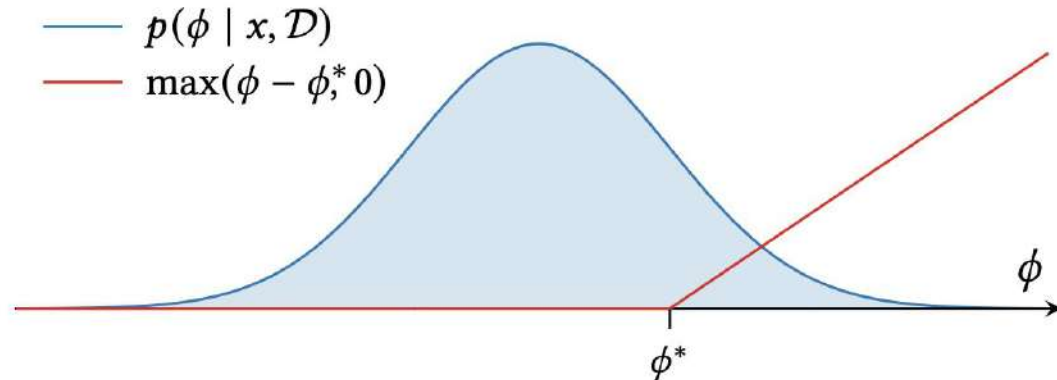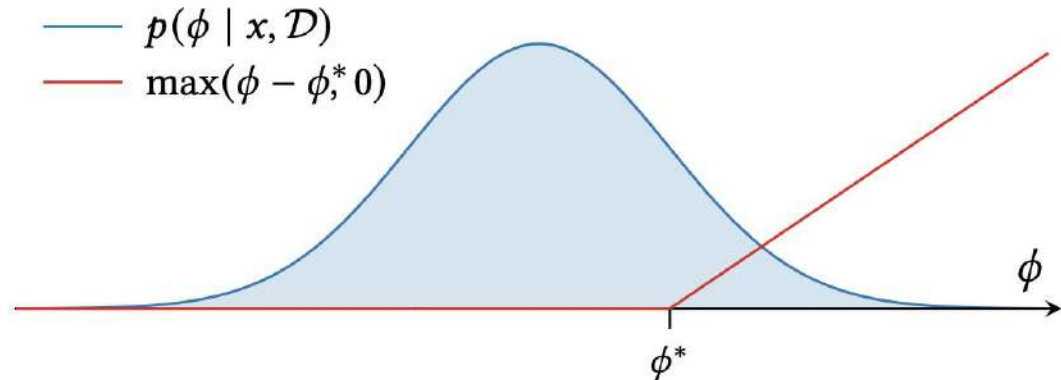
# Noiseless expected improvement

utility (best seen value):

$u(D) = \phi^* = \max \mathbf{f}$

marginal gain:

$\max(\phi - \phi^*, 0)$

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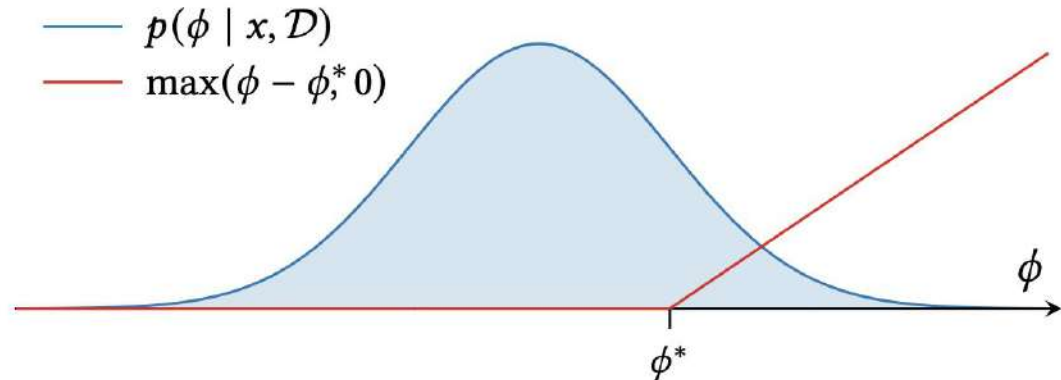


Legend:
— $p(\phi \mid x, \mathcal{D})$
— $\max(\phi - \phi^*, 0)$

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

# Noiseless expected improvement

expected utility easy to compute, has nice properties, etc.
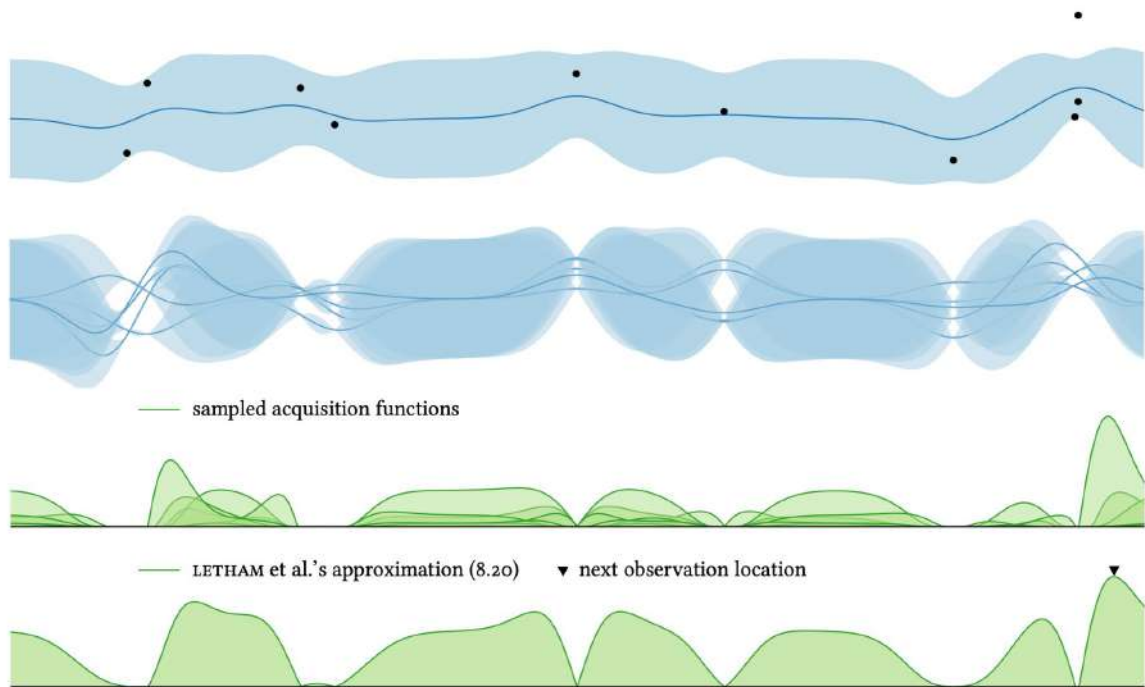
very tempting to start here and try to "fix" this!



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

# "Fixing" the expected utility

- plug-in estimators: use noiseless EI with "guess" of max **f**
- expectation of EI with respect to **f** (Letham, et al. 2019)



sampled acquisition functions

LETHAM et al.'s approximation (8.20)    ▼ next observation location
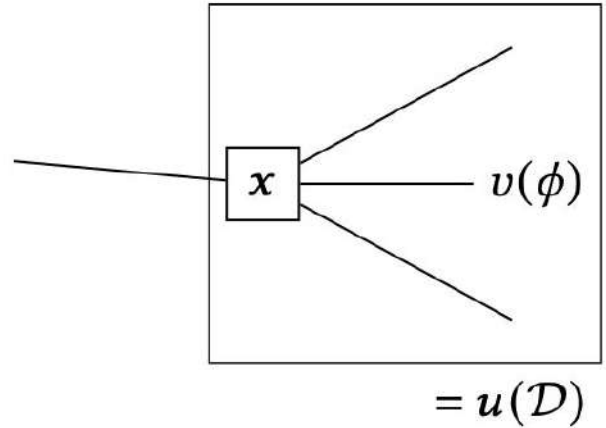
# Let's start with utility!

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 **x** = $u(D)$



$$= u(\mathcal{D})$$

# The noisy setting: Utility

maximum of posterior mean on **x** = $u(D)$

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



$$= u(\mathcal{D})$$

# 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

$$\sum_i a_i \big[ \Phi(c_{i+1}) - \Phi(c_i) \big] + b_i \big[ \phi(c_i) - \phi(c_{i+1}) \big]$$

# The result

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

# Alternative approaches



- plug-in estimate, $\phi^* \approx \max y$ (8.17)  ▾ next observation location
- LETHAM et al.'s approximation (8.20)  ▾ next observation location
- expected improvement (8.16)

# 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 $\theta$ 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 $\theta$ and integrate the hyprameter-conditional acquisition function against the hyperparameter posterior

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

blind to uncertainty in $\theta$!

# Standard approach

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

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

blind to uncertainty in $\theta$!

# Alternative approach

Define utility with respect to marginal model from the beginning!

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

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

# Example

- function is $f(x) = x$ or $f(x) = -x$
- knowledge gradient
- for standard approach, acquisition function is flat! (maximum of $\theta$-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?        Harold Kushner, 1964

expected improvement?              Jonas Mockus, 1972

upper confidence bound?            Cox and John, 1998

knowledge gradient?                Frazier, et al., 2009

# I was wrong!

probability of improvement?          Harold Kushner, 1964

expected improvement?               ~~Jonas Mockus, 1972~~

upper confidence bound?             ~~Cox and John, 1998~~

knowledge gradient?                 ~~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[1]

*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)

- 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)

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

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

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

is maximum.

# As far as I can tell...

upper confidence bound?          Harold Kushner, 1962

probability of improvement?      Harold Kushner, ~~1964~~ 1962

expected improvement?

knowledge gradient?

# Further Development (Kushner, 1964)

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

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[1]

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
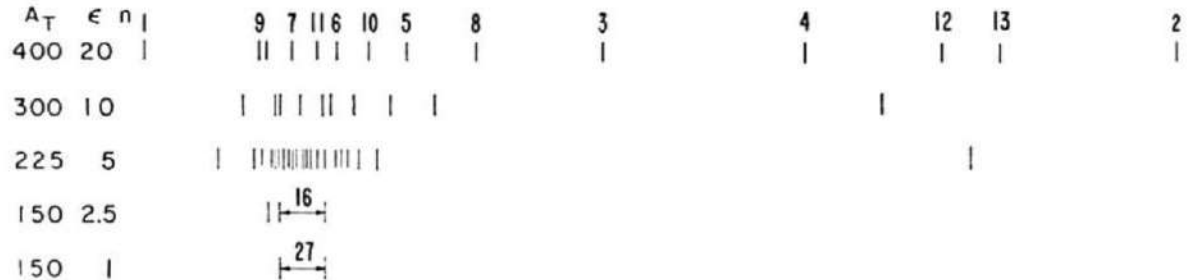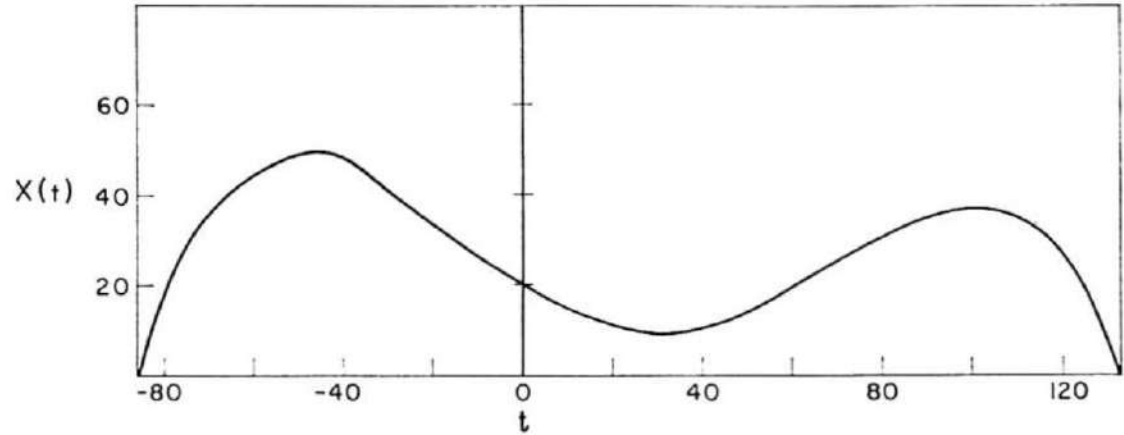- careful scheduling of improvement thresholds



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

# Very thoughtful! (Kushner, 1964)

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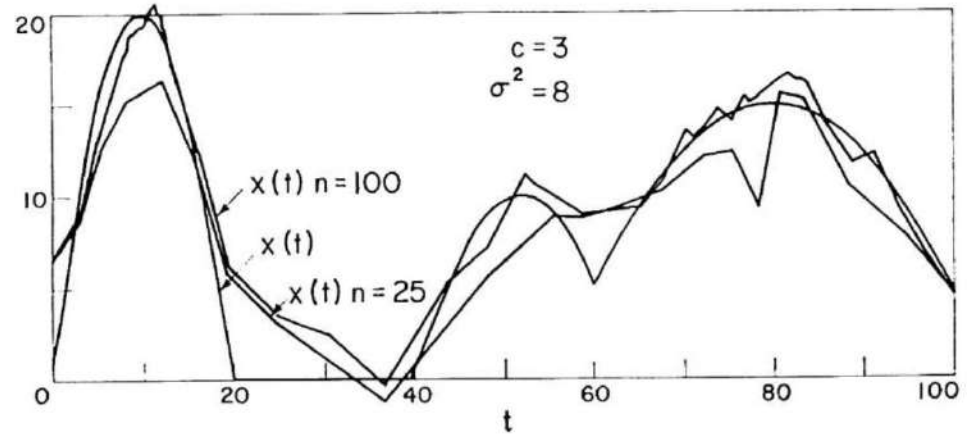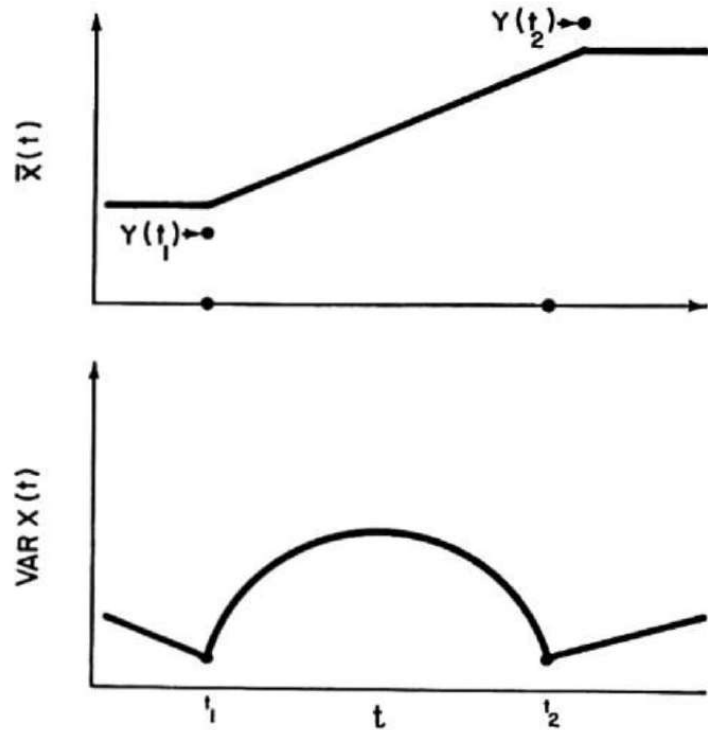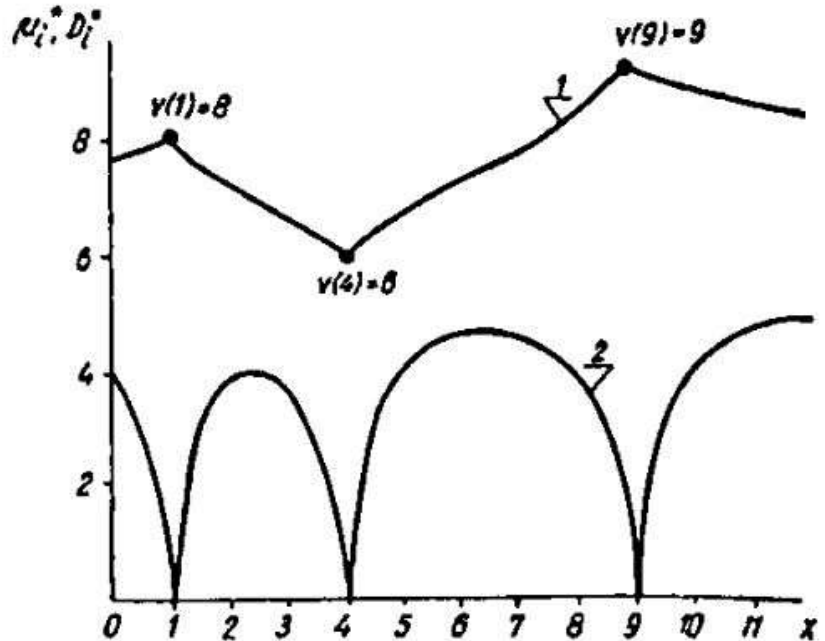- scheduling of improvement thresholds
- handling noise



Fig. 7 Experimental results with observation noise
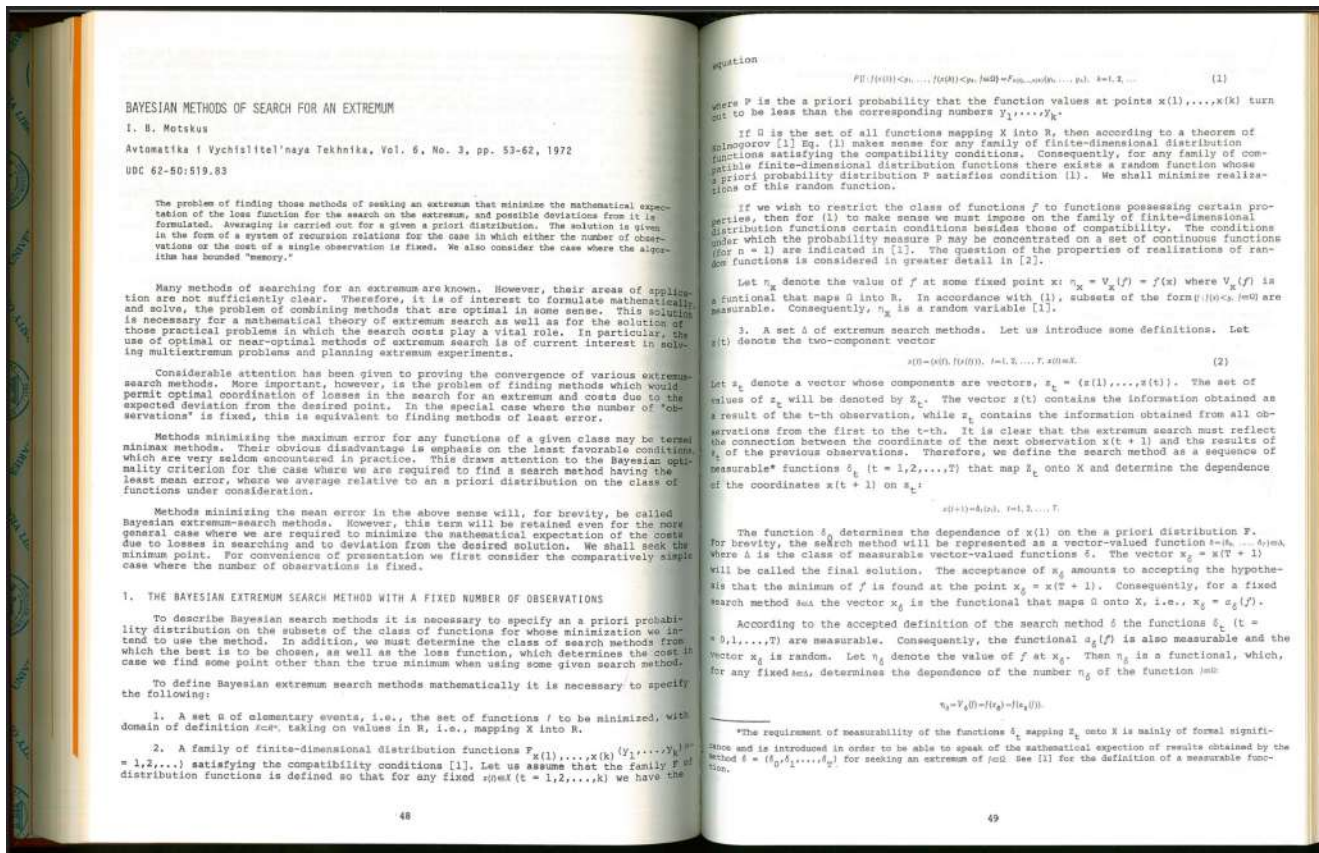
# Aside: Gauss-Markov models



Wiener process (Kushner)

OU process (Šaltyanis)

# Okay but this is EI right? (Mockus 1972)

# 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
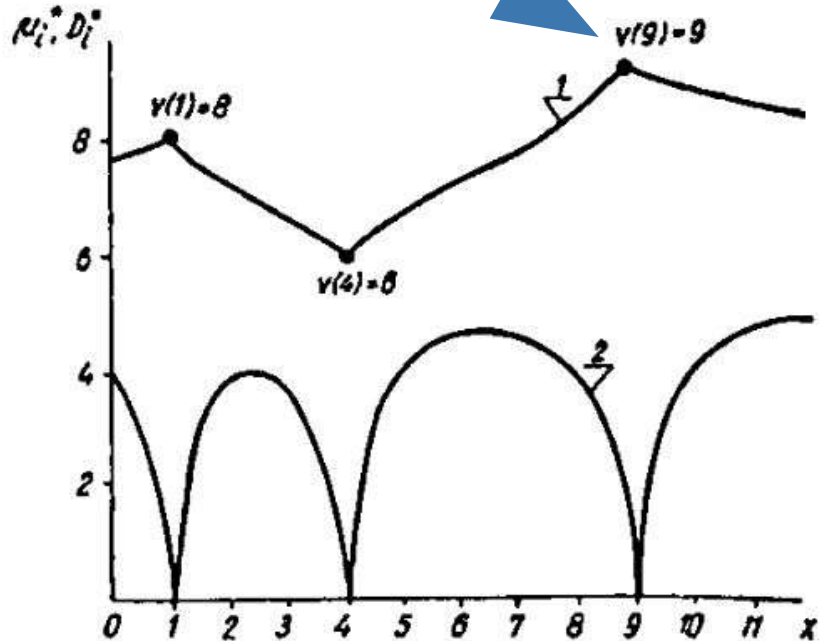
$$E\left\{u\left(z_n, f(x_{n+1}), x_{n+1}\right) | z_n\right\} = \boxed{\min_{x \in A}} E\left\{u\left(z_n, f(x), x\right) | z_n\right\}$$

where

$$u(z_{n+1}) = \boxed{\min_{x \in A}} E\left\{f(x) | z_{n+1}\right\}, \quad n = 0, \ldots, N.$$

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

# Equivalent to EI for G-M



Maximum of posterior mean occurs at observation location...

# As far as I can tell...

upper confidence bound?  Harold Kushner, 1962

probability of improvement?  Harold Kushner, ~~1964~~ 1962

expected improvement?

knowledge gradient?  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 effectiveness 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!
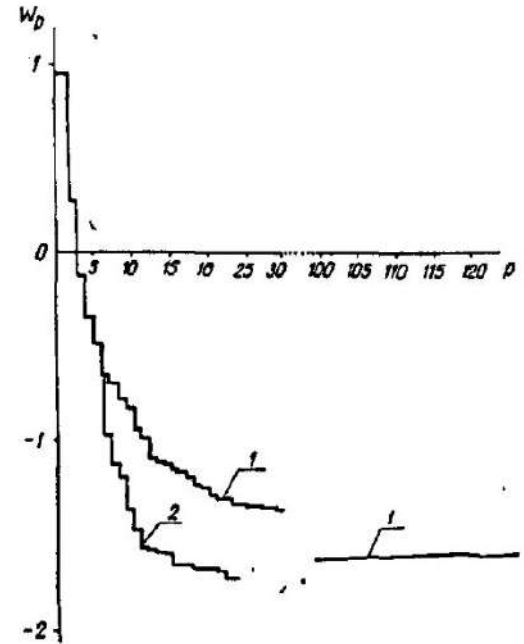- very familiar comparison to random search...



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

# As far as I can tell...

upper confidence bound?      Harold Kushner, 1962

probability of improvement?      Harold Kushner, ~~1964~~ 1962

expected improvement?      Šaltyanis, 1971

knowledge gradient?      Jonas Mockus, 1972

# Thank you!