

Continue : Bayesian optimization

Given a set of possible algorithms $A = \{A_1, A_2, \dots, A_n\}$

$\forall A_i \rightsquigarrow \Lambda^{(i)} \rightsquigarrow \underline{\lambda} \in \Lambda^{(i)}$

AutoML : solve the opt. problem.

$$A_{\lambda^*} = \underset{A, \lambda}{\operatorname{arg\,min}} \mathcal{L}(A, D_{\text{train}}, D_{\text{valid}})$$

very expensive opt problem.

evaluating $\mathcal{L}(\cdot)$ or finding $\nabla \mathcal{L}(\cdot)$

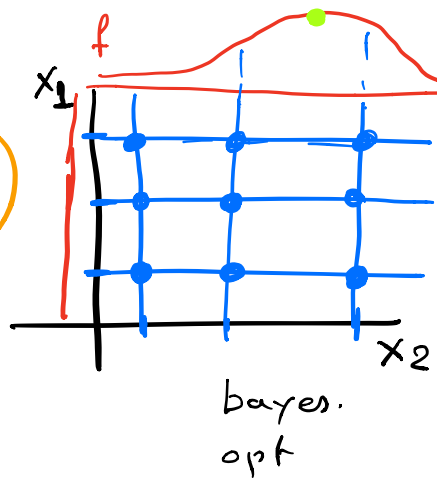
Blackbox optimization

Discretize ..

Every coordinate exponential in \bar{x} and perform exhaustive search

→ is not fine enough (I cannot find optimal solution)

grid search



$f(x)$

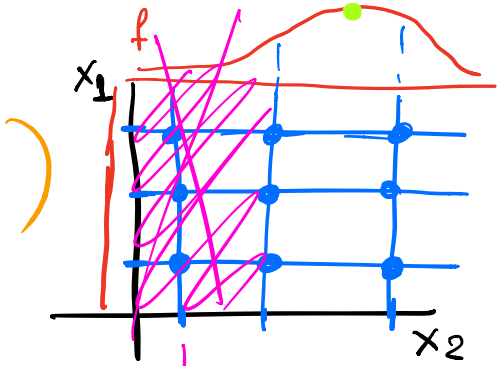
\bar{x} :

Some of the components are discrete; some are continuous; some are categorical.

find

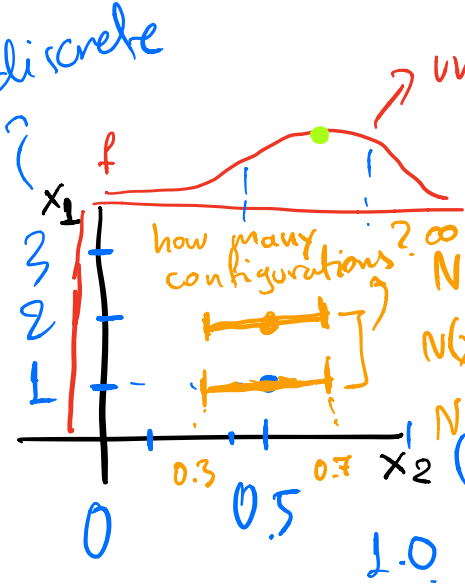
$x^* = \underset{x}{\operatorname{arg\,min}} f(x)$

Random Search



low value for f .

discrete



how many configurations? ∞

$N(1, 0.5)$?

$N(x_1) = \{x_1 - 1, x_1 + 1, x_1 + 3\}$

$N(x_2) = [x_2 \pm \delta]$

(continuous)

$x_1 \in \{1, 2, 3\}$

$x_2 \in [0, 1]$

(x_1^1, x_2^1)

$= (1, 0.5)$

Random search:
avoid evaluating low-value areas

→ start from a random (x_1^1, x_2^1)

→ get the neighborhood of (x_1^1, x_2^1)

→ evaluate f only

in that neighborhood

→ take (x_1^2, x_2^2) to

be the point with min f

exp. und wasteful

$N(1, 0.5) \rightarrow ?$ edit, reconfigure (x_1, x_2)

$0.5 \rightarrow 0.6$ (step?)

known

$$f(\bullet) = f(1, 0.5)$$

(2, 0.5)

(1, 0.6) x_2 is

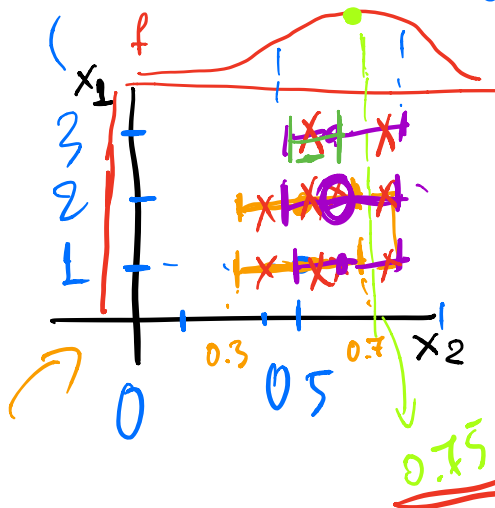
(1, 0.4) cont

$N(x_1, x_2)$: a "ball" around x_1, x_2
distance in the x_1 coordinate
" " " x_2 "

Take the product of the two sets of points that these two define.

Problem: the neighborhood $N(\bar{x})$ might require us to perform a large # of evaluations.

Attacks: (1) pick a random point in N
Idea 1: multiple random points (k) \leftarrow ?
evaluate \rightarrow convergence is not guaranteed.



step 2

Best guarantee local optimum.
w.o. knowing anything about f

(2) \rightarrow Adaptively change δ
 $\delta \rightarrow$ make it smaller.

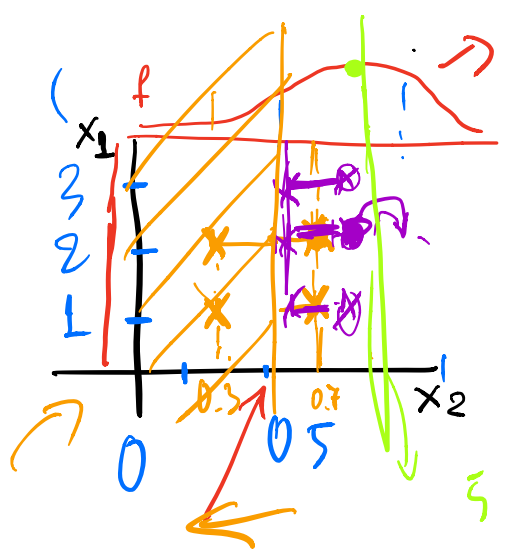
Question : what is the property that we should exploit? (monotonicity)
 if f is ^{locally} monotonic wrt x_i ?

Combine Ideas

- ① evaluate f on a subset of points in N
- ② adaptively decrease δ
- ③ hope that we locally monotonic \rightarrow evaluate boundary points.

New idea : (exploration) pick a new random point.

helps in this scenario



~~maximize~~

What does pruning exploit?

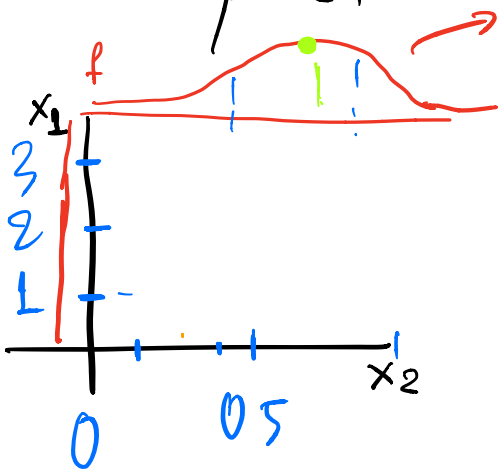
memory

Bayesian optimization

Key idea is "have memory"

we approximate f with some pdf and we will use inference to find the next best point to evaluate considering all previous evaluations.

assume a Gaussian (pdf) ^{smooth, global opt.}



what does this look like?

Gaussian.

$$x^* = \underset{x}{\operatorname{arg\,max}} f(x)$$

$x^* = ?$ assuming it's a Gaussian?

$$x^* = \mu$$

Steps of BO

Assume \rightarrow a model P such that P approximates

$$x^* = \underset{x}{\operatorname{arg\,max}} P(f(x) \mid f(x_1), f(x_2), \dots, f(x_n))$$

Posterior analysis.

① What is P?

→ fit this model (estimate P) (learning Problem)

x_1, x_2, \dots, x_n

→ $\operatorname{argmax}_x P(\cdot)$ vs $a(P(\cdot))$

option 1 is use $P(\cdot)$ directly

option 2 is use a function of $P(\cdot)$

BO allows us to fight local minima | acquisition → exploration

→ Consider a gaussian prior $\approx P$

→ as i am collecting evidence on the form of the real pdf

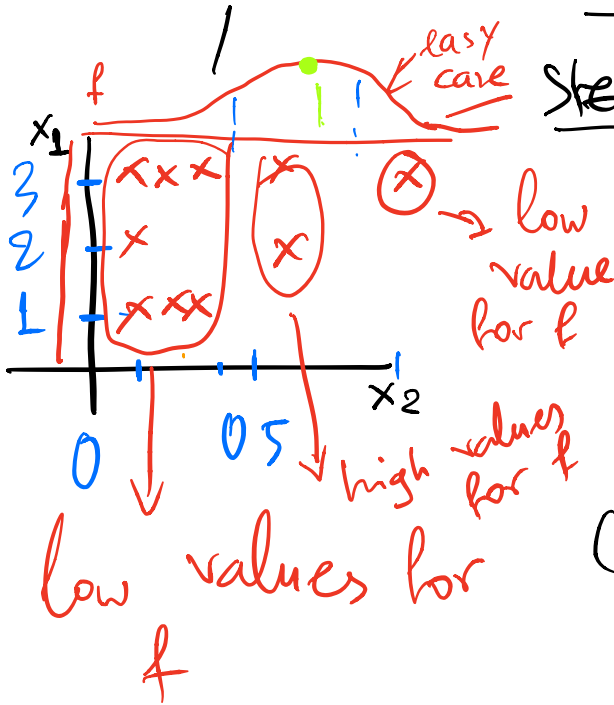
I am revising my belief on P

→ I need to use an acquisition function

a that takes as input my current belief on P and guides exploration of the

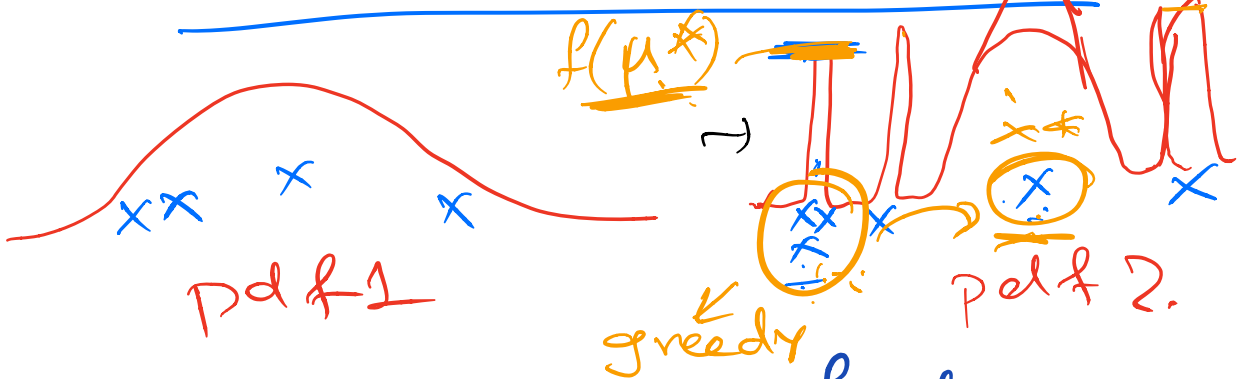
>pace for \bar{x}

Algorithm for BO



Step 1: Take 10 samples x_1, x_2, \dots, x_{10} and evaluate f over them.
(uniform)

One type of acquisition function
→ maximize exp. improvement



Acquisition function:
control exploration

→ Expected improvement. Y_{best} : running best solution

$$E \left[\min (f(x^*) - Y_{best}, 0) \right]$$

$$\frac{(f(x^*) - \mu(x^*))}{\sigma(x^*)} - \frac{Y_{best} - \mu(x^*)}{\sigma(x^*)}$$

→ An acquisition function that promotes exploration

→ for different regions of x
if they are not explored

confidence on what $f(x)$
is in reality is low (I have
a small
number of
samples)

assign a positive score
to these regions.

$$f(x_*) | f(x_1) = y_1, \dots, f(x_N) = y_N$$

$$N \left(\underbrace{k_*^T}_{\mu} \underbrace{\Sigma^{-1}}_{\sigma^2}, \underbrace{k(x_*, x_*)}_{\sigma^2} - \underbrace{k_*^T \Sigma^{-1} k_*}_{\sigma^2} \right)$$

$$k_* = \begin{bmatrix} k(x_1, x_*) \\ k(x_2, x_*) \\ \vdots \\ k(x_N, x_*) \end{bmatrix}$$

vector of evaluation.

k: kernel.
update this

Σ = Covariance

$$\begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) \\ \vdots & \vdots \\ k(x_N, x_1) & \dots \end{bmatrix}$$