## Navigating the Unknown: Intelligent and Efficient Solution Discovery via Bayesian Optimization

#### Wang MA

24 Summer Seminar - Week 1

maw6@rpi.edu

July 22, 2024

#### Contents

#### 1. Introduction

2. An Overview of Bayesian Optimization

3. A model of the function: Gaussian Process

4. Acquisition Functions

5. Conclusion

#### **Problem Setup**

Bayesian Optimization is a class of machine-learning-based optimization methods focused on sloving the problem

$$\hat{\mathbf{x}} = \max_{\mathbf{x} \in \mathcal{A}} f(\mathbf{x}),\tag{1}$$

where the feasible set and objective function typically follow some assumptions/properties.

- The input *x* is in  $\mathbb{R}^d$ , where typically  $d \leq 20$ .
- The feasible set *A* is a simple set, e.g., box constraints (hyper-rectangle) or a simplex.
- *f* is continuous but lacks special structure, e.g., non-convex.
- *f* is derivative-free, evaluations do not give gradient info.
- f is expensive to evaluate.
- f may be noisy.

### Scenarios Suitable for Bayesian Optimization

- Expensive function evaluations: the objective function is costly to evaluate, such as requiring significant computational resources or time
- Black-box functions: the objective function is a black-box with no explicit analytical form and no available gradient information.
- Hyperparameter Optimization: hyperparameter tuning in machine learning models, such as optimizing learning rates, regularization parameters, etc..

Optimization of expensive functions ariese in

- fitting machine learning models
- tuning algorithms via backtesting
- optimizing physics-based models
- drug and materials discovery

#### Contents

#### 1. Introduction

#### 2. An Overview of Bayesian Optimization

3. A model of the function: Gaussian Process

4. Acquisition Functions

5. Conclusion

## **Bayesian Optimization**

#### Algorithm 1 BayesOpt

Assume a Bayesian prior on f

(usually a Gaussian process prior, a probabilistic model of the function)

#### while budget is not exhausted do

Find *x* that maximizes **acquisition function** (*x*, *posterior*)

Sample *x* and observe f(x)

Update the posterior distribution on *f* 



7/48







#### Contents

1. Introduction

2. An Overview of Bayesian Optimization

#### 3. A model of the function: Gaussian Process

4. Acquisition Functions

5. Conclusion

### Gaussian Process: Definition

A GP is fully specified by its mean function m(x) and covariance function k(x, x'). When we model our function as  $f(x) \sim GP(m(x), k(x, x'))$ , we are saying that

- Mean function:  $m(x) = \mathbb{E}[f(x)]$
- Covariance function:  $k(x, x') = \mathbb{E}[(f(x) m(x))(f(x') m(x'))]$

The mean function m(x) often assumed to be constant or zero for simplicity, say, m(x) = 0. Covariance functions (kernels) decreases with |x - x'||, commonly used k(x, x'):

- Squared Exponential (RBF):  $k(x, x') = \sigma_f^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$
- Matern:  $\mathbf{k}(\mathbf{x}, \mathbf{x}') = \sigma_{\mathbf{f}}^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}|\mathbf{x}-\mathbf{x}'|}{\ell}\right)^{\nu} \mathbf{K}_{\nu} \left(\frac{\sqrt{2\nu}|\mathbf{x}-\mathbf{x}'|}{\ell}\right)$
- Rational Quadratic:  $k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \left(1 + \frac{(\mathbf{x} \mathbf{x}')^2}{2\alpha\ell^2}\right)^{-\alpha}$

#### **Gaussian Process**

#### **Gaussian Process**

A Gaussian Process is a collection of random variables, any finite number of which have a joint Gaussian distribution.

As a concrete example, let's choose:

$$\boldsymbol{m}(\boldsymbol{x}) = 0 \tag{2}$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^{T}(\mathbf{x} - \mathbf{x}')).$$
 (3)

Given observations  $\mathbf{f} = [f(x_1), f(x_2), ..., f(x_t)]$ , we would like to make a **prediction** at a new point  $x^*$ . According to the GP prior,  $f(x^*)$  is jointly normally distributed with  $\mathbf{f}$  so that

$$\Pr\left(\begin{bmatrix}\boldsymbol{f}\\\boldsymbol{f}^*\end{bmatrix}\right) = \operatorname{Norm}\left(\boldsymbol{0}, \begin{bmatrix}\boldsymbol{K}[\boldsymbol{X},\boldsymbol{X}] & \boldsymbol{K}[\boldsymbol{X},\boldsymbol{x}^*]\\\boldsymbol{K}[\boldsymbol{x}^*,\boldsymbol{X}] & \boldsymbol{K}[\boldsymbol{x}^*,\boldsymbol{x}^*]\end{bmatrix}\right).$$
(4)

#### Gaussian Process: Predicting

Given the jointly normal property, we have that

$$\Pr(\mathbf{f}^*|\mathbf{f}) = \operatorname{Norm}(\mu[\mathbf{x}^*], \sigma^2[\mathbf{x}^*]),$$
(5)

where

$$\mu[\mathbf{x}^*] = \mathbf{K}[\mathbf{x}^*, \mathbf{X}] \mathbf{K}[\mathbf{X}, \mathbf{X}]^{-1} \mathbf{f}$$
(6)

$$\sigma^{2}[\mathbf{x}^{*}] = \mathbf{K}[\mathbf{x}^{*}, \mathbf{x}^{*}] - \mathbf{K}[\mathbf{x}^{*}, \mathbf{X}]\mathbf{K}[\mathbf{X}, \mathbf{X}]^{-1}\mathbf{K}[\mathbf{x}^{*}, \mathbf{x}^{*}].$$
(7)

Using above formula, we can estimate the distribution at any new point  $\mathbf{x}^*$ . The best estimate of the funciton value is given by the mean  $\mu[\mathbf{x}]$ , and the uncertainty is given by the variance  $\sigma^2[\mathbf{x}]$ .

#### **Gaussian Process Model**



15/48

#### Contents

1. Introduction

- 2. An Overview of Bayesian Optimization
- 3. A model of the function: Gaussian Process

#### 4. Acquisition Functions

5. Conclusion

#### Basic Concepts

## **Acquisition Functions**

- Acquisition functions guide the selection of the next point to evaluate by balancing exploration and exploitation.
- Common acquisition functions include:
  - Probability of Improvement (PI)
  - Expected Improvement (EI)
  - Upper Confidence Bound (UCB)

## Upper Confidence Bound (UCB)

- UCB selects points based on an upper confidence bound on the surrogate model's predictions.
- Mathematically defined as:

$$\mathsf{UCB}[\mathbf{x}^*] = \mu[\mathbf{x}^*] + \kappa \sigma[\mathbf{x}^*],\tag{8}$$

where  $\kappa$  is a parameter that controls the trade-off between exploration and exploitation.

 This favors either (i) regions where is μ[x\*] large (for exploitation) or (ii) regions where σ[x\*] is large (for exploration). The positive parameter κ trades off these two tendencies.

## Probability of Improvement (PI)

- PI aims to maximize the probability that the next sample will improve over the current maximum.
- Mathematically defined as:

$$\mathsf{PI}[\mathbf{x}^*] = \int_{\mathbf{f}[\hat{\mathbf{x}}]}^{\infty} \mathsf{Norm}_{\mathbf{f}[\mathbf{x}^*]}[\mu[\mathbf{x}^*], \sigma[\mathbf{x}^*]] d\mathbf{f}[\mathbf{x}^*],$$
(9)

where  $f[\hat{x}]$  is the current maximum.

This acquisition function computes the likelihood that the function at  $x^*$  will return a result higher than current maximum. For each point  $x^*$ , we integrate the part of the associated normal distribution that is above the current maximum.

#### Basic Concepts

## Expected Improvement (EI)

- El balances exploration and exploitation by considering both the magnitude and the probability of improvement.
- Mathematically defined as:

$$\mathsf{EI}[x^*] = \mathsf{E}_n[(f[x^*] - f[\hat{x}])^+]$$

$$= \int_{f[\hat{x}]}^{\infty} (f[x^*] - f[\hat{x}]) \mathsf{Norm}_{f[x^*]}[\mu[x^*], \sigma[x^*]] df[x^*].$$
(10)
(11)

Expected Improvement (EI) takes into account how much the improvement will be, so that we can find the favorable larger improvement. A closed form of EI:

$$\mathsf{EI}[\mathbf{x}^*] = [\Delta(\mathbf{x}^*)]^+ + \sigma(\mathbf{x}^*)\varphi\left(\frac{\Delta(\mathbf{x}^*)}{\sigma(\mathbf{x}^*)}\right) - |\Delta(\mathbf{x}^*)|\Phi\left(-\frac{|\Delta(\mathbf{x}^*)|}{\sigma(\mathbf{x}^*)}\right),$$
(12)

where  $\Delta(\mathbf{x}^*) = \mu(\mathbf{x}^*) - \mathbf{f}[\hat{\mathbf{x}}].$ 

Basic Concepts

## Expected Improvement: Exploration ( $\sigma[x^*]$ ) vs. Exploitation ( $\Delta[x^*]$ )



## Acqusition Functions: Where should we sample next?



### Parallel Expected Improvement: Setting

We can parallelize EI:

$$\mathsf{EI}[\mathbf{x}_{1:q}] = \mathbb{E}_{n}[(\max(f[\mathbf{x}_{1}], ..., f[\mathbf{x}_{q}]) - f[\hat{\mathbf{x}}])^{+}]$$
(13)

How to maximize the parallel expected improvement?

- **Output** Construct an unbiased estimator of the gradient of  $EI[x_{1:q}]$ .
- **2** Use multistart stochastic gradient ascent to appoximately maximize  $EI[x_{1:q}]$ .

### PEI: Esitimate $\nabla EI$

Here's how we estimate  $\nabla EI$ :

- $Y = [f[x_1], ..., f[x_q]]$  is multivariate normal (GP prior)
- Let m = E[Y] and C = Chol(Cov[Y])
- then Y = m + CZ, where Z is a vector of independent standard normals
- $\mathsf{EI}[x_{1:q}] = \mathbb{E}[h(Y)]$ , where  $h(Y) = (\max[Y] f[\hat{x}])^+$
- Assume the problem is well-behaved, then we can switch derivative and expectation:

$$\nabla \mathsf{EI}[\mathbf{x}_{1:q}] = \mathbb{E}[\nabla h(\mathbf{m} + \mathbf{CZ})]$$

#### PEI: Esitimate $\nabla EI$

Here's how we estimate  $\nabla EI$ :

- Simulate a vector Z of independent standard normals
- Calculate  $m = \mathbb{E}[Y]$  and C = Chol(Cov[Y])
- Then the estimator of  $\nabla EI[x_1, ..., x_q]$  is  $\nabla h(m + CZ)$ , where  $h(Y) = (\max[Y] f[\hat{x}])^+$

## Approximately Maximize $EI[x_{1:q}]$

We use the previous estimator of  $\nabla EI$  in multistart stochastic gradient ascent:

- Select several starting points, uniformly at random
- From each starting point, iterate using the stochastic gradient method until convregence.

$$(\vec{x}_1,...,\vec{x}_q) \longleftarrow (\vec{x}_1,...,\vec{x}_q) + \alpha_n g(\vec{x}_1,...,\vec{x}_q,\omega),$$

where  $(\alpha_n)$  is a stepsize sequence,

- For each starting point, average the iterates to get an estimated stationary point. (Polyak-Ruppert averaging)
- Select the estimated stationary point with the best estimated value as the solution.







Here's a demonstration on a 6-dimensional Bayesian Optimization problem with up to 128 parallel evaluations.



**Bayesian Optimization** 

### **Knowledge Gradient**

- KG measures the expected increase in the maximum value of the objective function after a new observation.
- We will not select the previously evaluated point as the final solution: the exploration can be "unsuccessful"
- Formally, for a candidate point *x*, KG is defined as:

$$KG(\mathbf{x}) = \mathbb{E}[\mu_{n+1}^* - \mu_n^* | \text{observing} \mathbf{x}_{n+1}]$$

• Here,  $\mu_n^* = \max_x \mu_n[x] = \max_x \mathbb{E}_n[f(x)]$  is the best expected value of our ovjective under the posterior at tiem step *n*.

## Knowledge Gradient: Formula

- Let μ(x) and σ(x) be the posterior mean and standard deviation of the objective function at x.
- The KG acquisition function can be expressed as:

$$\mathbf{KG}(\mathbf{x}) = (\mu(\mathbf{x}) - \mu_n^*) \Phi\left(\frac{\Delta(\mathbf{x})}{\sigma(\mathbf{x})}\right) + \sigma(\mathbf{x}) \phi\left(\frac{\Delta(\mathbf{x})}{\sigma(\mathbf{x})}\right)$$

Where:

- $\mu_n^*$  is the current best observed value.
- $\Phi$  is the cumulative distribution function of the standard normal distribution.
- $\phi$  is the probability density function of the standard normal distribution.

#### Our approach for optimizing parallel EI also works for optimizing KG

**1.** Estimate  $\nabla KG(\mathbf{x}_{1:q})$  using infinitesimal perturbation analysis (IPA) & the envelope theorem:

$$\nabla \mathrm{KG} = \nabla \mu_{n+1}(\mathbf{x}^*; \mathbf{m}(\mathbf{x}_{1:q}) + \mathbf{CZ}, \mathbf{x}_{1:q}),$$

calculating the gradient holding  $x^*$  fixed.

**2.** Use multistart stochastic gradient ascent to maximize  $KG(\mathbf{x}_{1:q})$ 

## El May Make Poor Decisions













#### Incorporating noisy measurements

## **Incorporating Noisy Measurements**

We add an extra noise term to the expression for the Gaussian process covariance

$$\mathbb{E}[(\mathbf{y}(\mathbf{x}) - \mathbf{m}(\mathbf{x}))(\mathbf{y}(\mathbf{x}') - \mathbf{m}(\mathbf{x}'))] = \mathbf{k}(\mathbf{x}, \mathbf{x}') + \sigma_{\mathbf{n}}^{2},$$
(14)

where  $y(x) = f(x) + \epsilon$  is a noisy observation. And the joint normal becomes

$$\Pr\left(\begin{bmatrix}\boldsymbol{f}\\\boldsymbol{f}^*\end{bmatrix}\right) = \operatorname{Norm}\left(\boldsymbol{0}, \begin{bmatrix}\boldsymbol{K}[\boldsymbol{X},\boldsymbol{X}] + \sigma_n^2\boldsymbol{I} & \boldsymbol{K}[\boldsymbol{X},\boldsymbol{x}^*]\\\boldsymbol{K}[\boldsymbol{x}^*,\boldsymbol{X}] & \boldsymbol{K}[\boldsymbol{x}^*,\boldsymbol{x}^*]\end{bmatrix}\right).$$
(15)

Then We have the conditional distribution with noise:

$$\Pr(\mathbf{f}^*|\mathbf{f}) = \operatorname{Norm}(\mu[\mathbf{x}^*], \sigma^2[\mathbf{x}^*]),$$
(16)

where

$$\mu[\mathbf{x}^*] = \mathbf{K}[\mathbf{x}^*, \mathbf{X}] \left[ \mathbf{K}[\mathbf{X}, \mathbf{X}] + \sigma_n^2 I \right]^{-1} \mathbf{f}$$
(17)

$$\sigma^{2}[\boldsymbol{x}^{*}] = \boldsymbol{K}[\boldsymbol{x}^{*}, \boldsymbol{x}^{*}] - \boldsymbol{K}[\boldsymbol{x}^{*}, \boldsymbol{X}] \left[ \boldsymbol{K}[\boldsymbol{X}, \boldsymbol{X}] + \sigma_{\boldsymbol{n}}^{2} \boldsymbol{I} \right]^{-1} \boldsymbol{K}[\boldsymbol{x}^{*}, \boldsymbol{x}^{*}].$$
(18)

Incorporating noisy measurements

## Noisy Measurements Illustration

# a) No measurement noise b) Measurement noise $x^*$ $x^*$

#### Wang MA (SUSTech/RPI)

#### **Further Reading**

- Learning GP parameters: MLE, full Bayesian Approach
- Different Probabilistic Models: Random Forest (SMAC)
- Discrete Variables: Beta-Bernoulli Bandit
- Kernel Choices
- Tips, tricks, and limitations
  - Inducing points
  - Decomposing the kernel
  - Using random projections

#### Contents

1. Introduction

- 2. An Overview of Bayesian Optimization
- 3. A model of the function: Gaussian Process
- 4. Acquisition Functions
- 5. Conclusion

### Conclusion: Bayesian Optimization

#### Algorithm 2 BayesOpt

Assume a Bayesian prior on f

(usually a Gaussian process prior, a probabilistic model of the function)

#### while budget is not exhausted do

Find *x* that maximizes **acquisition function** (*x*, *posterior*)

Sample *x* and observe f(x)

Update the posterior distribution on *f* 

#### Conclusion: Gaussian Process Regression

$$\Pr(\mathbf{f}^*|\mathbf{f}) = \operatorname{Norm}(\mu[\mathbf{x}^*], \sigma^2[\mathbf{x}^*]), \tag{19}$$

where

$$\mu[\mathbf{x}^*] = \mathbf{K}[\mathbf{x}^*, \mathbf{X}]\mathbf{K}[\mathbf{X}, \mathbf{X}]^{-1}\mathbf{f}$$
(20)

$$\sigma^{2}[\boldsymbol{x}^{*}] = \boldsymbol{K}[\boldsymbol{x}^{*}, \boldsymbol{x}^{*}] - \boldsymbol{K}[\boldsymbol{x}^{*}, \boldsymbol{X}]\boldsymbol{K}[\boldsymbol{X}, \boldsymbol{X}]^{-1}\boldsymbol{K}[\boldsymbol{x}^{*}, \boldsymbol{x}^{*}].$$
(21)

#### **Conclusion: Acquisition Functions**

- Acquisition functions guide the selection of the next point to evaluate by balancing exploration and exploitation.
- Common acquisition functions include:
  - Probability of Improvement (PI)

$$\mathsf{PI}[\mathbf{x}^*] = \int_{\mathbf{f}[\hat{\mathbf{x}}]}^{\infty} \mathsf{Norm}_{\mathbf{f}[\mathbf{x}^*]}[\mu[\mathbf{x}^*], \sigma[\mathbf{x}^*]] d\mathbf{f}[\mathbf{x}^*],$$
(22)

Expected Improvement (EI)

$$EI[x^*] = E_n[(f[x^*] - f[\hat{x}])^+]$$

$$= \int_{f[\hat{x}]}^{\infty} (f[x^*] - f[\hat{x}]) \operatorname{Norm}_{f[x^*]}[\mu[x^*], \sigma[x^*]] df[x^*].$$
(23)
(24)

Upper Confidence Bound (UCB)

$$\mathsf{UCB}[\mathbf{x}^*] = \mu[\mathbf{x}^*] + \kappa \sigma[\mathbf{x}^*], \tag{25}$$

Wang MA (SUSTech/RPI)

**Bayesian Optimization** 

July 22, 2024 46/48

#### Review

- 1. Introduction
- 2. An Overview of Bayesian Optimization
- 3. A model of the function: Gaussian Process
- 4. Acquisition Functions
- 5. Conclusion

#### **Bayesian Optimization: END**

## Thank you!

#### Questions and Opinions are Welcome!