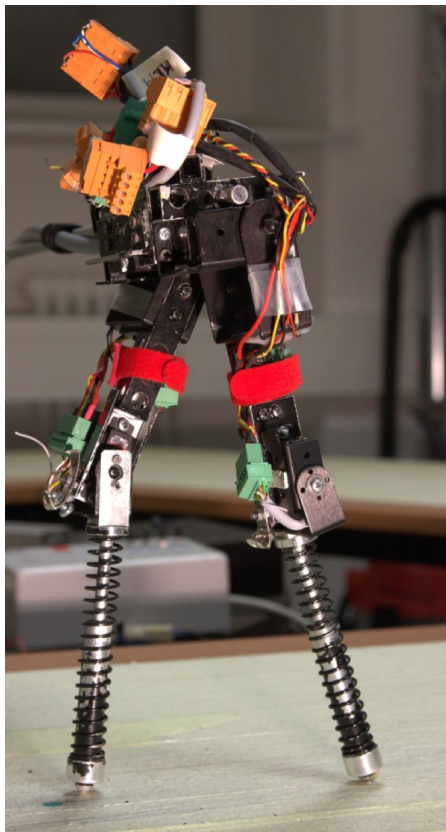


# A tutorial on Bayesian optimization

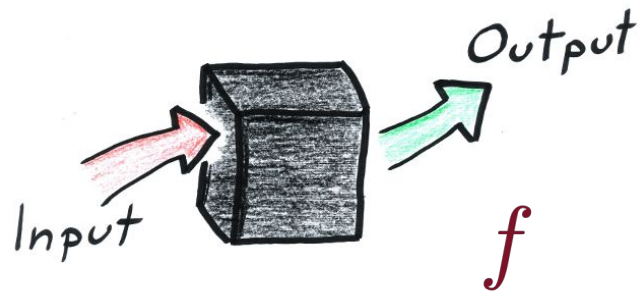
Zi Wang @ Google Brain



# Blackbox Function Optimization



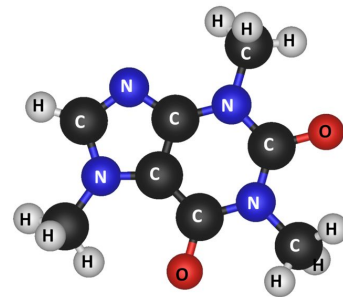
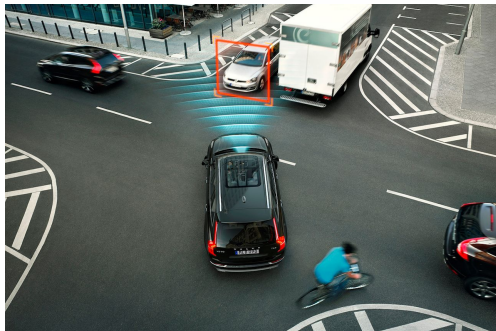
[Calandra et al., 2015]



Goal:

$$x_* = \operatorname{argmax}_{x \in \mathbb{R}^d} f(x)$$

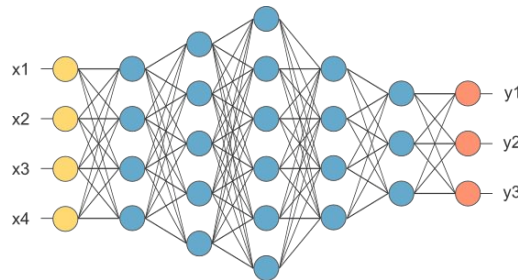
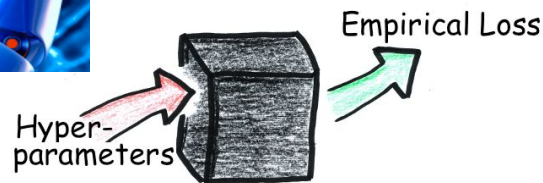
# Blackbox Function Optimization



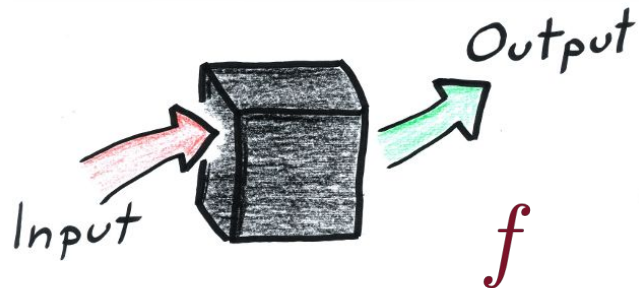
**Goal:**  $x_* = \operatorname{argmax}_{x \in \mathbb{R}^d} f(x)$

**Challenges:**

- $f$  is expensive to evaluate
- $f$  is multi-peak
- no gradient information
- evaluations can be noisy



# Blackbox Function Optimization

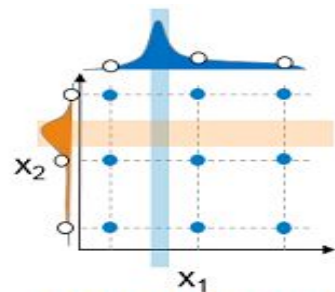


**Goal:**  $x_* = \operatorname{argmax}_{x \in \mathbb{R}^d} f(x)$

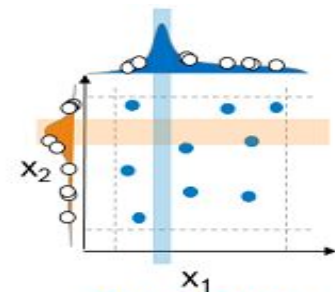
**Challenges:**

- $f$  is expensive to evaluate
- $f$  is multi-peak
- no gradient information
- evaluations can be noisy

Grid search?



Random search?



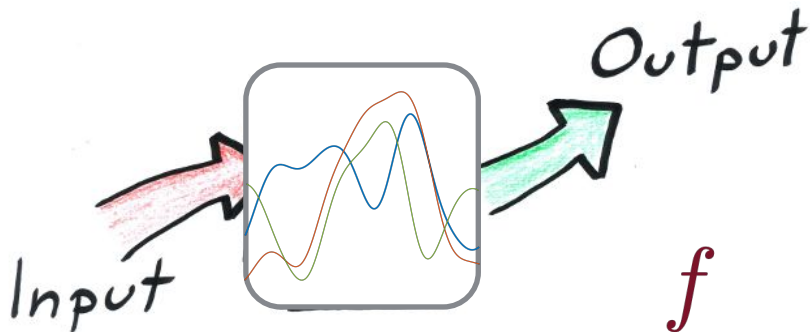
Many evaluations are wasted!

# Bayesian Optimization

Idea: build a **probabilistic model** of the function  $f$

## LOOP

- choose new query point(s) to evaluate  
*decision criterion: **acquisition function**  $\alpha_t(\cdot)$*
- update model



$$x_* = \operatorname{argmax}_{x \in \mathbb{R}^d} f(x)$$

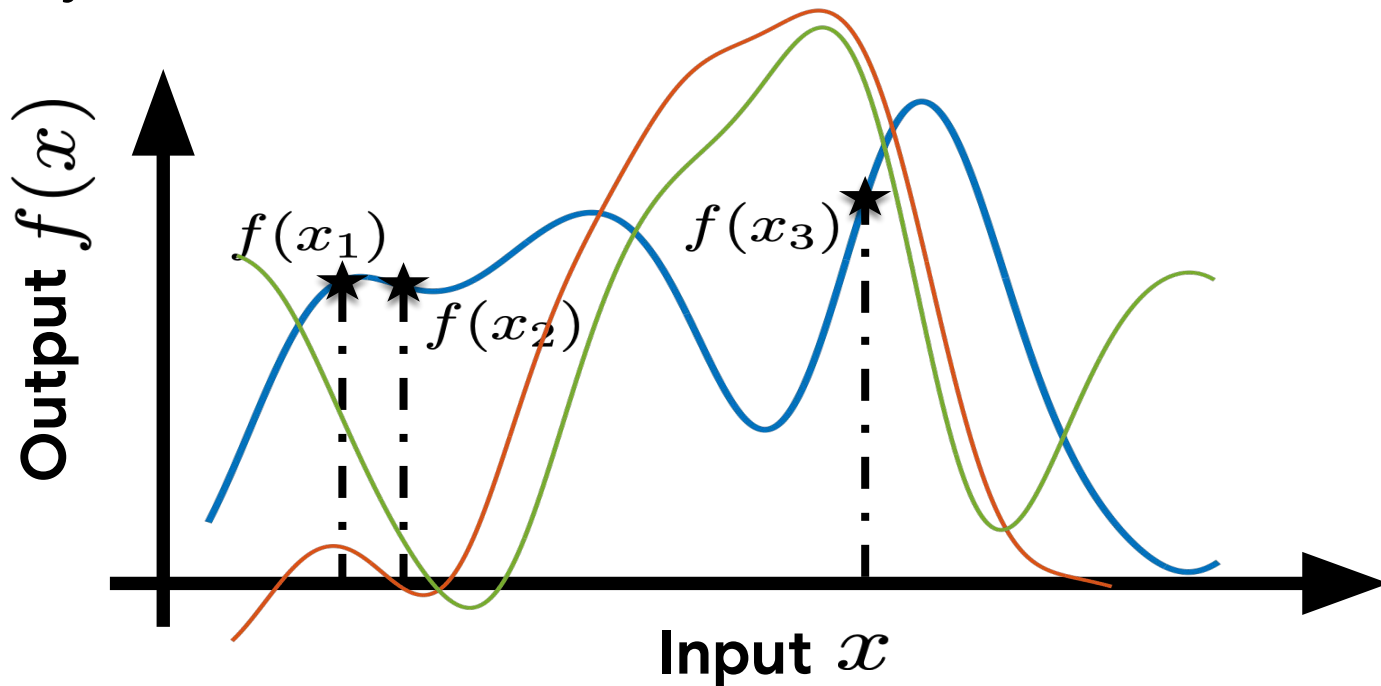


$$x_t = \operatorname{argmax}_{x \in \mathbb{R}^d} \alpha_t(x)$$

$$t = 1, \dots, T$$

# Gaussian Processes (GPs)

- probability distribution over functions
- any finite set of function values is a multi-variate Gaussian



# Gaussian Processes (GPs)

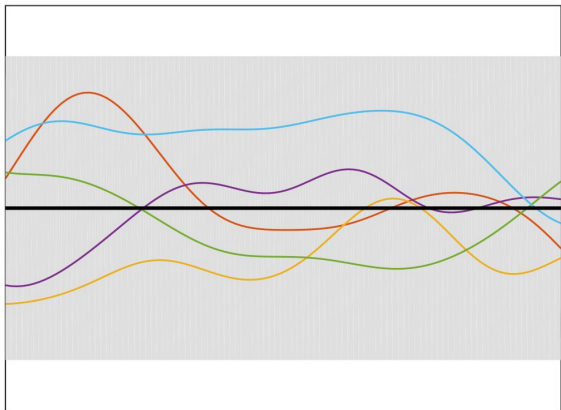
- probability distribution over functions
- any finite set of function values is a multi-variate Gaussian
- kernel function  $k(\cdot, \cdot)$ ; mean function  $\mu(\cdot)$

$$\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu(x_1) \\ \vdots \\ \mu(x_n) \end{bmatrix}, \begin{bmatrix} k(x_1, x_1), & \cdots, & k(x_1, x_n) \\ \vdots, & & \vdots \\ k(x_n, x_1), & \cdots, & k(x_n, x_n) \end{bmatrix} \right)$$

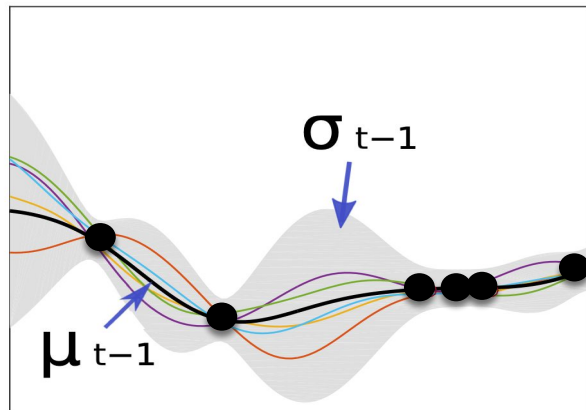
- function  $f \sim GP(\mu, k)$ ; observe noisy output at  $x_\tau$   
$$y_\tau = f(x_\tau) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

# Gaussian Processes (GPs)

Samples from the prior



Samples from the posterior



Given observations  $D_t = \{(x_\tau, y_\tau)\}_{\tau=1}^{t-1}$ , predict posterior mean and variance in **closed form** via conditional Gaussian

$$\mu_{t-1}(x) = k_{t-1}(x)^\top (K_{t-1} + \sigma^2 I)^{-1} y_{t-1}$$

$$\sigma_{t-1}(x)^2 = k(x, x) - k_{t-1}(x)^\top (K_{t-1} + \sigma^2 I)^{-1} k_{t-1}(x)$$

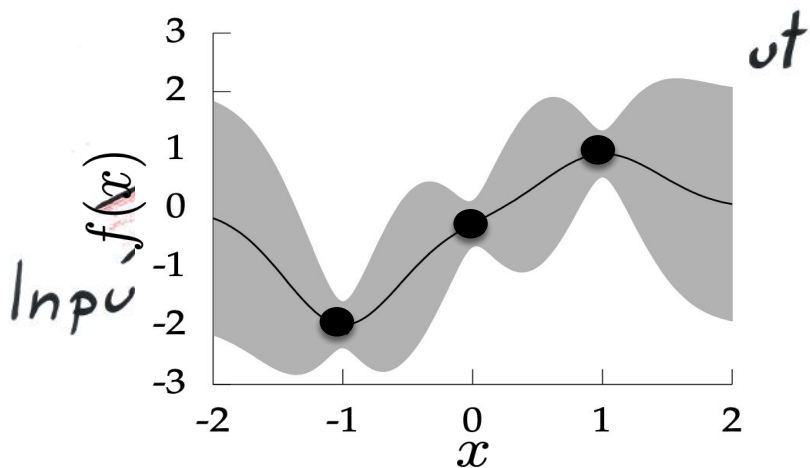


# Bayesian Optimization

Idea: build a **probabilistic model** of the function  $f$

## LOOP

- choose new query point(s) to evaluate  
*decision criterion: acquisition function*  $\alpha_t(\cdot)$
- update model



$$x_* = \operatorname{argmax}_{x \in \mathbb{R}^d} f(x)$$



$$x_t = \operatorname{argmax}_{x \in \mathbb{R}^d} \alpha_t(x)$$

$$t = 1, \dots, T$$

How to design acquisition functions?

# Acquisition functions in BayesOpt

# Examples of acquisition functions in BayesOpt

Prior:  $f \sim GP(\mu, k)$

At iteration  $t$ ,

- predict the posterior  $\mu_{t-1}(x)$  and  $\sigma_{t-1}^2(x)$
- pick an input by optimizing the acquisition function

$$x_t = \arg \max_x \alpha_t(x)$$

Upper confidence bounds, expected improvement, probability of improvement, entropy search methods...

[Auer, 2002; Srinivas et al., 2010; Kushner, 1964; Mockus, 1974; Hennig & Schuler, 2012; Hernandez-Lobato et al., 2014; Wang&Jegelka, 2017; Hoffman&Zoubin, 2015...]

# GP-UCB: an example of acquisition functions

[Auer, 2002; Srinivas et al., 2010]

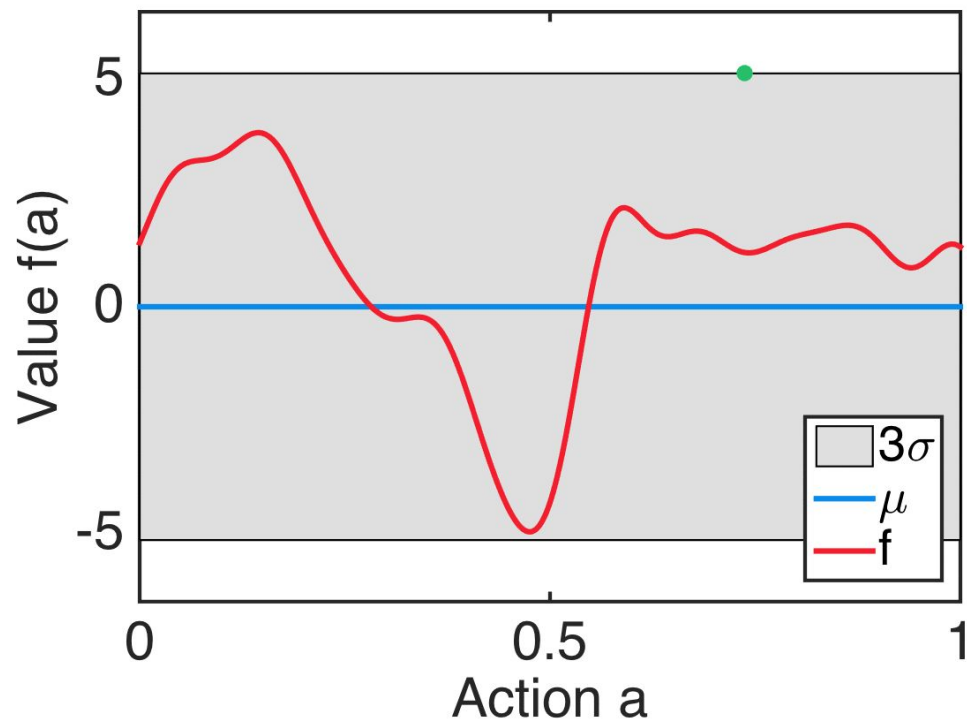
Prior:  $f \sim GP(\mu, k)$

At iteration  $t$ ,

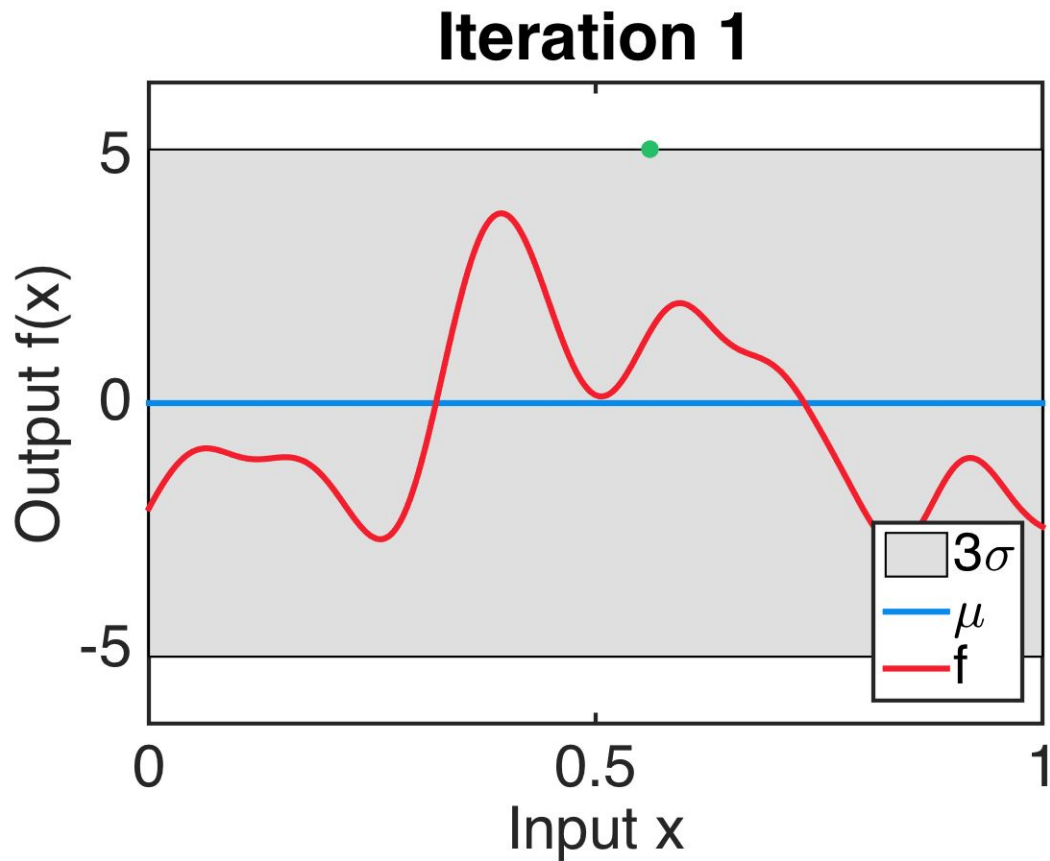
- predict the posterior  $\mu_{t-1}(x)$  and  $\sigma_{t-1}^2(x)$
- pick an input by optimizing the acquisition function

$$x_t = \arg \max \mu_{t-1}(x) + \beta \sigma_{t-1}(x)$$

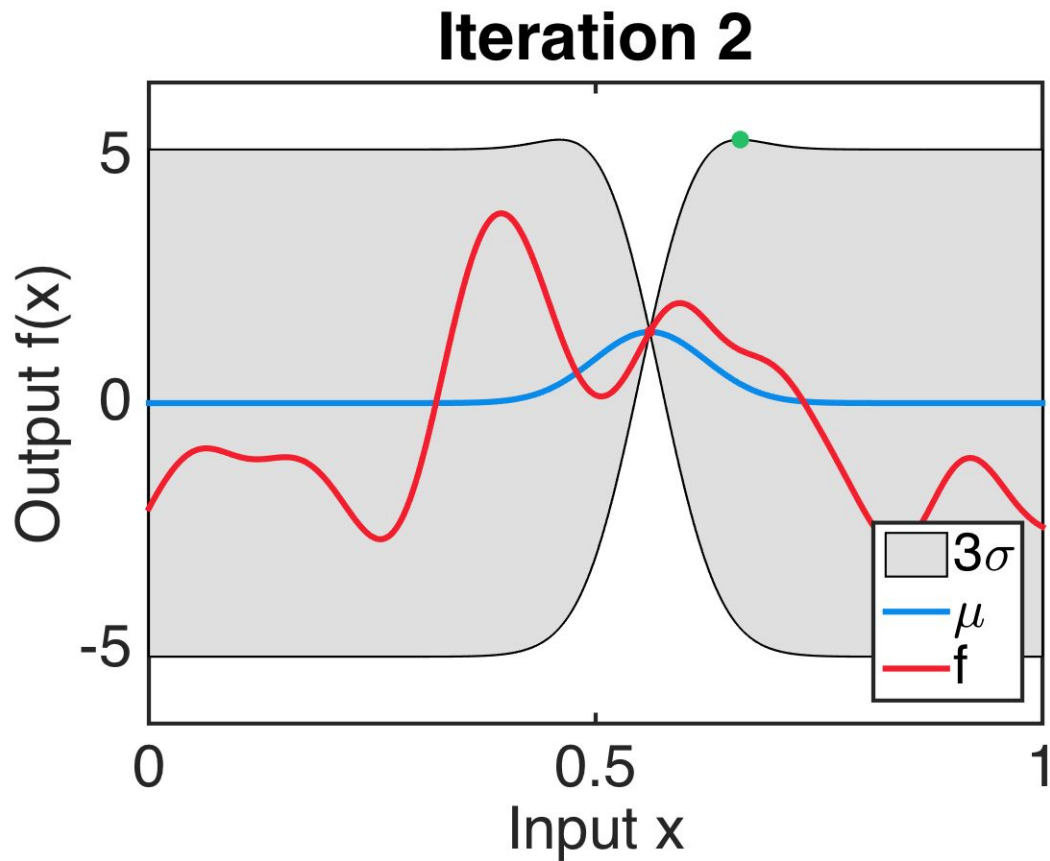
## Iteration 1



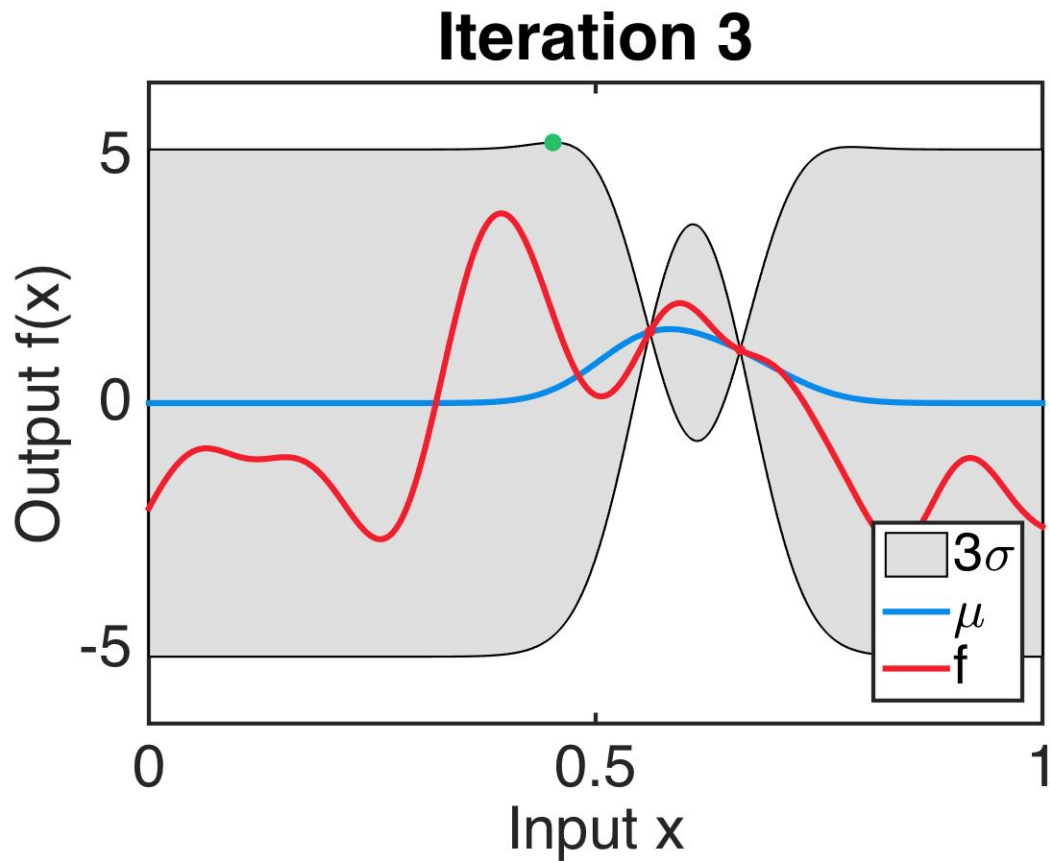
# GP-UCB: an example of acquisition functions



# GP-UCB: an example of acquisition functions

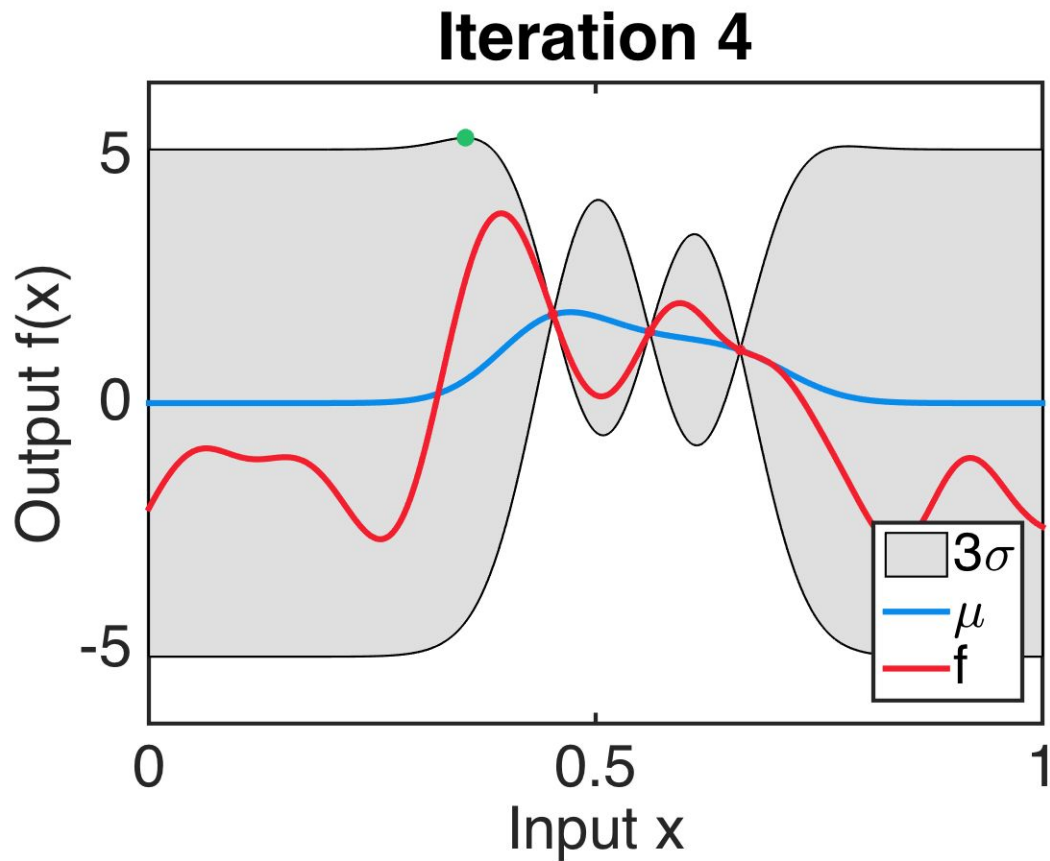


# GP-UCB: an example of acquisition functions

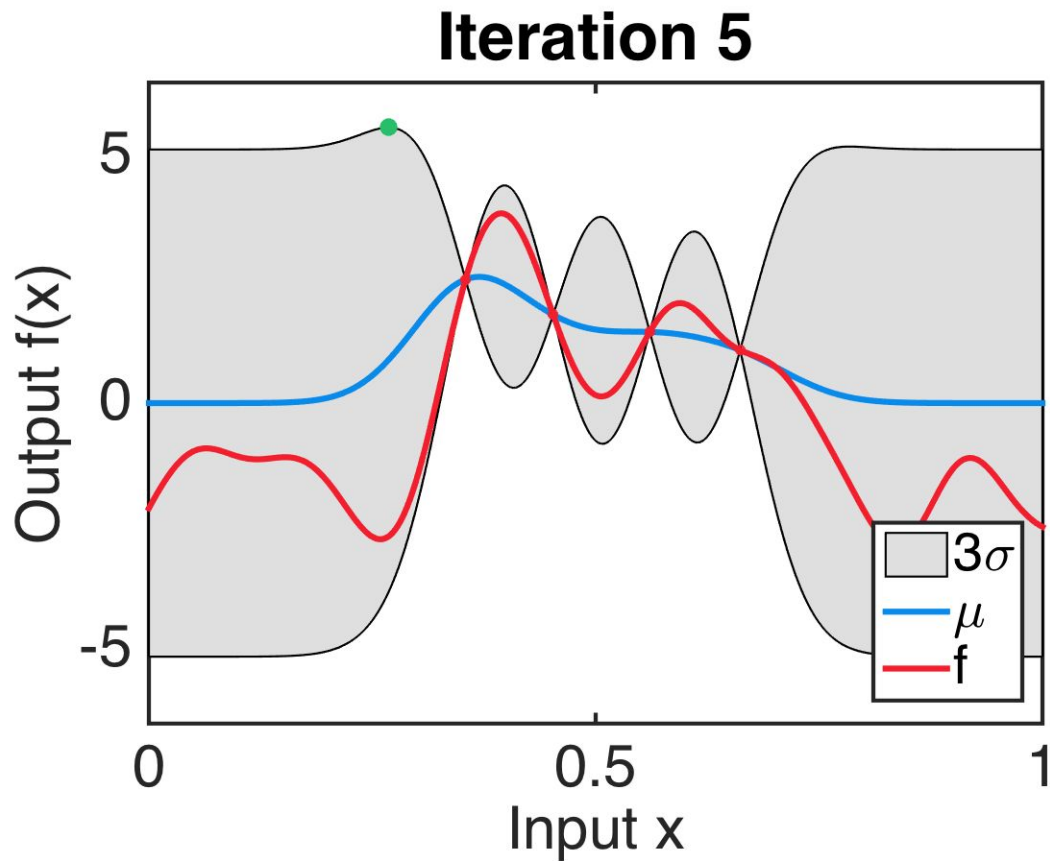




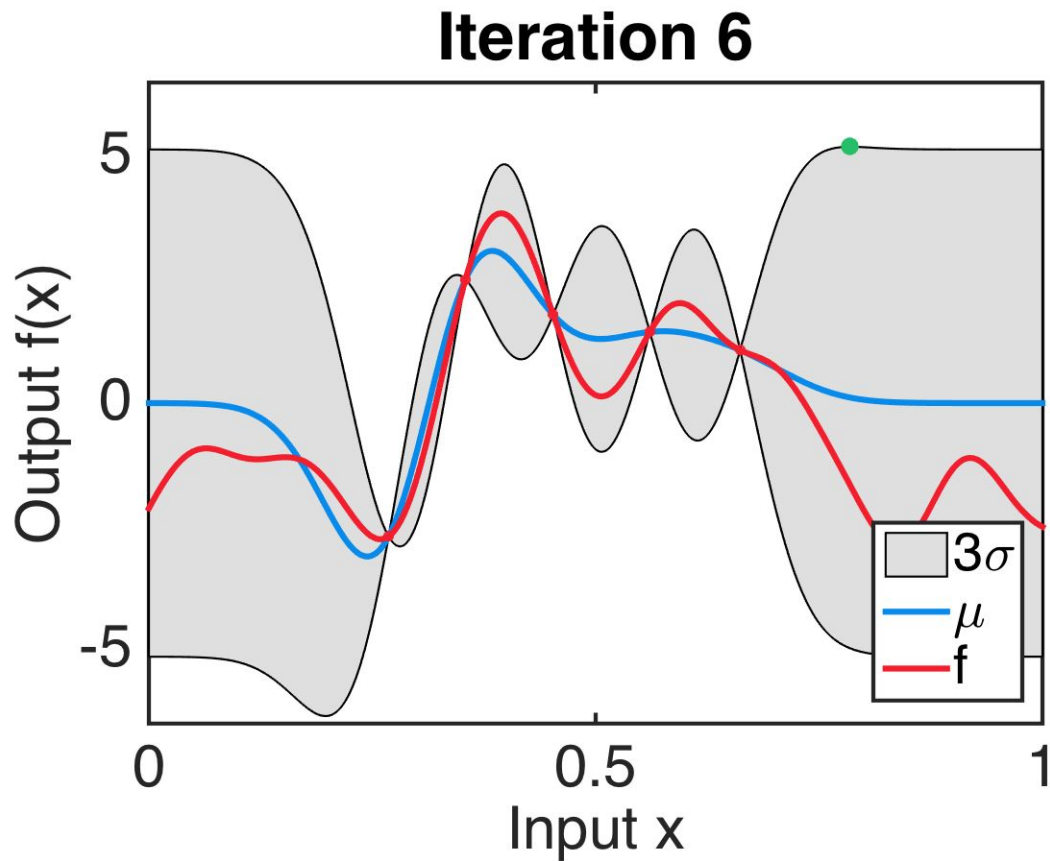
# GP-UCB: an example of acquisition functions



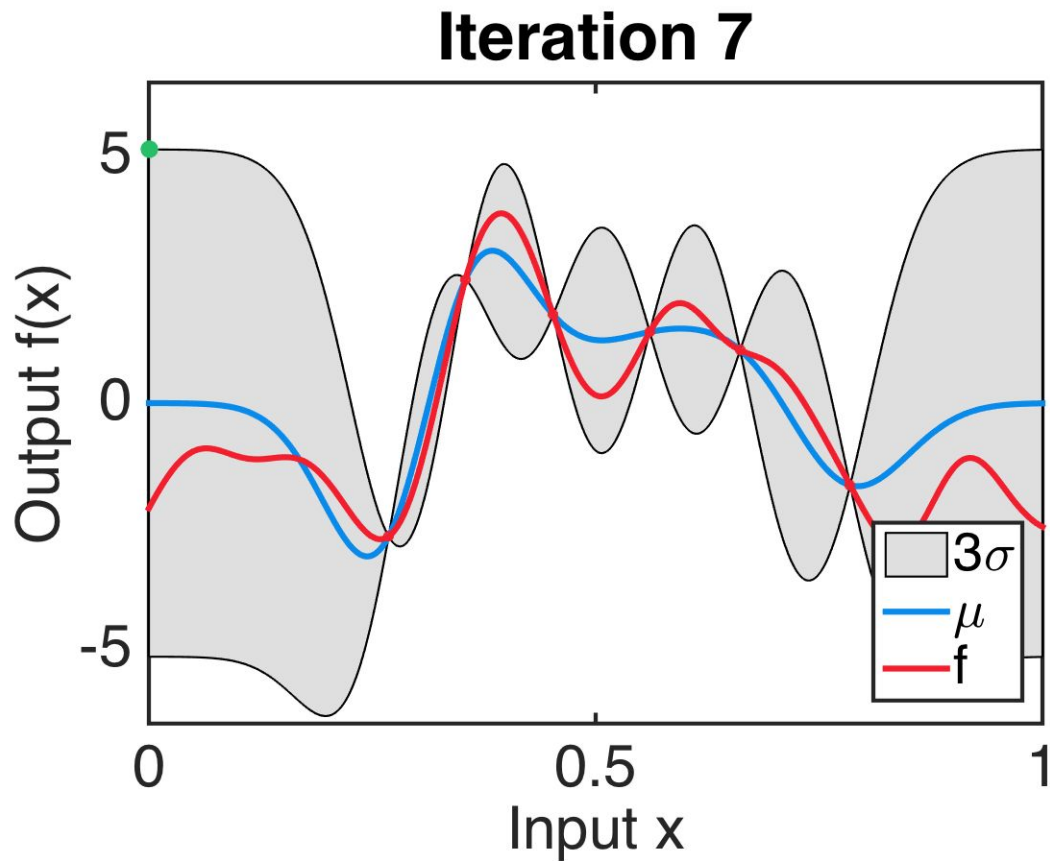
# GP-UCB: an example of acquisition functions



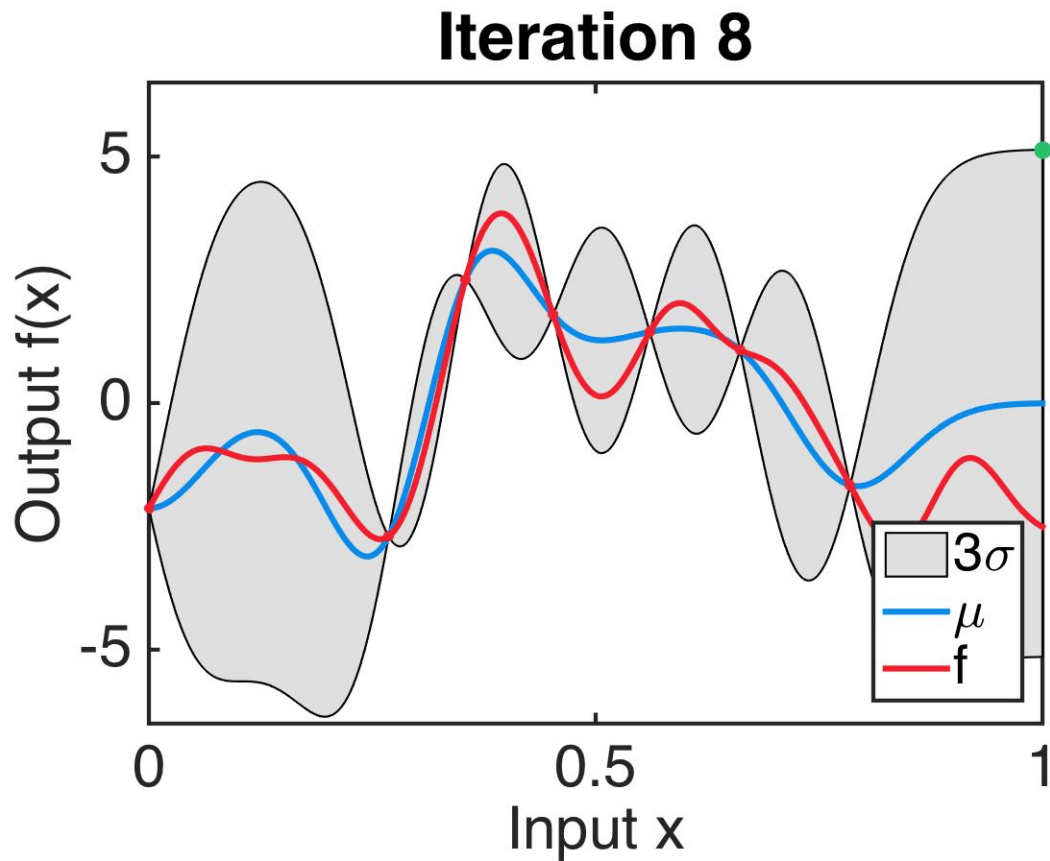
# GP-UCB: an example of acquisition functions



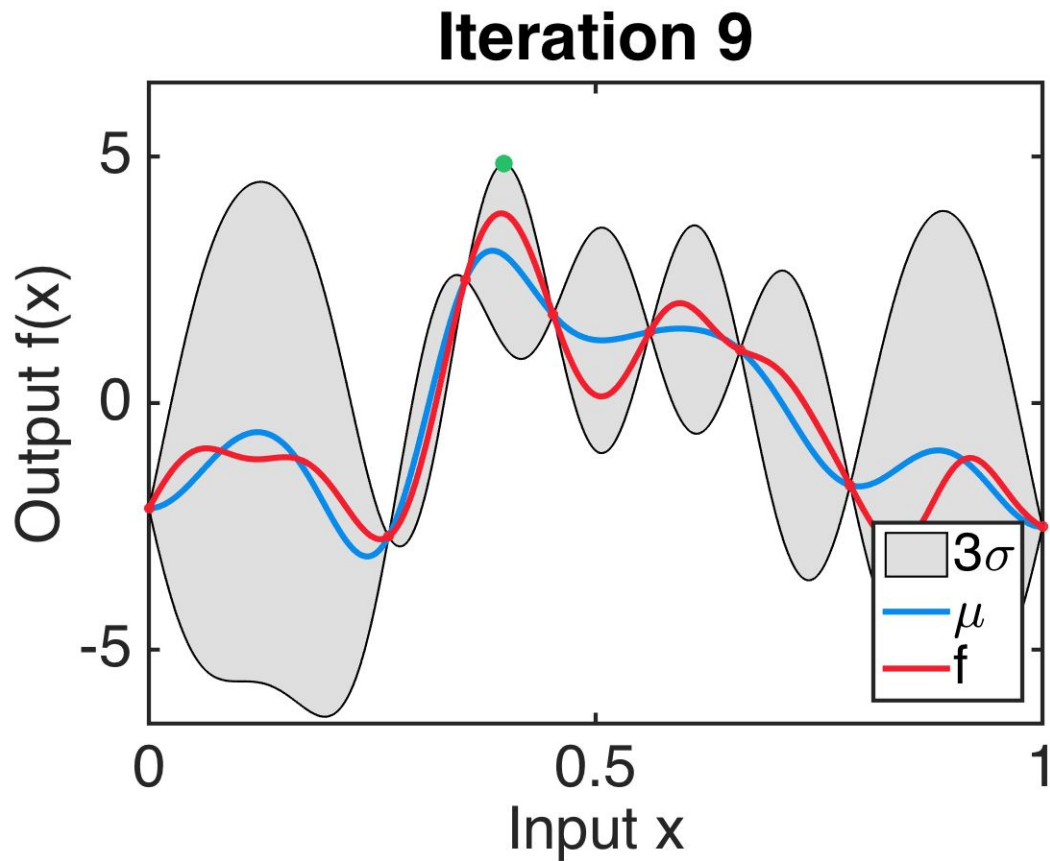
# GP-UCB: an example of acquisition functions



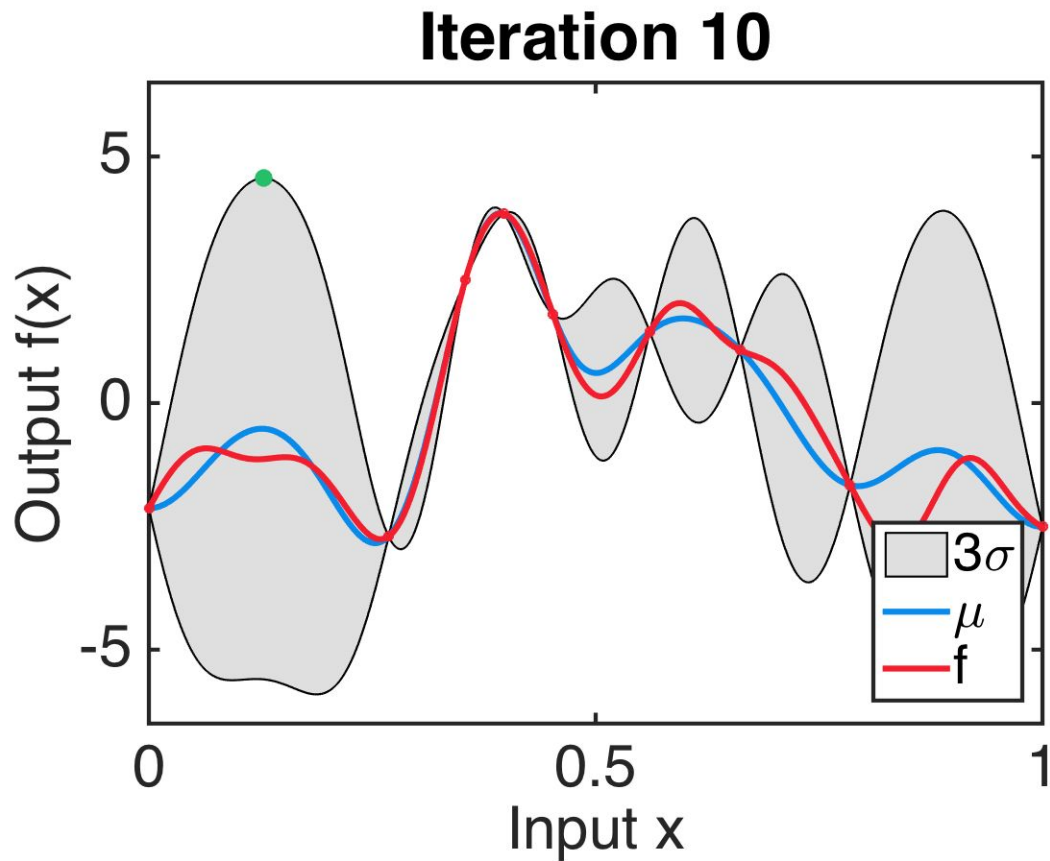
# GP-UCB: an example of acquisition functions



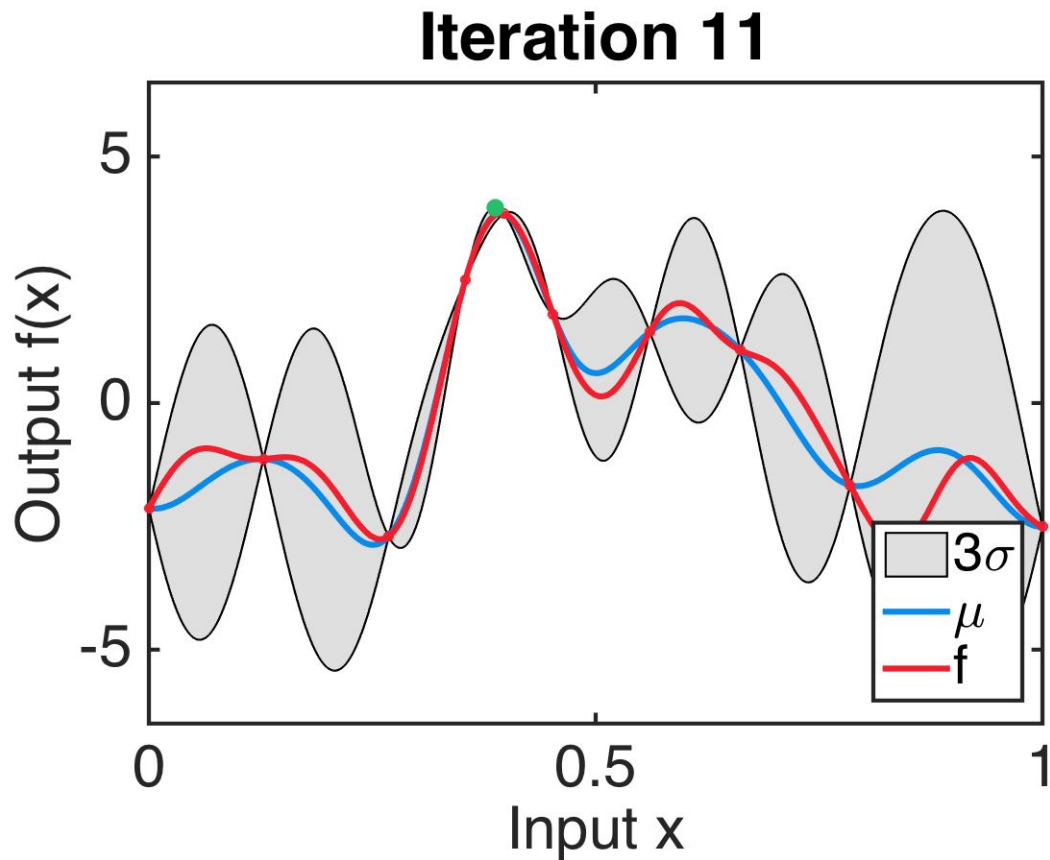
# GP-UCB: an example of acquisition functions



# GP-UCB: an example of acquisition functions

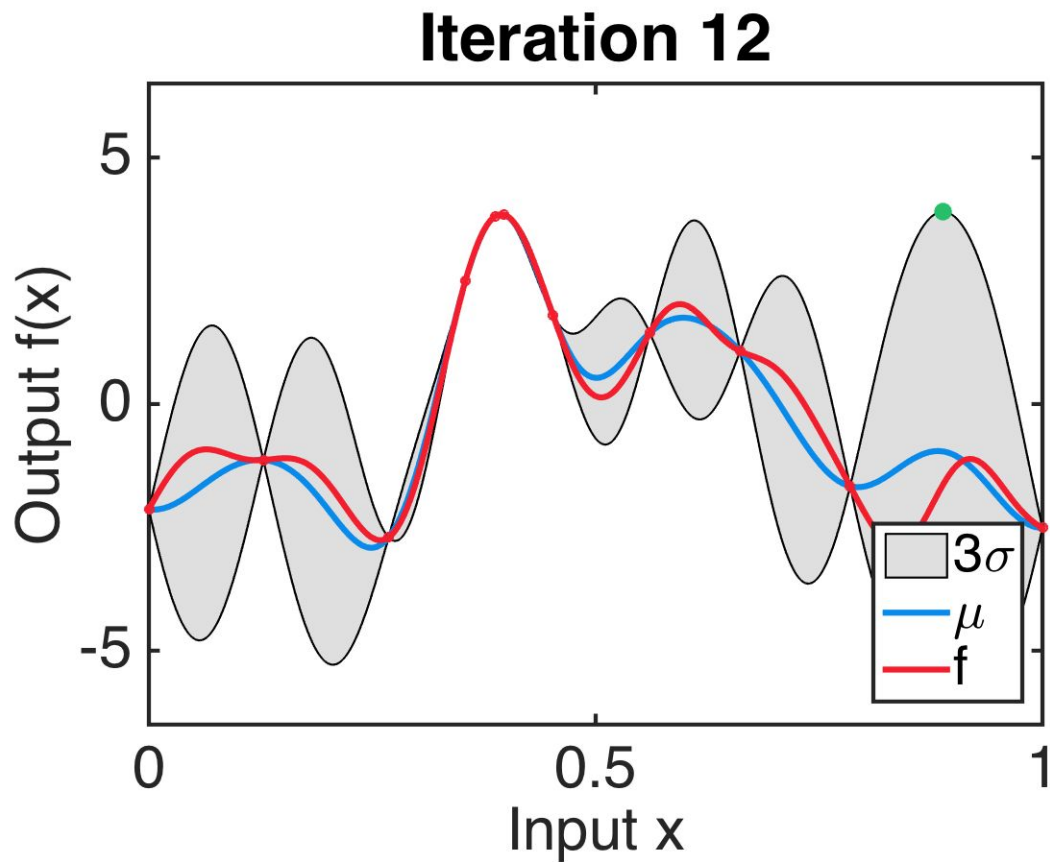


# GP-UCB: an example of acquisition functions

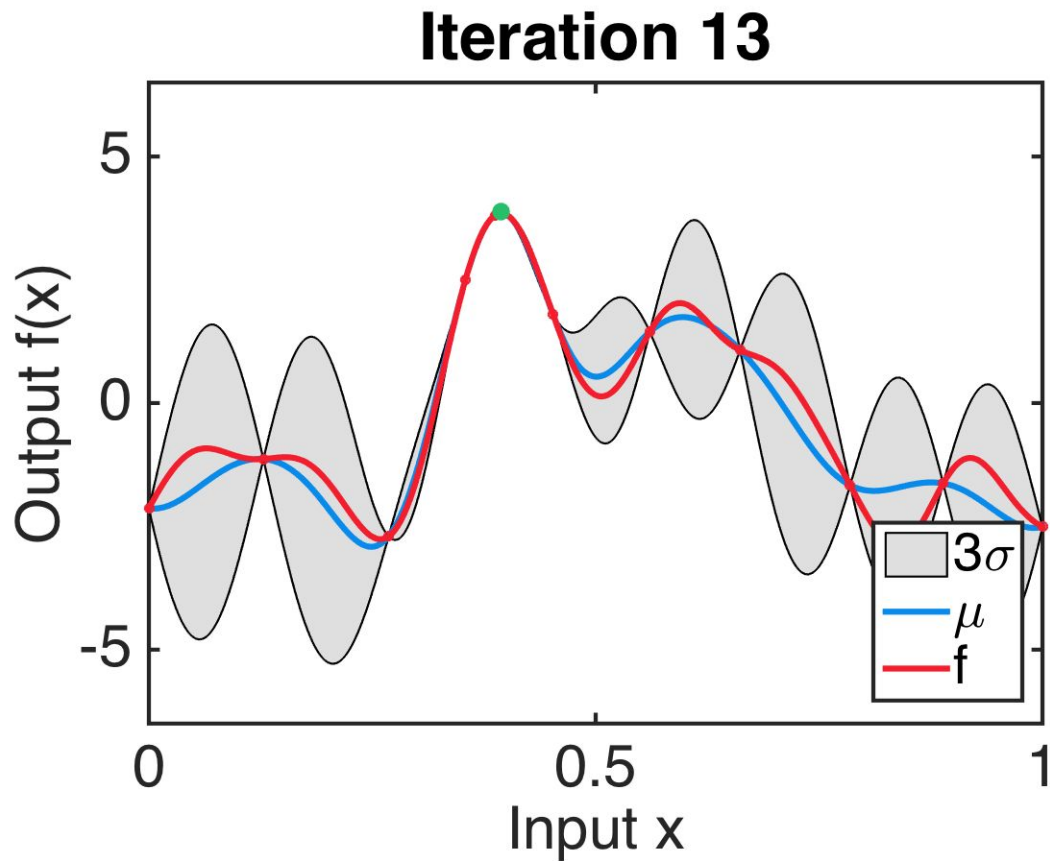




# GP-UCB: an example of acquisition functions



# GP-UCB: an example of acquisition functions

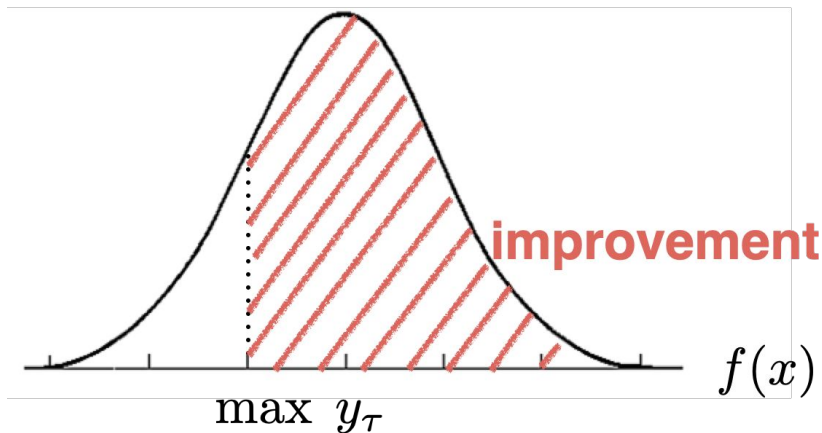


# Example of acquisition functions: PI

[Kushner, 1964]

## Probability of Improvement (PI)

- Observations:  $D_t = \{(x_\tau, y_\tau)\}_{\tau=1}^{t-1}$
- The best observation is  $\max y_\tau$
- for each  $x$ , predict the posterior mean and variance



$$\alpha_t(x) = \Pr[f(x) \geq \max y_\tau]$$

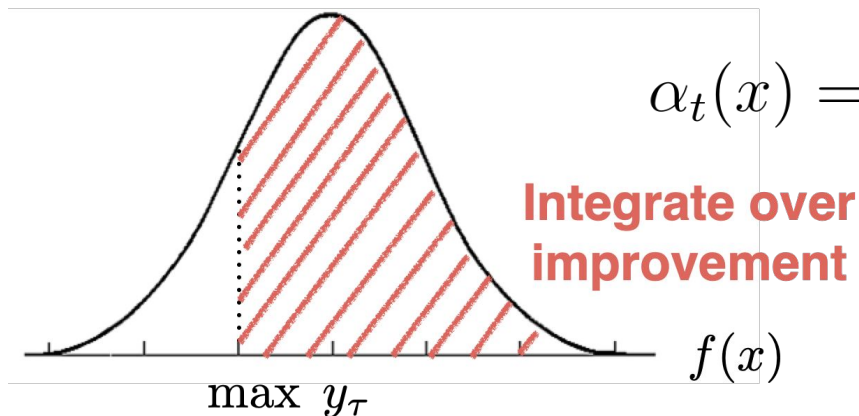
$$\alpha_t(x) = \Pr[f(x) \geq \max y_\tau + \epsilon]$$

# Example of acquisition functions: EI

[Kushner, 1964]

## Expected Improvement (EI)

- Observations:  $D_t = \{(x_\tau, y_\tau)\}_{\tau=1}^{t-1}$
- The best observation is  $\max y_\tau$
- for each  $x$ , predict the posterior mean and variance



$$\alpha_t(x) = \mathbb{E} [(f(x) - \max y_\tau)_+]$$

# Entropy Search and Predictive Entropy Search

$$\underset{x_t \in \mathcal{X}}{\text{maximize}} \alpha_t(x_t)$$

$$t = 1, \dots, T$$

Point to  
query

Location  
of global  
optimum

Observed  
Data

$$\begin{aligned} I(a; b) &= H(a) - H(a|b) \\ &= H(b) - H(b|a) \end{aligned}$$

$$\alpha_t(x) = I(\{x, y\}; x_* \mid D_t)$$

$$= H(p(x_* \mid D_t)) - \mathbb{E}_y[H(p(x_* \mid D_t \cup \{x, y\}))]$$

$$= H(p(y \mid D_t, x)) - \mathbb{E}_{x_*}[H(p(y \mid x_*, D_t, x))]$$

# Max-value Entropy Search

[Wang&Jegelka, 2017; Hoffman&Zoubin, 2015]

Point to query

Observed  
Data

$$\alpha_t(x) = I(\{x, y\}; x_* | D_t)$$

Location of global optimum

D-dimensional  
input space

$$\alpha_t(x) = I(\{x, y\}; y_* | D_t)$$

**Global max-value**

$$= H(p(y | D_t, x)) - \mathbb{E}_{y_*} [H(p(y | y_*, D_t, x))]$$

$$y \leq y_*$$

1-dimensional  
Output space

Gaussian

Truncated  
Gaussian

$$\approx \frac{1}{K} \sum_{y_* \in Y_*}$$

$$\left[ \frac{\gamma_{y_*}(x) \psi(\gamma_{y_*}(x))}{2\Psi(\gamma_{y_*}(x))} - \log(\Psi(\gamma_{y_*}(x))) \right]$$

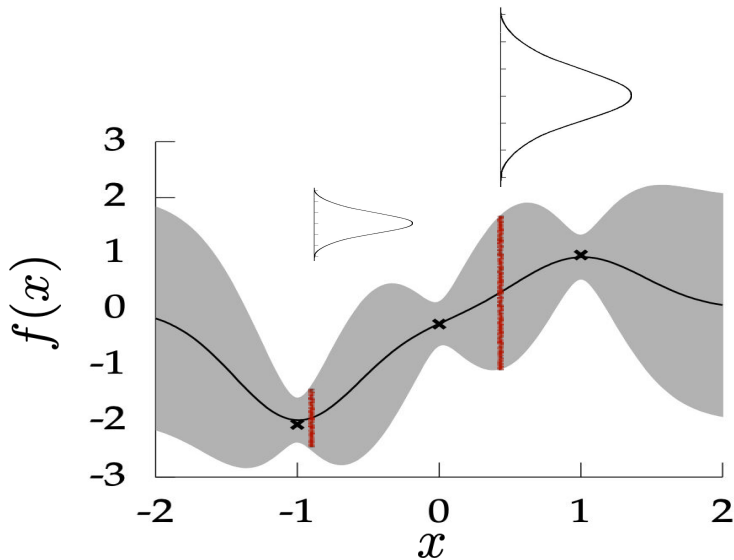
Something  
Closed-form

# Sample $y_*$ with a Gumbel Distribution

Intuition: each  $f(x)$  is a Gaussian

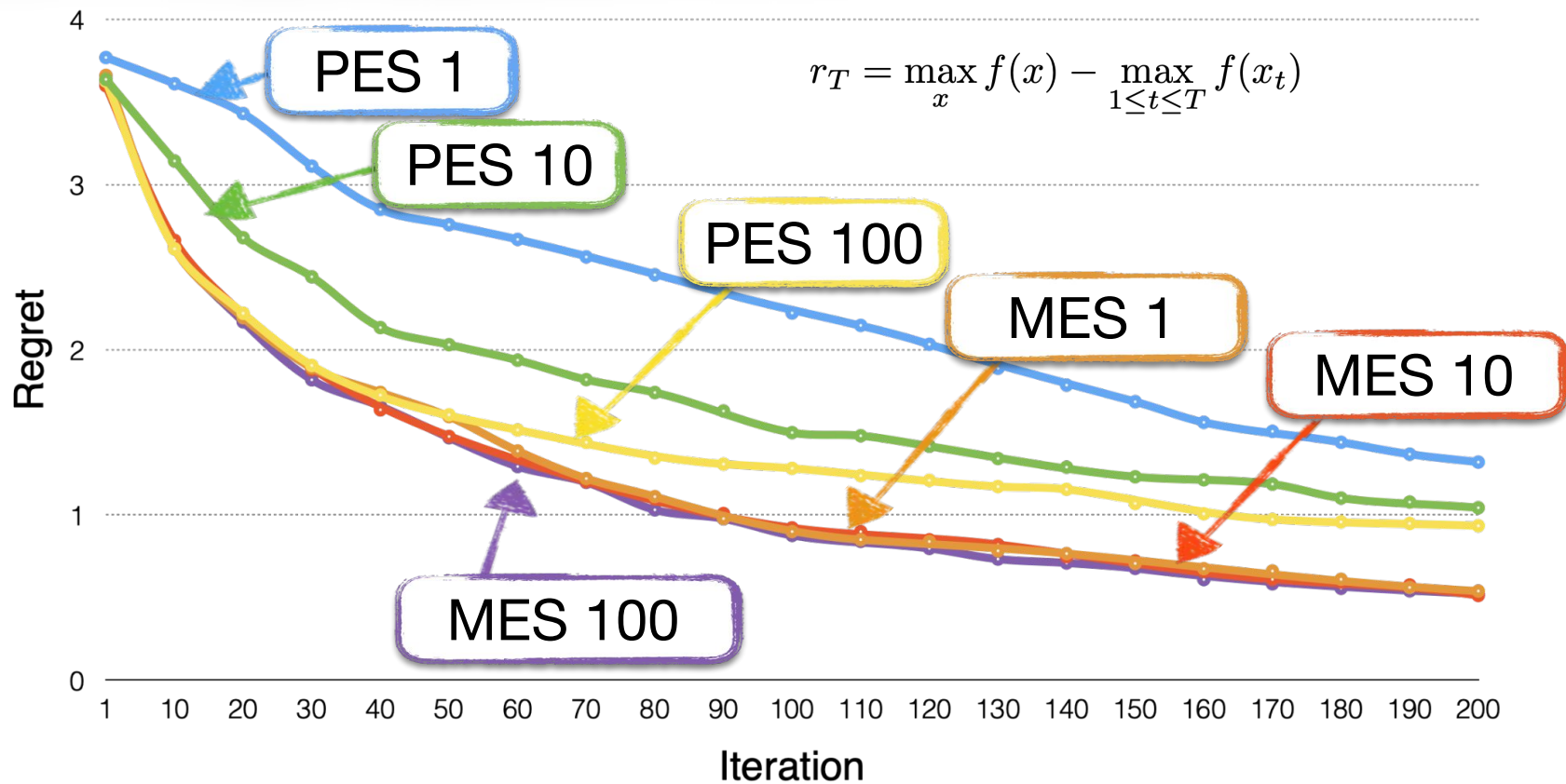
## Fisher-Tippett-Gnedenko Theorem

The maximum of a set of i.i.d. Gaussian variables is asymptotically described by a **Gumbel distribution**.



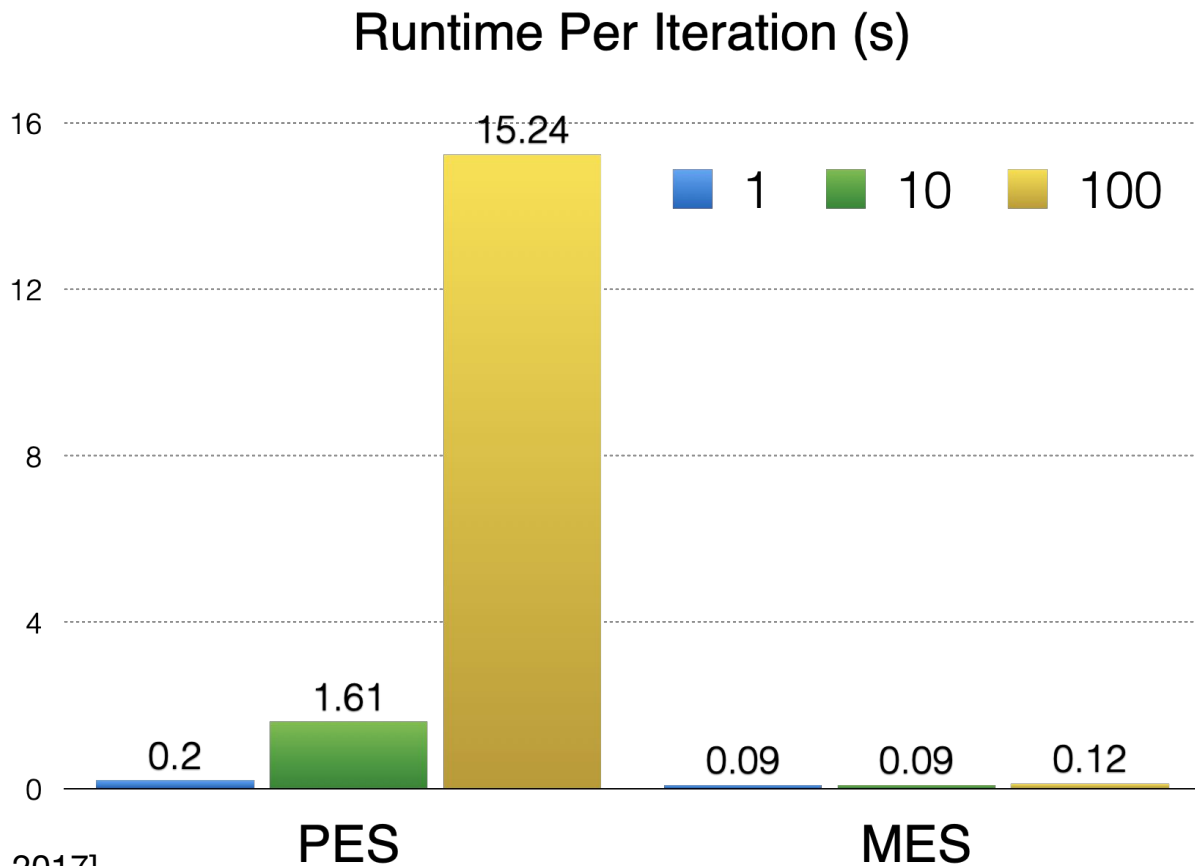
- Sample representative points
- Approximate the max-value of the representative points by a Gumbel distribution [Wang&Jegelka, 2017]

# MES gets faster and better empirical results than PES





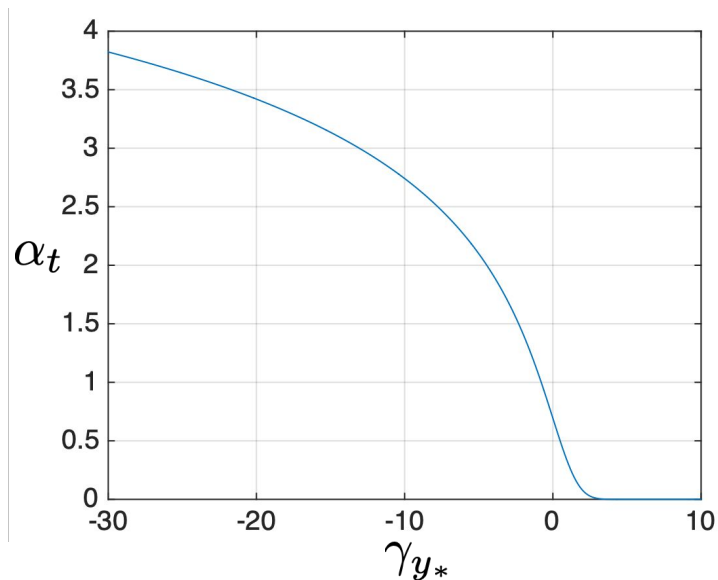
# MES gets faster and better empirical results than PES



# Understanding the acquisition function in MES

$$\alpha_t(x) \approx$$

$$\frac{\gamma_{y_*}(x)\psi(\gamma_{y_*}(x))}{2\Psi(\gamma_{y_*}(x))} \alpha_t(\gamma_{y_*}(x))$$



$$\gamma_{y_*}(x) = \frac{y_* - \mu_{t-1}(x)}{\sigma_{t-1}(x)}$$

So, maximize  $\alpha_t(x)$

is equivalent to

minimize  $\gamma_{y_*}(x)$ .

# Relations among GP-UCB, PI and MES

[Jones, 2001; Wang&Jegelka, 2017]

MES

$$\underset{x}{\text{minimize}} \gamma_{y_*}(x)$$

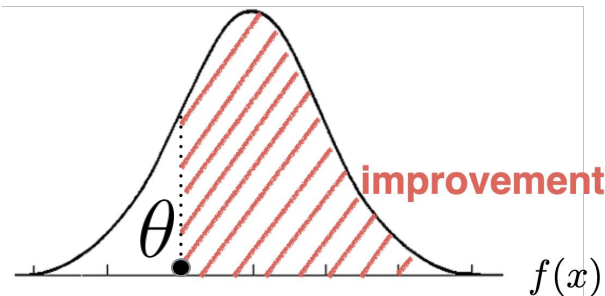
PI

$$\underset{x}{\text{minimize}} \gamma_{\max y_\tau + \epsilon}(x)$$

GP-UCB

$$\underset{x}{\text{minimize}} \gamma_{\max_{x'} \mu_{t-1}(x') + \beta \sigma_{t-1}(x')}(x)$$

$$\gamma_\theta(x) = \frac{\theta - \mu_{t-1}(x)}{\sigma_{t-1}(x)}$$



GP-UCB, PI and MES are equivalent to  $\underset{x}{\text{minimize}} \gamma_\theta(x)$  under special cases of  $\theta$ .

# Relations among GP-UCB, PI and MES

[Jones, 2001; Wang&Jegelka, 2017]

## MES

$$\underset{x}{\text{maximize}} \mu_{t-1}(x) + \sigma_{t-1}(x) \min_{x'} \gamma_{y_*}(x')$$

## PI

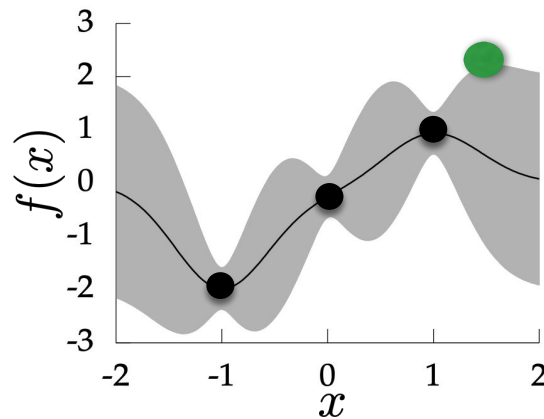
$$\underset{x}{\text{maximize}} \mu_{t-1}(x) + \sigma_{t-1}(x) \min_{x'} \gamma_{\max y_\tau + \epsilon}(x')$$

## GP-UCB

$$\underset{x}{\text{maximize}} \mu_{t-1}(x) + \beta \sigma_{t-1}(x)$$

They are also equivalent to  $\underset{x}{\text{maximize}} \mu_{t-1}(x) + \beta \sigma_{t-1}(x)$   
under special cases of  $\beta$ .

$$\gamma_\theta(x) = \frac{\theta - \mu_{t-1}(x)}{\sigma_{t-1}(x)}$$



# Regret bounds for GP-UCB and related methods

[Srinivas et al., 2010; Wang et al., 2016]

Define regret as:  $r_T = \max_x f(x) - \max_{1 \leq t \leq T} f(x_t)$

Key assumptions:  $f \sim GP(\mu, k)$

Mean function and kernel are both given.  
Optimize in a d-dimensional compact space.

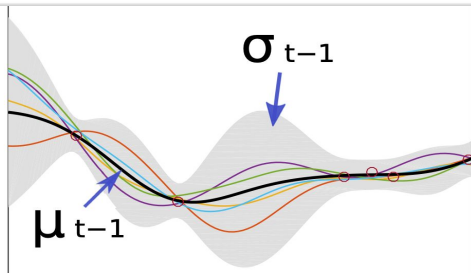
After T iterations, GP-UCB obtains  $r_T = O\left(\sqrt{\frac{d(\log T)^{d+2}}{T}}\right)$ .

MES obtains  $r_T = O\left(\sqrt{\frac{(\log T)^{d+2}}{T}} + \max_t \min_x \frac{y_* - \mu_{t-1}(x)}{\sigma_{t-1}(x)} \sqrt{\frac{(\log T)^{d+1}}{T}}\right)$ .

\* For simplicity we only show regret bounds for Gaussian kernels. Regret for other kernels may look different.

# Summary of how BayesOpt works

Posterior estimation



$t \leftarrow t + 1$

Evaluate  $f$  at  
 $x_t = \arg \max \alpha_t(x)$

Define an acquisition function

UCB:  $\mu_{t-1}(x) + \beta \sigma_{t-1}(x)$

El:  $\mathbb{E} [(f(x) - \max y_\tau)_+]$

Pl:  $\Pr[f(x) \geq \max y_\tau + \epsilon]$

ES:  $I(\{x, y\}; x_* \mid D_t)$   
 $I(\{x, y\}; y_* \mid D_t)$

and others...

Challenges, open problems and some attempts

# Selected topics in BayesOpt

- High dimensional search space
- Unknown GP prior
- Parallel evaluations
- Unknown constraints
- Applications in robotics

Challenges, open problems and some attempts

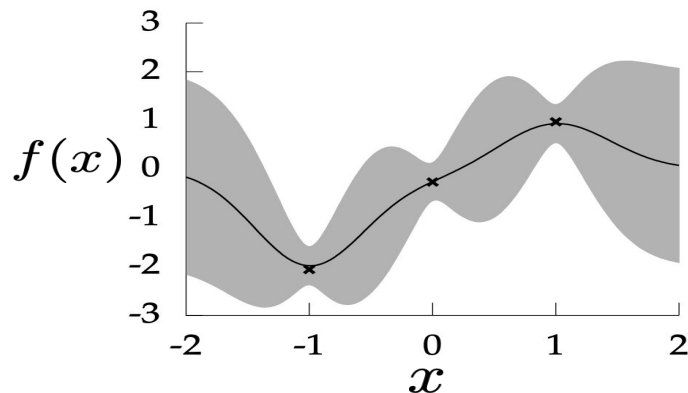
# High dimensional search space



# Challenges in high-dimensional BO

- optimizing multi-peak acquisition functions in high dimensions  
**computationally challenging**

- estimating a nonlinear function in high input dimensions: need more observations  
**statistically challenging**

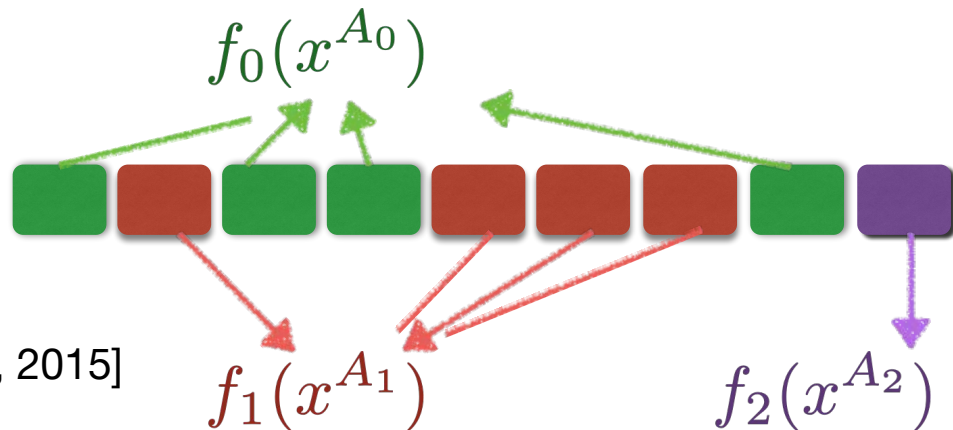


$$\text{regret } r_T \approx O\left(\sqrt{\frac{(\log T)^{d+2}}{T}}\right)$$

# Possible solution: additive Gaussian processes

$$f(x) = \sum_{m \in [M]} f_m(x^{A_m})$$

[Hastie&Tibshirani, 1990; Kandasamy et al., 2015]

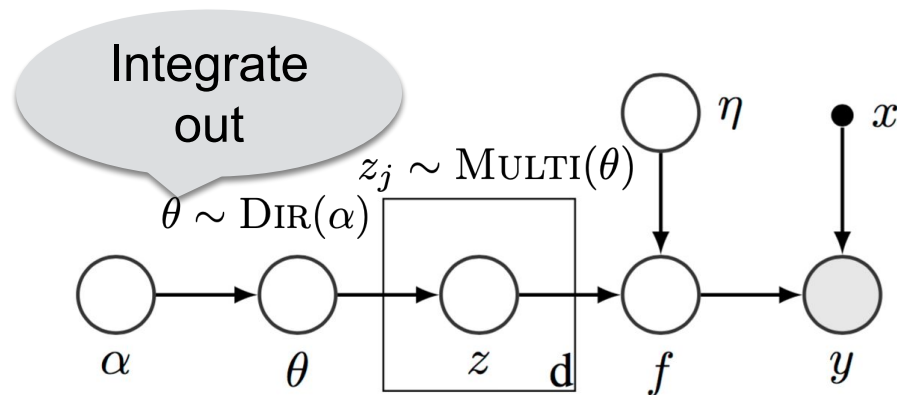
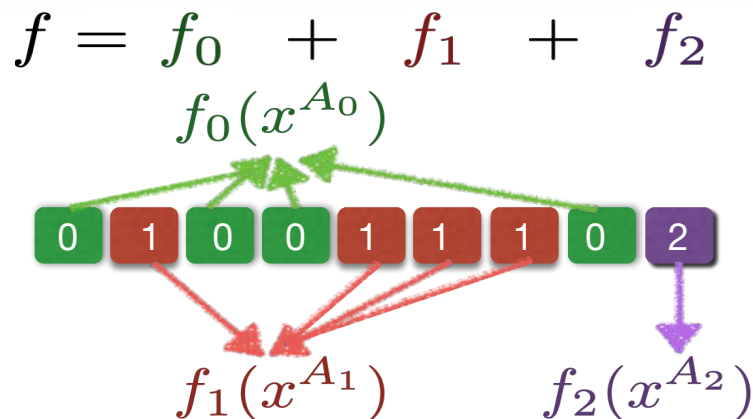


- optimize acquisition function block-wise  
**computational efficiency**
- lower-complexity functions  
**statistical efficiency**

**What is the additive structure?**

# Structural Kernel Learning (SKL)

[Wang et al., 2017]



*Decomposition indicator:*

$$z = [0, 1, 0, 0, 1, 1, 1, 0, 2]$$

Learn  $z$ !

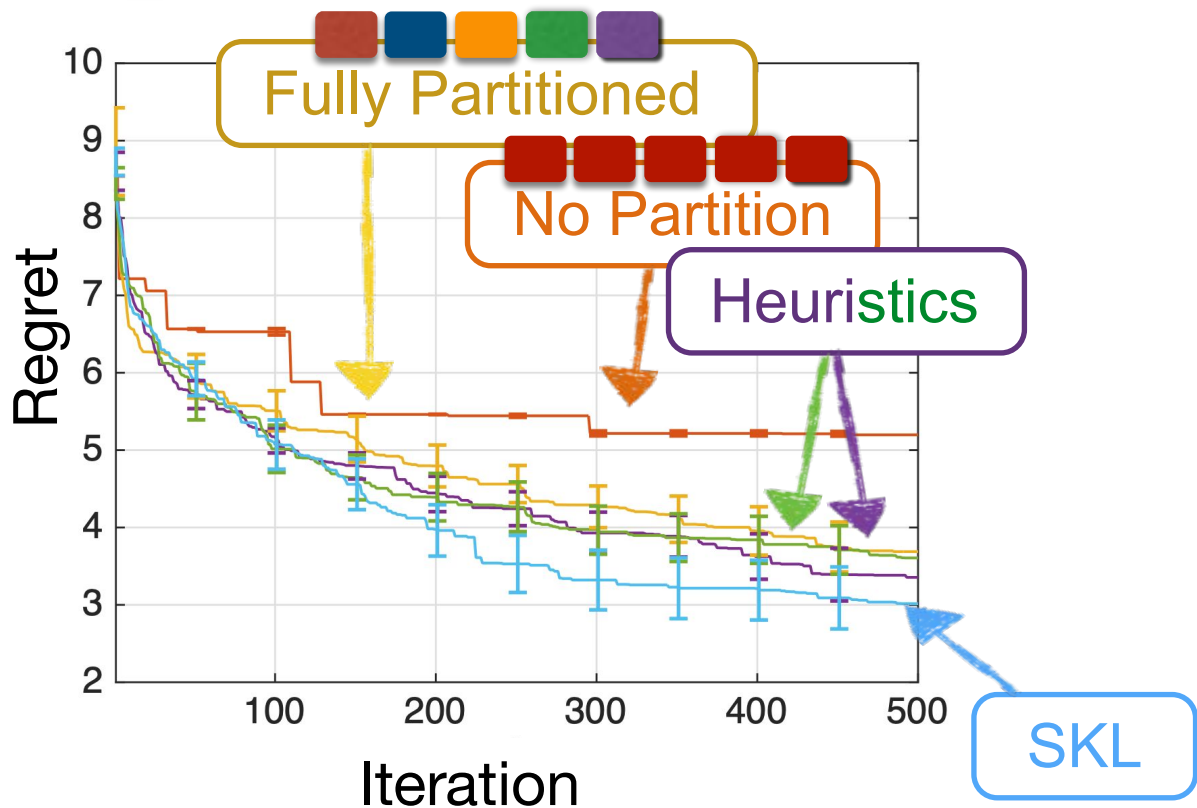
Learn posterior

$$p(z \mid D_n; \alpha)$$

via Gibbs sampling.

easy updates

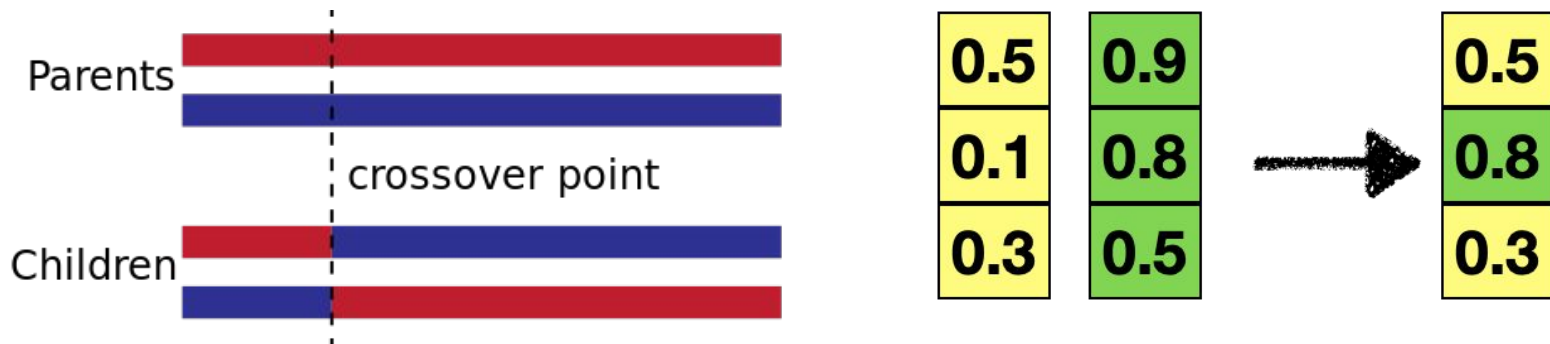
# Empirical results for Structural Kernel Learning (SKL)



# Connection to genetic algorithms?

## Evolutionary/Genetic algorithms:

- maintain ensemble of promising points
- new points from exchanging coordinates of good points randomly



# Connection to genetic algorithms?

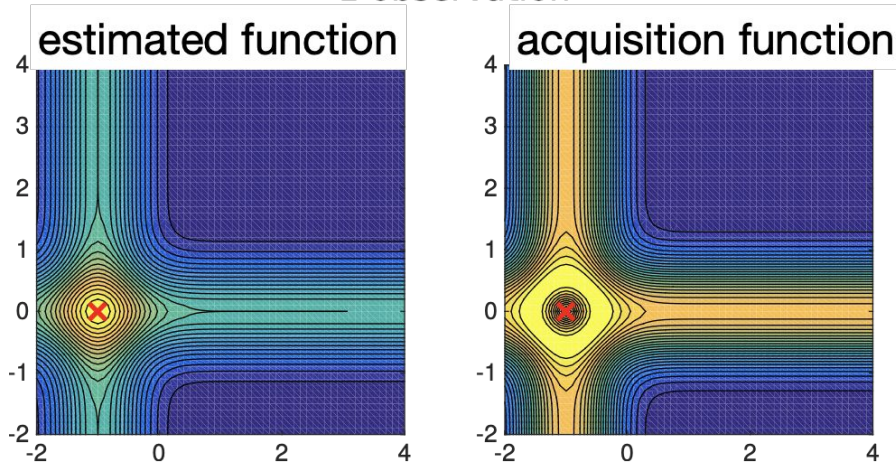
[Wang et al., 2018]

BayesOpt with additive GPs

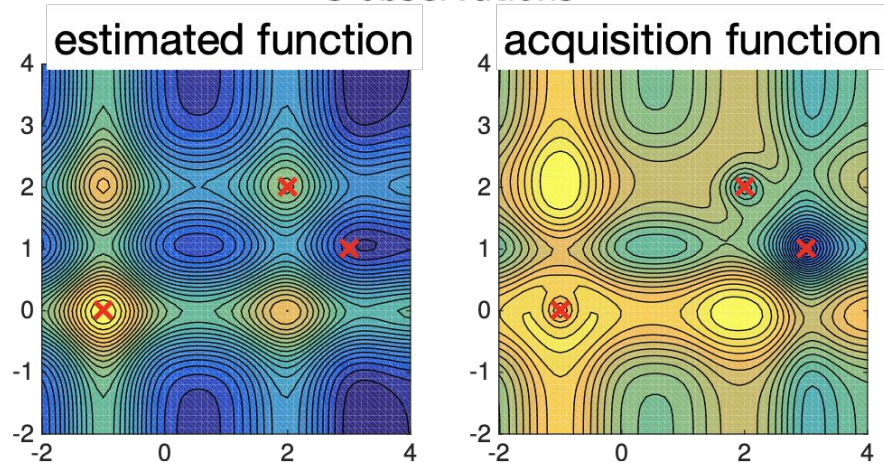
toy example: 2D



1 observation



3 observations



Observed good points:  $[-1, 0]$ ,  $[2, 2]$

Query points:  $[-1, 2]$ ,  $[2, 0]$



Learned instead of completely random coordinate partition.

# Other ideas to solve high-dim BayesOpt

- REMBO: low-dim embedding [Wang et al., JAIR 2016]
- BOCK: BO with cylindrical kernels [Oh et al., ICML 2018]
- Additive GPs with overlapping groups [Rolland et al., AISTATS 2018]
- .....

## Joint problems:

**Assume special structures of high-dim functions but with little data, it is difficult to verify if the assumptions are true.**

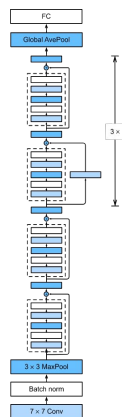
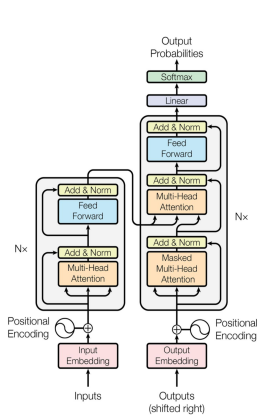
Challenges, open problems and some attempts

# Parallel evaluations



# BayesOpt with parallel compute resources

- GPUs running in parallel for hyperparameter tuning in deep learning;
- group of robots for offline learning of control parameter;
- parallel wet lab experiments for biology and chemistry applications; etc.



# Some ideas to propose a batch of queries

- Instead of optimizing one input over information gain, optimize  $Q$  inputs. [Shah&Ghahramani, 2015]

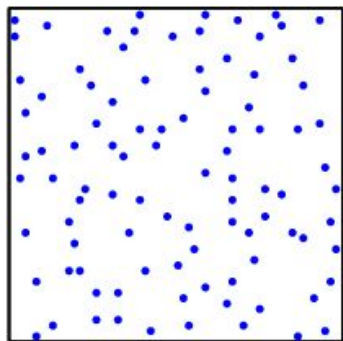
$$a_{\text{PPES}}(\mathcal{S}_t|\mathcal{D}) = H[p(\mathbf{x}^*|\mathcal{D})] - \mathbb{E}_{p(\{y_q\}_{q=1}^Q|\mathcal{D},\mathcal{S}_t)} \left[ H[p(\mathbf{x}^*|\mathcal{D} \cup \{\mathbf{x}_q, y_q\}_{q=1}^Q)] \right]$$

- Choose a new point based on expected acquisition function under all possible outcomes of pending evaluations. [Snoek et al., 2012]

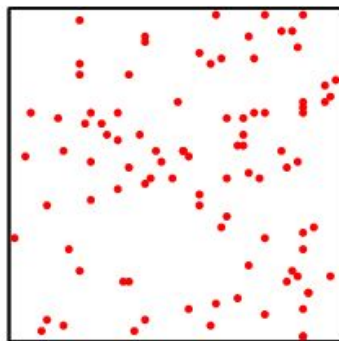
$$\hat{a}(\mathbf{x}; \{\mathbf{x}_n, y_n\}, \theta, \{\mathbf{x}_j\}) = \int_{\mathbb{R}^J} a(\mathbf{x}; \{\mathbf{x}_n, y_n\}, \theta, \{\mathbf{x}_j, y_j\}) p(\{y_j\}_{j=1}^J | \{\mathbf{x}_j\}_{j=1}^J, \{\mathbf{x}_n, y_n\}_{n=1}^N) dy_1 \cdots dy_J$$

# Some ideas to propose a batch of queries

- Use determinantal point process (DPP) to generate a diverse set of queries.  
[Kathuria et al., 2016]
- Use a Mondrian process to propose one query per partition. [Wang et al., 2018]

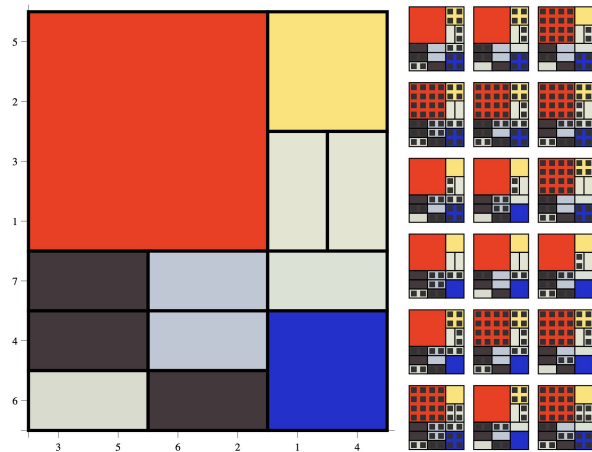


DPP



Independent

[Kulesza&Taskar, 2011]



[Roy&Teh, 2009]

# Potential issues with existing methods

- Computational cost is usually high.
- Not all adapt to asynchronous parallel BayesOpt settings.
- Difficult to debug especially in high-dimensional settings.
- Parallel BayesOpt typically co-occur with large scale high-dimensional problems, but a joint solution for these conditions is not yet satisfying.

Challenges, open problems and some attempts

# Unknown priors

# Bayesian optimization with an unknown prior

Which comes first?  
Data or prior?

Estimate “prior” from data

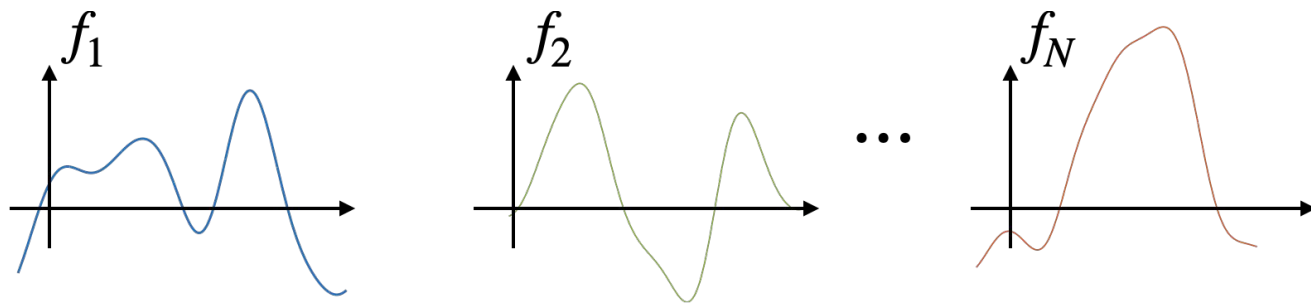
- maximum likelihood
- hierarchical Bayes
- Regret bounds exist only when prior is assumed given
- bad settings of priors make BO perform poorly and seem to be a bad approach



# Bayesian optimization with an unknown prior

meta / multi-task / transfer learning

- Our idea: learn the “prior” from past experience with similar functions
- Assumption: we can collect data on functions sampled from the same prior  $f_1, f_2, \dots, f_N \sim GP(\mu, k)$



# How to learn the GP prior?

$$|\mathcal{X}| = M$$

Use finite input space to illustrate; extensions to continuous case requires more assumptions. [Wang et al., 2018 + ongoing work]

Prior estimation with meta training data  $\{[(x_j, y_{ij})]_{j=1}^M\}_{i=1}^N$

Task 1	(x_1, y_11)	(x_2, y_12)	.....	(x_M, y_1M)
Task 2	(x_1, y_21)	(x_2, y_22)	.....	(x_M, y_2M)
.....	.....	.....	.....	.....
Task N	(x_1, y_N1)	(x_2, y_N2)	.....	(x_M, y_NM)
New Task	?	?	.....	?



# How to estimate the GP prior?

$$|\mathcal{X}| = M$$

Task 1	(x <sub>1</sub> , y <sub>11</sub> )	(x <sub>2</sub> , y <sub>12</sub> )	.....	(x <sub>M</sub> , y <sub>1M</sub> )
Task 2	(x <sub>1</sub> , y <sub>21</sub> )	$Y = \begin{bmatrix} y_{11} & \cdots & y_{1M} \\ \cdots & & \cdots \\ y_{N1} & \cdots & y_{NM} \end{bmatrix}$		(x <sub>M</sub> , y <sub>2M</sub> )
.....	.....			.....
Task N	(x <sub>1</sub> , y <sub>N1</sub> )			(x <sub>M</sub> , y <sub>NM</sub> )
New Task	?		.....	?

Unbiased prior estimator

$$\hat{\mu}(\mathcal{X}) = \frac{1}{N} Y^T \mathbf{1}_N \sim \mathcal{N}(\mu(\mathcal{X}), \frac{1}{N} (k(\mathcal{X}) + \sigma^2 \mathbf{I}))$$

$$\hat{k}(\mathcal{X}) = \frac{1}{N-1} (Y - \mathbf{1}_N \hat{\mu}(\mathcal{X})^T)^T (Y - \mathbf{1}_N \hat{\mu}(\mathcal{X})^T) \sim \mathcal{W}(\frac{1}{N-1} (k(\mathcal{X}) + \sigma^2 \mathbf{I}), N-1)$$

# How to estimate the GP posterior?

$$|\mathcal{X}| = M$$

Task 1	(x_1, y_11)	(x_2, y_12)	.....	(x_M, y_1M)	
Task 2	(x_1, y_21)	$Y = \begin{bmatrix} y_{11} & \cdots & y_{1M} \\ \cdots & & \cdots \\ y_{N1} & \cdots & y_{NM} \end{bmatrix}$			(x_M, y_2M)
.....	.....				.....
Task N	(x_1, y_N1)				(x_M, y_NM)
New Task	?	?	.....	?	

Unbiased posterior estimator

$$\hat{\mu}_t(x) = \hat{\mu}(x) + \hat{k}(x, \mathbf{x}_t) \hat{k}(\mathbf{x}_t, \mathbf{x}_t)^{-1} (y_t - \hat{\mu}(\mathbf{x}_t))$$

$$\hat{\sigma}_t^2(x, x') = \frac{N-1}{N-t-1} \left( \hat{k}(x, x') - \hat{k}(x, \mathbf{x}_t) \hat{k}(\mathbf{x}_t, \mathbf{x}_t)^{-1} \hat{k}(\mathbf{x}_t, x') \right)$$

# Regret bound without the knowledge of the GP prior

$$\text{regret: } r_T = \max_{x \in \mathcal{X}} f(x) - \max_{t \in [T]} f(x_t)$$

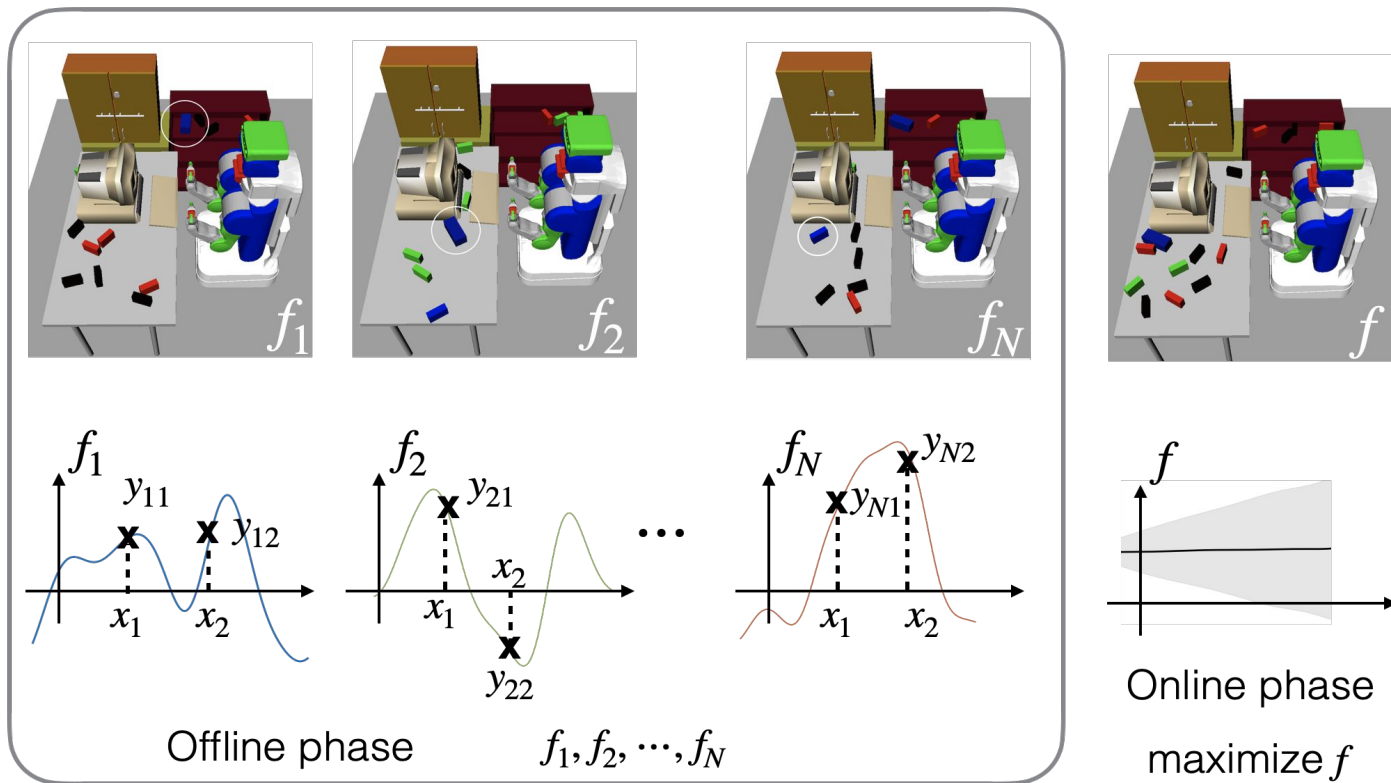
Important assumptions:

- functions are sampled from the same Gaussian process
- enough number of functions in offline phase  $N > T + 20$

Given  $T$  observations on the new function  $f$ , with probability  $1 - \delta$ ,

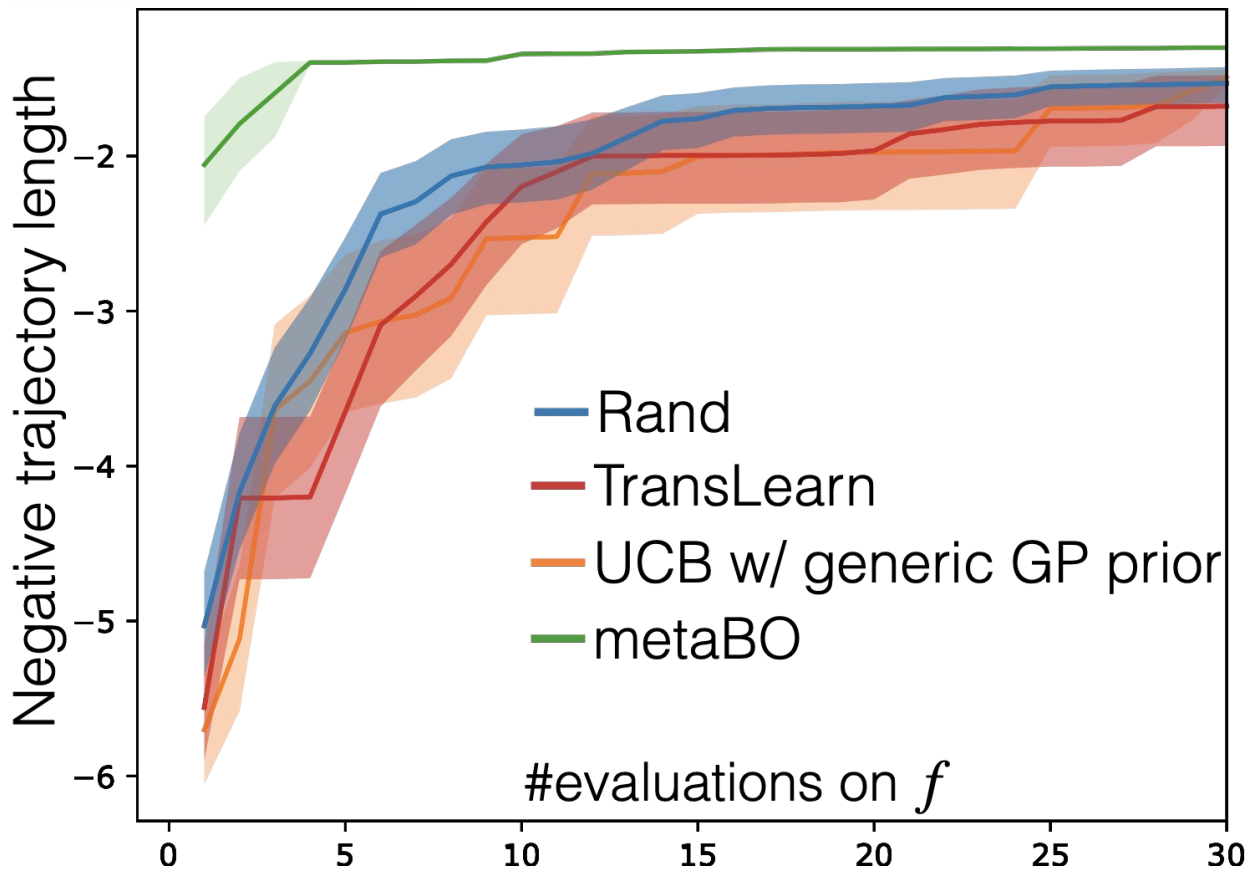
$$\text{regret } r_T \leq O \left( \left( \sqrt{\frac{1}{N-T}} + \underset{\substack{\text{constant} \\ \text{depending on } \delta}}{C} \right) \left( \sqrt{\frac{\log T}{T}} + \underset{\substack{\text{observation noise}}}{\sigma} \right) \right) \rightarrow C\sigma$$

# Empirical results on block picking and placing



# MetaBO gives better performance with fewer samples

[Wang et al., NeurIPS 2018]



Challenges, open problems and some attempts

# Applications in robotics

Type of tasks:

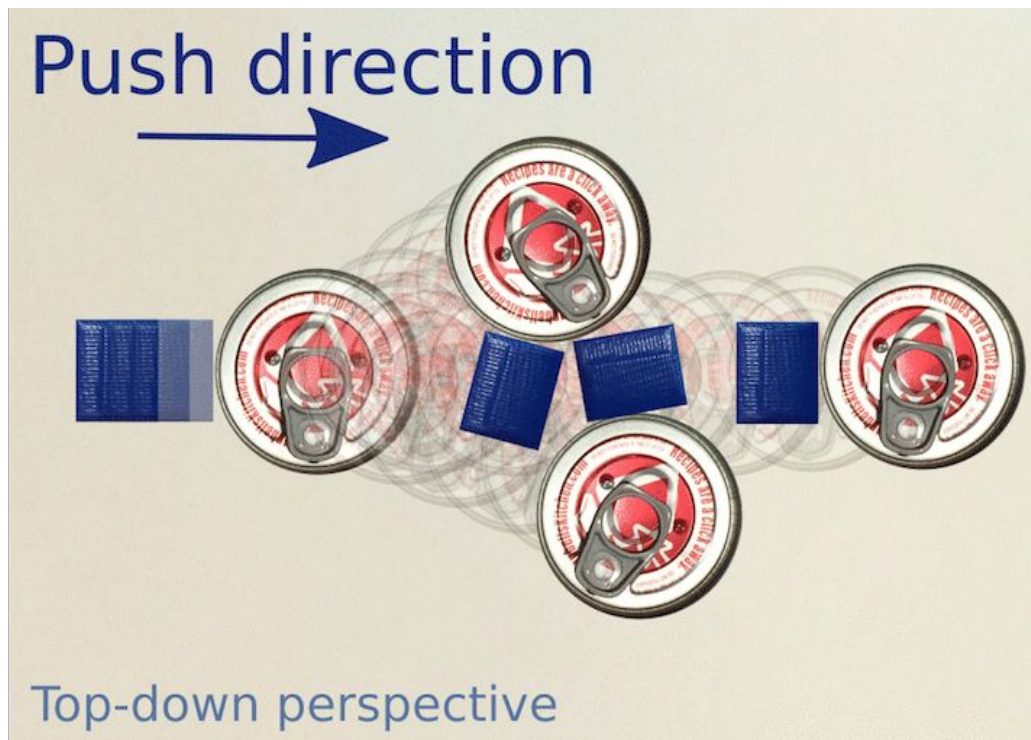
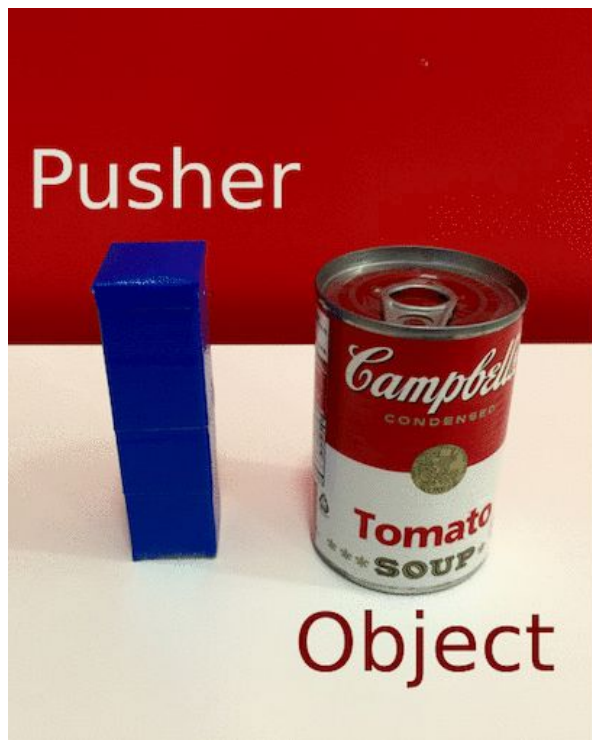
complex

long-horizon

stochastic

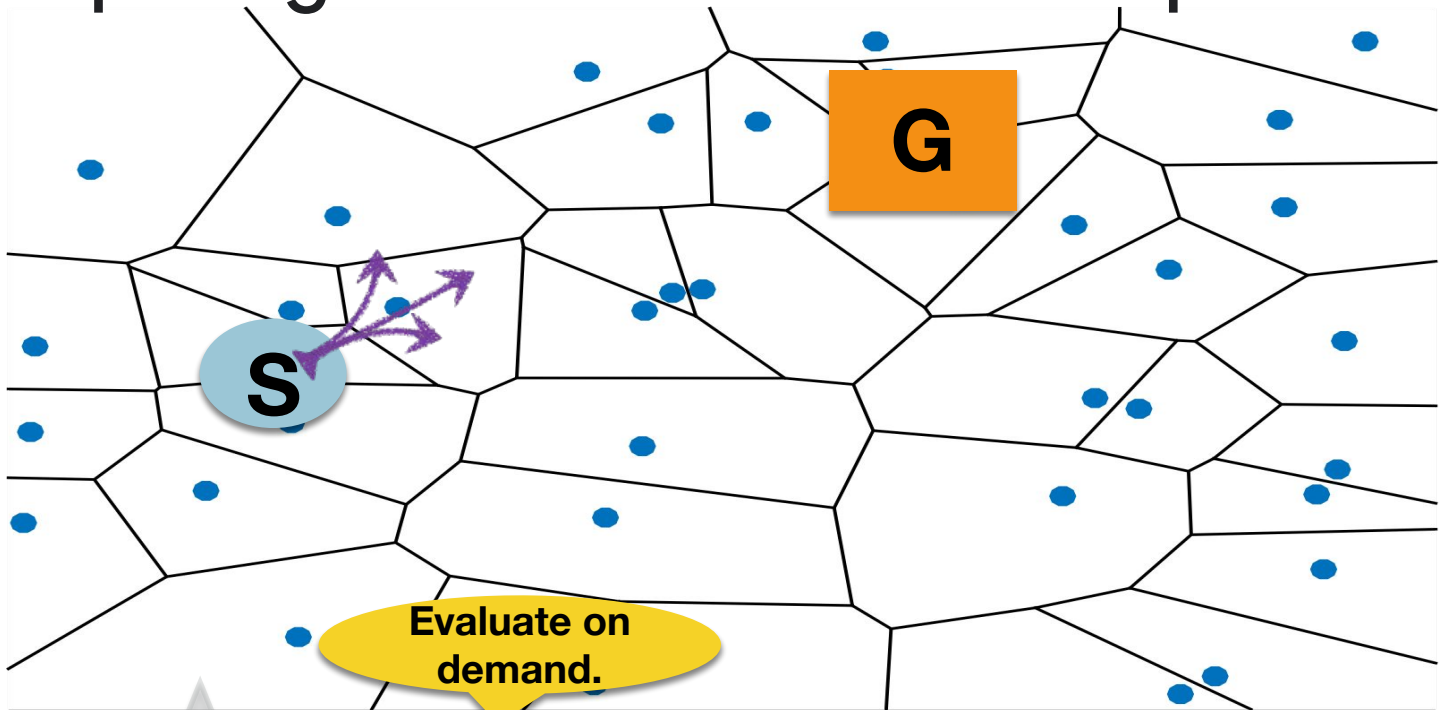
fully-observable

# Learning a pushing skill with multi-modal dynamics





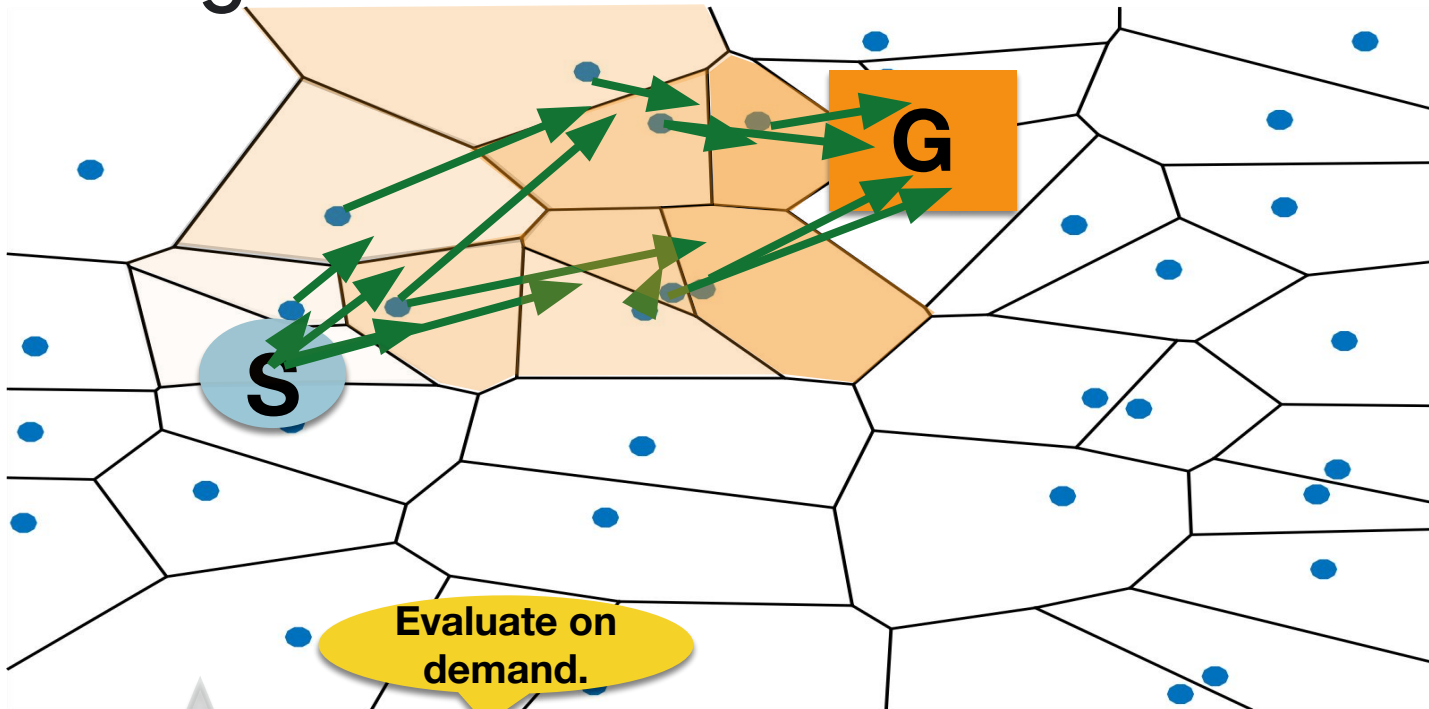
# Computing the action values is expensive



$$= \arg \max_a \sum_{s' \in \tilde{S}} P_{s'|s,a}(s' | s, a) (R(s' | s, a) + \gamma^{\Delta t} V(s'))$$

Solve via BO

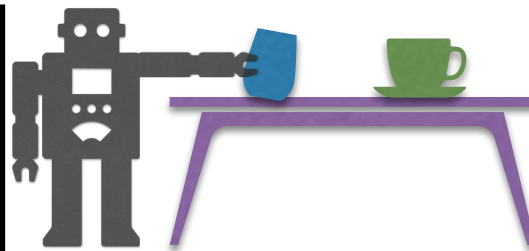
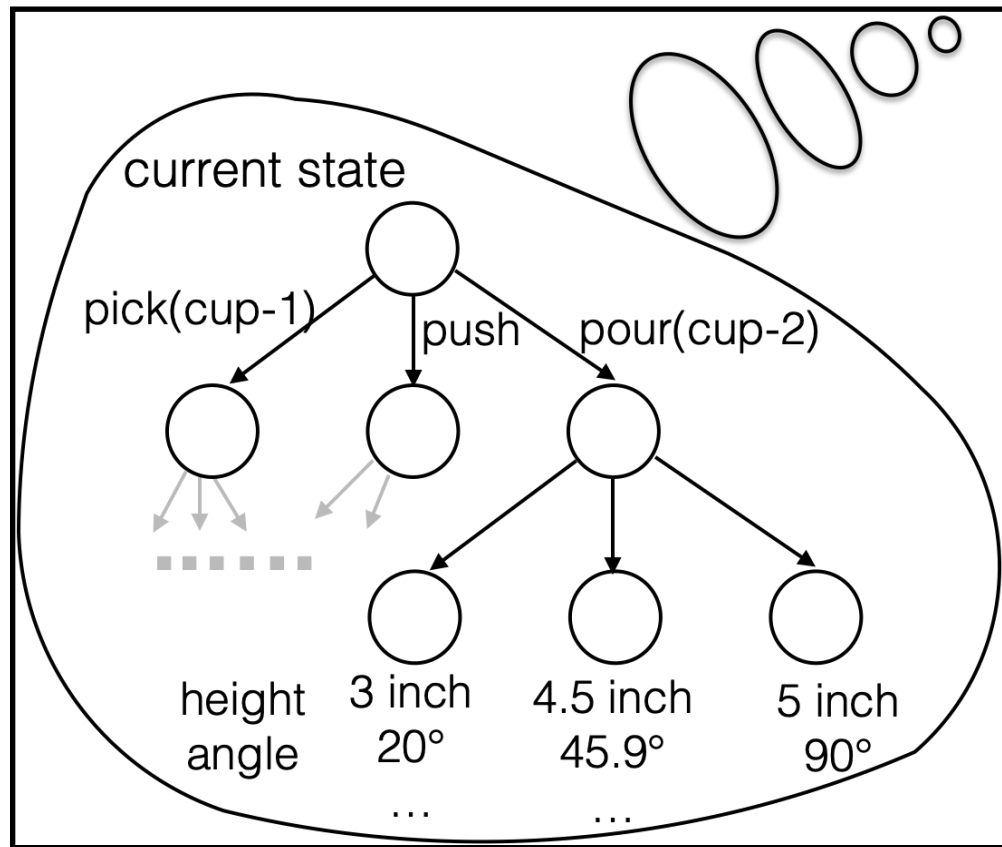
# Focusing on relevant states with RTDP



$$= \arg \max_a \sum_{s' \in \tilde{S}} P_{s'|s,a}(s' | s, a) (R(s' | s, a) + \gamma^{\Delta t} V(s'))$$

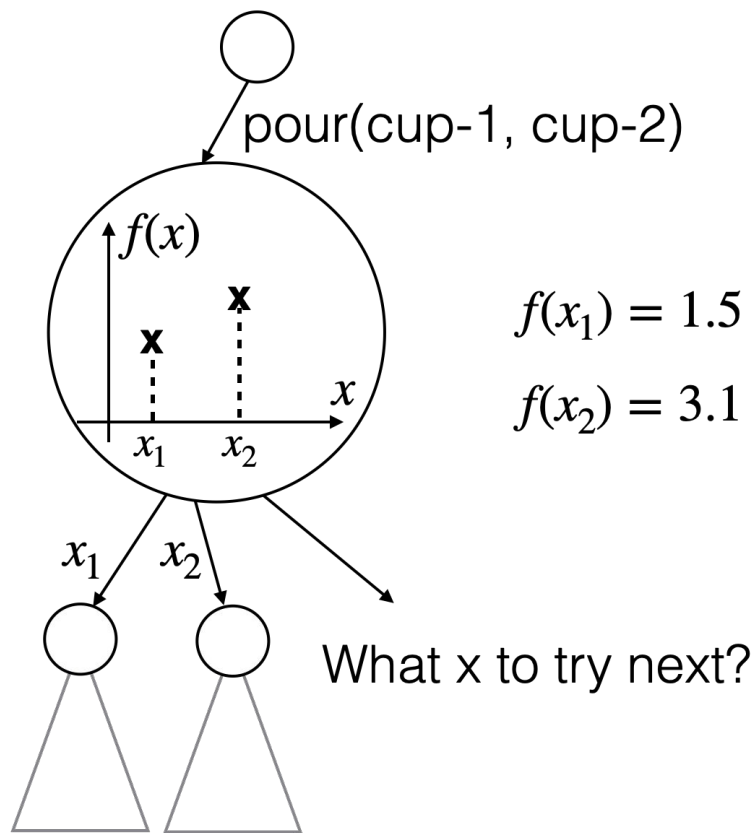
Solve via BO

# How to plan with learned skills?

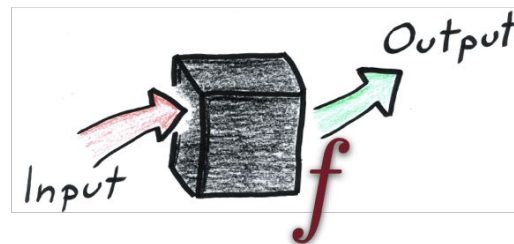


- Sample continuous skill parameters
- Tree search on both discrete and continuous variables

# How to sample skill parameters for the planner?



**Treat the problem of sampling skill parameters as a black-box function optimization problem**

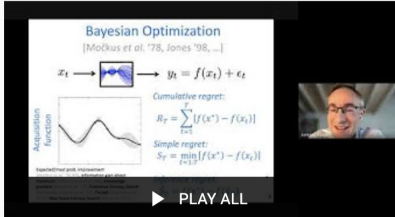


maximize  $f(x)$   
 $x \in \mathcal{X}$

# More talks on BayesOpt

## 2021 Google BayesOpt Speaker Series

now publicly available at [Google TechTalks](#) channel on YouTube




Bayesian Optimization  
[MoKus et al. '78, Jones '98, ...]  
 $x_t \rightarrow y_t = f(x_t) + \epsilon_t$   
Cumulative regret:  
 $R_T = \sum_{t=1}^T [f(x^*) - f(x_t)]$   
Simple regret:  
 $S_T = \min_{x \in \mathcal{X}} [f(x^*) - f(x_T)]$


PLAY ALL


2021 Google BayesOpt Speaker Series


4 videos • 339 views • Last updated on Jul 13, 2021


≡ ↗ ↶ ...

 Google TechTalks [SUBSCRIBE](#)

- 

1 Efficient Exploration in Bayesian Optimization – Optimism and Beyond by Andreas Krause  
Google TechTalks  
1:15:19
- 

2 Learning to Explore in Molecule Space by Yoshua Bengio  
Google TechTalks  
1:05:17
- 

3 Grey-box Bayesian Optimization by Peter Frazier  
Google TechTalks  
1:17:30
- 

4 Resource Allocation in Multi-armed Bandits by Kirthivasan Kandasamy  
Google TechTalks  
59:48

# Questions?