

# Bayesian Optimization for Hyperparameter Optimization

## Computationally Cheap Acquisition Functions

Bernd Bischl   Frank Hutter   Lars Kotthoff  
Marius Lindauer   Joaquin Vanschoren

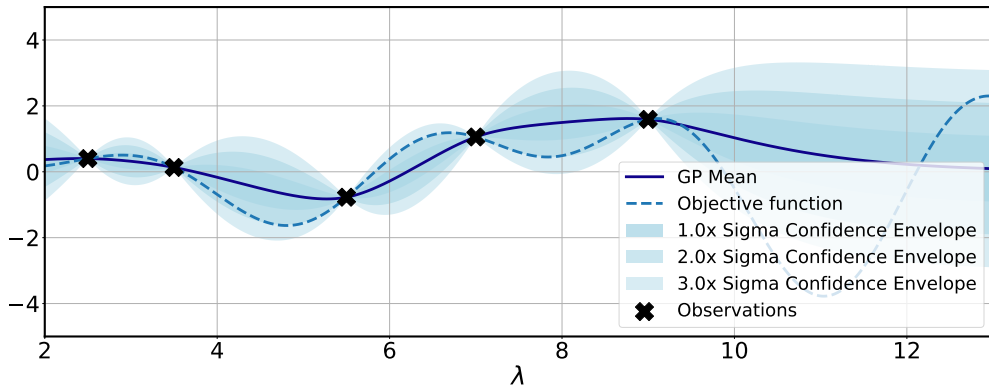
# Acquisition Functions: the Basics

- Given the surrogate model  $\hat{c}^{(t)}$  at the  $t$ -th iteration of BO, the acquisition function  $u(\cdot)$  judges the utility (or usefulness) of evaluating  $f$  at  $\lambda^{(t)} \in \Lambda$  next

# Acquisition Functions: the Basics

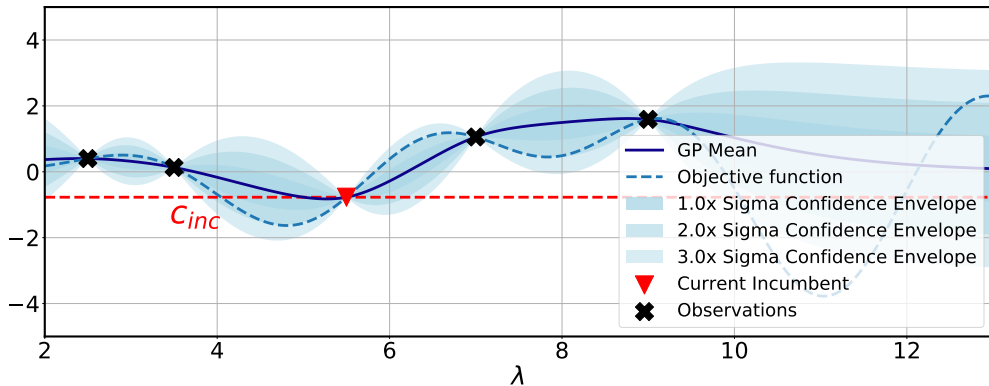
- Given the surrogate model  $\hat{c}^{(t)}$  at the  $t$ -th iteration of BO, the acquisition function  $u(\cdot)$  judges the utility (or usefulness) of evaluating  $f$  at  $\lambda^{(t)} \in \Lambda$  next
- The acquisition function needs to trade off exploration and exploitation
  - ▶ E.g., just picking the  $\lambda$  with lowest predicted mean would be too greedy
  - ▶ We also need to take into account the uncertainty of the surrogate model  $\hat{c}^{(t)}$  to explore

# Probability of Improvement (PI): Concept



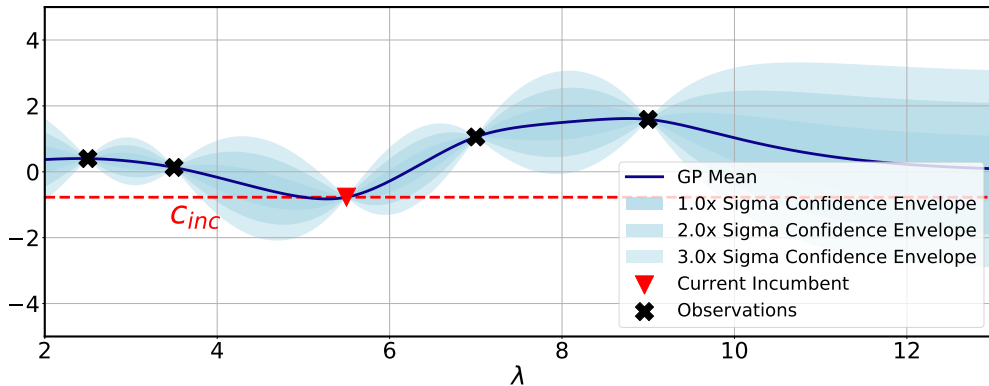
Given the surrogate fit at iteration  $t$

# Probability of Improvement (PI): Concept



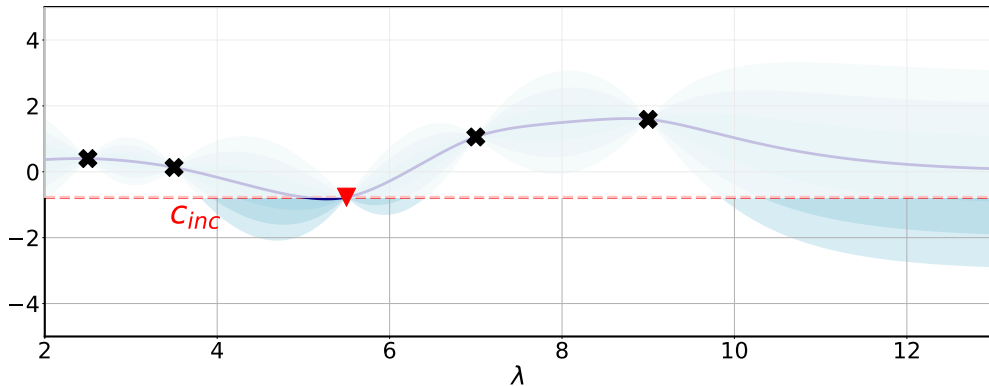
Current incumbent  $\hat{\lambda}$  and its observed cost  $c_{inc}$

# Probability of Improvement (PI): Concept



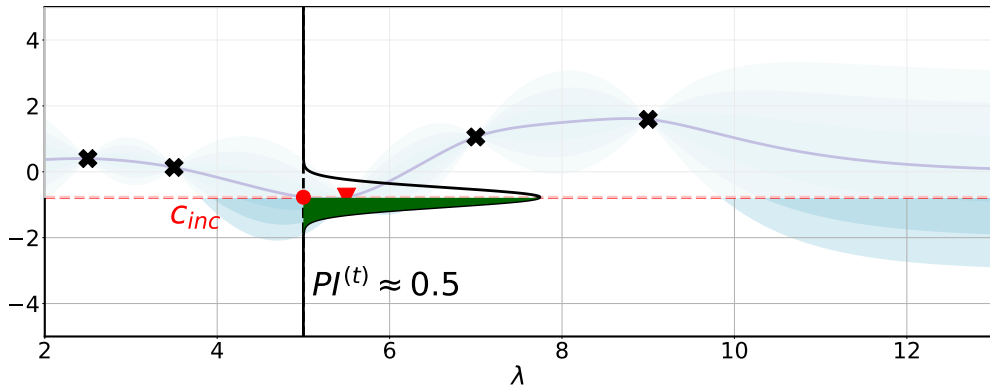
Now let's drop the objective function - it's unknown after all!

# Probability of Improvement (PI): Concept



Intuitively, we care about the probability of improving over the current incumbent

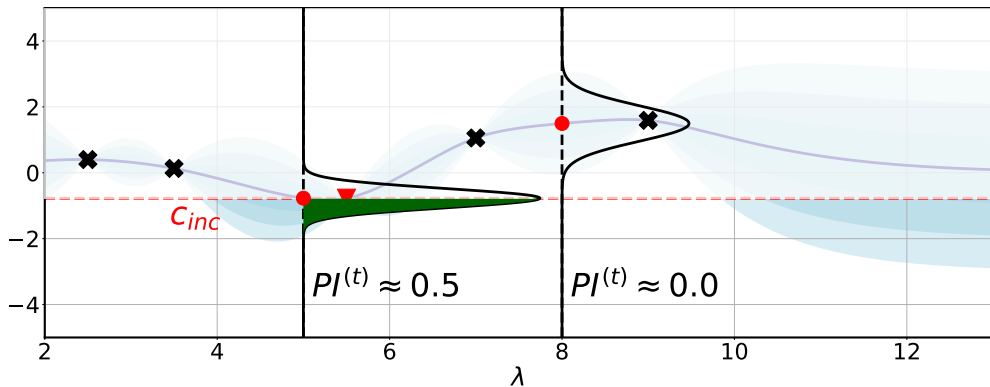
## Probability of Improvement (PI): Concept



PDF of a good candidate configuration. Only the green area is an improvement.



# Probability of Improvement (PI): Concept



PDF of a bad candidate configuration

## Probability of Improvement (PI): Formal Definition

- We define the **current incumbent** at time step  $t$  as:  $\hat{\lambda}^{(t-1)} \in \arg \min_{\lambda' \in \mathcal{D}^{(t-1)}} c(\lambda')$
- We write  $c_{inc}$  shorthand for the **cost of the current incumbent**:  $c_{inc} = c(\hat{\lambda}^{(t-1)})$
- The **probability of improvement**  $u_{PI}(\lambda)$  at a configuration  $\lambda$  is then defined as:

$$u_{PI}^{(t)}(\lambda) = P(c(\lambda) \leq c_{inc}).$$

# Probability of Improvement (PI): Formal Definition

- We define the **current incumbent** at time step  $t$  as:  $\hat{\lambda}^{(t-1)} \in \arg \min_{\lambda' \in \mathcal{D}^{(t-1)}} c(\lambda')$
- We write  $c_{inc}$  shorthand for the **cost of the current incumbent**:  $c_{inc} = c(\hat{\lambda}^{(t-1)})$
- The **probability of improvement**  $u_{PI}(\lambda)$  at a configuration  $\lambda$  is then defined as:

$$u_{PI}^{(t)}(\lambda) = P(c(\lambda) \leq c_{inc}).$$

- Since the predictive distribution for  $c(\lambda)$  is a Gaussian  $\mathcal{N}(\mu^{(t-1)}(\lambda), \sigma^{2(t-1)}(\lambda))$ , this can be written as:

$$u_{PI}^{(t)}(\lambda) = \Phi[Z], \quad \text{with } Z = \frac{c_{inc} - \mu^{(t-1)}(\lambda) - \xi}{\sigma^{(t-1)}(\lambda)},$$

where  $\Phi(\cdot)$  is the CDF of the standard normal distribution and  $\xi$  is an optional exploration parameter

# Probability of Improvement (PI): Formal Definition

- We define the **current incumbent** at time step  $t$  as:  $\hat{\lambda}^{(t-1)} \in \arg \min_{\lambda' \in \mathcal{D}^{(t-1)}} c(\lambda')$
- We write  $c_{inc}$  shorthand for the **cost of the current incumbent**:  $c_{inc} = c(\hat{\lambda}^{(t-1)})$
- The **probability of improvement**  $u_{PI}(\lambda)$  at a configuration  $\lambda$  is then defined as:

$$u_{PI}^{(t)}(\lambda) = P(c(\lambda) \leq c_{inc}).$$

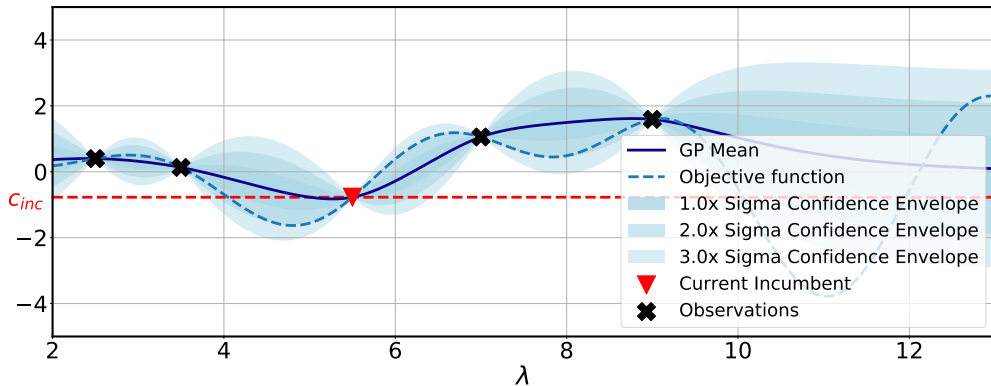
- Since the predictive distribution for  $c(\lambda)$  is a Gaussian  $\mathcal{N}(\mu^{(t-1)}(\lambda), \sigma^{2(t-1)}(\lambda))$ , this can be written as:

$$u_{PI}^{(t)}(\lambda) = \Phi[Z], \quad \text{with } Z = \frac{c_{inc} - \mu^{(t-1)}(\lambda) - \xi}{\sigma^{(t-1)}(\lambda)},$$

where  $\Phi(\cdot)$  is the CDF of the standard normal distribution and  $\xi$  is an optional exploration parameter

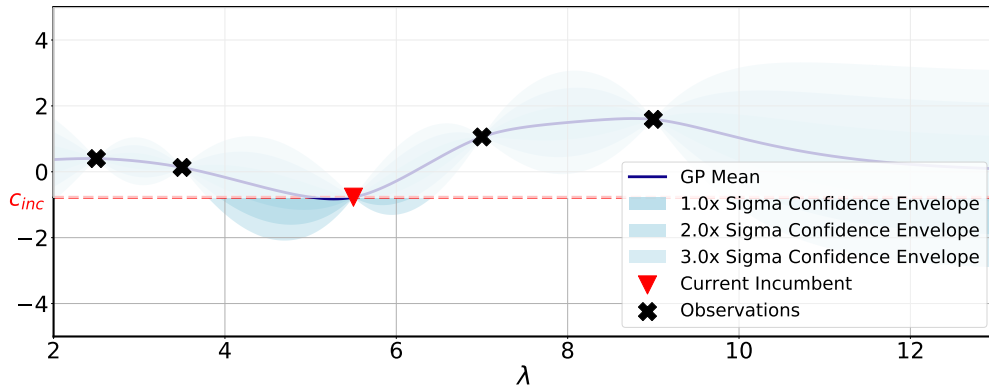
Choose  $\lambda^{(t)} \in \arg \max_{\lambda \in \Lambda} (u_{PI}^{(t)}(\lambda))$

# Expected Improvement (EI): Concept



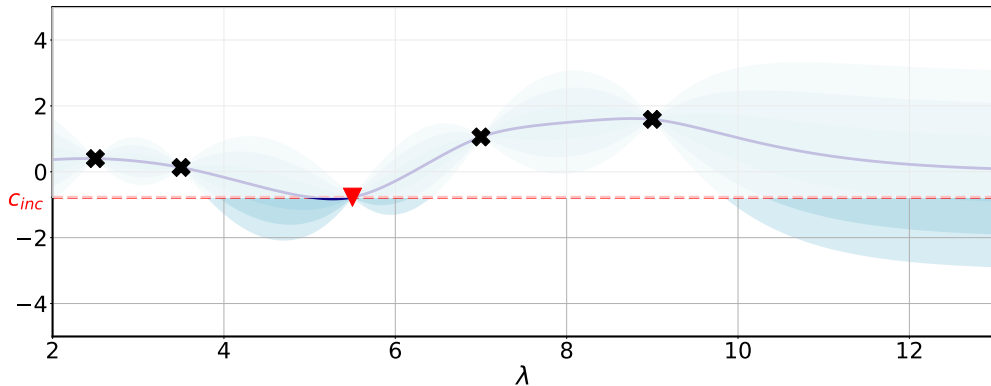
Given the surrogate fit at iteration  $t$

# Expected Improvement (EI): Concept



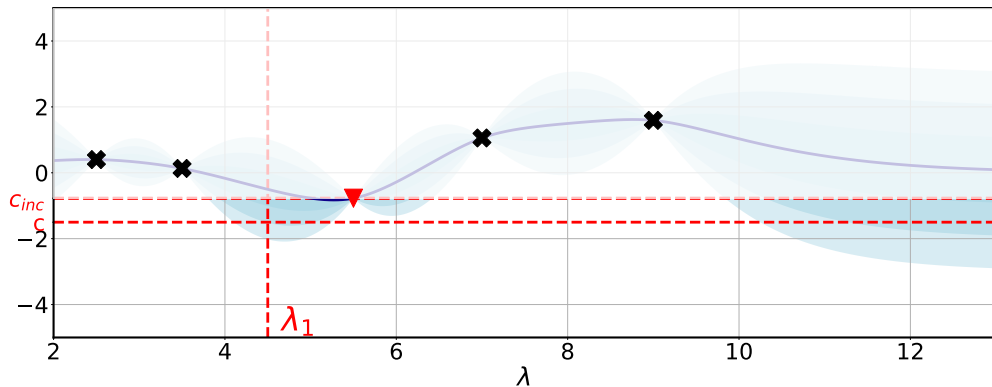
Region of probable improvement – but **how large** is the improvement?

# Expected Improvement (EI): Concept



Region of probable improvement – but **how large** is the improvement?

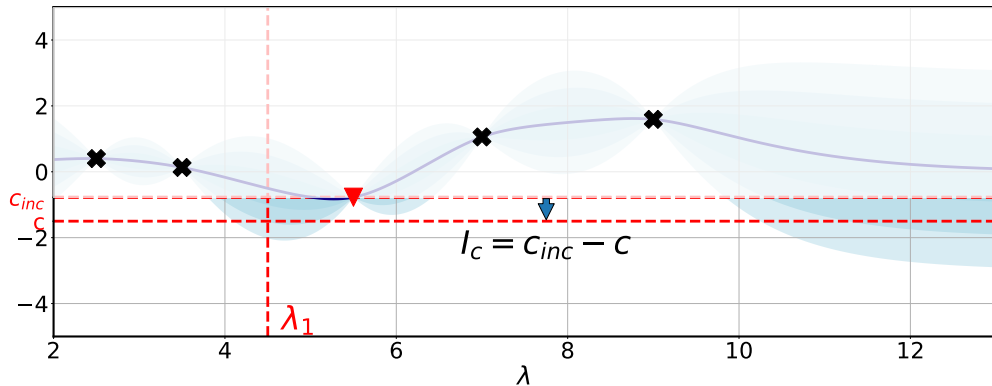
# Expected Improvement (EI): Concept



Hypothetical *real* cost  $c$  at a given  $\lambda$  - unknown in practice without evaluating

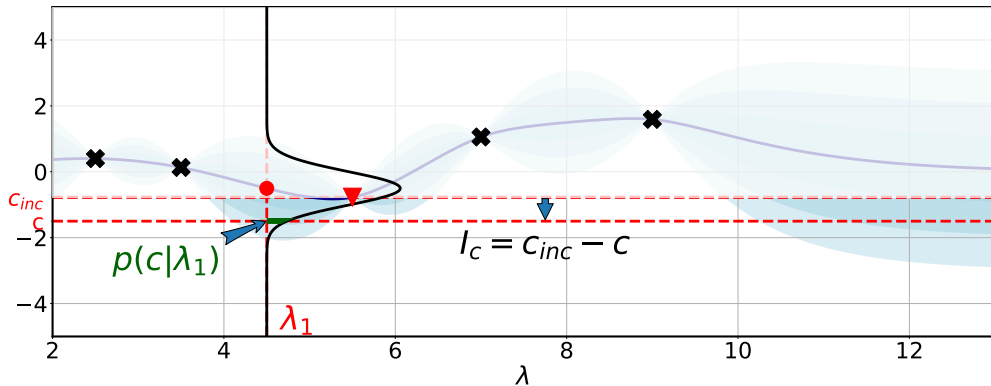


## Expected Improvement (EI): Concept



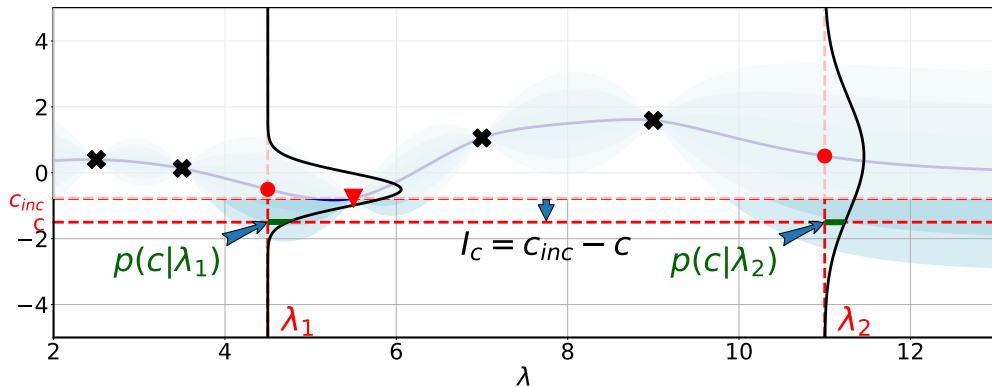
Given a hypothetical  $c$ , we can improve the improvement  $I_c(\lambda)$

# Expected Improvement (EI): Concept



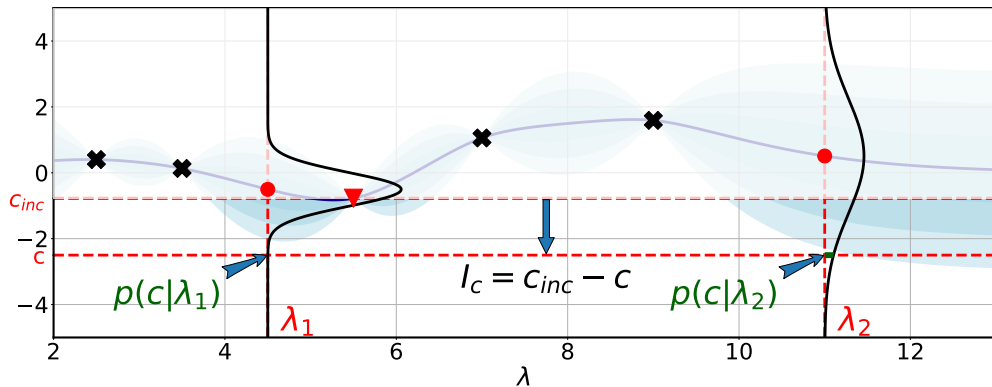
Given  $\hat{c}(\lambda) = \mathcal{N}(\mu(\lambda), \sigma^2(\lambda))$ , we can also compute  $p(c|\lambda)$ .

# Expected Improvement (EI): Concept



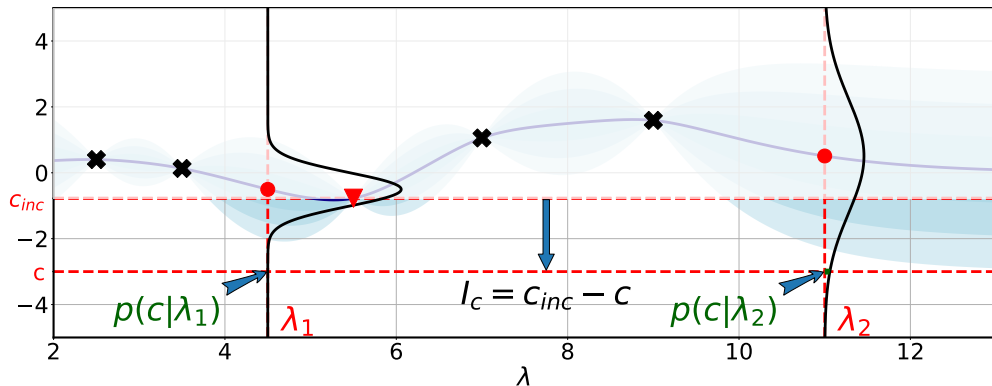
Compare the likelihood of a given improvement for two different configurations  $\lambda_1$  and  $\lambda_2$

## Expected Improvement (EI): Concept



Now consider the likelihood of a larger improvement.

# Expected Improvement (EI): Concept



Larger improvements are more likely in areas of high uncertainty.  
To compute  $\mathbb{E}[I(\lambda)]$ , intuitively, we sum  $p(c | \lambda) \times I_c$  over all possible values of  $c$ .

# Expected Improvement (EI): Formal Definition

- We define the one-step positive improvement over the current incumbent as

$$I^{(t)}(\boldsymbol{\lambda}) = \max(0, c_{inc} - c(\boldsymbol{\lambda}))$$

- Expected Improvement is then defined as

$$u_{EI}^{(t)}(\boldsymbol{\lambda}) = \mathbb{E}[I^{(t)}(\boldsymbol{\lambda})] = \int_{-\infty}^{\infty} p^{(t)}(c \mid \boldsymbol{\lambda}) \times I^{(t)}(\boldsymbol{\lambda}) \, dc.$$

# Expected Improvement (EI): Formal Definition

- We define the one-step positive **improvement over the current incumbent** as

$$I^{(t)}(\boldsymbol{\lambda}) = \max(0, c_{inc} - c(\boldsymbol{\lambda}))$$

- Expected Improvement is then defined as

$$u_{EI}^{(t)}(\boldsymbol{\lambda}) = \mathbb{E}[I^{(t)}(\boldsymbol{\lambda})] = \int_{-\infty}^{\infty} p^{(t)}(c \mid \boldsymbol{\lambda}) \times I^{(t)}(\boldsymbol{\lambda}) \, dc.$$

- Since the posterior distribution of  $\hat{c}(\boldsymbol{\lambda})$  is a Gaussian, EI can be computed in closed form (see exercise):

$$u_{EI}^{(t)}(\boldsymbol{\lambda}) = \begin{cases} \sigma^{(t)}(\boldsymbol{\lambda})[Z\Phi(Z) + \phi(Z)], & \text{if } \sigma^{(t)}(\boldsymbol{\lambda}) > 0 \\ 0 & \text{if } \sigma^{(t)}(\boldsymbol{\lambda}) = 0, \end{cases}$$

where  $Z = \frac{c_{inc} - \mu^{(t)}(\boldsymbol{\lambda}) - \xi}{\sigma^{(t)}(\boldsymbol{\lambda})}$  and  $\xi$  is an optional exploration parameter.

# Expected Improvement (EI): Formal Definition

- We define the one-step positive **improvement over the current incumbent** as

$$I^{(t)}(\boldsymbol{\lambda}) = \max(0, c_{inc} - c(\boldsymbol{\lambda}))$$

- Expected Improvement is then defined as

$$u_{EI}^{(t)}(\boldsymbol{\lambda}) = \mathbb{E}[I^{(t)}(\boldsymbol{\lambda})] = \int_{-\infty}^{\infty} p^{(t)}(c \mid \boldsymbol{\lambda}) \times I^{(t)}(\boldsymbol{\lambda}) \, dc.$$

- Since the posterior distribution of  $\hat{c}(\boldsymbol{\lambda})$  is a Gaussian, EI can be computed in closed form (see exercise):

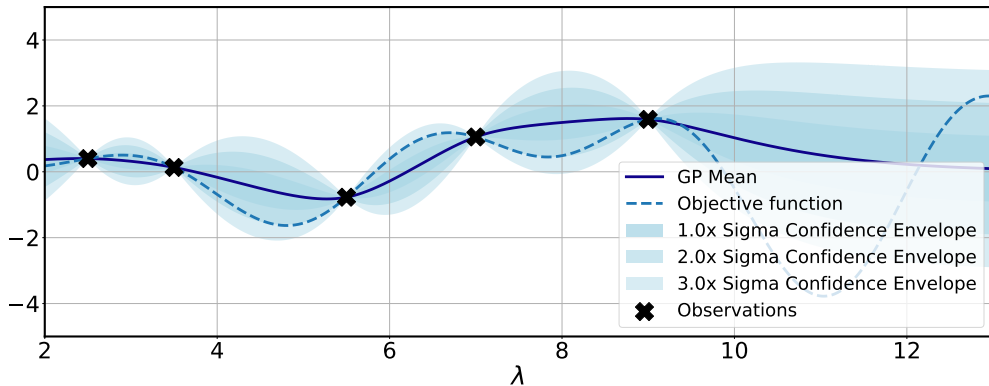
$$u_{EI}^{(t)}(\boldsymbol{\lambda}) = \begin{cases} \sigma^{(t)}(\boldsymbol{\lambda})[Z\Phi(Z) + \phi(Z)], & \text{if } \sigma^{(t)}(\boldsymbol{\lambda}) > 0 \\ 0 & \text{if } \sigma^{(t)}(\boldsymbol{\lambda}) = 0, \end{cases}$$

where  $Z = \frac{c_{inc} - \mu^{(t)}(\boldsymbol{\lambda}) - \xi}{\sigma^{(t)}(\boldsymbol{\lambda})}$  and  $\xi$  is an optional exploration parameter.

Choose  $\boldsymbol{\lambda}^{(t)} \in \arg \max_{\boldsymbol{\lambda} \in \Lambda} (u_{EI}^{(t)}(\boldsymbol{\lambda}))$

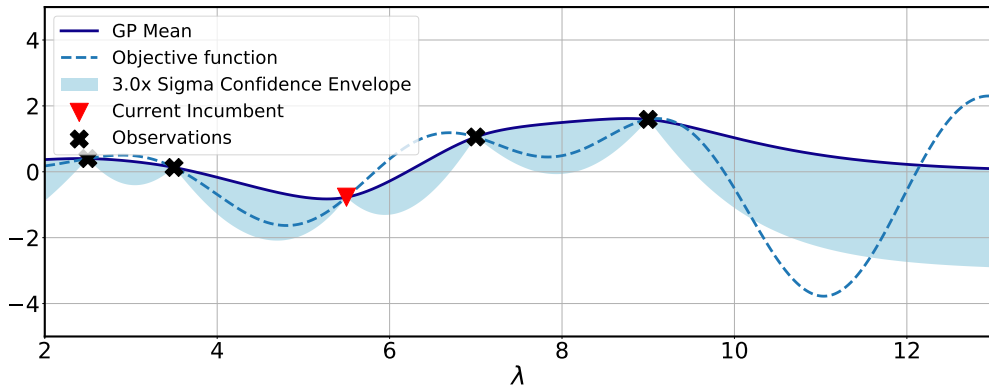


# Lower/Upper Confidence Bounds (LCB/UCB): Concept



Given the surrogate fit at iteration  $t$

# Lower/Upper Confidence Bounds (LCB/UCB): Concept



Lower Confidence Bound,  $\mu(\lambda) - \alpha\sigma(\lambda)$  (here, for  $\alpha = 3$ )

## Lower/Upper Confidence Bounds (LCB/UCB): Formal Definition

- We define the Lower Confidence Bound as

$$u_{LCB}^{(t)}(\boldsymbol{\lambda}) = \mu^{(t)}(\boldsymbol{\lambda}) - \alpha \sigma^{(t)}(\boldsymbol{\lambda}), \quad \alpha \geq 0$$

- One can schedule  $\alpha$  (e.g., increase it over time [Srinivas et al. 2009])

Choose  $\boldsymbol{\lambda}^{(t)} \in \arg \max_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \left( -u_{LCB}^{(t)}(\boldsymbol{\lambda}) \right)$

# Lower/Upper Confidence Bounds (LCB/UCB): Formal Definition

- We define the **Lower Confidence Bound** as

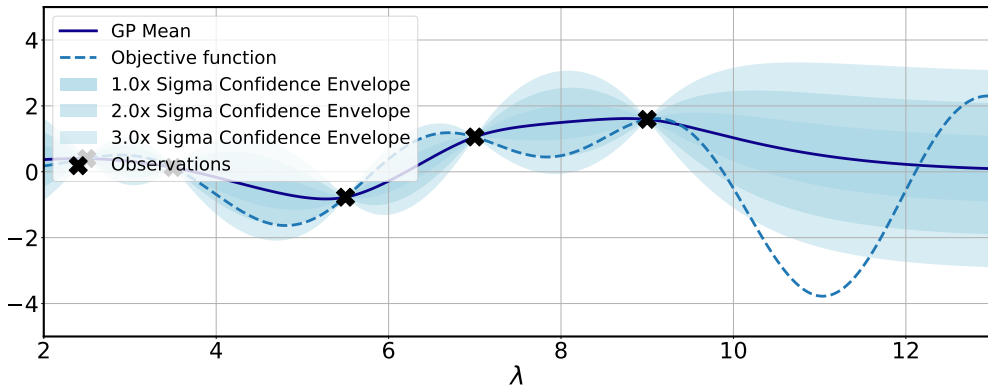
$$u_{LCB}^{(t)}(\boldsymbol{\lambda}) = \mu^{(t)}(\boldsymbol{\lambda}) - \alpha \sigma^{(t)}(\boldsymbol{\lambda}), \quad \alpha \geq 0$$

- One can schedule  $\alpha$  (e.g., increase it over time [Srinivas et al. 2009])

Choose  $\boldsymbol{\lambda}^{(t)} \in \arg \max_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \left( -u_{LCB}^{(t)}(\boldsymbol{\lambda}) \right)$

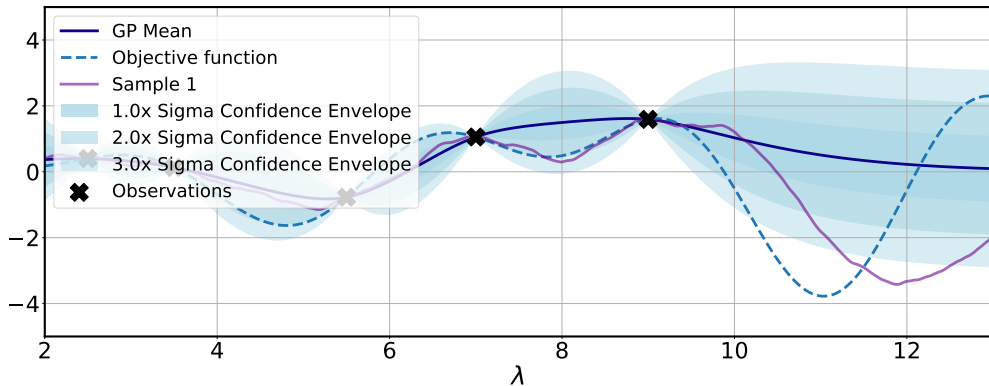
- Note: when one aims to **maximize** the objective function, one would use **UCB** instead
  - ▶  $u_{UCB}^{(t)}(\boldsymbol{\lambda}) = \mu^{(t)}(\boldsymbol{\lambda}) + \alpha \sigma^{(t)}(\boldsymbol{\lambda})$
  - ▶ For UCB, one would choose  $\boldsymbol{\lambda}^{(t)} \in \arg \max_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} (u_{UCB}^{(t)}(\boldsymbol{\lambda}))$

# Thompson Sampling (TS): Concept



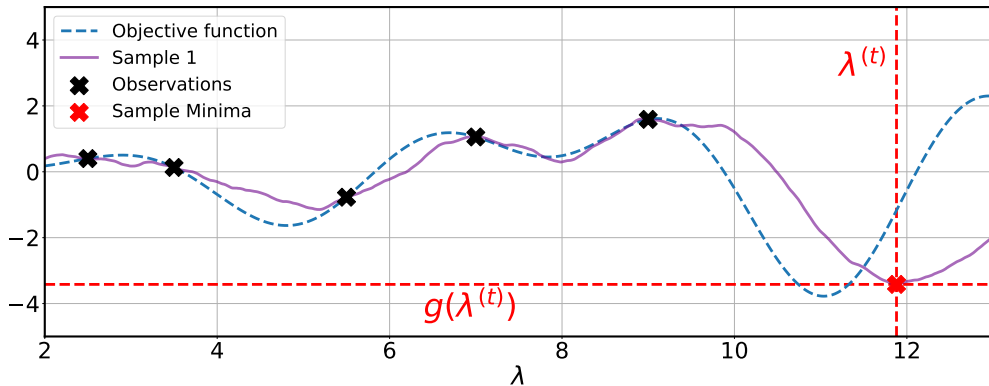
Given the surrogate at iteration  $t$  fit on dataset  $\mathcal{D}^{(t-1)}$

# Thompson Sampling (TS): Concept



Draw a sample  $g$  from the predictive surrogate model

# Thompson Sampling (TS): Concept



Then choose the minimum of this sample to evaluate at next

# Thompson Sampling (TS): Pseudocode

---

Bayesian Optimization using Thompson Sampling

---

**Require:** Search space  $\Lambda$ , cost function  $c$ , surrogate model  $\hat{c}$ , maximal number of function evaluations  $T$

**Result :** Best observed configuration  $\hat{\lambda}$  according to  $\mathcal{D}^{(T)}$  or  $\mathcal{G}$

- 1 Initialize data  $\mathcal{D}^{(0)}$  with initial observations
  - 2 **for**  $t = 1$  **to**  $T$  **do**
    - 3     Fit predictive model  $\hat{c}^{(t)}$  on  $\mathcal{D}^{(t-1)}$
    - 4     Sample a function from the surrogate:  $g \sim \hat{c}^{(t)}$
    - 5     Select next query point:  $\lambda^{(t)} \in \arg \min_{\lambda \in \Lambda} g(\lambda)$
    - 6     Query  $c(\lambda^{(t)})$
    - 7     Update data:  $\mathcal{D}^{(t)} \leftarrow \mathcal{D}^{(t-1)} \cup \{(\lambda^{(t)}, c(\lambda^{(t)}))\}$
-



## Questions to Answer for Yourself / Discuss with Friends

- **Discussion.** How would you set the exploration parameter  $\xi$  for PI if you want to avoid too incremental improvements?
- **Derivation.** Derive the closed form solution of expected improvement.
- **Discussion.** In which situations would EI perform substantially differently than PI?