

Bayesian Optimization for Hyperparameter Optimization

The Alternative Optimization Approach of the Tree-Parzen Estimator (TPE)

Bernd Bischl Frank Hutter Lars Kotthoff
Marius Lindauer Joaquin Vanschoren

Overview of TPE [Bergstra et al. 2011]

- Standard Bayesian optimization models the probability $p(y \mid \lambda)$ of observations y given configurations λ
- Instead, TPE fits kernel density estimators (KDEs) $l(\lambda \mid y \leq \gamma)$ and $g(\lambda \mid y \leq \gamma)$
 - ▶ These KDEs are for “good configurations” (leading to objective function values below a threshold γ) and “bad configurations”
 - ▶ By default, γ is set to the 15% quantile of the observations

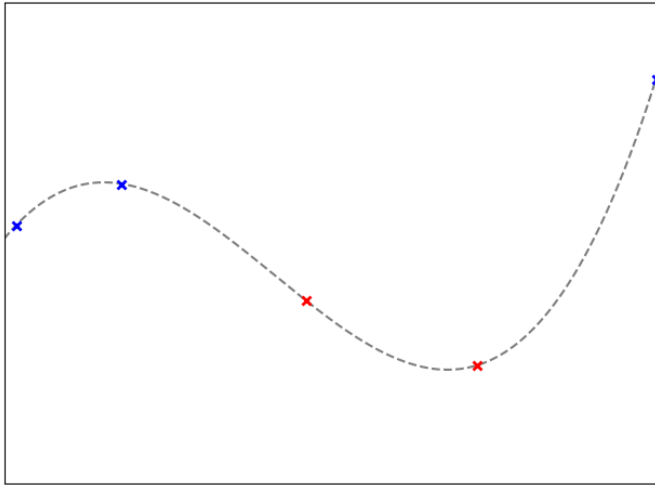
Overview of TPE [Bergstra et al. 2011]

- Standard Bayesian optimization models the probability $p(y \mid \boldsymbol{\lambda})$ of observations y given configurations $\boldsymbol{\lambda}$
- Instead, TPE fits kernel density estimators (KDEs) $l(\boldsymbol{\lambda} \mid y \leq \gamma)$ and $g(\boldsymbol{\lambda} \mid y \leq \gamma)$
 - ▶ These KDEs are for “good configurations” (leading to objective function values below a threshold γ) and “bad configurations”
 - ▶ By default, γ is set to the 15% quantile of the observations
- Optimizing $l(\boldsymbol{\lambda})/g(\boldsymbol{\lambda})$ is equivalent to optimizing standard expected improvement in Bayesian optimization [Bergstra et al. 2011]

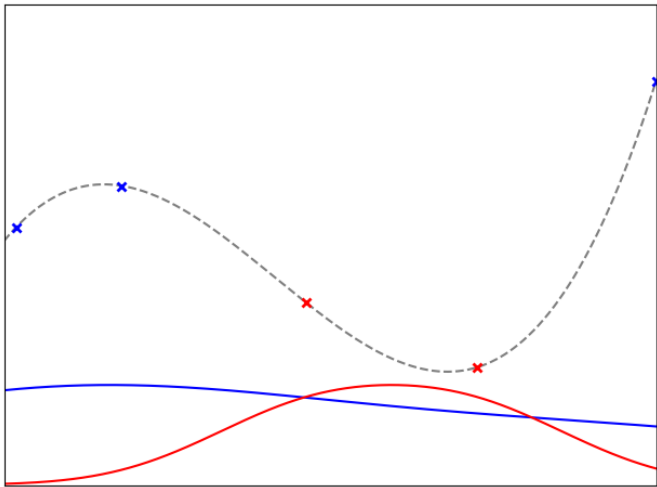
Overview of TPE [Bergstra et al. 2011]

- Standard Bayesian optimization models the probability $p(y \mid \lambda)$ of observations y given configurations λ
- Instead, TPE fits kernel density estimators (KDEs) $l(\lambda \mid y \leq \gamma)$ and $g(\lambda \mid y \leq \gamma)$
 - ▶ These KDEs are for “good configurations” (leading to objective function values below a threshold γ) and “bad configurations”
 - ▶ By default, γ is set to the 15% quantile of the observations
- Optimizing $l(\lambda)/g(\lambda)$ is equivalent to optimizing standard expected improvement in Bayesian optimization [Bergstra et al. 2011]
- Why is the technique called TPE?
 - ▶ The used KDEs are Parzen estimators
 - ▶ TPE can handle tree-structured search spaces

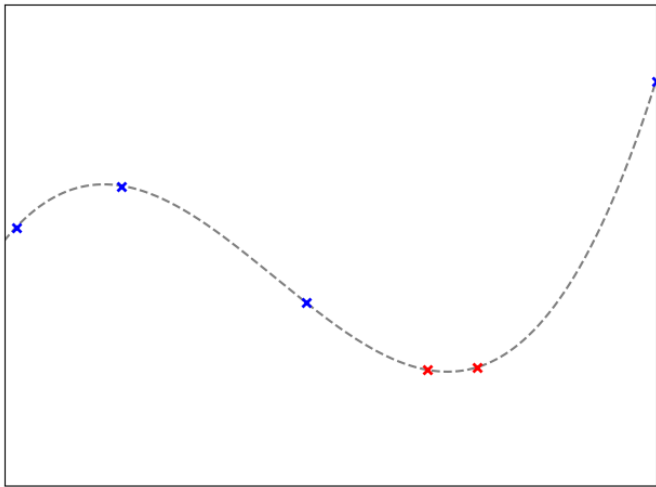
TPE Example



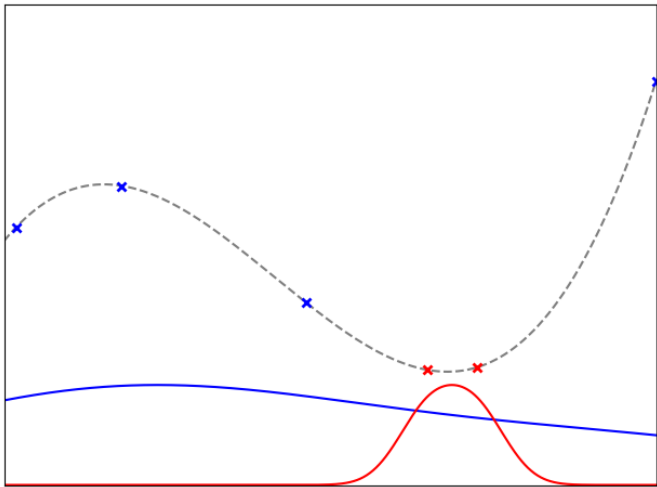
TPE Example



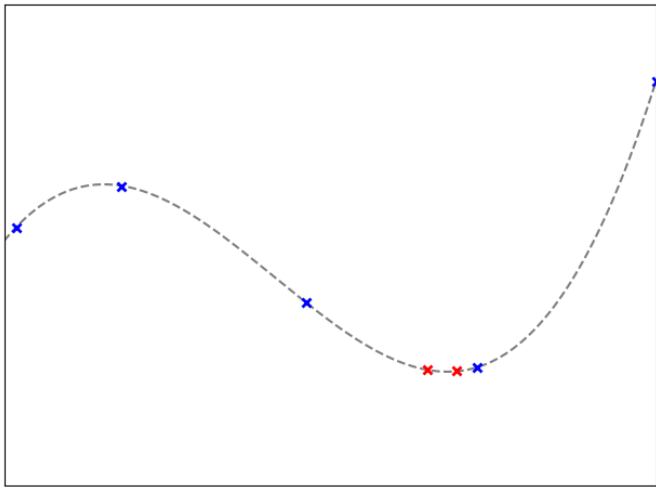
TPE Example



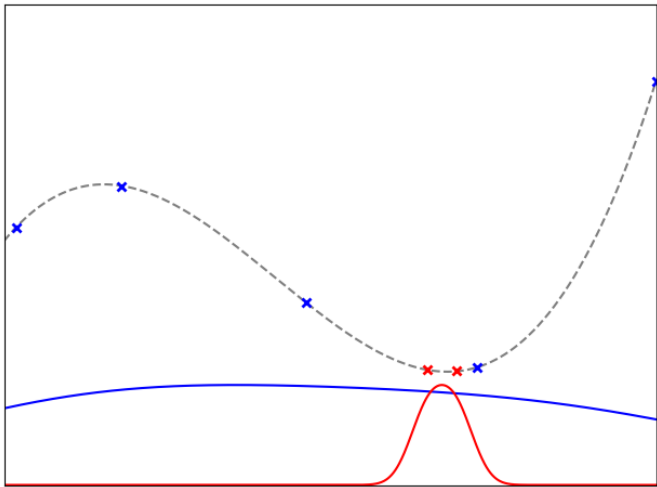
TPE Example



TPE Example



TPE Example



TPE Pseudocode

TPE loop

Require: Search space Λ , cost function c , percentile γ , maximal number of function evaluations T

Result : Best observed configuration λ according to $\mathcal{D}^{(T)}$

```
1  $\mathcal{D}^{(0)} \leftarrow \emptyset$ 
2 for  $t = 1$  to  $T$  do
3    $\mathcal{D}_{\text{good}}, \mathcal{D}_{\text{bad}} \leftarrow \text{split } \mathcal{D}^{(t-1)}$  according to quantile  $\gamma$ 
4    $l(\lambda), g(\lambda) \leftarrow \text{fit KDE on } \mathcal{D}_{\text{good}}, \mathcal{D}_{\text{bad}}$  respectively
5    $\Lambda_{\text{cand}} \leftarrow \text{draw samples from } l$ ;
6   Select next query point:  $\lambda^{(t)} \in \arg \max_{\lambda \in \Lambda_{\text{cand}}} l(\lambda)/g(\lambda)$ 
7   Query  $c(\lambda^{(t)})$ 
8    $\mathcal{D}^{(t)} \leftarrow \mathcal{D}^{(t-1)} \cup \{\langle \lambda^{(t)}, c(\lambda^{(t)}) \rangle\}$ 
9 end
```

Further Details

Remarks:

- TPE models $p(\boldsymbol{\lambda}|c(\boldsymbol{\lambda}))$
 - ▶ we can multiply it with a prior to add expert knowledge

Further Details

Remarks:

- TPE models $p(\boldsymbol{\lambda}|c(\boldsymbol{\lambda}))$
 - ▶ we can multiply it with a prior to add expert knowledge
- Performance of TPE depends on:
 - ▶ setting of γ to trade-off exploration and exploitation
 - ▶ bandwidth of the KDEs

Further Details

Remarks:

- TPE models $p(\boldsymbol{\lambda}|c(\boldsymbol{\lambda}))$
 - ▶ we can multiply it with a prior to add expert knowledge
- Performance of TPE depends on:
 - ▶ setting of γ to trade-off exploration and exploitation
 - ▶ bandwidth of the KDEs
- A successful tool implementing TPE is Hyperopt [Bergstra et al.]

Summary

Advantages

- Computationally efficient: $O(Nd)$
- Parallelizable
- Robust
- Can handle complex search spaces with priors

Summary

Advantages

- Computationally efficient: $O(Nd)$
- Parallelizable
- Robust
- Can handle complex search spaces with priors

Disadvantages

- Less sample-efficient than GPs

Questions to Answer for Yourself / Discuss with Friends

- **Disussion.** Is TPE really Bayesian optimization?
- **Disussion.** How does γ impact the optimization procedure?
- **Disussion.** Derive that optimizing $l(\boldsymbol{\lambda})/g(\boldsymbol{\lambda})$ is equivalent to optimizing expected improvement.