

Hyperparameters Optimization in Large Scale ML

Шугаепов Ильнур

VK.com
Performance Advertising Team

Higher School of Economics, 2020

Examples of hyperparameters

- ▶ Number of trees in GBDT
- ▶ Regularization params in LogReg
- ▶ ...

Why we need HPO?

Examples of hyperparameters

- ▶ Number of trees in GBDT
- ▶ Regularization params in LogReg
- ▶ ...

Goals of HPO

- ▶ Improve the performance of machine learning algorithms

Examples of hyperparameters

- ▶ Number of trees in GBDT
- ▶ Regularization params in LogReg
- ▶ ...

Goals of HPO

- ▶ Improve the performance of machine learning algorithms
- ▶ Reduce the human effort necessary for applying machine learning (AutoML)

- ▶ \mathcal{A} — algorithm with N hyperparameters
- ▶ Λ_n — domain of n -th hyperparameter
- ▶ $\Lambda = \Lambda_1 \times \Lambda_2 \times \dots \times \Lambda_N$ — configuration space
- ▶ \mathcal{A}_λ — algorithm with hyperparameters $\lambda \in \Lambda$

- ▶ \mathcal{A} — algorithm with N hyperparameters
- ▶ Λ_n — domain of n -th hyperparameter
- ▶ $\Lambda = \Lambda_1 \times \Lambda_2 \times \dots \times \Lambda_N$ — configuration space
- ▶ \mathcal{A}_λ — algorithm with hyperparameters $\lambda \in \Lambda$

HPO Problem

Given a dataset \mathcal{D} , our goal is to find

$$\lambda^* = \arg \min_{\lambda \in \Lambda} \mathbb{E}_{(\mathcal{D}_{train}, \mathcal{D}_{valid}) \sim \mathcal{D}} \mathbf{V}(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid}),$$

where $\mathbf{V}(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$ measure the loss of a model generated by \mathcal{A}_λ on training data \mathcal{D}_{train} and evaluated on validation data \mathcal{D}_{valid} .

Algorithm

1. For each $\Lambda_i, i = 1, \dots, N$ choose $L_i \subseteq \Lambda_i$ - set of values
2. $S = \prod_{i=1}^N L_i \subseteq \Lambda$ - set of trial points
3. Evaluate $V(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$ for each $\lambda \in S$ trial point
4. Report best result

Algorithm

1. For each $\Lambda_i, i = 1, \dots, N$ choose $L_i \subseteq \Lambda_i$ - set of values
2. $S = \prod_{i=1}^N L_i \subseteq \Lambda$ - set of trial points
3. Evaluate $V(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$ for each $\lambda \in S$ trial point
4. Report best result

Q: How to compare different HPO algorithms

Algorithm

1. For each $\Lambda_i, i = 1, \dots, N$ choose $L_i \subseteq \Lambda_i$ - set of values
2. $S = \prod_{i=1}^N L_i \subseteq \Lambda$ - set of trial points
3. Evaluate $V(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$ for each $\lambda \in S$ trial point
4. Report best result

Q: How to compare different HPO algorithms **A:** In terms of number of evaluations/wall-clock time

- ▶ Function evaluations can be extremely expensive

- ▶ Function evaluations can be extremely expensive
- ▶ The configuration space is often complex and high-dimensional

- ▶ Function evaluations can be extremely expensive
- ▶ The configuration space is often complex and high-dimensional
- ▶ The configuration space can contain *conditionality*

- ▶ Function evaluations can be extremely expensive
- ▶ The configuration space is often complex and high-dimensional
- ▶ The configuration space can contain *conditionality*
- ▶ We usually don't have access to a gradient of the loss function with respect to the hyperparameters

Challenges of HPO

Example

- ▶ XGBoost training takes on average about 10 hours

Challenges of HPO

Example

- ▶ XGBoost training takes on average about 10 hours
- ▶ Configuration space

```
1 space = {
2     'num_round' : choice(range(200, 1500, 100)), # 14
3     'eta' : quniform(0.05, 0.55, 0.05), # 11
4     'max_depth' : choice(range(4, 10)), # 6
5     'min_child_weight' : quniform(1, 100, 10), # 10
6     'subsample' : quniform(0.5, 1, 0.05), # 11
7     'gamma' : quniform(0, 1, 0.1), # 10
8     'colsample_bytree' : quniform(0.5, 1, 0.1), # 5
9     'alpha' : choice(range(0, 100, 10)), # 11
10    'grow_policy' : choice(['depthwise', 'lossguide']), # 2
11 }
```

Challenges of HPO

Example

- ▶ XGBoost training takes on average about 10 hours
- ▶ Configuration space

```
1 space = {  
2     'num_round' : choice(range(200, 1500, 100)), # 14  
3     'eta' : quniform(0.05, 0.55, 0.05), # 11  
4     'max_depth' : choice(range(4, 10)), # 6  
5     'min_child_weight' : quniform(1, 100, 10), # 10  
6     'subsample' : quniform(0.5, 1, 0.05), # 11  
7     'gamma' : quniform(0, 1, 0.1), # 10  
8     'colsample_bytree' : quniform(0.5, 1, 0.1), # 5  
9     'alpha': choice(range(0, 100, 10)), # 11  
10    'grow_policy' : choice(['depthwise', 'lossguide']), # 2  
11 }
```

- ▶ Number of evaluations required

$$14 \times 6 \times 10 \times 11 \times 10 \times 5 \times 11 \times 2 = 111804000$$

curse of dimensionality

Validation protocol

Reducing the evaluation time

- ▶ $V(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$ — holdout and cross-validation error
- ▶ Strategies for reducing the evaluation time
 - ▶ test algorithms on a subset of folds
 - ▶ only on a subset of data
 - ▶ for a small amount of iterations
 - ▶ auxiliary tasks

1. Random Search
 - Algorithm
 - Random vs Grid Search
2. Bayesian Optimization
 - Surrogate Models
 - Acquisition Functions
 - Noise
3. Multi-fidelity Optimization
 - Learning Curve-Based Methods
 - Hyperband
 - Multi-task Bayesian Optimization
4. HPO Systems
5. Resume

Algorithm

1. Set of trials $S \sim_{\text{iid}} \Lambda$
2. Evaluate $V(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{\text{train}}, \mathcal{D}_{\text{valid}})$ for each $\lambda \in S$ trial point
3. Report best result

Algorithm

1. Set of trials $S \sim_{\text{iid}} \Lambda$
2. Evaluate $V(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{\text{train}}, \mathcal{D}_{\text{valid}})$ for each $\lambda \in S$ trial point
3. Report best result

- ▶ In the simplest scenario $\Lambda = U(\Lambda_1) \times U(\Lambda_2) \times \dots \times U(\Lambda_N)$
- ▶ User can set distribution over Λ_i

Random vs Grid Search

- ▶ Given fixed budget B - number of evaluations
- ▶ The number of different values grid search can afford to evaluate for each of the N hyperparameters is only $B^{1/N}$, whereas random search will explore B different values for each

Random vs Grid Search

- ▶ Given fixed budget B - number of evaluations
- ▶ The number of different values grid search can afford to evaluate for each of the N hyperparameters is only $B^{1/N}$, whereas random search will explore B different values for each

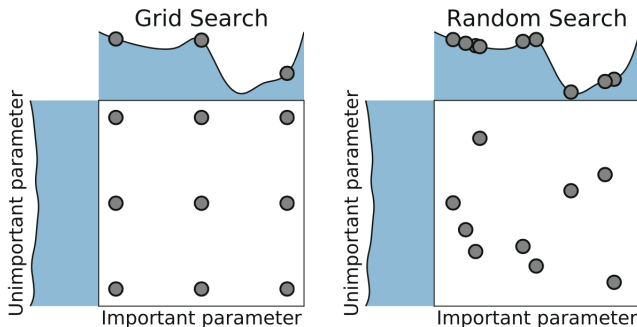


Рис.: Comparison of grid search and random search for minimizing a function with one important and one unimportant parameter. $f(x, y) = g(x) + h(y) \approx g(x)$ with *low effective dimensionality*. Above each square $g(x)$, left of each square $h(y)$

Grid Search is usually worse because:

- ▶ a small fraction of hyperparameters matter for any one data set, but
- ▶ different hyperparameters matter on different data sets

Remark

Failure of grid search is the rule rather than the exception in high dimensional hyperparameter optimization

How to prove that only a few hyperparameters matter? See [2]

Table of Contents

1. Random Search
2. Bayesian Optimization
 - Surrogate Models
 - Acquisition Functions
 - Noise
3. Multi-fidelity Optimization
4. HPO Systems
5. Resume



$$\max_{\mathbf{x} \in \Lambda \subset \mathbb{R}^N} f(\mathbf{x})$$

- ▶ $f: \Lambda \rightarrow \mathbb{R}$ — black-box true objective is costly to evaluate



$$\max_{\mathbf{x} \in \Lambda \subset \mathbb{R}^N} f(\mathbf{x})$$

- ▶ $f: \Lambda \rightarrow \mathbb{R}$ — black-box true objective is costly to evaluate

Algorithm 2.2 SMBO(f, M_0, T, S)

$\mathcal{D}_0 \leftarrow \emptyset$

▷ Observations

for $t = 1, \dots, T$ **do**

$\mathbf{x}^* \leftarrow \arg \min_{\mathbf{x}} S(\mathbf{x}, M_{t-1})$

 Evaluate $f(\mathbf{x}^*)$

▷ Expensive step

$\mathcal{D}_{1:t} \leftarrow \mathcal{D}_{1:t-1} \cup \{(\mathbf{x}^*, f(\mathbf{x}^*))\}$

 Fit a new model M_t to $\mathcal{D}_{1:t}$

return $\mathcal{D}_{1:T}$

▷ Observation history



$$\max_{\mathbf{x} \in \Lambda \subset \mathbb{R}^N} f(\mathbf{x})$$

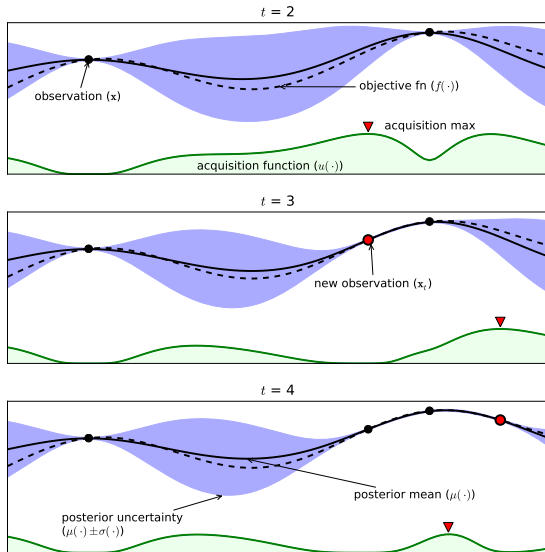
- $f: \Lambda \rightarrow \mathbb{R}$ — black-box true objective is costly to evaluate

Algorithm 2.3 SMBO(f, M_0, T, S)

$\mathcal{D}_0 \leftarrow \emptyset$ ▷ Observations
for $t = 1, \dots, T$ **do**
 $\mathbf{x}^* \leftarrow \arg \min_{\mathbf{x}} S(\mathbf{x}, M_{t-1})$
 Evaluate $f(\mathbf{x}^*)$ ▷ Expensive step
 $\mathcal{D}_{1:t} \leftarrow \mathcal{D}_{1:t-1} \cup \{(\mathbf{x}^*, f(\mathbf{x}^*))\}$
 Fit a new model M_t to $\mathcal{D}_{1:t}$
return $\mathcal{D}_{1:T}$ ▷ Observation history

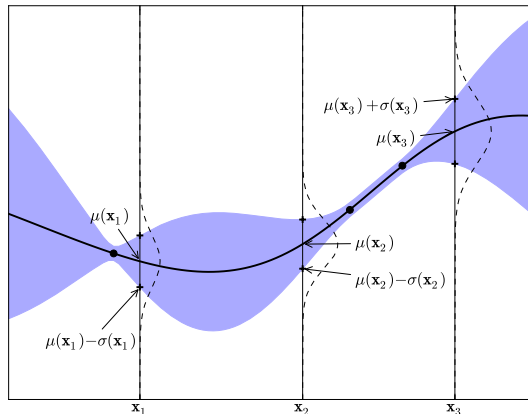
Key ingredients: a probabilistic *surrogate model* (M) and an *acquisition function* (S) to decide which point to evaluate next

- ▶ $P(f | \mathcal{D}_{1:t}) \propto P(\mathcal{D}_{1:t} | f)P(f)$
- ▶ *surrogate model* — Gaussian process (GP)



Surrogate Models

Gaussian process [10]



$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

$m(\mathbf{x})$ — mean function

$k(\mathbf{x}, \mathbf{x}')$ — covariance function

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{1}{2} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$$

Рис.: Simple 1D Gaussian process with three observations

► What we want

$$P(f \mid \mathcal{D}_{1:t}) \propto P(\mathcal{D}_{1:t} \mid f)P(f)$$

- What we want

$$P(f \mid \mathcal{D}_{1:t}) \propto P(\mathcal{D}_{1:t} \mid f)P(f)$$

- For GP holds the following (assume that $m(\mathbf{x}) = 0$)

$$\begin{bmatrix} \mathbf{f}_{1:t} \\ f_{t+1} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{k} \\ \mathbf{k}^T & k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) \end{bmatrix}\right),$$

where $\mathbf{k} = [k(\mathbf{x}_{t+1}, \mathbf{x}_1), k(\mathbf{x}_{t+1}, \mathbf{x}_2), \dots, k(\mathbf{x}_{t+1}, \mathbf{x}_t)]$ and $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))$ - kernel matrix.

- What we want

$$P(f \mid \mathcal{D}_{1:t}) \propto P(\mathcal{D}_{1:t} \mid f)P(f)$$

- For GP holds the following (assume that $m(\mathbf{x}) = 0$)

$$\begin{bmatrix} \mathbf{f}_{1:t} \\ f_{t+1} \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{k} \\ \mathbf{k}^T & k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) \end{bmatrix} \right),$$

where $\mathbf{k} = [k(\mathbf{x}_{t+1}, \mathbf{x}_1), k(\mathbf{x}_{t+1}, \mathbf{x}_2), \dots, k(\mathbf{x}_{t+1}, \mathbf{x}_t)]$ and $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))$ - kernel matrix.

GP surrogate

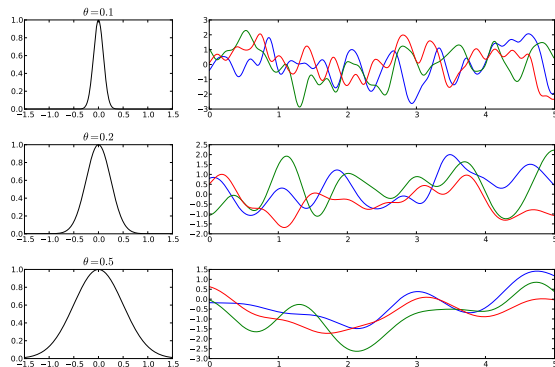
$$P(f_{t+1} \mid \mathcal{D}_{1:t}, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1})),$$

where

$$\mu_t(\mathbf{x}_{t+1}) = \mathbf{k}^T \mathbf{K}^{-1} \mathbf{f}_{1:t}$$

$$\sigma_t^2(\mathbf{x}_{t+1}) = k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T \mathbf{K}^{-1} \mathbf{k}$$

Hyperparameters of GP



$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{1}{2\theta^2} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right),$$

where θ controls the width of the kernel

Q: How to find optimal values of θ and other possible hyperparameters?

Q: How to find optimal values of θ and other possible hyperparameters?

A: MLE

Q: How to find optimal values of θ and other possible hyperparameters?

A: MLE

- Observations

$$\mathbf{f}_{1:t} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma),$$

where $\boldsymbol{\mu} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_t))$

- Log Likelihood

$$L = \log p(\mathbf{f}_{1:t} | \mathbf{x}_{1:t}, \theta) = -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{f}_{1:t} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{f}_{1:t} - \boldsymbol{\mu}) - \frac{t}{2} \log(2\pi)$$

- Find $\frac{\partial L}{\partial \theta_m}, \frac{\partial L}{\partial \theta_k}$, where θ_m, θ_k — hyperparameters of the m and k functions

Acquisition Function

The role of the acquisition function is to guide the search for the optimum

$$\mathbf{x}_{t+1} = \arg \max_{\mathbf{x}} u(\mathbf{x} \mid \mathcal{D}_{1:t}),$$

where u — acquisition function.

High acquisition corresponds to *potentially* high values of the objective function

Acquisition Function

The role of the acquisition function is to guide the search for the optimum

$$\mathbf{x}_{t+1} = \arg \max_{\mathbf{x}} u(\mathbf{x} \mid \mathcal{D}_{1:t}),$$

where u — acquisition function.

High acquisition corresponds to *potentially* high values of the objective function

Examples

► Probability of Improvement

$$\begin{aligned} \text{PI}(\mathbf{x}) &= P(f(\mathbf{x}) \geq f(\mathbf{x}^+)) \\ &= \Phi \left(\frac{\mu(\mathbf{x}) - f(\mathbf{x}^+)}{\sigma(\mathbf{x})} \right), \end{aligned}$$

where $\mathbf{x}^+ = \arg \max_{\mathbf{x}_i \in \mathbf{x}_{1:t}} f(\mathbf{x}_i)$

Unlike the original unknown objective function, $u(\cdot)$ can be cheaply sampled.

Instead of observing $f(\mathbf{x})$, we can often only observe a noisy transformation of $f(\mathbf{x})$.

Instead of observing $f(\mathbf{x})$, we can often only observe a noisy transformation of $f(\mathbf{x})$.

- ▶ Assume that $f(\mathbf{x})$ is corrupted with Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma_{noise}^2)$
- ▶ $y_i = f(\mathbf{x}_i) + \varepsilon_i$
- ▶ Since mean of noise is zero $\mathbf{K} = \mathbf{K} + \sigma_{noise}^2 \mathbf{I}$

Instead of observing $f(\mathbf{x})$, we can often only observe a noisy transformation of $f(\mathbf{x})$.

- ▶ Assume that $f(\mathbf{x})$ is corrupted with Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma_{noise}^2)$
- ▶ $y_i = f(\mathbf{x}_i) + \varepsilon_i$
- ▶ Since mean of noise is zero $\mathbf{K} = \mathbf{K} + \sigma_{noise}^2 \mathbf{I}$

GP surrogate with noise

$$P(f_{t+1} | \mathcal{D}_{1:t}, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1}) + \sigma_{noise}^2),$$

where

$$\mu_t(\mathbf{x}_{t+1}) = \mathbf{k}^T [\mathbf{K} + \sigma_{noise}^2 \mathbf{I}]^{-1} \mathbf{f}_{1:t}$$

$$\sigma_t^2(\mathbf{x}_{t+1}) = k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T [\mathbf{K} + \sigma_{noise}^2 \mathbf{I}]^{-1} \mathbf{k}$$

Instead of observing $f(\mathbf{x})$, we can often only observe a noisy transformation of $f(\mathbf{x})$.

- ▶ Assume that $f(\mathbf{x})$ is corrupted with Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma_{\text{noise}}^2)$
- ▶ $y_i = f(\mathbf{x}_i) + \varepsilon_i$
- ▶ Since mean of noise is zero $\mathbf{K} = \mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}$

GP surrogate with noise

$$P(f_{t+1} | \mathcal{D}_{1:t}, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1}) + \sigma_{\text{noise}}^2),$$

where

$$\mu_t(\mathbf{x}_{t+1}) = \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{f}_{1:t}$$

$$\sigma_t^2(\mathbf{x}_{t+1}) = k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{k}$$

Remark

Parameter σ_{noise}^2 of noise distribution could be learned while training GP

Table of Contents

1. Random Search

2. Bayesian Optimization

3. Multi-fidelity Optimization

- Learning Curve-Based Methods

- Hyperband

- Multi-task Bayesian Optimization

4. HPO Systems

5. Resume

A common technique to speed up manual tuning

- ▶ probe an hyperparameter configuration on a small subset of \mathcal{D}_{train}
- ▶ training it only for a few iterations
- ▶ running it on a subset of features

A common technique to speed up manual tuning

- ▶ probe an hyperparameter configuration on a small subset of \mathcal{D}_{train}
- ▶ training it only for a few iterations
- ▶ running it on a subset of features

Multi-fidelity Optimization

- ▶ low-fidelity (low cost) approximations of the actual loss function $V(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$
- ▶ tradeoff between optimization performance and runtime

Remark

MFO used to speedup discussed HPO algorithms

Predictive Termination [6]

Types of Learning Curves

1. Performance of an iterative algorithm as a function of its number of iterations
2. Performance of an algorithm as a function of the size of \mathcal{D}_{train}

Predictive Termination [6]

Approach

- ▶ When running SGD on DNNs we measure validation performance ($-\mathbf{V}(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$) in regular intervals
- ▶ $\mathbf{f}_{1:t}$ — observed performance values for the first t intervals

Predictive Termination [6]

Approach

- ▶ When running SGD on DNNs we measure validation performance ($-\mathbf{V}(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$) in regular intervals
- ▶ $\mathbf{f}_{1:t}$ — observed performance values for the first t intervals
- ▶ We observe $\mathbf{f}_{1:t}$ and aim to predict performance f_m after a large number of intervals $m \gg t$

Predictive Termination [6]

Approach

- ▶ When running SGD on DNNs we measure validation performance ($-\mathbf{V}(\mathcal{L}, \mathcal{A}_\lambda, \mathcal{D}_{train}, \mathcal{D}_{valid})$) in regular intervals
- ▶ $\mathbf{f}_{1:t}$ — observed performance values for the first t intervals
- ▶ We observe $\mathbf{f}_{1:t}$ and aim to predict performance f_m after a large number of intervals $m \gg t$

Predictive Termination Criteria

Probability that the network, after training for m intervals, will exceed the performance y^+ (current best).

$$P(f_m \geq f^+ \mid \mathbf{f}_{1:t}) \geq \delta$$

If this probability is above a threshold δ then training continues

Predictive Termination [6]

Experimental results

Predictive termination speeds up current hyperparameter optimization methods for DNNs by roughly a factor of two

Algorithm

1. For a given initial budget (amount of data), query all algorithms for that budget
2. Remove the half that performed worst
3. Double the budget
4. Repeat until only a single algorithm is left

Hyperband [8]

Successive halving

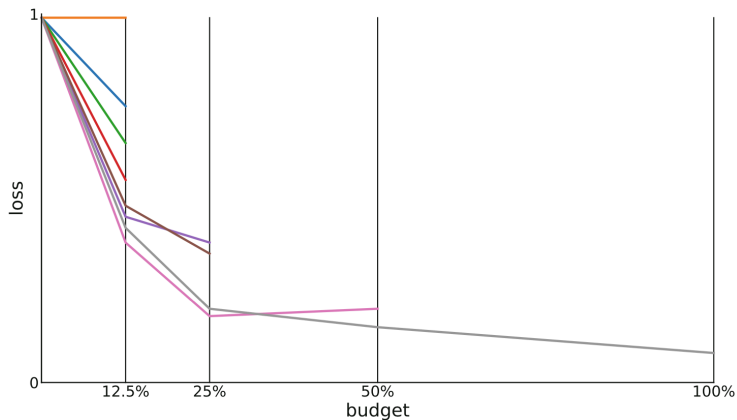


Рис.: Illustration of successive halving for eight algorithms/configurations. After evaluating all algorithms on 1/8 of the total budget, half of them are dropped and the budget given to the remaining algorithms is doubled

Hyperband [8]

Successive halving. Problem

User has to decide beforehand:

1. to try many configurations and only assign a small budget to each, or
2. to try only a few and assign them a larger budget.

Table of Contents

1. Random Search
2. Bayesian Optimization
3. Multi-fidelity Optimization
4. HPO Systems
5. Resume

- ▶ hyperopt¹ [3]
- ▶ Facebook/Ax²
- ▶ fmfnn/BayesianOptimization³
- ▶ Optuna⁴ [1]
- ▶ hyperband⁵ [8]
- ▶ tune⁶ [9]
- ▶ ...

¹<https://github.com/hyperopt/hyperopt>

²<https://github.com/facebook/Ax>

³<https://github.com/fmfnn/BayesianOptimization>

⁴<https://github.com/optuna/optuna>

⁵<https://github.com/zygmuntz/hyperband>

⁶<https://github.com/ray-project/ray>

Table of Contents

1. Random Search
2. Bayesian Optimization
3. Multi-fidelity Optimization
4. HPO Systems
5. Resume

- ▶ Why Random Search is better than Grid Search (low effective dimensionality)
- ▶ How to compare different HPO algorithms
- ▶ SOTA Bayesian Optimization (surrogate, acquisition)
- ▶ How to speedup Random Search and BO
 - ▶ Predictive Termination
 - ▶ Successive halving
 - ▶ Multi-task Bayesian Optimization

- [1] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2623–2631, 2019.
- [2] J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *Journal of machine learning research*, 13(Feb):281–305, 2012.
- [3] J. Bergstra, D. Yamins, and D. D. Cox. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. 2013.
- [4] J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pages 2546–2554, 2011.
- [5] E. Brochu, V. M. Cora, and N. De Freitas. A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *arXiv preprint arXiv:1012.2599*, 2010.

- [6] T. Domhan, J. T. Springenberg, and F. Hutter. Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.
- [7] F. Hutter, L. Kotthoff, and J. Vanschoren. *Automated Machine Learning*. Springer, 2019.
- [8] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *The Journal of Machine Learning Research*, 18(1):6765–6816, 2017.
- [9] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, and I. Stoica. Tune: A research platform for distributed model selection and training. *arXiv preprint arXiv:1807.05118*, 2018.
- [10] C. E. Rasmussen. Gaussian processes in machine learning. In *Summer School on Machine Learning*, pages 63–71. Springer, 2003.
- [11] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015.

- [12] K. Swersky, J. Snoek, and R. P. Adams. Multi-task bayesian optimization. In *Advances in neural information processing systems*, pages 2004–2012, 2013.