## Introduction to Machine Learning

# Hyperparameter Tuning Problem Definition





#### Learning goals

- Definition of HPO objective and components
- Understand its properties
- What makes tuning challenging

### HYPERPARAMETER OPTIMIZATION

Hyperparameters (HP)  $\lambda$  are parameters that are *inputs* to learner  $\mathcal{I}$  which performs ERM on training data set to find optimal **model parameters**  $\theta$ . HPs can influence the generalization performance in a non-trivial and subtle way.

Hyperparameter optimization (HPO) / Tuning is the process of finding a well-performing hyperparameter configuration (HPC)  $\lambda\in\tilde{\Lambda}$  for an learner  $\mathcal{I}_{\lambda}.$ 

× × ×

### **OBJECTIVE AND SEARCH SPACE**

Search space  $\tilde{\Lambda} \subset \Lambda$  with all optimized HPs and ranges:

 $ilde{\mathbf{\Lambda}} = ilde{\mathbf{\Lambda}}_1 imes ilde{\mathbf{\Lambda}}_2 imes \cdots imes ilde{\mathbf{\Lambda}}_l$ 

where  $\tilde{\Lambda}_i$  is a bounded subset of the domain of the i-th HP  $\Lambda_i$ , and can be either continuous, discrete, or categorical.

The general HPO problem is defined as:

$$\boldsymbol{\lambda}^* \in \argmin_{\boldsymbol{\lambda} \in \tilde{\boldsymbol{\Lambda}}} \boldsymbol{c}(\boldsymbol{\lambda}) = \argmin_{\boldsymbol{\lambda} \in \tilde{\boldsymbol{\Lambda}}} \widehat{\operatorname{GE}}(\mathcal{I}, \mathcal{J}, \rho, \boldsymbol{\lambda})$$

with  $\lambda^*$  as theoretical optimum, and  $c(\lambda)$  is short for estim. gen. error when  $\mathcal{I}$ , resampling splits  $\mathcal{J}$ , performance measure  $\rho$  are fixed.



#### **OBJECTIVE AND SEARCH SPACE**

$$\boldsymbol{\lambda}^* \in \argmin_{\boldsymbol{\lambda} \in \tilde{\boldsymbol{\Lambda}}} \boldsymbol{c}(\boldsymbol{\lambda}) = \argmin_{\boldsymbol{\lambda} \in \tilde{\boldsymbol{\Lambda}}} \widehat{\operatorname{GE}}(\mathcal{I}, \mathcal{J}, \rho, \boldsymbol{\lambda})$$



× × ×

• Evals are stored in archive

$$\mathcal{A} = ((\boldsymbol{\lambda}^{(1)}, \boldsymbol{c}(\boldsymbol{\lambda}^{(1)})), (\boldsymbol{\lambda}^{(2)}, \boldsymbol{c}(\boldsymbol{\lambda}^{(2)})), \dots),$$
 with  $\mathcal{A}^{[t+1]} = \mathcal{A}^{[t]} \cup (\boldsymbol{\lambda}^+, \boldsymbol{c}(\boldsymbol{\lambda}^+)).$ 

• We can define tuner as function  $\tau : (\mathcal{D}, \mathcal{I}, \tilde{\Lambda}, \mathcal{J}, \rho) \mapsto \hat{\lambda}$ 

### WHY IS TUNING SO HARD?

- Tuning is usually **black box**: No derivatives of the objective are availabe. We can only eval the performance for a given HPC via a computer program (CV of learner on data).
- Every evaluation can require multiple train and predict steps, hence it's **expensive**.
- Even worse: the answer we get from that evaluation is **not exact**, **but stochastic** in most settings, as we use resampling.
- Categorical and dependent hyperparameters aggravate our difficulties: the space of hyperparameters we optimize over can have non-metric, complicated structure.
- Many standard optimization algorithms cannot handle these properties.

× × ×