

Hyperparameter Tuning

Tuning: Finding the best Hyperparameters, also called HPO "Hyperpar. Optimization"

Many ML Learning Algorithms are quite sensitive to (good) setting of their hyperparams. These control flexibility/complexity of model, certain structural properties and computational aspects of training process. HPs can influence generalization performance in a subtle way.

Hyperparameters λ are not optimized during training! They must be specified in advance in order to give us a fixed Learner configuration which we then optimize w.r.t. model params θ on full-training set and evaluate on testset(s). [HPs are input of Learner t] Later we will embedd HP-tuning into a more complex training procedure, effectively removing HPs as input of Learner Algorithm. \Rightarrow self-tuning learner $t_{\lambda}: D \rightarrow y$

Search Space $\tilde{\Lambda} \subset \Delta$ with all optimizable HPs and ranges is a bounded subset of the hyperparameter space

certain λ s can only be optimized if we restrict their possible values

$\tilde{\Delta} = \tilde{\Delta}_1 \times \tilde{\Delta}_2 \times \dots \times \tilde{\Delta}_L$ with total L Hyperparameters to tune Search methods to find λ : Grid, Random Search, Bayesian Optimization, Hyperbands How choose limited amount of candidates?

Tuner T is a mapping which takes dataset D , inducer t , HP search space $\tilde{\Lambda}$, performance measure p_e and some Resampling strategy $J \Rightarrow$ returns Hyperparameter Configuration $\lambda \in \tilde{\Lambda}$

$$T(D, t, \lambda, p_e, J) = \lambda \approx \lambda^* \quad \text{theoretical "findable" optimum } \lambda^* \in \arg\min_{\lambda \in \tilde{\Lambda}} \hat{GE}(t, J, p_e, \lambda)$$

Evals are stored in archive $A = ((\lambda^{(1)}, c(\lambda^{(1)})), (\lambda^{(2)}, c(\lambda^{(2)})), \dots)$ where index refers to current iteration of HPO algorithm

- Tuning / HPO algorithm is usually a black box, partly that's because many standard optimization techniques (such as GD) don't work here!
No derivatives exist. We can only evaluate performance for given HPC via known resampling procedure
- Categorical HPs (split criterion for Class. Trees, distance measure for kNN) and dependency hierarchy of certain HPs (if we use KDE for Naive Bayes \rightarrow bandwidth?) means that we may have a very complicated mixed structure in our Hyperparam Space have to make choice after limited amount of evals
- Evaluation $c(\lambda)$ is not exact but stochastic due to resampling and each Eval requires multiple train & predict steps making it computationally expensive!

Tuning implies a two-level optimization process. First: $\arg\min \hat{GE}(t, J, p_e, \lambda)$ in regards to λ HPO

Second: For each λ do ERM on Inducer Level to obtain optimal model param θ (nested within first level)

↑ also includes eval. on test data, pessimistic bias due to $J_{train} \ll n$ which resampling method (such as k-fold CV) mitigates

How to combine both? HPO comes first, but if we used test data for 2nd lev. there is no unseen data left to evaluate final selected Learner config!

Specifically] The HP-optimal Learner config was found by choosing minimal \hat{GE} , which is calc. by evaluating different ERM-Models on different train-test splits and averaging over Error Measure for given λ . Chosen λ is the HP config that has min. \hat{GE} among all evaluated configs comparable & reproducible

But even though we use Resampling for every inspected λ -config: Overall by repeatedly doing this many times for different λ , the average or cleaner the aggregated Test Information "leaks into evaluation" (we choose λ with min of all those averaged GE-estimates), esp. if we use exact same splits $\forall \lambda$!

⇒ OVERTUNING This GE of final selected Learner config is optimistically biased $E[\min_i \hat{GE}(i)] \downarrow$ via Overfitting \leftarrow min. $E[\hat{GE}(\lambda)]$ actually want this!

In analogy to multiple testing we could also call this a multiple eval. problem. By searching & testing over more and more λ it becomes increasingly likely that one of these HP configs performs (too) good by random Luck GE-estims can be viewed as a random variable we sample from; imagine $\text{Bin}(n_{train}, \pi = MCE)$ for class. task

↓ for Performance Evaluation

Untouched Test Set / Unseen Test Data Principle needs to be adhered to consistently in the whole ML Pipeline!

For Model building steps that require resampling themselves this implies nested resampling (HP Tuning, Feature Selection etc.) or extra holdout testset

(initial) HP train | HP test | Test

Nested Resampling

Outer Resampling: Split all data into train | test \Rightarrow Retrain model with λ on all train data, evaluate Learner (config) on test data

for Tuning strictly separates model selection and performance evaluation

Inner Resampling: Split outer train data again into train | test \Rightarrow HP tuning as above, yielding $\lambda = \arg\min_{\lambda} c(\lambda)$

Now doing this k-times for k-fold CV or subsampling we will get k unbiased GE to average over

↳ aggregate $\hat{GE}_{\text{nested}} = \frac{1}{k} \sum_{i=1}^k \hat{GE}_{\text{outer}}^{(i)}$ estimates Performance of Learner / Learn Algorithm including HP-finding

can do same for other Learner types on identical splits & p-e measure ... \Rightarrow choose model / Learner type with Min. \hat{GE}_{nested} & train on full data

We have number for GE of Tuned Lear. Algo Process to report and compare (fair!)

Ind. Tuning with same strategy as Nested Resamplings inner loop on n! before final training our last resampling for HPs! (CV)
! Don't nest, new splits

In reality its actually about performance of the entire Learning Pipeline, not just Learner - HP space - Tuning Search Strategy

That entails: Pre-processing → Scaling / Standardizing, Transformations, Error Fix, Outliers, Value Imputation for N/A's etc.

BRUNNEN Feature-Engineering / Selection → Filter, Feature Encoding, Wrapper (RFE), Embedded, Combining Features, Interactions in Model, Dimension Reduction

ERM loss & Performance Measure But dont worry: Most of this is chosen upfront based on exp. and common knowledge [we can't break up Nested Res. inner loop too much]

Fix design choices → Thresholding, Early stopping in Trees?, Regularization implemented vs. overfitting etc.