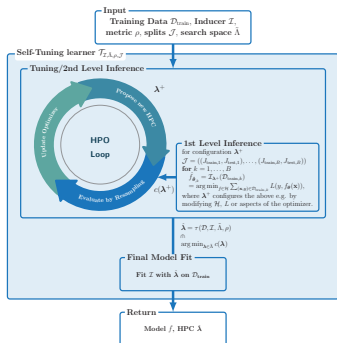


Introduction to Machine Learning

Hyperparameter Tuning - Practical Aspects

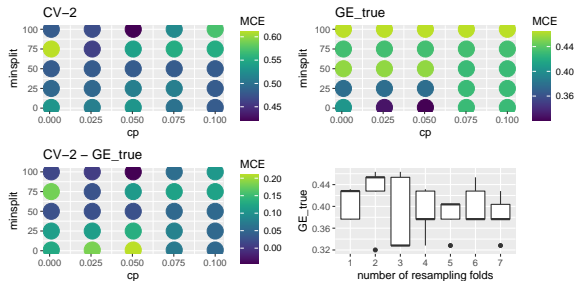


Learning goals

- Understand the possible design choices for HPO
- Know termination criteria of HPO

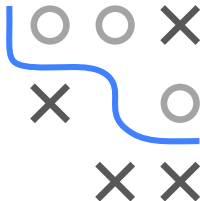
PRACTICAL ASPECTS OF HPO

- Choosing resampling
 - Nr of observations, i.i.d assumption for data sampling process
 - Higher resampling rates likely result in a better model; however they are computationally more expensive



PRACTICAL ASPECTS OF HPO_2

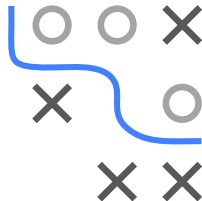
Tuning a CART on the `spirals` data with a k-fold CV (k=1 means here a 2/3 holdout split) using grid search and estimating the true GE with a very large test set (5 repetitions)



-
- The figure consists of two plots. The left plot is a heatmap showing the relationship between 'cd' (x-axis, 0.00 to 0.05) and 'misspelt' (y-axis, 0 to 100). The background is a heatmap where color represents 'accuracy', ranging from dark blue (0.76) to light blue (0.80). Points are plotted on this heatmap, colored by 'type' (red for $\hat{\lambda}^+$, grey for $\hat{\lambda}^-$) and shaped by 'accuracy' (circle for 0.80, triangle for 0.78, square for 0.76). The right plot is a line graph showing 'accuracy' (y-axis, 0.78 to 0.81) versus 'area' (x-axis, 1 to 5). It compares two types: $\hat{\lambda}^+$ (solid line with circles) and $\hat{\lambda}^-$ (dashed line with circles). Both types show an increase in accuracy as area increases, with $\hat{\lambda}^-$ consistently performing better.
- | area | accuracy ($\hat{\lambda}^+$) | accuracy ($\hat{\lambda}^-$) |
|------|--------------------------------|--------------------------------|
| 1 | 0.780 | 0.780 |
| 2 | 0.783 | 0.796 |
| 5 | 0.800 | 0.815 |

PRACTICAL ASPECTS OF HPO / 2

Tuning `cp` and `minsplit` for a CART on the `titanic` data over 3 increasing rectangular search spaces with random search (candidates number fixed) and comparing the result with the optimal model (found with exhaustive grid search)



PRACTICAL ASPECTS OF HPO

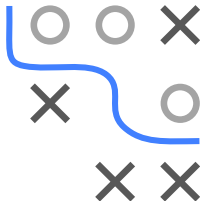
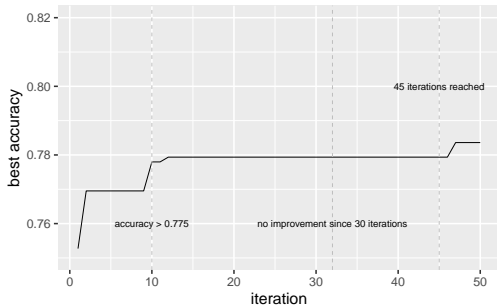
- Choosing HPO algorithm
 - For few HPS (1-3), grid search can be used
 - BO with GPs for upto 10 numeric HPs
 - BO with RFs handle mixed HP spaces
 - Random search and Hyperband work well as long as the “effective” dimension is low
 - EAs are somewhat in-between BO and RS, can handle very complex spaces, but less sample efficient than BO
 - **Also: use something that’s stable and robust! More an aspect of the implementation than the algo!**



PRACTICAL ASPECTS OF HPO

When to terminate HPO

- Specify a certain amount of runtime/budget beforehand
- Set a lower bound regarding \widehat{GE}
- Terminate if performance improvement stagnates



Different stopping points while tuning CART on the `titanic` data depending on which termination criterion is used

PRACTICAL ASPECTS OF HPO

- Warm starts
 - Evaluations (e.g., weight sharing of neural networks)
 - Optimization (initializing with HPCs that worked well before)
- Control of execution
 - Parallelizability of HPO algorithms differs strongly
 - HPO execution can be parallelized at different levels (outer resampling, iteration, evaluation, inner resampling, model fit)

