

# AutoML: Hyperparameter Optimization

## Practical Problems

Bernd Bischl   Frank Hutter   Lars Kotthoff  
Marius Lindauer

# Choice of Learning Algorithms

- A plethora of learners exists, for different data sets different models are likely needed.
- Studies and experience show:  
One these is often good – on tabular data:
  - ▶ Penalized regression, e.g. elastic net
  - ▶ Support vector machines
  - ▶ Gradient boosting
  - ▶ Random forests
  - ▶ Neural networks
- Example: Auto-Sklearn 2.0 [Feurer et al. 2020] uses:
  - ▶ Extra trees
  - ▶ Gradient boosting
  - ▶ Passive aggressive
  - ▶ Random forest
  - ▶ Linear regression with SGD
  - ▶ Multi-layer perceptron

# Choice of Search Space for a Learning Algorithm

Ranges often selected based on experience

- See other AutoML frameworks: e.g. Auto-Sklearn 2.0 [Feurer et al. 2020]
- Sensitivity analysis often does not exist for ML algorithms
- Check literature on specific ML algorithm

Algorithm	Hyperparameter	Type	Lower	Upper	Trafo
<b>glmnet</b>					
(Elastic net)	alpha	numeric	0	1	-
	lambda	numeric	-10	10	$2^x$
<b>rpart</b>					
(Decision tree)	cp	numeric	0	1	-
	maxdepth	integer	1	30	-
	minbucket	integer	1	60	-
	minsplit	integer	1	60	-
<b>kkm</b>					
(k-nearest neighbor)	k	integer	1	30	-
<b>svm</b>					
(Support vector machine)	kernel	discrete	-	-	-
	cost	numeric	-10	10	$2^x$
	gamma	numeric	-10	10	$2^x$
	degree	integer	2	5	-
<b>ranger</b>					
(Random forest)	num.trees	integer	1	2000	-
	replace	logical	-	-	-
	sample.fraction	numeric	0.1	1	-
	mtry	numeric	0	1	$x \cdot p$
	respect.unordered.factors	logical	-	-	-
	min.node.size	numeric	0	1	$n^x$
<b>xgboost</b>					
(Gradient boosting)	nrounds	integer	1	5000	-
	eta	numeric	-10	0	$2^x$
	subsample	numeric	0.1	1	-
	booster	discrete	-	-	-
	max_depth	integer	1	15	-
	min_child_weight	numeric	0	7	$2^x$
	colsample_bytree	numeric	0	1	-
	colsample_bylevel	numeric	0	1	-
	lambda	numeric	-10	10	$2^x$
	alpha	numeric	-10	10	$2^x$

Source: [Probst et al. 2019 ].

# Choice of Search Space for a Learning Algorithm

Parameter	Def.P	Def.O	Tun.P	Tun.O	q <sub>0.05</sub>	q <sub>0.95</sub>
glmnet			0.069	0.024		
alpha	1	0.403	0.038	0.006	0.009	0.981
lambda	0	0.004	0.034	0.021	0.001	0.147
rpart			0.038	0.012		
cp	0.01	0	0.025	0.002	0	0.008
maxdepth	30	21	0.004	0.002	12.1	27
minbucket	7	12	0.005	0.006	3.85	41.6
minsplit	20	24	0.004	0.004	5	49.15
kknn			0.031	0.006		
k	7	30	0.031	0.006	9.95	30
svm			0.056	0.042		
kernel	radial	radial	0.030	0.024		
cost	1	682.478	0.016	0.006	0.002	920.582
gamma	1/p	0.005	0.030	0.022	0.003	18.195
degree	3	3	0.008	0.014	2	4
ranger			0.010	0.006		
num.trees	500	983	0.001	0.001	206.35	1740.15
replace	TRUE	FALSE	0.002	0.001		
sample.fraction	1	0.703	0.004	0.002	0.323	0.974
mtry	$\sqrt{p}$	0.257	0.006	0.003	0.035	0.692
respect.unordered.factors	TRUE	FALSE	0.000	0.000		
min.node.size	1	1	0.001	0.001	0.007	0.513
xgboost			0.043	0.014		
nrounds	500	4168	0.004	0.002	920.7	4550.95
eta	0.3	0.018	0.006	0.005	0.002	0.355
subsample	1	0.839	0.004	0.002	0.545	0.958
boost	gbtree	gbtree	0.015	0.008		
max_depth	6	13	0.001	0.001	5.6	14
min_child_weight	1	2.06	0.008	0.002	1.295	6.984
colsample_bytree	1	0.752	0.006	0.001	0.419	0.864
colsample_bylevel	1	0.585	0.008	0.001	0.335	0.886
lambda	1	0.982	0.003	0.002	0.008	29.755
alpha	1	1.113	0.003	0.002	0.002	6.105

Table 3: Defaults (package defaults (Def.P) and optimal defaults (Def.O)), tunability of the hyperparameters with the package defaults (Tun.P) and our optimal defaults (Tun.O) as reference and tuning space quantiles (q<sub>0.05</sub> and q<sub>0.95</sub>) for different parameters of the algorithms.

Ranges often selected based on experience

- See other AutoML frameworks: e.g. Auto-Sklearn 2.0 [Feurer et al. 2020]
- Sensitivity analysis often does not exist for ML algorithms
- Check literature on specific ML algorithm

Options for automation:

- 1 Use huge search space to cover all possibilities (combine with meta-learning for good initial design for Bayesian optimization)
  - Use results of meta-experiments to obtain smaller search space that is estimated to work well.

# Choice of Search Space for a Learning Algorithm

Parameter	Def.P	Def.O	Tun.P	Tun.O	q <sub>0.05</sub>	q <sub>0.95</sub>
glmnet			0.069	0.024		
alpha	1	0.403	0.038	0.006	0.009	0.981
lambda	0	0.004	0.034	0.021	0.001	0.147
rpart			0.038	0.012		
cp	0.01	0	0.025	0.002	0	0.008
maxdepth	30	21	0.004	0.002	12.1	27
minbucket	7	12	0.005	0.006	3.85	41.6
minsplit	20	24	0.004	0.004	5	49.15
kknn			0.031	0.006		
k	7	30	0.031	0.006	9.95	30
svm			0.056	0.042		
kernel	radial	radial	0.030	0.024		
cost	1	682.478	0.016	0.006	0.002	920.582
gamma	1/p	0.005	0.030	0.022	0.003	18.195
degree	3	3	0.008	0.014	2	4
ranger			0.010	0.006		
num.trees	500	983	0.001	0.001	206.35	1740.15
replace	TRUE	FALSE	0.002	0.001		
sample.fraction	1	0.703	0.004	0.002	0.323	0.974
mtry	$\sqrt{p}$	0.257	0.006	0.003	0.035	0.692
respect.unordered.factors	TRUE	FALSE	0.000	0.000		
min.node.size	1	1	0.001	0.001	0.007	0.513
xgboost			0.043	0.014		
nrounds	500	4168	0.004	0.002	920.7	4550.95
eta	0.3	0.018	0.006	0.005	0.002	0.355
subsample	1	0.839	0.004	0.002	0.545	0.958
boost	gbtree	gbtree	0.015	0.008		
max_depth	6	13	0.001	0.001	5.6	14
min_child_weight	1	2.06	0.008	0.002	1.295	6.984
colsample_bytree	1	0.752	0.006	0.001	0.419	0.864
colsample_bylevel	1	0.585	0.008	0.001	0.335	0.886
lambda	1	0.982	0.003	0.002	0.008	29.755
alpha	1	1.113	0.003	0.002	0.002	6.105

Table 3: Defaults (package defaults (Def.P) and optimal defaults (Def.O)), tunability of the hyperparameters with the package defaults (Tun.P) and our optimal defaults (Tun.O) as reference and tuning space quantiles (q<sub>0.05</sub> and q<sub>0.95</sub>) for different parameters of the algorithms.

Ranges often selected based on experience

- See other AutoML frameworks: e.g. Auto-Sklearn 2.0 [Feurer et al. 2020]
- Sensitivity analysis often does not exist for ML algorithms
- Check literature on specific ML algorithm

Options for automation:

- 1 Use huge search space to cover all possibilities (combine with meta-learning for good initial design for Bayesian optimization)
  - Use results of meta-experiments to obtain smaller search space that is estimated to work well.
- 2 Start with a small space and increase bit by bit

# Choice of Resampling Strategy

For computation of generalization error / cost:

$$c(\boldsymbol{\lambda}) = \frac{1}{k} \sum_{i=1}^k \widehat{GE}_{\mathcal{D}_{\text{val}}^i} (\mathcal{I}(\mathcal{D}_{\text{train}}^i, \boldsymbol{\lambda}))$$

Rules of thumb:

- Default: 10-fold CV ( $k = 10$ )
- Huge datasets: holdout
- Tiny datasets: 10x10 repeated CV
- Stratification for imbalanced classes

# Choice of Resampling Strategy

For computation of generalization error / cost:

$$c(\boldsymbol{\lambda}) = \frac{1}{k} \sum_{i=1}^k \widehat{GE}_{\mathcal{D}_{\text{val}}^i} (\mathcal{I}(\mathcal{D}_{\text{train}}^i, \boldsymbol{\lambda}))$$

Rules of thumb:

- Default: 10-fold CV ( $k = 10$ )
- Huge datasets: holdout
- Tiny datasets: 10x10 repeated CV
- Stratification for imbalanced classes

Watch out for this:

- Small sample size because of imbalances
- Repeated measurements (leave-one-object out)
- Time dependencies
- A good AutoML system should let you customize resampling
- Meta-learn good resampling strategy [Feurer et al. 2020]

# Choice of Optimization Algorithm

Choose optimization algorithm based on ...

- complexity of search space / budget
- time-costs of evaluations

Complex search space

→ BO with RF surrogate, EA with exploratory character, TPE

---

<sup>1</sup>Still has its own hyperparameters [Lindauer et al. 2019]



# Choice of Optimization Algorithm

Choose optimization algorithm based on ...

- complexity of search space / budget
- time-costs of evaluations

Complex search space

→ BO with RF surrogate, EA with exploratory character, TPE

Numerical (lower-dim) search space and tight budget

→ BO with GP surrogate<sup>1</sup>

---

<sup>1</sup>Still has its own hyperparameters [Lindauer et al. 2019]

# Choice of Optimization Algorithm

Choose optimization algorithm based on ...

- complexity of search space / budget
- time-costs of evaluations

Complex search space

→ BO with RF surrogate, EA with exploratory character, TPE

Numerical (lower-dim) search space and tight budget

→ BO with GP surrogate<sup>1</sup>

Expensive evaluations

→ Hyperband, BOHB, DEHB

---

<sup>1</sup>Still has its own hyperparameters [Lindauer et al. 2019]

# Choice of Optimization Algorithm

Choose optimization algorithm based on . . .

- complexity of search space / budget
- time-costs of evaluations

Complex search space

→ BO with RF surrogate, EA with exploratory character, TPE

Numerical (lower-dim) search space and tight budget

→ BO with GP surrogate<sup>1</sup>

Expensive evaluations

→ Hyperband, BOHB, DEHB

Deep learning

→ common practice: Parameterize architectures, then HPO – better do it jointly!

→ one-shot models and gradient-based optimization

---

<sup>1</sup>Still has its own hyperparameters [Lindauer et al. 2019]

# Practical Problems: When to stop?

We need to specify a budget, e.g.

- walltime,
- function evaluations,
- performance threshold, or
- *stagnation* for a certain time.

Problems:

- Overtuning [Makarova et al. 2021]
- Missed opportunity
- Wasted computational resources

Ways out:

- Early stopping for BO [Makarova et al. 2021].
- Rules of thumb, maybe  $50 \times l$  to  $100 \times l$  (be careful and think for yourself!).
- Expert knowledge.

# Practical Problems: Stability

AutoML system should:

- Never fail to return a result.
- Terminate within a given time.
- Save intermediate results and allow to continue.

Failure points:

- Optimizer can crash.
- Pipeline training can crash.
- Training of a pipeline can run "forever".

Ways out:

- Encapsulate train/predict in separate process from HPO.
- Ressource limit time and memory of that process.
- If pipeline crashes, run robust fallback (e.g., constant predictor).
- Use random configuration if optimizer crashes.

# Practical Problems: Parallelization

Parallelization should allow:

- Multiple CPUs/GPUs on a single machine.
- Multiple machines / nodes.

Possible parallelization levels:

- Training of pipeline.
- Resampling.
- Evaluation of configurations (batch proposals or asynchronous).

Possible problems:

- Sequential nature of HPO algorithms (e.g. BO).
- Heterogeneous training times of pipelines can cause idling.
- Main memory or CPU-cache becomes bottleneck
- Communication between machine / nodes.

Way out: Use a robust framework for parallelization.

# Practical Problems: What to return?

What is the output of a an AutoML system, e.g.

- Pipeline with best validation error.
- Stacking, e.g. averaging, of top- $k$  pipelines.
- Pareto set for multi-objective optimization.
- "One-standard-error rule": Use the simplest model within one standard error of the performance of the best model [Hastie et al. 2009].

Ensure that simple but efficient pipelines have been tried out

- Baseline (Classification: Majority vote; Regression: Mean prediction).
- Linear Model.
- (untuned) Random Forest.
- ...

# Open Problems

- Most efficient HPO approach? Good benchmarks often missing.
- How to integrate human a-priori knowledge?
- Human-in-the-loop approaches for AutoML.
- How can we best (computationally) transfer “experience” into AutoML?
- Warmstarts, learned search spaces, etc.
- Multi-Objective goals, including model interpretability and fairness.
- AutoML as a process is too much of a black-box, hurts adoption.
- Incorporate Uncertainty quantification into AutoML.
- AutoML beyond supervised learning.
- ...

→ Lots of open research questions, feel free to approach us for if you are interested.