Towards Explainable AutoML: xAutoML

Marius Lindauer









* Slides available at automl.org/talks

M. Lindauer

ELLIS AutoML Seminar

Leibniz University Hannover

AutoML Process



Who is using AutoML?



Users of ML without any deep expertise in ML



ML experts / Data scientists ... (AutoML researchers)





Are ML-Experts using AutoML?



- Bouthillier and Varoguaux [2020] showed that authors of NeurIPS and ICLR papers:
 - a) they often optimize their hyperparameters (>75%) Ο
 - b) they often **do it manually** and don't use AutoML tools Ο
- Crisan and Fiore-Gartland [2021] interviewed data scientists and concluded:
 - a) experts don't necessarily trust AutoML Ο
 - b) visualization of results and processes Ο can help to increase the acceptance of AutoML results







Previous Approach-Agnostic xAutoML Methods





Note: Of course, approach-specific methods are also possible, e.g. [Ru et al. 2021]. AutoML.org

M. Lindauer

ELLIS AutoML Seminar

Leibniz University Hannover

General Setup



* can in principle be any kind of cost function.

M. Lindauer

6

0 -00-

-MA

Leibniz Universität Hannover

AutoML.org

Visualization of Pipelines



Source: [Ono et al. 2020]

ELLIS AutoML Seminar

- Visualization of sampled pipelines (incl. algorithms and hyperparameters) and their performance
- More descriptive analysis



10 -

00-

- M4

Leibniz Universität Hannover

Parallel Coordinate Plots



- Visualization of sampling in high-dimensional hyperparameter spaces [Golovin et al. 2017]
- Allows to identify:
 - optimization focus of AutoML optimiziers
 - well-performing combination of settings
 - Interaction effects between settings (to some degree)
- Rather qualitative, less quantitative analysis
- Follow up: Conditional parallel coordinate plots [Weidele et al. 2019]

AutoML.org

Ablation Studies



- AutoML can start from some default settings
 - a. Defined by the algorithm developer
 - b. User expertise
- Question: Which settings (λ) changes had the bigger impact on the performance?
- Characteristic:
 - a. Quantitative analysis
 - b. Hard to visualize in high-dimensional spaces
 - c. Rather local: subspace between default and incumbent setting
- Efficiency in high-dimensional spaces:
 - a. Greedy approach [Fawcett and Hoos 2016]
 - b. Use of surrogate models instead of expensive function evaluations (e.g., from Bayesian Optimization) [Biedenkapp et al. 2017]



0

Hannover

3

Leibniz Universität

ICE Curves and LPI



M. Lindauer

- ICE Curve [Goldstein et al. 2017]: individual effect of one feature for an individual observation
 - Slice through the space in one dimension at a given observation
- LPI: Local (Hyper-)Parameter Importance [Biedenkapp et al. 2018]
 - Same idea as ICE curves, but single ICE curve centered at the incumbent setting returned by an AutoML tool
 - Quantitative importance of hyperparameters:

$$LPI(h \mid \lambda) = \frac{Var_{v \in \Lambda_h} \hat{c}(\lambda[\lambda_h = v])}{\sum_{h' \in \mathcal{H}} Var_{w \in \Lambda_{h'}} \hat{c}(\lambda[\lambda_h = w])}$$



Leibniz Universität Hannover

fANOVA for Hyperparameter Importance

	fANOVA	LPI
discount	19.32	38.88
learning rate	3.70	35.4
batch size	15.77	21.5
# units 1	1.86	0.07
# units 2	0.39	0.01

PPO on cartpole Source: [Lindauer et al. 2019] Fraction of explained variance by main and interaction effects of hyperparameters can be quantified by

$$\mathbb{V}_{\mathcal{H}'\subset\mathcal{H}} = \frac{\frac{1}{||\Lambda_{\mathcal{H}'}||} \int \hat{c}(\lambda_{\mathcal{H}'})^2 d\lambda_{\mathcal{H}'}}{\frac{1}{||\Lambda||} \int (\hat{y}(\lambda) - \hat{c}_{\varnothing})^2 d\lambda}$$

- Efficient computation on a RF as surrogate model [Hutter et al. 2014]
- Allows to study importance across datasets [van Rijn & Hutter 2018, Sharma et al. 2019]

Leibniz Universität Hannover

Explaining HPO via Partial Dependence Plots

Julia Moosbauer¹, Julia Herbinger¹, Giuseppe Casalicchio¹, Marius Lindauer², Bernd Bischl¹









PDP: Partial Dependence Plots [Friedman 2001]

For, a subset *S* of the hyperparameters, the partial dependence function is:

$$c_S(\lambda_S) := \mathbb{E}_{\lambda_C} \left[c(\lambda) \right] = \int_{\Lambda_C} c(\lambda_S, \lambda_C) d\mathbb{P}(\lambda_C)$$

and can be approximated by Monte-Carlo integration:

$$\hat{c}_S(\lambda_S) = \frac{1}{n} \sum_{i=1}^n \hat{m}\left(\lambda_S, \lambda_C^{(i)}\right)$$

where $\left(\lambda_C^{(i)}\right)_{i=1} \sim \mathbb{P}(\lambda_C)$ and λ_S for a set of grid points.





Leibniz Universität Hannover

 \rightarrow Average of ICE curves.

[Hutter et al. 2014] showed how to do this efficiently for RFs as surrogate models.

Quantifying Uncertainties in PDPs

$$\hat{s}_{S}^{2}(\lambda_{S}) = \mathbb{V}_{\hat{c}} \left[\hat{c}_{S} \left(\lambda_{S} \right) \right]$$
$$= \mathbb{V}_{\hat{c}} \left[\frac{1}{n} \sum_{i=1}^{n} \hat{c} \left(\lambda_{S}, \lambda_{C}^{(i)} \right) \right]$$
$$= \frac{1}{n^{2}} \mathbf{1}^{\top} \hat{K} \left(\lambda_{S} \right) \mathbf{1}.$$

 \rightarrow requires a kernel correctly specifying the covariance structure (e.g., GPs).

Approximation:

$$\hat{s}_{S}^{2}(\lambda_{S}) \approx \frac{1}{n} \sum_{i=1}^{n} \hat{K}(\lambda_{S})_{i,i}$$

 \rightarrow Model-agnostic (local) approximation





Leibniz Universität Hannover

Problem: Biased Sampling

- PDPs assume that the data is i.i.d.
- Obviously not the case for efficient AutoML tools with a focus on high-performance regions



- BO with GPs and LCB 0
- Different exploration rate Ο for LCB to show different sampling bias

$$LCB(\lambda) = \mu(\lambda) + \beta \cdot \sigma(\lambda)$$







_eibniz Jniversität lannover



Impact of the Sampling Bias

- Simply using all observations from AutoML tools might lead to misleading PDPs
- Uncertainty estimates help to quantify the poor fits
- \rightarrow of course, sampling bias is wanted and the solution cannot be to change the sampling behavior



Leibniz Universität Hannover

AutoML.org

16

M. Lindauer

ELLIS AutoML Seminar

Partitioning of Space

Partition space to obtain interpretable subspaces \mathcal{N}^{\prime}

Uncertainty variation across all ICE estimates: $L(\lambda_{S}, \mathcal{N}') = \sum_{i \in \mathcal{N}} \left(\hat{s}^{2} \left(\lambda_{S}, \lambda_{C}^{(i)} \right) - \hat{s}_{S|\mathcal{N}'}^{2} \left(\lambda_{S} \right) \right)^{2}$ $\hat{s}_{S|\mathcal{N}'}^{2} \left(\lambda_{S} \right) := \frac{1}{|\mathcal{N}'|} \sum_{i \in \mathcal{N}'} \hat{s}^{2} \left(\lambda_{S}, \lambda_{C}^{(i)} \right)$

\rightarrow Uncertainty structure of ICE curves should maximally agree Split Loss = Aggregation over all grid points:

$$\mathcal{R}_{L2}(\mathcal{N}') = \sum_{g=1}^{G} L(\lambda_S^{(g)}, \mathcal{N}')$$

Note (i): Partition only along the marginalized dimensions











Empirical Results





M. Lindauer

19

ELLIS AutoML Seminar

Leibniz University Hannover

Effect of Splitting on an Artificial Function

Main Insights:

- For higher-dimensional problems, PDPs are potentially more uncertain
- Mean Confidence (MC) increases with deeper trees
- More to gain for high-sampling bias cases



AutoML.org



Explaining LCBench [Zimmer et al. 2021]

Name	Range	log	type
Number of layers	[1, 5]	no	int
Max. number of units	[64, 512]	yes	int
Batch size	[16, 512]	yes	int
Learning rate (SGD)	$[1e^{-4}, 1e^{-1}]$	yes	float
Weight decay	$[1e^{-5}, 1e^{-1}]$	no	float
Momentum	[0.1, 0.99]	no	float
Max. dropout rate	[0.0, 1.0]	no	float

Hyperparameter	δ MC (%)	$\delta \text{ OC } (\%)$	δ NLL (%)
Batch size	40.8 (14.9)	61.9 (13.5)	19.8 (19.5)
Learning rate	50.2 (13.7)	57.6 (14.4)	17.9 (20.5)
Max. dropout	49.7 (15.4)	62.4 (11.9)	17.4 (18.2)
Max. units	51.1 (15.2)	58.6 (12.7)	24.6 (22.0)
Momentum	51.7 (14.5)	58.3 (12.7)	19.7 (21.7)
Number of layers	30.6 (16.4)	50.9 (16.6)	13.8 (32.5)
Weight decay	36.3 (22.6)	61.0 (13.1)	11.9 (19.7)

Improvement of mean confidence (MC), confidence close to incumbent (OC), and negative log-likelihood (NLL) after 6 splits

Setting:

- Small configuration space of Auto-PyTorch Tabular
- Training of a RF as surrogate on LCBench with 2000 randomly sampled configurations
- Bayesian Optimization with 200 function evaluations

Take-away:

 \rightarrow the confidence of PDPs improves across all hyperparameters and metrics



AutoML.ord

Leibniz Universität Hannover

Explaining Auto-PyTorch (cont'd) [Zimmer et al. 2021]



Leibniz Universität Hannover

100

- M4





Future Work and Conclusion





ELLIS AutoML Seminar

Leibniz University Hannover

M. Lindauer

Take-Home

Summary:

- Explaining AutoML is important to create trust in them
- Common iML methods such as PDPs can be used to explain AutoML
- However, i.i.d assumptions might be violated
- PDPs can be extended to uncertainty estimates
- Split into subspaces with better interpretability

Future Work:

- Additional samples to efficiently reduce sampling bias
- Extension to multi-fidelity setting AutoML
- Other iML methods

Universität

Goal: Human-Centered AutoML



too 4 Hannover



M. Lindauer



Thank you!







M. Lindauer





AutoML.org