

AutoML: Interpretability

Overview: Automated Empirical Analysis

Bernd Bischl Frank Hutter Lars Kotthoff
Marius Lindauer Joaquin Vanschoren

- Big challenge of ML: Interpretability
 - ▶ In some applications, it is required to "understand" a prediction
 - ▶ Users have less trust in systems, they can't understand

- Big challenge of ML: Interpretability
 - ▶ In some applications, it is required to "understand" a prediction
 - ▶ Users have less trust in systems, they can't understand
- AutoML is even worse?
 - ▶ AutoML is a black-box that automates the design of another blackbox (ML)
 - ▶ Also ML-developers have an basic understanding of the design of their ML pipelines
- Automated empirical interpretability helps to
 - ▶ understand the finally returned ML system
 - ▶ understand the AutoML process

Approach

- Insights:

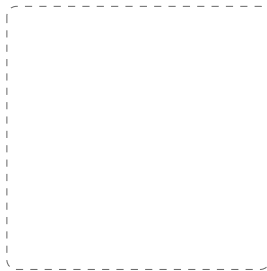
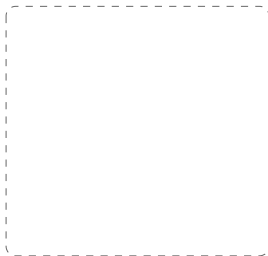
- ▶ AutoML is yet another optimization problem
- ▶ (Most) AutoML approach are iterative in nature

~> AutoML generates a lot of empirical data

Cost c

Budgets

Design Space Λ

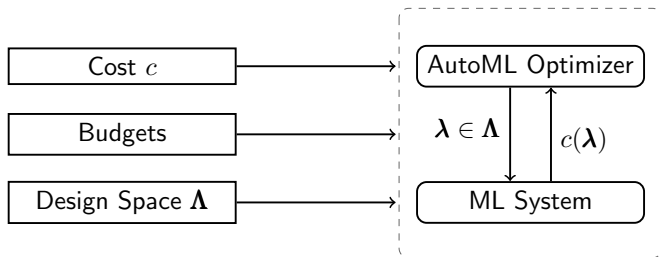


Approach

- Insights:

- ▶ AutoML is yet another optimization problem
- ▶ (Most) AutoML approach are iterative in nature

~> AutoML generates a lot of empirical data

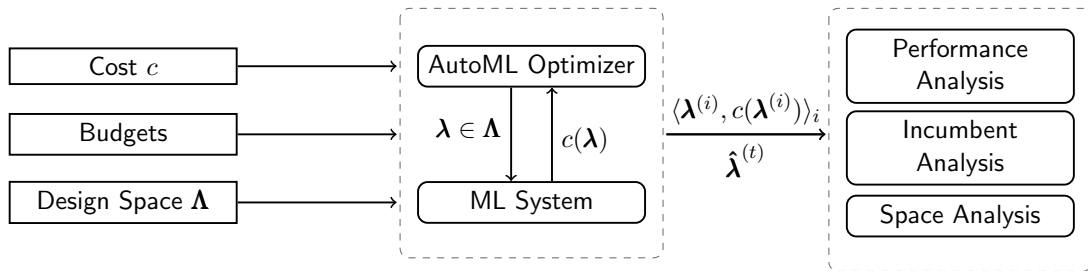


Approach

- Insights:

- ▶ AutoML is yet another optimization problem
- ▶ (Most) AutoML approach are iterative in nature

~> AutoML generates a lot of empirical data

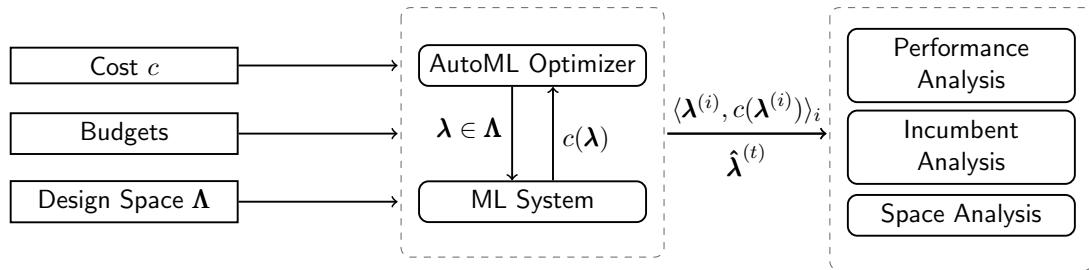


Approach

- Insights:

- ▶ AutoML is yet another optimization problem
- ▶ (Most) AutoML approach are iterative in nature

↪ AutoML generates a lot of empirical data



↪ Let's use this data to learn something about our AutoML problem

Basic Examples

- Visualize final incumbent $\hat{\lambda}$
 - ▶ ML pipeline with its components
 - ▶ Neural architecture

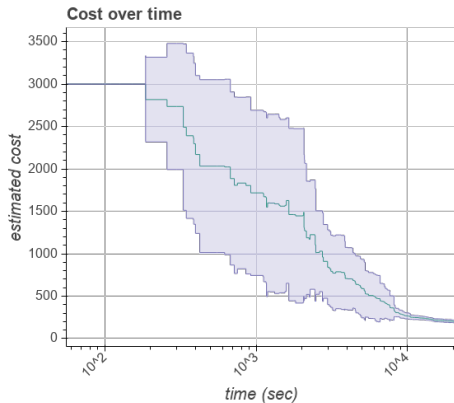
Basic Examples

- Visualize final incumbent $\hat{\lambda}$
 - ▶ ML pipeline with its components
 - ▶ Neural architecture
- Compare what changed between λ_{def} and $\hat{\lambda}$

Basic Examples

- Visualize final incumbent $\hat{\lambda}$
 - ▶ ML pipeline with its components
 - ▶ Neural architecture
- Compare what changed between λ_{def} and $\hat{\lambda}$
- Show $\hat{\lambda}$ on different budgets (if you used a multi-fidelity approach)

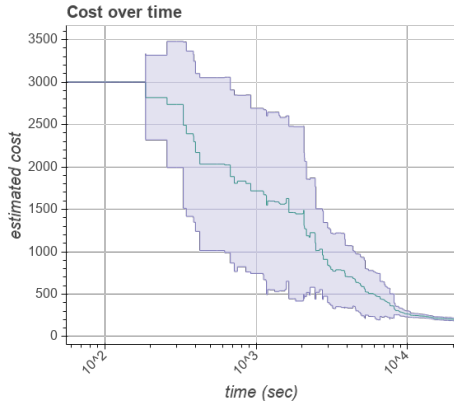
Cost Over Time



- Study how your AutoML tool improves cost (or loss) over time

Source: [Lindauer et al. 2019]

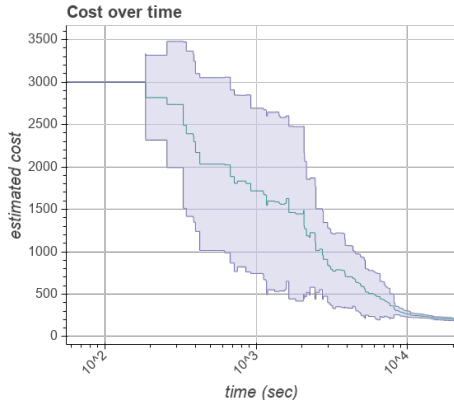
Cost Over Time



Source: [Lindauer et al. 2019]

- Study how your AutoML tool improves cost (or loss) over time
- Allows to identify whether
 - ▶ you need less time next time or

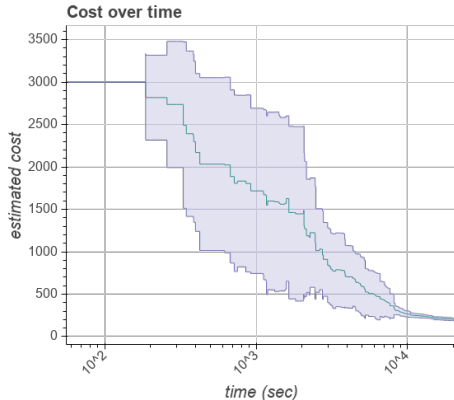
Cost Over Time



Source: [Lindauer et al. 2019]

- Study how your AutoML tool improves cost (or loss) over time
- Allows to identify whether
 - ▶ you need less time next time or
 - ▶ the AutoML system is still improving; so you should give it more time

Cost Over Time



Source: [Lindauer et al. 2019]

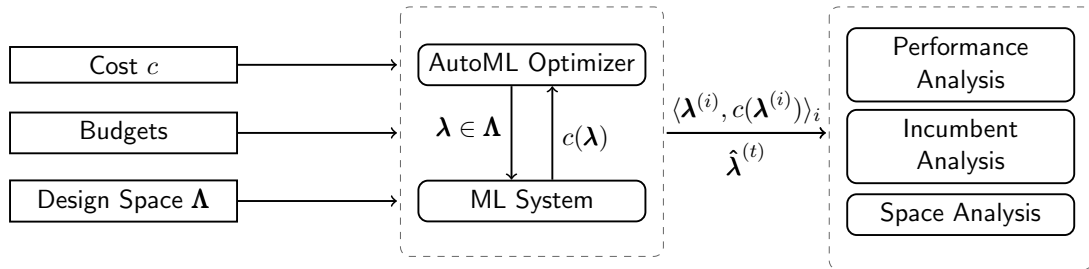
- Study how your AutoML tool improves cost (or loss) over time
- Allows to identify whether
 - ▶ you need less time next time or
 - ▶ the AutoML system is still improving; so you should give it more time
- Notes:
 - ▶ Plot on log-scale to see details in the beginning
 - ▶ If you done several runs, plot distribution (e.g., median and 25/75%-quartiles)

AutoML: Interpretability

Studying the AutoML Optimization Process

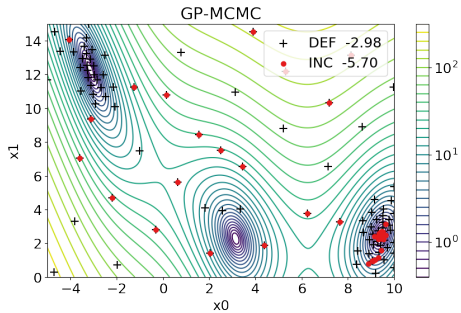
Bernd Bischl Frank Hutter Lars Kotthoff
Marius Lindauer Joaquin Vanschoren

Idea



~> focus on how the AutoML optimizer samples from the design space Λ

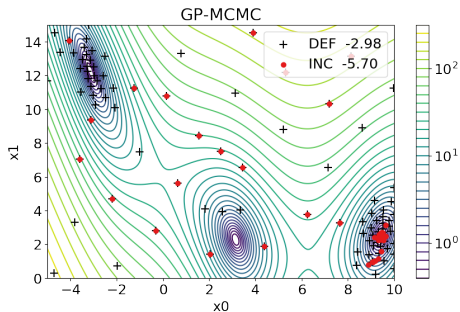
Analyzing the Sampling Behavior



- Plot of a 1D or 2D function

Source: [Lindauer et al. 2019]

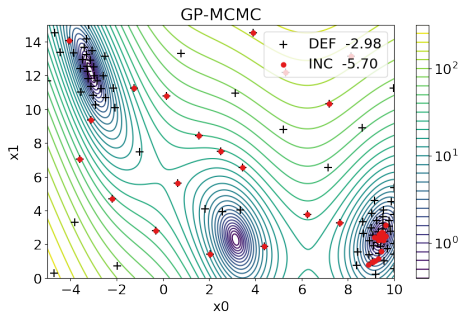
Analyzing the Sampling Behavior



- Plot of a 1D or 2D function
- Background shows the ground truth (real function values)

Source: [Lindauer et al. 2019]

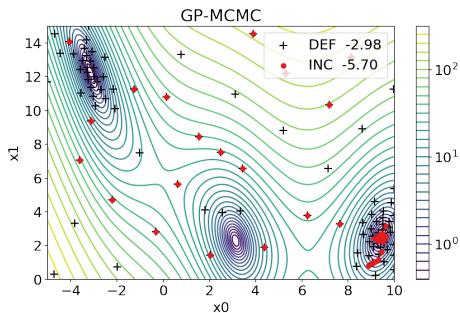
Analyzing the Sampling Behavior



Source: [Lindauer et al. 2019]

- Plot of a 1D or 2D function
- Background shows the ground truth (real function values)
- Dots are sampled points in the search space

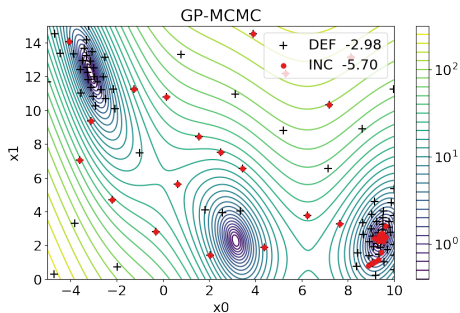
Analyzing the Sampling Behavior



Source: [Lindauer et al. 2019]

- Plot of a 1D or 2D function
- Background shows the ground truth (real function values)
- Dots are sampled points in the search space
- Typical approach in Bayesian Optimization community

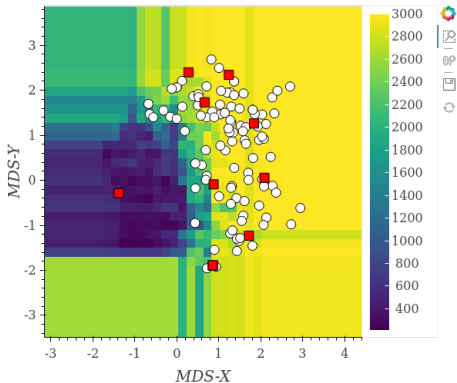
Analyzing the Sampling Behavior



Source: [Lindauer et al. 2019]

- Plot of a 1D or 2D function
 - Background shows the ground truth (real function values)
 - Dots are sampled points in the search space
 - Typical approach in Bayesian Optimization community
- ~> Impossible for higher dimensional problems?

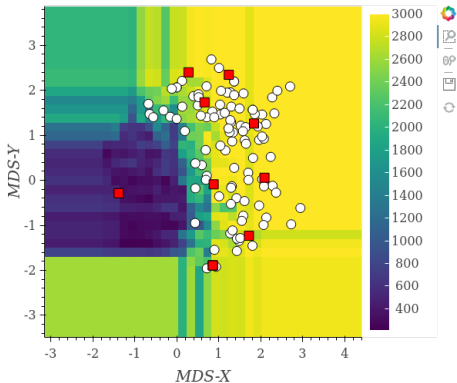
Analyzing the Sampling Behavior in N -D



- Same idea as before
but we have to project N -D into 2-D

Source: [Lindauer et al. 2019]

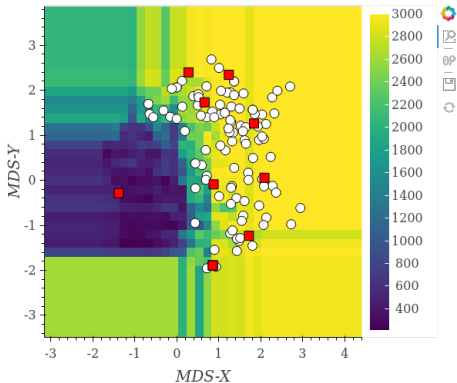
Analyzing the Sampling Behavior in N -D



Source: [Lindauer et al. 2019]

- Same idea as before
but we have to project N -D into 2-D
 - 1 Use an MDS to project down to 2-D

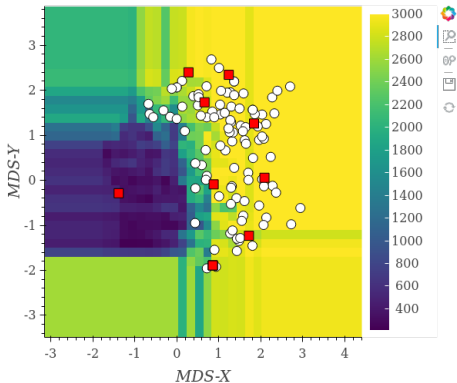
Analyzing the Sampling Behavior in N -D



Source: [Lindauer et al. 2019]

- Same idea as before
but we have to project N -D into 2-D
 - 1 Use an MDS to project down to 2-D
 - 2 Each dot is single hyperparameter configuration

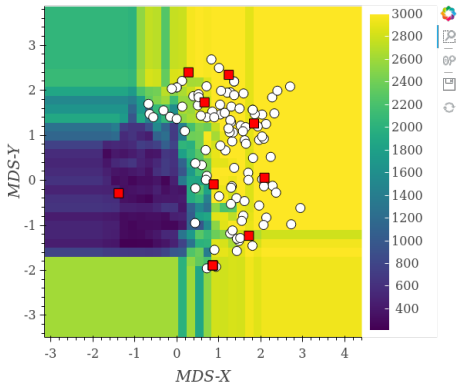
Analyzing the Sampling Behavior in N -D



Source: [Lindauer et al. 2019]

- Same idea as before
but we have to project N -D into 2-D
 - 1 Use an MDS to project down to 2-D
 - 2 Each dot is single hyperparameter configuration
 - 3 Red squares are intermediate incumbents

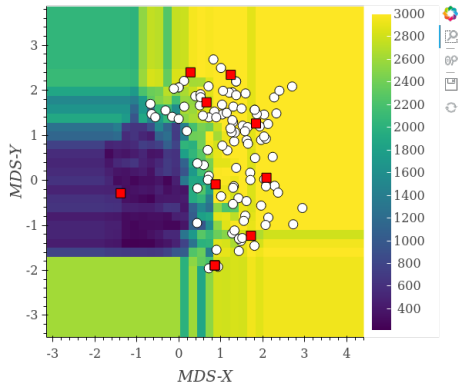
Analyzing the Sampling Behavior in N -D



Source: [Lindauer et al. 2019]

- Same idea as before
but we have to project N -D into 2-D
 - 1 Use an MDS to project down to 2-D
 - 2 Each dot is single hyperparameter configuration
 - 3 Red squares are intermediate incumbents
 - 4 The background is colored wrt a performance-estimate (e.g., reusing model fitted during BO)

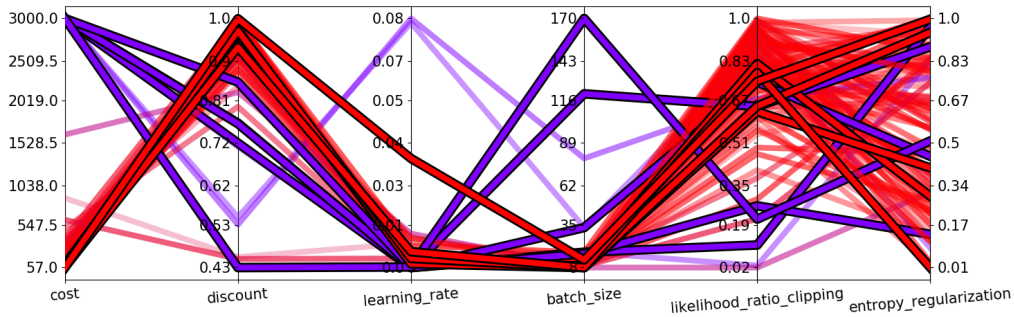
Analyzing the Sampling Behavior in N -D



Source: [Lindauer et al. 2019]

- Same idea as before
but we have to project N -D into 2-D
 - 1 Use an MDS to project down to 2-D
 - 2 Each dot is single hyperparameter configuration
 - 3 Red squares are intermediate incumbents
 - 4 The background is colored wrt a performance-estimate
(e.g., reusing model fitted during BO)
 - 5 Extension: Animation by showing
how points get added over time

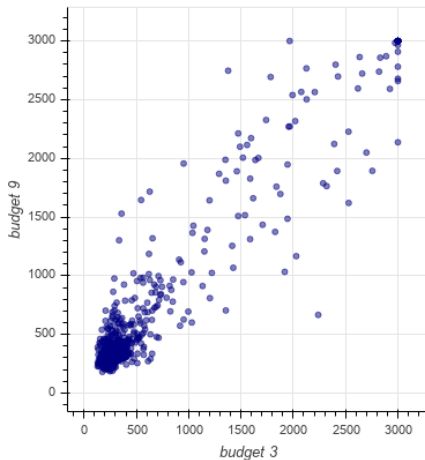
Parallel Coordinate Plot [Golovin et al. 2017]



Source: [Lindauer et al. 2019]

- Each coordinate is one hyperparameter;
- Except the most left one: cost or loss

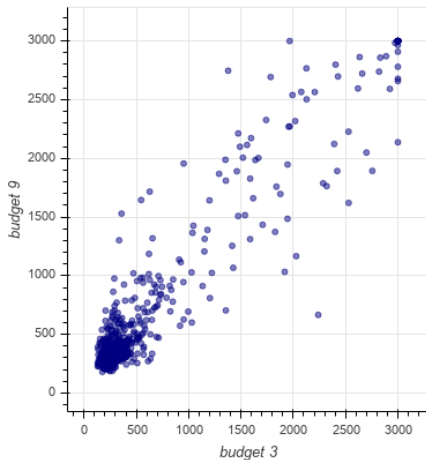
Multi-Fidelity Checks



- Challenge of multi-fidelity approaches:
 - ▶ How to choose the fidelities (a.k.a. budgets)

Source: [Lindauer et al. 2019]

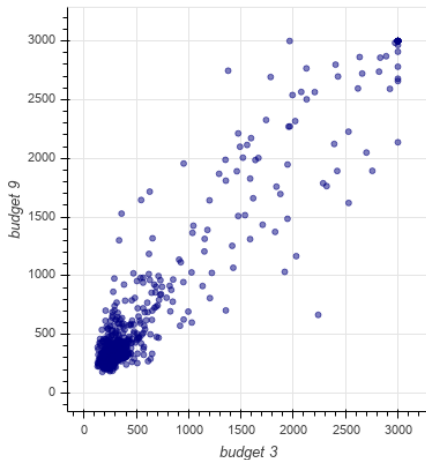
Multi-Fidelity Checks



- Challenge of multi-fidelity approaches:
 - ▶ How to choose the fidelities (a.k.a. budgets)
- Important Property:
 - ▶ Decisions on small budgets should be reasonable for higher budgets

Source: [Lindauer et al. 2019]

Multi-Fidelity Checks



Source: [Lindauer et al. 2019]

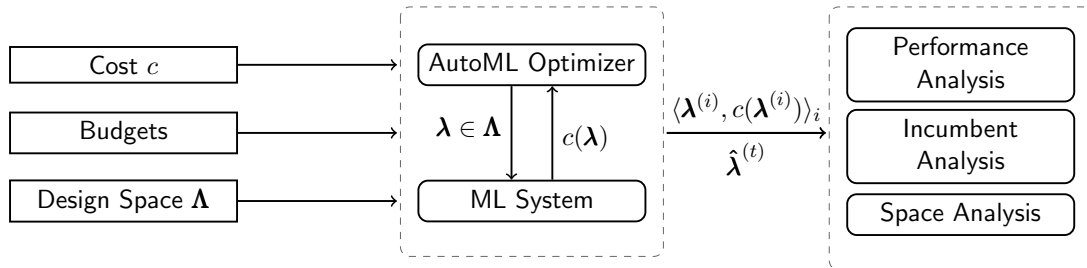
- Challenge of multi-fidelity approaches:
 - ▶ How to choose the fidelities (a.k.a. budgets)
- Important Property:
 - ▶ Decisions on small budgets should be reasonable for higher budgets
- Analysis:
 - 1 Scatter plot of performance on Budget X vs. Budget Y
 - 2 Each dot is sampled hyperparameter configuration
 - 3 Compute rank correlation (here: 0.69)

AutoML: Interpretability

Incumbent Analysis and Local Hyperparameter Importance

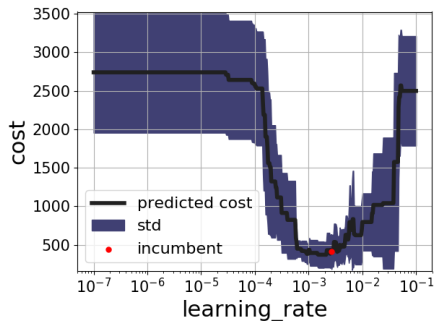
Bernd Bischl Frank Hutter Lars Kotthoff
Marius Lindauer Joaquin Vanschoren

Idea



~> focus on why is the eventually returned configuration a good choice

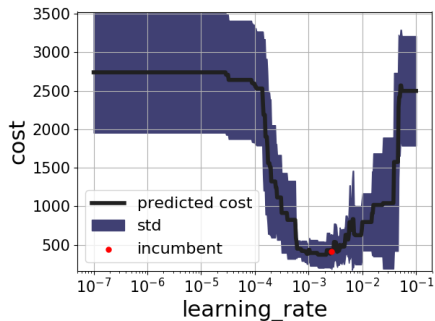
Local Importance [Biedenkapp et al. 2018]



- Typical question of users:
 - ▶ How would the performance change if we change hyperparameter λ_i ?

Source: [Lindauer et al. 2019]

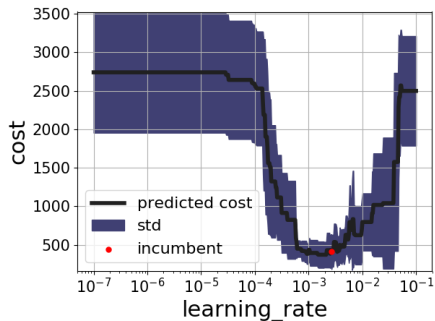
Local Importance [Biedenkapp et al. 2018]



Source: [Lindauer et al. 2019]

- Typical question of users:
 - ▶ How would the performance change if we change hyperparameter λ_i ?
- Problem: Running full study is often too expensive
 - ▶ Each run of an ML-system is potential expensive

Local Importance [Biedenkapp et al. 2018]



Source: [Lindauer et al. 2019]

- Typical question of users:
 - ▶ How would the performance change if we change hyperparameter λ_i ?
- Problem: Running full study is often too expensive
 - ▶ Each run of an ML-system is potential expensive
- Key Ideas:
 - ▶ Re-use probabilistic models as trained in BO
 - ▶ Plot performance change around $\hat{\lambda}^{(t)}$ along each dimension

$$\text{VAR}_{\boldsymbol{\lambda}}(i) = \sum_{v \in \boldsymbol{\Lambda}_i} (\mathbb{E}_{v \sim \boldsymbol{\Lambda}_i} [L(\boldsymbol{\lambda})] - L(\boldsymbol{\lambda}[\boldsymbol{\lambda}_i := v]))^2 \quad (1)$$

$$\text{VAR}_{\boldsymbol{\lambda}}(i) = \sum_{v \in \boldsymbol{\Lambda}_i} (\mathbb{E}_{v \sim \boldsymbol{\Lambda}_i} [L(\boldsymbol{\lambda})] - L(\boldsymbol{\lambda}[\boldsymbol{\lambda}_i := v]))^2 \quad (1)$$

$$\text{LPI}(i \mid \boldsymbol{\lambda}) = \frac{\text{VAR}_{\boldsymbol{\lambda}}(i)}{\sum_j \text{VAR}_{\boldsymbol{\lambda}}(j)} \quad (2)$$

$$\text{VAR}_{\boldsymbol{\lambda}}(i) = \sum_{v \in \boldsymbol{\Lambda}_i} (\mathbb{E}_{v \sim \boldsymbol{\Lambda}_i} [L(\boldsymbol{\lambda})] - L(\boldsymbol{\lambda}[\boldsymbol{\lambda}_i := v]))^2 \quad (1)$$

$$\text{LPI}(i \mid \boldsymbol{\lambda}) = \frac{\text{VAR}_{\boldsymbol{\lambda}}(i)}{\sum_j \text{VAR}_{\boldsymbol{\lambda}}(j)} \quad (2)$$

~> While fixing all other hyperparameters to the incumbent value,
the hyperparameter with the highest variance is the most important one

Ablation Study for Importance

- Users often start from some kind of default configuration
 - ① As given in the documentation
 - ② Or as always used in the last time

Ablation Study for Importance

- Users often start from some kind of default configuration
 - ① As given in the documentation
 - ② Or as always used in the last time
- **Key Idea:** Going from the default to the automatically optimized configuration, which choices were important?

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42]\end{aligned}$$

Ablation Study for Importance

- Users often start from some kind of default configuration
 - ① As given in the documentation
 - ② Or as always used in the last time
- **Key Idea:** Going from the default to the automatically optimized configuration, which choices were important?

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42]\end{aligned}$$

- Cheap approach: Assess $\lambda^{(\text{end})}$ with each hyperparameter value from $\lambda^{(\text{start})}$

Ablation Study for Importance

- Users often start from some kind of default configuration
 - ① As given in the documentation
 - ② Or as always used in the last time
- **Key Idea:** Going from the default to the automatically optimized configuration, which choices were important?

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42]\end{aligned}$$

- Cheap approach: Assess $\lambda^{(\text{end})}$ with each hyperparameter value from $\lambda^{(\text{start})}$
- Expensive approach: Try all mixtures of $\lambda^{(\text{end})}$ and $\lambda^{(\text{start})}$
 - ▶ Only feasible for small spaces and fairly cheap ML systems

Ablation Study for Importance

- Users often start from some kind of default configuration
 - ① As given in the documentation
 - ② Or as always used in the last time
- **Key Idea:** Going from the default to the automatically optimized configuration, which choices were important?

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42]\end{aligned}$$

- Cheap approach: Assess $\lambda^{(\text{end})}$ with each hyperparameter value from $\lambda^{(\text{start})}$
- Expensive approach: Try all mixtures of $\lambda^{(\text{end})}$ and $\lambda^{(\text{start})}$
 - ▶ Only feasible for small spaces and fairly cheap ML systems
- Trade-off: Find a way from $\lambda^{(\text{start})}$ to $\lambda^{(\text{end})}$ in a greedy fashion [Fawcett and Hoos. 2016]

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

1st Iteration:

$$\lambda^{(1)} = [0.98, 1, 0, 100] \quad L_1 = 19\%$$

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

1st Iteration:

$$\begin{aligned}\lambda^{(1)} &= [0.98, 1, 0, 100] & L_1 &= 19\% \\ \lambda^{(2)} &= [1, 2.42, 0, 100] & L_2 &= 20\%\end{aligned}$$

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

1st Iteration:

$$\begin{aligned}\lambda^{(1)} &= [0.98, 1, 0, 100] & L_1 &= 19\% \\ \lambda^{(2)} &= [1, 2.42, 0, 100] & L_2 &= 20\% \\ \lambda^{(3)} &= [1, 1, 1, 100] & L_3 &= 7\%\end{aligned}$$

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

1st Iteration:

$$\begin{aligned}\lambda^{(1)} &= [0.98, 1, 0, 100] & L_1 &= 19\% \\ \lambda^{(2)} &= [1, 2.42, 0, 100] & L_2 &= 20\% \\ \lambda^{(3)} &= [1, 1, 1, 100] & L_3 &= 7\% \\ \lambda^{(4)} &= [1, 1, 0, 42] & L_4 &= 16\%\end{aligned}$$

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

1st Iteration:

$$\begin{aligned}\lambda^{(1)} &= [0.98, 1, 0, 100] & L_1 &= 19\% \\ \lambda^{(2)} &= [1, 2.42, 0, 100] & L_2 &= 20\% \\ \lambda^{(3)} &= [1, 1, 1, 100] & L_3 &= 7\% \\ \lambda^{(4)} &= [1, 1, 0, 42] & L_4 &= 16\%\end{aligned}$$

\rightsquigarrow 1st step: λ_2 – flipping hyperparameter 3

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(s1)} &= [1, 1, \textcolor{blue}{1}, 100] & L &= 7\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

2nd Iteration:

$$\lambda^{(1)} = [\textcolor{blue}{0.98}, 1, 1, 100] \quad L_1 = 6\%$$

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(s1)} &= [1, 1, 1, 100] & L &= 7\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

2nd Iteration:

$$\begin{aligned}\lambda^{(1)} &= [0.98, 1, 1, 100] & L_1 &= 6\% \\ \lambda^{(2)} &= [1, 2.42, 1, 100] & L_2 &= 7\%\end{aligned}$$

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(s1)} &= [1, 1, 1, 100] & L &= 7\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

2nd Iteration:

$$\begin{aligned}\lambda^{(1)} &= [0.98, 1, 1, 100] & L_1 &= 6\% \\ \lambda^{(2)} &= [1, 2.42, 1, 100] & L_2 &= 7\% \\ \lambda^{(3)} &= [1, 1, 1, 42] & L_3 &= 5\%\end{aligned}$$

\rightsquigarrow 2nd step: λ_3 – flipping hyperparameter 4

Greedy Ablation Study

Given:

$$\begin{aligned}\lambda^{(\text{start})} &= [1, 1, 0, 100] & L_{\text{start}} &= 20\% \\ \lambda^{(s1)} &= [1, 1, 1, 100] & L &= 7\% \\ \lambda^{(s2)} &= [1, 1, 1, 42] & L &= 5\% \\ \lambda^{(\text{end})} &= [0.98, 2.42, 1, 42] & L_{\text{end}} &= 4\%\end{aligned}$$

3rd Iteration:

$$\begin{aligned}\lambda^{(1)} &= [0.98, 1, 1, 100] & L_1 &= 4\% \\ \lambda^{(2)} &= [1, 2.42, 1, 100] & L_2 &= 5\%\end{aligned}$$

\rightsquigarrow 2nd step: λ_3 – flipping hyperparameter 1

Greedy Ablation Study

Ablation Path:

$$\lambda^{(\text{start})} = [1, 1, 0, 100] \quad L_{\text{start}} = 20\%$$

$$\lambda^{(s1)} = [1, 1, 1, 100] \quad L = 7\%$$

$$\lambda^{(s1)} = [1, 1, 1, 42] \quad L = 5\%$$

$$\lambda^{(s3)} = [0.98, 1, 1, 42] \quad L = 4\%$$

$$\lambda^{(s4)} = [0.98, 2.42, 1, 42] \quad L = 4\%$$

$$\lambda^{(\text{end})} = [0.98, 2.42, 1, 42] \quad L_{\text{end}} = 4\%$$

Greedy Ablation Pseudo Code

Algorithm 1 Greedy Ablation

Input : Algorithm \mathcal{A} with configuration space Λ , start configuration $\lambda^{(\text{start})}$,
end configuration $\lambda^{(\text{end})}$, cost metric c

$\lambda \leftarrow \lambda^{(\text{start})};$

$P \leftarrow [] ;$

Greedy Ablation Pseudo Code

Algorithm 2 Greedy Ablation

Input : Algorithm \mathcal{A} with configuration space $\mathbf{\Lambda}$, start configuration $\lambda^{(\text{start})}$,
end configuration $\lambda^{(\text{end})}$, cost metric c

$\lambda \leftarrow \lambda^{(\text{start})};$

$P \leftarrow [];$

foreach $t \in \{1 \dots |\mathbf{\Lambda}|\}$ **do**

|

Greedy Ablation Pseudo Code

Algorithm 3 Greedy Ablation

Input : Algorithm \mathcal{A} with configuration space $\mathbf{\Lambda}$, start configuration $\lambda^{(\text{start})}$,
end configuration $\lambda^{(\text{end})}$, cost metric c

$\lambda \leftarrow \lambda^{(\text{start})};$

$P \leftarrow [];$

foreach $t \in \{1 \dots |\mathbf{\Lambda}|\}$ **do**

foreach $\delta \in \Delta(\lambda, \lambda^{(\text{end})})$ **do**

$\lambda'_\delta \leftarrow \text{apply } \delta \text{ to } \lambda;$

 evaluate $c(\lambda'_\delta);$

Greedy Ablation Pseudo Code

Algorithm 4 Greedy Ablation

Input : Algorithm \mathcal{A} with configuration space Λ , start configuration $\lambda^{(\text{start})}$,
end configuration $\lambda^{(\text{end})}$, cost metric c

$\lambda \leftarrow \lambda^{(\text{start})};$

$P \leftarrow [];$

foreach $t \in \{1 \dots |\Lambda|\}$ **do**

foreach $\delta \in \Delta(\lambda, \lambda^{(\text{end})})$ **do**

$\lambda'_\delta \leftarrow \text{apply } \delta \text{ to } \lambda;$

 evaluate $c(\lambda'_\delta);$

 Determine most important change $\delta^* \in \arg \min_{\delta \in \Delta(\lambda, \lambda^{(\text{end})})} c(\lambda_\delta);$

$\lambda \leftarrow \text{apply } \delta^* \text{ to } \lambda;$

$P.\text{append}(\delta^*);$

Greedy Ablation Pseudo Code

Algorithm 5 Greedy Ablation

Input : Algorithm \mathcal{A} with configuration space Λ , start configuration $\lambda^{(\text{start})}$, end configuration $\lambda^{(\text{end})}$, cost metric c

$\lambda \leftarrow \lambda^{(\text{start})};$

$P \leftarrow [];$

foreach $t \in \{1 \dots |\Lambda|\}$ **do**

foreach $\delta \in \Delta(\lambda, \lambda^{(\text{end})})$ **do**

$\lambda'_\delta \leftarrow \text{apply } \delta \text{ to } \lambda;$

 evaluate $c(\lambda'_\delta);$

 Determine most important change $\delta^* \in \arg \min_{\delta \in \Delta(\lambda, \lambda^{(\text{end})})} c(\lambda_\delta);$

$\lambda \leftarrow \text{apply } \delta^* \text{ to } \lambda;$

$P.\text{append}(\delta^*);$

return Ablation path P

Remarks on Ablation

- Even this greedy ablation requires $\mathcal{O}(n^2)$ steps

Remarks on Ablation

- Even this greedy ablation requires $\mathcal{O}(n^2)$ steps
- ~→ We can also speedup that up by using surrogate models
[Biedenkapp et al. 2017]

Remarks on Ablation

- Even this greedy ablation requires $\mathcal{O}(n^2)$ steps
- \rightsquigarrow We can also speedup that up by using surrogate models
[Biedenkapp et al. 2017]
- Common observations:
 - ① Some hyperparameters might not matter (λ_2 in the example)

Remarks on Ablation

- Even this greedy ablation requires $\mathcal{O}(n^2)$ steps
- \rightsquigarrow We can also speedup that up by using surrogate models
[Biedenkapp et al. 2017]
- Common observations:
 - 1 Some hyperparameters might not matter (λ_2 in the example)
 - 2 Often only a few of the hyperparameters have an big impact

Remarks on Ablation

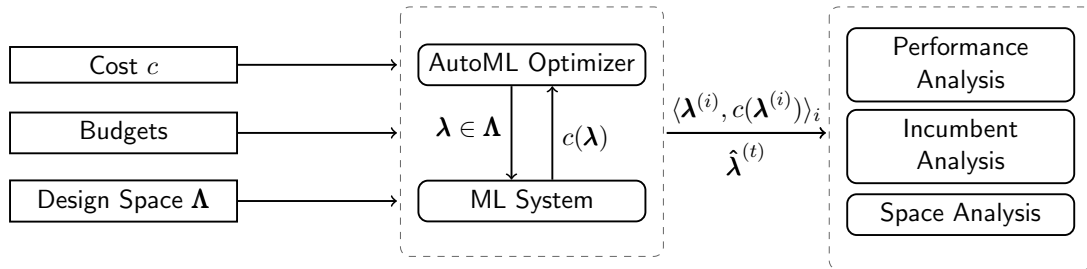
- Even this greedy ablation requires $\mathcal{O}(n^2)$ steps
- \rightsquigarrow We can also speedup that up by using surrogate models
[Biedenkapp et al. 2017]
- Common observations:
 - 1 Some hyperparameters might not matter (λ_2 in the example)
 - 2 Often only a few of the hyperparameters have an big impact
 - 3 You have plateaus in your ablation path because of interaction effects

AutoML: Interpretability

Global Hyperparameter Importance

Bernd Bischl Frank Hutter Lars Kotthoff
Marius Lindauer Joaquin Vanschoren

Idea



~> focus on which hyperparameters are important across the entire search space

Importance Analysis of Surrogate Model

- **Key Idea:** Surrogate models ($\Lambda \rightarrow \mathbb{R}$) learn from observations how to predict the performance of a **hyperparameter configuration** $\lambda \in \Lambda$

Importance Analysis of Surrogate Model

- **Key Idea:** Surrogate models ($\Lambda \rightarrow \mathbb{R}$) learn from observations how to predict the performance of a **hyperparameter configuration** $\lambda \in \Lambda$
- ~> These models can be used to figure out which hyperparameter was important

Importance Analysis of Surrogate Model

- **Key Idea:** Surrogate models ($\Lambda \rightarrow \mathbb{R}$) learn from observations how to predict the performance of a **hyperparameter configuration** $\lambda \in \Lambda$
- ~> These models can be used to figure out which hyperparameter was important
- For example:
 - ▶ Use forward selection [Hutter et al. 2013]
 - ▶ Use automatic feature relevance determination of the model (e.g., of a surrogate model based on random forest)

Importance Analysis of Surrogate Model

- **Key Idea:** Surrogate models ($\Lambda \rightarrow \mathbb{R}$) learn from observations how to predict the performance of a **hyperparameter configuration** $\lambda \in \Lambda$
- ~> These models can be used to figure out which hyperparameter was important
- For example:
 - ▶ Use forward selection [Hutter et al. 2013]
 - ▶ Use automatic feature relevance determination of the model (e.g., of a surrogate model based on random forest)
- Advantages:
 - ▶ Very cheap to do, since we only have to query the surrogate model several times

Importance Analysis of Surrogate Model

- **Key Idea:** Surrogate models ($\Lambda \rightarrow \mathbb{R}$) learn from observations how to predict the performance of a **hyperparameter configuration** $\lambda \in \Lambda$
- ~> These models can be used to figure out which hyperparameter was important
- For example:
 - ▶ Use forward selection [Hutter et al. 2013]
 - ▶ Use automatic feature relevance determination of the model (e.g., of a surrogate model based on random forest)
- Advantages:
 - ▶ Very cheap to do, since we only have to query the surrogate model several times
- Potential drawback:
 - ▶ The surrogate model might overfit to different subsets of the hyperparameters (if we don't provide sufficient data)

Global Importance Analysis

- **Key idea:** What is the importance of a hyperparameter by marginalizing over all other hyperparameter effects?

Global Importance Analysis

- **Key idea:** What is the importance of a hyperparameter by marginalizing over all other hyperparameter effects?
- **Key Insight:** We can use a surrogate model to compute these effects

Global Importance Analysis

- **Key idea:** What is the importance of a hyperparameter by marginalizing over all other hyperparameter effects?
- **Key Insight:** We can use a surrogate model to compute these effects

*f*ANOVA [Sobobl. 1993]

Write performance predictions as a sum of components:

$$\hat{y}(\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n) = \hat{f}_0 + \sum_{i=1}^n \hat{f}_i(\boldsymbol{\lambda}_i) + \sum_{i \neq j} \hat{f}_{ij}(\boldsymbol{\lambda}_i, \boldsymbol{\lambda}_j) + \dots$$

$\hat{y}(\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n) =$ average response + main effects +
2-D interaction effects + higher order effects

Global Importance Analysis

- **Key idea:** What is the importance of a hyperparameter by marginalizing over all other hyperparameter effects?
- **Key Insight:** We can use a surrogate model to compute these effects

*f*ANOVA [Sobobl. 1993]

Write performance predictions as a sum of components:

$$\begin{aligned}\hat{y}(\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n) &= \hat{f}_0 + \sum_{i=1}^n \hat{f}_i(\boldsymbol{\lambda}_i) + \sum_{i \neq j} \hat{f}_{ij}(\boldsymbol{\lambda}_i, \boldsymbol{\lambda}_j) + \dots \\ \hat{y}(\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n) &= \text{average response} + \text{main effects} + \\ &\quad \text{2-D interaction effects} + \text{higher order effects}\end{aligned}$$

Variance Decomposition

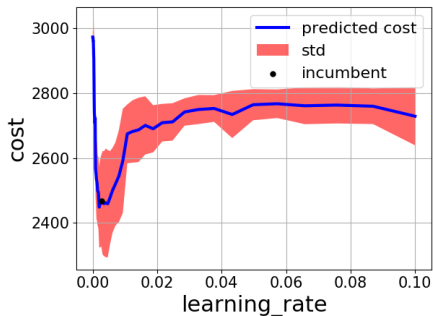
$$V = \frac{1}{||\boldsymbol{\Lambda}||} \int_{\boldsymbol{\lambda}_1} \dots \int_{\boldsymbol{\lambda}_n} [(\hat{y}(\boldsymbol{\lambda}) - \hat{f}_0)^2] d\boldsymbol{\lambda}_1 \dots d\boldsymbol{\lambda}_n$$

fANOVA Analysis

- The fANOVA and variance decomposition can be done efficiently in linear time if the surrogate model is a random forest [Hutter et al. 2014]

fANOVA Analysis

- The fANOVA and variance decomposition can be done efficiently in linear time if the surrogate model is a random forest [Hutter et al. 2014]

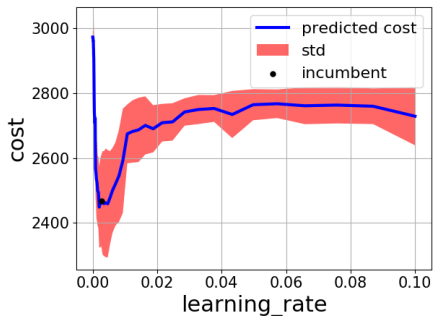


- predicted cost is marginalized over all other hyperparameter effects

Source: [Lindauer et al. 2019]

fANOVA Analysis

- The fANOVA and variance decomposition can be done efficiently in linear time if the surrogate model is a random forest [Hutter et al. 2014]



Source: [Lindauer et al. 2019]

- predicted cost is marginalized over all other hyperparameter effects
- **Warning:** The optimum on these curves does not have to be the global optimum across all hyperparameters

fANOVA Analysis

- How much of the variance can be explained by a hyperparameter (or combinations of hyperparameters) marginalized over all other parameters?

Table: Exemplary analysis of PPO on cartpole

Hyperparameter	Explained Variance
Discount rate	19.3 %
Batch size	15.7 %
Learning rate	3.7 %
Likelihood ration clipping	3.4%
...	

fANOVA Analysis

- How much of the variance can be explained by a hyperparameter (or combinations of hyperparameters) marginalized over all other parameters?

Table: Exemplary analysis of PPO on cartpole

Hyperparameter	Explained Variance
Discount rate	19.3 %
Batch size	15.7 %
Learning rate	3.7 %
Likelihood ration clipping	3.4%
...	
discount rate & batch size	10.4%
discount rate & likelihood ration clipping	4.4%
...	

Remarks on fANOVA

- Given compute higher-order interaction effects
 - ▶ Often too expensive for more than 2 or 3 dimensions

Remarks on fANOVA

- Given compute higher-order interaction effects
 - ▶ Often too expensive for more than 2 or 3 dimensions
- Implicit assumption: the surrogate model models the space fairly well

Remarks on fANOVA

- Given compute higher-order interaction effects
 - ▶ Often too expensive for more than 2 or 3 dimensions
- Implicit assumption: the surrogate model models the space fairly well
- Global analysis and local analysis of hyperparameter importance does not always agree
[Biedenkapp et al. 2018]

Remarks on fANOVA

- Given compute higher-order interaction effects
 - ▶ Often too expensive for more than 2 or 3 dimensions
- Implicit assumption: the surrogate model models the space fairly well
- Global analysis and local analysis of hyperparameter importance does not always agree
[Biedenkapp et al. 2018]
- ~> You should run both to get a good understanding of why an AutoML tool chose a configuration