

Auto-Sklearn: Automated Machine Learning in Python

Matthias Feurer



Katharina Eggensperger

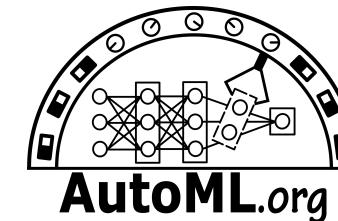


Eddie Bergman



Department of Computer Science

University of Freiburg, Germany



Find slides, notebooks and more here:
<https://github.com/automl/auto-sklearn-talks>

*Machine Learning for everyone
in 4 lines of code*

```
import autosklearn.classification
>>> cls = autosklearn.classification.AutoSklearnClassifier()
>>> cls.fit(X_train, y_train)
>>> predictions = cls.predict(X_test)
```

Our Vision

[automl / auto-sklearn](#) Public

Edit Pins ▾ Unwatch 214 Fork 1.2k Starred 6.5k

<> Code Issues 110 Pull requests 9 Discussions Actions Projects 1 Wiki ...

Contributors 78

Used by 296

+ 288

```
import autosklearn.classification
>>> cls = autosklearn.classification.AutoSklearnClassifier()
>>> cls.fit(X_train, y_train)
>>> predictions = cls.predict(X_test)
```



/automl/auto-sklearn

Goals for today & Outline

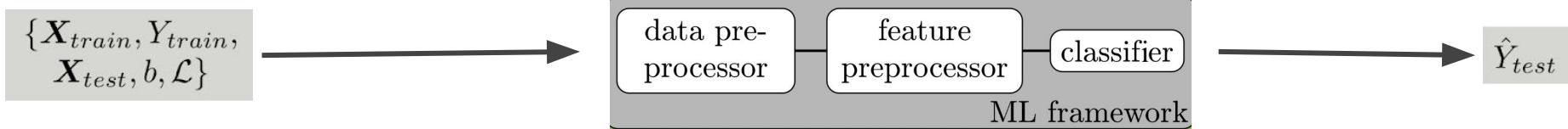
Goals

1. Understand how Auto-sklearn works
2. Apply Auto-Sklearn

Outline

1. Introduction to Auto-sklearn (25 mins)
2. **Task 1:** BYOP (10 mins)
3. **Task 2:** ASKL (20 mins)
4. **Task 3:** MO-ASKL (20 mins)
5. **+ Bonus Tasks** (? mins)
6. Outroduction (5mins)

Design Space: Traditional ML with scikit-learn



Design Space: Traditional ML with scikit-learn

 $\{\mathbf{X}_{train}, Y_{train}, \mathbf{X}_{test}, b, \mathcal{L}\}$ 

preprocessor	#λ
extreml. rand. trees prepr.	5
fast ICA	4
feature agglomeration	4
kernel PCA	5
rand. kitchen sinks	2
linear SVM prepr.	3
no preprocessing	-
nystroem sampler	5
PCA	2
polynomial	3
random trees embed.	4
select percentile	2
select rates	3
one-hot encoding	2
imputation	1
balancing	1
rescaling	1

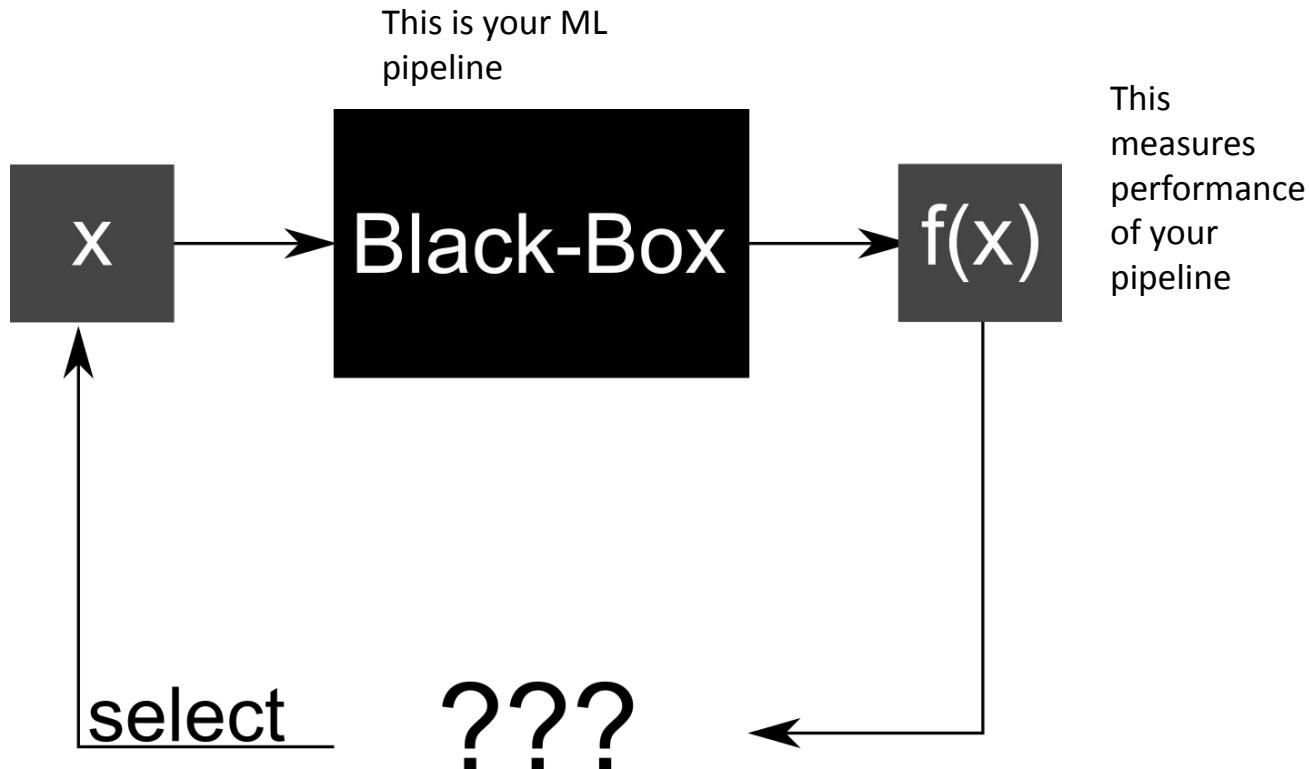


classifier	#λ
AdaBoost (AB)	4
Bernoulli naïve Bayes	2
decision tree (DT)	4
extreml. rand. trees	5
Gaussian naïve Bayes	-
gradient boosting (GB)	6
kNN	3
LDA	4
linear SVM	4
kernel SVM	7
multinomial naive Bayes	2
passive aggressive	3
QDA	2
random forest (RF)	5
Linear Class. (SGD)	10

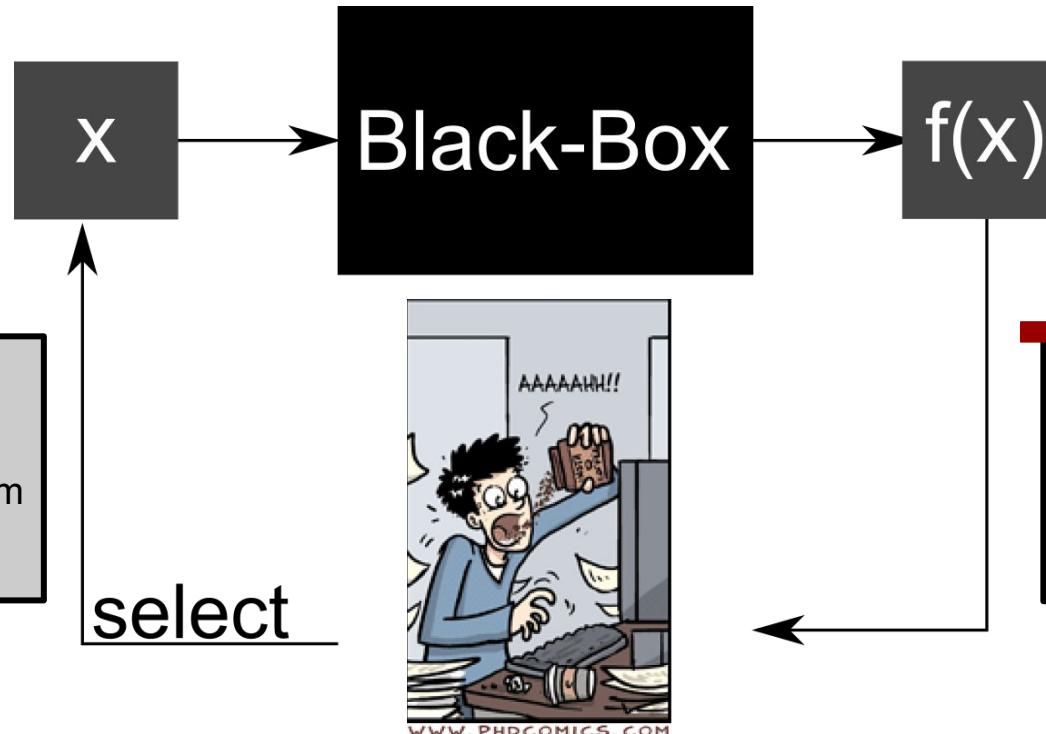
 \hat{Y}_{test}

Black Box Optimization

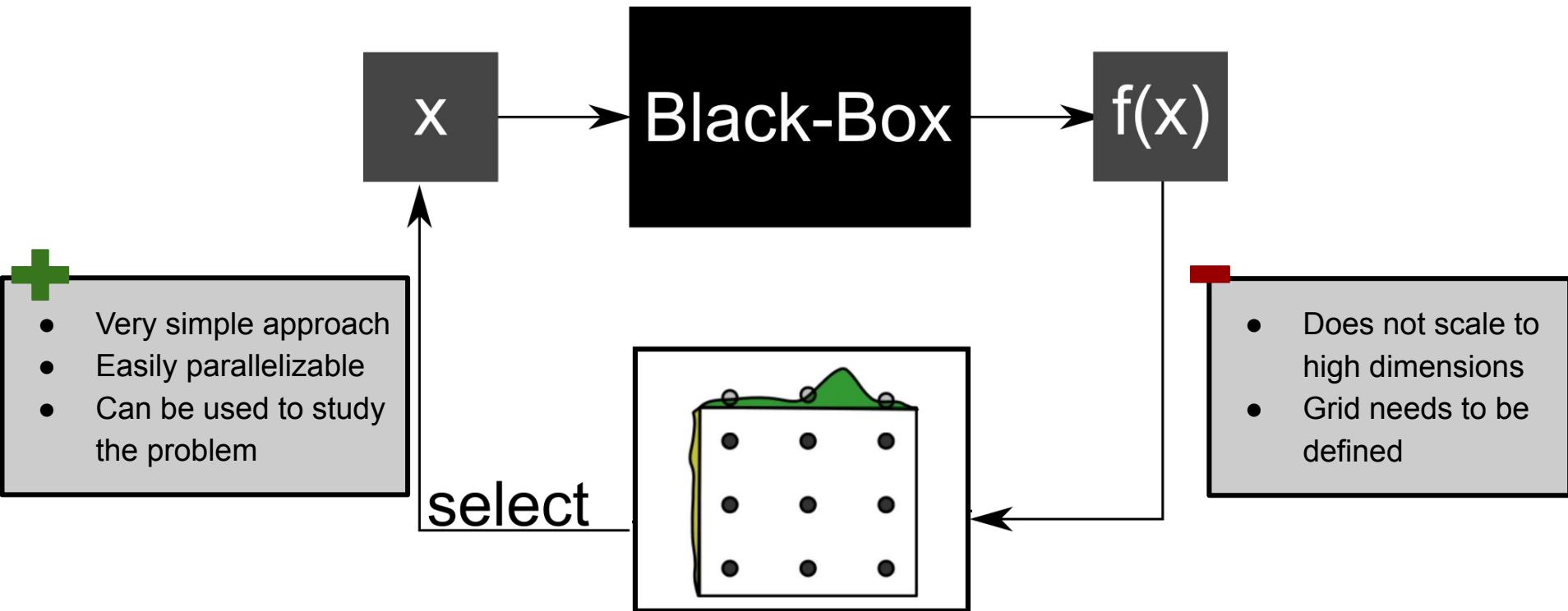
This is a configuration of a ML pipeline



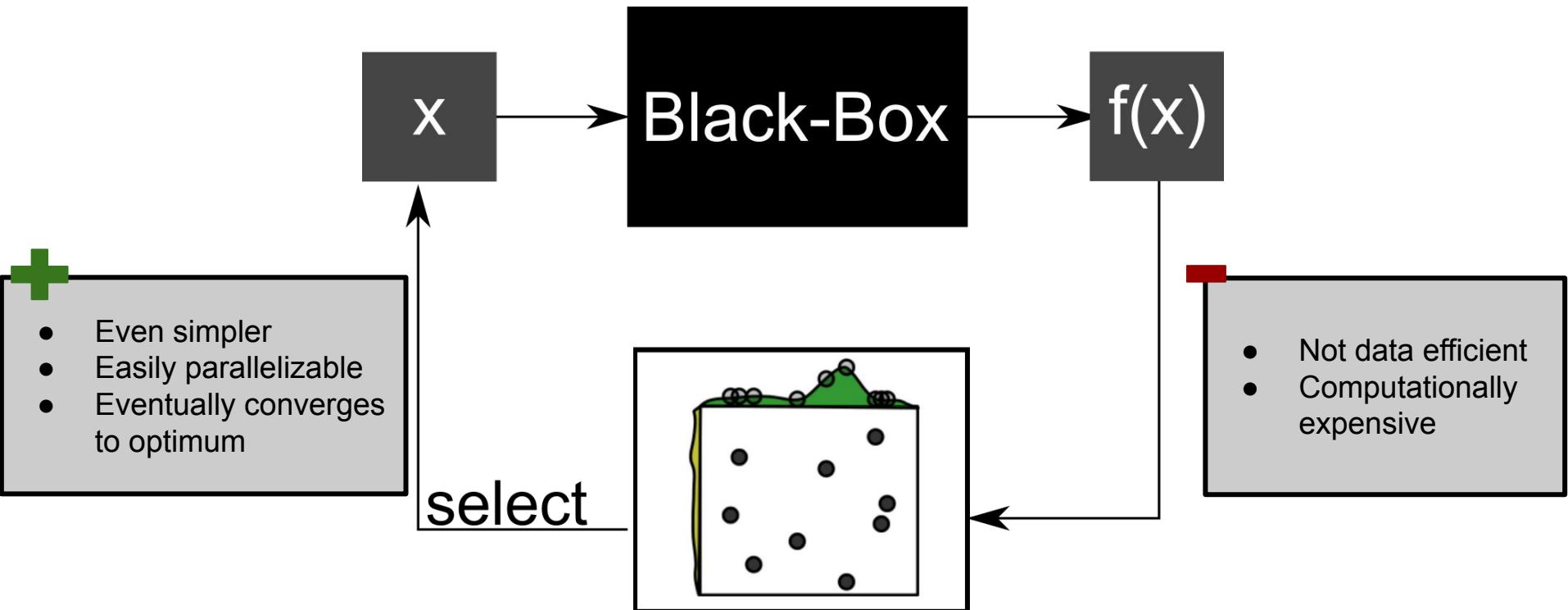
Black Box Optimization: The Human Optimizer



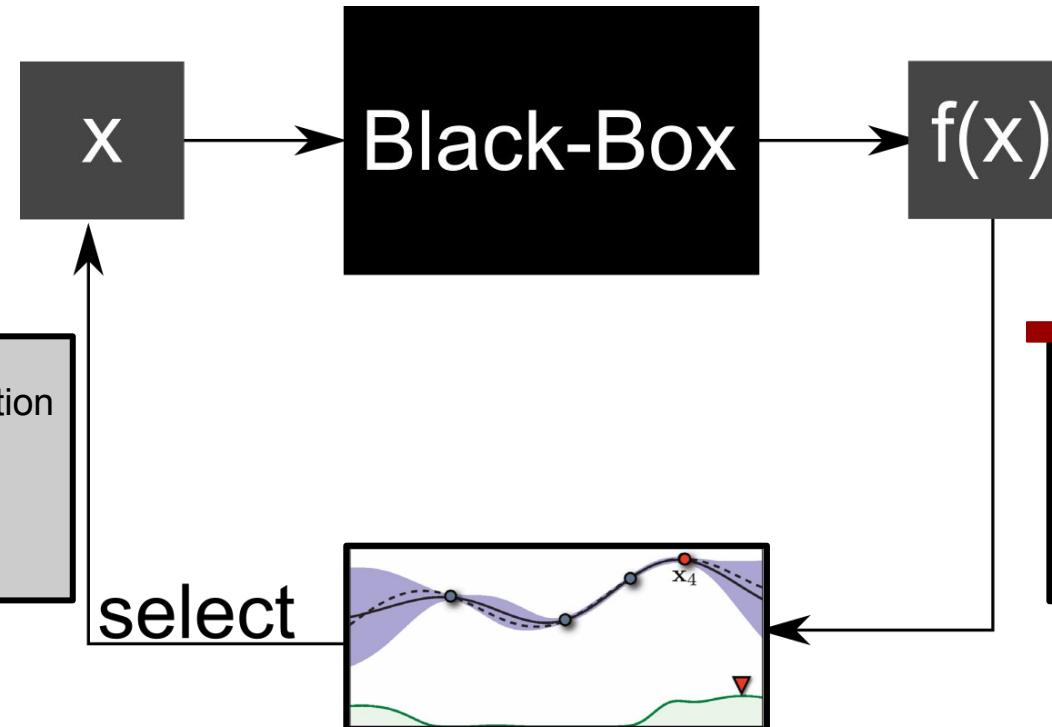
Black Box Optimization: Grid Search



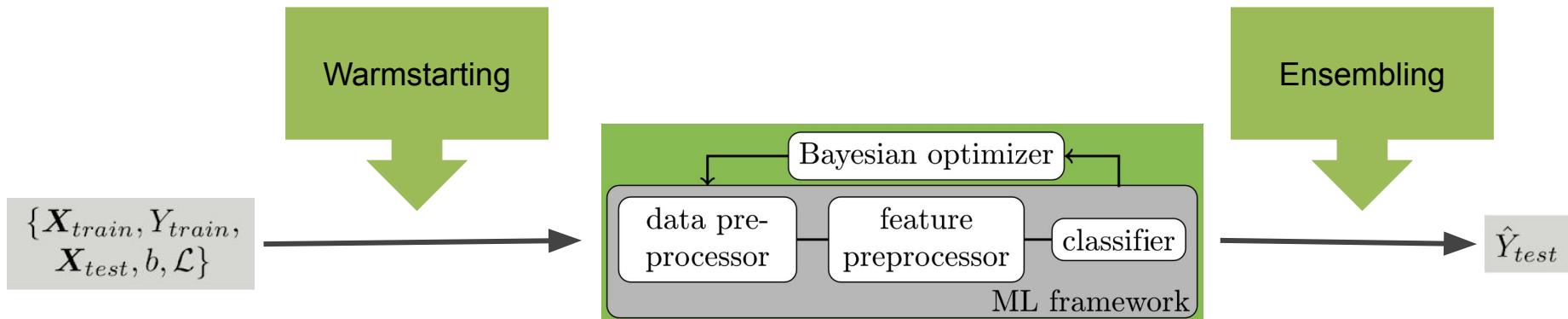
Black Box Optimization: Random Search



Black Box Optimization: Bayesian Optimization



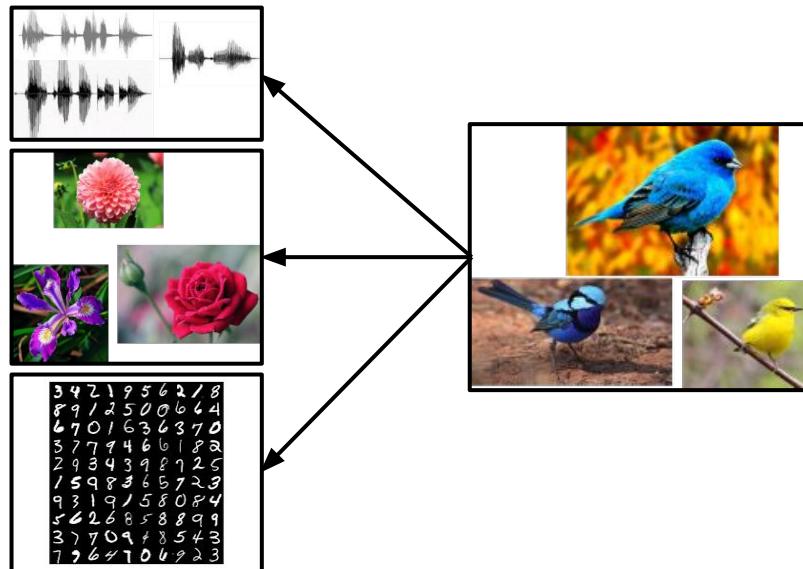
Design Space: Traditional ML with scikit-learn



More I: Meta-Learning

How to reuse previous experience?

→ Warmstart Bayesian Optimization



Offline / Before:

- 1) Collect >200 datasets
- 2) Find the best pipeline on each dataset

Online / For a new dataset:

- 1) Compute 38 meta-features, select 25 most similar previous datasets
- 2) Initialize optimization with best pipelines on those datasets

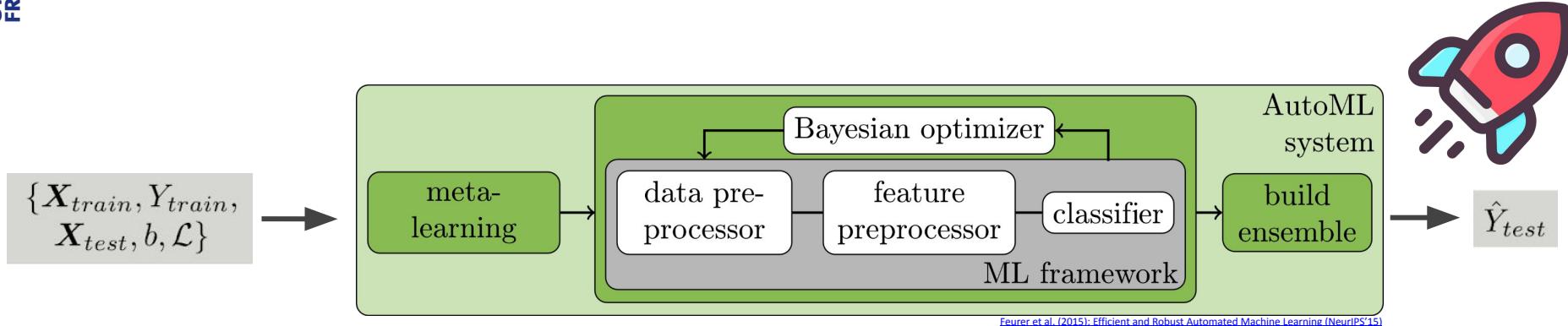
More II: Ensembling

How to get the best out of all evaluated models?

→ **Build an ensemble**



Auto-Sklearn 1.0



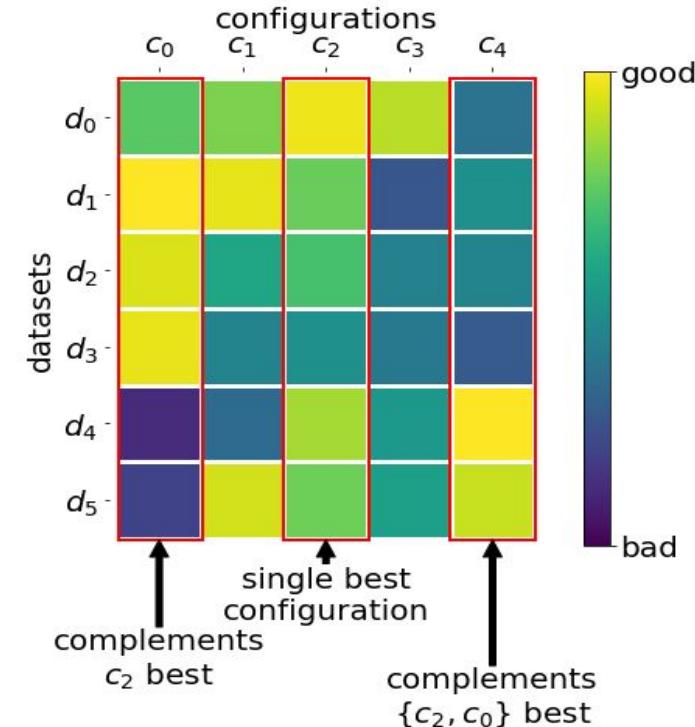
However, some things to be improved

- meta-features can be expensive to compute
- large datasets can be an issue

Even More I: Portfolios

Goal: Meta-Learning without meta-features

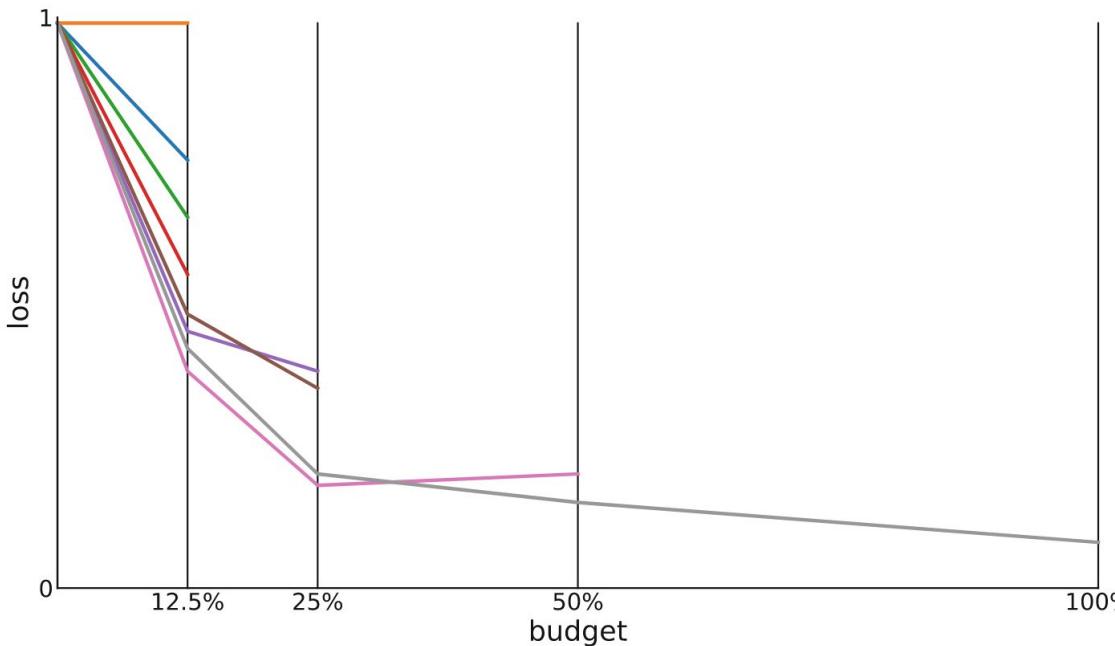
Idea: Construct a Portfolio
(a list of diverse pipelines)



Even More II: Successive Halving

Goal: Scale to large datasets.

Idea: Allocate more resources to promising pipelines



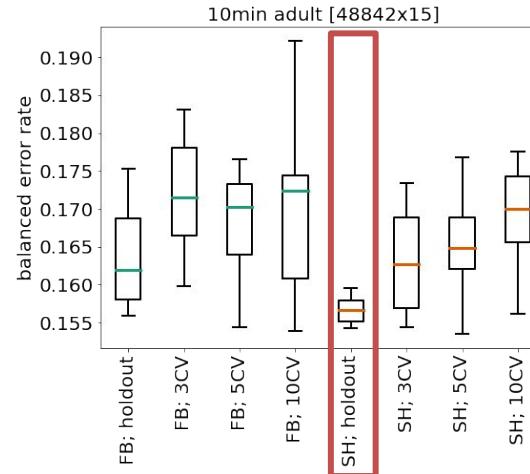
But what about
small datasets?

Image Credit - CC-BY

Matthias Feurer and Frank Hutter: *Hyperparameter optimization*

Automated Machine Learning, The Springer Series on Challenges in Machine Learning

Impact of the Optimization Strategy



But wait ... did we make it worse?

Can we automatically
select an optimization
policy?

→ Auto-sklearn 2.0

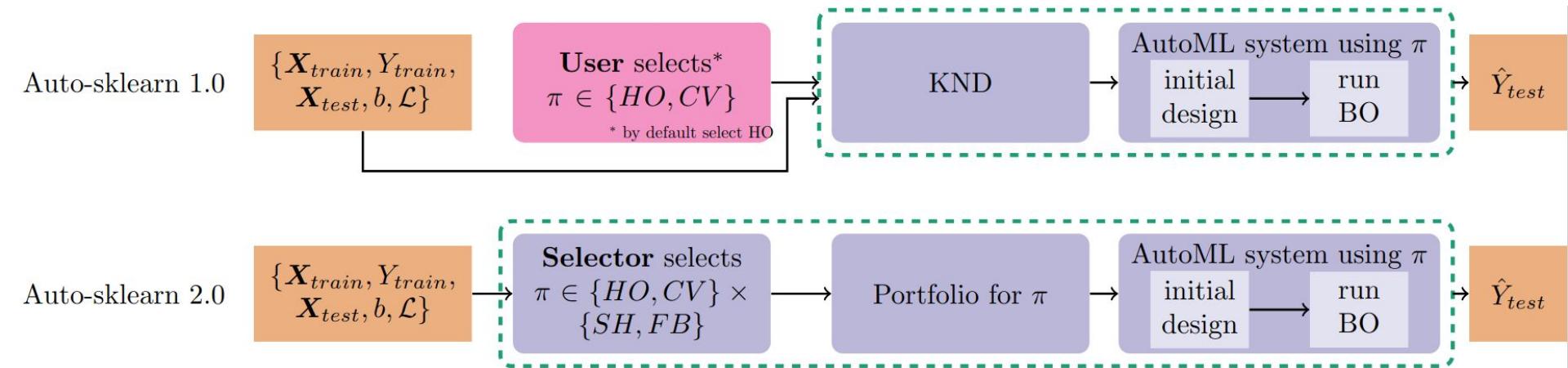


Image Credit - CC BY-NC-ND 2.0; by [Beagle Mama](#)

Yes, with a learned selector!

For more details see "[Feurer et al. \(2022\): Auto-Sklearn 2.0: Hands-free AutoML via Meta-Learning](#)"

Autosklearn 1.0 vs Auto-sklearn 2.0



Practical features

- Simplicity → follows scikit-learn API
- Parallelism → Uses Dask
- Extensibility → Simply add new algorithms
- Robustness → Limits on time and memory usage

... and many more:

- configurable (access to underlying SMAC)
- compatible with Pandas/numpy arrays
- support for text features
- ...

Summary

```
import autosklearn.classification
>>> cls = autosklearn.classification.AutoSklearnClassifier()
>>> cls.fit(X_train, y_train)
>>> predictions = cls.predict(X_test)
```

- based on **scikit-learn**; simple & familiar API
- integrates **latest research (>1K citations)**
- **>20K** downloads per month
- **BSD-3-Clause** License
- works best under **Linux**
- requires **Python>=3.7**



/automi/auto-sklearn



Matthias
Feurer
PhD Student @ ML Lab
Freiburg



Katharina
Eggensperger
PhD Student @ ML Lab
Freiburg



Edward
Bergman
Research Engineer @ ML Lab
Freiburg



Prof. Dr.
Marius Lindauer
Head of the ML Lab
Hannover

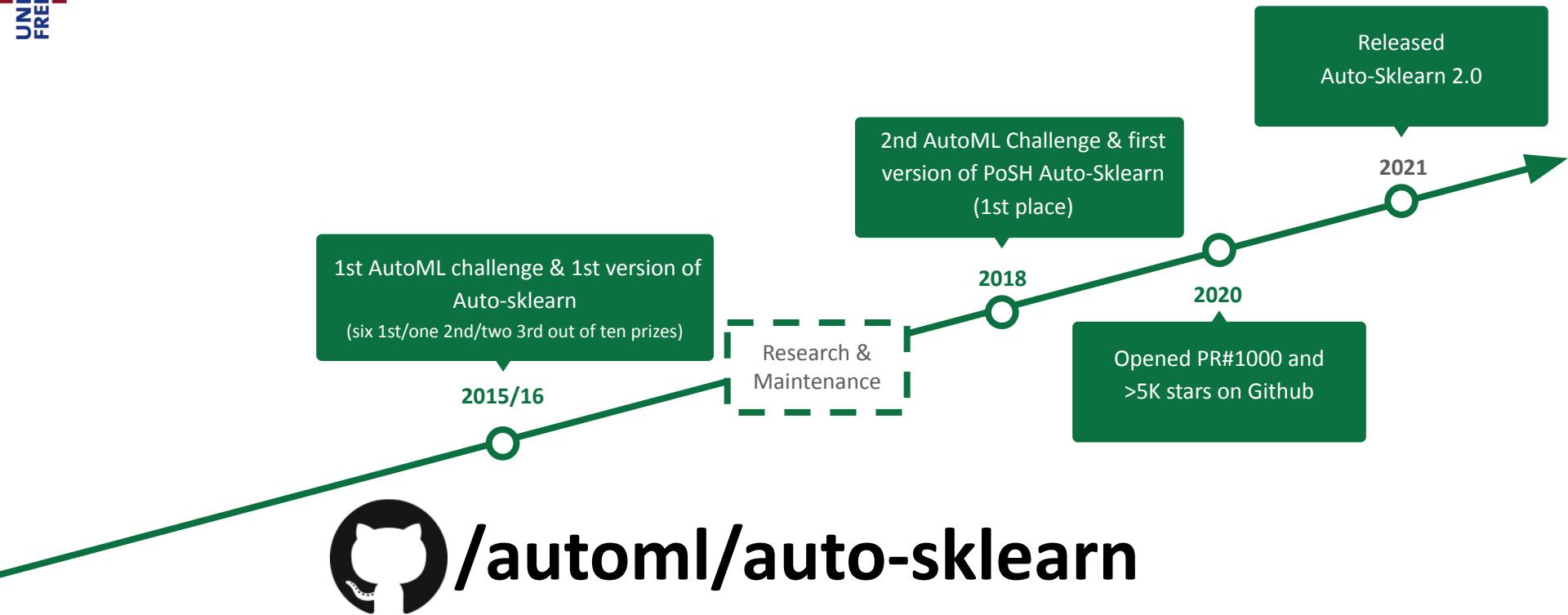


Prof. Dr.
Frank Hutter
Head of the ML Lab
Freiburg

Now: Hands-On Session

- Form groups of 3-5 people
- Copy and open the Colab notebook (see github.com/automl/auto-sklearn-talks)
- Work through
 - 1: BYOP
 - 2: ASKL
 - 3: Multi-objective ASKL
 - + Bonus tasks (if time)

Timeline & Success Stories



Matthias
Feurer
PhD Student @ ML Lab
Freiburg



Katharina
Eggensperger
PhD Student @ ML Lab
Freiburg



Edward
Bergman
Research Engineer @ ML Lab
Freiburg

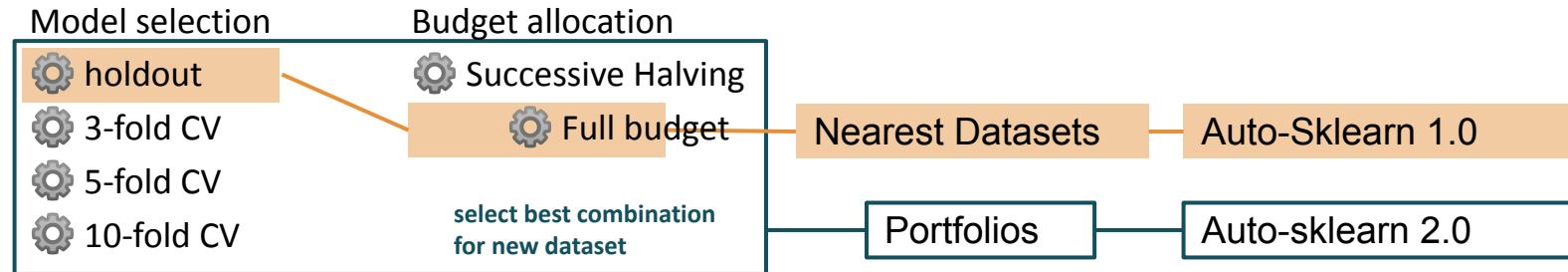


Prof. Dr.
Marius Lindauer
Head of the ML Lab
Hannover



Prof. Dr.
Frank Hutter
Head of the ML Lab
Freiburg

PoSH-Auto-sklearn



10MIN	60MIN
\emptyset	std

Auto-sklearn vs SMAC vs NASLib

	SMAC	NASLib	Auto-sklearn
Goal	<ul style="list-style-type: none">Algorithm configurationHyperparameter tuningBlackbox optimization	Neural Architecture Search research	<ul style="list-style-type: none"><i>ML in 4 lines of code</i>Automated Machine LearningDrop-in replacement for scikit-learn
Comments	Formed the basis for the winning entry to the NeurIPS 2020 Black-Box optimization competition!	Uses Blackbox and Greybox optimization under the hood	Uses SMAC under the hood
Alternatives	Hyperopt, Optuna, Ray, AX, DEHB, mlrMBO, etc...	None (so far)	Auto-PyTorch, AutoXGBoost, H2O, Auto-Gluon, etc...

Auto-sklearn - or why you should continue listening...

```
import autosklearn.classification
>>> cls = autosklearn.classification.AutoSklearnClassifier()
>>> cls.fit(X_train, y_train)
>>> predictions = cls.predict(X_test)
```

[automl / auto-sklearn](#) Public Edit Pins Unwatch 214 Fork 1.2k Starred 6.5k

<> [Code](#) (Issues 110) Pull requests 9 Discussions Actions Projects 1 Wiki ...

PyPI Stats [auto-sklearn](#)

Search [PyPI page](#) [Home page](#)
All packages [Author: Matthias Feurer](#)
Top packages [License: BSD3](#)
Track packages [Summary: Automated machine learning.](#)
[Latest version: 0.14.7](#)

Downloads last day: 1,671
Downloads last week: 9,213
Downloads last month: 40,289

Used by 296

Contributors 78

From github.com/automl/auto-sklearn 27