

AutoML with Auto-Sklearn

Matthias Feurer

 /__mfeurer__

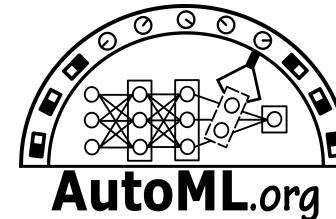
Katharina Eggensperger

 /KEggensperger

Eddie Bergman

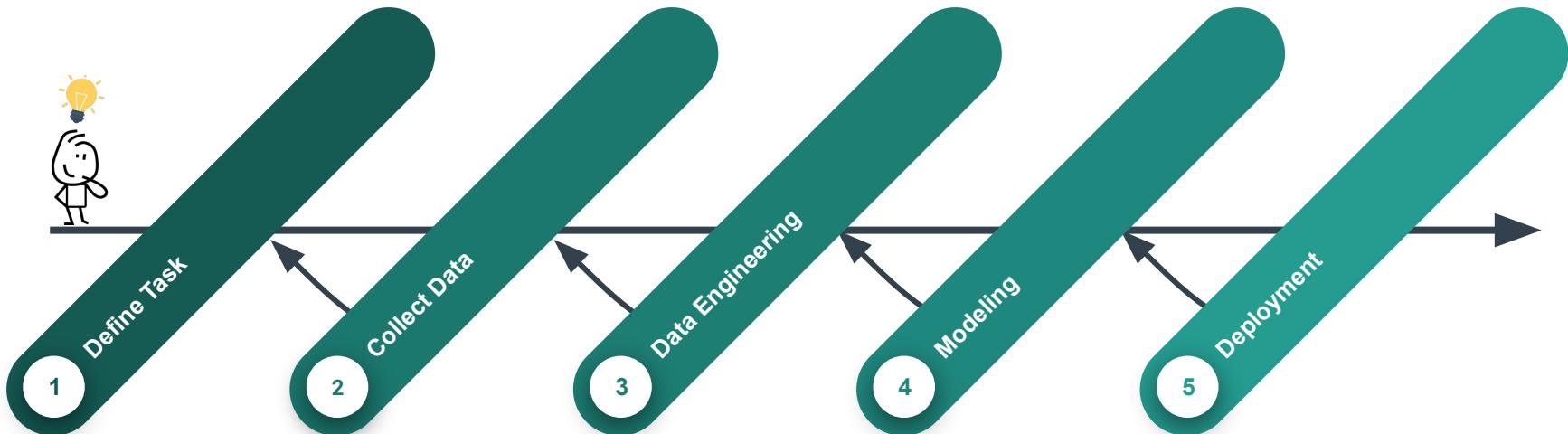
 /Edberg_Wardman

special guest today: Aron Bahram (**Hiwi working on the code**)

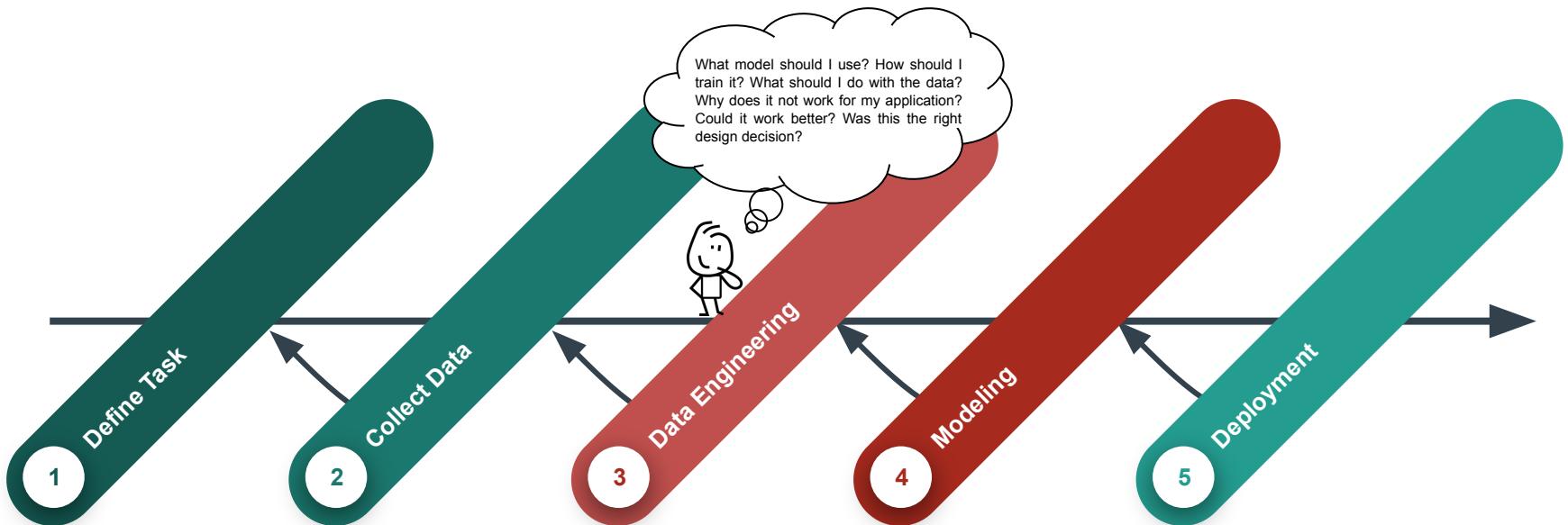


Find slides, notebooks and more here:
<https://github.com/automl/auto-sklearn-talks>

What is AutoML (and how can it help you)?



What is AutoML (and how can it help you)?



- automates parts of the ML workflow
- makes ML easy to use

*“Machine Learning
for everyone”*

What is AutoML (and how can it help you)?

Popular Approaches?

- Hyperparameter Tuning
- Neural Architecture Search
- **AutoML Systems**

Focus of today

Target User?

- non ML experts
- ML researchers
- ML practitioners

Use cases?

- fast prototyping
- strong baselines & fair comparisons
- avoid tedious manual tuning
- get the best out of a model
- study ML algorithms for new problems

*“Machine Learning
for everyone”*

*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)
```

How it started?



ChaLearn Automatic Machine Learning Challenge (AutoML)

Organized by automl.chalearn - Current server time: June 5, 2023, 7:14 a.m. UTC

Reward \$30,000

First phase

Tweakathon0

Dec. 8, 2014, midnight UTC

End

Competition Ends

June 25, 2016, midnight UTC



Joaquin Vanschoren @joavanschoren · Jun 25, 2016

Congratulations to the **AutoML Challenge** winners! @_mfeurer_@abhijithakur @narnars0 and Lisheng Sun

...



Comment 5 Like 8 Share



PAKDD 2018

AUTOML
Challenge 2018

Paradigm

CHALearn

AutoML 2018 challenge :: PAKDD2018

Organized by hugo.jair - Current server time: June 5, 2023, 7:21 a.m. UTC

Reward \$3,000

First phase

Feedback

Nov. 30, 2017, midnight UTC

End

Competition Ends

March 31, 2020, midnight UTC

Where we are now

automl / **auto-sklearn** Public

Edit Pins ▾ Unwatch 213 Fork 1.2k Starred 7k

< Code Issues 150 Pull requests 5 Discussions Actions Projects 1 Wiki ...

Contributors 80

/automl/auto-sklearn

Used by 425

7

Goals & Outline

Goals

Understand how Auto-sklearn works & Apply Auto-Sklearn

Outline

1. Interactive Discussion (15 mins)
2. Introduction to Auto-sklearn (45 mins)
3. Hands-On
 - **Task 1:** BYOP
 - **Task 2:** ASKL
 - **Task 3:** EXPLAIN

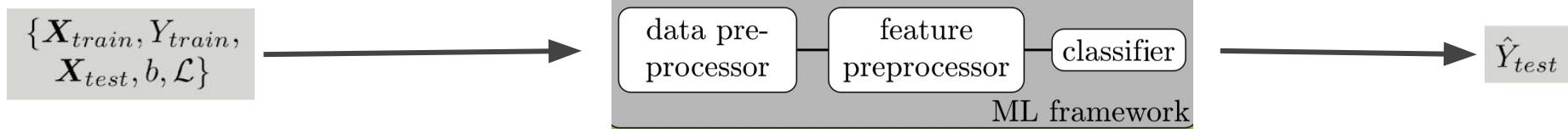
Bonus Tasks - CUSTOM METRICS, MOO, ASKL 2.0, CUSTOM MODEL
4. Outroduction (20 mins)

Now: Your input □

Back to

What is Auto-Sklearn?

Design Space: Traditional ML with scikit-learn



Design Space: Traditional ML with scikit-learn

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



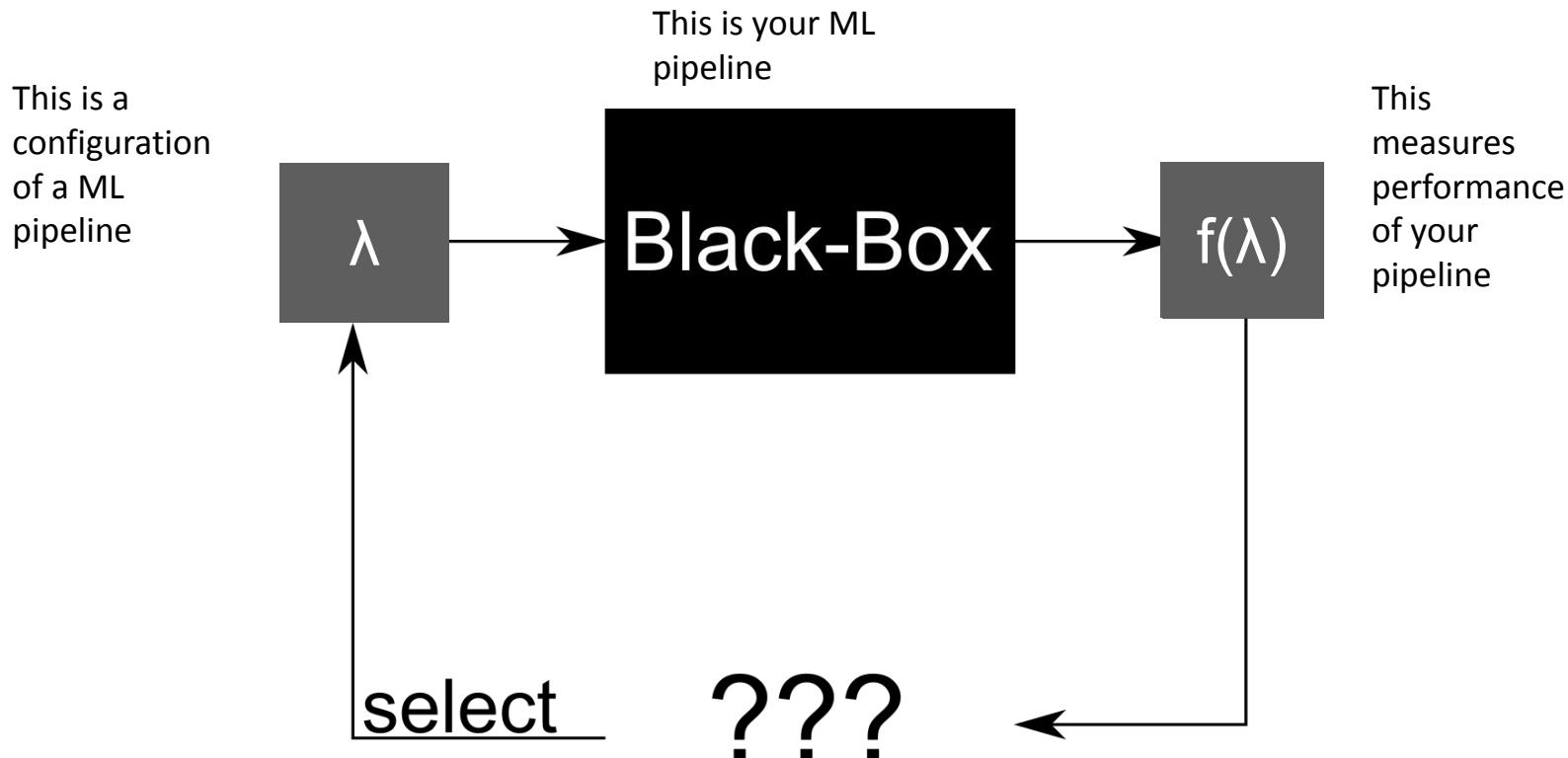
preprocessor	#λ
extreml. rand. trees prepr.	5
fast ICA	4
feature agglomeration	4
kernel PCA	5
rand. kitchen sinks	2
linear SVM prepr.	3
no preprocessing	-
nystroem sampler	5
PCA	2
polynomial	3
random trees embed.	4
select percentile	2
select rates	3
one-hot encoding	2
imputation	1
balancing	1
rescaling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



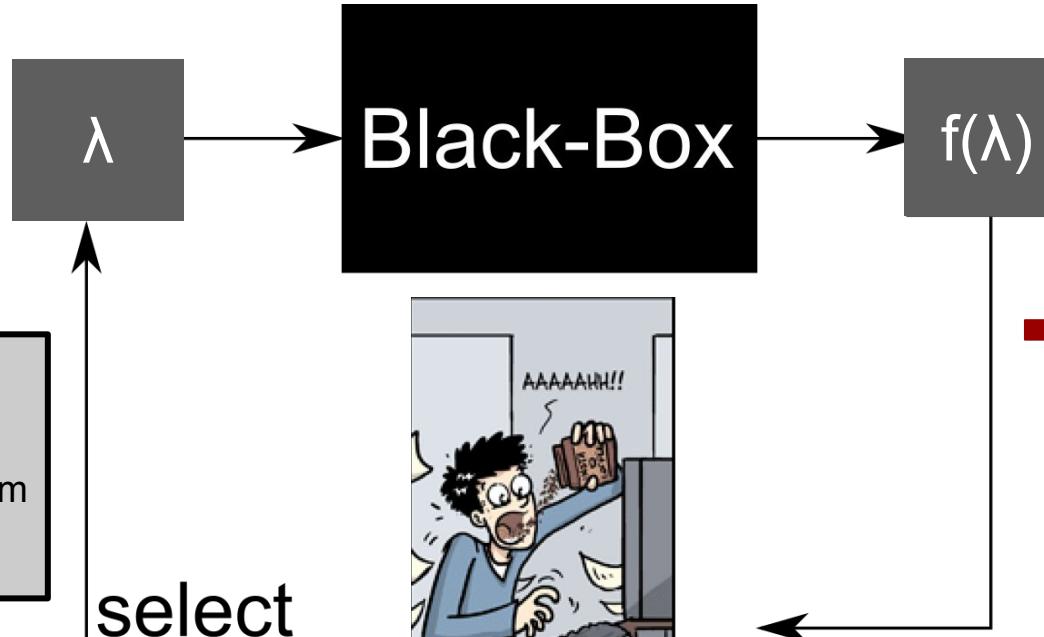
classifier	#λ
AdaBoost (AB)	4
Bernoulli naïve Bayes	2
decision tree (DT)	4
extreml. rand. trees	5
Gaussian naïve Bayes	-
gradient boosting (GB)	6
kNN	3
LDA	4
linear SVM	4
kernel SVM	7
multinomial naive Bayes	2
passive aggressive	3
QDA	2
random forest (RF)	5
Linear Class. (SGD)	10

\hat{Y}_{test}

Black Box Optimization

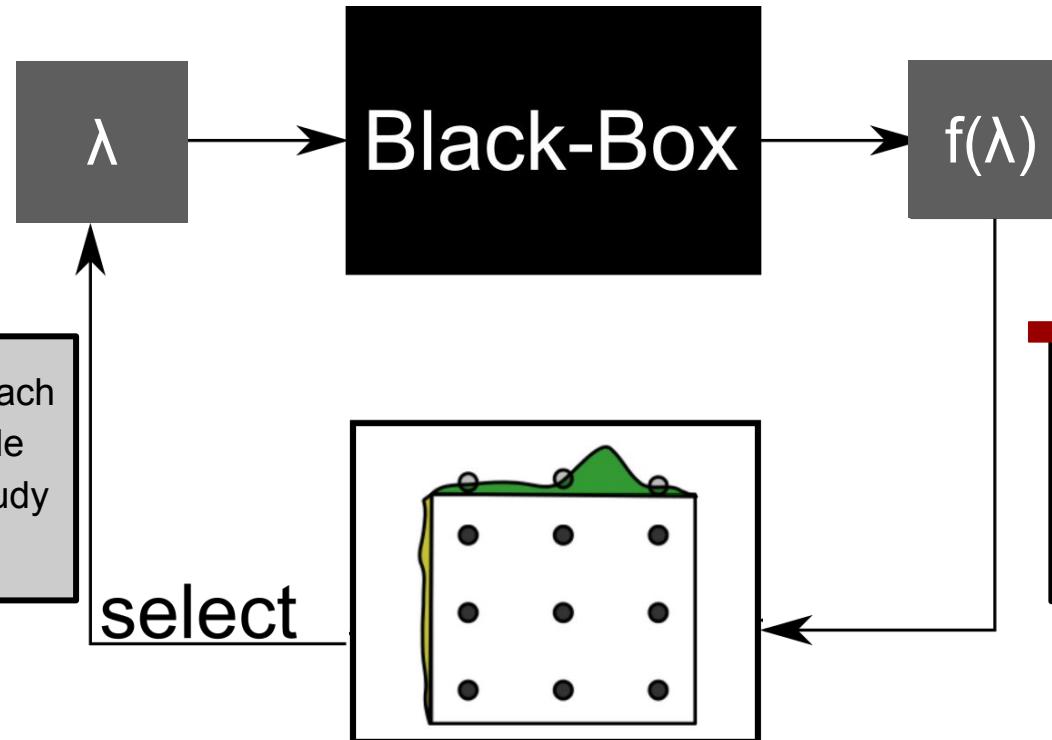


Black Box Optimization: The Human Optimizer



www.PHDCOMICS.COM

Black Box Optimization: Grid Search

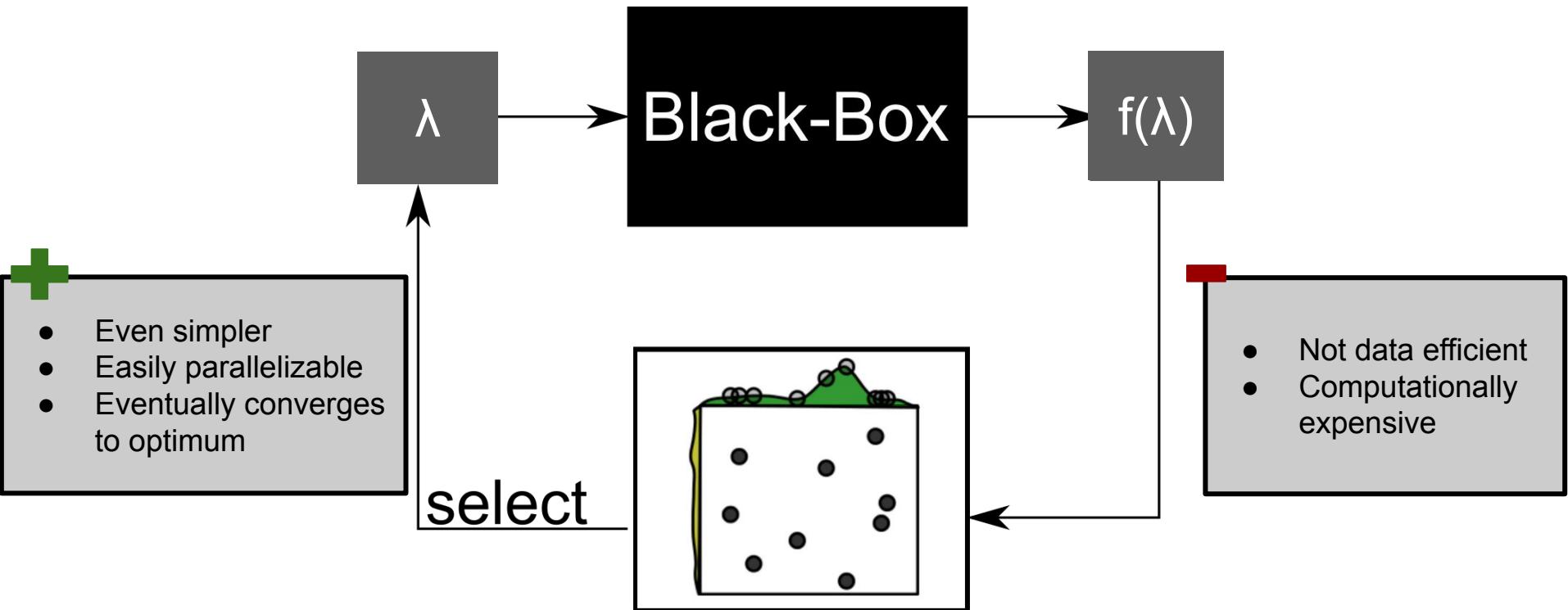


- Very simple approach
- Easily parallelizable
- Can be used to study the problem

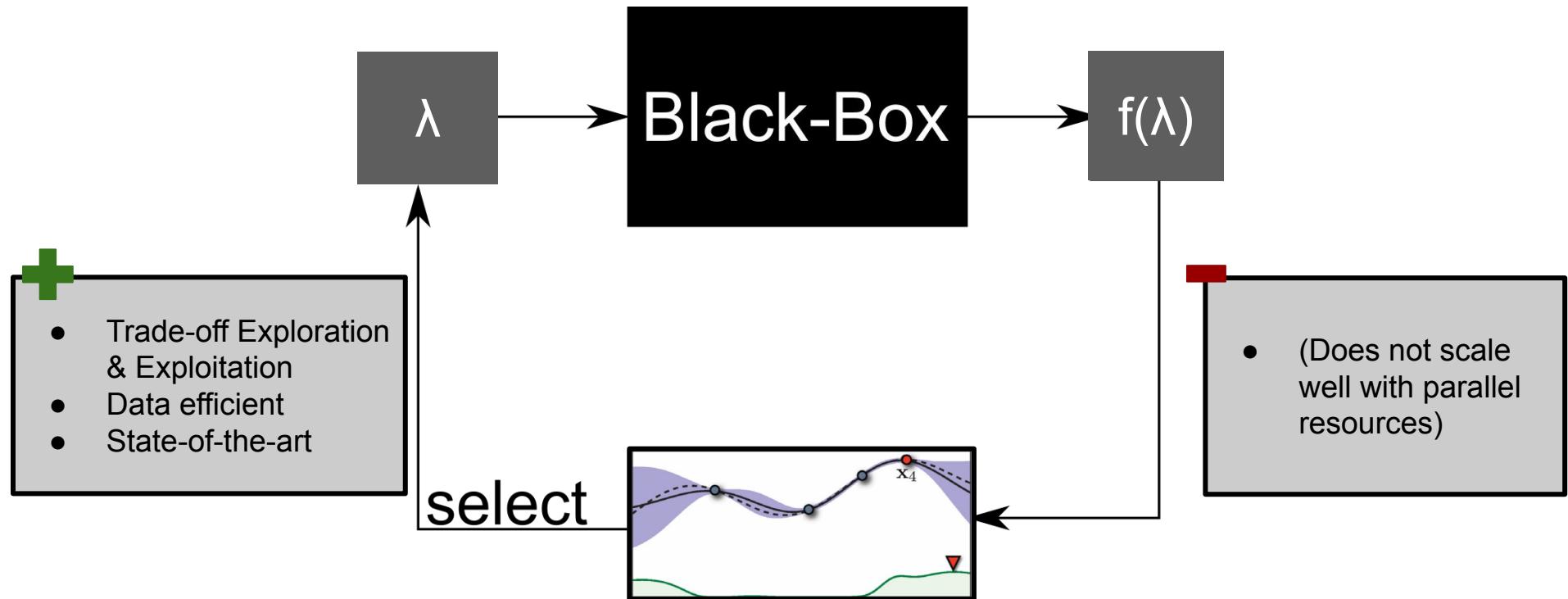
select

- Does not scale to high dimensions
- Grid needs to be defined

Black Box Optimization: Random Search

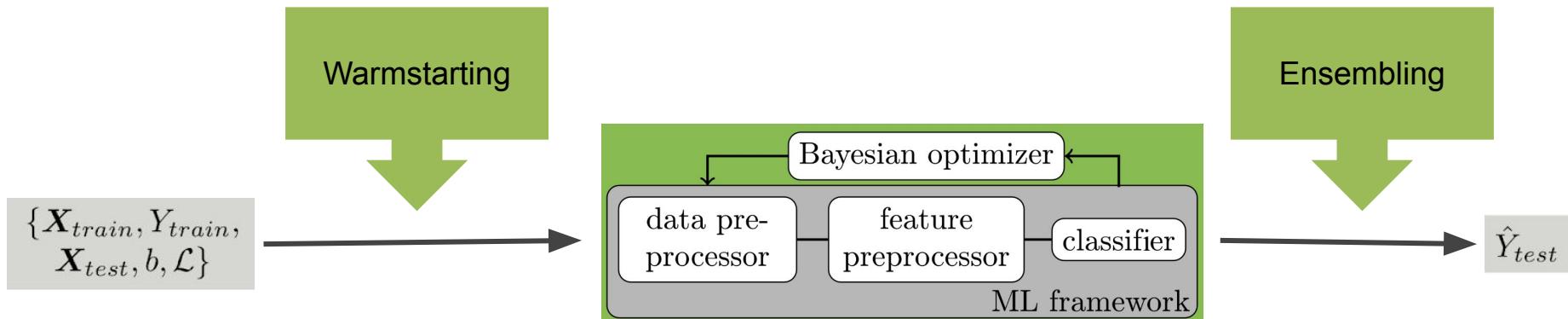


Black Box Optimization: Bayesian Optimization



A tiny intro to Bayesian Optimization:
<https://distill.pub/2020/bayesian-optimization/>

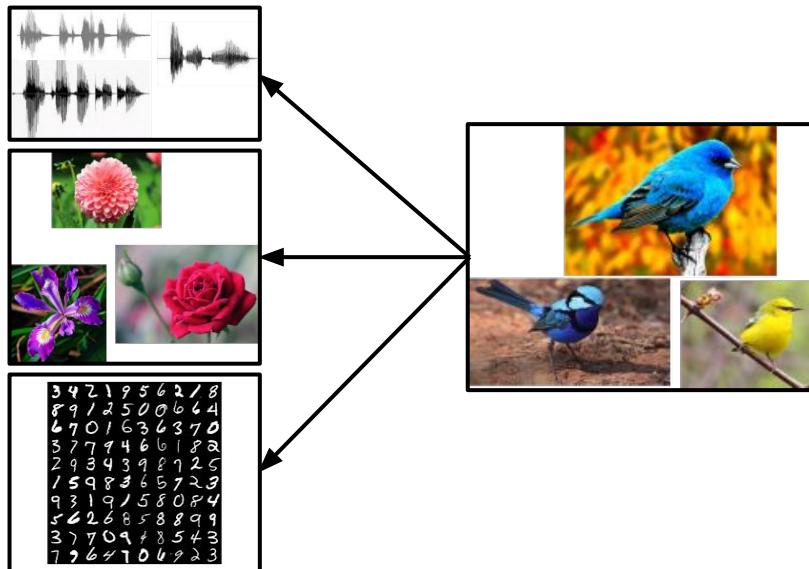
Design Space: Traditional ML with scikit-learn



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

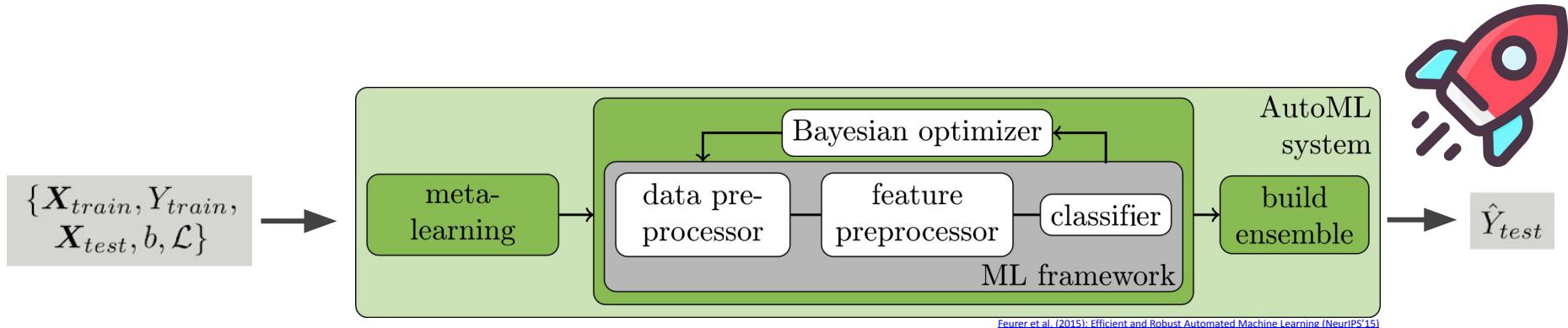
More II: Ensembling

How to get the best out of all evaluated models?

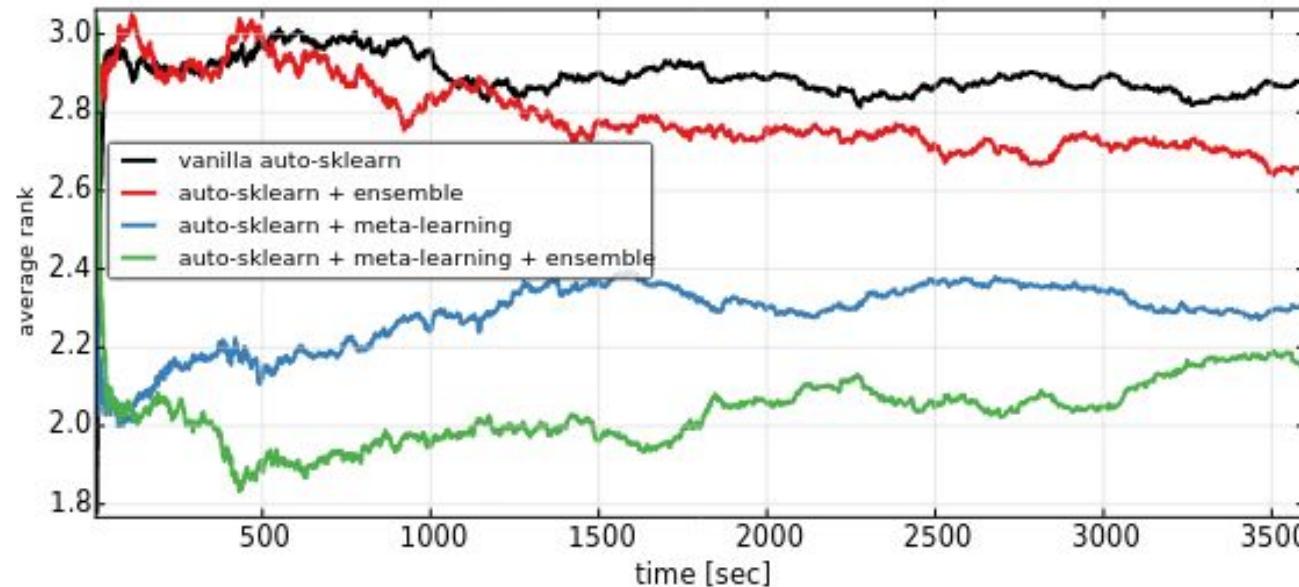
→ **Build an ensemble**



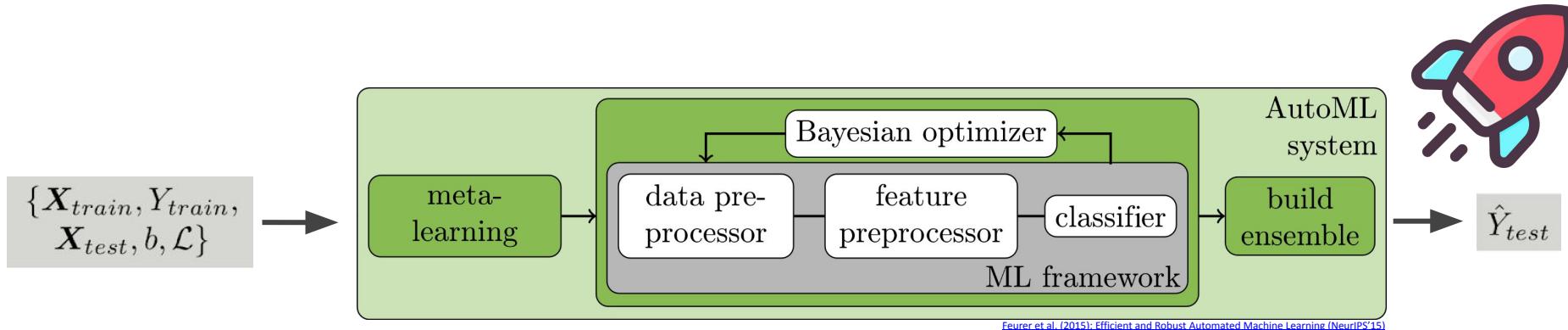
Auto-Sklearn 1.0



Impact of Meta-Learning and Ensembling



Auto-Sklearn 1.0



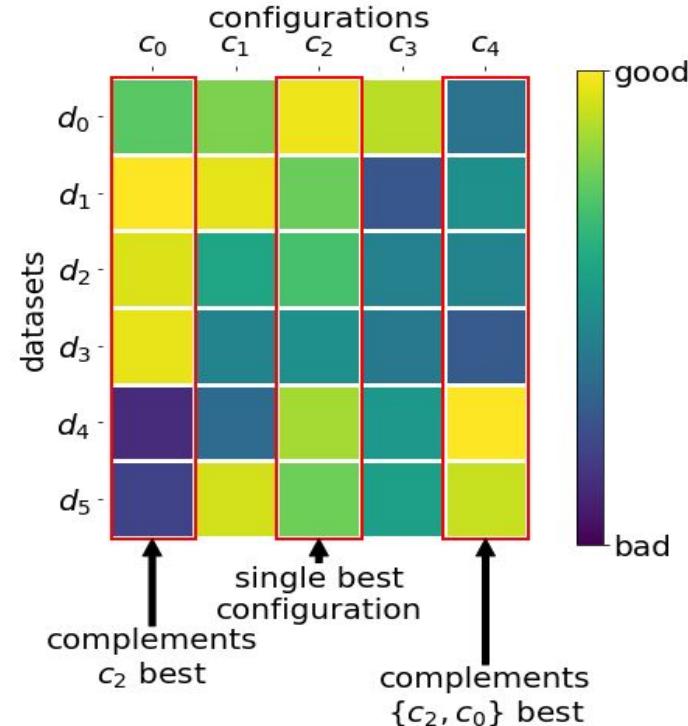
However, some things to be improved

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

Even More I: Portfolios

Goal: Meta-Learning without meta-features

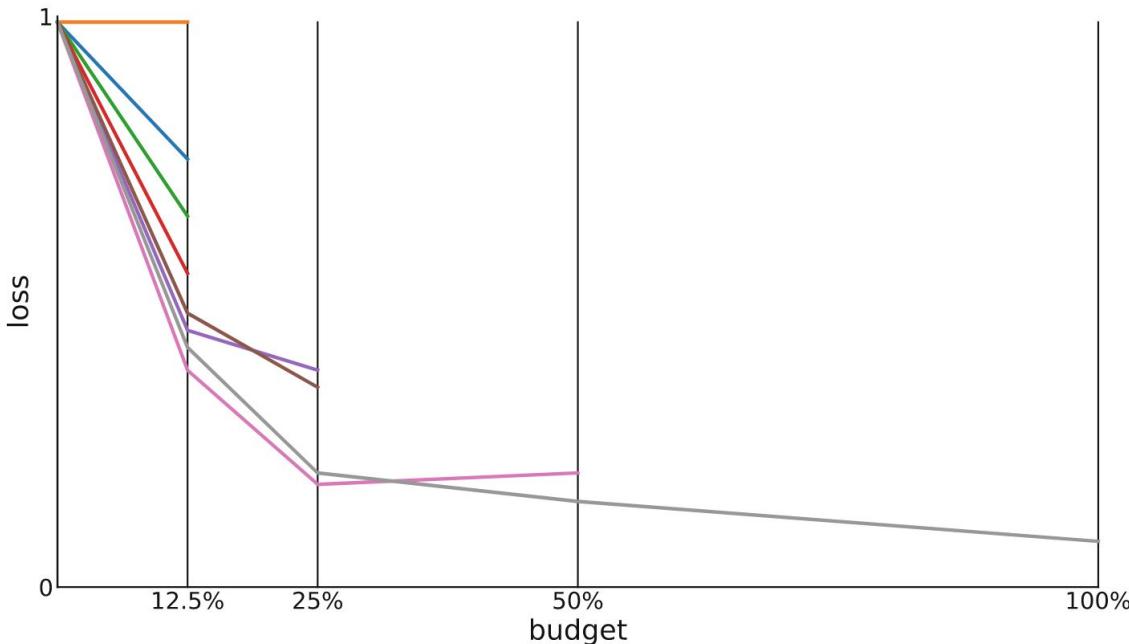
Idea: Construct a Portfolio
(a list of diverse pipelines)



Even More II: Successive Halving

Goal: Scale to large datasets.

Idea: Allocate more resources to promising pipelines



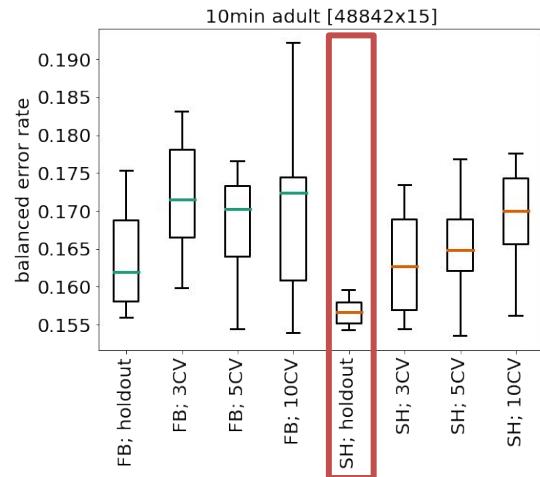
But what about
small datasets?

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

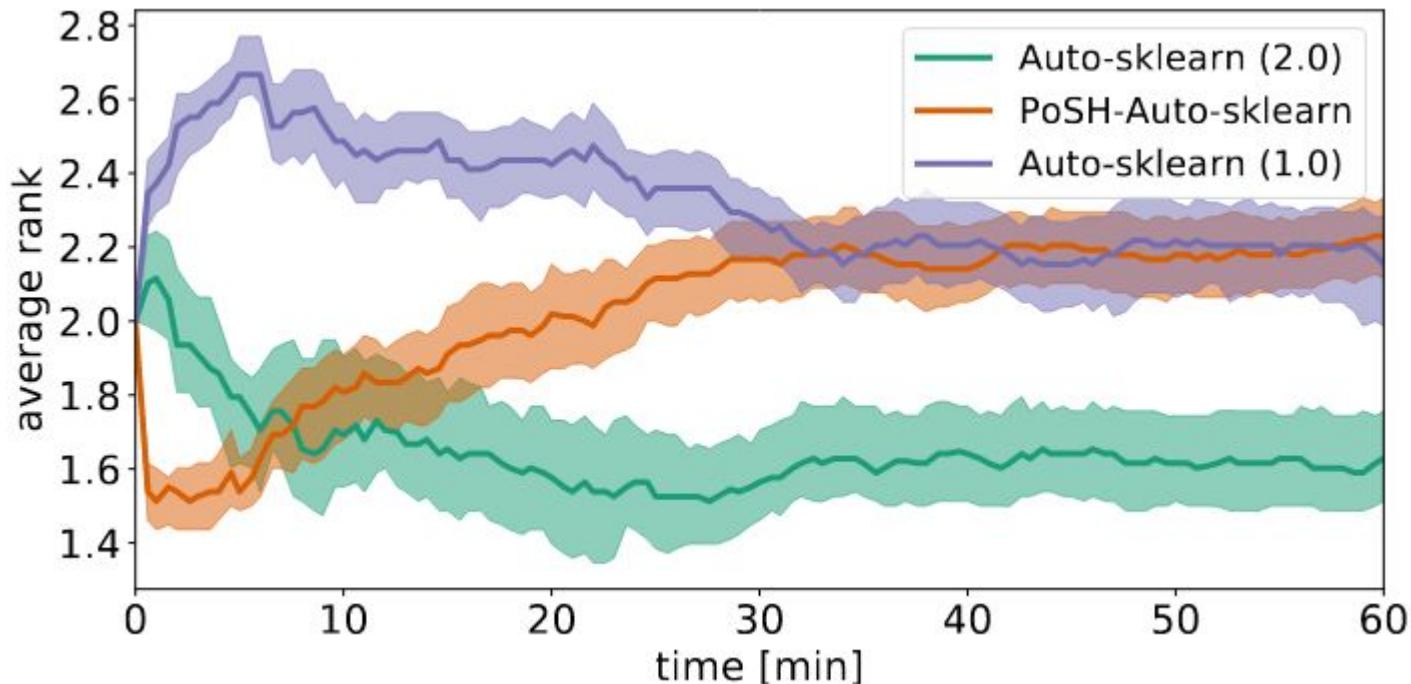


Image Credit - CC BY-NC-ND 2.0; by [Beagle Mama](#)

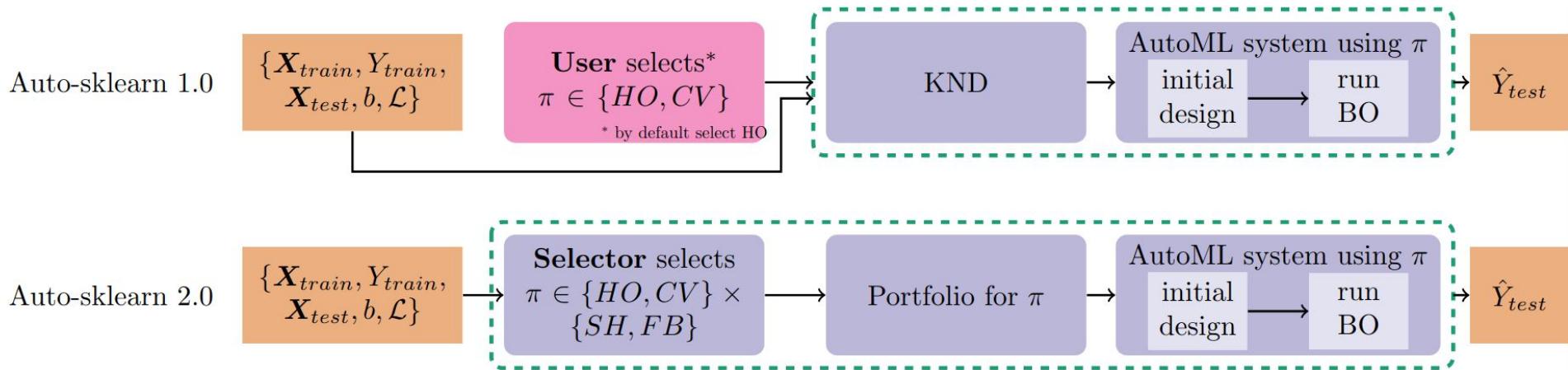
Yes, with a learned selector!

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

Improvement over Auto-sklearn 1.0



Autosklearn 1.0 vs Auto-sklearn 2.0



Practical features

- 💡 Simple → follows scikit-learn API
- ⚙️ Extensible → to your own algorithms
- 🎯 Efficient → parallelization using Dask
- 🔗 Robust → strict resource limitations

... and many more:

- configurable
- 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**



Aron Bahram
Hiwi @ ML Lab Freiburg



/automi/auto-sklearn



Matthias
Feurer
Thomas Bayes Fellow and
Substitute Professor @
MCML / LMU Munich



Katharina
Eggensperger
Group Leader @ Uni
Tübingen



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 tasks

Feedback Time!

What else?

Hyperparameter Tuning

Goal: Find the best performing configuration:

$$\lambda^* \in \arg \min_{\lambda \in \Lambda} f(\mathcal{A}_\lambda)$$

```
from ConfigSpace import ConfigurationSpace  
  
cs = ConfigurationSpace(  
    space={  
        "lr": (0.00001, 1.),      # UniformFloat  
        "num_neurons": (16, 512), # UniformInt  
        "opt": ["Adam", "SGD"],   # Categorical  
    },  
)
```

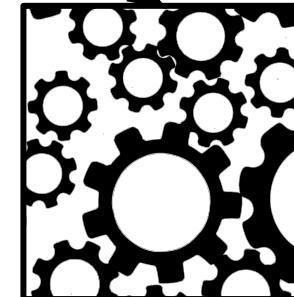
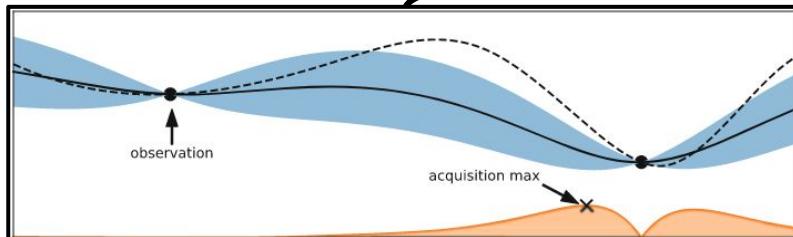
Hyperparameter Tuning

Goal: Find the best performing configuration:

$$\lambda^* \in \arg \min_{\lambda \in \Lambda} f(\mathcal{A}_\lambda)$$

λ_n

Bayesian optimization
(and its advanced variants)



target
algorithm A

$f(\mathcal{A}_{\lambda_n})$



/automl/SMAC3



Neural Architecture Search (NAS)

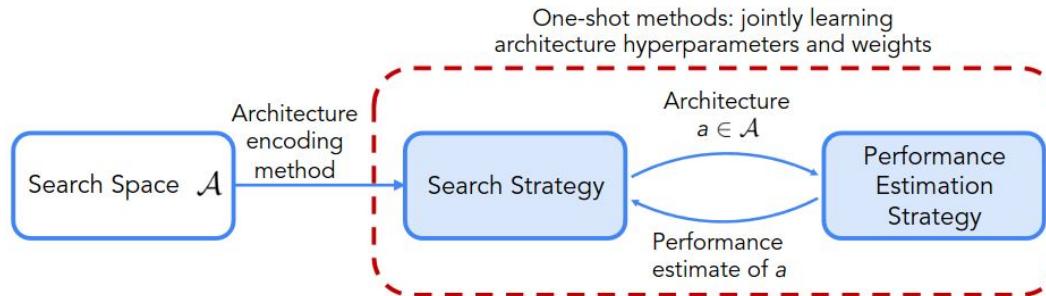


Image from "Colin White et al.: Neural Architecture Search: Insights from 1000 Papers. arXiv:2301.08727v2 [cs.LG]"

Want to build your own NAS optimizer?

 /automl/NASlib

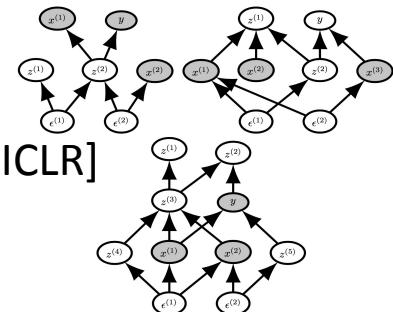
Want a network?

 /automl/AutoPytorch

TabPFNs for small data

PFN = Prior-data Fitted Networks [Müller et al. 2022; ICLR]

TabPFN = PFNs trained with a prior for **tabular** data [Hollman et al. 2023; ICLR]



Done once, offline

Sample synthetic datasets D_i from prior: $D_i \sim p(\mathbf{D})$

Train TabPFN q_θ on synthetic datasets $\{D_1, \dots, D_n\}$

Done per real-world dataset, online

Real-world training dataset D_{real} and test point x_{test}

Obtain $q_\theta(y_{test} | x_{test}, D_{real})$ with a single forward pass



Limitation: 1K samples;
100 features; 10 classes

Other Tools?

	search	searchspace
 AutoGluon	None	None
 AutoKeras	Hill-climbing	Keras
Auto-PyTorch	BO (SMAC)	PyTorch
Auto-sklearn	BO (SMAC)	scikit-learn
 AutoWEKA	BO (SMAC)	WEKA
FLAML	Custom	scikit-learn, XGBoost, CatBoost, LightGBM
 H2O.ai	RS	H2O + XGBoost
HPO tools	BO, Bandit, DE, etc...	everything :)

AutoML DIY - A modular AutoML Toolkit

→ Many AutoML systems exist. They're all awesome!

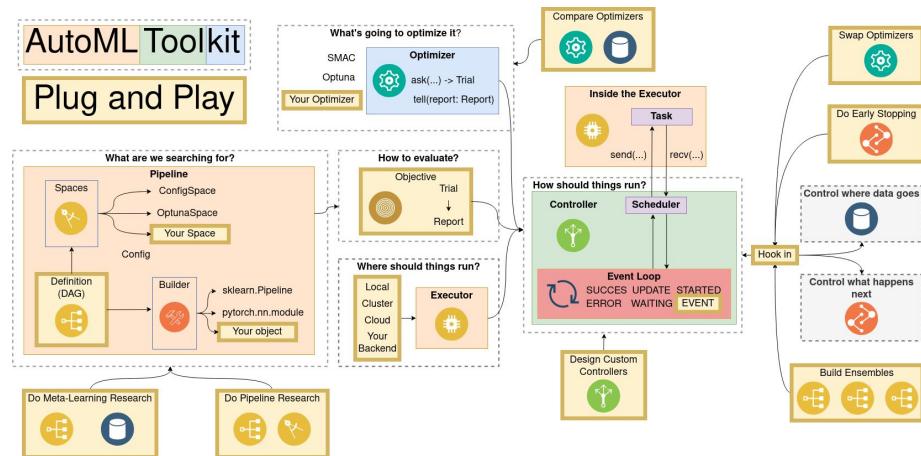
But: Research on AutoML systems is hard.

- No common code framework
- Software-engineering heavy

What we want? A Modular AutoML Framework.
(or “Build-your-own-system-in-a-few-lines-of-code”)

Who is this for? AutoML researcher & ML practitioners

Interested? Talk to us.



Thank you!

