# Using Meta Learning to Initialize Bayesian Optimization

Albert-Ludwigs-Universität Freiburg

### Matthias Feurer<sup>1</sup> Jost Tobias Springenberg<sup>2</sup> Frank Hutter<sup>1</sup>

<sup>1</sup>Research Group on Learning, Optimization, and Automated Algorithm Design <sup>2</sup>Machine Learning Lab

Department of Computer Science, University of Freiburg, Germany

Second Workshop on Configuration and Selection of Algorithms, 30 September 2014







Choose an algorithm based on dataset characteristics, e.g. for the Iris dataset this could be an SVM BURG

ZW



- Choose an algorithm based on dataset characteristics, e.g. for the Iris dataset this could be an SVM
- Manual tuning
   -> fiddling with
   hyperparameters.

COSEAL '14

iBURG



- Choose an algorithm based on dataset characteristics, e.g. for the Iris dataset this could be an SVM
- Manual tuning
   -> fiddling with
   hyperparameters.
- Better: Use automated methods like PSO, GA or SMBO

BUR



- Choose an algorithm based on dataset characteristics, e.g. for the Iris dataset this could be an SVM
- Manual tuning
   -> fiddling with
   hyperparameters.
- Better: Use automated methods like PSO, GA or SMBO
- Best: AutoWeka

BUR









 Manual tuning: Use experience and start from the parameters found on the Iris dataset



- Manual tuning: Use experience and start from the parameters found on the Iris dataset
- Automated methods -> start from scratch



- Manual tuning: Use experience and start from the parameters found on the Iris dataset
- Automated methods -> start from scratch
- $\rightarrow$  Cast *use experience* into an algorithm.

# Sequential Model-based Bayesian Optimization (SMBO)





Dataset D

# Sequential Model-based Bayesian Optimization (SMBO)





# Sequential Model-based Bayesian Optimization (SMBO)



# Metalearning-Initialized SMBO (MI-SMBO)



# Metalearning-Initialized SMBO (MI-SMBO)



UNI FREIBURG

### Metafeatures





- # training examples: 150
- # classes: 3
- # features: 4
- # numerical features: 4
- # categorical features: 0
- missing values? No

For a new dataset *D<sub>new</sub>*:

- Sort known datasets  $D_{1:N}$  by distance to  $D_{new}$ .
- For each of these datasets, extract the best known hyperparameter configuration λ<sup>\*</sup><sub>D</sub>.
- Initialize SMBO with the first k hyperparameter configurations from the sorted list.

# Similarity of Datasets





# Finding the nearest datasets (1)



# Finding the nearest datasets (2)



UNI FREIBURG

# Finding the nearest datasets (3)



# Finding the nearest datasets (3)



# Finding the nearest datasets (4)



### Distance metric (1)



#### Commonly used in literature, the $L_1$ norm:

$$d(D_{\text{new}}, D_j) = \sum_i |m_i^{\text{new}} - m_i^j|$$
(1)



■ 57 datasets from the OpenML repository



- 57 datasets from the OpenML repository
- 46 metafeatures from the literature:
  - Split into five different subsets, including landmarking [Pfahringer et al. 2000]



- 57 datasets from the OpenML repository
- 46 metafeatures from the literature:
  - Split into five different subsets, including landmarking [Pfahringer et al. 2000]
- Two case studies
  - Support Vector Machine with MI-Spearmint [Snoek et al. 2012]
  - AutoSklearn with MI-SMAC [Hutter et al. 2011]



- 57 datasets from the OpenML repository
- 46 metafeatures from the literature:
  - Split into five different subsets, including landmarking [Pfahringer et al. 2000]
- Two case studies
  - Support Vector Machine with MI-Spearmint [Snoek et al. 2012]
  - AutoSklearn with MI-SMAC [Hutter et al. 2011]
- Tried 5, 10, 20 and 25 initial configurations

- 57 datasets from the OpenML repository
- 46 metafeatures from the literature:
  - Split into five different subsets, including landmarking [Pfahringer et al. 2000]
- Two case studies
  - Support Vector Machine with MI-Spearmint [Snoek et al. 2012]
  - AutoSklearn with MI-SMAC [Hutter et al. 2011]
- Tried 5, 10, 20 and 25 initial configurations
- ran each instantiation 10 times on each dataset
  - ightarrow 26220 optimization runs

- 57 datasets from the OpenML repository
- 46 metafeatures from the literature:
  - Split into five different subsets, including landmarking [Pfahringer et al. 2000]
- Two case studies
  - Support Vector Machine with MI-Spearmint [Snoek et al. 2012]
  - AutoSklearn with MI-SMAC [Hutter et al. 2011]
- Tried 5, 10, 20 and 25 initial configurations
- ran each instantiation 10 times on each dataset
  - ightarrow 26220 optimization runs
- therefore, precomputed a dense grid for every dataset





[Auto-WEKA, Thornton et al. 2013]

Component	Hyperparameter	# Values
Main	$\lambda_{classifier}$	3
Main	preprocessing	2
SVM	$\log_2(C)$	21
SVM	$\log_2(\gamma)$	19
LinearSVM	$\log_2(C)$	21
LinearSVM	penalty	2
RF	min splits	5
RF	max features	10
RF	criterion	2
PCA	variance to keep	2

Component	Hyperparameter	# Values
Main	$\lambda_{classifier}$	3
Main	preprocessing	2
SVM	$\log_2(C)$	21
SVM	$\log_2(\gamma)$	19
LinearSVM	$\log_2(C)$	21
LinearSVM	penalty	2
RF	min splits	5
RF	max features	10
RF	criterion	2
PCA	variance to keep	2

#### 1623 hyperparameter configurations



UNI FREIBURG



UNI FREIBURG











#### COSEAL '14

REBURG









- Does MI-SMBO scale to larger configuration spaces?
- What if gridsearch is too expensive?
- Can the metalearning component be added directly into the SMBO procedure?



- SMBO can be substantially improved by providing good initial configurations.
- Metalearning provides a sound framework to find these configurations.
- MI-SMAC improves on state-of-the-art methods on a large configuration space, namely AutoSklearn.



Thank you for your attention.

Further questions: feurerm@cs.uni-freiburg.de

This presentation was partially supported by an *ECCAI Travel Award* and the *ECCAI sponsors*.







