

# Using Meta Learning to Initialize Bayesian Optimization

Albert-Ludwigs-Universität Freiburg



UNI  
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

ECAI-2014 Workshop on Meta-learning & Algorithm Selection, 19 August 2014

# Your task: Build an Iris classification system



The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0.

# Your task: Build an Iris classification system



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

The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0.

# Your task: Build an Iris classification system



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

The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0.

# Your task: Build an Iris classification system



- 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

The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0.

# Your task: Build an Iris classification system



- 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

The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0.

# Adding the Iris Japonica to the dataset



The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0; Bottom right: Iris Japonica is licensed by KENPEI under CC BY-SA 3.0

# Adding the Iris Japonica to the dataset



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

The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0; Bottom right: Iris Japonica is licensed by KENPEI under CC BY-SA 3.0



# Adding the Iris Japonica to the dataset



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

The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris virginica is licensed by C T Johansson under CC BY 3.0; Bottom right: Iris japonica is licensed by KENPEI under CC BY-SA 3.0

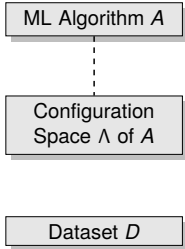
# Adding the Iris Japonica to the dataset



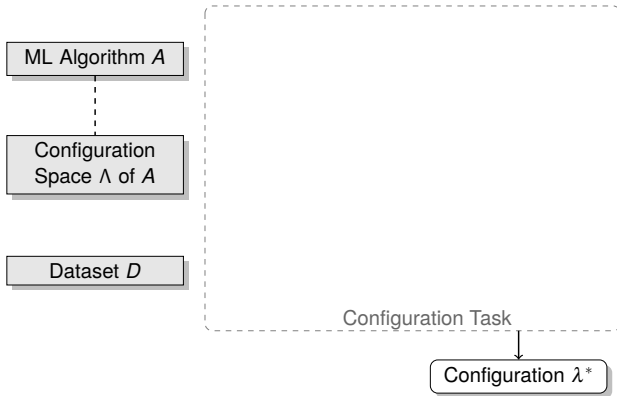
- Manual tuning:  
Use experience and start from the parameters found on the Iris dataset
- Automated methods  
→ start from scratch
- → *Cast use experience* into an algorithm.

The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0; Bottom right: Iris Japonica is licensed by KENPEI under CC BY-SA 3.0

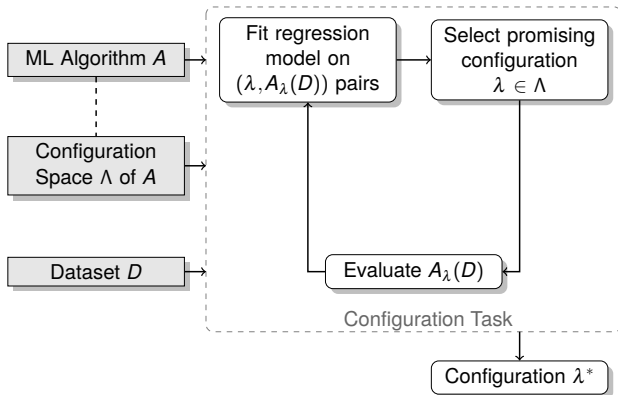
# Sequential Model-based Bayesian Optimization (SMBO)



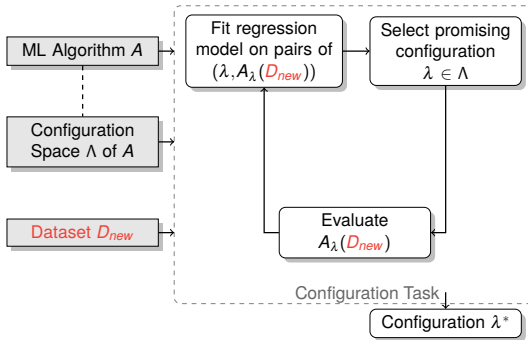
# Sequential Model-based Bayesian Optimization (SMBO)



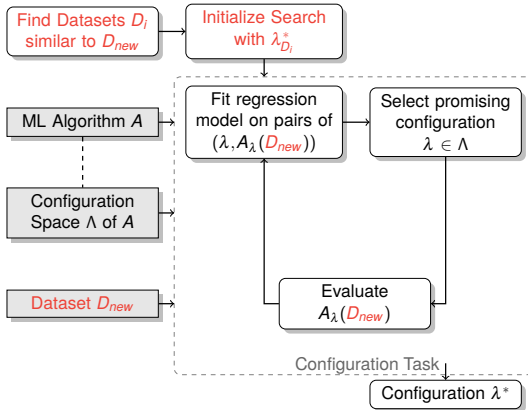
# Sequential Model-based Bayesian Optimization (SMBO)



# Metalearning-Initialized SMBO (MI-SMBO)



# Metalearning-Initialized SMBO (MI-SMBO)





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

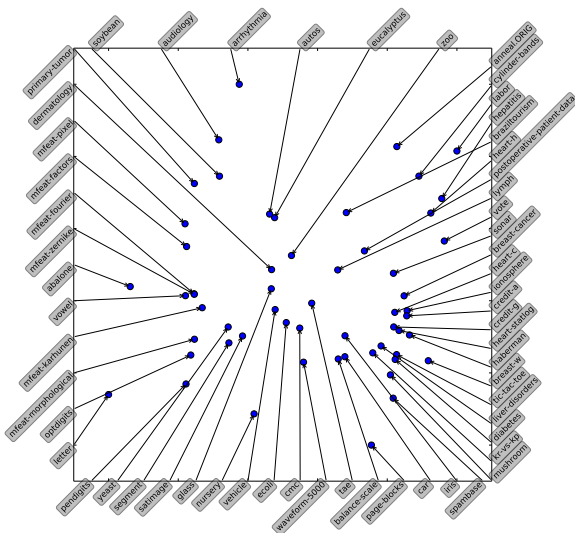
The iris pictures on this slide are from wikimedia commons and used under the following licenses: Top left: Iris Versicolor is public domain; Bottom left: Iris setosa is licensed by Radomil under CC BY-SA 3.0; Top right: Iris Virginica is licensed by C T Johansson under CC BY 3.0.



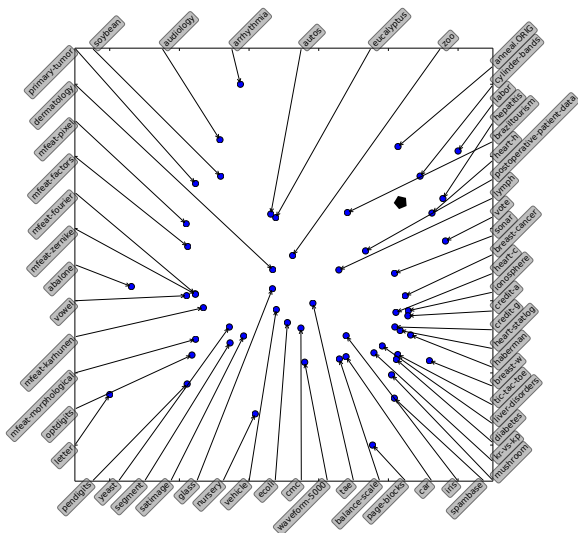
For a new dataset  $D_{new}$ :

- Sort known datasets  $D_{1:N}$  by distance to  $D_{new}$ .
- For each of these datasets, extract the best known hyperparameter configuration  $\lambda_{D_i}^*$ .
- Initialize SMBO with the first  $k$  hyperparameter configurations from the sorted list.

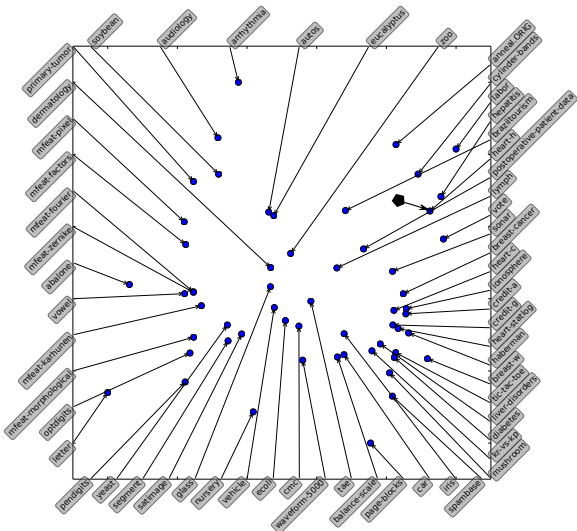
# Similarity of Datasets



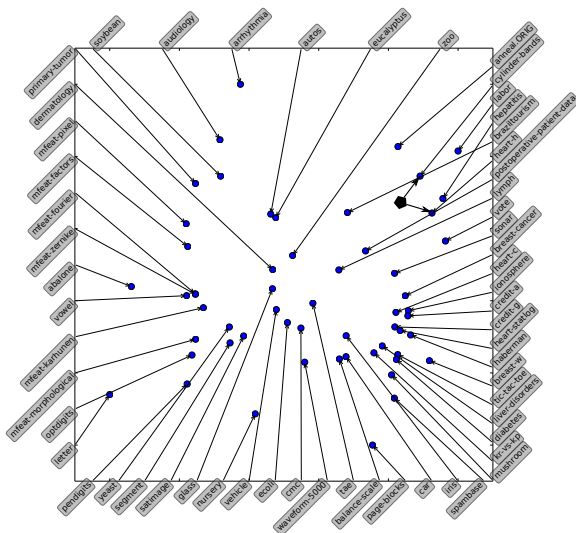
# Finding the nearest datasets (1)



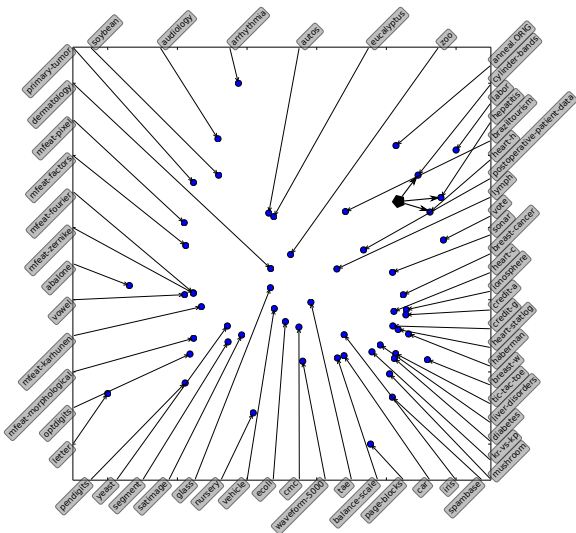
# Finding the nearest datasets (2)



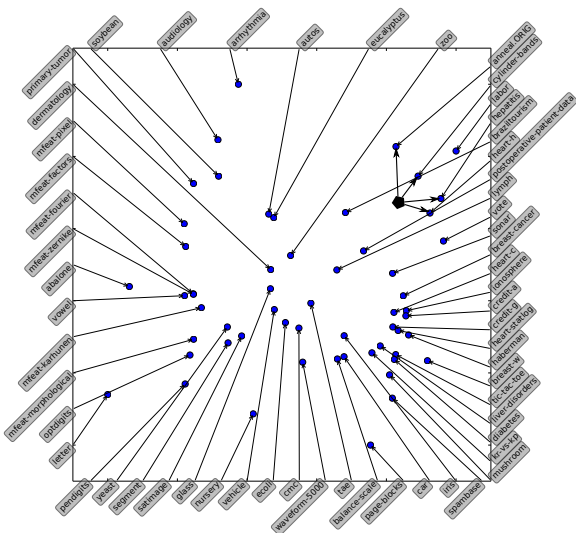
# Finding the nearest datasets (3)



# Finding the nearest datasets (3)



# Finding the nearest datasets (4)



Commonly used in literature, the  $L_1$  norm:

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



- 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]

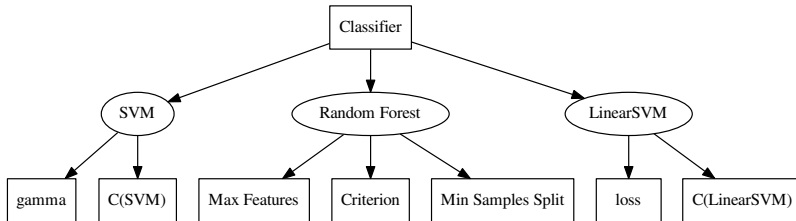
- 57 datasets from the OpenML repository
- 46 metafeatures from the literature:
  - Split into five different subsets, including landmarking [Pfahringer et al. 2000]
- Two case studies
  - Support Vector Machine with MI-Spearmint [Snoek et al. 2012]
  - AutoSklearn with MI-SMAC [Hutter et al. 2011]

- 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

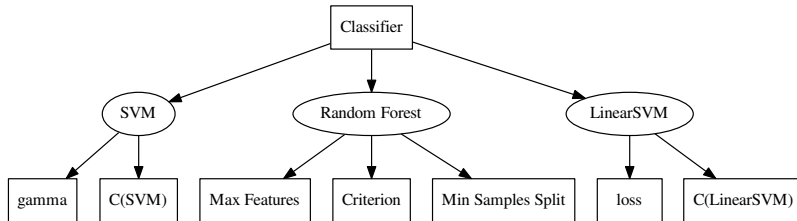
- 57 datasets from the OpenML repository
- 46 metafeatures from the literature:
  - Split into five different subsets, including landmarking [Pfahringer et al. 2000]
- Two case studies
  - Support Vector Machine with MI-Spearmint [Snoek et al. 2012]
  - AutoSklearn with MI-SMAC [Hutter et al. 2011]
-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

# Combined Algorithm Selection and Hyperparameter Optimization problem (CASH)



# Combined Algorithm Selection and Hyperparameter Optimization problem (CASH)



[Auto-WEKA, Thornton et al. 2013]



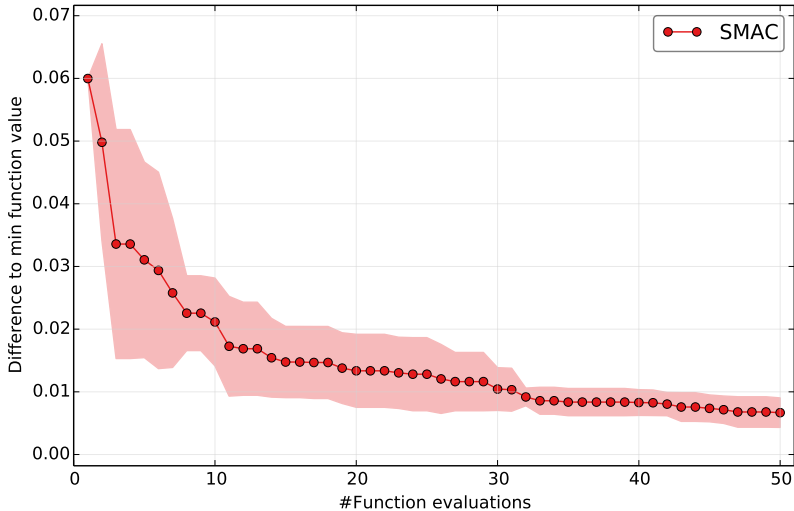
| Component | Hyperparameter                | # Values |
|-----------|-------------------------------|----------|
| Main      | $\lambda_{\text{classifier}}$ | 3        |
| Main      | preprocessing                 | 2        |
| SVM       | $\log_2(C)$                   | 21       |
| SVM       | $\log_2(\gamma)$              | 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      | $\lambda_{\text{classifier}}$ | 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        |

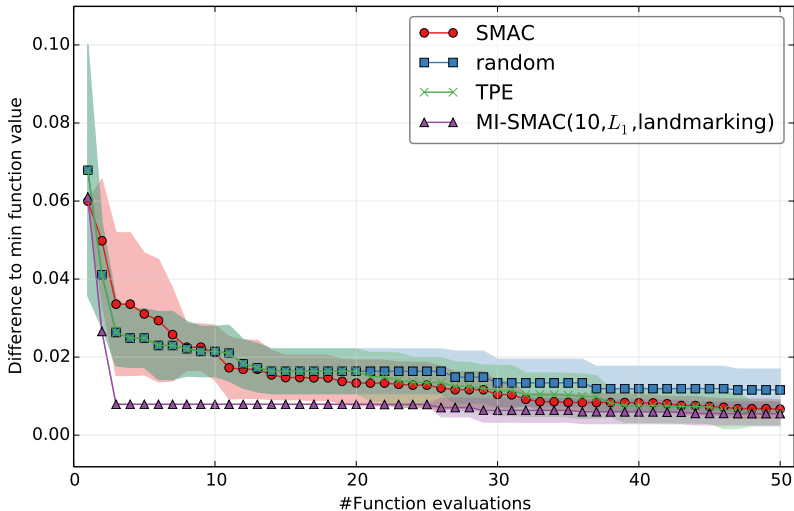
---

1623 hyperparameter configurations

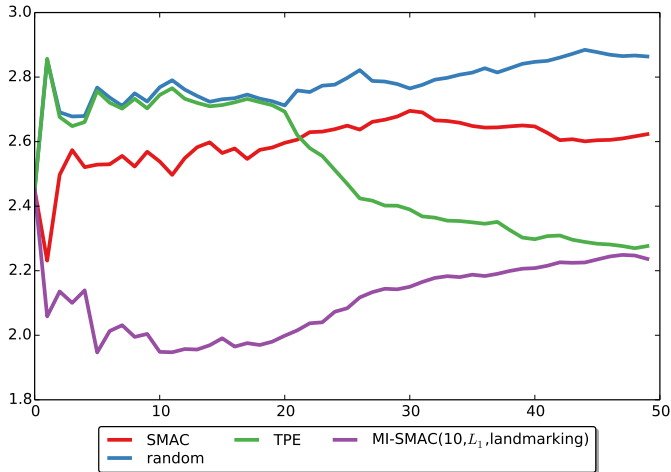
# AutoSklearn: Results (1)



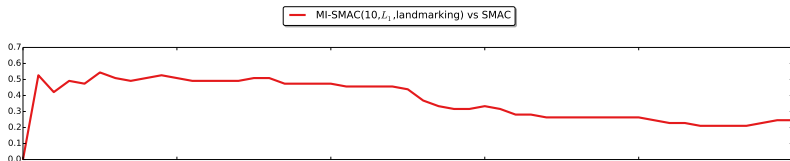
# AutoSklearn: Results (1)



# AutoSklearn: Results (2)



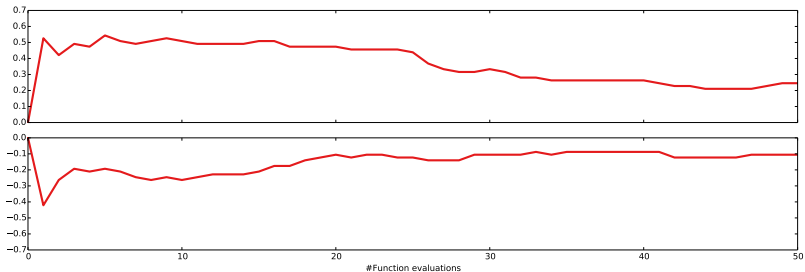
# AutoSklearn: Results (3)



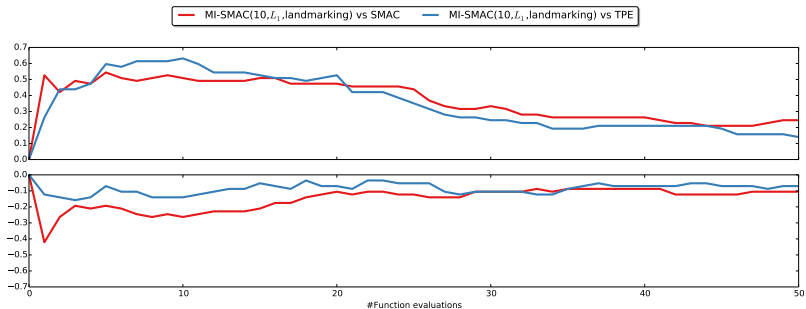
# AutoSklearn: Results (3)



MI-SMAC(10,  $I_1$ , landmarking) vs SMAC

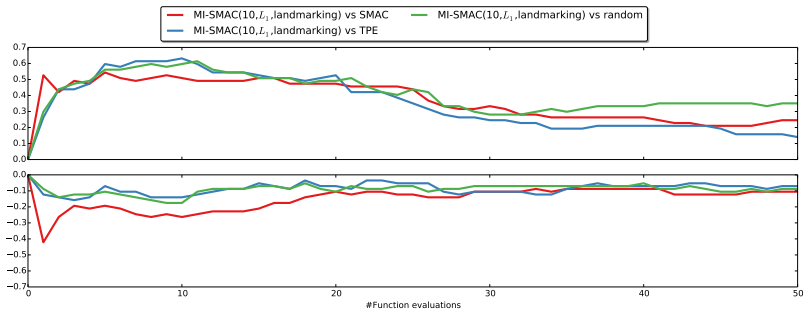


# AutoSklearn: Results (3)





# AutoSklearn: Results (3)



- 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.

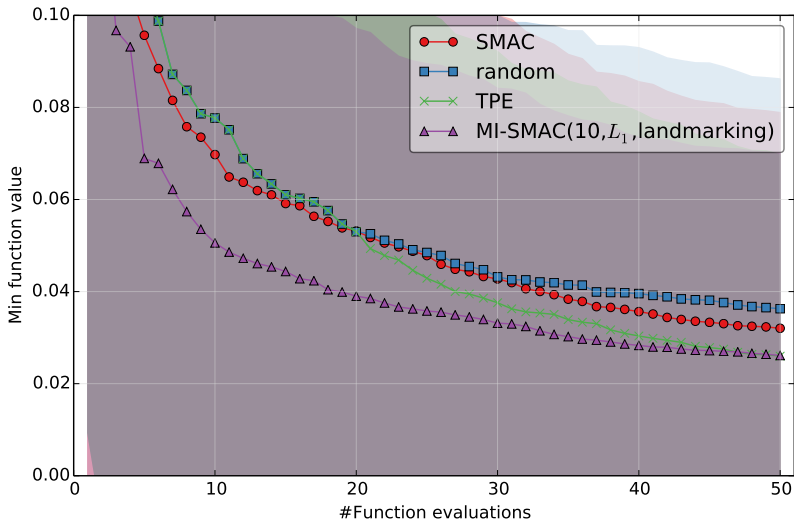
# The end

Thank you for your attention.

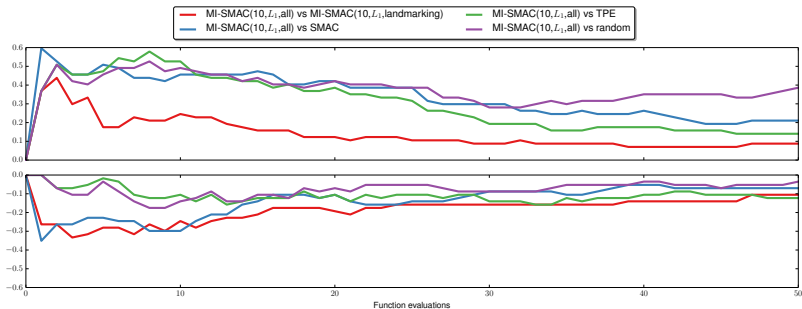
Further questions: `feurerm@cs.uni-freiburg.de`

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

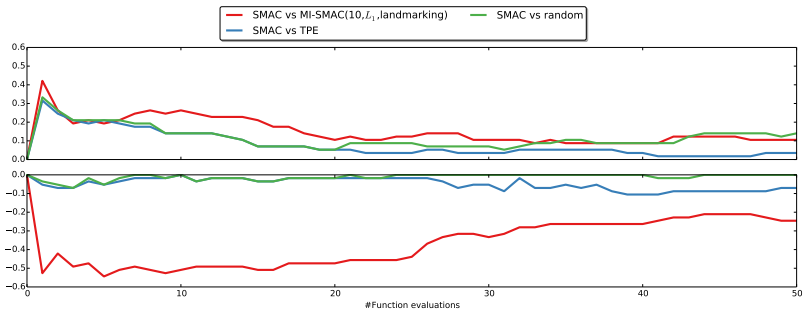
# AutoSklearn: Results (5)



# AutoSklearn: Results (7)



# AutoSklearn: Results (8)



# AutoSklearn: Results (9)

