

# MODEL-BASED OPTIMIZATION FOR EXPENSIVE BLACK-BOX PROBLEMS AND HYPERPARAMETER OPTIMIZATION

Bernd Bischl

Computational Statistics, LMU Munich

Dec 1st, 2017

SEQUENTIAL MODEL-BASED OPTIMIZATION

PARALLEL BATCH PROPOSALS

MULTICRITERIA SMBO

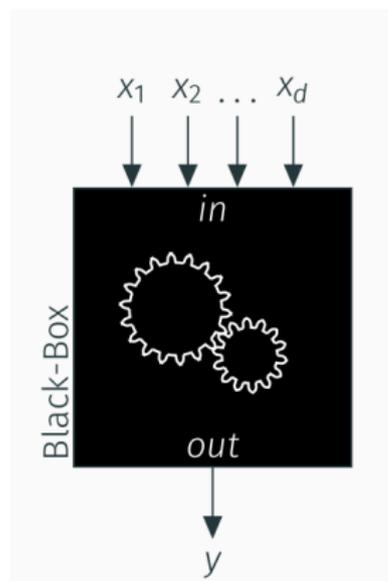
INTERESTING CHALLENGES

ML MODEL SELECTION AND HYPERPARAMETER  
OPTIMIZATION

## Section 1

# SEQUENTIAL MODEL-BASED OPTIMIZATION

# EXPENSIVE BLACK-BOX OPTIMIZATION



$$y = f(\mathbf{x}), \quad f : \mathbb{X} \rightarrow \mathbb{R} \quad (1)$$

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{X}} f(\mathbf{x}) \quad (2)$$

- ▶  $y$ , target value
- ▶  $\mathbf{x} \in \mathbb{X} \subset \mathbb{R}^d$ , domain
- ▶  $f(\mathbf{x})$  function with considerably long runtime
- ▶ Goal: Find optimum  $\mathbf{x}^*$

# SEQUENTIAL MODEL-BASED OPTIMIZATION

- ▶ Setting: Expensive black-box problem  $f : x \rightarrow \mathbb{R} = \min!$
- ▶ Classical problem: Computer simulation with a bunch of control parameters and performance output; or algorithmic performance on 1 or more problem instances; we often optimize ML pipelines
- ▶ Idea: Let's approximate  $f$  via regression!

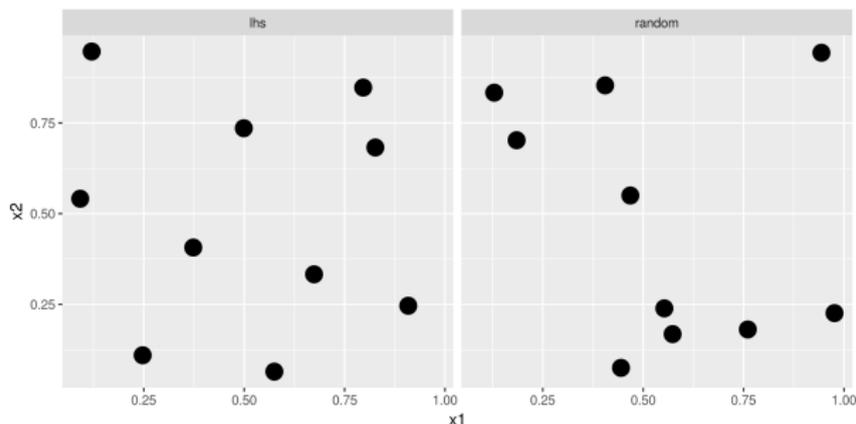
## GENERIC MBO PSEUDO CODE

- ▶ Create initial space filling design and evaluate with  $f$
- ▶ In each iteration:
  - ▶ Fit regression model on all evaluated points to predict  $\hat{f}(x)$  and uncertainty  $\hat{s}(x)$
  - ▶ Propose point via infill criterion

$$\text{EI}(x) \uparrow \iff \hat{f}(x) \downarrow \wedge \hat{s}(x) \uparrow$$

- ▶ Evaluate proposed point and add to design
- ▶ EGO proposes kriging (aka Gaussian Process) and EI  
*Jones 1998, Efficient Global Opt. of Exp. Black-Box Functions*

# LATIN HYPERCUBE DESIGNS

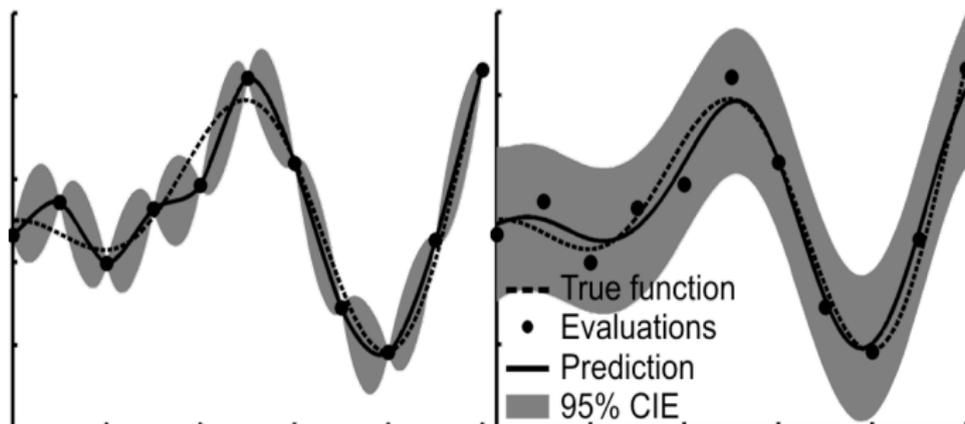


- ▶ Initial design to train first regression model
- ▶ Not too small, not too large
- ▶ LHS / maximin designs: Min dist between points is maximized
- ▶ But: Type of design usually has not the largest effect on MBO, and unequal distances between points could even be beneficial

# KRIGING AND LOCAL UNCERTAINTY PREDICTION

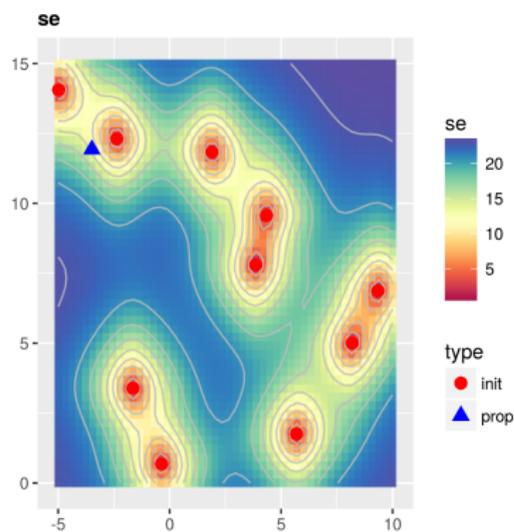
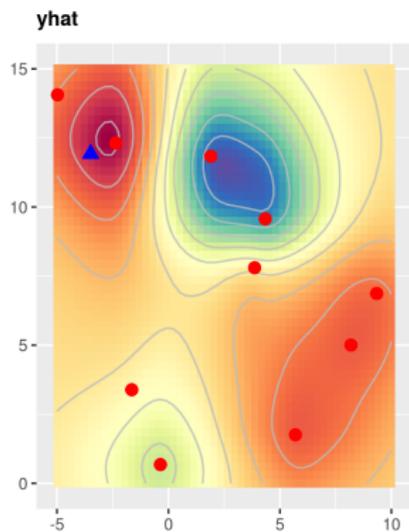
Model: Zero-mean GP  $Y(\mathbf{x})$  with const. trend and cov. kernel  $k_\theta(\mathbf{x}_1, \mathbf{x}_2)$ .

- ▶  $\mathbf{y} = (y_1, \dots, y_n)^T$ ,  $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1,\dots,n}$
- ▶  $\mathbf{k}_*(\mathbf{x}) = (k(\mathbf{x}_1, \mathbf{x}), \dots, k(\mathbf{x}_n, \mathbf{x}))^T$
- ▶  $\hat{\mu} = \mathbf{1}^T \mathbf{K}^{-1} \mathbf{y} / \mathbf{1}^T \mathbf{K}^{-1} \mathbf{1}$  (BLUE)
- ▶ Prediction:  $\hat{f}(\mathbf{x}) = E[Y(\mathbf{x}) | Y(\mathbf{x}_i) = y_i, i = 1, \dots, n] = \hat{\mu} + \mathbf{k}_n(\mathbf{x})^T \mathbf{K}^{-1} (\mathbf{y} - \hat{\mu} \mathbf{1})$
- ▶ Uncertainty:  $\hat{s}^2(\mathbf{x}) = \text{Var}[Y(\mathbf{x}) | Y(\mathbf{x}_i) = y_i, i = 1, \dots, n] = \sigma^2 - \mathbf{k}_n^T(\mathbf{x}) \mathbf{K}^{-1} \mathbf{k}_n(\mathbf{x}) + \frac{(1 - \mathbf{1}^T \mathbf{K}^{-1} \mathbf{k}_n^T(\mathbf{x}))^2}{\mathbf{1}^T \mathbf{K}^{-1} \mathbf{1}}$



# KRIGING / GP IS A SPATIAL MODEL

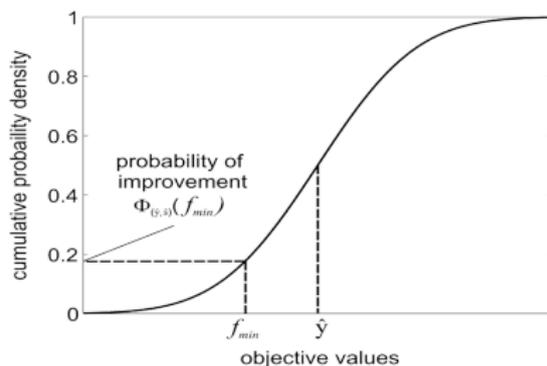
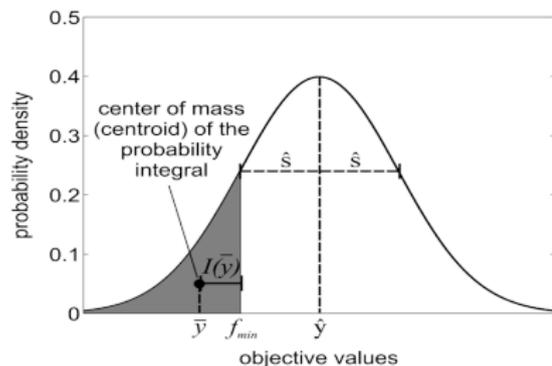
- ▶ Correlation between outcomes  $(y_1, y_2)$  depends on dist of  $x_1, x_2$   
E.g. Gaussian covar kernel  $k(x_1, x_2) = \exp\left(\frac{-\|x_1 - x_2\|}{2\sigma}\right)$
- ▶ Useful smoothness assumption for optimization
- ▶ Posterior uncertainty at new  $x$  increases with dist to design points
- ▶ Allows to enforce exploration



# INFILL CRITERIA: EXPECTED IMPROVEMENT

- ▶ Define improvement at  $x$  over best visited point with  $y = f_{min}$  as random variable  $I(x) = |f_{min} - Y(x)|^+$
- ▶ For kriging  $Y(x) \sim N(\hat{f}(\mathbf{x}), \hat{s}^2(\mathbf{x}))$  (given  $x = \mathbf{x}$ )
- ▶ Now define  $EI(x) = E[I(x)|x = \mathbf{x}]$
- ▶ Expectation is integral over normal density starting at  $f_{min}$
- ▶ Alternative: Lower confidence bound (LCB)  $\hat{f}(\mathbf{x}) - \lambda \hat{s}(\mathbf{x})$

$$\text{Result: } EI(x) = \left( f_{min} - \hat{f}(\mathbf{x}) \right) \Phi \left( \frac{f_{min} - \hat{f}(\mathbf{x})}{\hat{s}(\mathbf{x})} \right) + \hat{s}(\mathbf{x}) \phi \left( \frac{f_{min} - \hat{f}(\mathbf{x})}{\hat{s}(\mathbf{x})} \right)$$

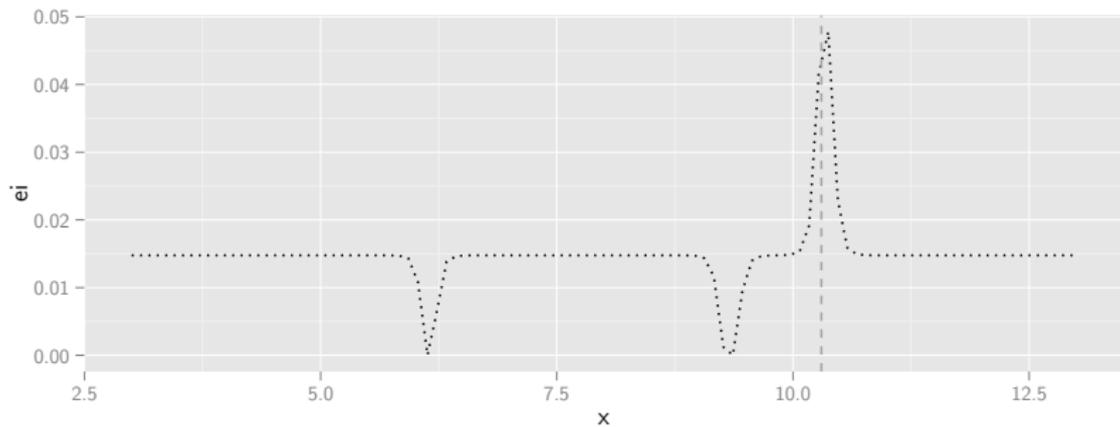
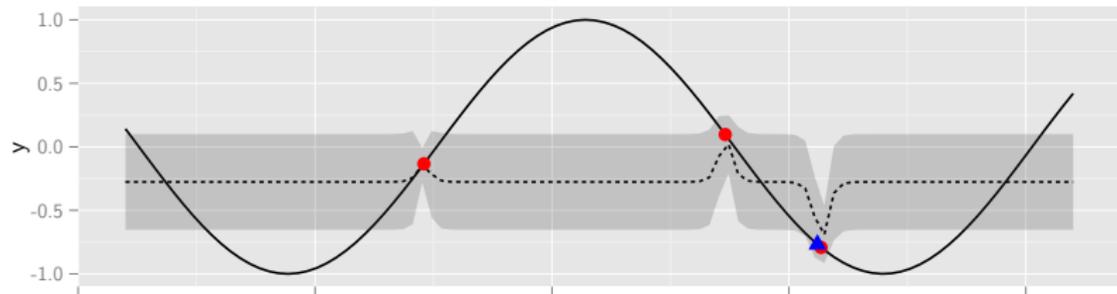


# FOCUSSEARCH

- ▶ El optimization is multimodal and not that simple
- ▶ But objective is now cheap to evaluate
- ▶ Many different algorithms exist, from gradient-based methods with restarts to evolutionary algorithms
- ▶ We use an iterated, focusing random search coined “focus search”
- ▶ In each iteration a random search is performed
- ▶ We then shrink the constraints of the feasible region towards the best point in the current iteration (focusing) and iterate, to enforce local convergence
- ▶ Whole process is restarted a few times
- ▶ Works also for categorical and hierarchical params

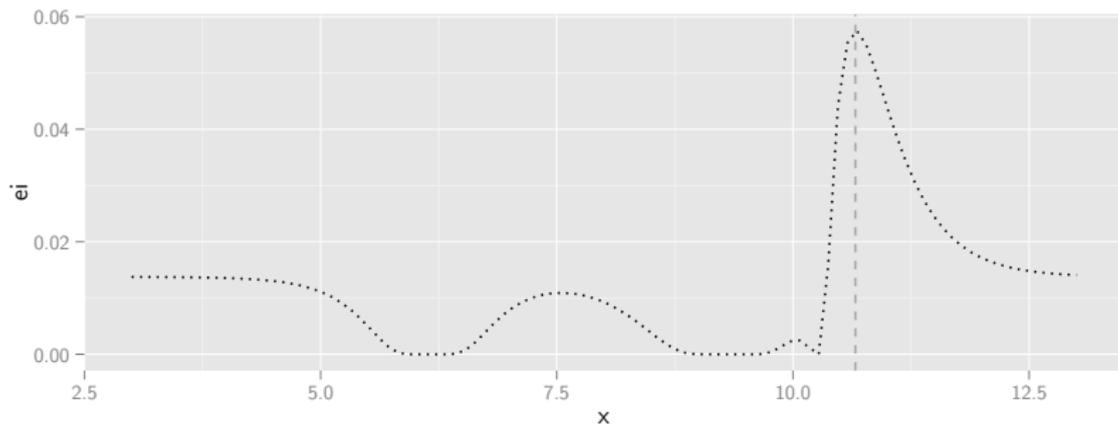
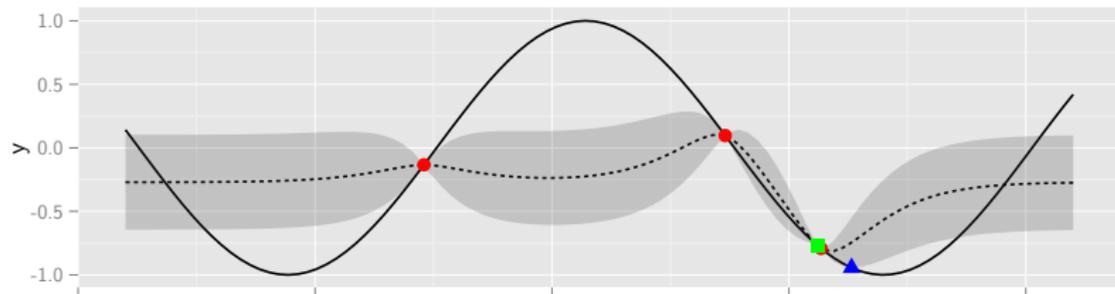
Iter = 1, Gap = 2.0795e-01

type — y - - - yhat type ● init ▲ prop



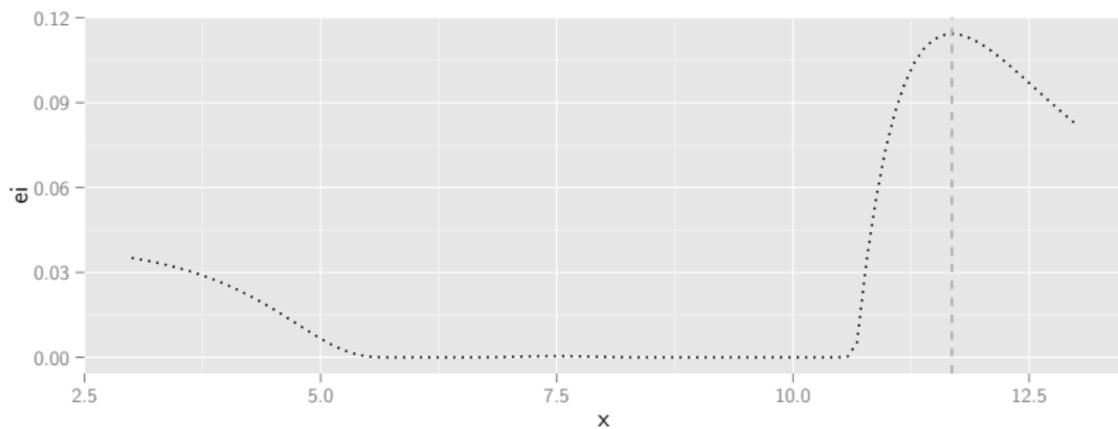
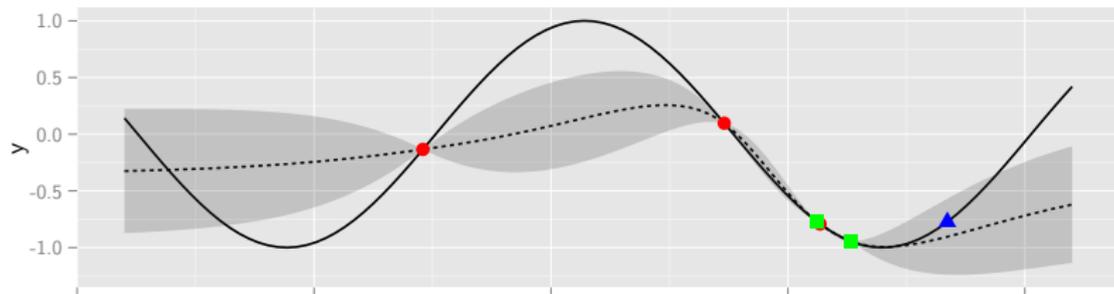
Iter = 2, Gap = 5.5410e-02

type — y ···· yhat    type ● init ▲ prop ■ seq



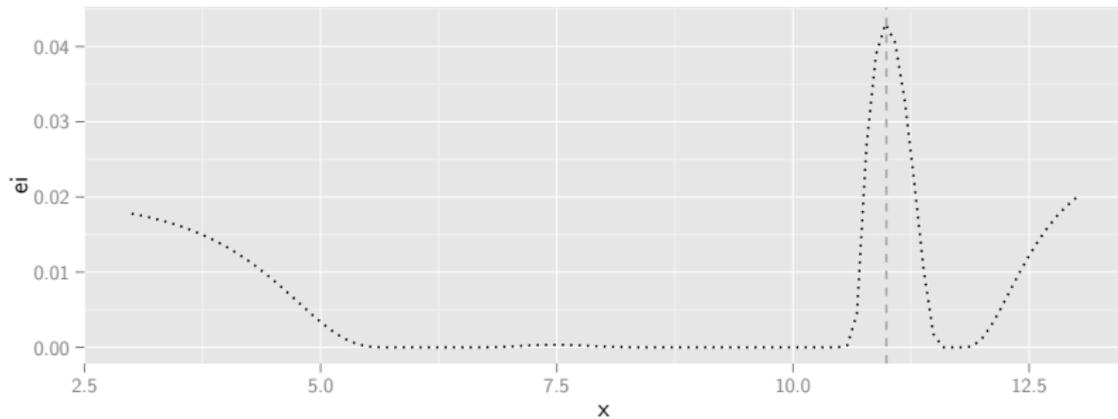
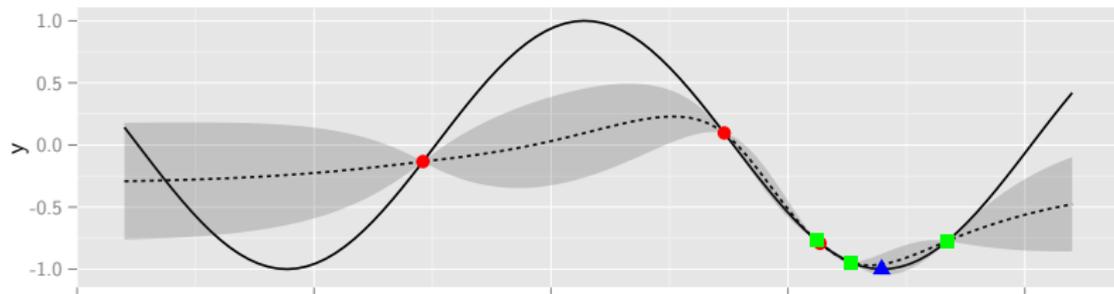
Iter = 3, Gap = 5.5410e-02

type — y ···· yhat    type ● init ▲ prop ■ seq



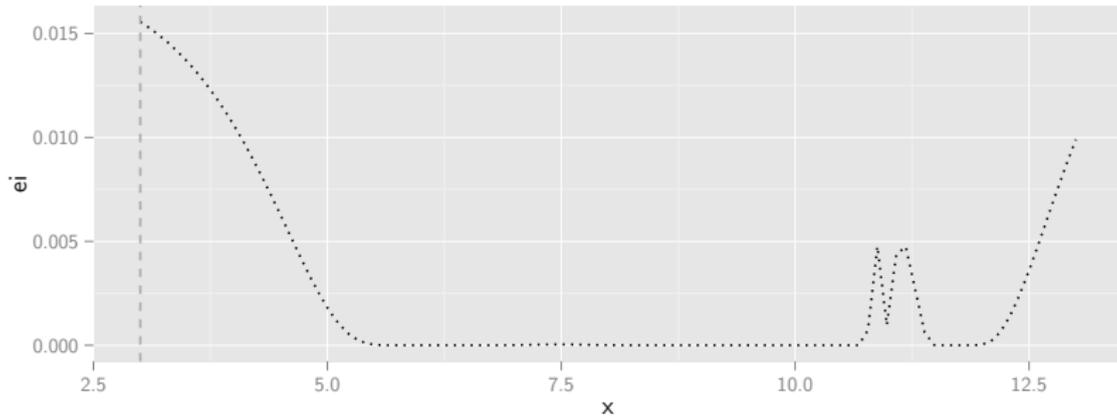
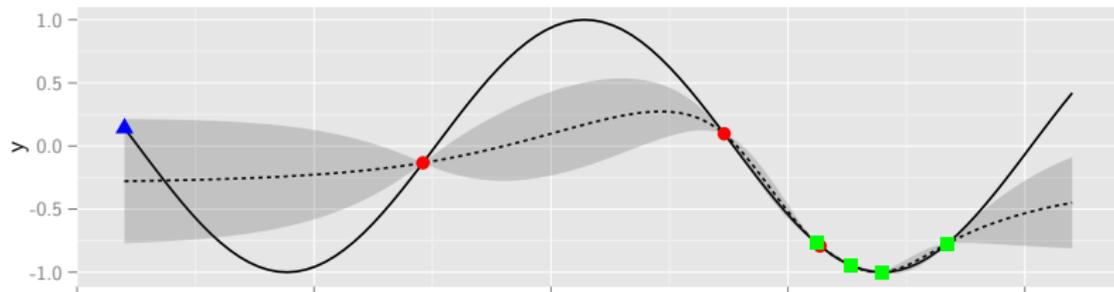
Iter = 4, Gap = 2.2202e-05

type — y ···· yhat type ● init ▲ prop ■ seq



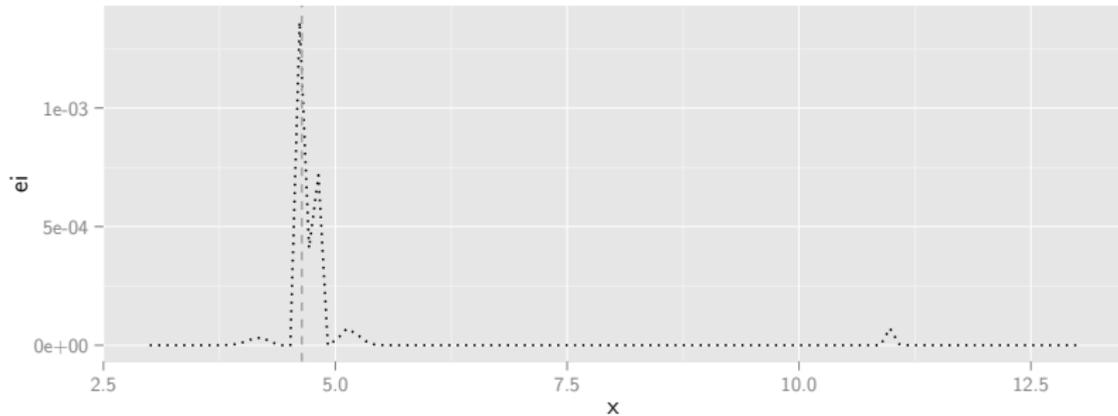
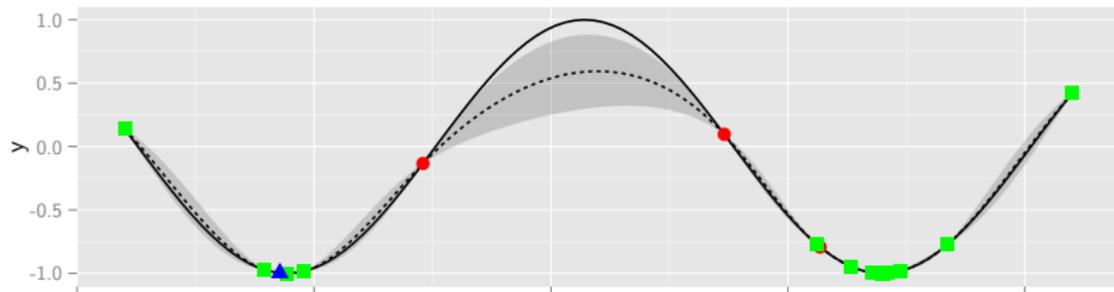
Iter = 5, Gap = 2.2202e-05

type — y ···· yhat    type ● init ▲ prop ■ seq



Iter = 15, Gap = 9.0305e-06

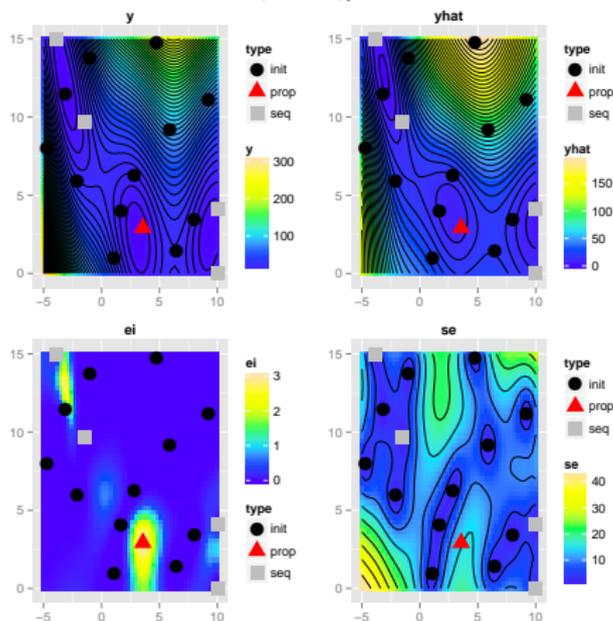
type — y ···· yhat    type ● init ▲ prop ■ seq



# MLRMBO: MODEL-BASED OPTIMIZATION TOOLBOX

Iter 5, x-axis: x1, y-axis: x2

- ▶ Any regression from mlr
- ▶ Arbitrary infill
- ▶ Single - or multi-crit
- ▶ Multi-point proposal
- ▶ Via parallelMap and batchtools runs on many parallel backends and clusters
- ▶ Algorithm configuration
- ▶ Active research

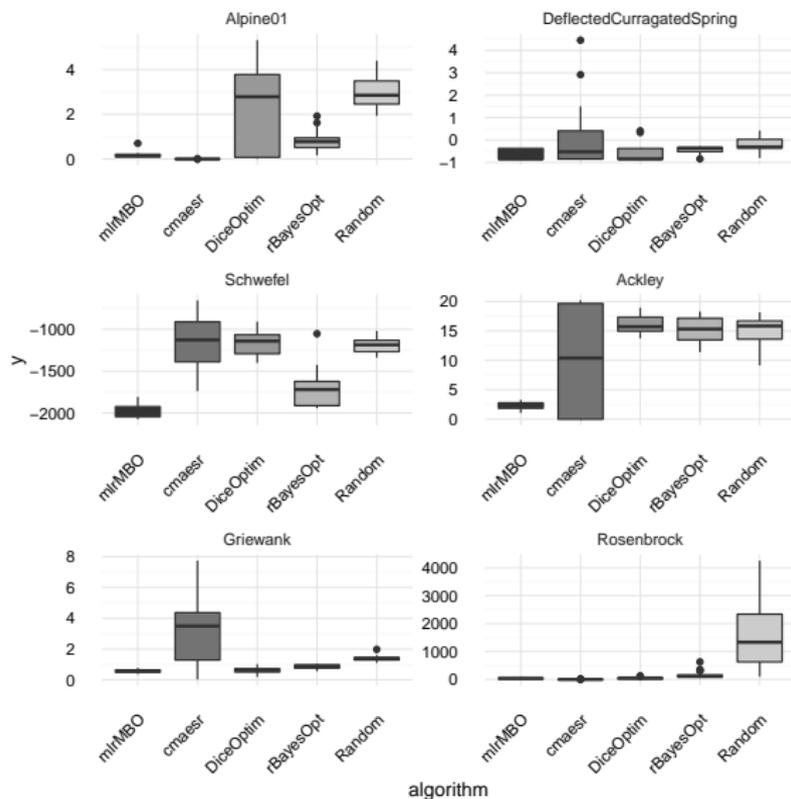


- ▶ mlr: <https://github.com/mlr-org/mlr>
- ▶ mlrMBO: <https://github.com/mlr-org/mlrMBO>
- ▶ mlrMBO Paper on arXiv (under review)  
<https://arxiv.org/abs/1703.03373>

# BENCHMARK MBO ON ARTIFICIAL TEST FUNCTIONS

- ▶ Comparison of mlrMBO on multiple different test functions
  - ▶ Multimodal
  - ▶ Smooth
  - ▶ Fully numeric
  - ▶ *Well known*
- ▶ We use GPs with
  - ▶ LCB with  $\lambda = 1$
  - ▶ Focusesearch
  - ▶ 200 iterations
  - ▶ 25 point initial design, created by LHS sampling
- ▶ Comparison with
  - ▶ Random search
  - ▶ CMAES
  - ▶ other MBO implementations in R

# MBO GP vs. COMPETITORS IN 5D



## Section 2

# PARALLEL BATCH PROPOSALS

# MOTIVATION FOR BATCH PROPOSAL

- ▶ Function evaluations expensive
- ▶ Often many cores available on a cluster
- ▶ Underlying  $f$  can in many cases not be easily parallelized
- ▶ Natural to consider batch proposal
- ▶ Parallel MBO: suggest  $q$  promising points to evaluate:  $\mathbf{x}_1^*, \dots, \mathbf{x}_q^*$
- ▶ We need to balance exploration and exploitation
- ▶ Non-trivial to construct infill criterion for this

# REVIEW OF PARALLEL MBO STRATEGIES

- ▶ **Constant liar:** (Ginsbourger et al., 2010)
  - ▶ Fit kriging model based on real data and find  $\mathbf{x}_1^*$  according to EI-criterion.
  - ▶ “Guess”  $f(\mathbf{x}_{i-1}^*)$ , update the model and find  $\mathbf{x}_i^*$ ,  $i = 2, \dots, q$
  - ▶ Use  $f_{min}$  for “guessing”
  
- ▶  **$q$ -LCB:** (Hutter et al., 2012)
  - ▶  $q$  times: sample  $\lambda$  from  $Exp(1)$  and optimize single LCB criterion
  - ▶  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} \text{LCB}(\mathbf{x}) = \arg \min_{\mathbf{x} \in \mathcal{X}} \hat{f}(\mathbf{x}) - \lambda \hat{s}(\mathbf{x})$ .

# MULTIOBJECTIVIZATION AND PROPOSED IDEA

## ► Multiobjectivization

- Originates from multi-modal optimization
- Add distance to neighbors for current set as artificial objective
- Use multiobjective optimization
- Select by hypervolume or first objective or ...

## ► Our approach

- Decouple  $\hat{f}(x)$  and  $\hat{s}(x)$  as objectives – instead of EI – to have different exploration / exploitation trade-offs
- Consider distance measure as potential extra objective
- Run multiobjective EA to select  $q$  well-performing, diverse points
- Distance is possible alternative if no or bad  $\hat{s}(x)$  estimator
- Decoupling  $y(x)$ ,  $\hat{s}(x)$  potential alternative when EI derivation does not hold for other model classes

Bischl, Wessing et al: *MOI-MBO: Multiobjective infill for parallel model-based optimization*, LION 2014

# EXPERIMENTAL SETUP

## Problem Instances

- ▶ All 24 test functions of the black-box optimization benchmark (BBOB) noise-free test suite
- ▶ Dimensions  $d \in \{5, 10\}$

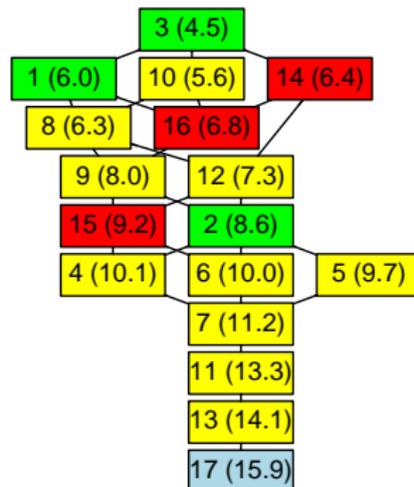
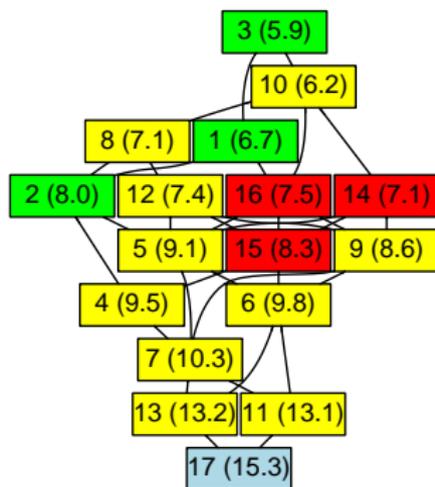
## Budget

- ▶ For every function 10 initial designs of size  $5 \cdot d$   
⇒ 10 statistical replications for each problem instance
- ▶  $40 \cdot d$  function evaluations on top of the initial design
- ▶ Parallel optimization: batches of size  $q = 5$

## Visualization: Preference relation graph

- ▶ Each node represents an approach (mean rank in braces)
- ▶ Two nodes are connected with an edge if one approach (the upper) is significantly better than the other (the lower) according to the sign test

# RESULT GRAPHS



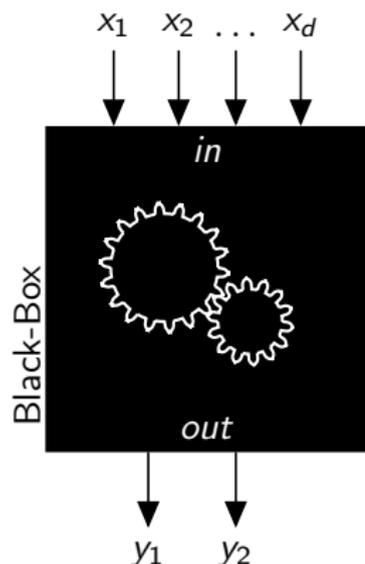
- 01: ego
- 02: ego\_ea\_ei
- 03: ego\_ea\_mean
- 04: moi\_ei.dist\_nb\_first
- 05: moi\_ei.dist\_nb\_hv
- 06: moi\_ei.dist\_nn\_first
- 07: moi\_ei.dist\_nn\_hv
- 08: moi\_mean.se.dist\_nb\_first
- 09: moi\_mean.se.dist\_nb\_hv

- 10: moi\_mean.se.dist\_nn\_first
- 11: moi\_mean.se.dist\_nn\_hv
- 12: moi\_mean.se\_nn\_first
- 13: moi\_mean.se\_nn\_hv
- 14: par\_cl
- 15: par\_cl\_ea
- 16: par\_lcb
- 17: random\_search

## Section 3

# MULTICRITERIA SMBO

# MODEL-BASED MULTI-OBJECTIVE OPTIMIZATION



$$\min_{\mathbf{x} \in \mathcal{X}} \mathbf{f}(\mathbf{x}) = \mathbf{y} = (y_1, \dots, y_m) \text{ with } \mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^m$$

►  $\mathbf{y}$  dominates  $\tilde{\mathbf{y}}$  if (3)

$$\forall i \in \{1, \dots, m\} : y_i \leq \tilde{y}_i \quad (4)$$

$$\text{and } \exists i \in \{1, \dots, m\} : y_i < \tilde{y}_i \quad (5)$$

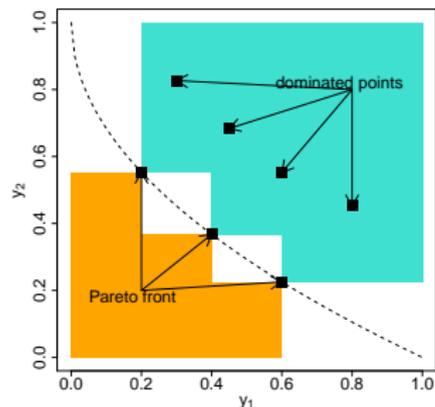
► Set of non-dominated solutions:

$$\mathcal{X}^* := \{\mathbf{x} \in \mathcal{X} \mid \nexists \tilde{\mathbf{x}} \in \mathcal{X} : \mathbf{f}(\tilde{\mathbf{x}}) \text{ dominates } \mathbf{f}(\mathbf{x})\}$$

► Pareto set  $\mathcal{X}^*$ , Pareto front  $\mathbf{f}(\mathcal{X}^*)$

► Goal: Find  $\hat{\mathcal{X}}^*$  of non-dominated points that estimates the true set  $\mathcal{X}^*$

# MODEL-BASED MULTI-OBJECTIVE OPTIMIZATION



$$\min_{\mathbf{x} \in \mathcal{X}} \mathbf{f}(\mathbf{x}) = \mathbf{y} = (y_1, \dots, y_m) \text{ with } \mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

- ▶  $\mathbf{y}$  dominates  $\tilde{\mathbf{y}}$  if (6)

$$\forall i \in \{1, \dots, m\} : y_i \leq \tilde{y}_i \quad (7)$$

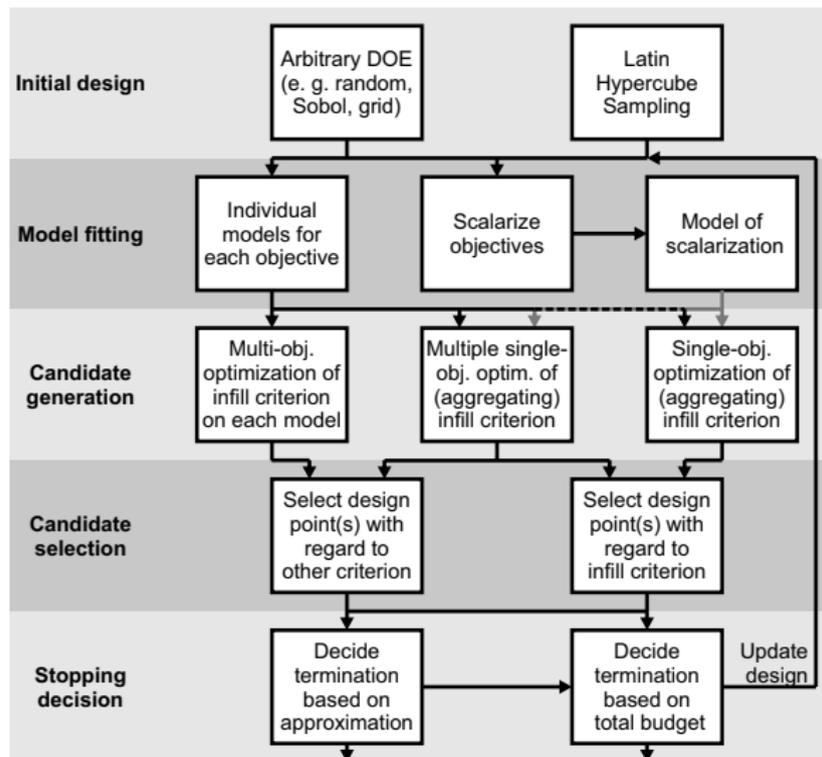
$$\text{and } \exists i \in \{1, \dots, m\} : y_i < \tilde{y}_i \quad (8)$$

- ▶ Set of non-dominated solutions:

$$\mathcal{X}^* := \{\mathbf{x} \in \mathcal{X} \mid \nexists \tilde{\mathbf{x}} \in \mathcal{X} : \mathbf{f}(\tilde{\mathbf{x}}) \text{ dominates } \mathbf{f}(\mathbf{x})\}$$

- ▶ Pareto set  $\mathcal{X}^*$ , Pareto front  $\mathbf{f}(\mathcal{X}^*)$
- ▶ Goal: Find  $\hat{\mathcal{X}}^*$  of non-dominated points that estimates the true set  $\mathcal{X}^*$

# TAXONOMY



Horn, Wagner, Bischl et al: *Model-based multi-objective optimization: Taxonomy, multi-point proposal, toolbox and benchmark*, EMO 2014

# BATCH PROPOSAL

- ▶ Most MBMO lack way to propose  $N > 1$  points (batch evaluation)
- ▶ Batch evaluations are essential for distributed computing
- ▶ We integrated such mechanism(s) for arbitrary MBMO
- ▶ Replaced single phases of the taxonomy

# PAREGO

1. Scalarize objectives using the augmented Tchebycheff norm

$$\max_{i=1,\dots,d} [w_i f_i(\mathbf{x})] + \rho \sum_{i=1}^d w_i f_i(\mathbf{x})$$

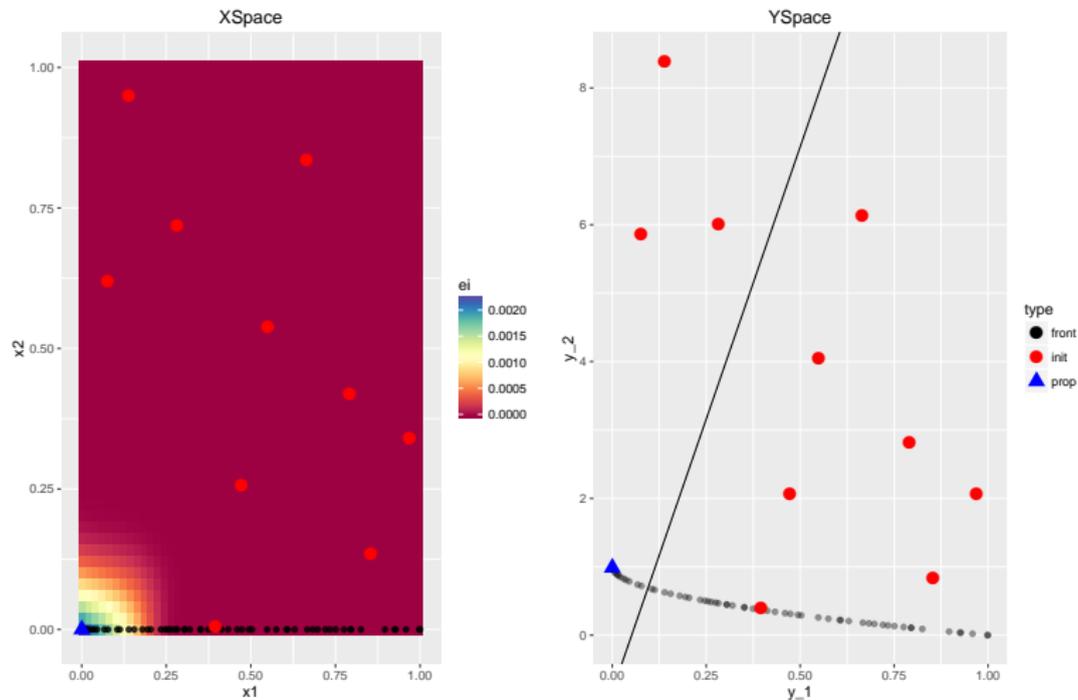
with uniformly distributed weight vector  $\mathbf{w}$  ( $\sum w_i = 1$ ) and fit surrogate model to the respective scalarization.

2. Single-objective optimization of EI (or LCB?)

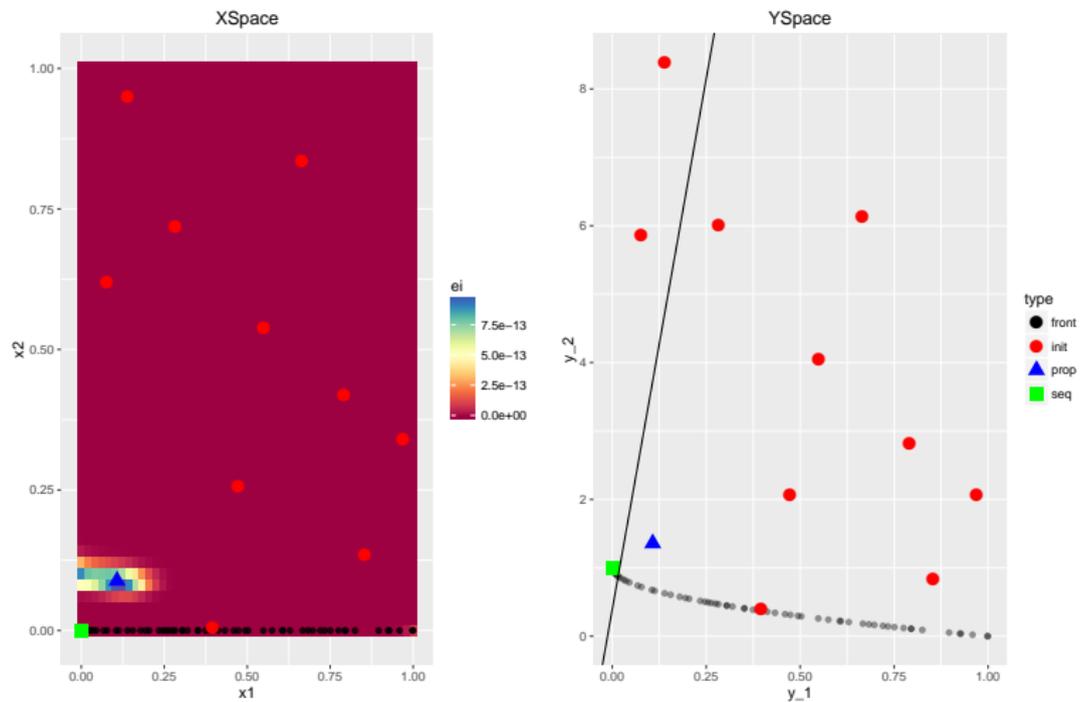
Batch proposal: **Increase the number and diversity of randomly drawn weight vectors**

- ▶ If  $N$  points are desired,  $cN$  ( $c > 1$ ) weight vectors are considered
- ▶ Greedily reduce set of weight vectors by excluding one vector of the pair with minimum distance
- ▶ Scalarizations implied by each weight vector are computed
- ▶ Fit and optimize models for each scalarization
- ▶ Optima of each model build the batch to be evaluated

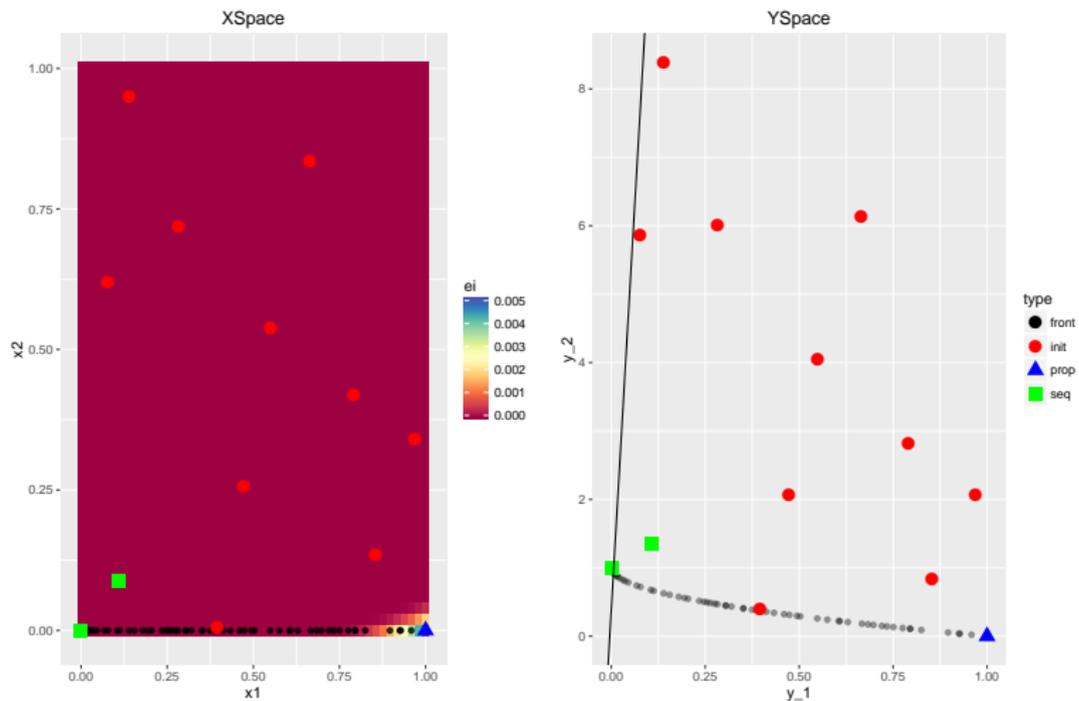
# ANIMATION OF PAREGO



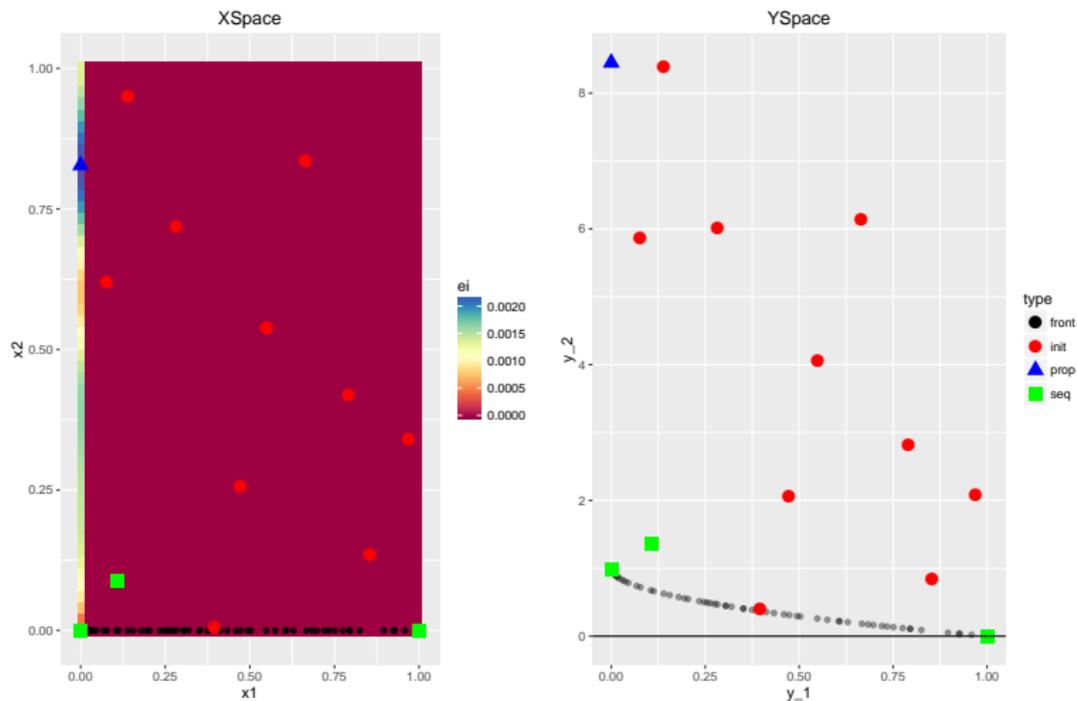
# ANIMATION OF PAREGO



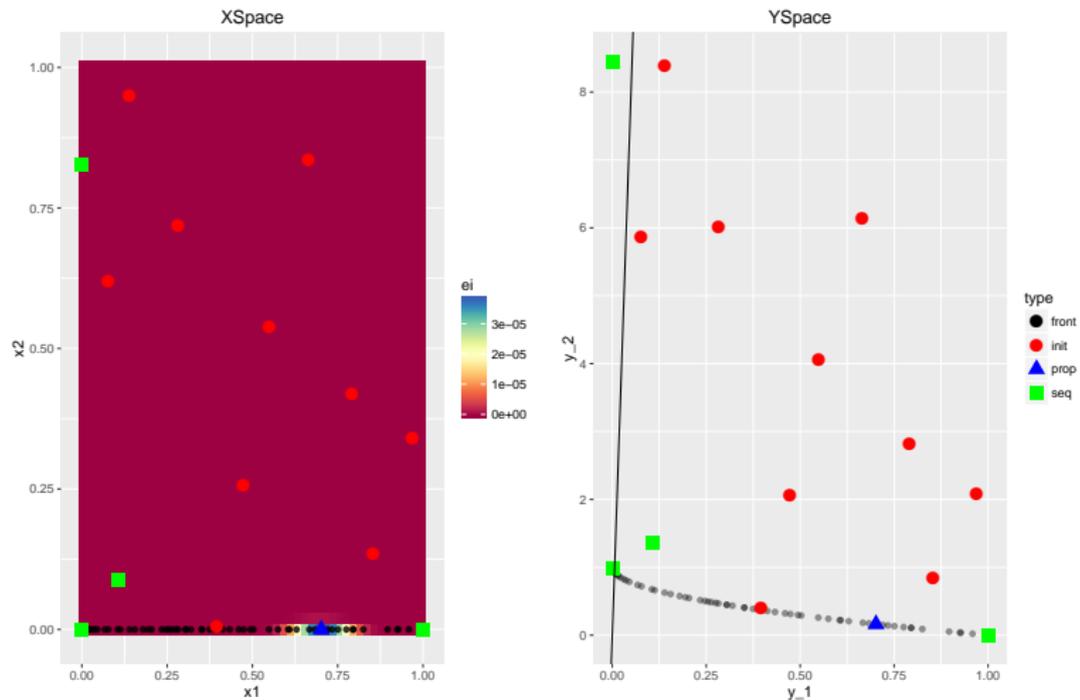
# ANIMATION OF PAREGO



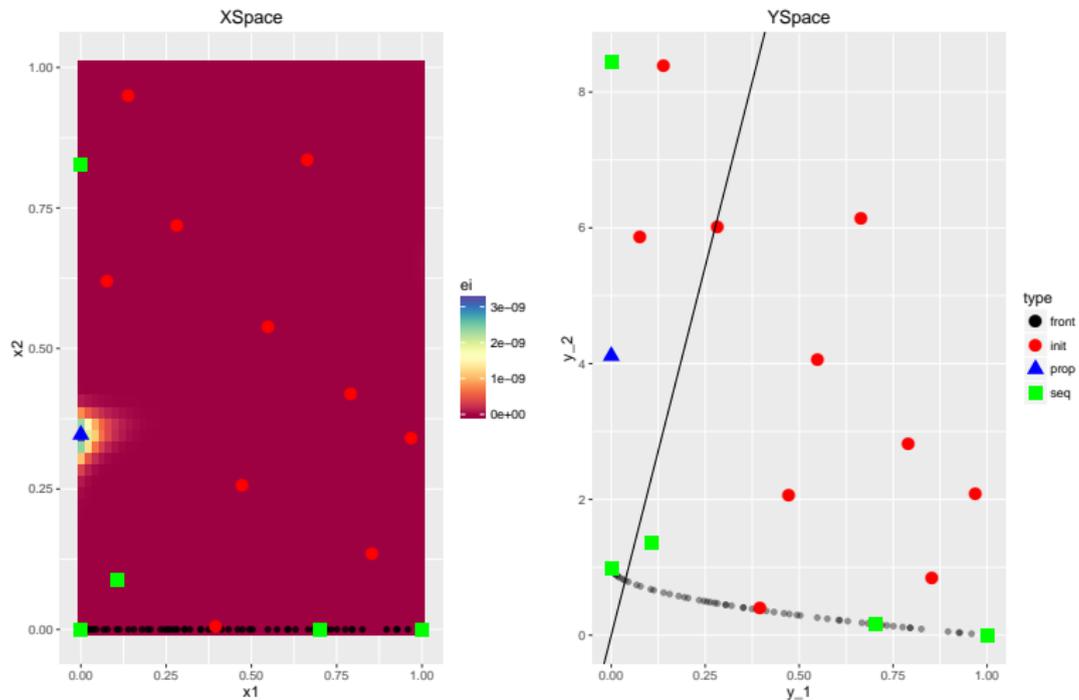
# ANIMATION OF PAREGO



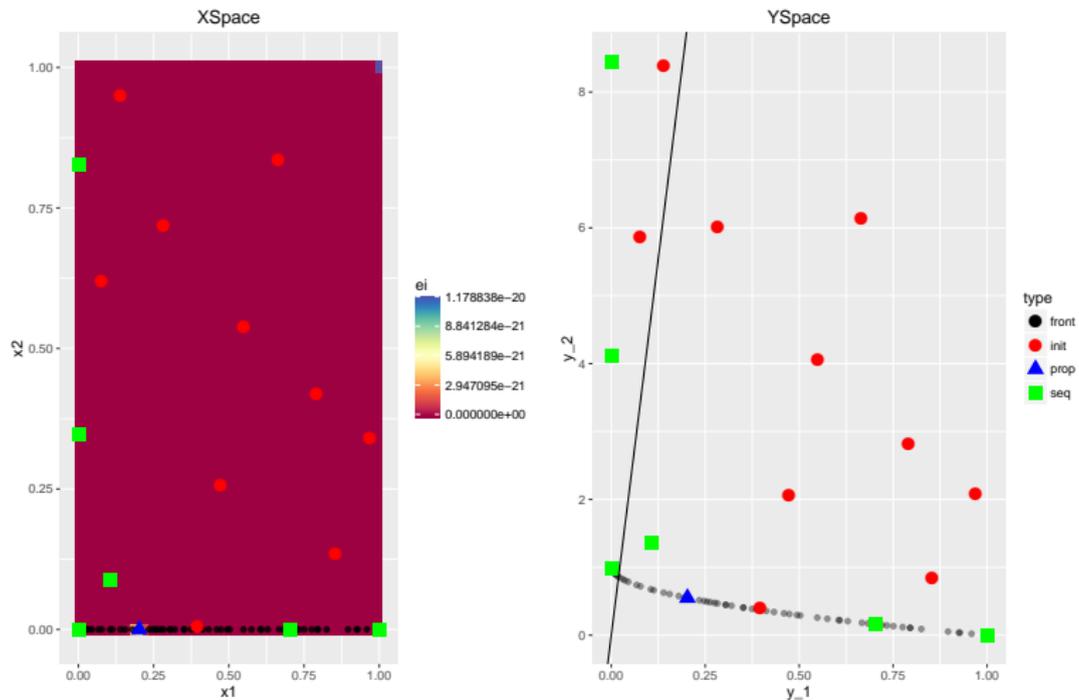
# ANIMATION OF PAREGO



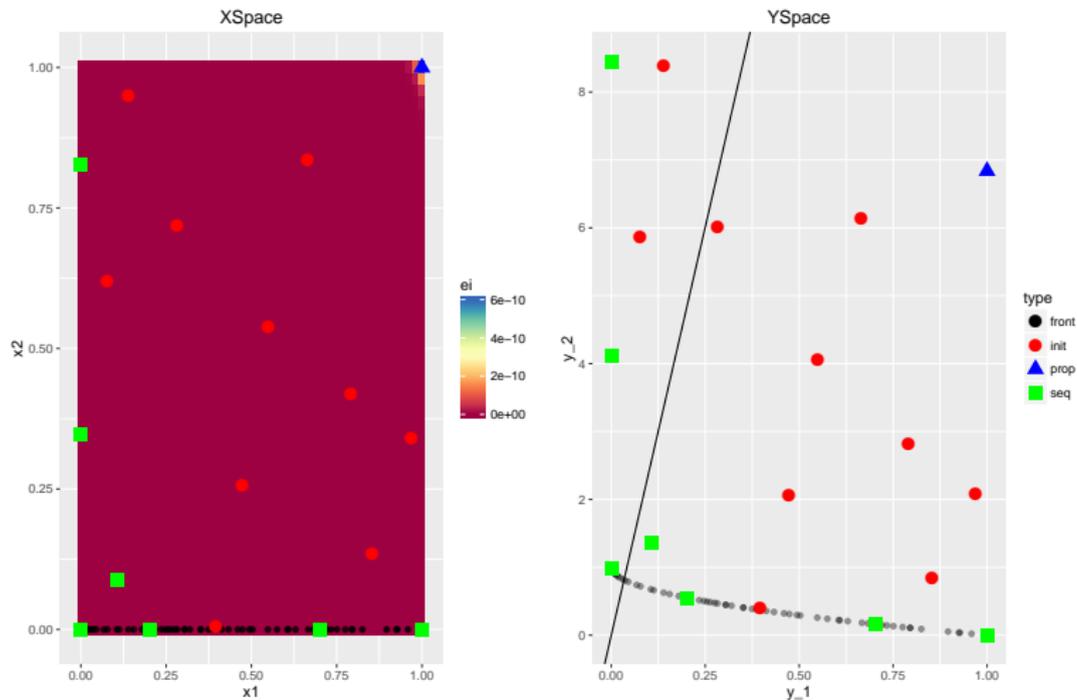
# ANIMATION OF PAREGO



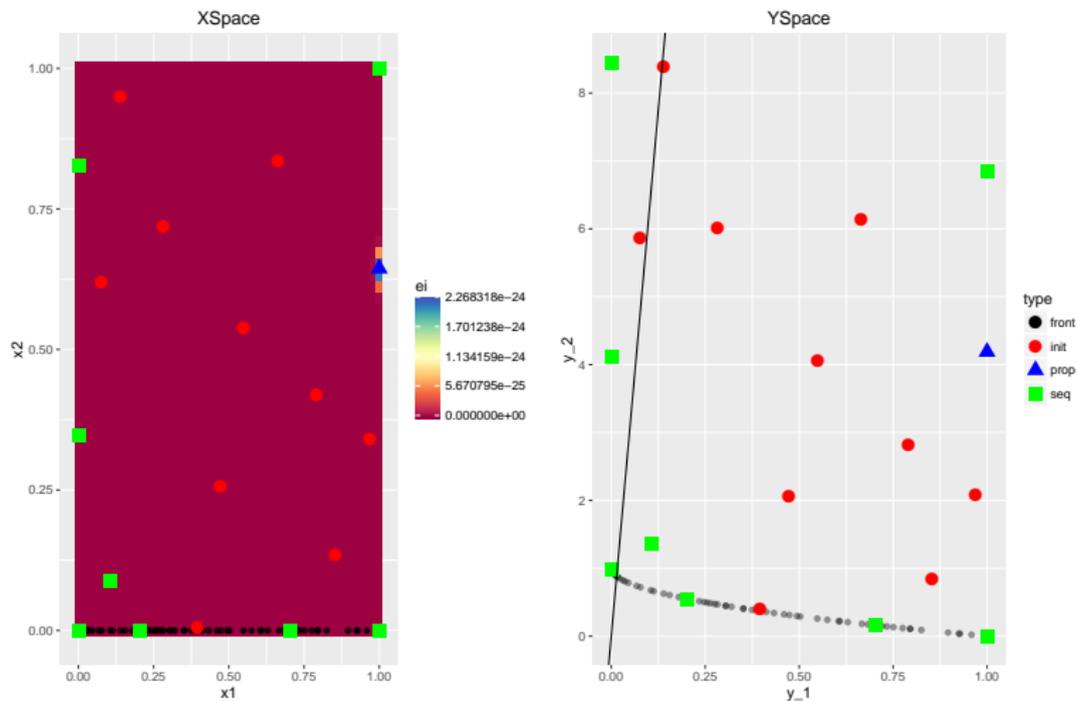
# ANIMATION OF PAREGO



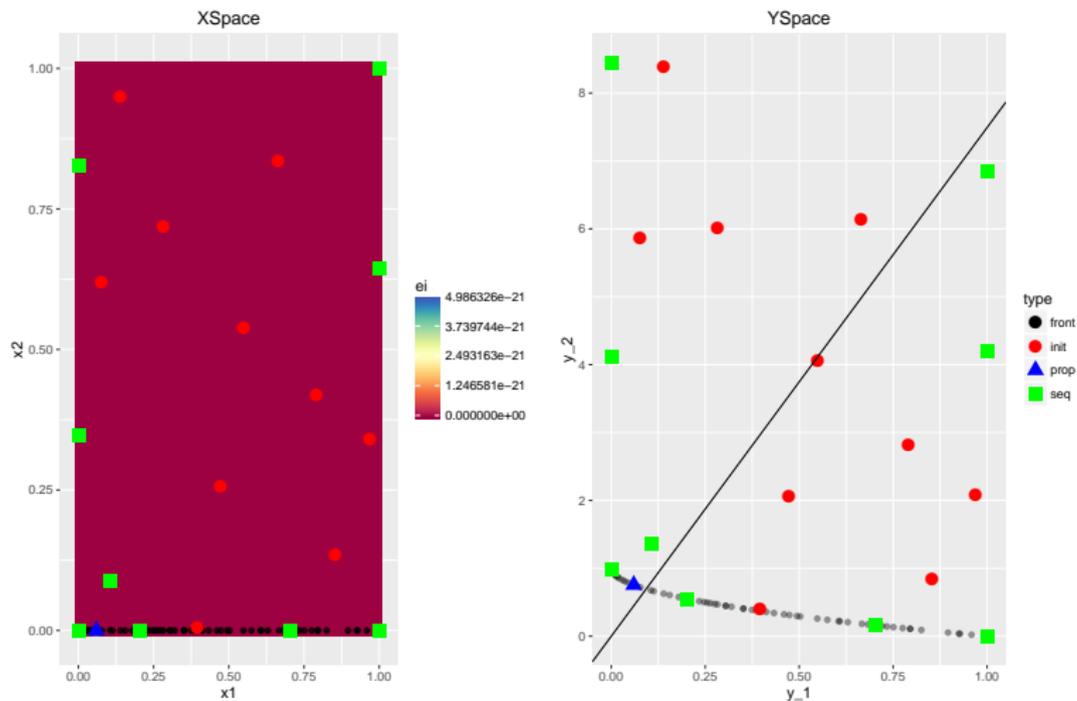
# ANIMATION OF PAREGO



# ANIMATION OF PAREGO

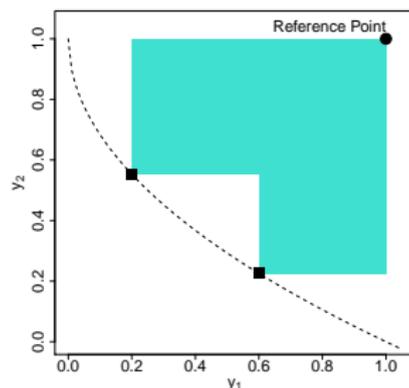


# ANIMATION OF PAREGO



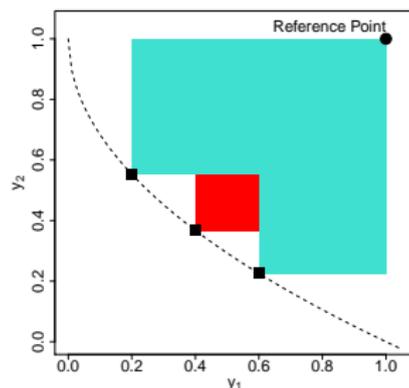
# SMS-EGO

- ▶ Individual models for each objective
- ▶ Single-objective optimization of aggregating infill criterion:  
Calculate contribution of the confidence bound of representative solution to the current front approximation
- ▶ Calculate LCB for each objective
- ▶ Measure contribution with regard to the hypervolume indicator
- ▶ For  $\varepsilon$ -dominated ( $\preceq_\varepsilon$ ) solutions, a penalty  
$$\Psi(\mathbf{x}) = -1 + \prod_{j=1}^m \left( 1 + (l(\mathbf{x}) - y_j^{(i)}) \right)$$
  
is added  
(Actually not needed for Focussearch.)



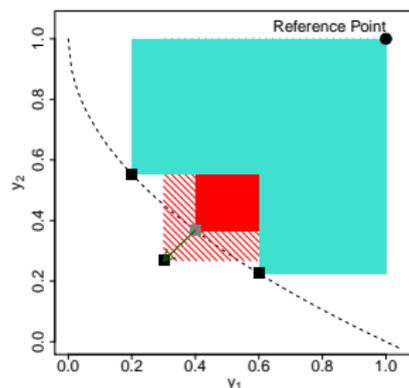
# SMS-EGO

- ▶ Individual models for each objective
- ▶ Single-objective optimization of aggregating infill criterion:  
Calculate contribution of the confidence bound of representative solution to the current front approximation
- ▶ Calculate LCB for each objective
- ▶ Measure contribution with regard to the hypervolume indicator
- ▶ For  $\varepsilon$ -dominated ( $\preceq_\varepsilon$ ) solutions, a penalty  
$$\Psi(\mathbf{x}) = -1 + \prod_{j=1}^m \left( 1 + (l(\mathbf{x}) - y_j^{(i)}) \right)$$
  
is added  
(Actually not needed for Focussearch.)



# SMS-EGO

- ▶ Individual models for each objective
- ▶ Single-objective optimization of aggregating infill criterion:  
Calculate contribution of the confidence bound of representative solution to the current front approximation
- ▶ Calculate LCB for each objective
- ▶ Measure contribution with regard to the hypervolume indicator
- ▶ For  $\varepsilon$ -dominated ( $\preceq_\varepsilon$ ) solutions, a penalty  
$$\Psi(\mathbf{x}) = -1 + \prod_{j=1}^m \left( 1 + (l(\mathbf{x}) - y_j^{(i)}) \right)$$
  
is added  
(Actually not needed for Focussearch.)



# SMS-EGO: BATCH PROPOSAL

Modification of phase candidate generation: **Use simulated evaluations for candidate generation**

- ▶ The proposed point  $\vec{x}^*$  is not directly evaluated, but the LCB  $l(\vec{x}^*)$  is added to the current approximation without refitting the model
- ▶ Repeat until  $N$  points for a batch evaluation have been found

# APPROXIMATIVE RBF-SVM TRAINING ALGORITHMS

---

SVM solver	Description
Pegasos	Stochastic Gradient Descent
BSGD	Budgeted Stochastic Gradient Descent
LLSVM	Low-rank kernel approximation + linear solver
<b>LIBSVM</b>	<b>"Exact" SMO solver</b>
LASVM	Online variant of SMO solver
LIBBVM/CVM	Minimum Enclosing Ball (only squared hinge loss)
SVMperf	Cutting Plane Algorithm

---

- ▶ What is the trade-off between training time and prediction error?
- ▶ Most solvers have 2 additional parameters on top of  $C$  and  $\gamma$
- ▶ Optimizing 2 expensive objectives in a 4-dim parameter space.
- ▶ Replace grid search with more sophisticated PAREGO-algorithm.

# APPROXIMATIVE SVM TRAINING ALGORITHMS

- ▶ **We expect:** Every solver has a trade-off between training time and prediction error: Given more time a solver (should) reach a better solution.
- ▶ **Our goal:** Analyze this trade-off! Solve the multi-criteria optimization problem with respect to the two objectives error and training time by varying the parameters.
- ▶ **The challenge:** Optimizing 2 expensive objectives in a 4-dimensional parameter space.
- ▶ **Our approach:** Replace standard grid search with more sophisticated PAREGO-algorithm.

# APPROXIMATIVE SVM TRAINING ALGORITHMS

The parameters  $(C, \gamma)$  of the SVM itself were optimized over  $2^{[-15,15]}$  respectively. Every solver has further approximation parameters:

SVM solver	Parameters	Optimization Space
Pegasos	#Epochs	$2^{[0,7]}$
BGSD	Budget size, #Epochs	$2^{[4,11]}$ , $2^{[0,7]}$
LLSVM	Matrix rank	$2^{[4,11]}$
LIBSVM	$\epsilon$ (Accuracy)	$2^{[-13,-1]}$
LASVM	$\epsilon$ (Accuracy), #Epochs	$2^{[-13,-1]}$ , $2^{[0,7]}$
LIBBVM/CVM	$\epsilon$ (Accuracy)	$2^{[-19,-1]}$
SVMperf	$\epsilon$ (Accuracy), #Cutting planes	$2^{[-13,-1]}$ , $2^{[4,11]}$

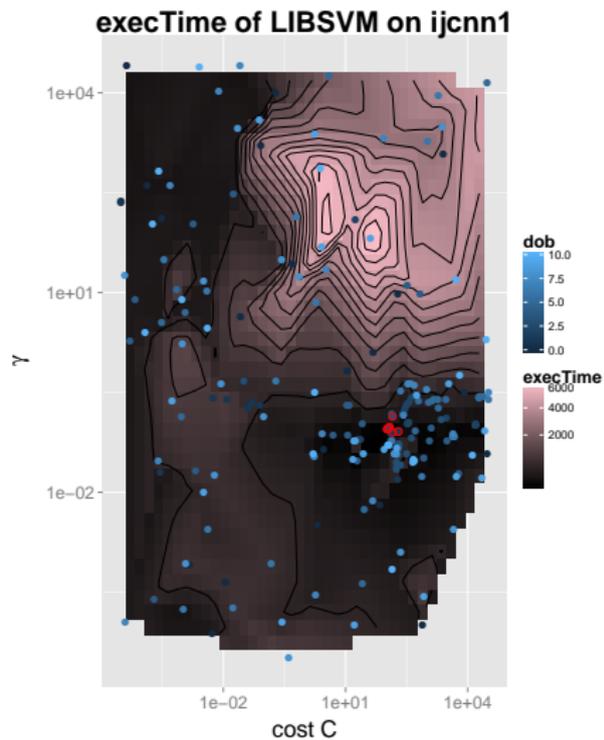
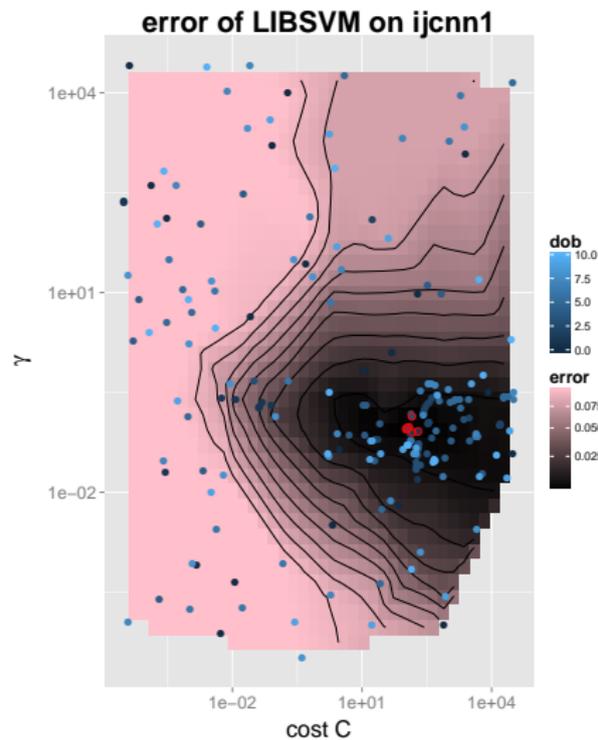
Additional parameters set to default values.

# DATASETS

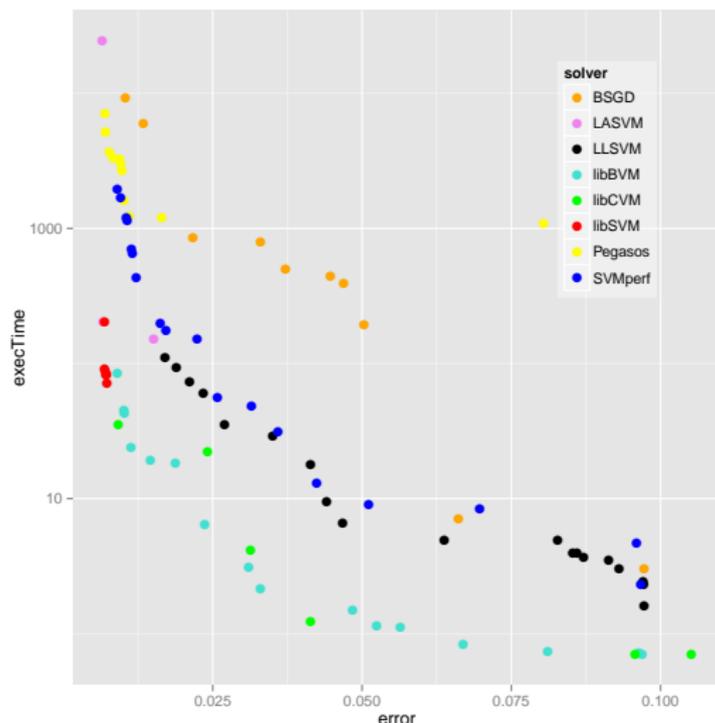
data set	# points	# features	class ratio	sparsity
wXa	34 780	300	34.45	95.19 %
aXa	36 974	123	3.17	88.72 %
protein	42 153	357	1.16	71.46 %
mnist	70 000	780	0.96	80.76 %
vehicle	98 528	100	1.00	0 %
shuttle	101 500	9	0.27	0.23 %
spektren	175 090	22	0.80	0 %
ijcnn1	176 691	22	9.41	40.91 %
arthrosis	262 142	178	1.19	0.01 %
cod-rna	488 565	8	2.00	0.02 %
covtype	581 012	54	1.05	78 %
poker	1 025 010	10	1.00	0 %

TABLE: Overview of the data sets.

# TEST ERROR LANDSCAPE (LIBSVM)



# ALL PARETO FRONTS FOR IJCNN1 DATASET



- ▶ 176 691 samples  
22 features  
9.41 class ratio
- ▶ We see:
  - LIBSVM: exact but slow
  - LIBBVM/CVM: good front - speed increase with small accuracy loss
  - SVMperf / LLSVM / BSGD: can be really fast, but higher accuracy loss
  - LASVM / Pegasos: less exact and even slower as LIBSVM

## Section 4

# INTERESTING CHALLENGES

## CHALLENGE: THE CORRECT SURROGATE?

- ▶ GPs are very much tailored to what we want to do, due to their spatial structure in the kernel and the uncertainty estimator.
- ▶ But GPs are rather slow. And (fortunately) due to parallelization (or speed-up tricks like subsampling) we have more design points to train on.
- ▶ Categorical features are also a problem in GPs (although methods exist, usually by changing the kernel)
- ▶ Random Forests handle categorical features nicely, are much faster. But they don't rely on a spatial kernel and the uncertainty estimation is much more heuristic / may not represent what we want.

# CHALLENGE: TIME HETEROGENEITY

- ▶ Complex configuration spaces across many algorithms results in vastly different runtimes in design points.
- ▶ Actually just the RBF-SVM tuning can result in very different runtimes.
- ▶ We don't care how many points we evaluate, we care about total walltime of the configuration.
- ▶ The option to subsample further complicates things.
- ▶ Parallelization further complicates things.
- ▶ Option: Estimate runtime as well with a surrogate, integrate it into acquisition function.

## Section 5

# ML MODEL SELECTION AND HYPERPARAMETER OPTIMIZATION

# AUTOMATIC MODEL SELECTION

## PRIOR APPROACHES:

- ▶ Looking for the silver bullet model  
    ↪ Failure
- ▶ Exhaustive benchmarking / search  
    ↪ Very expensive, often contradicting results
- ▶ Meta-Learning:  
    ↪ Good meta-features are hard to construct  
    ↪ IMHO: Gets more interesting when combined with SMBO

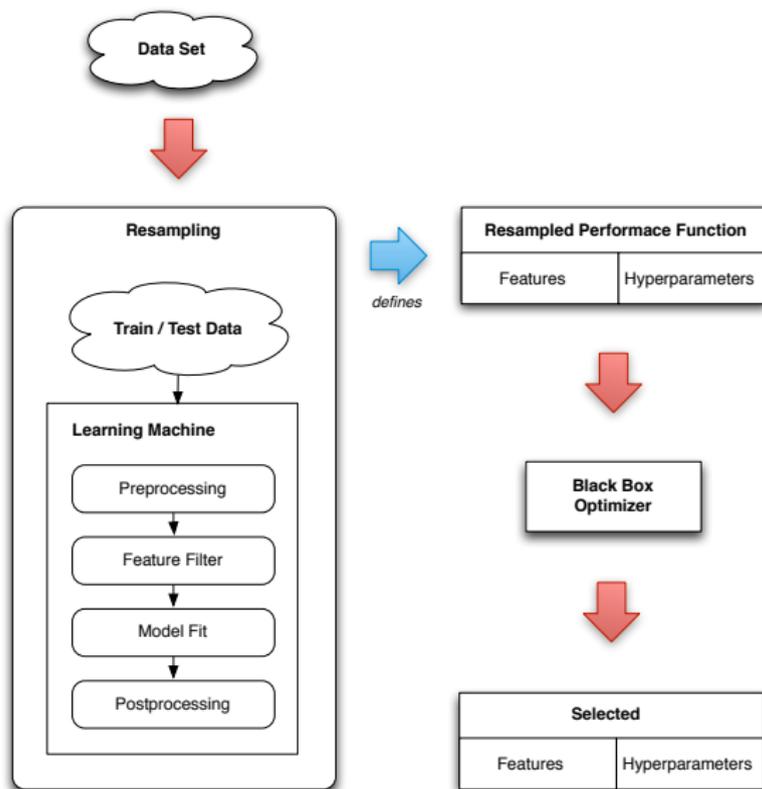
## GOAL FOR AUTOML:

- ▶ Data dependent
- ▶ Automatic
- ▶ Include every relevant modeling decision
- ▶ Efficient
- ▶ Learn on the model-settings level!

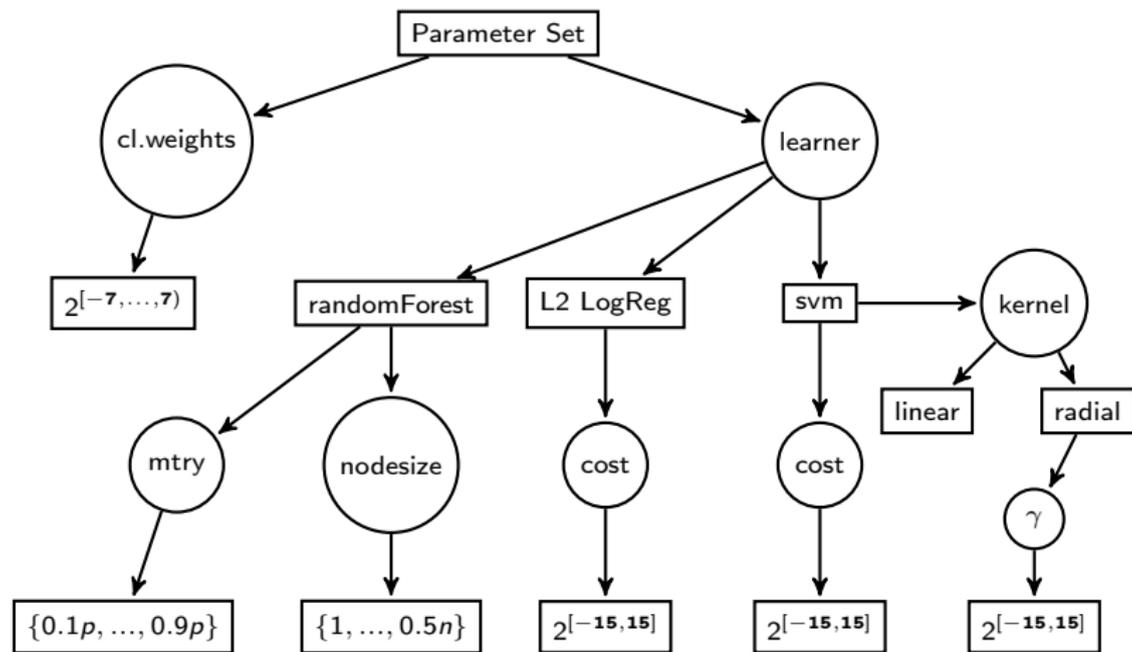
# FROM NORMAL SMBO TO HYPERPARAMETER TUNING

- ▶ Objective function is resampled performance measure
- ▶ Parameter space  $\theta \in \Theta$   
might be discrete and dependent / hierarchical
- ▶ No derivative for  $f(\cdot, \theta)$ , black-box
- ▶ Objective is stochastic / noisy
- ▶ Objective is expensive to evaluate
- ▶ In general we face a problem of algorithm configuration:
- ▶  $\rightsquigarrow$  Usual approaches: racing or model-based / bayesian optimization

# FROM NORMAL SMBO TO HYPERPARAMETER TUNING



# COMPLEX PARAMETER SPACE



# FROM NORMAL SMBO TO HYPERPARAMETER TUNING

- ▶ Initial design: LHS principle can be extended, or just use random
- ▶ Focus search: Can be (easily) extended, as it is based on random search. To zoom in for categorical parameters we randomly drop a category for each param which is not present in the currently best configuration.
- ▶ Few approaches for GPs with categorical params exist (usually with new covar kernels), not very established
- ▶ Alternative: Random regression forest (mlrMBO, SMAC)
- ▶ Estimate uncertainty / confidence interval for mean response by efficient bootstrap technique<sup>1</sup>, or jackknife, so we can define  $EI(x)$  for the RF
- ▶ Dependent params in mlrMBO: Imputation:
- ▶ Many of the current techniques to handle these problems are (from a theoretical standpoint) somewhat crude

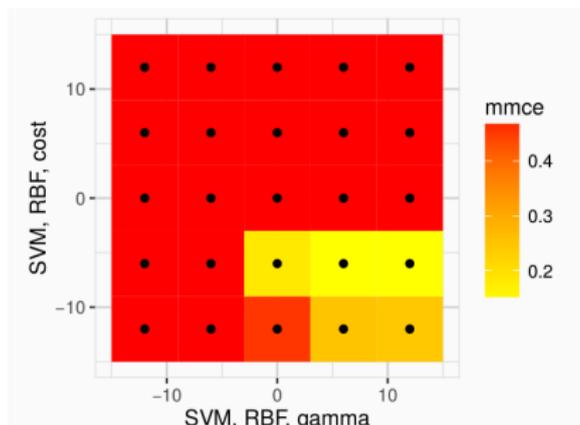
---

<sup>1</sup>Sexton et al, "Standard errors for bagged and random forest estimators, 2009."

# HYPERPARAMETER TUNING

- ▶ Still common practice: grid search  
For a SVM it might look like:
  - ▶  $C \in (2^{-12}, 2^{-10}, 2^{-8}, \dots, 2^8, 2^{10}, 2^{12})$
  - ▶  $\gamma \in (2^{-12}, 2^{-10}, 2^{-8}, \dots, 2^8, 2^{10}, 2^{12})$
  - ▶ Evaluate all  $13^2 = 169$  combinations  $C \times \gamma$
- ▶ Bad because:
  - ▶ optimum might be "off the grid"
  - ▶ lots of evaluations in bad areas
  - ▶ lots of costly evaluations
- ▶ How bad?

# HYPERPARAMETER TUNING



- ▶ Because of budget restrictions grid might even be smaller!
- ▶ Unpromising area quite big!
- ▶ Lots of costly evaluations!

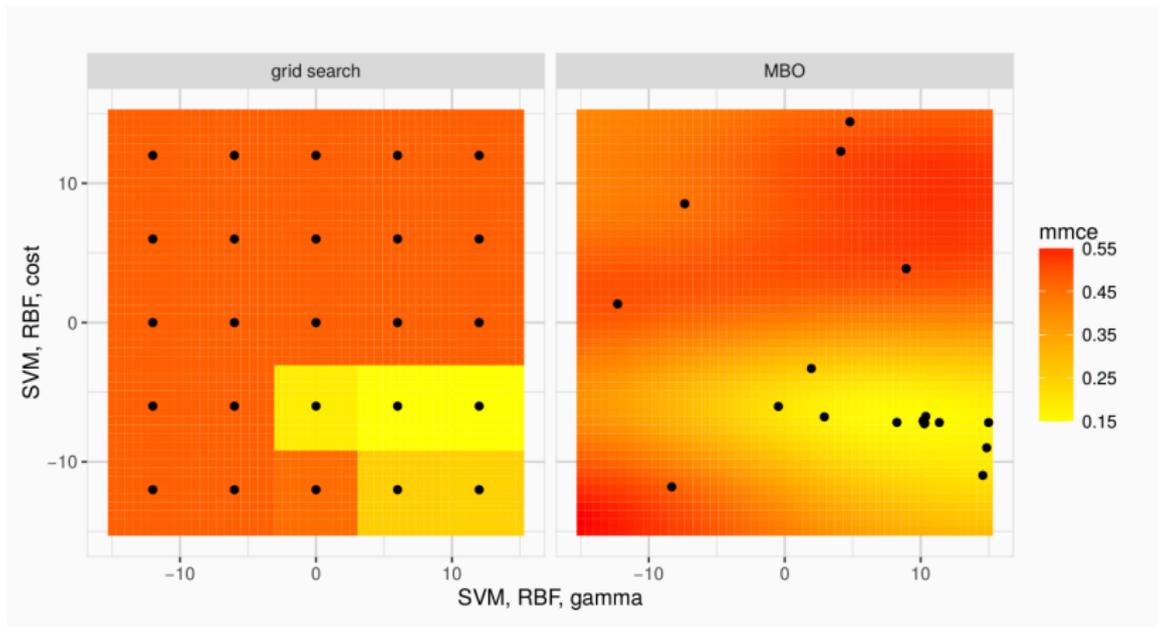
With **mlrMBO** it is not hard to do it better!

More interesting applications to time-series regression and cost-sensitive classification<sup>2</sup>

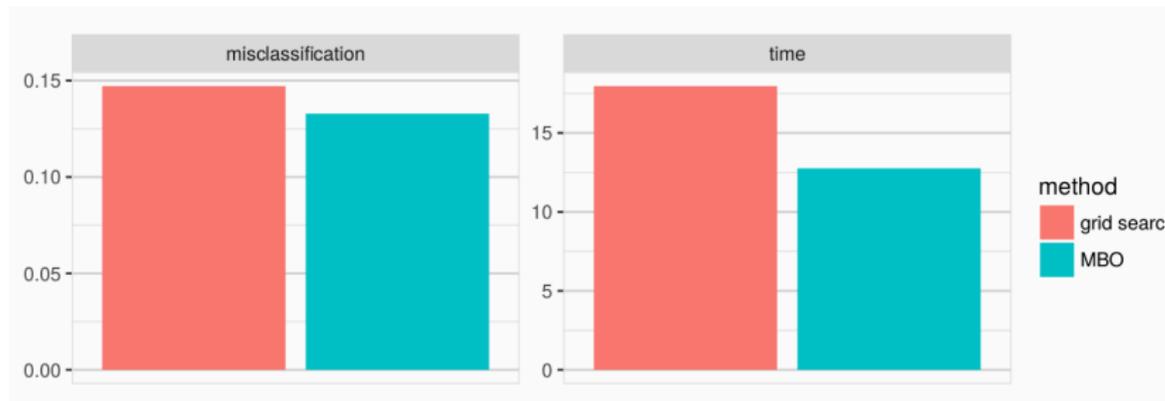
---

<sup>2</sup>Koch, Bischl et al: *Tuning and evolution of support vector kernels*, EI 2012

# HYPERPARAMETER TUNING



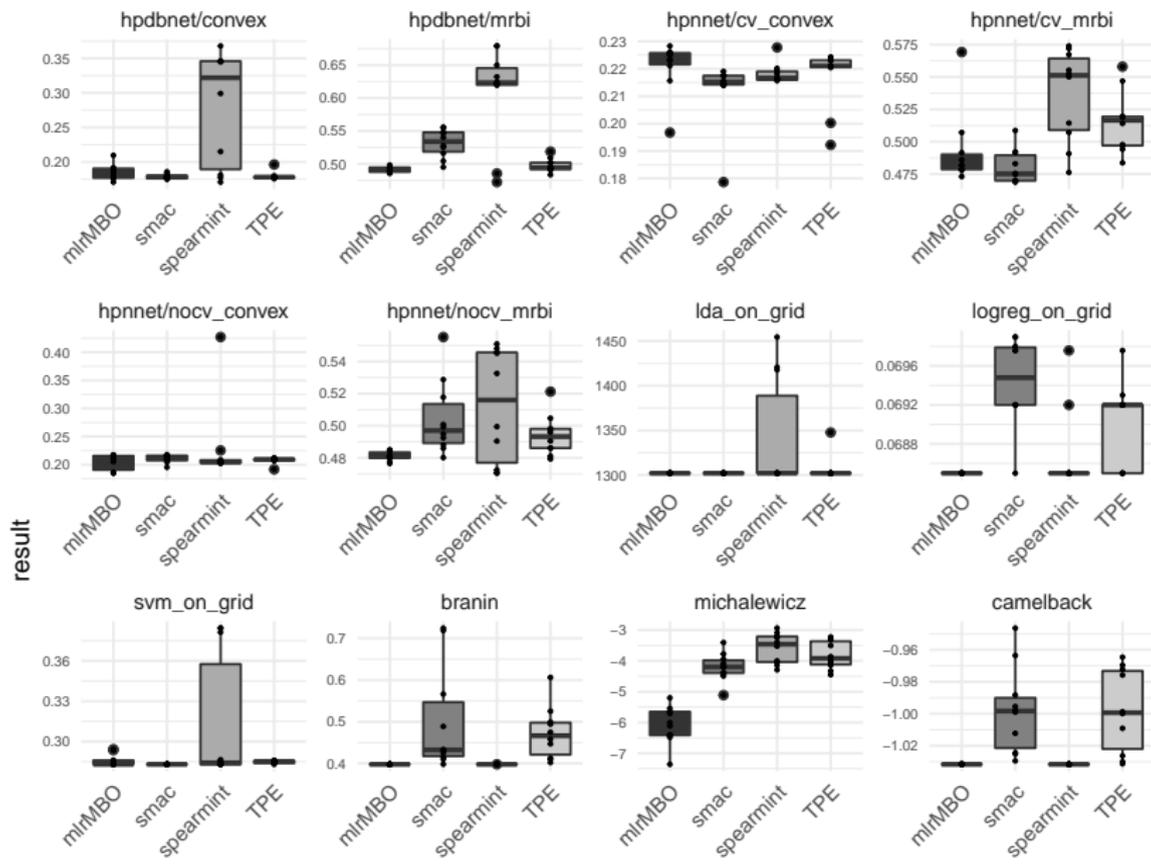
# HYPERPARAMETER TUNING



# HPOLIB

- ▶ HPOlib is a set of standard benchmarks for hyperparameter optimizer
- ▶ Allows comparison with
  - ▶ Spearmint
  - ▶ SMAC
  - ▶ Hyperopt (TPE)
- ▶ Benchmarks:
  - ▶ Numeric test functions (similar to the ones we've seen before)
  - ▶ Numeric machine learning problems (l1, SVM, logistic regression)
  - ▶ Deep neural networks and deep belief networks with 15 and 35 parameters.
- ▶ For benchmarks with discrete and dependent parameters (hpnnet, hpdbnet) a random forest with standard error estimation is used.

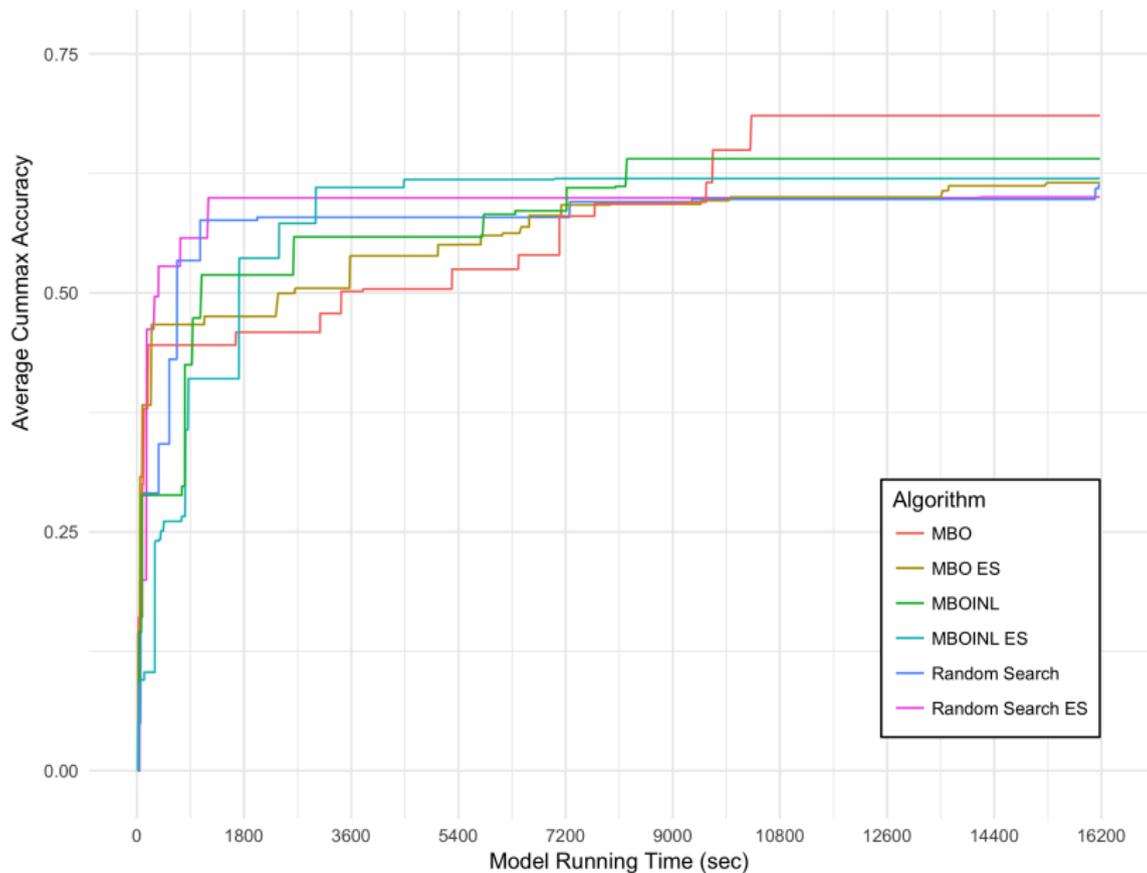
# MBO: HPOLIB



# DEEP LEARNING CONFIGURATION EXAMPLE

- ▶ Dataset: CIFAR-10 (60000 32x32 images with 3 color channels; 10 classes)
- ▶ Configuration of a deep neural network (mxnet)
- ▶ Size of parameter set: 30, including number of hidden layers, activation functions, regularization, convolution layer setting, etc.
- ▶ Split: 2/3 training set, 1/6 test set, 1/6 validation set
- ▶ Time budget per tuning run: 4.5h (16200 sec)
- ▶ Surrogate: Random forest
- ▶ Acquisition: LCB with  $\lambda = 2$

# DEEP LEARNING CONFIGURATION EXAMPLE



Thanks! Any  
comments or  
questions?