

MACHINE LEARNING IN R: PACKAGE MLR

Bernd Bischl
tiny.cc/cdfmay

WELCOME!

- Project home page
<https://github.com/mlr-org/mlr>
 - ▶ R documentation rendered in HTML
 - ▶ Tutorial for online viewing / download, including many examples
 - ▶ Don't hesitate to interrupt me
 - ▶ There will be a coffee break (I hope?)
- 8-10 main developers, quite a few contributors, 3 GSOC projects in 2015
- About 20K lines of code, 8K lines of unit tests

OVERVIEW

INTRODUCTION

WHY MLR?

BUILDING BLOCKS

BENCHMARKING AND MODEL COMPARISON

HYPERPARAMETER TUNING

FEATURE SELECTION

MLR LEARNER WRAPPERS

PARALLELIZATION

VISUALIZATIONS

OPENML

THE END

Section 1

INTRODUCTION

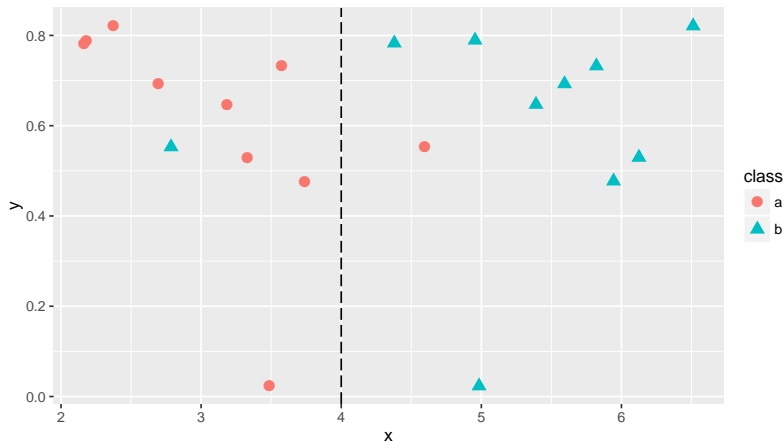
WHAT IS (SUPERVISED) MACHINE LEARNING?

- Learning structure in data
- The art of predicting stuff
- Model optimization
- Understanding of grey-box models

DISCLAIMER

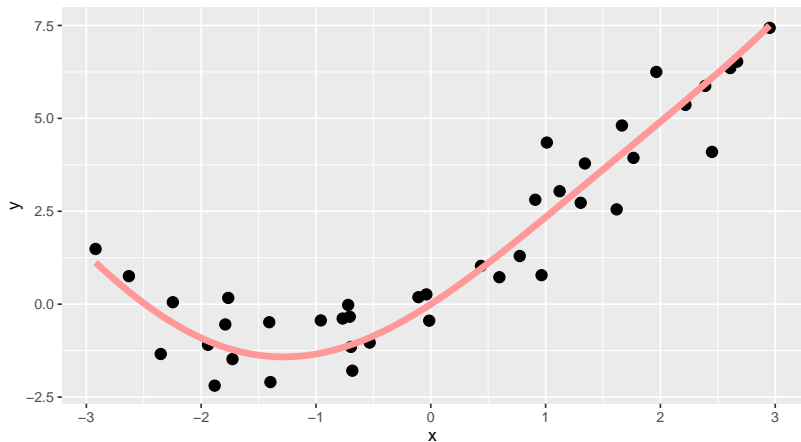
- The list is subjective and naively tailored to this talk
- ML is based on math and statistics, we will (mainly) talk about structure, software, and practical issues here

SUPERVISED CLASSIFICATION TASKS



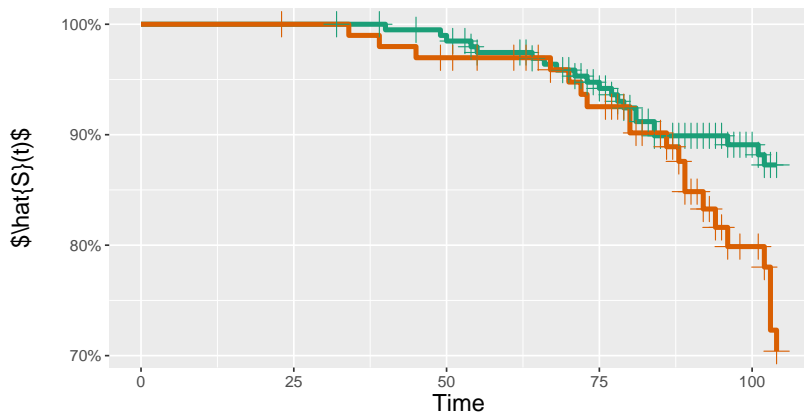
GOAL: Predict a class (or membership probabilities)

SUPERVISED REGRESSION TASKS



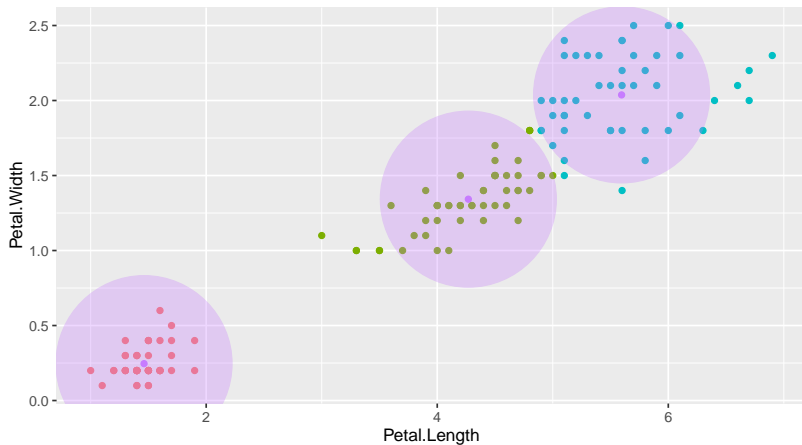
GOAL: Predict a continuous output

SUPERVISED SURVIVAL TASKS



GOAL: Predict a survival function $\hat{S}(t)$, i.e. the probability to survive to time point t

UNSUPERVISED CLUSTER TASKS



GOAL: Group data into similar clusters (or estimate fuzzy membership probabilities)

Section 2

WHY MLR?

MOTIVATION

THE GOOD NEWS

- CRAN serves hundreds of packages for machine learning (cf. CRAN task view machine learning)
- Many packages are compliant to the unwritten interface definition:

```
> model = fit(target ~ ., data = train.data, ...)  
> predictions = predict(model, newdata = test.data, ...)
```

MOTIVATION

THE BAD NEWS

- Some packages do not support the formula interface or their API is “just different”
- Functionality is always package or model-dependent, even though the procedure might be general
- No meta-information available or buried in docs (sometimes not documented at all)
- Many packages require the user to “guess” good hyperparameters
- Larger experiments lead to lengthy, tedious and error-prone code

Our goal: A domain-specific language for many machine learning concepts!

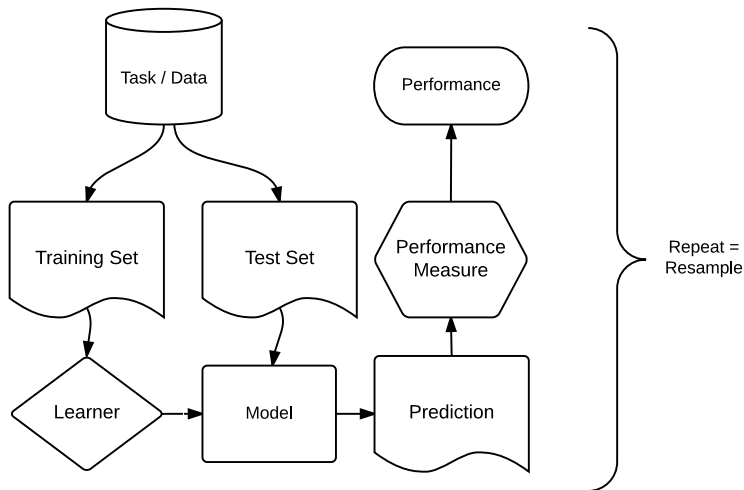
MOTIVATION: MLR

- Unified interface for the basic building blocks: tasks, learners, resampling, hyperparameters, ...
- Reflections: nearly all objects are queryable (i.e. you can ask them for their properties and program on them)
- The OO-structure allows many generic algorithms:
 - ▶ Bagging
 - ▶ Stacking
 - ▶ Feature Selection
 - ▶ ...
- Easily extensible via S3
 - ▶ Extension is not covered here, but explained in detail in the online tutorial
 - ▶ You do not need to understand S3 to use `mlr`
 - ▶ Wondering why we don't use S4? We care about code bloat and speed.

Section 3

BUILDING BLOCKS

BUILDING BLOCKS



- `mlr` objects: tasks, learners, measures, resampling instances.

TASK ABSTRACTION

- Tasks encapsulate data and meta-information about it
- Regression, classification, clustering, survival tasks
- Data is stored inside an environment to save memory

```
> task = makeClassifTask(data = iris, target = "Species")  
> print(task)
```

```
## Supervised task: iris  
## Type: classif  
## Target: Species  
## Observations: 150  
## Features:  
## numerics  factors  ordered  
##          4         0         0  
## Missings: FALSE  
## Has weights: FALSE  
## Has blocking: FALSE  
## Classes: 3  
##      setosa versicolor  virginica  
##      50         50         50  
## Positive class: NA
```


TASK ABSTRACTION: API I

```
> getTaskId(task)

## [1] "iris"

> str(getTaskData(task))

## 'data.frame': 150 obs. of  5 variables:
##  $ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
##  $ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
##  $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
##  $ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
##  $ Species      : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1
```

TASK ABSTRACTION: API II

```
> str(getTaskDescription(task))

## List of 11
## $ id           : chr "iris"
## $ type         : chr "classif"
## $ target       : chr "Species"
## $ size         : int 150
## $ n.feats      : Named int [1:3] 4 0 0
## $ ..- attr(*, "names")= chr [1:3] "numerics" "factors" "ordered"
## $ has.missings : logi FALSE
## $ has.weights  : logi FALSE
## $ has.blocking : logi FALSE
## $ class.levels : chr [1:3] "setosa" "versicolor" "virginica"
## $ positive     : chr NA
## $ negative     : chr NA
## $ - attr(*, "class")= chr [1:3] "TaskDescClassif" "TaskDescSupervised" "TaskD
```

TASK ABSTRACTION: API III

```
> getTaskSize(task)

## [1] 150

> getTaskFeatureNames(task)

## [1] "Sepal.Length" "Sepal.Width" "Petal.Length" "Petal.Width"

> getTaskTargetNames(task)

## [1] "Species"

> getTaskFormula(task)

## Species ~ .

> summary(getTaskTargets(task))

##      setosa versicolor  virginica
##      50          50          50
```

LEARNER ABSTRACTION I

- Internal structure of learners:
 - ▶ wrappers around `fit()` and `predict()` of the package
 - ▶ description of the parameter set
 - ▶ annotations
- Naming convention: `<tasktype>.<functionname>`

```
> makeLearner("classif.rpart")  
> makeLearner("regr.rpart")
```

- Adding custom learners is covered in the tutorial

LEARNER ABSTRACTION II

```
> lrn = makeLearner("classif.rpart")
> print(lrn)

## Learner classif.rpart from package rpart
## Type: classif
## Name: Decision Tree; Short name: rpart
## Class: classif.rpart
## Properties: twoclass,multiclass,missings,numerics,factors,ordered,prob,weigh
## Predict-Type: response
## Hyperparameters: xval=0
```

WHAT LEARNERS ARE AVAILABLE? I

CLASSIFICATION (72)

- LDA, QDA, RDA, MDA
- Trees and forests
- Boosting (different variants)
- SVMs (different variants)
- ...

CLUSTERING (8)

- K-Means
- EM
- DBscan
- X-Means
- ...

REGRESSION (53)

- Linear, lasso and ridge
- Boosting
- Trees and forests
- Gaussian processes
- ...

SURVIVAL (11)

- Cox-PH
- Cox-Boost
- Random survival forest
- Penalized regression
- ...

WHAT LEARNERS ARE AVAILABLE? II

We can explore them on the webpage – or ask `mlr`

WHAT LEARNERS ARE AVAILABLE? III

```
> # list all classification learners which can predict probabilities
> # and allow multiclass classification
> listLearners("classif",
+   properties = c("prob", "multiclass"))[1:5, c(-2, -5, -16)]
```

##		class	package	short.name	numerics	factors	ordered	
## 1		classif.avNNet	nnet	avNNet	TRUE	TRUE	FALSE	
## 2		classif.bdk	kohonen	bdk	TRUE	FALSE	FALSE	
## 3		classif.boosting	adabag,rpart	adabag	TRUE	TRUE	FALSE	
## 4		classif.cforest	party	cforest	TRUE	TRUE	TRUE	
## 5		classif.ctree	party	ctree	TRUE	TRUE	TRUE	
##		missings	weights	prob	oneclass	twoclass	multiclass	class.weights
## 1		FALSE	TRUE	TRUE	FALSE	TRUE	TRUE	FALSE
## 2		FALSE	FALSE	TRUE	FALSE	TRUE	TRUE	FALSE
## 3		TRUE	FALSE	TRUE	FALSE	TRUE	TRUE	FALSE
## 4		TRUE	TRUE	TRUE	FALSE	TRUE	TRUE	FALSE
## 5		TRUE	TRUE	TRUE	FALSE	TRUE	TRUE	FALSE

WHAT LEARNERS ARE AVAILABLE? IV

Get all applicable learners for a task

```
> listLearners(task)[1:5, c(-2, -5, -16)]
```

##		class		package	short.name	numerics	factors	ordered
## 1		classif.avNNet		nnet	avNNet	TRUE	TRUE	FALSE
## 2		classif.bdk		kohonen	bdk	TRUE	FALSE	FALSE
## 3		classif.boosting		adabag,rpart	adabag	TRUE	TRUE	FALSE
## 4		classif.cforest		party	cforest	TRUE	TRUE	TRUE
## 5		classif.ctree		party	ctree	TRUE	TRUE	TRUE
##		missings	weights	prob	oneclass	twoclass	multiclass	class.weights
## 1	FALSE	TRUE	TRUE	FALSE	TRUE	TRUE		FALSE
## 2	FALSE	FALSE	TRUE	FALSE	TRUE	TRUE		FALSE
## 3	TRUE	FALSE	TRUE	FALSE	TRUE	TRUE		FALSE
## 4	TRUE	TRUE	TRUE	FALSE	TRUE	TRUE		FALSE
## 5	TRUE	TRUE	TRUE	FALSE	TRUE	TRUE		FALSE

PARAMETER ABSTRACTION

- Extensive meta-information for hyperparameters available: storage type, constraints, defaults, dependencies
- Automatically checked for feasibility
- You can program on parameters!

```
> getParamSet(lrn)
```

##	Type	len	Def	Constr	Req	Tunable	Trafo
## minsplit	integer	-	20	1 to Inf	-	TRUE	-
## minbucket	integer	-	-	1 to Inf	-	TRUE	-
## cp	numeric	-	0.01	0 to 1	-	TRUE	-
## maxcompete	integer	-	4	0 to Inf	-	TRUE	-
## maxsurrogate	integer	-	5	0 to Inf	-	TRUE	-
## usesurrogate	discrete	-	2	0,1,2	-	TRUE	-
## surrogatestyle	discrete	-	0	0,1	-	TRUE	-
## maxdepth	integer	-	30	1 to 30	-	TRUE	-
## xval	integer	-	10	0 to Inf	-	FALSE	-
## parms	untyped	-	-	-	-	TRUE	-

LEARNER ABSTRACTION: API

```
> lrn$properties

## [1] "twoclass"      "multiclass" "missings"    "numerics"    "factors"
## [6] "ordered"      "prob"       "weights"

> getHyperPars(lrn)

## $xval
## [1] 0

> lrn = setHyperPars(lrn, cp = 0.3)
> lrn = setPredictType(lrn, "prob")
> lrn = setPredictThreshold(lrn, 0.7);
```

PERFORMANCE MEASURES

- Performance measures evaluate the predictions a test set and aggregate them over multiple in resampling iterations
- 22 classification, 10 regression, 5 cluster, 1 survival
- Internally: performance function, default aggregation function and annotations
- Adding custom measures is covered in the tutorial

```
> print(mmce)

## Name: Mean misclassification error
## Performance measure: mmce
## Properties: classif,classif.multi,req.pred,req.truth
## Minimize: TRUE
## Best: 0; Worst: 1
## Aggregated by: test.mean
## Note:
```

WHAT MEASURES ARE AVAILABLE?

We can explore them on the webpage – or ask `mlr`

```
> listMeasures("classif")
```

```
## [1] "timepredict"    "gmean"          "acc"
## [4] "auc"            "ber"            "fn"
## [7] "fp"            "fnr"            "gpr"
## [10] "featperc"       "ppv"            "fpr"
## [13] "mmce"           "timeboth"       "npv"
## [16] "timetrain"      "fdr"            "tnr"
## [19] "mcc"            "bac"            "tpr"
## [22] "tn"             "f1"             "tp"
## [25] "multiclass.auc" "brier"
```

```
> listMeasures(task)
```

```
## [1] "timepredict"    "acc"            "ber"
## [4] "featperc"       "mmce"           "timeboth"
## [7] "timetrain"      "multiclass.auc"
```

R EXAMPLE

Training and prediction

RESAMPLING ABSTRACTION I

- Procedure: Train, Predict, Eval, Repeat.
- Aim: Estimate expected model performance.
 - ▶ Hold-Out
 - ▶ Cross-validation (normal, repeated)
 - ▶ Bootstrap (OOB, B632, B632+)
 - ▶ Subsampling
 - ▶ Stratification
 - ▶ Blocking
- Instantiate it or not (= create data split indices)

```
> rdesc = makeResampleDesc("CV", iters = 3)
> rin = makeResampleInstance(rdesc, task = task)
> str(rin$train.inds)
```

```
## List of 3
## $ : int [1:100] 68 75 78 11 31 145 63 110 123 9 ...
## $ : int [1:100] 68 75 27 11 41 128 63 110 58 98 ...
## $ : int [1:100] 27 78 41 31 145 128 123 58 9 46 ...
```

RESAMPLING ABSTRACTION II

RESAMPLING A LEARNER

- Measures on test (or train) sets
- Returns aggregated values, predictions and some useful extra information

```
> lrn = makeLearner("classif.rpart")  
> rdesc = makeResampleDesc("CV", iters = 3)  
> measures = list(mmce, timetrain)  
> r = resample(lrn, task, rdesc, measures = measures)
```

- For the lazy

```
> r = crossval(lrn, task, iters = 3, measures = measures)
```


RESAMPLING ABSTRACTION III

```
> print(r)

## Resample Result
## Task: iris
## Learner: classif.rpart
## mmce.aggr: 0.08
## mmce.mean: 0.08
## mmce.sd: 0.04
## timetrain.aggr: 0.01
## timetrain.mean: 0.01
## timetrain.sd: 0.00
## Runtime: 0.0400295
```

RESAMPLING ABSTRACTION IV

```
> names(r)

## [1] "learner.id"      "task.id"          "measures.train"
## [4] "measures.test"   "aggr"             "pred"
## [7] "models"          "err.msgs"         "extract"
## [10] "runtime"

> r$measures.test

##   iter mmce timetrain
## 1    1 0.08    0.009
## 2    2 0.04    0.004
## 3    3 0.12    0.004

> r$aggr

##           mmce.test.mean timetrain.test.mean
##                0.080000                0.005667
```

RESAMPLING ABSTRACTION V

```
> head(as.data.frame(r$pred))
```

##	id	truth	response	iter	set
## 1	1	setosa	setosa	1	test
## 2	6	setosa	setosa	1	test
## 3	7	setosa	setosa	1	test
## 4	9	setosa	setosa	1	test
## 5	15	setosa	setosa	1	test
## 6	21	setosa	setosa	1	test

CONFIGURING THE PACKAGE

- What to do when training fails? error, warn, or be quiet?
 - You don't want to stop in complex loops like benchmark
 - `FailureModel` is created that predicts NAs
- Show verbose info messages?
- What if parameters are not described in learner?
- `?configureMlr` sets global flags and can be overwritten for individual learners

Section 4

BENCHMARKING AND MODEL COMPARISON

BENCHMARKING AND MODEL COMPARISON I

BENCHMARKING

- Comparison of multiple models on multiple data sets
- Aim: Find best learners for a data set or domain, learn about learner characteristics, ...

BENCHMARKING AND MODEL COMPARISON II

BENCHMARKING IN MLR

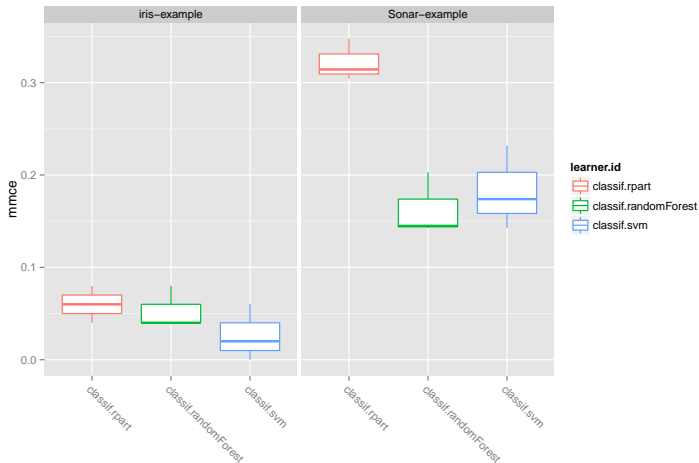
- Train and test sets are synchronized, i.e. all learners see the same data splits
- Can be done in parallel (see later)
- Can be combined with feature selection / tuning / nested resampling (see later)
- Results stored in well-defined container object, with getters and converters
- We are working on standard analysis tools

BENCHMARKING AND MODEL COMPARISON III

```
> library(mlr)
> # lets try a couple of methods on some (mlr example) tasks
>
> # these are predefined in mlr for toying around:
> tasks = list(iris.task, sonar.task)
>
> learners = list(
+   makeLearner("classif.rpart"),
+   makeLearner("classif.randomForest", ntree = 500),
+   makeLearner("classif.svm")
+ )
>
> rdesc = makeResampleDesc("CV", iters = 3)
> set.seed(1)
> br = benchmark(learners, tasks, rdesc)
```


BENCHMARKING AND MODEL COMPARISON IV

```
> plotBenchmarkResult(br)
```



BENCHMARKING AND MODEL COMPARISON V

```
> getBMRAggrPerformances(br, as.df = TRUE)
```

##	task.id	learner.id	mmce.test.mean
## 1	iris-example	classif.rpart	0.06000
## 2	iris-example	classif.randomForest	0.05333
## 3	iris-example	classif.svm	0.02667
## 4	Sonar-example	classif.rpart	0.32215
## 5	Sonar-example	classif.randomForest	0.16356
## 6	Sonar-example	classif.svm	0.18288

BENCHMARKING AND MODEL COMPARISON VI

```
> getBMRPerformances(br, as.df = TRUE)
```

##	task.id	learner.id	iter	mmce
## 1	iris-example	classif.rpart	1	0.0600
## 2	iris-example	classif.rpart	2	0.0400
## 3	iris-example	classif.rpart	3	0.0800
## 4	iris-example	classif.randomForest	1	0.0400
## 5	iris-example	classif.randomForest	2	0.0400
## 6	iris-example	classif.randomForest	3	0.0800
## 7	iris-example	classif.svm	1	0.0000
## 8	iris-example	classif.svm	2	0.0200
## 9	iris-example	classif.svm	3	0.0600
## 10	Sonar-example	classif.rpart	1	0.3478
## 11	Sonar-example	classif.rpart	2	0.3043
## 12	Sonar-example	classif.rpart	3	0.3143
## 13	Sonar-example	classif.randomForest	1	0.2029
## 14	Sonar-example	classif.randomForest	2	0.1449
## 15	Sonar-example	classif.randomForest	3	0.1429
## 16	Sonar-example	classif.svm	1	0.2319
## 17	Sonar-example	classif.svm	2	0.1739
## 18	Sonar-example	classif.svm	3	0.1429

BENCHMARKING AND MODEL COMPARISON VII

```
> head(getBMRPredictions(br, as.df = TRUE), 10)
```

##	task.id	learner.id	id	truth	response	iter	set
## 1	iris-example	classif.rpart	1	setosa	setosa	1	test
## 2	iris-example	classif.rpart	3	setosa	setosa	1	test
## 3	iris-example	classif.rpart	12	setosa	setosa	1	test
## 4	iris-example	classif.rpart	17	setosa	setosa	1	test
## 5	iris-example	classif.rpart	22	setosa	setosa	1	test
## 6	iris-example	classif.rpart	24	setosa	setosa	1	test
## 7	iris-example	classif.rpart	25	setosa	setosa	1	test
## 8	iris-example	classif.rpart	26	setosa	setosa	1	test
## 9	iris-example	classif.rpart	31	setosa	setosa	1	test
## 10	iris-example	classif.rpart	34	setosa	setosa	1	test

Section 5

HYPERPARAMETER TUNING

HYPERPARAMETER TUNING I

TUNING

- Used to find “best” hyperparameters for a method in a data-dependent way
- Essential for some methods, e.g. SVMs

TUNING IN MLR

- General procedure: Tuner proposes param point, eval by resampling, feedback value to tuner
- Multiple tuners through exactly the same interface
- All evals and more info is logged into `OptPath` object

HYPERPARAMETER TUNING II

GRID SEARCH

- Basic method: Exhaustively try all combinations of finite grid
- Inefficient, combinatorial explosion
- Searches large, irrelevant areas
- Reasonable for continuous parameters?
- Still often default method

RANDOM SEARCH

- Randomly draw parameters
- `mlr` supports all types and dependencies
- Scales better than grid search, easily extensible

R EXAMPLE

Tuning

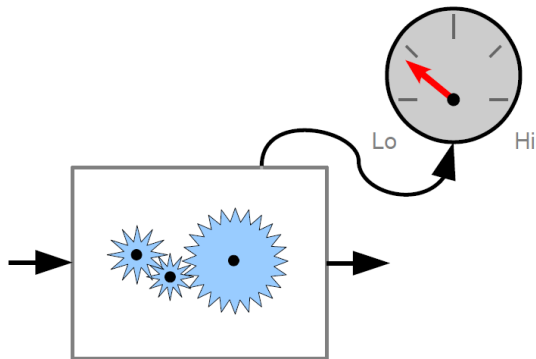
AUTOMATIC MODEL SELECTION

PRIOR APPROACHES:

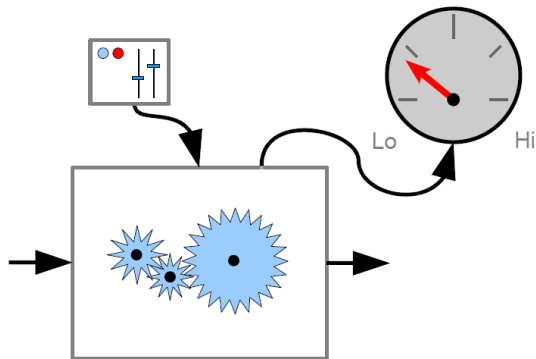
- Looking for the silver bullet model
 - ~ Failure
- Exhaustive benchmarking / search
 - ~ Per data set: too expensive
 - ~ Over many: contradicting results
- Meta-Learning:
 - ~ Failure
 - ~ Usually not for preprocessing / hyperparameters

GOAL: Data dependent + Automatic + Efficient

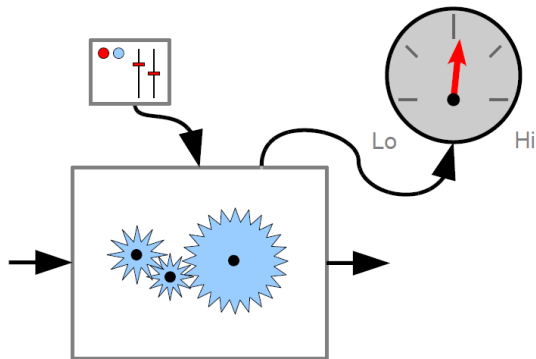
BLACK-BOX-PERSPECTIVE IN CONFIGURATION



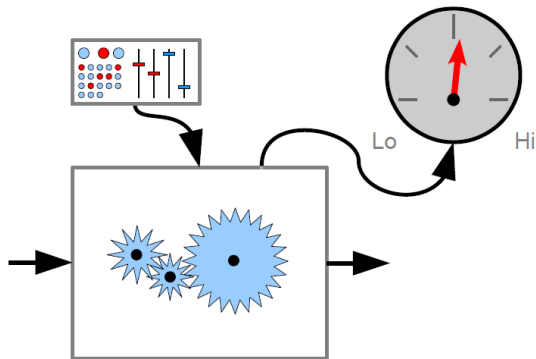
BLACK-BOX-PERSPECTIVE IN CONFIGURATION



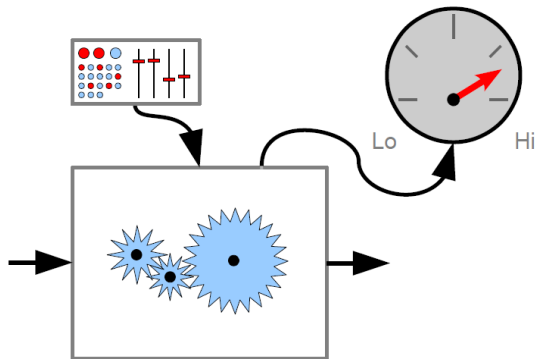
BLACK-BOX-PERSPECTIVE IN CONFIGURATION



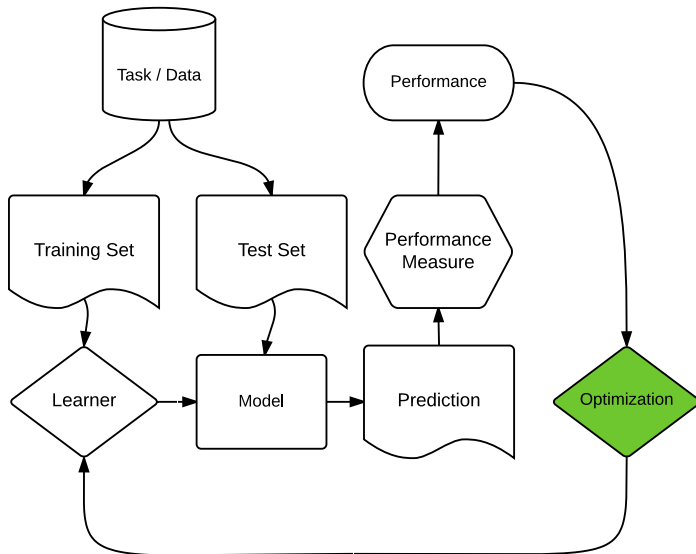
BLACK-BOX-PERSPECTIVE IN CONFIGURATION



BLACK-BOX-PERSPECTIVE IN CONFIGURATION



ADAPTIVE TUNING

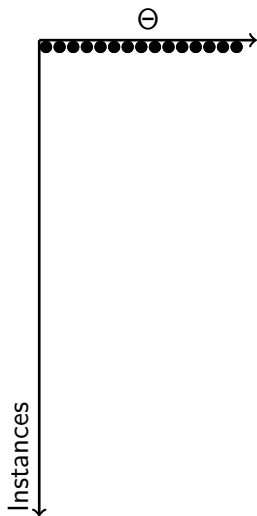


GENERAL ALGORITHM CONFIGURATION

- Assume a (parametrized) algorithm a
- Parameter space $\theta \in \Theta$
might be discrete and dependent / hierarchical
- Stochastic generating process for instances $i \sim P$, where we draw i.i.d. from.
- Run algorithm a on i and measure performance
 $f(i, \theta) = \text{run}(i, a(\theta))$
- Objective: $\min_{\theta \in \Theta} E_P[f(i, \theta)]$
- No derivative for $f(\cdot, \theta)$, black-box
- f is stochastic / noisy
- f is likely expensive to evaluate
- Consequence: very hard problem

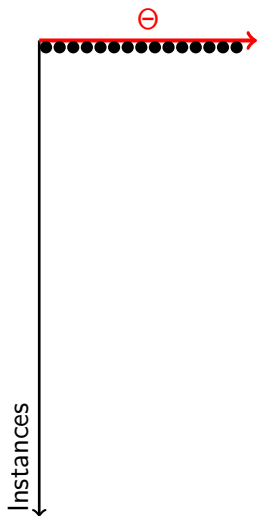
\leadsto RACING OR MODEL-BASED / BAYESIAN OPTIMIZATION

IDEA OF (F-)RACING



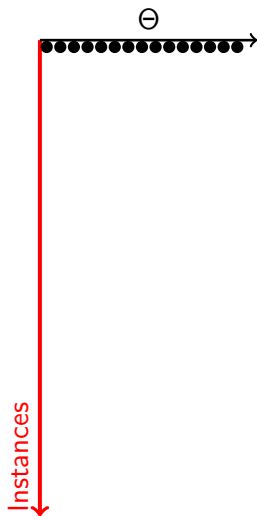
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



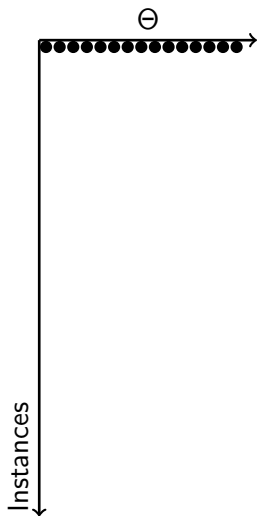
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



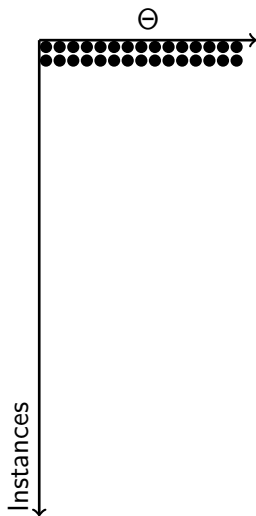
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



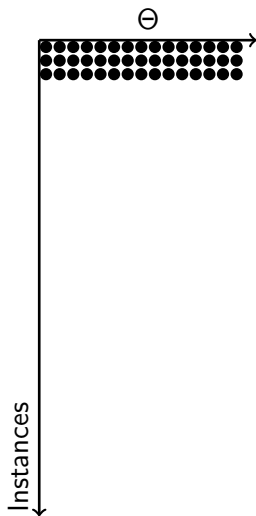
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



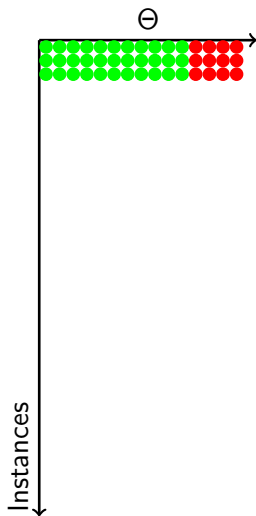
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



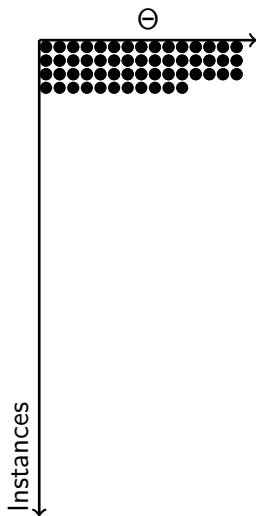
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



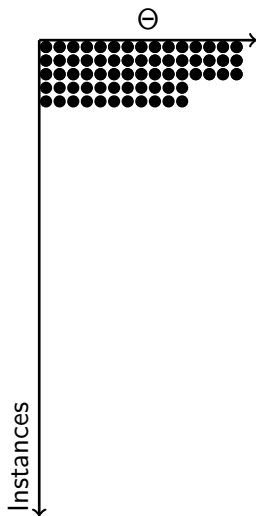
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



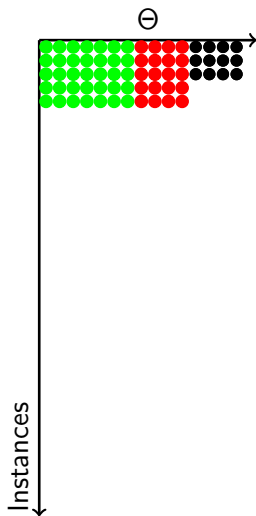
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



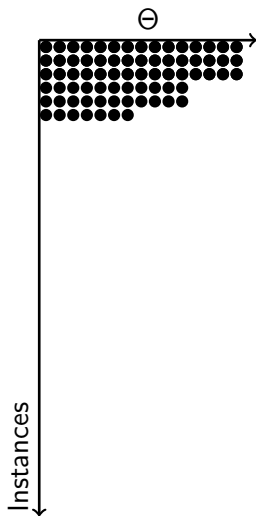
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



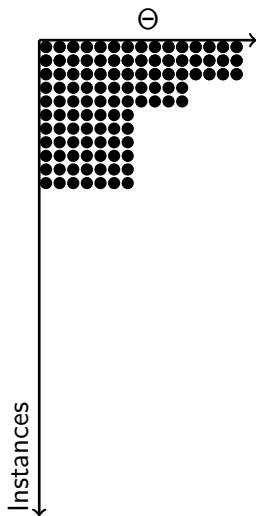
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



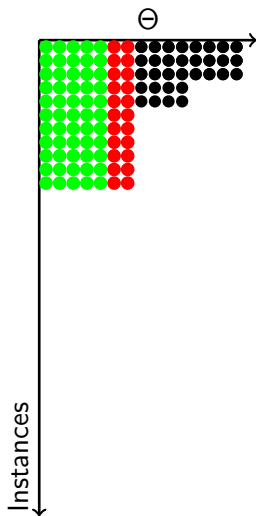
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



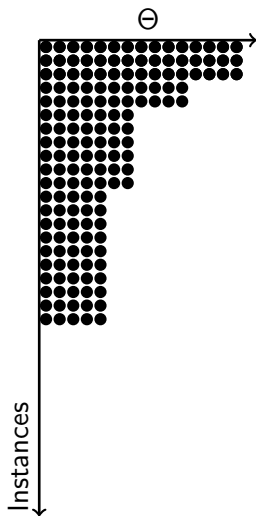
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



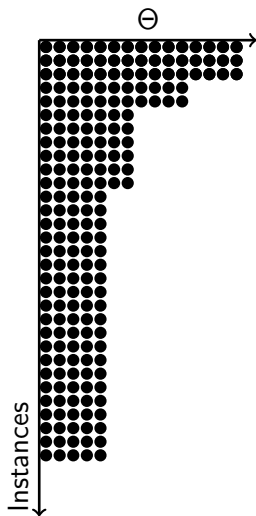
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



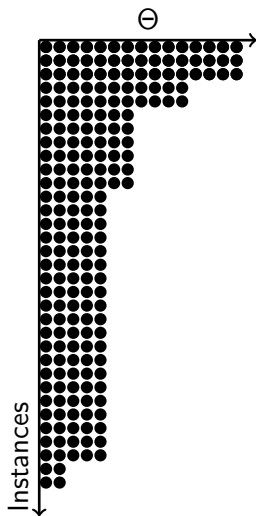
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



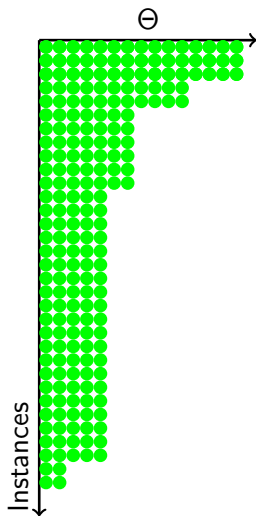
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



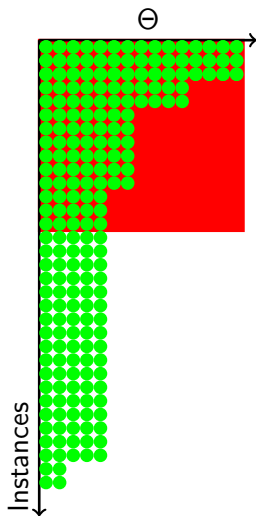
- Write down all candidate solutions
- Iterate the following till budget exhausted
 - One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



- Write down all candidate solutions
- Iterate the following till budget exhausted
- One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF (F-)RACING



- Write down all candidate solutions
- Iterate the following till budget exhausted
- One “generation”
 - ▶ Evaluate all candidates on an instance, and another, . . .
 - ▶ After some time, compare candidates via statistical test, e.g., Friedman test with post-hoc analysis for pairs
 - ▶ Remove outperformed candidates
- Output: Remaining candidates
- Yes, the testing completely ignores “sequentiality” and is somewhat heuristic.

IDEA OF ITERATED F-RACING

WHAT MIGHT BE PROBLEMATIC?

- We might have many or an infinite number of candidates

ITERATED RACING

- Have a stochastic model to draw candidates from in every generation
- For each parameter: Univariate, independent distribution (factorized joint distribution)
- Sample distributions centered at “elite” candidates from previous generation(s)
- Reduce distributions’ width / variance in later generations for convergence

IDEA OF ITERATED F-RACING

WHATS GOOD ABOUT THIS

- Very simple and generic algorithm
- Can easily be parallelized
- A nice R package exists: `irace`¹

WHAT MIGHT BE NOT SO GOOD

- Quite strong (wrong?) assumptions in the probability model
- Sequential model-based optimization is probably more efficient
(But be careful: Somewhat my personal experience and bias,
as not so many large scale comparisons exist)

¹Lopez-Ibanez et al, “The irace package, Iterated Race for Automatic Algorithm Configuration. Technical Report TR/IRIDIA/2011-004, IRIDIA, Université libre de Bruxelles, Belgium, 2011.”

Section 6

FEATURE SELECTION

FEATURE SELECTION I

- Reduce dimensionality, increase interpretability and predictive performance

- Concepts:

FILTER: Preliminary step, independent from model

WRAPPER: Wrapped around model fit which is iteratively scored

EMBEDDED: Model has feature selection embedded, e.g. lasso regression

mlr supports all of these, but we do not have enough time today for details.

Section 7

MLR LEARNER WRAPPERS

MLR LEARNER WRAPPERS I

WHAT?

- Extend the functionality of learners by adding an `mlr` wrapper to them
- The wrapper hooks into the `train` and `predict` of the base learner and extends it
- This way, you can create a new `mlr` learner with extended functionality
- Hyperparameter definition spaces get joined!

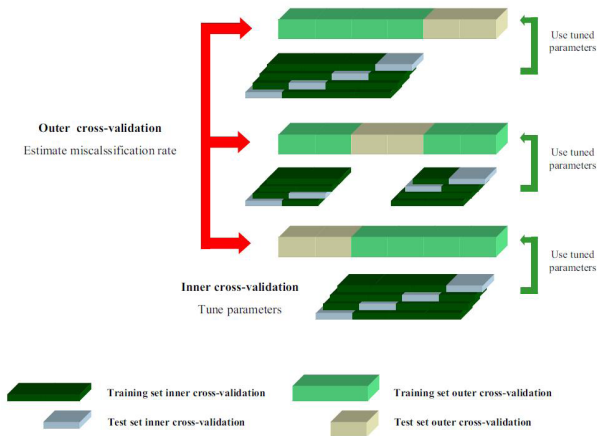
MLR LEARNER WRAPPERS II

AVAILABLE WRAPPERS

- PREPROCESSING: PCA, normalization (z-transformation)
- PARAMETER TUNING: grid, optim, random search, genetic algorithms, CMAES, iRace, MBO
- FILTER: correlation- and entropy-based, χ^2 -test, mRMR, ...
- FEATURE SELECTION: (floating) sequential forward/backward, exhaustive search, genetic algorithms, ...
- IMPUTE: dummy variables, imputations with mean, median, min, max, empirical distribution or other learners
- BAGGING to fuse learners on bootstrapped samples
- STACKING to combine models in heterogenous ensembles
- OVER- AND UNDERSAMPLING for unbalanced classification

NESTED RESAMPLING

- Using the TuningWrapper or FeatureSelectionWrapper allows to enable nested resampling
- Ensures **unbiased** results for model optimization
- Everything else is statistically unsound



R EXAMPLE

Complex tuning example

Section 8

PARALLELIZATION

PARALLELIZATION I

- We use our own package: `parallelMap`
- Initialize a backend with `parallelStart`
- Stop with `parallelStop`

```
> parallelStart("multicore")  
> benchmark(...)  
> parallelStop()
```

- Backends: `local`, `multicore`, `socket`, `mpi` and `BatchJobs`
- The latter means support for: makeshift SSH-clusters and HPC schedulers like SLURM, Torque/PBS, SGE or LSF
- The first loop which is marked as parallel executable will be automatically parallelized

PARALLELIZATION II

PARALLELIZATION LEVELS

- Which loop to parallelize depends on number of iterations
- Levels allow fine grained control over the parallelization
 - ▶ `mlr.resample`: Each resampling iteration (a train / test step) is a parallel job.
 - ▶ `mlr.benchmark`: Each experiment “run this learner on this data set” is a parallel job.
 - ▶ `mlr.tuneParams`: Each evaluation in hyperparameter space “resample with these parameter settings” is a parallel job. How many of these can be run independently in parallel depends on the tuning algorithm.
 - ▶ `mlr.selectFeatures`: Each evaluation in feature space “resample with this feature subset” is a parallel job.

PARALLELIZATION III

```
> lrns = list(makeLearner("classif.rpart"), makeLearner("classif.svm"))
> rdesc = makeResampleDesc("Bootstrap", iters = 100)
>
> parallelStart("multicore", 8)
```

```
## Starting parallelization in mode=multicore with cpus=8.
```

```
> benchmark(learners = lrns, tasks = iris.task, resamplings = rdesc)
```

```
## Mapping in parallel: mode = multicore; cpus = 8; elements = 2.
```

```
##      task.id    learner.id mmce.test.mean
## 1 iris-example classif.rpart      0.05689
## 2 iris-example  classif.svm      0.04234
```

```
> parallelStop()
```

```
## Stopped parallelization. All cleaned up.
```

PARALLELIZATION IV

Parallelize the bootstrap instead:

```
> parallelStart("multicore", 8, level = "mlr.resample")

## Starting parallelization in mode=multicore with cpus=8.

> benchmark(learners = lrns, tasks = iris.task, resamplings = rdesc)

## Mapping in parallel: mode = multicore; cpus = 8; elements = 100.
## Mapping in parallel: mode = multicore; cpus = 8; elements = 100.

##      task.id    learner.id mmce.test.mean
## 1 iris-example classif.rpart      0.05817
## 2 iris-example  classif.svm       0.04126

> parallelStop()

## Stopped parallelization. All cleaned up.
```


Section 9

VISUALIZATIONS

VISUALIZATIONS

- We use ggplot2 and interactive ggvis as a standard, if possible
- Some plots use Viper Charts as backend (cost curves, lift charts, ...)
- GSOC project 2015 with Zach Jones
 - ▶ Demo plots for models in teaching
 - ▶ ROC curves
 - ▶ Threshold vs. Performance
 - ▶ Partial dependency plot
 - ▶ Learning curves

Visualizations

Section 10

OPENML


OPENML-R-PACKAGE I

Caution: Work in progress

OPENML?

- Main idea: Make ML experiments reproducible and most parts computer-readable
- Share everything
- Enrich with meta-information
- Later: Mine the results, meta-learn on it

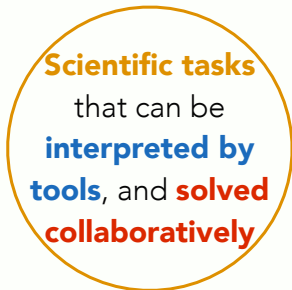
1 minute intro



Data from
various sources
**analysed and
organised online**
for easy access

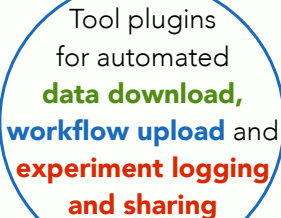
Scientists can **broadcast data**, explaining the challenge that needs to be addressed. OpenML will (for known data formats) **automatically analyze the data**, compute data characteristics, **annotate and index it for easy search**

1 minute intro



Tasks are **realtime (collaborative) data mining challenges**, allowing anyone to build on previous results. OpenML creates **machine-readable descriptions** so that tools can **automatically download data**, use the correct procedures, and **upload all results online**.


1 minute intro



Tool plugins
for automated
data download,
workflow upload and
experiment logging
and sharing

Flows are implementations of algorithms, workflows, or scripts **solving OpenML tasks**. OpenML keeps track of **flow details and versioning**, **organizes all their results** for easy comparison, even across tools.

1 minute intro



Experiments
auto-uploaded,
linked to **data, flows**
and **authors**, and
organised for easy
reuse

Runs contain the results that **flows** obtained on specific tasks. Runs are **fully reproducible**, linked to the underlying data, tasks, flows and authors. OpenML **organizes all results online** for **discovery, comparison and reuse**

OPENML-R-PACKAGE I

Let's visit website and project page

OPENML-R-PACKAGE II

`https://github.com/openml/r`

CURRENT API IN R

- Explore data and tasks
- Download data and tasks
- Register learners
- Upload runs
- Explore your own and other people's results

Already nicely connected to `mlr`!

OPENML: EXPLORE AND SELECT DATA I

```
> library(OpenML)
> listOMLDataSets()[1:3, 1:9]
```

Downloading from 'http://www.openml.org/api/v1/data/list' to '<mem>'

##	did	status	name	MajorityClassSize	MaxNominalAttDistinctValues	
## 1	1	active	anneal	684	10	
## 2	2	active	anneal	684	9	
## 3	3	active	kr-vs-kp	1669	3	
##			MinorityClassSize	NumBinaryAtts	NumberOfClasses	NumberOfFeatures
## 1			0	14	6	39
## 2			0	7	6	39
## 3			1527	34	2	37

OPENML: EXPLORE AND SELECT DATA II

```
> listOMLTasks()[1:3, c(1:5, 10:11)]
```

Downloading from 'http://www.openml.org/api/v1/task/list' to '<mem>'

##	task.id	task.type	did	status	name
## 1	1	Supervised Classification	1	active	anneal
## 2	2	Supervised Classification	2	active	anneal
## 3	3	Supervised Classification	3	active	kr-vs-kp
##	MajorityClassSize MaxNominalAttDistinctValues				
## 1		684		10	
## 2		684		9	
## 3		1669		3	

OPENML: DOWNLOAD A DATA SET

```
> # uses built in caching from disk
> getOMLDataSet(6)

## Data '6' file 'description.xml' found in cache.
## Data '6' file 'dataset.arff' found in cache.

##
## Data Set "letter" :: (Version = 1, OpenML ID = 6)
##   Default Target Attribute: class
```

OPENML: DOWNLOAD A TASK I

```
> # uses built in caching from disk
> oml.task = getOMLTask(1)

## Task '1' file 'task.xml' found in cache.
## Task '1' file 'datasplits.arff' found in cache.
## Data '1' file 'description.xml' found in cache.
## Data '1' file 'dataset.arff' found in cache.
```

OPENML: DOWNLOAD A TASK II

```
> print(oml.task)

##
## OpenML Task 1 :: (Data ID = 1)
##   Task Type           : Supervised Classification
##   Data Set             : anneal :: (Version = 2, OpenML ID = 1)
##   Target Feature(s)    : class
##   Tags                 : basic, study_1, study_7, under100k, under1m
##   Estimation Procedure : Stratified crossvalidation (1 x 10 folds)
```


OPENML: RUN A TASK

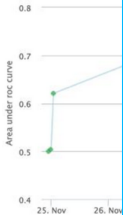
```
> lrn = makeLearner("classif.rpart")
> res = runTaskMlr(oml.task, lrn)

## Task: OpenML-Task-1, Learner: classif.rpart
## [Resample] cross-validation iter: 1
## [Resample] cross-validation iter: 2
## [Resample] cross-validation iter: 3
## [Resample] cross-validation iter: 4
## [Resample] cross-validation iter: 5
## [Resample] cross-validation iter: 6
## [Resample] cross-validation iter: 7
## [Resample] cross-validation iter: 8
## [Resample] cross-validation iter: 9
## [Resample] cross-validation iter: 10
## [Resample] Result: acc.test.mean=0.977, timetrain.test.sum=1.21, timepredict.test.sum=1.39
```

OPENML: UPLOAD LEARNER AND PREDICTIONS

```
> impl = createOpenMLImplementationForMlrLearner(lrn)
> uploadOpenMLImplementation(impl, session.hash = hash)
> uploadOpenMLRun(oml.task, lrn, impl, pred, hash)
```

Towards OpenML in education



frontier Joaquin Vanschoren
Olav Bunte Stephan Oostveen
Jorn Engelbart Mathijs van Lier
Stefan Majoer



Rogier Beckers

@RogierBeckers



Follow

Het bewijs dat ik studeer op zondag!
“@joavanschoren: #Machinelearning students
on a #collaborative data mining ”

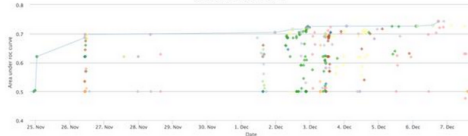
View translation

Lauradorp, Landgraaf



Contributions over time

every point is a task, click for details



frontier Joaquin Vanschoren Perry van Wesel Jose Melo Jos Mangnus Daan Peters Tom Becht Kevin Jacobs Koen Engelen
Olav Bunte Stephan Oostveen Ray van den Hurk Schwester Kogowski Ky-Anh Tran Edgar Salas Thomas Tiel Groenesteghe
Jorn Engelbart Mathijs van Lier Henry He Riche Brundensieker Hugo Lape Stanley Clark Christoforos Bousmoulas Rogier Beckers
Stefan Majoer

RETWEETS

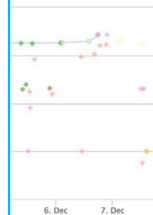
2

FAVORITES

2



9:48 PM - 7 Dec 2014



Kevin Jacobs Koen Engelen
Tiel Groenesteghe
pukouvalas Rogier Beckers

Section 11

THE END

THERE IS MORE ...

- Regular cost-sensitive learning (class-specific costs)
- Cost-sensitive learning (example-dependent costs)
- Multi-Label learning
- Model-based optimization
- Multi-criteria optimization
- OpenML
- ...

OUTLOOK

WE ARE WORKING ON

- Even better tuning system
- More interactive plots
- Large-Scale SVM ensembles
- Time-Series tasks
- Better benchmark analysis
- ...

Thanks!