

Benchmark Experiments – Comparing Statistical Learning Algorithms

Manuel J. A. Eugster

Institut für Statistik
Ludwig-Maximilians-Universität München

September 28th 2009, TU Dortmund

Benchmark experiments

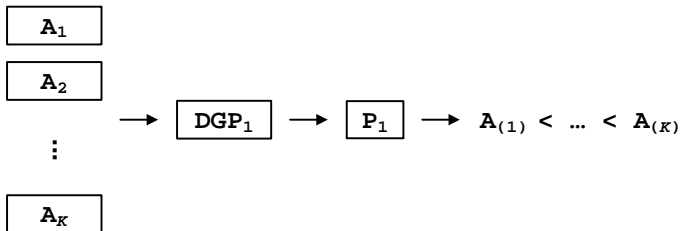
“Statistical decathlon”:

Empirical investigations with the aim of comparing and ranking (learning) algorithms with respect to certain performance measures.

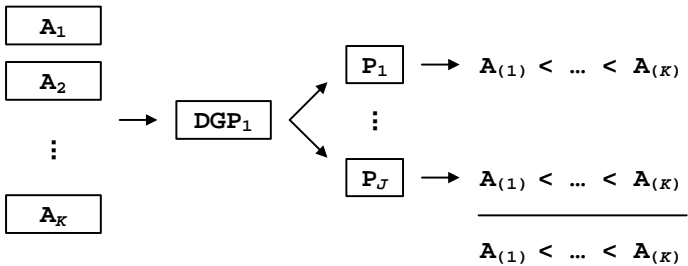
In essence:

Draw observations from theoretically intractable performance distributions of the learning algorithms, which are defined by either artificial data generating processes or empirical distribution functions from one or more data sets.

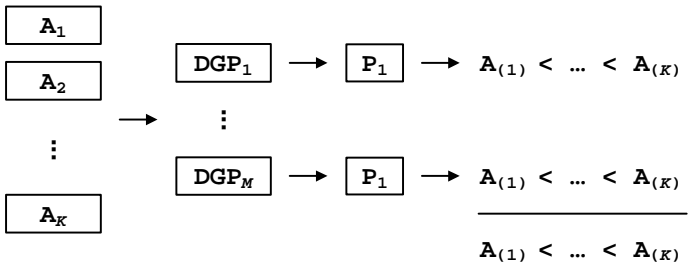
Scenario 1:



Scenario 2:



Scenario 3:



Abstract levels of benchmark experiments:

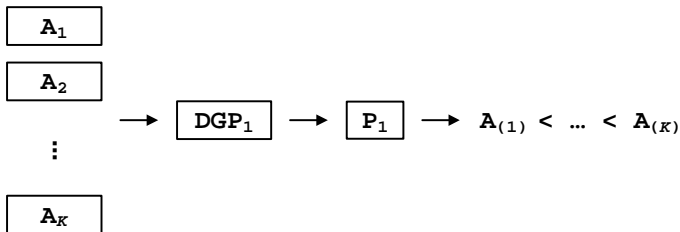
Setup: define design of a benchmark experiment; data sets, candidate algorithms, performance measures and resampling strategy.

Execution: execute Setup; parallel computation, monitoring, sequential and adaptive procedures.

Analysis: exploratory and inferential analyses of the raw performance measures; main objective is a statistically correct order of the candidate algorithms.

Scenario 1

“From the theoretical framework to a concrete order of some candidate algorithms”



General framework

Data generating process:

Given a data generating process DGP , we draw B independent and identically distributed learning samples:

$$\mathcal{L}^1 = \{z_1^1, \dots, z_n^1\} \sim DGP$$

\vdots

$$\mathcal{L}^B = \{z_1^B, \dots, z_n^B\} \sim DGP$$

Candidate algorithms:

There are $K > 1$ algorithms a_k ($k = 1, \dots, K$) with the function $a_k(\cdot | \mathcal{L}^b)$ the fitted model on the learning sample \mathcal{L}^b ($b = 1, \dots, B$):

$$a_k(\cdot | \mathcal{L}^b) \sim \mathcal{A}_k(DGP)$$

Performance measure:

Performance of algorithm a_k when provided with the learning sample \mathcal{L}^b is measured by a scalar function p :

$$p_{kb} = p(a_k, \mathcal{L}^b) \sim P_k = P_k(DGP)$$

Supervised learning problems

Observations are of the form $z = (y, x)$ where y denotes the response variable and x a vector of input variables.

Algorithms are learners $\hat{y} = a_k(x | \mathcal{L}^b)$ which, given the input variables, predict the unknown response variable.

Performance is defined by some functional μ of the distribution of a loss function $L(y, \hat{y})$ which measures the discrepancy between true response y and predicted value \hat{y} :

$$p_{kb} = p(a_k, \mathcal{L}^b) = \mu(L(y, a_k(x | \mathcal{L}^b))) \sim P_k(DGP)$$

Performance is defined by some functional μ of the distribution of a loss function $L(y, \hat{y})$ which measures the discrepancy between true response y and predicted value \hat{y} :

$$p_{kb} = p(a_k, \mathcal{L}^b) = \mu(L(y, a_k(x | \mathcal{L}^b))) \sim P_k(DGP)$$

Misclassification:

The loss function is

$$L(y, \hat{y}) = \begin{cases} 0, & y = \hat{y} \\ 1, & \text{otherwise} \end{cases}$$

and the functional μ is the expectation E :

$$p_{kb} = E_{a_k} E_{z=(y,x)} L(y, a_k(x | \mathcal{L}^b))$$

with $z = (y, x)$ drawn from the same data generating process as the observations in the learning sample \mathcal{L}^b .

If it is not possible to calculate μ analytically, the empirical analogue $\mu_{\mathfrak{T}}$ based on a test sample $\mathfrak{T} \sim DGP$ has to be used:

$$\hat{p}_{kb} = \hat{p}(a_k, \mathfrak{L}^b) = \mu_{\mathfrak{T}}(L(y, a_k(x | \mathfrak{L}^b))) \sim \hat{P}_k(DGP)$$

If it is not possible to calculate μ analytically, the empirical analogue $\mu_{\mathfrak{T}}$ based on a test sample $\mathfrak{T} \sim DGP$ has to be used:

$$\hat{p}_{kb} = \hat{p}(a_k, \mathfrak{L}^b) = \mu_{\mathfrak{T}}(L(y, a_k(x | \mathfrak{L}^b))) \sim \hat{P}_k(DGP)$$

Misclassification:

The empirical functional $\mu_{\mathfrak{T}}$ is the expectation E with respect to the test sample \mathfrak{T} :

$$\hat{p}_{kb} = E_{a_k} E_{z=(y,x)} L(y, a_k(x | \mathfrak{L}^b))$$

with $z = (y, x) \in \mathfrak{T}$.

Data generating process in a real world problem:

Given is one learning sample $\mathcal{L} \sim \mathcal{Z}_n$ of n observations from some distribution function \mathcal{Z} . Then the empirical distribution function $\hat{\mathcal{Z}}_n$ covers all knowledge about the data generating process:
 $DGP = \hat{\mathcal{Z}}_n$.

Learning samples are defined by bootstrapping: $\mathcal{L}^b \sim \hat{\mathcal{Z}}_n$.

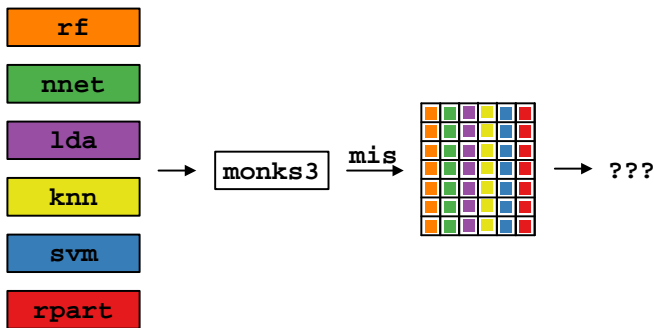
Test samples are defined by the out-of-bootstrap observations:
 $\mathcal{T}^b = \mathcal{L} \setminus \mathcal{L}^b$.

Empirical performance distributions:

For each candidate algorithm a_k the empirical performance distribution $\hat{P}_k(\mathcal{L})$ is estimated based on the \hat{p}_{kb} .

Exploratory data analysis tools and formal inference procedures to compare them and calculate a statistically correct order.

Exemplar benchmark experiment



(1) classification problem `monks3` (554 observations and 6 nominal attributes); (2) algorithms {`randomForest`, `nnet`, `lda`, `knn`, `svm`, `rpart`}; (3) misclassification performance measure; (4) 250 bootstrap learning samples; (5) out-of-bootstrap test samples;

Common practise

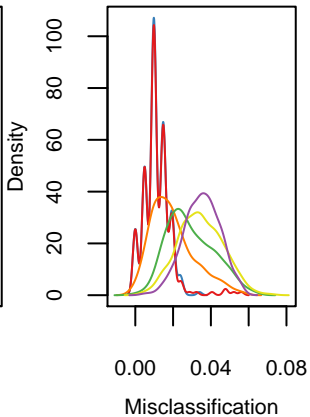
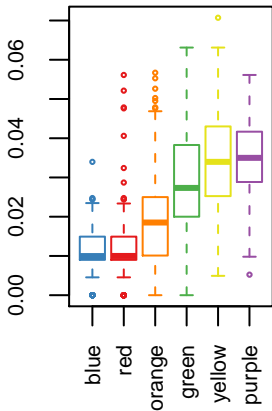
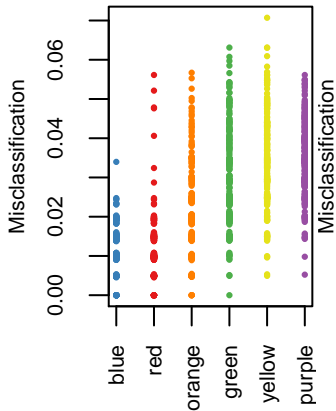
Comparison based on some summary statistics: algorithm a_1 is better than algorithm a_2 iff $\phi(\hat{P}_1) < \phi(\hat{P}_2)$.

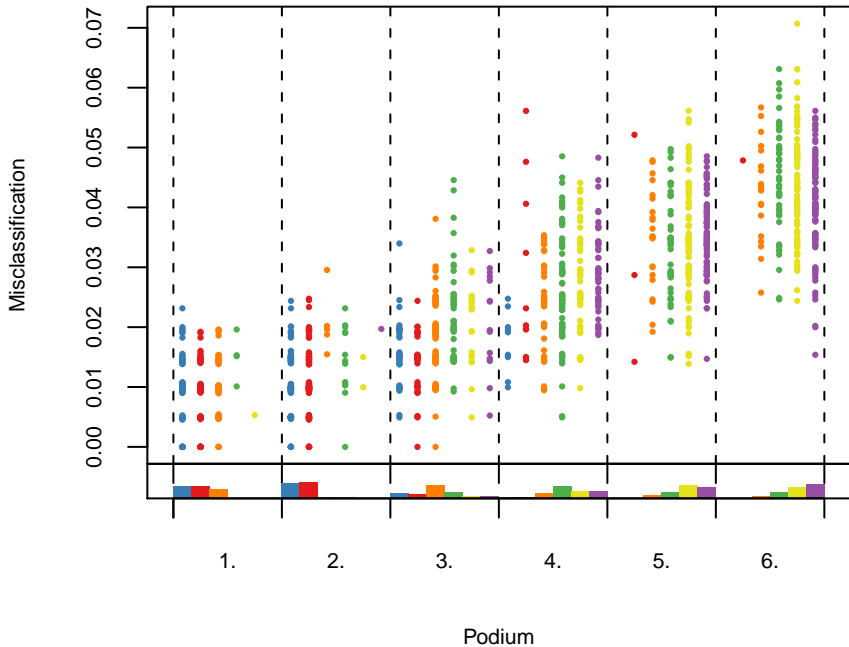
$\phi =$	Mean	SD	Median	Max
purple	0.0352	0.0094	0.0350	0.0561
orange	0.0197	0.0117	0.0185	0.0567
yellow	0.0344	0.0118	0.0340	0.0707
red	0.0116	0.0080	0.0100	0.0561
blue	0.0110	0.0059	0.0100	0.0340
green	0.0293	0.0123	0.0273	0.0631

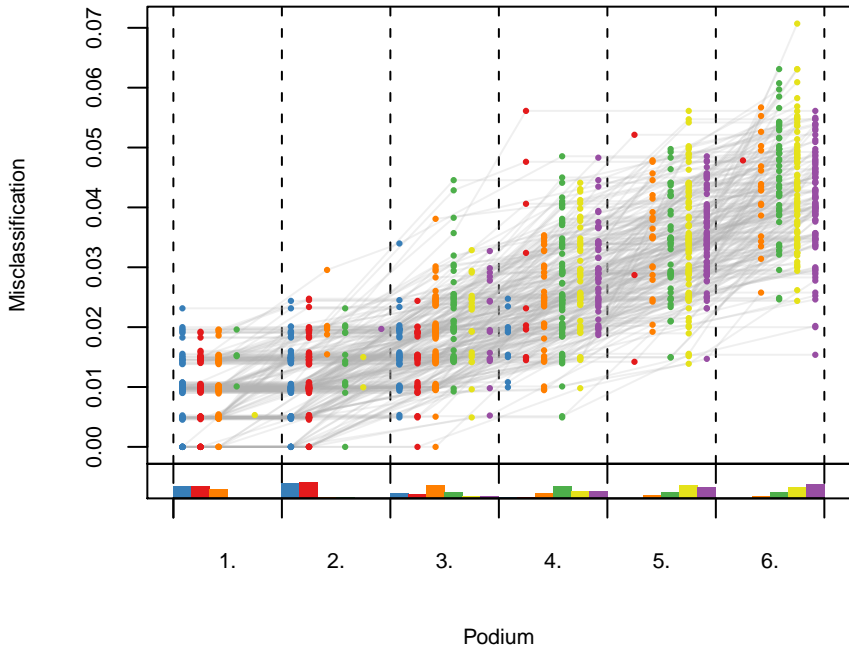
Order based on Mean:

blue < red < orange < green < yellow < purple

Exploratory analysis







Inferential analysis

Random block design:

The set of algorithms as experimental factor κ_j , a learning sample as the blocking factor b_i ; κ_0 the intercept which expresses the basic performance:

$$p_{ij} = \kappa_0 + \kappa_j + b_i + \epsilon_{ij},$$
$$i = 1, \dots, B, j = 1, \dots, (K - 1)$$

Hypothesis of interest:

$$H_0 : \kappa_1 = \cdots = \kappa_{K-1} = 0,$$

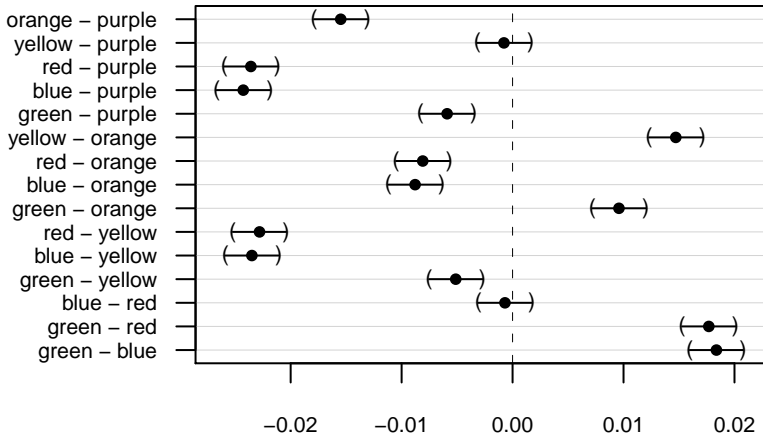
$$H_A : \exists j : \kappa_j \neq 0.$$

Testing using parametric and non-parametric methods (which take different assumptions on κ_j , b_i and ϵ_{ij}).

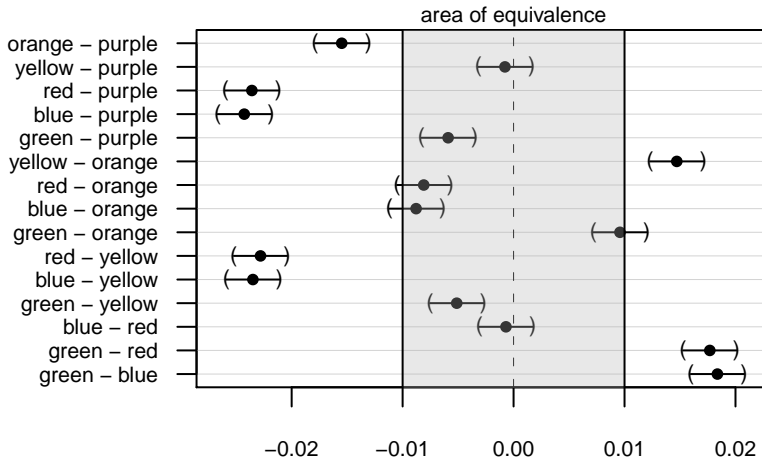
Mixed effects model:

k_j fixed effect, b_i random effect and $b_i \sim N(0, \sigma_b^2)$, $\epsilon_{ij} \sim N(0, \sigma^2)$.

Pairwise comparisons with Tukey contrasts; calculation of simultaneous confidence intervals (95%), which allow controlling the experimentwise error.



Pairwise comparisons with Tukey contrasts; calculation of simultaneous confidence intervals (95%), which allow controlling the experimentwise error.



Relations

Equivalence:

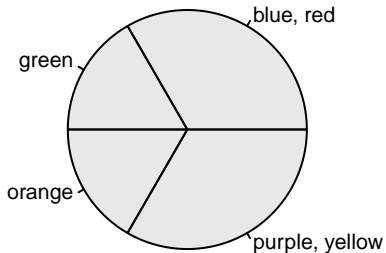
Based on the p -values: $a_i \approx a_j$ iff their difference is not significant (relevant).

Properties ($i, j, k = 1, \dots, K$):

1. reflexive: $a_i \approx a_i$
2. symmetric: $a_i \approx a_j \Rightarrow a_j \approx a_i$
- 3? transitive: $a_i \approx a_j \wedge a_j \approx a_k \Rightarrow a_i \approx a_k$

Equivalence classes (based on significance):

	purple	orange	yellow	red	blue	green
purple	1	0	1	0	0	0
orange	0	1	0	0	0	0
yellow	1	0	1	0	0	0
red	0	0	0	1	1	0
blue	0	0	0	1	1	0
green	0	0	0	0	0	1



Strict weak ordering:

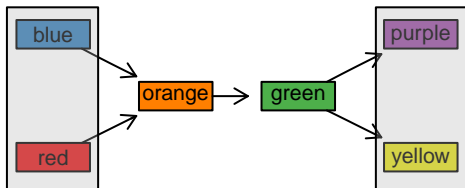
Based on the p -values and the directions of the tests: $a_i < a_j$ iff a_i is significantly (or relevantly) better than a_j .

Properties ($i, j, k = 1, \dots, K$):

1. irreflexive: $a_i \not< a_i$
2. asymmetric: $a_i < a_j \Rightarrow a_j \not< a_i$
- 3? transitive: $a_i \approx a_j \wedge a_j \approx a_k \Rightarrow a_i \approx a_k$
- 4? negatively transitive: $a_i \not< a_j \wedge a_j \not< a_k \Rightarrow a_i \not< a_k$

Equivalence classes and their order (based on significance):

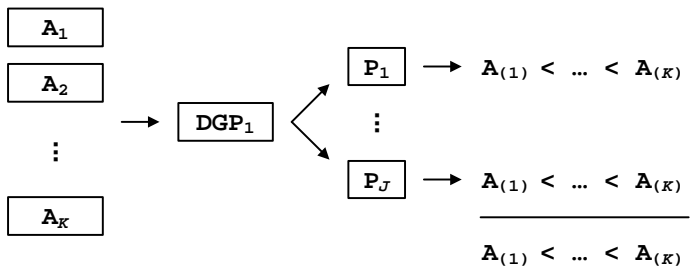
	purple	orange	yellow	red	blue	green
purple	0	0	0	0	0	0
orange	1	0	1	0	0	1
yellow	0	0	0	0	0	0
red	1	1	1	0	0	1
blue	1	1	1	0	0	1
green	1	0	1	0	0	0



$\text{blue} \approx \text{red} < \text{orange} < \text{green} < \text{purple} \approx \text{yellow}$

Scenario 2

“Aggregation of different performance measures”



Empirical performance distributions:

For each candidate algorithm a_k J empirical performance distributions $\hat{P}_k^j(\mathcal{L})$ are estimated based on \hat{p}_{kb}^j .

Relation ensemble:

The exploratory and inferential analysis leads to an order relation for each performance measure:

$$\mathcal{R} = \{R_1, \dots, R_J\}$$

Exemplar benchmark experiment

$\mathcal{R} = \{R_m, R_{mm}, R_c\}$ with R_m the relation based on the misclassification and the mixed effect model; R_{mm} the relation based on the minimax rule applied on the misclassification to control the worst-case scenario; and R_c the relation based on the computation time.

$R_m = \text{blue} \approx \text{red} < \text{orange} < \text{green} < \text{purple} \approx \text{yellow}$

$R_{mm} = \text{blue} < \text{purple} \approx \text{red} < \text{orange} < \text{green} < \text{yellow}$

$R_c = \text{red} < \text{purple} < \text{orange} < \text{yellow} < \text{green} < \text{blue}$

Consensus decision making

$R_m = \text{blue} \approx \text{red} < \text{orange} < \text{green} < \text{purple} \approx \text{yellow}$

$R_{mm} = \text{blue} < \text{purple} \approx \text{red} < \text{orange} < \text{green} < \text{yellow}$

$R_c = \text{red} < \text{purple} < \text{orange} < \text{yellow} < \text{green} < \text{blue}$

Aggregate the preferences of voters, i.e. the performance measures' orders of the candidate algorithms, in a way that "all voters can live with".

Borda count:

Relations as “voters ranks”: the number of points given to an algorithm equates to the number of related algorithms.

On relation R :

$$\text{borda}_R(a_i) = \#\{a_j \mid (a_i, a_j) \in R, j = 1, \dots, K\}$$

On the set of relations \mathcal{R} (w_R is a weighting factor):

$$\text{borda}_{\mathcal{R}}(a_i) = \sum_{R \in \mathcal{R}} w_R \cdot \text{borda}_R(a_i),$$

Final order R_{borda} is the order of the algorithms according to their Borda count.

Voters points:

	purple	orange	yellow	red	blue	green
R_m	0	3	0	4	4	2
R_{mm}	3	2	0	3	5	1
R_c	4	3	2	5	0	1
	7	8	2	12	9	4

Final ranking:

red	blue	purple	orange	green	yellow
1	2	3	4	5	6

$$R_{\text{borda}} = \text{red} < \text{blue} < \text{purple} < \text{orange} < \text{green} < \text{yellow}$$

Majority criterion:

If there is a single candidate preferred by a majority of voters to all other candidates, then that candidate should win.

Borda count fails:

51 voters: $a_1 < a_2 < a_3 < a_4$

5 voters: $a_2 < a_3 < a_4 < a_1$

23 voters: $a_3 < a_2 < a_4 < a_1$

21 voters: $a_4 < a_2 < a_3 < a_1$

a1	a2	a3	a4
153	205	151	91

$a_2 < a_1 < a_3 < a_4$

Condorcet methods:

Rank algorithm a_i above algorithm a_j iff the number of individual wins of a_i over a_j exceeds the number of losses.

Voting or Condorcet's paradoxon:

Collective preferences can be cyclic, i.e. not transitive, even if the preferences of individual voters are not.

Optimization approach:

Find the minimal solution of the problem

$$\sum_{R \in \mathcal{R}} w_R \cdot d_{\Delta}(R_{con}, R) \Rightarrow \min_{R_{con} \in \mathcal{C}}.$$

w_R is a weighting factor, and \mathcal{C} defines a suitable set of possible consensus relations, e.g., the set of linear orders or the set of complete preorders.

$d_{\Delta}(R_1, R_2)$ is the symmetric difference distance:

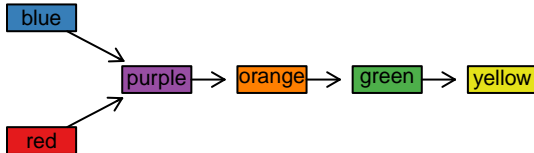
$$\#\{(a_i, a_j) \mid ((a_i, a_j) \in R_1) \oplus ((a_i, a_j) \in R_2), i, j = 1, \dots, K\}$$

$R_m = \text{blue} \approx \text{red} < \text{orange} < \text{green} < \text{purple} \approx \text{yellow}$

$R_{mm} = \text{blue} < \text{purple} \approx \text{red} < \text{orange} < \text{green} < \text{yellow}$

$R_c = \text{red} < \text{purple} < \text{orange} < \text{yellow} < \text{green} < \text{blue}$

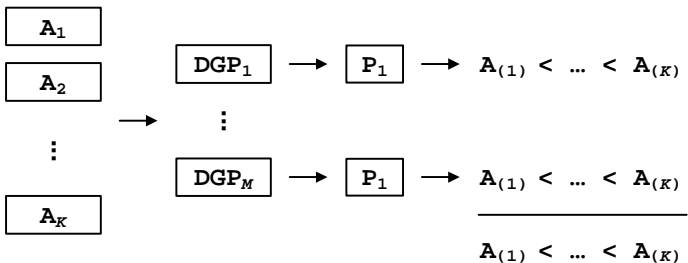
Optimization consensus over the set \mathcal{C} of partial orders:



$R_{con} = \text{blue} \approx \text{red} < \text{purple} < \text{orange} < \text{green} < \text{yellow}$

Scenario 3

“More than one data set”



Domain of interest:

Given is a set of data sets $\mathcal{D} = \{\mathcal{L}_1, \dots, \mathcal{L}_M\}$ consisting of M data sets representing the problem domain of interest.

For each data set \mathcal{L}_m ($m = 1, \dots, M$) the benchmark experiment is executed and results in $J = 1 \cdot K$ estimations of empirical performance distributions $\hat{\mathcal{P}}_k^j(\mathcal{L}_m)$ with $j = 1, \dots, J$ and $k = 1, \dots, K$.

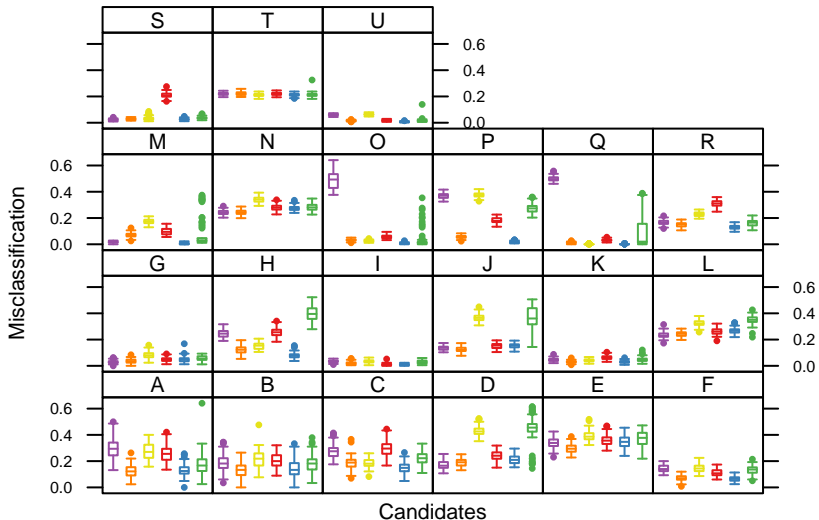
These raw results are aggregated to an order R_m of the candidate algorithms.

Exemplar benchmark experiment

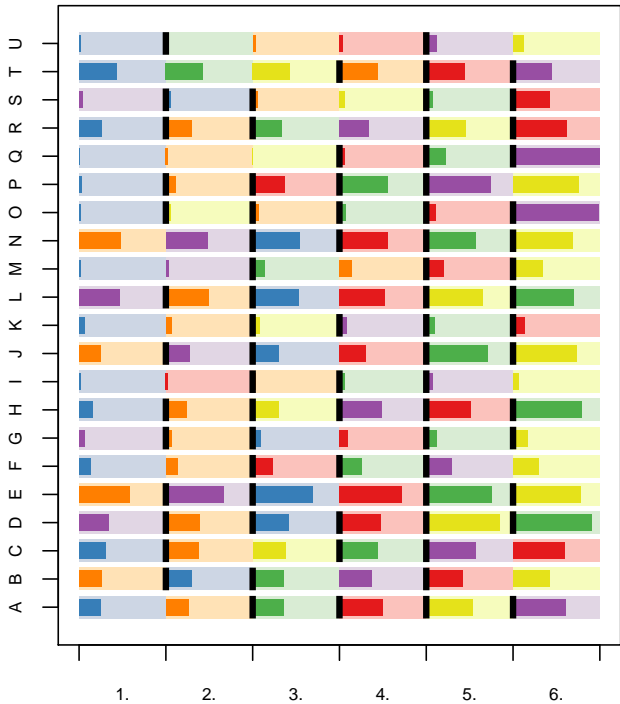
(1) 21 binary classification problems originated from the UCI Machine Learning repository: $\mathcal{D} = \{A, \dots, U\}$ (monks4 is letter I); **(2)** - **(5)** unchanged.

Execution results in 21×6 empirical performance distributions and 21 order relations.

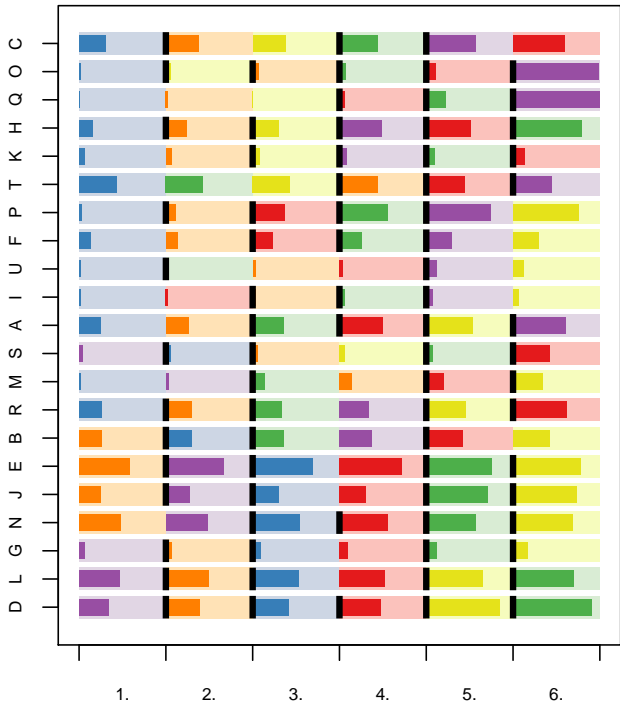
Exploratory analysis



Datasets



Datasets



Consensus ranking

Optimization consensus over the set \mathcal{C} of partial orders:

blue \approx orange $<$ purple $<$ green $<$ red $<$ yellow

blue \approx orange $<$ purple $<$ red $<$ green $<$ yellow

Final order:

blue \approx orange $<$ purple $<$ green \approx red $<$ yellow

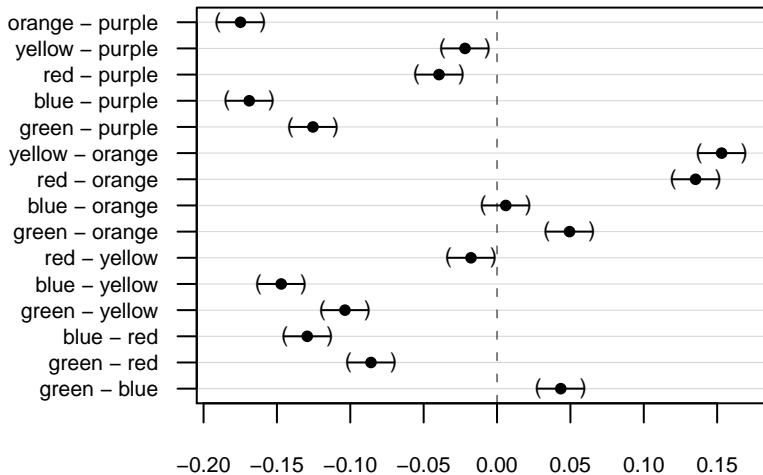
Inferential analysis

Design with two experimental factors, their interactions, and blocking factors at two levels:

$$p_{ijk} = \kappa_0 + \kappa_j + \gamma_k + \delta_{jk} + b_k + b_{ki} + \epsilon_{ijk}$$
$$i = 1, \dots, B, j = 1, \dots, (K - 1), k = 1, \dots, (M - 1).$$

The experimental factors are described by κ and γ , their interactions by δ . κ_0 represents the basic performance of each algorithm, κ_j the individual differences from the basic performance. γ_k represents the data set effect. The blocking factors are described by b_k , the data sets, and b_{ki} , the sampling within the data sets. ϵ_{ijk} describes the systematic error.

Mixed effects model:



Final order:

blue \approx orange < green < red < yellow < purple

and finally ...

Statistical correct orders:

Statistical correct orders within the domain represented by the 21 data sets $\{A, \dots, U\}$:

$R_{con} = \text{blue} \approx \text{orange} < \text{purple} < \text{green} \approx \text{red} < \text{yellow}$

$R_{mem} = \text{blue} \approx \text{orange} < \text{green} < \text{red} < \text{yellow} < \text{purple}$

Papers and Software

... from <http://statistik.lmu.de/~eugster/benchmark>.

R Package:

`benchmark` available from R-Forge.

Papers and reports:

Manuel J. A. Eugster and Friedrich Leisch. *Bench plot and mixed effects models: First steps toward a comprehensive benchmark analysis toolbox*. In Paula Brito, editor, *Compstat 2008-Proceedings in Computational Statistics*, pages 299-306. Physica Verlag, Heidelberg, Germany, 2008.

Manuel J. A. Eugster, Torsten Hothorn and Friedrich Leisch. *Exploratory and Inferential Analysis of Benchmark Experiments*. Technical Report 30, LMU Munich. Submitted.