

compboost

Fast and Flexible Way of bla

Daniel Schalk

July 4, 2019

LMU Munich

Working Group Computational Statistics



Use-Case

- We own a small booth at the city center that sells beer.
- As we are very interested in our customers' health, we only sell to customers who we expect to drink less than 110 liters per year.
- To estimate how much a customer drinks, we have collected data from 200 customers in recent years.
- These data include the beer consumption (in liter), age, sex, country of origin, weight, body size, and 200 characteristics gained from app usage (that have absolutely no influence).

Use-Case

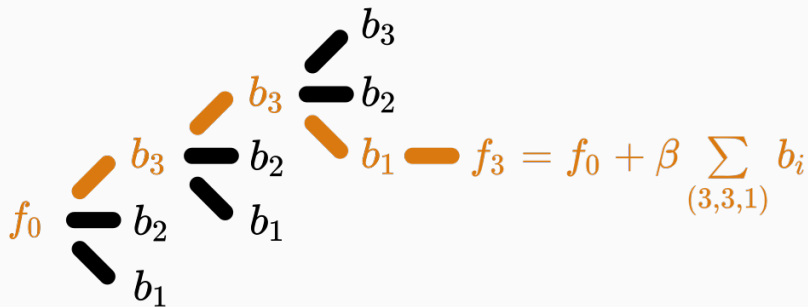
beer_consumption	gender	country	age	weight	height	app_usage1	app_usage2
106.5	m	Seychelles	33	87.2	173	0.168	0.606
85.5	f	Seychelles	52	89.4	200	0.808	0.938
116.5	f	Czechia	54	92.0	179	0.385	0.264
67.0	m	Australia	32	63.5	186	0.328	0.380
43.0	f	Australia	51	64.7	175	0.602	0.807
85.0	m	Austria	43	95.7	173	0.604	0.978
79.0	f	Austria	55	87.6	156	0.125	0.958
107.0	f	Austria	24	93.2	161	0.295	0.763
57.0	m	USA	55	76.3	183	0.578	0.510
89.0	m	USA	16	72.2	203	0.631	0.064

With these data we want to answer the following questions:

- Which of the customers' characteristics are important to be able to determine the consumption?
- How does the effect of important features look like?
- How does the model behave on unseen data?

What is Component-Wise Boosting?

The General Idea



- Inherent (unbiased) feature selection.
- Resulting model is sparse since important effects are selected first and therefore it is able to learn in high-dimensional feature spaces ($p \gg n$).
- Parameters are updated iteratively. Therefore, the whole trace of how the model evolves is available.

The Idea Behind Compboost

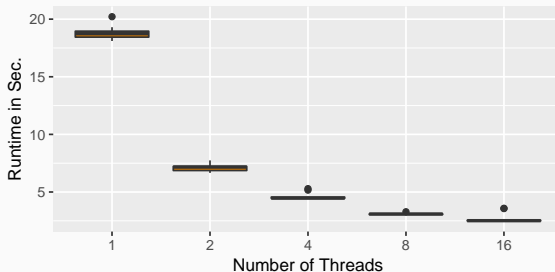
About Compboost

The `compboost` package is a fast and flexible framework for model-based boosting:

- With `mboost` as standard, we want to keep the modular principle of defining custom base-learner and losses.
- Completely written in C++ and exposed by `Rcpp` to obtain high performance and full memory control.
- R API is written in R6 to provide convenient wrapper.
- Major parts of the `compboost` functionality are unit tested against `mboost` to ensure correctness.

Runtime and Memory Considerations

- Matrices are stored (if possible) as sparse matrix.
- Take advantage of the matrix structure to speed up the algorithm by reducing the number of repetitive or too expensive calculations.
- Optimizer are parallelized via openmp:



A Short Demonstration

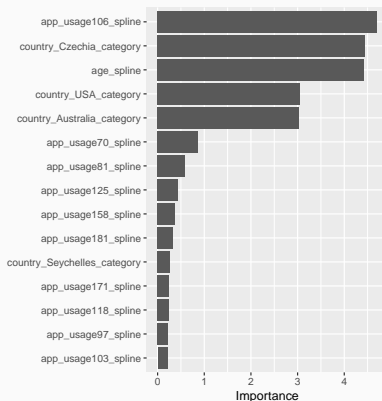
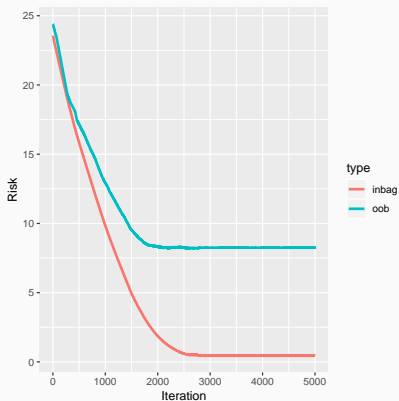
Using Convenience Wrapper

```
set.seed(618)
cboost = boostSplines(data = beer_data, target = "beer_consumption",
  loss = LossAbsolute$new(), learning_rate = 0.1, iterations = 5000L,
  penalty = 10, oob_fraction = 0.3, trace = 2500L)

##      1/5000    risk = 24  oob_risk = 24
## 2500/5000    risk = 0.6  oob_risk = 8.3
## 5000/5000    risk = 0.44 oob_risk = 8.3
##
##
## Train 5000 iterations in 18 Seconds.
## Final risk based on the train set: 0.44
```

Visualizing Results

```
gg1 = cboost$plotInbagVsOobRisk()  
gg2 = cboost$plotFeatureImportance()  
  
gridExtra::grid.arrange(gg1, gg2, ncol = 2L)
```

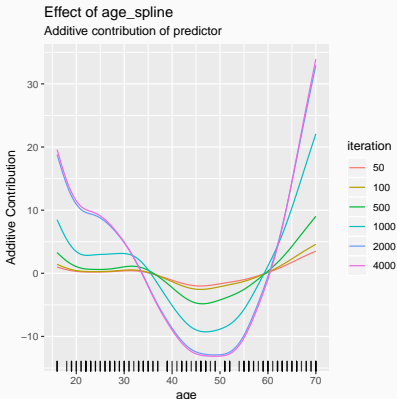
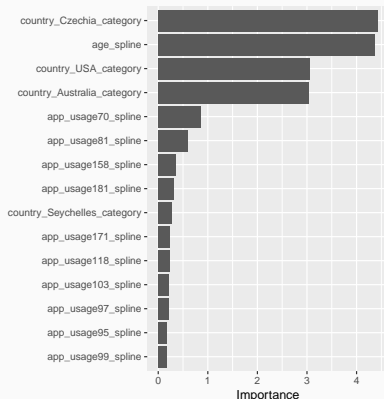


Visualizing Results

```
cboost$train(200L)

gg1 = cboost$plotFeatureImportance()
gg2 = cboost$plot("age_spline", iters = c(50, 100, 500, 1000, 2000, 4000))

gridExtra::grid.arrange(gg1, gg2, ncol = 2L)
```



Using the R6 Interface

```
cboost = Comboost$new(data = beer_data, target = "beer_consumption",
  loss = LossAbsolute$new(), learning_rate = 0.1, oob_fraction = 0.3,
  optimizer = OptimizerCoordinateDescent$new(nthreads = 8L))

cboost$addBaselearner("age", "spline", BaselearnerPSpline)
cboost$addBaselearner("country", "category", BaselearnerPolynomial)

cboost$addLogger(logger = LoggerTime, use_as_stopper = TRUE, logger_id = "time",
  max_time = 100000, time_unit = "microseconds")

cboost$train(10000, trace = 500)

##      1/10000   risk = 24  oob_risk = 24   time = 1
##    500/10000   risk = 14  oob_risk = 17   time = 25207
##   1000/10000   risk = 7.4  oob_risk = 9.4   time = 53724
##   1500/10000   risk = 3.5  oob_risk = 5    time = 88625
##
##
## Train 1637 iterations in 0 Seconds.
## Final risk based on the train set: 3.3
```


More Advanced Customizations

- Custom loss function and base-learner
- Advanced stopper for early stopping (e.g. time or performance based stopping)
- Parallelization via openmp is controlled by the optimizer, e.g.
`OptimizerCoordinateDescent$new(4L)`

What's Next?

What's Next?

- Research on computational aspects of the algorithm:
 - More stable base-learner selection process via resampling
 - Base-learner selection for arbitrary performance measures
 - Smarter and faster optimizer to select base-learner
- Greater functionality:
 - Functional data structures and loss functions
 - Unbiased feature selection
 - Effect decomposition into constant, linear, and non-linear
- Reducing the memory load by applying binning to numerical features.
- Python API

Questions?