# Scalability of Predictive Modeling Algorithms

Master's thesis presentation

Ing. Tomáš Frýda

Supervisor: Ing. Pavel Kordík, Ph.D.

# Outline

## FAKE GAME
### Base Models
### Ensembles
### Evolution

# Outline

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

# Outline

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

# Outline

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

# Outline

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

# Outline

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

# Outline

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
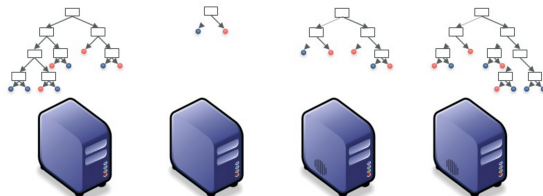Optimization
SMBO

Implementation

Experiments

Results

References

# Motivation

- ▶ Model selection usually does not depend solely on predictive performance
- ▶ I took in to account time, which gives me two basic use cases:
    - ▶ Good enough model trained using limited computational resources
    - ▶ Highly accurate model trained using as much computational resources as needed
- ▶ Make FAKE GAME usable on big data

# FAKE GAME

- Originally created for small data
- Base models
    - Decision tree, KNN, etc
    - Regression models
- Ensembles
    - Bagging
    - Boosting
    - Stacking, Cascade Correlation, ...
- Genetic programming-based ensemble creation

# Base Models

- ▶ regression models
  - ▶ linear
  - ▶ polynomial
  - ▶ sigmoid
  - ▶ sine, . . .
- ▶ regression models are used as discriminant functions for classification
- ▶ decision trees
- ▶ k-NN
- ▶ . . .

# Ensembles

- ► Arbitrating[1]
- ► Bagging
- ► Boosting
- ► Cascade Generalization
- ► Cascading[1]
- ► Delegating[1]
- ► Stacking

---

[1]used only for classification

# Evolution

Genetic programming used for evolving templates that can be expanded to hierarchical ensembles

# H2O

- ▶ framework for distributed machine learning based on MapReduce
- ▶ support for preprocessing and data manipulation
- ▶ RESTful API used by various language bindings (R, Python, . . . )

# Overview

Machine Learning algorithms included in H2O:

- ▶ Deep Learning
- ▶ Distributed Random Forest
- ▶ Gradient Boosting Machines
- ▶ Generalized Linear Model
- ▶ Naïve Bayes
- ▶ K-Means
- ▶ PCA
- ▶ GLRM
- ▶ . . .

# Architecture

- in-memory MapReduce
- uses distributed key-value storage
- tries to keep related data in the same or nearby node in order to minimize network usage
- columns are compressed and lazily decompressed just in time of usage in CPU registers
- parallel data load

See more at
http://blog.h2o.ai/2014/03/h2o-architecture/

# How are decision trees built in H2O?

**Implementation #1**

## Build independent trees per machine local data

- RVotes approach
- Each node builds a subset of forest



0xdata

Chawla, N., & Hall, L. (2004). *Learning ensembles from bites: A scalable and accurate approach.* The Journal of Machine Learning Research, 5, p421–451.

## Implementation #1

❤ Fast - trees are independent and can be built in parallel

🚫 Data have to fit into memory

🚫 Possible accuracy decrease if each node can see only
subset of data

0xdata

Implementation #2

Build a distributed tree over all data

Dataset
points

0xdata

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

## Implementation #2

### Tree is built per layer

• Histograms are reduced
  and a new layer
  is prepared



*Active
tree layer*

*In/Out of
bag
flags*

*Dataset
points*

0xdata

## Implementation #2

✔ Exact solution - no decrease of accuracy

✔ Elegant solution merging tree building and OOB scoring

🚫 More data transfers to exchange histograms

🚫 Can produce huge trees (since tree size depends on data)

0xdata

## Python

```python
import h2o

h2o.init()
data = h2o.import_file(
        path="data.csv")

# Create test/train split
s = data["Year"].runif()
train = data[s <= 0.75]
test  = data[s > 0.75]

# Create an estimator
dl=H2ODeepLearningEstimator(
    hidden=[10,10],
    epochs=5,
    balance_classes=True)

# Train an estimator
dl.train(
    x=myX,
    y=myY,
    training_frame=train,
    validation_frame=test)
```

## R

```r
library("h2o")

h2o.init()
dt <- h2o.importFile(
    path = "data.csv")

# Create test/train split
dt.split <- h2o.splitFrame(
    data = dt,
    ratios = 0.75)
train <- dt.split[[1]]
test <- dt.split[[2]]

# Create an estimator and
# train it
dl <- h2o.deeplearning(
  x = myX,
  y = myY,
  training_frame = train,
  validation_frame = test,
  hidden=c(10, 10))
```

# H2O Flow

# H2O Flow

Scalability of Predictive Modeling Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter Optimization
SMBO

Implementation

Experiments

Results

References

# H2O Flow

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
Usage
Related Software

Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

# H2O Flow

# H2O Flow

# Related Software

- ▶ Sparkling Water
- ▶ Deep Water
- ▶ Steam — H2O deployment

# Sparkling Water

# Deep Water

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution
H2O
Overview
Architecture
Usage
Related Software
Hyper-parameter
Optimization
SMBO
Implementation
Experiments
Results
References

# Hyper-Parameter Optimization

- ▶ Grid Search
- ▶ Random
  Search



- ▶ Bayesian optimization (SMAC)

# Sequential Model-based Bayesian Optimization (SMBO)

1. evaluate random configuration and add it to the probabilistic model

2. select promising configuration based on probabilistic model using an acquisition function[2]

3. evaluate the configuration

4. add the new configuration to the probabilistic model

5. go to step 2

---

[2]usually Expected Improvement $\mathbf{E}I(x) = \mathbf{E}(\max\{0, f_{t+1}(x) - f(x^+)\}|\mathsf{M}_t)$

# Instances of SMBO

- ► Gaussian Process based SMBO
  - ► no obvious way how to deal with categorical parameters
- ► Tree-structured Parzen Estimator (TPE)
- ► Sequential Model-based Algorithm Configuration (SMAC)
- ► Hyperband

# Tree-structured Parzen Estimator

$$p(x|y) = \begin{cases} l(x) & \text{if} \quad y < y^* \\ g(x) & \text{if} \quad y \geq y^* \end{cases}$$

► easy to sample space of promising values

► EI is proportional to $\left(\gamma + \frac{g(x)}{l(x)}(1-\gamma)\right)^{-1}$

---

$y^*$ is chosen as some quantile (e.g., $p(y < y^*) = 0.15 = \gamma$)

# Sequential Model-based Algorithm Configuration

- ▶ based upon ROAR — an racing optimization algorithm
- ▶ uses random forest as a probabilistic model
- ▶ usage of random forest makes it easy to use user-defined cost metric
- ▶ configuration to be evaluated is selected by following process
  1. take 10 best previously seen configurations
  2. initialize local search (using one-exchange neighbourhood for categorical, and four random neighbours for numerical variables)
  3. merge resulting 10 best configurations with 10 000 randomly sampled configurations
  4. sort by their EI
  5. interleave with uniformly sampled configurations

# Hyperband

- ▶ SMBO with enhanced selection and evaluation phase
- ▶ uses information from training phase of a model that is being optimized and eventually stops it if it doesn't converge well $\implies$ explores more space using the same amount of resources
- ▶ iteratively discards the worse half of evaluated configurations

# Implementation

- ▶ integration of FAKE GAME into H2O framework
- ▶ creation of benchmarking environment
  - ▶ written in Python
  - ▶ supports
    - ▶ all supervised machine learning algorithms in H2O
    - ▶ H2O Ensemble (implemented in R, based on SuperLearner package)
    - ▶ Hyper-Parameter optimization using Random Search and SMAC
  - ▶ configurable using YAML and Python

# Experiments

▶ 2 datasets with binomial response class
  ▶ Higgs
  ▶ Airline - 4 different scenarios
▶ 20+ models benchmarked on each of 5 scenarios
▶ Hyper-Parameter optimization on each dataset

# Overview of Results

Figure: Higgs dataset



Figure: Airline – predicting IsDepDelayed

# Decision Boundary on Airline dataset

# Decision Boundary on Airline dataset

**GBM**, probability of class

**GBM**, decision boundary

**Deep Learning**, probability of class

**Deep Learning**, decision boundary

# Hyper-Parameter Optimization

# Conclusion

- ▶ Successfully integrated FAKE GAME into H2O and created benchmarking environment
- ▶ Experiments took over 2000 hours (wall clock), used 12 CPUs and 16 GiB of RAM
- ▶ Experiments show that
  - ▶ newly integrated FAKE GAME can find better models than those previously present in H2O
  - ▶ H2O's auto-tuning yields good results by default
- ▶ Results of those experiments were submitted, as part of an article, to be published in Machine Learning

Scalability of
Predictive
Modeling
Algorithms

Ing. Tomáš Frýda

FAKE GAME
Base Models
Ensembles
Evolution

H2O
Overview
Architecture
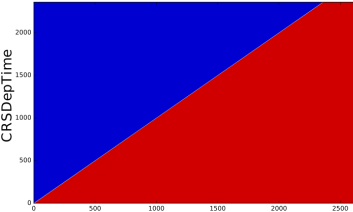Usage
Related Software

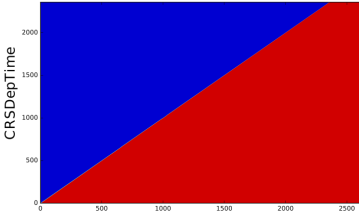Hyper-parameter
Optimization
SMBO

Implementation

Experiments

Results

References

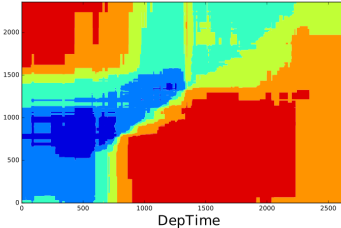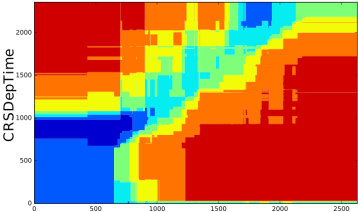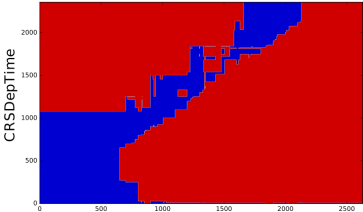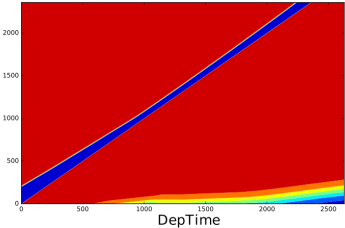0xdata,
`https://www.slideshare.net/0xdata/rf-brighttalk`
`https://www.slideshare.net/0xdata/`
`deep-water-gpu-deep-learning-for-h2o-arno-candel`

J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl,
"Algorithms for Hyper-Parameter Optimization,"
Adv. Neural Inf. Process. Syst., pp. 2546–2554, 2011.

F. Hutter, H. H. Hoos, and K. Leyton-Brown,
"Sequential model-based optimization for general algorithm
configuration",
Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif.
Intell. Lect. Notes Bioinformatics), vol. 6683 LNCS, pp. 507–523,
2011.

L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A.
Talwalkar,
"Efficient Hyperparameter Optimization and Infinitely Many
Armed Bandits,"
arXiv Prepr., 2016.