# Explainable AI:
# White Box Modeling by Symbolic Regression

**Stephan M. Winkler**
**FH Oberösterreich, Campus Hagenberg**

# FH Upper Austria, Campus Hagenberg



**FH-Prof. PD DI Dr. Stephan M. Winkler**

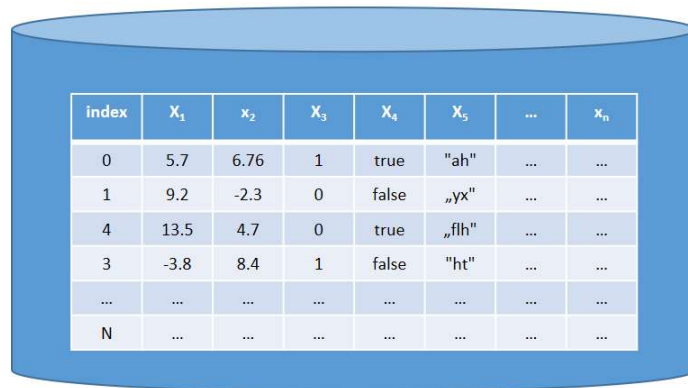**Professor for system identification, genetic programming, machine learning**

- **Teaching:**

  > Head of Department for Medical & Bioinformatics and Data Science & Engineering at FH Upper Austria, Campus Hagenberg

  > Interests: Machine learning, artificial intelligence, algorithm development, genetic programming, image analysis, …

- **Research:**

  > Lead of the Bioinformatics Research Group in Hagenberg

  > Member of the Heuristic and Evolutionary Algorithms Laboratory (HEAL)

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# Systems Identification and Machine Learning

| index | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | ... | $x_n$ |
|-------|-------|-------|-------|-------|-------|-----|-------|
| 0 | 5.7 | 6.76 | 1 | true | "ah" | ... | ... |
| 1 | 9.2 | -2.3 | 0 | false | „yx" | ... | ... |
| 4 | 13.5 | 4.7 | 0 | true | „flh" | ... | ... |
| 3 | -3.8 | 8.4 | 1 | false | "ht" | ... | ... |
| ... | ... | ... | ... | ... | ... | ... | ... |
| N | ... | ... | ... | ... | ... | ... | ... |

Supervised Machine Learning

$$y_1 = f_1(x_1, x_2, ..., x_n)$$
$$y_2 = f_2(x_1, x_2, ..., x_n)$$
$$y_3 = f_3(x_1, x_2, ..., x_n)$$

Models

– Systems identification

– Modeling

– Process analysis

– Prediction

– Optimization

UNIVERSITY OF APPLIED SCIENCES UPPER AUSTRIA
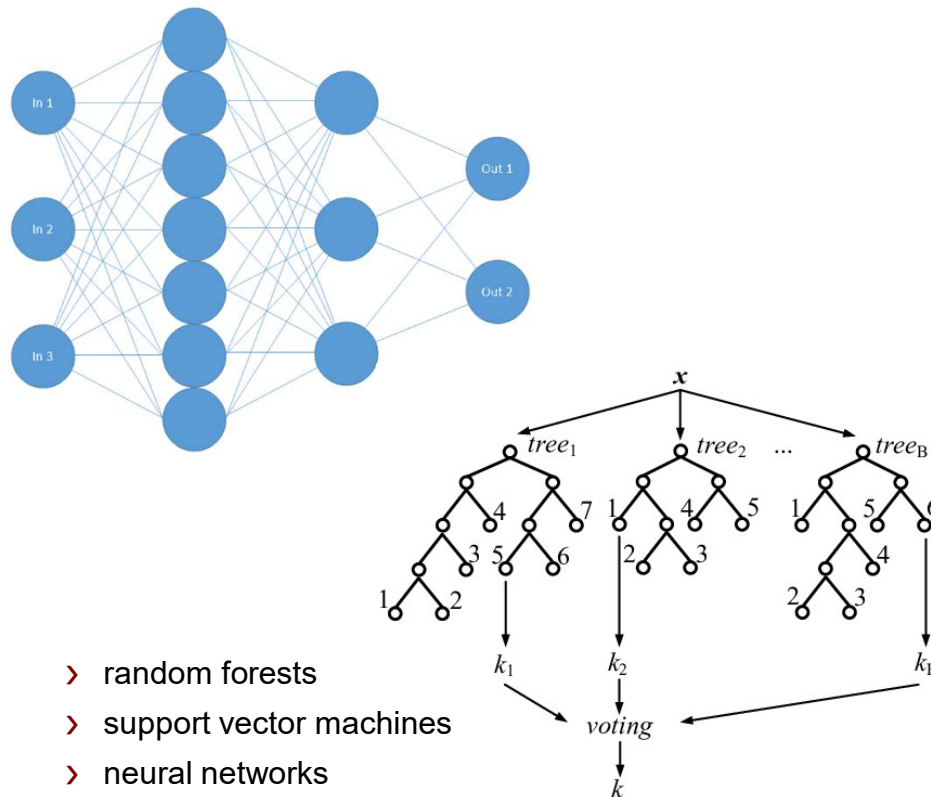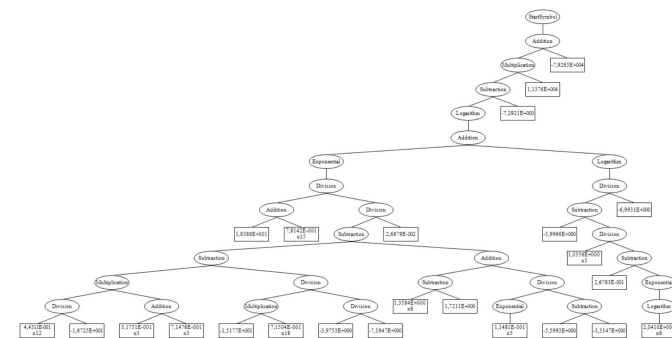
# Black-Box vs. White-Box Modeling

# Black-Box vs. White-Box Modeling

## Black Box

## White Box



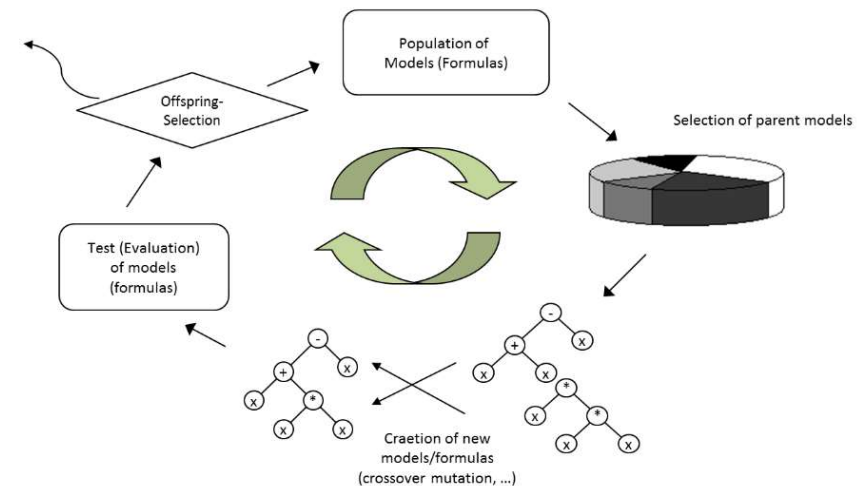$$y = \left(\left(\left(\left(c_0 \cdot \mathbf{x}31 + \left(\frac{(c_1 \cdot \mathbf{x}22 + c_2 \cdot \mathbf{x}17)}{c_3} - (c_4 \cdot c_5 + c_6 \cdot \mathbf{x}49)\right)\right) + (c_7 + c_8) \cdot \frac{c_9}{c_{10}}\right) + c_{11} \cdot \mathbf{x}14\right) \cdot c_{12} + c_{13}\right)$$

| | | | | | |
|---|---|---|---|---|---|
| $c_0 =$ | $-0.15366$ | $c_5 =$ | $-1.8194$ | $c_{10} =$ | $-11.55$ |
| $c_1 =$ | $1.4153$ | $c_6 =$ | $1.3732$ | $c_{11} =$ | $1.353$ |
| $c_2 =$ | $-0.090437$ | $c_7 =$ | $-5.5435$ | $c_{12} =$ | $-0.030191$ |
| $c_3 =$ | $18.821$ | $c_8 =$ | $4.6554$ | $c_{13} =$ | $-5.7891$ |
| $c_4 =$ | $-9.0013$ | $c_9 =$ | $-11.338$ | | |

> random forests

> support vector machines

> neural networks

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# White-Box Modeling



- **Genetic Programming**

  › evolutionary process

  › implicit feature selection
  › optimizes model structure and parameters

  › generates interpretable formulas

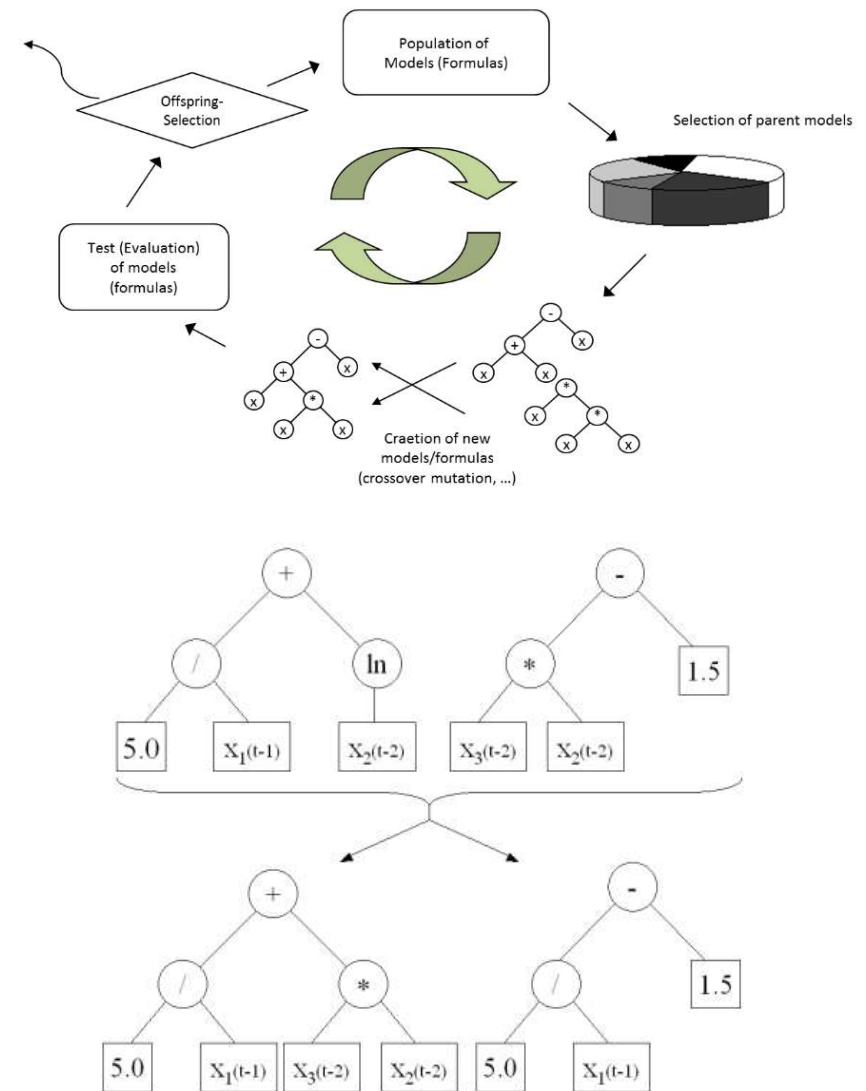  › results directly applicable

  › assessment of variable relevance

x1 ——→ [ ? ] ? ——→ y1
x2 ——→ ——→ y2

y1(t+1) = y2(t) / (x1+x2)
y2(t) = x1(t)+3*exp(x2(t-2)) ?

$$y = c_0 \cdot x_1 \cdot (\log(c_1 \cdot x_2) + c_2) + c_3$$
$$c_0 = 0.500331886126962$$
$$c_1 = 2.3702293890766$$
$$c_2 = 1.28570833399083$$
$$c_3 = 4.91856540837157$$

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# White-Box Modeling



- **Genetic Programming**

  › evolutionary process

  › implicit feature selection

  › optimizes model structure and parameters

  › generates interpretable formulas

  › results directly applicable

  › assessment of variable relevance

# White-Box Modeling

# HeuristicLab
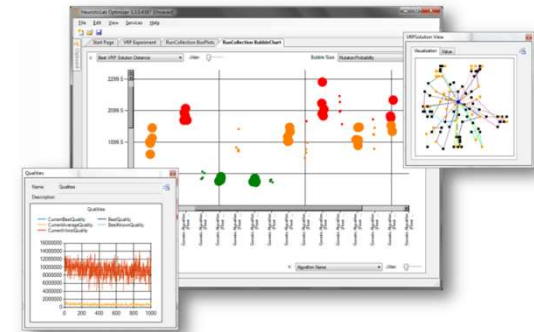
**Open Source Optimization Environment HeuristicLab**

- developed since 2002
- basis of many research projects and publications
- 2nd place at *Microsoft Innovation Award 2009*
- HeuristicLab 3.3.x since May 2010 under GNU GPL

## Motivation and Goals

- graphical user interface for interactive development, analysis and application of optimizations methods
- numerous optimization algorithms and optimization problems
- support for extensive experiments and analysis
- distribution through parallel execution of algorithms
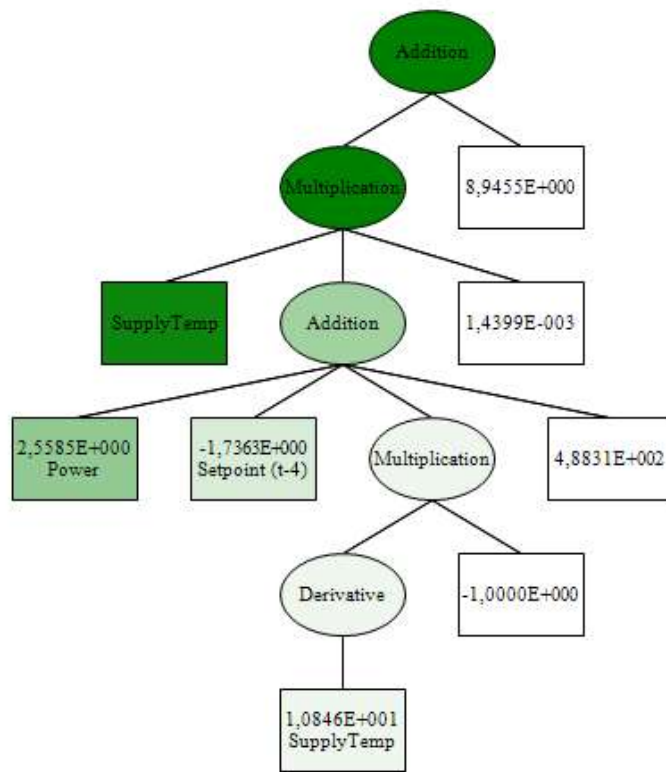- extensibility and flexibility (plug-in architecture)

## Distributed Computing with HeuristicLab Hive

- framework for distribution and parallel execution of HeuristicLab algorithms
- compute resources at Campus Hagenberg
    - 2006 – 2011: research cluster 1 (14 cores)
    - since 2009: research cluster 2 (112 cores, 448GB RAM)
    - since 2011: lab computers (100 PCs, on demand in the night)
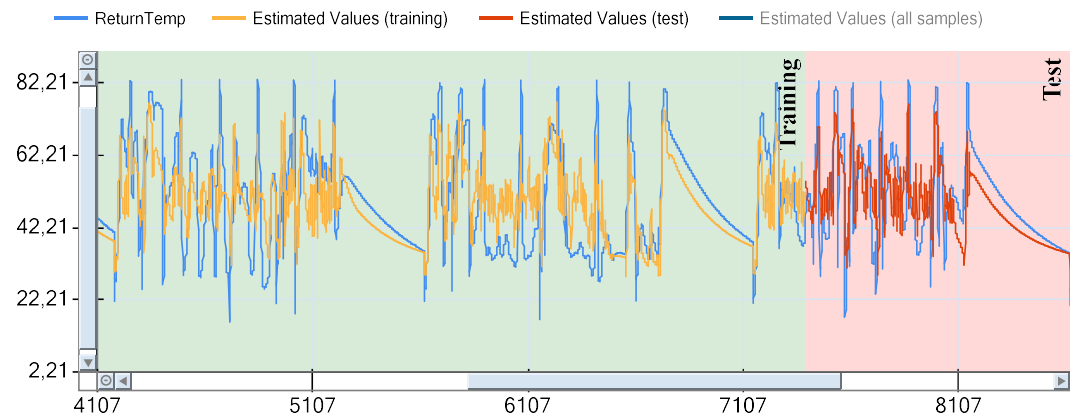    - since 2017: research cluster 3 (448 cores, 4TB RAM)

# White-Box Modeling

$$\text{SupplyTemp} \cdot \left( c_0 \cdot \text{Power} + c_1 \cdot \text{Setpoint}(t-4) + \frac{d(c_2 \cdot \text{SupplyTemp})}{dt} \cdot c_3 + c_4 \right) \cdot c_5 + c_6$$
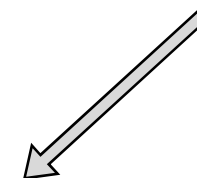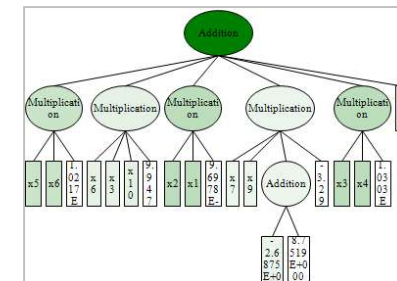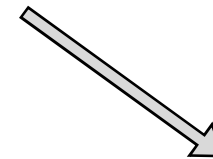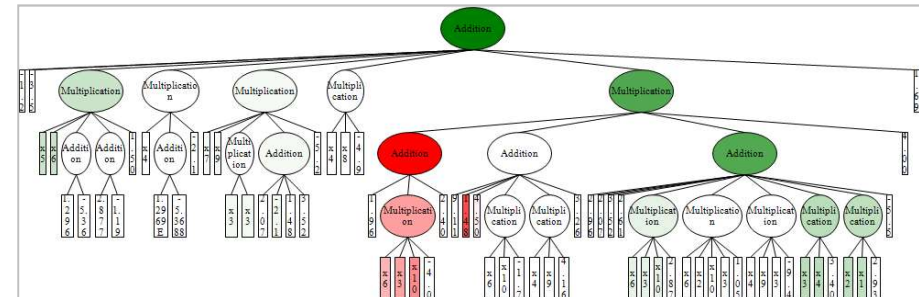
$$
\begin{aligned}
c_0 &= 2.5585 \\
c_1 &= -1.7363 \\
c_2 &= 10.846 \\
c_3 &= -1.0 \\
c_4 &= 488.31 \\
c_5 &= 0.0014399 \\
c_6 &= 8.9455
\end{aligned}
$$

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# Model Simplification
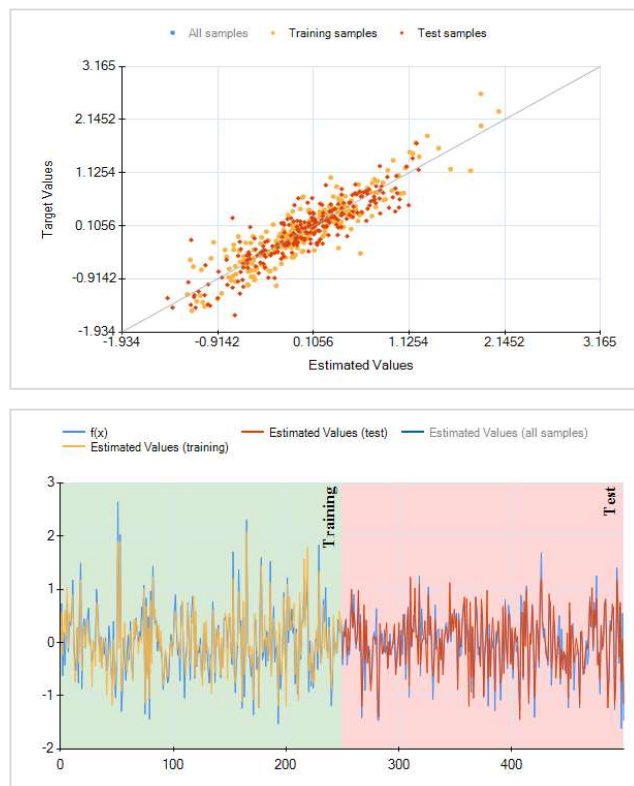
- Simplification Methods
  - mathematical transformation
  - remove nodes
  - constant optimization
  - external optimization

- Export
  - textual export
  - LaTeX, MATLAB
  - graphical export



$$y = x_1 \cdot x_2 + x_3 \cdot x_4 + x_5 \cdot x_6$$
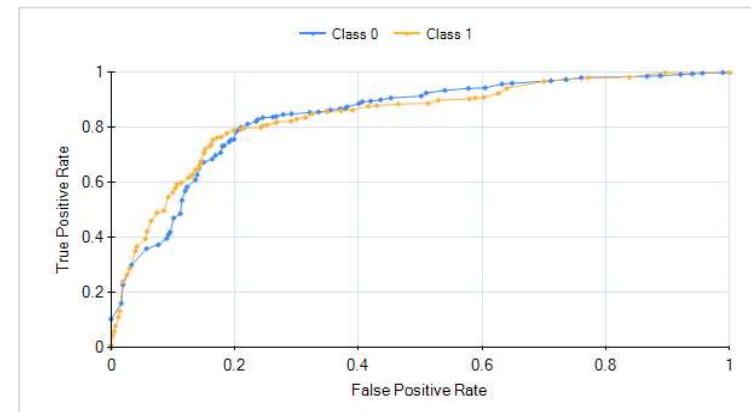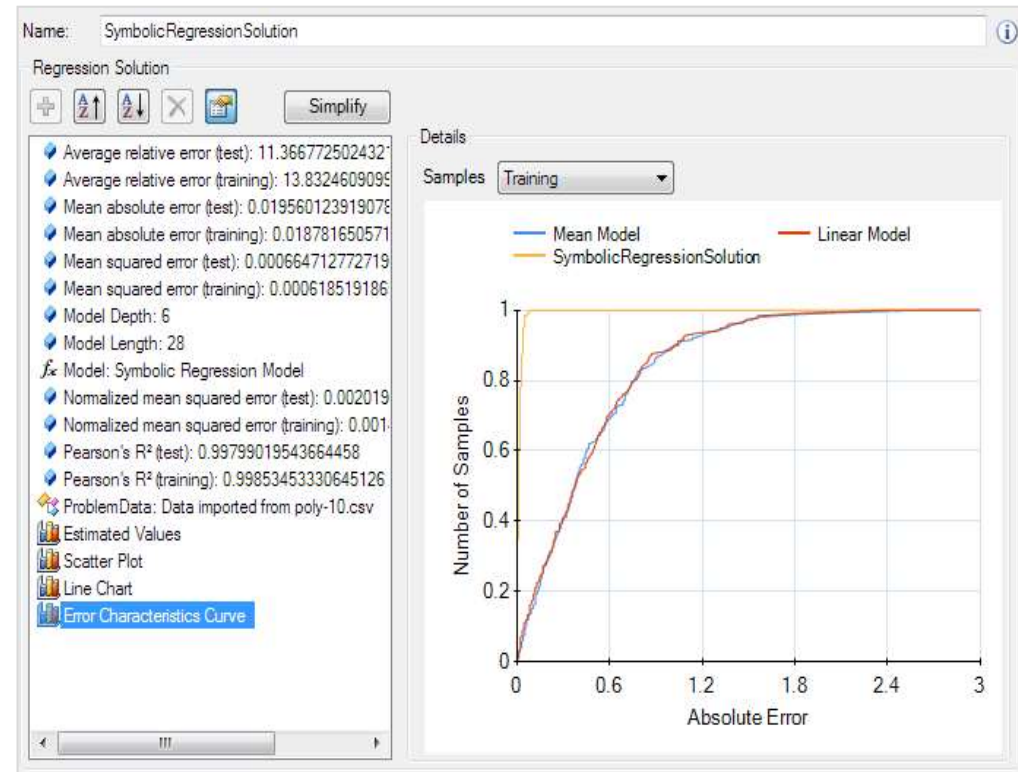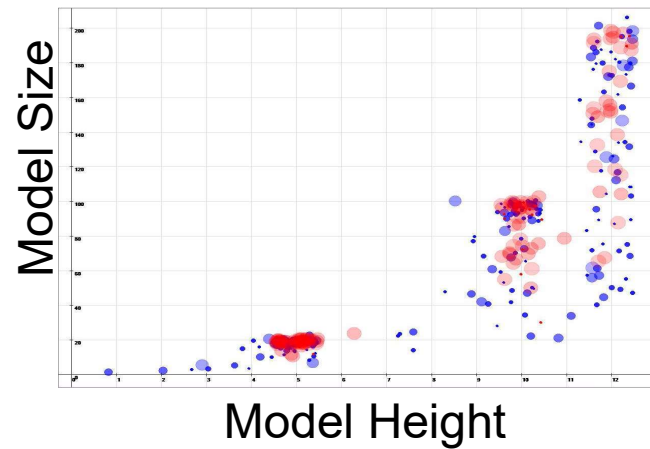$$+ x_1 \cdot x_7 \cdot x_9 + x_3 \cdot x_6 \cdot x_{10}$$

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# Model Evaluation

## Regression



## Classification

| | Id | Target Variable | Estimated Values (all) | Absolute Error (all) | Relative Error (all) |
|---|---|---|---|---|---|
| Row 1 | 0 | -0.051247964 | -0.244931259696236 | 0.193683295696236 | 0.790765931373734 |
| Row 2 | 1 | 0.727691161 | 0.566948971537046 | 0.160742189462954 | 0.283521441139877 |
| Row 3 | 2 | -0.623794992 | -0.235158714563106 | 0.388636277436894 | 1.65265522121487 |
| Row 4 | 3 | 0.184169363 | 0.312577120202989 | 0.128407757202989 | 0.410803443065828 |
| Row 5 | 4 | -0.425409255 | 0.607464911486624 | 1.03287416648662 | 1.70030259683463 |
| Row 6 | 5 | 0.13440877 | 0.135008413403134 | 0.000599643403133... | 0.00444152618358... |
| Row 7 | 6 | 0.723969158 | 1.02967884646345 | 0.305709688463453 | 0.296898095472629 |
| Row 8 | 7 | -0.175618484 | -0.096476538290749 | 0.079141945709251 | 0.820323232066462 |
| Row 9 | 8 | 0.412736644 | 0.559935700149158 | 0.147199056149158 | 0.262885642244183 |
| Row 10 | 9 | 0.321465414 | 0.391061335521024 | 0.0695959215210236 | 0.177966766845663 |
| Row 11 | 10 | 0.492008676 | 0.412907348968929 | 0.0791013270310709 | 0.191571613410599 |

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# Visual Model Exploration
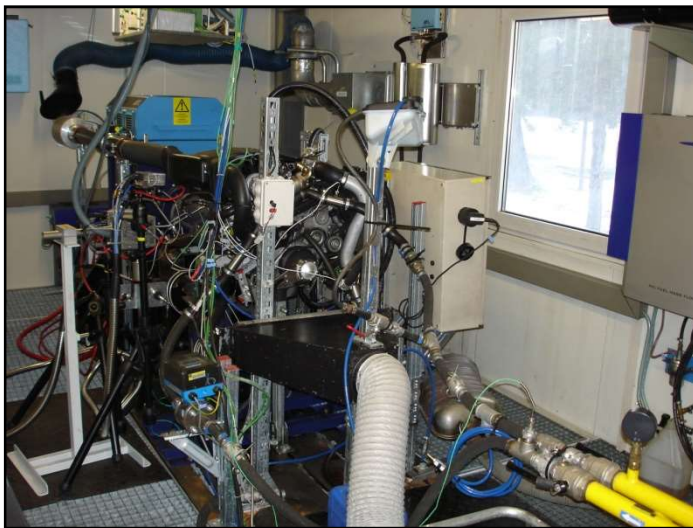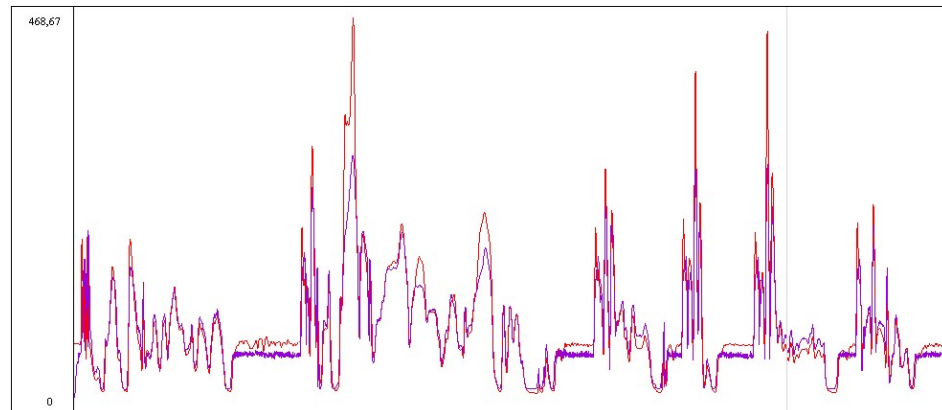
# Example: Virtual Sensors for Modeling Exhaust Gases

- high quality modeling of emissions (NOx and soot) of a diesel engine
- virtual sensors: (mathematical) models that mimic the behavior of physical sensors
- advantages: low cost and non-intrusive
- identify variable impacts: injected fuel, engine frequency, manifold air pressure, concentration of O2 in exhaustion etc.



$$NO_x(t) = f(x1_{(t-7)}, x2_{(t-2)}, \ldots)$$

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# Example: Virtual Sensors for Modeling Exhaust Gases





$$\left[ NO_x^*(t) \right] = \frac{2.696 m_f^*(t-10) + 2.618 m_f^*(t-7)}{\log\left(0.029 N^*(t-10)\right)}$$
$$+ \frac{\left[1.798 m_f^*(t-5)\right]\left[7.536 W^*(t-5)\right]}{\left[0.027 N^*(t-9)\right]\log\left(0.031 N^*(t-3)\right)}$$

Figure 2: Poor performance of the same neural network model at other operating points



Figure 6: Estimated and Measured NOx (System with EGR)

UNIVERSITY OF APPLIED SCIENCES UPPER AUSTRIA

# Example: Blast Furnace Modeling



Model → $f(x)$ → Prognosis

- results as formulas → domain experts can analyze, simplify and refine the models
- integration of prior physical knowledge into modeling process
- powerful data analysis tools: model simplification and variable impact analysis

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

# Example: Plasma Nitriding Modeling

- Motivation
  - › hardening of materials (e.g. transmission parts)
  - › process parameter settings based on expert knowledge

- Modeling Scenarios
  - a) prediction of quality values based on process parameters and material composition
  - b) propose process parameter settings to reach the desired material characteristics

# Example: Medical Data Analysis

- Virtual Tumor markers and cancer diagnosis prediction



Table 3: Modeling results for melanoma diagnosis

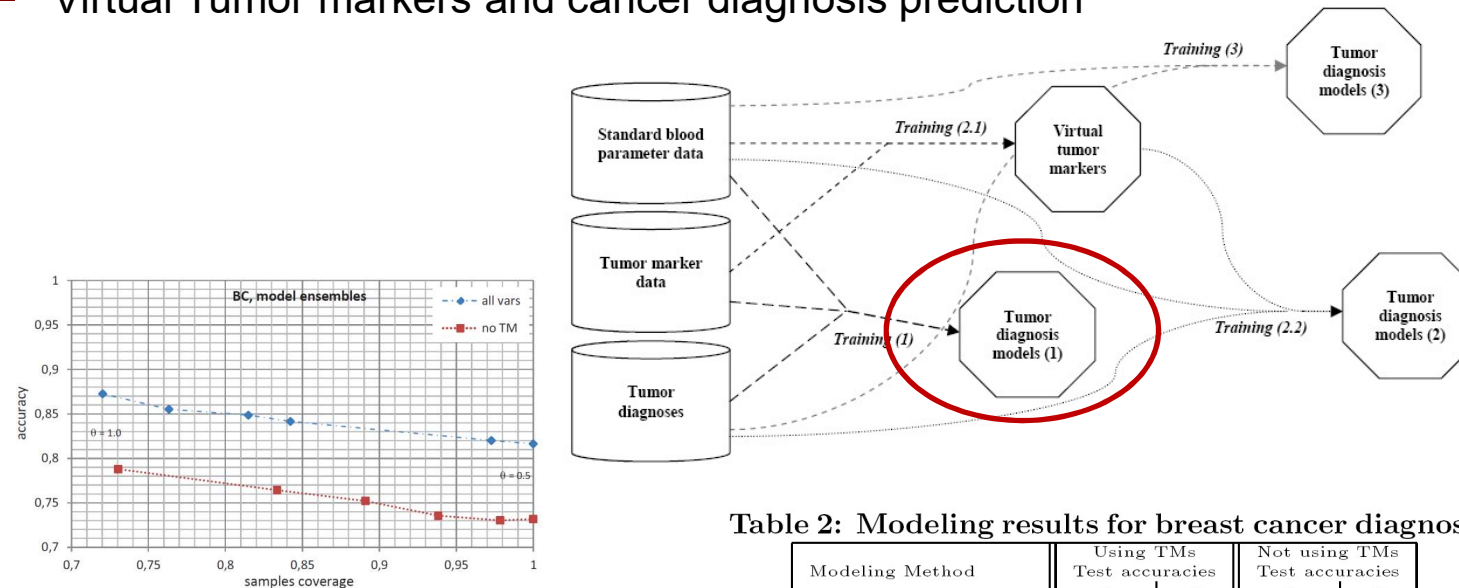| Modeling Method | Using TMs Test accuracies | | Not using TMs Test accuracies | |
|---|---|---|---|---|
| | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ |
| LR, full features set | 73.81% | 3.39 | 71.09% | 4.14 |
| OSGA + LR, $\alpha = 0.0$ | 72.45% | 4.69 | 72.36% | 2.30 |
| OSGA + LR, $\alpha = 0.1$ | 74.73% | 2.35 | 72.09% | 4.01 |
| OSGA + LR, $\alpha = 0.2$ | 73.85% | 2.54 | 72.70% | 2.02 |

Table 4: Modeling results for respiratory system cancer diagnosis

| Modeling Method | Using TMs Test accuracies | | Not using TMs Test accuracies | |
|---|---|---|---|---|
| | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ |
| LR, full features set | 91.32% | 0.37 | 85.97% | 0.27 |
| OSGA + LR, $\alpha = 0.0$ | 91.57% | 0.46 | 86.41% | 0.36 |
| OSGA + LR, $\alpha = 0.1$ | 91.16% | 1.18 | 85.80% | 0.45 |
| OSGA + LR, $\alpha = 0.2$ | 89.45% | 0.37 | 85.02% | 0.15 |
| OSGA + kNN, $\alpha = 0.0$ | 90.98% | 0.84 | 87.09% | 0.46 |
| OSGA + kNN, $\alpha = 0.1$ | 90.01% | 2.63 | 87.01% | 0.83 |
| OSGA + kNN, $\alpha = 0.2$ | 90.16% | 0.74 | 86.92% | 0.81 |
| OSGA + ANN, $\alpha = 0.0$ | 90.28% | 1.63 | 85.97% | 4.07 |
| OSGA + ANN, $\alpha = 0.1$ | 90.99% | 1.97 | 85.82% | 4.52 |
| OSGA + ANN, $\alpha = 0.2$ | 88.64% | 1.87 | 87.24% | 1.91 |
| OSGA + SVM, $\alpha = 0.0$ | 89.03% | 1.38 | 83.12% | 3.79 |
| OSGA + SVM, $\alpha = 0.1$ | 89.91% | 1.58 | 86.25% | 0.79 |
| OSGA + SVM, $\alpha = 0.2$ | 88.33% | 1.94 | 84.66% | 2.06 |
| OSGP, $ms = 50$ | 89.58% | 3.75 | 85.98% | 5.74 |
| OSGP, $ms = 100$ | 90.44% | 3.02 | 86.54% | 6.02 |
| OSGP, $ms = 150$ | 89.58% | 3.75 | 87.97% | 5.57 |

Table 2: Modeling results for breast cancer diagnosis

| Modeling Method | Using TMs Test accuracies | | Not using TMs Test accuracies | |
|---|---|---|---|---|
| | $\mu$ | $\sigma$ | $\mu$ | $\sigma$ |
| LR, full features set | 79.32% | 1.06 | 70.63% | 1.28 |
| OSGA + LR, $\alpha = 0.0$ | 81.78% | 0.21 | 73.13% | 0.36 |
| OSGA + LR, $\alpha = 0.1$ | 81.49% | 1.18 | 72.66% | 0.14 |
| OSGA + LR, $\alpha = 0.2$ | 81.44% | 0.37 | 71.40% | 0.57 |
| OSGA + kNN, $\alpha = 0.0$ | 79.21% | 0.78 | 74.22% | 2.98 |
| OSGA + kNN, $\alpha = 0.1$ | 78.99% | 0.57 | 75.55% | 0.87 |
| OSGA + kNN, $\alpha = 0.2$ | 78.33% | 1.04 | 74.50% | 0.20 |
| OSGA + ANN, $\alpha = 0.0$ | 81.41% | 1.14 | 75.60% | 2.47 |
| OSGA + ANN, $\alpha = 0.1$ | 80.19% | 1.68 | 72.38% | 6.08 |
| OSGA + ANN, $\alpha = 0.2$ | 79.37% | 1.17 | 70.54% | 6.10 |
| OSGA + SVM, $\alpha = 0.0$ | 81.23% | 1.10 | 73.90% | 2.36 |
| OSGA + SVM, $\alpha = 0.1$ | 80.46% | 1.80 | 72.19% | 0.94 |
| OSGA + SVM, $\alpha = 0.2$ | 77.43% | 3.55 | 71.89% | 0.70 |
| OSGP, $ms = 50$ | 79.72% | 1.80 | 75.32% | 0.45 |
| OSGP, $ms = 100$ | 75.50% | 4.95 | 71.63% | 2.75 |
| OSGP, $ms = 150$ | 79.20% | 6.60 | 75.75% | 2.16 |

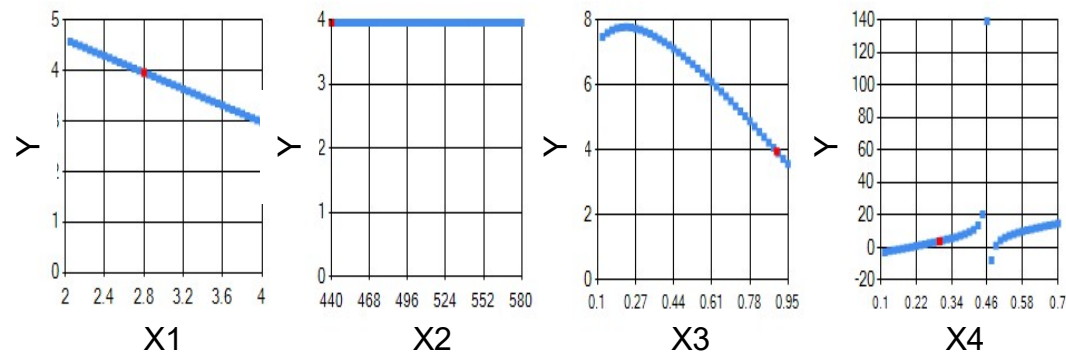UNIVERSITY OF APPLIED SCIENCES UPPER AUSTRIA

# Example: Medical Data Analysis

- Network Identification
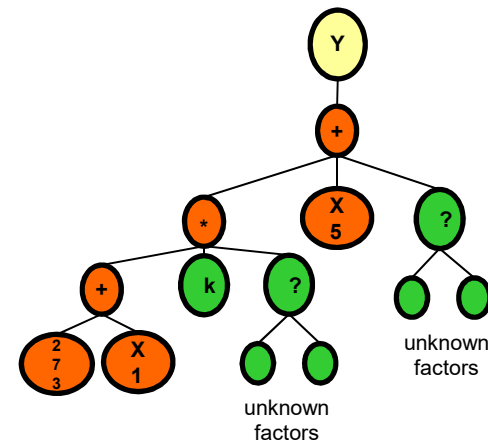- KUK, Prim. Stekel

# Integration of Expert Knowledge
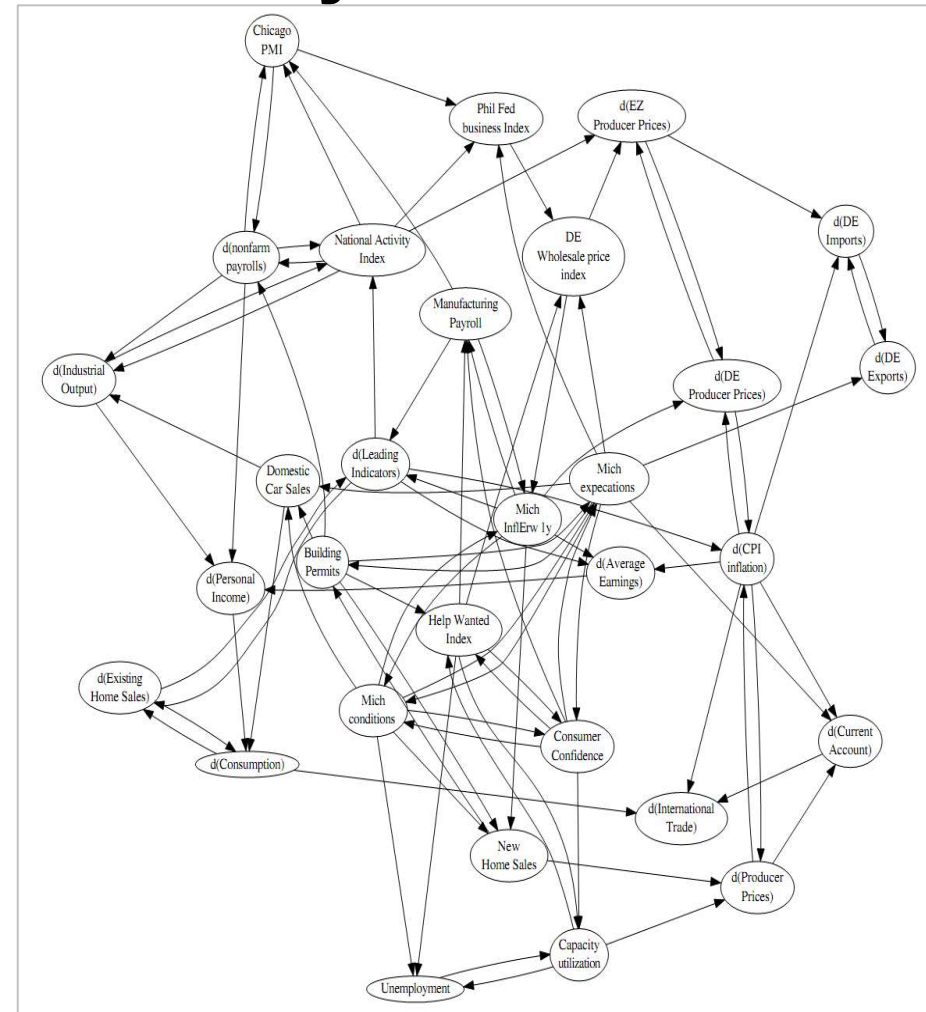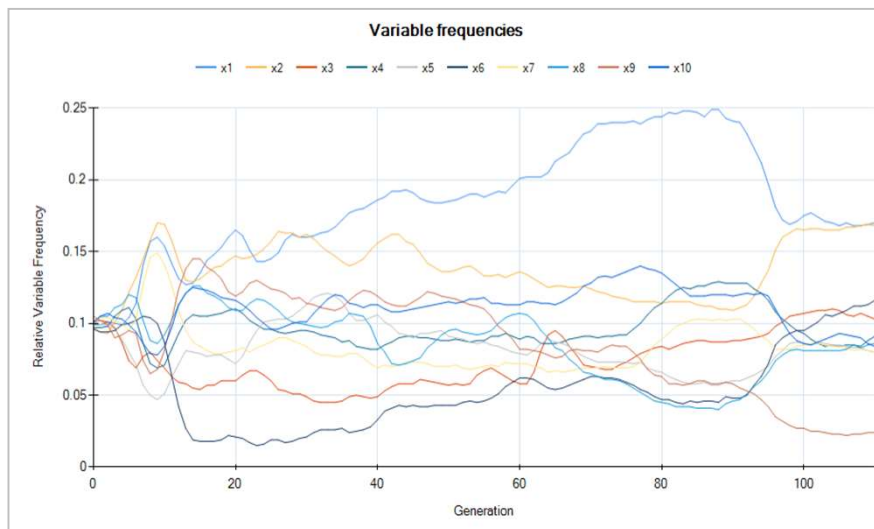
## Model Analysis



## Knowledge Integration

- specification of known correlations
- model extension through algorithm

# Holistic Knowledge Discovery

– Variable interaction networks
  › reveals non-linear correlations

– Variable frequencies
  › analyzed during the algorithm run

# Acknowledgements



**Bioinformatics Research Group**



**Heuristic and Evolutionary Algorithms Laboratory**

http://bioinformatics.fh-hagenberg.at/

https://heal.heuristiclab.com/

UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA