

# Symbolic Regression with Interaction-Transformation

---

Prof. Fabricio Olivetti de França

Federal University of ABC

Center for Mathematics, Computation and Cognition (CMCC)

Heuristics, Analysis and Learning Laboratory (HAL)

12 de Julho de 2018



1. Goals
2. Interaction-Transformation
3. Experiments
4. Conclusions

# Goals

---

## Find the function

Find a function  $f$  that minimizes the approximation error:

$$\begin{array}{ll} \underset{\hat{f}(\mathbf{x})}{\text{minimize}} & \|\epsilon\|^2 \\ \text{subject to} & \hat{f}(\mathbf{x}) = f(\mathbf{x}) + \epsilon. \end{array}$$

# Simple is better

- Ideally this function should be as simple as possible.
- Conflict of interests:
  - minimize approximation (use universal approximators)
  - maximize simplicity (walk away from generic approximators)

The Linear Regression:

$$\hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{w} \cdot \mathbf{x}.$$

- Very simple (and yet useful) model.
- Clear interpretation
- The variables may be non-linear transformation of the original variables.

The mean:

$$\hat{f} = \bar{f}(x)$$

- The average can lie!
- It's just for the sake of the pun

A deep chaining of non-linear transformations that just works!

fig

- Universal approximation
- Alien mathematical expression



# The best?

Despite its success with error minimization, it raises some questions:

- What does the answer mean?
- What if the data is wrong?

# Symbolic Regression

Searches for a function form and the correct parameters.

Hopefully a simple function

**Disclaimer:** I have large experience with evolutionary algorithms, but limited with Symbolic Regression. I have start studying that last year.

# Symbolic Regression

This was my first experience with GP:

$$\begin{aligned} & 6.379515826309025e - 3 + -0.00 * id(x_1^{-4.0} * x_2^{3.0} * x_3^{1.0}) \\ & + -0.00 * id(x_1^{-4.0} * x_2^{3.0} * x_3^{2.0}) - 0.01 * id(x_1^{-4.0} * x_2^{3.0} * x_3^{3.0}) \\ & - 0.02 * id(x_1^{-4.0} * x_2^{3.0} * x_3^{4.0}) + 0.01 * cos(x_1^{-3.0} * x_2^{-1.0}) + \\ & 0.01 * cos(x_1^{-3.0}) + 0.01 * cos(x_1^{-3.0} * x_3^{1.0}) + 0.01 * cos(x_1^{-3.0} * x_2^{1.0}) \\ & + 0.01 * cos(x_1^{-2.0} * x_2^{-2.0}) - 0.06 * log(x_1^{-2.0} * x_2^{-2.0}) \\ & + 0.01 * cos(x_1^{-2.0} * x_2^{-1.0}) + 0.01 * cos(x_1^{-2.0} * x_2^{-1.0} * x_3^{1.0}) \\ & + 0.01 * cos(x_1^{-2.0}) + 0.01 * cos(x_1^{-2.0} * x_3^{1.0}) \\ & + 0.01 * cos(x_1^{-2.0} * x_3^{2.0}) + 0.01 * cos(x_1^{-2.0} * x_2^{1.0}) \\ & + 0.01 * cos(x_1^{-2.0} * x_2^{1.0} * x_3^{1.0}) + -0.00 * id(x_1^{-2.0} * x_2^{2.0}) \\ & - 0.00 * sin(x_1^{-2.0} * x_2^{2.0}) + 0.01 * cos(x_1^{-2.0} * x_2^{2.0}) + \dots \end{aligned}$$

# Why?

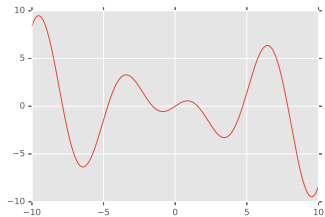
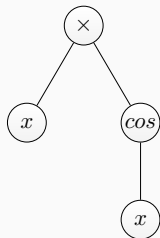
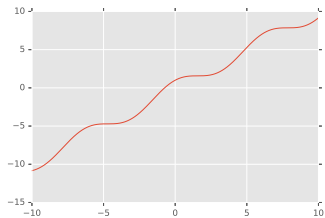
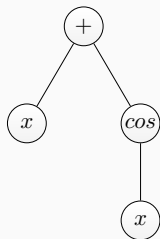
- Infinite search space
- Redundancy
- Rugged

$$f(x) = \frac{x^3}{6} + \frac{x^5}{120} + \frac{x^7}{5040}$$

$$f(x) = \frac{16x(\pi - x)}{5\pi^2 - 4x(\pi - x)}$$

$$\mathbf{f}(\mathbf{x}) = \sin(\mathbf{x}).$$

# Rugged space



# What I wanted

- A few additive terms (linear regression of transformed variables)
- Each term with as an interaction of a couple of variables
- Maximum of one non-linear function applied to every interaction (no chaining)



# Interaction-Transformation

---

Constrains the search space to what I want: a **linear combination** of the application of different **transformation functions** on **interactions** of the original variables.

Essentially, this pattern:

$$\hat{f}(x) = \sum_i w_i \cdot t_i(p_i(x))$$

$$p_i(x) = \prod_{j=1}^d x_j^{k_j}$$

$$t_i = \{id, \sin, \cos, \tan, \sqrt{\phantom{x}}, \log, \dots\}$$

Valid expressions:

- $5.1 \cdot x_1 + 0.2 \cdot x_2$
- $3.5 \sin(x_1^2 \cdot x_2) + 5 \log(x_2^3/x_1)$

Invalid expressions:

- $\tanh(\tanh(\tanh(w \cdot x)))$
- $\sin(x_1^2 + x_2)/x_3$

We can control the complexity of the expression by limiting the number of additive terms and the number of interactions:

$$\begin{aligned}\hat{f}(x) &= \sum_{i=1}^k w_i \cdot t_i(p_i(x)) \\ p_i(x) &= \prod_{j=1}^d x_j^{k_j} \\ \text{s.t. } |\{k_j \mid k_j \neq 0\}| &\leq n\end{aligned}$$

# Interaction-Transformation

Describing as an Algebraic Data Type can help us generalize to other tasks:

```
IT x      = 0 | Weight (Term x) `add` (IT x)
```

```
Term x    = Trans (Inter x)
```

```
Trans     = a -> a
```

```
Inter x:xs = 1 | x s `mul` Inter xs
```

The meaning of `add` and `mul` can lead us to boolean expressions, decision trees, program synthesis.

Simple search heuristic:

```
symtree x leaves | stop      = best leaves
                  | otherwise = symtree x leaves'

where
  leaves' = [expand leaf | leaf <- leaves]

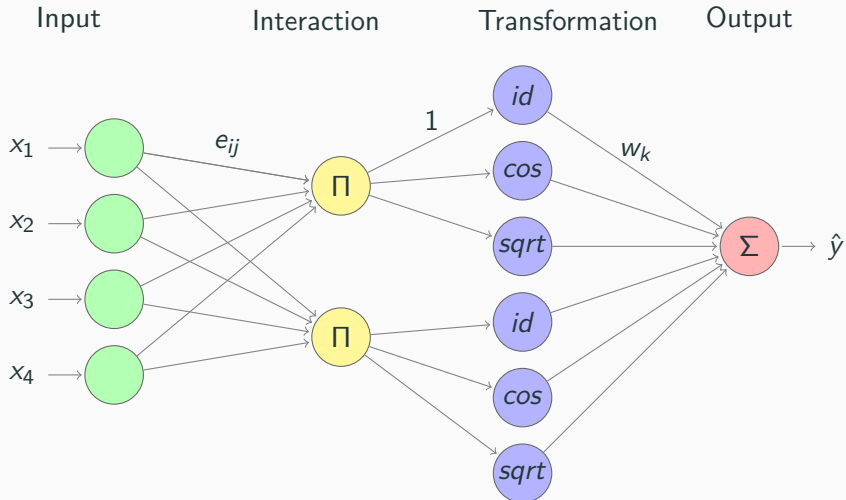
symtree x [linearRegression x]
```



```
expand leaf = expand' leaf terms
  where terms = interaction leaf U transformation leaf

expand' leaf terms = node : expand' leaf leftover
  where (node, leftover) = greedySearch leaf terms
```

Interaction-Transformation Extreme Learning Machine, it generates lots of random interactions, enumerates the transformations for each interaction and then adjust the weight of the terms using  $l_0$  or  $l_1$  regularization.



# Experiments

---

## Data sets

Data set	Features	5-Fold / Train-Test
Airfoil	5	5-Fold
Concrete	8	5-Fold
CPU	7	5-Fold
energyCooling	8	5-Fold
energyHeating	8	5-Fold
TowerData	25	5-Fold
wineRed	11	5-Fold
wineWhite	11	5-Fold
yacht	6	5-Fold
Chemical-I	57	Train-Test
F05128-f1	3	Train-Test
F05128-f2	3	Train-Test
Tower	25	Train-Test

For the sets with folds:

- Each algorithm was run 6 times per fold and the median of the RMSE of the test set is reported
- SymTree was run 1 time per fold (deterministic)

For the sets with train-test split:

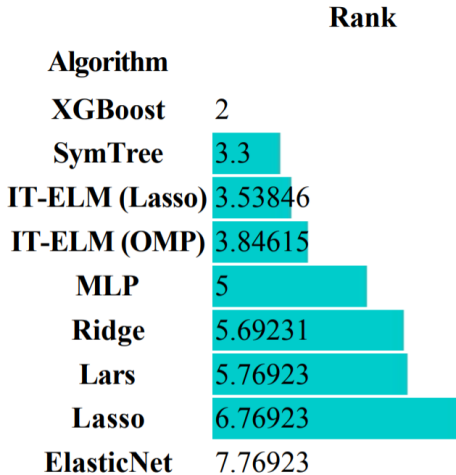
- Each algorithm was run 10 times and the median of the RMSE for the test set is reported
- SymTree was run 1 time per data set

For a complete table:

Binder

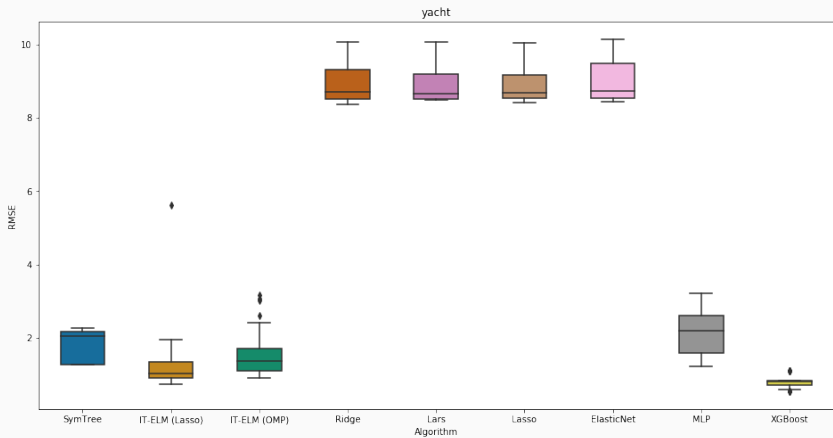
Cell -> Run All

# Results

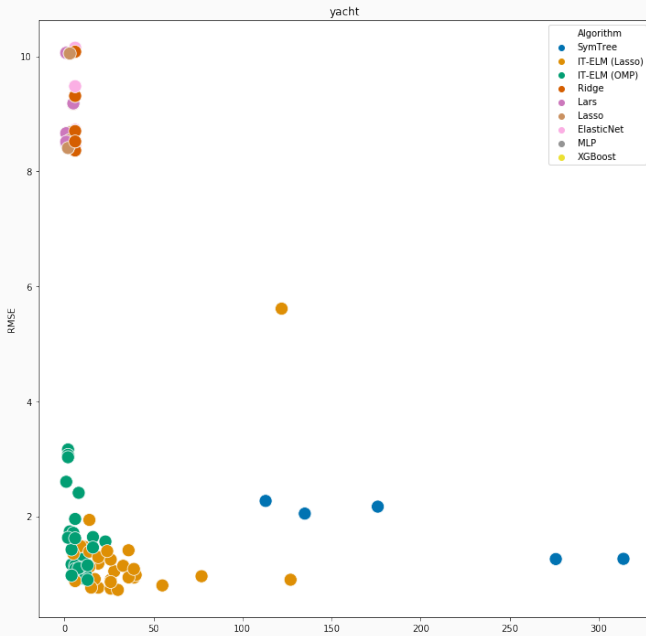




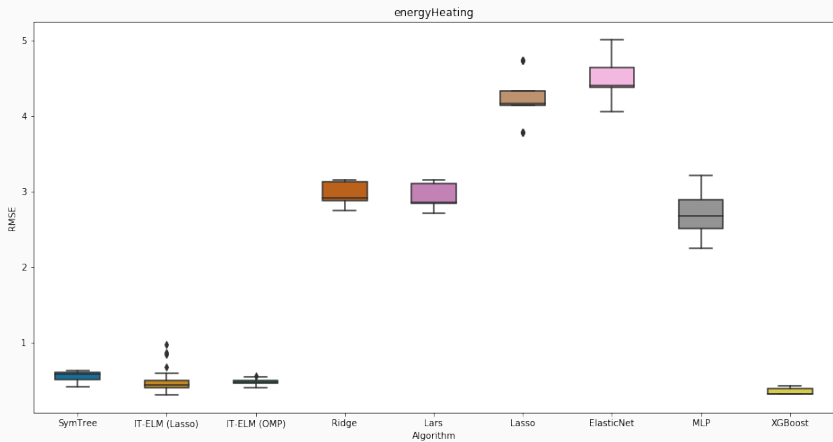
# Results



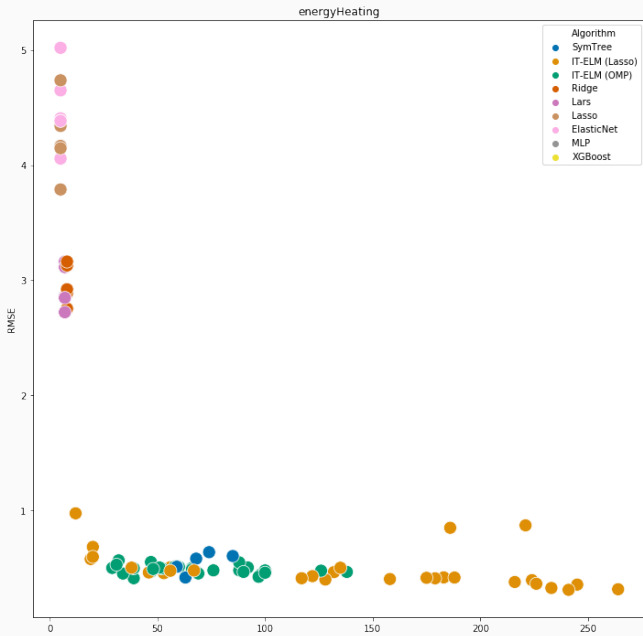
# Results



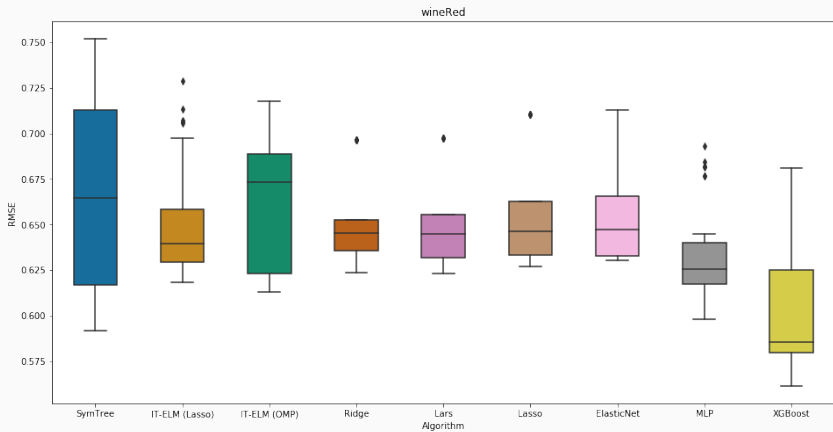
# Results



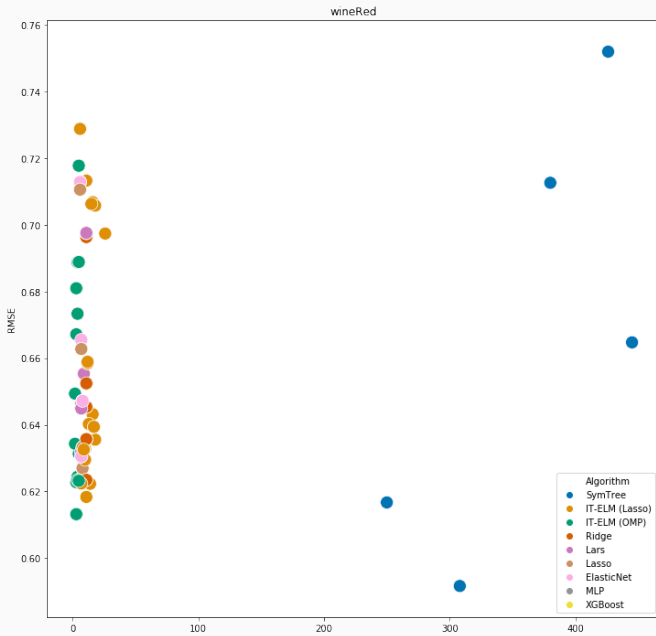
# Results



# Results



# Results



## Sample equation

CPU:

$$\approx 0.86 \cdot \text{cache} + 0.12 \cdot 10^{-6} \cdot \text{maxMem} \sqrt{\text{maxChan} \cdot \text{minMem}}$$

## Sample equation

CPU:

$$\approx 0.86 \cdot \text{cache} + 0.12 \cdot \text{maxMem}(\text{MB}) \sqrt{\text{maxChan} \cdot \text{minMem}(\text{MB})}$$



## Sample equation

CPU:

$$\approx 0.86 \cdot \text{cache} + 0.12 \cdot \text{maxMem}(\text{MB}) \sqrt{\text{maxChan}} \sqrt{\text{minMem}(\text{MB})}$$

More cache = more performance! (sounds about right)

## Sample equation

CPU:

$$\approx 0.86 \cdot \text{cache} + 0.12 \cdot \text{maxMem}(\text{MB}) \sqrt{\text{maxChan}} \sqrt{\text{minMem}(\text{MB})}$$

The original paper isn't clear about it, but it seems that max/min Mem refers to the range of machine tests with a given CPU. So the second term may represent the existence of not measured variables proportional to memory and channels of the experimented machines.

# Conclusions

---

- The Interaction-Transformation representation can help to eliminate *complicated* expressions from the symbolic regression search space.
- Two algorithms created so far: SymTree, a search-based heuristic, and IT-ELM, based on extreme learning machines.
- The results show a good compromise between model accuracy and simplicity.

- Generalization as a Algebraic Data Type
- Use this representation for classification, program synthesis, etc.
- Broaden the search space a little bit
- Explore other search heuristics (evolutionary based)

## Try it!

You can try a lightweight version of SymTree at:

<https://galdeia.github.io/>

It works even on midrange Smartphones!

## Some problems with the provided data sets

The *folds* data sets were provided by some authors that extensively used for GP performance comparison, but:

Forest Fire contains many samples with  $\text{target} = 0$ , because most of the time the forests did not caught fire.

CPU should use the last but one column as the target variable, the last column is just the predicted values from the original paper.