

OBERÖSTERREICH HeuristicLab

EuroCAST 2019

Data Aggregation for Reducing Training Data in Symbolic Regression

Lukas Kammerer, Gabriel Kronberger and Michael Kommenda

Josef Ressel Centre for Symbolic Regression (JRZ) University of Applied Sciences Upper Austria, Hagenberg (HEAL)







JOSEF RESSEL CENTER FOR SYMBOLIC REGRESSION

Contact:

Lukas Kammerer Heuristic and Evolutionary Algorithms Lab (HEAL) Softwarepark 11 4232 Hagenberg, Austria

E-mail: lukas.kammerer@fh-hagenberg.at

Web: http://heal.heuristiclab.com



Regression Algorithms

Cinear Regression

- Find linear relations in data
- Symbolic Regression
 - Find mathematical formulas of any structure
 - Search with genetic programming with strict offspring selection¹
 - Optimize constants with gradient descent
- Random Forest Regression
 - Build many decision trees with differently sampled training data



¹⁾ Affenzeller, M., Wagner, S., Winkler, S., & Beham, A. (2009). *Genetic Algorithms and Genetic Programming: Modern Concepts and Practical Applications*. Chapman and Hall/CRC.

Regression Algorithms

Linear Regression

• Find linear relations in data

Symbolic Regression

- Find mathematical formulas of any structure
- Search with genetic programming with strict offspring selection¹
- Optimize constants with gradient descent

Random Forest Regression

• Build many decision trees with differently sampled training data



Regression Algorithms

Linear Regression

- Find linear relations in data
- Symbolic Regression
 - Find mathematical formulas of any structure
 - Search with genetic programming with strict offspring selection¹
 - Optimize constants with gradient descent

Random Forest Regression

• Build many decision trees with differently sampled training data

Image: Ilan Reinstein, KDnuggets™, www.kdnuggets.com/2017/10/rand om-forests-explained.html



1) Affenzeller, M., Wagner, S., Winkler, S., & Beham, A. (2009). *Genetic Algorithms and Genetic Programming: Modern Concepts and Practical Applications*. Chapman and Hall/CRC.



Make Big Data small again!

GP with Big Data is impractical \rightarrow Problem: Evaluation

Can we reduce training data with without loss in accuracy?





Kugler et al.¹ suggest data clustering

- → Cancel out noise
- → Remove inconsistencies



Used undersampling methods:

- Random sampling
- K-means clustering
- Data binning

7



K-Means Clustering

Finds *K* centroids, so that the total (Euclidean) distance of all points to their closest centroid becomes minimal.

- \rightarrow Stochastic
- \rightarrow Scaling of variables
- → Mini-Batch K-Means¹ Much faster with similar accuracy¹



Image: Scikit-learn: *Machine Learning in Python*, Pedregosa et al., JMLR 12 pp. 2825-2830, 2011



K-Means Clustering

Result of K-Means are K centroids

 \rightarrow Use these centroids as training data!





Data Binning

- 1. Group the data along the target variable into bins of fixed width (similar to a histogram).
- 2. Use the average of each variable per bin as representative.
- **3.** Use these representatives as training data.
 - \rightarrow Deterministic
 - → Imposes n:1 mapping of features to target variable
 - → Should cancel out inconsistencies in process states





Data Binning

- 1. Group the data along the target variable into bins of fixed width (similar to a histogram).
- 2. Use the average of each variable per bin as representative.
- **3.** Use these representatives as training data.



Intuition









Loss in Quality with OSGA



Loss in Quality with OSGA













Conclusion

Results of k-means and random sampling comparable

Implications on... ٠

> hyperparameter optimization? statistical tests? etc.?

- Binning not usable
 - Unusable loss in accuracy in all cases
- Little loss of accuracy with >30% undersampling rate
- Speedup linear to undersampling rate



Discussion

EuroCAST 2019

Data Aggregation for Reducing Training Data in Symbolic Regression

Lukas Kammerer, Gabriel Kronberger and Michael Kommenda

Josef Ressel Centre for Symbolic Regression (JRZ) University of Applied Sciences Upper Austria, Hagenberg (HEAL)







JOSEF RESSEL CENTER FOR SYMBOLIC REGRESSION

Contact:

Lukas Kammerer Heuristic and Evolutionary Algorithms Lab (HEAL) Softwarepark 11 4232 Hagenberg, Austria

E-mail: lukas.kammerer@fh-hagenberg.at Web: http://heal.heuristiclab.com

Algorithm Settings

Offspring Selection Genetic Programming

Max. Tree Size	Depth: 30 Length:50	
Grammar	+ - × ÷ exp log square sqrt cbrt sin cos tan	
Population Size	300	
Mutation Rate	20 %	
Selector	GenderSpecific (Random & Proportional	



Random Forest

R	30 %
Μ	50 %
Number of Trees	50

Image: Affenzeller, M., Wagner, S., & Winkler, S. (2007, July). *Aspects of adaptation in natural and artificial evolution*. In Proceedings of the 9th annual conference companion on Genetic and evolutionary computation (pp. 2595-2602). ACM