

EVOFEAT: GENETIC PROGRAMMING BASED FEATURE ENGINEERING APPROACH TO TABULAR DATA CLASSIFICATION

HENGZHE ZHANG, QI CHEN, BING XUE, YAN WANG, AIMIN ZHOU, MENGJIE ZHANG

VICTORIA UNIVERSITY OF WELLINGTON

23/05/2023

- 1 Introduction
- 2 Related Work
- 3 Preliminaries
- 4 The Proposed Algorithm
- 5 Experiments

INTRODUCTION

- **Tabular Data Learning:** Widely used in recommendation systems ¹ and advertising ².
- **Goal:** Capture the relationship between explanatory variables $\{x_1, \dots, x_m\}$ and a response variable y .
- **Dataset Structure:** $\{(\{x_1^1, \dots, x_m^1\}, y^1), \dots, (\{x_1^n, \dots, x_m^n\}, y^n)\}$, where n is the number of instances.

Challenge

- Linear models assume linear relationships.
- Decision trees assume axis-parallel decision boundaries.
- Real-world data often violates these assumptions.

¹Ruoxi Wang et al., *Proceedings of the Web Conference 2021* (2021)

²Haizhi Yang et al., *Proceedings of the 30th ACM International Conference on Information & Knowledge Management* (2021)

- **Manual Design:** Based on domain knowledge.
- **Kernel Methods:** Use kernel tricks to transform data into higher dimensions.
- **Deep Learning:** Leverages neural networks to learn features automatically. ¹

Limitations

- Manual Design: Labor-intensive.
- Kernel Methods: Hard to integrate with tree-based methods.
- Deep Learning: Requires large datasets, effectiveness debatable for small, heterogeneous datasets ².

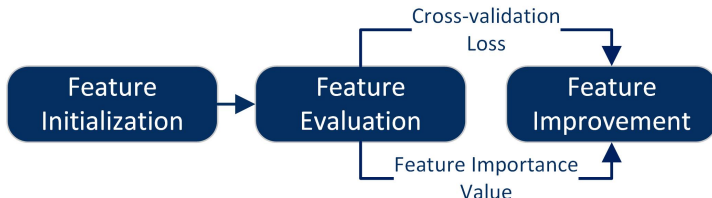
¹Jianxun Lian et al., *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2018)

²Yury Gorishniy et al., *Advances in Neural Information Processing Systems* (2021)

- **Objective:** Feature construction using genetic programming (GP).
- **GP Advantages:** Gradient-free, interpretable, and flexible.
- **Hypothesis:** GP-based feature engineering can outperform both traditional and deep learning methods on tabular data.

Our Approach: EvoFeat

- Constructs nonlinear features with GP.
- Enhances ensemble learning models.
- Uses cross-validation and feature importance for evaluation.



RELATED WORK

■ **Beam Search Methods:**

- ▶ Greedy, lacks strong mechanisms to prevent overfitting.

■ **Deep Learning Methods:**

- ▶ Effectiveness in comparison to tree-based methods is still debated ¹.

¹Yury Gorishniy et al., *Advances in Neural Information Processing Systems* (2021)

■ Iterative Feature Generation:

- ▶ Starts with low-order features.
- ▶ Generates higher-order features based on important low-order features ¹.

■ Evaluation:

- ▶ Uses logistic regression accuracy, or XGBoost feature importance.
- ▶ Sole reliance on training loss can lead to overfitting.

Key Limitation

Lack of effective mechanisms to prevent overfitting restricts feature construction capabilities.

¹Yuanfei Luo et al., *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2019)

■ High-Order Feature Construction:

- ▶ Cross Network in DCN.
- ▶ Field-wise feature cross in xDeepFM ¹.
- ▶ Attention mechanism in AutoInt ².

Effectiveness

- Effectiveness over fully connected NN is debatable ³.
- Lack of comprehensive studies comparing with XGBoost ⁴.

¹Jianxun Lian et al., *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2018)

²Weiping Song et al., *Proceedings of the 28th ACM International Conference on Information and Knowledge Management* (2019)

³Ruoxi Wang et al., *Proceedings of the Web Conference 2021* (2021)

⁴Yury Gorishniy et al., *Advances in Neural Information Processing Systems* (2021)

■ **Single Learner:**

- ▶ Traditionally, more focus on simple learner like single decision tree ¹.
- ▶ Gap in enhancing state-of-the-art algorithms.

■ **Ensemble-based Feature Construction:**

- ▶ Promising results in regression ².
- ▶ Requires adaptation for tabular classification.

Notes

Adapting evolutionary feature construction techniques for classification involves:

- Adapting loss functions.
- Using logistic regression models as base learners.

¹Binh Tran, Bing Xue, and Mengjie Zhang, *Pattern Recognition* (2019)

²Hengzhe Zhang, Aimin Zhou, and Hu Zhang, *IEEE Transactions on Evolutionary Computation* (2021)

PRELIMINARIES

■ Feature Initialization:

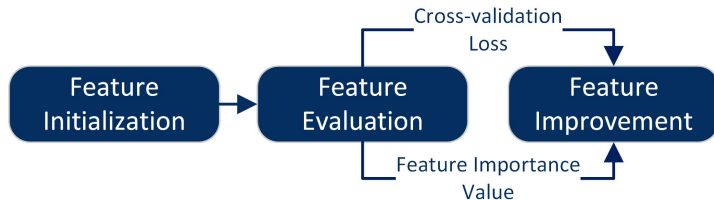
- ▶ Construct initial features based on domain knowledge or randomly.

■ Feature Evaluation:

- ▶ Evaluate features using cross-validation and calculate feature importance.

■ Feature Improvement:

- ▶ Discard ineffective features and replace with new ones derived from important features.



Feature engineering workflow.

■ **Cross-Validation:**

- ▶ Evaluates generalization performance.

■ **Feature Importance:**

- ▶ Identifies useful features.
- ▶ Risky to rely solely on feature importance.

Key Insight

Constructing multiple sets of features and evaluating them using cross-validation can provide better insights into their generalization capabilities.

THE PROPOSED ALGORITHM

■ **Symbolic Trees:**

- ▶ Each individual has k GP trees representing k new features.

■ **Tree Structure:**

- ▶ Non-leaf nodes: Functions (e.g., $+$, $-$, $*$, \log , \sin).
- ▶ Leaf nodes: Original Features.

■ **Base Learners:**

- ▶ Decision trees or linear regression models.

Initialization

Randomly initialize N individuals, each with k symbolic trees.

Evaluation

- Evaluate individuals using cross-validation loss.
- Calculate feature importance for each feature.

Selection

Use lexicase selection ¹ to select parent individuals based on cross-validation losses.

Generation

Generate new individuals using self-competitive crossover and guided mutation ².

Archive Update

Update archive with top-performing models using reduce-error pruning ³.

¹William La Cava et al., *Evolutionary Computation* (2019)

²Hengzhe Zhang et al., *IEEE Transactions on Evolutionary Computation* (2023)

³Rich Caruana et al., *Proceedings of the Twenty-First International Conference on Machine Learning* (2004)

■ Initialization Strategy:

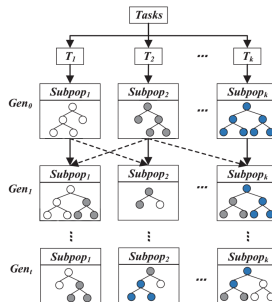
- ▶ Ramped-half-and-half for symbolic trees.
- ▶ Half full trees, half random depth.

■ Base Learner Assignment:

- ▶ Randomly assign **decision tree** or **linear regression model**.

- Three selection operators in EvoFeat:
 - ▶ Base Learner Selection
 - ▶ Individual Selection: Lexicase Selection
 - ▶ Feature Selection: Softmax Selection

- Divide population into two subgroups (decision trees, logistic regression).
- Random mating probability ($rpm = 0.5$):
 - ▶ 50%: Select parents from different subgroups.
 - ▶ 50%: Select parents from the same subgroup.



Inspired by multitask GP ¹

¹Fangfang Zhang et al., *IEEE Transactions on Cybernetics* (2021)

- Selects individuals based on a vector of cross-validation losses, one for each instance.
- Constructs filters based on each loss value ¹:

$$\tau_j = \min_i \mathcal{L}_j^i + \epsilon_j, \quad (1)$$

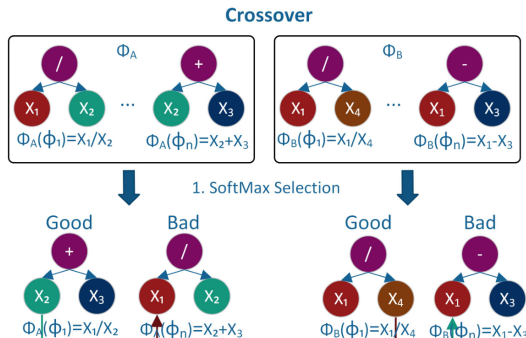
- Where:
 - ▶ τ_j is the threshold,
 - ▶ \mathcal{L}_j^i is the loss of the i -th individual on the j -th instance,
 - ▶ ϵ_j is the median absolute deviation.

¹William La Cava et al., *Evolutionary Computation* (2019)

- Select features based on importance values $\{\theta_1, \dots, \theta_k\}$.
- Uses softmax function:

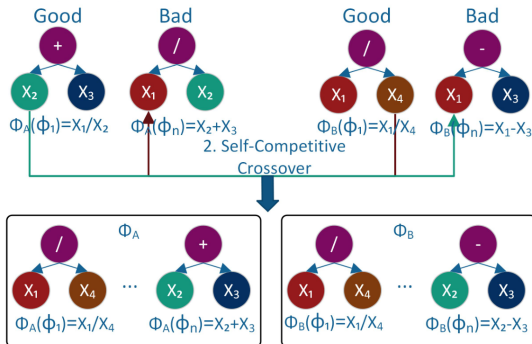
$$P(\theta_i) = \frac{e^{\theta_i/T}}{\sum_{j=1}^k e^{\theta_j/T}}, \quad (2)$$

- Good features sampled by $P(\theta_i)$, bad features by $P(-\theta_i)$.



■ Self-Competitive Crossover:

- ▶ Transfers beneficial material from good features to bad features.
- ▶ Biased crossover, only modifies bad features, preserving good features ¹.
- ▶ Ensures top-performing features are preserved.



¹Su Nguyen et al., *IEEE Transactions on Cybernetics* (2021)

■ Decision Tree:

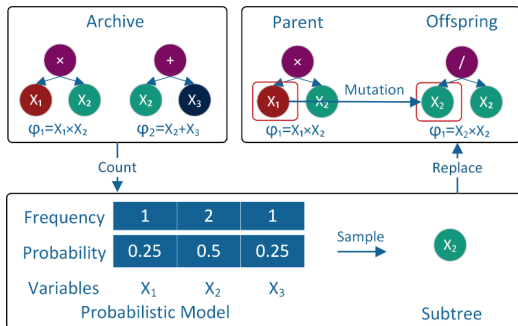
- ▶ Calculated by the total reduction of Gini impurity contributed by each feature ϕ .

■ Logistic Regression:

- ▶ Calculated by the absolute value of the model coefficients.
- ▶ Features are standardized to ensure equal influence on the coefficients.

■ Guided Mutation:

- Replaces subtree with a randomly generated subtree.
- Uses guided probability vector for terminal variable selection.
- Probability vector corresponds to the terminal usage of archived individuals.



■ Cross-Validation:

- ▶ Partition training set into five folds.
- ▶ Train on four folds, validate on one fold.

■ Loss Function:

- ▶ Cross entropy:

$$\sum_{c \in C} p_c * \log(q_c), \quad (3)$$

- ▶ Where p_c is the true probability, q_c is the predicted probability.

EXPERIMENTS

- **Objective:** Compare EvoFeat with popular machine learning and deep learning methods.
- **Datasets:** 130 datasets from DIGEN and PMLB benchmarks.
 - ▶ DIGEN ¹: Diverse synthetic datasets using genetic programming.
 - ▶ PMLB ²: Real-world datasets from OpenML.

¹Patryk Orzechowski and Jason H Moore, *Science Advances* (2022)

²Joseph D Romano et al., *Bioinformatics* (2022)

■ Evaluation Protocol:

- ▶ 80% training, 20% testing.
- ▶ 5-fold cross-validation on training set.
- ▶ Repeat experiments with 30 random seeds.

■ Hyperparameter Tuning:

- ▶ Use HEBO ¹ for tuning baseline algorithms.

¹Alexander I Cowen-Rivers et al., *Journal of Artificial Intelligence Research* (2022)

■ Machine Learning:

- ▶ XGBoost ¹, LightGBM ², Random Forest (RF), Decision Tree (DT), Logistic Regression (LR), K-Nearest Neighbors (KNN).

■ Deep Learning:

- ▶ Multilayer Perceptron (MLP), ResNet, DCN V2 ³, FT-Transformer ⁴.

¹Tianqi Chen and Carlos Guestrin, *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2016)

²Guolin Ke et al., *Advances in Neural Information Processing Systems* (2017)

³Ruoxi Wang et al., *Proceedings of the Web Conference 2021* (2021)

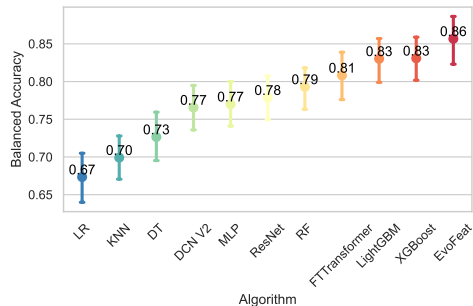
⁴Yury Gorishniy et al., *Advances in Neural Information Processing Systems* (2021)

■ Comparison:

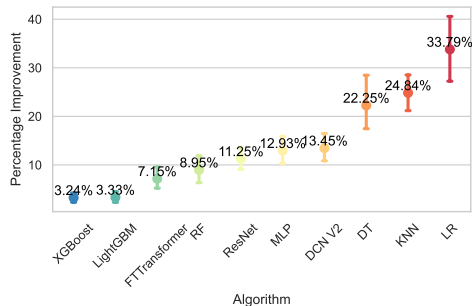
- ▶ Evaluate EvoFeat against traditional and deep learning methods.

■ Results:

- ▶ EvoFeat outperforms state-of-the-art methods on average accuracy.
- ▶ Demonstrates large improvements in predictive performance.



(a) Balanced testing accuracy.



(b) Improvement in accuracy.

- **Baseline:** XGBoost, LightGBM, RF, DT, LR, KNN.
- **Results:**
 - ▶ EvoFeat achieves better the best accuracy.
 - ▶ Significant improvements over XGBoost and LightGBM.

Statistical results of balanced testing accuracy on 90 PMLB and 40 DIGEN datasets.

	XGBoost	LightGBM	RF	LR	KNN	EvoFeat
DT	0/48/82	2/47/81	0/43/87	60/36/34	60/27/43	0/34/96
XGBoost	—	13/107/10	43/79/8	72/50/8	107/16/7	4/67/59
LightGBM	—	—	45/75/10	74/42/14	107/15/8	5/72/53
RF	—	—	—	73/47/10	102/20/8	7/62/61
LR	—	—	—	—	54/13/63	7/44/79
KNN	—	—	—	—	—	3/15/112

■ **Baseline:** MLP, ResNet, DCN V2, FT-Transformer.

■ **Results:**

- ▶ Deep learning methods perform comparably to RF.
- ▶ EvoFeat outperforms these deep learning methods significantly.

Statistical results of balanced testing accuracy on 90 PMLB and 40 DIGEN datasets.

	ResNet	DCN V2	FT-Transformer	EvoFeat
MLP	18/96/16	9/118/3	10/76/44	4/33/93
ResNet	—	8/99/23	46/73/11	3/32/95
DCN V2	—	—	45/79/6	2/35/93
FT-Transformer	—	—	—	4/34/92
EvoFeat	—	—	—	—

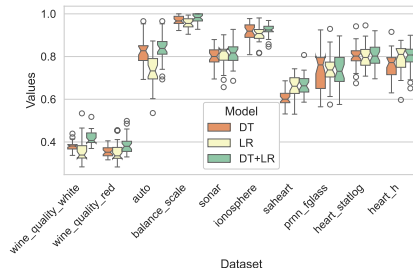
- **Objective:** Validate improvements from heterogeneous base learners and feature importance-guided search.
- **Components:**
 - ▶ **Heterogeneous base learners:** Compare EvoFeat with different combinations of base learners.
 - ▶ **Feature importance-guided search:** Evaluate the effectiveness of feature importance-guided operators.

- **Objective:** Compare heterogeneous base learners (DT+LR) with single base learners (DT, LR).
- **Results:**
 - ▶ DT+LR achieves better average performance.
 - ▶ Significant improvements over single learners.

Comparison of balanced testing accuracy across different base learners on 90 PMLB datasets.

	LR	DT+LR
DT	12(+)/47(~)/31(-)	0(+)/62(~)/28(-)
LR	—	5(+)/70(~)/15(-)

- **Objective:** Compare heterogeneous base learners (DT+LR) with single base learners (DT, LR).
- **Results:**
 - ▶ DT+LR achieves better average performance.
 - ▶ Significant improvements over single learners.



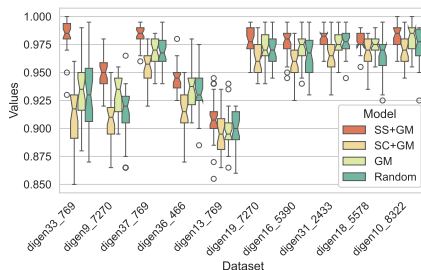
Balanced testing accuracy with different base learners.

- **Objective:** Evaluate effectiveness of feature importance-guided operators.
- **Methods:**
 - ▶ Compare random crossover and mutation (Random) with softmax-based self-competitive crossover and guided mutation (SS+GM).
- **Results:**
 - ▶ Feature importance-guided search achieves better performance.

Comparison of balanced testing accuracy across different selection operators on 40 DIGEN datasets.

	SC+GM	GM	Random
SS+GM	12(+)/26(~)/2(-)	5(+)/34(~)/1(-)	12(+)/28(~)/0(-)
SC+GM	—	0(+)/30(~)/10(-)	5(+)/30(~)/5(-)
GM	—	—	5(+)/35(~)/0(-)

- **Objective:** Evaluate effectiveness of feature importance-guided operators.
- **Methods:**
 - ▶ Compare random crossover and mutation (Random) with softmax-based self-competitive crossover and guided mutation (SS+GM).
- **Results:**
 - ▶ Feature importance-guided search achieves better performance.



Balanced testing accuracy with different selection operators.

■ **Summary:**

- ▶ EvoFeat outperforms state-of-the-art methods.
- ▶ Heterogeneous base learners and feature importance-guided search improve performance.

■ **Future Work:**

- ▶ Investigate modularization techniques for improved interpretability.
- ▶ Use diversity optimization to enhance ensemble performance.

THANKS FOR LISTENING!

EMAIL: HENGZHE.ZHANG@ECS.VUW.AC.NZ

GITHUB PROJECT: [HTTPS://GITHUB.COM/HENGZHE-ZHANG/EVOLUTIONARYFOREST/](https://github.com/hengzhe-zhang/evolutionaryforest/)