

May 6th, 2026  
First Dutch Workshop for  
Combatting Crime,  
Amsterdam,

# **A Hidden Formula for Fraud? Equation Learning for Fraud Detection**

Erman Acar

Informatics Institute  
Institute for Logic, Language and Computation



# TODAY'S ROADMAP

- Introduction
- What is Symbolic Regression?
- Two Research Papers
  - Explainable Fraud Detection through Deep Symbolic Classification (XAI 2024)
  - ECSEL: Explainable Classification via Signomial Equation Learning (ICML 2026)
- Final Takeaways



# A brief introduction

*“Who Am I? Ah, that’s the greatest puzzle”* — Lewis Carroll

Erman Acar (Assistant Professor for Safe & Explainable AI in Finance)

Affiliation: Informatics Institute & Institute of Logic, Language and Computation, at UvA

Research: XAI, Neurosymbolic Systems, Multi-agent Systems, Causality

Role: Leading Finesse Lab



**FINESSE LAB**  
Finance-inspired Neurosymbolic Systems  
for Safety & Explainability



**Adia Lumadjeng**  
PhD Student



**Andreas Sauter (VU)**  
PhD Student



**Angela van Sprang**  
PhD Student



**Arco van Breda**  
PhD Student



**Johannes Bendler (VU)**  
PhD Student



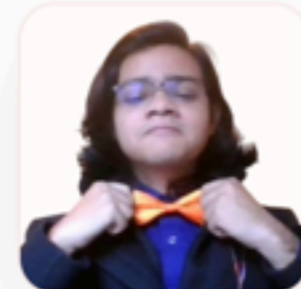
**Mayesha Tasnim**  
PhD Student



**Philip Wozny (VU)**  
PhD Student



**Raj Bhalwankar**  
PhD Student



**Satchit Chatterji**  
PhD Student



# FINESSE LAB

Finance-inspired Neurosymbolic Systems  
for Safety & Explainability

Qualities

Research Theme

Safety & Control

Circuit-Level Mechanisms

Explainability by design

- Equation Discovery
- Causality-driven AI/  
Fraud Detection
- Foundation Models  
Forecasting
- Multimodal Reasoning
- Multiagent Systems





# FINESSE LAB

Finance-inspired Neurosymbolic Systems  
for Safety & Explainability

*Promoting socially responsible, financially compliant AI Research.*

## For People

Starting from societal and stakeholder needs, we try to address fairness, transparency, accountability, and recourse.

## By People

We integrate human expert knowledge directly into AI through rules, logic, policy constraints, and stakeholder values.

## With People

Keeping humans in the loop, we see AI as a support for analysts, auditors, and citizens, not a black-box replacement.

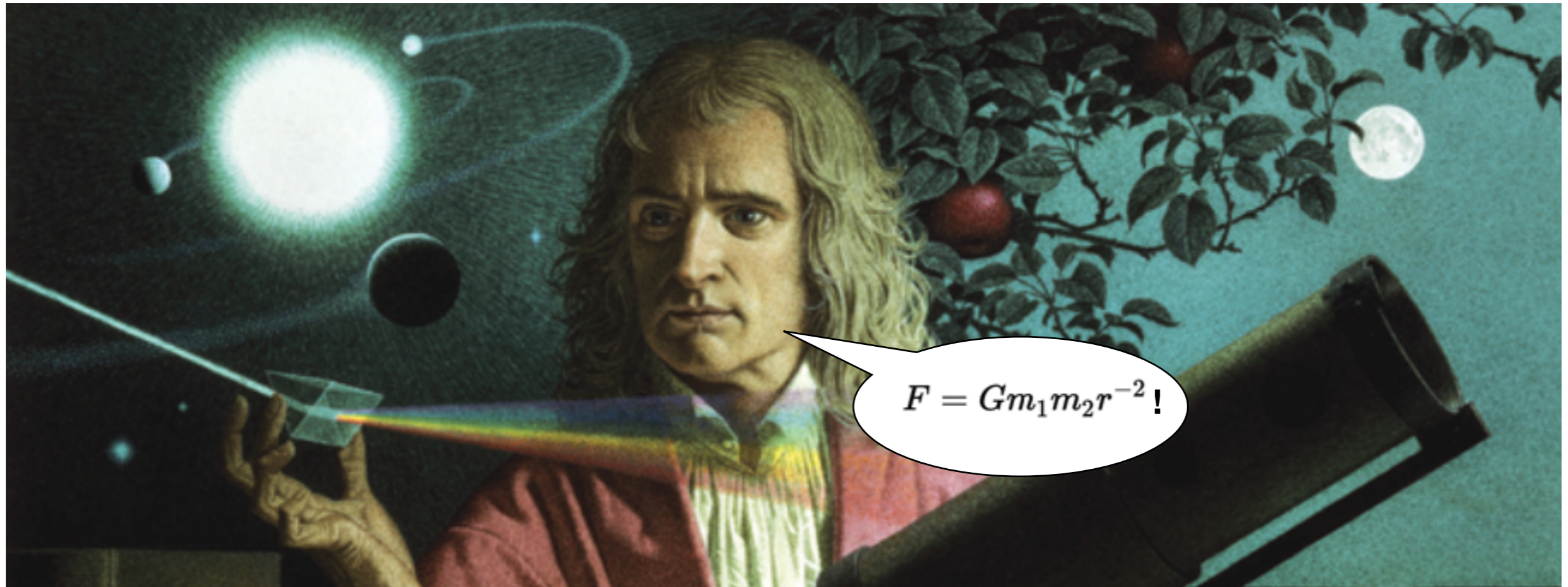
*Finesse is built on SIAS principles.*



# What is Symbolic Regression?

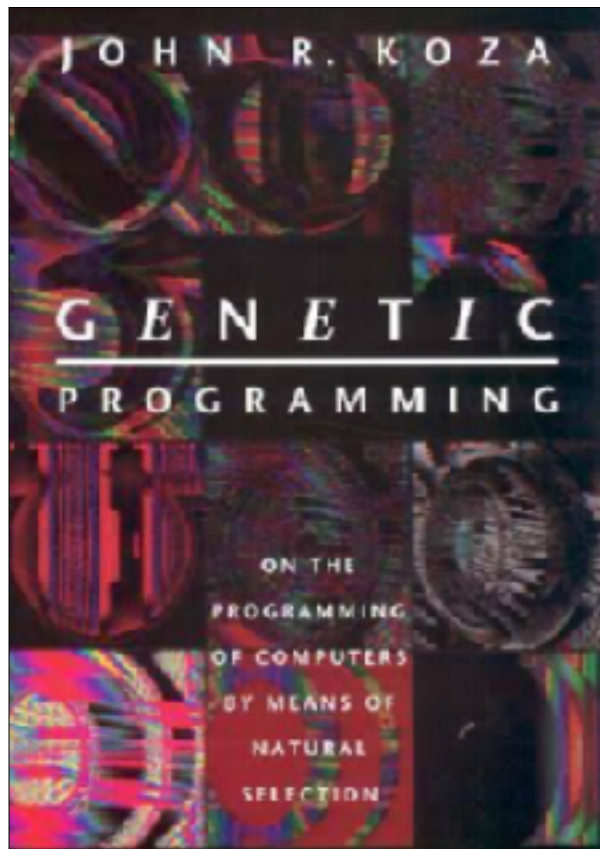


# What is the law behind the data?

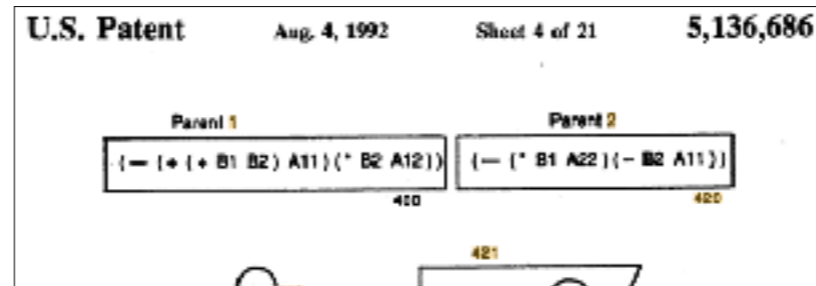


$G$	$m_1$ kg	$m_2$ kg	$r$ m	$F$ N
$6.674 \times 10^{-11}$	10	20	2	$3.337 \times 10^{-9}$
$6.674 \times 10^{-11}$	10	20	4	$8.343 \times 10^{-10}$
$1.000 \times 10^{-10}$	30	20	2	$1.500 \times 10^{-8}$
$5.000 \times 10^{-11}$	30	40	3	$6.667 \times 10^{-9}$
$8.000 \times 10^{-11}$	50	40	5	$6.400 \times 10^{-9}$
$4.000 \times 10^{-11}$	100	80	10	$3.200 \times 10^{-9}$

# Symbolic Regression: An old forgotten child of Machine Learning



Richard Feynman (1918-1988)



Feynman eq.	Equation	Solution time (s)
I.6.20a	$f = e^{-\theta^2/2} / \sqrt{2\pi}$	16
I.6.20	$f = e^{-\frac{\theta^2}{2\sigma^2}} / \sqrt{2\pi\sigma^2}$	2992
I.6.20b	$f = e^{-\frac{(\theta-\theta_1)^2}{2\sigma^2}} / \sqrt{2\pi\sigma^2}$	4792
I.8.14	$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$	544
I.9.18	$F = \frac{Gm_1m_2}{(x_2-x_1)^2 + (y_2-y_1)^2 + (z_2-z_1)^2}$	5975
I.10.7	$m = \frac{m_0}{\sqrt{1-\frac{v^2}{c^2}}}$	14
I.11.19	$A = x_1y_1 + x_2y_2 + x_3y_3$	184
I.12.1	$F = \mu N_n$	12
I.12.2	$F = \frac{q_1q_2}{4\pi\epsilon_0r^2}$	17
I.12.4	$E_f = \frac{q_1}{4\pi\epsilon_0r^2}$	12
I.12.5	$F = q_2E_f$	8
I.12.11	$F = q(E_f + Bv\sin\theta)$	19
I.13.4	$K = \frac{1}{2}m(v^2 + u^2 + w^2)$	22
I.13.12	$U = Gm_1m_2(\frac{1}{r_2} - \frac{1}{r_1})$	20
I.14.3	$U = mgz$	12
I.14.4	$U = \frac{k_{spring}x^2}{2}$	9
I.15.3x	$x_1 = \frac{x-ut}{\sqrt{1-u^2/c^2}}$	22

Expr

$x_1$	$x_2$	$x_3$	$y$
0.00	1.00	2.00	6.00
1.57	0.50	4.00	7.00
3.14	2.00	1.00	6.00
0.79	3.00	0.50	5.21



$$y = \sin(x_1) + 3x_2x_3$$

ed: 13 May 2019

ning black box machine learning models for decisions and use interpretable models

Science 1, 206-215 (2019) | [Cite this article](#)

13 Citations 524 Altmetric [Metrics](#)

## AI Feynman Dataset: An arena for Symbolic Regressors



# Blooming of Symbolic Regression

## Interpretable Machine Learning for Science with PySR and SymbolicRegression.jl

Miles Cranmer<sup>1,2</sup>

<sup>1</sup>Princeton University, Princeton, NJ, USA

<sup>2</sup>Flatiron Institute, New York, NY, USA

May 2, 2023

## GFN-SR: Symbolic Regression with Generative Flow Networks

Sida Li  
The University of Chicago  
l1star2000@uchicago.edu

Ioana Marinescu  
Princeton University  
ioanan@princeton.edu

Sebastian Musslick  
University of Osnabrück, Brown University  
sebastian.musslick@uos.de

Published as a conference paper at ICLR 2021

## DEEP SYMBOLIC REGRESSION: RECOVERING MATHEMATICAL EXPRESSIONS FROM DATA VIA RISK-SEEKING POLICY GRADIENTS

Brenden K. Peersen\*  
Lawrence Livermore National Laboratory  
Livermore, CA, USA  
bpeersi@llnl.gov

Mikel Landajuela Larra  
Lawrence Livermore National Laboratory  
Livermore, CA, USA  
landajue.m1a1@llnl.gov

T. Nathan Mundhenk  
Lawrence Livermore National Laboratory  
Livermore, CA, USA  
mundhenk1@llnl.gov

Claudio P. Santiago  
Lawrence Livermore National Laboratory  
Livermore, CA, USA  
santiago10@llnl.gov

Soo E. Kim  
Lawrence Livermore National Laboratory  
Livermore, CA, USA  
kim79@llnl.gov

Joanne T. Kim  
Lawrence Livermore National Laboratory  
Livermore, CA, USA  
kim102@llnl.gov

## SymbolicGPT: A Generative Transformer Model for Symbolic Regression

Mojtaba Valipour\*  
University of Waterloo  
mojtaba.valipour@uwaterloo.ca

Bowen You †  
University of Waterloo  
byyou@uwaterloo.ca

Maysum Panju †  
University of Waterloo  
mpanju@uwaterloo.ca

Ali Ghodsi †  
University of Waterloo

## Neural Symbolic Regression that scales

Luca Biggio, Tommaso Bendinelli, Alexander Neitz, Aurelien Lucchi, Giambattista Parascandolo  
Proceedings of the 38th International Conference on Machine Learning, PMLR 139:936-945, 2021.

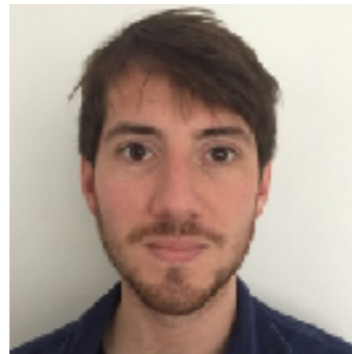
This is a **difficult problem**: *Symbolic Regression is NP-hard*, M. Virgolin, S. Pissis, 2022



**First Research Paper:**  
**Explainable Fraud Detection via  
Deep Symbolic Classification**  
(XAI 2024)



Samantha Visbeek  
(Zanders)

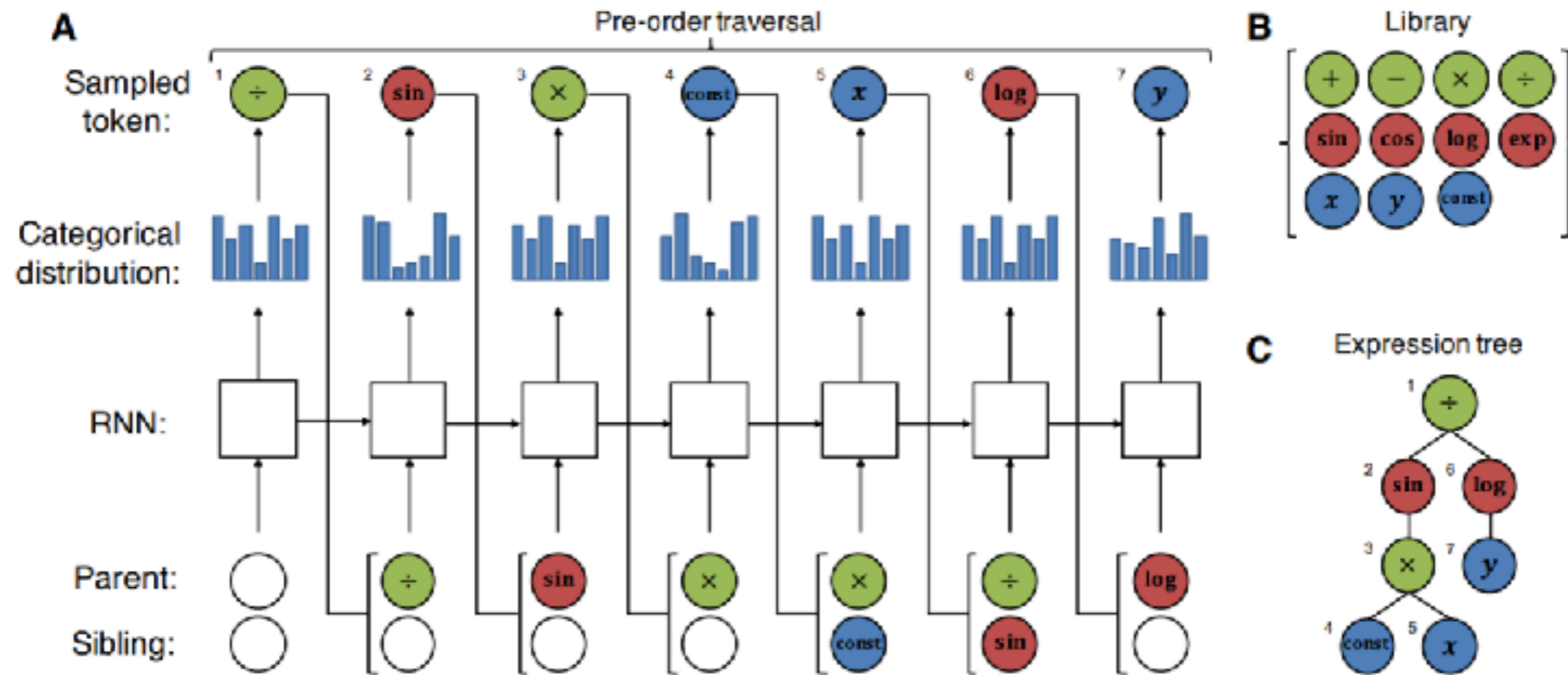


Floris den Hengst (ING & VU)



E. A.

# DEEP SYMBOLIC REGRESSION



- Elements are sampled from a categorical distribution emitted by the RNN (A), on the library  $L$  of elements in (B).
- The parent and sibling nodes of the next element are the next input to the RNN.
- The sampling process ends when all branches reach the leaf nodes.
- The resulting list  $[\div, \sin, \times, \text{constant}, x, \log, y]$  is the preorder traversal of the syntax tree that represents  $\sin(cx)/\log(y)$  (C).
- The syntax tree (C) that can be reconstructed from the preorder traversal list from (A).



# DATASET

## (PAYSIM)

### Synthetic Financial Datasets For Fraud Detection

Synthetic datasets generated by the PaySim mobile money simulator

[Data Card](#) [Code \(258\)](#) [Discussion \(29\)](#) [Suggestions \(0\)](#)

#### About Dataset

##### Context

There is a lack of public available datasets on financial services and specially in the emerging mobile money transactions domain. Financial datasets are important to many researchers and in particular to us performing research in the domain of fraud detection. Part of the problem is the intrinsically private nature of financial transactions, that leads to no publicly available datasets.

We present a synthetic dataset generated using the simulator called PaySim as an approach to such a problem. PaySim uses aggregated data from the private dataset to generate a synthetic dataset that resembles the normal operation of transactions and injects malicious behaviour to later evaluate the performance of fraud detection methods.

##### Content

PaySim simulates mobile money transactions based on a sample of real transactions extracted from one month of financial logs from a mobile money service implemented in an African country. The original logs were provided by a multinational company, who is the provider of the mobile financial service which is currently running in more than 14 countries all around the world.

#### Usability

8.82

#### License

[CC BY-SA 4.0](#)

#### Expected update frequency

Not specified

#### Tags

Finance

Crime

- PaySim is a data set of simulated transactions based on proprietary real transactions
- developed to provide researchers that exhibits statistical properties similar to a real payment transaction data set
- Contains ~ 6.3 million transactions over a period of a month
- Highly Imbalanced: with a fraudulent transaction rate of ~ 0.13%.

<https://www.kaggle.com/datasets/ealaxi/paysim1/data>



# RESULTS

method	accuracy	precision	recall	F1 score
RF + RUS	0.93	0.02	0.93	0.03
XGBoost + RUS	0.95	0.02	<b>0.94</b>	0.05
$k$ -NN + RUS	0.94	0.02	0.83	0.03
SVM + RUS	0.95	0.02	0.70	0.03
RF	<b>0.99</b>	<b>0.99</b>	0.67	0.81
XGBoost	<b>0.99</b>	0.98	0.70	<b>0.82</b>
DSC (average)	<b>0.99</b>	0.95 (.01)	0.67	0.78
DSC (best expression)	<b>0.99</b>	0.95	0.67	0.78

- Test set averaged over 5 runs: the baseline classification models with and without Random Undersampling (RUS)
- Split: 75% for training, 15% for validation, 10% for test.
- The reward function  $r_{F1}$  and a threshold value of 0.8 and the best expression obtained by DSC
- Undersampling (RUS) yields a loss in excessive amount of information and causes overfitting.
- $k$ -NN and SVM trained only with RUS. Otherwise, timeout after 50 hours due to high-class imbalance.



# PARETO FRONT OF PREDICTION VS. COMPLEXITY (OCCAM'S RAZOR)

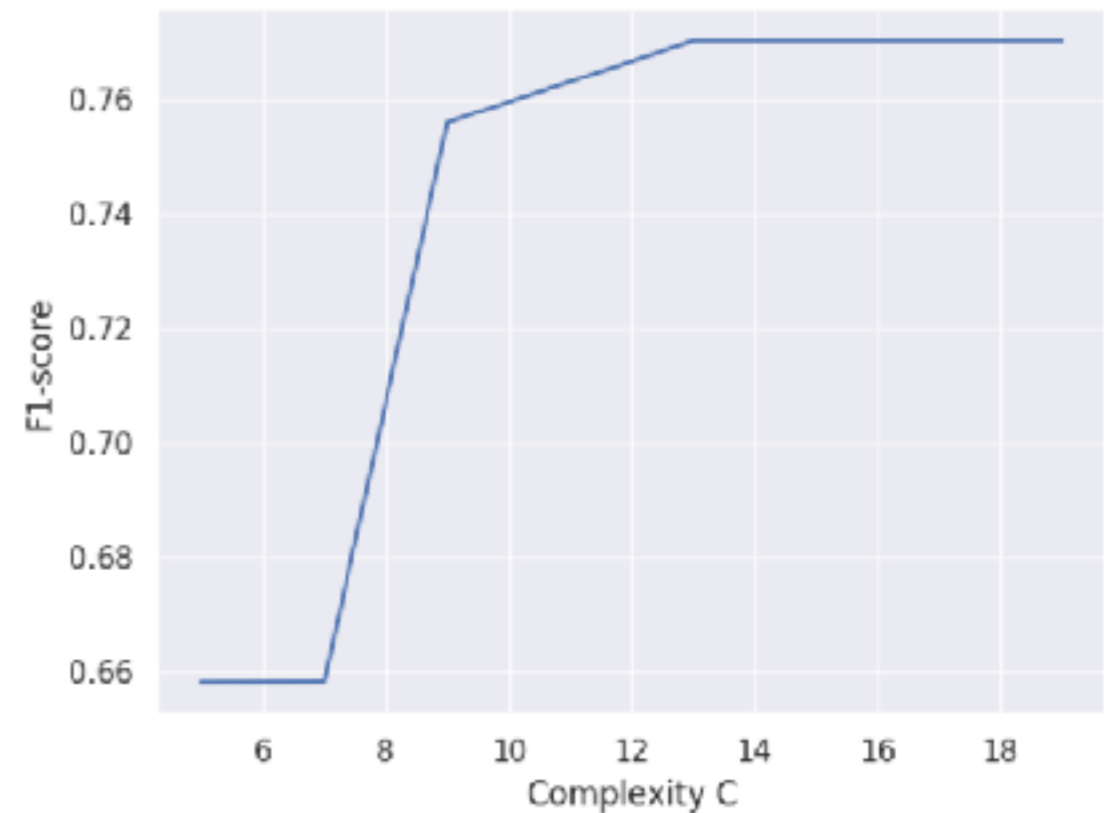
## Complexity of Tokens

token $\tau$	complexity $c$
+, -, $\times$ , feature, constant	1
$\div$ , square	2
sin, cos	3
exp, log, square-root	4

## Complexity of an Expression

$$C(f) = \sum_{i=0}^T c(\tau_i)$$

## Pareto Front of Performance vs. Complexity



$$t = 0. r_{F1} = \frac{2pr}{p+r}$$

The optimal expression selected w.r.t. adding more complexity to the expression does not result in a sufficient increase in the F1 score.

- prevents overfitting
- ensures that the expressions are not overly complex



# UNDERSTANDING THE BEST EXPRESSION

$$f = \sqrt{\text{externalDest} + \text{type\_cash-out} \cdot (\text{amount} - \text{maxDest7} + \text{type\_transfer})}$$

Complexity: 13

**A**

**B**

Boolean features: externalDest, type\_cash-out, type\_transfer

Non-negative numerical Features: amount, maxDest7

**Decision Rule:**  $\hat{y} = 1(\text{fraud})$ , if  $\sigma(f) > 0.7$ , (since  $t = 0.7$ )

$$\hat{y} = 1(\text{fraud}), \text{ if } f > 0.85. \quad \text{rewriting w.r.t. } \sigma(f) = 1/(1 + e^{-f})$$

Note that, (1)  $\text{amount} - \text{maxDest7} \leq 0$ , since maxDest7 includes current transaction

(2) type\_cash-out and type\_transfer are mutually exclusive (hence one-hot encoding)

**Two Scenarios: type\_cash-out = 1:**

Then, type\_transfer = 0, hence **B**  $\leq 0$ . But since, externalDest  $\in \{0,1\}$ ,  $f \leq 0$ , means no fraud.

**type transfer\_ = 1**

Then, type\_cash-out = 0, if externalDest = 0, **A** = 0, and no fraud.

if externalDest = 1, fraud = 1 iff  $f > 0.85$

iff **B**  $> 0.85$

iff  $\text{amount} - \text{maxDest7} > -0.15$

**Derived Decision Rule:** fraud if type = transfer  $\wedge$  externalDest = true  $\wedge$  amount - maxDest7  $> -0.15$ , legitimate otherwise.

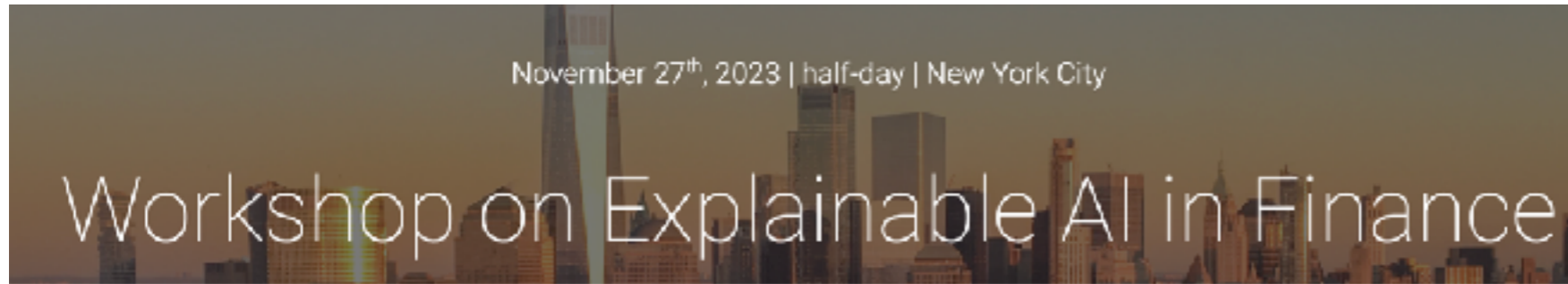
Complexity of Tokens

token $\tau$	complexity $c$
+, -, $\times$ , feature, constant	1
$\div$ , square	2
sin, cos	3
exp, log, square-root	4

Complexity of an Expression

$$C(f) = \sum_{i=0}^T c(\tau_i)$$

# IT WAS A MASTER THESIS



(Accepted and presented)



Samantha Visbeek

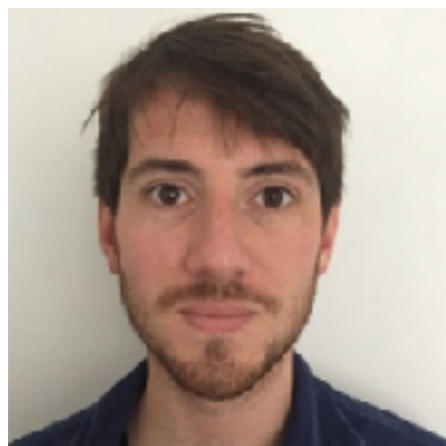
## amsterdam ai

### Call for Nominations is Open – Amsterdam AI Thesis Awards

The Amsterdam AI Thesis Awards aim to promote excellence in AI and Data Science from students at the Bachelor and Master level in Amsterdam-based Amsterdam-AI university partners (HvA, UvA, and VU).



The thesis won the Amsterdam AI Thesis Awards 2023



Co-supervised by Floris den Hengst (ING)



Published in XAI World conference 2024



# Second Research Paper:

## ECSEL: Explainable Classification via Signomial Equation Learning (ICML 2026 - to appear)



**Adia Lumadjeng**

PhD on Explainable AI for Fraud Detection  
Business Analytics (ABS)  
Socially Intelligent Artificial Systems Group (IVI)



**Ilker Birbil.**

Professor of AI & Optimisation Techniques for  
Business & Society  
Business Analytics (ABS)



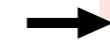
**E.A.**

# A Closer Look at AI Feynman data set



Richard Feynman (1918-1988)

Feynman eq.	Equation	Solution time (s)
I.6.20a	$f = e^{-\theta^2/2} / \sqrt{2\pi}$	16
I.6.20	$f = e^{-\frac{\theta^2}{2\sigma^2}} / \sqrt{2\pi\sigma^2}$	2992
I.6.20b	$f = e^{-\frac{(\theta-\theta_1)^2}{2\sigma^2}} / \sqrt{2\pi\sigma^2}$	4792
I.8.14	$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$	544
I.9.18	$F = \frac{Gm_1m_2}{(x_2-x_1)^2+(y_2-y_1)^2+(z_2-z_1)^2}$	5975
I.10.7	$m = \frac{m_0}{\sqrt{1-\frac{v^2}{c^2}}}$	14
I.11.19	$A = x_1y_1 + x_2y_2 + x_3y_3$	184
I.12.1	$F = \mu N_n$	12
I.12.2	$F = \frac{q_1q_2}{4\pi\epsilon r^2}$	17
I.12.4	$E_f = \frac{q_1}{4\pi\epsilon r^2}$	12
I.12.5	$F = q_2E_f$	8
I.12.11	$F = q(E_f + Bv \sin \theta)$	19
I.13.4	$K = \frac{1}{2}m(v^2 + u^2 + w^2)$	22
I.13.12	$U = Gm_1m_2(\frac{1}{r_2} - \frac{1}{r_1})$	20
I.14.3	$U = mgz$	12
I.14.4	$U = \frac{k_{spring}x^2}{2}$	9
I.15.3x	$x_1 = \frac{x^2-ut}{\sqrt{1-u^2/c^2}}$	22



I.12.2:  $F = \frac{1}{4\pi} \cdot q_1 \cdot q_2 \cdot \epsilon^{-1} \cdot r^{-2}$

I.13.4:  $K = \frac{1}{2}mv^2 + \frac{1}{2}mu^2 + \frac{1}{2}mw^2$



$$z(x) = \sum_{k=1}^K \alpha_k \prod_{j=1}^m x_j^{\beta_{k,j}}$$

*signomial function*

Udrescu, S.-M. And Tegmark, M. (2020): 'AI Feynman: A physics-inspired method for symbolic regression'. *Science Advances*, 6:eaay2631

45% are signomial functions!



# Not only Physics: Signomials are everywhere!

Field	Famous model	Signomial form
Economics	Cobb–Douglas production / utility	$Y = AK^\alpha L^\beta$ . This is a single monomial, so it is also a signomial. Used for output as a function of capital and labor; generalized versions use $\prod_i x_i^{\lambda_i}$ . <a href="#">Wikipedia</a>
Economics / optimization	Generalized cost or production functions	Sums such as $c_1 K^\alpha L^\beta + c_2 E^\delta + c_3 M^\epsilon$ . These are classic signomial-programming objects when coefficients/exponents are unrestricted. <a href="#">P/MC 41</a>
Biology	Kleiber's law / allometric scaling	$B = aM^{3/4}$ . A monomial relating metabolic rate to body mass; the exponent is debated, but the power-law form is canonical. <a href="#">P/MC 41</a>

**Theorem 3.1** (Universal Approximation for Signomials). *Let  $D \subset \mathbb{R}_{>0}^n$  be a compact subset of the positive orthant. Then the set of signomials*

$$\mathcal{S} = \left\{ S : S(x_1, \dots, x_n) = \sum_{k=1}^K \alpha_k \prod_{j=1}^n x_j^{\beta_{kj}}, \quad K \in \mathbb{N}, \alpha_k \in \mathbb{R}, \beta_{kj} \in \mathbb{R} \right\} \quad (9)$$

*is dense in  $C(D, \mathbb{R})$ . That is, for any continuous function  $f : D \rightarrow \mathbb{R}$  and any  $\epsilon > 0$ , there exists a signomial  $S \in \mathcal{S}$  such that*

$$\sup_{(x_1, \dots, x_n) \in D} |f(x_1, \dots, x_n) - S(x_1, \dots, x_n)| < \epsilon. \quad (10)$$

Physics	Newtonian gravity / inverse-square laws	$F = Gm_1 m_2 r^{-2}$ . A monomial signomial with a negative exponent. <a href="#">Wikipedia</a>
Physics / chemistry	Power-law rate laws	$r = kA^\alpha B^\beta$ . Monomial signomial; common in reaction kinetics and transport modeling.
Engineering	Drag or scaling laws	$F_D = cv^2, P = cv^3$ , or multi-term approximations such as $av^2 + bv + c$ . Polynomial-like models are signomials because integer powers are allowed.

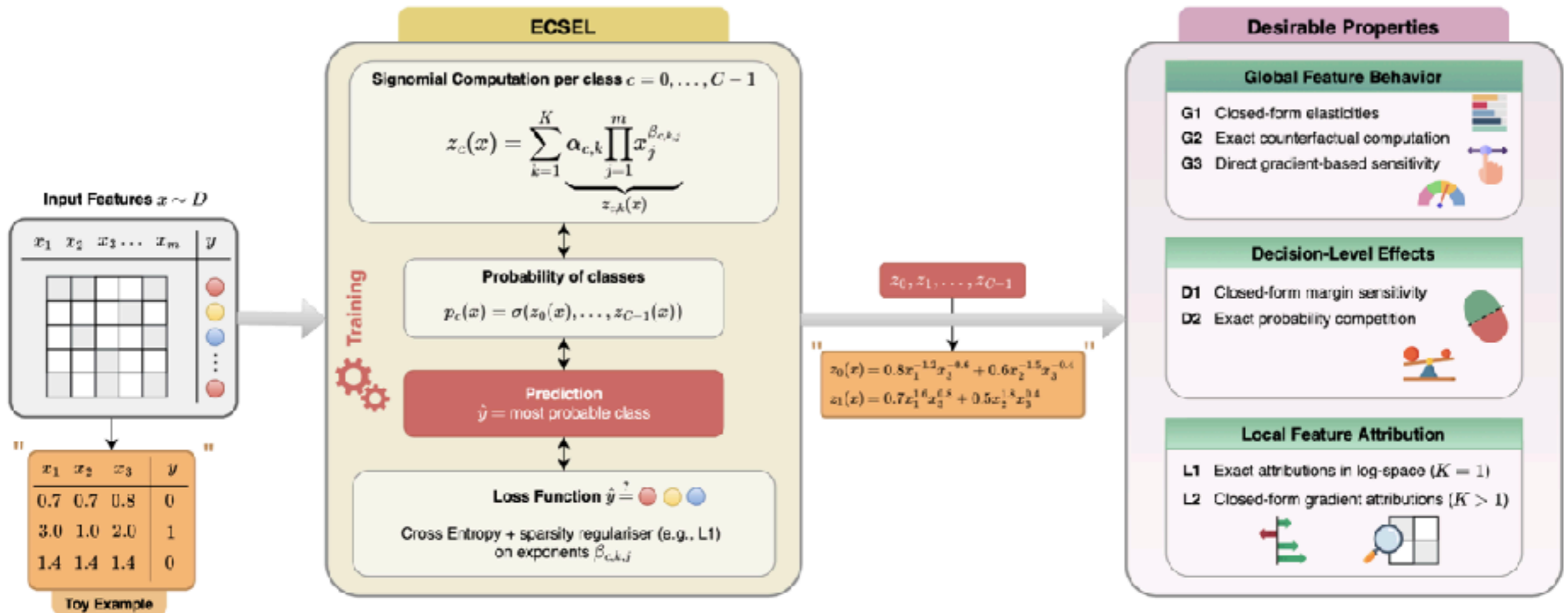
**Expressivity Result:**

**For a good reason: they are easy to compact, easy to interpret and expressive!**



# ECSEL

(Explainable classification via Signomial Equation Learning)



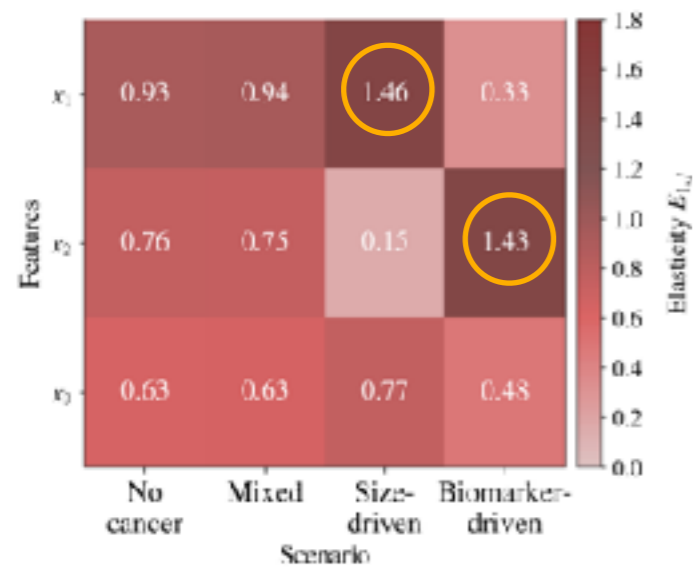
## Online Shopping Intention Example

$$z = 0.10 \cdot \frac{\text{PageValues}^{0.47} \cdot \text{Month}^{0.07} \cdot \text{PageValue per ExitRate}^{1.09} \cdot \text{ShopIntensity}^{0.66}}{\text{ExitRates}^{0.41} \cdot \text{Administrative}^{0.14} \cdot \text{IsReturn}^{0.04}}$$

REASONS I BUY THINGS:



# Interpretability Properties: Example



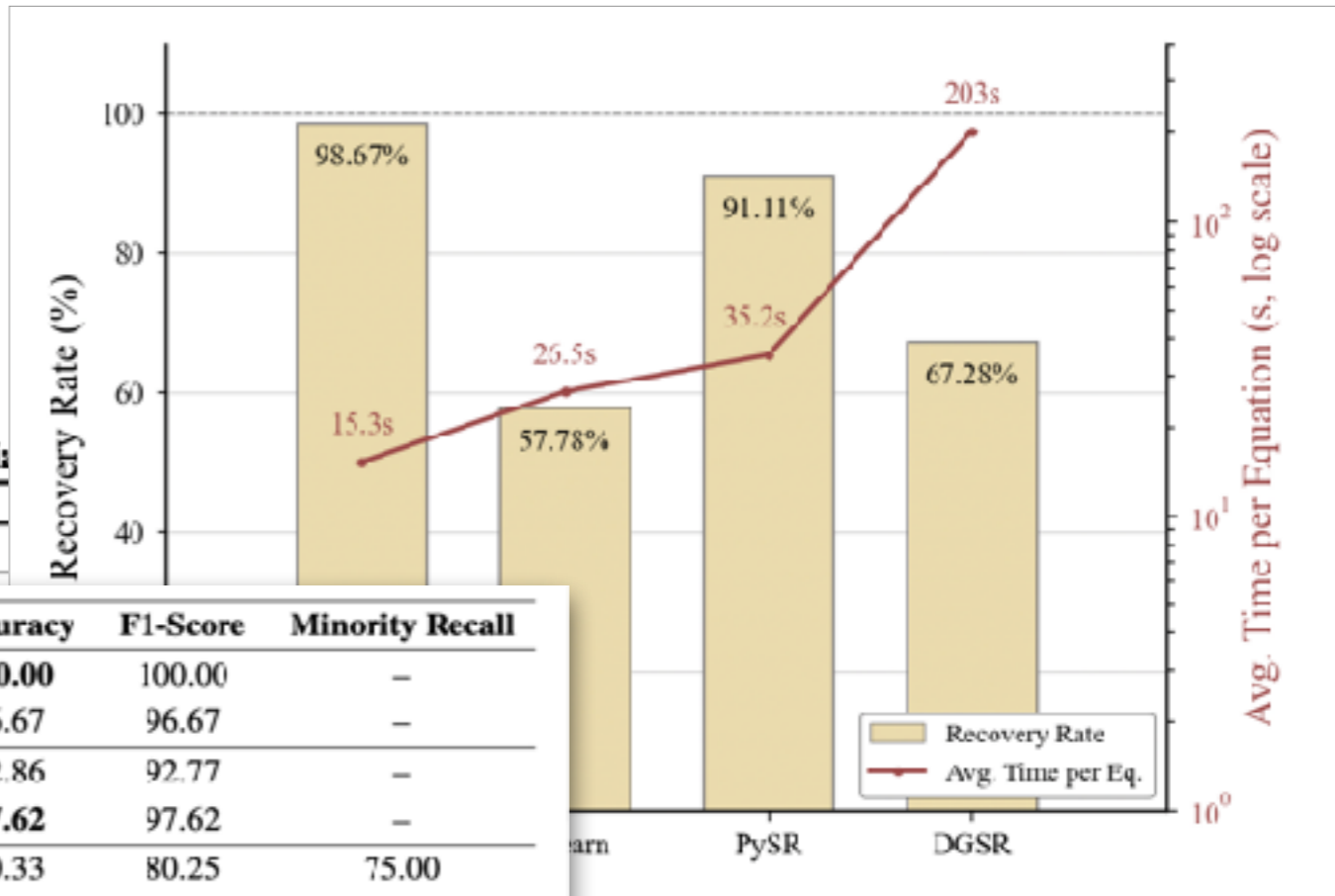
(a) Feature importance

# Symbolic Regression

## ECSEL vs. Generative Models

ECSEL: Explainable Classification via Signomix

ECSEL



Eq.	Expression
I.25.13	$V_c = \frac{2}{\omega C}$
I.29.4	$k = \frac{1}{\omega^2 a^2}$
I.32.5	$P = \frac{6\pi r^3}{q\omega B^3}$
I.34.8	$\omega = \frac{p}{E}$
I.34.27	$E = h\omega$
I.38.12	$r = \frac{4\pi\epsilon h^2}{3mq^2}$
I.39.10	$E = \frac{1}{2}pPV$
I.39.22	$P_F = \frac{nk_b T}{V}$
I.43.16	$w = \frac{\mu d r i f t q V_e}{\omega \epsilon k_b T}$
I.43.31	$D = \frac{d}{\omega \epsilon k_b T}$
I.47.23	$c = \sqrt{\frac{\rho}{P}}$
II.3.24	$F_B = \frac{4\pi r^2}{q}$
II.4.23	$V_e = \frac{4\pi\epsilon r^2}{3q}$
II.8.7	$E = \frac{5}{4\pi\epsilon r^2}$
II.8.31	$E_{dem} = \epsilon E_f^2$
II.11.20	$P_x = \frac{n_p P_d^2 E_f}{3k_b T}$
II.13.17	$B = \frac{4\pi\epsilon r^2}{2I}$
II.27.16	$F_B = \epsilon c E_f^2$
II.27.18	$E_{dem} = \epsilon E_f^2$
II.34.2A	$I = \frac{qv}{2\pi\epsilon_0 r}$
II.34.2	$\mu_M = \frac{2m}{9qB}$
II.34.11	$\omega = \frac{2m}{9qB}$

Dataset	Method	Accuracy	F1-Score	Minority Recall
IRIS	SVM	100.00	100.00	-
	ECSEL	96.67	96.67	-
SEEDS	LR	92.86	92.77	-
	ECSEL	97.62	97.62	-
HEARTS	LR	80.33	80.25	75.00
	ECSEL	83.61	83.61	82.14
ILPD	XGBoost	72.41	63.03	6.06
	ECSEL	75.86	74.39	42.42
TRANSFUSION	XGBoost	80.06	78.72	38.89
	ECSEL	79.33	77.95	41.67
CONTRACEPTIVE	XGBoost	60.00	59.14	-
	ECSEL	56.27	55.94	-
COMPAS	XGBoost	68.18	68.08	62.54
	ECSEL	68.47	68.36	62.82
DEFAULT	SVM	81.82	79.35	33.61
	ECSEL	81.74	79.34	34.06
SKINNONSKIN	XGBoost	99.96	99.96	99.97
	ECSEL	99.25	99.25	99.88
MAMMOGRAPHY	XGBoost	98.70	98.61	59.62
	ECSEL	98.66	98.57	59.62
LOAN	XGBoost	99.51	99.50	95.55
	ECSEL	99.22	99.21	92.97

100	0.55 ± 0.48	100	145.80 ± 32.03	100	25.51 ± 33.91
100	0.52 ± 0.45	100	121.01 ± 97.13	100	4.04 ± 16.88

-	-
0	250.22 ± 2.05
-	-
100	25.99 ± 1.29
0	260.26 ± 14.51
0	255.24 ± 1.02
0	262.06 ± 4.22
0	253.97 ± 5.09
100	23.10 ± 0.14
-	-
-	-
100	23.24 ± 0.81
100	0.60 ± 0.02
0	274.73 ± 6.11

# Back to the Fraud Detection (PaySim)

## Output

### Performance

Method	F1 (%)	Recall (%)	Precision (%)	ROC-AUC	Threshold	Time (s)
Logistic Regression	55.44	51.06	60.66	0.9593	1.000	1642.8
Random Forest	88.50	84.01	93.50	0.9992	0.417	204.0
XGBoost	<b>89.90</b>	<b>87.82</b>	92.09	<b>0.9998</b>	0.989	<b>8.5</b>
ECSEL	79.08	68.10	<b>94.27</b>	0.9914	0.904	960.0

### Equation ( $p > 0.904$ )

$$z = -0.07 \cdot \frac{A^{0.02} \cdot p^{0.03}}{exO^{0.03} \cdot exD^{0.16} \cdot CO^{0.14} \cdot T^{0.05} D^{0.03}} + 0.09 \cdot \frac{OBO^{1.42}}{NBO^{0.04} \cdot exD^{0.07} \cdot CO^{0.06} \cdot D^{0.06} p^{0.05}}$$

Feature name	Feature description	Type	Scaled	Notes
<b>Original numerical features</b>				
<i>step</i>	Discrete time step of the transaction (1 hour per step)	Continuous	Yes	Transaction timestamp
<i>amount</i>	Transaction amount	Continuous	Yes	Log-transformed
<i>oldbalanceOrig</i>	Origin account balance before transaction	Continuous	Yes	Log-transformed
<i>newbalanceOrig</i>	Origin account balance after transaction	Continuous	Yes	Log-transformed
<i>oldbalanceDest</i>	Destination account balance before transaction	Continuous	Yes	Log-transformed
<i>newbalanceDest</i>	Destination account balance after transaction	Continuous	Yes	Log-transformed
<b>Transaction type indicators (one-hot encoded)</b>				
<i>CashIn (CI)</i>	Cash-in transaction indicator	Binary	No	One-hot encoded
<i>CashOut (CO)</i>	Cash-out transaction indicator	Binary	No	One-hot encoded
<i>Debit (D)</i>	Debit transaction indicator	Binary	No	One-hot encoded
<i>Payment (P)</i>	Payment transaction indicator	Binary	No	One-hot encoded
<i>Transfer (T)</i>	Transfer transaction indicator	Binary	No	One-hot encoded
<b>Engineered features</b>				
<i>pct balance taken (PctBT)</i>	Fraction of the origin balance transferred in the transaction	Continuous	Yes	Clipped to [0, 1]
<i>externalOrig</i>	Indicator that the origin account is external (both balances equal zero)	Binary	No	Transaction-level rule

# FINAL TAKEAWAYS

- Symbolic Regressions opens domain-specific explainable AI expertise to explain data.
- Not only interesting for science, but also public sector, and many businesses both in regression and classification.
- So instead of a black box, you get a human-readable formula that explains exactly what the model is doing, by construction.
- Once the equation is obtained, easy to verify, on par with industry choice methods in accuracy.
- It is objective; so no human bias, but you can impose still the domain knowledge (part of future research).
- No post-hoc approximations, no proxy explanations, , no under or over sampling, the model *is* the explanation.
- The desirable properties imply more consistent and faster computation compared to post-hoc methods.
- The secret sauce: signomials, expressive enough to capture complex patterns, simple enough to read.
- More application on real-world data



# Thanks!

Thank you! Link to paper

