# DISCOVERING INTERPRETABLE PHYSICAL MODELS USING SR AND DEC

Simone Manti, Alessandro Lucantonio
Department of Mechanical and Production Engineering, Aarhus University

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT

# INTRODUCTION

**Problem:** use ML to discover new Physics

**Why**? Still unable to accurately model many physical phenomena (e.g. biological systems)

**How**? Many directions:

1.  Black-box methods (e.g. Neural Networks).

    Pros: many developed and tested models, more solid theory.

    Cons: difficult to interpret, requires large datasets.

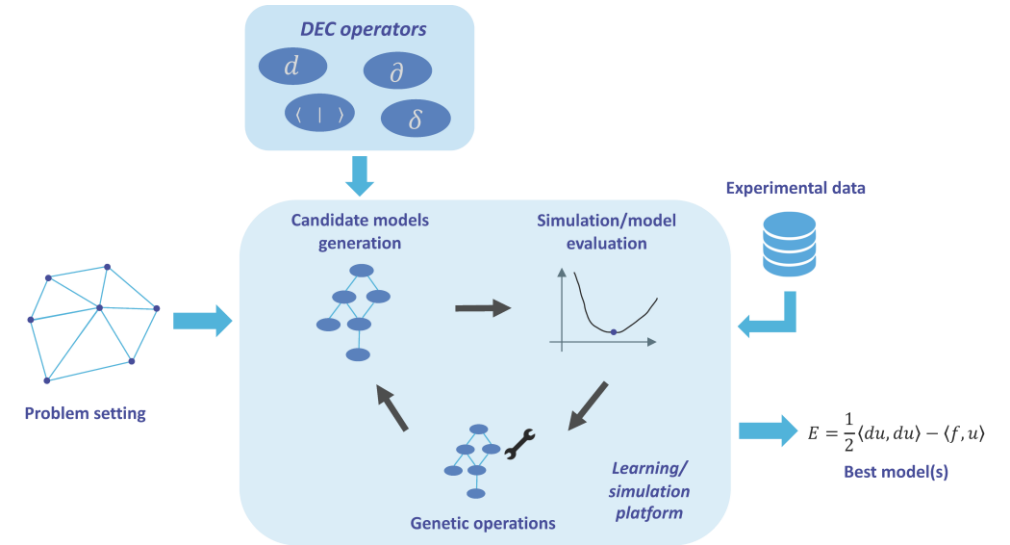2. Symbolic methods (e.g. Symbolic Regression models).

    Pros: interpretability + less data required (more constrained)

    Cons: algebraic equations (majority)

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT

# OUR CONTRIBUTION

**This work**: Develop a new symbolic method combining Symbolic Regression and Discrete Exterior Calculus to discover new physical models starting from data.

- General-purpose method designed for field problems

- The output of the method is an equation -> easy to interpret

- Working to extend it to face real-world and open problems

- Two open-source libraries: *dctkit* (to manage discrete mathematical tools) and *alpine* (to implement the learning strategy)

*dctkit*

*alpine*

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT

# STATE-OF-THE-ART COMPARISON

| | | Ours | PySINDy [1] | EQL [2] | Eureqa [3] | DSR [4] | AI Feynman [5] | PySR [6] |
|---|---|---|---|---|---|---|---|---|
| Field Problems | | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| | Domain source | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| | Stationary/ Non-stationary | ✓/✓ | ✗/✓ | ✗/✗ | ✗/✗ | ✗/✗ | ✗/✗ | ✗/✗ |
| | Variational/Non-variational | ✓/✓ | ✗/✓ | ✗/✗ | ✗/✗ | ✗/✗ | ✗/✗ | ✗/✗ |
| Dynamical Systems | | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| Algebraic Equations | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

[1] Kaptanoglu A A et al.,PySINDy: A comprehensive Python package for robust sparse system identification, *Journal of Open Source Software*

[2] Sahoo S, Lampert C and Martius G, Learning equations for extrapolation and control, *Int. Conf. on Machine Learning*

[3] Schmidt M and Lipson H, Distilling free-form natural laws from experimental data, *Science*

[4] Petersen B K, Landajuela M, Mundhenk T N, Santiago C P, Kim S K and Kim J T, Deep symbolic regression: recovering mathematical expressions from data via risk-seeking policy gradients (arXiv:1912.04871)

[5] Udrescu S-M and Tegmark M, AI Feynman: a physics-inspired method for symbolic regression, *Sci. Adv*

[6] Cranmer M, Interpretable machine learning for science with PySR and SymbolicRegression.jl (arXiv:2305.01582)

# DISCRETE EXTERIOR CALCULUS

- **Why?**

  - Discrete theory: no need for discretization schemes

  - Discrete geometric representation: suitable for field problems

  - Concise + effective set of operators: reduced search space

- **What?**

  - Discrete version of Exterior Calculus (differential forms in a manifold)

  - manifold <-> simplicial complex

  - field <-> form <-> cochain

  - grad, div, lap <-> coboundary (d) and hodge star ($\star$)



Figure from Desbrun et al., *Discrete Differential forms for Computational Modeling*



Figure from K.Crane, *Discrete Differential Geometry: An applied introduction*
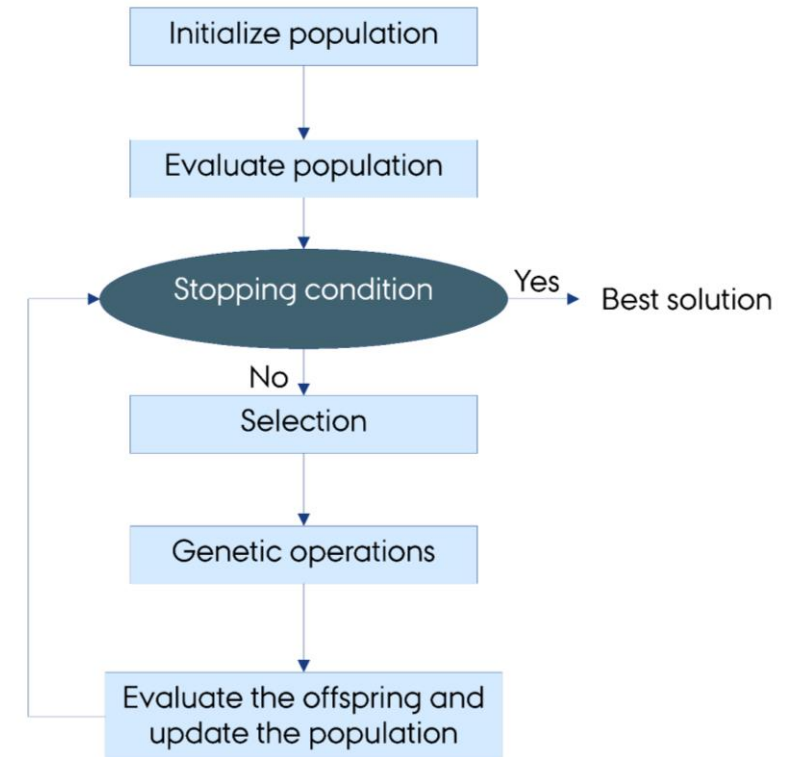
# SYMBOLIC REGRESSION

- **Why?**

  - Generate and manipulate candidate equations describing a given physical system;

  - interpretability;

  - use small dataset for training, validation and test.

- **What?**

  - Symbolic Regression -> find equations given data;

  - Genetic Programming -> evolutionary strategy that explores the space updating an initial population through genetic operations. The goal is to maximise a proper fitness function.



AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT

# THE METHOD

1.  Individual: potential energy or residual

2.  Pure GP does not work -> we cannot sum e.g. a cochain with a scalar.
    Sol: Strongly-Typed Genetic Programming -> type consistent trees

3.  Dimensionless variables -> every generated expression is physically meaningful

4.  Each individual is minimized (for the residual-> its norm) according to initial and
    boundary conditions. Then, we compare the solution with the true data, maximizing

$$F(I) = -(\alpha \mathrm{MSE}(I) + \eta R(I))$$

fixed scaling factor (problem dependent)

Regularization hyperparameter

Regularization function

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024    |    SIMONE MANTI
21 AUGUST 2024    |    PHD STUDENT

# RESULTS - PRELIMINARY INFO

- 3 different benchmarks in variational form: *Poisson, Elastica, Linear Elasticity* equations

- Data are split in training, validation and test (double hold-out).

- 50 final model discovery runs to compute recovery rate or MSE mean ± std



Training fitness

Best individual so far

A look inside *alpine*

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT

# RESULTS – 2D POISSON

Dirichlet energy using DEC operators:

$$\mathcal{E}_{\mathrm{p}}(u) := \frac{1}{2}\langle du, du \rangle - \langle u, f \rangle = \frac{1}{2}\langle \delta du, u \rangle - \langle u, f \rangle.$$

- Task: learn the correct energy
- Full dataset (training + validation + test): 12 (u,f) pairs without noise
- Dirichlet BCs
- Recovery rate: 66% (best hyperparameters)



True solution    Model solution

| # | Energy | Training Fit. | Test Fit. | Test MSE |
|---|--------|---------------|-----------|----------|
| 1 | $\langle u, \delta(d(u/2) - f)\rangle$ | 0.9 | 0.9 | $9.8 \cdot 10^{-10}$ |
| 2 | $\langle u, \delta(du) - 2f\rangle$ | 0.9 | 0.9 | $9.8 \cdot 10^{-10}$ |
| 3 | $\langle \star(\star u), \delta(du) - 2f\rangle$ | 1.1 | 1.1 | $9.8 \cdot 10^{-10}$ |
| 4 | $\langle du, du\rangle - \langle f, 2u\rangle$ | 1.1 | 1.1 | $9.8 \cdot 10^{-10}$ |

# RESULTS – EULER'S ELASTICA

$$\mathcal{E}_{\text{el}}(u) := \frac{1}{2}\langle \mathbb{1}_{\text{int}} \odot \star d^{\star}u, \mathbb{1}_{\text{int}} \odot \star d^{\star}u \rangle - \langle f\mathbb{1}, \sin u \rangle$$

- Task: learn the correct energy and the best related B

- f = PL$^2$/B, where P is the vertical component of the load and B is the *bending stiffness* of the rod

- Full dataset: 10 pairs (u, PL$^2$) perturbed with uniform noise

- To autotune the constant B we solve at each time

$$\min_{f \geq 0} \quad ||u_f - \bar{u}||^2 \qquad \text{s.t. } u_f \in \underset{u\,:\,u(0)=0}{\arg\min}\, \mathcal{E}(u, f)$$



| # | Energy | Training Fit. | Test Fit. | Test MSE | $B$ |
|---|--------|---------------|-----------|----------|-----|
| 1 | $\langle \star\mathbb{1}_{\text{int}}, (d^{\star}u)^2 \rangle - \langle \sin u, f\mathbb{1} \rangle$ | 0.196 | 0.1939 | 0.0084 | 37.0312 Nm$^2$ |
| 2 | $\langle \arccos(-1)\sin u + \delta^{\star}(\sin(\sin(d^{\star}u))), u - f\mathbb{1} \rangle$ | 0.2123 | 0.2247 | 0.0075 | 7.1779 Nm$^2$ |
| 3 | $\langle u - f\mathbb{1}, \delta^{\star}(\sin(\sin(d^{\star}u))) + 1/\exp(-1)\sin u \rangle$ | 0.214 | 0.2168 | 0.0067 | 6.8262 Nm$^2$ |
| 4 | $\langle \star\mathbb{1}_{\text{int}}, (d^{\star}u)^2 \rangle - \langle f\mathbb{1}, \sin u \rangle + 1/2$ | 0.216 | 0.2139 | 0.0084 | 37.0312 Nm$^2$ |

AARHUS UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION ENGINEERING

HAMLET 2024
21 AUGUST 2024

SIMONE MANTI
PHD STUDENT

# RESULTS – 2D LINEAR ELASTICITY

$$\mathcal{E}_{LE}(\boldsymbol{F}) := \frac{1}{2}\langle \boldsymbol{F} + \boldsymbol{F}^T - 2\boldsymbol{I} + \frac{\lambda_-}{\mu_-}\operatorname{tr}\left(\frac{1}{2}(\boldsymbol{F} + \boldsymbol{F}^T) - \boldsymbol{I}\right)\boldsymbol{I}, \frac{1}{2}(\boldsymbol{F} + \boldsymbol{F}^T) - \boldsymbol{I}\rangle$$

- Task: learn the energy as a function of **F**

- Full dataset: 20 node coordinates in the deformed configuration (homogeneous deformations)

- Implemented an internal filter (during training) for frame indifference

- Recovery rate: 92% with the best hyperparameter set

- Some of the learned energies are equivalent to the correct one only if **F** is constant



| # | Energy | Training MSE | Test MSE |
|---|--------|--------------|----------|
| 1 | $\langle \boldsymbol{I} - \boldsymbol{F}^T, \boldsymbol{I}\rangle^2 - 0.1(\langle \boldsymbol{I} - \boldsymbol{F}^T, \operatorname{sym}(\boldsymbol{F}) - \boldsymbol{I}\rangle + \langle \boldsymbol{I} - \operatorname{sym}(\boldsymbol{F}), \operatorname{sym}(\boldsymbol{F}) - \boldsymbol{I}\rangle)$ | 0 | $1.1672 \cdot 10^{-16}$ |
| 2 | $\langle -0.5\boldsymbol{I} + 0.5\operatorname{sym}(\boldsymbol{F}), \operatorname{sym}(\boldsymbol{F}) - \boldsymbol{I} + (\operatorname{tr}(\boldsymbol{F}) - \operatorname{tr}(\boldsymbol{I}))(2\operatorname{tr}(\boldsymbol{I})\boldsymbol{I} + \boldsymbol{I})\rangle$ | 0 | $2.0785 \cdot 10^{-16}$ |
| 3 | $\langle \boldsymbol{I} - \boldsymbol{F}^T, \boldsymbol{I} - \operatorname{sym}(\boldsymbol{F}) + 5\langle \boldsymbol{I} - \boldsymbol{F}, \boldsymbol{I}\rangle\boldsymbol{I}\rangle$ | 0 | $3.7309 \cdot 10^{-16}$ |
| 4 | $\langle \operatorname{sym}(\boldsymbol{F}) - \boldsymbol{I}, \operatorname{sym}(\boldsymbol{F}) - \boldsymbol{I}\rangle + 5\langle \boldsymbol{I} - \boldsymbol{F}, \boldsymbol{I}\rangle^2$ | 0 | $8.7033 \cdot 10^{-16}$ |

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024          SIMONE MANTI
21 AUGUST 2024       PHD STUDENT

# Thanks for your attention!



*Scan for the full paper*

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT

# PRIMITIVES

| Primitive names | Input types | Return type | Problem | | |
|---|---|---|---|---|---|
| | | | *Poisson* | *Elastica* | *Linearelasticity* |
| `Add, Sub, MulF, Div` | `(float, float)` | `float` | ✓ | ✓ | ✓ |
| `SinF, ArcsinF, CosF, ArccosF, ExpF, LogF InvF` | `float` | `float` | ✓ | ✓ | ✗ |
| `SqrtF, SquareF` | `float` | `float` | ✓ | ✓ | ✓ |
| `dXJS, delXJS, SinXJS, ArcsinXJS, CosXJS, ArccosXJS, ExpXJS, LogXJS` | `CochainXJS` | `CochainXJS` | ✓ | ✓ | ✗ |
| `StJR, InvStJR` | `CochainXJR` | `CochainXJR` | ✓ | ✓ | ✓ |
| `SqrtXJR, SquareXJR` | `CochainXJR` | `CochainXJR` | ✓ | ✓ | ✗ |
| `tranXJT, symXJT` | `CochainXJT` | `CochainXJT` | ✗ | ✗ | ✓ |
| `trXJT` | `CochainXJT` | `CochainXJS` | ✗ | ✗ | ✓ |
| `MulXJR` | `(CochainXJR, float)` | `CochainXJR` | ✓ | ✓ | ✓ |
| `MulvXJ` | `(CochainXJS, CochainXJT)` | `CochainXJT` | ✗ | ✗ | ✓ |
| `InvMulXJS` | `(CochainXJS, float)` | `CochainXJS` | ✓ | ✓ | ✗ |
| `InnXJR` | `(CochainXJR, CochainXJR)` | `float` | ✓ | ✓ | ✓ |
| `AddCXJR, SubCXJR` | `(CochainXJR, CochainXJR)` | `CochainXJR` | ✓ | ✓ | ✓ |
| `CochMulXJS` | `(CochainXJS, CochainXJS)` | `CochainXJS` | ✓ | ✓ | ✗ |

# BEST HYPERPARAMETERS

| Hyperparameter | Value | | |
|---|---|---|---|
| | *Poisson* | *Elastica* | *Linear elasticity* |
| Number of individuals ($\mu$) | 2000 | 2000 | 2000 |
| Crossover/mutation probabilities | (0.2, 0.8) | (0, 1) | (0.2,0.8) |
| Mixed mutation probabilities | (0.8, 0.2, 0) | (0.8, 0.2, 0) | (0.7, 0.2, 0.1) |
| Stochastic tournament probabilities | (0.7, 0.3) | (1, 0) | (0.7,0.3) |
| Regularization factor ($\eta$) | 0.1 | 0.01 | 0 |

# 2D LINEAR ELASTICITY – STANDARD GP



4% recovery rate

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT

# 2D LINEAR ELASTICITY – NON HOMOG.

AARHUS
UNIVERSITY
DEPARTMENT OF MECHANICAL AND PRODUCTION
ENGINEERING

HAMLET 2024 | SIMONE MANTI
21 AUGUST 2024 | PHD STUDENT