

Multiview Symbolic Regression

How to learn laws from examples

<u>Etienne Russeil</u> - *LPC Université Clermont Auvergne, France* **Fabricio Olivetti** - *CMCC Federal University of ABC, Brazil* **Konstantin Malanchev** - *University of Illinois Urbana–Champaign, USA* **Emille Ishida** - *LPC Université Clermont Auvergne, France* **Emmanuel Gangler** - *LPC Université Clermont Auvergne, France*

PLAN

- Motivation
- Symbolic Regression
- Multiview Symbolic regression
- Scientific applications
- Conclusion

Laws in science

The importance of laws in science

Science's goal is to describes and predicts phenomena.

The ultimate form of that is a mathematical equation

Einstein matter energy equivalence law

$$E = mc^2$$

Maxwell's equations

 $\nabla \cdot E = 0 \quad \nabla \times E = -\frac{1}{c} \frac{\partial H}{\partial t}$

$$\nabla \cdot H = 0 \quad \nabla \times H = \frac{1}{c} \frac{\partial E}{\partial t}$$

Galileo Galilei

"Mathematics is the language with which God has written the universe."



Newton's gravitation law

$$F = G \frac{m_1 m_2}{d^2}$$

Pythagoras theorem

 $a^2 + b^2 = c^2$

Shannon's information theory

$$H = -\sum p(x)\log p(x)$$

Equations come from observations : Kepler's laws



Tycho Brahe

The planets orbit around the sun and I will prove it with measurements ! Nova MVN DANI STSTEMATIS HYPOTYPOSIS ab Authore nuper adinuenta, qua tum vetus illa Ptolemaica redundantia & inconcinnitas, tum etiam recens Coperniana in motu Terra Physica absurditas, excluduntur, omniag, Apparentiis Cælestibus aptissime correspondent.

Equations come from observations : Kepler's laws



Johannes Kepler

Thank you for the dataset





Tabula TYCHONIS BRAHE ob matarum & computatarum oppolitionum MARTIS cum linea men Imotus Solis, ejusque examen.

Igitur tabula, e i qua fupra, fuit ifta.

T	Lon	ong. obs.re-				ati	tuc	lo	Long. obs. re							
fpec.circuli@											ob	s.	Au ecliptic			
Anni	Menf.	D	H	M	G	M	S		G	M	S		G	M		
1580	Novemb.	17	9	40	6	50	10	п	I	40	0	B	6	:46		
1582	Decembr.	28	I 2.	16	16	51	30	59	4	6	0	B	16	45		
1585	Januarii	31	19	35	21	9	50	ณ	4	32	10	B	2. I	IO		
1587	Martii	7	17	22	25	5	10	np	3	38	12	B	25	10		
1589	Aprilis	15	13	34	3	54	35	m	I	6	45	B	3	58	1	
1591	Junii	8	16	25	26	40 41	30	PF	3	59	0	M	26	32		
1593	Augusti	24	2	13	I 2	35	0	х	6	3	0	M	I 2	43	-	
1595	Octobris	29	21	22	17	56	5	8	0	5	15	B	17	56	1	
1597	Decemb.	13	13	35	2.	34	0	20	3	33	0	B	2	28		
1600	Januarii	19	9	40	8	18	45	ຄ	4	30	50	B	8	18		

Planetæ er motus in fuo eccentrico e certis obfervationibusan anthis per annos xx(ab Lxxx usque M D c)fedulo per noftra instrumenta habitis respectu variarum dispositionus min subjecta tabula patet, accurata restitutio.

-																		
22		Differentia			Simpl.Long.~				Apog. o				Præcefs.æ- quin.noftra			Ρ̈́Ρ		
T	I	MS		S	G	M	S	S	G	M	S	G	M	S	G	M	S	
		410	A	0	27	29	46	.3	25	21	40	27	58	50	6	50	40	
		1520	A	2	II	34	56	3	25	22	17	28	0	38	16	5Ì	26	
		036	S	3	22	37	46	3	25	22	55	28	2	25	21	9	4 I	
		510	S	5	3	27	46	3	25	23	32	28	4	10	25	:4	50	
1		335	S	6	16	53	7	3	25	24	10	28	5	55	3	54	33	
1		1020	A	8	7	47	30	3	25	24	48	2.8	7	47	26	40	23	
I		845	S	10	IO	53	50	3	25	25	26	28	9	40	12	34	36	
		012	A	Ó	8	26	47	3	25	27	35	28	, I I	27	17	57	14	
1	1	60	A	I	24	.55	47	3	25	29	5	28	13	20	2	32	20	
-	1	045	S	3	6	46	16	3	25	130	6	2.8	IŚ	5	8	19	57	

Equations come from observations : Kepler's laws

1: The orbit of every planet is an ellipse with the Sun at one of the two foci

2: A line joining a planet and the Sun sweeps out equal areas during equal intervals of time $\frac{dA}{dt} = \frac{\pi \times a \times b}{T}$

3: The ratio of the square of an object's orbital period with the cube of the semi-major axis of its orbit is the same for all objects orbiting the same primary.

 $\frac{R^3}{T^2} = constant$

Johannes Kepler

What if we could automatise this discovery process ?

Symbolic Regression



Caution: I will present a standard method based on Genetic Programming but it is not the only one ! In practice each implementation varies







As a first step the algorithm will randomly generate many different equations













Many different implementations of SR available !









https://github.com/heal-research/operon

DATA SET



does not output a general law







I would have found f(X) = A X + B





MultiView Symbolic Regression











MvSR in a nutshell

- (1) Receive multiple datasets as input.
- (2) Perform a minimization of the parameters independently
- for each dataset.
- (3) Use an aggregation function to compute an overall loss.
- (4) Allow parameters to be repeated.
- (5) Control the maximum number of parameters.
- (6) Penalise solutions based on the number of parameters used





[Submitted on 6 Feb 2024] Multi-View Symbolic Regression

Etienne Russeil, Fabrício Olivetti de França, Konstantin Malanchev, Bogdan Burlacu, Emille E. O. Ishida, Marion Leroux, Clément Michelin, Guillaume Moinard, Emmanuel Gangler

Symbolic regression (SR) searches for analytical expressions representing the relationship between a set of explanatory and response variables. Current SR methods assume a single dataset extracted from a single experiment. Nevertheless, frequently, the researcher is confronted with multiple sets of results obtained from experiments conducted with different setups. Traditional SR methods may fail to find the underlying expression since the parameters of each experiment can be different. In this work we present Multi-View Symbolic Regression (MvSR), which takes into account multiple datasets simultaneously, mimicking experimental environments, and outputs a general parametric solution. This approach fits the evaluated expression to each independent dataset and returns a parametric family of functions f(x; \theta) simultaneously capable of accurately fitting all datasets. We demonstrate the effectiveness of MvSR using data generated from known expressions, as well as real-world data from astronomy, chemistry and economy, for which an a priori analytical expression is not available. Results show that MvSR obtains the correct expression more frequently and is robust to hyperparameters change. In real-world data, it is able to grasp the group behaviour, recovering known expressions from the literature as well as promising alternatives, thus enabling the use SR to a large range of experimental scenarios.

Comments: Submitted to GECCO-2024. 10 pages, 6 figures

Subjects: Machine Learning (cs.LG); Instrumentation and Methods for Astrophysics (astro-ph.IM); Applications (stat.AP)

Cite as: arXiv:2402.04298 [cs.LG]

(or arXiv:2402.04298v1 [cs.LG] for this version)

On arXiv since yesterday: https://arxiv.org/abs/2402.04298

Scientific applications



 ZTF DR OID or SNAD name
 633207400004730 or 202
 co

 Coordinates or name
 00h00m00s +00d00m00s or M57
 radius, arcsec
 1
 Coordinates

717101200000722



No minor planets found in 15"

https://ztf.snad.space/dr17/view/717101200000722








MvSR recovers the literature !



It can also generate more complexe and better solution



As well as unexpected but very effective forms

Let's try other use cases !





S&P500

Stock market of the 500 biggest american companies

S&P500

Stock market of the 500 biggest american companies

Usually aggregated for analysis



S&P500

NO randomies

Stock market of the 500 biggest american companies

490 to check

	Recover literature		
Models	Equation f(x)	med(MSE)	
Gaussian [2, 5]	$A \cdot e^{-\frac{x^2}{B}}$	0.363	
Laplace [17]	$A \cdot e^{-B x }$	0.342	
Cauchy [20]	$A \cdot B^2 / (x^2 + B^2)$	0.305	
Linear-Laplace	$(A - Bx) \cdot e^{-C x }$	0.327	
Exp-Laplace	$A \cdot e^{Bx - C x }$	0.328	
Power-Laplace	$A \cdot e^{B x ^C}$	0.246	

Find new models

	Recover literature		
Models	Equation f(x)	med(MSE)	
Gaussian [2, 5]	$A \cdot e^{-\frac{x^2}{B}}$	0.363	
Laplace [17]	$A \cdot e^{-B x }$	0.342	
Cauchy [20]	$A \cdot B^2 / (x^2 + B^2)$	0.305	
Linear-Laplace	$(A - Bx) \cdot e^{-C x }$	0.327	
Exp-Laplace	$A \cdot e^{Bx - C x }$	0.328	
Power-Laplace	$A \cdot e^{B x ^C}$	0.246	

Find new models



And another one?











General Beer Lambert's law:





Conclusion

Conclusion

- MvSR is working, have a look at the <u>arXiv</u>
- It has potential to be used in every science
- If you see a potential use case for you science come and talk to me
- Still need some work on our side for a proper full implementation



Multiview Symbolic Regression (MvSR)





Multiview Symbolic Regression (MvSR)

Toy data illustration



61





Point mutations





Point mutations



Hoist mutations



Point mutations

Х





sin

7





Create a new population from the previous best candidates



$$f_1(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$f_2(x) = \sin(\theta_0 x_0 x_1) + \theta_1 (x_2 - \theta_2)^2 + \theta_3 x_3 + x_4$$

$$f_3(x) = (\theta_0 x_0^2 + (\theta_1 x_1 x_2 - \frac{\theta_2}{(\theta_3 x_1 x_3 + 1)})^2)^{0.5}$$



	View 1	View 2	View 3	View 4	Partial view
θ_0	2	0	0	2	2
θ_1	2	2	0	0	-2
θ_2	0	2	2	0	2
θ_3	0	0	2	2	2

Parameter	Value
population size	1000
number of evaluations	10000000
pool size	5
error metric	MSE
prob. cx	1.0
prob. mut.	0.25
max depth.	10
optim. iterations	100
aggregation function	max
operators	add, sub, mul, div, square, exp, sqrt,
	$sin(f_2 only)$

Models	Equation f(x)	med(MSE)	MSE _{S&P}
Gaussian [2, 5]	$A \cdot e^{-\frac{x^2}{B}}$	0.363	0.260
Laplace [17]	$A \cdot e^{-B x }$	0.342	0.084
Cauchy [20]	$A \cdot B^2 / (x^2 + B^2)$	0.305	0.079
Linear-Laplace	$(A - Bx) \cdot e^{-C x }$	0.327	0.065
Exp-Laplace	$A \cdot e^{Bx - C x }$	0.328	0.063
Power-Laplace	$A \cdot e^{B x ^C}$	0.246	0.075
