

PySINDy: A comprehensive Python package for robust sparse system identification

Alan A. Kaptanoglu¹, Brian M. de Silva², Urban Fasel³, Kadierdan Kaheman³, Jared Callahan³, Charles B. Delahunt², Kathleen Champion², Jean-Christophe Loiseau⁴, J. Nathan Kutz², and Steven L. Brunton³

¹ Department of Physics, University of Washington ² Department of Applied Mathematics, University of Washington ³ Department of Mechanical Engineering, University of Washington ⁴ Arts et Métiers Institute of Technology, CNAM, DynFluid, HESAM Université

DOI: [10.21105/joss.03994](https://doi.org/10.21105/joss.03994)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Sebastian Benthall](#) ↗

Reviewers:

- [@henrykironde](#)
- [@tuelwer](#)

Submitted: 21 October 2021

Published: 14 December 2021

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Automated data-driven modeling, the process of directly discovering the governing equations of a system from data, is increasingly being used across the scientific community. PySINDy is a Python package that provides tools for applying the sparse identification of nonlinear dynamics (SINDy) approach to data-driven model discovery. In this major update to PySINDy, we implement several advanced features that enable the discovery of more general differential equations from noisy and limited data. The library of candidate terms is extended for the identification of actuated systems, partial differential equations (PDEs), and implicit differential equations. Robust formulations, including the integral form of SINDy and ensembling techniques, are also implemented to improve performance for real-world data. Finally, we provide a range of new optimization algorithms, including several sparse regression techniques and algorithms to enforce and promote inequality constraints and stability. Together, these updates enable entirely new SINDy model discovery capabilities that have not been reported in the literature, such as constrained PDE identification and ensembling with different sparse regression optimizers.

Statement of need

Traditionally, the governing laws and equations of nature have been derived from first principles and based on rigorous experimentation and expert intuition. In the modern era, cheap and efficient sensors have resulted in an unprecedented growth in the availability of measurement data, opening up the opportunity to perform automated model discovery using data-driven modeling. These data-driven approaches are also increasingly useful for processing and interpreting the information in these large datasets. A number of such approaches have been developed in recent years, including the dynamic mode decomposition ([Kutz et al., 2016](#); [Schmid, 2010](#)), Koopman theory ([Steven L. Brunton et al., 2021](#)), nonlinear autoregressive algorithms ([Billings, 2013](#)), neural networks ([Pathak et al., 2018](#); [M. Raissi et al., 2019](#); [Vlachas et al., 2018](#)), Gaussian process regression ([Maziar Raissi et al., 2017](#)), operator inference and reduced-order modeling ([Benner et al., 2015](#); [Peherstorfer & Willcox, 2016](#); [Qian et al., 2020](#)), genetic programming ([Bongard & Lipson, 2007](#); [Schmidt & Lipson, 2009](#)), and sparse regression ([Steven L. Brunton et al., 2016](#)). These approaches have seen many variants and improvements over the years, so data-driven modeling software must be regularly updated to remain useful to the scientific community. The SINDy approach has experienced

particularly rapid development, motivating this major update to aggregate these innovations into a single open-source tool that is transparent and easy to use for non-experts or scientists from other fields.

The original PySINDy code (de Silva et al., 2020) provided an implementation of the traditional SINDy method (Steven L. Brunton et al., 2016), which assumes that the dynamical evolution of a state variable $\mathbf{q}(t) \in \mathbb{R}^n$ follows an ODE described by a function \mathbf{f} ,

$$\frac{d}{dt}\mathbf{q} = \mathbf{f}(\mathbf{q}). \quad (1)$$

SINDy approximates the dynamical system \mathbf{f} in Eq. (1) as a sparse combination of terms from a library of candidate basis functions $\boldsymbol{\theta}(\mathbf{q}) = [\theta_1(\mathbf{q}), \theta_2(\mathbf{q}), \dots, \theta_p(\mathbf{q})]$

$$\mathbf{f}(\mathbf{q}) \approx \sum_{k=1}^p \theta_k(\mathbf{q}) \boldsymbol{\xi}_k, \quad \text{or equivalently} \quad \frac{d}{dt}\mathbf{q} \approx \boldsymbol{\Theta}(\mathbf{q}) \boldsymbol{\Xi}, \quad (2)$$

where $\boldsymbol{\Xi} = [\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p]$ contain the sparse coefficients. In order for this strategy to be successful, a reasonably accurate approximation of $\mathbf{f}(\mathbf{q})$ should exist as a sparse expansion in the span of $\boldsymbol{\theta}$. Therefore, background scientific knowledge about expected terms in $\mathbf{f}(\mathbf{q})$ can be used to choose the library $\boldsymbol{\theta}$. To pose SINDy as a regression problem, we assume we have a set of state measurements sampled at time steps t_1, \dots, t_m and rearrange the data into the data matrix $\mathbf{Q} \in \mathbb{R}^{m \times n}$,

$$\mathbf{Q} = \begin{bmatrix} q_1(t_1) & q_2(t_1) & \cdots & q_n(t_1) \\ q_1(t_2) & q_2(t_2) & \cdots & q_n(t_2) \\ \vdots & \vdots & \ddots & \vdots \\ q_1(t_m) & q_2(t_m) & \cdots & q_n(t_m) \end{bmatrix}. \quad (3)$$

A matrix of derivatives in time, \mathbf{Q}_t , is defined similarly and can be numerically computed from \mathbf{Q} . In this case, Eq. (2) becomes $\mathbf{Q}_t \approx \boldsymbol{\Theta}(\mathbf{Q}) \boldsymbol{\Xi}$ and the goal of the SINDy sparse regression problem is to choose a sparse set of coefficients $\boldsymbol{\Xi}$ that accurately fits the measured data in \mathbf{Q}_t . We can promote sparsity in the identified coefficients via a sparse regularizer $R(\boldsymbol{\Xi})$, such as the l_0 or l_1 norm, and use a sparse regression algorithm such as SR3 (Champion et al., 2020) to solve the resulting optimization problem,

$$\operatorname{argmin}_{\boldsymbol{\Xi}} \|\mathbf{Q}_t - \boldsymbol{\Theta}(\mathbf{Q}) \boldsymbol{\Xi}\|^2 + R(\boldsymbol{\Xi}). \quad (4)$$

The original PySINDy package was developed to identify a particular class of systems described by Eq. (1). Recent variants of the SINDy method are available that address systems with control inputs and model predictive control (MPC) (Fasel et al., 2021; Kaiser et al., 2018), systems with physical constraints (Kaptanoglu, Morgan, et al., 2021; Loiseau & Brunton, 2018), implicit ODEs (Kaheman et al., 2020; Mangan et al., 2016), PDEs (Rudy et al., 2017; Schaeffer, 2017), and weak form ODEs and PDEs (Messenger & Bortz, 2021; Reinbold et al., 2020; Schaeffer & McCalla, 2017). Other methods, such as ensembling and sub-sampling (Delahunt & Kutz, 2021; Maddu et al., 2019; Reinbold et al., 2021), are often vital for making the identification of Eq. (1) more robust. In order to incorporate these new developments and accommodate the wide variety of possible dynamical systems, we have extended PySINDy to a more general setting and added significant new functionality. Our code¹ is thoroughly documented, contains extensive examples, and integrates a wide range of functionality, some of which may be found in a number of other local SINDy implementations². In contrast to

¹<https://github.com/dynamicslab/pysindy>

²<https://github.com/snagcliffs/PDE-FIND>, <https://github.com/eurika-kaiser/SINDY-MPC>, <https://github.com/dynamicslab/SINDy-PI>, <https://github.com/SchatzLabGT/SymbolicRegression>, https://github.com/dynamicslab/databook_python, <https://github.com/sheadan/SINDy-BVP>, <https://github.com/sethbirsh/BayesianSindy>, <https://github.com/racdale/sindyr>, <https://github.com/SciML/DataDrivenDiffEq.jl>, https://github.com/MathBioCU/WSINDy_PDE, https://github.com/pakreinbold/PDE_Discovery_Weak_Formulation, <https://github.com/ZIB-IOL/CINDy>

74 some of these existing implementations, PySINDy is completely open-source, professionally-
 75 maintained (for instance, providing unit tests and adhering to PEP8 stylistic standards), and
 76 minimally dependent on non-standard Python packages.

77 New features

78 Given spatiotemporal data $\mathbf{Q}(\mathbf{x}, t) \in \mathbb{R}^{m \times n}$, and optional control inputs $\mathbf{u} \in \mathbb{R}^{m \times r}$ (note m
 79 has been redefined here to be the product of the number of spatial measurements and the
 80 number of time samples), PySINDy can now approximate algebraic systems of PDEs (and
 81 corresponding weak forms) in up to 3 spatial dimensions. Assuming the system is described
 82 by a function g , we have

$$g(\mathbf{q}, \mathbf{q}_t, \mathbf{q}_x, \mathbf{q}_y, \mathbf{q}_{xx}, \dots, \mathbf{u}) = 0. \quad (5)$$

83 ODEs, implicit ODEs, PDEs, and other dynamical systems are subsets of Eq. (5). We can
 84 accommodate control terms and partial derivatives in the SINDy library by adding them as
 85 columns in $\Theta(\mathbf{Q})$, which becomes $\Theta(\mathbf{Q}, \mathbf{Q}_t, \mathbf{Q}_x, \dots, \mathbf{u})$.

86 In addition, we have extended PySINDy to handle more complex modeling scenarios, includ-
 87 ing trapping SINDy for provably stable ODE models for fluids (Kaptanoglu, Callaham, et al.,
 88 2021), models trained using multiple dynamic trajectories, and the generation of many mod-
 89 els with sub-sampling and ensembling methods for cross-validation and probabilistic system
 90 identification. In order to solve Eq. (5), PySINDy implements several different sparse regres-
 91 sion algorithms. Greedy sparse regression algorithms, including step-wise sparse regression
 92 (SSR) (Boninsegni et al., 2018) and forward regression orthogonal least squares (FROLS)
 93 (Billings, 2013), are now available. Figure 1 illustrates the PySINDy code structure, changes,
 94 and high-level goals for future work.

95 PySINDy includes extensive Jupyter notebook tutorials that demonstrate the usage of various
 96 features of the package and reproduce nearly the entirety of the examples from the original
 97 SINDy paper (Steven L. Brunton et al., 2016), trapping SINDy paper (Kaptanoglu, Callaham,
 98 et al., 2021), and the PDE-FIND paper (Rudy et al., 2017). We include an extended example
 99 for the quasiperiodic shear-driven cavity flow (Callaham et al., 2021). As a simple illustration
 100 of the new functionality, we demonstrate how SINDy can be used to identify the Kuramoto-
 101 Sivashinsky (KS) PDE from data. We train the model on the first 60% of the data from Rudy
 102 et al. (Rudy et al., 2017), which in total contains 1024 spatial grid points and 251 time steps.
 103 The KS model is identified correctly and the prediction for $\dot{\mathbf{q}}$ on the remaining testing data
 104 indicates strong performance in Figure 2. Lastly, we provide a useful flow chart in Figure 3
 105 so that users can make informed choices about which advanced methods are suitable for their
 106 datasets.

107 Conclusion

108 The goal of the PySINDy package is to enable anyone with access to measurement data to
 109 engage in scientific model discovery. The package is designed to be accessible to inexperienced
 110 users, adhere to scikit-learn standards, include most of the existing SINDy variations in
 111 the literature, and provide a large variety of functionality for more advanced users. We hope
 112 that researchers will use and contribute to the code in the future, pushing the boundaries of
 113 what is possible in system identification.

114 Acknowledgments

115 PySINDy is a fork of [sparsereg](#) (Quade, 2018). SLB, AAK, KK, and UF acknowledge support
116 from the Army Research Office (ARO W911NF-19-1-0045). JLC acknowledges support from
117 funding support from the Department of Defense (DoD) through the National Defense Science
118 & Engineering Graduate (NDSEG) Fellowship Program.

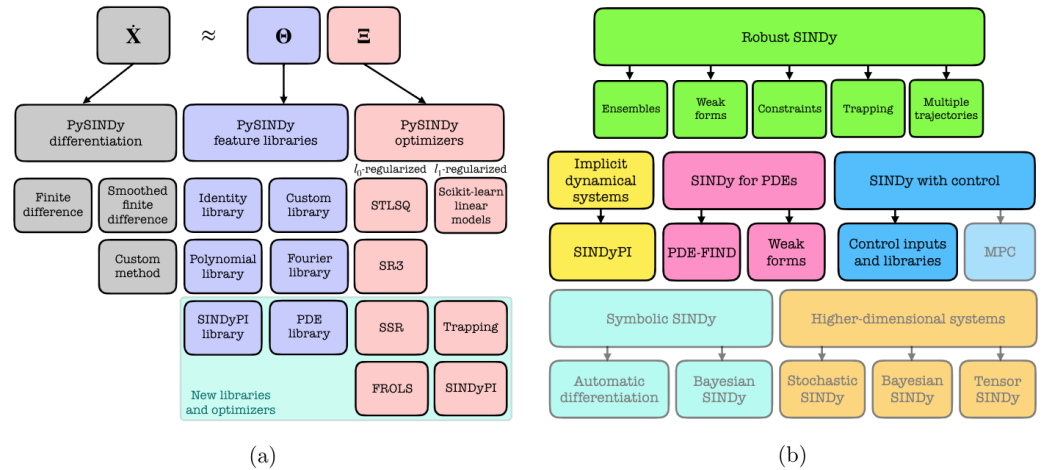


Figure 1: Summary of SINDy features organized by (a) PySINDy structure and (b) functionality. (a) Hierarchy from the sparse regression problem solved by SINDy, to the submodules of PySINDy, to the individual optimizers, libraries, and differentiation methods implemented in the code. (b) Flow chart for organizing the SINDy variants and functionality in the literature. Bright color boxes indicate the features that have been implemented through this work, roughly organized by functionality. Semi-transparent boxes indicate features that have not yet been implemented.

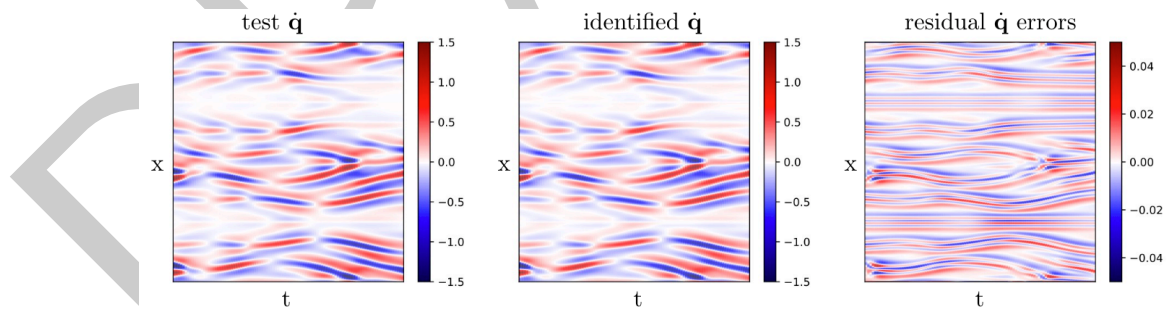


Figure 2: PySINDy can now be used for PDE identification; we illustrate this new capability by accurately capturing a set of testing data from the Kuramoto-Sivashinsky system, described by $q_t = -qq_x - q_{xx} - q_{xxx}$. The identified model is $q_t = -0.98qq_x - 0.99q_{xx} - 1.0q_{xxx}$.

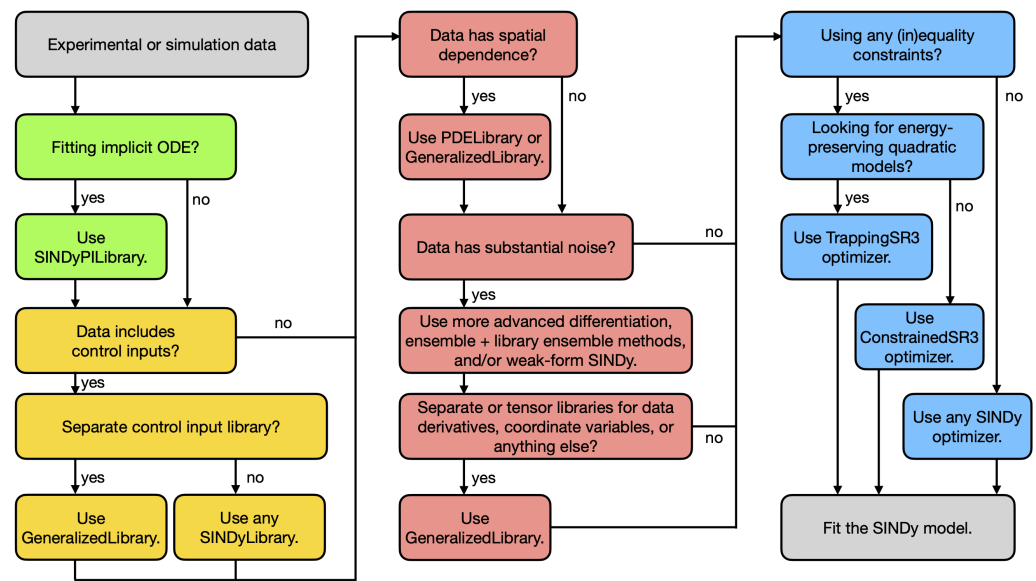


Figure 3: This flow chart summarizes how PySINDy users can start with a dataset and systematically choose the proper candidate library and sparse regression optimizer that are tailored for a specific scientific task. The GeneralizedLibrary class allows for tensoring, concatenating, and otherwise combining many different candidate libraries.

References

- Benner, P., Gugercin, S., & Willcox, K. (2015). A survey of projection-based model reduction methods for parametric dynamical systems. *SIAM Review*, 57(4), 483–531. <https://doi.org/10.1137/130932715>
- Billings, S. A. (2013). *Nonlinear system identification: NARMAX methods in the time, frequency, and spatio-temporal domains*. John Wiley & Sons.
- Bongard, J., & Lipson, H. (2007). Automated reverse engineering of nonlinear dynamical systems. *Proc. Natl. Acad. Sciences*, 104(24), 9943–9948. <https://doi.org/10.1073/pnas.0609476104>
- Boninsegna, L., Nüske, F., & Clementi, C. (2018). Sparse learning of stochastic dynamical equations. *The Journal of Chemical Physics*, 148(24), 241723. <https://doi.org/10.1063/1.5018409>
- Brunton, Steven L., Budišić, M., Kaiser, E., & Kutz, J. N. (2021). Modern Koopman theory for dynamical systems. *arXiv Preprint arXiv:2102.12086*.
- Brunton, Steven L., Proctor, J. L., & Kutz, J. N. (2016). Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15), 3932–3937. <https://doi.org/10.1073/pnas.1517384113>
- Callahan, J. L., Brunton, S. L., & Loiseau, J.-C. (2021). On the role of nonlinear correlations in reduced-order modeling. *arXiv Preprint arXiv:2106.02409*.
- Champion, K., Zheng, P., Aravkin, A. Y., Brunton, S. L., & Kutz, J. N. (2020). A unified sparse optimization framework to learn parsimonious physics-informed models from data. *IEEE Access*, 8, 169259–169271. <https://doi.org/10.1109/access.2020.3023625>
- de Silva, B., Champion, K., Quade, M., Loiseau, J.-C., Kutz, J. N., & Brunton, S. (2020). PySINDy: A Python package for the sparse identification of nonlinear dynamical systems

- 144 from data. *Journal of Open Source Software*, 5(49), 1–4. [https://doi.org/10.21105/joss.](https://doi.org/10.21105/joss.02104)
145 02104
- 146 Delahunt, C. B., & Kutz, J. N. (2021). A toolkit for data-driven discovery of governing
147 equations in high-noise regimes. *arXiv Preprint arXiv:2111.04870*.
- 148 Fasel, U., Kaiser, E., Kutz, J. N., Brunton, B. W., & Brunton, S. L. (2021). SINDy with
149 control: A tutorial. *arXiv Preprint arXiv:2108.13404*.
- 150 Kaheman, K., Kutz, J. N., & Brunton, S. L. (2020). SINDy-PI: A robust algorithm for parallel
151 implicit sparse identification of nonlinear dynamics. *Proceedings of the Royal Society A*,
152 476(2242), 20200279. <https://doi.org/10.1098/rspa.2020.0279>
- 153 Kaiser, E., Kutz, J. N., & Brunton, S. L. (2018). Sparse identification of nonlinear dynamics
154 for model predictive control in the low-data limit. *Proceedings of the Royal Society of*
155 *London A*, 474(2219). <https://doi.org/10.1098/rspa.2018.0335>
- 156 Kaptanoglu, A. A., Callahan, J. L., Aravkin, A., Hansen, C. J., & Brunton, S. L. (2021).
157 Promoting global stability in data-driven models of quadratic nonlinear dynamics. *Phys.*
158 *Rev. Fluids*, 6, 094401. <https://doi.org/10.1103/PhysRevFluids.6.094401>
- 159 Kaptanoglu, A. A., Morgan, K. D., Hansen, C. J., & Brunton, S. L. (2021). Physics-
160 constrained, low-dimensional models for magnetohydrodynamics: First-principles and data-
161 driven approaches. *Phys. Rev. E*, 104, 015206. [https://doi.org/10.1103/physreve.104.](https://doi.org/10.1103/physreve.104.015206)
162 015206
- 163 Kutz, J. N., Brunton, S. L., Brunton, B. W., & Proctor, J. L. (2016). *Dynamic mode*
164 *decomposition: Data-driven modeling of complex systems*. SIAM.
- 165 Loiseau, J.-C., & Brunton, S. L. (2018). Constrained sparse Galerkin regression. *Journal of*
166 *Fluid Mechanics*, 838, 42–67. <https://doi.org/10.1017/jfm.2017.823>
- 167 Maddu, S., Cheeseman, B. L., Sbalzarini, I. F., & Müller, C. L. (2019). Stability selection
168 enables robust learning of partial differential equations from limited noisy data. *arXiv*
169 *Preprint arXiv:1907.07810*.
- 170 Mangan, N. M., Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2016). Inferring biological
171 networks by sparse identification of nonlinear dynamics. *IEEE Transactions on Molecu-*
172 *lar, Biological and Multi-Scale Communications*, 2(1), 52–63. [https://doi.org/10.1109/](https://doi.org/10.1109/tmbmc.2016.2633265)
173 [tmbmc.2016.2633265](https://doi.org/10.1109/tmbmc.2016.2633265)
- 174 Messenger, D. A., & Bortz, D. M. (2021). Weak SINDy for partial differential equations.
175 *Journal of Computational Physics*, 110525. <https://doi.org/10.1016/j.jcp.2021.110525>
- 176 Pathak, J., Hunt, B., Girvan, M., Lu, Z., & Ott, E. (2018). Model-free prediction of large
177 spatiotemporally chaotic systems from data: A reservoir computing approach. *Physical*
178 *Review Letters*, 120(2), 024102. <https://doi.org/10.1103/physrevlett.120.024102>
- 179 Peherstorfer, B., & Willcox, K. (2016). Data-driven operator inference for nonintrusive
180 projection-based model reduction. *Computer Methods in Applied Mechanics and Engi-*
181 *neering*, 306, 196–215. <https://doi.org/10.1016/j.cma.2016.03.025>
- 182 Qian, E., Kramer, B., Peherstorfer, B., & Willcox, K. (2020). Lift & Learn: Physics-informed
183 machine learning for large-scale nonlinear dynamical systems. *Physica D: Nonlinear Phe-*
184 *nomena*, 406, 132401. <https://doi.org/10.1016/j.physd.2020.132401>
- 185 Quade, M. (2018). *Sparsereg - collection of modern sparse regression algorithms*. <https://doi.org/10.5281/zenodo.1173754>
- 186
- 187 Raissi, M., Perdikaris, P., & Karniadakis, G. (2019). Physics-informed neural networks: A
188 deep learning framework for solving forward and inverse problems involving nonlinear partial
189 differential equations. *Journal of Computational Physics*, 378, 686–707. [https://doi.org/](https://doi.org/10.1016/j.jcp.2018.10.045)
190 [10.1016/j.jcp.2018.10.045](https://doi.org/10.1016/j.jcp.2018.10.045)

- 191 Raissi, Maziar, Perdikaris, P., & Karniadakis, G. E. (2017). Machine learning of linear differen-
192 tial equations using Gaussian processes. *Journal of Computational Physics*, 348, 683–693.
193 <https://doi.org/10.1016/j.jcp.2017.07.050>
- 194 Reinbold, P. A., Gurevich, D. R., & Grigoriev, R. O. (2020). Using noisy or incomplete
195 data to discover models of spatiotemporal dynamics. *Physical Review E*, 101(1), 010203.
196 <https://doi.org/10.1103/physreve.101.010203>
- 197 Reinbold, P. A., Kageorge, L. M., Schatz, M. F., & Grigoriev, R. O. (2021). Robust
198 learning from noisy, incomplete, high-dimensional experimental data via physically con-
199 strained symbolic regression. *Nature Communications*, 12(1), 1–8. [https://doi.org/10.](https://doi.org/10.1038/s41467-021-23479-0)
200 [1038/s41467-021-23479-0](https://doi.org/10.1038/s41467-021-23479-0)
- 201 Rudy, S. H., Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2017). Data-driven discovery of
202 partial differential equations. *Science Advances*, 3(e1602614). [https://doi.org/10.1126/](https://doi.org/10.1126/sciadv.1602614)
203 [sciadv.1602614](https://doi.org/10.1126/sciadv.1602614)
- 204 Schaeffer, H. (2017). Learning partial differential equations via data discovery and sparse
205 optimization. *Proceedings of the Royal Society a*, 473, 20160446. [https://doi.org/10.](https://doi.org/10.1098/rspa.2016.0446)
206 [1098/rspa.2016.0446](https://doi.org/10.1098/rspa.2016.0446)
- 207 Schaeffer, H., & McCalla, S. G. (2017). Sparse model selection via integral terms. *Physical*
208 *Review E*, 96(2), 023302. <https://doi.org/10.1103/physreve.96.023302>
- 209 Schmid, P. J. (2010). Dynamic mode decomposition of numerical and experimental data.
210 *Journal of Fluid Mechanics*, 656, 5–28. <https://doi.org/10.1017/s0022112010001217>
- 211 Schmidt, M., & Lipson, H. (2009). Distilling free-form natural laws from experimental data.
212 *Science*, 324(5923), 81–85. <https://doi.org/10.1126/science.1165893>
- 213 Vlachas, P. R., Byeon, W., Wan, Z. Y., Sapsis, T. P., & Koumoutsakos, P. (2018). Data-driven
214 forecasting of high-dimensional chaotic systems with long short-term memory networks.
215 *Proc. R. Soc. A*, 474(2213), 20170844. <https://doi.org/10.1098/rspa.2017.0844>