

---

# Sparse Regression and Adaptive Feature Generation for the Discovery of Dynamical Systems

---

Chinmay S. Kulkarni chinmayk@mit.edu  
Department of Mechanical Engineering, MIT  
6.860: Statistical Learning Theory and Applications, Fall 2018, MIT

## Abstract

We study the performance of sparse regression methods and propose new techniques to distill the governing equations of dynamical systems from data. We first look at the generic methodology of learning interpretable equation forms from data, proposed by Brunton et al. [3], followed by performance of LASSO for this purpose. We then propose a new algorithm that uses the dual of LASSO optimization for higher accuracy and stability. In the second part, we propose a novel algorithm that learns the candidate function library in a completely data-driven manner to distill the governing equations of the dynamical system. This is achieved via sequentially thresholded ridge regression (STRidge [17]) over a orthogonal polynomial space. The performance of the three discussed methods is illustrated by looking the Lorenz 63 system and the quadratic Lorenz system.

## 1 Introduction and Overview

The role of data today is no longer limited to the verification of first-principles-derived-models but it is also being used to learn such models. This is particularly important in non-autonomous nonlinear dynamical systems that describe a multitude of problems from science and engineering. Recent methods leverage the fact that most dynamical equations governing physical systems contain a few terms, making them sparse in high-dimensional nonlinear function space [3, 17]. By constructing an appropriate feature library based on the data coordinates, one can apply sparse regression to discover the governing equations of the dynamical system. Initial work in this field has mainly focused on the behavior of the approach with regards to noise, multi-fidelity data etc., however few have tried to improve upon the sparse regression algorithm at the core of the approach.

This is exactly the first focus area of the current work. We first look at the sparse regression method most commonly employed in this field: LASSO [18]. We point out the theoretical underpinnings of LASSO along with accuracy bounds. Although LASSO works well for uncorrelated features, for highly correlated it tends to choose a feature at random from each of the correlated groups [9]. When the feature library is built from data, higher order features are more correlated with each other, where LASSO potentially runs into difficulties. To alleviate these difficulties, we propose to solve the dual of LASSO to learn the governing equations. Even in the case of correlated features, the dual LASSO has a unique solution, which then allows us to correctly choose the present features. The second part of this work deals with the case when the exact function blocks that describe the dynamical system are not present in the feature library. We propose and investigate an approach to handle such cases by using an appropriate family of orthogonal functional basis to span the feature library. As the result may not be sparse in terms of the orthogonal functions, we resort to sequentially thresholded least squares (STRidge) algorithm [17] that does not impose sparsity, but chooses the dominant features.

These approaches are demonstrated on the Lorenz 63 system [15] and the quadratic Lorenz system [8]. We compare the performance of the LASSO and dual LASSO, and then use the approach of data-driven orthogonal regression to learn the feature library.

## 1.1 General Methodology

Let us assume that we have  $n$  state space parameters  $(x_1, \dots, x_n)$ , with measurements for  $x_i$  and  $\dot{x}_i = dx/dt$  at times  $t = 1, \dots, T$  (denoted by a superscript). If only state observations are available, the rate parameters can be computed using finite difference. This is followed by constructing a non-linear library of features using the state space parameters. The span of these features now describes the feature space. Typically we would construct this feature space through a class of functions that are dense in the space that our dynamical system lives in. In this work we assume a polynomial feature library. After constructing the feature library (say  $X$ ), we formulate the regression problem as  $\dot{X} = XW + \varepsilon$ , where  $\dot{X}_{(t,j)} = \dot{x}_j^t$  and  $W$  are the unknown weights, with  $\varepsilon$  being the noise.

Often, dynamical models have a functional form where the terms live in a smaller manifold in the feature space. We thus use sparse regression to pick the relevant features. These features, with their corresponding coefficients describe the functional form of the governing equations. This procedure is pictorially illustrated in figure 1. For noisy observations, we carry out data filtering and smoothing at the input step and perform multiple sparse regression solves until convergence. However we do not focus on this, and we assume that we have noise free observations.

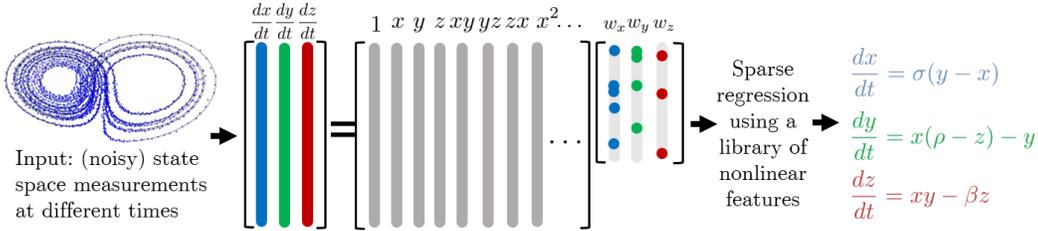


Figure 1: Schematic of the generic methodology used.

## 2 Sparse regression over fixed feature space

In this case, we assume that the feature library is fixed, and that we wish to find either the exact sparse equation form from this library or the closest approximation to the governing equation only from the terms in the library. The highest polynomial degree in the feature space ( $X$ ) is  $p$ . Then, the feature space contains terms of the form  $(x_1^t)^{p_1} \dots (x_n^t)^{p_n}$ , such that  $p_1 + \dots + p_n \leq p$ . The number of solutions to this system (*i.e.* the number of terms in the feature library) is  $m = \binom{n+p}{n} = \frac{(n+p)!}{n!p!}$ . Whereas, empirically the number of distinct terms in the governing equations is  $\mathcal{O}(n)$ . Thus even for small enough  $p$ , the terms in the feature library are much more in number than those to be chosen, which justifies sparse regression to select the features. Let us denote the coefficient matrix, obtained from the sparse regression solve by  $W$ . The optimization problem with some form of penalty ( $\mathcal{P}$ ) is:

$$\min_W \mathcal{L}(W) = \left[ \left( \dot{X} - XW \right)^2 + \mathcal{P}(W) \right] \quad \text{where } \dot{X} \in \mathbb{R}^{T \times n}, X \in \mathbb{R}^{T \times m} \text{ and } W \in \mathbb{R}^{m \times n} \quad (1)$$

To further select the features appropriately, we use our knowledge of the underlying physics of the dynamical system. Often, we have knowledge about the general scales of the state space parameters (for example, in the Lorenz system, the state space parameters represent the rate of convection, horizontal temperature variation and the vertical temperature variation). Thus, we know the typical scales of these parameters from physics. Let us denote the scale of  $x_i$  by  $L_i$ , then  $(x_1^t)^{p_1} \dots (x_n^t)^{p_n} \approx L_1^{p_1} \dots L_n^{p_n}$ . Thus, instead of selecting features based on their coefficient values, we select features by looking at their net magnitude, by replacing the state space parameters with their respective scales. We refer to this as ‘scale based thresholding’.

### 2.1 Sparsity through $L_1$ regularization (LASSO)

We now look at the use of LASSO for imposing sparsity. As is well-known, the LASSO penalty is  $\mathcal{P}(W) = \lambda \|W\|_1$ , which serves as a convex counterpart to the non-convex  $L_0$  norm. Kakade et al. [12] showed that the Rademacher averages for linear regressors with  $L_1$  penalty are bounded by  $\mathcal{R}_n(\mathcal{F}_W) \leq XW_{\max} \sqrt{\frac{2 \log(m)}{n}}$ , where  $W_{\max} = \max \|W\|_1$ . We can then use the contraction

property, assuming the Lipschitz constant of the LASSO problem to be  $L$  (more details and methods to compute it may be found in [14, 13]) to obtain the bound of Rademacher averages of the loss function class  $\mathcal{G}$ . Using the symmetrization lemma, this yields a bound on the expected maximum error (uniform deviation), as given by Eq. (2).

$$\mathbb{E} \max_{g \in \mathcal{G}} \left[ \mathbb{E} g(x_i) - \frac{1}{n} \sum_{i=1}^n g(x_i) \right] \leq 2L X W_{\max} \sqrt{\frac{2 \log(m)}{n}} \quad (2)$$

Note that even though we cannot explicitly compute the right hand side, we have  $m \gg \gg n$ , which renders this accuracy bound impractical. Thus, LASSO in the case of exceedingly high number of features does not provide good accuracy bounds. To alleviate this problem, we first use the SAFE bounds provided by Ghaoui et al. [10] to remove the irrelevant features and then apply LASSO over the remaining library. Even after weeding out the absent features, another important task that remains is to choose the hyperparameter  $\lambda$ . Literature widely suggests the order of  $\lambda$  to be  $\lambda \approx \sqrt{T \log(m)}$  [2], however choosing this parameter independent of the degree of correlation is not recommended [5]. It is further shown that the larger the correlation amongst the features, the smaller the value of  $\lambda$  [11]. Hebiri and Lederer [11] also prove that LASSO can achieve fast rate bounds (*i.e.* error bounded by  $\lambda^2$ ) even for correlated features. Even though reassuring, in practice this requires significant amount of tuning and cross validation, which is a major issue in the case at hand, as typically we only have sparse and limited observations of the physical states. Further, even though the LASSO fit is unique in case of high correlations, they are problematic for parameter estimation and feature selection [16, 19]. The pitfalls of LASSO (even after removing the irrelevant features using the SAFE rules) are that it requires significant hyperparameter tuning and it is extremely sensitive to  $\lambda$  for correlated features (observed empirically). These motivate us to instead formulate a new approach to solve the sparse regression problem.

## 2.2 Sparsity through dual LASSO regression

To overcome the difficulties in the application of LASSO (along with the SAFE rules), we formulate and solve its dual problem. In order for the LASSO solution (that is, the feature and parameter choices) to be unique, the feature matrix must satisfy the irrepresentability condition (IC) of Zhao and Yu [22] and beta-min condition of Bühlmann and Van De Geer [4]. The feature library violates the IC for highly correlated columns, leading to an unstable feature selection. However, even for highly correlated features, the corresponding dual LASSO solution is always unique [19]. The dual problem is given by Eq. (3) [16], which is strictly convex in  $\theta$  (implying a unique solution).

$$\max_{\theta} \mathcal{D}(\theta) = \|\dot{X}\|_2^2 - \|\theta - \dot{X}\|_2^2 \quad \text{such that} \quad \|X^T \theta\|_{\infty} \leq \lambda \quad (3)$$

Now, let  $\hat{W}$  be a solution of Eq. (1) with LASSO penalty and  $\hat{\theta}$  be the unique solution to the corresponding dual problem Eq. (3). Then stationarity condition implies  $\hat{\theta} = \dot{X} - X\hat{W}$ . Even though LASSO does not have a unique  $\hat{W}$ , the fitted value  $X\hat{W}$  is unique, as the optimization problem Eq. (1) is strongly convex in  $XW$  for  $\mathcal{P}(W) = \lambda \|W\|_1$ . We make use of this by first computing a solution to the primal LASSO problem and then computing the unique dual solution by using the primal fitted value and Eq. (2). Once we have the unique dual solution  $\hat{\theta}$ , we do feature selection by using the dual active set, which is same as the primal active set with high probability under the IC [9]. It is also shown that features discarded by SAFE rules are the inactive constraints over the optimal dual solution [20], which justifies omitting them. KKT conditions imply:

$$\hat{\theta}^T X_i \begin{cases} = \text{sign}(\hat{W}_i) & \text{if } \hat{W}_i \neq 0 \\ \in (-1, 1) & \text{if } \hat{W}_i = 0 \end{cases} \quad (4)$$

Eq. (4) gives us a direct way to compute the active dual set (which is the same as the active primal set with high probability) once we have  $\hat{\theta}$ . We discard the features for which  $\hat{\theta}^T X_i \in (-1, 1)$  and retain the others. This does not give us a good fit of the solution, so to compute the coefficients accurately, we perform ridge regression ( $\mathcal{P}(W) = \lambda_2 \|W\|_2^2$ ) over these active features. This pseudocode for this is given in algorithm 1. We refer to this algorithm as ‘dual LASSO’.

## 3 Sparse Regression Over Adaptive Feature Space

In this section, we look at cases where the feature library is not known. If we have no prior belief over the form of the equations, we may not be able to construct an efficient feature library. In such

---

**Algorithm 1** Sparse regression using dual LASSO

---

**Require:** state parameters:  $\mathbf{x} = x_i^t, \hat{\mathbf{x}} = \hat{x}_i^t$ ; LASSO penalty  $\lambda$ , ridge penalty  $\lambda_2$

Compute the primal LASSO solution:  $\hat{W} = \min_W \left[ \left( \dot{X} - XW \right)^2 + \lambda \|W\|_1 \right]$

Compute the unique dual solution  $\hat{\theta} = \dot{X} - X\hat{W}$

Compute dual active set (same as primal active with h.p.)  $S^d = \{1 \leq j \leq m : \hat{\theta}^T X_j \notin (-1, 1)\}$

Construct reduced feature matrix  $\tilde{X}$  by only considering features whose indices are in  $S^d$

Solve ridge regression  $W^* = \min_W \left[ \left( \dot{X} - XW \right)^2 + \lambda_2 \|W\|_2^2 \right]$

---

situations, learning this library from data might be the most advantageous choice. One approach is to make the use of orthogonal functions of some parametric family, such as polynomials to construct this library. This ensures no correlation between the features and affords maximum freedom to the algorithm. The drawback in this case is that the regressor may not be sparse over this feature library.

### 3.1 Adaptive feature space growth through orthogonal regression

We propose the following approach for this purpose: Starting with an empty library, we recursively add a feature to it and compute the corresponding loss function of the resulting fit by using STRidge (Sec. 3.2). If the loss function decreases by more than a certain fraction, we keep this feature. Otherwise we discard it and look at the next orthogonal feature. Once every few of the addition timesteps, we perform a removal step. That is, we discard the feature(s) that do not result in a significant increase in the loss function (if any). This ensures that we do not keep lower order functions that may not be required to describe the equations as higher order functions are added. Our algorithm is inspired by previous greedy feature development algorithms by Efron et al. [7] (LARS), Zhang [21] (FoBa) etc. However these require pre-determined full possible feature space, whereas we construct new features on the fly. Once the equations are obtained in terms of these orthogonal polynomials, we distill their sparse forms by using symbolic equation simplification [1].

### 3.2 Sequentially thresholded ridge regression

To compute regressors over the orthogonal feature space, we use sequentially thresholded ridge regression (STRidge), proposed by Rudy et al. [17]. The idea is simple: we iteratively compute the ridge regression solution with decreasing penalty proportional to the condition number of  $X$ , and discard the components using scale based thresholding (Sec. 2). As the feature matrix is orthonormal by construction, the analytical solution is  $W = (1 + \lambda)^{-1} X^T \dot{X}$ . This algorithm is effective in choosing optimal features without necessitating sparsity. The overall pseudocode for learning the governing equations through adaptively growing the feature library is given by algorithm 2, and the corresponding results are presented in Sec. 4.

## 4 Results

For the first part (Sec. 2), our testbed will be the Lorenz 63 system ( $n = 3$ ), given by Eq. (5).

$$\dot{x} = 10(yz - x); \quad \dot{y} = x(28 - z); \quad \dot{z} = xy - 2.667z \quad (5)$$

For the first part, we consider three polynomial feature libraries with  $p = 3, 10$  and  $20$  ( $m = 20, 286$  and  $1771$ ). The idea behind considering larger orders ( $p$ ) is that it highlights the poor performance of LASSO for highly correlated features. Fig. (2a) shows the residual after the LASSO and the dual LASSO model fits. We can see that the residuals for both methods are comparable for all of the feature libraries, which empirically shows that the model fit even for a high degree of correlations is good for both the methods. Fig. (2b) plots the number of non-zero features in the equations for different  $p$  values. LASSO has a much higher number of non-zero terms, and this number increases significantly with  $p$  (and  $m$ ), indicating instability of the solution. Dual LASSO performs very well, and the number of present features does not change for the most part with  $p$ . Fig. (2c) plots the absolute weights for the components for the  $p = 3$  case for the  $\dot{y}$  equation. Dual LASSO retrieves the

---

**Algorithm 2** Learning the governing equations through adaptive growth of the feature library
 

---

**Require:** state parameters:  $\mathbf{x} = x_i^t, \dot{\mathbf{x}} = \dot{x}_i^t$ ; orthogonal family  $F_j(\bullet)$ ; feature addition / removal thresholds:  $r_a (\leq 1), r_r (\geq 1), \lambda_0$ ; removal step frequency  $k_r$

Initialize:  $X = \emptyset, W = \mathbf{0}, t = 0, \mathcal{L} = \infty$

**while** True **do**

$X_t = \text{append}(X, F_k(\mathbf{x}))$

    Solve the STRidge problem:  $W_t = \text{STRidge}(\dot{X}, X_t, \lambda_0)$

    Compute the loss  $\mathcal{L}_t = (\dot{X} - X_t W_t)^2$

**if**  $\mathcal{L}_t \leq r_a \mathcal{L}$  **then**

$X = X_t; W = W_t$

**if**  $\text{mod}(k, k_r) == 0$  **then**

**for**  $i = 1, \dots, X.\text{shape}[1]$  **do** (number of columns of  $X$ )

$X_t = \text{append}(X[:, 1 : i - 1], X[:, i + 1 : \text{end}])$  (ignore the  $i^{\text{th}}$  column of  $X$ )

            Solve the STRidge problem:  $W_t = \text{STRidge}(\dot{X}, X_t, \lambda_0)$

            Compute the loss  $\mathcal{L}_t = (\dot{X} - X_t W_t)^2$

**if**  $\mathcal{L}_t \leq r_r \mathcal{L}$  **then**

$X = X_t; W = W_t$

$k = k + 1.$

**break** if no change in feature space over multiple iterations.

Perform symbolic simplification of  $\dot{X} = XW$  to obtain the final form of the equations

---

correct features (with accurate weights), while LASSO detects the correct features but also detects high order features that have low weights and are highly correlated to each other. This serves as a great validation of the superiority of dual LASSO over conventional LASSO for model discovery.

Now we look at the results for Sec. 3, on the quadratic Lorenz system (Eq. (6)) [8].

$$\dot{x} = 10(yz - x); \quad \dot{y} = x(28 - z); \quad \dot{z} = (xy)^2 - 2.667z \quad (6)$$

We start with an empty feature library and  $W = \mathbf{0}$  and iteratively grow the feature space using algorithm 2 using Legendre polynomials (denoted by  $\mathbb{L}_p(\bullet)$ ),  $r_a = 0.75, r_r = 1.25, \lambda_0 = 1$  and removal step working every 10 addition steps ( $k_r = 10$ ). The final obtained model is given in Eq. (7). We use symbolic simplification [1] to simply this model. Note that the model is not unique before symbolic simplification (as  $\mathbb{L}_0(x) = \mathbb{L}_0(y) = \mathbb{L}_0(z)$ ). We also do scale based thresholding (Sec. 2) once after the symbolic simplification to remove any terms that may remain due to approximate factorization. The results before and after scale based thresholding are given in Eq. (8). We can see that this algorithm does a good job at identifying the correct governing equations by iteratively building the feature library, and does not include any incorrect features.

$$\begin{aligned} \dot{x} &= 9.93\mathbb{L}_1(y)\mathbb{L}_1(z) - 9.89\mathbb{L}_1(x) & \dot{z} &= 0.43\mathbb{L}_2(x)\mathbb{L}_2(y) + 0.22\mathbb{L}_2(x) + 0.21\mathbb{L}_2(y) \\ \dot{y} &= 27.66\mathbb{L}_1(x) - 1.04\mathbb{L}_1(x)\mathbb{L}_1(z) & & - 2.62\mathbb{L}_1(z) + 2.09\mathbb{L}_0(x) - 0.22\mathbb{L}_0(y) - 1.95\mathbb{L}_0(z) \end{aligned} \quad (7)$$

$$\begin{aligned} \dot{x} &= 9.93yz - 9.89x \\ \dot{y} &= 27.66x - 1.04xz \\ \dot{z} &= 0.97(xy)^2 + 0.007(x^2 - y^2) \\ &\quad - 2.62z + 0.027 \end{aligned} \quad \implies \quad \begin{aligned} \dot{x} &= 9.93yz - 9.89x \\ \dot{y} &= 27.66x - 1.04xz \\ \dot{z} &= 0.9675(xy)^2 - 2.62z \end{aligned} \quad (8)$$

after symbolic simplification
after scale based thresholding

## 5 Conclusions and Future Work

In this work, we investigated LASSO, proposed dual LASSO and data driven feature learning approaches to solve the problem of discovering governing equations only from state parameter data. Future work directions involve extending the ideas of feature library building where one can construct the function to be added through a mix of a larger family of orthogonal functions. Approaches to involve kernel compositionality [6] to build these libraries can also be investigated.

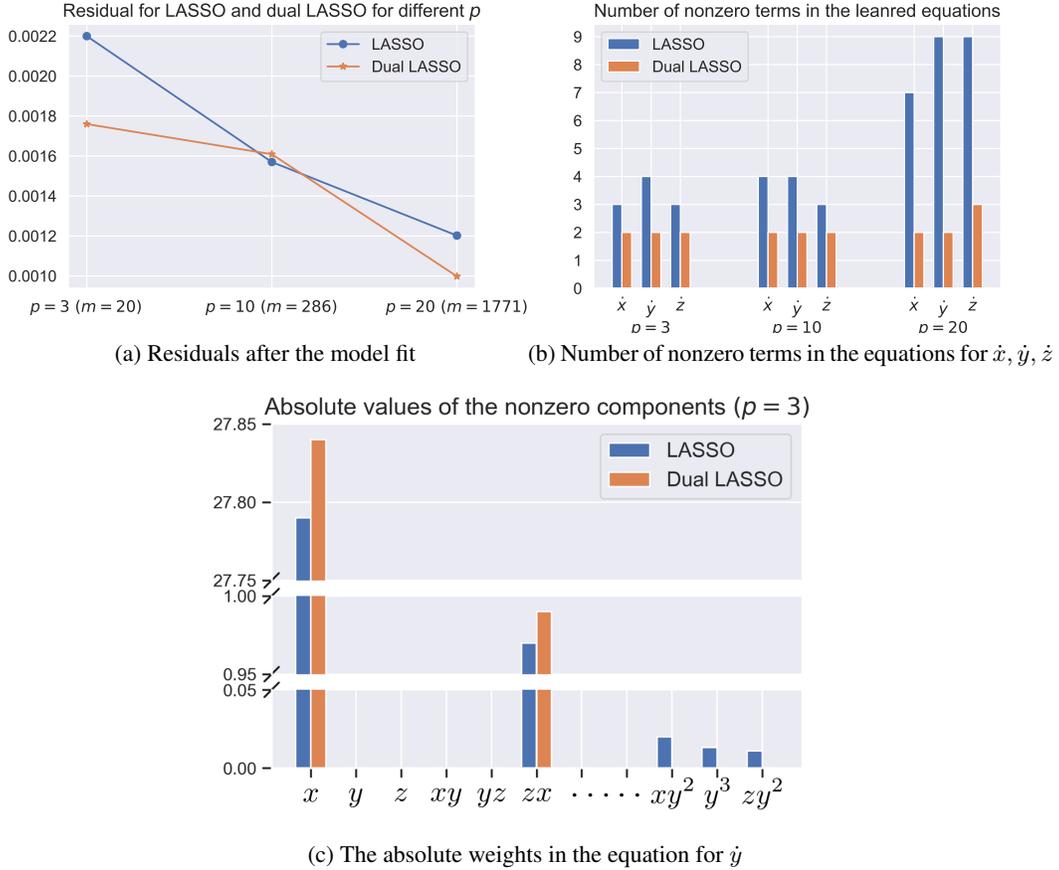


Figure 2: Normalized (combined) residual after model fit, corresponding number of non-zero coefficients in the equation for  $\dot{y}$ , and the absolute weights for the components in the  $\dot{y}$  equation ( $p=3$ ).

## References

- [1] David H Bailey, Jonathan M Borwein, and Alexander D Kaiser. Automated simplification of large symbolic expressions. *Journal of Symbolic Computation*, 60:120–136, 2014.
- [2] Peter J Bickel, Ya’acov Ritov, Alexandre B Tsybakov, et al. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- [3] Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–3937, 2016.
- [4] Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.
- [5] Arnak S Dalalyan, Mohamed Hebiri, Johannes Lederer, et al. On the prediction performance of the lasso. *Bernoulli*, 23(1):552–581, 2017.
- [6] David Duvenaud, James Robert Lloyd, Roger Grosse, Joshua B Tenenbaum, and Zoubin Ghahramani. Structure discovery in nonparametric regression through compositional kernel search. *arXiv preprint arXiv:1302.4922*, 2013.
- [7] Bradley Efron, Trevor Hastie, Iain Johnstone, Robert Tibshirani, et al. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [8] Buğçe Eminağa, Hatice Aktöre, and Mustafa Riza. A modified quadratic lorenz attractor. *arXiv preprint arXiv:1508.06840*, 2015.

- [9] Niharika Gauraha. Dual lasso selector. *arXiv preprint arXiv:1703.06602*, 2017.
- [10] Laurent El Ghaoui, Vivian Viallon, and Tarek Rabbani. Safe feature elimination for the lasso and sparse supervised learning problems. *arXiv preprint arXiv:1009.4219*, 2010.
- [11] Mohamed Hebiri and Johannes Lederer. How correlations influence lasso prediction. *IEEE Transactions on Information Theory*, 59(3):1846–1854, 2013.
- [12] Sham M Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *Advances in neural information processing systems*, pages 793–800, 2009.
- [13] Jakub Konecny and Peter Richtárik. Semi-stochastic gradient descent methods. *arXiv preprint arXiv:1312.1666*, 2013.
- [14] Jakub Konecny, Jie Liu, Peter Richtárik, and Martin Takac. Mini-batch semi-stochastic gradient descent in the proximal setting. *IEEE Journal of Selected Topics in Signal Processing*, 10(2):242–255, 2016.
- [15] Edward N Lorenz. Deterministic nonperiodic flow. *Journal of the atmospheric sciences*, 20(2):130–141, 1963.
- [16] Michael R Osborne, Brett Presnell, and Berwin A Turlach. On the lasso and its dual. *Journal of Computational and Graphical statistics*, 9(2):319–337, 2000.
- [17] Samuel H Rudy, Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Data-driven discovery of partial differential equations. *Science Advances*, 3(4):e1602614, 2017.
- [18] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- [19] Ryan J Tibshirani et al. The lasso problem and uniqueness. *Electronic Journal of Statistics*, 7: 1456–1490, 2013.
- [20] Jie Wang, Jiayu Zhou, Peter Wonka, and Jieping Ye. Lasso screening rules via dual polytope projection. In *Advances in Neural Information Processing Systems*, pages 1070–1078, 2013.
- [21] Tong Zhang. Adaptive forward-backward greedy algorithm for sparse learning with linear models. In *Advances in Neural Information Processing Systems*, pages 1921–1928, 2009.
- [22] Peng Zhao and Bin Yu. On model selection consistency of lasso. *Journal of Machine learning research*, 7(Nov):2541–2563, 2006.