

# Constrained-Based Sparse Identification of Cardiac Electrophysiology Models Using PySINDy

Mariana A S de Carvalho<sup>1</sup>, Rodrigo W dos Santos<sup>1</sup>, Bernardo M Rocha<sup>1</sup>

<sup>1</sup> Universidade Federal de Juiz de Fora, Juiz de Fora, Brasil

## Abstract

*Complex models in cardiac electrophysiology are commonly used to simulate cardiac action potential and cardiac dynamics. However, their high complexity can lead to significant computational costs and challenges in large-scale simulations. Therefore, it could be interesting to find simplified or alternative versions of cardiac cell models that preserve the essential characteristics of the original system, enabling efficient simulations. The identification of these models must strike a balance between accuracy and complexity, ensuring that critical aspects of electrophysiology are preserved while maintaining simplicity and reducing computational costs. In this work, we explore the constrained formulation in PySINDy to identify cardiac electrophysiology models as a proof of concept. The constrained formulation allows the imposition of physically based conditions on the identified model terms and coefficients, ensuring that certain principles or mathematical properties are preserved. In this work, we demonstrate that polynomial formulations of the gating equations in the classical Hodgkin-Huxley model can be derived using the methodology developed in this work.*

## 1. Introduction

The Hodgkin-Huxley (HH) model is an example of a biophysical model that describes how action potentials in neurons are initiated. It consists of a set of nonlinear ordinary differential equations (ODEs) that represent ionic currents across the neuronal membrane based on voltage-dependent conductance, and it is widely used to study the electrical behavior of excitable cells [1]. Today, the HH model serves as a foundational framework for numerous cardiac electrophysiology cell models.

Data-driven approaches can be used to discover governing equations directly from time series data. Among these, the Sparse Identification of Nonlinear Dynamical Systems (SINDy) [2] method can identify compact and interpretable models by assuming that the dynamics can be expressed as a sparse combination of candidate functions and coefficients. PySINDy is a Python library that imple-

ments this approach from input data and time derivatives and a function library, which then applies a sparse regression technique to discover an equivalent system by assuming the dynamics are governed by a sparse set of active terms [3].

The original SINDy formulation does not allow the user to incorporate prior knowledge about the system, so this can be problematic in cases where parts of the system's dynamics are already consolidated by physical laws, such as in biophysical models like Hodgkin-Huxley. To address these limitations in this work, a constrained version of the SINDy method was used, including prior structural or physical knowledge by fixing or constraining some terms during the regression process.

This work aims to apply the SINDy-constrained method using the SR3 optimization algorithm with  $\ell_2$ -norm regularization to recover the Hodgkin-Huxley model from noisy data. In our approach, we fix the equation for the transmembrane potential, and identify alternative forms for the gating variable equations. By doing so, we ensure that the identified model remains biophysically plausible while leveraging the flexibility of SINDy to discover the remaining components of the system. This represents an initial step toward the development of simplified and/or reduced, data-driven models for cardiac electrophysiology that retain interpretability and physiological relevance.

## 2. The Hodgkin-Huxley Model

The Hodgkin-Huxley model consists of a set of nonlinear differential equations that simulate the electrical behavior of neurons and cardiac muscle cells. The model is represented by the following system of four ordinary differential equations:

$$\dot{V} = \frac{1}{C_m} (I_{\text{ion}} - I_{\text{Na}} - I_{\text{K}} - I_{\text{L}} + I_{\text{app}}), \quad (1)$$

$$\dot{m} = \alpha_m(V)(1 - m) - \beta_m(V)m, \quad (2)$$

$$\dot{h} = \alpha_h(V)(1 - h) - \beta_h(V)h, \quad (3)$$

$$\dot{n} = \alpha_n(V)(1 - n) - \beta_n(V)n, \quad (4)$$

where  $V$  represents the transmembrane potential, and  $\alpha_k$  and  $\beta_k$  for  $k = \{m, n, h\}$  are the following functions of  $V$ :

$$\alpha_m = 0.1 \frac{25 - V}{\exp\left(\frac{25 - V}{10}\right) - 1}, \quad (5)$$

$$\beta_m = 4 \exp\left(\frac{-V}{18}\right), \quad (6)$$

$$\alpha_h = 0.07 \exp\left(\frac{-V}{20}\right), \quad (7)$$

$$\beta_h = \frac{1}{\exp\left(\frac{30 - V}{10}\right) + 1}, \quad (8)$$

$$\alpha_n = 0.01 \frac{10 - V}{\exp\left(\frac{10 - V}{10}\right) - 1}, \quad (9)$$

$$\beta_n = 0.125 \exp\left(\frac{-V}{80}\right). \quad (10)$$

In the model  $m$ ,  $h$ , and  $n$  are auxiliary variables for sodium ( $Na^+$ ), and potassium ( $K^+$ ) currents, and  $I_{app}$  denotes an externally applied current. In addition, the currents are defined as:

$$I_{Na} = g_{Na} m^3 h (V - E_{Na}), \quad (11)$$

$$I_K = g_K n^4 (V - E_K), \quad (12)$$

$$I_L = g_L (V - E_L), \quad (13)$$

which correspond to the sodium, potassium, and leak currents. The following parameters were used in this work:  $g_{Na} = 120.0 \text{ mS/cm}^2$ ,  $g_K = 36.0 \text{ mS/cm}^2$ ,  $g_L = 0.3 \text{ mS/cm}^2$ ,  $E_{Na} = 50.0 \text{ mV}$ ,  $E_K = -77.0 \text{ mV}$ ,  $E_L = -54.4 \text{ mV}$ , and  $I_{app} = 10.0 \text{ } \mu\text{A/cm}^2$ .

### 3. Methodology

Here, we present some fundamental concepts of the SINDy method and its computational implementation. In particular, we will focus on the constrained version of the algorithm, known as SINDy-constrained, which allows for the incorporation of prior knowledge and physical constraints into the model identification process. This approach will be applied to identify the dynamics of the Hodgkin-Huxley model.

#### 3.1. Sparse Identification of Nonlinear Dynamics

SINDy is a data-based sparse identification method introduced in [3] that combines sparse regression techniques with a library of candidate nonlinear functions to identify the underlying equations governing a dynamical system directly from data. As input, the method requires a set of time series data from the variables of interest over time,

and as output, it provides a set of sparse differential equations.

The method works by assuming that the temporal evolution of a system can be described by a sparse combination of candidate functions (such as polynomials or trigonometric functions). SINDy identifies which terms are most relevant to capturing the essential dynamics, based on the premise that many dynamical systems can be written as

$$\frac{d}{dt} \mathbf{x} = \mathbf{f}(\mathbf{x}) \quad (14)$$

have dynamics  $\mathbf{f}$  with only a few active terms and can be approximated by the following linear combination:

$$\mathbf{f}(\mathbf{x}) \approx \sum_{k=1}^p \theta_k(\mathbf{x}) \xi_k = \Theta(\mathbf{x}) \boldsymbol{\xi}, \quad (15)$$

where  $\Theta(\mathbf{x})$  is a library of candidate functions depending on the state  $\mathbf{x}$  and the input  $\mathbf{u}$ , and  $\boldsymbol{\xi}$  are the coefficients. The sparse regression method aims to identify a model with the fewest possible active terms in the set  $\Theta$ , that is, with only a few nonzero coefficients in  $\boldsymbol{\xi}$ .

To evaluate  $\Theta$ , a set of  $m$  measurements of the variables  $\mathbf{x}$  over time is taken and organized into matrices as follows:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}(t_1) & \mathbf{x}(t_2) & \cdots & \mathbf{x}(t_m) \end{bmatrix}^T, \quad (16)$$

In addition to these data, the method also requires the time derivatives of  $\mathbf{x}$ , which are organized in matrix form as:

$$\dot{\mathbf{X}} = \begin{bmatrix} \dot{\mathbf{x}}(t_1) & \dot{\mathbf{x}}(t_2) & \cdots & \dot{\mathbf{x}}(t_m) \end{bmatrix}. \quad (17)$$

These time derivatives can be obtained directly from the measurements, if available, or approximated using a numerical differentiation method. Thus, the problem in Equation (14), can be written in terms of data matrices as:

$$\dot{\mathbf{X}} \approx \Theta(\mathbf{X}) \boldsymbol{\Xi}, \quad (18)$$

where  $\boldsymbol{\Xi}$  is the coefficient matrix representing the model, and  $\Theta(\mathbf{X})$  is a matrix of candidate functions. This library can be constructed from the data in  $\mathbf{X}$ , as follows:

$$\Theta(\mathbf{X}) = \begin{bmatrix} \mathbf{1} & \mathbf{X} & \dots & \mathbf{X}^d & \dots & \sin(\mathbf{X}) \end{bmatrix}, \quad (19)$$

where  $\mathbf{X}^d$  denotes the matrix with column vectors composed of all available time series of polynomial degree  $d$ . The matrix  $\boldsymbol{\Xi}$  stores the coefficients that determine the magnitude of the candidate functions. That is, each column  $\xi_k$  of  $\boldsymbol{\Xi}$  is a vector of coefficients that determines the active terms in the  $k$ -th row of Equation (18).

To determine which model best fits the data, a sparse regression is performed using the SR3-constrained algorithm (Sparse Relaxed Regularized Regression) through the use of a  $\ell_2$  norm, aiming to minimize the difference between

the right and left-hand sides of Equation (18) [4]. For this, we solve the following optimization problem:

$$\min_{\Xi, \mathbf{W}} \frac{1}{2} \|(\dot{\mathbf{X}} - \Theta(\mathbf{X})\Xi)\|^2 + \lambda \mathbf{R}(\mathbf{W}) + \frac{1}{2\nu} \|\Xi - \mathbf{W}\|^2, \quad (20)$$

$$\text{subject to } \mathbf{C}\xi = \mathbf{d}, \quad \xi = \Xi(\cdot), \quad (21)$$

where  $\mathbf{C}$  is the constraint matrix and  $\mathbf{d}$  is a vector containing the prescribed values for the constraints, where  $\lambda$  is the parameter that promotes sparsity. The auxiliary variable  $\mathbf{W}$  is introduced to relax the constrained problem of estimating  $\xi$ , leading to a cost function that balances data fit, sparsity, and a penalty (scaled by  $\nu$ ) enforcing  $\mathbf{W} \approx \xi$  [5, 6].

### 3.2. Constraints

In this work, we employed the following types of constraints for model identification: one related to the  $V$  equation, which effectively fixes it to the original Hodgkin-Huxley model equation (see equation (1)), and another related to the gating variables equations  $m$ ,  $h$  and  $n$ . The constraints related to  $V$  were used to essentially remove the identification of the  $V$  equation from the process. The constraints for the gating variables were such that in each gating variable equation, only the corresponding gating variable (either  $m$ ,  $h$ , or  $n$ ) and  $V$  terms appear in the equation, replicating the structure from the original model as observed in equations (2)-(4). Also, we only allowed the gating variable up to first power, to reproduce the structure of the original formulation.

## 4. Results

The numerical experiments were performed using PySINDy with the SR3 optimizer configured to enforce equality constraints and applying an  $\ell_2$ -norm thresholding. Synthetic data were generated from the Hodgkin-Huxley model, and Gaussian noise was added to mimic experimental conditions, with different amplitudes applied to membrane potential and gating variables.

A noise amplitude of 1 mV was applied to the membrane potential  $V$ , and 0.01 to the gating variables. The time derivatives were calculated using high-order finite differences, while applying local smoothing to reduce noise in the data.

A polynomial library of degree five was used to construct candidate functions to match the original model, where the equation  $V$  is also a fifth-order equation. The constraints matrix was defined to restrict the active terms in each identified equation, ensuring that prior structural knowledge was incorporated directly into the model discovery process.

With respect to noise, it was observed in the experiments that, when no noise was provided, the method was not able

to identify any model capable of representing the original system. When the noise level was increased to 20%, the algorithm was able to discover models, but the responses exhibited a noticeable delay compared to the original system, leading to a time-shifted behavior. At 50% noise, the identification process did provide a model. However, the simulated response did not follow the curve of the original system, and the overall fit was unsatisfactory. Therefore, a noise level of 10% was adopted in the reported results.

To find a suitable model, the parameter  $\lambda$  was tested in a range of  $[0.0, 10.0]$ , shortening the subdivisions until the best model was obtained, which was  $\lambda = 0.505$ , for the regression problem defined by equation (20). With this choice for  $\lambda$ , the following model was obtained:

$$\dot{V} = -6.383 - 0.303V + 6060.000m^3h \quad (22)$$

$$- 2799.720n^4 - 121.200Vm^3h - 36.360Vn^4$$

$$\dot{m} = 3.858 + 0.099V - 3.965m + 0.001V^2 \quad (23)$$

$$- 0.092Vm - 0.001V^2m$$

$$\dot{h} = -0.002 + 0.001V - 1.013h - 0.010Vh \quad (24)$$

$$\dot{n} = 0.477 + 0.007V - 0.500n - 0.004Vn. \quad (25)$$

In the equations above, we report only the coefficients with magnitudes greater than  $10^{-4}$ , as smaller-magnitude terms were also present but could not be removed without compromising the stability of the identified model. This behavior warrants further investigation.

Figure 1 shows the numerical simulations of the action potential and gating variables from both the identified and original model (with noise), demonstrating a good agreement between the identified model by SINDy and the original HH model. We also conducted qualitative validations of key action potential properties (all-or-none response, refractory period, and sensitivity to initial conditions), confirming that the identified model provides an adequate representation of the original system. Figure 2 shows the response for different initial conditions.

## 5. Conclusions

In this work, we presented a sparse identification approach applied to the HH model for electrophysiology. By constraining some terms, the identification process becomes more stable and was able to discover a polynomial model for the gating variables. It also ensured that the resulting model remains consistent with the known structure of dynamics. The results showed that the SR3-constrained method was effective in accurately identifying the system.

This work was conducted as a proof of concept for applying SINDy's data-driven capabilities to electrophysiology models, and the results encourage further exploration toward model simplification, reduction, and discovery.

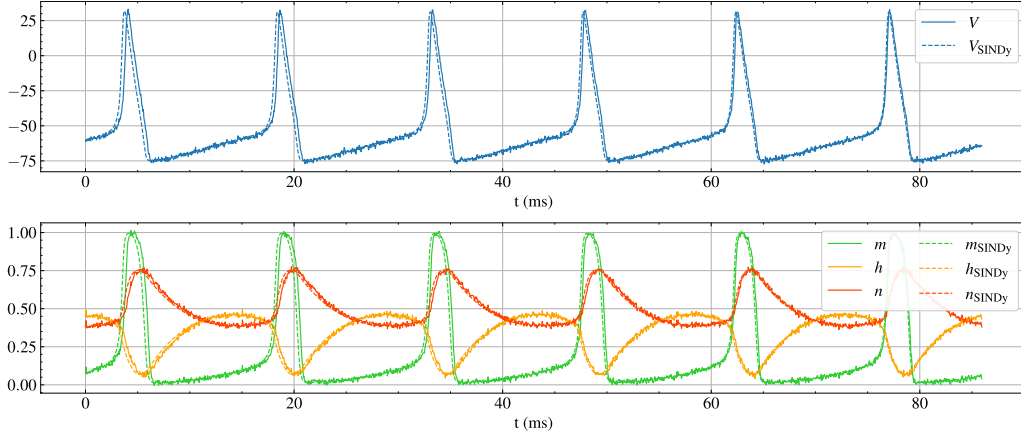


Figure 1. Simulations of the identified model (SINDy) (dashed) compared to the original HH data with noise (solid).

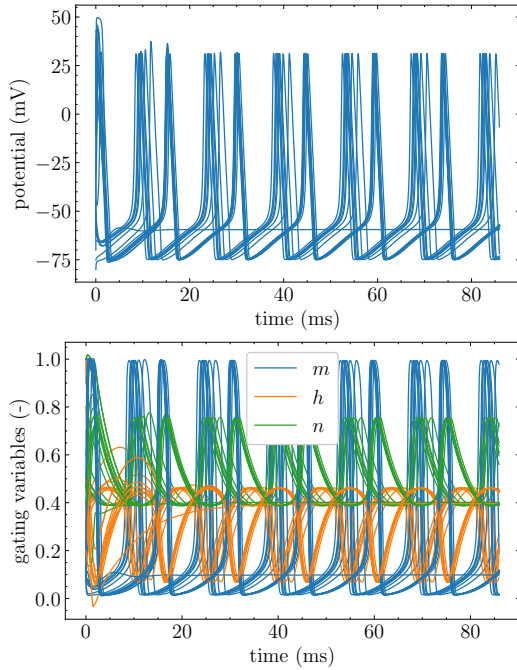


Figure 2. Response of the discovered SINDy model under different initial conditions.

## Acknowledgments

The authors would like to express their thanks to Wellcome Trust fellowship (214290/Z/18/Z), the EPSRC project CompBioMedX (EP/X019446/1), CompBioMed2 grant agreements No.675451 and No.823712, Minas Gerais State Research Support Foundation (FAPEMIG) - PCE-00048-25; APQ-02752-24, APQ-02445-24, APQ-02513-22, FINEP (SOS Equipamentos 2021 AV02 0062/22), "Conselho Nacional de Desenvolvimento Científico e Tecnológico"(CNPq) 423278/2021-5 and

310722/2021-7; "Coordenação de Aperfeiçoamento de Pessoal de Nível Superior" (CAPES), "Empresa Brasileira de Serviços Hospitalares" (Ebserh), SINAPAD Santos-Dumont, Federal University of Juiz de Fora (UFJF) for funding this work.

## References

- [1] Hodgkin AL, Huxley AF. A quantitative description of membrane current and its application to conduction and excitation in nerve. *The Journal of Physiology* 1952;117(4):500.
- [2] Champion K, Lusch B, Kutz JN, Brunton SL. Data-driven discovery of coordinates and governing equations. *Proceedings of the National Academy of Sciences* 2019; 116(45):22445–22451.
- [3] Brunton SL, Proctor JL, Kutz JN. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proceedings of The National Academy of Sciences* 2016;113(15):3932–3937.
- [4] Kaptanoglu AA, de Silva BM, Fasel U, Kaheman K, Goldschmidt AJ, Callahan JL, Delahunt CB, Nicolaou ZG, Champion K, Loiseau JC, et al. PySINDy: A comprehensive Python package for robust sparse system identification. *arXiv preprint arXiv:211108481* 2021;.
- [5] Brunton SL, Kutz JN. *Data-driven science and engineering: Machine learning, dynamical systems, and control*. Cambridge University Press, 2022.
- [6] Champion K, Zheng P, Aravkin AY, Brunton SL, Kutz JN. A unified sparse optimization framework to learn parsimonious physics-informed models from data. *IEEE Access* 2020; 8:169259–169271.

Address for correspondence:

Mariana Carvalho  
Programa de Pós-graduação em Modelagem Computacional/UFJF  
Rua José Lourenço Kelmer, Juiz de Fora, MG, CEP 36036-330  
marianaasdcavvalho@gmail.com