

## ABSTRACT

Biological networks are complex and constructing the dynamics for large systems can prove difficult. We will explore utilizing data-driven techniques to uncover the dynamics of a biological network from data provided. In particular, we are interested in inferring the dynamics for biological networks containing conservation laws or other special structures which can result in a singularity. These properties will be discussed in detail and their effect of recovering the dynamics will be investigated. Several numerical dynamics identification algorithms will be presented and their strengths and weaknesses will be discussed. In addition, the question of whether reconstructed networks exhibit adaptation properties based on recently developed analytical criteria will be investigated.

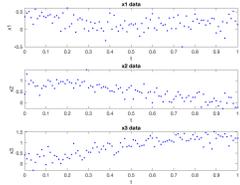
## METHODOLOGY

### We want to:

- Recover the dynamics from biological data

**Given:** Data from a biological network

**Goal:** Utilizing the data, recover a series of differential equations which accurately portray the data



$$x'_1 = f_1(x)$$

$$x'_2 = f_2(x)$$

$$x'_3 = f_3(x)$$

where  $x = [x_1, x_2, x_3]$  is the vector of proteins of interest.

- Recreate network structure from dynamics
- Using recovered dynamics, test for properties of interest such as:
  - Conservation Laws [1]
  - Special network structures, such as compounds in the network
  - Response to input/stimuli
  - Adaptation properties [1]

### Known problems:

- Rational functions are typically present in biological models (Michaelis-Menten kinetics, etc.)
- Systems can contain properties which could interfere with recovering dynamics
- May not have access to derivative information or it may not be accurate enough

## FUNDING

This research is part of the Industrial Immersion Program funded by the George Mason Provost. Both AWM and AMS have provided travel funding to present this work.

## INFERRING DYNAMICS

We begin by examining the structure of the Sparse Identification of Nonlinear Dynamics (SINDy) Method [2]:

Consider a dynamical system:

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

where we wish to approximate  $f(x)$  by a generalized linear model:

$$f_k(x) \approx \Theta(x)\xi_k$$

with the fewest nonzero terms in  $\xi_k$  as possible.

Let  $X = [x(t_1) \dots x(t_m)]^T$  be the time series data collected from the system. Then, we can develop a library of candidate nonlinear functions  $\Theta(X)$  constructed from the data in  $X$ :

$$\Theta(X) = [1 \quad X \quad X^2 \quad \dots \quad X^d \quad \dots \quad \sin(X) \quad \dots]$$

Now the dynamical system can be represented in the form:

$$\dot{X} = [\dot{x}(t_1) \dots \dot{x}(t_m)]^T = \Theta(X) [\xi_1 \dots \xi_k \dots]$$

We wish to minimize:

$$\xi_k = \arg \min \|\dot{X}_k - \Theta(X)\xi'_k\|_2 + \lambda \|\xi'_k\|_1$$

where  $\lambda$  weights the sparsity constraint.

**Alternative Framework:** For our purposes, we will use the ideas presented in the SINDy Method and apply them utilizing a different solver.

Consider the dynamical system

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

for which we would like to recover from data provided.

We will consider an approximation of  $f(x)$  using a generalized linear model (as in the SINDy Method):

$$f_k(x) \approx \begin{bmatrix} \Theta(x) \\ -\Theta(x) \end{bmatrix} \omega_k$$

where  $\Theta$  is a library of candidate nonlinear functions constructed from the data and  $\omega_k$  contains the fewest nonzero terms possible.

In order to find the optimal  $\xi_k$ , we will utilize the time series data in the **Non-negative Least Squares (NNLS)** algorithm [3] as follows:

$$\omega_k = \arg \min \left\| \begin{bmatrix} \Theta(X) \\ -\Theta(X) \end{bmatrix} \omega'_k - \dot{X}_k \right\|_2$$

where the top entries of  $\omega_k$  will correspond to positive coefficients in the recovered dynamics and the bottom entries are the negative.

## INFERRING DYNAMICS WITH RATIONAL FUNCTIONS

We wish to discover the dynamics of:

$$\frac{d}{dt}x_k(t) = \frac{f_{N,k}(x)}{f_{D,k}(x)}$$

from time series data of the state  $x(t) = [x_1(t), \dots, x_n(t)]^T$  and where  $f_{N,k}(x)$  and  $f_{D,k}(x)$  represent the numerator and denominator polynomials in the state variable  $x_k$ .

We will proceed as in Mangan et al. [4] by multiplying both sides by the denominator:

$$f_{N,k}(x) - f_{D,k}(x)\dot{x}_k = 0$$

As before, we wish to approximate  $f_{N,k}(x)$  and  $f_{D,k}(x)$  by generalized linear models:

$$f_{N,k} \approx \Theta_N(x)\omega_{N,k} \quad f_{D,k} \approx \Theta_D(x)\omega_{D,k}$$

where  $\Theta_N(x)$ ,  $\Theta_D(x)$  are the candidate function libraries and

$\omega_{N,k}$ ,  $\omega_{D,k}$  are the corresponding coefficients.

We will consider the library of candidate linear functions constructed from the data:

$$\Theta_N(X) = \Theta_D(X) = [1 \quad X \quad X^2 \quad \dots \quad X^d \quad \dots]$$

Thus, we have:

$$\Theta_N(x)\omega_{N,k} - \Theta_D(x)\omega_{D,k}\dot{x}_k = 0$$

In order to apply NNLS, we will assume the coefficient associated with the  $x_k\dot{x}_k$  term is 1 and thus solve:

$$x_k\dot{x}_k = \Theta_N(x)\omega_{N,k} - \tilde{\Theta}_D(x)\omega_{D,k}\dot{x}_k$$

where  $\tilde{\Theta}_D$  is the  $\Theta_D$  matrix with the column corresponding to  $x_k$  removed. The optimal  $\omega_{N,k}$  and  $\omega_{D,k}$  can now be found using a similar method as before.

## RECOVER DYNAMICS

$$\dot{x}_1 = k_3x_1x_3 - k_1x_1x_2$$

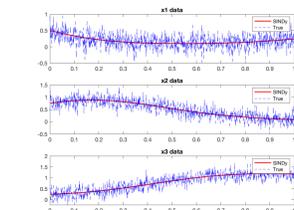
$$\dot{x}_2 = k_1x_1x_2 - k_2x_2x_3$$

$$\dot{x}_3 = k_2x_2x_3 - k_3x_1x_3$$

$$C_{tot} = x_1 + x_2 + x_3$$

$$k_1 = 6, k_2 = 4, k_3 = 3$$

$$x_0 = (.5, .75, .25)$$

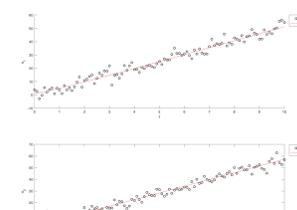


The noisy data and the solution found using SINDy. Similar results were obtained using NNLS.

## RECOVER RATIONAL FUNCTIONS

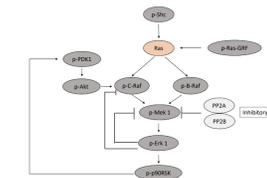
$$\frac{dx_1}{dt} = \frac{2x_1 + 3x_2}{1 + x_1}$$

$$\frac{dx_2}{dt} = \frac{2 + x_1 + 5x_2}{1 + x_2}$$



The noisy data and the solution found using NNLS to recover dynamics for systems with rational functions.

## THE MAPK PATHWAY



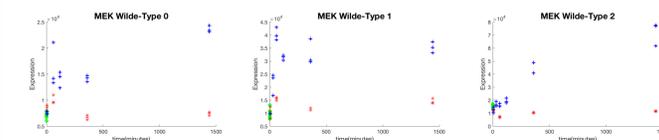
### Properties:

- Plays a role in regulation of gene expression and cellular growth
- Dysregulated MAPK signaling is implicated in a wide range of cancers
- Ras protein is prone to mutation
- Data for the proteins in this pathway is a measure of the protein activation

Above is a diagram for the MAPK Pathway provided by Dr. Pierobon.

### Data Properties

- Data from 8 different cell lines over a period of 24 hours with two different solutions applied, Dimethyl sulfoxide (DMSO) and the MEK Inhibitor Selumetinib.
- 3 different states for the KRAS protein:
  - Wilde-Type = 0: 0 normal copies of gene (mutated)
  - Wilde-Type = 1: 1 normal copies of gene (mutated)
  - Wilde-Type = 2: 2 normal copies of gene (unmutated)



MEK protein expression over time where  $\circ$  represent baseline value,  $*$  are DMSO Control, and  $+$  are Selumetinib (MEK Inhibitor).

## FUTURE RESEARCH

There are several avenues for future research:

- Recover a model for the MAPK pathway
  - Model differences between DMSO application and MEK inhibitor application
  - Model differences for mutated cell lines
  - For all models, identify any interesting properties of the dynamics, such as conservation laws, properties of adaptation, etc.
- Continue to develop the theory behind the recovery process
- What are the limitations for recovery?

## REFERENCES

- T. Oellerich, M. Emelianenko, L. A. Liotta, and R. P. Araujo. Biological networks with singular jacobians: their origins and adaptation criteria. *Submitted*.
- S. L. Brunton, J. L. Proctor, and J. N. 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.
- Charles L. Lawson and Richard J. Hanson. *Solving Least Squares Problems*. SIAM, 1995.
- N.M. Mangan, S.L. Brunton, J.L. Proctor, and J.N. Kutz. Inferring biological networks by sparse identification of nonlinear dynamics. *IEEE Transactions on Molecular, Biological and Multi-Scale Communications*, 2, 2016.