Playing Newton

Ilya Nemenman

Departments of Physics and Biology Emory University

JAMES S. MCDONNELL FOUNDATION







Thanks

- Bryan Daniels (Arizona State)
- Andre Levchenko (Yale)
- Will Ryu (Toronto)





Why? (paraphrasing Richard Hamming)

- 1. What are the important problems in your field?
 - 1. What are the important problems in your field?
 - 2. What important problems are you working on?
 - 3. Why are the answers to (1) and (2) different?

So:

What are the important problems in theoretical biophysics?



Of exactitude in science

...In that Empire, the craft of Cartography attained such Perfection that the Map of a Single province covered the space of an entire City, and the Map of the Empire itself an entire Province. In the course of Time, these Extensive maps were found somehow wanting, and so the College of Cartographers evolved a Map of the Empire that was of the same Scale as the Empire and that coincided with it point for point. Less attentive to the Study of Cartography, succeeding Generations came to judge a map of such Magnitude cumbersome, and, not without Irreverence, they abandoned it to the Rigours of sun and Rain. In the western Deserts, tattered Fragments of the Map are still to be found, Sheltering an occasional Beast or beggar; in the whole Nation, no other relic is left of the Discipline of Geography.

> From Travels of Praiseworthy Men (1658) by J. A. Suarez Miranda (a fictional reference). By Jorge Luis Borges and Adolfo Bioy Casares. English translation quoted from J. L. Borges, A Universal History of Infamy, Penguin Books, London, 1975.



At a recent meeting...

• A common opinion: "The final theory of biological systems will be a large multiscale computational model. We need more and more experimental data to specify details of these models."

openworm

A simulation platform to build digital in silico living systems -- starting with a c. elegans worm virtual organism simulation

- There's something wrong with this statement.
 - The "final" theory?
 - Do we need the theory of
 "everything" in any biological (or physical) system?
- The best material model of a cat is another, or preferably the same, cat. (*Philosophy of Science*, Wiener and Rosenblueth, 1945)





Physics analogy

- What is the final, complete theory of the chair you are sitting in?
 - How does it fall from the second floor?
 - How does the cloth seat age and tear?
 - How much weight would the chair hold before it breaks?
 - How does it conduct electricity?
 - How much food can I cook when I burn it?

- ...

- There's no such thing as "the full theory of the chair".
 - We build models tailored to answer **specific questions**.
 - The complete theory that answers every question would need to include quarks, superstrings...
 - Each modeling level needs its own effective degrees of freedom
 - "Don't model bulldozers with quarks." (Goldenfeld and Kadanoff, 1999)
- Models must loose details. Otherwise, just use the same cat...



So...

- Are there phenomenological, coarse-grained, and yet functionally accurate representations of (some) biological dynamics, or are we forever doomed to every detail mattering?
 - And, of course, these models would not answer *every* question, but specific questions on coarse scales.
 - E.g., not What is a position of this particular atom in the cell? But What is the whole system doing?
 - (Parenthetically): Without such coarse-grained descriptions, if everything is equally important, modern biology ceases to be a Western science.



What do I mean by this?

- Western tradition:
 - There are laws (of nature, of god, whatever). A rock is a rock everywhere. It falls the same in Pisa and in Atlanta.
 - There're causes and there are effects.
 - There is "useless information" (Oscar Wilde).
 - But this belief requires closing one's eyes to minor discrepancies
 - Two balls dropped from the Leaning Tower didn't actually land simultaneously.
 - "If we had the STM in the 1920s, there wouldn't be the Debye theory of solids." (H. Levine)
- Non-western tradition, e.g., buddhism
 - Pratityasamutpada: dependent origination: "Pratitya samutpada is sometimes called the teaching of cause and effect, but that can be misleading, because we usually think of cause and effect as separate entities, with cause always preceding effect, and one cause leading to one effect. According to the teaching of Interdependent Co-Arising, cause and effect co-arise (samutpada) and everything is a result of multiple causes and conditions... " Thich Nhat Hanh



Where would new laws come from?

- More is different! (PW Anderson)
 - The law of large numbers produces universalities if the right questions are asked (e.g., about *large-scale* quantities).
- Important to use the right level of description.
 - Coarse-graining: Each modeling level needs its own effective degrees of freedom.
 - "Don't model bulldozers with quarks." (Goldenfeld and Kadanoff, 1999)
- This is already common in your every-day life, not just physics
 - Which level of description is better for driving to a local school?





A side note: This has worked before! Hodgkin-Huxley



A good theory! (and only roughly correct)



Cellular networks: complex beasts Large *N* is baked in

A culture's icons are a window onto its soul. Few would disagree that, in the culture of molecular biology that dominated much of the life sciences for the last third of the 20th century, the dominant icon was the double helix. In the present, post-modern, 'systems biology' era, however, it is, arguably, the hairball.

A.D. Lander. BMC Biology 2010, 8:40







And yet, for typical inputs, their dynamics are rather simple



- A handful of parameters (time scales, amplitudes) describe responses of networks to most experimentally accessible perturbations.
- Relation of phenomenological to mechanistic parameters often unclear.
- Do we need complex networks to describe simple dynamics?



The field (and we) has worked on this for a while

- Bottom-up methods, reducing a known microscopic, mechanistic network.
 - with Sinitsyn et al., 2006-2010; with Munsky, Bel et al., 2009-2013; with Merchan, 2016.
 - Schwab and Mehta 2015
 - Large number of publications in chemical physics (Petzold et al.)
 - **Problem:** need to know and start with the microscopic model.
- Can we instead build phenomenological models topdown, from data directly, and without reconstructing a mechanistic network as an intermediate step?
 - Purpose: **predict responses** to exogenous signals.
 - Purpose: drive all of us, theorists, out of work?



First steps...

- Assume that dynamics of cellular networks is given by local ordinary differential equations.
 - Do not fit curves; fit dynamics.
- Neglect stochasticity, and spatial structure for now

$$\begin{cases} \frac{dx_1}{dt} = f_1(x_1, x_2, \dots, x_n) \\ \dots \\ \frac{dx_n}{dt} = f_n(x_1, x_2, \dots, x_n) \end{cases}$$

- Can we automatically fit these functions f_i from data?
 - How do we enumerate the set of all possible multivariate functions?
 - How do we search through this list?
 - How do we not overfit?



Prior (and posterior) art

- The full search approach for an exact model
 - Small systems dynamics search for all possible models using S-systems formalism (Voit et al, Theor Biol Med Model 2006).
 - Searching for a control model from a (small) set of a priori allowed models (Lillacci and Khammash, PLoS CB 2010).
 - Searching for a stochastic model from a (small) set of a priori allowed models (Munsky, et al., MSB 2009, Science 2013).
 - Eureqa: exhaustive genetic algorithm search through all possible elementary function combinations, with selection of new experiments to optimize discriminability among models (Lipson et al., Science 2009, Phys Biol 2011).
 - SINDy: Compressed sensing approaches (Brunton et al., 2016).
 - AutomatedStatistician.com: compositional tree structures for stochastic data.
 - More at the 2017 MM meeting sessions on "Robot-scientist".
- Phenomenological search (Crutchfield and McNamara, Compl Syst 1987).
- Problems (limiting the analysis to only a few variables)
 - data/computing demands explode with the number of variables;
 - cannot handle unobserved variables.



Test Model: Yeast Glycolytic Oscillator



- 7 species, 28 parameters
- Complex rational dynamical laws





Eureqa performance: Yeast Glycolytic Oscillator



• 7 species, 28 variables

Automatically inferred system

Complex rational dynamical laws

Amazing accuracy!

Original system

$$\begin{array}{ll} \frac{dS_1}{dt} = 2.5 - \frac{100*A_3S_1}{1+13.68*A_3^4} \\ \frac{dS_2}{dt} = \frac{200*A_3S_1}{1+13.68*A_3^4} - 6*S_2 - 6*S_2N_2 \\ \frac{dS_3}{dt} = 6*S_2 - 6*N_2S_2 - 64*S_3 + 16*A_3S_3 \\ \frac{dS_4}{dt} = 64*S_3 - 16*A_3S_3 - 13*S_4 - 100*N_2S_4 \\ +13*S_5 \\ \frac{dN_2}{dt} = 6*S_2 - 18*N_2S_2 - 100*N_2S_4 \\ \frac{dA_3}{dt} = -1.28*A_3 - \frac{200*A_3S_1}{1+13.68*A_3^4} + 128*S_3 + 32*A_3S_3 \\ \frac{dS_5}{dt} = 1.3*S_4 - 3.1*S_5 \end{array}$$

$$\begin{array}{l} \frac{dS_1}{dt} = 2.53 - \frac{98.79\cdotA_3S_1}{1+12.66\cdotA_3^4} \\ \frac{dS_2}{dt} = 200.23\cdotA_3S_1 - 6.87\cdot S_2 \\ \frac{dS_3}{dt} = 6.00\cdot S_2 - 6.00\cdot N_2 \\ \frac{dS_4}{dt} = 64.04\cdot S_3 - 16.03 + 100*N_2S_4 \\ \frac{dA_3}{dt} = -1.28*A_3 - \frac{200*A_3S_1}{1+13.68*A_3^4} + 128*S_3 + 32*A_3S_3 \\ \frac{dS_5}{dt} = 1.3*S_4 - 3.1*S_5 \end{array}$$

$$\frac{dS_1}{dt} = 2.53 - \frac{98.79 \cdot A_3 S_1}{1+12.66 \cdot A_3^4} \frac{dS_2}{dt} = \frac{200.23 \cdot A_3 S_1}{1+13.80 \cdot A_3^4} - 6.87 \cdot S_2 - 6.87 \cdot N_2 + 0.95 \frac{dS_3}{dt} = 6.00 \cdot S_2 - 6.00 \cdot N_2 S_2 - 64.16 \cdot S_3 + 16.08 \cdot A_3 S_3 \frac{dS_4}{dt} = 64.04 \cdot S_3 - 16.03 \cdot A_3 S_3 - 13.03 \cdot S_4 - 100.11 \cdot N_2 S_4 + 15.21 \cdot S_5 \frac{dN_2}{dt} = -0.051 + 5.99 \cdot S_2 - 17.94 \cdot N_2 S_2 - 98.82 \cdot N_2 S_4 \frac{dA_3}{dt} = -1.12 \cdot A_3 - \frac{192.24 \cdot A_3 S_1}{1+12.50 \cdot A_3^4} + 124.92 \cdot S_3 + 31.69 \cdot A_3 S_3$$

Schmidt et al., Phys Biol 2011



But at the same time: Need derivatives, and...



- Astronomical computation times -- exhaustive search.
 - **Overfitting** -- need astronomical sample sizes.
- Two exponential costs: **selecting** the best model family, **fitting** the best family with the model.





SINDy (Brunton et al., 2016) (But needs derivatives and correct basis)





Can we avoid the exhaustive search and the need for the correct basis?

• Do not need an exhaustive search or exact fits when fitting dimensional curves with progressively increasing complexity



- Use **nested**, **complete** model families, e.g., Taylor series.
- Use Bayesian model selection to limit the complexity of the search space (the value of maximum *K*).

Schwartz, Ann Stat 1978; MacKay, Neural Comp, 1992 Balasubramanian, Neural Comp1996; Nemenman, Neural Comp, 2005



Bayesian Model selection

• For large sample size *N*, averages done in the Laplace (saddle point) limit.

$$P(K|\{x_i\}) = \int d^K \vec{\alpha} P(\vec{\alpha}|\{x_i\}) = \int d^K \vec{\alpha} \frac{P(\{x_i\}|\vec{\alpha})\mathcal{P}(\alpha)}{P(\{x_i\})}$$
$$= \int d^k \vec{\alpha} \exp(-N\mathcal{L})$$

 $\log P(K|\{x_i\}) = \log P(\{x_i\}|\vec{\alpha}_{\mathrm{ML}}) - \frac{1}{2}\log \det N\mathcal{F} + O(N^0)$

- Penalty for model complexity (the log term) "selects" the best model family.
- Not that simple in detail, but this description is roughly accurate.
- Beautiful consistency properties for nested, complete model families.
 MacKay 1992, Balasubramanian 1996,

EMOR UNIVERSIT Nemenman 2005

Why is fitting dynamics so hard?



- Hidden degrees of freedom and nonlinearities breaks nestedness -- no consistency.
- Choose any (reasonable) **complete** path through the model space
 - Good choice good fits with few data; Bad choice not worse than exhaustive search.



Two types of model families

- Both nested and complete.
- Account for nonlinearities and hidden variables.



Daniels and Nemenman, Nature Comm 2015; PLoS ONE 2015



Algorithm

- Specify a particular hierarchy of model families.
- For given data:
 - Choose a model family within the hierarchy.
 - Fit for the best model within the family.
 - Calculate the posterior likelihood of the family using modified Bayesian criterion.
 - Choose more complex family and terminate when the modified likelihood starts to decrease.
- Algorithmic improvements to ensure that no complete re-fitting is done when move to the next family, or increase data set size.
- Two exponential complexities: search of a model family, and fitting a model within the family.
 - This only solves the **first**.
 - In practice works OK for both.



Finding laws that we already know: An automated Sir Isaac (Sirlsaac on GitHub)



Daniels and Nemenman, Nature Comm 2015, PLoS ONE 2015



Simple dynamics from a complex network: Combinatorial multisite phosphorylation

• Effective models fit better than the true model for finite data





- Rates depend on occupancy of the nearby sites, about 50 parameters total.
- Caricature of some of the most combinatorially complex signaling models.
- Typically more parameters than data.

2.0

1.6

1.8

Effective, reduced model of multi-site phosphorylation





The yeast glycolytic oscillations: Complex dynamics needing complex structure

- Observe only 3/7 of variables; add 10%
- Data: N samples of structure
 - Initial condition of the 3 species;
 - Some random time later;
 - The value of these 3 species at that time.





Results

2.6

 $S_1 \ 1.3$

0.0

0.0

 $\mathbf{2}$

0

0.0

0 1 $\mathbf{2}$ 3 4

 A_3

 S_4^{ex} 0.1

Condition 1 ...

Condition N

-300-350

-400

-450

1.0

0.8

0.4

0.2

0.0

0

8 $\mathbf{24}$

0

1 $\mathbf{2}$ 3

Computational effort

 $(\times 10^8 \text{ model evals})$

1.0

0.8

0.6

0.4

0.20.0

20

4

Out-of-samp

Mean out-of-sample correlation

 $0\ 1\ 2\ 3\ 4\ 5$

Time (minutes)

correlation 0.6

Daniels and Nemenman, Nature Comm 2015; PLoS ONE 2015

- ~100x fewer evaluations for the same accuracy compared to full search.
 - ~1000x fewer data points than full search.
 - Better accuracy than curve fitting.
- Linear scaling with the amount of data and with the number of variables.

 $2\ 3\ 4$

Time (minutes)

 $0\ 1\ 2\ 3\ 4\ 5$

Condition 1 ...

2.6

0.0

0.0

0.0

1 -4

0

-4

-2

0 1

 $X_{6} - 1$

 $S_1 \,\, 1.3$

 $S_2 \,\, 1.3$

 $S_3 \,\, 0.2$

 X_4

 X_5

40 60

40

Number of parameters

Num, measurements N

Condition M

Learning new science: Ca++ oscillations in beta cells

- Beta cells: cells in the pancreatic islets. Primary function: to store and release insulin.
- Insulin is released in pulses in synchrony with pulsatile activity of Ca++ spikes.
- Incorrect frequency / irregular pulsing of Ca++ / insulin measured in the portal vein, which connects pancreas to liver, prevents the liver



from absorbing glucose from blood and is a precursor of type 2 diabetes.

 Protein Kinase A (PKA), cAMP, and Ca++ are a part of a synchronized oscillatory circuit.



Calcium-PKA oscillatory dynamics in beta cells



- Doesn't account for decreasing amplitude of oscillations.
- Different model for every cell.

Ni et al, NCB, 2010



Calcium-PKA oscillatory dynamics in beta cells: automatic inference

Daniels, IN, Levchenko, in prep.



Fits every cell with <10 parameters individualized..



Modeling *C. elegans* temperature nociception escape response



among behavioral ral states?



Properties of the escape behavior







Worm nociception: Data and fits



Leung, Mohhamadi, Ryu, IN, 2016



Worm nociception: phase space



Daniels, Ryu, IN, in prep.



Worm nociception: phase space



 Settle the debate: Do stimuli bias transition probabilities among behavioral states, or do they create new behavioral states? In some sense, both.



Daniels, Ryu, IN, in prep.



Summary

- Focus on refining phenomenological dynamics.
- **Complete, nested** model families of dynamics allow to use Bayesian model selection to adapt model complexity to the available data.
- Such phenomenological models make accurate predictions in the undersampled regime, where true models overfit.
- Why do this?
 - The duck test: If it looks like a duck, swims like a duck, and quacks like a duck, then it probably is a duck (Nemenman, Physics Today 2015).
 - Find new phenomenological laws of nature
 - Repeat Hookean approach in biology: build effective models of similar systems and look for patterns (e.g., chemotaxis in *C. elegans* and *E. coli*).



