# Parameter-free statistical model invalidation for biochemical reaction networks

Kenneth L. Ho

Courant Institute, New York University

NYU COB Colloquium (2012 Oct 23)

## Driving problem

Given observed data and multiple candidate models for the process generating that data, which is the most appropriate model for that process?

Example (multisite phosphorylation):

$$K + S_0 \xrightarrow[\kappa_{02}]{\kappa_{02}} S_1 \qquad K + S_1 \xrightarrow[\kappa_{12}]{\kappa_{12}} S_2$$

If  $\kappa_{02} = 0$ , the kinase is distributive; otherwise, processive.



Aoki et al. (2011) PNAS

Basic setting:

- Data:  $x_1, \ldots, x_m$
- Models:  $f_1, \ldots, f_n$

Typically, each model also depends on some parameters, e.g.,  $f = f(\theta)$ .

Standard approaches:

- ▶ If the parameters are known, just simulate and compare, e.g., compute  $\hat{x}_i = f(\theta)$  and check  $\varepsilon = ||x \hat{x}||$ .
- If the parameters are unknown:
  - Optimize parameters and check best-case error (parameter estimation).
  - Average over parameters according to priors (Bayesian):

$$\Pr(f|x) \propto \int \Pr(x|f(\theta)) \Pr(f(\theta)) d\theta$$

For unknown parameters, most methods essentially involve some form of optimization or exploration over the parameter space.

In biology, the true parameters are hardly ever known (inability to measure, uncertainties, etc.). Thus, some type of parameter optimization is often required.

But this optimization can be very difficult:

- Nonlinearity of objective function
- High dimensionality of parameter space



#### Can we get by without parameter optimization?

In biology, the true parameters are hardly ever known (inability to measure, uncertainties, etc.). Thus, some type of parameter optimization is often required.

But this optimization can be very difficult:

- Nonlinearity of objective function
- High dimensionality of parameter space



#### Can we get by without parameter optimization?

In this talk, we present a statistical model invalidation technique that does not depend on any parameters, i.e., only on the model structure/topology.

Philosophically related:

- Chemical reaction network theory (Jackson, Horn, Feinberg)
- Stoichiometric network analysis (Clarke)
- Flux balance analysis (Palsson, Edwards)

#### Chemical reaction networks

To be concrete, consider a chemical reaction network

$$\sum_{j=1}^{N} r_{ij} X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij} X_j, \quad i = 1, \dots, R$$

with mass-action kinetics

$$\dot{x}_j = \sum_{i=1}^R \kappa_i \left( p_{ij} - r_{ij} \right) x^{r_i}, \quad j = 1, \dots, N,$$

where  $x^{r_i} = x_1^{r_{i1}} \cdots x_N^{r_{iN}}$ .

#### Example

$$\begin{array}{ll} X+Y \xrightarrow{\kappa_1} 2X & \dot{x} = \kappa_1 xy - 3\kappa_2 x^3 + \kappa_3 z \\ 3X \xrightarrow{\kappa_2} Y + 2Z & \dot{y} = -\kappa_1 xy + \kappa_2 x^3 + \kappa_4 z \\ Z \xrightarrow{\kappa_3} X , \ Z \xrightarrow{\kappa_4} Y & \dot{z} = 2\kappa_2 x^3 - (\kappa_3 + \kappa_4) z \end{array}$$

## Model compatibility

The dynamics  $\dot{x}$  provide a quantitative description of the model and can, in principle, be used to test its compatibility with observed data.

In practice, however, not all variables can be measured:

- Velocities are often difficult, so consider only the steady state  $\dot{x} = 0$
- Experimentally inaccessible species must be eliminated

Back to our example:

$$\kappa_1 xy - 3\kappa_2 x^3 + \kappa_3 z = 0$$
  
$$-\kappa_1 xy + \kappa_2 x^3 + \kappa_4 z = 0$$
  
$$2\kappa_2 x^3 - (\kappa_3 + \kappa_4) z = 0$$

If Z cannot be measured, eliminate  $z = \frac{2\kappa_2 x^3}{\kappa_3 + \kappa_4}$ . Compatibility conditions:

$$\kappa_1 xy - 3\kappa_2 x^3 + \left(\frac{2\kappa_2\kappa_3}{\kappa_3 + \kappa_4}\right) x = 0$$
$$-\kappa_1 xy + \kappa_2 x^3 + \left(\frac{2\kappa_2\kappa_4}{\kappa_3 + \kappa_4}\right) x = 0$$

How to test compatibility without knowing parameters in advance?

Easy case: suppose all variables can be measured. Then the compatibility conditions  $f(x; \kappa) = 0$  are linear in  $\kappa$ , e.g.,

$$\begin{bmatrix} xy & -3x^3 & z \\ -xy & x^3 & z \\ & 2x^3 & -z & -z \end{bmatrix} \begin{bmatrix} \kappa_1 \\ \kappa_2 \\ \kappa_3 \\ \kappa_4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

How to test compatibility without knowing parameters in advance?

Easy case: suppose all variables can be measured. Then the compatibility conditions  $f(x; \kappa) = 0$  are linear in  $\kappa$ , e.g.,

$\begin{bmatrix} x_1y_1 \end{bmatrix}$	$-3x_{1}^{3}$	$z_1$				
$-x_1y_1$	$x_{1}^{3}$	<i>z</i> <sub>1</sub>		[ [m]	[0]	
	$2x_1^3$	$-z_1$	$-z_1$	$\begin{bmatrix} \kappa_1 \\ \kappa_2 \end{bmatrix}_{-}$	0	
x <sub>2</sub> y <sub>2</sub>	$-3x_{2}^{3}$	<i>z</i> <sub>2</sub>		$ \kappa_3  =$	0	•
$-x_2y_2$	$x_{2}^{3}$	<i>z</i> <sub>2</sub>		$\lfloor \kappa_4 \rfloor$		
	$2x_2^3$	- <i>z</i> <sub>2</sub>	$-z_2$			

How to test compatibility without knowing parameters in advance?

Easy case: suppose all variables can be measured. Then the compatibility conditions  $f(x; \kappa) = 0$  are linear in  $\kappa$ , e.g.,

$$\begin{bmatrix} x_1y_1 & -3x_1^3 & z_1 \\ -x_1y_1 & x_1^3 & z_1 \\ \vdots & 2x_1^3 & -z_1 & -z_1 \\ \vdots & x_2y_2 & -3x_2^3 & z_2 \\ -x_2y_2 & x_2^3 & z_2 \\ \vdots & 2x_2^3 & -z_2 & -z_2 \end{bmatrix} \begin{bmatrix} \kappa_1 \\ \kappa_2 \\ \kappa_3 \\ \kappa_4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Compatibility with the data requires that there exist  $\kappa$  (possibly with restrictions) satisfying this matrix equation.

In general, there is a nonlinear matrix function  $\Phi$  such that the compatibility conditions can be expressed as  $\Phi(X) \cdot \kappa = 0$ , where X is the data matrix.

- Compatible if and only if κ ∈ null(Φ(X))
- Necessary condition: dim(null(Φ(X))) > 0
- Compute SVD and check smallest singular value
  - Reject model if σ<sub>min</sub>(Φ(X)) > 0
  - Criterion depends only on data and so is parameter-free

Geometric perspective:

- Rows of Φ(X), considered as points, are coplanar
- $\sigma_{\min}$  quantifies the deviation from coplanarity

Necessary but not sufficient:

- Cannot demonstrate correctness
- Contrasts with parameter optimization



(2008) Biophys J

How close to zero is close enough?

- ▶ Assume i.i.d. Gaussian noise in X
- Estimate propagated noise in  $\Phi(X)$ 
  - $\Phi$  depends only on model topology and hence is known
  - First-order expansion in noise magnitude using  $\nabla \Phi$

• Rescale rows of  $\Phi(X)$  so that all noise components have variance  $\leq 1$ :

$$\sigma_{\min} \leq \|z\|, \quad z_i \sim \mathcal{N}(0, \mu_i), \quad |\mu_i| \leq 1$$

Compare with chi distribution for p-value

This provides a means to reject the null hypothesis that  $\Phi(X)$  is coplanar.

Other bounds are also possible:

Courant-Fisher-Weyl, Wielandt-Hoffman, concentration of measure, etc.

In the general case, we must eliminate all variables that cannot be measured.

- How to do this systematically?
- Elimination also destroys linearity in  $\kappa$ 
  - Recall example:

$$\kappa_1 xy - 3\kappa_2 x^3 + \left(\frac{2\kappa_2\kappa_3}{\kappa_3 + \kappa_4}\right) x = 0$$
$$-\kappa_1 xy + \kappa_2 x^3 + \left(\frac{2\kappa_2\kappa_4}{\kappa_3 + \kappa_4}\right) x = 0$$

• How to linearize?

Our solution:

- Algebraic geometry, Gröbner bases
- Lifting procedure by relaxing nonlinearities



$$f_i(x) = \sum_{j=1}^p a_{ij} x^{b_{ij}} = 0, \quad i = 1, ..., n$$
  
 $x = (x_1, ..., x_N)$ 



- Algebraic geometry: zeros of polynomial equations
- Algebraic variety:  $V = \{x \mid f_1(x) = \cdots = f_n(x) = 0\}$
- Gröbner bases: Gaussian elimination for multivariate polynomial systems
- ▶ Form polynomial ring  $\mathbb{Q}[a] = \left\{ \sum_i c_i a^{b_i} \mid c_i \in \mathbb{Q} \right\}$  and let

$$\mathbb{K} = \mathsf{Frac}(\mathbb{Q}[a]) = \left\{ rac{p}{q} \mid p, q 
eq \mathsf{0} \in \mathbb{Q}[a] 
ight\}$$

- Allows symbolic computation over a
- Construct ideal  $I = \langle f_1, \ldots, f_n \rangle = \{ \sum_i f_i h_i \mid h_i \in \mathbb{K}[x] \}$ 
  - Contains all elements of  $\mathbb{K}[x]$  that vanish on V

$$f_i(x) = \sum_{j=1}^p a_{ij} x^{b_{ij}} = 0, \quad i = 1, ..., n$$
  
 $x = (x_1, ..., x_N)$ 



To eliminate  $x_1, \ldots, x_k$ , consider the elimination ideal  $I_k = I \cap \mathbb{K}[x_{k+1}, \ldots, x_N]$ .

# Elimination property

If  $g = (g_1, \ldots, g_m)$  is a Gröbner basis for I over  $\mathbb{K}$  under the lexicographic ordering  $x_1 > \cdots > x_N$ , then  $I_k = \langle g_k \rangle$  for  $g_k = g \cap \mathbb{K}[x_{k+1}, \ldots, x_N]$ .

The basis polynomials  $g_k$  all vanish on V and depend only on  $x_{k+1}, \ldots, x_N$ . We call the elements  $\sum_i h_i(a) x^{p_i}$  of  $g_k$ , where  $h_i(a) \in \mathbb{K}$ , steady-state invariants.

- Properties can depend delicately on the monomial ordering
- Can be computed using standard computer algebra packages
- No reasonable bounds on computing time or storage

# Algorithm

Given dynamics  $\dot{x}$ , identified observables  $x_{obs}$ , and steady-state data  $\hat{x}_{obs,1}, \ldots, \hat{x}_{obs,m}$ :

- Compute steady-state invariants in x<sub>obs</sub> using Gröbner bases
- For each invariant  $\sum_{i=1}^{n} a_i(\kappa) x_{obs}^{p_i}$ :
  - Linearize by writing as  $\sum_{i=1}^{n} b_i y_i$ , where  $b_i = a_i(\kappa)$  and  $y_i = \varphi_i(x_{obs}) = x_{obs}^{p_i}$
  - Test coplanarity of  $Y = \Phi(\hat{X}_{obs}) \in \mathbb{R}^{m \times n}$  with respect to effective parameters b
- Reject model if any invariant does not induce a transformation to coplanarity



lifted

Many caveats remain (enlarged lifted space, Gröbner basis issues, b = 0, etc.); however, the method is still "surprisingly" effective.

#### Example: two-site phosphorylation



Kinase/phosphatase: distributive/processive Four models: PP, PD, DP, DD

- ► Variable ordering: (ks<sub>00</sub>, ks<sub>01</sub>, ks<sub>10</sub>, fs<sub>01</sub>, fs<sub>10</sub>, fs<sub>11</sub>, k, f, s<sub>00</sub>, s<sub>01</sub>, s<sub>10</sub>, s<sub>11</sub>)
- Only kinase mechanism is discriminative
- Can reject DP/DD models on the basis of PP data
- Results strongly dependent on ordering, e.g., reversing x<sub>obs</sub> makes phosphatase mechanism discriminative instead





Xohs

Harrington et al. (2012) PNAS

# Example: cell death signaling



- Crosslinking model: sequential Fas recruitment
- Cluster model: scaffold for Fas clustering, capable of bistability
- Can reject crosslinking model from cluster model data



Ho and Harrington (2010) PLoS Comput Biol, Harrington et al. (2012) PNAS

#### Remarks

Summary:

- Parameter-free statistical model invalidation
- Very simple yet can still be reasonably effective
- Cheap compared to parameter optimization
- Use as preprocessor to thin out model space
- Hierarchy of methods:

 $\mathsf{parametric} \to \mathsf{Bayesian} \to \mathsf{coplanarity}$ 

 Can probably generalize to periodic systems using integrated variables

Primary limitation: Gröbner bases are unreliable

- Does not always work
- Often requires manual intervention



# Beyond Gröbner bases

Can we eliminate without using Gröbner bases?

- Nonlinear methods can be somewhat brittle
- ▶ How far can we get using only linear (i.e., Gaussian) elimination?

# Beyond Gröbner bases

Can we eliminate without using Gröbner bases?

- Nonlinear methods can be somewhat brittle
- ▶ How far can we get using only linear (i.e., Gaussian) elimination?

Chemical reaction network models: nonlinear ODEs

Need to restrict to classes for which linear elimination is possible

#### Appeal to chemical reaction network theory

- Banaji, Conradi, Craciun, Dickenstein, Feinberg, Gunawardena, Horn, Jackson, Pantea, Pérez Millán, Shinar, Shiu, Sontag, and many others
- Qualitative dynamics, algebraic structure of chemical reaction networks
- Fundamental insight: there is a lot of hidden linearity

# Beyond Gröbner bases

Can we eliminate without using Gröbner bases?

- Nonlinear methods can be somewhat brittle
- ▶ How far can we get using only linear (i.e., Gaussian) elimination?

Chemical reaction network models: nonlinear ODEs

Need to restrict to classes for which linear elimination is possible

#### Appeal to chemical reaction network theory

- Banaji, Conradi, Craciun, Dickenstein, Feinberg, Gunawardena, Horn, Jackson, Pantea, Pérez Millán, Shinar, Shiu, Sontag, and many others
- Qualitative dynamics, algebraic structure of chemical reaction networks
- Fundamental insight: there is a lot of hidden linearity

What follows is very much a work in progress (with Heather Harrington). Any comments, thoughts, or connections are very much appreciated.

#### Chemical reaction network theory

$$\sum_{j=1}^{N} r_{ij}X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij}X_j, \quad i = 1, \dots, R$$

$$\stackrel{\mathbb{R}^C}{\longrightarrow} \frac{A_{\kappa}}{\mathbb{R}^C} \xrightarrow{\mathbb{R}^C} \psi$$

$$\stackrel{\chi_j}{\longrightarrow} \stackrel{\uparrow \psi}{\longrightarrow} \psi$$

$$\hat{x}_j = \sum_{i=1}^{R} \kappa_i (p_{ij} - r_{ij}) x^{r_i}, \quad j = 1, \dots, N$$

$$\stackrel{\mathbb{R}^S}{\longleftarrow} \frac{f}{\mathbb{R}^S} \xrightarrow{f} \mathbb{R}^S$$

$$\dot{x} = f(x) = YA_{\kappa}\Psi(x)$$

- Species:  $S = \{X_j \mid j = 1, \dots, N\}$
- Complexes:  $C = \left\{ \sum_{j=1}^{N} r_{ij} X_j, \sum_{j=1}^{N} p_{ij} X_j \mid i = 1, \dots, R \right\}$
- $\Psi$  : nonlinear species-to-complex map
- A<sub>κ</sub>: complex-to-complex rate matrix
- Y : complex-to-species stoichiometric matrix

Note that  $A_{\kappa}$  is linear; hence the dynamics in complex space is linear. Furthermore,  $A_{\kappa}$  is a Laplacian matrix (complex conservation).

#### **CRNT** example

$$X + Y \xrightarrow{\kappa_1} 2X, \qquad 3X \xrightarrow{\kappa_2} Y + 2Z, \qquad Z \xrightarrow{\kappa_3} X$$

S = {X, Y, Z}
C = {X, Y, Z, 2X, X + Y, 3X, Y + 2Z}

			x		у		z	$x^2$	xy	<i>x</i> <sup>3</sup>	yz <sup>2</sup>
	$\begin{bmatrix} x \end{bmatrix}$	x	Γ		$\kappa_3$						]
	y	У			$\kappa_4$						
		z		—	$\kappa_3$ –	$-\kappa_4$					
$\Psi(x) =$	$ x^2 $ ,	$A_{\kappa} = x^2$							$\kappa_1$		
	xy	xy							$-\kappa_1$		
	$ x^3 $	$x^{3}$								$-\kappa_2$	
	$yz^2$	yz <sup>2</sup>	L							$\kappa_2$	
					2		3	2			
		- ×	у	z	X <sup>2</sup>	хy	x	yz²	-		
		x   1	0	0	2	1	3	0			
		$Y = y \mid 0$	1	0	0	1	0	1			
		z [ 0	0	1	0	0	0	2			

Linear elimination in complex space (complex-linear invariants)

- Elimination on  $YA_{\kappa}$  (complex-to-species map)
  - Karp et al. (2012) J Theor Biol
  - · Same algorithm as before: check smallest singular value
  - · Complexity bounds, quite general, no ordering issues
  - Still can be difficult to understand
- Elimination on  $A_{\kappa}$  (complex-to-complex map)
  - Restrict to complex-balanced networks
  - Exploit Laplacian structure (off-diagonal non-negativity, diagonal dominance)
  - · Much more powerful results, can basically understand everything

## Definition

A chemical reaction network is complex-balanced if  $A_{\kappa}\Psi(x) = 0$  at any steady state  $x \in \mathbb{R}^{S}$ . A network is unconditionally complex-balanced if it is complex-balanced for all parameters  $\kappa$ .

## Definition

A chemical reaction network is complex-balanced if  $A_{\kappa}\Psi(x) = 0$  at any steady state  $x \in \mathbb{R}^{S}$ . A network is unconditionally complex-balanced if it is complex-balanced for all parameters  $\kappa$ .

Properties:

- Completely specified by complex reaction graph (directed, acyclic)
- Precludes "interesting" behavior, e.g., no multistationarity
- Sufficient graph-theoretic condition: deficiency zero (Feinberg)

Examples:

$$A + B \xrightarrow{\sim} AB \qquad A \xrightarrow{\sim} B \longrightarrow C \xrightarrow{\sim} D$$
$$A + B \xrightarrow{\sim} AB \longrightarrow A + C \xrightarrow{\sim} AC$$

Choose from C an arbitrary subset  $C^*$  of observable complexes (comprising only observable species). We first consider how to compute invariants in  $C^*$ .

- Invariants are useful in their own right beyond model selection
- Absolute concentration robustness, e.g.,

$$\left(\frac{\kappa_1\kappa_3}{\kappa_2}\right)xy - (\kappa_4 + \kappa_5)x^2y = 0$$
 implies  $x = \frac{\kappa_1\kappa_3}{\kappa_2(\kappa_4 + \kappa_5)}$  if  $y \neq 0$ 

Choose from C an arbitrary subset  $C^*$  of observable complexes (comprising only observable species). We first consider how to compute invariants in  $C^*$ .

- Invariants are useful in their own right beyond model selection
- Absolute concentration robustness, e.g.,

$$\left(\frac{\kappa_1\kappa_3}{\kappa_2}\right)xy - (\kappa_4 + \kappa_5)x^2y = 0$$
 implies  $x = \frac{\kappa_1\kappa_3}{\kappa_2(\kappa_4 + \kappa_5)}$  if  $y \neq 0$ 

#### Result: Gaussian elimination never breaks

In other words, if  $A_{\kappa}$  is block partitioned as

$$A_{\kappa} = rac{\mathcal{C}^*}{\mathcal{C}\setminus\mathcal{C}^*} \left[ egin{array}{ccc} \mathcal{A} & \mathcal{C}\setminus\mathcal{C}^* \ \mathcal{A} & \mathcal{B} \ \mathcal{C} & \mathcal{D} \end{array} 
ight],$$

then the Schur complement  $A - BD^{-1}C$  for  $C^*$  always exists (and provides invariant coefficients).

It is possible that the resulting invariants  $\sum_i b_i y_i$  have b = 0, in which case y need not be coplanar. How to guarantee unconditionally nontrivial invariants?

Necessary and sufficient graph-theoretic conditions (closed systems):

- Exists  $c \in C^*$  in a non-terminal SCC
- Exists distinct  $c, c' \in \mathcal{C}^*$  in the same terminal SCC

Proofs are standard (induct on terminal SCCs, diagonal dominance).



Remarks:

Intuition:

- Has a sink, concentration goes to zero
- Proportional concentrations by equilibrium constant
- Similar statements for open systems (synthesis and degradation)
- Punchline: can determine if complexes are coplanar by inspection

- Closely related to Feinberg's results on ker A<sub>k</sub>
- ► In fact, our approach can be used to prove and extend to open systems
  - Matrix approach is quite easy
- Assume no constitutive synthesis without degradation somewhere
- Main results:
  - If C is in a non-terminal SCC, then  $x^{C} = 0$
  - If C is in a terminal SCC without syn/deg, then dim(span( $x^{C}$ )) = 1
  - If C is in a terminal SCC with syn/deg, then  $x^{C} = \chi > 0$  fixed
- Characterizes concentration robustness for complex-balanced networks
- Much more stringent than coplanarity:
  - Test for zero, constant, and rank-one ( $\sigma_2 = 0$ )



Given a chemical reaction network:

- Check complex-balancing (deficiency zero)
- Determine all steady-state properties by graph inspection
  - Zero, constant, rank-one
  - Only fast, scalable graph algorithms required
- Test all steady-state properties statistically
- Control rejections with FDR

## Conclusion

- Parameter-free statistical model invalidation
- Quantitative "qualitative" biology
- Current disconnect:
  - Complex-balanced: quite limited, but know everything
  - Everything else: don't know very much of anything
  - How to bridge the gap?
- Fundamental idea: detect low-dimensional representations
  - Specific representation may be parametric, but low dimensionality is not
  - Can exploit in other ways besides coplanarity, rank-zero, rank-one
- Success perhaps attributable to biological robustness
- ► Generalizations: Laplacian dynamics, reverse engineering, design

## Conclusion

- Parameter-free statistical model invalidation
- Quantitative "qualitative" biology
- Current disconnect:
  - Complex-balanced: quite limited, but know everything
  - Everything else: don't know very much of anything
  - How to bridge the gap?
- Fundamental idea: detect low-dimensional representations
  - Specific representation may be parametric, but low dimensionality is not
  - Can exploit in other ways besides coplanarity, rank-zero, rank-one
- Success perhaps attributable to biological robustness
- ► Generalizations: Laplacian dynamics, reverse engineering, design

Acknowledgements:

- ► Heather Harrington, Tom Thorne, Michael Stumpf, Anne Shiu
- ICL Theoretical Systems Biology group