Parameter-free statistical model invalidation for biochemical reaction networks

Kenneth L. Ho (Stanford)

Joint work with Heather Harrington (Oxford)

Theranos, Apr. 2015

- Model selection: observed data, multiple models; which model is 'best'?
- Example (sequential phosphorylation):



• Two models: distributive ($\kappa_{02} = 0$), processive ($\kappa_{02} > 0$)



Closely related to model invalidation

- Data x, model $f = f(\kappa)$ with parameters κ
- How to tell if model is incompatible with data?
- **Known parameters**: compute $\hat{x} = f(\kappa)$ and check $||x \hat{x}||$
- Unknown parameters: fit parameters and check best-case error

- Data x, model $f = f(\kappa)$ with parameters κ
- How to tell if model is incompatible with data?
- **Known parameters**: compute $\hat{x} = f(\kappa)$ and check $||x \hat{x}||$
- Unknown parameters: fit parameters and check best-case error

Parameter problem: in biology, parameters are hardly ever known

- ▶ Technical limitations, uncertainties, etc.
- Partial data: experimentally inaccessible species
- Nonlinear, high-dimensional optimization often required



- Data x, model $f = f(\kappa)$ with parameters κ
- How to tell if model is incompatible with data?
- **Known parameters**: compute $\hat{x} = f(\kappa)$ and check $||x \hat{x}||$
- Unknown parameters: fit parameters and check best-case error

Parameter problem: in biology, parameters are hardly ever known

- ► Technical limitations, uncertainties, etc.
- Partial data: experimentally inaccessible species
- Nonlinear, high-dimensional optimization often required



What can be done?

• Optimization 'tricks': random seeding, simulated annealing, etc.

- > Optimization 'tricks': random seeding, simulated annealing, etc.
- Convex relaxation + SDP: lower bound for best-case error
 - Polynomial-time search through parameter space

- > Optimization 'tricks': random seeding, simulated annealing, etc.
- Convex relaxation + SDP: lower bound for best-case error
 - Polynomial-time search through parameter space

This talk: (quantitative) parameter-free methods

- Like SDP but no dependence on parameters
- Based only on model structure/topology
- Not necessarily 'better' but new framework



- > Optimization 'tricks': random seeding, simulated annealing, etc.
- Convex relaxation + SDP: lower bound for best-case error
 - Polynomial-time search through parameter space

This talk: (quantitative) parameter-free methods

- Like SDP but no dependence on parameters
- Based only on model structure/topology
- Not necessarily 'better' but new framework

Philosophically related (qualitative):

- Chemical reaction network theory
- Stoichiometric network analysis



- Reactions:
- Mass-action dynamics:

$$\sum_{j=1}^{N} r_{ij} X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij} X_j, \qquad i = 1, \dots, R$$
$$\dot{x}_j = \sum_{i=1}^{R} \kappa_i \left(p_{ij} - r_{ij} \right) \prod_{k=1}^{N} x_k^{r_{ik}}, \quad j = 1, \dots, N$$

Reactions:

$$\sum_{j=1}^{N} r_{ij}X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij}X_j, \qquad i = 1, \dots, R$$
$$\dot{x}_j = \sum_{i=1}^{R} \kappa_i (p_{ij} - r_{ij}) \prod_{k=1}^{N} x_k^{r_{ik}}, \quad j = 1, \dots, N$$

Example:



$$\begin{split} \dot{x}_1 &= \kappa_1 x_1 x_2 - 3 \kappa_2 x_1^3 + \kappa_3 x_3 \\ \dot{x}_2 &= -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3 \\ \dot{x}_3 &= 2 \kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3 \end{split}$$

Reactions:

$$\sum_{j=1}^{N} r_{ij} X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij} X_j, \qquad i = 1, \dots, R$$
$$\dot{x}_j = \sum_{i=1}^{R} \kappa_i \left(p_{ij} - r_{ij} \right) \prod_{k=1}^{N} x_k^{r_{ik}}, \quad j = 1, \dots, N$$

Example:

 $X_1 + X_2 \xrightarrow{\kappa_1} 2X_1$



$$\begin{split} \dot{x}_1 &= \kappa_1 x_1 x_2 - 3\kappa_2 x_1^3 + \kappa_3 x_3 \\ \dot{x}_2 &= -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3 \\ \dot{x}_3 &= 2\kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3 \end{split}$$

Reactions:

$$\sum_{j=1}^{N} r_{ij} X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij} X_j, \qquad i = 1, \dots, R$$
$$\dot{x}_j = \sum_{i=1}^{R} \kappa_i \left(p_{ij} - r_{ij} \right) \prod_{k=1}^{N} x_k^{r_{ik}}, \quad j = 1, \dots, N$$

Example:

 $X_1 + X_2 \xrightarrow{\kappa_1} 2X_1$

 $3X_1 \xrightarrow{\kappa_2} X_2 + 2X_3$



$$\begin{split} \dot{x}_1 &= \kappa_1 x_1 x_2 - 3\kappa_2 x_1^3 + \kappa_3 x_3 \\ \dot{x}_2 &= -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3 \\ \dot{x}_3 &= 2\kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3 \end{split}$$

Reactions:

$$\sum_{j=1}^{N} r_{ij} X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij} X_j, \qquad i = 1, \dots, R$$
$$\dot{x}_j = \sum_{i=1}^{R} \kappa_i \left(p_{ij} - r_{ij} \right) \prod_{k=1}^{N} x_k^{r_{ik}}, \quad j = 1, \dots, N$$

Example:



 $\begin{aligned} \dot{x}_1 &= \kappa_1 x_1 x_2 - 3\kappa_2 x_1^3 + \kappa_3 x_3 \\ \dot{x}_2 &= -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3 \\ \dot{x}_3 &= 2\kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3 \end{aligned}$

Reactions:

$$\sum_{j=1}^{N} r_{ij} X_j \xrightarrow{\kappa_i} \sum_{j=1}^{N} p_{ij} X_j, \qquad i = 1, \dots, R$$
$$\dot{x}_j = \sum_{i=1}^{R} \kappa_i \left(p_{ij} - r_{ij} \right) \prod_{k=1}^{N} x_k^{r_{ik}}, \quad j = 1, \dots, N$$

Example:



 $\begin{aligned} \dot{x}_1 &= \kappa_1 x_1 x_2 - 3\kappa_2 x_1^3 + \kappa_3 x_3 \\ \dot{x}_2 &= -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3 \\ \dot{x}_3 &= 2\kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3 \end{aligned}$

ODE: quantitative test for model compatibility

$$\begin{split} \dot{x}_1 &= \kappa_1 x_1 x_2 - 3 \kappa_2 x_1^3 + \kappa_3 x_3 \\ \dot{x}_2 &= -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3 \\ \dot{x}_3 &= 2 \kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3 \end{split}$$

- ODE: quantitative test for model compatibility
- Problem: partial data

$$\begin{split} \dot{x}_1 &= \kappa_1 x_1 x_2 - 3 \kappa_2 x_1^3 + \kappa_3 x_3 \\ \dot{x}_2 &= -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3 \\ \dot{x}_3 &= 2 \kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3 \end{split}$$

- ODE: quantitative test for model compatibility
- ▶ Problem: partial data
 - Derivatives unreliable; assume steady state

$$0 = \dot{x}_1 = \kappa_1 x_1 x_2 - 3\kappa_2 x_1^3 + \kappa_3 x_3$$
$$0 = \dot{x}_2 = -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3$$
$$0 = \dot{x}_3 = 2\kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3$$

- ODE: quantitative test for model compatibility
- Problem: partial data
 - Derivatives unreliable; assume steady state
 - Eliminate experimentally inaccessible species

$$0 = \dot{x}_1 = \kappa_1 x_1 x_2 - 3\kappa_2 x_1^3 + \kappa_3 x_3$$

$$0 = \dot{x}_2 = -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3$$

$$0 = \dot{x}_3 = 2\kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3$$

- ODE: quantitative test for model compatibility
- Problem: partial data
 - Derivatives unreliable; assume steady state
 - Eliminate experimentally inaccessible species

Example: $0 = \dot{x}_1 = \kappa_1 x_1 x_2 - 3\kappa_2 x_1^3 + \kappa_3 x_3$ $0 = \dot{x}_2 = -\kappa_1 x_1 x_2 + \kappa_2 x_1^3 + \kappa_4 x_3$ $0 = \dot{x}_3 = 2\kappa_2 x_1^3 - (\kappa_3 + \kappa_4) x_3$

• Observe x_1, x_2 ; eliminate $x_3 = \frac{2\kappa_2 x_1^3}{\kappa_3 + \kappa_4}$

- ODE: quantitative test for model compatibility
- Problem: partial data
 - Derivatives unreliable; assume steady state
 - Eliminate experimentally inaccessible species

$$\begin{split} \mathbf{0} &= \kappa_1 x_1 x_2 + \left(\frac{2\kappa_3}{\kappa_3 + \kappa_4} - 3\right) \kappa_2 x_1^3 \\ \mathbf{0} &= -\kappa_1 x_1 x_2 + \left(\frac{2\kappa_4}{\kappa_3 + \kappa_4} + 1\right) \kappa_2 x_1^3 \end{split}$$

• Observe x_1, x_2 ; eliminate $x_3 = \frac{2\kappa_2 x_1^3}{\kappa_3 + \kappa_4} \implies$ steady-state invariants

- ODE: quantitative test for model compatibility
- Problem: partial data
 - Derivatives unreliable; assume steady state
 - Eliminate experimentally inaccessible species

$$\begin{split} \mathbf{0} &= \kappa_1 x_1 x_2 + \left(\frac{2\kappa_3}{\kappa_3 + \kappa_4} - 3\right) \kappa_2 x_1^3 \\ \mathbf{0} &= -\kappa_1 x_1 x_2 + \left(\frac{2\kappa_4}{\kappa_3 + \kappa_4} + 1\right) \kappa_2 x_1^3 \end{split}$$

- Observe x_1, x_2 ; eliminate $x_3 = \frac{2\kappa_2 x_1^3}{\kappa_3 + \kappa_4} \implies$ steady-state invariants
- In general, use computational algebraic geometry (Gröbner bases):

$$0 = \sum_{i=1}^{n} \alpha_i(\kappa) \varphi_i(x)$$

[Manrai/Gunawardena]

$$\sum_{i=1}^{n} lpha_i(\kappa) arphi_i(x^{(1)}) = 0$$

 \vdots
 $\sum_{i=1}^{n} lpha_i(\kappa) arphi_i(x^{(m)}) = 0$

[Manrai/Gunawardena]

$$\sum_{i=1}^{n} \alpha_i(\kappa) \varphi_i(x^{(1)}) = 0$$
$$\vdots$$
$$\sum_{i=1}^{n} \alpha_i(\kappa) \varphi_i(x^{(m)}) = 0$$

• Let
$$y^{(k)} = (\varphi_1(x^{(k)}), \dots, \varphi_n(x^{(k)})) \in \mathbb{R}^n$$

• **Geometry**: $y^{(1)}, \ldots, y^{(m)}$ are coplanar

Depends only on data, hence parameter-free



$$\sum_{i=1}^{n} lpha_i(\kappa) arphi_i(x^{(1)}) = 0$$

 \vdots
 $\sum_{i=1}^{n} lpha_i(\kappa) arphi_i(x^{(m)}) = 0$

• Let
$$y^{(k)} = (\varphi_1(x^{(k)}), \dots, \varphi_n(x^{(k)})) \in \mathbb{R}^n$$

• **Geometry**: $y^{(1)}, \ldots, y^{(m)}$ are coplanar

Depends only on data, hence parameter-free

Linear algebra: $Y\alpha = 0$

- Compatible $\implies \exists \alpha \in \operatorname{null}(Y) \implies \operatorname{dim}(\operatorname{null}(Y)) > 0$
- Compute SVD, reject if $\sigma_{\min}(Y) > 0$



- ▶ Null hypothesis: $\sigma_{\min}(Y) = \min_{\|\alpha\|=1} \|Y\alpha\| = 0 \implies \exists \alpha \text{ such that } Y\alpha = 0$
- Assume Gaussian noise in $x^{(k)}$, estimate noise in $y^{(k)} = \varphi(x^{(k)})$
- To first-order, Gaussian noise in $Y \implies z = Y \alpha$ Gaussian
- Rescale rows: $DY \implies (Dz)_i \sim \mathcal{N}(0, \sigma_i^2), \ \sigma_i^2 \leq 1$
- ► Tail bound: $\Pr(\sigma_{\min}(DY) > t) \le \Pr(\|Dz\| > t) \le \Pr(\chi_m > t)$
- ▶ Other bounds possible: Weyl, Wielandt-Hoffman, concentration of measure, etc.

Algorithm



- Test coplanarity for each invariant
- Reject model if any invariant fails
- Main costs: elimination, $\sigma_{\min}(\mathbb{R}^{m \times n})$



Kinase/phosphatase: distributive/processive

obs

- ► Four models: PP, PD, DP, DD
- 12 species, 22 parameters
- ▶ Variable ordering: (*ks*₀₀, *ks*₀₁, *ks*₁₀, *fs*₀₁, *fs*₁₀, *fs*₁₁, *k*, *f*, *s*₀₀, *s*₀₁, *s*₁₀, *s*₁₁)
- ▶ Kinase is discriminative, can reject DP/DD models on basis of PP data
- Can discriminate instead on phosphatase by reversing variable ordering







- Extrinsic pathway
- FasL/Fas interactions
- Measure activated Fas
- Crosslinking model: sequential Fas recruitment (8 species, 2 parameters)
- Cluster model: scaffold for Fas clustering, bistable (6 species, 9 parameters)
- Can reject crosslinking model from cluster data





[Harrington/Ho/Thorne/Stumpf, Ho/Harrington, Lai/Jackson]

- Parameter-free statistical model invalidation
- Detection of non-parametric linear structure
- Very efficient once invariants have been computed

Broader perspective: parameter-independent model properties

- Structure, topology, robustness, modularity
- Algebraic systems biology

Limitations: nonlinear elimination, steady-state data, necessary but not sufficient

- CRNs: nonlinear ODEs \implies nonlinear elimination
- Fundamental insight of CRNT: hidden linearity
- Complex-balanced networks: underlying Laplacian dynamics
- Study properties of Laplacian matrices
- Steady state: kernel = zero + constant + rank-one
- No elimination: decomposition based on graph connectivity
- Test σ₁, σ₂, etc.



- CRNs: nonlinear ODEs \implies nonlinear elimination
- Fundamental insight of CRNT: hidden linearity
- Complex-balanced networks: underlying Laplacian dynamics
- Study properties of Laplacian matrices
- Steady state: kernel = zero + constant + rank-one
- No elimination: decomposition based on graph connectivity
- Test σ₁, σ₂, etc.



[Harrington/Ho]

- CRNs: nonlinear ODEs \implies nonlinear elimination
- Fundamental insight of CRNT: hidden linearity
- Complex-balanced networks: underlying Laplacian dynamics
- Study properties of Laplacian matrices
- Steady state: kernel = zero + constant + rank-one
- No elimination: decomposition based on graph connectivity
- Test σ₁, σ₂, etc.



- CRNs: nonlinear ODEs \implies nonlinear elimination
- Fundamental insight of CRNT: hidden linearity
- Complex-balanced networks: underlying Laplacian dynamics
- Study properties of Laplacian matrices
- Steady state: kernel = zero + constant + rank-one
- No elimination: decomposition based on graph connectivity
- Test σ₁, σ₂, etc.



[Harrington/Ho]

- Differential elimination: dynamical invariants involving derivatives
- ► Example: Lotka-Volterra $\dot{x} = ax - bxy$ $\dot{y} = -cy + dxy$ $\Rightarrow acx^2 + ax\dot{x} - bcx^3 - bx^2\dot{x} = -\dot{x}^2 + x\ddot{x}$
- Estimate derivatives using Gaussian processes



- Goal: solve global nonlinear optimization
- Exploit polynomial structure, numerical algebraic geometry
- Find closest intersection between model and data varieties
- Maximum-likelihood parameter estimation, model invalidation, model selection



- Mathematical modeling of cell signaling networks
- Automated all-atom protein crystal structure refinement
- Fast multipole methods, direct solvers, matrix factorizations



[Harrington/Ho/Ghosh/Tung, Ho/Harrington]

- Mathematical modeling of cell signaling networks
- Automated all-atom protein crystal structure refinement
- ► Fast multipole methods, direct solvers, matrix factorizations



- Mathematical modeling of cell signaling networks
- Automated all-atom protein crystal structure refinement
- Fast multipole methods, direct solvers, matrix factorizations



[Greengard/Ho/Lee, Ho/Greengard, Ho/Ying, Li/Yang/Martin/Ho/Ying, Minden/Damle/Ho/Ying]

- Mathematical modeling of cell signaling networks
- Automated all-atom protein crystal structure refinement
- ► Fast multipole methods, direct solvers, matrix factorizations



[Greengard/Ho/Lee, Ho/Greengard, Ho/Ying, Li/Yang/Martin/Ho/Ying, Minden/Damle/Ho/Ying]

References

Parameter-free invalidation:

- E. Gross, B. Davis, K.L. Ho, D.J. Bates, H.A. Harrington. Numerical algebraic geometry for model selection. In preparation.
- H.A. Harrington, K.L. Ho. Parameter-free statistical model invalidation for weakly complex-balanced chemical reaction networks. In preparation.
- H.A. Harrington, K.L. Ho, N. Meshkat. Model rejection using differential algebra techniques. In preparation.
- H.A. Harrington, K.L. Ho, T. Thorne, M.P.H. Stumpf. Parameter-free model discrimination criterion based on steady-state coplanarity. Proc. Natl. Acad. Sci. U.S.A. 109 (39): 15746–15751, 2012.

Other:

- J.A. Bell, K.L. Ho, R. Farid. Significant reduction in errors associated with nonbonded contacts in protein crystal structures: automated all-atom refinement with PrimeX. Acta Cryst. D68: 935–952, 2012.
- L. Greengard, K.L. Ho, J.-Y. Lee. A fast direct solver for scattering from periodic structures with multiple material interfaces in two dimensions. J. Comput. Phys. 258: 738–751, 2014.
- H.A. Harrington, K.L. Ho, S. Ghosh, KC Tung. Construction and analysis of a modular model of caspase activation in apoptosis. Theor. Biol. Med. Model. 5: 26, 2008.
- K.L. Ho. Fast direct methods for molecular electrostatics. Ph.D. thesis, New York University, 2012.
- K.L. Ho, L. Greengard. A fast direct solver for structured linear systems by recursive skeletonization. SIAM J. Sci. Comput. 34 (5): A2507–A2532, 2012.
- K.L. Ho, L. Greengard. A fast semidirect least squares algorithm for hierarchically block separable matrices. SIAM J. Matrix Anal. Appl. 35 (2): 725–748, 2014.
- K.L. Ho, H.A. Harrington. Bistability in apoptosis by receptor clustering. PLoS Comput. Biol. 6 (10): e1000956, 2010.
- K.L. Ho, L. Ying. Hierarchical interpolative factorization for elliptic operators: differential equations. Preprint, arXiv:1307.2895 [math.NA], 2013. To appear, Comm. Pure Appl. Math.
- K.L. Ho, L. Ying. Hierarchical interpolative factorization for elliptic operators: integral equations. Preprint, arXiv:1307.2666 [math.NA], 2013. To appear, Comm. Pure Appl. Math.
- Y. Li, H. Yang, E. Martin, K. Ho, L. Ying. Butterfly factorization. Preprint, arXiv:1502.01379 [math.NA], 2015.
- V. Minden, A. Damle, K.L. Ho, L. Ying. A technique for updating hierarchical factorizations of integral operators. Preprint, arXiv:1411.5706 [math.NA], 2014.