

Discerning protein activity regulation mechanism using "inverse" machine learning

Sameer Varma

University of South Florida, Tampa

Protein function



.

Protein function regulation



Post-translational regulation of protein function

Regulators

- (a) non-covalent interactions
- (b) covalent modifications
- (c) pH
- (d) salt concentrations
- (e)
- () combination of any of the above



Mechanism of protein function

Understanding how a protein works requires a systematic assessment of relationships between the structure, dynamics and energetics of its various states



Pioneers: Molecular simulation methods for biomolecules



The 1970's revolutionary vision...

Since energy drives structure, dynamics and activity, shouldn't we be able to predict them directly, provided we truly understand the energetics of underlying interactions?









Michael Levitt

Arieh Warshel

Direct estimation of activity from energy

$$H = \sum \frac{p_i^2}{2m_i} + V$$

Define V

Example:
$$V = \sum V_{ij}^{bond} + \sum V_{ijk}^{angle} + \sum V_{ijkl}^{dihedral} + \sum V_{ijkl}^{dihedral} + \sum V_{ij}^{electrostatics} + \sum V_{ij}^{vdw}$$

Solve Hamilton's equation of motion to get $r_i(t)$

$$\frac{d\vec{r_i}}{dt} = +\frac{\partial H}{\partial \vec{p}_i} \quad , \quad \frac{d\vec{p}_i}{dt} = -\frac{\partial H}{\partial \vec{r}_i}$$





How well do we understand inter-atomic interactions?

Reliability test: ab initio protein folding



Pioneers

Nobel Prizes and Laureates

÷ < 2013 >

Chemistry Prizes

 About the Nobel Prize in Chemistry 2013
 Summary
 Prize Announcement
 Press Release
 Advanced Information
 Popular Information
 Greetings

Martin Karplus

- Michael Levitt
- Arieh Warshel

All Nobel Prizes in Chemistry All Nobel Prizes in 2013



The Nobel Prize in Chemistry 2013 Martin Karplus, Michael Levitt, Arieh Warshel

The Nobel Prize in Chemistry 2013



© Nobel Media AB Martin Karplus



Photo: Keilana via Wikimedia Commons Michael Levitt



Photo: Wikimedia Commons

Arieh Warshel

The Nobel Prize in Chemistry 2013 was awarded jointly to Martin Karplus, Michael Levitt and Arieh Warshel *"for the development of multiscale models for complex chemical systems"*.

Does this imply that we can now blindly apply molecular simulations to any protein?

NO!

While there is a lot that we can model, which we couldn't just a decade ago,

there remains a lot more that we still cannot...



Our interests



fundamental atomically-detailed principles underlying

- **1)** Selective ion transport by channels
- 2) Protein assembly

3) Allosteric regulation of protein activity

Allostery

Basic problem



Allostery

General problem



The regulatory signal

 $\Delta \mathbb{R}_{signal} \subset \Delta \mathbb{R}$

Allosteric regulation involving large structural changes

For many proteins, their regulatory models can be understood in terms of how their 3D structures, or conformations, change during transition

(a) Negligible overlap between conformational ensembles of 2 states



Many methods are available to study allosteric signaling in such proteins

Allosteric regulation involving small structural changes

For several major pharmaceutical proteins, distinguishing their states **requires** consideration of their finite temperature conformational ensembles

(b) Discernible overlap between conformational ensembles of 2 states



Understanding mechanisms requires accounting for information on both structure and dynamics from multiple states.

Role of dynamic allostery - examples

Fusion regulation in paramyxoviruses



GPCR regulation by natural compounds and drug molecules



Regulation of Immune response



Heat shock response



- 1) A method to quantify differences between conformational ensembles of two states
- **Traditional approach**



Problem with comparing summary statistics: Need to know the "right set" of summary statistics that differentiate ensembles. But how does one identify such summary statistics beforehand?

2) A method to relate $\Delta \mathbb{R}$ to regulatory signals

Current strategy

- Construct correlation matrices
- Determine efficient signaling paths from correlation matrices.



2) A method to related $\Delta \mathbb{R}$ to regulatory signals

Assumptions

(i) structural changes are large, and so conformational fluctuations can be ignored.

<u>OR</u> (ii) ignore $\Delta \mathbb{R}$ and focus only on \mathbb{R}' of only the active state.

Problem: These assumptions break down for proteins regulated by dynamic allostery

Method	Description	Predictions	tot. num.	тр	TN	FP	FN	AUC (Fig. 3A)	AUC (Fig. 3B)	TPR	FPR
Fuentes et al 2004	Experimental measurements	L18, I20, V22, V26, V30, A39, V40, V61, V64, L66, A69, L78, T81, V85	14	n. a.	n. a.	n.a.	n. a.				
Kong and Karplus 2009	Molecular Dynamics (MD) approach	A12, L18, V26, I41, A45, A46, V58, L59, L66, A69, A74, L78, T81, L89	14	6	6	5	8	n. a.	n. a.	0.43	0.45
Serek and Ozkan 2011	Perturbation Response Scanning method	L11, L18, I20, V22, T23, I35, V37, A39, V40, I41, A45, A46, V58, L59, A60, V61, L66, A69, A74, V75, T77, L78, T81, V85, L87	26	11	3	8	3	n. a.	n. a.	0.79	0.73

Cilia et al. PLoS Comp Bio. 2012

- 11 11 11

New method the quantify differences in ensembles

Instead of comparing summary statistics, compare ensembles directly against each other





Leighty and Varma, JCTC 2013 Varma, Botlani, Leighty, Proteins 2014

 $\eta \in$

Works even for multi-modal distributions

$$\eta = 1 - \left\| \sum_{i=1}^{n} c_i f_i \cap \sum_{j=1}^{n} c'_j f'_j \right\| = 1 - \left\| \sum_{i,j=1}^{n} c_i f_i \cap c'_j f'_j \right\|$$



Dutta, Siddiqui, Botlani, Varma. BJ 2016

Machine Learning

Typical scenario: learning from data

- Given data set $\,x\,$ and labels $\,y\,$

(generated by some joint probability distribution p(x,y))



- LEARN/INFER underlying unknown mapping

$$y = f(x)$$

Learn f from examples such that **R**isk of <u>prediction</u>, is minimized

$$R[f] = \int \left| f(x) - y \right|^2 dP(x, y)$$



"Inverse" machine learning

Typical scenario: learning from data

- Given data set $\,x\,$ and labels $\,y\,$

(generated by some joint probability distribution p(x,y))



- LEARN/INFER underlying unknown mapping

y = f(x)

Learn f from examples such that **R**isk of <u>prediction</u>, is minimized

Instead of training f for prediction

- Construct and train f in an appropriate Hilbert space such that
- it can be used to derive causalities in physical space

A support vector machine is a binary classifier that is trained on a set of instances $\{x_i\}$ for which their corresponding group identities $y_i = \pm 1$ are known.

It is constructed by defining two hyperplanes

$$y_i(\mathbf{w}\cdot\mathbf{x}-b)=1$$

that divide the instances into two groups,

$$y_i(\mathbf{w} \cdot \mathbf{x} - b) \ge 1$$

which generates a classification function

$$F(\mathbf{x}) = \sum_{i=1}^{2m} \alpha_i y_i k(\mathbf{x}_i, \mathbf{x})$$



The optimization task: maximize the distance (2/||w||) between the two hyperplanes.

This constrained optimization problem can be cast in terms of Lagrangian multipliers as

$$L = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^{2m} \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i - b) - 1]$$

where $0 \le \alpha_i \le C$ are the Lagrange multipliers, and the limit *C* is a regularization parameter

Set $\nabla L = 0$

$$\frac{\partial L}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^{2m} \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^{2m} \alpha_i y_i = 0$$



$$L = \sum_{i=1}^{2m} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

where $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$



 $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i) . \phi(\mathbf{x}_j) \rangle_R = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$$

Optimization of

$$L = \sum_{i=1}^{2m} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

yields two sets of $\alpha_{i:}$



The number of support vectors are bounded, that is, for type of instance $2 \le s \le 2m$

Use $\{\alpha_i\} > 0$ to construct the classification function:

$$F(\mathbf{x}) = \sum_{i=1}^{2m} \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) \longrightarrow \text{use for prediction}$$



Optimization of

$$L = \sum_{i=1}^{2m} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

yields two sets of $\alpha_{i:}$

The number of support vectors are bounded, that is, for type of instance $2 \le s \le 2m$.

It can be expected that

Higher the number of support vectors (s), the less discriminable are the two groups in a given Hilbert space.

Is it possible to define discriminability between two groups as $\eta = 1 - \frac{s}{2m}$



Selection of *C* and γ

Possibility?

Can we select C and γ such that

 $\eta = 1 - s/2m$

=1-||Overlap||

 $= 1 - \left| \left| \mathbb{R} \cap \mathbb{R}' \right| \right|$





Such a definition of discriminability satisfies two conditions

(i)
$$|\eta(\mathbb{R} \to \mathbb{R}')| = |\eta(\mathbb{R}' \to \mathbb{R})|$$

(ii) If $|\eta(\mathbb{R} \to \mathbb{R}')| = |\eta(\mathbb{R}' \to \mathbb{R}'')|$, then it does not

necessarily imply that $\left|\eta(\mathbb{R} o\mathbb{R'})
ight| = \left|\eta(\mathbb{R} o\mathbb{R''})
ight|$

Leighty and Varma, JCTC 2013

Selection of *C* and γ

Possibility?

Can we select C and γ such that

 $\eta = 1 - s/2m$

=1-||Overlap||

 $= 1 - ||\mathbb{R} \cap \mathbb{R'}||$





Selection of *C* and γ



Leighty and Varma, JCTC 2013 Dutta, Siddiqui, Botlani, Varma. BJ 2016

Performance



Performance

Works for multi-modal distributions

$$\eta = 1 - \left| \left| \sum_{i=1}^{n} c_i f_i \cap \sum_{j=1}^{n} c'_j f'_j \right| \right| = 1 - \left| \left| \left(\sum_{i,j=1}^{n} c_i f_i \cap c'_j f'_j \right) \right| \right|$$







Simplifies data representation and analysis

Effect of oxidation on phosphatase structure and dynamics



New applications enabled

Permits ranking of population shifts



Leighty and Varma, JCTC 2013



ing the folding and unfolding processes of proteins as a function of temperature is a major charteins as a function of temperature is a major charteins as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens as a function of temperature is a major chartens of temperature is a function of temperature is a major chartens of temperature is a major chartens of temperature is a major cool proteins within tens of temperature is a function of temperature is a major cool protein within tens of temperature is a function of the temperature is a major cool protein within tens of temperature is a major cool protein within tens of temperature is a major condition. Our result and cool proteins within tens of temperature is a major condition.

Enables comparison of multiple population shifts



Varma, et al. Proteins 2014

Enables comparison of multiple population shifts

Effect of mutation on ligand-induced population shifts



Set of residues affected by mutation satisfy at least one of the three conditions

$$|\eta - \eta^{m}| > 2 \times \text{MAE},$$

 $\eta_{\text{apo}} > \operatorname{erf}\left(1/\sqrt{2}\right), \text{ and}$
 $\eta_{\text{bnd}} > \operatorname{erf}\left(1/\sqrt{2}\right).$

Dutta, et al. BJ 2016

Source code and tutorials @ SimTK



Version 2.0.4. Website design by Viewfarm. Icons created by SimTK team using art by GraphBerry from www.flaticon.com under a CC BY 3.0 license. Forked from FusionForge 5.3.2.

Dutta, Siddiqui, Botlani, Varma. BJ 2016

Sincerely thankful to...

Contributing lab members



Mohsen Botlani





Ahnaf Siddiqui



Collaborators

Machine learning



Klaus Muller (TU Berlin)

Protein activity & regulation



Matteo Porotto

(Columbia U)



Stan Stevens (USF)





Priyanka Dutta

Thank you.