

BIOMENG 261

TISSUE AND BIOMOLECULAR ENGINEERING

Module I: Reaction kinetics and systems biology

Oliver Maclarens
oliver.maclarens@auckland.ac.nz

LECTURE 8 SYSTEMS BIOLOGY CONT'D

- Intro to parameter estimation
- Intro to Flux Balance Analysis (FBA)

1

3

MODULE OVERVIEW

Reaction kinetics and systems biology (*Oliver Maclarens*)
[11-12 lectures/3 tutorials/2 labs]

1. Basic principles: modelling with reaction kinetics [5-6 lectures]

Physical principles: conservation, directional and constitutive. Reaction modelling. Mass action. Enzyme kinetics. Enzyme regulation. Mathematical/graphical tools for analysis and fitting.

2. Systems biology I: signalling and metabolic systems [3 lectures]

Overview of systems biology. Modelling signalling systems using reaction kinetics. Introduction to parameter estimation. Modelling metabolic systems using reaction kinetics. Flux balance analysis and constraint-based methods.

3. Systems biology II: genetic systems [3 lectures]

Modelling genes and gene regulation using reaction kinetics. Gene regulatory networks, transcriptomics and analysis of microarray data.

FORWARD AND INVERSE PROBLEMS

Typical modelling deals with (so-called) *forward problems*:

*Given parameters and initial conditions,
predict data*

2

4

TRADE-OFFS

The key to solving ill-posed problems is recognising the *trade-offs* involved

- Fit vs complexity
- Overfitting vs underfitting
- Training vs test
- Efficiency vs stability
- Bias vs variance
- Etc

These trade-offs are closely related and *recur throughout science, statistics and engineering*

5

7

FORWARD AND INVERSE PROBLEMS

However, in science and engineering we are usually confronted with *inverse problems*:

Given data, estimate parameters and/or initial conditions and predict future data

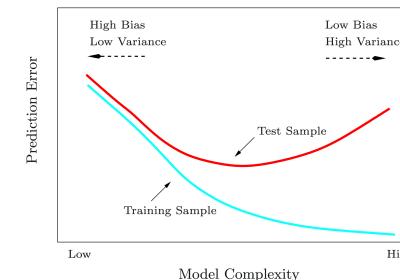
TRADE-OFFS

ILL-POSED PROBLEMS

In contrast to forward problems, inverse problems can have

- No solution
- Many solutions
- Unique but unstable solutions

These are called *ill-posed* (c.f. 'well-posed') problems
(Hadamard, 1902)



From Hastie et. al 'Elements of Statistical Learning: Data Mining, Inference and Prediction'.

available at: <http://web.stanford.edu/~hastie/ElemStatLearn/>. See also: 'An Introduction to Statistical Learning' (simplified version of above), available at: <http://www-bcf.usc.edu/~gareth/ISL/>

6

8

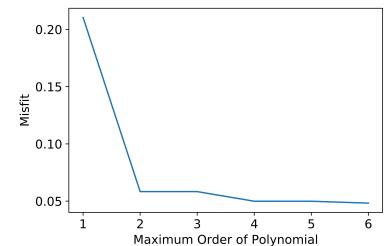
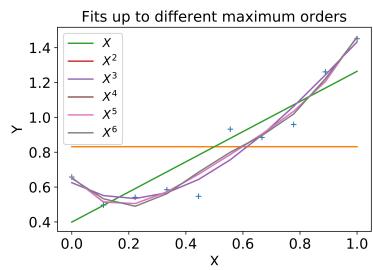
EXAMPLE: SIMPLE POLYNOMIAL FITTING

Suppose we have enzyme data and we *want to fit a curve to the double-reciprocal plot*

- Usually use linear (first order) regression (relates to MM)

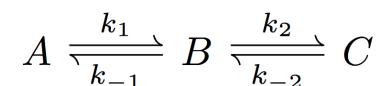
What if we tried to fit higher order polynomials?

EXAMPLE: SIMPLE POLYNOMIAL FITTING



EXAMPLE: SIMPLE REACTION

Suppose we have the reaction:



Exercise: write down the differential equations (assuming mass action)!

9

11

SIMPLE REACTION: SYNTHETIC DATA

Simulate 'synthetic' data assuming

$$k_1 = 0.6, k_2 = 0.4, k_{-1} = 0.02, k_{-2} = 0.01$$

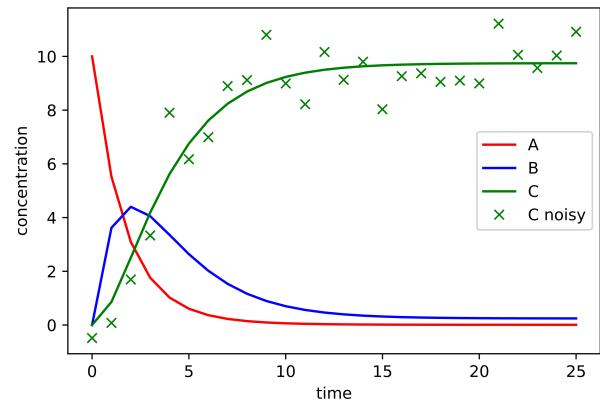
- Add Gaussian noise
- Assume can only measure $[C]$ (others unknown)

*Q: can we *recover* the original parameter values?*

10

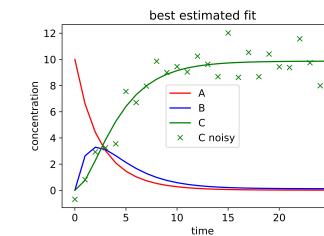
12

SIMPLE REACTION: SYNTHETIC DATA



13

SIMPLE REACTION: BEST FIT AND SIMPLEST FIT

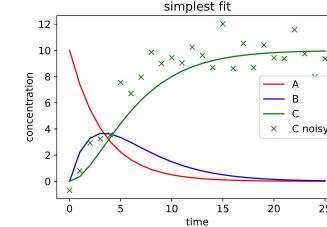


Simplest: $k_1 = 0.3, k_2 = 0.3, k_{-1} = 0.0, k_{-2} = 0.00$

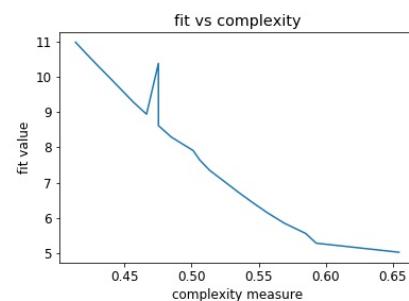
Best: $k_1 = 0.55, k_2 = 0.38, k_{-1} = 0.015, k_{-2} = 0.00$

True: $k_1 = 0.6, k_2 = 0.4, k_{-1} = 0.02, k_{-2} = 0.01$

15



SIMPLE REACTION: FIT VS COMPLEXITY CURVE



Note: not (quite) monotonically decreasing since get stuck in *local* minima in general.

14

COMMENTS

- Can *recover* parameters from synthetic data OK
 - Some better than others
 - *Real* data much noisier and 'true' model unknown/non-existent
- Parameters can *trade-off* against each other
 - E.g. maybe only combinations are identifiable
- Need *efficient algorithms* for solving minimisation problems!

In general, parameter estimation requires care and awareness of trade-offs involved

16

HOW TO DEAL WITH PARAMETERS FOR WHOLE-CELL MODELLING?

INTERFACE

rsif.royalsocietypublishing.org

Review



Cite this article: Babtie AC, Stumpf MPH. 2017 How to deal with parameters for whole-cell modelling. *J. R. Soc. Interface* 14: 20170237.
<http://dx.doi.org/10.1098/rsif.2017.0237>

Received: 30 March 2017
Accepted: 22 June 2017

How to deal with parameters for whole-cell modelling

Ann C. Babtie and Michael P. H. Stumpf

Department of Life Sciences, Imperial College London, London, UK

DOI: MPPHS, 0000-0002-3577-1222

Dynamical systems describing whole cells are on the verge of becoming a reality. But as models of reality, they are only useful if we have realistic parameters for the molecular reaction rates and cell physiological processes. There is currently no suitable framework to reliably estimate hundreds, let alone thousands, of reaction rate parameters. Here, we map out the relative weaknesses and promises of different approaches aimed at redressing this issue. While suitable procedures for estimation or inference of the whole (vast) set of parameters will, in all likelihood, remain elusive, some hope can be drawn from the fact that much of the cellular behaviour may be explained in terms of smaller sets of parameters. Identifying such parameter sets and assessing their behaviour is now becoming possible even for very large systems of equations, and we expect such methods to become central tools in the development and analysis of whole-cell models.

METABOLIC MODELS AS NETWORKS

Usually viewed as a *large network* of interacting pathways

Can reconstruct via genome sequencing

17

19

EXAMPLE: METABOLIC MODELS

Recall: *metabolism*

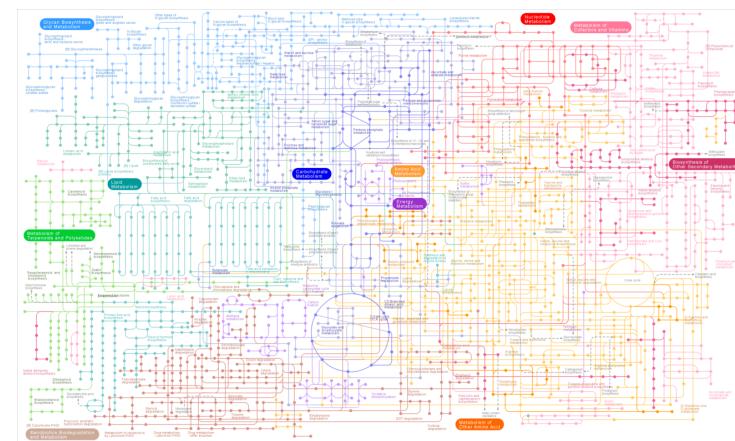
The consumption and production of chemical substances and energy to sustain life

- Catabolism: breakdown
- Anabolism: build up

Food → Life

18

EXAMPLE METABOLIC MODELS/NETWORKS

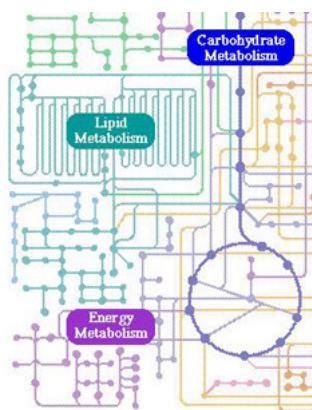


From: http://www.genome.jp/kegg-bin/show_pathway?map01100

20

EXAMPLE METABOLIC MODELS/NETWORKS

Zoomed in:



FLUX BALANCE ANALYSIS

More from Babtie and Stumpf (2017):

present in a system. In cell biology, for example, there are now numerous attempts at modelling aspects of metabolism, gene regulation and signalling at cellular level [17–24]. Perhaps the best established are metabolic models, where a powerful set of tools, based around *flux balance analysis* (FBA) [25], allows us to explore metabolic phenotypes *in silico* at a genomic level for an increasing range of organisms (and some individual cell types) [24,26,27]. However, such models are stoichiometric and thus give us information about biochemical reaction schemes and fluxes, but not details about the system dynamics.

21

23

APPROACHES

- *Full dynamic model* and search for smaller sets of important parameters (e.g. introduce complexity trade-off)
- See previous. *Difficult to scale up* to this size!
- *Try something else!*

MOTIVATING EXAMPLE

Consider

$$A \xrightleftharpoons[J_1^-]{J_1^+} B, B \xrightleftharpoons[J_2^-]{J_2^+} C, C \xrightleftharpoons[J_3^-]{J_3^+} A$$

or, equivalently in terms of *net* fluxes:

$$A \xrightarrow{J_1} B, B \xrightarrow{J_2} C, C \xrightarrow{J_3} A$$

where $J_1 = J_1^+ - J_1^-$ etc, and the arrows represent *net* fluxes

22

24

STOICHIOMETRIC MATRIX

MATRIX/VECTOR FORM

Derivation: see handout.

Result:

$$\frac{d\mathbf{C}}{dt} = \mathbb{S}\mathbf{J}$$

where \mathbf{C} is a *vector* of concentrations/metabolites, \mathbf{J} is a *vector* of fluxes and... \mathbb{S} is the...

...we write these as a *matrix*

$$\mathbb{S} = \begin{bmatrix} \beta_{11} - \alpha_{11} & \dots & \beta_{1N} - \alpha_{1N} \\ \dots & \beta_{ij} - \alpha_{ij} & \dots \\ \beta_{M1} - \alpha_{M1} & \dots & \beta_{MN} - \alpha_{MN} \end{bmatrix}$$

Note

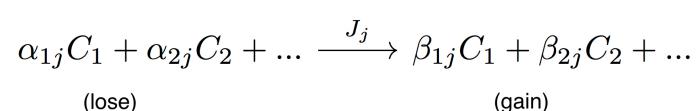
- Rows: species
- Columns: fluxes
- Entries: gain minus loss (stoichiometric coefficients)

25

27

STOICHIOMETRIC MATRIX \mathbb{S}

Given M species/metabolites and N reactions, each of the form



MATHEMATICAL FRAMEWORK OF FLUX BALANCE ANALYSIS

Rather than the dynamic model, we aim to solve the *steady-state* equation

$$\mathbb{S}\mathbf{J} = \mathbf{0}$$

for the vector of *fluxes* \mathbf{J} , *here treated as unknown*.

- No constitutive equations/no rate parameters involved here.
- We don't need to know the metabolite concentrations, just solve for fluxes

26

28

Biomeng 261 Lecture 8 :

Systems Biology cont'd

- Intro to parameter estimation for relatively simple dynamic models
- Scaling up to much larger systems using (steady state) Flux Balance Analysis (FBA).

Example paper (see Canvas).

INTERFACE

rsif.royalsocietypublishing.org

Review



Cite this article: Babtie AC, Stumpf MPH. 2017 How to deal with parameters for whole-cell modelling. *J. R. Soc. Interface* 14: 20170237. <http://dx.doi.org/10.1098/rsif.2017.0237>

Received: 30 March 2017
Accepted: 22 June 2017

How to deal with parameters for whole-cell modelling ↗ WCM.

Ann C. Babtie and Michael P. H. Stumpf

Department of Life Sciences, Imperial College London, London, UK

DOI: [MPHS, 0000-0002-3577-1222](https://doi.org/10.1098/rsif.2017.0237)

Dynamical systems describing whole cells are on the verge of becoming a reality. But as models of reality, they are only useful if we have realistic parameters for the molecular reaction rates and cell physiological processes. There is currently no suitable framework to reliably estimate hundreds, let alone thousands, of reaction rate parameters. Here, we map out the relative weaknesses and promises of different approaches aimed at addressing this issue. While suitable procedures for estimation or inference of the whole (vast) set of parameters will, in all likelihood, remain elusive, some hope can be drawn from the fact that much of the cellular behaviour may be explained in terms of smaller sets of parameters. Identifying such parameter sets and assessing their behaviour is now becoming possible even for very large systems of equations, and we expect such methods to become central tools in the development and analysis of whole-cell models.

State of the art :

techniques, and instead start to require computer simulations to explore their behaviour.

1.1. From simple to complex models

Some models aim to capture the essential hallmarks of life—such as metabolism, nutrient uptake, gene expression regulation and replication—but in a simplified representation that does not aim to replicate the true complexity of a whole organism [11–16]. These *coarse-grained* models have shown great promise and allow us to integrate molecular, cellular and population level/scale processes into a coherent—and analytically tractable—modelling framework. While real cells will be much more complicated, these simple model systems have successfully provided insight into fundamental cellular physiology, e.g. processes affecting microbial growth rates [12,13,15,16].

Increasingly there is interest in generating more realistic and complicated models that, rather than aiming to provide abstract representations of key features, incorporate extensive details of known components and interactions (or reactions) present in a system. In cell biology, for example, there are now numerous attempts at modelling aspects of metabolism, gene regulation and signalling at cellular level [17–24]. Perhaps the best established are metabolic models, where a powerful set of tools, based around flux balance analysis (FBA) [25], allows us to explore metabolic phenotypes *in silico* at a genomic level for an increasing range of organisms (and some individual cell types) [24,26,27]. However, such models are stoichiometric and thus give us information about biochemical reaction schemes and fluxes, but not details about the system dynamics.

Advances in both high-throughput experimentation and computational power have opened up the possibility of creating and analysing more complex dynamic models of biological systems, including many which represent processes occurring at different scales [28,29]. Numerous models now face the challenge of being *large* (in terms of numbers of species and parameters represented), multi-scale and/or *hybrid* in nature (incorporating multiple different mathematical representations) [23,28–30]. The most ambitious models to date—the WCMs—aim to provide faithful *in silico* representations of real biological cells, including all major cellular processes and components, and are both very large scale and hybrid (figure 1) [31–33]. There are several potential uses for such WCMs:

- (1) To gain mechanistic insights, by serving as an *in silico* ‘blueprint’ through which we study the behaviour of real cells.
- (2) As a rational screening and predictive tool, to explore *in silico* what might be hard or impossible to study *in vivo*.
- (3) To drive new biological discoveries, by showing where we lack sufficient understanding, and identifying promising future directions to pursue experimentally.
- (4) To study *emergent* phenomena which are only apparent when we consider a system as a whole.
- (5) To integrate heterogeneous datasets and amalgamate our current knowledge into a single modelling framework.
- (6) Perhaps eventually to study, via virtual competitions between different cell architectures, evolutionary dynamics in unprecedented detail (but at enormous, currently crippling, computational cost).
- (7) In the meantime, as the community strives to develop viable WCMs, the technological, computational and

Here, we focus on the inference and statistical modelling challenges inherent to developing WCMs (and other complex models), in terms of model construction, parameter estimation, uncertainty and sensitivity analyses, and model validation and refinement. Some of these are generic modelling challenges—but worth reiterating—while others are specific to large-scale, multi-scale and hybrid models.



State of art cont'd.

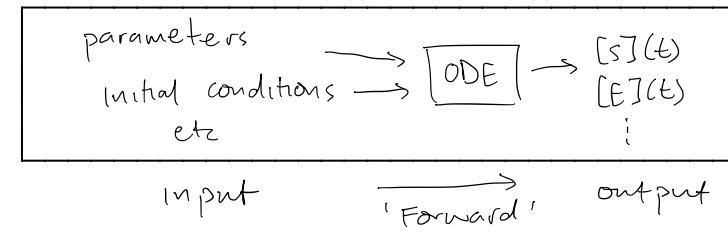
At present, none of the statistical inference methods outlined above are applicable at the scale of WCMs. However, smaller subsystems, such as individual pathways, regulatory motifs, receptor complexes or systems comprising small sets of metabolic reactions and the associated regulatory processes can be effectively parametrized using such methods [72]. For such systems, we can often estimate parameters, including uncertainty; and we are frequently able to assess parameter sensitivity (typically measured as the change in some model output, e.g. predicted protein abundance, in response to varying a single parameter). In some cases, experimental measurements of species concentrations may allow us to effectively decompose our models into smaller modules for efficient parameter estimation [73]. Bayesian inference methods in particular are limited in terms of scale and are generally only feasible for models with up to tens to hundreds of species and parameters [74,75]. Some optimization approaches are much more scalable though, with recent advances allowing parametrization of ODE models comprising hundreds to thousands of species and parameters [65,76,77]. As always, however, the chance of being trapped in local optima is high for such large-dimensional problems.

A combination of both inference and experimental estimation will probably be needed to parametrize complex biological models. It is currently impractical to use inference techniques within the context of a full WCM, unless considering very small pre-defined subsets of the parameters and, even then, the computational costs are enormous [67]. We can, however, make use of scalable inference techniques [78] to help us parametrize the component submodels, using experimental information where available as prior knowledge for the inference procedures. This will allow us to avoid some of the potential pitfalls outlined above of experimental estimates, and generate parameter estimates that take into account—to the best of our ability—the influences of cellular and system context, and make use of the most appropriate *in vivo* datasets. Crucially, rigorous statistical inference also enables us to explore the relationships between model parameters and start to understand and quantify the uncertainties inherent to any mathematical model.

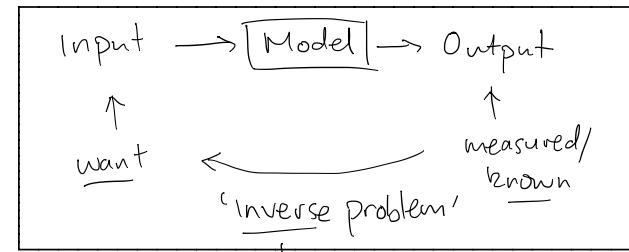
- Let's do
- estimation for 'simple' } today.
models
 - scaling up to large } today/
systems via FBA. } tomorrow

Parameter estimation : Forward vs inverse problems

Usual modelling procedure :

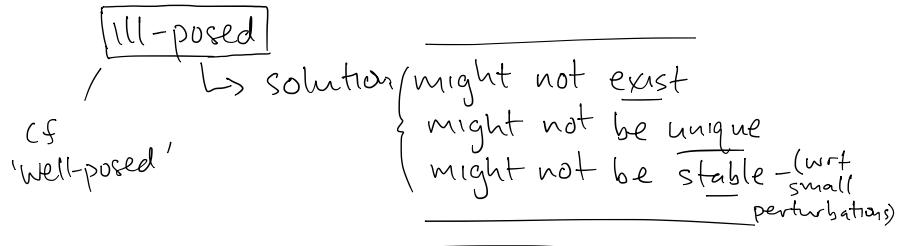


Problem : in reality we have data like (noisy) measurements of S, E concentrations (outputs) while we don't know inputs eg k_1, k_2 etc, ie



What's the problem?

→ in contrast to 'forward' problems,
inverse problems are usually:



Simple example

$$y = f(x) \quad \text{eg } y = x^2 \quad \left. \begin{array}{l} \text{unique output} \\ \text{for each input} \end{array} \right\}$$

$$x \mapsto x^2$$

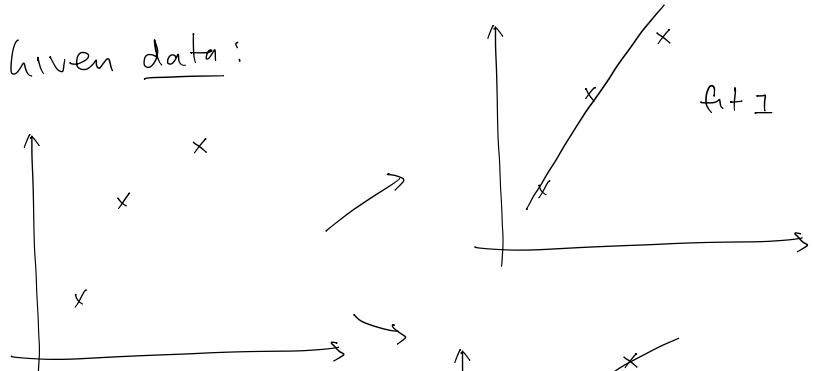
- Observe output $y = 4$
 - What was input x ?
- no unique solution!

Solution set: $\{-2, 2\}$

Importantly:
parameter fitting is 'ill-posed'

Illustration: curve fitting

Given data:



Which is better?

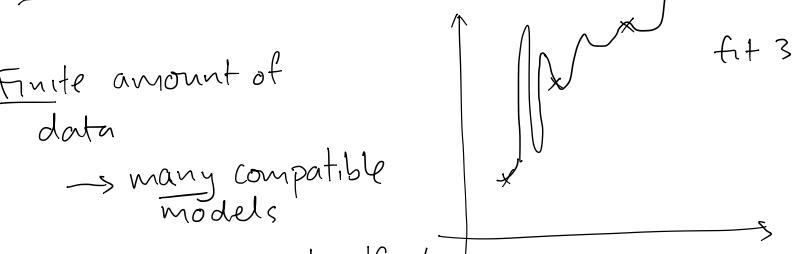
Finite amount of data

→ many compatible models

→ more complex 'fits'
given better

↳ but are often unstable
&/or fit future data worse

→ simpler ≈ more understandable?



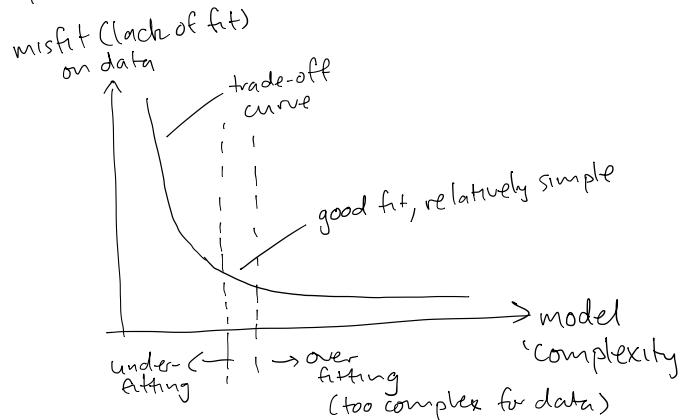
Approach: Trade-offs are key.

* Key Page *

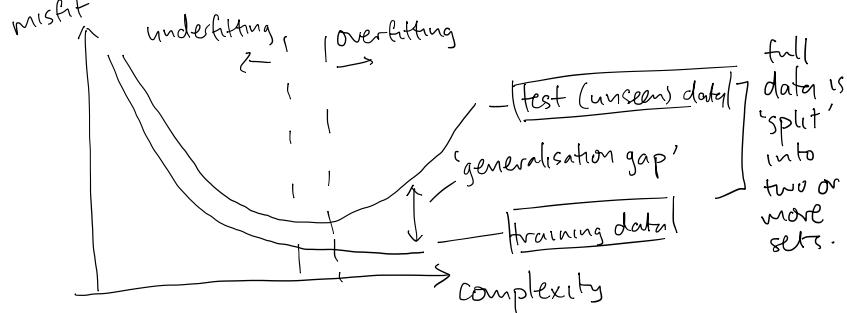
- (A) - make more complex } conflicting
- (B) - make simpler } → you need to decide!

- (A) - Fit to given ('training') data } 'predictive' or 'empirical' view
- (B) - stability / fit to future ('test') data } 'unseen'

Simple vs complex:



Training/Test ('Predictive' or 'machine learning' view)

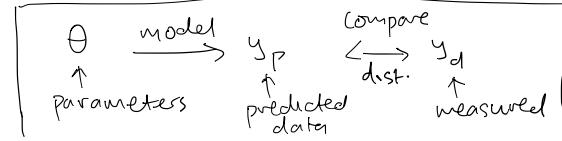


Measuring data fit / misfit (e.g. non-negative, symmetric etc.)

do in lab

- Need a norm, distance, metric, cost function etc relating data and model

- Evaluates model by its predicted data



e.g. $d(\underline{y}_p, \underline{y}_d)$ (e.g. vector $\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \end{pmatrix}$)

& $\underline{y}_p = M(\theta)$ leads to

$'\text{misfit}'(\theta; \underline{y}_d) = d(\underline{y}_d, M(\theta)) = d(\underline{y}_d, \underline{y}_p(\theta))$

how good are parameters θ at 'predicting' data \underline{y}_d

Typical 'distances' or 'cost' measures:

$$d(\underline{x}, \underline{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \leftarrow \sqrt{\text{sum of squares}}$$

$$d(\underline{x}, \underline{y}) = \sum_{i=1}^n (x_i - y_i)^2 \leftarrow \text{sum of squared diff's.}$$

$$d(\underline{x}, \underline{y}) = \sum_{i=1}^n |x_i - y_i| \leftarrow \text{sum of absolute differences}$$

Measuring 'complexity'?

→ Harder!

- Simple idea is to use a 'norm' (size) of model/parameters or 'distance' from a 'default' or 'null' model
- clearer for linear models (less so for nonlinear)

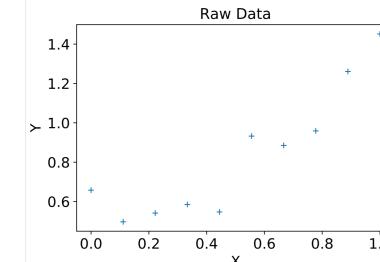
$$\text{Eg } \|\theta\|_d = d(\theta, 0) \stackrel{\text{e.g.}}{=} \sum_{i=1}^n \theta_i^2$$

'size' or
'complexity'
of parameters 'distance from' or
'cost relative to'
zero (or other ref)

(Machine learning:
see 'VC dimension')

Examples

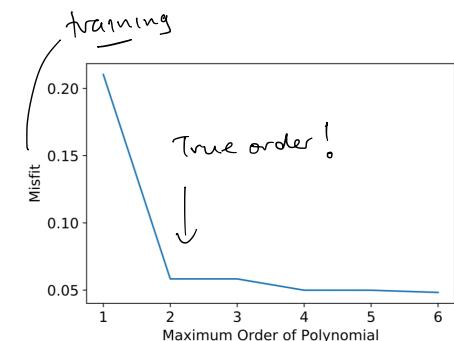
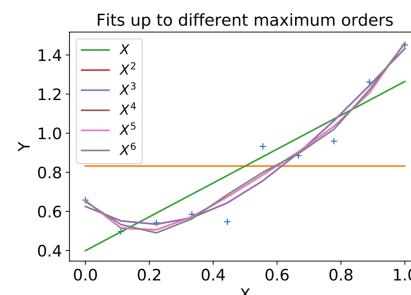
→ Polynomial:



True model

$$Y = X^2 + 0.5 + \text{noise}$$

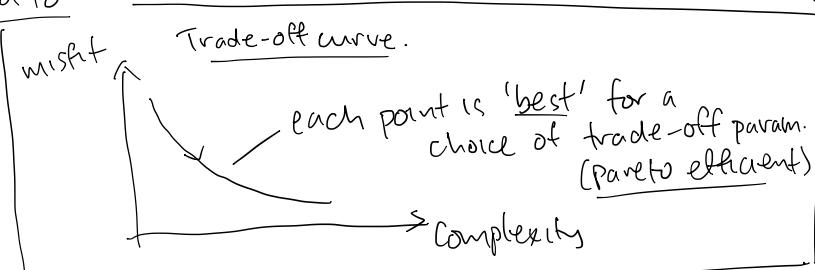
Fits & trade-offs:



Combining: optimal trade-offs

- minimise complexity
subject to acceptable data fit
 - minimise data misfit
subject to acceptable complexity
 - minimise data misfit + model complexity
cost
cost
- can show are equivalent
→ need to choose relative importance (tradeoff parameter) though.

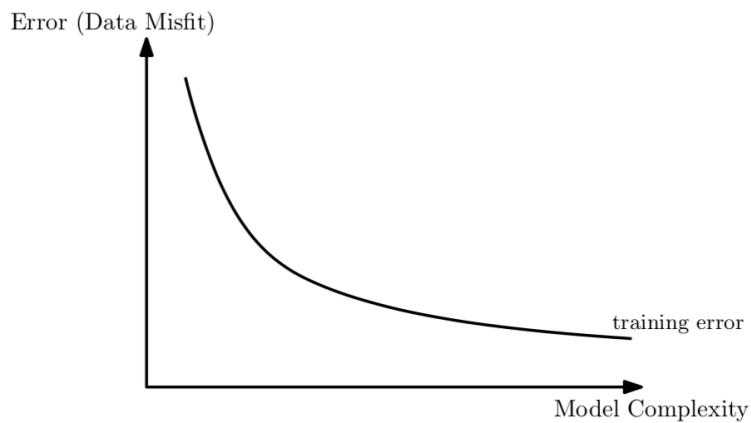
All lead to:



ODE? see slides

Example test question:

- c) Given the following figure illustrating the training error vs. model complexity for a family of models fitted to a given dataset, add
- A curve illustrating what you would expect the *test error* to look like as a function of model complexity
 - A vertical line indicating what would be a 'good' choice of model complexity to use, given your answer to i.



Flux Balance Analysis (Part I).

◦ Motivation & formulation.

→ common for metabolic models

Recall : Metabolism (cellular)

'all the chemical processes
keeping the cell alive'

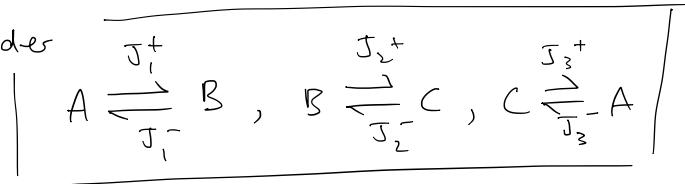
- large systems! (eg 100-1000s of ODEs, each with many parameters)
 - ↳ can be genome scale!
 - ↳ many interacting pathways'
- naive parameter estimation doesn't work well
- dynamic data often not available

Alternative trade-off :

- i.e. mass balance → {
 - fluxes only
 - stoichiometric matrices
- but no constitutive equations
- (◦ can include some overall thermodynamic info)

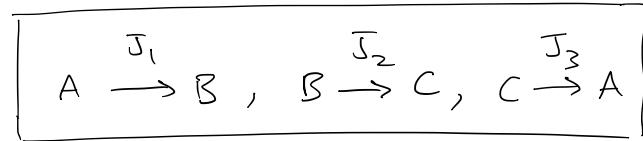
FBA: simple motivating example.

Consider



In terms of net fluxes, where $\bar{J}_i = J_i^+ - J_i^-$ etc

can write



where $\xrightarrow{\bar{J}_1}$ represents $\xrightleftharpoons[\bar{J}_1^-]{\bar{J}_1^+}$ etc.

Conservation of mass is then

$$\boxed{\begin{aligned} \frac{d[A]}{dt} &= -J_1 + J_3 \\ \frac{d[B]}{dt} &= +J_1 - J_2 \\ \frac{d[C]}{dt} &= +J_2 - J_3 \end{aligned}}$$

FBA: simple motivating example.

Now, we write as matrix/vector equation:

$$\bar{c} = \begin{bmatrix} [A] \\ [B] \\ [C] \end{bmatrix}, \quad \frac{d\bar{c}}{dt} = \begin{bmatrix} d[A] \\ \frac{d[B]}{dt} \\ d[C] \end{bmatrix}$$

$$\bar{J} = \begin{bmatrix} \bar{J}_1 \\ \bar{J}_2 \\ \bar{J}_3 \end{bmatrix}$$

↳ also use v_i & \bar{v} instead of J_i & \bar{J} sometimes

Gives: (overbar: vector
underbar: matrix)

$$\frac{d\bar{c}}{dt} = \begin{bmatrix} -1 & 0 & +1 \\ +1 & -1 & 0 \\ 0 & +1 & -1 \end{bmatrix} \bar{J}$$

S matrix

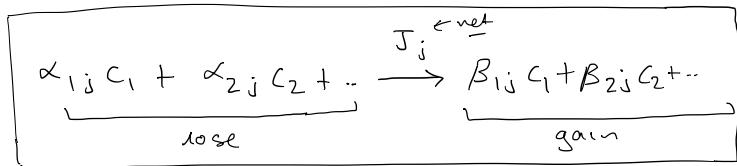
$$\text{i.e. } \boxed{\frac{d\bar{c}}{dt} = S \bar{J}}$$

Note: we are not using a constitutive equation $\bar{J} = f(\bar{c})$ so we do not have a 'closed' system of eqns

$$\frac{d\bar{c}}{dt} = F(\bar{c})$$

FBA: Stoichiometric matrices S

Given M species & N reactions, each of the form



we define the Stoichiometric matrix

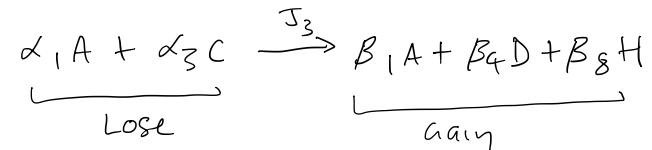
$$S = \begin{matrix} & c_1 & c_2 & \cdots & c_i & \cdots & c_M \\ & \begin{matrix} J_1 & \cdots & J_j & \cdots & J_N \end{matrix} & & & & & \\ & \begin{bmatrix} \beta_{11} - \alpha_{11} & \cdots & \beta_{1N} - \alpha_{1N} \\ \vdots & \ddots & \vdots \\ \beta_{ij} - \alpha_{ij} & \cdots & \beta_{NN} - \alpha_{NN} \end{bmatrix} & & & & & \\ & \text{is just to help fill in matrix} & & & & & \\ & \text{not part of it!} & & & & & \end{matrix}$$

Note $+\beta$, $-\alpha$

→ sign is determined by choice of sign for net flux.

FBA: Stoichiometric matrices

Example: 3rd reaction in some system is



Then:

$$S = \begin{matrix} & J_1 & J_2 & J_3 & \cdots \\ A & \beta_1 - \alpha_1 & 0 & 0 & \cdots \\ B & 0 & \cdots & \cdots & \cdots \\ C & -\alpha_3 & \cdots & \cdots & \cdots \\ D & \beta_4 & 0 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ H & 0 & \cdots & \cdots & \cdots \end{matrix}$$

↑
3rd reaction: 3rd column

As mentioned, we are working in terms of pure fluxes J (ie conservation of mass)

- no constitutive eq's
- no rate parameters
- no 'closure' eg $\frac{dc}{dt} = F(c)$.

FBA: Steady states

To avoid needing constitutive equations &/or needing dynamic data, in FBA we usually

- just consider steady states
- treat fluxes as unknowns to determine.

→ makes sense when thinking about eg overall metabolism & homeostasis
 — hence popular in this area

→ often have eg 'metabolic network' maps available to help build models

So : New goal!

Solve

$$\underline{S} \bar{J} = \bar{O}$$

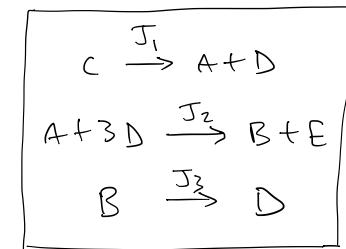
for fluxes \bar{J}

↑
not rate
constants.

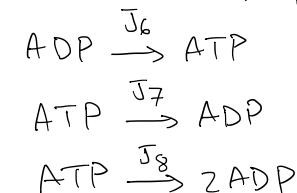
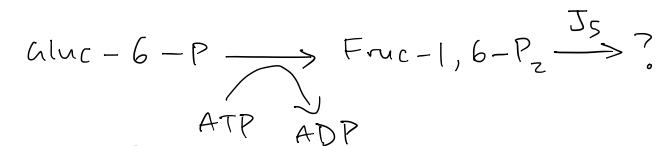
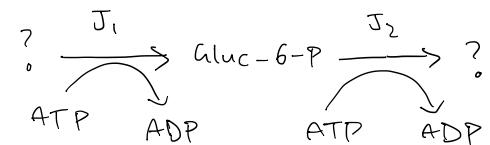
Tomorrow: FBA
 cont'd!

FBA: Exercises!

A. Determine \underline{S} for :



B. Determine \underline{S} for the system :



Notes:

- $\begin{array}{c} A \xrightarrow{J} C \\ \downarrow \\ B \xrightarrow{J} D \end{array}$ is just $\begin{array}{c} A+B \xrightarrow{J} C+D \end{array}$ concentrations
- we allow unseen metabolites via '?'
 → no row in \underline{S} , but still include reaction col!