

BIOMENG 261

TISSUE AND BIOMOLECULAR ENGINEERING

Module I: Reaction kinetics and systems biology

Oliver Maclarens

oliver.maclarens@auckland.ac.nz

LECTURE 12: GENE EXPRESSION DATA AND GRNS

- *Larger systems*
 - Gene space and gene regulatory networks (GRNs)
- *Brief overview of microarray data*
 - Experiment types
 - Data organisation and expression matrices
- *Analysis types*
 - Clustering, distance matrices and dendrograms
 - Control analysis and regulatory matrices

Note: there are many images stolen from the internet in what follows...

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.

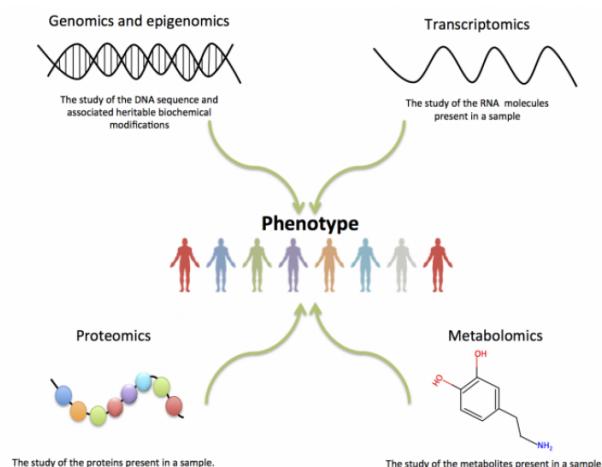
2. Systems biology I: signalling and metabolic systems [3 lectures]

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.

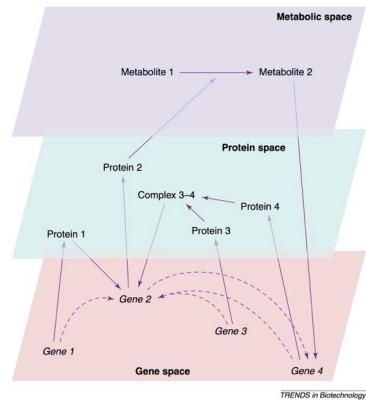
MUCH LARGER SYSTEMS - 'OMICs'



2

4

GENE SPACE



See: Brazhnik et al. (2002) 'Gene networks - how to put the function in genomics' (on Canvas)

TRANSCRIPTOMICS

- A subfield of *functional genomics*
 - Functional genomics: study of how genes and intergenic regions contribute to biological function
 - The focus is on *gene expression*
 - In particular, via *measuring mRNA* (the transcripts)

See: Lowe et al. (2017) 'Transcriptomics technologies' (on Canvas)

EXPRESSION ANALYSIS

- *Microarrays*
 - Mature technology
 - Relatively well-established data analysis methods
 - *RNA-seq*
 - Newer technology, rapidly overtaking microarrays
 - Less standardisation of analysis methods
 - Much more computationally/storage intensive

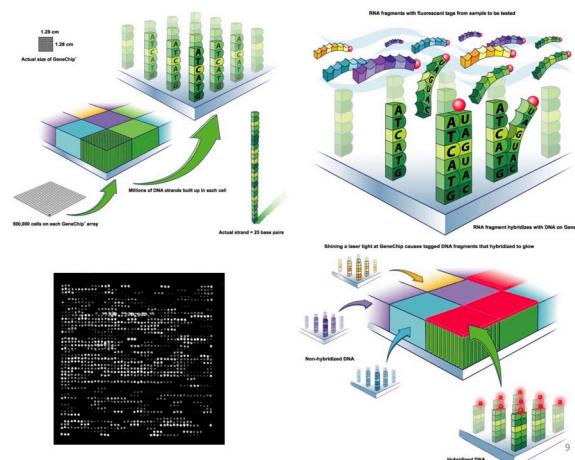
But: *microarrays still relevant and useful*: we will consider these (easier and better understood)

MICROARRAYS



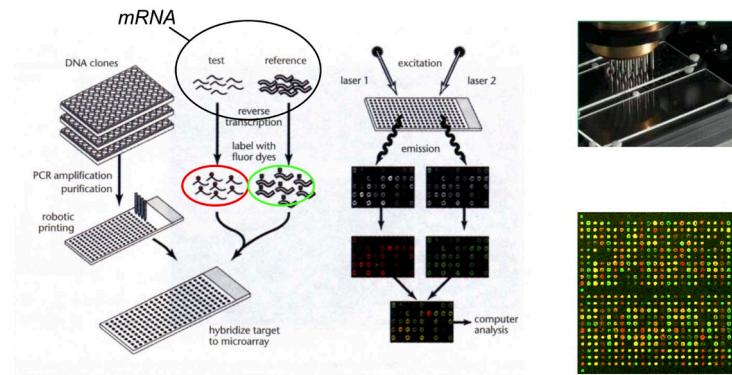
For a video intro: see e.g. <http://www.youtube.com/watch?v=VNsThMNIKhM>

MICROARRAYS



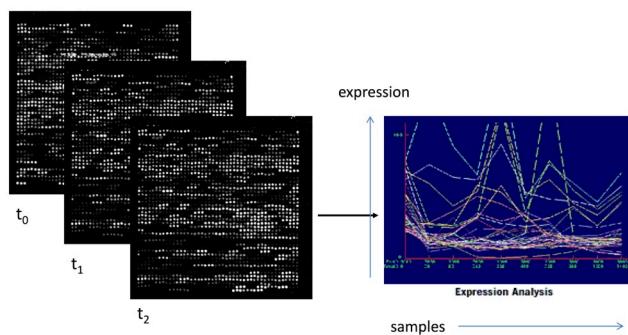
9

MICROARRAYS: COMPARATIVE EXPRESSION



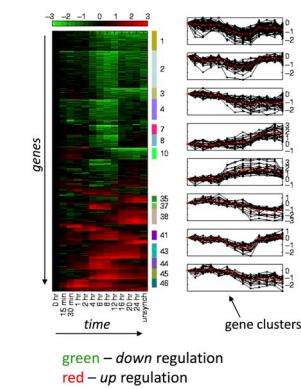
11

MICROARRAYS: TIME SERIES



10

MICROARRAYS: RELATIVE EXPRESSION OVER TIME



12

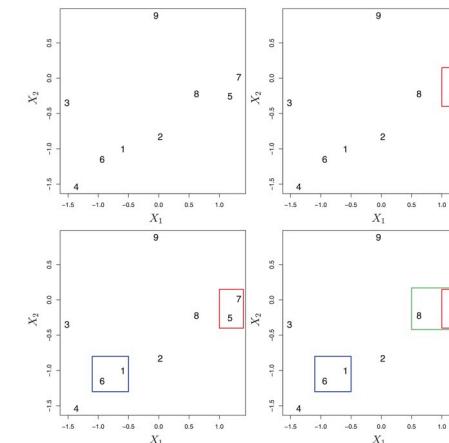
DATA ANALYSIS: STATISTICAL/MACHINE LEARNING

Clustering, unsupervised and supervised learning etc. see:

- James et al. 'An Introduction to Statistical Learning'
 - Available at: <http://www-bcf.usc.edu/~gareth/ISL/>
- Hastie et. al 'Elements of Statistical Learning: Data Mining, Inference and Prediction'
 - Available at:
<http://web.stanford.edu/~hastie/ElemStatLearn/>

13

HIERARCHICAL CLUSTERING EXAMPLE (JAMES ET AL.)



15

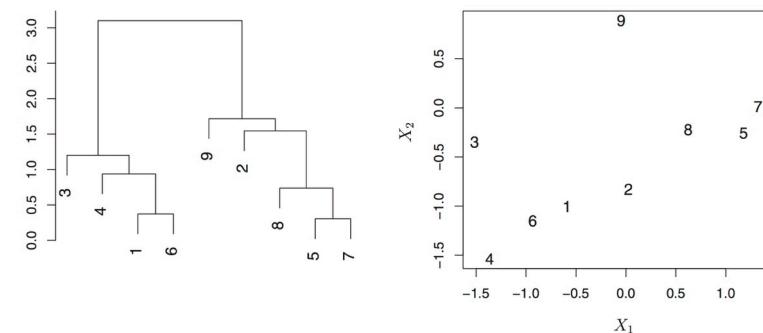
CLUSTERING

- An *unsupervised learning* method for *pattern discovery*
- Two popular algorithms are
 - K-means
 - Hierarchical clustering

See James et al. Chapter 10 for detailed algorithms. We will look at *hierarchical clustering* here.

14

HIERARCHICAL CLUSTERING: DENDROGRAMS



16

PERTURBATION APPROACH FOR INFERRING REGULATORY MATRICES/NETWORKS

- Perturb *transcription rates* for *each gene in turn*
- Measure changes in *steady-state expression levels* for all genes (including self)
- Gives indication of underlying *regulatory network*
- Summarise in *regulatory strength matrix or network diagram*.

See de la Fuente et al. (2002) 'Linking the genes' (on Canvas)

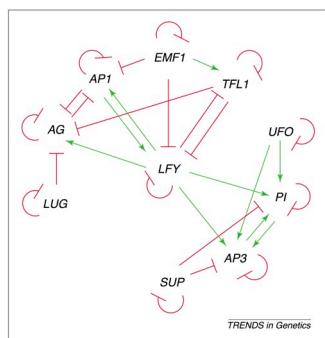
17

PERTURBATION APPROACH FOR INFERRING REGULATORY MATRICES/NETWORKS

Example: flower morphogenesis (see de la Fuente et al. 2002 for details):

LUG	AG	API	EMF1	TFL1	LFY	SUP	ASP	PI	UFO	LUG
-1	0	0	0	0	0	0	0	0	0	0
-0.579	-1.14	-0.184	-0.002	-0.12	0.114	0	0	0	0	AG
-0.009	-0.894	-1.14	-0.109	-0.026	0.124	0	0	0	0	EMF1
0	0	0	0.094	-1.09	-0.979	0	0	0	0	TFL1
-0.002	-0.005	0.053	-0.103	-0.107	-0.09	0	0	0	0	LFY
0	0	0	0	0	0	0	0	0	0	SUP
0	0	0	-0.001	-0.001	-0.138	-0.119	-1.04	-1.04	0.109	ASP
0	0	0	-0.001	-0.001	0.138	-0.119	0.088	0.088	0.109	PI
0	0	0	0	0	0	0	0	0	0	UFO

TRENDS in Genetics



18

EXAMPLE PAST QUESTIONS

- (b) DNA microarrays are used to simultaneously measure the expression of many different genes in a sample. Explain briefly the difference between *time series* and *comparative* microarray measurements. (4 marks)
- (c) In a series of experiments, the amounts of mRNA for different genes are perturbed and the changes in mRNA of all genes of interest (namely A, B and C) are measured. A set of 'co-control coefficients' was calculated and organised into the 'Regulatory Strength' matrix, R_d , given below

$$R_d = \begin{bmatrix} A & B & C \\ -0.5 & 0 & 0 \\ 0.1 & 0 & 1.2 \\ 0.8 & -0.1 & 1 \end{bmatrix}$$

Using the information in R_d sketch a qualitative regulatory network showing how each gene regulates the expression of all genes (including itself). Use arrows (\rightarrow) to show positive regulation and blunt arrows ($\overline{\longrightarrow}$) to indicate negative regulation. (5 marks)

19

EXAMPLE PAST QUESTIONS

- 5) a) In a microarray experiment, the expression of 8 genes was measured as a function of time and the data were analysed to create the following dendrogram:



State if the following statements are TRUE, FALSE or INDETERMINATE.

i) Genes C and I show similar expression patterns at the different time points.

(1 mark)

ii) The Euclidean distance between genes C and I is less than the Euclidean distance between genes I and B.

(1 mark)

20

Biomeng 261 Lecture 12

- Large genetic systems
 - 'Gene space' & GRNs
(genetic regulatory networks)
 - Expression analysis methods
- Takeaways: know & interpret:
- Matrix types
 - ↳ Gene expression matrix } main
 - ↳ gives { distance matrix ①
 - ↳ regulatory strength matrix ②
 - Dendograms
 - ① ↳ summarise distance matrix / clustering algorithm
 - Regulatory networks
 - ↳ summarise regulatory strength matrix ②

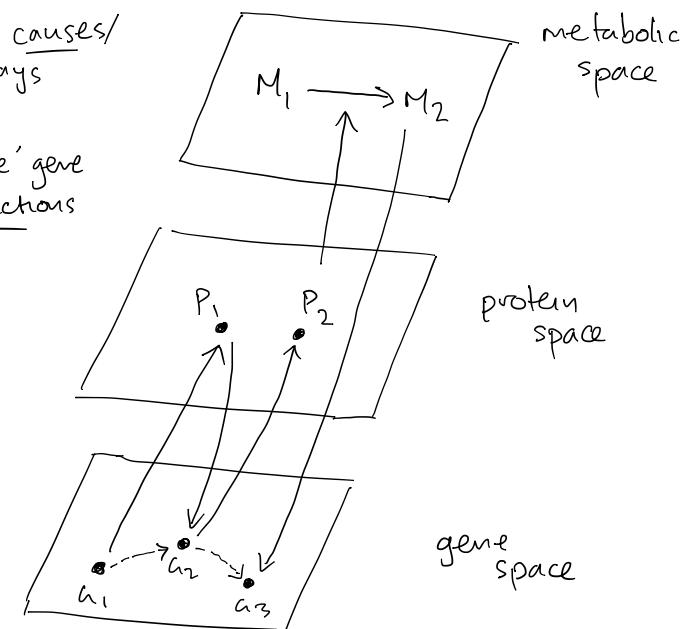
Background: Gene space

- A way of 'projecting all the action' into interactions between genes

Arrow types:

→ 'actual' causes/
pathways

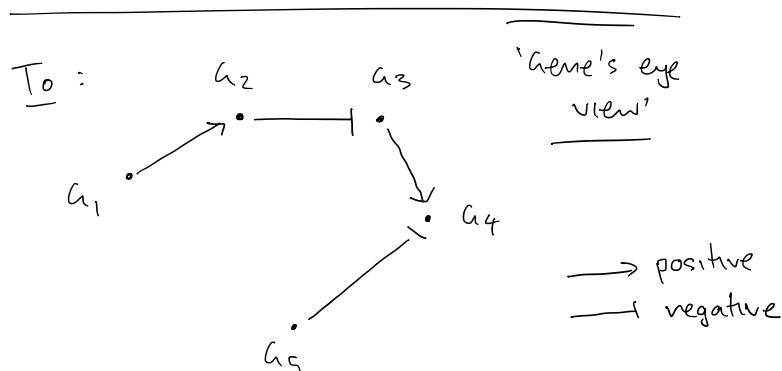
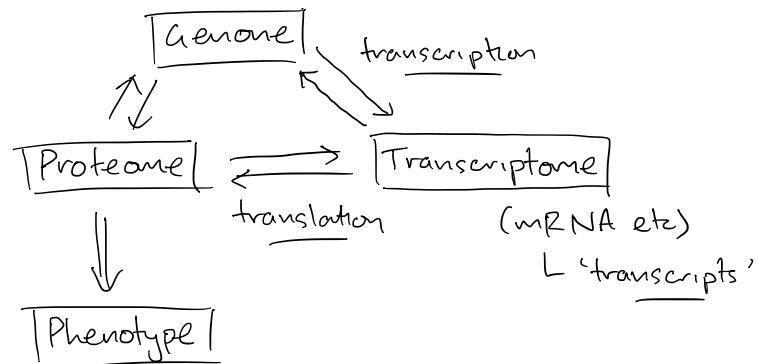
→ 'effective' gene
interactions



(Background)

Leads to : Gene Regulatory Networks

From :



Warning:
we also loosely
used this term
for 'genetic
regulation'
models in gene
networks - gene's eye
view.

"Gene regulatory network"
(GRN)

(Background)

Transcriptomics

('omics': as a whole)

- Subfield of 'functional genomics'

[The study of how
genes & intergenic
regions contribute to
biological function]

- Transcriptomics focuses on
gene expression levels
as measured by the levels
of the Transcripts (mRNA)
associated with genes

↳ mRNA typically easier to
measure etc than proteins,
but see 'proteomics'!

Key idea

under
'treatment':

- does expression go up or down?
- do groups of genes go up or down together?

Expression Analysis

Two key approaches:

- Microarrays

- L uses 'probes' (eg cDNA)
complementary DNA
- L samples 'hybridise' to probes complementary to probes
- L amplify & quantify using qPCR
(PCR: polymerase chain reaction)

- RNA-seq

- L direct sequencing of transcripts
- L 'next gen' high-throughput sequencing

We will focus on microarrays!

- well understood
 - more mature, easier to analyse
 - still used / still useful
- but RNA-seq is 'rapidly overtaking'

Microarrays: (more) background (see slides & videos etc)

Idea: want to measure mRNA levels (hence gene expression)

- Microarrays can measure 1000s of mRNA levels at a time

↳ i.e. 1000s of genes at a time

→ consist of a grid of 'spots' containing cDNA (complementary DNA)

→ mRNA samples preferentially bind to ('hybridise' with) corresponding cDNA

↳ 'reverse transcription'!
(DNA ← mRNA)

Also:

- use fluorescent tags to distinguish samples
- amplify levels of products via PCR/qPCR.

Preprocessing

Preprocessing is crucial

BUT : we won't go into here

Typical considerations :

- artifacts
- absolute vs relative expression levels
- log transformations
- etc

usually work with

$$\boxed{\log(\text{relative expression})}$$

From now: take measure as 'given'
& just use numbers.

Data & Experiment Types

We consider two types of experiment

1. Perturbation (or comparative)

- ↳ effect of treatment
 - ↳ different tissue types etc
- } multiple sample types, same time

2. Time Series

- ↳ same cell/tissue etc
 - studied over time
- } one sample type, multiple times

Idea : o time has natural order

- ↳ continuous or categorical ordinal

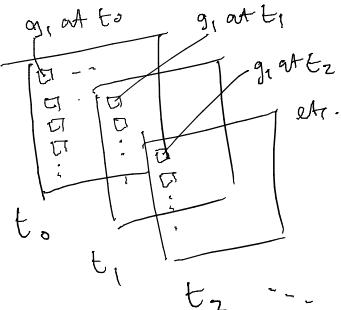
o treatment/class not necessarily ordered

- ↳ cancer or not etc

[BUT] essentially

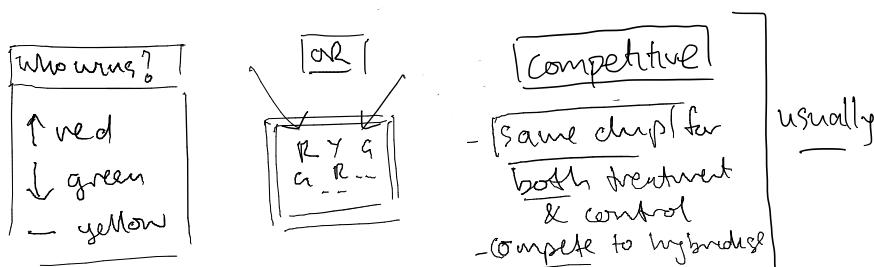
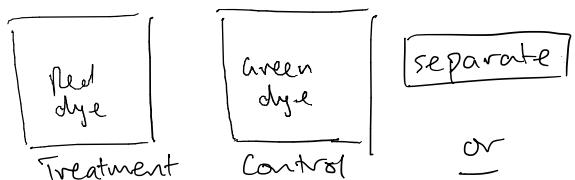
- just different independent vars,
same basic ideas
- also, often want to combine,
eg compare two time series
from different treatments
at same times

Time Series



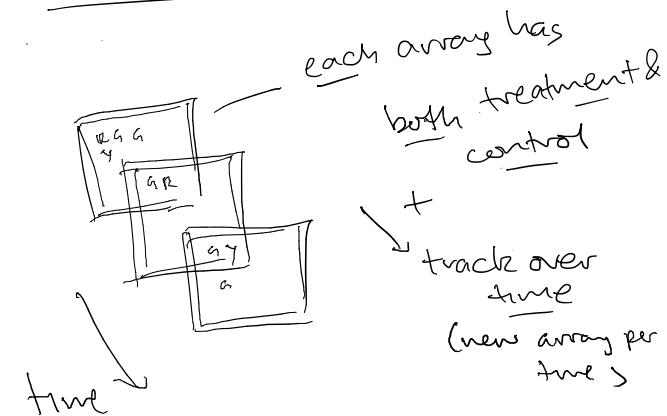
new chip for each time

Perturbation/comparative



↳ direction of treatment relative to control

Combined



Idea: genes that co-express (correlated expression) are potentially co-regulated.

Two key goals

- try to cluster (group) similar
- gene-gene effects
 - ↳ perturb each gene and see effects

But first →

Gene expression matrices

A general format that we can put both time series & perturb./comparative into.

	Exp 1	Exp 2	Exp 3
Gene 1			
Gene 2			
Gene 3			
:			
:			

Exp:
experiment
(not expression)

- e.g. exp 1: time t_0 , control $\}^{experiment}$
- exp 2: time t_1 , control
- :
- exp n: time t_n , treatment
- exp n+1: time t_{n+1} , treatment
- etc.

(warning: actually the transpose of many 'data array' formats
e.g. R data frames)

Gene expression matrices

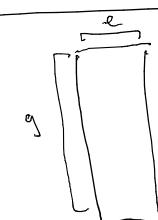
Column: an experiment, i.e. an array

(gene chip/
microarray)

	Exp 1	Exp 2	Exp 3
Gene 1			
Gene 2			
Gene 3			
:			
:			

Row: a gene expression profile for a given gene

Typically: \gg genes \gg experiments



- o many more genes than experiments
- o problem for statistical inference!
- ⇒ focus on smaller sets or clusters of genes

Some typical analyses

1. Clustering: Similarity analysis
(eg over time)

2. Perturbations: Linear control
based analysis

Clustering

- a form of unsupervised learning for pattern discovery
- Need notion of 'distance', however.
↳ Recall prev. lectures.

Note: more precise distinctions

Distance

- i) $d(x, y) \geq 0$
- ii) $d(x, y) = d(y, x)$
- iii) $d(x, x) = 0$

Metric

- i-iii)
- &
- iv) $d(x, y) = 0$ iff $x = y$
- v) $d(x, y) + d(y, z) \geq d(x, z)$

Dissimilarity

- i-iii)
- &
- iii) $d(x, y) \neq d(x, y)$
increases monotonically as $x \neq y$
more dissimilar (subjective)

(Terminology)

Unsupervised? Note on 'learning types'

Supervised

$$X \rightarrow Y$$

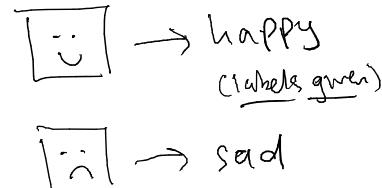
Learn function $X \rightarrow Y$

Unsupervised

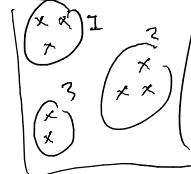
$$X \hookrightarrow X$$

Pattern discovery in X

↳ Train: supervisor examples



Eg: Find clusters



group 'similar'

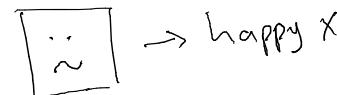
→ no labels given

→ goal: discover labels/clusters etc

→ do need 'distance' / 'similarity'
 $d(x_1, x_2)$

→ Can do some eval/val.
eg via consistency/stability
on training/test sets
↳ but in general, harder to validate

Test: predict on new set



'Reasonably straightforward' to evaluate/validate.

Example : clustering expression profiles
→ across time/exp - (profiles)

expression
matrix

	1	2	3	...	8
gene A	-1	-2	2	...	-2
gene B	-1	-2	-1	-	-
gene I	-1	-2	0	-	-

which genes have 'similar'
profiles?

Distance ? Here : Euclidean (for simplicity)



Distance between two profiles

$$d(A, B) = \sqrt{\sum_{i=1}^n (y_i^{(A)} - y_i^{(B)})^2}$$

Euclidean:
square root of
sum of
square
diff.

where
n : number of experiments,

$y_i^{(A)}$: expression level of
gene A in experiment
i

Example :

A	1	1	-1	1	2	1
B	1	1	1	1	1	2
	↓	↓	↓	↓		

squared
diffs: 0² 0² 3² 2² sum squares
 i.e. 0 0 9 4 [13]

$$\text{so } d(A, B) = \sqrt{13}$$

Distance matrices

We can summarise all pairwise differences in a

+ distance matrix | note: is
(not an expression matrix)

	A	B	C	D	..
A	$d(A, A)$	\dots	$d(A, D)$		
B		$d(B, A)$			
C			$d(C, A)$	$d(C, C)$	\ddots
D				$d(D, A)$	\dots

[ie $d(\text{from}, \text{to}) = d(\text{row}, \text{col})$]

Note: • symmetric so don't need upper part i.e

$$d(B, A) = d(A, B)$$

• diagonals are zero

Distance-based clustering

Goal: group together of close.

+ Two popular algorithms |

o K-means

o hierarchical

see over → for pseudo code

Here: consider hierarchical



Hierarchical clustering

Begin with n observations

Find pairwise distances

For $i = n, n-1, \dots, 2$:

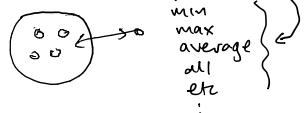
Find smallest distance

'Fuse' together (^{so:} $n \rightarrow n-1$
Observ.)

Recompute distances to
'fused' cluster

\hookrightarrow need notion of

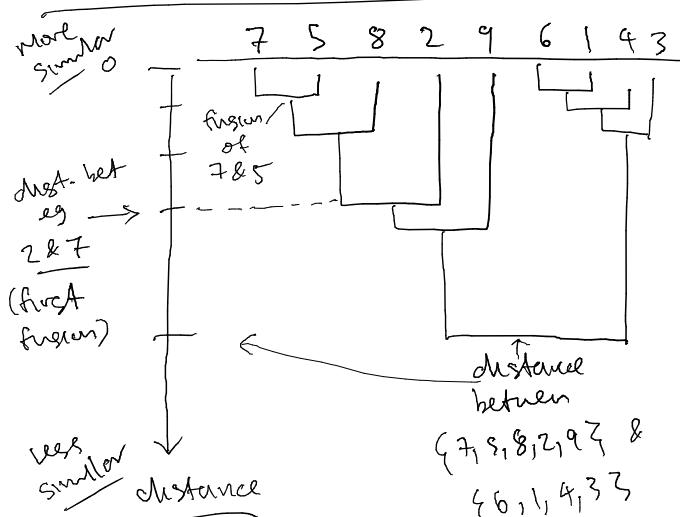
distance to a
cluster ('linkage')



Dendograms

- Summarise a hierarchical
clustering algorithm

- Indicate distance between various clusters in 'tree'-style diagram
- height of 'fusion' indicates distance
- for any two observations, can find where they first fused



2. 'Control' analysis via individual perturbations

Consider gene expression
matrix again:

		Experiment			
		1	2	3	...
gene 1	gene 1				
	gene 2				
	gene 3				
	:				

Goal: how do genes 'affect'
other genes?

(can we deduce 'network'?)



Effect of perturbations

→ Make each experiment
a perturbation of
a single gene process

↳ we increase the
transcription rates of each gene
in turn.

↳ measure the change
in steady state
concentrations of all genes

↳ get (co-) 'control coeff'

↳ summarise regulatory strengths matrix



⇒ see de la Fuente (2002)
'Linking the genes'
for details

Effect of perturbations of gene transcription rates on S.S. concentrations

upshot

- Can summarise in regulatory strengths matrix of R_d
- Can use to draw potential gene regulatory network

$T R_d =$

		Experiment:		
		Rate ↑ of gene ¹	gene ²	...
Δ gene ¹ level	Δ gene ¹ level	-1	0.5	-
	Δ gene ² level	-0.2	1	-
:	:	:	:	:

Example

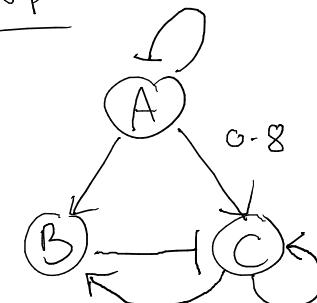
$$R_d = \begin{bmatrix} A & B & C \\ A & -0.5 & 0 & 0 \\ B & 0.1 & 0 & 1.2 \\ C & 0.8 & -0.1 & 1 \end{bmatrix}$$

↑ rate increase
↓ SS. conc change

Network

Interp:

↑ neg
← post.



negative autoreg : A

positive autoreg : C