

# 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)  
[12 lectures/3 tutorials/2 labs]

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

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

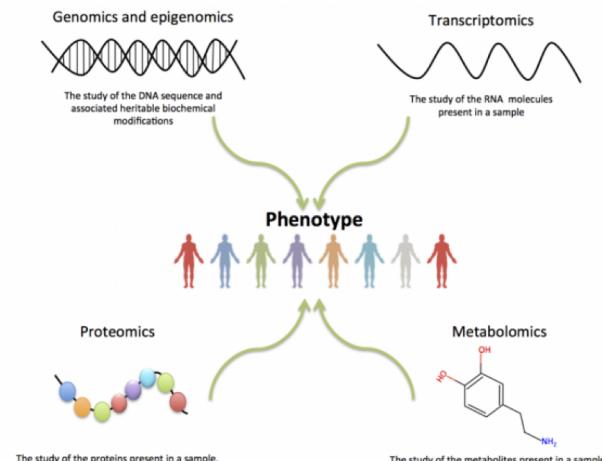
### 2. Systems biology I: overview, signalling and metabolic systems [3 lectures]

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.

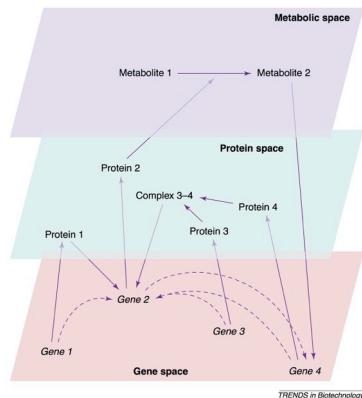
## MUCH LARGER SYSTEMS - 'OMICS'



2

4

## GENE SPACE



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

5

## 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)

7

## 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)

6

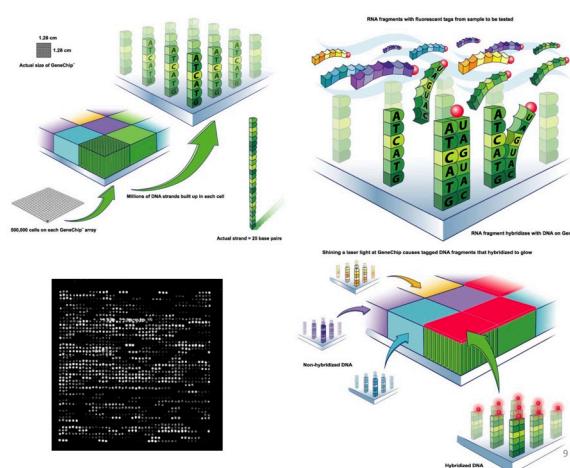
## MICROARRAYS



For video intros: see e.g. <https://youtu.be/0ATUjAxNf6U> or <https://youtu.be/VNsThMNjKhM>

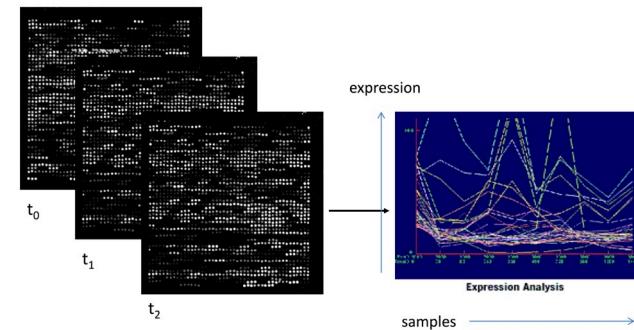
8

## MICROARRAYS



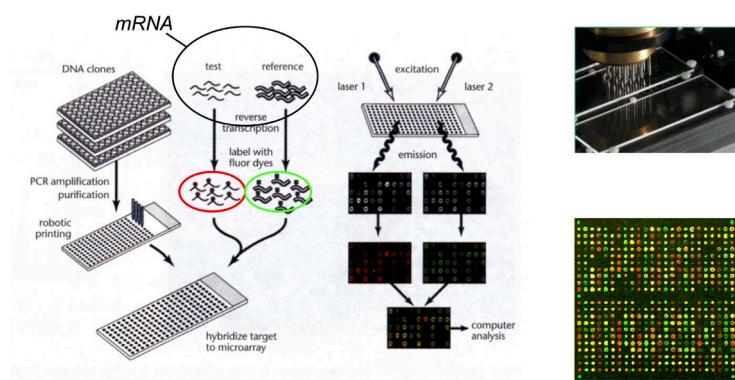
9

## MICROARRAYS: TIME SERIES



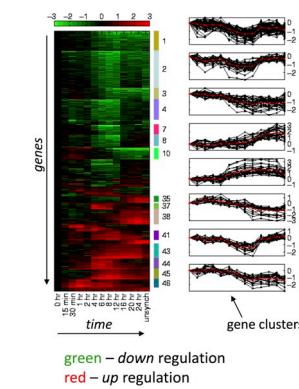
11

## MICROARRAYS: COMPARATIVE EXPRESSION



10

## MICROARRAYS: RELATIVE EXPRESSION OVER TIME



12

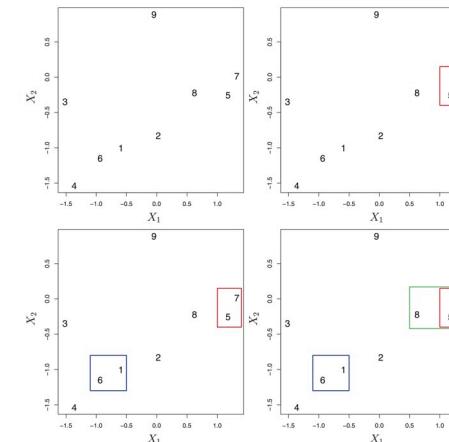
## DATA ANALYSIS: STATISTICAL/MACHINE LEARNING

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

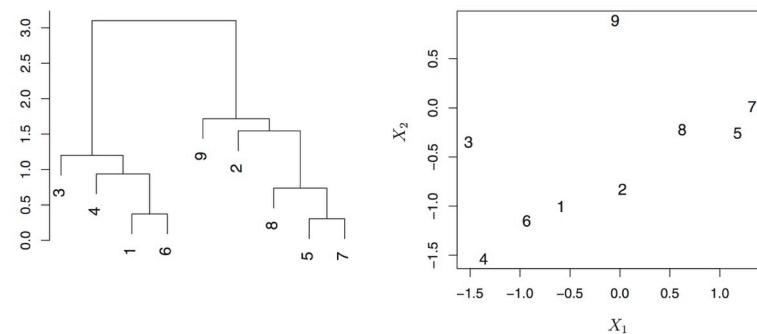
## CLUSTERING

- 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

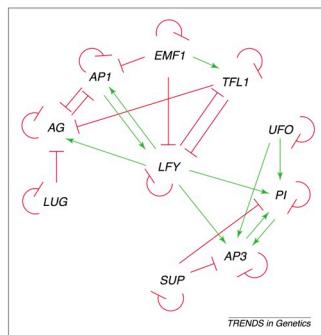
- Perturb *transcription rates* for *each gene in turn*
  - Measure changes in *steady-state expression levels* for all genes (including self)
  - Gives indication of underlying *regulatory network*
  - Summarise in *regulatory strength matrix or network diagram*.

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):



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)

(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 *Rd* 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)

## 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 mank

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

*(1 mark)*

20

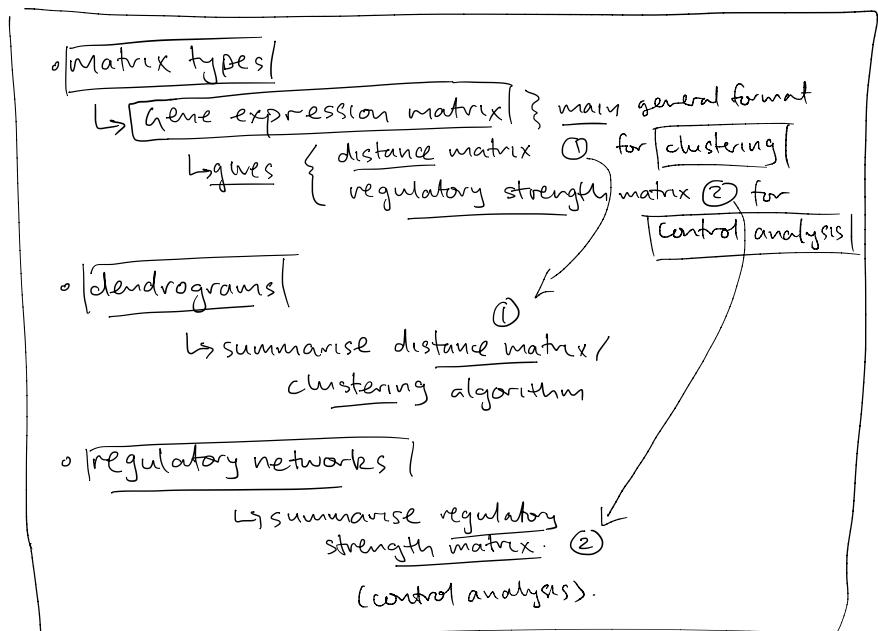
## Biomeng 261 Lecture 12

### Large genetic systems

- 'Gene space' & GRNs  
(genetic regulatory networks)

### Expression analysis methods

Falzeaways **Upshot**: know/interpret/draw



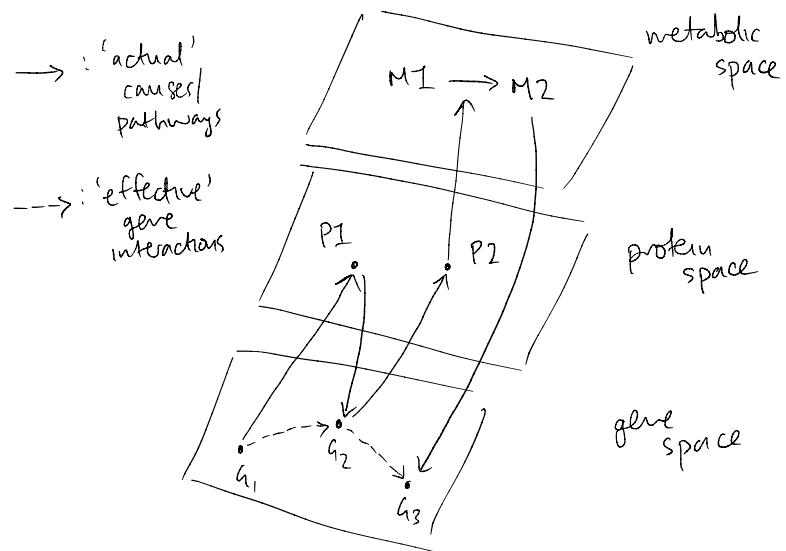
Setting: even more complex networks  
than last time

→ eg 1000s of genes (or more!)

(Background):

Gene Space

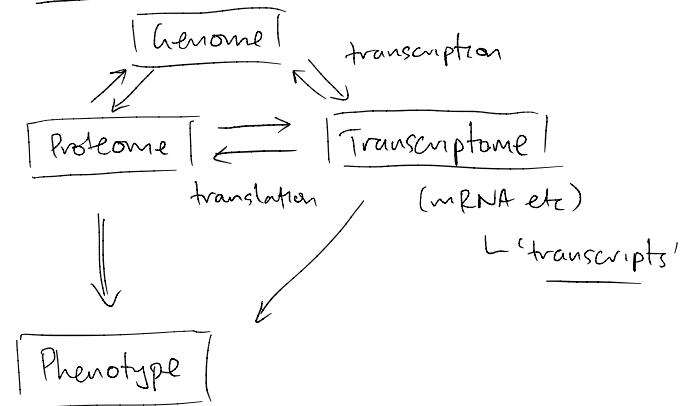
- A way of 'projecting all the action'  
down into interactions between genes



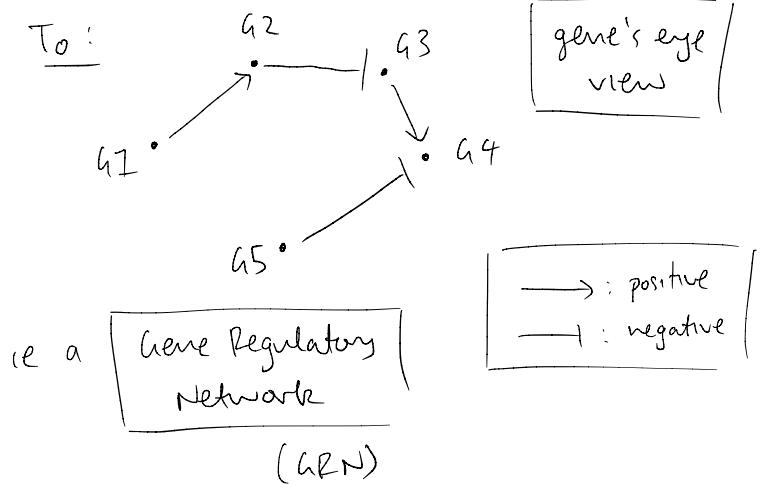
(Background)

### Gene Regulatory Networks

From:



To:



### Transcriptomics

- A subfield of 'functional genomics'

i.e. the study of how genes & intergenic (between gene) regions contribute to biological function

- Transcriptomics focuses on

gene expression levels

In particular as measured via the levels of the transcripts (mRNA) associated with genes

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

Ideas: - does expression go up or down (under treatment)?

- do groups of genes go up/down together?

## Expression Analysis

Two key approaches :

### - Microarrays

- ↳ uses (known) 'probes' (eg cDNA)
- ↳ samples 'hybridise' if complementary to probes
- ↳ amplify & quantify via qPCR
- ↳ PCR: polymerase chain (see lab lectures)

need to know roughly what looking for

### - RNA-seq

- ↳ direct sequencing of transcripts
- ↳ 'next gen', high-throughput sequencing

don't need pre-chosen 'probes'

\* we will discuss microarrays \*

- well understood
- more mature & easier to analyse
- still used & useful
- ... but RNA-seq overtaking!

(see Lowe et al. 2017  
'Transcriptomics technologies' )

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 different 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

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

From now: take measure as 'given'

& just use numbers.

→ Vinod knows more about experimental side!

## 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

eg treated vs untreated cancer or not etc.

### 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

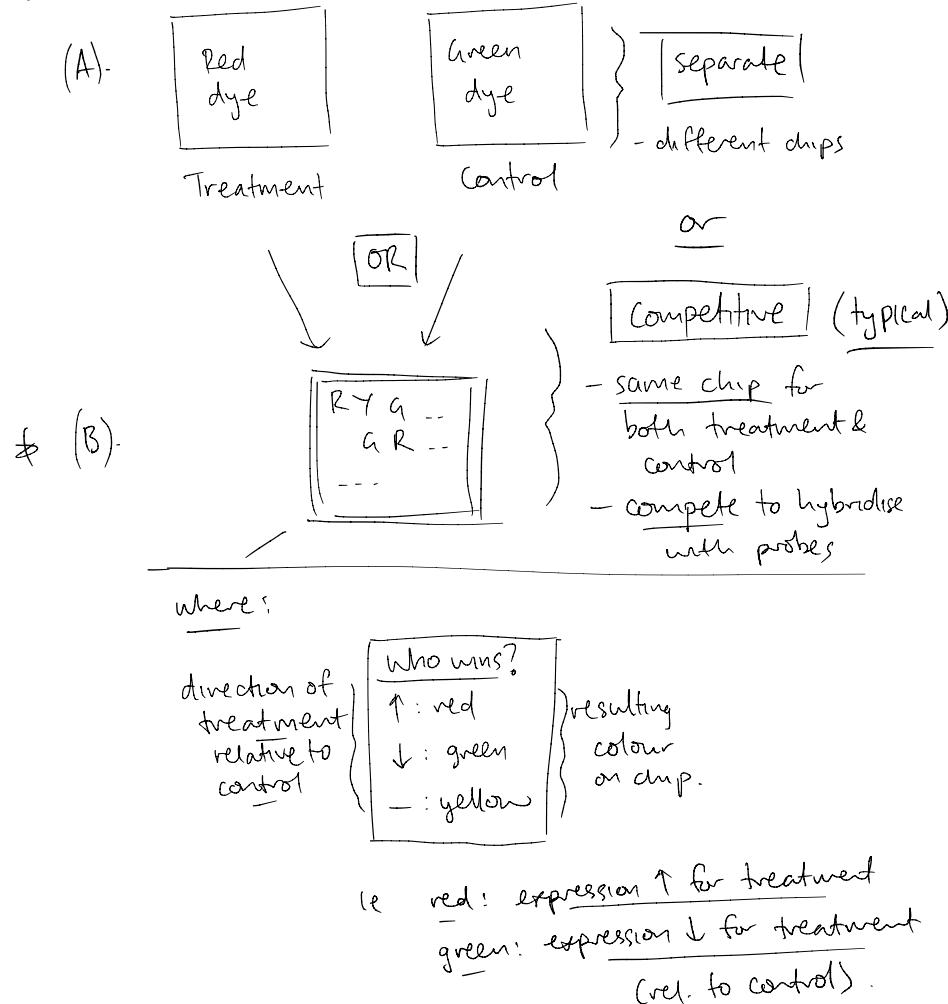
— just different independent vars, same basic ideas

— also, often want to combine,  
eg compare two time series from different treatments at same times

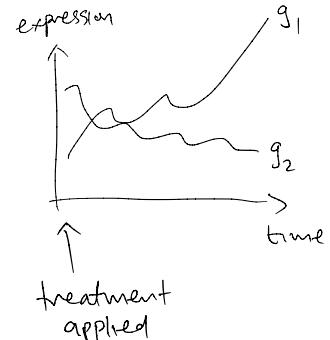
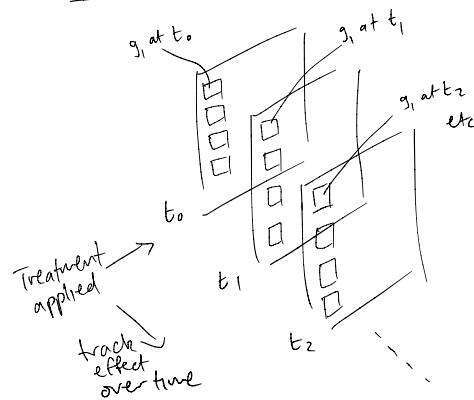
## Perturbation/comparative

- Have both a treatment & control (& at single time)

### Variations:

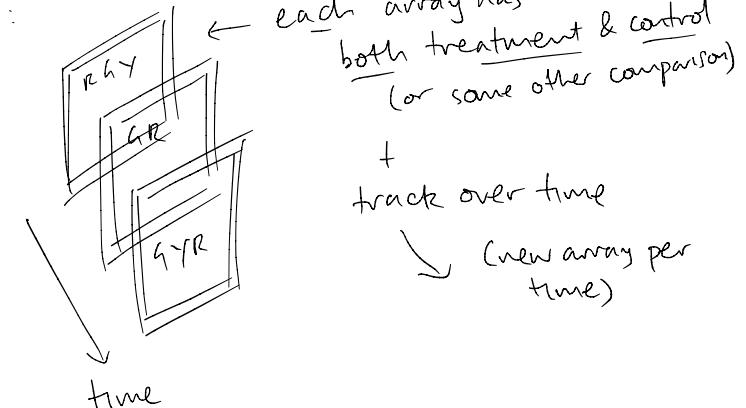


## Time Series



- treat one sample type at t<sub>0</sub> (say)
- new chip for each time
- track resulting expression patterns over time.

### Combined:



## Data format: Gene expression matrices

→ A general format that we can put both time series & perturbation/comparative into:

	Exper.1	Exper.2	Exper.3	...
Gene1				
Gene2				
Gene3				
:				
:				

Eg experiment 1: time  $t_0$ , control

experiment 2: time  $t_1$ , control

:

experiment  $n+1$ : time  $t_0$ , treatment

experiment  $n+2$ : time  $t_1$ , treatment etc.

(Warning: many 'data array' formats in stat./  
eg R, Python etc are the transpose  
of above: genes are col, exp. are row)

## Gene expression matrices

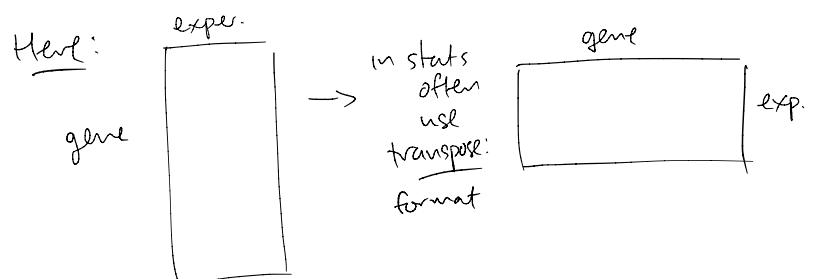
Column: a particular experiment, i.e an array (gene chip/microarray) of all genes.

	Exper.1	Exper.2	Exper.3	...
Gene1				
Gene2				
Gene3				
:				
:				

row: a gene expression profile over experiments for a given gene.

Typically: ~~genes~~  $\rightarrow$  experiments

unobserved data



## Questions / Problems

Q: which genes 'cause' difference between control & treatment / two samples for comparison?

→ many more genes than experiments

→ overfitting issues if try to explain via single genes!

eg are difference (cancer or not) → 1000 possible explanations!

→ problem for naive/traditional stat. inference!

↳ motivation for many modern stat/ML methods -

One way to tackle:

- focus on sets / clusters of genes

eg



→ smaller number of these groups

→ also relates to idea that genes work together & co-express

## Typical analyses

1. Clustering : similarity analysis (eg over time)

2. Perturbation : linear control analysis based (ss).

### 1. Clustering

- a form of unsupervised learning for pattern discovery

- need a notion of 'distance' however  
↳ user input (see eg prev. lectures & below)

### For fun (not examinable):

Defining idea of 'distance' & related:

#### distance

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

#### metric

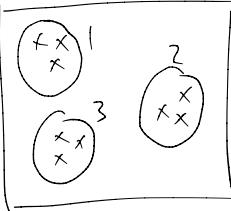
- i-iii) &
- iv)  $d(x, y) = 0$  iff  $x = y$
- v)  $d(x, y) + d(y, z) \geq d(x, z)$

#### dissimilarity

- i-ii) &
- iii)\*:  $d(x, y)$  increases monotonically as  $x$  &  $y$  more 'dissimilar' (subjective)

## (Terminology)

Unsupervised? Note on 'learning types'

<u>Supervised</u>	<u>Unsupervised</u>
$X \rightarrow Y$	$X \curvearrowright$
learns function $X \rightarrow Y$	pattern discovery (in $X$ )
<u>Train</u> : supervisor gives examples	Eg. Find <u>clusters</u>
→ Happy } Labels given → Sad }	
<u>Test</u> : predict on new (unseen) set	<ul style="list-style-type: none"> <li>→ group 'similar'</li> <li>→ <u>no labels given</u></li> <li>→ <u>do need a distance/</u> similarity measure</li> <li>→ <u>In general harder to evaluate</u> (some methods exist --&gt;)</li> </ul>
→ Happy X	
→ Happy ✓	

Example: Clustering expression profiles

→ across time &/or experiments

(ie profiles of genes: )

Expression matrix:

	Experiment							
	1	2	3	...	...	8		
gene A	-1	-2	2	---	---	2		

gene B	-1	-2	-1	---	---		
	1	1	1	1	1	1	1

gene I	-1	-2	0	---	---		
	1	1	1	1	1	1	1

Q: which genes have 'similar' profiles?

→ define distance --- eg Euclidean  
(since easy)

Distance between two profiles A, B :

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

( square root of sum of squared differences)

where

$\left\{ \begin{array}{l} n : \text{number of experiments} \\ y_i^{(A)} : \text{expression level of} \\ \text{gene A in experiment } i \end{array} \right.$

Example

A  $\begin{bmatrix} 1 & -1 & 2 & 0 \end{bmatrix}$

} two gene profiles.

B  $\begin{bmatrix} 1 & -1 & -1 & 2 \end{bmatrix}$

Squared  
differences:

$0^2$	$0^2$	$3^2$	$2^2$	Sum squares!
0	0	9	4	

13

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

Distance matrices

- We can summarise all pairwise  
differences between profiles

in a  $\boxed{\text{distance matrix}}$  (note: this  
isn't an  
expression matrix)

$$\begin{array}{cccc} A & B & C & D \\ \begin{bmatrix} d(A, A) \\ d(B, A) \\ d(C, A) \\ d(D, A) \end{bmatrix} & \ddots & \ddots & \ddots \\ \begin{bmatrix} d(A, B) \\ d(B, B) \\ d(C, B) \\ d(D, B) \end{bmatrix} & \ddots & d(C, C) & \ddots \\ \begin{bmatrix} d(A, C) \\ d(B, C) \\ d(C, C) \end{bmatrix} & \ddots & \ddots & d(D, C) \\ \begin{bmatrix} d(A, D) \\ d(B, D) \\ d(C, D) \\ d(D, D) \end{bmatrix} & \ddots & \ddots & \ddots \end{array}$$

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

$$\begin{array}{cccc} A & B & C & D \\ \begin{bmatrix} 0 \\ d(B, A) \\ d(C, A) \\ \vdots \end{bmatrix} & \ddots & \ddots & \ddots \\ \begin{bmatrix} 0 \\ 0 \\ d(C, B) \\ \vdots \end{bmatrix} & \ddots & 0 & \ddots \\ \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} & \ddots & \ddots & 0 \\ \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} & \ddots & \ddots & 0 \end{array}$$

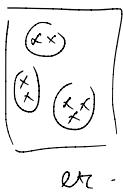
ignore

Note: - symmetric, so don't need  
explicit upper part  
ie  $d(B, A) = d(A, B)$

- diagonals are zero.  
 $d(A, A) = 0$  etc.

## Distance-based clustering

Goal: group together if 'close' eg:



etc.

Two popular algorithms

- K-means

- Hierarchical

Here: will consider Hierarchical

### Pseudocode (Hierarchical)

Begin with n observations

Find all pairwise distances

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

Find smallest distance

'Fuse' or 'group' together (so  $n$  obs  $\rightarrow n-1$  obs)

Recompute distances to 'fused' cluster \*

(see below)

\* Distance to cluster?

- multiple types
  - $\min$
  - $\max$
  - $\text{average}$
  - $\text{etc} \dots$

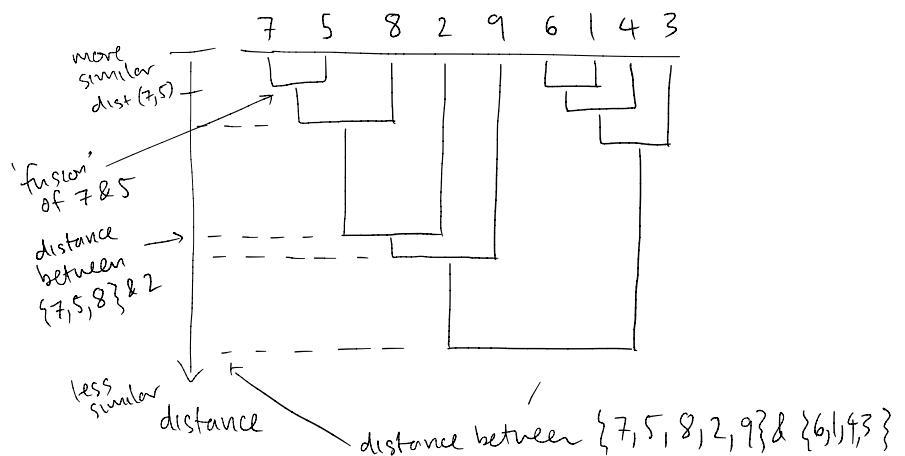
?

## Dendograms

- summarise results of hierarchical clustering

• Indicate distance between various clusters in 'tree-style' diagram

Example:



- height (y axis) of 'fusion' indicates distance

- for any two clusters, can find where they first fused

## Other analysis type.

### 2. 'Control' analysis via individual perturbations

Consider gene expression matrix again:

		Experiment			
		1	2	3	---
gene 1	1				
	2				
	3				

Goal: how do genes 'affect  
other genes?

Can we deduce 'control network'?

### Effect of perturbations

→ Make each experiment a perturbation of a single gene

- increase transcription/rates of each gene in turn
- measure the change in steady state (concentrations / expressions of mRNA) for all genes (incl. self).

→ get (co-) ['control coefficients'] (eg +1.5 or -0.5)  
→ summarise in a [regulatory strength matrix]

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

Effect of perturbations of gene transcription rates on ss. concentrations/expressions levels

Result:

- can summarise in regulatory strength matrix  $R_d$

- can use to draw potential gene regulatory networks

$R_d$		Experiment: rate ↑ of...		
		gene 1	gene 2	...
Δ gene 1 level	-1	0.5	-	-
Δ gene 2 level	-0.2	1	-	-
⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮

↑  
normalised  
change in  
expression level

Example

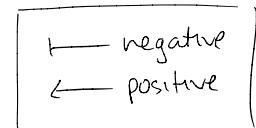
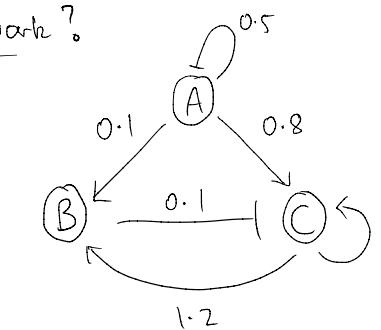
$$R_d = \begin{bmatrix} A & B & C \\ A & -0.5 & 0 & 0 \\ B & 0.1 & 0 & 1.2 \\ C & 0.8 & -0.1 & 1 \end{bmatrix}$$

rate increase  
↓  
s.s. conc. change  
↑

Interp: effect of -- on [ ]

Q: Network?

A:



negative auto (self) regulation: A  
positive auto regulation: C