

Blomeng 261 Microarrays

- Basic biology (very brief)
- Data / Experiment types
- Basic analysis methods

Upshot: Know & interpret:

- o matrix types
 gene expression matrix } main
 ↳ gives { distance matrix
 regulatory strength matrix
- o dendrograms
 ↳ summarise distance matrix / clustering alg.
- o regulatory networks
 ↳ summarise regulatory strength matrices

Background Biology (see slides & video link)

Idea

Want to measure mRNA levels (gene expression)

o microarrays can measure 1000s of mRNA levels at a time

↳ 1000s of genes at a time

↳ a grid of 'spots' containing cDNA (complementary DNA)

reverse transcription
(get DNA from RNA)

↳ 'Hybridisation'

↳ mRNA preferentially binds to corresp. cDNA

also → use fluorescent tags to distinguish samples
→ amplify via PCR/qPCR

Preprocessing

Preprocessing is crucial

BUT: we won't go into

-
- Qs: - artifacts
- absolute vs relative expression
- log transformations
etc.

Typically work with

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

- see slides for an example conversion
- from now, take as given

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'

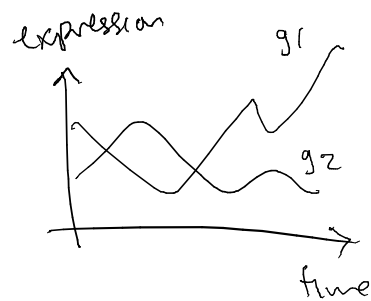
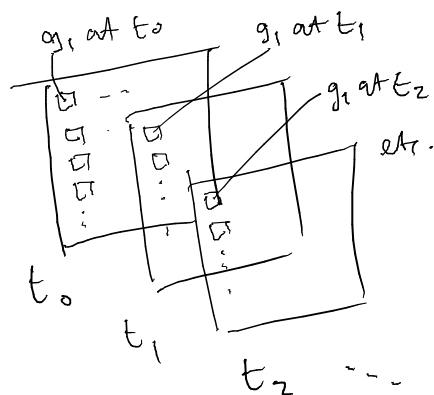
2. - Time-series
 ↳ same cell/tissue etc studied over time } one sample type, multiple times

idea: time has natural order (continuous or categ. ordinal)
• treatment/class not rec. ordered
 - has cancer, doesn't, etc.

BUT essentially

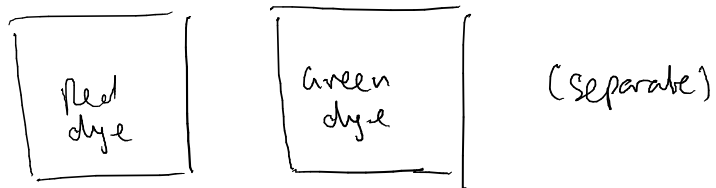
- ⇒ just different independent vars, same basic ideas
- ⇒ also, often want to combine:
 eg compare two time series from diff. treatments at same times

Time Series

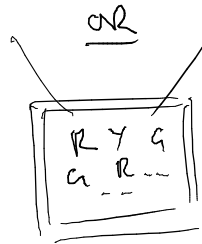
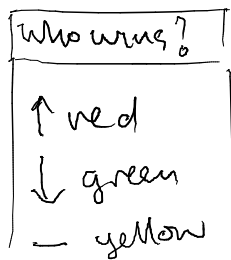


new chip for each time

Perturb. / comparative



treatment I control (more commonly)



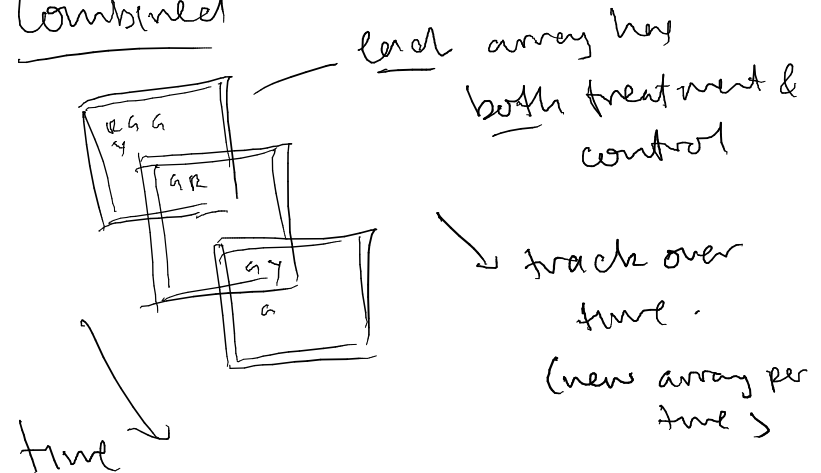
Competitive

↳ same chip for both treatment & control

↳ compete to hybridise

↳ direction of treatment rel. to control.

Combined



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

Two key goals

- try to cluster 'similar' &/or
- 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.

	Exp1	Exp2	Exp3
gene1			
gene2			
gene3			
:			
:			

Exp:
experiment
(not expression)

→ eg exp1: time t_0 , control } experiment 1.
exp2: time t_1 , control
:
expn: time t_0 , treatment
expn+1: time t_1 , treatment

etc.

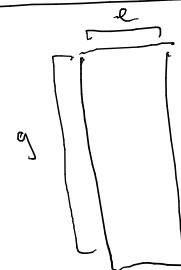
Gene expression matrices

Column: an experiment, ie an array
(gene chip/microarray)

	Exp1	Exp2	Exp3
gene1			
gene2			
gene3			
:			
:			

Row: a gene expression profile for a given gene

Typically: $\# \text{ genes} \gg \# \text{ experiments}$



• many more genes than experiments

• 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 analysis

1. Clustering

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

Note: more precise distinctions

Distance

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

Metric

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

Dissimilarity

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

Unsupervised? Note on 'learning types'

Supervised

$X \rightarrow Y$


learn function $X \rightarrow Y$

Unsupervised

$X \rightarrow X$


Pattern discovery in X


↳ Train: Supervisor examples

 \rightarrow happy
(labels given)

 \rightarrow sad

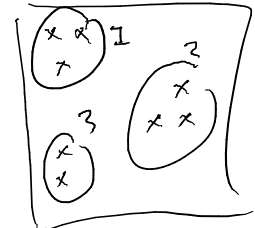
Test: predict on new set

 \rightarrow happy X

 \rightarrow happy ✓

'Reasonably straightforward' to evaluate/validate.

Eg: Find clusters



group 'similar'

\rightarrow no labels given

\rightarrow goal: discover labels/clusters etc

\rightarrow do need distance / similarity $d(x_1, x_2)$

\rightarrow Can do some eval/valid. eg via consistency/stability on training/test sets
↳ but in general, harder to 'validate'

Example : clustering expression profiles

→ across time/exp - (profiles)

Expression
matrix

	Experiment					
	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
similarity)

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
squared
diff.

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

Example :

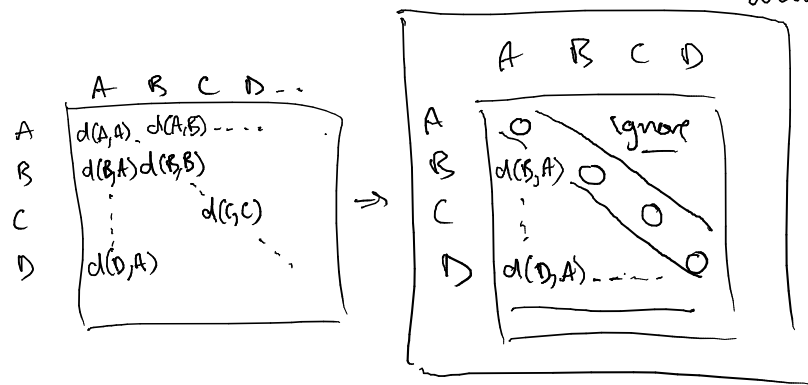
A	1	-1	2	1	0	1
	↓	↓	↓	↓		
B	1	1	-1	-1	2	1
squared diffs:						sum squares
12 0 0 9 4						13

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

Distance matrices

We can summarise all pairwise differences in a

distance matrix (note: is not an expression matrix)



Note: • symmetric so don't need upper part if

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

• diagonals are zero

Example: slides

Distance-based clustering

Goal: group together if close.

Two popular algorithms

• K-means

• hierarchical

see slides for pseudocode

Here: consider hierarchical



Hierarchical clustering

Begin with n observations

Find pairwise distances

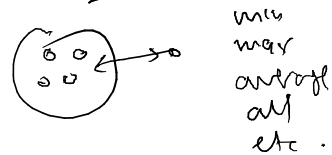
For $i = n, n-1, \dots, 2$:

Find smallest distance

'Fuse' together

Recompute distances to
'fused' cluster

↳ need notion of
distance to
cluster ('linkage')

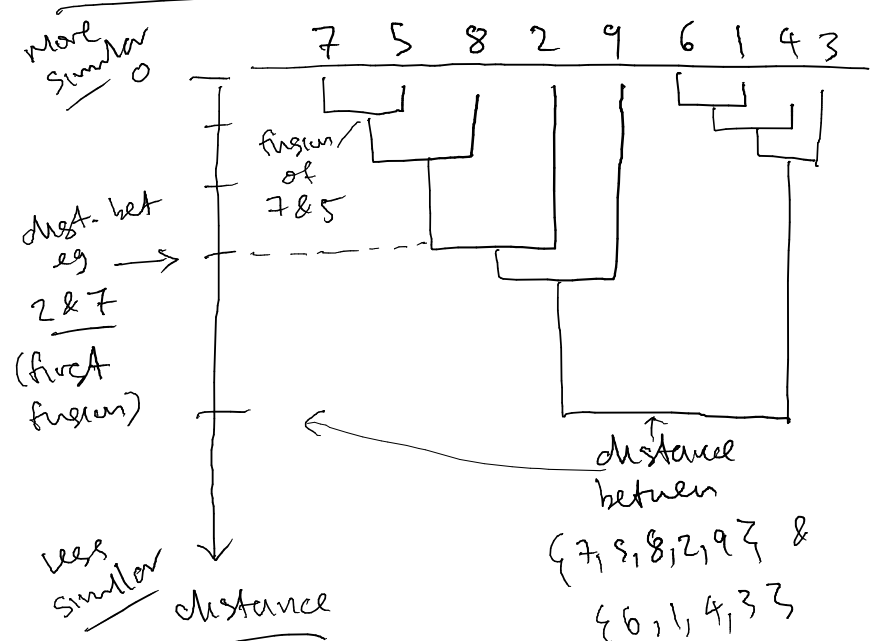


pics: see slides

Dendrograms

- 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 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 (w-) 'control cell'

↳ 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 R_d

Can use to draw potential gene regulatory network

$$|R_d| =$$

Normalised change in exp. level.

Experiment:
Rate ↑ of
gene 1 gene 2 ...

Δ Gene 1 level	-1	0.5	...
Δ Gene 2 level	-0.2	1	...
...

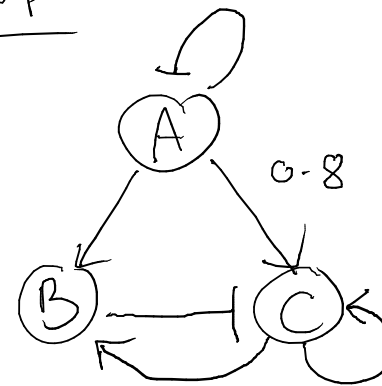
Example

rate increase

SS. conc change

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

Network Interp:



T-bar neg
← pos.

negative autoreg : A
positive autoreg : C