

Gene Expression & Annotation

Clare Foyle, Nick Fieller &
AstraZeneca R & D

University of Sheffield & Alderley Park

October 2005



Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05



Outline

- ◆ Basic Biology:–
 - DNA, Genes, Proteins,
- ◆ Measuring Gene Expression
- ◆ Gene Expression Data
- ◆ Annotation
- ◆ Mesh Terms
- ◆ MeSH DAG
- ◆ Quantifying Annotation
- ◆



Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

2

Basic Biology

- ◆ Genes \Rightarrow Proteins \Rightarrow Bioactivity
 - bioactivity includes disease (e.g. cancer)
 Aim is to design drugs to interfere with bioactivity by interrupting pathways of protein activity
- ◆ DNA: [linear] sequence of ~3 billion base pairs of nucleotides (C G A T)
 - Genes: subsequence of few hundred base pairs
 - recognised by key start & end sequences
 - genes don't vary [much] between people
 - most base pairs are identical across the population
 - ~30,000 genes in total
- ◆ ~ 10% of DNA consists of genes ('apparently')



Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05



3

Genes \Rightarrow Proteins

- ◆ Mechanism involves genes 'expressing' RNA which in turn 'constructs' proteins
- ◆ Genes do not 'express' continuously but only in relation to biological events
 - i.e. if 'switched on'
- \Rightarrow **key interest:** –
 - determine which genes are (or are not) expressed in various circumstances
 - e.g. cancer tissue may have different gene expression profile from other types of tissue
 - differences may reveal processes involved in that type of cancer



Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

4

Measuring Gene Expression

- 3 Basic technologies:
 - ◆ Affymetrix® GeneChips® oligonucleotide assays
 - expression levels of each gene in sample
 - ◆ Red/green spot assays
 - differential expression levels in 2 samples
 - ◆ Radionucleotide assays
 - differential expression levels in 2 samples
- All involve 'labelling' RNA from tissue
 - e.g. fluorescent markers



Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05



5

Affymetrix GeneChips

- ◆ Sequences of [all] genes is known
 - Can be manufactured to order
 - Just need to stick individual molecules of C G T & A together in right order
 - possible to do with very great precision
 - Subsequence of ~25 bases unique within the genome to that gene is a **Probe Set** for that gene
- ◆ Affymetrix GeneChips have many thousands of probe sets placed at known coordinates on a chip



Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

6



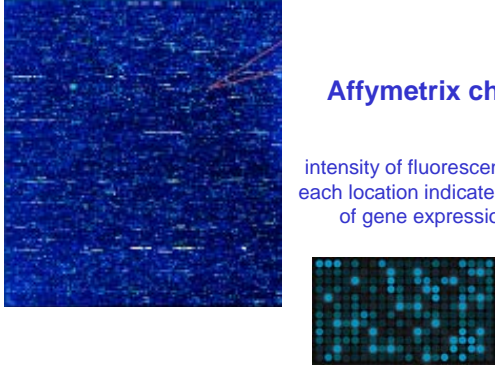
7

- ◆ Chip has single strands of DNA probe sets
- ◆ Tissue containing mRNA processed by reverse transcription to cDNA incorporating fluorescent marker
 - (only contains DNA from 'expressing genes')
- ◆ Denatured to produce a single strand
 - (from the double strand helix)
- ◆ Solution flushed over chip and strands in sample bind to single strands of probe set
 - fluorescence measured

⇒ measures of expression of genes in sample

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

8



Affymetrix chip

intensity of fluorescence at each location indicates level of gene expression

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

9

- **Points to Consider:–**
 - ◆ Each chip has several copies of each probe set
 - result is mean and variance of measures
 - Affymetrix software provides mean level & a p-value
 - (Better ways of initial processing widely investigated)
 - ◆ May be several probe sets for single gene
 - e.g. subsequences from beginning/middle/end
 - may happen that only some probe sets from a gene will indicate expression (& not others)
 - ◆ genes ⇒ proteins
 - but single gene may produce several different proteins under different circumstances
 - different genes may produce same protein

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

10

- **Summary so far**
 - ◆ Gene expression data:–
 - typically 'gene expression data sets' have relatively few subjects, little or no replication and quantitative measures of many thousands of gene [or probe sets]
 - data are typically of 'poor quality'
 - many outliers, high variability, 'contamination'
 - data are 'expensive' to obtain per person
 - each chip ~\$100
 - one chip can be used for only one sample

⇒ data sets typically are "n < p" with no replication but there may be covariates

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

11

- **Need some dimensionality reduction**
- Most common is to proceed in 2 stages:
 - ◆ initial selection of 'interesting genes'
 - ◆ standard dimensionality reduction on these
 - ◆ Stage 1:
 - alternative 1: select the N genes with highest expression levels
 - alternative 2: select the N genes with greatest inter-group variability
 - groups defined by covariates (e.g. cancer 'level')
 - » c.f. "differential expression"
 - N=500 and N=100 were tried
 - ◆ Stage 2: PCA, LDA, PLS,.....

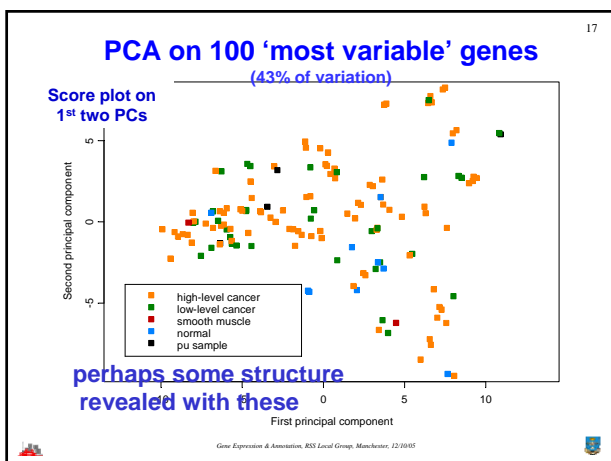
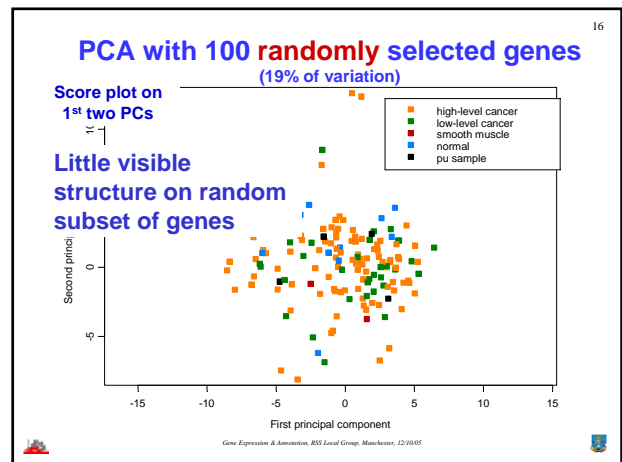
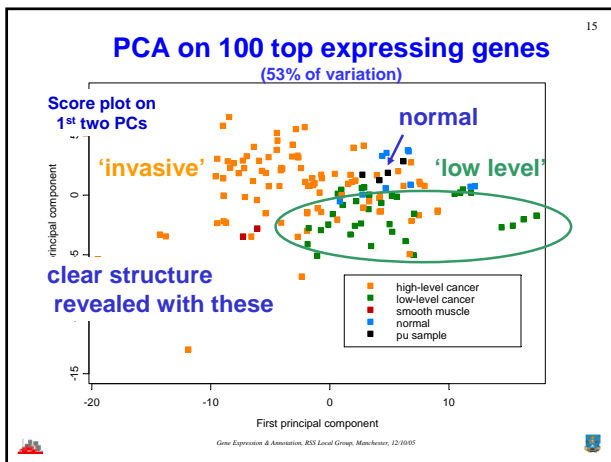
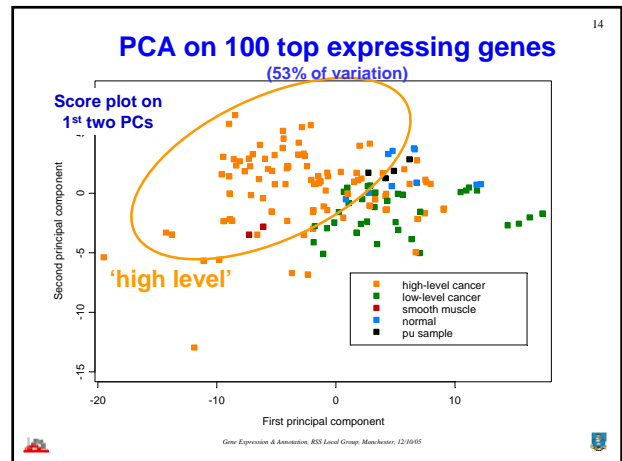
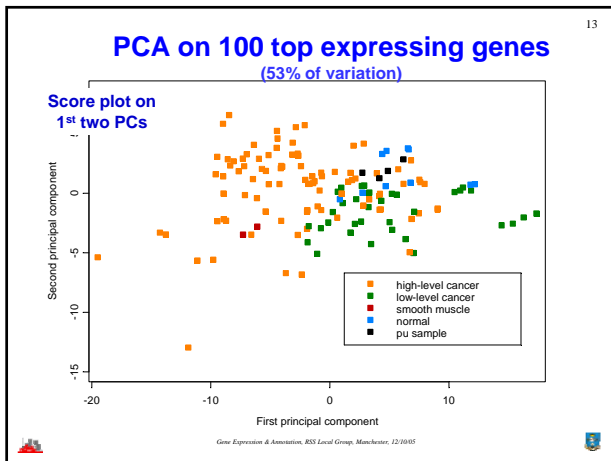
Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

12

- ◆ **example:–**
 - ~60 cancer subjects (+ some controls),
 - ~13,000 probe sets (coded)
 - one sample from each, measured twice,
 - (note these are not full replicates)
 - + background information sex/severity/age/...
- ◆ **Aside: 3rd party confidentiality meant we had only code numbers for probe sets**
 - no information on which probe sets were from same gene
 - inhibits some avenues of analysis

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05





- 18
- ◆ Analysis with top 100 'most variable' gives some structure on 1st two PCs but not as clear as top expressing
 - 3rd & 4th PCs needed (+16% variation)
 - (not shown)
 - ◆ Analysis with 500 genes gives less clear results on high order PCs than with only first 100
 - (not shown)
 - ◆ Similar pictures and conclusions from other categorizations of samples
 - e.g. demographic features
 - (not shown)
- Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05



19

- Results suggest that initial selection methods capture at least some of available information in genes
 - (*actually probe sets in this case*)
 - suggests may be able to reduce variability by 'averaging' over subjects with common covariates — perhaps literally or by modelling
- What else is known?

ANNOTATION

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

20

ANNOTATION

- 'Annotation' — information on genes
 - typically text, possibly published, free form
 - position on chromosome
 - role in intracellular pathway
 - citation in published paper or tech report
 -
 -
- may be different types & quantities of information on different genes
 - very extensive

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

21

- Annotation clearly important in **follow-up** to analysis when 'statistically important' genes have been identified
 - in which pathway does it have a role?
 - which proteins does it 'produce'?
 -
- Can annotation be used in the analysis?
 - requires some form of **quantification** of the annotation information
 - (I presume...)

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

22

- Initial ideas on annotation**

structural information on genes (annotation)

Sample × Gene data matrix of expression levels

structural information on samples (stage, demographics etc)

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

23

- Initial ideas on annotation**

inherent symmetry:
Gene × Sample
or
Sample × Gene matrix

structural information on samples (stage, demographics etc)

structural information on genes (annotation)

Gene × Sample data matrix of expression levels

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

24

- Initial ideas on annotation**

invites iterative analysis of genes × samples & samples × genes to update selection of genes related to structure (work in progress)

can also compare structure with annotation on randomly selected genes

Gene × Sample data matrix of expression levels

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05



25

- A simple approach
 - » (simplistic?)
 - ◆ underlying idea is the genes associated with same [or related] 'biofunctions' are themselves 'related'
 - ◆ start by looking for co-citations of genes with specific biofunctions
 - i.e. how many papers refer to both gene X and biofunction A
 - need extensive text mining literature search
 - unrealistic to try for all 13,000 genes so restrict [initially] to 'preselected interesting' ones

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

26

- Next step:
 - ◆ determine which groups of genes are linked to same biofunction term
 - regard these as 'related' genes
- Problem:
 - ◆ very sparse information
 - SO, link together biofunctions
 - genes linked to 'related biofunctions' are [distantly] related
- Finally
 - ◆ adjust gene × sample analysis by 'averaging' over related genes
 - c.f. 'averaging' over related samples

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

27

- biofunctions???
- ◆ need definitive list
- ◆ need declared hierarchy to provide relationships
- Two (+) possibilities
 - ◆ GO
 - Gene Ontology Consortium
 - » <http://www.geneontology.org/>
 - ◆ MeSH
 - Medical Subject Headings
 - US National Library of Medicine
 - » <http://www.nlm.nih.gov/mesh/meshhome.html>

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

28

- biofunctions???
- ◆ need definitive list
- ◆ need declared hierarchy to provide relationships
- Two (+) possibilities
 - ◆ GO
 - Gene Ontology Consortium
 - » <http://www.geneontology.org/>
 - ◆ MeSH ←
 - Medical Subject Headings
 - US National Library of Medicine
 - » <http://www.nlm.nih.gov/mesh/meshhome.html>

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

29

- MeSH
 - Medical Subject Headings
 - US National Library of Medicine
 - » <http://www.nlm.nih.gov/mesh/meshhome.html>
 - Terms under 15 broad headings organised in hierarchical tree structure (DAG)
 - Selected 2 headings
 - Diseases
 - Biological Sciences
 - scanned literature for co-citation of each of our 'interesting' genes with each term under these headings
 - i.e. two [overlapping] sets of 500 of probe sets
 - » (< 500 distinct genes)
 - ~6000 terms

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

MeSH Tree Structures - 2005

[Return to Entry Page](#)

1. [\[A\] Anatomy \[A\]](#)
2. [\[B\] Organisms \[B\]](#)
3. [\[C\] Diseases \[C\]](#)
4. [\[D\] Chemicals and Drugs \[D\]](#)
5. [\[E\] Analytical, Diagnostic and Therapeutic Techniques and Equipment \[E\]](#)
6. [\[F\] Psychiatry and Psychology \[F\]](#)
7. [\[G\] Biological Sciences \[G\]](#)
8. [\[H\] Physical Sciences \[H\]](#)
9. [\[I\] Anthropology, Education, Sociology and Social Phenomena \[I\]](#)
10. [\[J\] Technology and Food and Beverages \[J\]](#)
11. [\[K\] Humanities \[K\]](#)
12. [\[L\] Information Science \[L\]](#)
13. [\[M\] Persons \[M\]](#)
14. [\[N\] Health Care \[N\]](#)
15. [\[Z\] Geographic Locations \[Z\]](#)

[Return to Entry Page](#)

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05



MeSH Tree Structures - 2005

[Return to Entry Page](#)

1. Anatomy [A]
2. Organisms [B]
3. Diseases [C]
 - o [Bacterial Infections and Mycoses \[C01\]](#) +
 - o [Virus Diseases \[C02\]](#) +
 - o [Parasitic Diseases \[C03\]](#) +
 - o [Neoplasms \[C04\]](#) +
 - o [Musculoskeletal Diseases \[C05\]](#) +
 - o [Digestive System Diseases \[C06\]](#) +
 - o [Stomatognathic Diseases \[C07\]](#) +
 - o [Respiratory Tract Diseases \[C08\]](#) +
 - o [Otorhinolaryngologic Diseases \[C09\]](#) +
 - o [Nervous System Diseases \[C10\]](#) +
 - o [Eye Diseases \[C11\]](#) +
 - o [Urologic and Male Genital Diseases \[C12\]](#) +
 - o [Female Genital Diseases and Pregnancy Complications \[C13\]](#) +
 - o [Cardiovascular Disease](#)
 - o [Hemic and Lymphatic](#)

&c., &c., &c..... until

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

32

2005 MeSH
MeSH Descriptor Data

[Return to Entry Page](#)

MeSH Heading	Gift Giving
Tree Number	F01.145.813.208
Scope Note	The bestowing of tangible or intangible benefits, voluntarily and usually without expectation of anything in return. However, gift giving may be motivated by feelings of ALTRUISM or gratitude, by a sense of obligation, or by the hope of receiving something in return.
Entry Term	Financial Contributions
Entry Term	Gifts, Financial
Entry Term	Giftgiving
Entry Term	Gifts
See Also	Altruism
See Also	Charities
See Also	Financial Support
See Also	Fund Raising
Allowable Qualifiers	IS
History Note	2003, for GIFTS, FINANCIAL use FUND RAISING 1978-2002
Unique ID	D057921

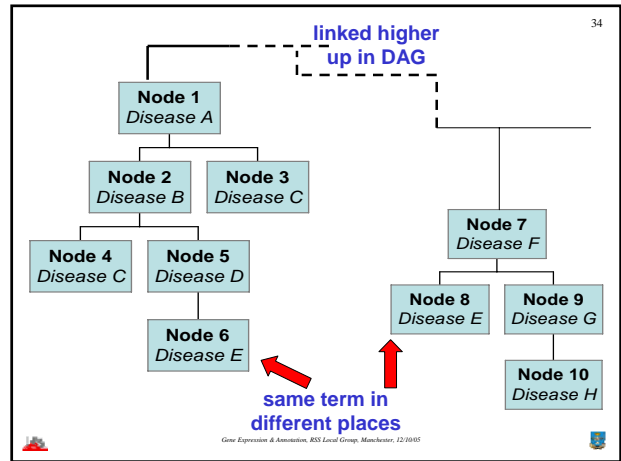
Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

33

- Some MeSH terms can appear in several places in the DAG

numbers give sequence of nodes in DAG

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05

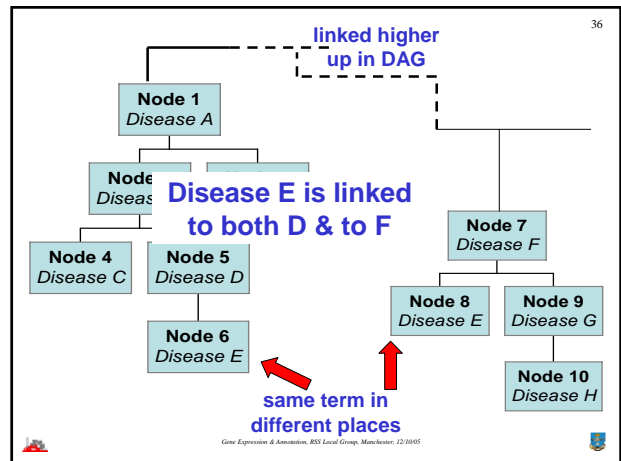


35

- Adjacent terms in DAG are 'linked'

- ◆ Disease D is linked to both B and E
- ◆ 'Up-propagation' involves combining nodes with categories one level higher in DAG
- ◆ complication with multi-occurrence of terms

Gene Expression & Annotation, RSS Local Group, Manchester, 12/10/05



37

Algebraic Manipulations

- ◆ Annotation information, A , $n \times p$ matrix
 - with n probesets and p annotation terms.
 - gives number of associations found linking an annotation term to a probe set
 - **convert to binary form:**
- 1 link between annotation & probe set
0 is if no link has been found
- **(as yet!)**

Gene Expression & Annotation, RSI Local Group, Manchester, 12/10/05

38

AA^T gives the co-occurrence (or 'similarity') matrix of the probe sets

- ◆ Typically large & sparse
 - provides few links between genes of interest
 - to provide further (but weaker) links up-propagation of the MeSH DAG is used to establish further relationships
- ◆ Initially, genes linked via mother-daughter pairs of descriptors can be determined
 - Repetition of process allows successively weaker links to be determined via successively less closely related descriptors

Gene Expression & Annotation, RSI Local Group, Manchester, 12/10/05

39

Up-propagation:

- ◆ Up-propagation matrix, P
 - $p \times p$ matrix of the DAG
 - indicates which nodes are directly linked
- ◆ AA^T gives probe sets directly linked via co-occurrence with a single MeSH term
- ◆ APA^T gives the probe sets related by links to one-generation related MeSH descriptors
- ◆ AP^rA^T gives those linked via r-generation related descriptor

Gene Expression & Annotation, RSI Local Group, Manchester, 12/10/05

40

Issues:-

- ◆ Probe Sets / Genes
 - data in probe sets, annotation re genes
 - (+ have to go through 3rd party for annotation on probe sets)
- ◆ multiple occurrence of MeSH terms
- ◆ both require 'fiddling' with rows/columns in matrix manipulation
- ◆ large sparse matrices (10,000 \times 500)
 - computational care required
- ◆ co-occurrence may be a false measure
 - "This gene has nothing whatsoever to do with that disease"
- ◆ only have annotation for genes which are 'interesting' & preselected – not on all genes
- ◆

Gene Expression & Annotation, RSI Local Group, Manchester, 12/10/05

41

■ <http://www.shef.ac.uk/paspgr/foyle>

■ <http://www.shef.ac.uk/nickfieller>



Gene Expression & Annotation, RSI Local Group, Manchester, 12/10/05

