Analysis of microarray data and interpretation of gene signature

Agenda

How to pre-process and perform statistical analysis of microarray data
 How to interpret gene signature, once it is obtained

- Importance of learning about microarray in RNA-Seq era With respect to dealing with expression values, many aspects of microarray analysis can be applied to RNA-Seq analysis.
- □ Large volume of publicly available microarray data



Seungwoo Hwang, <u>swhwang@kribb.re.kr</u>

Korean Bioinformation Center (KOBIC)

Korea Research Institute of Bioscience and Biotechnology (KRIBB)

2013/06/17 The 11th CJK Bioinformatics Training Course

Two-channel array vs. One-channel array

Channel: The number of fluorescent dyes used in the experiment



cDNA chip: an outdated technology. Just ignore it. Affymetrix GeneChip, etc.

Preprocessing of Affymetrix array: Robust Multi-Chip Average (RMA) Method

Scanned chip image



An easy reading: <u>http://www.plexdb.org/modules/documentation/RMAexplained.pdf</u>

Quantile normalization

- What it does: Makes the replicate arrays to have equal distribution.
- Rationale: It can be assumed that **overall distribution** of expression values does not vary much across all samples under most types of studies.
- How it works: Values for each column are ranked, then the average per rank is taken and is reattributed to each column according to the original rank.

	Chip1	Chip2			Chip1	Chip2
Probe A	10	14 ·		Rank1	4 (probe C)	8 (probe C)
Probe B	12	12	Sort the data from each	Rank2	8 (probe D)	10 (probe D)
Probe C	4	8	chip	Rank3	10 (probe A)	12 (probe B)
Probe D	8	10	• P	Rank4	12 (probe B)	14 (probe A)
				F ↓ F	For each rank he data with t	, substitute heir average
	Chip1	Chip2			Chip1	Chip2
Probe A	11	13	◀	Rank1	6 (probe C)	6 (probe C)
Probe B	13	11	Re-sort to the original	Rank2	9 (probe D)	9 (probe D)
Probe C	6	6	probe order	Rank3	11 (probe A)	11 (probe B)
Probe D	9	9		Rank4	13 (probe B)	13 (probe A)

Through quantile normalization, all chips in a dataset get the same distribution. That is,

- The highest intensity values in each chip become identical.
- The second highest intensity values in each chip become identical.
- ...
- The lowest intensity values in each chip become identical.

Normalize all the samples in a dataset together



□ All group normalization: presumes that the sample groups are similar enough.

- □ Separate group normalization: presumes that the sample groups are very different.
- In most situations (e.g., tumor vs. normal // treated vs. untreated), all group normalization is to be used.
- □ In some rare situations (e.g., brain vs. liver), do something else.

Gene-level collapsing of duplicate probe sets' intensities

□ Some probe sets correspond to identical gene. For example, in Affymetrix U133A, there are ~21,000 probe sets that have Entrez Gene annotation, but they are mapped to only ~13,000 Entrez Gene IDs. Thus ~8,000 probe sets are duplicates.

□ Many of the signature interpretation methods require the genes in a signature to be unique.



Probeset-level data table

Collapsing methods: Miller (2011) BMC Bioinformatics PMID:21816037

Simple averaging (most popular): Mean or Median

□ <u>Choosing a representative</u>: Probe set with highest overall intensity, or largest variance

Also it is the most important to use the most recent probe set annotation

Unpaired *t*-test

Test whether, on average, gene expression levels are different between control group and case group

log expression level from *i*th control individual



Paired t-test

Test whether gene expression levels are different, for example, before and after treatment for each and every individual.





Multiple testing correction



3. How we do

False-discovery rate (FDR) correction: FDR of 0.05 means 5% of false positives are expected among those identified as positives. Benjamini-Hochberg's procedure is the norm in bioinformatics.
 Family-wise error rate (FWER) correction: FWER of 0.05 means 0.05 genes are expected to be false positive. Bonferroni's procedure is well known but too strict for most bioinformatics tasks.

An apparent paradox of multiple testing correction

Microarray experiment

- Suppose raw P-value of 0.01 was obtained for TP53 gene from microarray.
- The TP53 gene was not called significant since multiple testing corrected P-value was not small enough, even though raw P-value was small.

RT-PCR on a single gene

- Suppose that, upon measuring only TP53 with RT-PCR, raw P-value of 0.01 was obtained.
- This time, TP53 was called significant since raw P-value was small.

Apparent paradox

Why TP53 was penalized only in the large-scale experiment even though they showed the same result in both small- and large-scale experiments?

Large-scale experiment is a screening experiment

Here the scientist did not have any hypothesis.

So the scientist should recognize the possibility of the presence of many false findings in his result, which can occur by chance alone, and do the best to prevent these false leads from getting into scientific community, through multiple testing correction.

Small-scale experiment is a <u>confirmatory</u> experiment

 Here the scientist had a clear hypothesis. So he already narrowed down many possibilities before the experiment, and performed the experiment only on that possibility to confirm his hypothesis.

• So the scientist is able to be confident on his finding (*"TP53 is differentially expressed, just as I expected*!")

R for microarray data analysis

EMA - A R package for Easy Microarray data

analysis A combined R package of many individual R packages. Easy to use.

Nicolas Servant^{1,2,3*†}, Eleonore Gravier^{1,2,3,4†}, Pierre Gestraud^{1,2,3}, Cecile Laurent^{1,2,3,6,7,8}, Caroline Paccard^{1,2,3}, Anne Biton^{1,2,3,5}, Isabel Brito^{1,2,3}, Jonas Mandel^{1,2,3}, Bernard Asselain^{1,2,3}, Emmanuel Barillot^{1,2,3}, Philippe Hupé^{1,2,3,5}

Run R code at the command line



Avoid using R interactive mode. Make a code and run it at UNIX command line.

How to interpret microarray data and gene signature



Functional enrichment analysis

□ Also called gene set analysis and GO analysis

It is a way to see the forest (gene sets—GO terms, KEGG pathways) through the trees (individual genes)

	 Functional enrichment analysis What are the prevalent biological themes in the gene list? Which gene sets are differentially regulated? A way to distill the gene list down to a more digestible level. 	Short list of prevalent gene sets
Long list of genes		

□ Two types of functional enrichment analysis (and other interpretation approaches as well)

	Cutoff-based methods	Cutoff-free method
Input	Unordered list of selected genes above a selection cutoff (e.g., p- value<0.05, fold change>2-fold)	Ordered list of all genes, along with their statistics (e.g., p-value, fold change)
Advantage	Simpler and intuitive	Can detect subtle signals

Cutoff-based vs. Cutoff-free functional enrichment analysis



Softwares

Cutoff-based enrichment analysis: DAVID [Huang (2008) Nat Protoc; PMID:19131956]

Cutoff-free enrichment analysis: GSEA [Subramanian (2007) Bioinformatics; PMID:17644558]

GeneTrail [Backes (2007) Nucleic Acid Res; PMID:17526521]

Cutoff-based enrichment analysis (Fisher's exact test)



Cutoff-free functional enrichment analysis (GSEA)



Step 2: Calculation of *p*-value by permutation test



GSEA can detect subtle but coordinate expression changes

OPEN O ACCESS Freely available online

PLos one

Gene Expression Pattern in Transmitochondrial Cytoplasmic Hybrid Cells Harboring Type 2 Diabetes-Associated Mitochondrial DNA Haplogroups

Seungwoo Hwang ¹⁹ , Soo Heon Kwak ²⁹ , Jong Bhak ³ , Hae Sun Kang ² , You Ri Lee ² , Bo Kyung Koo ² , Kyor	۱g
Soo Park ² , Hong Kyu Lee ⁴ , Young Min Cho ² *	

- ❑ Other than the mitochondrial SNPs, the two groups were identical
- □ Therefore, expression profile differences between the two groups were very subtle



Removing redundancy from GO analysis result (GO-Module)



GO-Module: functional synthesis and improved interpretation of Gene Ontology patterns

Xinan Yang¹, Jianrong Li¹, Younghee Lee¹ and Yves A. Lussier^{1,2,*}

Removing redundancy from GO analysis result (GO-Module)



Removing redundancy from GO analysis result (GO-Module)



Identifying differentially expressed PPI subnetworks

Overall scheme



Why identify differentially expressed subnetworks?

- □ Current set of pre-defined pathways (e.g., GO, KEGG) is based on current biological knowledge, which is far from complete.
- Only part of the pathway is usually altered during biological processes.

□ Solution: To identify differentially expressed subnetworks de novo.

Functional enrichment analysis on pre-defined pathways is not sufficient

Interactome mapping of DEGs versus DE subnetworks



Primary databases of PPI

Coverage of human PPIs on major primary databases [De Las Rivas (2010) PLoS Comput Biol; PMID:20589078]



Need a collective database of all these primary databases of PPI \rightarrow Called a meta-database (or consolidated database)

Meta-databases of PPI

MiMI (Michigan Molecular Interactions)

Michigan molecular interactions r2: from interacting proteins to pathways

V. Glenn Tarcea, Terry Weymouth, Alex Ade, Aaron Bookvich, Jing Gao, Vasudeva Mahavisno, Zach Wright, Adriane Chapman, Magesh Jayapandian, Arzucan Özgür, Yuanyuan Tian, Jim Cavalcoli, Barbara Mirel, Jignesh Patel, Dragomir Radev, Brian Athey, David States and H. V. Jagadish*

I2D (Interologous Interaction Database)

Unequal evolutionary conservation of human protein interactions in interologous networks

Kevin R Brown^{*†} and Igor Jurisica^{*†‡}

iRefWeb (Interaction Reference Web)

iRefWeb: interactive analysis of consolidated protein interaction data and their supporting evidence

Brian Turner¹, Sabry Razick^{2,3}, Andrei L. Turinsky¹, James Vlasblom⁴, Edgard K. Crowdy⁵, Emerson Cho¹, Kyle Morrison¹, Ian M. Donaldson^{2,6} and Shoshana J. Wodak^{1,4,7,*}

Identifying differentially expressed subnetworks with GiGA

Graph-based iterative Group Analysis enhances microarray interpretation

Rainer Breitling*^{1,2}, Anna Amtmann¹ and Pawel Herzyk^{2,3}

Step 1: Assign the ranks to all the nodes in PPI



Step 2: Find local minimum nodes



Local minima:

nodes that have a higher rank than their direct neighbors

❑ serve as seed nodes for subnetwork extension

Identifying differentially expressed subnetworks with GiGA

Step 3: Subnetwork extension from rank 1 seed node



Identifying differentially expressed subnetworks with GiGA

Step 4: Calculation of *p*-value of subnetwork by hypergeometric distribution



Step 6: Repeat the procedure from next seed node

Signature comparison using signature databases



Signature comparison using signature databases

Gene s	signature tables from pap	pers	DB of signatures		
Paran pr $P-va$ pr $p-va$ pr pr pr 0.0002 pr H 0.0006 $H1$ $1.20E-1$ M $p < m$ $H1$ $2.40E-05$ $CGI-69$ $CGI-69$	De II. Top 30 Up-regulated Genes Distinguishing N Predictor genes roteasome 26S subunit, ATPase, 5 ytochrome c oxidase subunit VIa polypeptide 1 haperonin containing TCP1, subunit 3 rohibitin uman D9 splice variant B mRNA roteasome subunit, β , type 4 ydroxyacyl-coenzyme A dehydrogenase, type II eptidylprolyl isomerase A denosine deaminase, RNA-specific icCN5-like 1 htochondrial ribosomal protein L12 100 1.425 -1.5 1.5	MD from <u>P value</u> 2.774 1.992 1.983 1.803 1.753 1.733 1.697 1.662 1.654 1.591 1.493 35	Advantages of publication-derived signature DBs Advantages of publication-derived signature DBs Directly import the end results from expert original analysis of individual studies Can always obtain signatures from papers Can obtain signatures from complex experimental design beyond simple two-class comparison (e.g., survival analysis, tissue		



GeneSigDB: a manually curated database and resource for analysis of gene expression signatures

Aedín C. Culhane^{1,2,*}, Markus S. Schröder¹, Razvan Sultana¹, Shaita C. Picard¹, Enzo N. Martinelli¹, Caroline Kelly¹, Benjamin Haibe-Kains^{1,2}, Misha Kapushesky³, Anne-Alyssa St Pierre¹, William Flahive¹, Kermshlise C. Picard¹, Daniel Gusenleitner¹, Gerald Papenhausen¹, Niall O'Connor¹, Mick Correll¹ and John Quackenbush^{1,3,4,*}

Liverome: a curated database of liver cancer-related gene signatures with self-contained context information

specificity analysis)

Langho Lee¹, Kai Wang², Gang Li², Zhi Xie², Yuli Wang², Jiangchun Xu², Shaoxian Sun², David Pocalyko², Jong Bhak³, Chulhong Kim³, Kee-Ho Lee⁴, Ye Jin Jang⁵, Young II Yeom⁵, Hyang-Sook Yoo^{6*}, Seungwoo Hwang^{1*}

29

Comparison of two signatures



OrderedList—a bioconductor package for detecting similarity in ordered gene lists

Claudio Lottaz^{1,*}, Xinan Yang^{1,2}, Stefanie Scheid¹ and Rainer Spang¹

Calculating similarity score between signatures and *p*-value

Sorted list 1 Sorted list 2		Rank i	Cumulative no. of matches at rank <i>i</i>	Weight at rank <i>i</i> $W_i = 1/e^{\alpha \cdot i}$	Cumulative weighted no. of matches at rank <i>i</i>	
COL1A2	MMP14	1	0	$1/\exp(\alpha)$	$0 \times 1/\exp(\alpha)$	
	A2M	2	1	$1/\exp(2\alpha)$	$1 \times 1/\exp(2\alpha)$	
SPARC		3	2	$1/\exp(3\alpha)$	$2 \times 1/\exp(3\alpha)$	
VIM 🚺 /	MST1	4	2	$1/\exp(4\alpha)$	$2 \times 1/\exp(4\alpha)$	
A2M	COL1A2	5	4	1/exp(5α)	4 x 1/exp(5α)	
MST1	VIM	N	Up-match sum	1/exp(<i>Να</i>)	Sum x 1/exp(<i>Nα</i>)	
MT3	SOD1	N	Down-match sum	1/exp(<i>Nα</i>)	Sum x 1/exp(<i>Na</i>)	
					🛉 🚺	
RHOB —	RHOB	5	2	1/exp(5α)	2 x 1/exp(5α)	
RND3	SERPINE1	4	1	1/exp(4α)	1 x 1/exp(4α)	
SERPINE1	GAPDH	3	0	1/exp(3α)	0 x 1/exp(3α)	
SOD1	MT3	2	0	1/exp(2α)	0 x 1/exp(2α)	
VTN	PLG	1	0	1/exp(α)	$0 \times 1/\exp(\alpha)$	
					ĮĻ	
Obtain overall similarity score between signatures 1 and 2						
Calculate <i>p</i> -value of the						
similarity score by permutation						
List of major common genes that contribute to the score the most						

Deriving a consensus signature from several signatures



RankAggreg, an R package for weighted rank aggregation Vasyl Pihur, Susmita Datta and Somnath Datta*

Rank aggregation methods Shili Lin*

Book recommendations: Microarray





Ο

Statistics and Data Analysis for Microarrays <u>using R and Bioconductor</u> (**2011, 2nd Edition**) X Data Analysis Tools for DNA Microarrays (2003, 1st Edition)

Book recommendations: R

Teach how to program in R (Useful for bioinformaticians)



Teach how to do statistics in R (Not appropriate for bioinformaticians)



... And may others