Introduction of Bioinformatics Fall 2007

Part II: Microarray Data Analysis

Biological Principle Behind Microarray Experiments



Base complementary A ↔ T, U
C ↔ G

- Gene expression experiments measure the amount of mRNA to see which genes are being expressed in the cell.
- Measuring protein levels directly is also possible, but is currently harder.

Transcription



From DNA to mRNA

Reverse Transcription

Clone cDNA strands, complementary to the mRNA





Microarray Technology

- Microarray technology takes advantage of hybridization properties of nucleic acid
 - probes: complementary molecules (e.g. cDNAs) attached to a solid surface
 - target: the specific nucleic acid transcripts of interest



 A specific scanner is used to measure the amount of hybridized target at each probe, which is reported as an intensity.

Microarray Platforms

- Two main classes of platforms:
 - High-density oligonucleotide array
 - Two-color spotted array

Two-color Spotted Arrays

- Two-color (two-channel) spotted platforms:
 - Two colors represent the two samples (experiment and reference) competitively hybridized
 - Each spot has "red" and "green" measurements associated with it.

Two-color Spotted Arrays



Two-color Spotted Arrays



<u>E</u> ie <u>E</u> dit	<u>Y</u> iew <u>I</u> nse	art F <u>o</u> rmat	<u>T</u> ools <u>D</u> a	ata S-PULG	<u>W</u> indow	RExcel <u>H</u> a	ap <u>A</u> rrayTo	iols		
🖻 🖪 🔒	e 🔁 🛃	🗟 🌄 🐰		🤮 Σ -	24 🛍 (?) ? Ari	-l	- 10	- B 2	<u> </u>
A1	-	& NAME								
A	В	С	D	E	F	G	Н		J	þ
NAME	TYPE	ACC	CLID	SPOT	svec77	syde77	svec77	svec78	sveic78	svec7
					CH1D	CH2D	FLAG	CH1D	CH2D	FLAG
zinc finger	cDNA	AA406467	753234	1	1 B46	208B	0	6660	10328	
small proli	cDNA	AA447835	B13614	577	352	492	0	527	275	
zinc finger	cDNA	T57959	71626	1153	2601	1883	0	2677	1685	
486544	cDNA	AA043334	486544	1729	3527	2745	0	2059	1431	
zinc finger	cDNA	H17047	50794	2	3444	3039	0	6830	9445	
small indu	cDNA	H52985	206633	578	B42	1292	0	1330	1514	
Human PC	cDNA	AA425602	768644	1154	548	665	0	1275	860	
amall indu	cDNA	AA425102	768561	1730	4951	1797	Ū	2714	1065	
ESTs, Higi	cDNA	W16724	302190	3	2170	1897	0	6442	13371	
Sjogren sy	cDNA	H29484	49970	579	3710	2356	0	10049	2749	
zind finger	cDNA	AA088564	511814	1155	3106	189D	0	3429	2074	
signal reco	:cDNA	AA411407	754998	1731	3538	3169	0	3523	1513	
Down synd	cDNA	H19439	51408	4	1694	1256	0	6505	6343	
sex hormo	cDNA	T69346	82871	580	387	676	0	502	406	
alpha thala	cDNA	AA410435	753430	1156	94	185	0	929	798	
selectin L	cDNA	H00756	149910	1732	1338	949	0	900	764	
wingless-ty	cDNA	W49672	324901	5	6690	3050	0	8633	5541	
selectin E	cDNA	H3956D	186132	581	207	517	0	254	954	
wingless-ty	cDNA	T99653	122762	1157	272	2075	0	365	1404	
sarcoglyca	cDNA	AA234982	866829	1733	476	588	0	654	529	
von Hippel-	cDNA	H73054	234856	6	2517	2255	0	3252	2177	
steraid sul	cDNA	H15215	49591	582	5075	532B	0	7069	5692	
visinin-like	cDNA	H65066	210575	1158	1 D98	666	0	1046	448	
SRY (sex-	cDNA	AA400739	753184	1734	5D26	707D	0	2180	3533	
ESTs, Moo	cDNA	H69834	213260	2305	112	374	0	420	274	
phosphoin	cDNA	AA464765	B10372	2881	B79	409	1	867	87	
1.1	THE R. L. P.		400040			4040			1	

High-density Oligonucleotide Array

- 寡聚核苷酸晶片 ~ 1.28 cm²
- Contain one set of probe-level data per microarray; some probes for specific finding and others for nonspecific finding



High-density Oligonucleotide Array



*http://keck.med.yale.edu/affymetrix/technology.htm

High-density Oligonucleotide Array

- **PM (Perfect Match):** The perfect match probe has a sequence exactly complimentary to the particular gene.
- **MM (Mismatch):** The mismatch probe differs from the perfect match probe by a single base substitution at the center base position.
- $\Sigma w_k(PM_k-MM_k)^+$ quantifies the expression level of a gene (MAS 5.0).

Affymetrix GeneGhips

.dat file: a huge image file .cel file: cell intensity file

Microsoft Excel - Project.xls											
:	檔案(E) 編輯	(E) 檢視(Ÿ)	插入①	格式(2)) 工具(<u>T</u>)	資料(D) 前	見窗(W) RE:	xcel Stanford	. Tools 說明	(<u>H) A</u> rrayTo	ols
1	💕 🛃 👌 (🛋 🛍 🖻	- 12) -	Σ -	🛍 🕜 SAN	4 SAM Contro	ller 🚆 新細	明體	- 12	- B 2	Ū∣≣
: 🛅 🖄 🖄 🕼 🍋 🖄 📨 🏷 🤰 📭 📭 🖤 回覆變更(C) 結束檢閱(U)											
	К1	-	fx								
		А			В	С	D	E	F	G	Н
	Click to	display the	data	c	chipC-	chipC-	chipC-	Missing	P-Value	Rank	Variance
1	Probe set			, r	epl 🔽	^{rep2} 🔽	^{rep3} 🔽	F) 🔽) 🕞	
2	141200_at				11.71562	11.73581	11.73581	0			
3	141201_at				10.70887	10.66056	10.80053	0			
4	141202_at				11.62488	11.62488	11.62488	0			
5	141203_at				9.96371	9.954425	9.954425	0			
6	141204_at				9.01742	9.034827	9.034827	0			
- 7	141205_at				12.04037	12.21155	12.04684	0			
8	141206 at				11.2604	11.25606	11.2604	0			

Microarray Result

- Probe-level data: the intensities read for each of the components.
- Genomic-level data: the measures being used in real research.

Preprocessing

 Preprocessing is the process to account for various sources of variation and make the probe-level data to genomic-level data

• Statistical analyses will be applied to genomic-level data.

Outline for Part II

- 1. Introduction of Microarry
- 2. Statistical Analyses:
 - Distance Measures in DNA Microarray Data Analysis
 - Cluster Analysis of Genomic Data
 - Analysis of Differential Gene Expression Studies
 - Multiple Testing Procedures and Applications to Genomics
 - Machine Learing Concepts and Tools for Statistical Genomics
- 3. Data Visualization
- 4. Preprocessing
 - Preprocessing High-density Oligonucleotide and Two-Color Spotted Arrays
 - Preprocessing SELDI-TOF Mass Spectrometry Protein Data