

Basics in Bioinformatics

Heping Zhang

What is Bioinformatics?

The use of techniques, including applied mathematics, informatics, statistics, computer science, artificial intelligence, chemistry, and biochemistry, to solve biological problems usually on the molecular level.

- ## Important Topics in Bioinformatics
- ◆ Sequence alignment
 - ◆ Sequencing and sequence assembly
 - ◆ Microarray data analysis

- ## Sequence Alignment
- ◆ A sequence alignment is a way of arranging the primary sequences of DNA, RNA, or protein to identify regions of similarity that may be a consequence of functional, structural, or evolutionary relationships between the sequences
 - ◆ Finding sequence similarities with genes of known function is a common approach to infer a newly sequenced gene's function
 - ◆ In 1984 Russell Doolittle and colleagues found similarities between cancer-causing gene and normal growth factor (PDGF) gene

Alignment of Two Sequences

Aligned sequences of nucleotide or amino acid residues are typically represented as rows within a matrix. Gaps are inserted between the residues so that residues with identical or similar characters are aligned in successive columns.

V = ATCTGATG n = 8
W = TGCATAC m = 7

V	A	T	-	C	-	T	G	A	T	G
W	-	T	G	C	A	T	-	A	-	C

Labels: match (A-T, C-A), mismatch (T-G), insertion (G in W), deletion (A in V), indels (gaps).

4 matches
1 mismatches
2 insertions
2 deletions

Alignment : $2 * k$ matrix ($k > m, n$)

- ## Global Alignment vs. Local Alignment
- Calculating a **global alignment** is a form of global optimization that "forces" the alignment to span the entire length of all query sequences.
 - **Local alignments** identify regions of similarity within long sequences that are often widely divergent overall.

Needleman and Wunsch Algorithm

Step 1: Dot Matrix

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	1												
J					1								
C			1					1		1			
N				1									
R						1						1	
C			1					1		1			
K													
C			1					1		1			
R						1						1	
B		1											
P													1

9/24/2007 <http://www2.cs.uh.edu/~fofanov/> 7

Needleman and Wunsch Algorithm

Step 2: Dynamic Programming

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	1												
J					1								
C			1					1		1			
N				1									
R						1		4	3	3	2	2	0
C			1					3	4	3	3	1	0
K								3	3	3	3	2	1
C			1					2	2	3	2	3	1
R								2	1	1	1	1	2
B		1						1	1	1	1	1	0
P								0	0	0	0	0	1

$$M_{ij} = M_{ij} + \max_{k>j} (M_{ij+1}, M_{i+1k})$$

9/24/2007 <http://www2.cs.uh.edu/~fofanov/> 8

Needleman and Wunsch Algorithm

Step 3: Trace Back

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	8	7	6	6	5	4	4	3	3	2	1	0	0
J	7	7	6	6	6	4	4	3	3	2	1	0	0
C	6	6	7	6	5	4	4	4	3	3	2	0	0
J	6	6	6	5	6	4	4	3	3	2	1	0	0
N	5	5	5	6	5	4	4	3	3	2	1	0	0
R	4	4	4	4	4	4	4	3	3	2	2	0	0
C	3	3	4	3	3	3	3	4	3	3	1	0	0
K	3	3	3	3	3	3	3	3	3	2	1	0	0
C	2	2	3	2	2	2	2	3	2	3	1	0	0
R	2	1	1	1	1	2	1	1	1	1	2	0	0
B	1	2	1	1	1	1	1	1	1	1	1	0	0
P	0	0	0	0	0	0	0	0	0	0	0	1	0

Start from the biggest element in upper row or Left column.

At each step move one row and column to the lower right and pick maximum in this row or column.

9/24/2007 <http://www2.cs.uh.edu/~fofanov/> 9

Needleman and Wunsch Algorithm

Results

A	B	C	N	J	R	Q	C	L	C	R	P	M	
A	J	C		J	N	R		C	K	C	R	B	P

A	B	C		N	J	R	Q	C	L	C	R	P	M
A	J	C	J	N		R		C	K	C	R	B	P

9/24/2007 <http://www2.cs.uh.edu/~fofanov/> 10

Smith and Waterman Algorithm

Step 1: Dot Matrix

	H	E	A	G	A	W	G	H	E	E
P										
A										
W										
H										
E										
A										
E										

Scoring system

- Match: 1
- Mismatch: -1/3
- Gap penalty: -(1+1/3*k)

9/24/2007 <http://www2.cs.uh.edu/~fofanov/> 11

Smith and Waterman Algorithm

Step 2: Dynamic Programming

Three possible ways to end an alignment at current cell:

- extend from last stage adding one base from both sequences
- sliding along first sequence causing gap(s) on second sequence
- sliding along second sequence causing gap(s) on first sequence

	H	E	A	G	A	W	G	H	E	E
P	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
A	0.00	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00
W	0.00	0.00	0.00	0.00	0.00	2.00	0.67	0.33	0.00	0.00
H	0.00	0.00	0.00	0.00	0.00	0.67	1.67	?		
E										
A										
E										

$$H_{ij} = \max\{H_{i-1,j-1} + s(a_i b_j), \max\{H_{i-k,j} - P_k\}, \max\{H_{i,j-l} - P_l\}, 0\}$$

9/24/2007 <http://www2.cs.uh.edu/~fofanov/> 12

Smith and Waterman Algorithm

Step 3: Trace Back

	H	E	A	G	A	W	G	H	E	E
P	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
A	0.00	0.00	1.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00
W	0.00	0.00	0.00	0.00	0.00	2.00	0.67	0.33	0.00	0.00
H	0.00	0.00	0.00	0.00	0.33	0.67	1.67	-1.67	0.33	0.00
E	0.00	1.00	0.00	0.00	0.00	0.33	0.33	1.33	2.67	1.33
A	0.00	0.00	2.00	0.67	1.00	0.00	0.00	0.00	1.33	2.33
E	0.00	1.00	0.67	1.67	0.33	0.67	0.00	0.00	1.00	2.33

Best local alignment: **AWGHE**
AW-HE

9/24/2007

<http://www2.cs.uh.edu/~fofanov/>

Multiple Alignments

- A **multiple sequence alignment** is a sequence alignment of three or more biological sequences
- Reason to do multiple sequence alignment
 - The sequences *may* share a common origin - a common ancestor sequence. If the similarity is sufficiently convincing or if we have additional evidence for an evolutionary relationship, then we say that the sequences are homologous.
 - The sequences *may* have the same or related structure and function

9/24/2007

<http://www2.cs.uh.edu/~fofanov/>

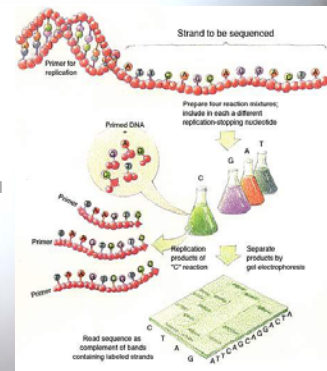
DNA Sequencing

- The term **DNA sequencing** encompasses biochemical methods for determining the order of the nucleotide bases, adenine (A), guanine (G), cytosine (C), and thymine (T), in a DNA oligonucleotide.
- Techniques
 - Chain-termination methods
 - Large scale sequencing strategies
 - Sequencing by hybridization

9/24/2007

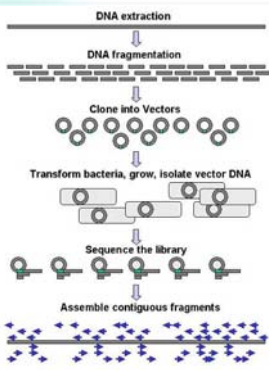
Sanger Method

- Start at primer (restriction site)
- Grow DNA chain
- Include ddNTPs
- Stops reaction at all possible points
- Separate products by length, using gel electrophoresis



9/24/2007

Shotgun Sequencing



- Suitable for longer sequences
- DNA is broken up randomly into numerous small segments, which are sequenced using the chain termination method to obtain *reads*.
- Multiple overlapping reads for the target DNA are obtained by performing several rounds of this fragmentation and sequencing
- Computer programs then use the overlapping ends of different reads to assemble them into a contiguous sequence.

9/24/2007

http://en.wikipedia.org/wiki/DNA_sequencing

Fragment Assembly

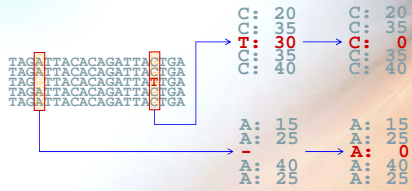
- Assemble individual short fragments (reads) into a single genomic sequence
- The **“OLC Framework”**
 - Overlap**: find all the overlaps between the reads that satisfy certain quality criteria.
 - Layout**: given the set of overlap relationships between the reads, determine a consistent layout of the reads, i.e., find a consistent tiling of all the reads that preserves most of the overlap constraints
 - Consensus**: given a tiling of reads determined in the layout stage, determine the most likely DNA sequence (the consensus sequence) that can be explained by the tiling

9/24/2007

Shotgun Sequence Assembly MIHAI POP TIGR

Finding Overlapping Reads

- Correct errors using multiple alignment



- Score alignments
- Accept alignments with good scores

9/24/2007

19

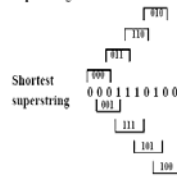
Layout

- A variation of the shortest common superstring problem

The Shortest Superstring problem

Set of strings: {000, 001, 010, 011, 100, 101, 110, 111}

Concatenation Superstring 000 001 010 011 100 101 110 111



– Problem: Given a set of strings, find a shortest string that contains all of them

– Complexity: NP-complete

9/24/2007

20

SSP to TSP (Traveling Salesman Problem)

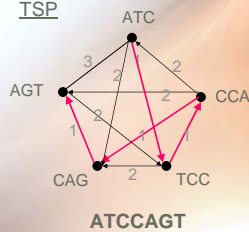
Create an overlap graph G where each vertex represents a string. Edge between string s_1 and s_2 will have weight equal to the length of s_1 minus the overlap of s_1 with s_2 . (Edge weight is not symmetric) The path visiting all the vertices of minimum total weight defines the shortest common superstring.

$S = \{ ATC, CCA, CAG, TCC, AGT \}$

SSP

AGT
CCA
ATC
ATCCAGT
TCC
CAG

TSP

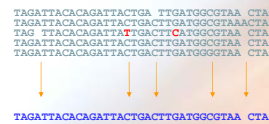


9/24/2007

21

Derive Consensus Sequence

Derive each consensus base by weighted voting



9/24/2007

22

Sequencing of the Human Genome

◆ Aims

- Sequence the entire 3 billion DNA bases
- Dissect the code of estimated 25,000 genes that determine the physical characters of the human body

9/24/2007

23

The Human Genome Project

- Funded by the National Institutes of Health in the United States, and the UK charity, the Wellcome Trust, and numerous other groups from around the world
 - Total cost \$3 billion
 - Started in 1990
- The genome was broken into smaller pieces
 - Approximately 150,000 base pairs in length
 - These pieces are called "bacterial artificial chromosomes", or BACs, because they can be inserted into bacteria where they are copied by the bacterial DNA replication machinery
 - Each of these pieces was then sequenced separately as a small "shotgun" project and then assembled
 - The larger, 150,000 base pairs go together to create chromosomes
 - This is known as the "hierarchical shotgun" approach
 - The genome is first broken into relatively large chunks, which are then mapped to chromosomes before being selected for sequencing
- Paper published in *Nature*

9/24/2007

24

Celera Genomics

- ◆ A similar, privately funded quest was launched by the American researcher Craig Venter and his firm Celera Genomics
 - ◆ Total cost: \$300 million
 - ◆ Started in 1998, intended to proceed at a faster pace and at a fraction of the cost of the publicly funded project.
- ◆ Used a riskier technique called whole genome shotgun sequencing

Whole genome shotgun had been used to sequence bacterial genomes of up to six million base pairs in length, but not for anything nearly as large as the three thousand million base pair human genome
- ◆ Draft genome published at the same time as the public effort did
- ◆ Paper published in *Science*

9/24/2007

25

Sequencing by Hybridization (SBH)

• History

- 1988: SBH suggested as an alternative sequencing method. Nobody believed it would ever work
- 1991: Light directed polymer synthesis developed by Steve Fodor and colleagues.
- 1994: Affymetrix develops first 64-kb DNA microarray

First microarray prototype (1989)



First commercial DNA microarray prototype w/16,000 features (1994)



500,000 features per chip (2002)



9/24/2007

26

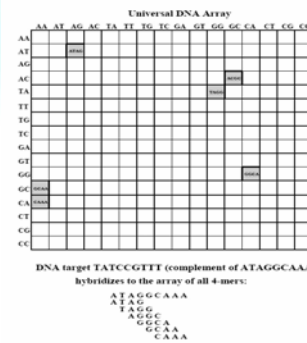
How SBH Works

- Attach all possible DNA probes of length l to a flat surface, each probe at a distinct and known location. This set of probes is called the DNA array.
- Apply a solution containing fluorescently labeled DNA fragment to the array.
- The DNA fragment hybridizes with those probes that are complementary to substrings of length l of the fragment.
- Using a spectroscopic detector, determine which probes hybridize to the DNA fragment to obtain the l -mer composition of the target DNA fragment.
- Apply the combinatorial algorithm (below) to reconstruct the sequence of the target DNA fragment from the l -mer composition.

9/24/2007

27

Hybridization on DNA Array



9/24/2007

28

l -mer composition

- **Spectrum** (s, l) - unordered multiset of all possible $(n - l + 1)$ l -mers in a string s of length n
- The order of individual elements in $\text{Spectrum}(s, l)$ does not matter
- For $s = \text{TATGGTGC}$ all of the following are equivalent representations of $\text{Spectrum}(s, 3)$:
 - {TAT, ATG, TGG, GGT, GTG, TGC}
 - {ATG, GGT, GTG, TAT, TGC, TGG}
 - {TGG, TGC, TAT, GTG, GGT, ATG}

We usually choose the lexicographically maximal representation as the canonical one.

9/24/2007

29

The SBH Problem

- Goal: Reconstruct a string from its l -mer composition
- Input: A set S , representing all l -mers from an (unknown) string s
- Output: String s such that $\text{Spectrum}(s, l) = S$

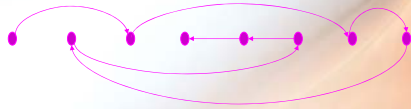
9/24/2007

30

SBH: Hamiltonian Path Approach

$S = \{ \text{ATG AGG TGC TCC GTC GGT GCA CAG} \}$

H ATG AGG TGC TCC GTC GGT GCA CAG



ATGCAGGTCC

Path visited every VERTEX once

9/24/2007

31

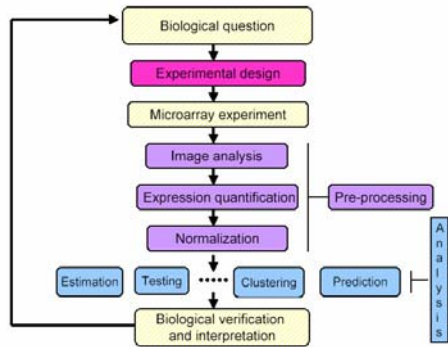
Some Difficulties with SBH

- ◆ **Fidelity of Hybridization:** difficult to detect differences between probes hybridized with perfect matches and 1 or 2 mismatches
- ◆ **Array Size:** Effect of low fidelity can be decreased with longer k -mers, but array size increases exponentially in k . Array size is limited with current technology.
- ◆ **Alternative:** SBH is not competitive due to the above reason
 - ◆ Positional SBH
 - ◆ Shotgun SBH

9/24/2007

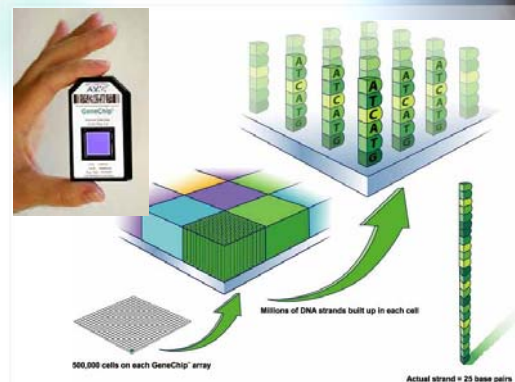
32

Microarray Data Analysis



9/24/2007

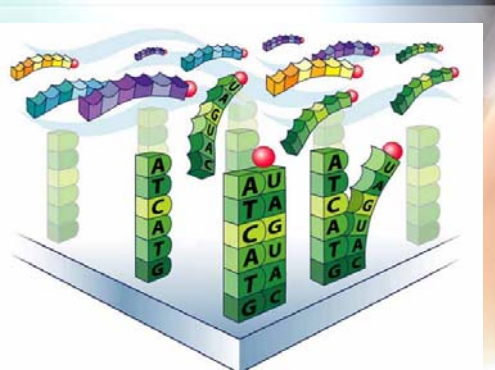
Affymetrix GeneChip



9/24/2007

33

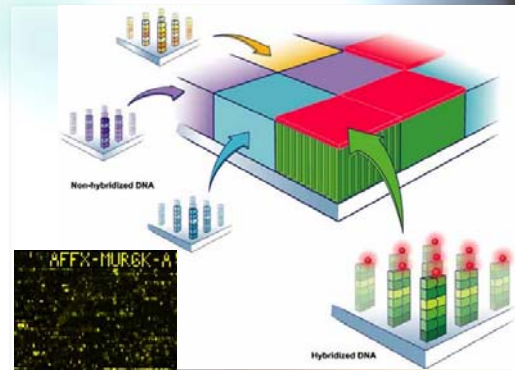
Hybridization



9/24/2007

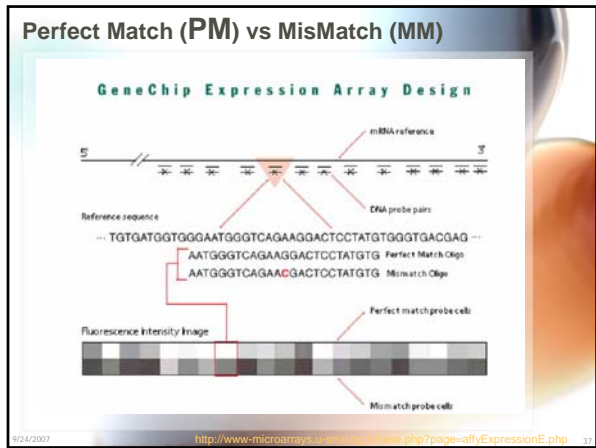
34

Microarray Images

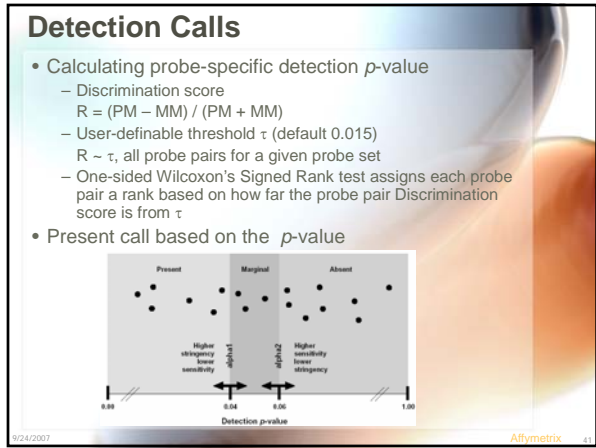
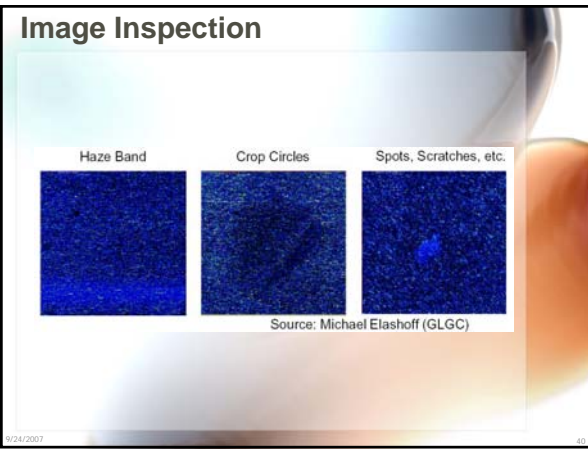
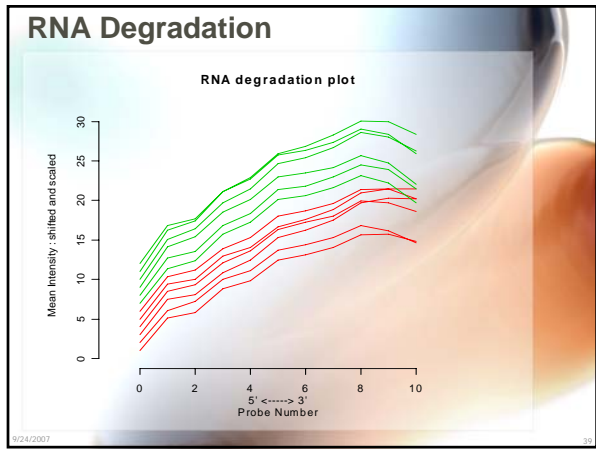


9/24/2007

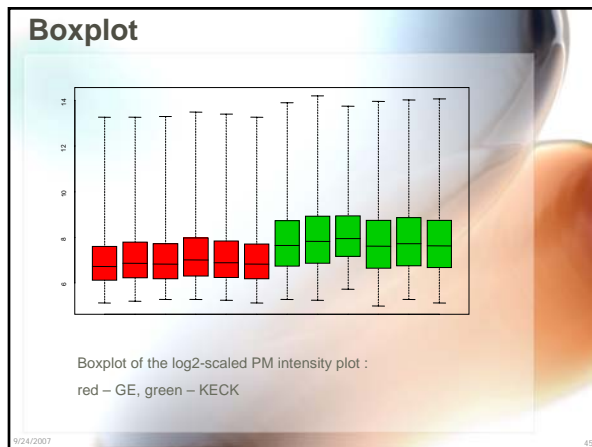
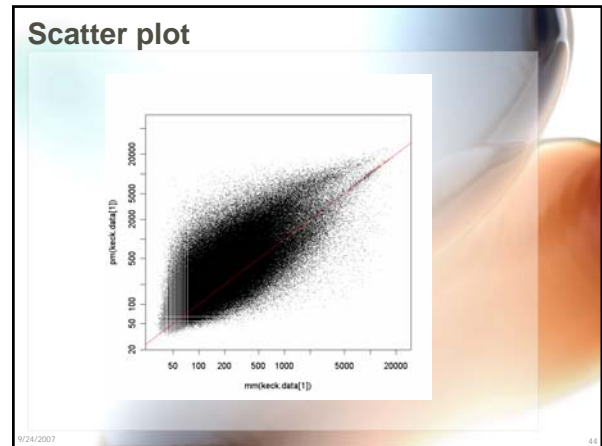
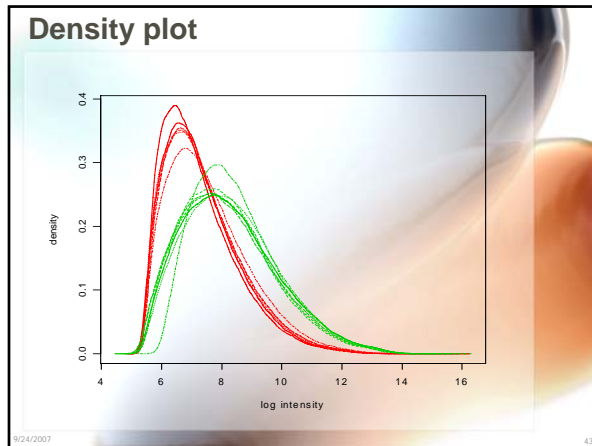
35



- ### Quality Assessment and Control
- RNA sample quality control (e.g., RNA degradation plot)
 - Array hybridization quality control
 - Probe array image inspection
 - Control genes
 - PM/MM present calls
 - Scaling Factors
 - Statistical diagnostics plots
- 9/24/2007 38



- ### Diagnostics plots
- Image plot
 - PM/MM intensity plot
 - Histogram plot
 - Density plot
 - Scatter Plot
 - Boxplot
- 9/24/2007 42



- ### Pre-analysis
- Background correction
 - Normalization
 - Artifacts (outliers) detection and management
- 9/24/2007 46

- ### Background correction
- Why: to eliminate the low levels of noise that are present on any microarray
 - Different methods:
 - Local neighborhood detection
 - Negative control
 - Mismatch probes
 - Model-based: RMA and MBEI
- 9/24/2007 <http://www.genome.gov/cancer/emory.edu/analysis/normal/vmsr> 47

- ### Normalization
- Why normalization?
 - Minimizing non-biological factor
 - Reducing unwanted variance across chips
 - Different methods:
 - All genes on a slide (global normalization)
 - Constantly expressed genes (Housekeeping)
 - Set of control genes
 - e.g., Mouse430 chip, 1415670_at to 1415769_at
 - Rank-invariant gene
 - Quantile: giving each chip the same empirical distribution; reducing variance w/o introducing drastic bias effects
- 9/24/2007 48

Artifacts & Outliers Detection

- Different methods apply different algorithms to detect outliers (artifacts) and take different actions;
- For example, Li-Wong model-based (dChip) method: identify extreme residuals, remove them, re-fit, ..., converge.

9/24/2007

49

Low-level analysis

- Generally, refer to **probe-level analysis** for Affymetrix Chips: how to extract gene expressions from probe data
- Three common approaches:
 - **MAS 5.0** (MicroArray Suite Version 5)
 - **Li-Wong model-based analysis** (dChip)
 - **RMA** (Robust multi-chip analysis)

9/24/2007

50

Identifying differentially expressed genes

- Fold change detection;
- Student's *t* test;
- Mann-Whitney U test;
- Multiple comparisons adjusted P-values and confidence intervals.

9/24/2007

51

Microarray High-level Analysis

- **Clustering** (unsupervised learning)
Grouping objects based on their similarity in feature space, e.g., identifying groups of co-regulated genes;
- **Classification** (supervised learning)
Training machine and assigning new cases into known classes, i.e., differentiating tumor / normal cells;
- **Network analysis**
Inferring and building regulatory networks.

9/24/2007

52

Clustering

- Two categories
 - Hierarchical methods
 - Divisive: successively splitting larger clusters, top-down
 - Agglomerative: successively merging smaller clusters, bottom-up
 - Partitional methods
 - Determines all clusters at once
 - e.g., K-means, SOM, etc.

9/24/2007

53

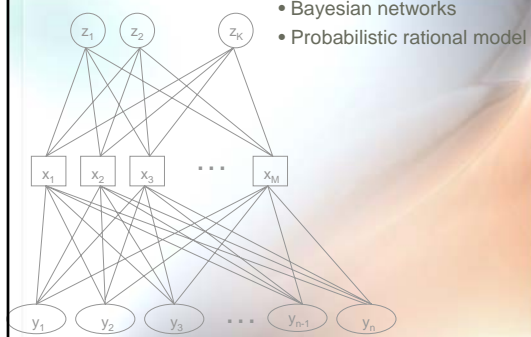
Classification

- k-nearest neighbors (KNN)
- Classification Tree
- Linear discriminant analysis (LDA)
- Bayesian Regression
- Support vector machines (SVM)
- Artificial neural networks (ANN)

9/24/2007

54

Pathway analysis



9/24/2007

55

References

Many slides are from www.bioalgorithms.info

• Simons, Robert W. *advanced Molecular Genetics Course*, UCLA (a00a).

<http://www.mimg.ucla.edu/bobs/Ca59/Presentations/Benz.pdf>

Batzoglou, S. *Computational Genomics Course*, Stanford University (a004).

<http://www.stanford.edu/class/csa6a/handouts.html>

9/24/2007

56