An Introduction to Knowledge Discovery Applications and Challenges in Life Sciences

Limsoon Wong





ICDE'07, Istanbul, Turkey, 16-20 April 2007

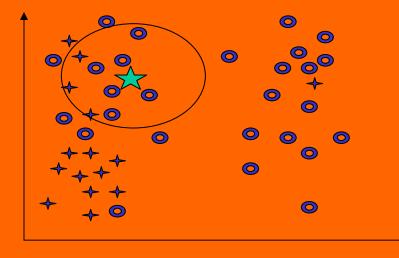
Plan



- Quick introduction to knowledge discovery (10 min)
- Protein function inference (40 min)
- Gene feature recognition (50 min)
- Disease diagnosis, treatment, and understanding (70 min)
- Advancing knowledge discovery (10 min)

10 min

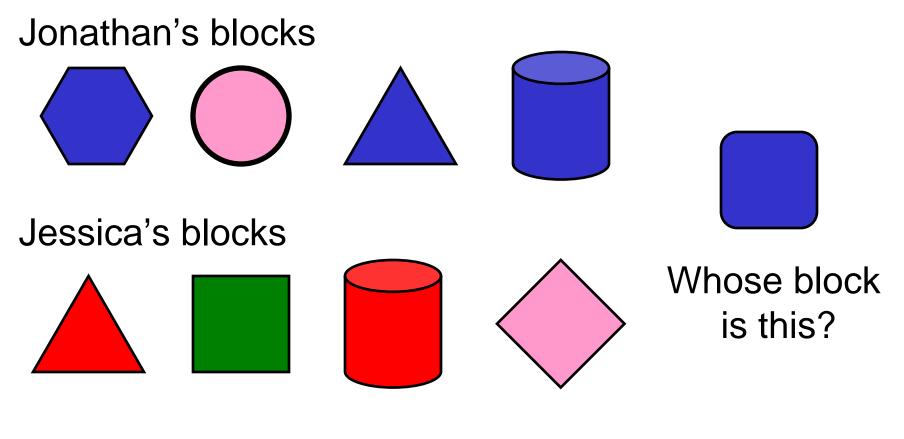
Quick Intro to Knowledge Discovery







What is Knowledge Discovery?



Jonathan's rules Jessica's rules : Blue or Circle : All the rest



What is Knowledge Discovery?









Question: Can you explain how?



Main Steps of Knowledge Discover

- Training data gathering
- Feature generation
 - k-grams, colour, texture, domain know-how, ...
- Feature selection
 - Entropy, χ 2, CFS, t-test, domain know-how...
- Feature integration

- SVM, ANN, PCL, CART, C4.5, kNN, ...

classifier/

methods

Some

ICDE'07, Istanbul, Turkey, 16-20 April 2007

An Example Classifier: kNN (k=8)

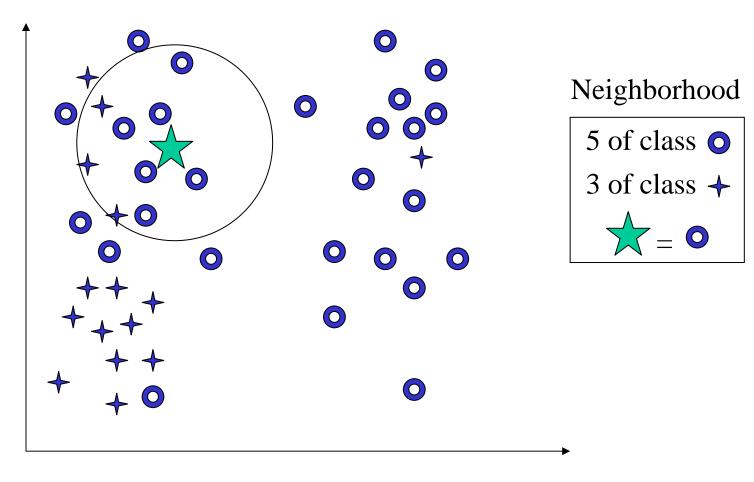
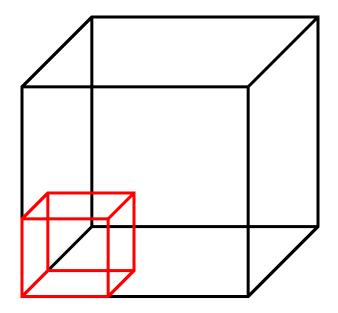


Image credit: Zaki

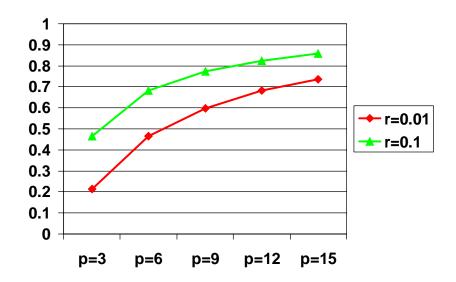


Curse of Dimensionality

• How much of each dimension is needed to cover a proportion r of total sample space?



- Calculate by $e_p(r) = r^{1/p}$
- So, to cover 1% of a 15-D space, need 85% of each dimension!



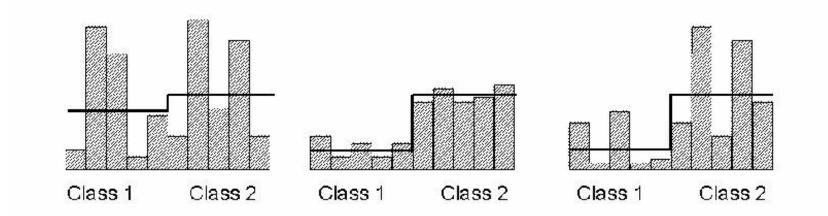
Tackling the Curse: Signal Selection

- Choose a feature w/ low intra-class distance
- Choose a feature w/ high inter-class distance

The t-state of a signal is defined as

$$t = rac{|\mu_1 - \mu_2|}{\sqrt{(\sigma_1^2/n_1) + (\sigma_2^2/n_2)}}$$

where σ_i^2 is the variance of that signal in class *i*, μ_i is the mean of that signal in class *i*, and n_i is the size of class *i*.



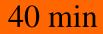


Self-fulfilling Oracle

- Construct artificial dataset with 100 samples, each with 100,000 randomly generated features and randomly assigned class labels
- select 20 features with the best t-statistics (or other methods)

- Evaluate accuracy by cross validation using only the 20 selected features
- The resultant estimated accuracy can be ~90%
- But the true accuracy should be 50%, as the data were derived randomly

What went wrong?



Protein Function Inference





ICDE'07, Istanbul, Turkey, 16-20 April 2007



 Doolittle et al. (Science, July 1983) searched for platelet-derived growth factor (PDGF) in his own DB. He found that PDGF is similar to v-sis oncogene

PDGF-2 1 SLGSLTIAEPAMIAECKTREEVFCICRRL?DR?? 34 p28sis 61 LARGKRSLGSLSVAEPAMIAECKTRTEVFEISRRLIDRTN 100

• A seq alignment maximizes the number of positions that are in agreement in two sequences

Sequence Alignment: Poor Example

Poor seq alignment shows few matched positions
 The two proteins are not likely to be homologous

Alignment by FASTA of the sequences of amloyanin and domain 1 of ascorbate oxidase

60 70 80 100Amicyanin **MPHNVHFVAGVLGEAALKGPMMKKEQAYSLTFTEAGTYDYHCTPHPFMRGKVVVE** 11. 11 Ascorbate Oxidase ILQRGT9WADGTASISQCAINPGETFFYNFTVDNPGTFFYHGHLGMQRSAGLYGSLI 7080 90 100 110120 No obvious match between Amicyanin and Ascorbate Oxidase

ICDE'07, Istanbul, Turkey, 16-20 April 2007

Sequence Alignment: Good Example

- Good alignment usually has clusters of extensive matched positions
- \Rightarrow The two proteins are likely to be homologous

D >gil13476732|ref|NP_108301.1|
 unknown protein [Mesorhizobium loti]
 gil14027493|dbj|BAB53762.1|
 unknown protein [Mesorhizobium loti]
 Length = 105

Score = 105 bits (262), Expect = 1e-22 Identities = 61/106 (57%), Positives = 73/106 (68%), Gaps = 1/106 (0%)

Query: 1 MKPGRLASIALAIIFLPMAVPAHAATIEITMENLVISPTEVSAKVGDTIRWVNKDVFAHT 60 MK G L ++ MA PA AATIE+T++ LV SP V AKVGDTI WVN DV AHT Sbjct: 1 MKAGALIRLSWLAALALMAAPAAAATIEVTIDKLVFSPATVEAKVGDTIEWVNNDVVAHT 60

> good match between Amicyanin and unknown M. loti protein

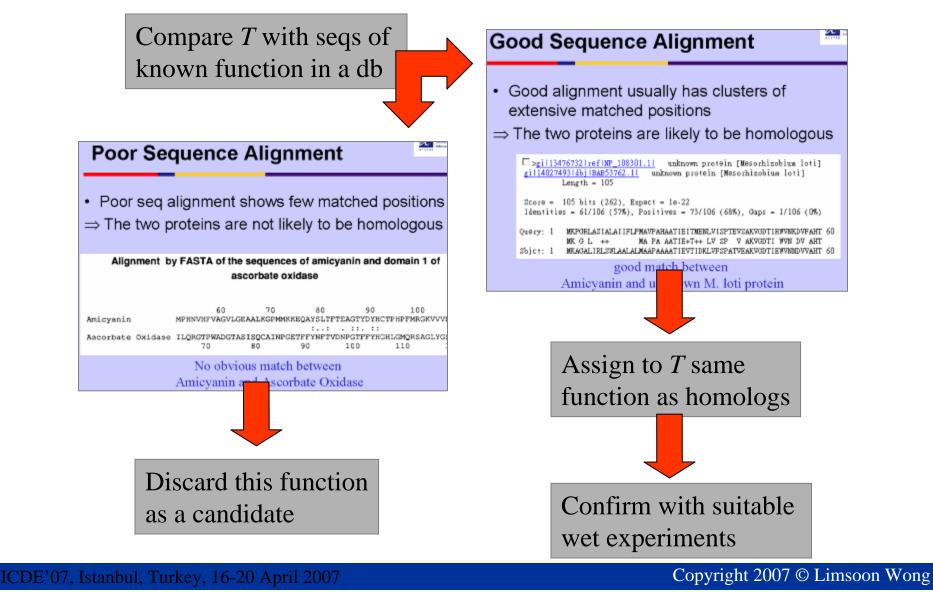


SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR YVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE VT

• How do we attempt to assign a function to a new protein sequence?



Guilt-by-Association

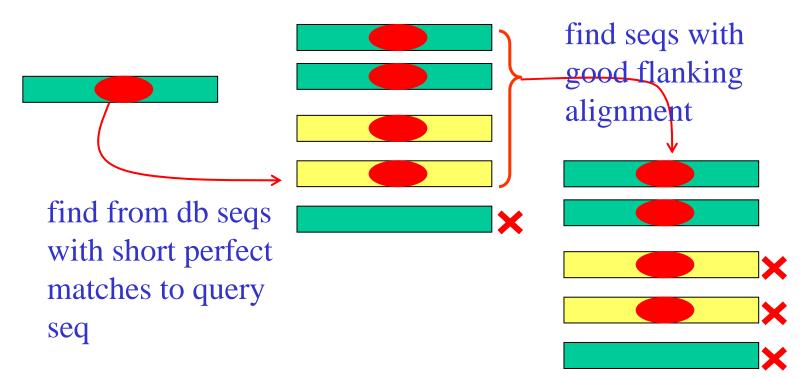


10

BLAST: How It Works Altschul et al., *JMB*, 215:403--410, 1990



 BLAST is one of the most popular tool for doing "guilt-by-association" sequence homology search



G Back • 🕑 •	😰 🚱 🏸 Search 🌟 Favorites 🤣 🖾 • 🧽 🖃 • 🖵 🖏 🦄				
Google -	National University				
Google	G Search 🔹 🧭 🌍 🥵 PageRank 🕸 6 blocked 👋 Check 🔹 🛝 AutoLink 💞 of Singapore				
S NCBI					
Nucleotide	Protein Translations Retrieve results for an RID				
Search NRYVNILPYDHSRVHLTPVEGVPDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWR MIWEQNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFC IQQVGDVTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHC SAGVGRTGTFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLE HYLYGDTELE					
<u>Set subsequence</u>	From: To:				
<u>Choose database</u>	nr 💌				
Do CD-Search					
Now:	BLAST! or Reset query Reset all				
Options Limit by entrez guery	for advanced blasting or select from: All organisms				



Homologs obtained by BLAST

	Score	E
Sequences producing significant alignments:	(bits)	Value
gi 14193729 gb AAK56109.1 AF332081_1 protein tyrosin phosph	<u>62:</u>	e-177
gi 126467 sp P18433 PTRA_HUMAN Protein-tyrosine phosphatase	<u>621 L</u>	e-177
gi 4506303 ref NP_002827.1 protein tyrosine phosphatase, r	<u>621 L</u>	e-176
<u>gi 227294 prf 1701300A</u> protein Tyr phosphatase	<u>620</u>	e-176
gi 18450369 ref NP_543030.1 protein tyrosine phosphatase,	<u>621 L</u>	e-176
gi 32067 emb CAA37447.1 tyrosine phosphatase precursor [Ho	<u>61:</u>	e-176
gi 285113 pir JC1285 protein-tyrosine-phosphatase (EC 3.1	<u>619</u>	e-176
gi 6981446 ref NP_036895.1 protein tyrosine phosphatase, r	<u>61;</u>	e-176
gi 2098414 pdb 1YFO A Chain A, Receptor Protein Tyrosine Ph	<u>61</u> S	e-174
qi 32313 emb CAA38662.1 protein-tyrosine phosphatase [Homo	<u>61</u>	e-174
<pre>gi 450583 qb AAB04150.1 protein tyrosine phosphatase >gi 4</pre>	<u>605</u>	e-172
gi 6679557 ref NP 033006.1 protein tyrosine phosphatase, r	<u>60.</u>	e-172
qi 483922 qb AAA17990.1 protein tyrosine phosphatase alpha	599	e-170

• Thus our example sequence could be a protein tyrosine phosphatase α (PTP α)



Example Alignment with $PTP\alpha$

Score = 632 bits (1629), Expect = e-18U Identities = 294/302 (97%), Positives = 294/302 (97%)

- Sbjct: 202 SPSTNRKYPPLPVDKLEEEINRRMADDNKLFREEFNALPACPIQATCEAASKEENKEKNR 261
- Query: 61 YVNILPYDHSRVHLTPVEGVFDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE 120 YVNILPYDHSRVHLTPVEGVFDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE
- Shjet: 2.62. YVNILPYDHSRVHLTPVEGVEDSDYINASFINGYQEKNKFIAAQGPKEETVNDFWRMIWE 32.1
- Query: 121 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD 180 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD
- Sbjet: 322 QNTATIVMVTNLKERKECKCAQYWPDQGCWTYGNVRVSVEDVTVLVDYTVRKFCIQQVGD 381
- Query: 181 VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG 240 VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG
- Sbjet: 382 VTNRKPQRLITQFHFTSWPDFGVPFTPIGMLKFLKKVKACNPQYAGAIVVHCSAGVGRTG 441
- Query: 241 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE 300 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE
- Sbjct: 442 TFVVIDAMLDMMHSERKVDVYGFVSRIRAQRCQMVQTDMQYVFIYQALLEHYLYGDTELE 501

ICDE'07, Istanbul, Turkey, 16-20 April 2007



Guilt-by-Association: Caveats

- Ensure that the effects of database size and composition have been accounted for
- Ensure that the function of the homology is not derived via invalid "transitive assignment"
- Ensure that the target sequence has all the key features associated with the function, e.g., active site and/or domain



Law of Large Numbers

- Suppose you are in a room with 365 other people
- Q: What is the prob that a specific person in the room has the same birthday as you?
- A: 1/365 = 0.3%

- Q: What is the prob that there is a person in the room having the same birthday as you?
- A: $1 (364/365)^{365} = 63\%$
- Q: What is the prob that there are two persons in the room having the same birthday?
- A: 100%



Interpretation of P-value

- Seq. comparison progs, e.g. BLAST, often associate a P-value to each hit
- P-value is interpreted as prob that a random seq has an equally good alignment

- Suppose the P-value of an alignment is 10⁻⁶
- If database has 10⁷ seqs, then you expect 10⁷ * 10⁻⁶ = 10 seqs in it that give an equally good alignment
- ⇒ Need to correct for database size if your seq comparison prog does not do that!



Lightning Does Strike Twice!

- Roy Sullivan, a former park ranger from Virgina, was struck by lightning 7 times
 - 1942 (lost big-toe nail)
 - 1969 (lost eyebrows)
 - 1970 (left shoulder seared)
 - 1972 (hair set on fire)
 - 1973 (hair set on fire & legs seared)
 - 1976 (ankle injured)
 - 1977 (chest & stomach burned)
- September 1983, he committed suicide



Cartoon: Ron Hipschman Data: David Hand

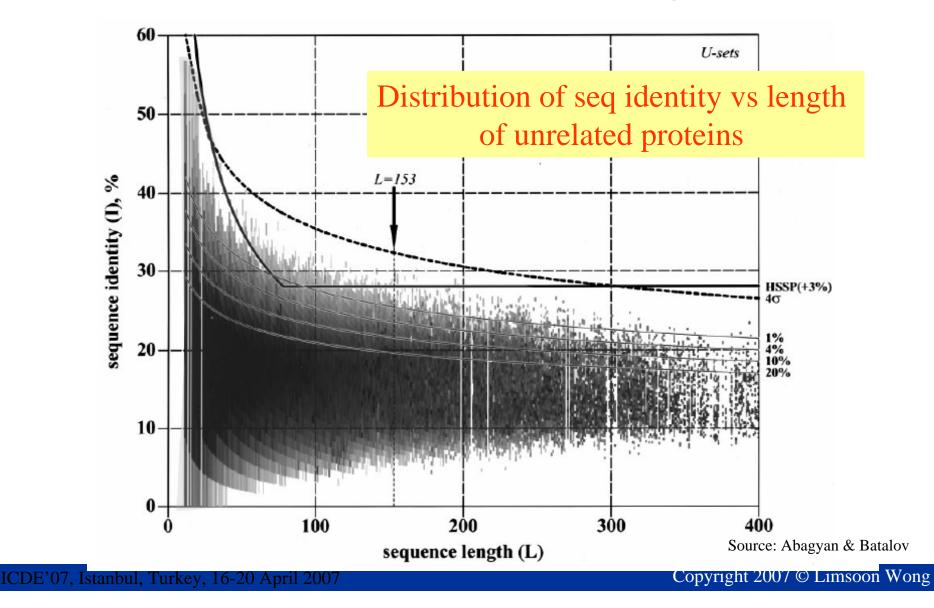
Effect of Seq Compositional Bias

- One fourth of all residues in protein seqs occur in regions with biased amino acid composition
- Alignments of two such regions achieves high score purely due to segment composition
- While it is worth noting that two proteins contain similar low complexity regions, they are best excluded when constructing alignments
- BLAST employs the SEG algorithm to filter low complexity regions from proteins before executing a search

Source: NCBI



Effect of Sequence Length





Important Unsolved Challenges

- What if there is no useful sequence homolog?
- Guilt by other types of association!
 - Domain modeling (e.g., HMMPFAM)
 - Similarity of dissimilarities (e.g., SVM-PAIRWISE)
 - Similarity of phylogenetic profiles
 - Similarity of subcellular co-localization & other physico-chemico properties(e.g., PROTFUN)
 - Similarity of gene expression profiles
 - Similarity of protein-protein interaction partners
 - Fusion of multiple types of info



28

Similarity of Dissimilarities

	orange ₁	banana ₁	
apple ₁	Color = red vs orange Skin = smooth vs rough Size = small vs small Shape = round vs round	Color = red vs yellow Skin = smooth vs smooth Size = small vs small Shape = round vs oblong	
apple ₂	Color = red vs orange Skin = smooth vs rough Size = small vs small Shape = round vs round	Color = red vs yellow Skin = smooth vs smooth Size = small vs small Shape = round vs oblong	
orange ₂	Color = orange vs orange Skin = rough vs rough Size = small vs small Shape = round vs round	Color = orange vs yellow Skin = rough vs smooth Size = small vs small Shape = round vs oblong	••



SVM-Pairwise Framework

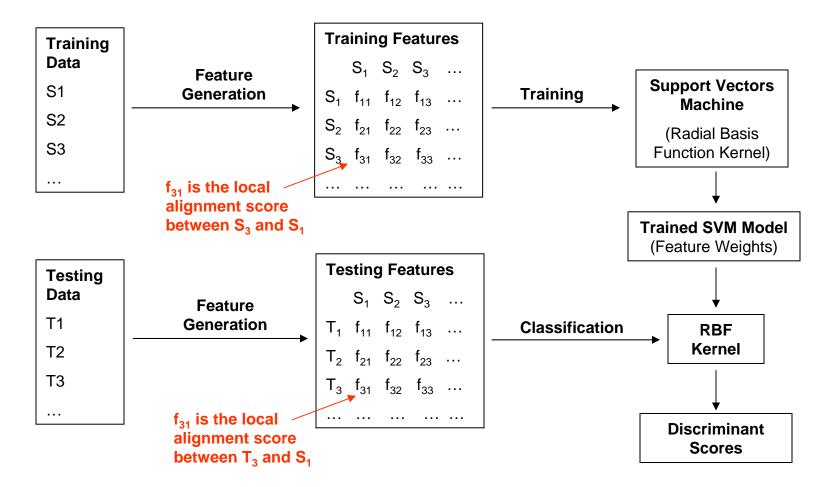
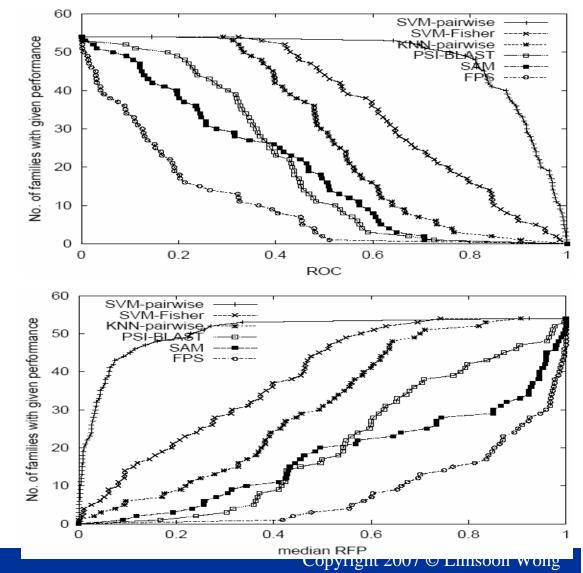


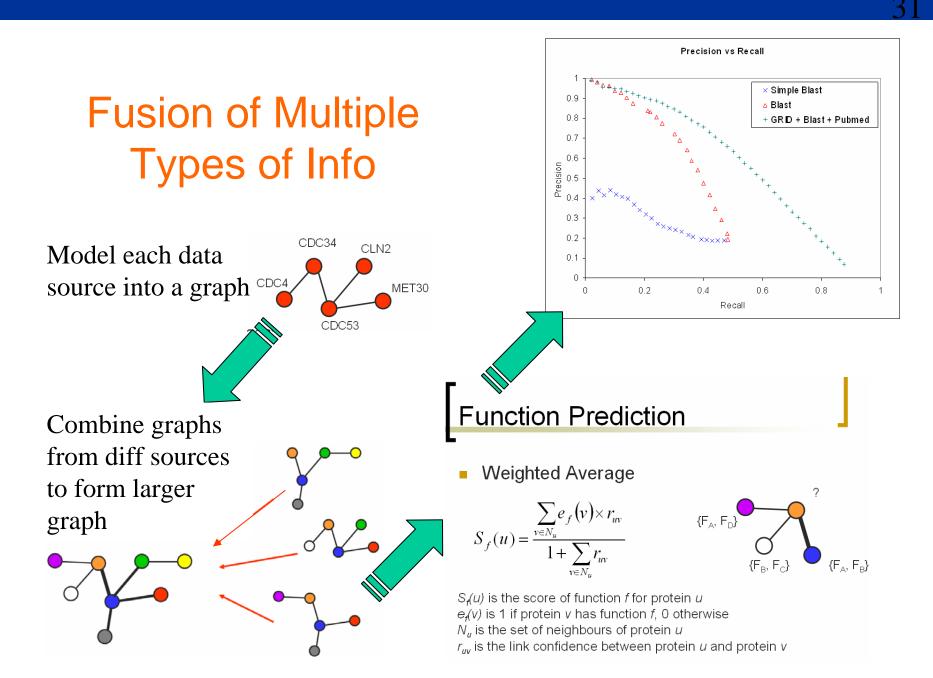
Image credit: Kenny Chua



Performance of SVM-Pairwise

- Receiver Operating
 Characteristic (ROC)
 - The area under the curve derived from plotting true positives as a function of false positives for various thresholds.
- Rate of median False Positives (RFP)
 - The fraction of negative test examples with a score better or equals to the median of the scores of positive test examples.







Recommended Readings

- Wong, *The Practical Bioinformatician*, 2004, ICP. Chapters 10 & 19
- Chua et al., Exploiting indirect neighbours and topological weight to predict protein function from protein-protein interactions. *Bioinformatics*, 22:1623-1630, 2006.
- Bateman et al., The Pfam protein families database. NAR, 32:D138-D141, 2004
- Jaakkola et al., A discriminative framework for detecting remote protein homologies. *JCB*, 7:95-114, 2000
- S.F.Altschul et al. "Gapped BLAST and PSI-BLAST: A new generation of protein database search programs", *NAR*, 25(17):3389--3402, 1997
- Pellegrini et al., Assigning protein functions by comparative genome analysis: Protein phylogenetic profiles. *PNAS*, 96:4285-4288, 1999

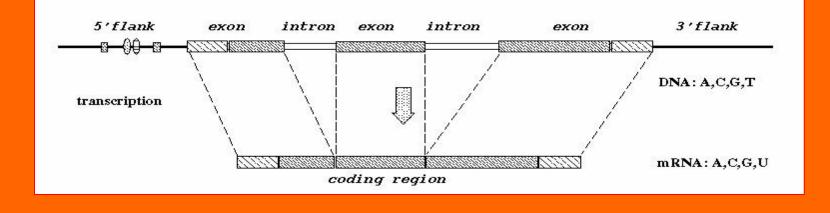
5-Minute Break?



ICDE'07, Istanbul, Turkey, 16-20 April 2007

60 min

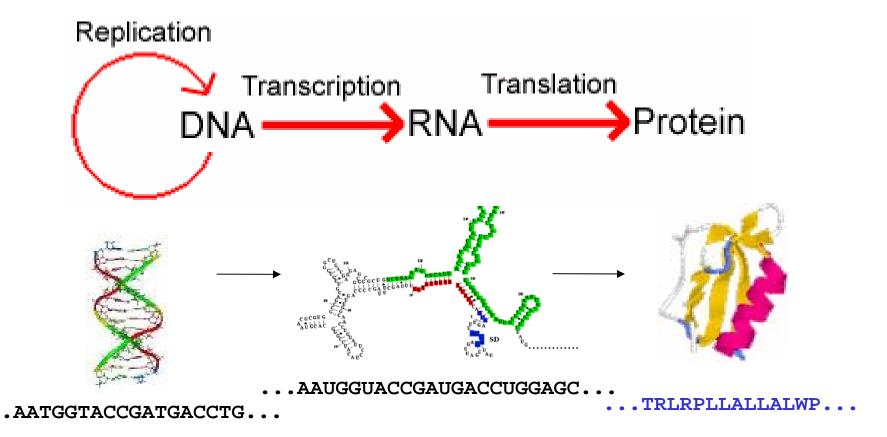
Gene Feature Recognition





Central Dogma

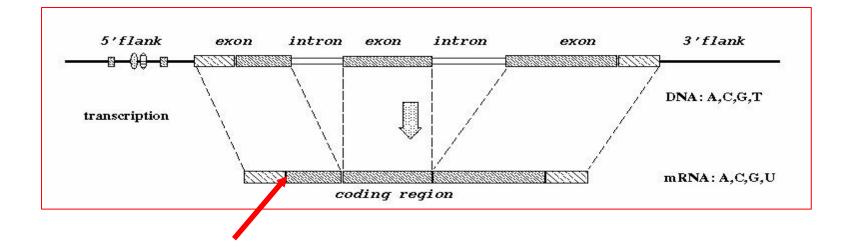






36

Translation Initiation Site





A Sample cDNA

299 HSU27655.1 CAT U27655 Homo sapiens	
CGTGTGTGCAGCAGCCTGCAGCTGCCCCAAGCC <u>ATG</u> GCTGAACACTGACTCCCAGCTGTG	80
CCCAGGGCTTCAAAGACTTCTCAGCTTCGAGC <u>ATG</u> GCTTTTGGCTGTCAGGGCAGCTGTA	160
GGAGGCAG <mark>ATG</mark> AGAAGAGGGAG <mark>ATG</mark> GCCTTGGAGGAAGGGAAGGGGCCTGGTGCCGAGGA	240
CCTCTCCTGGCCAGGAGCTTCCTCCAGGACAAGACCTTCCACCCAACAAGGACTCCCCT	
	80
ieeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeee	160
EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE	240
EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE	

• What makes the second ATG the TIS?

Copyright 2007 © Limsoon Wong



Approach

- Training data gathering
- Signal generation
 - k-grams, distance, domain know-how, ...
- Signal selection
 - Entropy, χ 2, CFS, t-test, domain know-how...
- Signal integration
 - SVM, ANN, PCL, CART, C4.5, kNN, ...



Training & Testing Data

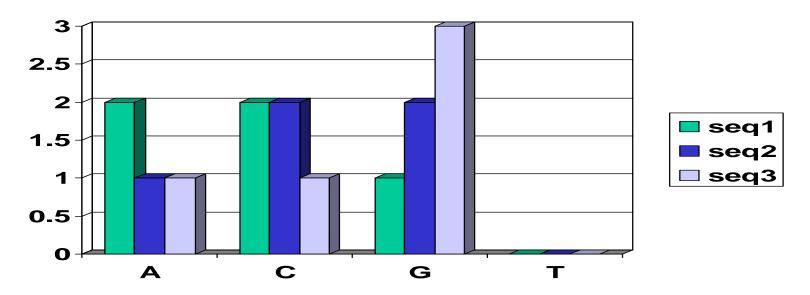
- Vertebrate dataset of Pedersen & Nielsen [ISMB'97]
- 3312 sequences
- 13503 ATG sites
- 3312 (24.5%) are TIS
- 10191 (75.5%) are non-TIS
- Use for 3-fold x-validation expts



Signal Generation

- K-grams (ie., k consecutive letters)
 - $K = 1, 2, 3, 4, 5, \dots$
 - Window size vs. fixed position
 - Up-stream, downstream vs. any where in window

- In-frame vs. any frame



Copyright 2007 © Limsoon Wong



- Window = ± 100 bases
- In-frame, downstream
 - GCT = 1, TTT = 1, ATG = 1...
- Any-frame, downstream
 - GCT = 3, TTT = 2, ATG = 2...
- In-frame, upstream

- GCT = 2, TTT = 0, ATG = 0, ...

Too Many Signals

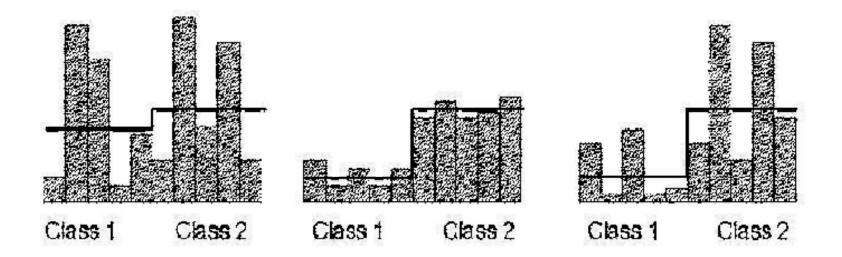


- For each value of k, there are 4^k * 3 * 2 k-grams
- If we use k = 1, 2, 3, 4, 5, we have 24 + 96 + 384 + 1536 + 6144 = 8184 features!
- This is too many for most machine learning algorithms



Signal Selection (Basic Idea)

- Choose a signal w/ low intra-class distance
- Choose a signal w/ high inter-class distance

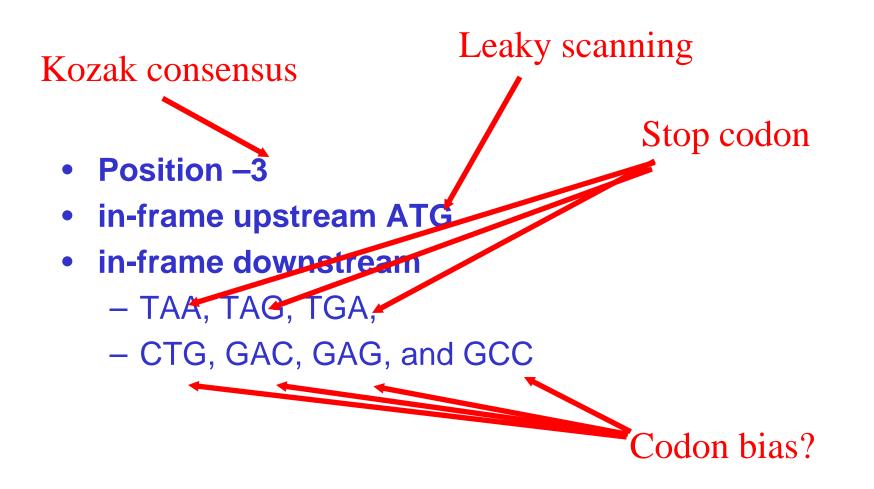




Signal Selection (e.g., CFS)

- Instead of scoring individual signals, how about scoring a group of signals as a whole?
- CFS
 - Correlation-based Feature Selection
 - A good group contains signals that are highly correlated with the class, and yet uncorrelated with each other





Copyright 2007 © Limsoon Wong



Signal Integration

- kNN
 - Given a test sample, find the k training samples that are most similar to it. Let the majority class win
- SVM
 - Given a group of training samples from two classes, determine a separating plane that maximises the margin of error
- Naïve Bayes, ANN, C4.5, ...



Results (3-fold x-validation)

	predicted as positive	predicted as negative
positive		FN
negative	FP	TN

	TP/(TP + FN)	TN/(TN + FP)	TP/(TP + FP)	Accuracy
Naïve Bayes	84.3%	86.1%	66.3%	85.7%
SVM	73.9%	93.2%	77.9%	88.5%
Neural Network	77.6%	93.2%	78.8%	89.4%
Decision Tree	74.0%	94.4%	81.1%	89.4%

ICDE'07, Istanbul, Turkey, 16-20 April 2007

Copyright 2007 © Limsoon Wong



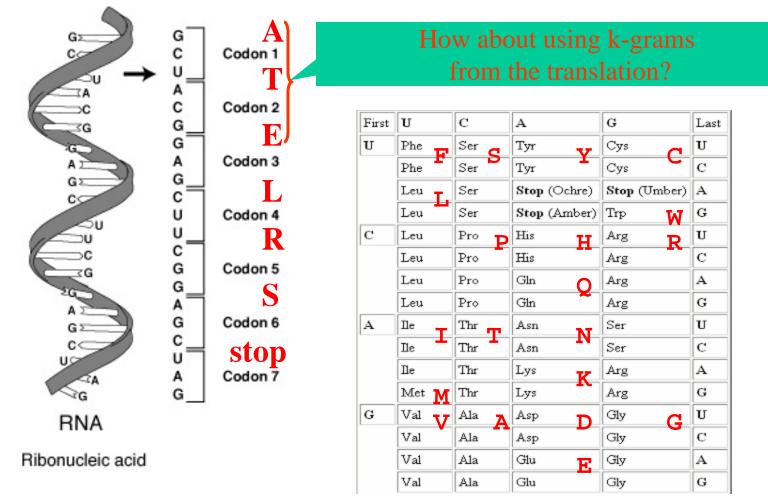
Improvement by Scanning

- Apply Naïve Bayes or SVM left-to-right until first ATG predicted as positive. That's the TIS
- Naïve Bayes & SVM models were trained using TIS vs. Up-stream ATG

	TP/(TP + FN)	TN/(TN + FP)	TP/(TP + FP)	Accuracy
NB	84.3%	86.1%	66.3%	85.7%
SVM	73.9%	93.2%	77.9%	88.5%
NB+Scanning	87.3%	96.1%	87.9%	93.9%
SVM+Scanning	88.5%	96.3%	88.6%	94.4%

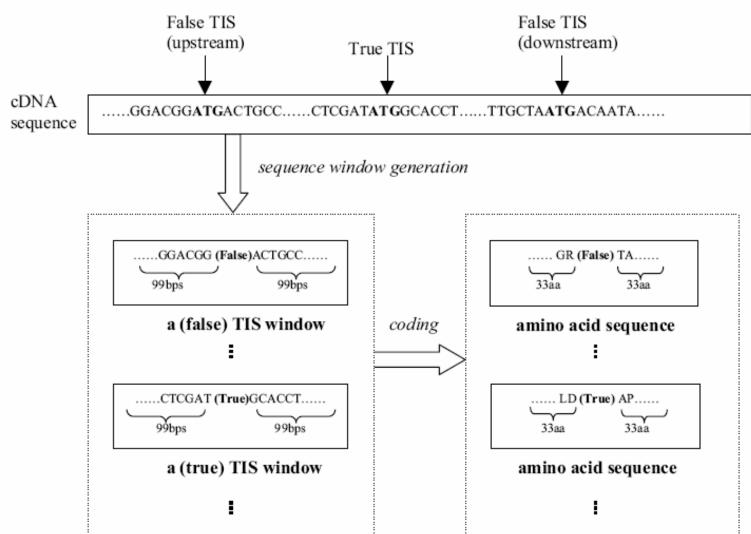


mRNA→Protein





Amino-Acid Features



Amino-Acid Features

	(upstream) T	TIS	(downstream)					
cDNA sequence	GGACGGATGACTGCCCTCGA	TATGGCACCT	TTGCTAATGACAATA					
	sequence window generation							
			GR (False) TA 334a 33aa					
	a (false) TIS window	coding	amino acid sequence					
	:	\Longrightarrow						
	998ps 998ps		LD (True) AP 33aa 33aa					
	a (true) TIS window		amino acid sequence					
	I		ጚ፟፟					

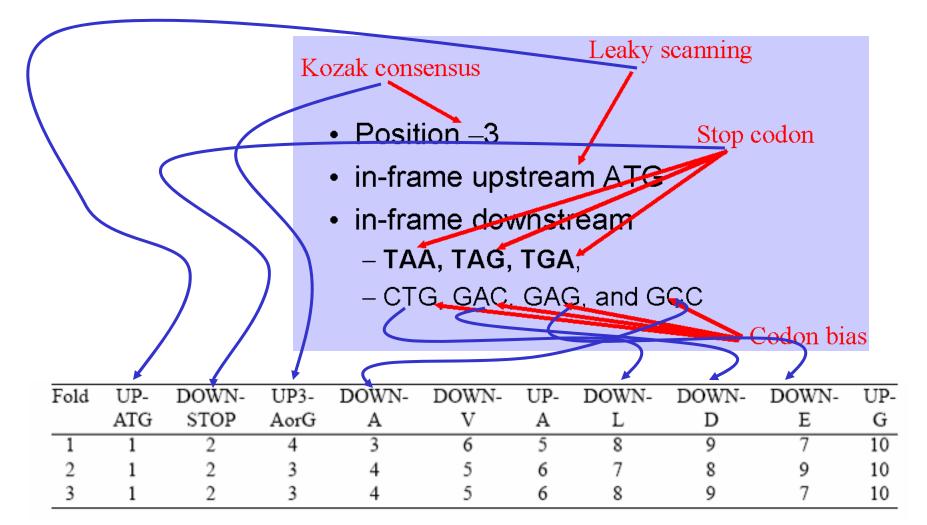
False TIS

False TIS



New feature space (total of 927 features + class label)								
42 1-gram amino acid patterns882 2-gram amino acid patterns3 bio-know- ledge patternsclass label								
UP-A, UP-R,UP-AA, UP-AR,,,UP-N, DOWN-UP-NN, DOWN-AA,A, DOWN-R,,DOWN-AR,,DOWN-NDOWN-NN(numeric type)(numeric type)								
Frequency as values								
1, 3, 5, 0, 4, 6, 2, 7, 0, 5, N, N, N,								
2, 0, 3, 10, 0,	Y, Y, Y,	True						
	882 2-gram amino acid patterns UP-AA, UP-AR,, UP-NN, DOWN-AA, DOWN-AR,, DOWN-NN (numeric type) Frequency as val 6, 2, 7, 0, 5,	882 2-gram amino acid patterns3 bio-know- ledge patternsUP-AA, UP-AR,, UP-NN, DOWN-AA, DOWN-AR,, DOWN-AR,, DOWN-NN (numeric type)DOWN4-G UP3-AorG, UP-ATG (boolean type, Y or N)Frequency as values6, 2, 7, 0, 5,N, N, N,						

Amino Acid K-grams Discovered (by Entropy)



ICDE'07, Istanbul, Turkey, 16-20 April 2007

Copyright 2007 © Limsoon Wong

of Singapore



Independent Validation Sets

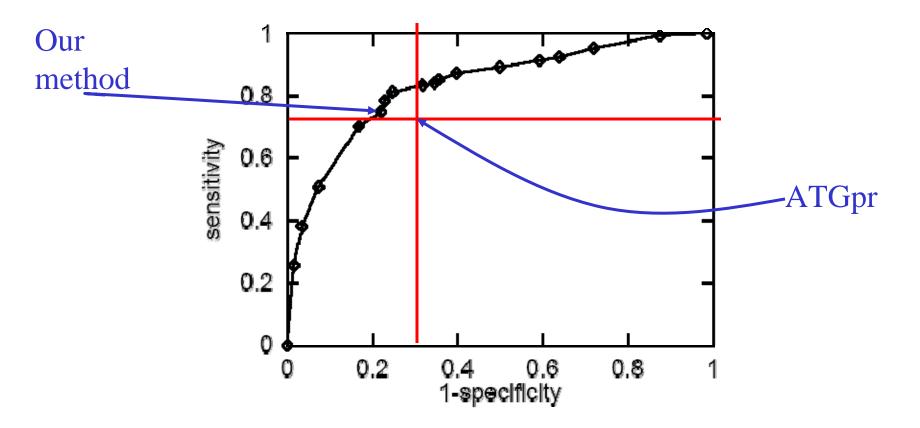
- A. Hatzigeorgiou:
 - 480 fully sequenced human cDNAs
 - 188 left after eliminating sequences similar to training set (Pedersen & Nielsen's)
 - 3.42% of ATGs are TIS
- Our own:
 - well characterized human gene sequences from chromosome X (565 TIS) and chromosome 21 (180 TIS)

Validation Results (on Hatzigeorgious)

Algorithm	Sensitivity	Specificity	Precision	Accuracy
SVMs(linear)	96.28%	89.15%	25.31%	89.42%
SVMs(quad)	94.14%	90.13%	26.70%	90.28%
Ensemble Trees	92.02%	92.71%	32.52%	92.68%
2018-878-82-241 - X	A# A40/	AA #14/	A 1 2001	AA AAA/

 Using top 100 features selected by entropy and trained on Pedersen & Nielsen's dataset

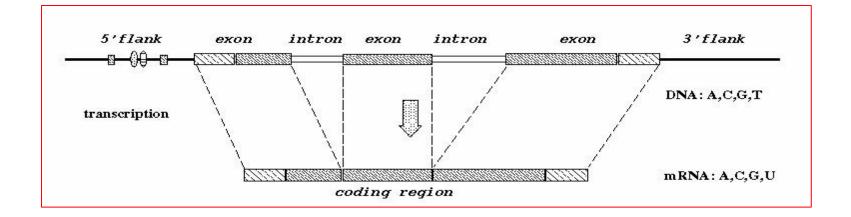
Validation Results (on Chr X and Chr 215)



• Using top 100 features selected by entropy and trained on Pedersen & Nielsen's

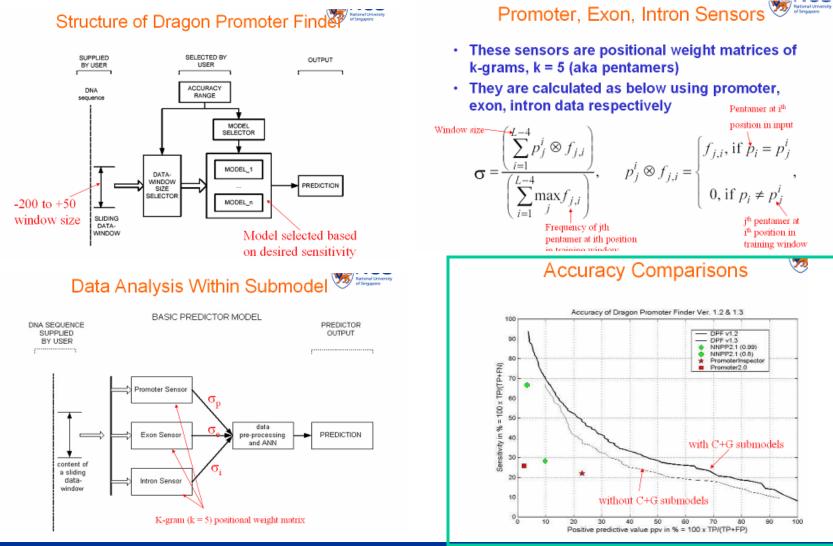
Important Unsolved Challenges

• Other gene features: TSS, PolyA, Splicing, ...



- Gene regulation: TFBS, EREs, miRNA targets, ...
- "Dark matters": miRNA & non-coding genes

TSS: Is State of the Art Good Enough?



ICDE'07, Istanbul, Turkey, 16-20 April 2007

Copyright 2007 © Limsoon Wong



TFs and TFBSs

- Recognize TFBS
 - Short & degenerate
 - Incomplete erroneous data
- Characterize type of regulation, e.g., suppressor vs enhancer
- Identify synergistic TF pairs
 - Genes that are coordinately bound are coordinately expressed
 - Requires data over many time points
 - TFs may interact only under certain conditions

- A method to identify synergistic TF pairs in cell cycle
 - Study if a pair of TFs are assoc with the same genes more often than random
 - For each significant pair, test whether there is a phase in cell cycle where the common associated genes are significantly diff in other cell cycle phases
- ⇒ Quite primitive. Lots more work needed!

miRNA: New Frontier in Molecular Biology

- Impt functions of miRNA:
 - Control of cell fate and fat metabolism in flies
 - Neuronal patterning in nematodes
 - Modulation of hematopoietic lineage differentiation in mammals
 - control of leaf and flower development in plants
- By regulating gene expr in binding targets:
 - Repression/inhibition of translation
 - Degradation of mRNA

- miRNA identification thru direct cloning has limitations:
 - miRNAs may be expressed only in certain tissues or at certain times
 - Expr levels vary greatly
 - Degradation products from mRNAs and other endogenous non-coding RNAs coexist w/ miRNAs and are sometimes dominant in small RNA molecule samples extracted from cells



Recommended Readings

- Wong, *The practical Bioinformatician*, 2004, ICP. Chapters 4, 5, 6, & 7
- Liu & Wong, Data mining tools for biological sequences. *JBCB*, 1:139-168, 2003
- Tsai et al., Statistical methods for identifying yeast cell cycle transcription factors. *PNAS*, 102:13532-13537, 2005
- Yang et al., Identification of microRNA precursors via SVM. APBC 2006, pages 267-276
- Zheng et al., Exploring Essential Attributes for Detecting MicroRNA Precursors from Background Sequences. *VLDB 2006 Workshop on Data Mining in Bioinformatics*, pages 131--145.

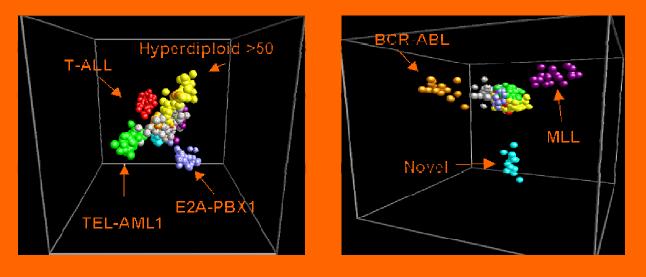
5-Minute Break?



ICDE'07, Istanbul, Turkey, 16-20 April 2007

70 min

Disease Diagnosis, Treatment, & Understanding

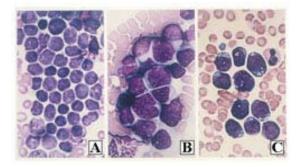




Childhood Acute Lymphoblastic Leukernia

- Major subtypes: T-ALL, E2A-PBX, TEL-AML, BCR-ABL, MLL genome rearrangements, Hyperdiploid>50
- Diff subtypes respond differently to same Tx
- Over-intensive Tx
 - Development of secondary cancers
 - Reduction of IQ
- Under-intensiveTx
 - Relapse

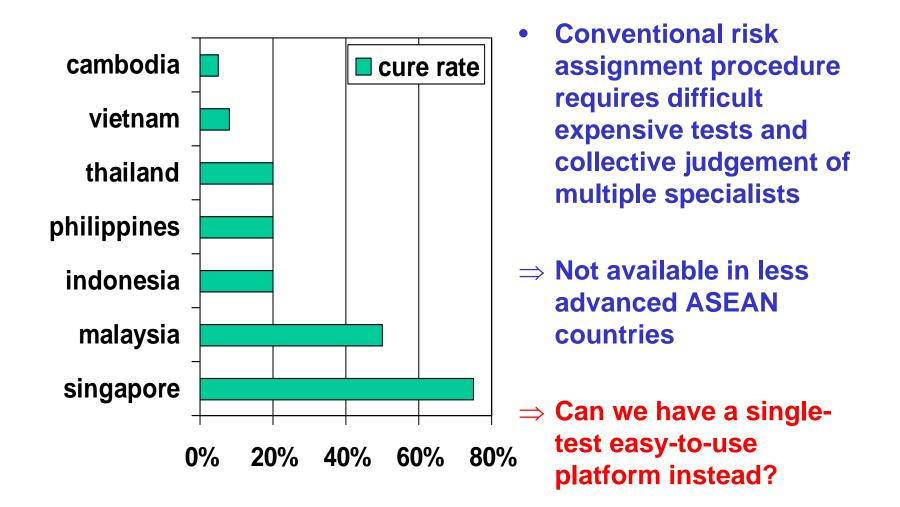
The subtypes look similar



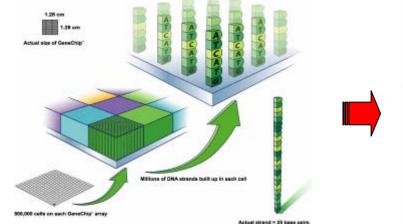
- Conventional diagnosis
 - Immunophenotyping
 - Cytogenetics
 - Molecular diagnostics



Childhood ALL Cure Rates

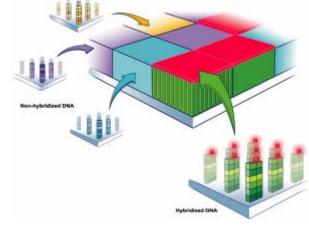






(II) Intra-class distance is too large

Class A Class B

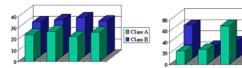


Shining a laser light at GeneChip causes tagged DNA fragments that hybridized to gio

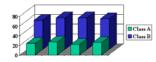


	00-0586-09	00-0586-U	0D-0586-U	00-0586-U	0D-0586-U	Descriptions
	Positive	Negative	Pairs InAw	Avg Diff	Abs Call	
AFFX-Murl	5	2	19	297.5	A	M16762 Mouse int
AFFX-Murl	3	2	19	654.2	A	M37897 Mouse int
AFFX-Murl	4	2	19	308.6	A	M25892 Mus musi
AFFX-Murl	1	3	19	141	A	MB3649 Mus musi
AFFX-BioE	13	1	19	9340.6	P	JD4423 E coli bioB
AFFX-BioE	15	0	19	12862.4	P	JD4423 E coli bioB
AFFX-BioE	12	0	19	8716.5	P	JD4423 E coli bioB
AFFX-BioC	17	0	19	25942.5	P	JD4423 E coli bioC
AFFX-BioC	16	0	20	28838.5	P	JD4423 E coli bioC
AFFX-BioD	17	0	19	25785.2	P	JD4423 E coli bioD
AFFX-BioD	19	0	20	140113.2	P	JD4423 E coli bioD
AFFX-Cre>	20	0	20	280036.6	P	X03453 Bacterioph
AFFX-CreX	20	0	20	401741.8	P	X03453 Bacterioph
AFFX-BioE	7	5	18	-483	A	JD4423 E coli bicB
AFFX-BioE	5	4	18	313.7	A	JD4423 E coli bicB
AFFX-BioE	7	6	20	-1016.2	A	JD4423 E coli bicB

(I) Inter-class distance is too small



(III) Inter- and intra-class distances of a good signal



ICDE'07, Istanbul, Turkey, 16-20 April 2007

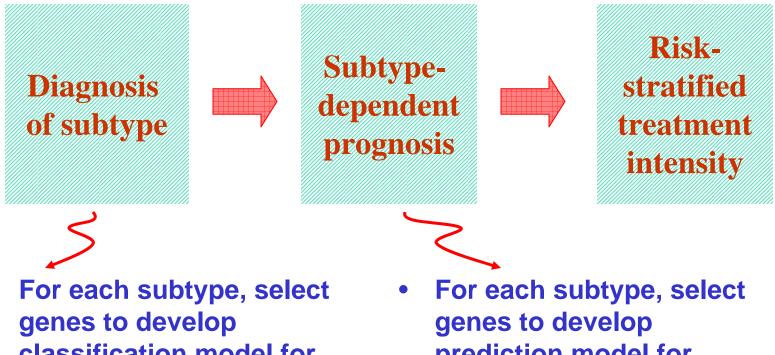


A Sample Affymetrix GeneChip Data File (U95A)

	00-0586-U	00-0586-U	(00-0586-U	00-0586-U	00-0586-U	Descriptions
	Positive	Negative	Pairs InAv	Avg Diff	Abs Call	
AFFX-Murl	5	2	19	297.5	A	M16762 Mouse interleukin 2 (IL-2) gene, exon 4
AFFX-Murl	3	2	19	554.2	A	M37897 Mouse interleukin 10 mRNA, complete cds
AFFX-Murl	4	2	19	308.6	A	M25892 Mus musculus interleukin 4 (II-4) mRNA, comp
AFFX-Murf	1	3	19	141	A	M83649 Mus musculus Fas antigen mRNA, complete i
AFFX-BioE	13	1	19	9340.6	Р	J04423 E coli bioB gene biotin synthetase (-5, -M, -3 r
AFFX-BioE	15	0	19	12862.4	Р	J04423 E coli bioB gene biotin synthetase (-5, -M, -3 r
AFFX-BioE	12	0	19	8716.5	Р	J04423 E coli bioB gene biotin synthetase (-5, -M, -3 r
AFFX-BioC	17	0	19	25942.5	Р	J04423 E coli bioC protein (-5 and -3 represent transcr
AFFX-BioC	16	0	20	28838.5	Р	J04423 E coli bioC protein (-5 and -3 represent transcr
AFFX-BioD	17	0	19	25765.2	Р	J04423 E coli bioD gene dethiobiotin synthetase (-5 ar
AFFX-BioD	19	0	20	140113.2	Р	J04423 E coli bioD gene dethiobiotin synthetase (-5 ar
AFFX-CreX	20	0	20	280036.6	Р	X03453 Bacteriophage P1 cre recombinase protein (-5
AFFX-CreX	20	0	20	401741.8	Р	X03453 Bacteriophage P1 cre recombinase protein (-5
AFFX-BioE	7	5	18	-483	A	J04423 E coli bioB gene biotin synthetase (-5, -M, -3 r
AFFX-BioE	5	4	18	313.7	A	J04423 E coli bioB gene biotin synthetase (-5, -M, -3 r
AFFX-BioE	7	6	20	-1016.2	A	J04423 E coli bioB gene biotin synthetase (-5, -M, -3 r



Overall Strategy



- classification model for diagnosing that subtype
- prediction model for prognosis of that subtype

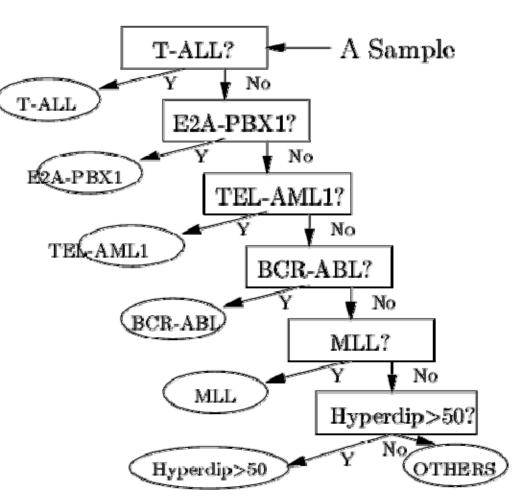


Subtype Diagnosis by PCL

- Gene expression data collection
- Gene selection by $\chi 2$
- Classifier training by emerging pattern
- Classifier tuning (optional for some machine learning methods)
- Apply classifier for diagnosis of future cases by PCL

Childhood ALL Subtype Diagnosis Workflow







National University of Singapore



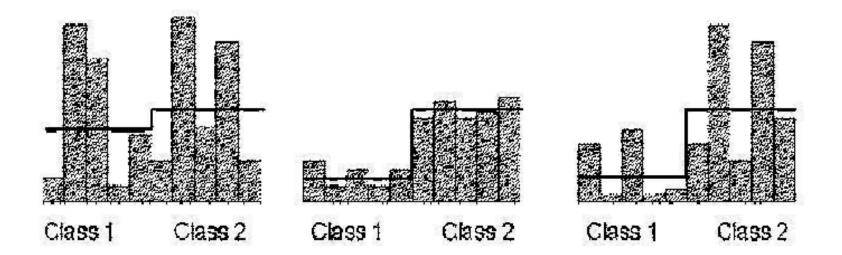
Training and Testing Sets

Paired datasets	Ingredients	Training	Testing
T-ALL vs	$OTHERS1 = \{E2A-PBX1, TEL-AML1, $	$28 \ \mathrm{vs} \ 187$	15 vs 97
OTHERS1	BCR-ABL, Hyperdip>50, MLL, OTHERS}		
E2A-PBX1 vs	$OTHERS2 = \{TEL-AML1, BCR-ABL$	18 vs 169	9 vs 88
OTHERS2	Hyperdip>50, MLL, OTHERS}		
TEL-AML1 vs	$OTHERS3 = \{BCR-ABL$	52 vs 117	27 vs 61
OTHERS3	Hyperdip>50, MLL, OTHERS}		
BCR-ABL vs	$OTHERS4 = \{Hyperdip > 50,$	9 vs 108	6 vs 55
OTHERS4	MLL, OTHERS}		
MLL vs	$OTHERS5 = {Hyperdip>50, OTHERS}$	14 vs 94	6 vs 49
OTHERS5			
Hyperdip>50 vs	$OTHERS = \{Hyperdip47-50, Pseudodip, \}$	42 vs 52	22 vs 27
OTHERS	Hypodip, Normo}		



Signal Selection Basic Idea

- Choose a signal w/ low intra-class distance
- Choose a signal w/ high inter-class distance





Signal Selection by $\chi 2$

The \mathcal{X}^2 value of a signal is defined as:

$$\mathcal{X}^2 = \sum_{i=1}^{m} \sum_{j=1}^{k} \frac{(A_{ij} - E_{ij})^2}{E_{ij}},$$

where *m* is the number of intervals, *k* the number of classes, A_{ij} the number of samples in the *i*th interval, *j*th class, R_i the number of samples in the *i*th interval, C_j the number of samples in the *i*th interval, C_j the number of samples in the *j*th class, *N* the total number of samples, and E_{ij} the expected frequency of A_{ij} ($E_{ij} = R_i * C_j/N$).



Emerging Patterns

- An emerging pattern is a set of conditions
 - usually involving several features
 - that most members of a class satisfy
 - but none or few of the other class satisfy
- A jumping emerging pattern is an emerging pattern that
 - some members of a class satisfy
 - but no members of the other class satisfy
- We use only jumping emerging patterns



Examples

Patterns	Frequency (P)	Frequency(N)			
{9, 36}	38 instances	0			
{9, 23}	38	0			
$\{4, 9\}$	38	0			
{9, 14}	38	0 Easy interpretation			
<i>{</i> 6 <i>,</i> 9 <i>}</i>	38	0			
{7, 21}	0	36			
{7, 11}	0	35			
{7, 43}	0	35			
{7, 39}	0	34			
{24, 29}	0	34			

Reference number 9: the expression of gene 37720_at > 215 Reference number 36: the expression of gene 38028_at <= 12

75



PCL: Prediction by Collective Likelihood

- Let EP_1^P, \ldots, EP_i^P be the most general EPs of D^P in descending order of support.
- Suppose the test sample T contains these most general EPs of D^P (in descending order of support):

$$EP_{i_1}^P, EP_{i_2}^P, \cdots, EP_{i_x}^P$$

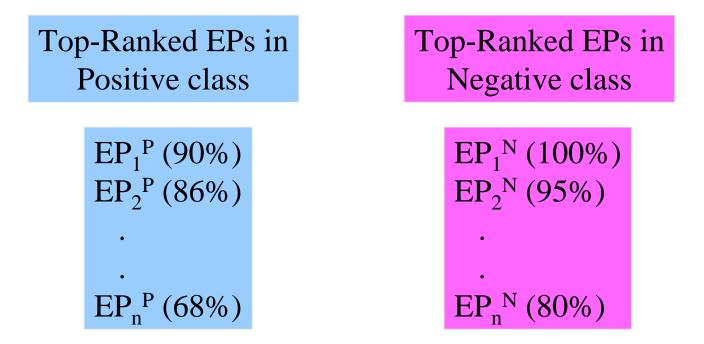
• Use k top-ranked most general EPs of D^P and D^N . Define the score of T in the D^P class as

$$score(T, D^P) = \sum_{m=1}^{k} \frac{frequency(EP^P_{i_m})}{frequency(EP^P_m)}$$

- Ditto for $score(T, D^N)$.
- If $score(T, D^P) > score(T, D^N)$, then T is class P. Otherwise it is class N.

PCL Learning





The idea of summarizing multiple top-ranked EPs is intended to avoid some rare tie cases

Copyright 2007 © Limsoon Wong



PCL Testing

Most freq EP of pos class in the test sample

Score^P =
$$\stackrel{\checkmark}{\text{EP}_1^{P'}} / \stackrel{\text{EP}_1^{P}}{\uparrow} + \dots + \stackrel{\text{EP}_k^{P'}}{\uparrow} / \stackrel{\text{EP}_k^{P}}{\uparrow}$$

Most freq EP of pos class

Similarly, Score^N = $EP_1^{N'} / EP_1^{N} + ... + EP_k^{N'} / EP_k^{N}$

If Score^P > Score^N, then positive class, Otherwise negative class

Accuracy of PCL (vs. other classifiers)

Testing Data	Error rate of different models					
	C4.5	SVM	NB	\mathbf{PCL}		
T-ALL vs OTHERS1	0:1	0:0	0:0	0:0		
E2A-PBX1 vs OTHERS2	0:0	0:0	0:0	0:0		
TEL-AML1 vs OTHERS3	1:1	0:1	0:1	1:0		
BCR-ABL vs OTHERS4	2:0	3:0	1:4	2:0		
MLL vs OTHERS5	0:1	0:0	0:0	0:0		
Hyperdiploid>50 vs OTHERS	2:6	0:2	0:2	0:1		
Total Errors	14	6	8	4		

The classifiers are all applied to the 20 genes selected by $\chi 2$ at each level of the tree



Understandability of PCL

• E.g., for T-ALL vs. OTHERS, one ideally discriminatory gene 38319_at was found, inducing these 2 EPs

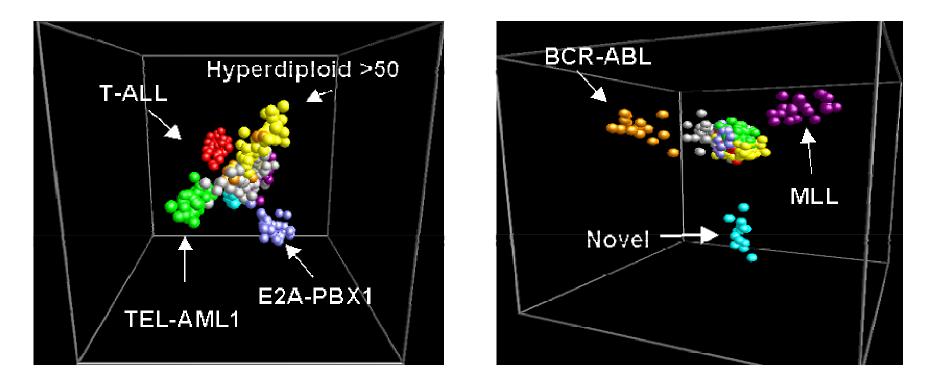
{ $gene_{(38319_at)}@(-\infty, 15975.6)$ } and { $gene_{(38319_at)}@[15975.6, +\infty)$ }.

• These give us the diagnostic rule

If the expression of 38 319_*at* is less than 15 975.6, then this ALL sample must be a T-ALL. Otherwise it must be a subtype in OTHERS1.



Multidimensional Scaling Plot for Subtype Diagnosis



Obtained by performing PCA on the 20 genes chosen for each level

Impact



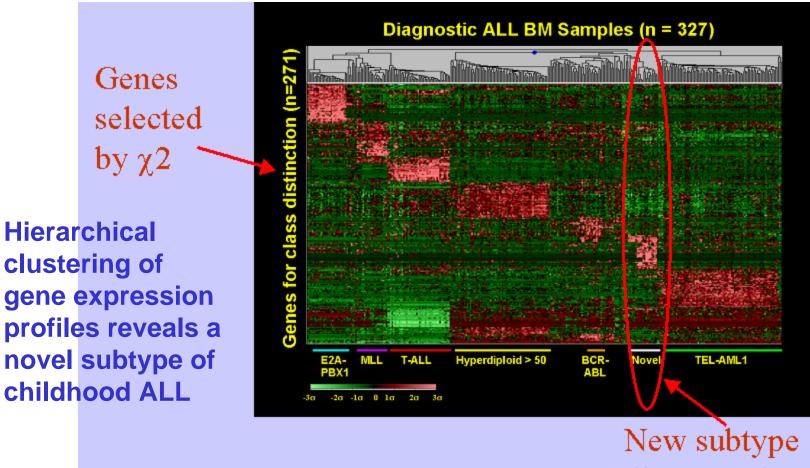
Childhood ALL in ASEAN Countries **Conventional Tx:** (2000 new cases per year) • intermediate intensity to 🗆 cure rate cambodia everyone vietnam Childhood ALL \Rightarrow 10% suffers relapse Patients Profile thailand \Rightarrow 50% suffers side effects 🗆 High philippines 10% Low indonesia \Rightarrow costs US\$150m/yr 50% Inter malaysia 40% singapore 60% 0% 20% 40% 80% ►High cure rate of 80% • Less relapse **Our optimized Tx:** • Less side effects • high intensity to 10% • Save US\$51.6m/yr intermediate intensity to 40% low intensity to 50% costs US\$100m/yr

ICDE'07, Istanbul, Turkey, 16-20 April 2007

Copyright 2007 © Limsoon Wong



Is there a new subtype?



discovered

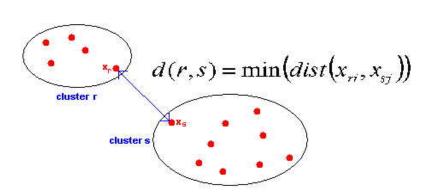
Copyright 2007 © Limsoon Wong

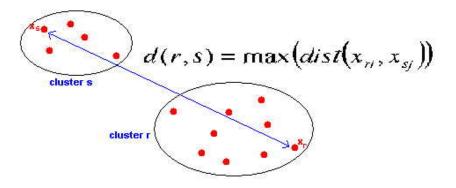


Hierarchical Clustering

- Assign each item to its own cluster
 - If there are N items initially, we get N clusters, each containing just one item
- Find the "most similar" pair of clusters, merge them into a single cluster, so we now have one less cluster
 - "Similarity" is often defined using
 - Single linkage
 - Complete linkage
 - Average linkage
- Repeat previous step until all items are clustered into a single cluster of size N







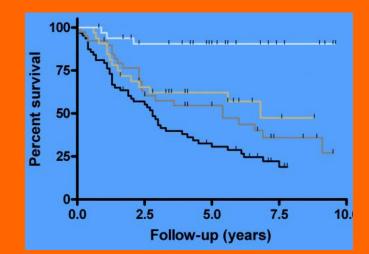
Single linkage defines distance betw two clusters as min distance betw them

Complete linkage defines distance betw two clusters as max distance betw them

Exercise: Give definition of "average linkage"

Image source: UCL Microcore Website

Selection of Patient Samples and Genes for Disease Prognosis





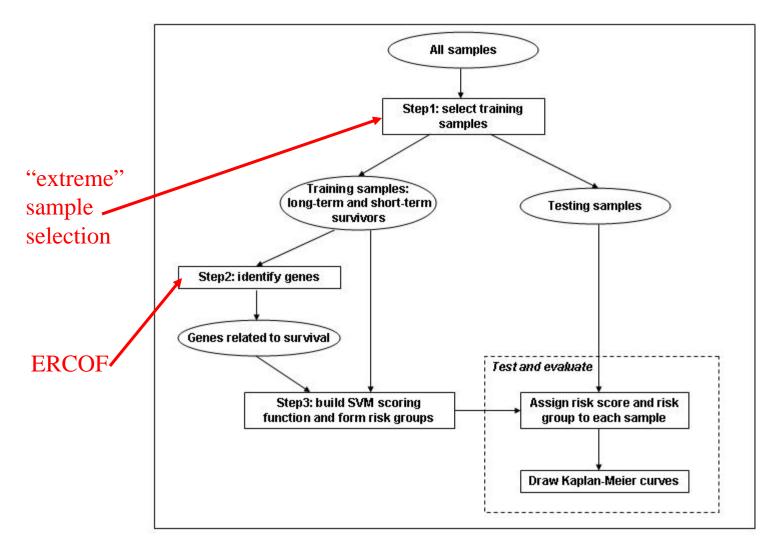


Gene Expression Profile + Clinical Data ⇒ Outcome Prediction

- Univariate & multivariate Cox survival analysis (Beer et al 2002, Rosenwald et al 2002)
- Fuzzy neural network (Ando et al 2002)
- Partial least squares regression (Park et al 2002)
- Weighted voting algorithm (Shipp et al 2002)
- Gene index and "reference gene" (LeBlanc et al 2003)
-



Another Approach ...





Extreme Sample Selection

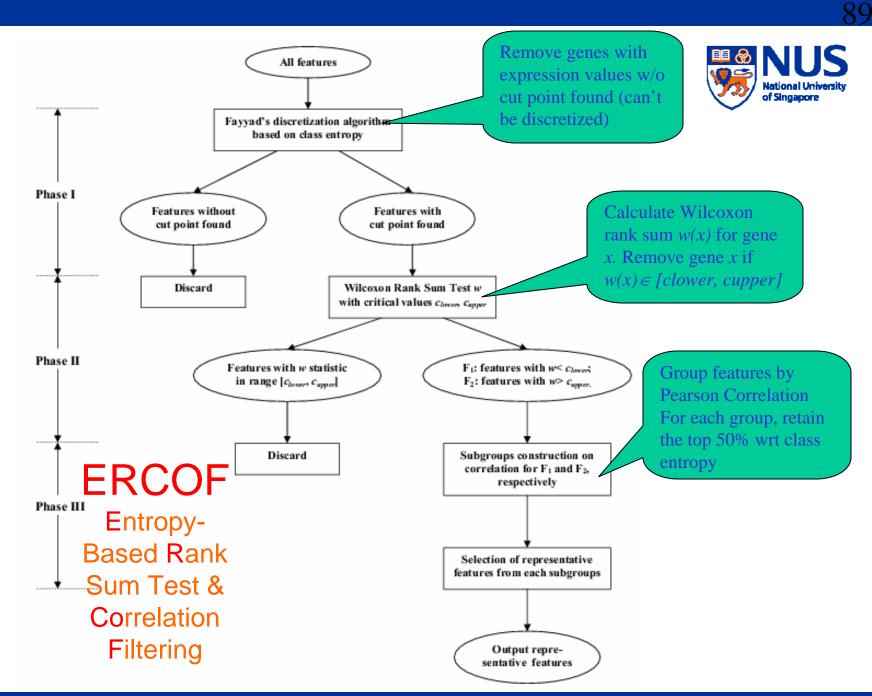
Short-term Survivors v.s. Long-term Survivors

Short-term survivors who died within a short period ↓ Long-term survivors who were alive after a *long* follow-up time

 $F(T) < c_1$ and E(T) = 1

 $F(T) > C_2$

T: sample F(T): follow-up time E(T): status (1:unfavorable; 0: favorable) c_1 and c_2 : thresholds of survival time



ICDE'07, Istanbul, Turkey, 16-20 April 2007

Copyright 2007 © Limsoon Wong



Risk Score Construction

Linear Kernel SVM regression function $G(T) = \sum_{i} a_{i} y_{i} K(T, x(i)) + b$

T: test sample, x(i): support vector, y_i : class label (1: short-term survivors; -1: long-term survivors)

Transformation function (posterior probability) $S(T) = \frac{1}{1 + e^{-G(T)}} \quad (S(T) \in (0,1))$

S(T): *risk score* of sample T



Diffuse Large B-Cell Lymphoma

- DLBC lymphoma is the most common type of lymphoma in adults
- Can be cured by anthracycline-based chemotherapy in 35 to 40 percent of patients
- ⇒ DLBC lymphoma comprises several diseases that differ in responsiveness to chemotherapy

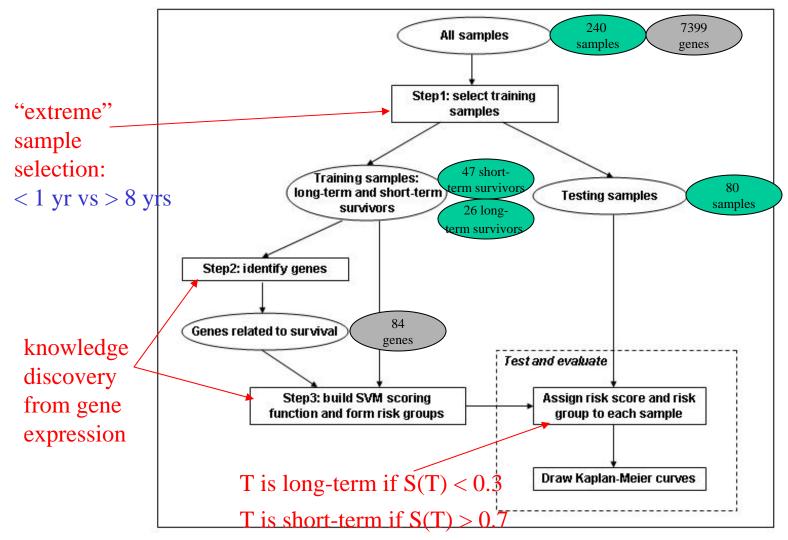
- Intl Prognostic Index (IPI)
 - age, "Eastern Cooperative Oncology Group" Performance status, tumor stage, lactate dehydrogenase level, sites of extranodal disease, ...
- Not very good for stratifying DLBC lymphoma patients for therapeutic trials
- ⇒ Use gene-expression profiles to predict outcome of chemotherapy?



Rosenwald et al., *NEJM* 2002

- 240 data samples
 - 160 in preliminary group
 - 80 in validation group
 - each sample described by 7399 microarray features
- Rosenwald et al.'s approach
 - identify gene: Cox proportional-hazards model
 - cluster identified genes into four gene signatures
 - calculate for each sample an outcome-predictor score
 - divide patients into quartiles according to score

Knowledge Discovery from Gene of Singapore **Expression of "Extreme" Samples**



ICDE'07, Istanbul, Turkey, 16-20 April 2007

Copyright 2007 © Limsoon Wong

ational University



Discussions: Sample Selection

Application	Data set	Sta	Total	
		Dead	Alive	
DLBCL	Original	88	72	160
	Informative	47+1(*)	25	73

Number of samples in original data and selected informative training set. (*): Number of samples whose corresponding patient was dead at the end of follow-up time, but selected as a long-term survivor.

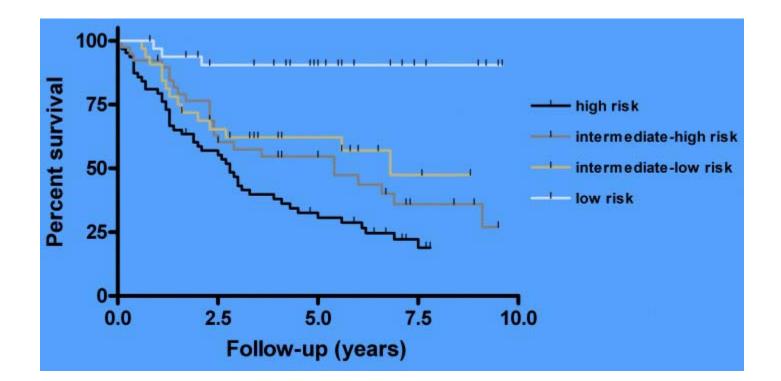


Discussions: Gene Identification

Gene selection	DLBCL
Original	4937(*)
Phase I	132(2.7%)
Phase II	84(1.7%)

Number of genes left after feature filtering for each phase. (*): number of genes after removing those genes who were absent in more than 10% of the experiments.

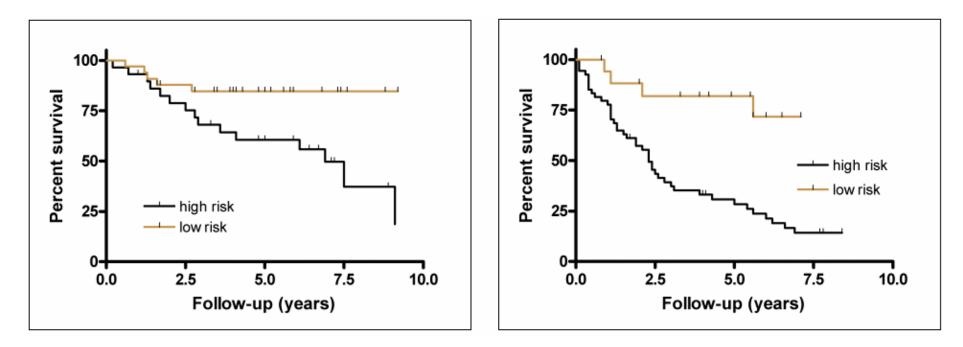




p-value of log-rank test: < 0.0001 Risk score thresholds: 0.7, 0.3 96



Improvement Over IPI

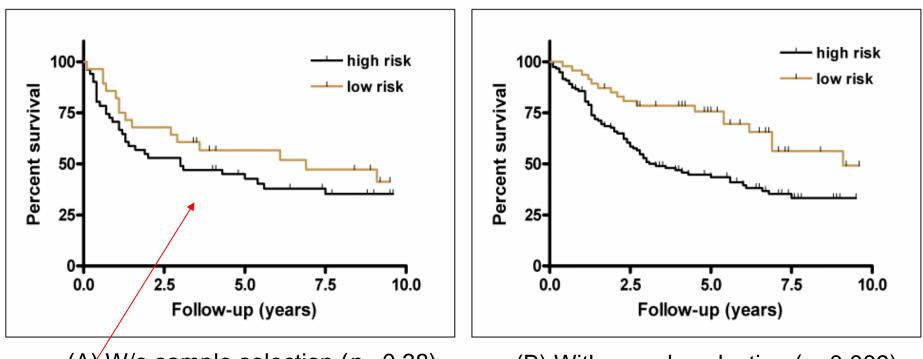


(A) IPI low, p-value = 0.0063 (B) IPI intermediate, p-value = 0.0003



98

Merit of "Extreme" Samples



(A) W/o sample selection (p = 0.38)

(B) With sample selection (p=0.009)

No clear difference on the overall survival of the 80 samples in the validation group of DLBCL study, if **no training sample selection conducted**

Is ERCOF Useful? Observations from 1000+ Expts



- Feature selection methods considered
 - All use all features
 - All-entropy select features whose value range can be partitioned by Fayyad & Irani's entropy method
 - Mean-entropy select features whose entropy is better than the mean entropy
 - Top-number-entropy select the top 20, 50, 100, 200 genes by their entropy
 - ERCOF at 5% significant level for Wilcoxon rank sum test and 0.99 Pearson correlation coeff threshold

- Data sets considered
 - Colon tumor
 - Prostate cancer
 - Lung cancer
 - Ovarian cancer
 - DLBC lymphoma
 - ALL-AML
 - Childhood ALL
- Learning methods considered
 - C4.5
 - Bagging, Boosting, CS4
 - SVM, 3-NN



ERCOF vs All-Entropy

Experiment	SVM	3-NN	Bagging	AdaBoostM1	RandomForests	CS4	
ColonTumor	С	A,C	С	С	С	С	
Prostate	С	С	A,C	A,C	С	С	
Lung test	С	A,C	Α	A,C	С	A,C	
Lung	A,C	С	A,C	С	С	С	
Ovarian	A,C	С	A,C	С	С	A,C	
DLBCL	С	С	Α	С	А	A,C	
ALLAML test	A,C	С	A,C	A,C	С	С	
ALLAML	A,C	A,C	A,C	С	A,C	A,C	
		Peo	liatric ALL	data — test			
T-ALL	A,C	A,C	A,C	A,C	A,C	A,C	
E2A-PBX1	A,C	A,C	A,C	A,C	A,C	A,C	All-entropy
TEL-AML1	A,C	A,C	A,C	A,C	A,C	С	
BCR-ABL	A,C	С	A,C	A,C	С	A,C	wins 4 times
MLL	A,C	A,C	С	A,C	С	С	1
Hyperdip>50	A,C	Α	A,C	С	С	С	
	Pedi)-fold cross valid	ation		
T-ALL	A,C	С	A,C	A,C	A,C	С	
E2A-PBX1	С	С	A,C	A,C	С	С	/ ERCOF
TEL-AML1	С	С	С	С	С	С	wins 60 times
BCR-ABL	С	С	С	С	С	A,C	whils oo times
MLL	A,C	С	С	С	С	C	
Hyperdip>50	С	С	A,C	С	С	A,C	
Sum <	A:0	A:1	A:2	A:0	A:1	A:0	
	<c:8< td=""><td>C:12</td><td>C:5</td><td>C:10</td><td>C:14</td><td>C:11</td><td></td></c:8<>	C:12	C:5	C:10	C:14	C:11	
	Tie:12	Tie:7	Tie:13	Tie:10	Tie:5	Tie:9	



ERCOF vs Mean-Entropy

Experiment	SVM	3-NN	Bagging	AdaBoostM1	RandomForests	CS4	
ColonTumor	С	С	B,C	С	С	С	
Prostate	С	B,C	С	В	С	B,C	
Lung test	B,C	B,C	В	B,C	С	B,C	
Lung	B,C	С	B,C	В	В	B,C	
Ovarian	B,C	С	В	С	B,C	С	
DLBCL	B,C	С	В	B,C	С	B,C	
ALLAML test	B,C	С	B,C	B,C	С	B,C	
ALLAML	B,C	В	В	С	В	B,C	
		Pe	diatric ALL	data — test			
T-ALL	B,C	B,C	B,C	B,C	B,C	B,C	
E2A-PBX1	B,C	B,C	B,C	B,C	B,C	B,C	
TEL-AML1	B,C	B,C	B,C	В	С	С	
BCR-ABL	С	B,C	В	B,C	B,C	B,C	
MLL	B,C	B,C	B,C	B,C	B,C	B,C	
Hyperdip>50	B,C	В	В	B,C	С	B,C	
Pediatric ALL data — 10-fold cross validation							
T-ALL	B,C	В	B,C	B,C	B,C	B,C	
E2A-PBX1	С	С	B,C	B,C	С	С	
TEL-AML1	С	С	B,C	С	С	C /	
BCR-ABL	С	С	С	В	В	в	
MLL	B,C	B,C	С	С	В	с /	
Hyperdip>50	С	С	С	B,C	С	с	
Sum	B :0	B:3	B:6	B:4	B:4	B:1	
	< <u>C:7</u>	C:9	C:4	C:5	C:10	C:7	
	Tie:13	Tie:8	Tie:10	Tie:11	Tie:6	Tie:12	

Mean-entropy wins 18 times

ERCOF wins 42 times



Effectiveness of ERCOF

Table 5.32: A summary of the total winning times (including tie cases) of each classifier (under different feature selection methods) across the 20 validation tests on the six gene expression profiles and one proteomic data set. The number with bold font in each row indicates the feature selection method that owns most winning times for the relevant classifier. In the brackets, there is the total number of misclassified samples across the same 20 validation tests. Similarly, the figure with bold font in the brackets in each row is the minimum number of total misclassified samples among feature selection methods for the classifier.

Classifier	All	All-entropy	Mean-entropy	Top-number-entropy			ERCOF	
				20	50	100	200	
SVM	4(100)	9(52)	11(48)	6(76)	6(74)	11(52)	11(59)	16(38)
3-NN	1(187)	5(87)	8(77)	6(88)	4(81)	6(77)	5(73)	12(61)
Bagging	7(123)	5(117)	8(115)	11(123)	11(122)	7(122)	9(114)	8(112)
AdaBoostM1	5(191)	8(181)	8(166)	11(138)	10(144)	10(157)	9(162)	10(154)
RandomForests	0(228)	5(111)	5(93)	6(96)	7(83)	8(96)	5(90)	9(80)
CS4	5(87)	6(77)	6(76)	7(101)	10(81)	9(74)	8(74)	12(66)
Total wins	22	38	46	47	48	51	47	67



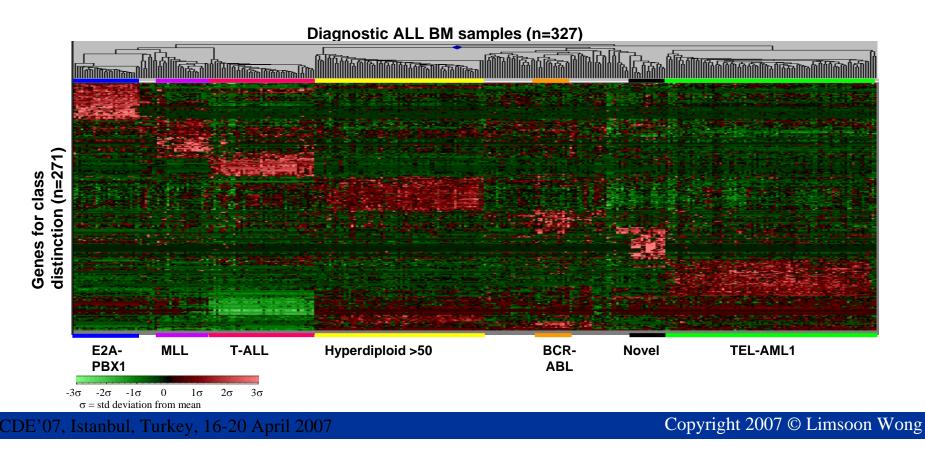


- Selecting extreme cases as training samples is an effective way to improve patient outcome prediction based on gene expression profiles and clinical information
- ERCOF is very suitable for SVM, 3-NN, CS4, Random Forest, as it gives these learning algos highest no. of wins
- ERCOF is suitable for Bagging also, as it gives this classifier the lowest no. of errors
- ⇒ ERCOF is a systematic feature selection method that is very useful

Beyond Classification of Gene Expression Profiles



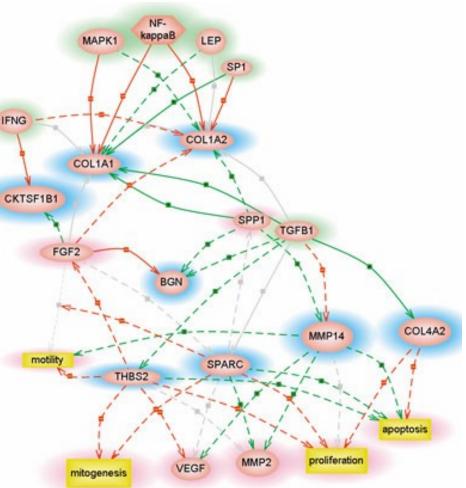
 After identifying the candidate genes by feature selection, do we know which ones are causal genes and which ones are surrogates?





Gene Regulatory Circuits

- Genes are "connected" in "circuit" or network
- Expr of a gene in a network depends on expr of some other genes in the network
- Can we "reconstruct" the gene network from gene expression and other data?



Source: Miltenyi Biotec



Recommended Readings

- Wong, *The Practical Bioinformatician*, 2004, ICP. Chapter 14
- Li & Wong, Identifying good diagnostic genes or gene groups from gene expression data by using the concept of emerging patterns. *Bioinformatics*, 18:725-734, 2002
- Miller et al., Optimal gene expression analysis by microarrays. *Cancer Cell*, 2:353-361, 2002
- Liu et al,. Use of Extreme Patient Samples for Outcome Prediction from Gene Expression Data. *Bioinformatics*, 21:3377-3384, 2005
- Eisen et al., Cluster analysis and display of genome-wide expression patterns. *PNAS*, 8:14863-14868, 1998
- Tibshirani et al., Diagnosis of multiple cancer types by shrunken centroids of gene expression. *PNAS*, 99:6567-6572, 2002

5 min



Advancing Knowledge Discovery



ICDE'07, Istanbul, Turkey, 16-20 April 2007



Some of those "techniques" frequently needed in analysis of biomedical data are insufficiently studied by current data mining researchers

- Recognizing what samples are relevant and what are not
- Recognizing what features are relevant and what are not & handling missing or incorrect values
- Recognizing trends, changes, and their causes

Any Question?



ICDE'07, Istanbul, Turkey, 16-20 April 2007



Acknowledgements

- Disease Treatment Planning & Understanding
 - Jinyan Li, Huiqing Liu
 - Allen Yeoh
- Gene Feature Recognition
 - Huiqing Liu, Rajesh Chowdhary, Vladimir Bajic
 - Fanfan Zeng, Roland Yap
- Protein Function Inference
 - Hon Nian Chua, Wing-Kin Sung