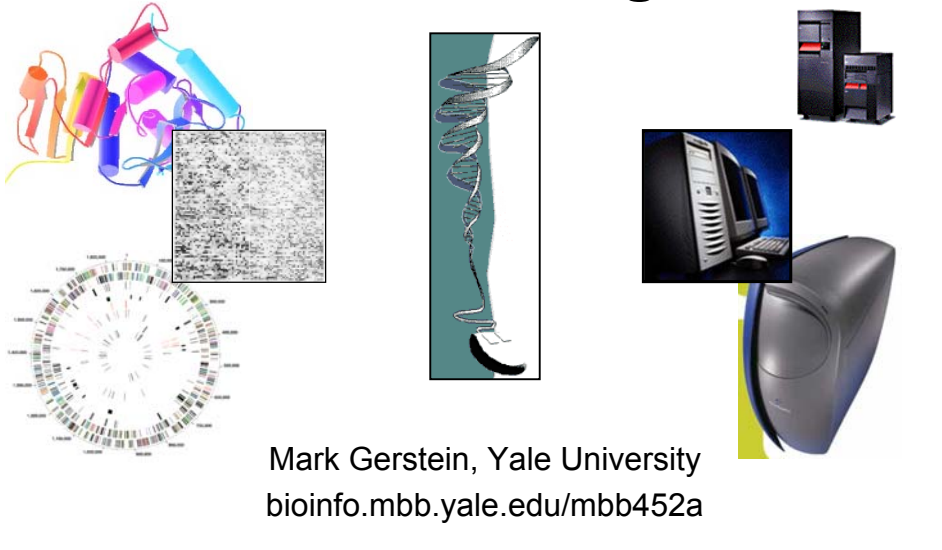


BIOINFORMATICS

Datamining



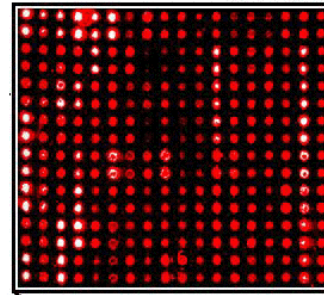
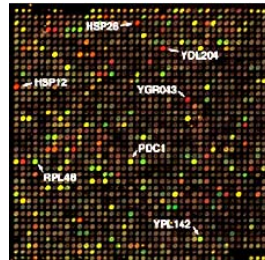
Mark Gerstein, Yale University
bioinfo.mbb.yale.edu/mbb452a

Large-scale Datamining

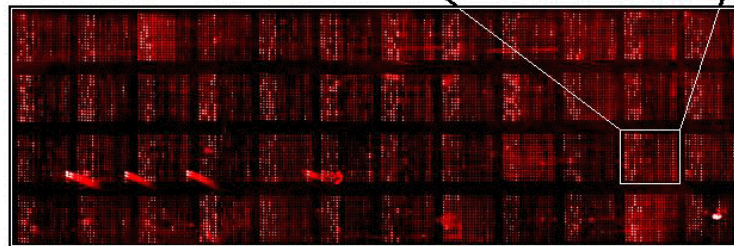
- Gene Expression
 - ◇ Representing Data in a Grid
 - ◇ Description of function prediction in abstract context
- Unsupervised Learning
 - ◇ clustering & k-means
 - ◇ Local clustering
- Supervised Learning
 - ◇ Discriminants & Decision Tree
 - ◇ Bayesian Nets
- Function Prediction EX
 - ◇ Simple Bayesian Approach for Localization Prediction

The recent advent and subsequent onslaught of microarray data

1st generation, Expression Arrays (Brown)



2nd gen., Proteome Chips (Snyder)



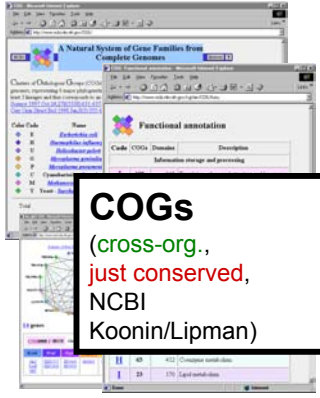
3 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Gene Expression Information and Protein Features

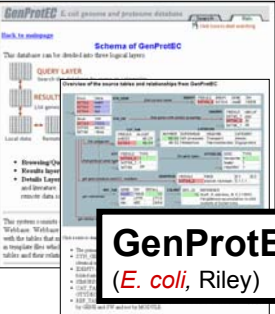
Yeast Gene ID	Basics																Predictors																	
	Sequence Features										Genomic Features						Cell cycle timecourse																	
	seq. length	Amino Acid Composition					How many times does the sequence have these motif features?					Abs. expr. Level (mRNA copies / cell)	Prot. Abundance	Cell cycle timecourse																				
Sequence	A	C	D	S	W	Y	Start site	NLS	Index motif	muc2	signalp	tmms1	Gene-Chip exp. from RY Lab	sage tag freq.	(1000 copies /cell)	†#0	†#1	†#2	†#3	†#4	†#5	†#6	†#7	†#8	†#9	†#10	†#11	†#12	†#13	†#14	†#15	†#16		
YAL001C	1160	.08	.02	.06	.01	.04	0	1	0	1	0	0	0.3	0	?	?	5	3	4	4	5	5	5	5	5	5	5	5	5	5	5	5	5	
YAL002W	1176	.09	.02	.06	.01	.04	0	0	0	0	0	1	0.2	?	?	8	4	2	3	3	4	4	4	4	4	4	4	4	4	4	4	4	4	
YAL003W	206	.08	.02	.06	.01	.04	0	0	0	0	0	0	19.1	19	?	23	70	73	91	68	105	52	112	88	64	168	108	104	75	103	140	98	126	
YAL004W	215	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	?	?	18	12	9	5	3	3	3	3	3	3	3	3	3	3	3	3	3	3
YAL005C	641	.08	.02	.06	.01	.04	0	0	0	0	0	1	13.4	16	?	17	39	38	30	13	17	8	11	8	7	8	6	8	8	7	9	8	14	
YAL007C	193	.08	.02	.06	.01	.04	0	0	0	0	1	4	2.2	8	?	?	15	21	32	20	21	19	29	19	16	22	26	28	23	22	25	16	17	
YAL008W	198	.08	.02	.06	.01	.04	0	0	0	0	0	3	1.2	?	?	6	2	4	3	3	3	4	4	4	4	4	4	4	4	4	4	4	4	4
YAL009W	259	.08	.02	.06	.01	.04	0	2	0	0	0	3	0.6	?	?	11	6	4	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
YAL010C	483	.08	.02	.06	.02	.04	0	0	0	0	0	1	0.3	?	?	11	6	4	3	3	3	4	4	4	4	4	4	4	4	4	4	4	4	4
YAL011W	616	.08	.02	.06	.01	.04	0	0	0	1	0	0	0.4	?	?	6	5	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
YAL012W	393	.08	.02	.06	.01	.04	0	0	0	0	0	1	8.9	4	?	6.7	29	26	25	27	53	26	43	36	25	28	23	28	31	28	34	23	28	
YAL013W	362	.08	.02	.06	.01	.04	0	0	0	0	0	0	0.6	?	?	7	9	6	5	14	6	12	14	10	9	9	10	8	8	6	10			
YAL014C	202	.08	.02	.06	.01	.04	0	0	0	0	0	0	1.1	?	?	12	13	10	8	10	10	12	13	12	14	11	11	11	10	11	9	12		
YAL015C	399	.08	.02	.06	.01	.04	0	1	0	0	0	0	0.7	0	1	19	18	14	10	14	12	17	17	14	13	11	13	16	11	14	12	13		
YAL016W	635	.08	.02	.06	.01	.04	0	0	0	0	0	1	3.3	5	?	15	20	20	102	20	20	30	22	18	19	19	20	21	21	23	16	16		
YAL017W	1356	.08	.02	.06	.01	.04	0	0	0	0	0	0	0.4	?	?	14	3	3	4	8	5	6	6	5	5	8	9	10	6	5	4	7		
YAL018C	325	.08	.02	.06	.01	.04	0	0	0	0	0	4	?	?	?	4	2	2	2	1	1	2	2	2	2	2	2	2	2	2	2	2	2	

4 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Functional Classification

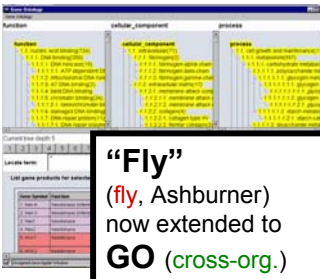


COGs
(cross-org.,
just conserved,
NCBI
Koonin/Lipman)




GenProtEC
(*E. coli*, Riley)

ENZYME
(SwissProt
Bairoch/
Apweiler,
just enzymes,
cross-org.)



"Fly"
(fly, Ashburner)
now extended to
GO (cross-org.)



MIPS/PEDANT
(yeast, Mewes)

Also:
Other
SwissProt
Annotation
WIT, KEGG
(just pathways)
TIGR EGAD
(human ESTs)
SGD (yeast)

5 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Core Prediction of Function on a Genomic Scale from Array Data & Sequence Features

Gene

Gene	Array Experiments		Sequence Features	
	Expr. Level	Proteome Chip	Length	Amino Acid Composition
YAL001C	0.5	0.3	5	3
YAL002W	0.2	0.2	5	4
YAL003W	10.1	0.909	70	72
YAL004W	0.632			
YAL005C	13.4	0.338	38	38

Gene

Gene	Array Experiments		Sequence Features	
	Expr. Level	Proteome Chip	Length	Amino Acid Composition
1160	0.8	0.2	0.04	1
1170	0.8	0.2	0.04	0
206	0.8	0.2	0.04	0
215	0.8	0.2	0.04	0
641	0.8	0.2	0.04	0

"Function" Description

Gene	Array Experiments		Sequence Features		"Function" Description		
	Expr. Level	Proteome Chip	Length	Amino Acid Composition	phrase description	Std. func. # (from MIPS)	Complex #
TFIIC (transcription initiation)	4.2	0	6	6			
ribosomal protein S11	6.8	0	6	6			
heat shock protein of HSP70	1.1	no	14	14			
heat shock protein of HSP71	4.8	no	4	4			

6000+

Different Aspects of function: molecular action, cellular role, phenotypic manifestation
Also: localization, interactions, complexes

6 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Arrange data in a tabulated form, each row representing an example and each column representing a feature, including the dependent experimental quantity to be predicted.

	predictor1	Predictor2	predictor3	predictor4	response
G1	A(1,1)	A(1,2)	A(1,3)	A(1,4)	Class A
G2	A(2,1)	A(2,2)	A(2,3)	A(2,4)	Class A
G3	A(3,1)	A(3,2)	A(3,3)	A(3,4)	Class B

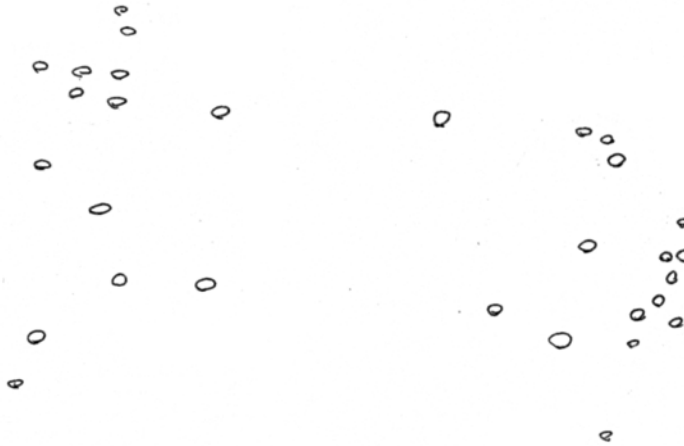
(adapted from Y Kluger)

Typical Predictors and Response for Yeast

Basics		Predictors												Response															
		Sequence Features						Genomic Features						Function	Localization														
Yeast Gene ID	Sequence	seq. length	Amino Acid Composition						How many times does the sequence have these motif features?	Abs. expr. Level (mRNA copies / cell)	Prot. Abundance	Cell cycle timecourse				function ID(s) (from MIPS)	function description	5-compartment											
			A	C	D	E	K	R	S	T	V	W	Y	farn site	NLS	total motif	myc2	signalp	tmms1	Gene-Chip expt. from RY Lab	sage tag freq. (1000 copies /cell)	E=0	E=1	E=15	E=16				
YAL001C	MNIFEMLR	1160	.08	.02	.06	.01	.04	0	1	0	0	0	0	0	0	0	1	0	0	0.3	0	?	5	3	4	04.01.01.04.03	TFIIIC (transcription initi	N	
YAL002W	KVFGRCLELA	1176	.09	.02	.06	.01	.04	0	0	0	0	0	0	0	0	1	0	0	1	0.2	?	8	4	4	06.04.08.13	vacuolar sorting protein	C		
YAL003W	RMLQFNLRW	206	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	0	0	19.1	19	23	70	73	98	05.04.30.03	translation elongation fac	N	
YAL004W	RPDFCLEPF	215	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	0	0	?	0	?	18	12	4	01.01.01	0	N	
YAL005C	VINTFDGVA	641	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	0	1	13.4	16	17	39	38	8	06.01.06.04.08	heat shock protein of HS	????	
YAL007C	KKAVINGEQ	190	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	1	4	2.2	8	?	15	20	16	99	????	????	
YAL008W	HPETLVKVK	198	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	3	1.2	?	?	9	6	2	3	99	????	????	
YAL009W	PTLEWFLSH	259	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	3	0.6	?	?	6	2	3	03.10.03.13	meiotic protein	????		
YAL010C	MEQRTILKD	493	.08	.02	.06	.02	.04	0	0	0	0	0	0	0	0	0	0	1	0.3	?	?	11	6	6	6	30.16	involved in mitochondrial	????	
YAL011W	KSFPEVGR	616	.08	.02	.06	.01	.04	0	8	0	0	0	0	0	0	0	1	0	0	0.4	?	?	6	5	5	30.16.99	protein of unknown func	????	
YAL012W	GVQVEITSE	393	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	1	8.9	4	6.7	29	26	23	20	01.01.01.30.03	cystathionine gamma-ly	C	
YAL013W	RTDCYGRNV	362	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	0	0.6	?	?	7	9	6	10	01.06.10.30.03	regulator of phospholipid	N	
YAL014C	GDVGRKKI	202	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	0	1.1	?	?	12	13	9	12	99	????	?	
YAL015C	MTPAVTTYK	399	.08	.02	.06	.01	.04	0	1	0	0	0	0	0	0	0	0	0	0	0.7	0	1	19	18	12	13	11.01.11.04	DNA repair protein	N
YAL016W	KKPLTQQL	635	.08	.02	.06	.01	.04	0	0	0	0	0	0	0	0	0	0	1	3.3	5	?	15	20	16	18	03.01.03.04.03	ser/thr protein phosphat	????	

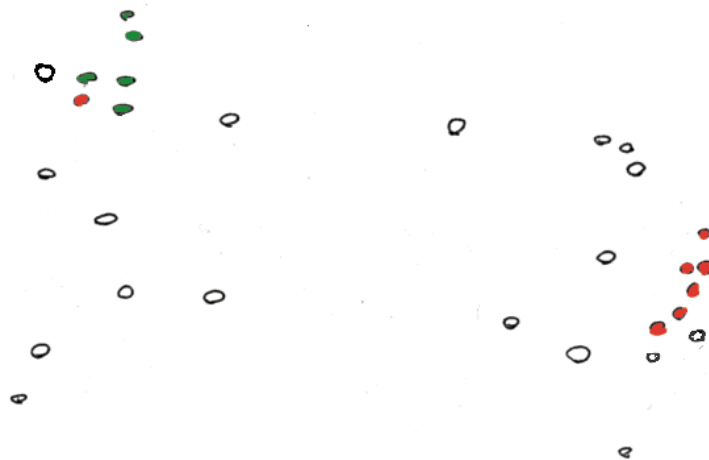
Represent predictors in abstract high dimensional space

Core

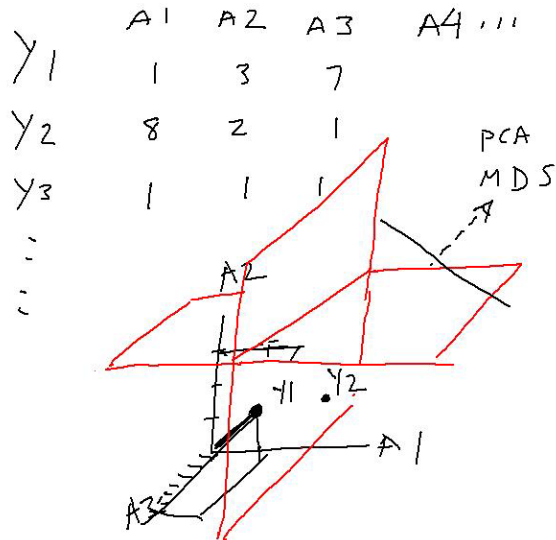


“Tag” Certain Points

Core



Abstract high-dimensional space representation



11 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Large-scale Datamining

- Gene Expression
 - ◇ Representing Data in a Grid
 - ◇ Description of function prediction in abstract context
- Unsupervised Learning
 - ◇ clustering & k-means
 - ◇ Local clustering
- Supervised Learning
 - ◇ Discriminants & Decision Tree
 - ◇ Bayesian Nets
- Function Prediction EX
 - ◇ Simple Bayesian Approach for Localization Prediction

12 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

“cluster” predictors

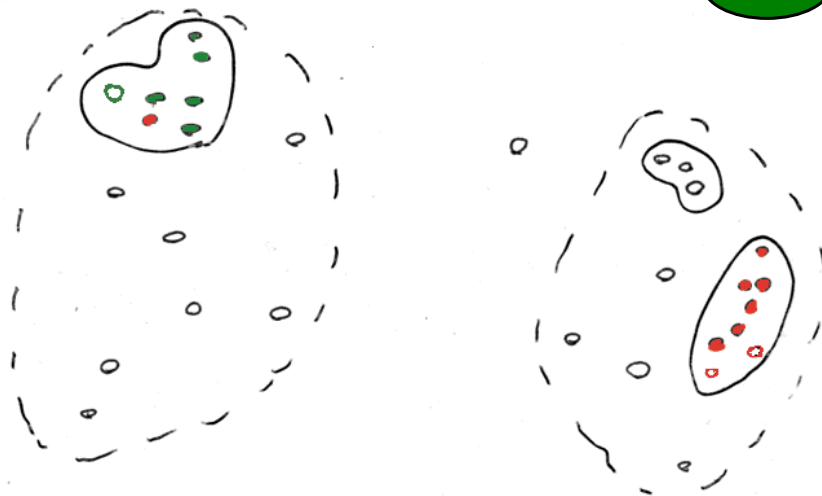
Core



13 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Use clusters to predict Response

Core



14 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

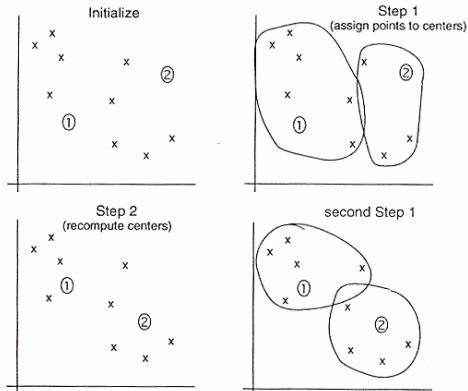
K-means



Heuristics Research, Inc.

Core

K-means algorithm in 2-D clustering



6

© Copyright 1995

15 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

K-means

Top-down vs. Bottom up

Top-down when you know how many subdivisions

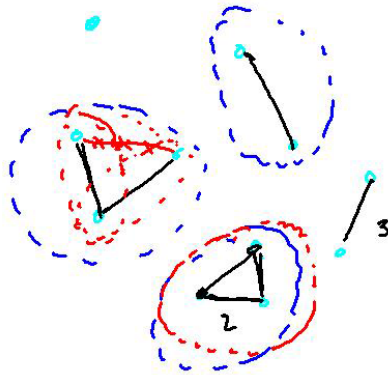
k-means as an example of top-down

- 1) Pick ten (i.e. k ?) random points as putative cluster centers.
- 2) Group the points to be clustered by the center to which they are closest.
- 3) Then take the mean of each group and repeat, with the means now at the cluster center.
- 4) I suppose you stop when the centers stop moving.

16 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Bottom up clustering

Core



SINGLE LINK
MULTI LINK

(THRESHOLD)



BOTTOM UP

TOP-DOWN

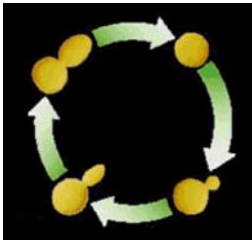


K MEANS

Large-scale Datamining

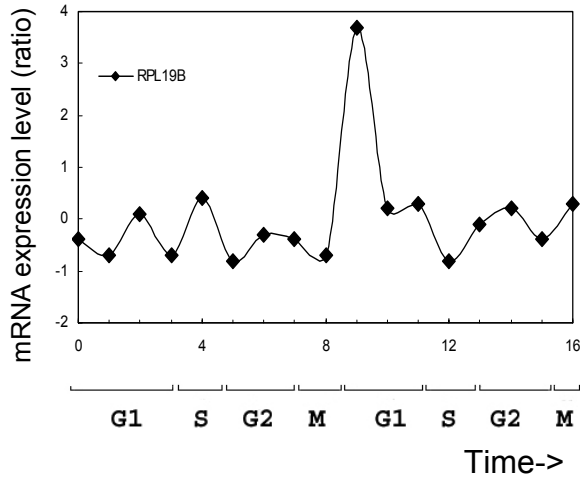
- Gene Expression
 - ◇ Representing Data in a Grid
 - ◇ Description of function prediction in abstract context
- Unsupervised Learning
 - ◇ clustering & k-means
 - ◇ Local clustering
- Supervised Learning
 - ◇ Discriminants & Decision Tree
 - ◇ Bayesian Nets
- Function Prediction EX
 - ◇ Simple Bayesian Approach for Localization Prediction

Clustering
the
yeast cell
cycle to
uncover
interacting
proteins



Extra

[Brown, Davis]

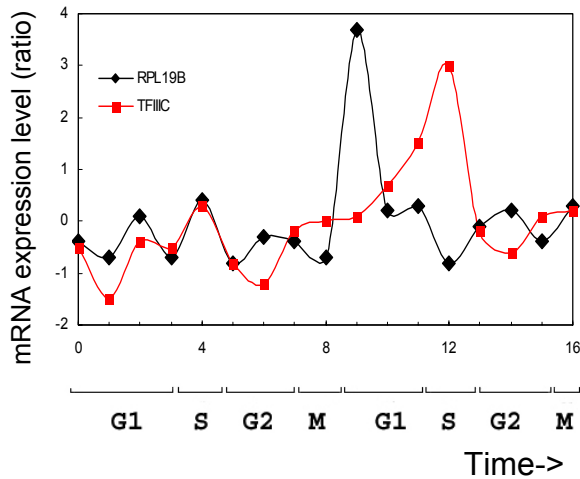


Microarray timecourse of
1 ribosomal protein

Clustering
the
yeast cell
cycle to
uncover
interacting
proteins



Extra



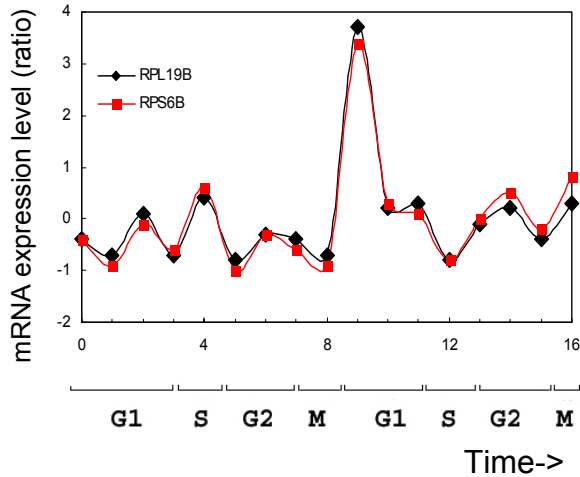
Random relationship from ~18M

Clustering
the
yeast cell
cycle to
uncover
interacting
proteins



Extra

[Botstein; Church, Vidal]

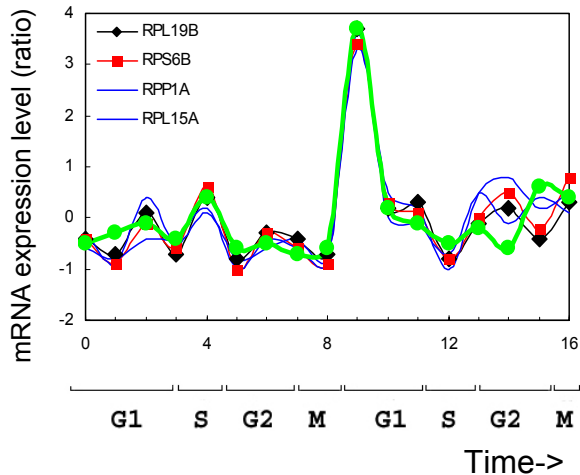


Close relationship from 18M
(2 Interacting Ribosomal Proteins)

Clustering
the
yeast cell
cycle to
uncover
interacting
proteins



Extra



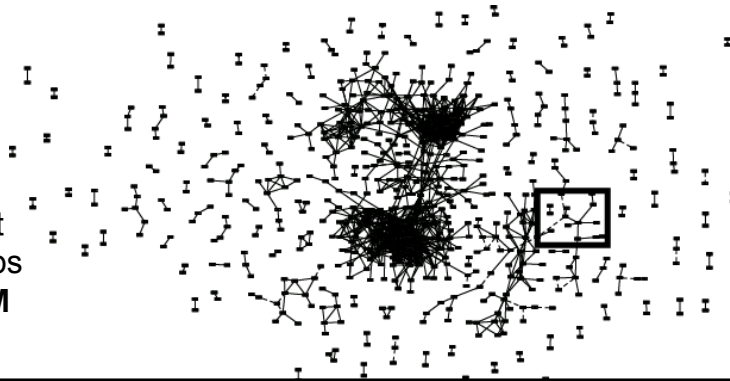
Predict Functional Interaction of
Unknown Member of Cluster



Global Network of Relationships

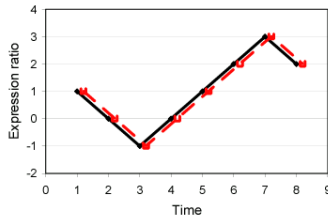
Core

~470K significant relationships from ~18M possible

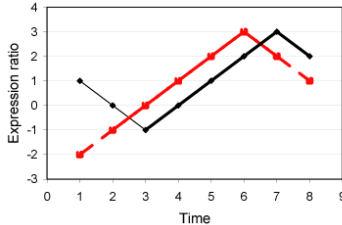


Simultaneous

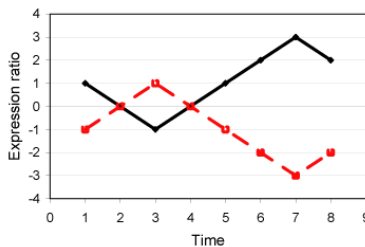
Traditional Global Correlation



Time-Shifted



Inverted



Local Clustering algorithm identifies further (reasonable) types of expression relationships

Core

[Church]

Local Alignment

Suppose there are n (1, 2, ..., n) time points:

➤ The expression ratio is normalized in “Z-score” fashion;

➤ Score matrix: $S_{i,j} = S(x_i, y_j) = x_i \cdot y_j$;

Qian J. et al. Beyond Synexpression Relationships: Local clustering of Time-shifted and Inverted Gene Expression Profiles Identifies New, Biologically Relevant Interactions. J. Mol. Biol. (2001) 314, 1053-1066

Local Alignment

Suppose there are n (1, 2, ..., n) time points:

➤ Sum matrices $E_{i,j}$ and $D_{i,j}$:

$$E_{i,j} = \max(E_{i-1,j-1} + S_{i,j}, 0);$$

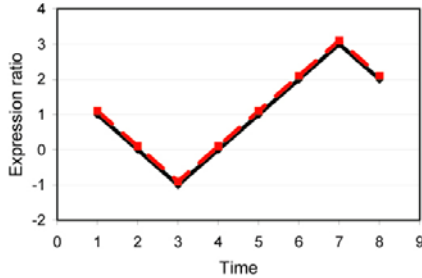
$$D_{i,j} = \max(D_{i-1,j-1} - S_{i,j}, 0);$$

➤ Match Score = $\max(E_{i,j}, D_{i,j})$

Qian J. et al. Beyond Synexpression Relationships: Local clustering of Time-shifted and Inverted Gene Expression Profiles Identifies New, Biologically Relevant Interactions. J. Mol. Biol. (2001) 314, 1053-1066

Local Alignment

Simultaneous



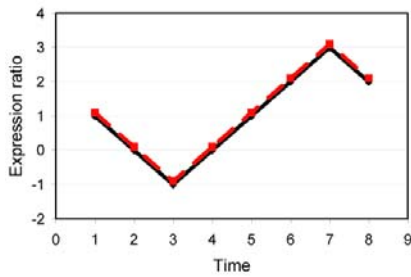
		1	0	-1	0	1	2	3	2
0	0	0	0	0	0	0	0	0	0
1	0	1	0	-1	0	1	2	3	2
0	0	0	0	0	0	0	0	0	0
-1	0	-1	0	1	0	-1	-2	-3	-2
0	0	0	0	0	0	0	0	0	0
1	0	1	0	-1	0	1	2	3	2
2	0	2	0	-1	0	2	4	6	4
3	0	3	0	-3	0	3	6	9	6
2	0	1	0	-2	0	2	4	6	4

$$S_{i,j} = x_i \cdot y_j$$

Qian J. et al. Beyond Synexpression Relationships: Local clustering of Time-shifted and Inverted Gene Expression Profiles Identifies New, Biologically Relevant Interactions. J. Mol. Biol. (2001) 314, 1053-1066

Local Alignment

Simultaneous



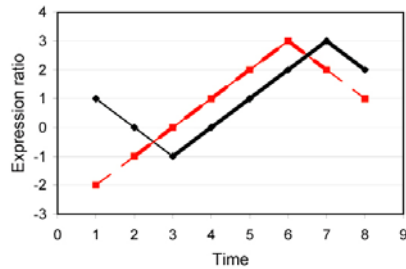
		1	0	-1	0	1	2	3	2
0	0	0	0	0	0	0	0	0	0
1	0	1	0	0	0	1	2	3	2
0	0	0	1	0	0	0	1	2	3
-1	0	0	0	2	0	0	0	0	0
0	0	0	0	0	2	0	0	0	0
1	0	1	0	0	0	3	2	3	2
2	0	2	1	0	0	2	7	8	7
3	0	3	2	0	0	3	8	16	14
2	0	2	3	0	0	2	7	14	20

$$E_{i,j} = \max(E_{i-1,j-1} + x_i \cdot y_j, 0)$$

Qian J. et al. Beyond Synexpression Relationships: Local clustering of Time-shifted and Inverted Gene Expression Profiles Identifies New, Biologically Relevant Interactions. J. Mol. Biol. (2001) 314, 1053-1066

Local Alignment

Time-Shifted



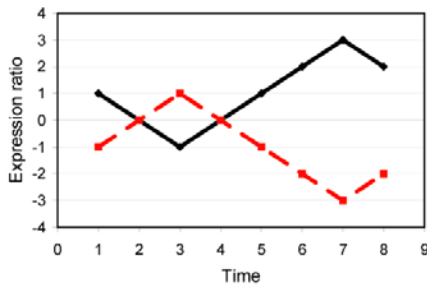
		-2	-1	0	1	2	3	2	1
0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	1	2	3	2	1
0	0	0	0	0	0	1	2	3	2
-1	0	2	1	0	0	0	0	0	2
0	0	0	2	1	0	0	0	0	0
1	0	0	0	2	2	2	3	2	1
2	0	0	0	0	4	6	8	7	4
3	0	0	0	0	3	10	15	14	10
2	0	0	0	0	2	7	16	19	16

$$E_{i,j} = \max(E_{i-1,j-1} + x_i \cdot y_j, 0)$$

Qian J. et al. Beyond Synexpression Relationships: Local clustering of Time-shifted and Inverted Gene Expression Profiles Identifies New, Biologically Relevant Interactions. J. Mol. Biol. (2001) 314, 1053-1066

Local Alignment

Inverted

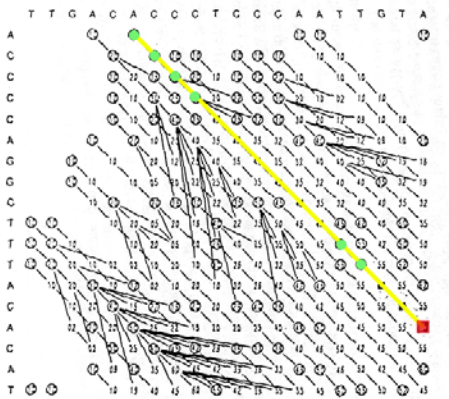


		-1	0	1	0	-1	-2	-3	-2
0	0	0	0	0	0	0	0	0	0
1	0	1	0	0	0	1	2	3	2
0	0	0	1	0	0	0	1	2	3
-1	0	0	0	2	0	0	0	0	2
0	0	0	0	0	2	0	0	0	0
1	0	1	0	0	0	3	2	3	2
2	0	2	1	0	0	2	7	8	7
3	0	3	2	0	0	3	8	16	14
2	0	2	3	0	0	2	7	14	20

$$D_{i,j} = \max(D_{i-1,j-1} - x_i \cdot y_j, 0)$$

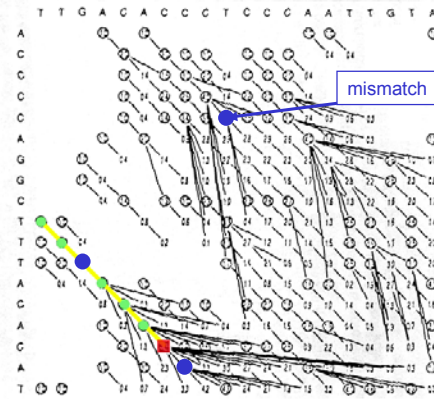
Qian J. et al. Beyond Synexpression Relationships: Local clustering of Time-shifted and Inverted Gene Expression Profiles Identifies New, Biologically Relevant Interactions. J. Mol. Biol. (2001) 314, 1053-1066

Global (NW) vs Local (SW) Alignments



TTGACACCCTCCCAATTGTA...
 |||| | | |
 A C C C C A G G C **TTTACAC**A T
 123444444456667

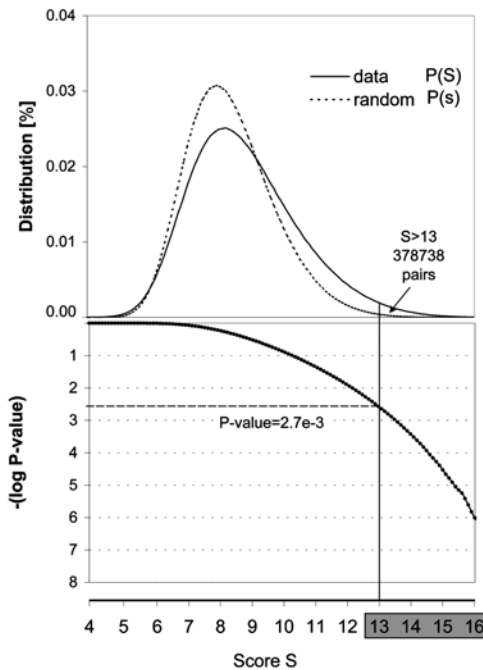
Match Score = +1
 Gap-Opening = -1.2, Gap-Extension = -.03
 for local alignment Mismatch = -0.6



T T G A C A C C . . .
 | | - | | | | -
 T T T A C A C A . . .
 1 2 1 2 3 4 5 4
 0 0 4 4 4 4 4 8

Adapted from D J States & M S Boguski, "Similarity and Homology," Chapter 3 from Gribskov, M. and Devereux, J. (1992). Sequence Analysis Primer. New York, Oxford University Press. (Page 133)

Statistical Scoring



Examples time-shifted relationships

Suggestive

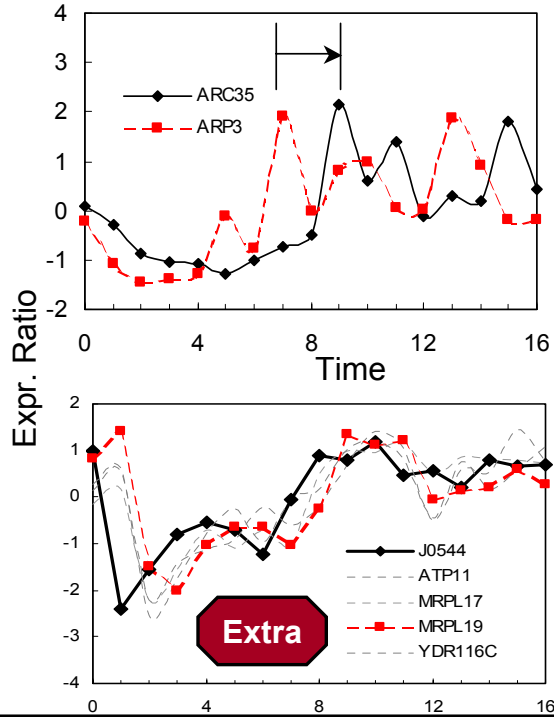
ARP3 : in actin
remodelling cplx.

ARC35 : in same cplx.
(required late in
cell cycle)

Predicted

J0544 : **unknown**
function

MRPL19: mito.ribosome



Examples time-shifted relationships

Suggestive

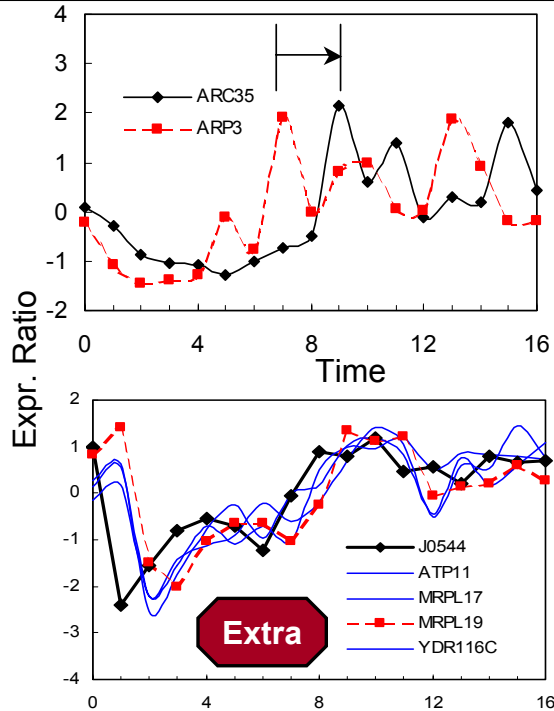
ARP3 : in actin
remodelling cplx.

ARC35 : in same cplx.
(required late in
cell cycle)

Predicted

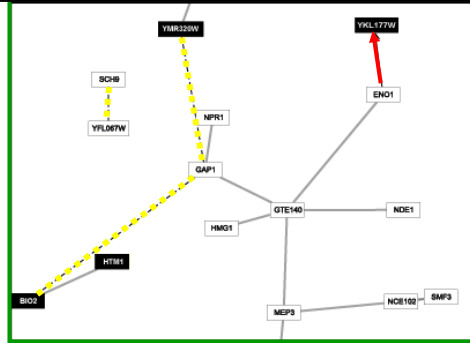
J0544 : **unknown**
function

MRPL19: mito.ribosome



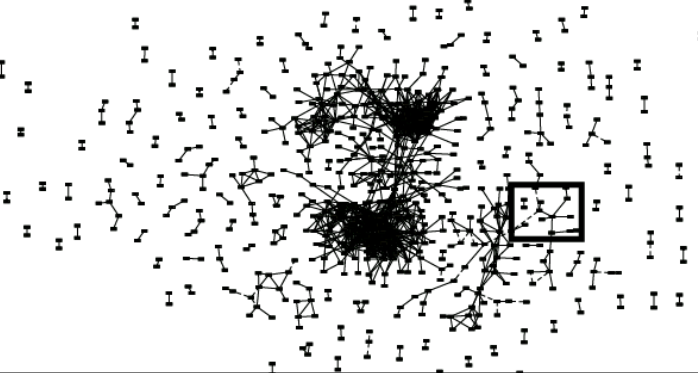
Global Network of 3 Different Types of Relationships

Simultaneous —
Inverted
Shifted →



Extra

~470K
 significant
 relationships
 from ~18M
 possible

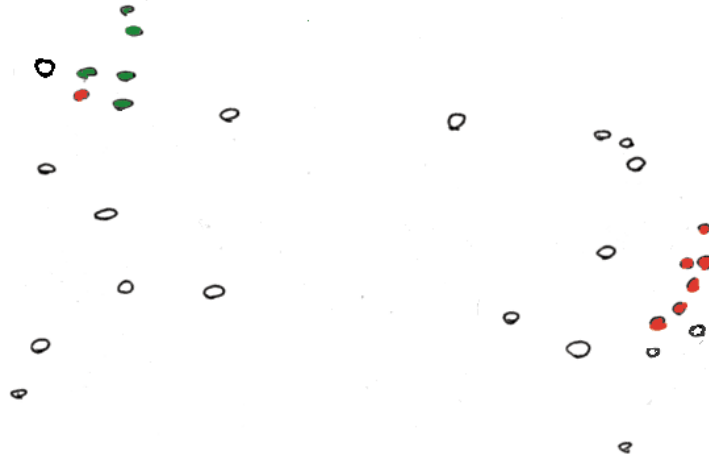


Large-scale Datamining

- Gene Expression
 - ◇ Representing Data in a Grid
 - ◇ Description of function prediction in abstract context
- Unsupervised Learning
 - ◇ clustering & k-means
 - ◇ Local clustering
- Supervised Learning
 - ◇ Discriminants & Decision Tree
 - ◇ Bayesian Nets
- Function Prediction EX
 - ◇ Simple Bayesian Approach for Localization Prediction

"Tag" Certain Points

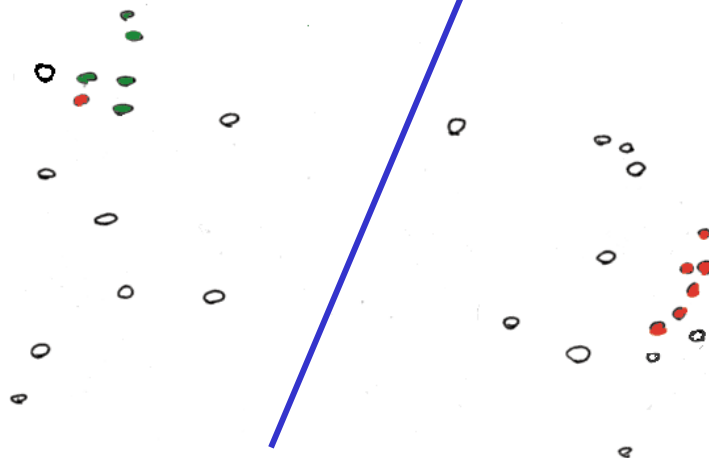
Core



37 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

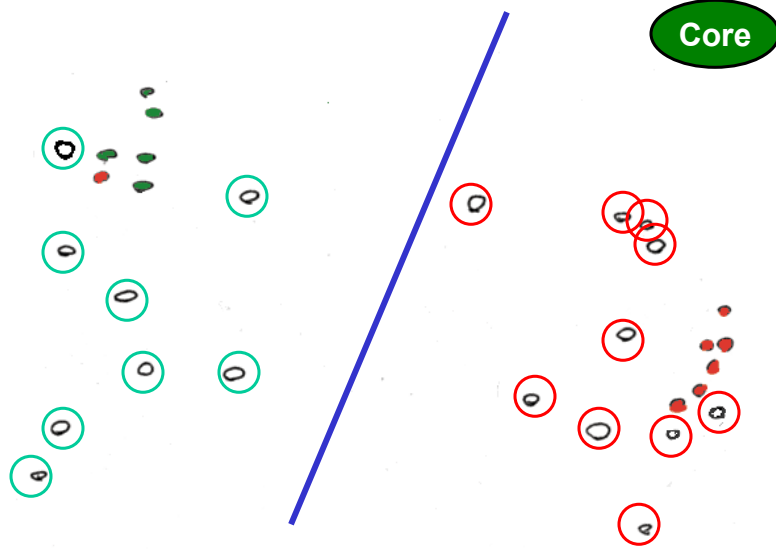
Find a Division to Separate Tagged Points

Core



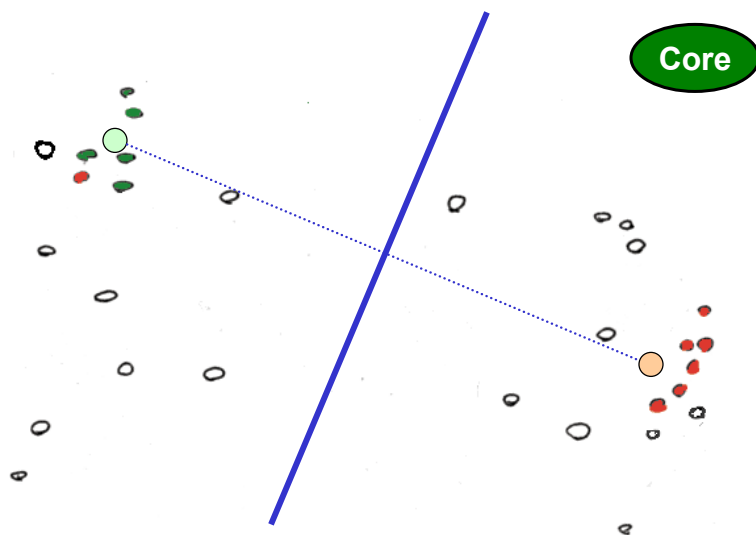
38 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Extrapolate to Untagged Points



39 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Discriminant to Position Plane

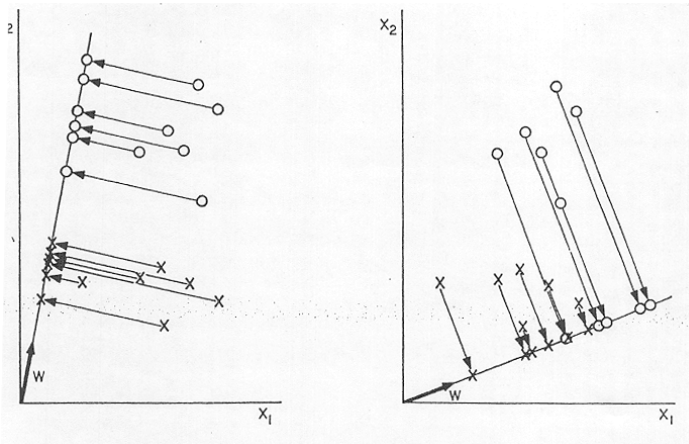


40 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Fisher discriminant analysis

- Use the training set to reveal the structure of class distribution by seeking a linear combination
- $y = w_1x_1 + w_2x_2 + \dots + w_nx_n$ which maximizes the ratio of the separation of the class means to the sum of each class variance (within class variance). This linear combination is called the first linear discriminant or first canonical variate. Classification of a future case is then determined by choosing the nearest class in the space of the first linear discriminant and significant subsequent discriminants, which maximally separate the class means and are constrained to be uncorrelated with previous ones.

Fischer's Discriminant



(Adapted from ???)

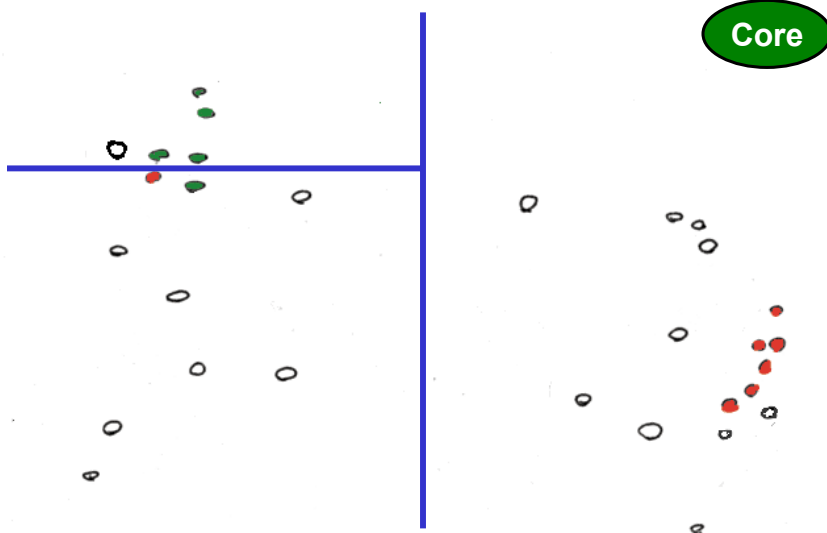
Fisher cont.

$$m_i = \underline{w} \cdot \underline{m}_i \quad s_i^2 = \sum_{y \in Y_i} (y - m_i)^2$$

Solution of 1st
variate

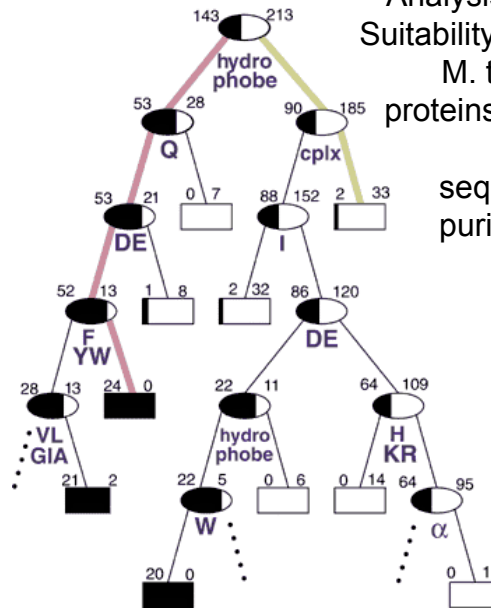
$$\underline{w} = S_W^{-1} (\underline{m}_1 - \underline{m}_2)$$

Find a Division to Separate Tagged Points



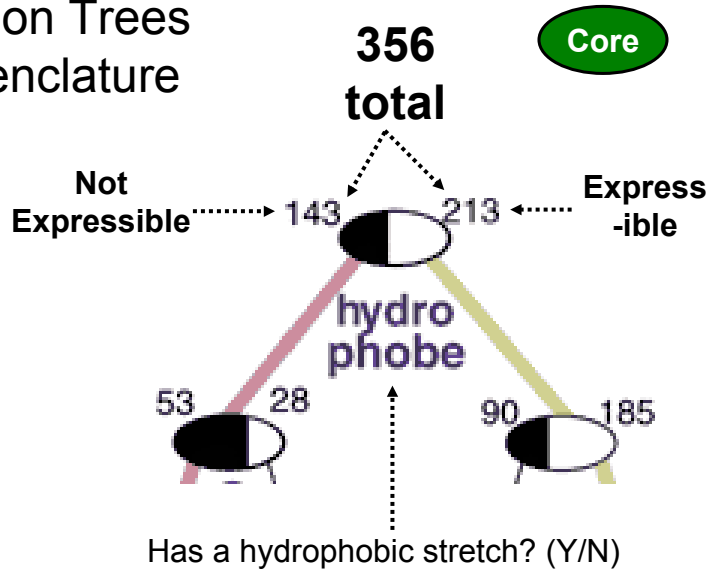
Retrospective Decision Trees

Analysis of the Suitability of 500 M. thermo. proteins to find optimal sequences purification

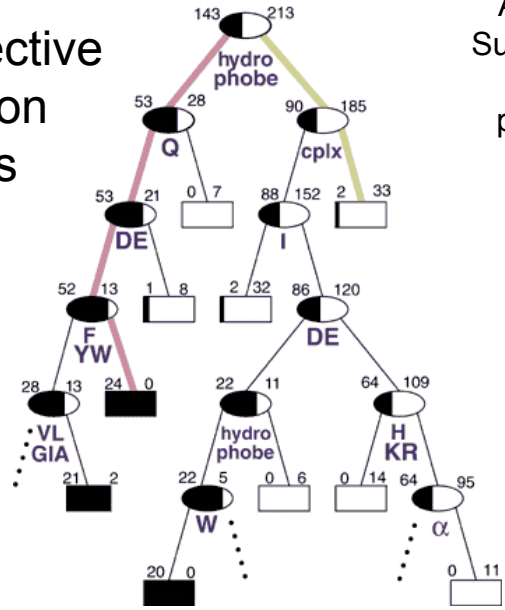


Core

Retrospective Decision Trees Nomenclature



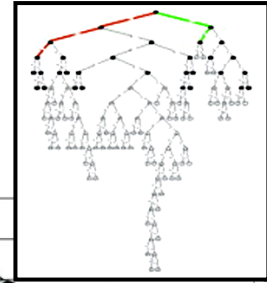
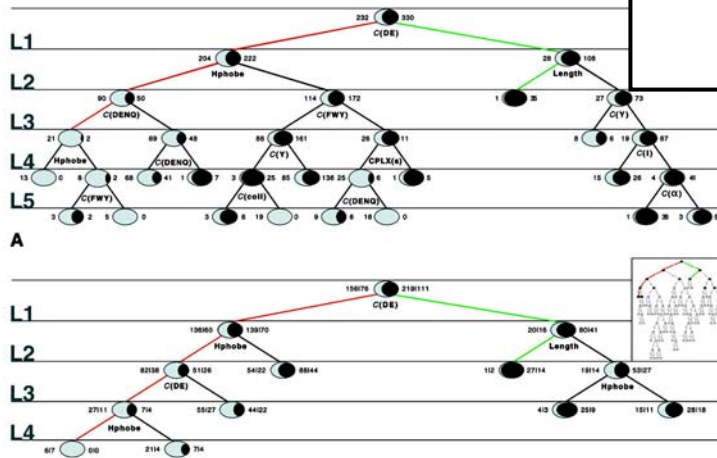
Retrospective Decision Trees



Analysis of the Suitability of 500 M. thermo. proteins to find optimal sequences purification

Not Express Express

Overfitting, Cross Validation, and Pruning



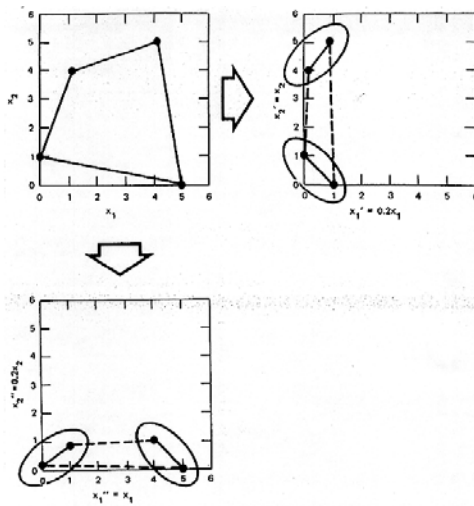
Decision Trees

- can handle data that is not linearly separable.
- A decision tree is an upside down tree in which each branch node represents a choice between a number of alternatives, and each leaf node represents a classification or *decision*. One classifies instances by sorting them down the tree from the root to some leaf nodes. To classify an instance the tree calls first for a test at the root node, testing the feature indicated on this node and choosing the next node connected to the root branch where the outcome agrees with the value of the feature of that instance. Thereafter a second test on another feature is made on the next node. This process is then repeated until a leaf of the tree is reached.
- Growing the tree, based on a training set, requires strategies for (a) splitting the nodes and (b) pruning the tree. Maximizing the decrease in average impurity is a common criterion for splitting. In a problem with noisy data (where distribution of observations from the classes overlap) growing the tree will usually over-fit the training set. The strategy in most of the cost-complexity pruning algorithms is to choose the smallest tree whose error rate performance is close to the minimal error rate of the over-fitted larger tree. More specifically, growing the trees is based on splitting the node that maximizes the reduction in deviance (or any other impurity-measure of the distribution at a node) over all allowed binary splits of all terminal nodes. Splits are *not* chosen based on misclassification rate. A binary split for a continuous feature variable v is of the form $v < \text{threshold}$ versus $v > \text{threshold}$ and for a "descriptive" factor it divides the factor's levels into two classes. Decision tree-models have been successfully applied in a broad range of domains. Their popularity arises from the following: Decision trees are easy to interpret and use when the predictors are a mix of numeric and nonnumeric (factor) variables. They are invariant to scaling or re-expression of numeric variables. Compared with linear and additive models they are effective in treating missing values and capturing non-additive behavior. They can also be used to predict nonnumeric dependent variables with more than two levels. In addition, decision tree models are useful to devise prediction rules, screen the variables and summarize the multivariate data set in a comprehensive fashion. We also note that ANN and decision tree learning often have comparable prediction accuracy [Mitchell p. 85] and SVM algorithms are slower compared with decision tree. These facts suggest that the decision tree method should be one of our top candidates to "data-mine" proteomics datasets. C4.5 and CART are among the most popular decision tree algorithms.

Optional: not needed for Quiz

(adapted from Y Kluger)

Effect of Scaling



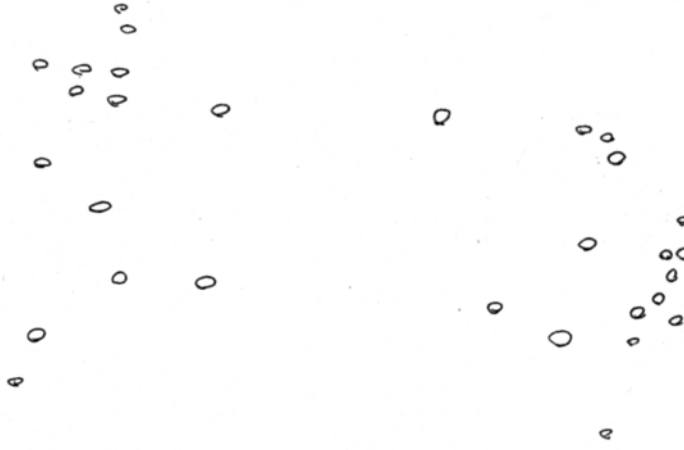
(adapted from ref?)

**End of class 2002,12.01
(Bioinfo-13)
[started at beg. of datamining]**

Large-scale Datamining

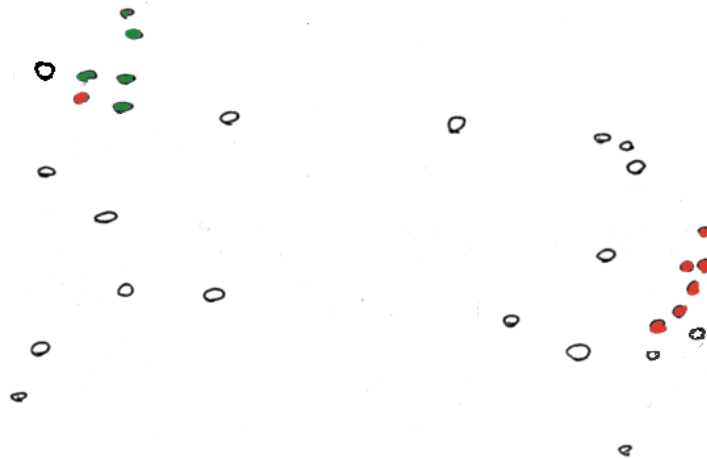
- Gene Expression
 - ◇ Representing Data in a Grid
 - ◇ Description of function prediction in abstract context
- Unsupervised Learning
 - ◇ clustering & k-means
 - ◇ Local clustering
- Supervised Learning
 - ◇ Discriminants & Decision Tree
 - ◇ Bayesian Nets
- Function Prediction EX
 - ◇ Simple Bayesian Approach for Localization Prediction

Represent predictors in abstract high dimensional space



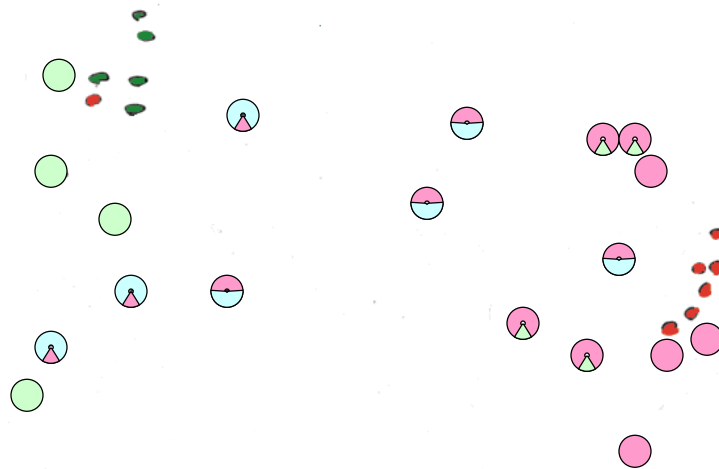
53 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Tagged Data



54 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Probabilistic Predictions of Class



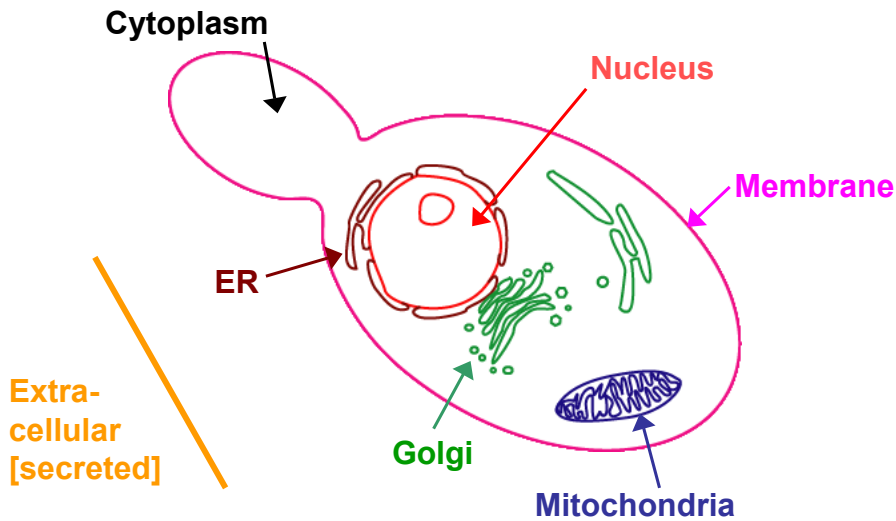
55 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Large-scale Datamining

- Gene Expression
 - ◇ Representing Data in a Grid
 - ◇ Description of function prediction in abstract context
- Unsupervised Learning
 - ◇ clustering & k-means
 - ◇ Local clustering
- Supervised Learning
 - ◇ Discriminants & Decision Tree
 - ◇ Bayesian Nets
- Function Prediction EX
 - ◇ Simple Bayesian Approach for Localization Prediction

56 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

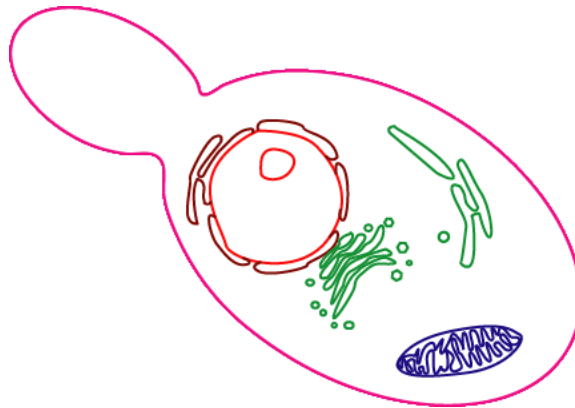
Subcellular Localization, a standardized aspect of function



57 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

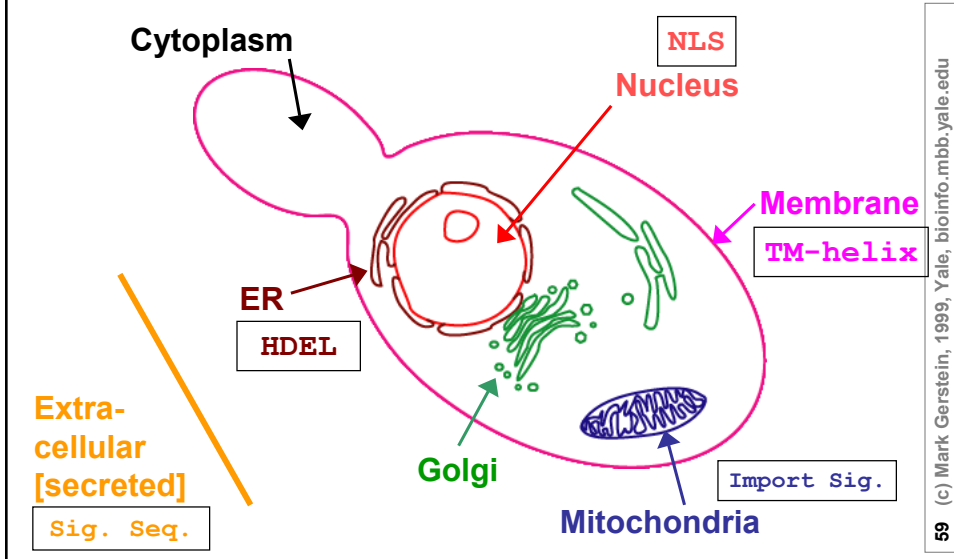
Subcellular Localization, Provides a simple goal for genome-scale functional prediction

Determine how many of the ~6000 yeast proteins go into each compartment

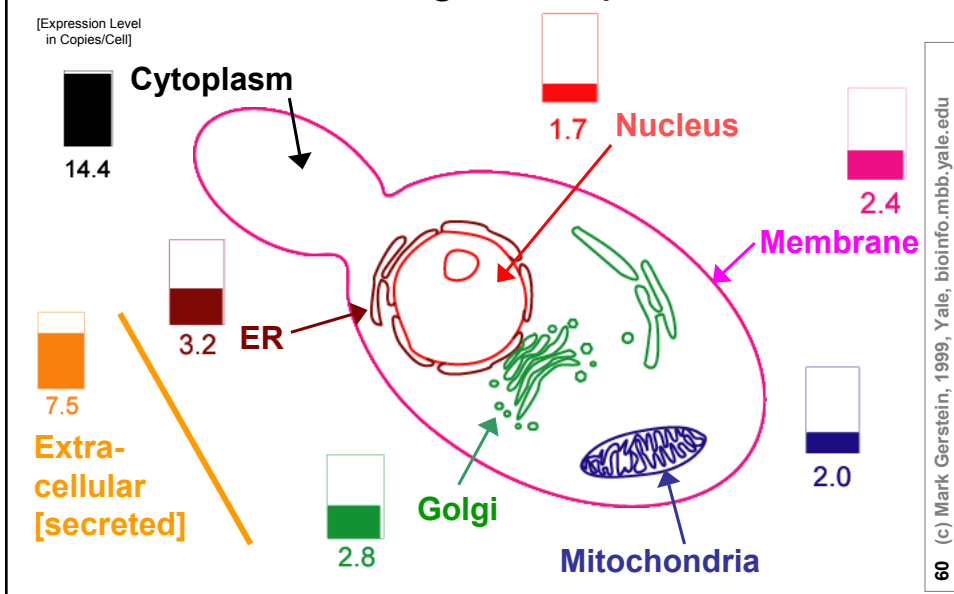


58 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

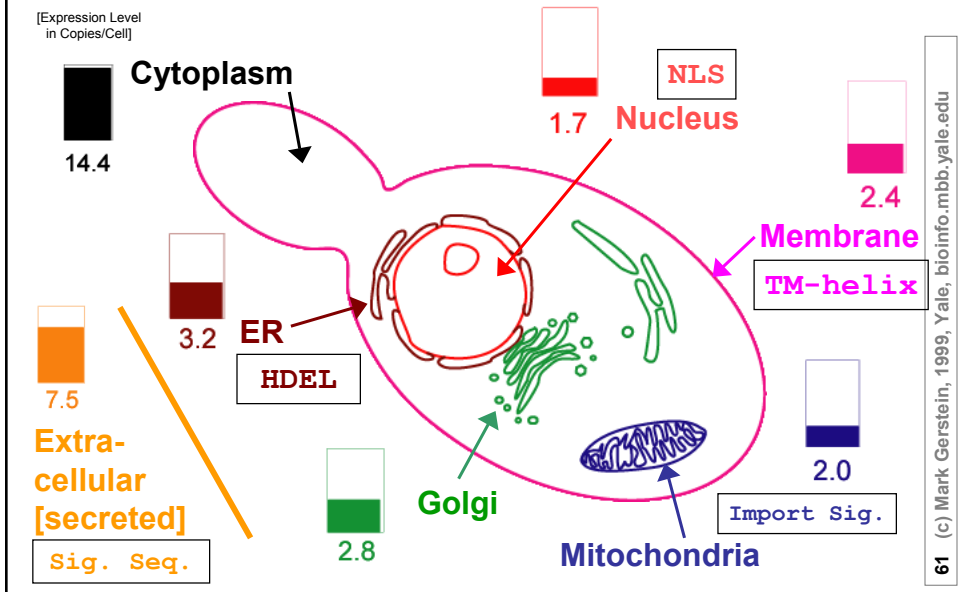
"Traditionally" subcellular localization is "predicted" by sequence patterns



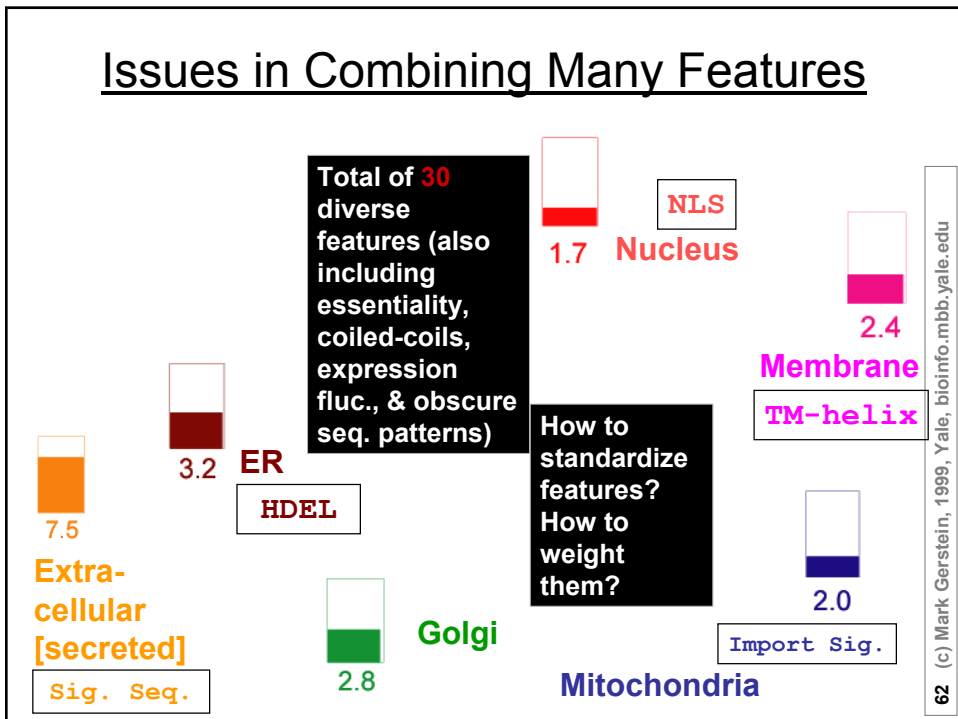
Subcellular localization is associated with the level of gene expression



Combine Expression Information & Sequence Patterns to Predict Localization



Issues in Combining Many Features



Bayesian System for Localizing Proteins

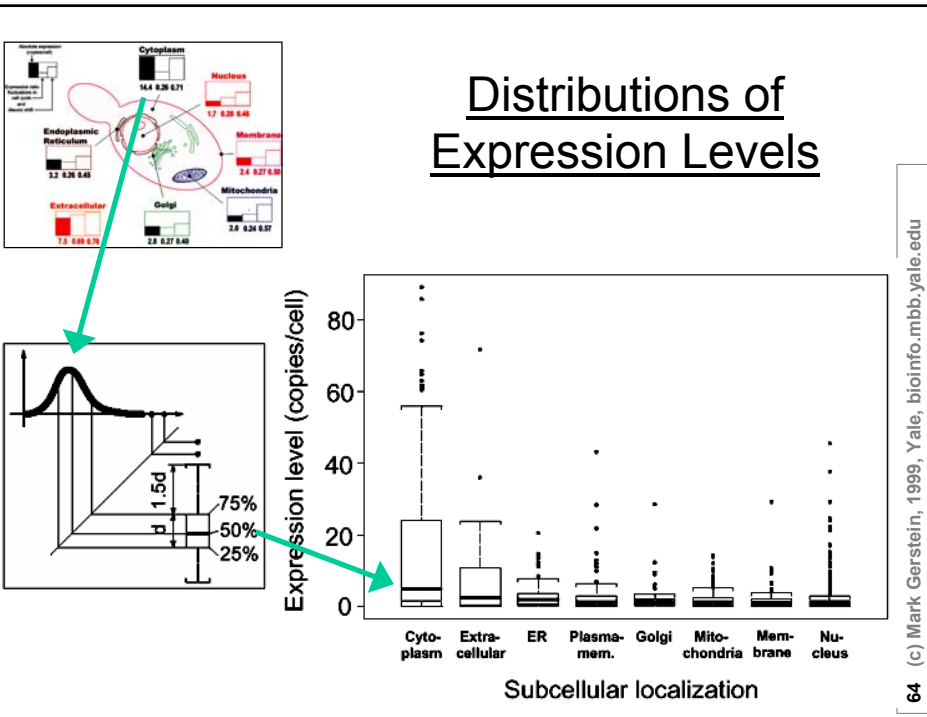
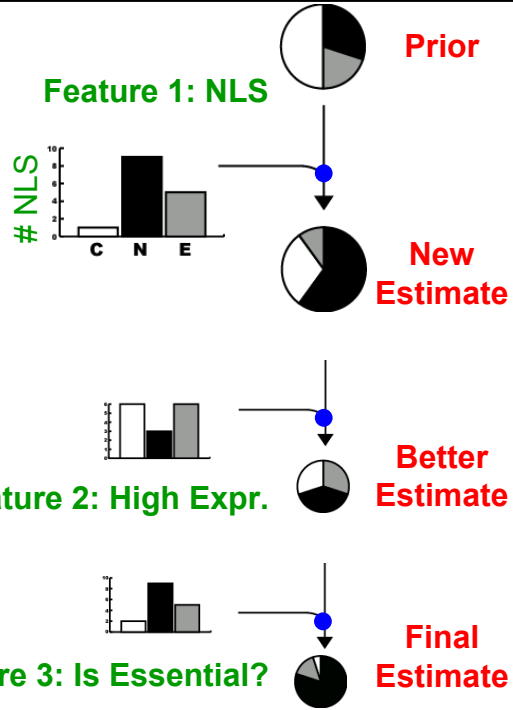
a) Everything in **standard probabilistic terms**

(Handles indefinite proteins, 50% cyt., 50% nuc.)

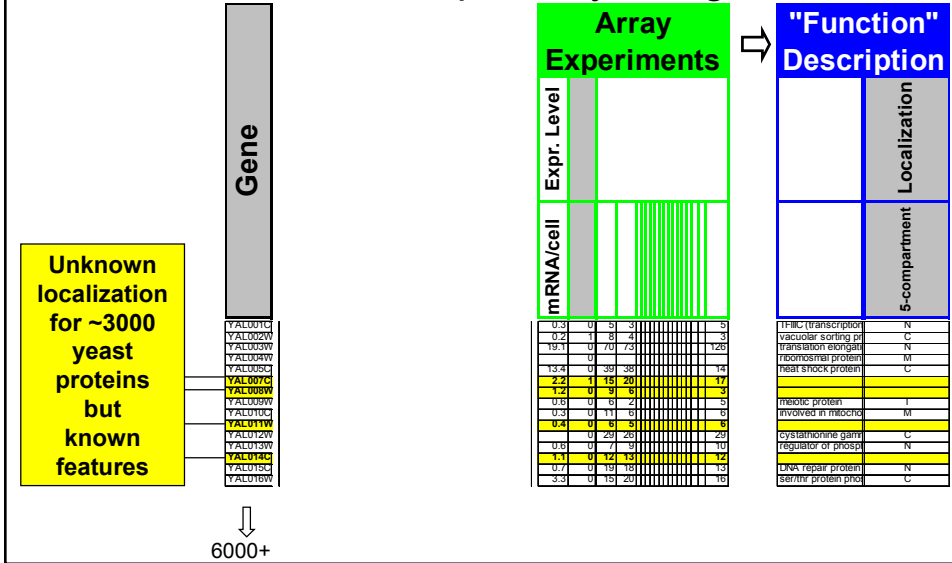
b) Sequentially **combine features** using **Bayes Rule**

(**Feature** x **Prior** / Normalization)

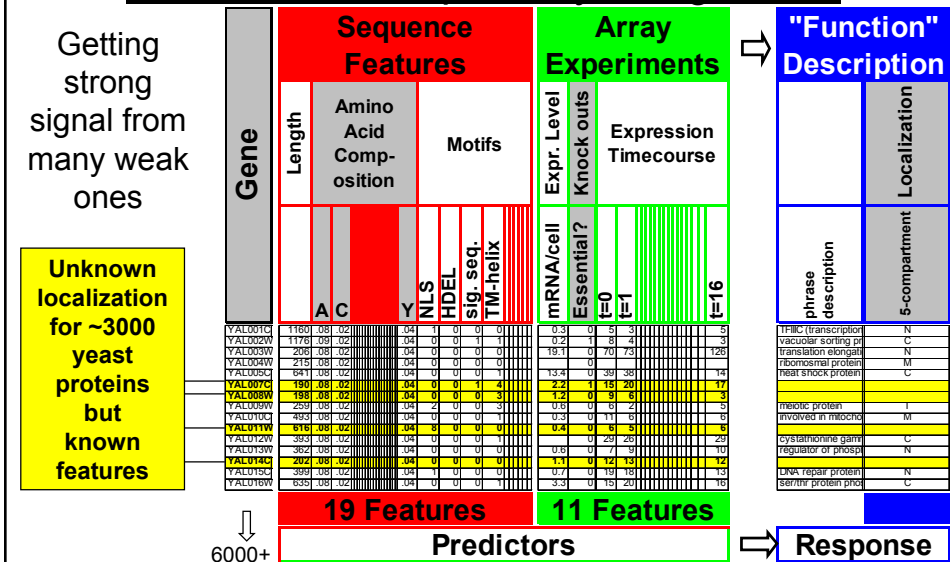
c) Final estimate naturally **weights features**



Integrate heterogeneous set of 30 diverse features to predict localization for uncharacterized part of yeast genome

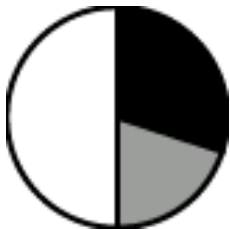
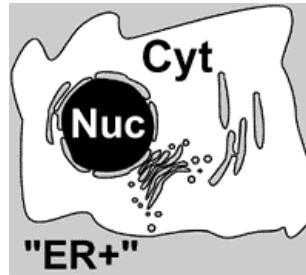


Integrate heterogeneous set of 30 diverse features to predict localization for uncharacterized part of yeast genome



Bayesian System for Localizing Proteins:
Prior

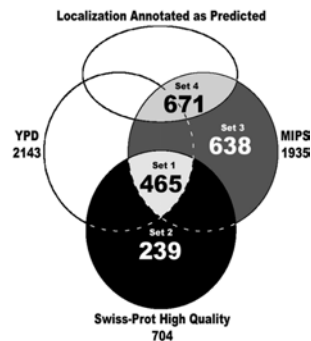
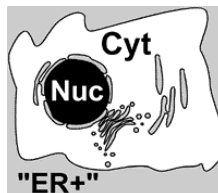
Simplified Cell: 3 (5) compartment



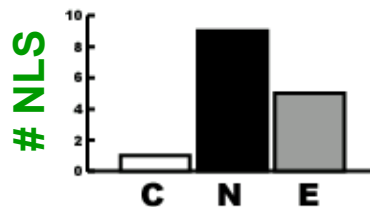
Prior probability distribution for a protein to be in each compartment

Bayesian System for Localizing Proteins:
Features

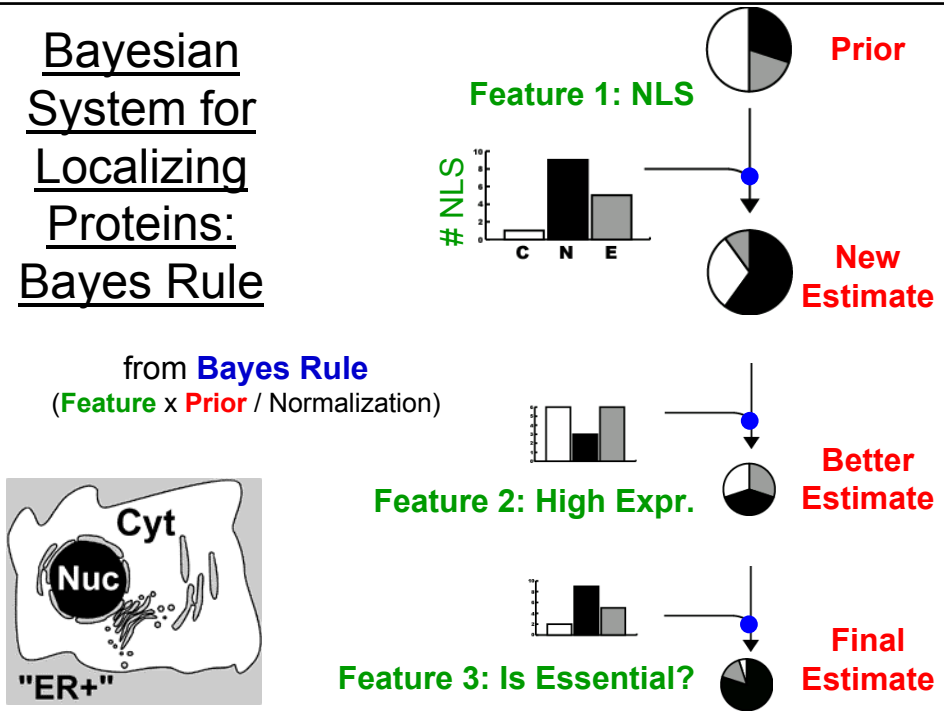
Training Data: 1342 proteins with known localizations



Tabulate occurrence of **feature** across compartments in training data for all 30 features



Bayesian System for Localizing Proteins: Bayes Rule



$$P(c|F) = P(F|c) P(c) / P(F)$$

P(c|F): Probability that protein is in class c given it has feature F

P(F|c): Probability in training data that a protein has feature F if it is class c

P(c): Prior probability that that protein is in class c

P(F): Normalization factor set so that sum over all classes c and ~c is 1 – i.e. $P(c|F) + P(\sim c|F) = 1$

This formula can be iterated with

P(c) [at iter. **i+1**] \leq **P(c|F)** [at iter. **i**]

$$C_{MAP} = \arg \max_{C_j \in \{C_1, C_2\}} P(c_j) \prod_{i=1}^n P(x_i | c_j)$$

Bayes Rule

Graphical Models - Microsoft Internet Explorer

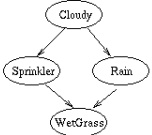
Address: C:\My Documents\Bayesian_summary.html

example, in which all nodes are binary, i.e., have two possible values, which we will denote by T (true) and F (false).

BN

$$\frac{P(C=F) \quad P(C=T)}{0.5 \quad 0.5}$$

C	P(S=F)	P(S=T)
F	0.8	0.2
T	0.2	0.8



C	P(R=F)	P(R=T)
F	0.5	0.5
T	0.9	0.1

S	R	P(W=F)	P(W=T)
F	F	1.0	0.0
T	F	0.1	0.9
F	T	0.1	0.9
T	T	0.01	0.99

We see that the event "grass is wet" ($W=true$) has two possible causes: either the water sprinkler is on ($S=true$) or it is raining ($R=true$). The strength of this relationship is shown in the table. For example, we see that $Pr(W=true | S=true, R=false) = 0.9$ (second row), and hence, $Pr(W=false | S=true, R=false) = 1 - 0.9 = 0.1$, since each row must sum to one. Since the C node has no parents, its CPT specifies the prior probability that it is cloudy (in this case, 0.5).

The simplest conditional independence relationship encoded in a Bayesian network can be stated as follows: a node is independent of its ancestors given its parents.

Optional: not needed for Quiz

71

Yeast Tables for Localization Prediction

Basics	Predictors										Response	Bayesian Localization																		
	Sequence Features					Genomic Features						Localization	State Vector giving localization prediction			Collapsed Prediction														
seq. length	Amino Acid Composition				How many times does the sequence have these motif features?	Abs. expr. Level (mRNA copies / cell)	Cell cycle timecourse			5-compartment	C		N	M	T	E	Training	Extrapolation												
Yeast Gene ID	A	C	D	E	W	Y	Yarn site	NLS	Index motif		nuc2	signalp	times1	Gene-Chip expt. from RY Lab	sage tag freq.	le=0	le=1	le=15	le=16	le=17	le=18	le=19	le=20	C	N	M	T	E	Training	Extrapolation
YAL001C	1160	08	02	06	01	04	0	1	0	1	0	0	0.3	0	5	3	4	5	0	1	N	0%	100%	0%	0%	0%	0%	N		
YAL002W	1176	09	02	06	01	04	0	0	0	0	0	1	0.2?	0	8	4	4	3	0	6	C	95%	3%	2%	0%	0%	0%	C		
YAL003W	206	08	02	06	01	04	0	0	0	0	0	0	19.1	19	70	73	98	126	0	1	N	67%	33%	0%	0%	0%	0%	C		
YAL004W	215	08	02	06	01	04	0	0	0	0	0	0	?	0	18	12	4	6	0	0	N	41%	59%	0%	0%	0%	0%	N		
YAL005C	641	08	02	06	01	04	0	0	0	0	0	0	1	13.4	16	39	38	8	14	0	6	68%	32%	0%	0%	0%	0%		C	
YAL007C	190	08	02	06	01	04	0	0	0	0	0	1	4	2.2	8	15	20	16	17	7	7	26%	43%	31%	0%	0%	0%		-	
YAL008W	198	08	02	06	01	04	0	0	0	0	0	0	3	1.2?	7	9	6	2	3	7	7	37%	60%	3%	0%	0%	0%		-	
YAL009W	259	08	02	06	01	04	0	2	0	0	0	0	3	0.6?	7	6	2	3	5	0	1	2%	98%	0%	0%	0%	0%		N	
YAL010C	493	08	02	06	02	04	0	0	0	0	0	0	1	0.3?	7	11	6	6	6	1	1	6%	90%	4%	0%	0%	0%		N	
YAL011W	616	08	02	06	01	04	0	0	0	0	0	1	0	0.4?	7	6	5	5	6	3	6	28%	62%	10%	0%	0%	0%		N	
YAL012W	393	08	02	06	01	04	0	0	0	0	0	0	1	8.9	4	29	26	23	29	0	0	C	92%	5%	4%	0%	0%	0%	C	
YAL013W	362	08	02	06	01	04	0	0	0	0	0	0	0	0.6?	7	7	9	6	10	0	1	N	0%	98%	0%	0%	0%	1%	N	
YAL014C	202	08	02	06	01	04	0	0	0	0	0	0	0	1.1?	7	12	13	9	12	7	7	N	1%	96%	4%	0%	0%	0%	N	
YAL015C	399	08	02	06	01	04	0	0	0	0	0	0	0	0.7	0	19	18	12	13	1	1	N	4%	96%	0%	0%	0%	0%	N	
YAL016W	635	08	02	06	01	04	0	0	0	0	0	0	1	3.3	5	15	20	16	16	0	1	74%	26%	0%	0%	0%	0%		C	
YAL017W	1356	08	02	06	01	04	0	0	0	0	0	0	0	0.4?	7	14	3	4	7	7	7	0%	1%	99%	0%	0%	0%		M	
YAL018C	325	08	02	06	01	04	0	0	0	0	0	0	4?	?	4	2	4	2	2	1	7	0%	100%	0%	0%	0%	0%		N	

Large-scale Datamining

- **Gene Expression**
 - ◇ Representing Data in a Grid
 - ◇ Description of function prediction in abstract context
- **Unsupervised Learning**
 - ◇ clustering & k-means
 - ◇ Local clustering
- **Supervised Learning**
 - ◇ Discriminants & Decision Tree
 - ◇ Bayesian Nets
- **Function Prediction EX**
 - ◇ Simple Bayesian Approach for Localization Prediction