



Computational Biology Lecture #9: Transcriptomics

Bud Mishra
Professor of Computer Science, Mathematics, & Cell Biology
Nov 14 2005

11/14/2005

© Bud Mishra, 2005

L7-1



Interrupted Genes

11/14/2005

© Bud Mishra, 2005

L7-2



Interrupted Genes:

- ◇ An open reading frame (containing a gene) consists of
 - INTRONS: Intervening sequences a Noncoding regions
 - EXONS: Protein coding regions
- ◇ Introns are abundant in eukaryotes and certain animal viruses.

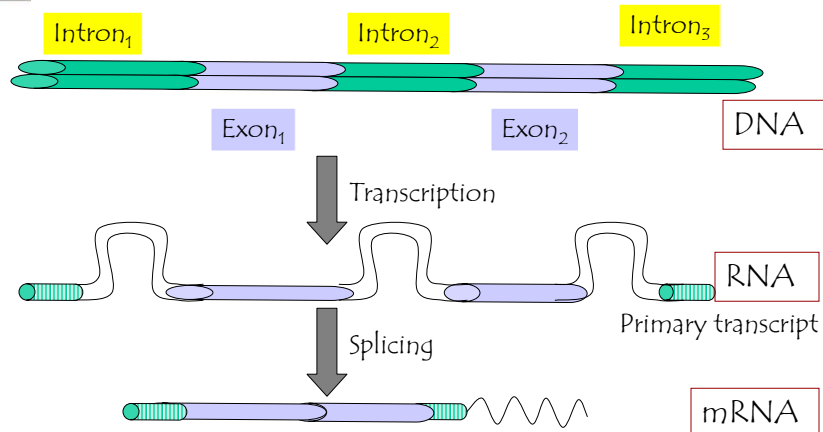
11/14/2005

© Bud Mishra, 2005

L7-3



Interrupted Genes:



11/14/2005

© Bud Mishra, 2005

L7-4



Some Genes...

Gene Product	Organism	Exon Length	#Introns	Intron Length
Adenosine deaminase	Human	1500	11	30,000
Apolipoprotein B	Human	14,000	28	29,000
Erythropoietin	Human	582	4	1562
Thyroglobulin	Human	8500	= 40	100,000
α -interferon	Human	600	0	0
Fibroin	Silk Worm	18,000	1	970
Phaseolin	French Bean	1263	5	515

11/14/2005

© Bud Mishra, 2005

L7-5



Organization of Genetic Information

◇ Bacterial Genome:

- Genes are closely spaced along the DNA.
- The sequences of genes may overlap.
- Related genes (encoding enzymes whose functions are part of the same pathway or whose activities are related) are linked as a single transcription unit.

11/14/2005

© Bud Mishra, 2005

L7-6



Organization of Genetic Information

◇ Eukaryotic Genome:

- Genes are separated by long stretches of noncoding DNA sequences.
- Multiple genes in a single transcription unit is extremely rare.
- Multiple chromosomes \mapsto Linear
- Chloroplasts and mitochondria \mapsto Circular
- Genes appearing on the same chromosome are syntenic.

11/14/2005

© Bud Mishra, 2005

L7-7



Location of Some Genes on Human Chromosome.

Genes	chromosomes	Genes	chromosomes
α -globin cluster	16	Insulin	11
β -globin cluster	11	Galactokinase	11
Immunoglobulin		Viral oncogene homologues	
κ (light chain)	2	C-sis	22
λ (light chain)	22	C-mos	8
Heavy Chain	14	C-Ha-Ras-1	11
Pseudogenes	9,32,15,18	C-myb	6
Growth Hormone gene cluster	17	Interferons	
Thymidine kinase	17	α & β cluster	9
		γ	12

11/14/2005

© Bud Mishra, 2005

L7-8



Eukaryotic Genome

- ◇ **Multiple copies of the same gene**
 - Solve "supply problem"
 - There are several hundred ribosomal RNA genes in mammals
- ◇ **Pseudogenes**
 - Nonfunctional copies of genes...(Deletions or alterations in the DNA sequence)
 - Number of pseudo genes for a particular gene varies greatly...Different from one organism to another.

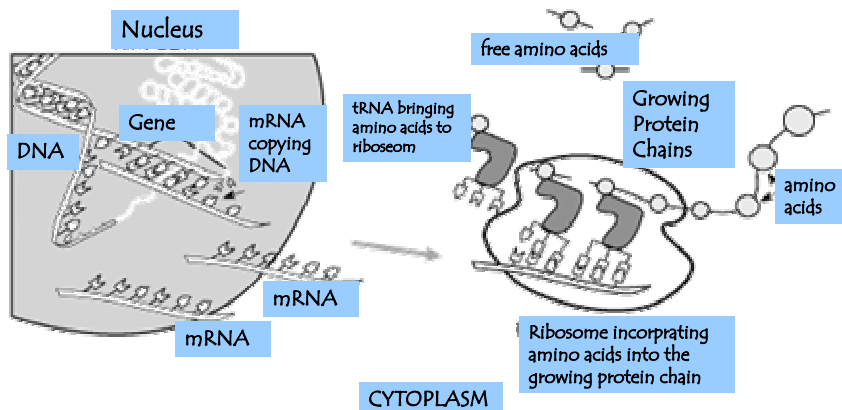
11/14/2005

© Bud Mishra, 2005

L7-9



Gene Expression



11/14/2005

© Bud Mishra, 2005

L7-10



Transcription

- ◇ A gene consists of a coding region and a regulatory region.
 - The coding region is the part that is transcribed into an mRNA and translated into a finished protein.
 - The regulatory region is the part of the DNA that contributes to the control of the gene.
- ◇ The regulatory region contains
 - Binding sites for transcription factors (TF), which act by binding to the DNA (directly or with other transcription factors in a small complex);

11/14/2005

© Bud Mishra, 2005

L7-11



Regulatory Regions

- ◇ In prokaryotes,
 - The regulatory region is short (10-100 bases) and contains binding sites for small number of TFs.
- ◇ In eukaryotes
 - The regulatory regions can be very long (up to 10,000 or 100,000 bases), and contains binding sites for multiple TFs.
 - TFs may act positively or negatively.
 - Another input mechanism is phosphorylation or dephosphorylation of a bound TF by other proteins.

11/14/2005

© Bud Mishra, 2005

L7-12



Terminology

- ◇ TFs are sometimes called *trans-regulatory elements*, and DNA sites where TFs bind are called *cis-regulatory elements*.
 - The collection of cis-regulatory elements upstream of the coding region is called the **promoter**.

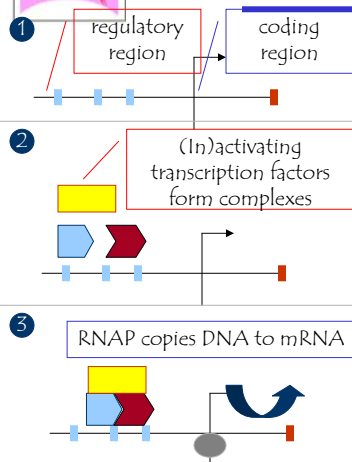
11/14/2005

© Bud Mishra, 2005

L7-13



Transcription Initiation



11/14/2005

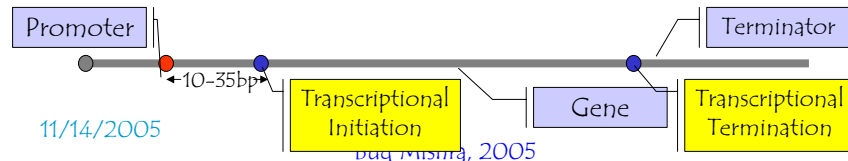
© Bud Mishra, 2005

L7-14



Regulation of Gene Expression

- ◇ **Motifs (short DNA sequences) that regulate transcription:**
 - Promoter & Terminator
 - The rate of transcription varies according to experimental conditions
- ◇ **Motifs that modulate transcription**
 - Repressor, Activator, Antiterminator

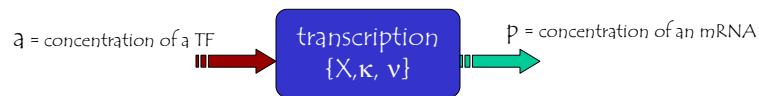


11/14/2005

© Bud Mishra, 2005



Model of transcription



- ν = Cooperativity coefficient
- κ = Concentration of a at which transcription of m is "half-maximally" activated.

- ◇ $dp/dt = \Phi(a, \kappa, \nu) = V a^\nu / [\kappa^\nu + a^\nu]$
- ◇ A graph of function Φ = Sigmoid Function
- ◇ If $\nu = 1$ then, the transcription activation function resembles the classical Michaelis-Menten!

11/14/2005

© Bud Mishra, 2005

L7-16



Regulatory Networks

11/14/2005

© Bud Mishra, 2005

L7-17



Regulatory Networks

◇ **Variations among Cells:**

- All cells in an organism have the same genomic data, but the proteins synthesized in each vary according to cell type, time and environmental factors
- There are network of interactions among various biochemical entities in a cell (DNA RNA, protein, small molecules)

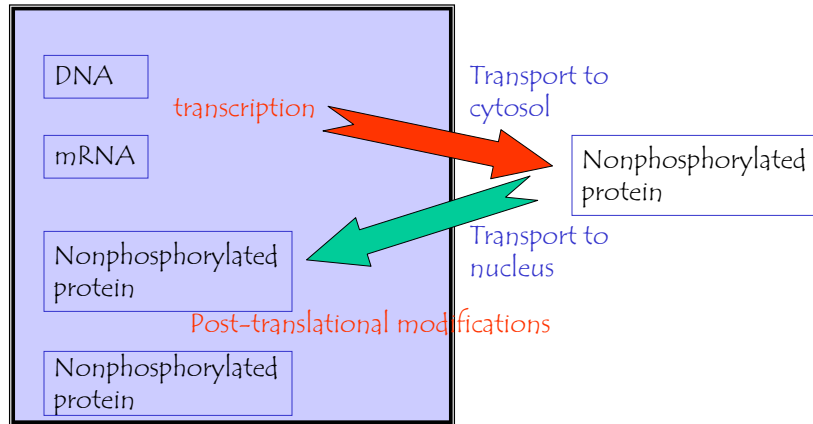
11/14/2005

© Bud Mishra, 2005

L7-18



Gene Regulation



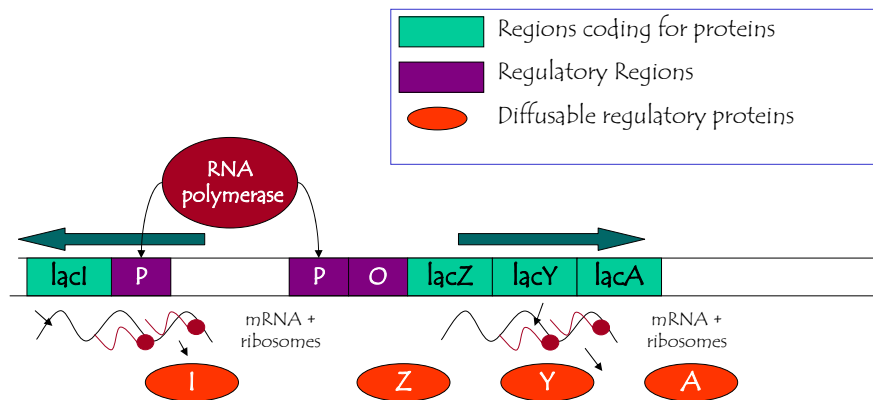
11/14/2005

© Bud Mishra, 2005

L7-19



Transcriptional Regulation: Example: The lac Operon



11/14/2005

© Bud Mishra, 2005

L7-20



The lac Operon

- ◇ **Regulates utilization of lactose by the bacterium *E. coli*.**
 - Lactose is not generally available to *E. coli* as a food substrate, so the bacterium does not usually synthesize the enzymes necessary for its metabolic use.
- ◇ **There is an operon, called the lac operon, normally turned off, that codes for three enzymes:**
 - β -galactoside permease, β -galactosidase and β -thiogalactoside acetyl transferase.

11/14/2005

© Bud Mishra, 2005

L7-21



Activation of the lac operon

- ◇ **An autocatalytic reaction..**
 - If the bacterium is exposed to lactose, these enzymes work together to (1) transport lactose into the cell and (2) isomerizes lactose into allolactose (an allosteric isomer of lactose).
 - ❖ The allolactose binds with a repressor molecule to keep it from repressing the production of mRNA.
 - ❖ Production of allolactose turns on the production of mRNA, which then leads to production of more enzyme, enabling production of more lactose to allolactose...

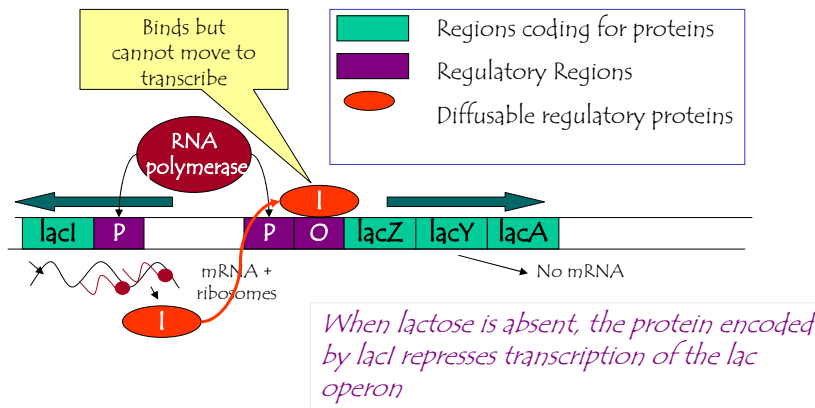
11/14/2005

© Bud Mishra, 2005

L7-22



Transcriptional Regulation: Example: The lac Operon



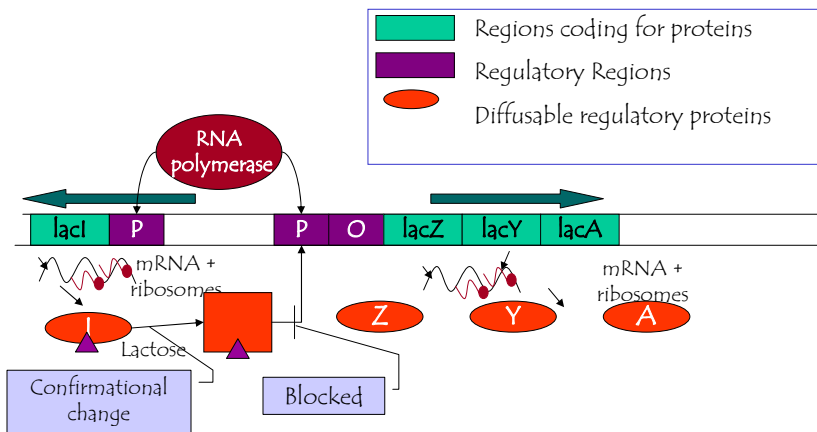
11/14/2005

© Bud Mishra, 2005

L7-23



Transcriptional Regulation: Example: The lac Operon



11/14/2005

© Bud Mishra, 2005

L7-24



Mathematical Model



- ◇ Production of enzyme is turned on by m molecules of the product allolactose P ...

- G =Inactive state of the gene
- X =Active state of the gene
- In a large population of genes, the percentage of active genes is given by the chemical equilibrium:

$$p = [P]^m / (k_{eq}^m + [P]^m)$$

11/14/2005

© Bud Mishra, 2005

L7-25



Production of mRNA

- ◇ The differential equation governing the (average) production of mRNA

$$dM/dt = M_0 + k_1 [P]^m / (k_{eq}^m + [P]^m) - k_2 M,$$

- ◇ where M is the concentration of mRNA that codes for the enzyme.

- Production of the enzymes (responsible for transforming into allolactose substrate):

$$dE_1/dt = c_1 M - d_1 E_1;$$

$$dE_2/dt = c_2 M - d_2 E_2.$$

11/14/2005

© Bud Mishra, 2005

L7-26



Lactose states

- S_0 = Concentration of the lactose, exterior to the cell.
- S = Concentration of the lactose interior to the cell.
- ◇ **[P] = Concentration of allolactose.**

$$dS_0/dt = -\sigma_0 E_1 S_0 / (k_0 + S_0)$$

$$dS/dt = \sigma_0 E_1 S_0 / (k_0 + S_0) - \sigma_1 E_2 S / (k_s + S)$$

$$d[P]/dt = \sigma_1 E_2 S / (k_s + S) - \sigma_2 E_2 [P] / (k_p + [P])$$

11/14/2005

© Bud Mishra, 2005

L7-27



Simplification

- ◇ **Assume:**

- mRNA is in quasi-steady state:

$$M = (k_1/k_2) [P]^m / (k_{eq}^m + [P]^m) + M_0/k_2;$$

- $d_1 = d_2$. Degradation is slow compared to cell growth.
Also, $E_1 = E_2$.

$$dE_1/dt = c_1 M_0/k_2 + (c_1 k_1/k_2) [P]^m / (k_{eq}^m + [P]^m) - d_1 E_1;$$

- No delay in conversion from lactose to allolactose:

$$d[P]/dt = \sigma_0 E_1 S_0 / (k_0 + S_0) - \sigma_2 E_1 [P] / (k_p + [P]).$$

11/14/2005

© Bud Mishra, 2005

L7-28



Dimensionless Form

◇ Dimensionless variables:

- $S_0 = k_0 s$, $[P] = k_p p$, $E_1 = e_0 e$, and $t = t_0 \tau$...

$$de/d\tau = m_0 + p^m/(\kappa^m + p^m) - \varepsilon e,$$

$$dp/d\tau = \mu e[s/(s+1) - \lambda p/(p+1)],$$

$$ds/d\tau = -e s/(s+1),$$

- where $e_0^2 = c_1 k_0 k_1 / (\sigma_0 k_2)$, $t_0 = k + O / (e_0 \sigma_0)$,

$$\lambda = \sigma_2 / \sigma_0, \mu = k_0 / k_p, \kappa = k / k_p, m_0 = M_0 / k_1,$$

$$\text{and } \varepsilon = t_0 d_1 \dots$$

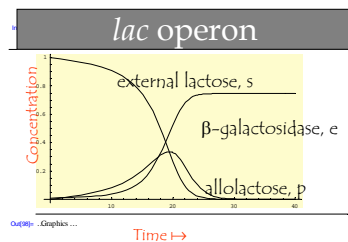
11/14/2005

© Bud Mishra, 2005

L7-29



The *lac* operon



- ◇ If the amount of lactose is too small, then the lactose is gradually depleted, although there is no increase in enzyme concentration.
- ◇ However, if the lactose dose is sufficiently large, then there is an autocatalytic response, as the lac operon is turned on and enzyme is produced.
- ◇ The production of enzyme shuts down when the lactose stimulus is consumed, and the enzyme concentration gradually declines...

11/14/2005

© Bud Mishra, 2005

L7-30



Example of Competition

- ◇ **The mutant Lac repressor X186:**
 - This mutant represses transcription of the *lac* genes in the presence of lactose...
 - The mutant binds DNA so tightly that, in the absence of inducer (allolactose), it is sequestered on non-operator DNA sites.
 - The inducer weakens the binding of the mutant repressor; thus, allowing it bind to the *lac* operon.

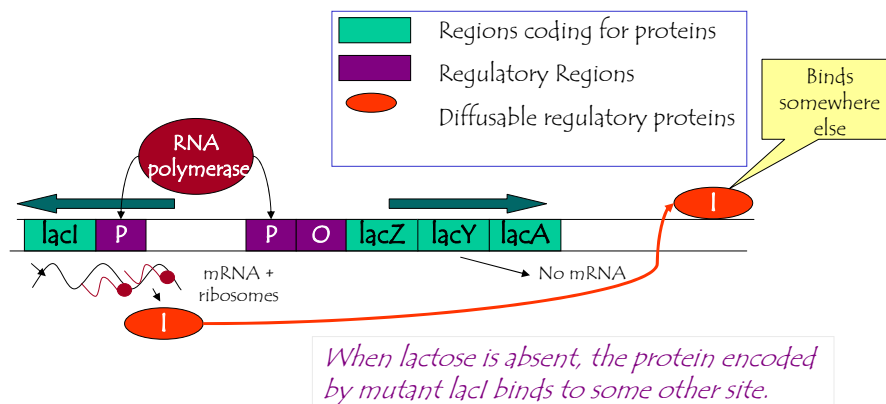
11/14/2005

© Bud Mishra, 2005

L7-31



Lac repressor X 186



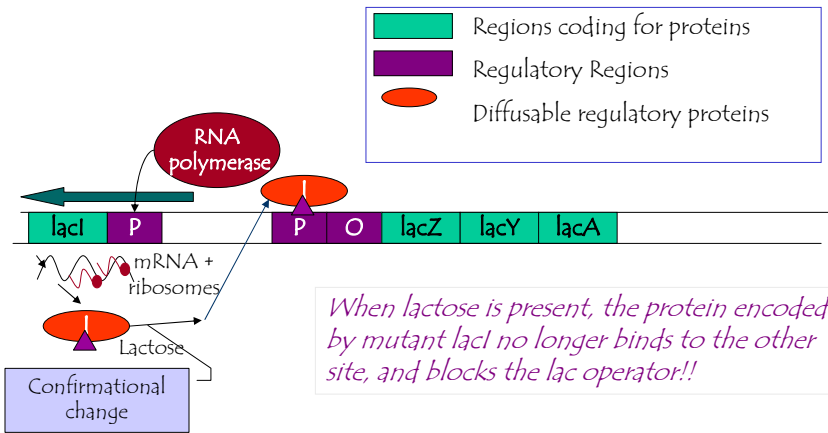
11/14/2005

© Bud Mishra, 2005

L7-32



Lac repressor X 186



11/14/2005

© Bud Mishra, 2005

L7-33



S-Systems

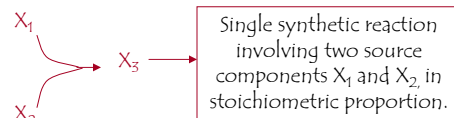
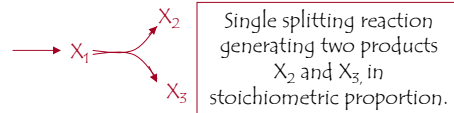
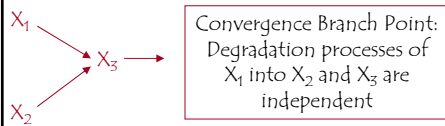
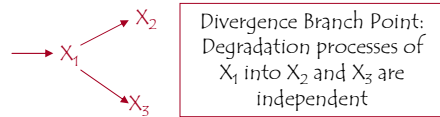
11/14/2005

© Bud Mishra, 2005

L7-34



Graphical Representation



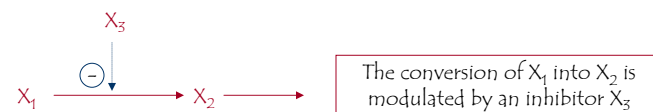
11/14/2005

© Bud Mishra, 2005

L7-35



Graphical Representation



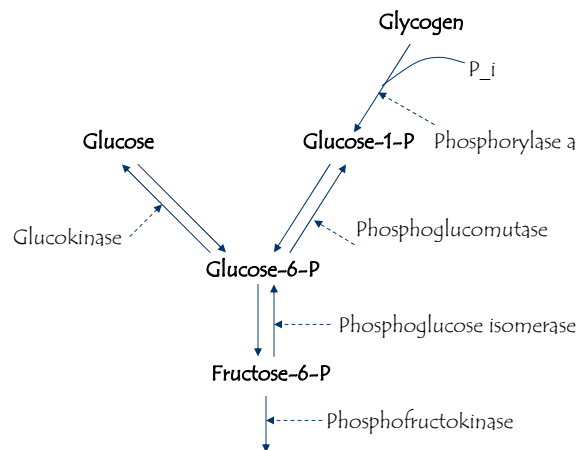
11/14/2005

© Bud Mishra, 2005

L7-36



Glycolysis



11/14/2005

© Bud Mishra, 2005

L7-37



S-Systems

- ◇ **Dependent Variables:** $X_i(t)$, $i=1, \dots, n$, O & t .
- ◇ **System is described in terms of the temporal changes in dependent variables:**
 - E.g., Instantaneous product formation in response to changes in the exogenous substrate, inhibitor or enzyme concentration...
 - Kinetic Laws: Relate a reaction rate to concentrations.
 - Reaction Rate = Instantaneous temporal rate of change in concentration of substrate or product.
- ◇ **Is this information sufficient to deduce the dynamics of a biochemical system? Yes.**

11/14/2005

© Bud Mishra, 2005

L7-38



Systems of Differential Equations

- dX_i/dt = (instantaneous) rate of change in X_i at time t = Function of substrate concentrations, enzymes, factors and products:

$$dX_i/dt = f(S_1, S_2, \dots, E_1, E_2, \dots, F_1, F_2, \dots, P_1, P_2, \dots)$$

- E.g. Michaelis-Menten for substrate S & product P :

1. $dS/dt = -V_{\max} S/(K_M + S)$

2. $dP/dt = V_{\max} S/(K_M + S)$

11/14/2005

© Bud Mishra, 2005

L7-39



Cell Informatics

11/14/2005

© Bud Mishra, 2005

L7-40



The dynamics of cell:

- ◊ The cell cycle) the set of events that occur within a cell between its birth by mitosis and its division into daughter cells again by mitosis
 - interphase period when DNA is synthesized and
 - mitotic phase
 - ❖ The cell division by mitosis (into 2 daughter cells) and meiosis (into 4 gametes from germ-line cells);

11/14/2005

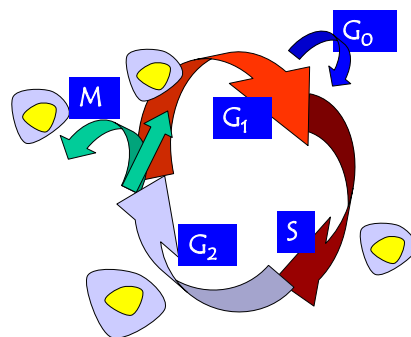
© Bud Mishra, 2005

L7-41



The Cell Cycle:

- ◊ In growing cells, the four phases proceed successively, taking from 10-20 hrs.
- ◊ Interphase: comprises the G_1 , S , and G_2 phases. DNA is synthesized in S and other cellular macromolecules are synthesized throughout interphase, roughly doubling cell's mass.
- ◊ During G_2 the cell is prepared for mitotic (M) phase when the genetic material is evenly proportioned and the cell divides.
- ◊ Nondividing cells exit the normal cycle, entering the quiescent G_0 state.



11/14/2005

© Bud Mishra, 2005

L7-42



Differentiation

- ◇ Cellular dynamics controls how a cell changes (or differentiates) to carry out a specialized functions
 - Structural or morphological changes (muscles, neural, skin..)
 - Immune systems: Many cell types come together in organized tissues designed to let the body distinguish self from non-self.

11/14/2005

© Bud Mishra, 2005

L7-43



Suicide

- ◇ Programmed Cell Death/Apoptosis:
 - Condensation of the nucleus.
 - Fragmentation of the DNA.
 - Morphological changes followed by consumption by macrophages.

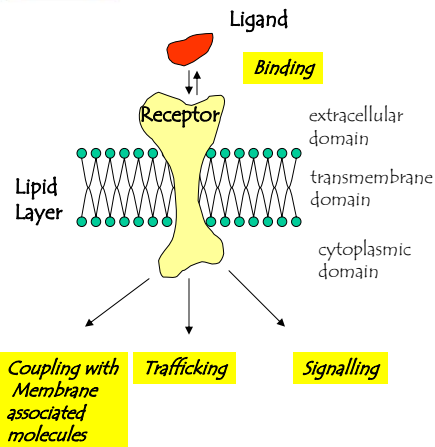
11/14/2005

© Bud Mishra, 2005

L7-44



Cell Talk



- ◇ **Cell Surface Receptors**
 - Extracellular domain for binding ligands (e.g., growth factors, adhesion molecules, etc.)
 - Transmembrane domain
 - Intracellular cytoplasmic domain
- ◇ **Receptor driven cellular behavior are extremely important**
 - E.g., Growth, Secretion, Contraction, Motility and Adhesion

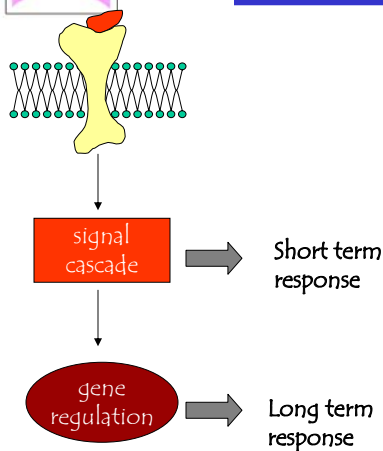
11/14/2005

© Bud Mishra, 2005

L7-45



Receptors and Gene Regulation



- ◇ **Ligands bind to receptors at the cell surface.**
- ◇ **Bound receptors activate various intracellular enzymes and initiate entire cascades of intracellular reactions**
 - Some of these regions trigger short term (of the order of milliseconds to minutes) responses.
 - Some eventually trigger long-term responses..e.g., requiring protein synthesis and additional molecular interactions

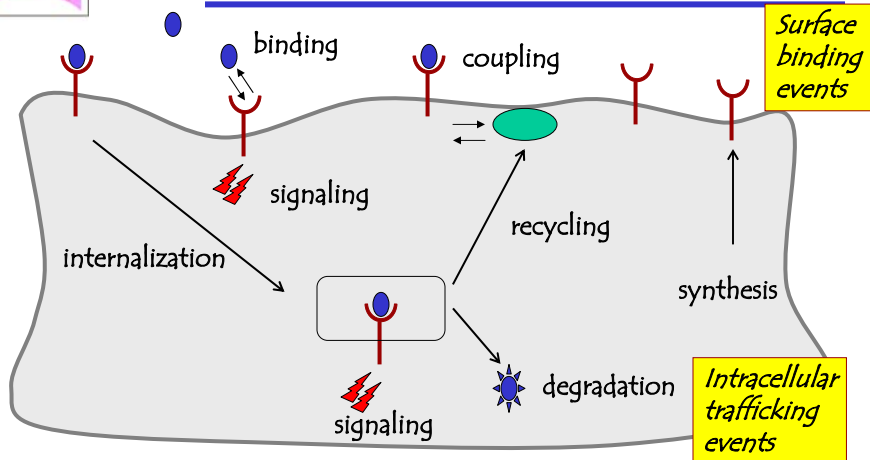
11/14/2005

© Bud Mishra, 2005

L7-46



A Complex Picture



11/14/2005

© Bud Mishra, 2005

L7-47



A Complex Picture

◇ Trafficking

- Receptor population undergoes many complex events of coupling with other cell surface molecules
- Internalization (RME: receptor-mediated endocytosis)
- Recycling
- Degradation
- Synthesis

11/14/2005

© Bud Mishra, 2005

L7-48



Gene Expression Data

- ◇ **Microarrays enable one**
 - To simultaneously measure the activity of up to 30,000 ($\gg 10^4$ — 10^5) genes.
 - In particular, the amount of mRNA for each gene in a given sample (or a pair of samples) can be measured.

11/14/2005

© Bud Mishra, 2005

L7-49



Gene Expression Data

- ◇ **Microarrays provide a tool for answering a wide variety of questions:**
 - In which cells is each gene active?
 - Under what environmental conditions is each gene active?
 - How does the activity level of a gene change under different conditions?
 - ❖ Stage of a cell cycle? Environmental conditions? Diseases?

11/14/2005

© Bud Mishra, 2005

L7-50



Gene Expression Data

- ◇ **Functional genomics with microarrays:**
 - What genes can be inferred to be regulated together?
 - What happens to the expression level of every gene when a (candidate) gene is mutated?
 - What can be inferred about the regulatory structure?

11/14/2005

© Bud Mishra, 2005

L7-51



The Computational Tasks

- ◇ **Clustering Genes:**
 - Genes co-regulated together
- ◇ **Classifying Genes:**
 - Functional class a particular gene fall into?
- ◇ **Classifying Gene Expressions:**
 - Disease classification from the set of all mRNA expressed in a cell
- ◇ **Inferring Regulatory Networks:**
 - The "circuitry" of the cell

11/14/2005

© Bud Mishra, 2005

L7-52



Microarrays

- ◇ Two general types currently popular...
 - Spotted Arrays (Pat Brown, Stanford)
 - Oligonucleotide Arrays (Affymetrix)
 - Other variations (Agilent, Incyte, NGS, ...)
- ◇ The key idea is to query a genome for a particular pattern by complementary hybridization.

11/14/2005

© Bud Mishra, 2005

L7-53



Complementary Hybridization

AGCGTTCGAATACC
UCGCAAGCUUAUGG
ATCGGTACGTTAACG
CCGAAAATAGCCAG

mRNA only hybridizes here

- Due to Watson-Crick base pairing, an mRNA molecule will hybridize to a complementary DNA molecule.

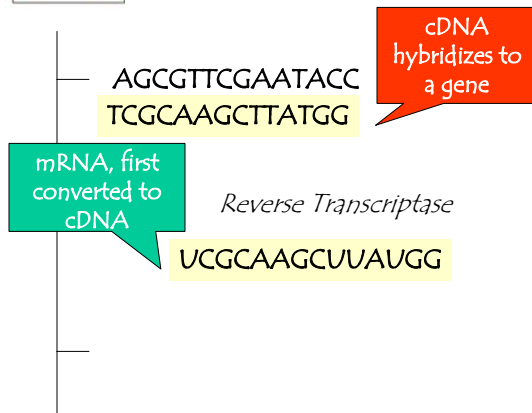
11/14/2005

© Bud Mishra, 2005

L7-54



Complementary Hybridization



Practical implementation:

- Put the actual gene sequence on array
- Convert mRNA to cDNA using reverse transcriptase
- Hybridize cDNA to the sequence on the array

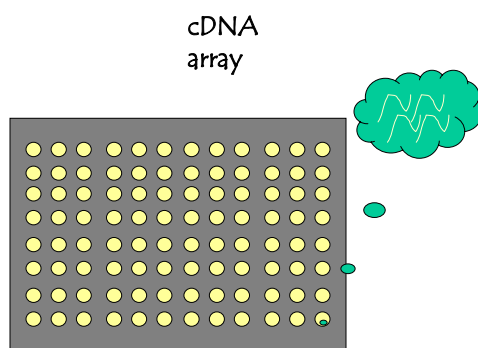
11/14/2005

© Bud Mishra, 2005

L7-55



Spotted Array



Robots array microscopic sized spots of DNA on glass slides

- Each spot is DNA analog (cDNA) of one of the mRNA's we wish to measure...

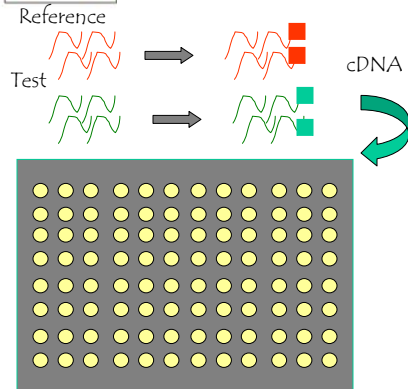
11/14/2005

© Bud Mishra, 2005

L7-56



Spotted Arrays



- Two samples (reference and test) of mRNA are converted to cDNA, labeled with fluorochrome dyes and allowed to hybridize to the array.

11/14/2005

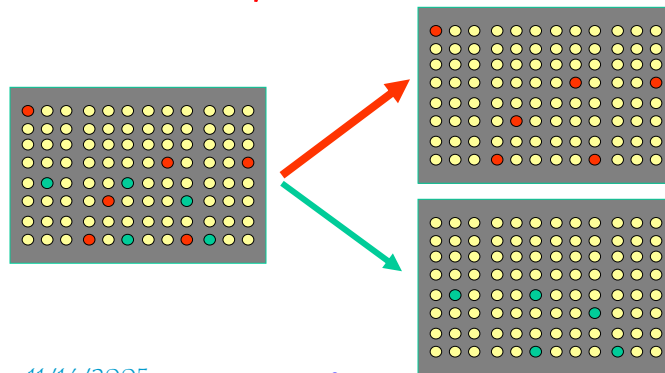
© Bud Mishra, 2005

L7-57



Spotted Arrays

- Lasers applied to the arrays yield an emission for each fluorescent dye.



11/14/2005

© Bud Mishra, 2005

L7-58



Oligonucleotide Arrays

◇ "Gene Chips"

- Instead of putting entire genes on array, put sets of DNA 25-mers (synthesized oligonucleotides)
- Produced using a photolithography process similar to the ones used to create semiconductor chips
- mRNA samples are processed separately instead of in pairs (of reference/control and test)

11/14/2005

© Bud Mishra, 2005

L7-59



Oligonucleotide Arrays



- ◇ Given a gene to be queried/measured, select a large number (» 20) 25-mers for that gene.
- ◇ Selection criteria
 - Specificity
 - Hybridization properties
 - Ease of manufacturing

11/14/2005

© Bud Mishra, 2005

L7-60



Oligonucleotide Arrays

- ◇ Each of these probes is put on the chip
 - Additionally a slight variant (that differs only at the 13th base) of each oligo is put next to it.
 - This helps factor out false hybridization (pm [perfect match] vs. mm [mismatch])
- ◇ The measurement for a gene is derived from these 40 separate measurements.

11/14/2005

© Bud Mishra, 2005

L7-61



Genome-wide Cluster Analysis

- ◇ Put all genes (» 6200) of *S. cerevisiae* (yeast) on a single microarray
- ◇ Measure experiment across m independent experiments
- ◇ Group together genes that have similar expression profiles.
 - Eisen et al. PNAS 1998

11/14/2005

© Bud Mishra, 2005

L7-62



Genome-wide Cluster Analysis

- ◇ Each measurement G_i represents

$$\log(\text{red}_i/\text{green}_i)$$

Where red is the test expression level and green is the reference expression level for gene G in the i^{th} experiment.

- ◇ The expression profile of a gene is the vector of measurements across all experiments:

$$\mathbf{h}_{G_1, \dots, G_m} \mathbf{i}$$

11/14/2005

© Bud Mishra, 2005

L7-63



The Data

- ◇ 79 measurements for each of 2467 genes
- ◇ Data collected at various times during
 - Diauxic shift (shutting down genes for metabolizing sugar, activating genes for metabolizing ethanol)
 - Mitotic cell division cycle
 - Sporulation
 - Temperature shock
 - Reducing Shock

11/14/2005

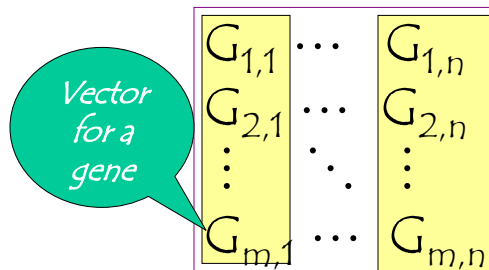
© Bud Mishra, 2005

L7-64



The Data

- ◇ n genes measured in m experiments:



11/14/2005

© Bud Mishra, 2005

L7-65



The Task

- ◇ **Given**
 - Expression profiles for a set of genes.
- ◇ **Compute**
 - An organization of genes into clusters such that genes within a cluster have similar profiles.

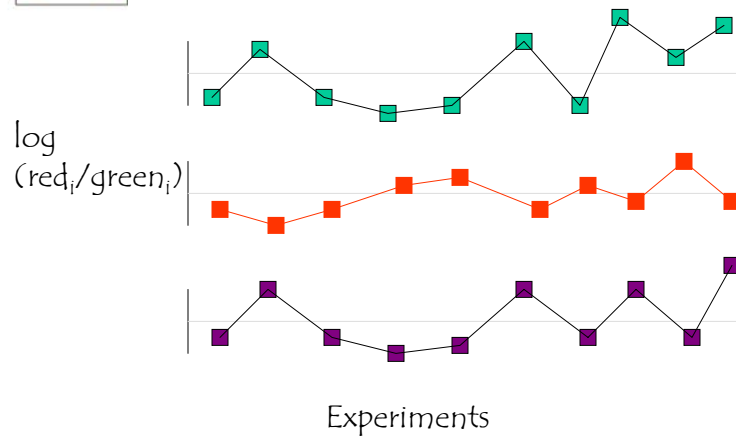
11/14/2005

© Bud Mishra, 2005

L7-66



The Task



11/14/2005

© Bud Mishra, 2005

L7-67



Approaches

- ◇ Eisen et al.: Hierarchical clustering.
- ◇ Other clustering methods have been applied to this gene expression data:
 - EM with Gaussian Clusters [Mjolsness et al. '99]
 - Self Organizing Maps [Tamayo et al. '99]
 - Graph Theory Algorithms [Ben-Dor & Yakhini '98, Hartuv et al. '99]

11/14/2005

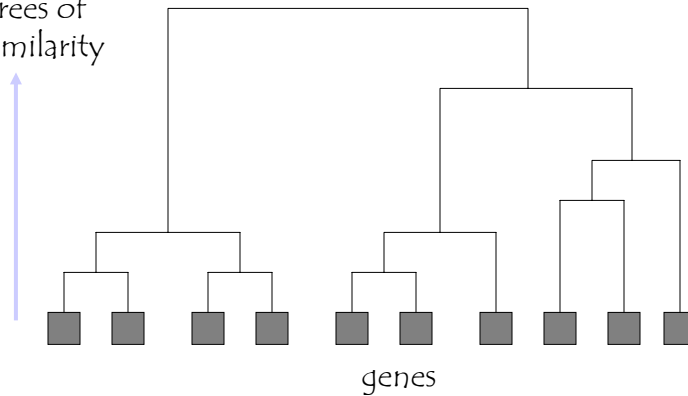
© Bud Mishra, 2005

L7-68



Hierarchical Clustering

Degrees of
dissimilarity



11/14/2005

© Bud Mishra, 2005

L7-69



Hierarchical Clustering

- ◇ P = set of genes
- ◇ While more than one subtree in P
 - Pick the most similar pair i, j in P
 - Define a new subtree k joining i and j
 - Remove i and j from P and insert k

11/14/2005

© Bud Mishra, 2005

L7-70



Gene Similarity Metric

- ◇ Similarity between two genes: X and Y

$S(X, Y) =$

$$(1/N) \sum_{i=1}^N (X_i - X_{\text{offset}} / \Phi_X) (Y_i - Y_{\text{offset}} / \Phi_Y)$$

where

$$\Phi_G = [\sum_{i=1}^N (G_i - G_{\text{offset}})^2 / N]^{1/2}$$

11/14/2005

© Bud Mishra, 2005

L7-71



Gene Similarity Metric

- Since there is an assumed reference state (the gene's expression level did not change), G_{offset} is set to 0 for all genes
- ◇ $S(X, Y)$
 - $= (1/N) \sum_{i=1}^N [X_i / \{ \sum_{i=1}^N X_i^2 / N \}^{1/2}] [Y_i / \{ \sum_{i=1}^N Y_i^2 / N \}^{1/2}]$
 - $= (1/N) \{ \sum_{i=1}^N X_i Y_i \} / \{ SD(X) SD(Y) \}$
 - $= (1/N) \{ X \cdot Y \} / \{ SD(X) SD(Y) \}$

11/14/2005

© Bud Mishra, 2005

L7-72



Results

- ◇ **Redundant representations of genes cluster together.**
 - But individual genes can be distinguished from related genes by subtle differences in expression.

11/14/2005

© Bud Mishra, 2005

L7-73



Results

- ◇ **Genes of similar function cluster together.**
 - E.g., 126 genes were found strongly down-regulated in response to stress.
 - ❖ 112 of these genes encode ribosomal and other proteins related to translation
 - ❖ The result agrees with previously known result that yeast responds to favorable growth conditions by increasing the production of ribosomes.

11/14/2005

© Bud Mishra, 2005

L7-74



Molecular Classification of Cancer

- ◇ Measure activity of 6817 genes in 38 leukemia patients
- ◇ Two tasks:
 - ❖ Class Prediction
 - ❖ Class Discovery
 - Golub et al., Science '99.
 - Slonim et al. '99.

11/14/2005

© Bud Mishra, 2005

L7-75



Cancer Class Prediction

- ◇ Learning Task
 - Given: Expression profiles of leukemia patients
 - Compute: A model distinguishing disease classes (e.g., AML vs. ALL patients) from expression data.
- ◇ Classification Task
 - Given: Expression profile of a new patient + A learned model (e.g., one computed in a learning task)
 - Determine: The disease class of the patient (e.g., whether the patient has AML or ALL)

11/14/2005

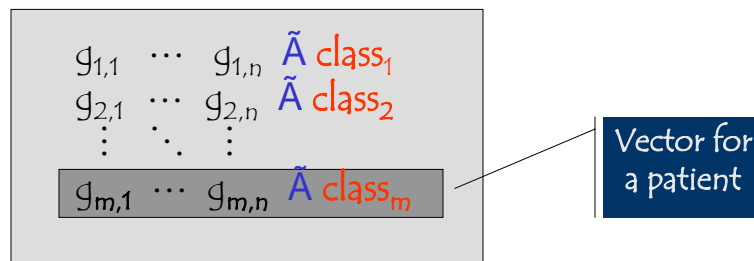
© Bud Mishra, 2005

L7-76



Cancer Class Prediction

◇ n genes measured in m patients



11/14/2005

© Bud Mishra, 2005

L7-77



Cancer Class Prediction Approach

- Rank genes by their correlation with class variable (AML/ALL)
- Select subset of "informative" genes
- Have these genes do a weighted vote to classify a previously unclassified patient.

11/14/2005

© Bud Mishra, 2005

L7-78



Ranking Genes

- ◇ Rank genes by how predictive they are (individually) of the class...

$g_{1,1}$...	$g_{1,n}$	\tilde{A} class ₁
$g_{2,1}$...	$g_{2,n}$	\tilde{A} class ₂
\vdots	\vdots	\vdots	
$g_{m,1}$...	$g_{m,n}$	\tilde{A} class _m

11/14/2005

© Bud Mishra, 2005

L7-79



Ranking Genes

- Split the expression values for a given gene g into two pools – one for each class (AML vs. ALL)
- Determine their mean μ and standard deviation σ of each pool
- ◇ Rank genes by

$$P(g, \text{class}) = (\mu_{\text{ALL}} - \mu_{\text{AML}}) / (\sigma_{\text{ALL}} + \sigma_{\text{AML}})$$

11/14/2005

© Bud Mishra, 2005

L7-80



Selecting Genes

- Select the kALL top ranked genes (highly expressed in ALL) and the kAML bottom ranked genes (highly expressed in AML)

$$P(g, \text{class}) = (m_{\text{ALL}} - m_{\text{AML}}) / (s_{\text{ALL}} + s_{\text{AML}})$$

11/14/2005

© Bud Mishra, 2005

L7-81



Weighted Voting

- ◇ Given a new patient to classify,
 - Each of the selected genes casts a weighted vote for only one class.
 - The class that gets the most vote is the prediction.

11/14/2005

© Bud Mishra, 2005

L7-82



Weighted Voting

- Suppose that x is the expression level measured for gene g in the patient

$$V = P(g, \text{class}) \times (x - [\mu_{\text{ALL}} + \mu_{\text{AML}}]/2)$$

Weight for
gene g

Distance from the
measurement to the
class boundary

11/14/2005

© Bud Mishra, 2005

L7-83



Prediction Strength

- Can assess the "strength" of a prediction as follows:

$$PS = (V_{\text{winner}} - V_{\text{loser}}) / (V_{\text{winner}} + V_{\text{loser}})$$

where V_{winner} is the summed vote from the winning class,
and V_{loser} is the summed vote for the losing class

11/14/2005

© Bud Mishra, 2005

L7-84



Prediction Strength

- ◇ When classifying new cases, the algorithm ignores those cases where the strength of the prediction is below a threshold...
- ◇ Prediction =
 - [ALL, if $V_{ALL} > V_{AML} \wedge PS > \theta$
 - [AML, if $V_{AML} > V_{ALL} \wedge PS > \theta$
 - [No-call, otherwise.

11/14/2005

© Bud Mishra, 2005

L7-85



Experiments

- ◇ Cross validation with the original set of patients
 - For $i = 1$ to 38
 - ❖ Hold the i^{th} gene aside
 - ❖ Use the other 37 genes to determine weights
 - ❖ With this set of weights, make prediction on the i^{th} gene
- ◇ Testing with another set of 34 patients...

11/14/2005

© Bud Mishra, 2005

L7-86



Results

- ◇ **Cross-validation experiments**
 - All trials that used at least 3 genes had 0 prediction error, with 1—4 no-calls.
- ◇ **Using the 50 gene model on a test set of 34 additional patients**
 - 29 correct predictions
 - 5 no-calls.

11/14/2005

© Bud Mishra, 2005

L7-87



Comments on Molecular Classification of Cancer

- ◇ **Gene expression profiling appears to be a promising tool for molecular medicine**
 - Screening
 - Diagnosis
 - Prognosis
 - Highly targeted genome-based therapy

11/14/2005

© Bud Mishra, 2005

L7-88



Cancer Class Discovery

- ◇ Given
 - Expression profiles of leukemia patients
- ◇ Do
 - Cluster the profiles, leading to discovery of the subclasses of leukemia represented by the set of patients

11/14/2005

© Bud Mishra, 2005

L7-89



Cancer Class Discovery Experiment

- ◇ Cluster the expression profiles of 38 patients in the training set
 - Using self-organizing maps with a predefined number of clusters (say, k)
- ◇ Run with $k = 2$
 - Cluster 1 contained 1 AML, 24 ALL
 - Cluster 2 contained 10 AML, 3 ALL

11/14/2005

© Bud Mishra, 2005

L7-90



Cancer Class Discovery Experiment

- ◇ Run with $k = 4$
 - Cluster 1 contained mostly AML
 - Cluster 2 contained mostly T-cell ALL
 - Cluster 3 contained mostly B-cell ALL
 - Cluster 4 contained mostly B-cell ALL
- ◇ It is unlikely that the clustering algorithm was able to discover the distinction between T-cell and B-cell ALL cases

11/14/2005

© Bud Mishra, 2005

L7-91



Comments on Cancer Class Discovery

- ◇ Potential to discover unknown but clinically significant classes
 - One may still be able to take advantage of class labels to guide subclass discovery
 - Room for novel statistical algorithms...

11/14/2005

© Bud Mishra, 2005

L7-92



To be continued...

...

11/14/2005

© Bud Mishra, 2005

L7-93



Hidden Markov Models

HMM

11/14/2005

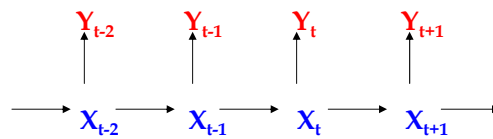
© Bud Mishra, 2005

L7-94



Hidden Markov Model

- ◇ Bayesian network structure for both
 - Hidden Markov Model
- ◇ Important Independence Assumptions:
 - Current state X_t depends only on the past state X_{t-1}
 - Current output Y_t only depends on the state X_t



11/14/2005

© Bud Mishra, 2005

L7-95



Hidden Markov Models (HMM)

- ◇ Defined by an alphabet S ,
 - A set of (hidden) states Q ,
 - A matrix of state transition probabilities A ,
 - and a matrix of emission probabilities E .

11/14/2005

© Bud Mishra, 2005

L7-96



States

- S = An alphabet of symbols
- Q = A set of states that emit symbols from the alphabet S
- $A = (a_{kl}) = |Q| \times |Q|$ matrix of state transition probabilities
- $E = (e_{k(B)}) = |Q| \times |S|$ matrix of emission probabilities

11/14/2005

© Bud Mishra, 2005

L7-97



A Path in the HMM

- ◇ $\pi = \pi_1 \pi_2 \cdots \pi_n$
= a sequence of states $\in Q^*$ in the hidden Markov model M
- $x \in \Sigma^*$ = sequence generated by the path π , determined by the model M
- $P(x | \pi) = P(\pi_1) \prod_{i=1}^{n-1} P(x_i | \pi_i) P(\pi_i | \pi_{i+1})$

11/14/2005

© Bud Mishra, 2005

L7-98



A Path in the HMM

- $P(x|\pi) = [\prod_{i=1}^n P(x_i | \pi_i) P(\pi_i | \pi_{i+1})] P(\pi_1)$
- $P(x_i | \pi_i) = e_{\pi_i}(x_i)$
- $P(\pi_i | \pi_{i+1}) = a_{\pi_i, \pi_{i+1}}$
- π_0 = Initial state "begin"
- π_{n+1} = Final state "end"
- ◇ $P(x|\pi)$
 - = $a_{\pi_0, \pi_1} e_{\pi_1}(x_1) a_{\pi_1, \pi_2} e_{\pi_2}(x_2) \dots e_{\pi_n}(x_n) a_{\pi_n, \pi_{n+1}}$
 - = $a_{\pi_0, \pi_1} \prod_{i=1}^n e_{\pi_i}(x_i) a_{\pi_i, \pi_{i+1}}$

11/14/2005

© Bud Mishra, 2005

L7-99



Decoding Problem

- ◇ For a given sequence x , and a given path π ,
 - The model (Markovian) defines the probability
 - $P(x|\pi)$
 - The dealer knows π and x
 - The player knows x but not π
"The path of x is hidden."
- ◇ Decoding Problem: Find an optimal path π^* for x such that $P(x|\pi)$ is maximized.

$$\pi^* = \arg \max_{\pi} P(x|\pi)$$

11/14/2005

© Bud Mishra, 2005

L7-100



Dynamic Programming Approach

- ◇ Principle of Optimality:
- ◇ **Optimal path for the (i+1)-prefix of x**

$$x_1 \cdots x_{i+1}$$
 - uses a path for an i-prefix of x that is optimal among the paths ending in an (unknown) state $\pi_i = k \in Q$

11/14/2005

© Bud Mishra, 2005

L7-101



Dynamic Programming Approach

- ◇ $s_k(i) =$
 - The probability of the most probable path for the i-prefix ending in state k.

$$s_k(i+1) = e_i(x_{i+1}) \cdot \max_{k' \in Q} [s_{k'}(i) \cdot a_{k'k}]$$

11/14/2005

© Bud Mishra, 2005

L7-102



Dynamic Programming

◇ $i=0$

$$s_{\text{begin}}(0) = 1, s_k(0) = 0, \mathbf{8}_{k \neq \text{begin}}$$

◇ $0 < i \leq n$

$$s_i(i+1) = e_i(x_{i+1}) \cdot \max_{k \in Q} [s_k(i) \cdot a_{kl}]$$

◇ $i = n+1$

$$P(x | \pi^*) = \max_{k \in Q} s_k(n) a_{k, \text{end}}$$

11/14/2005

© Bud Mishra, 2005

L7-103



Viterbi Algorithm

◇ **Dynamic Programming**

- with log-score function

$$S_i(i) = \log s_i(i)$$

- Space complexity = $O(n |Q|)$
- Time complexity = $O(n |Q|)$
- $S_i(i+1) = \log e_i(x_{i+1}) + \max_{k \in Q} [S_k(i) + \log a_{kl}]$

11/14/2005

© Bud Mishra, 2005

L7-104



Estimating the i^{th} State

- ◇ $P(\pi_i = k | x) =$
 - Given a sequence $x \in \Sigma^*$, the probability that the HMM was in state k at instant i .

11/14/2005

© Bud Mishra, 2005

L7-105



Forward Estimate:

- ◇ $f_k(i) = P(x_1 \dots x_i, \pi_i = k) =$
 - Probability of emitting the prefix $x_1 \dots x_i$ and reaching the state $\pi_i = k$
- ◇ $f_k(i) = e_k(x_i) \sum_{l \in Q} f_l(i-1) a_{lk}$

11/14/2005

© Bud Mishra, 2005

L7-106



Backward Estimate:

- ◇ $b_k(i) = P(x_{i+1} \dots x_n, \pi_i = k) =$
 - Probability of being at the state $\pi_i = k$ and emitting the suffix $x_{i+1} \dots x_n$.
- ◇ $b_k(i) = \sum_{l \in Q} e_k(x_{i+1}) \phi b_l(i+1) a_{kl}$

11/14/2005

© Bud Mishra, 2005

L7-107



Applying Bayes' Rule

- ◇ $P(\pi_i = k | x) = (1/P(x))$
 - $P(x_1 \dots x_i, \pi_i = k)$
 - $P(x_{i+1} \dots x_n, \pi_i = k)$
- ◇ $P(x) = \sum_{\pi} P(x | \pi)$

11/14/2005

© Bud Mishra, 2005

L7-108



Sequence Motifs

- ◇ A Sequence of patterns of biological significance.
- ◇ Examples:
 - DNA: Protein binding sites
 - ❖ (e.g. promoters, regulatory sequences)
 - Protein: sequences corresponding to conserved pieces of structure
 - ❖ (e.g. Local features, At various scales: blocks, domains & families)

11/14/2005

© Bud Mishra, 2005

L7-109



MEME Algorithm

- Uses EM (Expectation Minimization) algorithm to find multiple motifs in a set of sequences.
- ◇ Description of a motif:
 - W = (Fixed) width of a motif
 - $P = (p_{lc})_{l \in \Sigma, c \in 1..W}$
= Matrix of probabilities that letter l occurs at position c
= $|\Sigma| \times W$ matrix

11/14/2005

© Bud Mishra, 2005

L7-110



Example

$$P_\rho = \begin{matrix} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} A \\ T \\ C \\ G \end{matrix} & \left\{ \begin{matrix} 0.1 & 0.5 & 0.2 \\ 0.2 & 0.2 & 0.1 \\ 0.4 & 0.2 & 0.1 \\ 0.3 & 0.1 & 0.6 \end{matrix} \right\} \end{matrix}$$

◇ DNA motif of width

- $W = 3$,
- $\Sigma = \{A, T, C, G\}$

◇ $\rho = \text{motif}$) $W_\rho = 3$,

◇ $P_\rho = 4 \times 3$ stochastic matrix

11/14/2005

© Bud Mishra, 2005

L7-111



Computational Problem

◇ Given:

- A set of sequences, G
- A width parameter W

◇ Find:

- Motifs of width W common to sequences G and present their probabilistic descriptions.
- Note that motif start sites in each sequence are unknown (hidden).

11/14/2005

© Bud Mishra, 2005

L7-112



Basic EM Approach

	[1]	[2]	[3]	[4]
$seq1$	0.1	0.1	0.2	0.6
$seq2$	0.4	0.2	0.1	0.3
$seq3$	0.3	0.1	0.5	0.1
$seq4$	0.1	0.5	0.1	0.3

- ◇ $\Gamma =$
 - Training sequences.
- ◇ $|\Gamma| =$
 - Total number of sequences = m ;
 - Minimum length of a sequence in $\Gamma = l$.
- ◇ $Z = m \times l$ matrix of probabilities
 - z_{ij} = Probability that the motif starts in position j in sequence i .

11/14/2005

© Bud Mishra, 2005

L7-113



EM Algorithm

- ◇ Set initial values for P
- ◇ do
 - Re-estimate Z from P
 - Re-estimate P from Z
- ◇ until change in $P < \epsilon$
- ◇ return P

11/14/2005

© Bud Mishra, 2005

L7-114



EM Algorithm

◇ Maximize the likelihood of a motif in the training sequence:

- $S_i \in \Gamma$ i^{th} sequence
- $I_{ij} = \{1, \text{if motif starts at posn. } j \text{ in seq. } i\}$
 $\{0, \text{otherwise.}\}$
- $I_k = \text{the char. at posn. } j+k-1 \text{ in seq. } S_i$

◇ $\Pr(S_i | I_{ij} = 1, \rho) = \prod_{k=1}^W \rho_{\{I_k, k\}}$

11/14/2005

© Bud Mishra, 2005

L7-115



Example

$S_i = \text{AGGCTGTAGACAC}$

$P_\rho =$

	1	2	3
A	0.1	0.5	0.2
T	0.2	0.2	0.1
C	0.4	0.2	0.1
G	0.3	0.1	0.6

◇ $\Pr(S_i = \text{TGT} | I_{i5} = 1, \rho)$
 $= \rho_{T,1} \rho_{G,2} \rho_{T,3}$
 $= 0.2 \times 0.1 \times 0.1$
 $= 2 \times 10^{-3}$

11/14/2005

© Bud Mishra, 2005

L7-116



Estimating Z

- ◇ $z_{ij} = \Pr(I_{ij} = 1 \mid \rho, S_i)$
= Estimates the starting position in S_i Γ .
- ◇ $z_{ij} = \Pr(I_{ij} = 1 \mid \rho, S_i)$
= $\Pr(S_i, I_{ij} = 1 \mid \rho) / \Pr(S_i \mid \rho)$
= $\Pr(S_i \mid I_{ij} = 1, \rho) \Pr(I_{ij} = 1) / \sum_k \Pr(S_i \mid I_{ik} = 1, \rho) \Pr(I_{ik} = 1)$
= $\Pr(S_i \mid I_{ij} = 1, \rho) / \sum_k \Pr(S_i \mid I_{ik} = 1, \rho)$
- ◇ Follows from an application of the Bayes' rule and the assumption that "it is equally likely that the motif will start in any position."
 $\sum_k \Pr(I_{ik} = 1) = \Pr(I_{ij} = 1)$

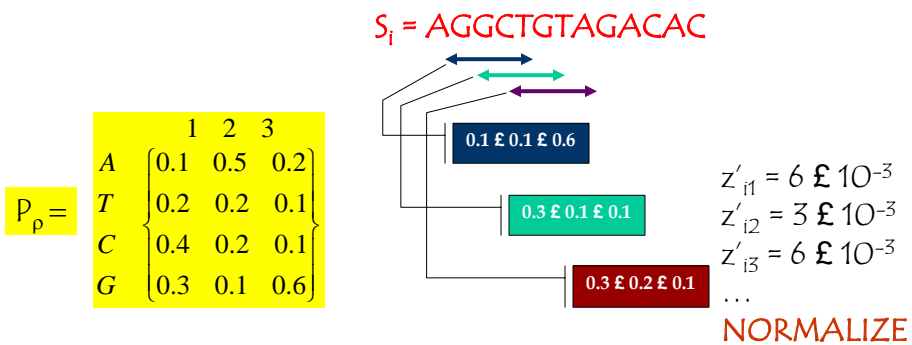
11/14/2005

© Bud Mishra, 2005

L7-117



Example



11/14/2005

© Bud Mishra, 2005

L7-118



Estimating P_ρ

- Given Z , estimate the probability that the character c occurs at the k^{th} position of a motif.
- $n_{ck} = \sum_{s(i)} \sum_{\Gamma; \{j \mid s(i, j+k-1) = c\}} Z_{ij}$
- Expected number of occurrences of the character c at the k^{th} position of a motif ρ (assuming that the motif "start position" is known.)
- $p_{ck} = (n_{ck} + 1) / \sum_d (n_{dk} + 1)$

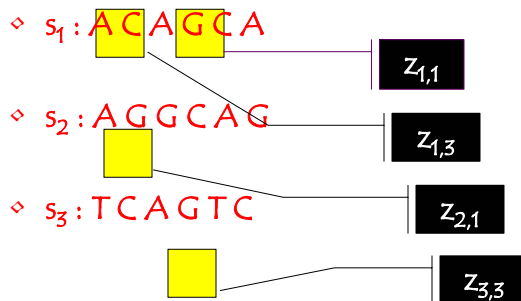
11/14/2005

© Bud Mishra, 2005

L7-119



Example



$$p_{A,1} = \frac{(z_{11} + z_{13} + z_{21} + z_{33} + 1)}{(z_{11} + z_{12} + \dots + z_{33} + z_{34} + 4)}$$

11/14/2005

© Bud Mishra, 2005

L7-120



Meme

- ◇ Uses the basic EM approach
 - Try many starting points.
 - Allow multiple occurrences of a motif per sequence
 - Allows multiple motifs to be learned simultaneously.

11/14/2005

© Bud Mishra, 2005

L7-121



Meme

- ◇ Initial set of possible motifs:
 - Take every distinct subsequences of length W in the training set
 - Derive an initial matrix P
 - $p_{ck} = \begin{cases} \alpha & \text{if } c \text{ occurs in position } k \text{ in} \\ & \text{the subsequence} \\ \{(1-\alpha) / (|\Sigma| - 1) & \text{Otherwise} \end{cases}$

11/14/2005

© Bud Mishra, 2005

L7-122



Example

- ◇ $W = 3,$
- ◇ $\rho = TAT,$
- ◇ $\alpha = 0.5$

$P_\rho =$		1	2	3
A	}	0.17	0.50	0.17
T		0.17	0.17	0.17
C		0.17	0.17	0.17
G		0.50	0.17	0.50

- Choose the motif model with the highest likelihood.
- Run EM to convergence

11/14/2005

© Bud Mishra, 2005

L7-123