

dhp HST

Making sense of microarrays

Isaac S Kohane

dhp HST

RNA expression detection chips

Schena M, et al. Proc Natl Acad Sci USA; 93: 10614 (1996).
Entire issue. Nature Genetics, 21: supplement (Jan 1999).
How I learned to love Bioinformatics...

dhp HST

Two Biochip technologies

Affymetrix

Fig. 2 Gene expression monitoring with oligonucleotide arrays. A single 1.28x1.28 cm array containing probe sets for approximately 40,000 human genes and ESTs. This array contains features smaller than 22x22 μm and only four probe pairs per gene or EST. Expression probe and array design. Oligonucleotide probes are chosen based on uniqueness criteria and composition design rules. For eukaryotic organisms, probes are chosen typically from the 3' end of the gene or transcript (closer to the poly(A) tail) to reduce problems that may arise from the use of partially degraded mRNA. The use of the PM minus MM difference averaged across a set of probes greatly reduces the contribution of background and cross hybridization and increases the quantitative accuracy and reproducibility of the measurements.

P. Brown / Stanford

How I learned to love E

dhp HST

What is a microarray?

- Low cost. The cost should be such that at least hundreds of samples be measurable within a typical NIH RO1 budget
- Commodity level workflow. The microarray should be commoditized such that a routine set of procedures requiring no scientific judgment can be performed using standard equipment to obtain the needed measurement.
- Automation. The process of data acquisition should be completely automated so that after the biomaterial whether it is protein, RNA or DNA is loaded into the analytic pipeline, most of the steps are fully automated and those that are not automated can be done by a non-specialized technician.
- Form factor. The equipment required should easily fit into a standard laboratory bench format and not require its own room.

How I learned to love Bioinformatics...

dhp HST

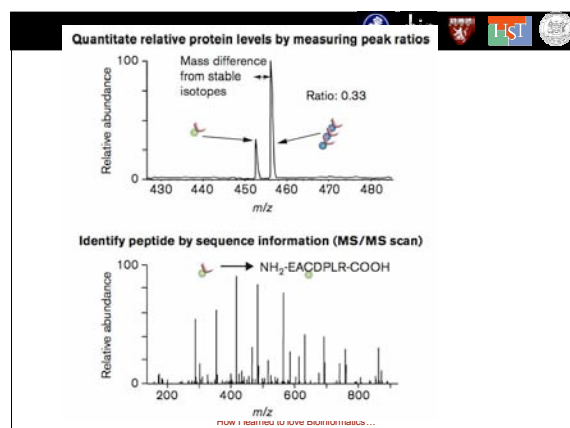
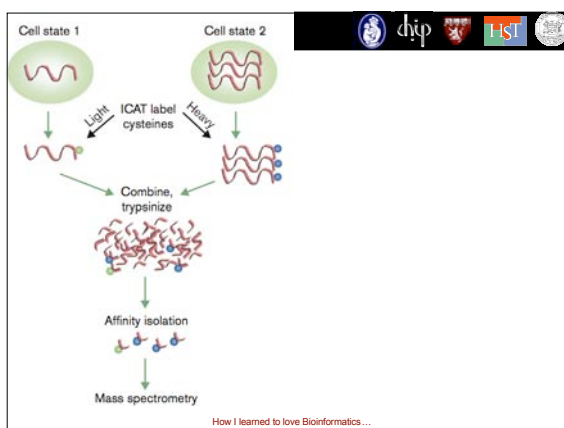
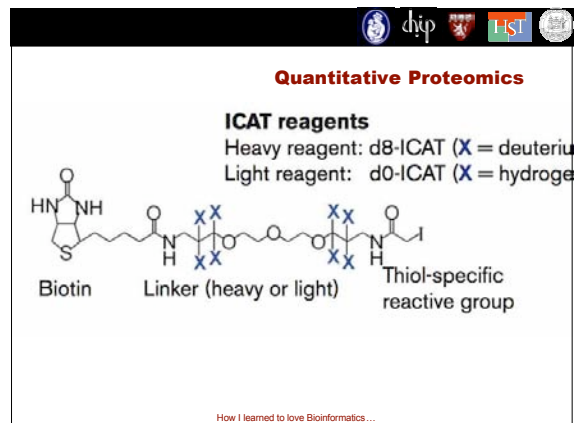
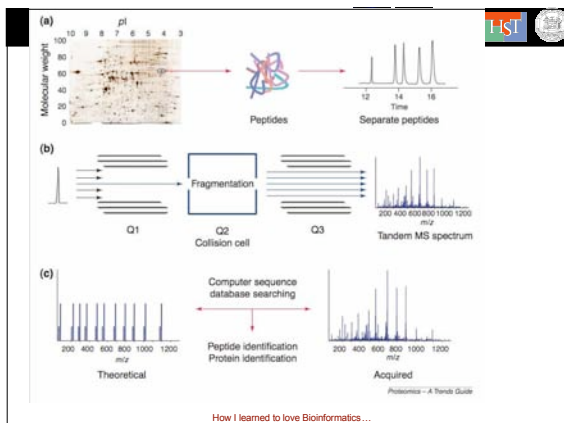
What is a microarray? (II)

- Translational friendliness. A clinical investigator does not have to understand molecular biology techniques in order to be able to provide the necessary materials for the acquisition of the microarray data.
- Identifiability. All items identified by the microarray technology whether they are proteins or RNA species should be automatically identified against standard reference nomenclatures.
- High Throughput. Hundred of patient samples can be processed within days.
- Commodity level priced infrastructure. The technology equipment and budget should be available to most biological and clinical investigational laboratories.
- Massively parallel measurements of the relevant analytes. That is, the members of transcriptome, the members of the proteome, the metabolome or any other comprehensive measure of molecular physiology.

How I learned to love Bioinformatics...

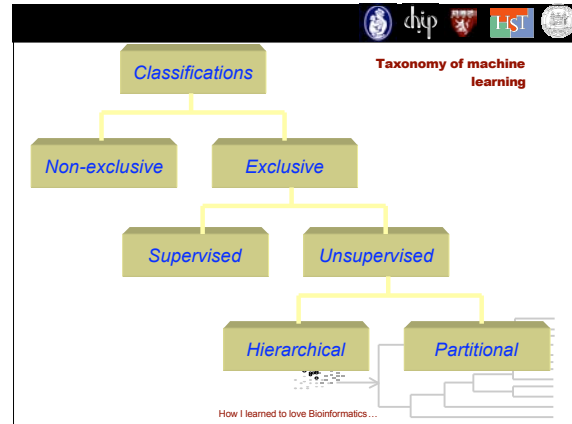
dhp HST

Can we build microarrays for proteomics?



So, do we have proteomic microarrays?

- Making sense of all the data**
- The underlying dogma
 - What data may be included in the data sets?
 - ✓ Beyond genomic data
 - ✓ Multiple scales and non-numeric measures



Phylogenetic-type tree

	RNA Expr Gene 1	RNA Expr Gene 2	RNA Expr Gene 3
Experiment 1	0.7	0.3	7.3
Experiment 2	1.2	1.9	6.5
Experiment 3	1.1	0.9	8.1

Correlation Coefficient

	Gene 1	Gene 2	Gene 3
Gene 1			
Gene 2	.88		
Gene 3	-.19	-.62	

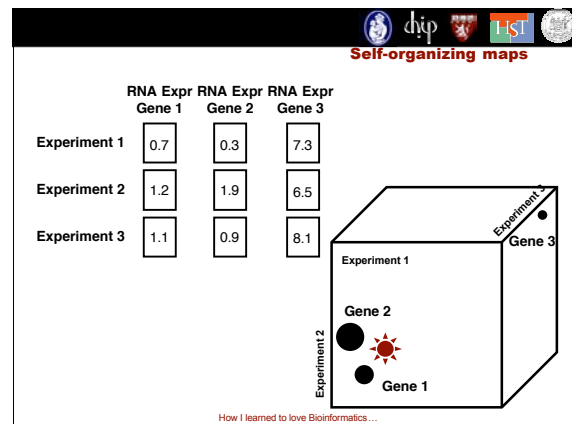
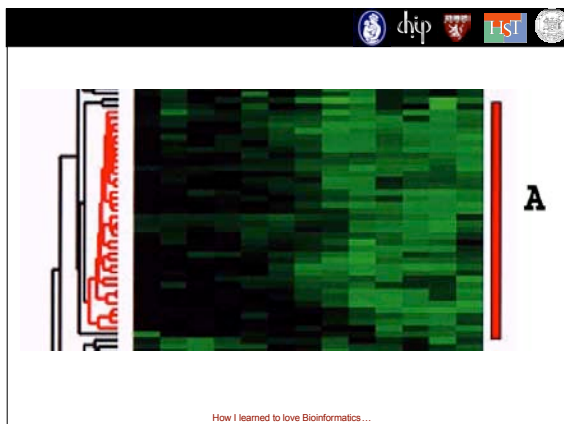
How I learned to love Bioinformatics...

Phylogenetic-type tree / Correlation coefficient

- Similarity score computed for two genes over the same conditions, similar to Pearson's correlation coefficient
- Found redundant representations and similarly functioning genes cluster together (for *S. cerevisiae*)
- Also found 10 temporal clusters in 8613 genes in the response of human fibroblasts to serum
- Suggests role for over 200 genes with previously unknown function

Eisen MB, PNAS 1998;95(25):14863-8.
Iyer VR, Science 1999;283(25):83-7.

How I learned to love Bioinformatics...



Relevance Networks

- Several algorithms have already been developed for knowledge discovery and data-mining of RNA expression data sets
- We are interested in finding networks of genes that are functionally clustered with little or no *a priori* knowledge (unsupervised learning)
- Relevance Networks** are an approach to analyze these data sets
- Previously validated in the clinical laboratory result domain

Children's Hospital, Patent Pending.
Butte A.J. Kohane IS, Unsupervised Knowledge Discovery in Medical Databases Using Relevance Networks, Symposia AMIA, 1999.

How I learned to love Bioinformatics...



Construction of Relevance Networks 1

- Patients and cell lines are analyzed as cases
- Clinical parameters, laboratory tests, RNA expression, and susceptibility to anti-cancer agents are all example features of those cases

Patient, Cell Line, Time, etc.	Lab Test	Lab Test 2	Clinical Param 1	RNA Expr J02923	Susceptibility to Anti-cancer Agent 169517
138	3.7	105	0.7	8.1	
134	4.5	99		2.1	
132	5.3	102	7.4	3.3	

How I learned to love Bioinformatics...

Construction of Relevance Networks 2

- For all pairs of features, we take overlapping values over the cases and make a scatter plot of values

How I learned to love Bioinformatics...

Construction of Relevance Networks 3

- Perform a pairwise comparison between all features
- For each scatter plot, we fit a linear model and stored \checkmark Correlation coefficient r

How I learned to love Bioinformatics...

Construction of Relevance Networks 4

- $r^2 = r^2 * r / \text{abs}(r)$
- Choose r^2 network
- Drop under threshold
- Breaks connected islands where connections are stronger than threshold
- Islands are what we call "relevance networks"
- Display graphically, with thick lines representing strongest links

How I learned to love Bioinformatics...

NCI 60 data set

- NCI 60 is set of 60 cancer cell lines against which the National Cancer Institute has tested over 50,000 chemicals to find anti-cancer agents
- NCI recently provided us with consistent, validated data for 5,084 agents, representing the concentration of each agent needed to cause 50% growth inhibition compared to control (GI50) for each cell line in NCI 60
- Data was provided as $-\log_{10}(\text{GI50})$
 - ✓ Lower number means higher GI50, or less sensitivity to an agent
- Unfortunately, few of these anti-cancer agents have documented mechanisms of action (or even a name listed)

Weinstein, et al. Science; 258: 447 (1992).
<http://www.dtp.nci.gov>

How I learned to love Bioinformatics...

RNA expression data set

- RNA expression in NCI 68 cell lines was determined using Affymetrix HU6000 arrays
 - ✓ 5,223 known genes
 - ✓ 1,193 expressed sequence tags
- The RNA expression data set and Anti-cancer susceptibility data set were merged, using the 60 cell lines the two tables had in common

How I learned to love Bioinformatics...

Distribution of r^2

- 11,692 features
- 68,345,586 total associations
 - ✓ 22 M between genes
 - ✓ 12 M between agents
 - ✓ 33 M between a gene and an agent

How I learned to love Bioinformatics...

Genes and Anti-Cancer Agents

- Threshold r^2 was 0.8
- 202 networks
- 834 features out of 11,692 (7.1%)
- 1,222 links out of 68,345,586 (.0018%)
- Only one link between a gene and anti-cancer agent

How I learned to love Bioinformatics...

Taxonomy of Network Links

- Identity / Synonymy
- Functional Similarity
- Derivation
- Biological Relationship

How I learned to love Bioinformatics...

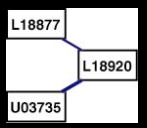
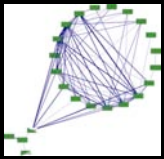
Identity / Synonymy Gene Networks

- 15 networks demonstrated synonymy associations
 - ✓ Most linked endogenous or spiked controls
- 5 / 15 were linked genes with multiple GenBank entries
 - ✓ L10838 and D28423: SRP20
 - ✓ M19267 and Z24727: Tropomyosin alpha chain
 - ✓ X17567 and X52979: snRNP B
 - ✓ U08021 and U51010: Nicotinamide N-methyltransferase
 - ✓ M14199 and U43901: Laminin receptor precursor and laminin receptor mRNA

How I learned to love Bioinformatics...

Functional Similarity Networks

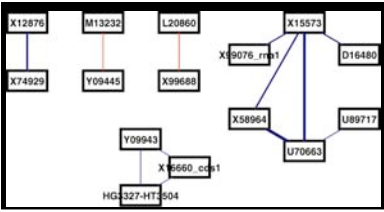
- Melanoma-associated antigens 2, 3 and 12
- Caldesmon 1 and alternative splicing products 3 and 4
- Two sequences from MHC class I (D32129 and X12432)
- 3088 (chlorambucil), 344007 (piperazine alkylator), 34462 (uracil nitrogen mustard), 48034 (aziridine), 6396 (thio-tepa), and 9706 (triethylenemelamine) are all alkylating agents
- 243928 (ethanesulfonamide derivative) and 249992 (amsacrine) are both active against topoisomerase II

How I learned to love Bioinformatics...

Biological Associations

- Keratin 8 linked to keratin 18
- Glycoprotein Ib beta linked to *psd*
- Others were not easily explained using the medical literature



How I learned to love Bioinformatics...

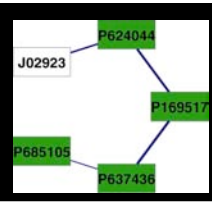
Derivative Networks

- 112167 (19-norlanosta-1) was derived from 112166 (curcubitacin)
- 274555 (gold derivative) was derived from 306388
- 295500, 374028, 606497, 606499, 610456, 610457, and 610459 are all camptothecin derivatives
- 302325 and 382034 are pyrimidinediamine derivatives
- 351710, 352299, and 627505 are methylleptidinium derivatives
- 603071 and 629971 are 9-aminocamptothecin derivatives
- 634785 and 634786 are 4-piperidinone derivatives

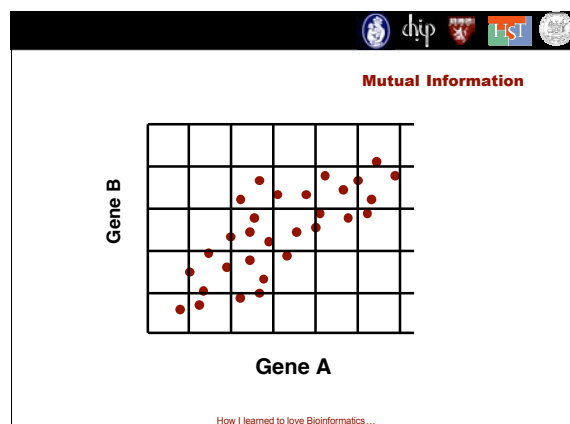
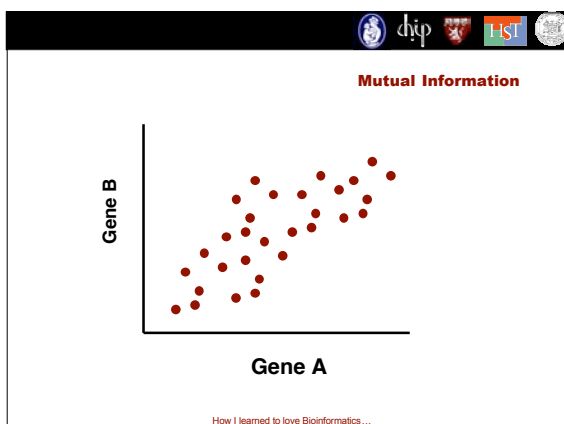
How I learned to love Bioinformatics...

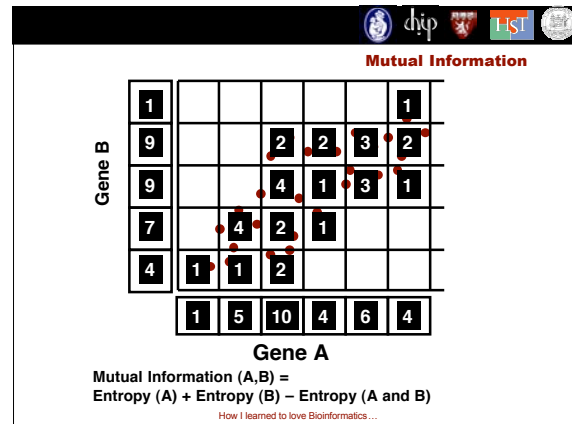
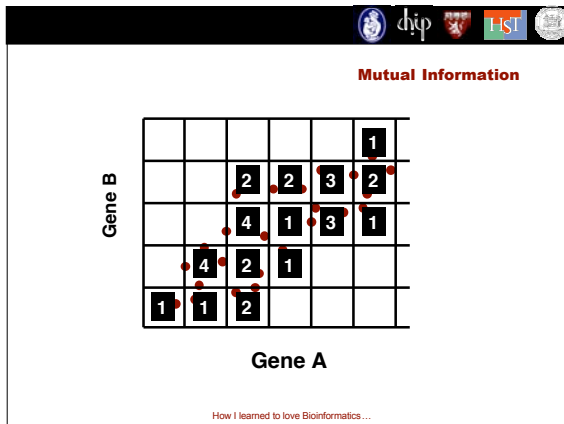
Genes and Anti-Cancer Agents

- Elevated levels of J02923 (lymphocyte cytosolic protein-1, LCP1, L-plastin, pp65) is associated with increased sensitivity to 624044
- Agent 624044 is 4-Thiazolidinecarboxylic acid, 3-[[6-[2-oxo-2-(phenylthio)ethyl]-3-cyclohexen-1-yl]acetyl]-2 thioxo-, methyl ester, [1R-[1a(R*),6a]]- (9CI)
- LCP1 is an actin-binding protein involved in leukocyte adhesion
- A role for LCP1 has been previously postulated
- Low level expression thought to occur in cancer cell lines
- Other thiazolidine derivatives are known to inhibit tumor cell growth



How I learned to love Bioinformatics...



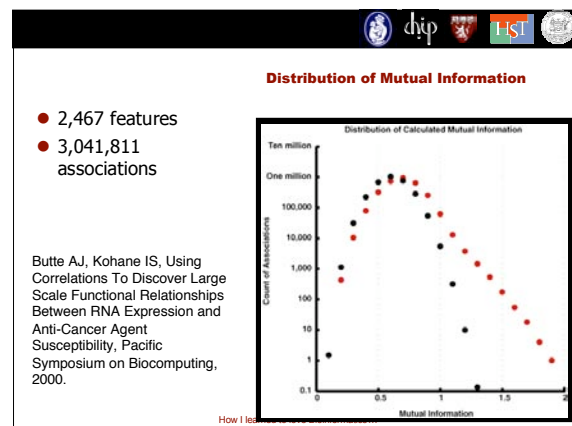


Global Regulatory Pathways Using Relevance Networks

- Stanford provides a yeast RNA expression data set for public use
- Data set of 2,467 open reading frames (ORF) measured under 79 conditions
- Several experiments with various time points in
 - ✓ Diauxic shift
 - ✓ Mitotic cell division cycle
 - ✓ Sporulation
 - ✓ Temperature shocks
 - ✓ Reducing shocks

Eisen MB, Proc Natl Acad Sci U S A 1998;95(25):14863-8.

How I learned to love Bioinformatics...



Synonymy Networks

- Threshold MI was 1.2
- 22 networks
- 199 features out of 2,467 (8.1%)
- Two networks with synonymy associations
 - ✓ Copper metallothionein
 - ✓ L-asparaginase II

How I learned to love Bioinformatics...

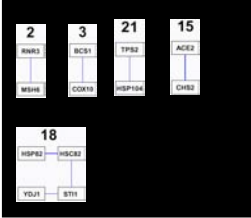
Similar Function Networks

- Nine networks with associations linking genes with similar functions
 - ✓ Histones
 - ✓ Acid phosphatases
 - ✓ Ribosomal proteins
 - ✓ Translation initiation
 - ✓ 70 kDa heat shock proteins
 - ✓ Hexose transporters
 - ✓ Mitochondrial ribosomal proteins

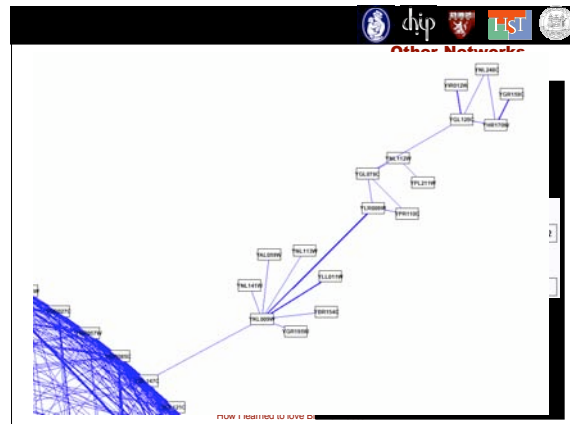
How I learned to love Bioinformatics...

Biological Pathway Networks

- Five networks with associations linking genes known to be related in the same biological pathway
 - ✓ Base pair mismatch repair
 - ✓ Cytochrome complex assembly
 - ✓ Trehalose-6-phosphate and chaparone
 - ✓ Chitinase expression regulator and chitin synthase II
 - ✓ Isoforms of chaparones

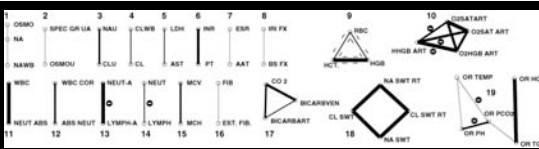


How I learned to love Bioinformatics...



Another domain: Clinical Laboratory

- Matrix of 642 lab tests with 28,566 overlapping results
- Threshold r^2 was 0.6, n was 50



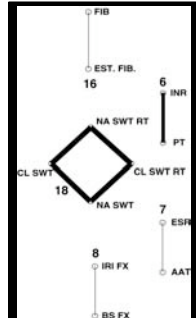
- 48 lab tests out of 642 (7%), 36 links out of 205,761 (.017%)

Butte AJ, Kohane IS, Unsupervised Knowledge Discovery in Medical Databases Using Relevance Networks, Symposia AMIA, 1999.

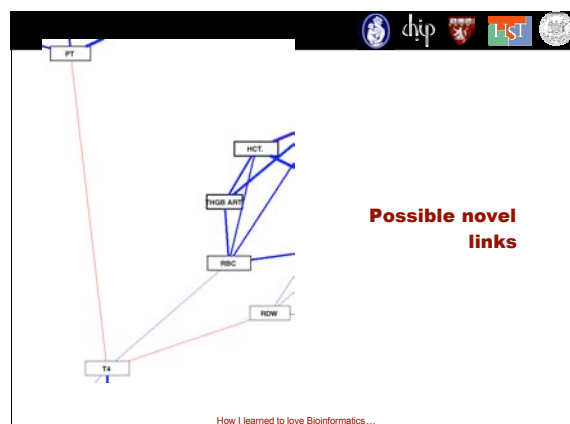
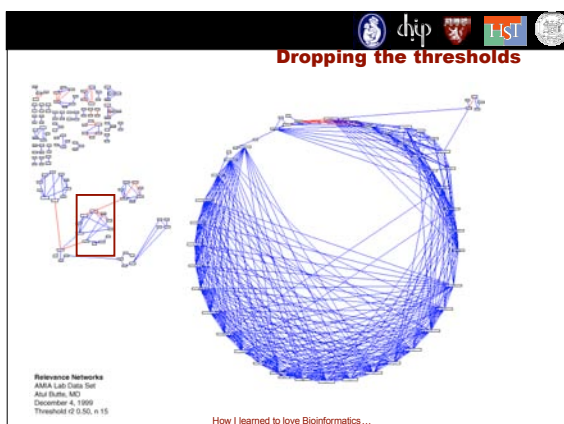
How I learned to love Bioinformatics...


Taxonomy of links

- Identity or synonymy
 - ✓ Serum fibrinogen and estimated fibrinogen
- Mathematical
 - ✓ Prothrombin time and International Normalized Ratio
- Physiologic
 - ✓ Sodium and chloride in sweat
- Pathologic
 - ✓ Erythrocyte sedimentation rate and alpha-1 antitrypsin
- Causal
 - ✓ Serum blood sugar and serum insulin level




How I learned to love Bioinformatics...





Supervised learning


- Several techniques available
- Decision trees
 - ✓ Y. Sun et al.
- Support Vector Machines
 - ✓ M. P. S. Brown et al., *Proc Natl Acad Sci U S A* **97**, 262-267 (2000).
- Neural Networks
 - ✓ C. Wu, S. Shivakumar, *Nucleic Acids Research* **22**, 4291-4299 (1994).

- 


Leukemia Classification

Morphology does not distinguish leukemias very well

Acute lymphoblastic leukemia (ALL)



Acute Myelogenous Leukemia (AML)



Courtesy P. Tamayo

How I learned to love Bioinformatics ...

How I learned to love Bioinformatics ...

Leukemia Initial (training) Dataset

Lymphoma Initial (training) Dataset

Normalized Expression

Color scale: -3 (blue) to 3 (red)

Genes (Leukemia Dataset):

- Proteinase 3 (C53612)
- Mitf 1 (U02292)
- C-myc (U03622)
- Myb (U03622)
- Bcl-2 (U03622)
- NF- κ B (U03622)
- IL-6 (U03622)
- IL-1 (U03622)
- IL-2 (U03622)
- IL-3 (U03622)
- IL-4 (U03622)
- IL-5 (U03622)
- IL-6 (U03622)
- IL-7 (U03622)
- IL-8 (U03622)
- IL-9 (U03622)
- IL-10 (U03622)
- IL-11 (U03622)
- IL-12 (U03622)
- IL-13 (U03622)
- IL-14 (U03622)
- IL-15 (U03622)
- IL-16 (U03622)
- IL-17 (U03622)
- IL-18 (U03622)
- IL-19 (U03622)
- IL-20 (U03622)
- IL-21 (U03622)
- IL-22 (U03622)
- IL-23 (U03622)
- IL-24 (U03622)
- IL-25 (U03622)
- IL-26 (U03622)
- IL-27 (U03622)
- IL-28 (U03622)
- IL-29 (U03622)
- IL-30 (U03622)
- IL-31 (U03622)
- IL-32 (U03622)
- IL-33 (U03622)
- IL-34 (U03622)
- IL-35 (U03622)
- IL-36 (U03622)
- IL-37 (U03622)
- IL-38 (U03622)
- IL-39 (U03622)
- IL-40 (U03622)
- IL-41 (U03622)
- IL-42 (U03622)
- IL-43 (U03622)
- IL-44 (U03622)
- IL-45 (U03622)
- IL-46 (U03622)
- IL-47 (U03622)
- IL-48 (U03622)
- IL-49 (U03622)
- IL-50 (U03622)
- IL-51 (U03622)
- IL-52 (U03622)
- IL-53 (U03622)
- IL-54 (U03622)
- IL-55 (U03622)
- IL-56 (U03622)
- IL-57 (U03622)
- IL-58 (U03622)
- IL-59 (U03622)
- IL-60 (U03622)
- IL-61 (U03622)
- IL-62 (U03622)
- IL-63 (U03622)
- IL-64 (U03622)
- IL-65 (U03622)
- IL-66 (U03622)
- IL-67 (U03622)
- IL-68 (U03622)
- IL-69 (U03622)
- IL-70 (U03622)
- IL-71 (U03622)
- IL-72 (U03622)
- IL-73 (U03622)
- IL-74 (U03622)
- IL-75 (U03622)
- IL-76 (U03622)
- IL-77 (U03622)
- IL-78 (U03622)
- IL-79 (U03622)
- IL-80 (U03622)
- IL-81 (U03622)
- IL-82 (U03622)
- IL-83 (U03622)
- IL-84 (U03622)
- IL-85 (U03622)
- IL-86 (U03622)
- IL-87 (U03622)
- IL-88 (U03622)
- IL-89 (U03622)
- IL-90 (U03622)
- IL-91 (U03622)
- IL-92 (U03622)
- IL-93 (U03622)
- IL-94 (U03622)
- IL-95 (U03622)
- IL-96 (U03622)
- IL-97 (U03622)
- IL-98 (U03622)
- IL-99 (U03622)
- IL-100 (U03622)
- IL-101 (U03622)
- IL-102 (U03622)
- IL-103 (U03622)
- IL-104 (U03622)
- IL-105 (U03622)
- IL-106 (U03622)
- IL-107 (U03622)
- IL-108 (U03622)
- IL-109 (U03622)
- IL-110 (U03622)
- IL-111 (U03622)
- IL-112 (U03622)
- IL-113 (U03622)
- IL-114 (U03622)
- IL-115 (U03622)
- IL-116 (U03622)
- IL-117 (U03622)
- IL-118 (U03622)
- IL-119 (U03622)
- IL-120 (U03622)
- IL-121 (U03622)
- IL-122 (U03622)
- IL-123 (U03622)
- IL-124 (U03622)
- IL-125 (U03622)
- IL-126 (U03622)
- IL-127 (U03622)
- IL-128 (U03622)
- IL-129 (U03622)
- IL-130 (U03622)
- IL-131 (U03622)
- IL-132 (U03622)
- IL-133 (U03622)
- IL-134 (U03622)
- IL-135 (U03622)
- IL-136 (U03622)
- IL-137 (U03622)
- IL-138 (U03622)
- IL-139 (U03622)
- IL-140 (U03622)
- IL-141 (U03622)
- IL-142 (U03622)
- IL-143 (U03622)
- IL-144 (U03622)
- IL-145 (U03622)
- IL-146 (U03622)
- IL-147 (U03622)
- IL-148 (U03622)
- IL-149 (U03622)
- IL-150 (U03622)
- IL-151 (U03622)
- IL-152 (U03622)
- IL-153 (U03622)
- IL-154 (U03622)
- IL-155 (U03622)
- IL-156 (U03622)
- IL-157 (U03622)
- IL-158 (U03622)
- IL-159 (U03622)
- IL-160 (U03622)
- IL-161 (U03622)
- IL-162 (U03622)
- IL-163 (U03622)
- IL-164 (U03622)
- IL-165 (U03622)
- IL-166 (U03622)
- IL-167 (U03622)
- IL-168 (U03622)
- IL-169 (U03622)
- IL-170 (U03622)
- IL-171 (U03622)
- IL-172 (U03622)
- IL-173 (U03622)
- IL-174 (U03622)
- IL-175 (U03622)
- IL-176 (U03622)
- IL-177 (U03622)
- IL-178 (U03622)
- IL-179 (U03622)
- IL-180 (U03622)
- IL-181 (U03622)
- IL-182 (U03622)
- IL-183 (U03622)
- IL-184 (U03622)
- IL-185 (U03622)
- IL-186 (U03622)
- IL-187 (U03622)
- IL-188 (U03622)
- IL-189 (U03622)
- IL-190 (U03622)
- IL-191 (U03622)
- IL-192 (U03622)
- IL-193 (U03622)
- IL-194 (U03622)
- IL-195 (U03622)
- IL-196 (U03622)
- IL-197 (U03622)
- IL-198 (U03622)
- IL-199 (U03622)
- IL-200 (U03622)
- IL-201 (U03622)
- IL-202 (U03622)
- IL-203 (U03622)
- IL-204 (U03622)
- IL-205 (U03622)
- IL-206 (U03622)
- IL-207 (U03622)
- IL-208 (U03622)
- IL-209 (U03622)
- IL-210 (U03622)
- IL-211 (U03622)
- IL-212 (U03622)
- IL-213 (U03622)
- IL-21



Resources

- MeV - TIGR
- Genecluster 2.0 (Broad Institute)

How I learned to love Bioinformatics...

