

## Computational Biology Machine Learning Approaches to Biology

**Computational Science Colloquium** 

## **Sepp Hochreiter**

Institute of Bioinformatics Johannes Kepler University, Linz, Austria

#### **Bioinformatics and Mathematics**





Hardy's 1908 paper:

To the Editor of Science: I am reluctant to intrude in a discussion concerning matters of which I have no expert knowledge, and I should have expected the <u>very simple point which I wish to</u> <u>make to have been familiar to biologists</u>. However, some remarks of Mr. Udny Yule, to which Mr. R. C. Punnett has called my attention, suggest that it may still be worth making...

Hardy–Weinberg principle: both allele and genotype frequencies in a population remain constant from generation to generation

## **Computational Biology** $\subset$ **Bioinformatics**



http://www.ncbi.nlm.nih.gov/About/ primer/bioinformatics.html:



N ational C enter for B iotechnology I nformation

#### **Bioinformatics:**

the field of bioinformatics has evolved such that the most pressing task now involves the *analysis and interpretation* of various types of data, including nucleotide and amino acid sequences, protein domains, and protein structures

Computational biology: actual process of analyzing and interpreting data

#### **Overview**

BIOINF

- virus drug targets
- single nucleotide polymorphism: schizophrenia
- microarray preprocessing
- classification: therapy outcome from gene expression profiles
- feature selection  $\rightarrow$  gene selection
- diagnosis based on peptide arrays
- copy number variations for complex diseases (bipolar)
- sequence classification: recurrent neural networks
- next generation sequencing
- nucleosome positions
- WE NEED HELP





Group: Molecular Libraries Dr. Volkmer

Peptide arrays



School of Chemistry Prof. Dek Woolfson



## U.S. Human Cases of Swine Flu Infection (As of April 28, 2009, 11:00 AM ET) State # of laboratory confirmed cases

California

Kansas

New York City

Ohio

Texas

#### TOTAL COUNT

10 cases

2 cases

45 cases

1 case

6 cases

64 cases



Welt

28.04.2009 um 23:33 Uhr

# Schweinegrippe: Erster Fall in Wien bestätigt

Im Krankenhaus in Steyr (OÖ) wurde heute Vormittag eine junge Frau eingeliefert, die ersten Angaben zufolge die Symptome des Schweinegrippe-Virus aufweist. Die 28-Jährige war vor kurzem von einem Mexiko-Urlaub heimgekehrt und hat sich offensichtlich dort infiziert. Die mexikanischen Behörden versuchen indes, die Herkunft des Virus ausfindig zu machen - die Spur führt zu einer Schweinefarm im Norden des Landes.











#### Coiled-coil domain prediction

#### a-b-c-d-e-f-g

(a) and (d)  $\rightarrow$  hydrophobic amino acids  $\rightarrow$  interaction between helices

(e) and (g)  $\rightarrow$  charged amino acids  $\rightarrow$  inter- and intra-helical salt bridges

(b), (c), and (f)  $\rightarrow$  polar amino acids  $\rightarrow$  outer surface



#### RMKQLEDKVEELLSKNYHLENEVARLKKLVG abcdefgabcdefgabcdefgabcdefgab

















Charité



Group: Dept. Psychiatry and Psychotherapy Prof. Andreas Heinz

Group: Dept. Nephrology and Internal Intensive Care Prof. Petra Reinke



## **Cleveland Clinic** Cleveland, Ohio, United States

Group:

Experimental Haematology and Hematopoiesis Taussig Cancer Center Prof. Dr. J. Maciejewski





single nucleotide polymorphism (SNP - pronounced snip)

- →Variation in the DNA at the same position in at least 1% of the population
- →SNPs occur all 100 to 300 base pairs
- →Current research relate diseases to SNPs
- →Each position exist twice in the human genome: heterozygous











## BIOINF

#### Lactose SNP



### Single Nucleotide Polymorphism: Schizophrenia





Schizophrenia is associated with genetic factors

#### Rates of Schizophrenia Among Relatives of Schizophrenic Patients\*



## Single Nucleotide Polymorphism: Schizophrenia





- medial prefrontal cortex and posterior cingulate cortex: associated with self-reflection
- increased connectivity between these regions:
   lost in their own world →trouble performing tasks →schizophrenic
- "People normally suppress this default system when they perform challenging tasks, but we found that patients with schizophrenia don't do this" (John D. Gabrieli)

## Single Nucleotide Polymorphism: Schizophrenia





Computational Science Colloquium, 29.04.2009

BIOIN

Johnson-Johnson



Beerse, Belgium













#### University Hasselt Statistics



## BIOINF

#### **Gene Network Science**



http://www.gnsbiotech.com/static\_content/johannes-kepler.html

"GNS is collaborating with researchers at Johannes Kepler University, Linz, by way of genomics data processing algorithms developed by the researchers and licensed by GNS."



Affymetrix Fluidics station Wash / Stain

Affymetrix Scanner







#### Fluorescence intensity image





Our new approach based on factor analysis: FARMS (Factor Analysis for Robust Microarray Summarization)









$$x = \lambda z + \varepsilon$$

Generative model:

- $\rightarrow$  z: factor N(0,1)
- $\rightarrow$   $\epsilon$ : noise  $N(0, \Psi)$
- $\rightarrow$   $\lambda$ : loading vector
- $\rightarrow$  x: data  $N(0, \lambda\lambda' + \Psi)$
- $\rightarrow \Psi, \lambda \rightarrow \text{EM-algorithm}$

 $\Psi$  is diagonal Covariance x:  $\lambda\lambda' + \Psi$ 

#### Correlations between probes can only be explained by hidden factor

Feature selection: minimal  $\operatorname{var}(z \mid \boldsymbol{x}) = (1 + \boldsymbol{\lambda}^T \boldsymbol{\Psi}^{-1} \boldsymbol{\lambda})^{-1}$ 

Maximum a posterior:

$$riangle$$
 Data:  $\{oldsymbol{x}\}\ =\ \{oldsymbol{x}_1,\ldots,oldsymbol{x}_N\}$ 

Posterior:  $p(\boldsymbol{\lambda}, \boldsymbol{\Psi} \mid \{\boldsymbol{x}\}) \propto p(\{\boldsymbol{x}\} \mid \boldsymbol{\lambda}, \boldsymbol{\Psi}) p(\boldsymbol{\lambda}, \boldsymbol{\Psi})$ 

• Likelihood:  

$$p(\{ m{x}\} \mid m{\lambda}, m{\Psi}) = \prod_{i=1}^N \mathcal{N} \left( m{0} \ , \ m{\lambda} m{\lambda}^T \ + \ m{\Psi} 
ight) (m{x}_i)$$

Prior:

 $p(\boldsymbol{\lambda}, \boldsymbol{\Psi})$ 



E-step  

$$\begin{aligned} \mu_{z_i | \boldsymbol{x}_i} &= (\boldsymbol{x}_i)^T \left( \boldsymbol{\Lambda} \ \boldsymbol{\Lambda}^T + \boldsymbol{\Psi} \right)^{-1} \boldsymbol{\Lambda} \\ E_{z_i | \boldsymbol{x}_i} &(z_i^2) &= \mu_{z_i | \boldsymbol{x}_i}^2 + \sigma_{z_i | \boldsymbol{x}_i}^2 \end{aligned} \qquad \begin{aligned} \mu_{z_i | \boldsymbol{x}_i} &= 1 - \boldsymbol{\Lambda}^T \left( \boldsymbol{\Lambda} \ \boldsymbol{\Lambda}^T + \boldsymbol{\Psi} \right)^{-1} \boldsymbol{\Lambda} \\ &= 1 - \boldsymbol{\Lambda}^T \left( \boldsymbol{\Lambda} \ \boldsymbol{\Lambda}^T + \boldsymbol{\Psi} \right)^{-1} \boldsymbol{\Lambda} \end{aligned} = \\ & \left( 1 + \boldsymbol{\Lambda}^T \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda} \right)^{-1} \end{aligned}$$

$$\begin{split} & \mathsf{M}\text{-step} \\ \lambda_{j}^{\text{Gauss}} = \left(\frac{1}{N}\sum_{i=1}^{N} x_{ij} \operatorname{E}_{z_{i}|\boldsymbol{x}_{i}}(z_{i}) + \frac{1}{N}\frac{\mu_{\Lambda} \Psi_{jj}^{\text{new}}}{\sigma_{\Lambda}^{2}}\right) \left(\frac{1}{N}\sum_{i=1}^{N} \operatorname{E}_{z_{i}|\boldsymbol{x}_{i}}(z_{i}^{2}) + \frac{1}{N}\frac{\Psi_{jj}^{\text{new}}}{\sigma_{\Lambda}^{2}}\right)^{-1} \\ \lambda_{j}^{\text{new}} = \begin{cases} \lambda_{j}^{\text{Gauss}} & \text{for } \lambda_{j}^{\text{Gauss}} > 0 \\ 0 & \text{for } \lambda_{j}^{\text{Gauss}} \leq 0 \\ 0 & \text{for } \lambda_{j}^{\text{Gauss}} \leq 0 \end{cases} \\ \Psi_{jj}^{\text{new}} = \left[\operatorname{diagvect}\left(\frac{1}{N}\sum_{i=1}^{N} \boldsymbol{x}_{i}\boldsymbol{x}_{i}^{T}\right)\right]_{j} - \lambda_{j}^{\text{new}}\left[\frac{1}{N}\sum_{i=1}^{N} \operatorname{E}_{z_{i}|\boldsymbol{x}_{i}}(z_{i}) \boldsymbol{x}_{i}\right]_{j} + \frac{1}{N}\frac{\Psi_{jj}^{\text{new}}}{\sigma_{\Lambda}^{2}}\lambda_{j}^{\text{new}}\left(\mu_{\Lambda} - \lambda_{j}^{\text{new}}\right) \end{split}$$



International competition: "Affycomp" (http://affycomp.biostat.jhsph.edu/)

N	Method / Submitter	1	2	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>	<u>9</u>	<u>10</u>	11	12	<u>13</u>	<u>14</u>
0	(perfection)	0.00	0.00	0.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
89	<u>VST-RMA / dupan</u>	0.04	0.07	0.23	0.80	0.49	0.12	0.51	0.46	0.47	0.15	0.54	0.96	0.93	0.64
88	<u>en.FARMS / clevert</u>	0.00	0.00	0.00	0.85	0.63	0.16	0.72	0.62	0.65	0.25	1.00	1.00	1.00	1.00
81	<u>dfwfcb / zhongxue</u>	0.00	0.00	0.00	0.63	0.31	0.09	0.32	0.26	0.29	0.12	1.00	1.00	1.00	1.00

AUC (most relevant criterion)

Participants from Berkely, Affymetrix, EBI, Roche, etc.

Johnson & Johnson tested it on 30 internal data sets

 $\rightarrow$  Now their default microarray normalization method

BIOINF

Extensions FARMS:

- Few arrays differ → few factors differ → factor is Laplace distributed Laplace FARMS
  - integral cannot be analytically computed
  - variational approach (from physics, "calculus of variations"  $\rightarrow$  functionals)
  - mixture of Gaussians ("independent factor analysis", Attias 99)

#### Laplace

Gauss

variational:  $\alpha_i$  is variance of the local Gaussian approximation of the factor

#### E-step

$$\sigma_{z_i \mid \boldsymbol{x}_i}^2 = \left( lpha_i^{-1} + \boldsymbol{\Lambda}^T \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda} 
ight)^{-1}$$

$$\sigma_{z_i|\boldsymbol{x}_i}^2 = \left(1 + \boldsymbol{\Lambda}^T \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda}\right)^{-1}$$

#### M-step

$$lpha_i = \sqrt{\mu_{z_i|m{x}_i}^2 + \sigma_{z_i|m{x}_i}^2}$$

BIOINF

Extensions FARMS:

- Conditional FARMS (factor depends on external input)

Update Rules:  $\boldsymbol{\Lambda} = \left(\frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}_{i} \operatorname{E}_{\boldsymbol{z}_{i} \mid \boldsymbol{x}_{i}, \boldsymbol{u}_{i}} (\boldsymbol{z}_{i}^{T})\right) \left(\frac{1}{N} \sum_{i=1}^{N} \operatorname{E}_{\boldsymbol{z}_{i} \mid \boldsymbol{x}_{i}, \boldsymbol{u}_{i}} (\boldsymbol{z}_{i} \boldsymbol{z}_{i}^{T})\right)^{-1}$  $oldsymbol{B} = \left( rac{1}{N} \sum_{i=1}^{N} \mathbb{E}_{oldsymbol{z}_i \mid oldsymbol{x}_i, oldsymbol{u}_i}(oldsymbol{z}_i) oldsymbol{u}_i^T 
ight) \left( rac{1}{N} \sum_{i=1}^{N} oldsymbol{u}_i oldsymbol{u}_i^T 
ight)^{-1}$  $\boldsymbol{\Sigma} = \frac{1}{N} \sum_{\boldsymbol{z}_i \mid \boldsymbol{x}_i, \boldsymbol{u}_i}^{N} \left( \boldsymbol{z}_i \ \boldsymbol{z}_i^T \right) - \boldsymbol{B} \frac{1}{N} \sum_{\boldsymbol{z}_i \mid \boldsymbol{x}_i, \boldsymbol{u}_i}^{N} \left( \boldsymbol{z}_i^T \right)$  $\boldsymbol{\Psi} \;=\; rac{1}{N}\;\sum_{i}^{N} oldsymbol{x}_{i} oldsymbol{x}_{i}^{T} \;-\; oldsymbol{\Lambda}\;rac{1}{N}\;\sum_{i}^{N} \mathrm{E}_{oldsymbol{z}_{i}|oldsymbol{x}_{i},oldsymbol{u}_{i}}\left(oldsymbol{z}_{i}
ight)oldsymbol{x}_{i}^{T}$ 



#### Given

- →Tissue samples (Tumor, Blood, …)
- →Therapy outcome, clinical features, tumor type or state

#### Goals:

- →Diagnosis (tumor type)
- →Prognosis of the therapy outcome (alternative therapy, medication)
- →Identification of relevant genes (drug design)

#### Working hypothesis

Gene expression profile  $\rightarrow$  cell state  $\rightarrow$  tumor state  $\rightarrow$  prediction

Prediction: Support Vector Machine and Feature Selection



#### Classification



Angiogenesis and

Metastasis








#### SVM model selection

primal problem	Lagrangian ${\cal L}$ with multipliers $lpha_i$
$egin{array}{ll} \min & rac{1}{2} \ oldsymbol{w}\ ^2 \ { m s.t.} & y^i ~ig(oldsymbol{w}^Toldsymbol{x}^i ~+~ big) ~\geq~ 1 \end{array}$	$ \begin{aligned} \mathcal{L}(\boldsymbol{w}, b, \boldsymbol{\alpha}) &= \frac{1}{2} \ \boldsymbol{w}\ ^2 - \sum_{i=1}^{l} \alpha_i \left( y^i \left( \boldsymbol{w}^T \boldsymbol{x}^i + b \right) - 1 \right) \\ \frac{\partial \mathcal{L}(\boldsymbol{w}, b, \boldsymbol{\alpha})}{\partial \boldsymbol{w}} &= \boldsymbol{w} - \sum_{i=1}^{l} \alpha_i y^i \boldsymbol{x}^i = \boldsymbol{0} \end{aligned} $
dual problem	$\Rightarrow w = \sum_{i=1}^{\infty} \alpha_i y^i x^i$ with respect to the data to classify
$ \min_{\boldsymbol{\alpha}}  \frac{1}{2} \sum_{i,j} \alpha_i \ \alpha_j \ y^i \ y^j \ \left(\boldsymbol{x}^j\right)^T \boldsymbol{x}^i \ - \ \sum_{i=1}^l \alpha_i $ s.t. $ \alpha_i \ \ge \ 0 $ $ \sum_{i=1}^l \alpha_i \ y^i \ = \ 0 $	$\frac{\partial \mathcal{L}(\boldsymbol{w}, b, \boldsymbol{\alpha})}{\partial b} = \boxed{\sum_{i=1}^{l} \alpha_i \ y^i} = 0$

not linear separable

- class –1
- class +1







#### non-linear support vector machine

- ( )
- $\bigcirc$ 
  - support vectors



#### Theorem 1 (Mercer)

Let the kernel k be symmetric and from  $L_2(X \times X)$  defining a Hilbert-Schmidt operator

$$T_k(f)(\boldsymbol{x}) = \int_X k(\boldsymbol{x}, \boldsymbol{x}') f(\boldsymbol{x}') d\boldsymbol{x}'$$

If  $T_k$  is positive definite, i.e. for all  $f \in L_2(X)$  $\int_{X \times X} k(\boldsymbol{x}, \boldsymbol{x}') \ f(\boldsymbol{x}) \ f(\boldsymbol{x}') \ d\boldsymbol{x} \ d\boldsymbol{x}' \ge 0$ ,

then  $T_k$  has eigenvalues  $\lambda_j \geq 0$  with associated eigenfunctions  $\psi_j \in L_2(X)$ . Further

$$egin{array}{lll} (\lambda_1,\lambda_2,\ldots)\in\ell_1 \ k(oldsymbol{x},oldsymbol{x}') \ = \ \sum_j\lambda_j \ \psi_j(oldsymbol{x}) \ \psi_j(oldsymbol{x}') \ , \end{array}$$

where  $\ell_1$  is the space of vectors with finite one-norm and the last sum converges absolutely and uniformly for almost all  $\boldsymbol{x}$  and  $\boldsymbol{x}'$ .





#### Problem

- →Many genes (several 10,000 features) but few samples (about 100 examples)
- →SVM theory requires more examples (*L*) then features (*n*) for a low error bound on future data (determined by L/n).

#### Solution

→ Feature selection: Selection of relevant genes and, therefore, decreasing the input dimension n







**Theorem 1 (Singular Value Expansion)** Let  $\alpha$  be from  $H_1 := L^2(\mathcal{Z})$  and let k be a kernel from  $H_3 := L^2(\mathcal{X}, \mathcal{Z})$  which defines a Hilbert-Schmidt operator  $T_k : H_1 \to H_2$   $(H_2 := L^2(\mathcal{X}))$   $(T_k \alpha)(x) = f(x) = \int_{\mathcal{Z}} k(x, z) \alpha(z) dz$ . Then  $||f||^2_{H_2} = \langle T_k^* T_k \alpha, \alpha \rangle_{H_1}$   $(T_k^* adjoint operator of T_k)$  and  $k(x, z) = \sum_n s_n e_n(z) g_n(x)$ 

which converges in the  $L^2$ -sense. The  $s_n \ge 0$  are the singular values of  $T_k$ , and  $e_n \in H_1, g_n \in H_2$  are the corresponding orthonormal functions. If  $\int_{\mathcal{Z}} (k(x,z))^2 dz \le K^2$  for all  $x \in \mathcal{X}$ , then the following sum convergences absolutely and uniformly:

$$f(x) = \sum_{n} s_n \langle \alpha, e_n \rangle_{H_1} g_n(x)$$



$$egin{array}{lll} \min & rac{1}{2} \left( oldsymbol lpha^+ \, - \, oldsymbol lpha^- 
ight)^ op K^ op K \left( oldsymbol lpha^+ \, - \, oldsymbol lpha^- 
ight) \ & - oldsymbol y^ op K \left( oldsymbol lpha^+ \, - \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- 
ight) \ & + \epsilon \ oldsymbol 1^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^- \, + \, oldsymbol a^ op \left( oldsymbol lpha^+ \, + \, oldsymbol lpha^+ \, + \, oldsymbol a^ op \left( oldsymbol lpha^+ \, + \, oldsymbol a^ op \left( oldsymbol lpha^+ \, + \, oldsymbol a^ op \left( oldsymbol a^ op \left( oldsymbol lpha^+ \, + \, oldsymbol a^ op \left( oldsymbol a^$$

 $m{y}$  is vector of labels  $m{K}$  is kernel matrix 1 is vector of ones  $m{lpha} = m{lpha}^+ - m{lpha}^$  $m{w} = m{Z} \m{lpha} = \sum_{j=1}^l \alpha_j \m{z}$ 

Fast optimization through a new sequential minimal optimization (SMO) technique: only box constraints!

$$f(x) = w^T x + b = \sum_{j=1}^{l} \alpha_j z^T x + b = \sum_{j=1}^{l} \alpha_j K_{(x)j} + b$$

#### Task

Brain tumor (medulloblastoma) patients respond differently to the therapy (chemo, radiation)

- →Negative prognosis: alternative therapy or more intensive control
- →Positive prognosis: reduced medication

60 patients and 7129 genes

S. L. Pomeroy, P. Tamayo, M. Gaasenbeek, L. M. Sturla, M. Angelo, M. E. McLaughlin, J. Y. H. Kim, L. C. Goumnerova and P. M. Black, C. Lau, J. C. Allen, D. Zagzag, J. M. Olson, T. Curran, C. Wetmore, J. A. Biegel, T. Poggio, S. Mukherjee, R. Rifkin, A. Califano, G. Stolovitzky, D. N. Louis, J. P. Mesirov, E. S. Lander, T. R. Golub Prediction of central nervous system embryonal tumour outcome based on gene expression Nature 415(687):436-442, 2002





## **Microarray: Feature Selection / Classification**

#### **Classification results**

Standard		New method				
Method	F	Error	Method	Features	Error	
TrkC (1 gene)	1	20	SVM	40 / 45 / 50	5 / 4 / 5	
SVM		15	SVM	40 / 45 / 50	5/5/5	
TrkC & SVM		14	P-SVM	40 / 45 / 50	4 / 4 / 5	
KNN	8	13				
KNN & SVM		12				

Standard feature selection with signal-to-noise- and *t*-statistics

# BIOINF

#### Task

Breast cancer: 70-80% of the patients do not need a treatment because metastasis does not appear

- → Prediction of metastasis: therapy selection
- →Alternative treatment or individual medication (toxicity)

78 patients and 25,000 genes

L. J. van't Veer, H, Dai, M. J. van de Vijver, Y. D. He, A. A. M. Hart, M. Mao, H. L. Peterse, K. van der Kooy, M. J. Marton, A. T. Witteveen, G. J. Schreiber, R. M. Kerkhoven, C. Roberts, P. S. Linsley, R. Bernards, S. H. Friend Gene expression profiling predicts clinical outcome of breast cancer Nature 415: 530-536, 2002



## **Microarray: Feature Selection / Classification**

#### **Classification results**

Standard feature selection				Ne					
Method	F	Error	ROC	Test	Method	F	Error	ROC	Test
weighted voting	70	20	0.77	2	SVM	30	12	0.88	2

Standard feature selection with signal-to-noise-statistics

## **Microarray: Feature Selection / Classification**

Signature Diagnostics AG Prof. Dr. André Rosenthal, CEO



Press release (<u>http://www.signature-diagnostics.de/press10.htm</u>):

Signature Diagnostics AG announced today that the company has developed a new <u>gene expression classifier</u>, that predicts the risk of recurrence and metastasis (progression of disease) in stage I and II colorectal cancer.

Note: 55% employee patent fee hold by Djork-Arne Clevert, JKU







Group: Institute of theoretical Biology Systems Immunology Prof. Michal Or-Guil

Humboldt University



Computational Science Colloquium, 29.04.2009



Mice with and without parasite (*H. polygyrus*, worm)



BIOINF

Given: serum for each mouse containing antibodies

- Goal: predict whether the mouse has a parasite
- Assumption: Immune system produces antibodies against parasite proteins





Antibody detection with random arrays and ML (feature selection and classification)



## BIOINF

#### Minimal number of peptide features

	lgM	lgG	lgM	lgG	lgM	lgG
Sensitivity	1	1	0.93	0.80	0.93	0.86
Specificity	0.96	0.93	1	0.87	1	1
Diagnostic reliability (%)	97.7	90.7	96.7	83.4	95.3	90.7
Number of features	2	2	6	4	5	4
Significance	<0.001	<0.001	0.004	0.014	<0.001	<0.001

Josef Penninger (Vienna): cancer cells activate the immune system (verified at mice)

 $\rightarrow$  Diagnosis of cancer by peptide arrays

## MERCK Serono



Computational Science Colloquium, 29.04.2009



Geneva, Switzerland







Darmstadt, Germany







complex diseases

e.g.: bipolar disorder

Computational Science Colloquium, 29.04.2009

Copy-number variant (CNV):

- →1 kb or larger (50 kb) DNA sequences
- →variable copy-number compared to a reference
- ⇒insertions, deletions and duplication
- → 2004 first publications: Sebat et al., Science Iafrate et al., Nature Genetics











BIOINF



Computational Science Colloquium, 29.04.2009

BIOINF







Smooth scatter plots of an Affymetrix SNP 6.0 array for Chr. 6 (using FARMS)

*Left*: defective amplifications in a sample stemming from lung cancer tissue

*Right*. the same region corresponding to a healthy cell







#### The problem

→ classify a protein given as sequence of amino acids into functional or structural classes





Myohemerythrin





Tobacco mosaic coat protein (a) Predominantly a helix

Immunoglobulin, V. domain

(b) Predominantly β sheet

Hexokinase, domain 2 (c) Mixed a helix and B sheet

APSRKFFVGGNWKMNGRKQSL GELTGTLNAAKVPADTEVVCA PPTAYIDFARQKLDPKIAVAA QNCYKVTNGAFTGEISPGMIK **DCGATWVVLGHSERRHVFGES** DELIGQKVAHALAEGLGVIAC IGEKLDEREAGITEKVVFEQT KVTADNVKDWSKVVI.AYEPVW AIGTGKTATPQQAQEVHEKLR GWLKSNVSDAVAQSTRIIYGG SVTGATCKELASQPDVDGFLV GGASLKPEFVDIINAKQ

Computational Science Colloquium, 29.04.2009

BIOINF

Long Short-Term Memory (LSTM)











#### MTYKLILNGKTLKGETTTEAVDAATAEKVFKQYANDNGVDGEWTYDDATKTFTVTE



#### **Next Generation Sequencing**

massive parallel sequencing = high-throughput sequencing = next-generation sequencing

→Roche's 454:
 400,000 reads (sequenced fragments)
 of 250–400bp length in up to 16 lanes

→Illumina's Solexa:
 30–60 million reads in one run in 8 lanes
 (a lane may be a sample) of 25-72bp long
 ⇒ higher coverage



FIGURE 1: ILLUMINA GENOME ANALYZER FLOW CELL



BIOINF

An amplification (vertical line) in chromosome 19 detected by BAC arrays

1.0 **Density Difference** 0.5 0.0 10635276 21270551 31905826 42541101 53176376 6381165<sup>-</sup> 0 Location

chr19 of Hapmap NA18947

## **Next Generation Sequencing**

BIOINF

Unexplained



#### chr2 of Hapmap NA18947

#### **Nucleosome Position**



Prof. Douglas Murray Institute for Advanced Biosciences





Institute for Advanced Biosciences, Keio University




#### **Nucleosome Position**





Division of Molecular Biology Chromatin- and Epigenetics Laboratory a. Prof. Alexandra Lusser

#### MEDIZINISCHE UNIVERSITÄT Innsbruck







Institute for Theoretical Chemistry Prof. Dr. Peter Schuster Rainer Machne

#### **Nucleosome Position**





#### **Nucleosome Position**



Isolation of natural nucleosome DNAs





# We need help from other disciplines in computational science

Next generation sequencing:

Task1: read mapping hg18 against hg18

- about 100,000,000,000 reads of length 32 bp
- 64 runs
- one run: memory of 20GB, 50 processes simultaneously, 20h
- 53 days of computation

Task2: read mapping 1000Genomes data

- 10,000,000 to 50,000,000 reads per data set
- 836 data sets
- one data set: memory of 2GB, processes simultaneously, 8h
- 278 days of computation







Find substrings in strings

can results of polynom factorization help?

**HELP** 

→ Symbolic Computation



Estimate posterior for latent variable models:

$$p(\boldsymbol{z} \mid \{\boldsymbol{x}\}) = \frac{p(\{\boldsymbol{x}\} \mid \boldsymbol{z}) \ p(\boldsymbol{z})}{\int p(\{\boldsymbol{x}\} \mid \boldsymbol{z}) \ p(\boldsymbol{z}) \ d\boldsymbol{z}}$$

#### Evaluate

$$\int p(\{\boldsymbol{x}\} \mid \boldsymbol{z}) \ p(\boldsymbol{z}) \ d\boldsymbol{z}$$

to compute moments of

 $p(\boldsymbol{z} \mid \{\boldsymbol{x}\})$ 

HELP

→Statistics (sampling)
 →Computational Mathematics



BIOINF

lower bound on likelihood

$$\frac{\ln \mathcal{L}(\{\boldsymbol{x}\};\boldsymbol{w})}{\log Q(\boldsymbol{z} \mid \{\boldsymbol{x}\})} = \ln p(\{\boldsymbol{x}\};\boldsymbol{w}) = \ln \int_{Z} p(\{\boldsymbol{x}\},\boldsymbol{z};\boldsymbol{w}) d\boldsymbol{z} =$$

$$\ln \int_{Z} \frac{Q(\boldsymbol{z} \mid \{\boldsymbol{x}\})}{Q(\boldsymbol{z} \mid \{\boldsymbol{x}\})} p(\{\boldsymbol{x}\},\boldsymbol{z};\boldsymbol{w}) d\boldsymbol{z} \ge$$

$$\int_{Z} Q(\boldsymbol{z} \mid \{\boldsymbol{x}\}) \ln \frac{p(\{\boldsymbol{x}\},\boldsymbol{z};\boldsymbol{w})}{Q(\boldsymbol{z} \mid \{\boldsymbol{x}\})} d\boldsymbol{z} =$$

$$\int_{Z} Q(\boldsymbol{z} \mid \{\boldsymbol{x}\}) \ln p(\{\boldsymbol{x}\},\boldsymbol{z};\boldsymbol{w}) d\boldsymbol{z} - \int_{Z} Q(\boldsymbol{z} \mid \{\boldsymbol{x}\}) \ln Q(\boldsymbol{z} \mid \{\boldsymbol{x}\}) d\boldsymbol{z}$$

#### maximize

$$h(\boldsymbol{w}) = \int_{Z} Q(\boldsymbol{z} \mid \{\boldsymbol{x}\}) \ln p(\{\boldsymbol{x}\}, \boldsymbol{z}; \boldsymbol{w}) d\boldsymbol{z}$$
  

$$\ln p(\{\boldsymbol{x}\}, \boldsymbol{z}; \boldsymbol{w}) =$$
  

$$\phi(\|\boldsymbol{x} - \rho(\boldsymbol{w}, \boldsymbol{z})\|) + \ln g(\boldsymbol{w})$$
  
HELP  

$$\rightarrow \text{Computational Mathematics}$$



