

# Fisher Discriminant Analysis

P. Agius – Spring 2008

## Contents

- What is Fisher's discriminant?
- Kernelizing Fisher's discriminant
- Results

P. Agius – Spring 2008

## Some history

- Introduced in 1998 by Tommi Jaakkola
- Named in honor of Sir Ronald Fisher who was an English statistician, evolutionary biologist and geneticist (analysis of variance, maximum likelihood)



P. Agius – Spring 2008

## Fisher discriminant analysis

The Fisher discriminant is a classification function

$$f(x) = \text{sgn}(\langle \phi(x), w \rangle + b)$$

where the weight  $w$  is chosen to maximize the quotient

$$J(w) = \frac{(\mu_w^+ - \mu_w^-)^2}{(\sigma_w^+)^2 + (\sigma_w^-)^2}$$

where  $\mu_w^+$  is the mean of the projection of the positive examples onto the direction  $w$ ,  $\mu_w^-$  the mean for the negative examples, and  $\sigma_w^+$ ,  $\sigma_w^-$  the corresponding standard deviations.

P. Agius – Spring 2008

## FLD with ref to Mika's paper – Fisher Discriminant Analysis with Kernels

$$\begin{aligned}
 \mathcal{X}_1 &= \{x_1^1, \dots, x_{\ell_1}^1\} & \mathcal{X}_2 &= \{x_1^2, \dots, x_{\ell_2}^2\} \\
 \mathcal{X} &= \mathcal{X}_1 \cup \mathcal{X}_2 = \{x_1, \dots, x_\ell\}
 \end{aligned}$$

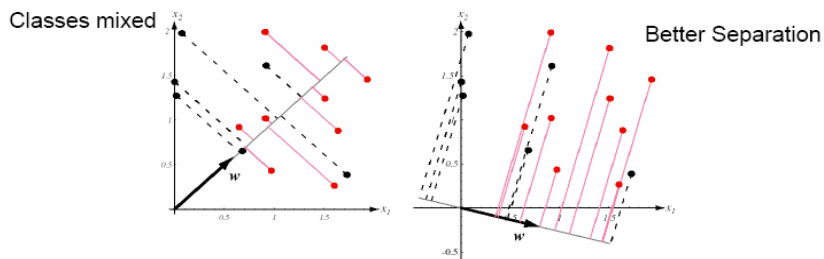
$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

$$\begin{aligned}
 S_B &:= (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \text{ and} & \mathbf{m}_i &:= \frac{1}{\ell_i} \sum_{j=1}^{\ell_i} \mathbf{x}_j^i \\
 S_W &:= \sum_{i=1,2} \sum_{\mathbf{x} \in \mathcal{X}_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T
 \end{aligned}$$

Different notation .... same concept

P. Agius – Spring 2008

### Fisher Liner Discriminant: two-dimensional example



Projection of same set of two-class samples onto two different lines in the direction marked  $w$ .

- Unique solution :
- Easy to compute
- Probabilistic interpretation
- Kernelizable
- Naturally extends to k-class problems

## Probabilistic Interpretation

The Fisher discriminant is the Bayes optimal classifier for two normal distributions with equal covariance.

$$p(y|\omega, \sigma, x) \propto \exp\left(-\frac{1}{\sigma^2}(\omega^T x - y)^2\right) \quad p(\omega) \propto \exp(-c \omega^T \omega)$$

Fisher discriminant analysis can be shown to:

$$\begin{aligned} \text{maximize}_{\omega} \quad & \sum_i \log p(y_i, \omega | x_i, \sigma) & \equiv & \text{minimize}_{\omega, \delta} \quad \|\delta\|_2^2 + c \|\omega\|_2^2 \\ \text{subject to} \quad & \sum_{i \in \{i|y_i=1\}} (\omega^T x_i - y_i) = 0 & & \text{subject to} \quad \delta_i = \omega^T x_i - y_i, \forall i \\ & \sum_{i \in \{i|y_i=-1\}} (\omega^T x_i - y_i) = 0 & & \sum_{i \in \{i|y_i=1\}} \delta_i = 0 \quad \sum_{i \in \{i|y_i=-1\}} \delta_i = 0 \end{aligned}$$

## Colonel Trick

From the theory of reproducing kernels, we know that  $w$  must lie in the *span* of the training examples. Therefore, we can write  $w$  as

$$w = \sum_{i=1}^{\ell} \alpha_i \Phi(x_i)$$

Recall that

$$m_i^{\Phi} := \frac{1}{\ell_i} \sum_{j=1}^{\ell_i} \Phi(x_j^i)$$

Therefore

$$\begin{aligned} w^T m_i^{\Phi} &= \frac{1}{\ell_i} \sum_{j=1}^{\ell} \sum_{k=1}^{\ell_i} \alpha_j k(x_j, x_k^i) \\ &= \alpha^T M_i \quad \text{where} \quad (M_i)_j := \frac{1}{\ell_i} \sum_{k=1}^{\ell_i} k(x_j, x_k^i) \end{aligned}$$

Recall 
$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_B^\Phi \mathbf{w}}{\mathbf{w}^T S_W^\Phi \mathbf{w}}$$

The numerator and denominator now become

$$\mathbf{w}^T S_B^\Phi \mathbf{w} = \boldsymbol{\alpha}^T M \boldsymbol{\alpha} \quad \mathbf{w}^T S_W^\Phi \mathbf{w} = \boldsymbol{\alpha}^T N \boldsymbol{\alpha}$$

$$M := (M_1 - M_2)(M_1 - M_2)^T \quad N := \sum_{j=1,2} K_j (I - \mathbf{1}_{\ell_j}) K_j^T$$

$$(M_i)_j := \frac{1}{\ell_i} \sum_{k=1}^{\ell_i} k(\mathbf{x}_j, \mathbf{x}_k^i) \quad (K_j)_{nm} := k(\mathbf{x}_n, \mathbf{x}_m^j)$$

Therefore 
$$J(\boldsymbol{\alpha}) = \frac{\boldsymbol{\alpha}^T M \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T N \boldsymbol{\alpha}}$$

P. Agius – Spring 2008

## Kernel Fisher Discriminant

Solve the problem

$$J(\boldsymbol{\alpha}) = \frac{\boldsymbol{\alpha}^T M \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T N \boldsymbol{\alpha}}$$

by finding the leading eigenvector of  $N^{-1}M$

**One regularization method** that leads to more stable results when solving for the leading eigenvector involves replacing  $N$  by

$$N_\mu := N + \mu I$$

where  $I$ =identity matrix  
and  
 $\mu$  is a positive parameter  
inducing a positive semi-  
definite matrix

P. Agius – Spring 2008

## Another Regularized Formulation

Choose  $w$  that solves the following optimization problem:

$$\max_w J(w) = \frac{(\mu_w^+ - \mu_w^-)^2}{(\sigma_w^+)^2 + (\sigma_w^-)^2 + \lambda \|w\|^2}$$

Kernel Methods for Pattern Analysis, Ch 5

P. Agius – Spring 2008

## Results – Mika paper

	RBF	AB <sup>-</sup>	AB <sub>R</sub>	SVM	KFD
Banana	<b>10.8±0.6</b>	12.3±0.7	10.9±0.4	11.5±0.7	<b>10.8±0.5</b>
B.Cancer	27.6±4.7	30.4±4.7	26.5±4.5	26.0±4.7	<b>25.8±4.6</b>
Diabetes	24.3±1.9	26.5±2.3	23.8±1.8	23.5±1.7	<b>23.2±1.6</b>
German	24.7±2.4	27.5±2.5	24.3±2.1	<b>23.6±2.1</b>	23.7±2.2
Heart	17.6±3.3	20.3±3.4	16.5±3.5	<b>16.0±3.3</b>	16.1±3.4
Image	3.3±0.6	<b>2.7±0.7</b>	<b>2.7±0.6</b>	3.0±0.6	4.8±0.6
Ringnorm	1.7±0.2	1.9±0.3	1.6±0.1	1.7±0.1	<b>1.5±0.1</b>
F.Sonar	34.4±2.0	35.7±1.8	34.2±2.2	<b>32.4±1.8</b>	33.2±1.7
Splice	10.0±1.0	10.1±0.5	<b>9.5±0.7</b>	10.9±0.7	10.5±0.6
Thyroid	4.5±2.1	4.4±2.2	4.6±2.2	4.8±2.2	<b>4.2±2.1</b>
Titanic	23.3±1.3	22.6±1.2	22.6±1.2	<b>22.4±1.0</b>	23.2±2.0
Twonorm	2.9±0.3	3.0±0.3	2.7±0.2	3.0±0.2	<b>2.6±0.2</b>
Waveform	10.7±1.1	10.8±0.6	<b>9.8±0.8</b>	9.9±0.4	9.9±0.4

RBF = RBF classifier

AB = Adaboost, ABR = regularized Adaboost

KFD = Kernel Fisher Discriminant

P. Agius – Spring 2008

Using the Fisher kernel method to detect  
remote protein homologies  
(Jaakkola, 1999)

Goal:

- Detect remote protein homologies
- Sequence-based tools
- SCOP database
- Use HMM for each family/superfamily

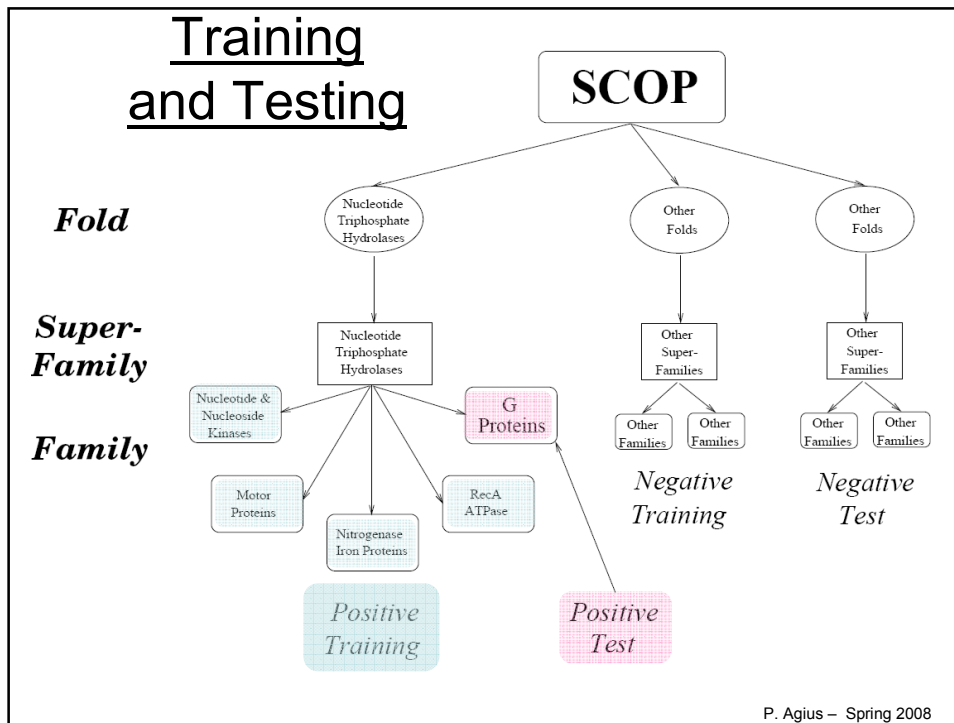
P. Agius – Spring 2008

## Jaakkola 1999

Method:

- Derive a kernel function using HMMs
- Kernel function specifies a similarity score for any pair of protein sequences (in contrast to the HMM likelihood which is a measure of the closeness of a sequence to the HMM model)
- Derive a Fisher score

P. Agius – Spring 2008



## Methodology details

$X = [x_1, \dots, x_n]$  is a protein sequence (amino acids).

Build a HMM  $H_1$  for a particular family of proteins

Then  $P(X|H_1)$  is the probability of  $X$  under  $H_1$  and  $P(X|H_0)$  is the probability of  $X$  under the null model.

The score used in the database search is defined by the log likelihood of  $X$

$$\begin{aligned} \mathcal{L}(X) &= \log \frac{P(X|H_1)P(H_1)}{P(X|H_0)P(H_0)} \\ &= \log \frac{P(X|H_1)}{P(X|H_0)} + \log \frac{P(H_1)}{P(H_0)} \end{aligned}$$

Positive  $\mathcal{L}(x) \rightarrow x$  is a member of the family

P. Agius – Spring 2008



## The Fisher kernel method

$$\begin{aligned}\mathcal{L}(X) &= \log P(H_1|X) - \log P(H_0|X) & (2) \\ &= \sum_{i: X_i \in H_1} \lambda_i K(X, X_i) - \sum_{i: X_i \in H_0} \lambda_i K(X, X_i)\end{aligned}$$

The sign of the discriminant function determines the assignment of the sequences into hypothesis classes.

The Fisher score is  $U_X = \nabla_{\theta} \log P(X|H_1, \theta)$

Set of model parameters

where  $U_X$  is the derivative of the loglikelihood for sequence  $X$  with respect to a particular parameter (hence determines contribution of parameter)

The kernel function is then defined to be

$$K(X, X') = e^{-\frac{1}{2\sigma^2}(U_X - U_{X'})^T (U_X - U_{X'})}$$

P. Agius – Spring 2008

## Summary of the method

- Derive HMM trained on positive examples
- Use HMM to map a sequence  $X$  to a feature vector (its Fisher score)
- Compute the kernel function
- The resulting discriminant function is:

$$\mathcal{L}(X) = \sum_{i: X_i \in H_1} \lambda_i K(X, X_i) - \sum_{i: X_i \in H_0} \lambda_i K(X, X_i),$$

This is the SVM-Fisher method

P. Agius – Spring 2008

## Other comparative methods

- BLAST – family pairwise search; database is queried with each positive training sequence
- SAMT-98 – builds HMM for a SCOP domain sequence, then iteratively selects positive training examples from among potential homologs and refines the model
- BLAST + SAM
- SVM-Fisher

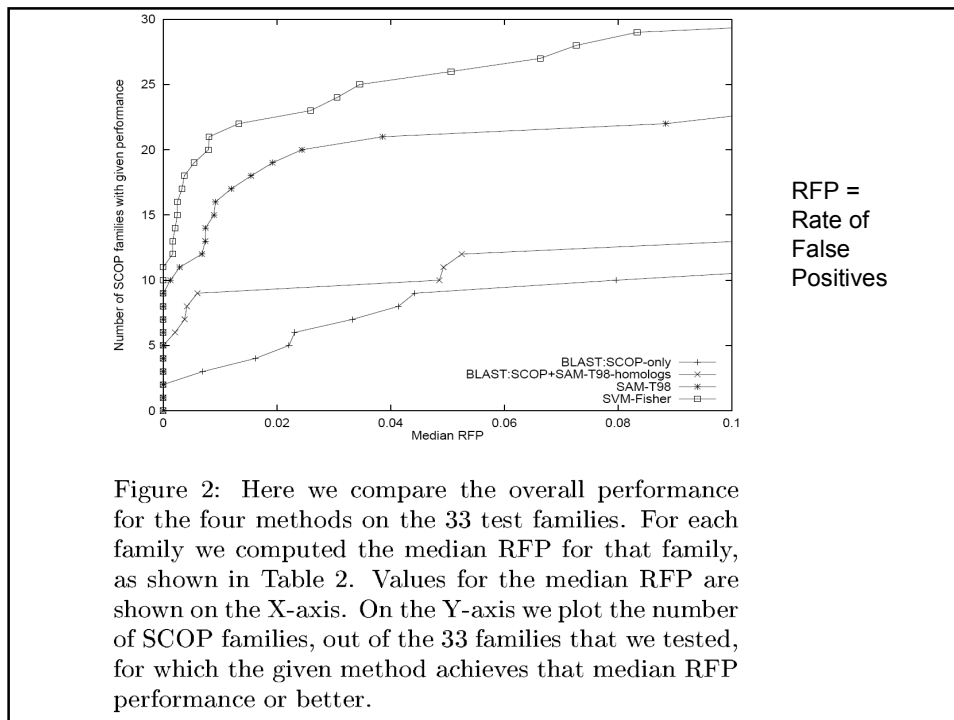
P. Agius – Spring 2008

## Results

Sequence	BLAST	B-Hom	S-T98	SVM-F
5p21	0.043	0.010	0.001	0.000
1guaA	0.179	0.031	0.000	0.000
1etu	0.307	0.404	0.428	0.038
1hurA	0.378	0.007	0.007	0.000
1eft(3)	0.431	0.568	0.041	0.051
1dar(2)	0.565	0.391	0.289	0.019
1tadA(2)	0.797	0.330	0.004	0.000
1gia(2)	0.867	0.421	0.017	0.000

Table 1: Rate of false positives for *G proteins* family. BLAST = BLAST:SCOP-only, B-Hom = BLAST:SCOP+SAMT-98-homologs, S-T98 = SAMT-98, and SVM-F = SVM-Fisher method.

P. Agius – Spring 2008



## Future work and other

- Which sequences should be used to estimate the HMM parameters?
- Which sequences should be left for the discriminative function?
- Can build generic models by building HMMs on short amino acid sequences that map to conserved regions in the proteins
- Extend the method to identify multiple domains within large protein sequences
- Fisher's discriminant for multiple classes (Ch 4, Chris Bishop)