

A Hybrid Kernel Machine

for

Protein Secondary Structure Prediction

Yann Guermeur

LORIA

Université Nancy I

Alain Lifchitz

LIP6

CNRS

Régis Vert

LRI

Université Paris XI

<http://www.loria.fr/~guermeur/>

Overview

Multi-class support vector machines (M-SVMs)

- One single architecture
- Several training algorithms
- M-SVMs and the SRM inductive principle

M-SVMs for protein secondary structure prediction

- Improving the generalization performance of Spro2
- Hybrid and modular architecture : M-SVM + IHMM
- Experimental results

Multi-class pattern recognition

Hypotheses : empirical data characterizing a joint probability distribution

- Q -category discrimination problem
- $Z = (X, Y)$: random variable on a probability space (Ω, \mathcal{B}, P)
- $X(\Omega) = \mathcal{X}$: input space (set of descriptions), $Y(\Omega) = \mathcal{Y}$: finite set of categories
- P : joint probability distribution function on $\mathcal{X} \times \mathcal{Y}$, fixed but unknown
- $s = \{(x_1, y_1), \dots, (x_m, y_m)\} \subset (\mathcal{X} \times \mathcal{Y})^m$, learning set : observations i.i.d. according to P
- \mathcal{H} : family of vector-valued functions $h = [h_k], (1 \leq k \leq Q)$, from \mathcal{X} into \mathbb{R}^Q

Goal : for a given pattern, find its category

Find in \mathcal{H} a function associated with the lowest expected risk (generalization error)

$$R(h) = R(f) = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{I}_{\{f(x) \neq y\}} dP(x, y)$$

f : discriminant function corresponding to h , obtained by choosing the category associated with the index of the highest output

Multi-class Support Vector Machines

Architecture

Functions $h = [h_k]$ computed by the architecture are defined by :

$$\forall x \in \mathcal{X}, \forall k \in \{1, \dots, Q\}, h_k(x) = h_k(x) + b_k = \langle w_k, \Phi(x) \rangle + b_k$$

where

- Φ is either a nonlinear map into a feature space or identity
- $\tilde{\mathcal{H}} = \{ \tilde{h} = [\tilde{h}_k] \}$ is the product of Q reproducing kernel Hilbert spaces (RKHS)

- K , the kernel associated with $\tilde{\mathcal{H}}$, is related to Φ through :

$$\langle \Phi(x_{(1)}), \Phi(x_{(2)}) \rangle = K(x_{(1)}, x_{(2)}) \quad \forall x_{(1)}, x_{(2)} \in \mathcal{X}_2$$

Training algorithm of M-SVMs : general principle

Let $s_m = \{(x_1, C(x_1)), \dots, (x_m, C(x_m))\}$ be the training set

Primal formulation

Problem 1

$$\min_{h \in \mathcal{H}} \left\{ \frac{1}{2} \| \tilde{h} \|_{\mathcal{H}}^2 + C \sum_{i=1}^m l(h(x_i), C(x_i)) \right\}$$

Dual formulation

Representer theorems establish that training amounts to finding the values of the coefficients β_{ik} in :

$$\forall k \in \{1, \dots, Q\}, h_k(x) = \sum_{i=1}^m \beta_{ik} K(x_i, x) + b_k$$

Training algorithms of M-SVMs : primal formulation

Problem 2 (M-SVM1 (Vapnik & Blanz 98, Weston & Watkins 98,...))

$$\min_{h \in \mathcal{H}} \left\{ \frac{1}{Q} \sum_{k=1}^K \|w_k\|_2 + C \sum_{m=1}^m \sum_{k=1}^K \xi_{ik} \right\}$$

$$\left. \begin{aligned} & \langle w_{C(x_i)} - w_k, \Phi(x_i) \rangle + b_{C(x_i)} - b_k \geq 1 - \xi_{ik}, \quad (1 \leq i \leq m), (1 \leq k \neq C(x_i) \leq Q) \\ & \xi_{ik} \geq 0, \quad (1 \leq i \leq m), (1 \leq k \leq Q) \end{aligned} \right\} \text{s.t.}$$

Problem 3 (M-SVM2 (Guerneur 02))

$$\min_{h \in \mathcal{H}} \left\{ \frac{1}{2} t^2 + C \sum_{m=1}^m \sum_{k=1}^K \xi_{ik} \right\}$$

$$\left. \begin{aligned} & \|w_k - w_l\|_2 \leq t^2, \quad (1 \leq k < l \leq Q) \\ & \text{Constraints of Problem 2} \end{aligned} \right\} \text{s.t.}$$

Training algorithms of M-SVMs : primal formulation

Problem 4 (M-SVM3 (Crammer & Singer 01))

$$\min_{h \in \mathcal{H}} \left\{ \frac{1}{Q} \sum_{k=1}^K \|w_k\|_2 + C \sum_{i=1}^m \xi_i \right\}$$

$$s.t. \langle w_{C(x_i)} - w_k, \Phi(x_i) \rangle + b_{C(x_i)} - b_k + \delta_{C(x_i), k} \geq 1 - \xi_i, (1 \leq i \leq m), (1 \leq k \leq Q)$$

Problem 5 (M-SVM4 (Lee, Lin & Wahba 01))

$$\min_{h \in \mathcal{H}} \left\{ \frac{1}{Q} \sum_{k=1}^K \|w_k\|_2 + C \sum_{i=1}^m \sum_{k=1}^Q \xi_{ik} \right\}$$

$$\left. \begin{aligned} \sum_{k=1}^Q w_k = 0, \quad \sum_{k=1}^Q b_k = 0 \\ \langle w_k, \Phi(x_i) \rangle + b_k \leq -1/(Q-1) + \xi_{ik}, \quad (1 \leq i \leq m), (1 \leq k \neq C(x_i) \leq Q) \\ \xi_{ik} \geq 0, \quad (1 \leq i \leq m), (1 \leq k \leq Q) \end{aligned} \right\} s.t.$$

Behaviour of a M-SVM

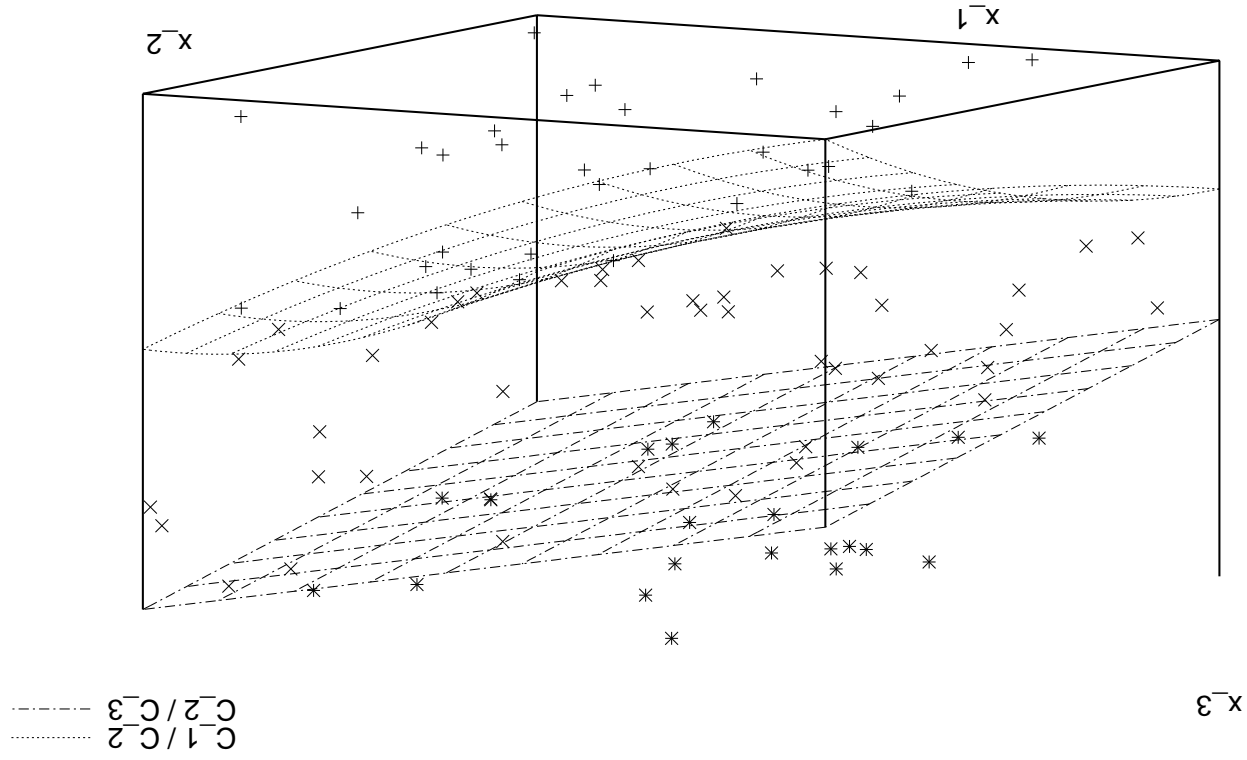


FIG. 1 – 3 categories non-linearly separable in 3D

Behaviour of a M-SVM

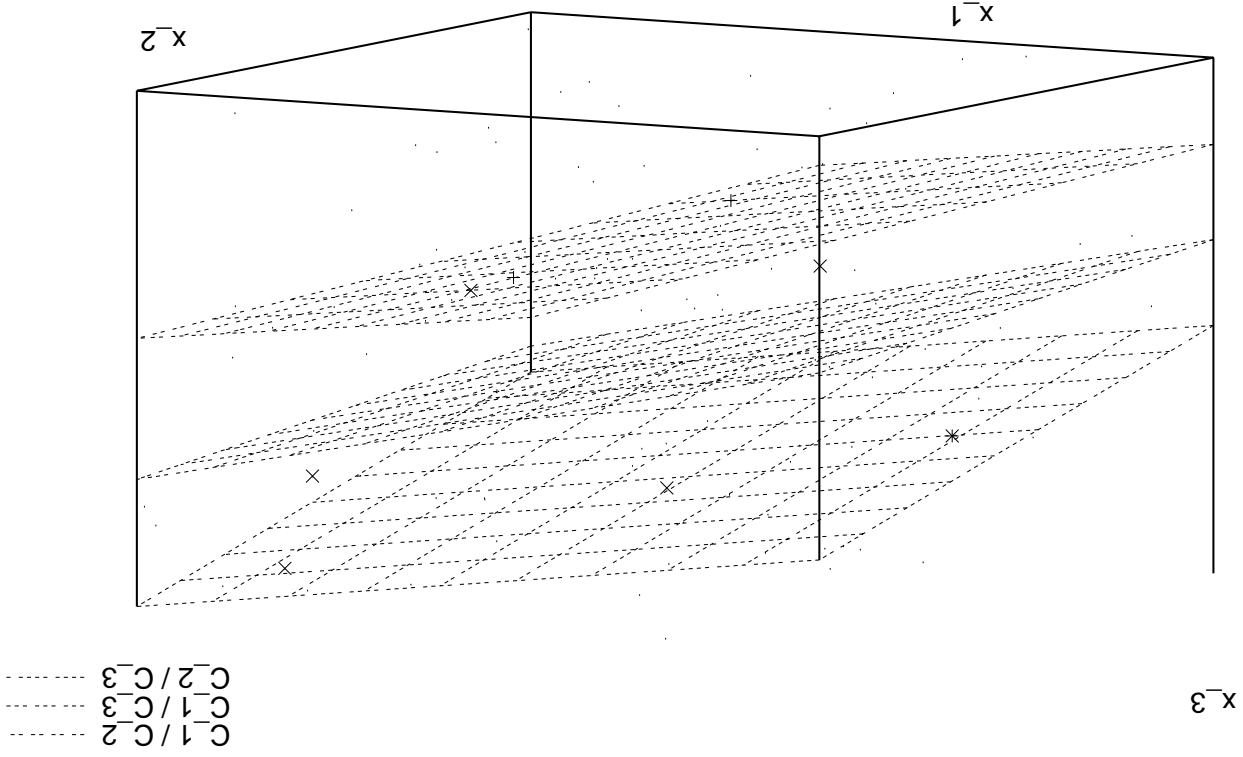


FIG. 2 – Hyperplanes and support vectors of a linear M-SVM

Empirical margin risk and uniform convergence result : multi-class case

For $\gamma \in (0, 1]$, let $\pi_\gamma : \mathbb{R} \rightarrow [-\gamma, \gamma]$ be the piecewise-linear squashing function defined as

$$\pi_\gamma(x) = \begin{cases} x & \text{if } |x| \leq \gamma \\ \gamma \cdot \text{sign}(x) & \text{otherwise} \end{cases}$$

Definition 1 (Canonical function) Let $h = [h_k]$ be a function from \mathcal{X} . Its canonical function $\Delta h = [\Delta h_k]$, ($1 \leq k \leq Q$), is the function from \mathcal{X} into \mathbb{R}^Q satisfying

$$\Delta h_k(x) = \frac{1}{2} \left\{ h_k(x) - \max_{l \neq k} h_l(x) \right\}$$

Definition 2 (Empirical margin risk (Elisseeff & al. 99)) The empirical risk with margin $\gamma \in (0, 1]$ of h on a set $s_m = \{(x_1, C(x_1)), \dots, (x_m, C(x_m))\}$ of size m is

$$R_\gamma^{s_m}(h) = \frac{1}{m} \left| \{(x_i, C(x_i)) \in s_m / \Delta h_{C(x_i)}(x_i) > \gamma\} \right|$$

Empirical margin risk and uniform convergence result : multi-class case

Theorem 1 (Elisseeff & al. 99) *Let s_m be a m -sample of examples drawn independently from P . With probability at least $1 - \delta$, for every value of γ in $(0, 1]$, the risk $R(h)$ of a function h computed by a numerical Q -class discriminant model \mathcal{H} is bounded above by :*

$$R(h) \leq \underbrace{R_{\gamma}^{s_m}(h)} + \sqrt{\frac{1}{2m} \left(\ln(2N_{\infty, \infty}(\gamma/2, \Delta\mathcal{H}_{\gamma}, 2m)) + \ln\left(\frac{\gamma\delta}{2}\right) \right)} + \frac{1}{m}$$

where $\Delta h_{\gamma} = [\pi_{\gamma} \circ \Delta h_k]$, $(1 \leq k \leq Q)$, $\Delta\mathcal{H}_{\gamma} = \{\Delta h_{\gamma} / h \in \mathcal{H}\}$

$$\forall s_m \in \mathcal{X}_m, \forall (h_{(1)}, h_{(2)}) \in \mathcal{H}_2, d_{l_{\infty, l_{\infty}}(s_m)}(h_{(1)}, h_{(2)}) = \max_{x_i \in s_m} \max_{k \in \{1, \dots, Q\}} \left| h_{k(1)}(x_i) - h_{k(2)}(x_i) \right|$$

$$N_{\infty, \infty}(\gamma/2, \Delta\mathcal{H}_{\gamma}, 2m) = \max_{s_{2m} \in \mathcal{X}_{2m}} N(\gamma/2, \Delta\mathcal{H}_{\gamma}, d_{l_{\infty, l_{\infty}}(s_{2m})})$$

Extended notions of VC dimension

Definition 3 (Fat-shattering dimension (Kearns & Schapire 90)) Let \mathcal{H} be a set of real-valued functions on a set \mathcal{X} . For $\gamma > 0$, a subset $s_m = \{x_i\}$, $(1 \leq i \leq m)$ of \mathcal{X} is said to be γ -shattered by \mathcal{H} if there is a vector $v_b = [b_i] \in \mathbb{R}^m$ such that, for each binary vector $v_y = [y_i] \in \{-1, 1\}^m$, there is a function $h_y \in \mathcal{H}$ satisfying

$$(h_y(x_i) - b_i) y_i \geq \gamma, \quad (1 \leq i \leq m)$$

The vector v_b is then said to witness the γ -shattering of s_m by \mathcal{H} . The fat-shattering dimension $\text{fat}_\gamma \mathcal{H}$ of the set \mathcal{H} is a function from the positive real numbers to the integers which maps a value γ to the size of the largest set γ -shattered by functions of \mathcal{H} , if this size is finite, or to infinity otherwise.

Definition 4 (Graph dimension (Dudley 87, Natarajan 89)) Let \mathcal{H} be a set of functions on a set \mathcal{X} taking their values in a countable set. For any $h \in \mathcal{H}$, the graph \mathcal{G} of h is $\mathcal{G}(h) = \{(x, h(x)) \mid x \in \mathcal{X}\}$ and the graph space of \mathcal{H} is $\mathcal{G}(\mathcal{H}) = \{\mathcal{G}(h) \mid h \in \mathcal{H}\}$. Then the graph dimension of \mathcal{H} is defined to be the VC dimension of the space $\mathcal{G}(\mathcal{H})$.

M-fat-shattering dimension

Definition 5 (M-fat-shattering dimension) Let \mathcal{H} be a set of functions on a set X taking their values in \mathbb{R}^Q . For $\gamma > 0$, a subset $s_m = \{x_i\}$, $(1 \leq i \leq m)$ of X is said to be M- γ -shattered by \mathcal{H} if there is a vector $v_i = [b_i] \in \mathbb{R}^m$ and a vector $v_c = [c_i] \in \{1, \dots, Q\}^m$ such that, for each binary vector $v_y = [y_i] \in \{-1, 1\}^m$, there is a function $h_y = [h_{y_k}]$, $(1 \leq k \leq Q) \in \mathcal{H}$ satisfying

$$(h_{y_{c_i}}(x_i) - b_i) y_i \geq \gamma, \quad (1 \leq i \leq m)$$

The couple (v_i, v_c) is then said to witness the M- γ -shattering of s_m by \mathcal{H} . The M-fat-shattering dimension $M\text{-fat}_{\gamma}(\mathcal{H})$ of the set \mathcal{H} is a function from the positive real numbers to the integers which maps a value γ to the size of the largest set M- γ -shattered by functions of \mathcal{H} , if this size is finite, or to infinity otherwise.

M-fat-shattering dimension : extension of the fat-shattering dimension to the multivariate case and scale-sensitive version of the graph dimension

Generalization of Sauer's lemma : multi-class case

Theorem 2 Let \mathcal{H} be a set of functions from X into \mathbb{R}^Q . For every value of γ in $(0, 1]$ and every value of m in \mathbb{N}^* , the following bound is true :

$$N_{\infty, \infty}(\gamma/2, \Delta\mathcal{H}^\gamma, 2m) \leq 2(2mQ9^Q)^{d \log_2(18emQ/d)}$$

where $d = M\text{-fat}_{\Delta\mathcal{H}^\gamma}(\gamma/8)$.

Theorem 3 Let \mathcal{H} be a set of vector-valued functions $h = [h_k]$, $(1 \leq k \leq Q)$, from a set X into \mathbb{R}^Q . Let \mathcal{H}_k , $(1 \leq k \leq Q)$, be the sets of real-valued functions h_k corresponding to the different components of the functions h . Then the following bound holds true for all positive value of ϵ :

$$M\text{-fat}_{\mathcal{H}^\epsilon}(\epsilon) \leq \sum_{k=1}^Q \text{fat}_{\mathcal{H}_k}(\epsilon)$$

Main difficulty : $M\text{-fat}_{\Delta\mathcal{H}^\epsilon}$ cannot be bounded in terms of $M\text{-fat}_{\mathcal{H}^\epsilon}$

Dependence of the capacity on the control term

Theorem 4 Let \mathcal{H} be the set of vector-valued functions $h = [h_k]$, ($1 \leq k \leq Q$), computed by a M -SVM. For k in $\{1, \dots, Q\}$, let $\Delta\mathcal{H}_k$ be the set of real-valued functions Δh_k corresponding to the k^{th} component of the functions Δh :

$$\Delta h_k(x) = \frac{1}{2} \min_{l \neq k} \{ \langle w_k - w_l, \Phi(x) \rangle + b_k - b_l \}$$

Suppose that $\Phi(x)$ is included in a ball of radius $\Delta\Phi(x)$ and that the vectors w_k satisfy $\max_{1 \leq k < l \leq Q} \|w_k - w_l\|_2 \leq \Delta_w$. Then

$$\forall k \in \{1, \dots, Q\}, \text{fat}_{\mathcal{H}_k}(\epsilon) = O \left(\left(\frac{\epsilon}{\Delta\Phi(x)\Delta_w} \right)^2 \right)$$

Simple pathway to bound the covering numbers

Theorem 5 Let \mathcal{H} be a set of functions from X into \mathbb{R}^Q . For every value of ϵ and γ satisfying $0 < \epsilon < \gamma$, the following bounds hold true :

$$N_{\infty, \infty}(\epsilon, \Delta \mathcal{H}^\gamma, m) \leq N_{\infty, \infty}(\epsilon, \Delta \mathcal{H}, m) \leq N_{\infty, \infty}(\epsilon, \mathcal{H}, m)$$

Theorem 6 Let \mathcal{H} be a set of functions computed by a M-SVM. Suppose that the biases b_k all belong to the interval $[-b, b]$ Then :

$$N_{\infty, \infty}(\epsilon, \mathcal{H}, m) \leq \left\lceil \frac{\epsilon}{2b} \right\rceil N_{\infty, \infty}(2\epsilon, \tilde{\mathcal{H}}, m)$$

Covering numbers and entropy numbers

Definition 6 (Entropy numbers) Let (E, ρ) be a pseudo-metric space and H a subset of E . Then, for n in \mathbb{N}^* , the n^{th} entropy number of H , $\epsilon_n(H)$, is :

$$\epsilon_n(H) = \inf \{ \epsilon > 0 / \mathcal{N}(\epsilon, H, \rho) \leq n \}$$

Definition 7 (Entropy numbers of a bounded linear operator) Let E_H and F_H be Hilbert spaces endowed with the norms $\|\cdot\|_{E_H}$ and $\|\cdot\|_{F_H}$. Let $\mathcal{L}(E_H, F_H)$ be the set of all bounded linear operators between $(E_H, \|\cdot\|_{E_H})$ and $(F_H, \|\cdot\|_{F_H})$. Let $T \in \mathcal{L}(E_H, F_H)$.

$$\|T\| = \sup_{x \in U_{E_H}} \|Tx\|_F \quad \epsilon_n(T) = \epsilon_n(T(U_{E_H}))$$

Theorem 7 (Maurey) Let $T \in \mathcal{L}(E_H, \mathbb{R}^m)$. Then, there exists a constant c such that, for all $n \in \mathbb{N}^*$,

$$\epsilon_n(T) \leq c \|T\| \left((\log(n) + 1)^{-1} \log \left(1 + \frac{\log(n) + 1}{m} \right) \right)^{1/2}$$

Covering numbers and entropy numbers

Theorem 8 Let $\tilde{\mathcal{H}} = \{h\}$ be defined as above, with the additional constraint that each

function $h = [w_k], (1 \leq k \leq Q)$, satisfies : $\max_k \|w_k\|_2 \leq \Lambda_w$

Let T_{s_m} be the linear operator given by $T_{s_m} = \Lambda_w S_{s_m}$ with :

$$S_{s_m} : l_2 \rightarrow l_2^\infty : w \mapsto [\langle w, \Phi(x_1) \rangle, \dots, \langle w, \Phi(x_m) \rangle, \dots, \langle w, \Phi(x_m) \rangle]$$

Then

$$\epsilon_n(T_{s_m}) \leq \epsilon_0 \iff \mathcal{N}_{\infty, \infty}(\epsilon_0, \tilde{\mathcal{H}}, m) \leq n$$

Constraints on the hyperplanes : $\|w_k\|_2^2$ or $\|w_k - w_l\|_2^2$?

Multi-class SVM	Objective function	Add. const.
Vapnik & Blanz 98	$J_1(w, b, \xi) = \sum_{k=1}^Q \ w_k\ _2^2 + C_1 1^T \xi$	-
Weston & Watkins 98	$J_1(w, b, \xi) = \sum_{k=1}^Q \ w_k\ _2^2 + C_1 1^T \xi$	-
Bredensteiner & al. 99	$J_2(w, b, \xi) = \sum_{Q^{k>l}} \ w_k - w_l\ _2^2 + \sum_{k=1}^Q \ w_k\ _2^2 + C_2 1^T \xi$	-
Guerneur & al. 00	$J_3(w, b, \xi) = \sum_{Q^{k>l}} \ w_k - w_l\ _2^2 + C_3 1^T \xi$	$\sum_{k=1}^Q w_k = 0^p$

Objective function	Add. const.	C	Solution
$J_1(w, b, \xi)$	-	C_1	$\left(w_{(1)}, b_{(1)}, \xi_{(1)}, \alpha_{(1)}, \beta_{(1)} \right)$
$J_2(w, b, \xi)$	-	$C_1 + \hat{Q}$	$\left(w_{(1)}, b_{(1)}, \xi_{(1)}, \alpha_{(1)}, \beta_{(1)}, \hat{Q} \right)$
$J_3(w, b, \xi)$	$\sum_{k=1}^Q w_k = 0^p$	$\hat{Q} C_1$	$\left(w_{(1)}, b_{(1)}, \xi_{(1)}, \alpha_{(1)}, \beta_{(1)}, \hat{Q}, 0^p \right)$

The same set of primal variables generates solutions for the three problems \implies All these multi-class SVMs are equivalent

Protein secondary structure prediction

Basic notions about proteins

- Proteins : macromolecules made up of amino acids
- 20 amino acids, each of them represented by a letter (A, R, N, D, C, E, ...)

Hierarchical description of the conformation

- Primary structure (sequence of amino acids) \Rightarrow sequencing
- Secondary structure (sequence of structural elements) \Rightarrow circular dichroism
- Tertiary structure (three-dimensional structure) \Rightarrow X-ray, NMR
- ...

Importance of structure prediction

- Knowing the structure is a prerequisite to gain a thorough understanding of the function
- Difficulty to determine the structure experimentally

Modular and hierarchical approach of the prediction

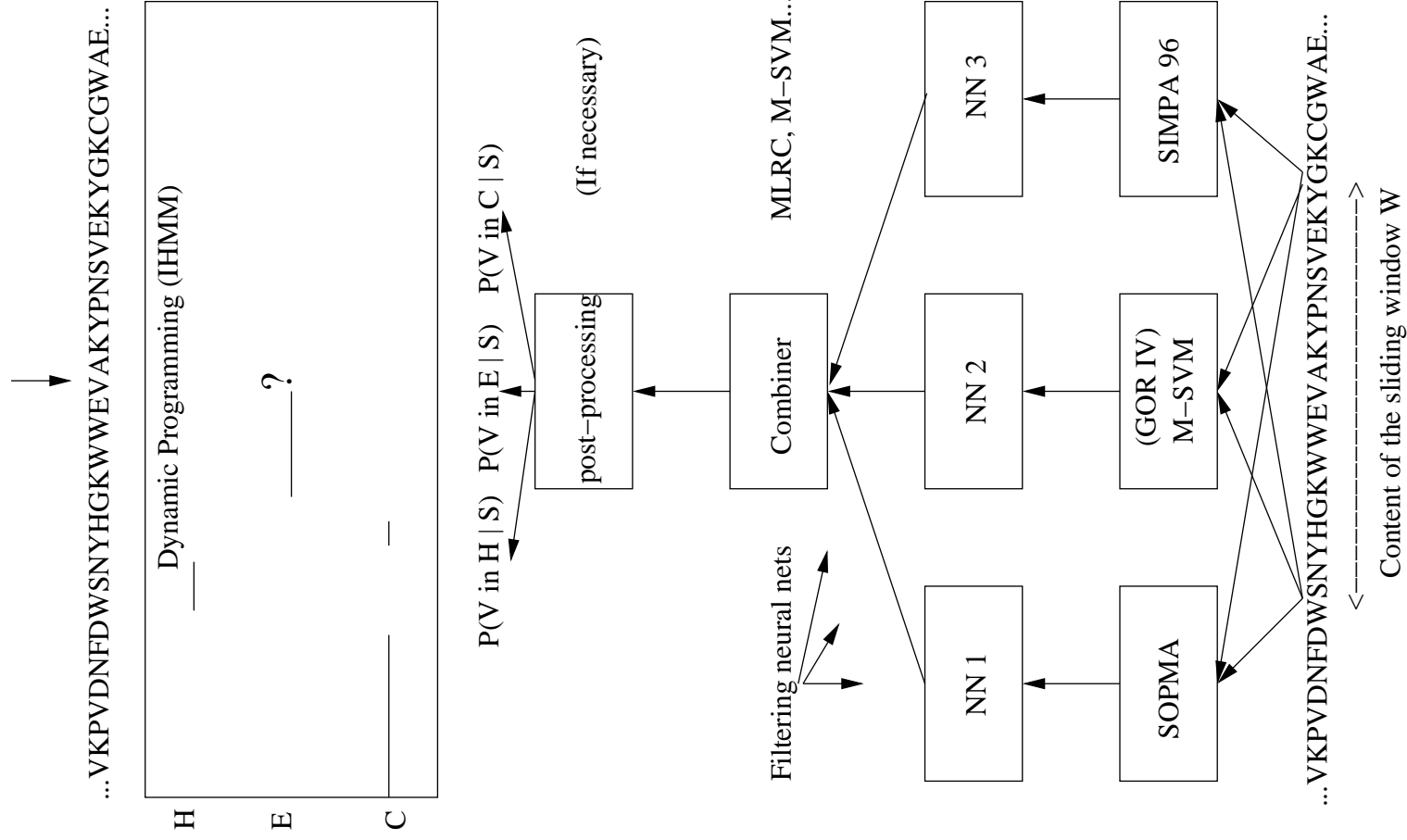


Fig. 3 – Hierarchical architecture for protein secondary structure prediction

Combination of the BRNNs of SSpro2

Collaboration with the biocomputing group of Prof. P. Baldi at the UCI

SSpro2 (Pollastri & al. 02)

- Data : profiles of alignments derived from PSI-BLAST
- Experts : 11 'Bidirectional Recurrent Neural Networks (BRNNs)
- SSpro2 : unweighted average on the outputs of the BRNNs

Goal

Improve performance by implementing a superior ensemble method \Rightarrow M-SVMs!

Experimental protocol

- New database of **1096 protein sequences (255551 amino acids)** exhibiting no homology with the training set of the BRNNs
- Training of the M-SVM + post-processing (perceptron with *softmax* units) + IHMM : *stacked generalization*

Different measures of prediction accuracy

Q_3 : recognition rate at the residue level

Pearson's/Matthews' correlation coefficients

$$C_i = \frac{\sqrt{(p_i + u_i)(d_i + o_i)(n_i + u_i)(n_i + o_i)}}{d_i n_i - u_i o_i}$$

Root mean square deviation (r.m.s.d.)

$$\sigma_i = \sqrt{\frac{1}{n_s} \sum_{j=n_s}^{j=1} (obs_{ij} - pred_{ij})^2}$$

Sov coefficients (Rost & al. 94, Zemla & al. 99)

$$Sov(\delta) = \frac{\sum_{S_1} \frac{1}{n_{S_1}} \left\{ \sum_{S_2/S_1 \cup S_2 \neq \emptyset} \frac{1}{n_{S_2}} \right.}{\left. \frac{\min(end(S_1), end(S_2)) - \max(beg(S_1), beg(S_2)) + 1 + \delta}{\max(end(S_1), end(S_2)) - \min(beg(S_1), beg(S_2)) + 1 + \delta} \right\}}{len(S_1)}$$

Relative prediction accuracy of combiners

	Av.	MLP	$SV M_{1^{u,r}}$	M-SVM1	M-SVM2	M-SVM3
Q_3	76.94	76.91	77.01	77.09	77.12	77.03
Q_α	86.7	86.7	86.8	86.8	86.9	86.8
Q_β	87.7	87.6	87.3	87.8	87.9	87.6
Q^c	79.5	79.6	79.6	79.6	79.6	79.5
C_α	0.72	0.72	0.71	0.72	0.72	0.72
C_β	0.62	0.63	0.62	0.63	0.63	0.63
C^c	0.58	0.58	0.58	0.58	0.58	0.58
Sov	72.2	72.2	72.0	72.4	72.5	72.3
Sov^α	75.6	76.1	76.0	76.1	76.0	76.0
Sov^β	67.1	69.0	67.3	68.9	69.0	68.6
Sov^c	69.0	67.6	68.1	68.5	68.5	68.4

TAB. 1 – Combination of the 11 BRNNs of SSpro2 with two-class and multi-class SVMs

Taking the PSI-BLAST profiles into account improves performance

	M-LP	$SV M^{\alpha+\beta+c}$	M-SVM1
Q_3	77.02	77.06	77.26
Q_α	86.9	86.9	87.1
Q_β	87.6	87.3	87.8
Q_c	79.5	79.7	79.6
C_α	0.72	0.73	0.73
C_β	0.63	0.62	0.64
C_c	0.58	0.58	0.58
Sov	72.3	72.2	72.5
Sov_α	74.8	74.8	74.6
Sov_β	69.0	68.7	69.4
Sov_c	68.8	68.5	68.7

TAB. 2 – Combination of the 11 BRNNs of SS_{pro2} and PSI-BLAST derived profiles

M-SVM as basic classifier - Shortcomings of standard kernels

Canonical coding of the content of the sliding window

$$- x = [x_{-n}, \dots, x_i, \dots, x_n]^T \in \{0, 1\}^{(2n+1).22} \quad (|W| = 2n + 1)$$

- x_i : canonical coding of the amino acid in position i (binary vector)

Function computed by the kernel

$$K(x, x') = \exp \left(- \frac{2\sigma^2}{\|x - x'\|_2^2} \right) = \exp \left(- \frac{2\sigma^2}{\sum_{i=-n}^n \delta_{x_i, x'_i} - (2n + 1)} \right)$$

The kernel only computes the Hamming distance!

The possibility of insertions/deletions, the nature of the substitutions and the position in the window are not taken into account

M-SVM as basic classifier - Optimization of the kernel

Kernel target alignment : multi-class case (Vert 02)

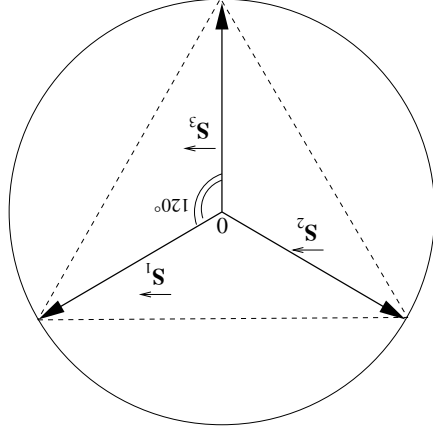


FIG. 5 - Optimal choice of the representatives of the categories

$$K_t(x, x') = \begin{cases} 1 & \text{if } C(x) = C(x') \\ -1/(Q-1) & \text{otherwise} \end{cases}$$

Training algorithm

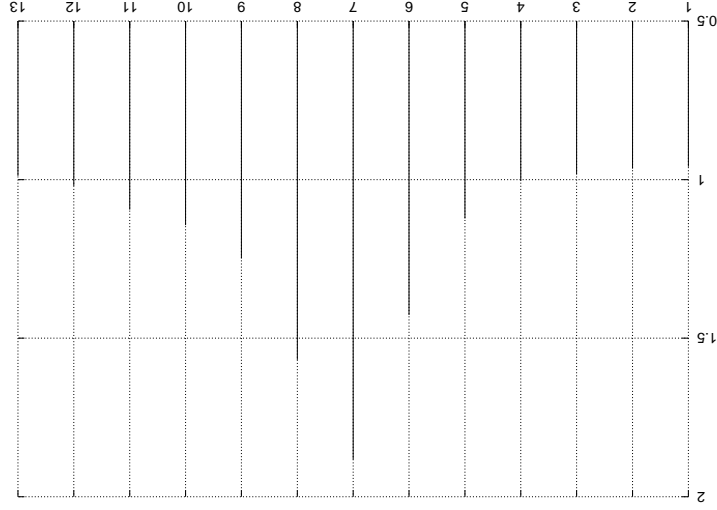
Stochastic gradient descent

Dedicated kernel

Analytical expression

$$K_{\theta}(x, x') = \exp \left(- \frac{\sum_{i=1}^n \theta_i^2 (1 - \langle x_i, x'_i \rangle)}{2\sigma^2} \right)$$

Influence of the position in the window

FIG. 6 – θ^* maximizing the alignment of K_{θ}

Prediction accuracy of a M-SVM without post-processing

	$SV M_{1v,r}$	M-SVM1
Q_3	68.4	69.0
C_α	0.53	0.54
C_β	0.50	0.51
C_γ	0.50	0.52

TAB. 3 – Prediction accuracy on a set of **1096** non-homologous protein sequences

Post-processing of the posterior probability estimates
with a Dynamic Programming algorithm

Underlying Hidden Markov Model : IHMM (Ramesh & Wilpon 92)

- One state for each structural state

- Parameters : $\lambda = (A, B, \Pi)$

$$A = \{a_{kl}(d)\}, (1 \leq k, l \leq 3), (1 \leq d \leq D_k)$$

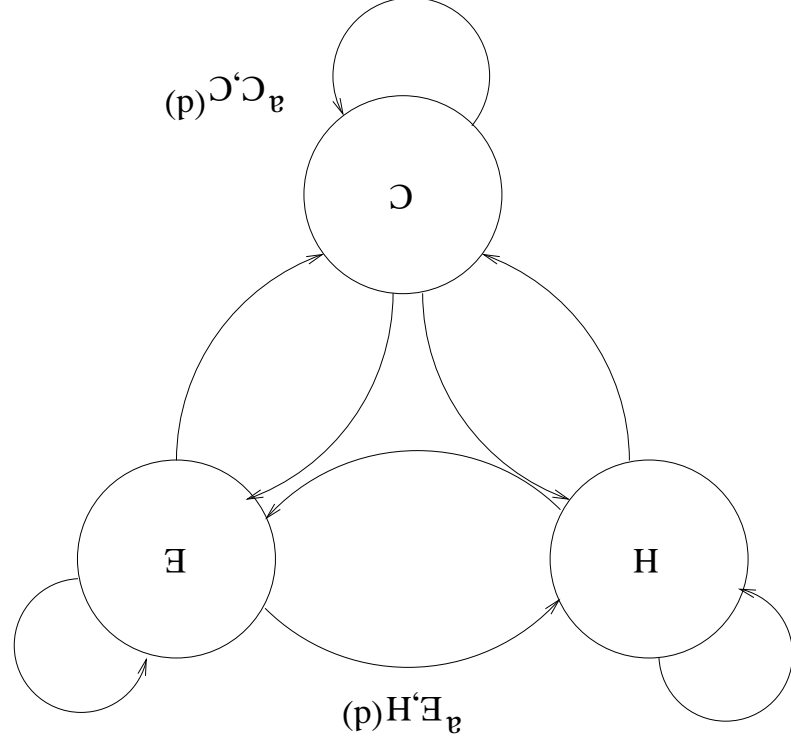
$a_{kl}(d)$: transition probability including a duration model

Prediction

“Best” sequence of states obtained with a variation of Viterbi’s algorithm

Implementation of *N-Best* algorithms

Top of the hierarchy : higher-level treatments



Training algorithm

Transition probabilities $a_{kl}(d)$: iteratively adapted with a Monte-Carlo and simulated annealing like algorithm so that the corresponding frequencies before and after post-processing become similar

Lengths of the conformational segments, observed and predicted



Fig. 7 – Distributions for the observed (left) and predicted (right) structures

Conclusions

Theory of M-SVMs

- New pathway to bound the risk of multi-class discriminant models
- New justifications of the training algorithms of all the M-SVMs proposed so far
- Possibility to develop new machines

M-SVMs for secondary structure prediction

- As ensemble methods, provide good generalization performance
- Little difference in generalization performance
- As basic classifiers, should prove superior to standard connectionist architectures

Work in progress

Theory and implementation of M-SVMs

- Comparison with studies involving *data dependent capacity measures* (Boucheron & al. 99, Bartlett & al. 02, Bousquet 02)
- Specification of optimization methods devoted to the new machines

M-SVMs for secondary structure prediction / biological sequence processing

- Additional work on the design and (efficient) implementation of specific kernels
- Fusion of additional knowledge sources provided by the biologists
- Global optimization of all the components of the hierarchical architecture