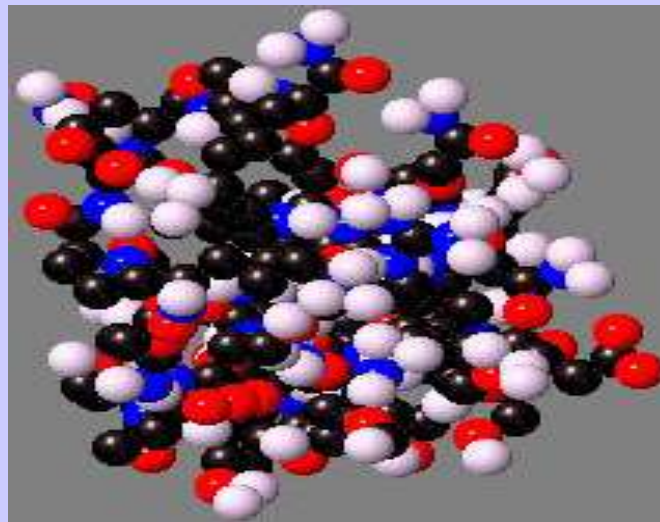


Apprentissage semi-supervisé asymétrique et estimations d'affinités locales dans les protéines



Christophe N. Magnan (Doctorant, LIF, Marseille)
CAP 2005, 03/06/2005



UNIVERSITÉ
DE PROVENCE
Aix-Marseille I

Plan

- ◆ Apprentissage semi-supervisé asymétrique
- ◆ Identifiabilité
- ◆ Expériences sur données artificielles
- ◆ Prédiction des ponts disulfures
- ◆ Conclusions et perspectives



Christophe N. Magnan (Doctorant, LIF, Marseille)
CAP 2005, 03/06/2005



Apprentissage semi-supervisé asymétrique

- ◆ variante de l'apprentissage semi-supervisé
 $S_{\text{lab}} = \{ (x_i, y_i) \}, S_{\text{unl}} = \{ x_i \}, x_i \in X, y_i \in Y$
- ◆ classification binaire: $Y = \{ 0, 1 \}$
- ◆ données: $S_{\text{pos}} = \{ (x_i, 1) \}, S_{\text{unl}} = \{ x_i \}$
- ◆ S_{pos} distribué selon $P(x \mid y=1)$



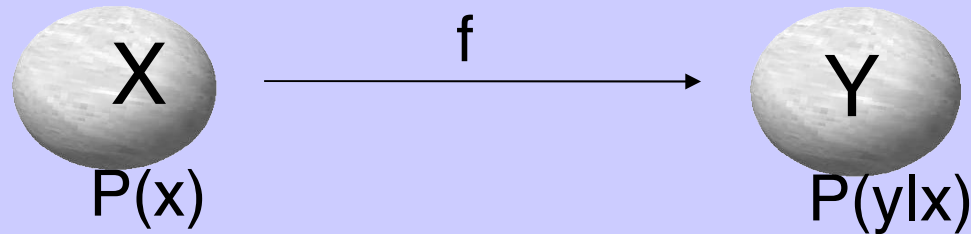
Apprentissage semi-supervisé asymétrique

- ◆ F. Denis: « PAC learning from positive statistical queries » (ALT 1998)
- ◆ F. Denis, F. De Comité, R. Gilleron, F. Letouzey: « Positive and unlabeled examples help learning » (ALT 1999)
- ◆ F. Denis, R. Gilleron, A. Laurent, M. Tommasi: « Text classification and cotraining from positive and unlabeled examples » (ICML 2003)
- ◆ B. Liu, X. Li: « Learning to classify text using positive and unlabeled data » (IJCAI 2003)

- ◆ Problème: identifiabilité?



Identifiabilité



- ◆ Données: $P(x)$ et $P(x \mid y=1)$
- ◆ Définition: une distribution \mathcal{D} est identifiable ssi les données permettent de déterminer cette distribution.
- ◆ Objectif: déterminer $P(y=1 \mid x) \forall x$
- ◆ Mais $P(y=1 \mid x)$ est-elle déterminée par les données?



Identifiabilité (suite)

$$P(y=1 | x) = \frac{P(x | y=1) \cdot P(y=1)}{P(x)}$$

- ◆ Propriété: $P(y=1 | x)$ déterminée ssi $P(y=1)$ est déterminé
- ◆ Propriété: \mathcal{D} identifiable ssi $P(y=1)$ identifiable
- ◆ Les données positives et non étiquetées permettent-elles d'identifier $P(y=1)$?



Identifiabilité (suite)

- ◆ NON sans hypothèse supplémentaire

$\Rightarrow \forall \lambda \leq \inf \{ P(x)/P(x | y=1) \}$, on peut construire une distribution P où $\lambda = P(y=1)$

- ◆ Mais pour certaines classes de distributions, $P(y=1)$ identifiable à partir de $P(x)$ et $P(x | y=1)$.



Exemple trivial de distribution identifiable

- ◆ Modèles déterministes:

$\forall x \in X, P(y=1|x)=1$ ou $P(y=1|x)=0$

$$P(y=1) = \sum_x P(x)P(y=1|x)$$



Distributions du classifieur naïf de Bayes

◆ $C_{\text{Bayes}}(\mathbf{x}) = \operatorname{argmax}_y P(y|\mathbf{x}) = \operatorname{argmax}_y P(\mathbf{x}|y) \cdot P(y)$

=> optimal mais incalculable dans la plupart des problèmes réels

◆ Hypothèse naïve de Bayes : les attributs sont indépendants conditionnellement à chaque classe

◆ $C_{\text{NaiveBayes}}(\mathbf{x}) = \operatorname{argmax}_y P(y) \prod P(x^i|y)$

=> nombre de paramètres raisonnable, hypothèse naïve.



Identifiabilité des modèles naïfs de Bayes

- ◆ **Théorème:** pour les distributions qui respectent l'hypothèse naïve de Bayes, $P(y=1)$ est déterminé par $P(x)$ et $P(x \mid y=1)$ dès 2 attributs *.
- ◆ * **Condition:** il existe x^i et x^j pour lesquels $P(x^i \mid y=1) \neq P(x^i)$ et $P(x^j \mid y=1) \neq P(x^j)$.
- ◆ Soient x^i et x^j tels que $P(x^i \mid y=1) \neq P(x^i)$ et $P(x^j \mid y=1) \neq P(x^j)$, alors:

$$P(y=1) = \frac{P(x^i \cap x^j) - P(x^i)P(x^j)}{P(x^i \mid y=1)P(x^j \mid y=1) - P(x^i \mid y=1)P(x^j) - P(x^j \mid y=1)P(x^i) + P(x^i \cap x^j)}$$

Estimateur consistant

- ◆ Estimateur consistant de $P(y=1)$:

$$(*) P(y=1) = \frac{\sum_i |P(x^i \cap x^j) - P(x^i)P(x^j)|}{\sum_i |P(x^i | y=1)P(x^j | y=1) - P(x^i | y=1)P(x^j) - P(x^j | y=1)P(x^i) + P(x^i \cap x^j)|}$$

- ◆ Algorithme sous-jacent:

- ◆ **Entrée:** S_{pos}, S_{unl}
 - ◆ 1) Estimer $P(x)$ et $P(x|y=1)$
 - ◆ 2) Calculer $P(y=1)$ avec (*)
 - ◆ 3) Calculer $P(y=1|x) \forall x$
- ◆ **Sortie:** un modèle θ

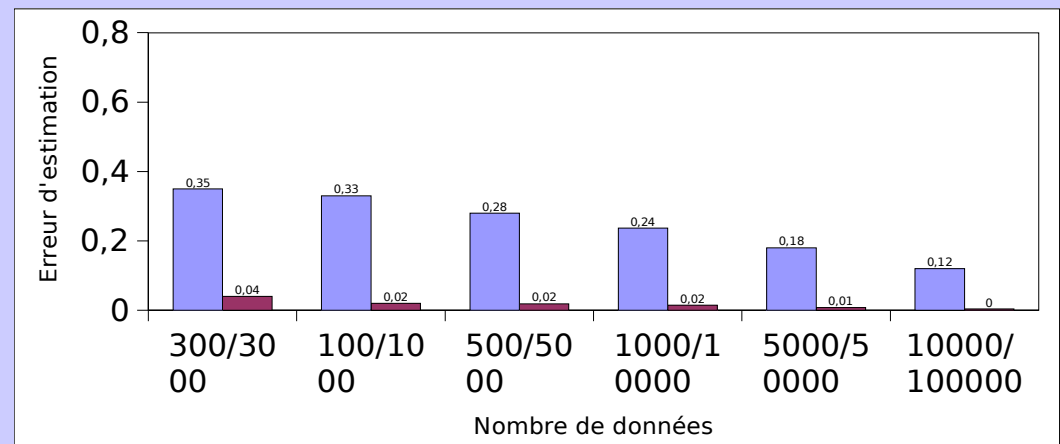


Identifiabilité des modèles naïfs de Bayes à partir de $P(\cdot)$

- ◆ 2001:D. Geiger, D. Heckerman, H. King, C. Meek:
« Stratified exponential families: Graphical models and model selection » (The Annals of statistics, 29(2), 505-529)
- ◆ Identifiabilité des modèles naïfs de Bayes dès 3 attributs et à permutation de classes près à partir de $P(\cdot)$
- ◆ Convergence très lente, inutilisable en pratique

Erreur d'estimation de $P(y=1)$ en fonction du nombre de données:

$$\left(E(\hat{P}(y=1) - P(y=1))^2 \right)^{1/2}$$



Estimations analytiques

VS

Principe du maximum de vraisemblance

◆ Vraisemblance:

◆ $L(\theta, S_{lab}) = \prod P((x_i, y_i) | \theta), i \in \{1, \dots, |S|\}$

◆ $L(\theta, S_{unl}) = \prod [P(y=1)P(x_i | y=1, \theta) + (1 - P(y=1))P(x_i | y=0, \theta)]$

◆ $L(\theta, S_{lab}, S_{unl}) = L(\theta, S_{lab}) \cdot L(\theta, S_{unl})$

◆ Principe : maximiser L



Méthode E.M. et classifieur naïf de Bayes

- ◆ 1999, A. Mc Callum, S. Thrun, T. Mitchell:
« Text classification from labeled and unlabeled documents using e.m. »
- ◆ Application de E.M. à l'inférence des paramètres des modèles naïfs de Bayes en contexte semi-supervisé
- ◆ Classification de textes: gain significatif de performances grâce aux données de S_{unl}



Adaptation au cas asymétrique

◆ Entrée: S_{pos} , S_{unl}

- ◆ 1) Estimer $P(\cdot)$ et $P(\cdot|y=1)$
- ◆ 2) Estimer $P(y=1)$ avec (*)
- ◆ 3) Model initial θ^0
- ◆ 4) $\forall x \in S_{\text{unl}}$, calculer $P(y=1|x, \theta^k)$ et $P(y=0|x, \theta^k)$
- ◆ 5) Calculer les paramètres du modèle θ^{k+1}
- ◆ 6) Itérer à l'étape 4) jusqu'à convergence

◆ Sortie: un modèle θ

◆ Pas de résultats théoriques sur la vitesse de convergence



Résultats expérimentaux sur données artificielles

Nombre d'attributs binaires			20
Performances modèle cible			0.035 (0.025)
S	Algorithme	Ens. apprentissage	Performances
$ S_{lab} =100$ $ S_{unl} =1000$	NB semi-sup	$S_{lab} + S_{unl}$	0.048 (0.097) 24,1
	NB asym analytique	$S_{pos} + S_{unl}$	0.082 (0.051)
	NB asym + EM	$S_{pos} + S_{unl}$	0.042 (0.031) 24
	NB unl	S_{unl}	0.204 (0.167)
$ S_{lab} =500$ $ S_{unl} =5000$	NB semi-sup	$S_{lab} + S_{unl}$	0.037 (0.027) 21,6
	NB asym analytique	$S_{pos} + S_{unl}$	0.044 (0.027)
	NB asym + EM	$S_{pos} + S_{unl}$	0.037 (0.027) 23,6
	NB unl	S_{unl}	0.126 (0.124)



Christophe N. Magnan (Doctorant, LIF, Marseille)
CAP 2005, 03/06/2005



Résultats expérimentaux sur données artificielles

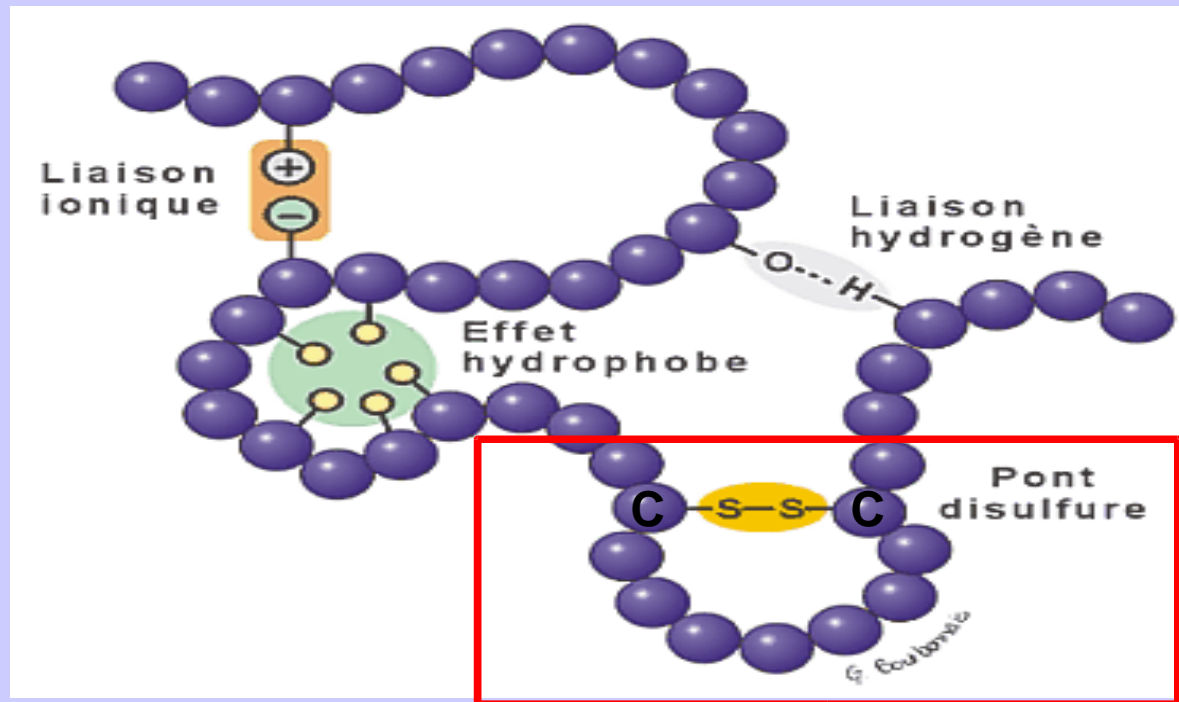
Nombre d'attributs binaires			20
Performances modèle cible			0.035 (0.025)
S	Algorithme	Ens. apprentissage	Performances
$ S_{lab} =100$ $ S_{unl} =1000$	NB semi-sup	$S_{lab} + S_{unl}$	0.048 (0.097) 24,1
	NB asym analytique	$S_{pos} + S_{unl}$	0.082 (0.051)
	NB asym + EM	$S_{pos} + S_{unl}$	0.042 (0.031) 24
	NB unl	S_{unl}	0.204 (0.167)
$ S_{lab} =500$ $ S_{unl} =5000$	NB semi-sup	$S_{lab} + S_{unl}$	0.037 (0.027) 21,6
	NB asym analytique	$S_{pos} + S_{unl}$	0.044 (0.027)
	NB asym + EM	$S_{pos} + S_{unl}$	0.037 (0.027) 23,6
	NB unl	S_{unl}	0.126 (0.124)



Christophe N. Magnan (Doctorant, LIF, Marseille)
CAP 2005, 03/06/2005

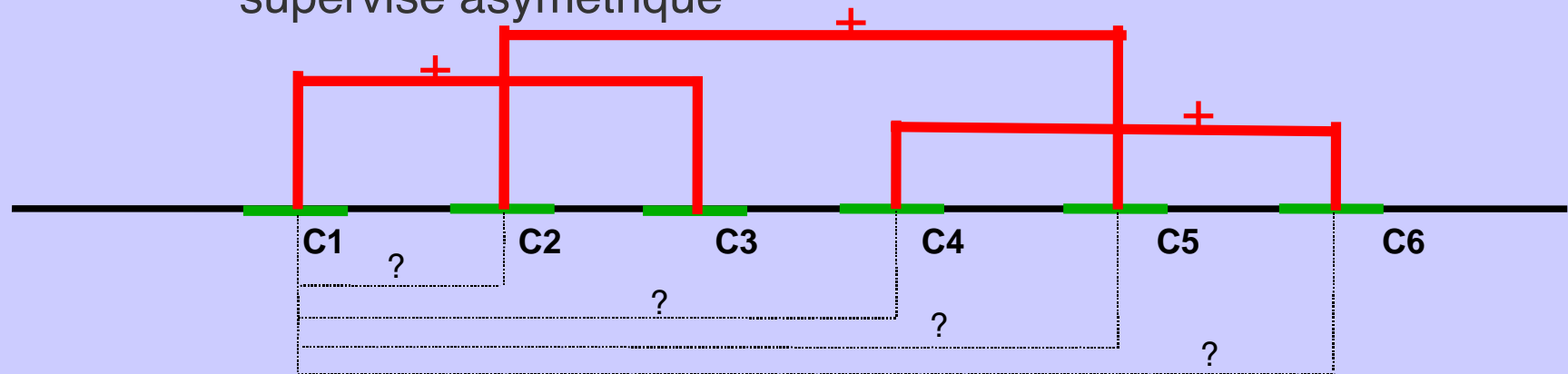


Application aux ponts disulfures



Motivations

- ◆ Etudier la contribution de l'environnement local aux cystéines à l'appariement de celles-ci : notion d'affinité locale
- ◆ Comment extraire au mieux cette information
- ◆ Etudier le statut des paires de cystéines non appariées:
 - ◆ exemples négatifs : apprentissage supervisé
 - ◆ exemples de classe non déterminée : apprentissage semi-supervisé asymétrique



Méthodes utilisées

- ◆ Apprentissage supervisé:
 - ◆ algorithme naïf de Bayes
- ◆ Apprentissage semi-supervisé asymétrique:
 - ◆ algorithme naïf de Bayes avec E.M. et estimateur de $P(y=1)$

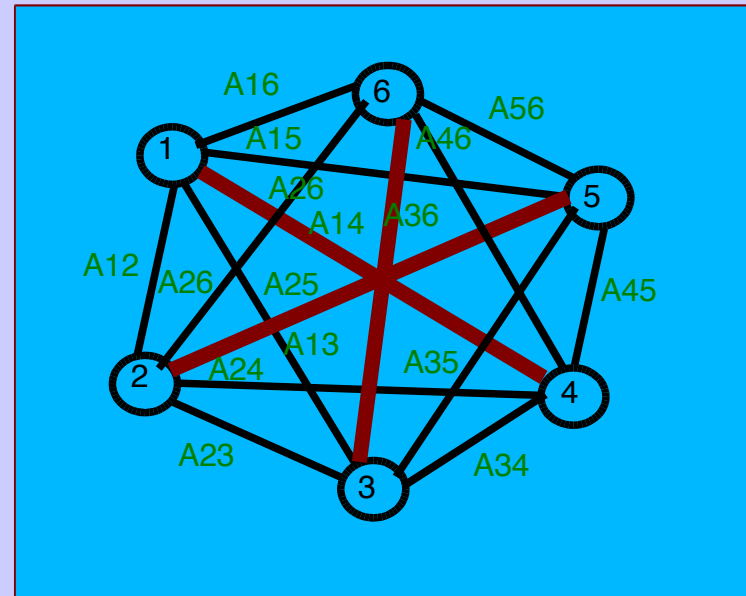
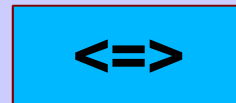
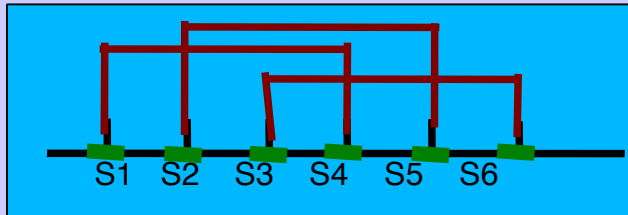


Représentation des données

- ◆ Soient f et f' deux segments, centrés sur des cystéines, de rayon n :
 - ◆ $f = x_{-n}, \dots, x_{-1}, C, x_1, \dots, x_n$
 - ◆ $f' = x'_{-n}, \dots, x'_{-1}, C, x'_1, \dots, x'_n$
- ◆ $(f, f') = \{ x_i x'_j \}, i, j \in \{ -n, \dots, n \}$
- ◆ Alphabet de taille 231



Protocole de test



■ = portion de 15 résidus (x_k)
centrée sur une cystéine

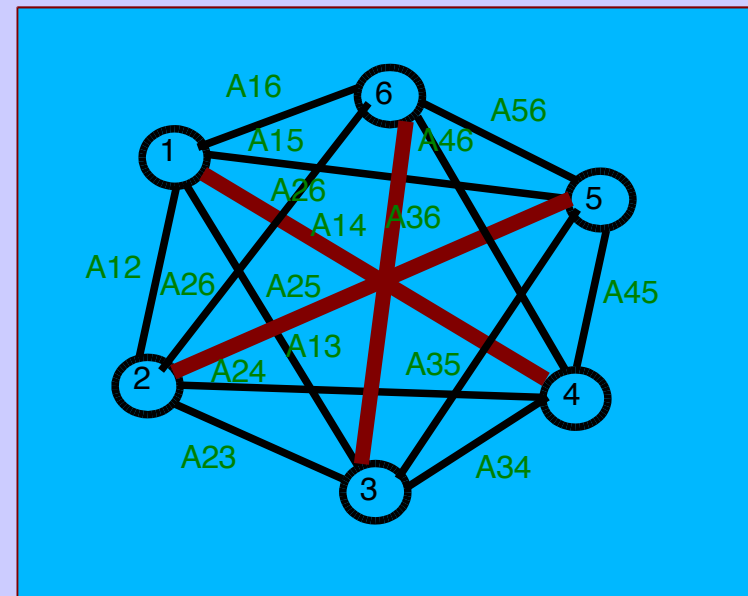
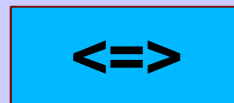
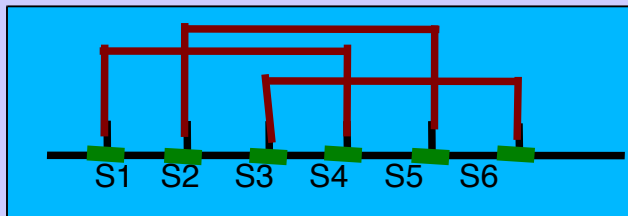
$$A_{ij} = P(y=1 \mid f_i, f_j)$$



Christophe N. Magnan (Doctorant, LIF, Marseille)
CAP 2005, 03/06/2005



Protocole de test



■ = portion de 15 résidus (xk)
centrée sur une cystéine

$$A_{ij} = P(y=1 \mid f_i, f_j)$$

⇒ couplage parfait de poids maximal



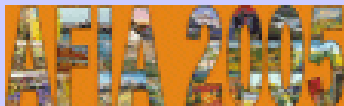
Christophe N. Magnan (Doctorant, LIF, Marseille)
CAP 2005, 03/06/2005



Résultats

- ◆ Données issues de la Protein Data Bank (PDB)
C. Geourjon (IBCP, Lyon), ACI GENOTO3D
- ◆ Comparaison des résultats sur le nombre de ponts correctement prédits

Nombre de ponts	2	3
Nombre de protéines	51	50
Nombre de ponts total	102	150
Aléa	34 (33,33%)	30 (20%)
Naïve Bayes supervisé	41 (40,2%)	26,25 (17,5%)
Naïve Bayes asym + EM	60 (58,8%)	50,1 (33,4%)



Conclusions et perspectives

Apprentissage asymétrique

- ◆ Identifiabilité des modèles naïfs de Bayes
- ◆ Algorithmes
- ◆ Résultats expérimentaux très satisfaisants
- ◆ Perspectives:
 - ◆ Autres classes de modèles identifiables
 - ◆ Méthodes pour identifier ces modèles



Christophe N. Magnan (Doctorant, LIF, Marseille)
CAP 2005, 03/06/2005



Conclusions et perspectives

Ponts disulfures

- ◆ Présence probable d'une information locale autour des cystéines
- ◆ L'hypothèse de représentation des données semble valide
- ◆ Perspectives:
 - ◆ Utiliser plus d'informations
 - ◆ Intégrer ces résultats dans des méthodes plus sophistiquées
 - ◆ Stratégies locales/globales



Références apprentissage semi-supervisé asymétrique

- ◆ F. Denis: « PAC learning from positive statistical queries » (ALT 1998)
- ◆ F. Denis, F. De Comit e, R. Gilleron, F. Letouzey: « Positive and unlabeled examples help learning » (ALT 1999)
- ◆ F. Denis, R. Gilleron, A. Laurent, M. Tommasi: « Text classification and cotraining from positive and unlabeled examples » (ICML 2003)
- ◆ B. Liu, X. Li: « Learning to classify text using positive and unlabeled data » (IJCAI 2003)



Références sur la prédiction des ponts disulfures

- ◆ P. Fariselli, R. Casadio: « Prediction of disulfide connectivity in proteins » (2001)
- ◆ P. Fariselli, P. Martelli, R. Casadio: « A neural network-based method for predicting the disulfide connectivity in proteins » (2002)
- ◆ A. Vullo, P. Frasconi: « A recursive connectionist approach for predicting disulfide connectivity in proteins » (2003)
- ◆ A. Vullo, P. Frasconi: « Disulfide connectivity prediction using recursive neural networks and evolutionary information » (2004)

