

Versions récursive et adaptative d'une mesure d'indépendance.

Pierre-Olivier AMBLARD^{1,2 *}

¹GIPSA-lab, DIS(UMR CNRS 5216), ENSE3-BP 46 38402 Saint Martin d'Hères Cedex

²Dept. Math&Stat, The University of Melbourne, Parkville, VIC3010, Australie.

bidou.amblard@gipsa-lab.inpg.fr

Résumé – L'objet de ce papier est de développer des algorithmes récursifs et adaptatifs pour mesurer l'indépendance entre deux suites de variables aléatoires. Nous considérons la mesure HSIC fondée sur le plongement des probabilités dans des espaces de Hilbert à noyau reproduisant (RKHS). Nous donnons deux versions récursive et adaptative d'un algorithme d'estimation de cette mesure d'indépendance. Une première version récursive permet le calcul de HSIC rapidement en utilisant peu de mémoire, autorisant ainsi son utilisation pratique sur des échantillons de grandes tailles. Une deuxième version récursive incorpore une représentation parcimonieuse (alternative aux approximations de rang faible) qui permet de traiter des données en-ligne et qui se prête à la mise sous forme adaptative.

Abstract – Recursive and adaptive algorithms to measure independence between two series of random variables are developed. We consider the HSIC measure based on embeddings of probability measures in reproducing kernel Hilbert spaces. We give two algorithms for the recursive calculation of HSIC. The first one uses low memory and is adapted to large data sets. The second one supports a sparsification front-end and can easily be turned into an adaptive algorithm.

1 Introduction

De nombreux problèmes de traitement du signal font appel explicitement à des mesures d'indépendance. Un exemple fameux est la séparation de source, qui fût initialement résolue en utilisant des mesures d'indépendance. Un autre exemple où l'indépendance intervient est la modélisation graphique de signaux multivariés. En neuroscience, la notion d'indépendance est utilisée pour rendre compte de la connectivité fonctionnelle entre deux aires cérébrales. Dans ce contexte, des études actuelles examinent l'évolution au cours du temps de la connectivité. Rendre dépendant du temps des mesures d'indépendance est donc une question pertinente. Dans [10] par exemple, une mesure de causalité fondée sur une mesure d'indépendance est évaluée au cours du temps en utilisant une fenêtre glissante.

Dans cet article, nous examinons une mesure utilisant les plongements de mesures de probabilité dans des espaces de Hilbert à noyau reproduisant (RKHS pour reproducing kernel Hilbert space). Implicitement, le plongement est effectuée à l'aide d'une fonction non linéaire φ qui envoie l'espace des observations dans un espace de Hilbert \mathcal{H} . Si φ est bien choisie, l'espace est un RKHS, et les produits scalaires s'y calculent remarquablement aisément, puisque φ est associée à une fonction définie positive symétrique k , le noyau, qui vérifie $k(x, y) = \langle \varphi(x) | \varphi(y) \rangle_{\mathcal{H}}$. Les données sont transformées non linéairement, mais plongées dans un espace linéaire, dans lequel les traitements linéaires (utilisant le produit scalaire) se calculent facilement. Autrement dit, il est possible de réaliser des traitements non linéaires à l'aide d'outils linéaires ! Dans cet ordre d'idée, on peut se demander s'il est possible de mesurer l'indépendance entre deux variables en utilisant des mesures de corrélation entre des transformées non linéaires des variables. Cette heuristique fût utilisée par exemple

par Jutten&Hérault dans leur algorithme de séparation aveugle de sources. La formalisation récente de cette idée repose sur des idées anciennes de Rényi, et a conduit à utiliser les normes d'opérateurs de covariance entre espaces de Hilbert à noyau reproduisant. Une mesure emblématique de cette approche est appelée HSIC pour Hilbert-Schmidt Independence Criterion.

Dans la suite de cette introduction, nous rappelons l'essentiel relatif à cette mesure. Nous en donnerons ensuite deux versions récursives. Ces versions permettent un calcul rapide (mais toujours en $O(n^2)$) de la mesure sans considérer explicitement les matrices de Gram. La deuxième version autorise alors une implantation approchée utilisant une procédure parcimonieuse (initialement proposée par C. Richard et ses collaborateurs dans un contexte tout différent [6]). Enfin, on propose une version adaptative permettant de détecter des changements dans la structure de dépendance entre deux grandeurs. Le dernier paragraphe sera dédié à des illustrations des différents algorithmes proposés.

Quelques rappels sur HSIC. HSIC a été introduite par Gretton et ses collaborateurs durant la dernière décennie [4, 8]. Soient deux variables aléatoires X, Y définies sur un espace probabilisé (Ω, \mathcal{F}, P) et prenant leurs valeurs dans respectivement \mathbb{X} and \mathbb{Y} (dotés de tribus adéquates). Soient deux noyaux (fonctions symétriques définies positives) $k_x : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$ et $k_y : \mathbb{Y} \times \mathbb{Y} \rightarrow \mathbb{R}$ (typiquement des noyaux gaussiens dans tout le résumé). Les espaces de Hilbert à noyau reproduisant correspondants sont notés \mathcal{H}_x et \mathcal{H}_y et sont supposés ici séparables [9, 7]. L'opérateur de covariance entre X et Y est l'unique opérateur linéaire borné de \mathcal{H}_x dans \mathcal{H}_y tel que $\text{Cov}[f(X), g(Y)] = \langle g | \Sigma_{YX} f \rangle_{\mathcal{H}_y}$. Cet opérateur est bien défini moyennant $E[k_x(X, X)] < +\infty$ et $E[k_y(Y, Y)] < +\infty$. L'opérateur est alors de Hilbert-Schmidt, *i.e.* sa norme de Hilbert-Schmidt (HS) définie par $\|\Sigma_{YX}\|^2 := \sum_i \|\Sigma_{YX} \varphi_i\|_{\mathcal{H}_y}^2$ est finie, pour toute base orthonormée $\{\varphi_i\}_{i \in \mathbb{N}}$ de \mathcal{H}_x [3].

*P.O.A. est soutenu par une bourse Marie Curie International Outgoing Fellowship de l'Union Européenne.

Un théorème prouvé par Gretton étend un résultat de Rényi stipulant l'équivalence entre indépendance et $\sup_{f,g} \text{Cov}[f(X), g(Y)] = 0$ pour des fonctions continues bornées f et g [5]. L'extension proposée par Gretton consiste à choisir f et g dans un espace à noyau reproduisant qui doit être suffisamment riche pour que la propriété de Rényi reste valide. Cette richesse est conférée par la propriété d'universalité des noyaux. Un noyau est universel si le RKHS associé est dense dans l'ensemble des fonctions continues bornées. Le théorème de Gretton stipule alors que $\sup_{f \in \mathcal{U}_x, g \in \mathcal{U}_y} \langle g | \Sigma_{YX} f \rangle_{\mathcal{H}_y} = 0$ est équivalent à l'indépendance entre X et Y . \mathcal{U} est le sous-ensemble des vecteurs de \mathcal{H} de norme inférieure ou égale à 1. L'hypothèse d'universalité permet d'approcher n'importe quelle fonction continue bornée par une suite de fonctions du RKHS, et donc de satisfaire le théorème de Rényi. Cette hypothèse est vérifiée par de nombreux noyaux, dont le noyau gaussien très largement utilisé. Par contre un noyau de type $(1 + \langle x|y \rangle)^q$ génère un espace vectoriel de dimension finie et n'est pas universel.

La magie du résultat de Gretton réside dans le fait que $\sup_{f \in \mathcal{U}_x, g \in \mathcal{U}_y} \langle g | \Sigma_{YX} f \rangle_{\mathcal{H}_y}$ n'est rien d'autre que la norme usuelle de l'opérateur de covariance ! La norme usuelle est inférieure ou égale à la norme de Hilbert-Schmidt, de sorte que le résultat de Gretton reste valide avec la norme HS. Or cette norme est très simple à évaluer à partir de données [4]. Les calculs de l'estimée peuvent être effectués en $O(n^2)$ pour un échantillon de taille n , une complexité élevée qui peut être réduite en utilisant des approximations de rang faible pour les matrices de Gram. Une version largement utilisée dans le contexte, voir [4, 8], est la transformée de Cholesky incomplète [2].

Dans la suite, des grandeurs ont x ou y en indice. Lorsque les définitions ou les expressions sont communes, ou nous ne les écrivons que pour un des deux indices, ou nous supprimons l'indice ; il est implicitement entendu que la définition ou l'expression en jeu est valide pour les deux indices.

2 Algorithmes

Première forme récursive. L'implantation proposée permet de calculer en des temps raisonnables la mesure pour des échantillons de taille très grande (10^4 à 10^5 échantillons), chose impossible en utilisant des matrices. De nombreuses multiplications peuvent être éliminées et les nécessités de stockage réduites en ne manipulant que des vecteurs. La version empirique que nous utilisons ici n'est seulement qu'asymptotiquement non biaisé. Les développements peuvent se réaliser sur la version non biaisée de [8].

Considérons n données X_i, Y_i . L'estimateur \hat{H}_n de HSIC $\|\Sigma_{YX}\|^2$ est obtenu à partir des matrices de Gram K_x and K_y , dont les (i, j) èmes éléments sont respectivement $k_x(X_i, X_j)$ et $k_y(Y_i, Y_j)$. Soit $C_n := I_n - \mathbf{1}\mathbf{1}^\top/n$ la matrice de centrage ($\mathbf{1}$ est un vecteur de 1 de dimension appropriée). Alors l'estimateur de HSIC est simplement $\hat{H}_n = n^{-2} \text{Tr}(\tilde{K}_x \tilde{K}_y)$ où $\tilde{K} = C_n K C_n$. K^n est la matrice de Gram de taille $n \times n$ obtenue après observation des n premières données. Comme $\text{Tr}(ABC) = \text{Tr}(BCA)$ et $C_n^2 = C_n$, une forme plus simple pour \tilde{K} est $\tilde{K} = K C_n$, définition que nous adoptons. Notons toutefois que $\tilde{K} = K C_n$ n'est plus symétrique.

Pour obtenir une version récursive, les matrices de Gram sont partitionnées. Soient \tilde{k}^n et $\tilde{\ell}^n$ deux vecteurs de dimension $n-1$, \tilde{K}^{n-1} une matrice carrée de dimension $n-1$ et $\tilde{\kappa}^n$ une constante. Alors

$$\tilde{K}_x \tilde{K}_y^n = \begin{pmatrix} \tilde{K}_x^{n-1} & \tilde{k}_x^n \\ (\tilde{\ell}_x^n)^\top & \tilde{\kappa}_x^n \end{pmatrix} \begin{pmatrix} \tilde{K}_y^{n-1} & \tilde{k}_y^n \\ (\tilde{\ell}_y^n)^\top & \tilde{\kappa}_y^n \end{pmatrix} \text{ et} \\ n^2 \hat{H}_n = \left(\text{Tr}(\tilde{K}_x^{n-1} \tilde{K}_y^{n-1}) + (\tilde{\ell}_x^n)^\top \tilde{k}_y^n + (\tilde{\ell}_y^n)^\top \tilde{k}_x^n + \tilde{\kappa}_x^n \tilde{\kappa}_y^n \right) (1)$$

Nous cherchons alors une version récursive pour chacun des termes. Soient \mathbf{k}_x^n le vecteur de dimension $(n-1)$ d'éléments $k_x(X_i, X_n), i = 1, \dots, n-1$, $\kappa_x^n = k_x(X_n, X_n)$ (et les mêmes définitions en y). Nous avons $\tilde{K}^n = K^n C_n$, matrices dont les partitions s'écrivent

$$\begin{pmatrix} K^{n-1} & \mathbf{k}^n \\ (\mathbf{k}^n)^\top & \kappa^n \end{pmatrix} \begin{pmatrix} I_{n-1} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top & -\frac{1}{n} \mathbf{1} \\ -\frac{1}{n} \mathbf{1}^\top & 1 - \frac{1}{n} \end{pmatrix}$$

Comme $I_{n-1} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top = C_{n-1} + \frac{1}{n(n-1)} \mathbf{1}\mathbf{1}^\top$, on a par identification

$$\tilde{K}^{n-1} = K^{n-1} C_{n-1} - \frac{1}{n} (\mathbf{k}^n - \frac{1}{n-1} K^{n-1} \mathbf{1}) \mathbf{1}^\top \quad (2)$$

$$\tilde{\mathbf{k}}^n = \mathbf{k}^n - \frac{1}{n} (\mathbf{k}^n + \frac{1}{n} K^{n-1} \mathbf{1}) \quad (3)$$

$$\tilde{\kappa}^n = \kappa^n - \frac{1}{n} (\kappa^n + \frac{1}{n} \mathbf{1}^\top \mathbf{k}^n) \quad (4)$$

$$\tilde{\ell}^n = \mathbf{k}^n - \frac{1}{n} (\kappa^n + \mathbf{1}^\top \mathbf{k}^n) \mathbf{1} \quad (5)$$

Posons $\mathbf{m}^n = K^n \mathbf{1}$. En utilisant la matrice de Gram en partitions, \mathbf{m}^n , (3), (4) s'écrivent

$$\mathbf{m}^n = \begin{cases} \mathbf{m}^{n-1} + \mathbf{k}^n \\ \mathbf{1}^\top \mathbf{k}^n + \kappa^n \end{cases} ; \begin{pmatrix} \tilde{\mathbf{k}}^n \\ \tilde{\kappa}^n \end{pmatrix} = \begin{pmatrix} \mathbf{k}^n \\ \kappa^n \end{pmatrix} - \frac{1}{n} \mathbf{m}^n \quad (6)$$

qui avec (5), permet d'obtenir une récursion pour $(\tilde{\ell}_x^n)^\top \tilde{\mathbf{k}}_y^n + (\tilde{\ell}_y^n)^\top \tilde{\mathbf{k}}_x^n + \tilde{\kappa}_x^n \tilde{\kappa}_y^n$ dans (1). Le dernier terme est $\text{Tr}(\tilde{K}_x^{n-1} \tilde{K}_y^{n-1})$. Remarquons que $\mathbf{k}^n - \frac{1}{n-1} K^{n-1} \mathbf{1}$ dans (2) est égal à $\mathbf{k}^n - \frac{1}{n-1} \mathbf{m}^{n-1}$ et peut être interprété comme une estimation de $\tilde{\mathbf{k}}^n$ à partir des données passées (comparer à 6). On note ce terme $\tilde{\mathbf{k}}^{n|n-1}$. Nous avons alors $\tilde{K}^{n-1} = K^{n-1} C_{n-1} - \frac{1}{n} \tilde{\mathbf{k}}^{n|n-1} \mathbf{1}^\top$, et par suite la trace de $\tilde{K}_x^{n-1} \tilde{K}_y^{n-1}$ s'écrit $\text{Tr}(K_x^{n-1} C_{n-1} K_y^{n-1} C_{n-1}) - \frac{1}{n} \mathbf{1}^\top K_y^{n-1} C_{n-1} \tilde{\mathbf{k}}_x^{n|n-1} - \frac{1}{n} \mathbf{1}^\top K_x^{n-1} C_{n-1} \tilde{\mathbf{k}}_y^{n|n-1} + \frac{1}{n^2} \mathbf{1}^\top \tilde{\mathbf{k}}_x^{n|n-1} \mathbf{1}^\top \tilde{\mathbf{k}}_y^{n|n-1}$. En utilisant $\mathbf{m}^n = K^n \mathbf{1}$ et quelques calculs, on obtient :

HSIC récursif :

$$\mathbf{m}_x^n = \begin{cases} \mathbf{m}_x^{n-1} + \mathbf{k}_x^n \\ (\mathbf{k}_x^n)^\top \mathbf{1} + \kappa_x^n \end{cases} ; \begin{pmatrix} \tilde{\mathbf{k}}_x^n \\ \tilde{\kappa}_x^n \end{pmatrix} = \begin{pmatrix} \mathbf{k}_x^n \\ \kappa_x^n \end{pmatrix} - \frac{1}{n} \mathbf{m}_x^n$$

$$\tilde{\ell}_x^n = \mathbf{k}_x^n - \frac{1}{n} (\kappa_x^n + \mathbf{1}^\top \mathbf{k}_x^n) \mathbf{1} ; \tilde{\mathbf{k}}_x^{n|n-1} = \mathbf{k}_x^n - \frac{1}{n-1} \mathbf{m}_x^{n-1}$$

$$c_{xy}^{n|n-1} = \frac{1}{n} (\mathbf{m}_x^{n-1})^\top \tilde{\mathbf{k}}_y^{n|n-1} - \frac{1}{n(n-1)} \mathbf{1}^\top \mathbf{m}_x^{n-1} \mathbf{1}^\top \tilde{\mathbf{k}}_y^{n|n-1}$$

$$\hat{H}_n = \hat{H}_{n-1} + \frac{1}{n^2} \left((\tilde{\ell}_x^n)^\top \tilde{\mathbf{k}}_y^n + (\tilde{\ell}_y^n)^\top \tilde{\mathbf{k}}_x^n + \tilde{\kappa}_x^n \tilde{\kappa}_y^n - c_{xy}^{n|n-1} - c_{yx}^{n|n-1} + \frac{1}{n^2} \mathbf{1}^\top \tilde{\mathbf{k}}_x^{n|n-1} \mathbf{1}^\top \tilde{\mathbf{k}}_y^{n|n-1} - (2n-1) \hat{H}_{n-1} \right)$$

Les vecteurs \mathbf{m} , \mathbf{k} voient leur dimension augmenter au cours du temps, rendant évidemment inutilisable cet algorithme de manière permanente. Cette difficulté est à l'origine de la deuxième forme récursive, qui se prête aisément à un calcul approché fondé sur une notion de parcimonie, offrant une alternative aux approximations de rang faible.

Deuxième forme et parcimonie. Pour réduire la taille mémoire et accélérer le calcul, nous avons recours à une stratégie de parcimonie, comme pour les algorithmes de régression en-ligne. L'idée force

est celle de [6] d'un dictionnaire cohérent \mathcal{D}_n . Ce dictionnaire est constitué de données passées et se construit récursivement. Intuitivement, une nouvelle donnée n'est incluse dans le dictionnaire que si elle ne peut être correctement estimée à partir des données déjà incluses.

Soit \mathcal{D}_n un sous-ensemble de $\{1, \dots, n\}$, $\mathcal{D}_1 = \{1\}$. Soit $Z_i = (X_i, Y_i)$. Le dictionnaire est mis à jour par la règle $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\} \iff \max_{\alpha \in \mathcal{D}_{n-1}} |k_{xy}^c(Z_\alpha, Z_n)| < \mu_0$ où $0 < \mu_0 \leq 1$, et k_{xy}^c est un noyau normalisé sur l'espace produit $\mathbb{X} \times \mathbb{Y}$ ($k^c(x, y) = k(x, y)(k(x, x)k(y, y))^{-1/2}$). Une donnée nouvelle est incluse dans le dictionnaire si elle n'est pas suffisamment cohérente avec les précédentes. Si $\mu = 1$ les données sont systématiquement agrégées : on obtient alors HSIC usuel. La règle de cohérence simplifie le critère ALD (approximate linear dependence) proposé dans [1] qui repose sur la prédiction linéaire. Comme montré par Richard dans [6], pour $\mu_0 < 1$, la taille du dictionnaire reste finie si les données vivent dans un sous-espace compact. Ceci interdit l'utilisation de la technique pour des données non bornées (variables gaussiennes par exemple) puisque la taille du dictionnaire ne peut que croître. En pratique toutefois, l'approche peut être appliquée avec une grande confiance, la croissance du dictionnaire étant très lente pour des données raisonnables.

Pour appliquer le principe à un algorithme récursif, chaque mise à jour est précédée d'un test d'inclusion dans le dictionnaire. Dans le cas où une nouvelle donnée n'est pas incluse, nous ajoutons par rapport aux travaux de [6] une nouvelle information. Le dictionnaire réalise une quantification vectorielle des données. Pour des données à support compact, lorsque la taille finale du dictionnaire est atteinte, chaque membre du dictionnaire est le centre d'une cellule, appelée *cellule de cohérence*, dont la taille et la forme sont définies par le critère de cohérence. Une nouvelle donnée ne sera jamais incluse dans le dictionnaire puisque sa taille finale est atteinte. Par contre la nouvelle donnée appartient évidemment à une des cellules. Ne pas prendre en compte cette information provoque une perte d'information. Pour la minimiser, nous suggérons d'utiliser la mesure empirique des données calculée sur la partition induite par le critère de cohérence. Pour $\alpha \in \mathcal{D}_n$, on note $V_\alpha = \{Z, |k_{xy}^c(Z_\alpha, Z)| \geq \mu_0\}$ la cellule de cohérence de centre Z_α . On pose $\pi_n(\alpha) = \sum_{k=1}^n \mathbf{1}_{V_\alpha}(Z_k)$.

Dans la deuxième version, nous allons utiliser le critère de cohérence pour ne garder qu'un nombre restreint d'éléments représentatifs des variables analysées. Pour obtenir la forme de l'algorithme, on ne travaille que sur les données conservées par le critère de cohérence : à une date n donnée, on n'utilise que les centres de cellules Z_α et la mesure empirique $\pi_n(\alpha)$. On peut alors établir un algorithme récursif exact sur ce jeu de données. Les estimateurs empiriques de $m_x(\cdot) = E[k_x(\cdot, X)]$ et $M_{yx} = E[k_x(\cdot, X) \otimes k_y(\cdot, Y)]$ s'écrivent sous la forme générique e^n

$$e^n = \frac{1}{n} \sum_{\alpha \in \mathcal{D}_n} \pi_n(\alpha) f(X_\alpha, Y_\alpha) = \frac{n-1}{n} e^{n-1} + \frac{1}{n} f(X_a, Y_a) \quad (7)$$

où $f(X, Y) = k_x(\cdot, X) \otimes k_y(\cdot, Y)$ pour $e = M_{yx}$ et $f(X, Y) = k_x(\cdot, X)$ pour m_x . Rappelons que le produit tensoriel de deux vecteurs f, g est défini par $(f \otimes g)(u, v) = \langle u | f \rangle \langle v | g \rangle$. De plus, $a = n$ si la nouvelle donnée est incluse dans le dictionnaire, et $a = \arg \max_{\alpha \in \mathcal{D}_{n-1}} |k_{xy}^c(Z_\alpha, Z_n)|$ est l'indice de la cellule de cohérence de la nouvelle donnée si elle n'est pas incluse. En jouant avec les deux formes des estimateurs données en (7), on obtient un algorithme récursif. Il faut évaluer la norme HS de $C_{yx}^n = M_{yx}^n - m_y^n \otimes m_x^n$, c'est-à-dire calculer $\sum_l \langle C_{yx}^n \varphi_l | C_{yx}^n \varphi_l \rangle_{\mathcal{H}_y}$ où $\{\varphi_l\}_{l \in \mathbb{N}}$ est une base orthonormée quelconque de \mathcal{H}_x . On obtient alors $\|C_{yx}^n\|^2 = \|M_{yx}^n\|^2 + \|m_x^n\|^2 \|m_y^n\|^2 - 2c_{yx}^n$ où

$c_{yx}^n = \langle M_{yx}^n m_x^n | m_y^n \rangle$. Le calcul des normes est assez aisé, une fois remarquées les expressions $\langle M_{yx}^{n-1} k_x(\cdot, X_a) | k_y(\cdot, Y_a) \rangle = \frac{1}{n-1} \pi_{n-1}^\top \mathbf{k}_{x_a}^n \circ \mathbf{k}_{y_a}^n$ et $\langle m_x^{n-1} | k_x(\cdot, X_a) \rangle = \frac{1}{n-1} \pi_{n-1}^\top \mathbf{k}_{x_a}^n$, où $\mathbf{k}_{x_a}^n$ est le vecteur contenant les $k_x(X_\alpha, X_a)$, $\alpha \in \mathcal{D}_{n-1}$ et π_n le vecteur de dimension $|\mathcal{D}_n|$ contenant les $\pi_n(\alpha)$. Le produit scalaire c_{yx}^n est un peu plus délicat. On introduit $v_x^n(\beta) = \sum_{\gamma \in \mathcal{D}_n} \pi_n(\gamma) k_x(X_\gamma, X_\beta)$ pour $\beta \in \mathcal{D}_n$ et \mathbf{v}_x^n le vecteur correspondant. On a alors $c_{yx}^n = n^{-3} \pi_n^\top \mathbf{v}_x^n \circ \mathbf{v}_y^n$. Une version récursive de $v_x^n(\beta)$ est disponible, mais dépend de l'inclusion ou non de la nouvelle donnée dans le dictionnaire. Si $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$, la nouvelle donnée est $X_a = X_n$ puisqu'une nouvelle cellule est créée, le vecteur π_{n-1} voit sa dimension augmenter de 1 et on a la récursion $\pi_n = (\pi_{n-1}^\top, 1)^\top$. Il vient $v_x^n(\beta) = \sum_{\gamma \in \mathcal{D}_{n-1}} \pi_{n-1}(\gamma) k_x(X_\gamma, X_\beta) + k_x(X_a, X_\beta)$, et \mathbf{v}_x^n suit la récursion $\mathbf{v}_x^n = ((\mathbf{v}_x^{n-1} + \mathbf{k}_{x_a}^n)^\top, \pi_{n-1}^\top \mathbf{k}_{x_a}^n + \kappa_{x_a}^n)^\top$, où $\kappa_{x_a}^n = k_x(X_a, X_a)$. Si la nouvelle donnée n'est pas incluse dans le dictionnaire, elle est remplacée par le centre de sa cellule de cohérence, à savoir Z_a où $a = \arg \max_{\alpha \in \mathcal{D}_{n-1}} |k_{xy}^c(Z_\alpha, Z_n)|$. Dans ce cas $\mathcal{D}_n = \mathcal{D}_{n-1}$, le vecteur π_n reste égal π_{n-1} excepté pour l'élément $\pi_n(a)$ qui est incrémenté de 1, et on a $\mathbf{v}_x^n = \mathbf{v}_x^{n-1} + \mathbf{k}_{x_a}^n$ qui ne change pas de dimension. En collectionnant les divers éléments on obtient une forme récursive exacte pour l'estimateur n'utilisant que la quantification proposée. D'ailleurs, les données arrivant sont classées dans leur cellule et remplacée par le centre a de cette cellule. Pour terminer et améliorer l'algorithme, il suffit alors non pas d'utiliser Z_a comme nouvelle donnée, mais simplement Z_n systématiquement. Dans ce qui précède, on remplace alors \mathbf{k}_{x_a} par \mathbf{k}_x , vecteur qui contient les $k_x(X_n, X_\alpha)$, $\forall \alpha \in \mathcal{D}_{n-1}$. En résumé, on a obtenu :

HSIC parcimonieux :

$$\begin{aligned} \widehat{H}_n^{\mu_0} &= \|M_{yx}^n\|^2 + \|m_x^n\|^2 \|m_y^n\|^2 - 2c_{yx}^n \\ \|M_{yx}^n\|^2 &= \frac{(n-1)^2}{n^2} \|M_{yx}^{n-1}\|^2 + \frac{2}{n^2} \pi_{n-1}^\top \mathbf{k}_x^n \circ \mathbf{k}_y^n + \frac{\kappa_x^n \kappa_y^n}{n^2} \\ \|m_x^n\|^2 &= \frac{(n-1)^2}{n^2} \|m_x^{n-1}\|^2 + \frac{2}{n^2} \pi_{n-1}^\top \mathbf{k}_x^n + \frac{\kappa_x^n}{n^2} \\ c_{yx}^n &= \frac{1}{n^3} \pi_n^\top \mathbf{v}_x^n \circ \mathbf{v}_y^n \text{ où :} \end{aligned}$$

1. Si $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$

$$\mathbf{v}_x^n = \begin{pmatrix} \mathbf{v}_x^{n-1} + \mathbf{k}_x^n \\ \pi_{n-1}^\top \mathbf{k}_x^n + \kappa_x^n \end{pmatrix} \quad \text{et} \quad \pi_n = \begin{pmatrix} \pi_{n-1} \\ 1 \end{pmatrix}$$

2. Si $\mathcal{D}_n = \mathcal{D}_{n-1}$: $a = \arg \max_{\alpha \in \mathcal{D}_{n-1}} |k_{xy}^c(Z_\alpha, Z_n)|$,

$$\mathbf{v}_x^n = \mathbf{v}_x^{n-1} + \mathbf{k}_x^n \quad \text{et} \quad \pi_n = \pi_{n-1} + \delta_{a\alpha}$$

Dans ces expressions, on a $\kappa_x^n = k_x(X_n, X_n)$, π_n le vecteur de dimension $|\mathcal{D}_n|$ contenant les $\pi_n(\alpha)$ (initialisé par $\pi_1(1) = 1$).

Une version adaptative. La version adaptative s'obtient de manière analogue que la forme précédente, mais repose sur une moyenne exponentielle.

HSIC adaptatif parcimonieux :

$$\begin{aligned} \widehat{H}_n^{\mu_0} &= \|M_{yx}^n\|^2 + \|m_x^n\|^2 \|m_y^n\|^2 - 2c_{yx}^n \\ \|M_{yx}^n\|^2 &= (1-\gamma)^2 \|M_{yx}^{n-1}\|^2 + 2\gamma(1-\gamma) \pi_{n-1}^{\top} \mathbf{k}_x^n \circ \mathbf{k}_y^n + \gamma^2 \kappa_x^n \kappa_y^n \\ \|m_x^n\|^2 &= (1-\mu)^2 \|m_x^{n-1}\|^2 + 2\mu(1-\mu) \pi_{n-1}^{\top} \mathbf{k}_x^n + \mu^2 \kappa_x^n \\ c_{yx}^n &= \pi_n^{\top} \mathbf{v}_{x,n}^n \circ \mathbf{v}_{y,n}^n \text{ où :} \end{aligned}$$

HSIC parcimonieux (suite) :

1. Si $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{n\}$

$$\mathbf{v}_{x,n}^\mu = \begin{pmatrix} (1-\mu)\mathbf{v}_{x,n-1}^\mu + \mu\mathbf{k}_x^n \\ (1-\mu)\boldsymbol{\pi}_{n-1}^{\mu T}\mathbf{k}_x^n + \mu\mathbf{v}_x^n \end{pmatrix} \text{ et } \boldsymbol{\pi}_n^\eta = \begin{pmatrix} (1-\eta)\boldsymbol{\pi}_n^\eta \\ \eta \end{pmatrix}$$

2. Si $\mathcal{D}_n = \mathcal{D}_{n-1} : a = \arg \max_{\alpha \in \mathcal{D}_{n-1}} |k_{xy}^c(Z_\alpha, Z_n)|$,

$$\mathbf{v}_{x,n}^\mu = (1-\mu)\mathbf{v}_{x,n-1}^\mu + \mu\mathbf{k}_x^n \text{ et } \boldsymbol{\pi}_n^\eta = (1-\eta)\boldsymbol{\pi}_{n-1}^\eta + \eta\delta_{a\alpha}$$

où $\eta = \gamma, \mu$ ou ζ . Les expressions en x sont valables pour y en remplaçant x par y et μ par ζ . Les vecteurs $\boldsymbol{\pi}_n^\eta$ dépendent des facteurs d'oubli $\eta = \gamma, \mu$ ou ζ , et réalisent en quelque sorte une estimation adaptative des mesures de probabilités des variables X, Y et du couple.

3 Illustrations

Pour illustrer le comportement des algorithmes et la perte de performance limitée de HSIC parcimonieux, nous travaillons sur des données bidimensionnelles : la première composante X étant distribuée suivant une gaussienne centrée réduite, et la deuxième Y suivant une loi exponentielle bilatérale de paramètre 1. X et Y sont indépendantes, et la dépendance est obtenue par une rotation. Les statistiques suivantes sont calculées par méthode Monte-Carlo en utilisant 500 réalisations des processus, pour des tailles d'échantillons de $n = 3000$ (sauf les temps de calculs qui sont évalués jusque $n=10^4$ échantillons en moyennant sur 100 réalisations). L'ensemble des résultats est synthétisé par la figure (1). Nous utilisons le noyau gaussien $\exp(-\|x-y\|^2/1.2)$ pour lequel la valeur 1.2 est choisie empiriquement. HSIC parcimonieux est évalué pour $\mu_0 = 0.85, 0.9, 0.95$. La figure A montre des traces typiques et un zoom après convergence, qui illustre la faible perte de performance en fonction de μ_0 quand on compare HSIC parcimonieux à HSIC. La perte décroît avec μ_0 croissant, ce qui est attendu puisque $\mu_0 = 1$ correspond à HSIC usuel. Pour quantifier cette perte, on trace l'erreur quadratique moyenne (la valeur vraie de HSIC étant inconnue, on choisit la valeur de HSIC estimée). Au delà de la confirmation d'une convergence au taux usuel $1/n$, la figure (1B) montre encore une fois une très faible perte pour $\mu_0 = 0.95$. Ceci est confirmé dans l'insert où nous montrons les valeurs finales moyennes estimées μ_0 de 0.85 à 1, ainsi que des intervalles $\pm 2\sigma$ (de confiance à 95 % sous hypothèse gaussienne). Les moyennes sont très proches de la moyenne de HSIC, et les variances très comparables.

Ces observations prennent toute leur importance avec C et D. En C, nous traçons les temps de calculs (en unités arbitraires) en fonction de n pour $\mu_0 = 0.85, 0.9, 0.95, 0.99$ et 1 (HSIC). En D nous traçons la taille moyenne des dictionnaires obtenus en fonction de μ_0 (et les intervalles $\pm 2\sigma$). Pour $\mu_0 = .99$, le gain en temps de calcul est comparable à HSIC (même plus mauvais pour les n petits). Mais dès $\mu_0 \leq 0.95$, les temps sont divisés par au moins 2 pour des tailles de données très grandes. De plus, en examinant D, on voit que la la taille du dictionnaire est très réduite, puisque le taux de compression est de 25 pour $\mu_0 = 0.95$!

Pour illustrer l'algorithme adaptatif, on l'applique pour détecter un saut brutal dans la structure de dépendance et pour suivre une évolution continue de cette structure. Dans le premier cas, les variables aléatoire utilisées précédemment deviennent soudainement indépendantes, puis à nouveau dépendantes avec une rotation toutefois différente. Pour le deuxième cas, l'angle de rotation suit une courbe gaussienne sur l'intervalle d'observation, son maximum étant atteint au milieu de l'intervalle. La figure (1E) illustre la poursuite, et la figure

(1F) le comportement pour les ruptures. Les facteurs d'oubli choisis sont $\gamma = 0.05, \mu = \zeta = 0.1$ et le paramètre de cohérence $\mu_0 = 0.95$. Ces figures illustrent les comportements usuels de l'adaptivité : compromis vitesse de convergence-variance, retard en poursuite. L'algorithme reste à être étudié.

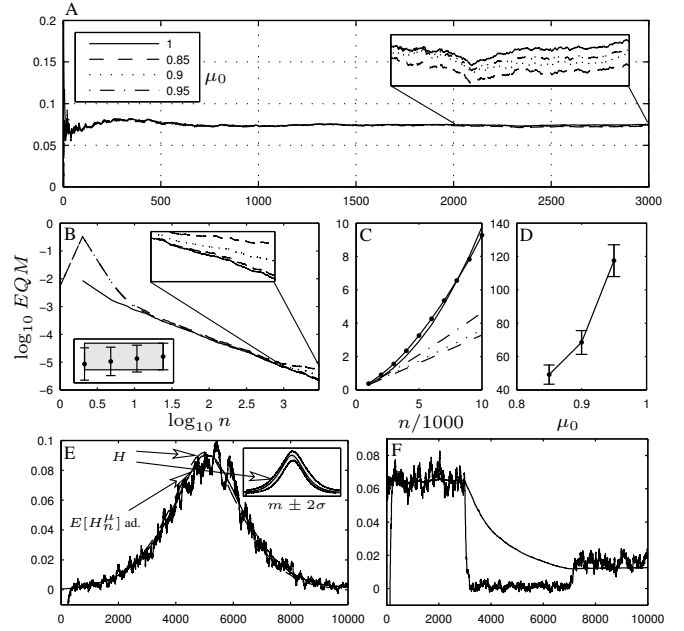


FIGURE 1 – A : traces typiques pour différentes valeurs de μ_0 et un zoom. B : EQM en échelle loglog et un zoom pour différentes valeurs de μ_0 . L'insert bas montre la moyenne finale $\pm 2\sigma$ pour les trois valeurs de μ_0 et HSIC ($\mu_0 = 1$). C : temps de calcul (u. a.). Les points épais correspondent à $\mu_0 = 0.99$. D : taille moyenne $\pm 2\sigma$ du dictionnaire pour $n = 3000$. E : Illustration en poursuite. Insert : bande moyenne $\pm 2\sigma$ et paramètre H à poursuivre. F : Versions adaptative et récursive face à des changements brusques.

Références

- [1] Y. Engel, S. Mannor, and R. Meir. The kernel recursive least-squares algorithm. *IEEE Trans. on Signal Processing*, 52(8) :2275–2284, 2004.
- [2] S. Fine and K. Scheinberg. Efficient svm training using low-rank kernel representations. *Journal of Machine Learning Research*, 2 :243–264, 2001.
- [3] R. Fortet. *Vecteurs, fonctions et distributions aléatoires dans les espaces de Hilbert*. Hermès, 1995.
- [4] A. Gretton, R. Herbrich, A. Smola, O. Bousquet, and B. Schölkopf. Kernel methods for measuring independence. *Journal of Machine Learning Research*, 6 :2075–2129, 2005.
- [5] A. Rényi. On measures of dependence. *Acta Math. Acad. Sci. Hungar.*, 10(2) :441–451, september 1959.
- [6] C. Richard, J. C. M. Bermudez, and P. Honeine. Online prediction of time series data with kernels. *IEEE Trans. on Signal Processing*, 57(3) :1058–1067, Mar 2009.
- [7] B. Schölkopf and A. J. Smola. *Learning with kernels*. MIT Press, Cambridge, Ma, USA, 2002.
- [8] L. Song, A. Smola, A. Gretton, J. Bedo, and K. Borgwardt. Feature selection via dependence maximization. *Journal of Machine Learning Research*, 13 :1393–1434, 2012.
- [9] I. Steinwart and A. Christmann. *Support vector machines*. Springer, 2008.
- [10] R. Vicente, M. Wibrals, M. Lindner, and G. Pipa. Transfer entropy—a model-free measure of effective connectivity for the neurosciences. *Journal of Computational Neuroscience*, 30(1) :45–67, 2011.