

# Prétraitement de données spatialement imprécises pour une classification supervisée basée sur les images satellitaires

Jannai Tokotoko\*, Frédéric Flouvat\*, Claire Goiran\*, Laetitia Hédouin\*\*, Antoine Collin\*\*\*, Nazha Selmaoui-Folcher \*

\*ISEA - New Caledonia University, BP R4, 98851, Nouméa, Nouvelle-Calédonie  
prénom.nom@univ-nc.nc,  
<http://pages.univ-nc.nc/~nom>

\*\*USR 3278 CNRS EPHE UPVD CRIOBE, 98729 Papetoai, Moorea, Polynésie Française  
laetitia.hedouin@criobe.pf

\*\*\*Laboratoire de Géomorphologie et Environnement Littoral,  
CNRS UMR 8586 Prodig, Ecole Pratique des Hautes Etudes, Dinard, France  
antoine.collin@ephe.sorbonne.fr

**Résumé.** Dans un problème de classification supervisée, les données d'apprentissage proviennent souvent d'inventaires acquis sur le terrain par des experts du domaine. Toutefois, la localisation de ces inventaires est approximative (en raison de la précision intrinsèque des GPS portables utilisés). Cette imprécision spatiale est particulièrement problématique lorsque ces données sont utilisées pour entraîner un classifieur sur des images satellitaires très haute résolution (THR). En effet, la précision spatiale des inventaires peut être dans certains cas bien inférieure à celles de ces images. Dans ce papier, nous proposons trois approches visant à améliorer la précision spatiale des données terrain via des pré-traitements. Le principe est d'exploiter les images satellitaires THR disponibles pour corriger spatialement les données terrain. Nos expérimentations mettent en avant l'intérêt de ces pré-traitements sur un jeu de données constitué de 24 inventaires d'habitats coralliens et une image satellitaire THR (WorldView-2).<sup>1</sup>

## 1 Introduction

Le suivi et la protection de l'environnement sont des enjeux majeurs en réponse aux changements climatiques globaux. Dans ce contexte, l'intérêt de l'imagerie satellitaire pour effectuer un suivi au long terme de grands espaces, à un coût relativement faible, n'est plus à démontrer. Ce type d'approches s'appuie sur des inventaires réalisés sur le terrain (données terrain) pour identifier les habitats et construire une carte de distribution de ceux-ci en utilisant un classifieur supervisé. Nous sommes dans le cadre d'une classification *multi-target* (appelée aussi multidimensionnelle ou multi-objectifs). Les instances sont les pixels de l'image caractérisés par leurs valeurs radiométriques (une valeur pour chaque bande spectrale). Les classes

---

1. Ce travail a été soutenu par le Labex "CORAIL" et la fondation DigitalGlobe (*Satellite image courtesy of the DigitalGlobe Foundation*)

## Prétraitement de données inventaires spatialement imprécises

à prédire sont les habitats ou les groupes d'habitats (un pixel pouvant représenter une zone avec plusieurs habitats). Chaque classe peut avoir plusieurs valeurs (pourcentage d'un habitat donné dans le pixel). Pour entraîner ces classifieurs, des données terrain sont nécessaires. Une approche classique pour de tels inventaires écologiques est le LIT ("Line Intercept Transect"). Elle répertorie les différents habitats rencontrés le long d'un "transect" (un segment de droite matérialisé par un décamètre sur le terrain). La figure 1 montre comment cette approche est utilisée pour un récif corallien.

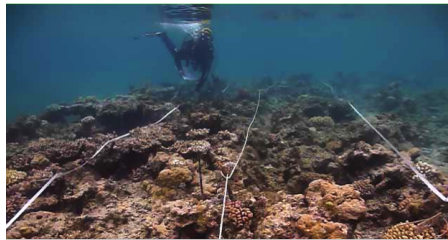


FIG. 1: Exemple d'inventaire réalisé sur un récif corallien.

Cependant, la position géographique de ces inventaires n'est pas exacte. Elle dépend de la précision des GPS utilisés sur le terrain. Cette imprécision spatiale peut être très importante au regard de la résolution des images satellitaires (p.ex. GPS avec une précision à 5 m, et image avec une résolution 0. m). Elle génère ainsi un biais important lors de l'entraînement du classifieur (Mustière, 2014). Des méthodes ont été développées pour traiter des données spatialement imprécises, souvent dans la phase de classification. Par exemple, la logique floue a été utilisée pour construire des classifieurs et les comparer (Hagen-Zanker et al., 2005; Fritz et See, 2005; Pontius Jr et Cheuk, 2006). (Blewitt et Taylor, 2002) s'intéressent au recalage de points localisés par un GPS mais toujours par rapport à des données de références. A notre connaissance, il n'existe pas de méthode directe et non supervisée pour corriger ce type d'imprécision spatiale. Cela reste un problème ouvert pour les experts et pour l'analyse de ces données (Hedley et al., 2016).

Si la position des inventaires était précise, chaque transect serait associé à une séquence de pixels de l'image. Dans notre cas, en raison de l'imprécision spatiale, plusieurs milliers de séquences de pixels (appelés transects candidats dans la suite de l'article) peuvent correspondre à un transect fait sur le terrain. Notre objectif est de sélectionner la séquence correspondant le mieux aux habitats rencontrés. Pour cela, nous devons évaluer si les valeurs spectrales sont cohérentes par rapport aux habitats. Malheureusement, même si nous savons que chaque habitat à une réponse spectrale spécifique, nous ne connaissons pas cette valeur. Nous ne pouvons donc qu'étudier l'homogénéité des valeurs spectrales et leur distribution.

Dans cet article, nous proposons des pré-traitements visant à améliorer cette précision spatiale des inventaires terrain à partir de l'image satellitaire Très Haute Résolution (THR) utilisée pour la classification. Tout d'abord, tous les transects candidats (i.e. les localisations possibles) sont extraits de l'image à partir d'une méthode classique de traitement d'images. Ensuite, pour chaque séquence de pixels extraite, nous générons la séquence de valeurs radiométriques associée, et calculons la séquence d'habitats correspondante à partir des données terrain. Pour finir, nous comparons ces deux séquences et sélectionnons la meilleur pour entraîner le classi-

fiur. Trois méthodes ont été étudiées pour comparer ces deux séquences. La première méthode fait un clustering des pixels (par rapport à leurs valeurs spectrales) et compare simplement les transitions d’habitats avec celles des clusters. La deuxième fait aussi un clustering des pixels, mais considère en plus les habitats afin de caractériser chaque cluster extrait. Au final, elle compare les transitions d’habitats mais aussi leur composition. La dernière méthode considère que le clustering des pixels est donné par les données terrain, et étudie simplement la qualité de celui-ci (l’inertie inter et intra classe).

Nos expérimentations montrent l’impact positif de ces pré-traitements sur les performances de 15 classifieurs mis à disposition par l’outil *MEKA* (Read et al., 2016).

## 2 Données et Problématique

Une base de données terrain est un ensemble structuré d’informations spatiales incluant des objets géographiques (p.ex. des trajectoires, des zones, ...), des habitats (p.ex. rocher, sable, corail branchu, ...) et la distribution spatiale de ces habitats par rapport aux objets géographiques. Elle peut être définie comme un triplet  $BD = (D_s, D_h, D_d)$  où  $D_s$  est la dimension spatiale,  $D_h$  est la dimension décrivant les habitats et  $D_d$  est la dimension précisant la distribution spatiale des habitats. Dans cet article, pour faciliter la compréhension, nous ne considérerons que des segments comme objets géographiques, mais notre approche est généralisable à tout type de zones (un polygone étant un ensemble fini non vide de segments de droites). Ces segments représentent des trajectoires suivies par les experts sur le terrain et sur lesquelles sont inventoriés les différents habitats.

La dimension spatiale est associée à un domaine de valeurs dénoté  $dom(D_s) = \{T_1, \dots, T_n\}$  où  $T_i$  pour  $i \in [1..n]$  est un *transect* (un segment de droite) délimité par deux points  $(x, y)$  et  $(x', y')$ , i.e.  $T_i = [(x, y), (x', y')]$ . Chacun de ces points est associé à une erreur de mesure  $\delta$ , i.e.  $x_{real} \in [x_{field} - \delta, x_{field} + \delta]$  et  $y_{real} \in [y_{field} - \delta, y_{field} + \delta]$ , où  $(x_{real}, y_{real})$  sont les coordonnées réelles,  $(x_{field}, y_{field})$  sont les coordonnées mesurées par le GPS, et  $\delta$  est la précision de celui-ci. La dimension décrivant les habitats est associée à un domaine de valeurs dénoté  $dom(D_h) = \{h_1, h_2, \dots, h_m\}$  où  $h_i$  pour  $i \in [1..m]$  est un *type d’habitats*. La dimension précisant la distribution spatiale des habitats est associée à un domaine de valeurs dénoté  $dom(D_d) = \mathbb{R}$  représentant la distance à laquelle à lieu le changement d’habitat.

Transect	Habitat	Distance par rapport au début du transect
$T_1$	rck	30
$T_1$	ca	380
$T_1$	rck	0
$T_1$	r	550
$T_1$	rck	570
$T_1$	mc	575
$T_1$	rck	620
$T_1$	ce	630
$T_1$	r	660
...	...	...

TAB. 1: Base de données terrain répertoriant des habitats coralliens le long d’un transect

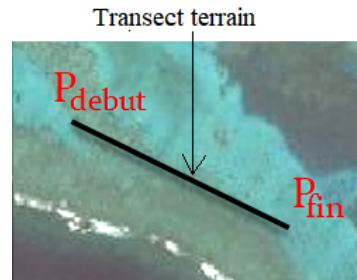


FIG. 2: Transect  $T_1$  suivi par les experts lors de leur inventaire.

## Prétraitement de données inventaires spatialement imprécises

Pour illustrer ces définitions, nous utilisons l'exemple du tableau 1. Il représente un inventaire réalisé sur un récif corallien. Seulement une partie des données du transect  $T_1$  est présentée dans le tableau. Sa localisation est quant à elle affichée dans la figure 2. Dans ce tableau,  $D_s = \{T_1\}$ ,  $D_h = \{rck, ca, r, mc, ce\}$  et  $D_d = \{30, 380, 0, 550, 570, 575, 620, 630, 660\}$ . Les habitats  $rck$ ,  $ca$ ,  $r$ ,  $mc$  et  $ce$  correspondent respectivement à des rochers, des algues corallines, des débris de coraux, des coraux massifs et des coraux encroûtants. Ce tableau indique que des rochers sont présents sur les 30 premiers centimètres du transect, suivi d'algues corallines sur une longueur de 350 cm ( $380 - 30 = 350$  cm), suivi de débris de coraux sur une longueur de 80 cm ( $0 - 380 = 80$  cm), ...

L'image satellitaire associée à cette base de données est une image dont la résolution est beaucoup plus élevée que la précision du GPS utilisé sur le terrain (p.ex. 0. m contre une précision à 5 m pour le GPS). Chaque pixel  $P$  de cette image est associé à un n-uplet  $(b_1, b_2, \dots, b_s)$  où  $b_i \in \mathbb{R}$  est la valeur radiométrique associée à la  $i$ ème bande spectrale de l'image (p.ex. rouge, vert, bleu, proche infrarouge, etc). Il est également associé à une coordonnée géographique  $(x_P, y_P)$  (celle de son centre). Nous considérons que l'image satellitaire THR et la base de données terrain sont acquises à des dates proches (et aucun évènement important entre les deux).

**Problème** Soit  $BD = (D_s, D_h, D_d)$  une base de données terrain. Soit  $T$  un transect de cette base de données, i.e.  $T \in \text{dom}(D_s)$ , et soit  $\delta$  l'imprécision spatiale du GPS utilisé sur le terrain pour le localiser. Soit  $I$  une image satellitaire multi-spectrale THR géo-référencée telle que  $T \subset I$  et  $\text{resol}_I \ll \delta$ , où  $\text{resol}_I$  est la taille des pixels de l'image (i.e. sa résolution spatiale). L'objectif est d'améliorer la précision spatiale des données terrain en exploitant l'image  $I$ , i.e. de trouver la localisation de  $T$  à une précision de  $\pm \text{resol}_I$ .

### 3 Prétraitement des données terrain à partir d'une image satellitaire

Dans cette section, nous allons introduire la chaîne de prétraitements développée pour recalculer des données terrain spatialement imprécises par rapport à une image satellitaire THR. Pour chaque transect  $T \in D_s$ , ce processus est composé des étapes suivantes :

1. génération de tous les emplacements possibles de  $T$  dans l'image  $I$ , noté  $E_{cand}(T)$ , en fonction de l'imprécision spatiale  $\delta$ .
2. pour chaque emplacement possible (transect candidat),  $T_{cand} \in E_{cand}(T)$  :
  - (a) génération de la séquence de pixels correspondante,  $Pix_I(T_{cand})$ , et de sa séquence de valeurs radiométriques,  $Rad_I(T_{cand})$
  - (b) génération d'une séquence  $Hab_T(T_{cand})$  représentant la distribution de l'habitats des données terrain dans chaque pixel de  $Pix_I(T_{cand})$
  - (c) estimation de la similarité entre  $Rad_I(T_{cand})$  et  $Hab_T(T_{cand})$
3. classement des transects candidats en fonction de la similarité et sélection du meilleur candidat.

### 3.1 Identification des transects candidats dans l'image

Pour un transect terrain donné, cette étape consiste à extraire les séquences de pixels correspondant aux positions possibles du transect dans l'image satellitaire.

En raison de l'imprécision du GPS utilisé sur le terrain, la localisation d'un inventaire est en réalité approximative. Comme le montre la figure 3, les coordonnées des points de départ et d'arrivée du transect ( $P_{debut}$  et  $P_{fin}$  de la figure 3) sont données avec une précision de  $\pm \delta$ . Les cercles  $E_{debut}$  et  $E_{fin}$  représentent cette zone d'incertitude autour des extrémités du transect. Comme la résolution du pixel de l'image est bien inférieure, il existe une multitude de positions possibles pour le transect considéré. Nous notons  $E_{cand}(T)$  l'ensemble des positions possibles du transect  $T$ .

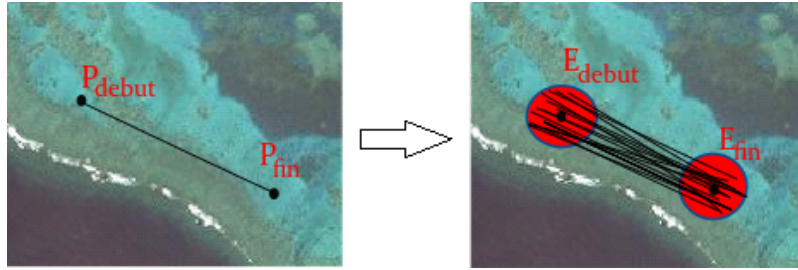


FIG. 3: Ensemble des transects candidats pour un transect  $T$  et une imprécision spatiale  $\delta$ .

Pour pouvoir comparer les transects candidats à l'inventaire réalisé, il est nécessaire d'étudier plus en détail les valeurs spectrales des pixels, et donc d'extraire la séquence de pixels associée à chaque transect candidat. Pour cela, il faut tracer le segment de droite passant par les deux pixels associés aux points de départ et d'arrivée du transect candidat, et extraire la séquence de pixels correspondants. Ce traitement n'est pas trivial car il faut approximer un segment de droite continu dans un plan discret. Nous nous sommes appuyés sur l'algorithme de Bresenham (1965) de tracé de segments dans un espace discret. La figure 4 montre un exemple de séquence de pixels associée à un segment.

Suite à ce traitement, nous obtenons donc pour chaque transect candidat  $T_{cand} \in E_{cand}(T)$  d'un transect inventorié  $T \in D_s$ , une séquence de  $k$  pixels  $Pix_I(T_{cand}) = \langle P_1, P_2, \dots, P_k \rangle$  avec  $P_1 \in E_{debut}$  et  $P_k \in E_{fin}$ . Chaque pixel étant associé à un vecteur de valeurs spectrales, chaque séquence  $Pix_I(T_{cand})$  est liée à une séquence de valeurs spectrales  $Rad_I(T_{cand}) = \langle R_1, R_2, \dots, R_k \rangle$  où  $R_i \in R^s$  est le vecteur de valeurs spectrales du  $i$ -ème pixel ( $s$  est le nombre de bandes spectrales de l'image).

### 3.2 Redécoupage des données de l'inventaire en fonction des transects candidats

Une image est une représentation de la réalité dans un espace discret alors qu'un inventaire est exprimée de façon continue (un segment avec des transitions d'habitats). Afin de pouvoir comparer l'inventaire réalisé et un transect candidat associé à des pixels de l'image, il est nécessaire de découper l'inventaire en fonction des pixels considérés. Pour cela, il faut d'abord

connaître l'endroit précis où chaque pixel "coupe" le segment considéré. La figure 4 illustre un exemple de ce découpage d'habitats suivant les pixels d'un candidat. Le pixel  $P_2$  contient la portion du transect délimitée par  $d1$  et  $d2$ . Une fois  $d1$  et  $d2$  calculées, on peut déduire les habitats (et leurs proportion) associés au pixel considéré en explorant la base de données terrain.

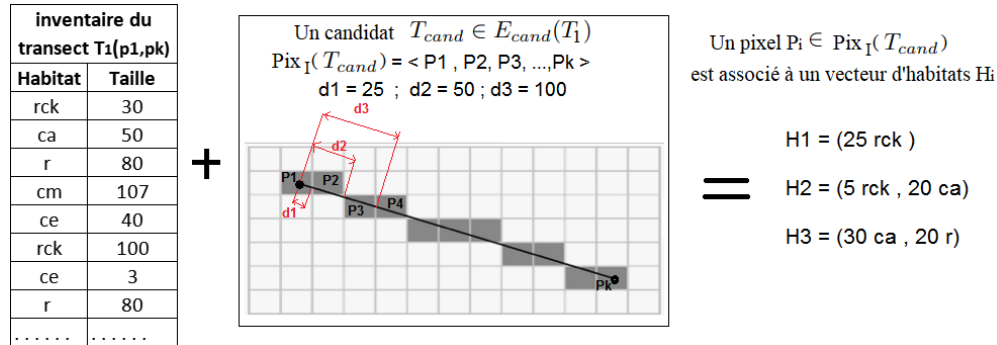


FIG. 4: Exemple illustrant la proportion d'habitats associée à chaque pixel.

Ce traitement permet pour chaque transect candidat  $T_{cand} \in E_{cand}(T)$  de générer une séquence d'habitats  $Hab_T(T_{cand}) = \langle H_1, H_2, \dots, H_k \rangle$ , où  $H_i$  est un vecteur décrivant la proportion de chaque habitat observé sur le terrain dans le pixel  $P_i$  de  $Pix_I(T_{cand})$ ,  $\forall i \in [1..k]$ . Soit  $H_i(j)$  la proportion du j-ème habitat de  $D_h$  dans le pixel  $P_i$ .

Par exemple, dans le Fig. 4,  $H_2 = (5rck; 20ca)$  correspond aux habitats associés au pixel  $P_2$ . Il y a 5 cm de roche dans ce pixel (puisque  $d1 = 25cm$  et la taille du premier habitat, roche, est de 30 cm), mais aussi 20 cm d'algues corallines (puisque  $d2 = 50cm$ ).  $H_2(1) = 5$  et  $H_2(2) = 20$  puisque  $D_h = \{rck, ca, r, mc, ce, \dots\}$ .

### 3.3 Estimation de la similarité entre transects candidats et inventaire : trois approches

Dans cette section, nous proposons trois méthodes pour estimer la similarité en la séquence de valeurs radiométriques  $Rad_I(T_{cand})$  et la séquence d'habitats  $Hab_T(T_{cand})$ , et ainsi mesurer la similarité entre un transect candidat et l'inventaire réalisé sur le terrain.

#### 3.3.1 Méthode basée uniquement sur les transitions d'habitats (M1)

Cette première méthode fait un clustering des pixels (par rapport à leurs valeurs spectrales) et compare simplement les transitions d'habitats avec celles des clusters.

En télédétection, un habitat (ou groupement d'habitats) est supposé avoir une signature spectrale particulière. Un clustering des pixels en fonction de leurs valeurs radiométriques génère des groupes de pixels avec des valeurs radiométriques similaires, c.à.d. des habitats supposés. Dans ce travail, nous utilisons la méthode des k-means pour trouver les changements d'habitats, avec  $k$  le nombre de changements dans les données terrain. Nous transfor-

mons donc la séquence  $Rad_I(T_{cand})$  en séquence de clusters  $S_{clusters} = \langle cluster(P_1), \dots, cluster(P_k) \rangle$ , où  $cluster(P)$  est l'identifiant du cluster auquel est affecté le pixel  $P$ .

En théorie, s'il y a une bonne correspondance entre les clusters et les habitats réels, la transition des clusters de  $S_{clusters}$  devrait ressembler à la transition des habitats de  $Hab_T(T_{cand})$ . Nous transformons donc la séquence  $S_{clusters}$  et la séquence  $Hab_T(T_{cand})$  en vecteurs binaires représentant les changements de clusters et les changements d'habitats. Soit  $V_{clusters} = (V_1, V_2, \dots, V_k)$  (resp.  $V_{habitats}$ ) le vecteur de changements de clusters associé à  $S_{clusters}$  (resp.  $Hab_T(T_{cand})$ ). Nous avons donc  $V_i = 0$  si  $cluster(P_i) = cluster(P_{i+1})$  (resp.  $H_i = H_{i+1}$ ), et  $V_i = 1$  sinon,  $\forall i \in [1 \dots k - 1]$ . Le tableau 2 illustre cette transformation sur un exemple. La première et la deuxième colonne du tableau représentent  $Hab_T(T_{cand})$  (sans les proportions) et  $V_{habitats}$ . La troisième colonne et la dernière représentent  $S_{clusters}$  et  $V_{clusters}$ .

Habitats	Vecteur Habitat	Cluster	Vecteur Cluster
mc+dc	0	cluster_1	0
r	1	cluster_2	1
r	0	cluster_2	0
acb+r	0	cluster_4	1
r	0	cluster_4	0
mc	0	cluster_3	1
mc+dc	0	cluster_3	0
acb	1	cluster_1	0
acb	1	cluster_4	1
acd		cluster_4	

TAB. 2: Exemple de distribution habitats et distribution clusters

Pour comparer ces deux vecteurs de changement d'habitats, nous utilisons une mesure de similarité de type II (Rifqi, 2010) telle que celle de Rogers et Tanimoto (pour information, plusieurs mesures de type II ont été testées et fournissent les mêmes résultats). Ces mesures ont l'avantage de prendre en compte les différences entre les deux vecteurs, mais aussi les valeurs identiques (que soit un 1 ou un 0). Soit  $X$  et  $Y$  deux vecteurs binaires de  $\{0, 1\}^p$ . On note  $a$  le nombre de 1 en commun,  $b$  le nombre de 1 dans  $X$  mais pas dans  $Y$ ,  $c$  le nombre de 0 dans  $X$  mais pas dans  $Y$ , et  $d$  le nombre de 0 en commun. La mesure de similarité de Rogers et Tanimoto entre les vecteurs binaires  $X$  et  $Y$  est  $\frac{a+d}{a+d+2(b+c)}$ . Au final, le transect candidat ayant la similarité la plus élevée par rapport aux données terrain est conservé.

### 3.3.2 Méthode basée sur la réponse spectrale des habitats et leur distribution (M2)

Cette deuxième méthode fait aussi un clustering des pixels, mais considère en plus les habitats afin de caractériser chaque cluster extrait. Au final, elle compare les transitions d'habitats mais aussi leur composition.

Elle utilise aussi la méthode des k-means pour identifier des clusters (habitats supposés) et transformer  $Rad_I(T_{cand})$  en  $S_{clusters}$ . Le nombre de clusters en entrée de la méthode des k-means est cette fois égal au nombre d'habitats différents dans la dimension  $D_h$ . Les clusters représentent ainsi des habitats supposés. On connaît leur signature spectrale mais pas encore le détail des habitats qui la compose.  $Hab_T(T_{cand})$  précise la répartition des habitats observés sur le terrain en fonction de chaque pixel. Pour chaque pixel  $P_i \in Pix_I(T_{cand})$ , nous connaissons donc les valeurs spectrales  $R_i$  de chaque  $cluster(P_i) \in S_{clusters}$ , mais aussi le détail des habitats  $H_i \in Hab_T(T_{cand})$  qui composent  $P_i$ . Par conséquent, il suffit de parcourir tous les

pixels et de calculer la proportion totale d'habitats associés à chaque  $cluster(P_i)$  (proportion normalisée).

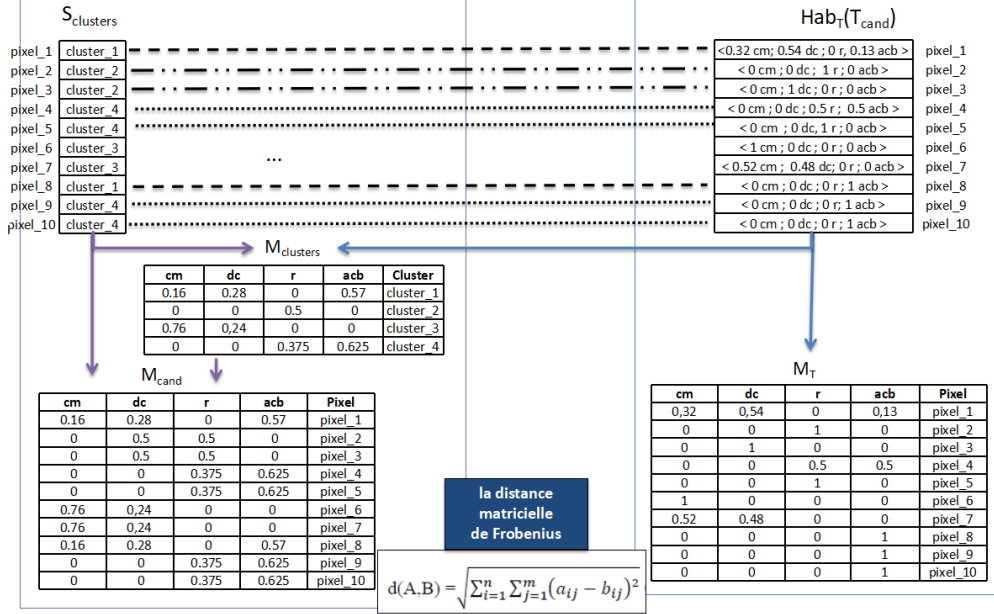


FIG. 5: Exemple illustrant le calcul des matrices  $M_{cand}$  et  $M_T$ .

Soit  $M_{clusters}$  cette matrice clusters/proportions d'habitats. La proportion du  $j$ -ème habitat de  $D_h$  associée au cluster  $C_l \in S_{clusters}$  est  $M_{clusters}(l, j) = \frac{\sum_{P_i \in C_l} H_i(j)}{\sum_j \sum_{P_i \in C_l} H_i(j)}$  (matrice normalisée), où  $H_i(j)$  est la proportion de l'habitat  $j$  qui apparaît dans le pixel  $P_i$  (ce pixel étant dans le cluster  $C_l$ ). Cette matrice permet d'associer au transect candidat une matrice  $M_{cand}$  précisant pour chaque pixel/cluster la proportion d'habitats supposés. Cette matrice peut ensuite être comparée à la matrice "terrain", notée  $M_T$ , dérivée de la séquence  $Hab_T(T_{cand})$ . La distance matricielle utilisée ici est celle de Frobenius (Horn et Johnson, 2012). La figure 5 illustre la construction de ces deux matrices et leur comparaison par la distance de Frobenius.

### 3.3.3 Méthode basée sur l'homogénéité des valeurs spectrales des habitats (M3)

Cette dernière méthode considère que le clustering des pixels est donné par les données terrain, et étudie simplement la qualité de celui-ci. Elle évalue l'homogénéité des valeurs spectrales des pixels associées à un même habitat dans le transect candidat. Nous croisons pour cela  $Hab_T(T_{cand})$  et  $Pix_I(T_{cand})$ , et mesurons l'inertie intra-classe et inter-classe des pixels de  $Pix_I(T_{cand})$  groupés en fonction des habitats de  $Hab_T(T_{cand})$ . Plus l'inertie intra-classe est petite (inertie inter-classe est grande) plus la classification habitat/pixel est bonne. Dans notre cas, les pixels  $P_i \in Pix_I(T_{cand})$  sont les individus, les classes sont les différents habitats (ou groupement d'habitats) représentés dans  $Hab_T(T_{cand})$ , les caractéristiques (attributs)



sont les 8 bandes spectrales des pixels dont les valeurs sont données par  $Rad_I(T_{cand})$ . Soit  $m$  le nombre d'habitats différents (les classes),  $s$  le nombre de valeurs spectrales (les attributs), et  $k$  le nombre de pixels de  $Pix_I(T_{cand})$  (les individus). On note  $R_i = (r_{i_1}, \dots, r_{i_s})$  le vecteur correspondant aux valeurs radiométriques du pixel  $P_i \in Pix_I(T_{cand})$ . Soit  $C_l$  une classe d'habitats et  $n_l$  le nombre de pixels de classe  $C_l$ . Son centre de gravité est  $\mu_l = \frac{1}{n_l} \sum_{i=1}^{n_l} R_i$ , et son inertie est  $I_l = \sum_{P_i \in C_l} d^2(R_i, \mu_l)$  avec  $d^2(R_i, \mu_l) = \sum_{j=1}^s (r_{i_j} - \mu_{l_j})^2$ . Plus l'inertie intra-classe, notée  $J_a = \sum_{i=1}^m I_i$ , est faible plus les clusters sont homogènes et éloignés les uns des autres. De même, on pourra calculer l'inertie inter-classe notée  $J_b = \sum_l n_l d^2(\mu, \mu_l)$  avec  $\mu = \frac{1}{k} \sum_{i=1}^k R_i$ . Plus l'inertie inter-classe est grande et plus les clusters sont éloignés.

Le meilleur transect candidat est celui pour lequel  $J_a$  est minimum ( $J_b$  est maximum).

## 4 Expérimentations

**Données utilisées** Dans nos expérimentations, nous étudions les données issues du suivi d'un récif corallien en Nouvelle-Calédonie. Les données terrain répertorient 33 habitats selon 24 transects. Les biologistes ont utilisé la méthode du LIT (Line Intercept Transect) de English et al. (1997) pour effectuer cet inventaire. Les positions de départ et d'arrivée des transects ont été acquises avec un GPS Garmin dont la précision spatiale est 5 mètres. En plus de ces données, une image satellitaire THR WorldView-2 a été utilisée. Cette image a été acquise 3 jours après l'inventaire. Elle possède 8 bandes spectrales (proche infrarouge 1, rouge, vert, bleu, red-edge, jaune, bande côtière, et proche infrarouge 2) et sa résolution spatiale est 0. m. Cette image a subi les corrections classiques en télédétection (géométrique, atmosphérique, radiométrique, et scintillement des vagues) et a été fusionnée avec l'image panchromatique (*pansharpening*).

**Protocol expérimental** La chaîne de prétraitements proposée a été implémentée dans les logiciels *ENVI* et *KNIME* (Berthold et al., 2009). Le logiciel de traitements d'images *ENVI* est utilisé pour extraire les transects candidats de l'image (via un script *IDL*) et le logiciel d'analyse de données *KNIME* est utilisé pour faire les autres traitements (via un workflow).

Afin de valider nos approches, nous avons étudié l'impact des trois méthodes précédentes sur 15 algorithmes de classification supervisée *multi-target* implémentés dans l'outil *MEKA* (Read et al., 2016) : BCC (*Bayesian Classifier Chains*), CC (*Classifier Chains*), CR (*Class-Relevance*), NSR (*Nearest Set Replacement*) et RAKELd (Tsoumakas et al., 2011). Chaque méthode s'appuie sur des méthodes de classification mono-label (J48 par défaut). Nous avons aussi testé en entrée de ces algorithmes NB (*Naive Bayes*) et deux algorithmes de classification floue : ENORA de Jiménez et al. (2014) et FURIA de Hühn et Hüllermeier (2009).

Dans ces expérimentations, les instances sont les pixels (800 environs) caractérisés par leurs valeurs radiométriques (8), et les classes sont des habitats coralliens ou des combinaisons d'habitats (puisque un pixel peut être constitué de plusieurs habitats). Chaque classe peut prendre plusieurs valeurs (une proportion d'habitat dans le pixel considéré). En effet, nous avons utilisé *KNIME* pour discrétiser les proportions d'habitats en 12 catégories (0, ]0, 10%], ]10%, 20%], ..., 100%) représentant le pourcentage du pixel "occupé" par l'habitat.

Plusieurs mesures sont proposées dans *MEKA* pour évaluer la qualité de la prédiction (Read et al., 2011). Dans nos expérimentations, nous présenterons les résultats obtenus avec *Ham-*

*ming Score* (HS) et *Exact Match* (EM). La première mesure indique le pourcentage de bonnes prédictions sur le nombre total de classes. La deuxième mesure *Exact Match* (EM) indique le pourcentage d'échantillons dont les étiquettes ont été exactement prédites.

multi-target	Hamming Score				Exact Match			
	terrain	M1	M2	M3	terrain	M1	M2	M3
BCC - J48	0.97	0.97	0.97	0.97	0.047	0.13	0.055	0.103
BCC - NB	0.975	0.889	0.89	0.871	0.006	0.006	0.0	0.01
BCC - ENORA	0.909	0.929	0.912	0.87	0.003	0.0	0.0	0.0
CC - J48	0.964	0.969	0.969	0.97	0.051	0.109	0.052	0.108
CC - ENORA	0.975	0.968	0.969	0.971	0.004	0.0	0.0	0.0
CR - J48	0.975	0.969	0.97	0.97	0.049	0.125	0.055	0.097
CR - ENORA	0.975	0.891	0.891	0.956	0.001	0.0	0.0	0.0
NSR - J48	0.939	0.96	0.953	0.958	0.133	0.314	0.164	0.22
NSR - FURIA	0.958	0.956	0.953	0.958	0.111	0.197	0.148	0.216
NSR - ENORA	0.961	0.956		0.958		0.197		0.216
RAkELd - J48	0.974	0.967	0.969	0.97	0.058	0.15	0.056	0.106
<b>Moyenne</b>	<b>0.961</b>	<b>0.947</b>	<b>0.945</b>	<b>0.947</b>	<b>0.0</b>	<b>0.11</b>	<b>0.053</b>	<b>0.098</b>

TAB. 3: Résultats de la classification sur les données brutes et les données générées par nos 3 approches

**Impact des prétraitements sur les classifieurs** Le tableau 3 affiche les performances des algorithmes de classification étudiés sur les données brutes (*terrain*) et sur les données pré-traitées avec nos méthodes (méthode basée sur les vecteurs de transitions *M1*, méthode basée sur les matrices d'habitats *M2* et méthode basée sur l'inertie *M3*). La mesure de Hamming Score (HS) est assez élevée en raison du nombre important de classes. En effet, nous avons en moyenne 50 classes (habitats ou regroupement d'habitats) pour 800 individus (pixels) en moyenne. Malgré ce nombre de classes élevé, nous constatons que la précision augmente pour la plus part des algorithmes. La valeur EM (qui est une mesure stricte) double dans certains cas avec nos prétraitements (p.ex. *M3*) par rapport aux classifications effectuées sur les données brutes.

La figure 6 montre les performances moyennes des classifieurs et classifieurs flous en fonction des prétraitements effectués sur les données. On remarque que les performances diminuent légèrement lors de l'utilisation des méthodes floues, et ceci même sur les données brutes. On constate aussi que nos prétraitements améliorent globalement les résultats des classifieurs. Les méthodes *M2* et *M3* semblent plus particulièrement efficaces. Ces résultats montrent aussi que nos prétraitements sont complémentaires des algorithmes de classification floue pour traiter des données spatialement imprécises.

## 5 Conclusion et perspectives

Dans cet article, nous nous sommes intéressés au problème de l'imprécision spatiale des données collectées sur le terrain et son impact sur la classification supervisée d'images satellitaires. Nous avons proposé des prétraitements pour améliorer cette précision spatiale en exploitant une image satellitaire THR. Ce processus s'appuie sur une identification des transects candidats dans l'image, suivi d'un redécoupage des données terrain en segments d'habitats correspondant à ces pixels. Un clustering est aussi fait pour détecter des habitats supposés à

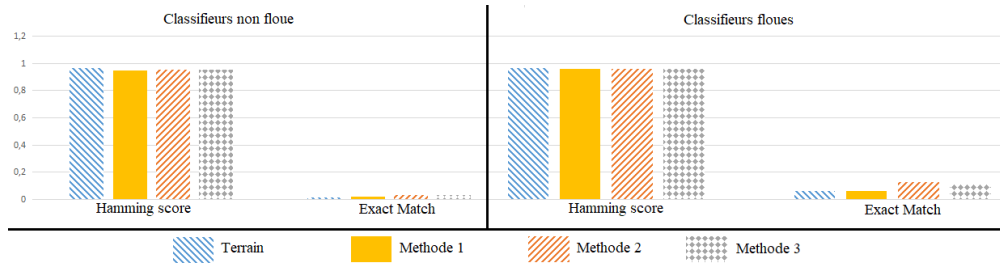


FIG. 6: Comparaison de la performance moyenne des classifieurs et classifieurs floues.

partir des valeurs spectrales associées aux transects candidats. Trois approches ont été proposées et comparées. Ces approches permettent d’extraire la séquence de pixels ressemblant le plus aux données terrain. Des expérimentations ont été réalisées sur un jeu de données réelles afin de mesurer l’impact de nos prétraitements sur la qualité d’une classification supervisée. Nous avons pour cela étudié les performances obtenues par 15 algorithmes de classification *multi-target* de la littérature. Ces résultats montrent que notre approche permet de doubler le nombre d’*Exact-Match*. Ils mettent également en avant la limite des mesures existantes (p.ex. *Hamming-Loss* ou distance de *Levenshtein*) pour évaluer la qualité d’une classification multi-classes avec des valeurs nominales (ou numériques discrétisées).

Les perspectives à ce travail sont nombreuses. Il serait par exemple intéressant d’étudier plus précisément l’impact du nombre de classes, de la discrétisation et du *pansharpening* sur nos prétraitements et les performances des classifieurs. Pour cela, le développement d’une mesure de similarité adaptée à des classes avec des valeurs numériques discrétisées semble important. Il serait également intéressant d’utiliser des indices radiométriques calculés ou des produits dérivés (p.ex. une estimation de la bathymétrie) en plus des valeurs spectrales des images. Une autre perspective serait de combiner nos prétraitements avec des algorithmes de régression, et d’en étudier les performances.

## Références

- Berthold, M. R., N. Cebron, F. Dill, T. R. Gabriel, T. Kötter, T. Meinl, P. Ohl, K. Thiel, et B. Wiswedel (2009). Knime-the konstanz information miner : version 2.0 and beyond. *AcM SIGKDD explorations Newsletter* 11(1), 26–31.
- Blewitt, G. et G. Taylor (2002). Mapping dilution of precision (mdop) and map-matched gps. *International Journal of geographical information science* 16(1), 55–67.
- Bresenham, J. E. (1965). Algorithm for computer control of a digital plotter. *IBM Systems journal* 4(1), 25–30.
- English, S. S., C. C. Wilkinson, et V. V. Baker (1997). *Survey manual for tropical marine resources*. Australian Institute of Marine Science.
- Fritz, S. et L. See (2005). Comparison of land cover maps using fuzzy agreement. *International Journal of Geographical Information Science* 19(7), 787–807.

- Hagen-Zanker, A., B. Straatman, et I. Uljee (2005). Further developments of a fuzzy set map comparison approach. *International Journal of Geographical Information Science* 19(7), 769–785.
- Hedley, J. D., C. M. Roelfsema, I. Chollett, A. R. Harborne, S. F. Heron, S. Weeks, W. J. Skirving, A. E. Strong, C. M. Eakin, T. R. Christensen, et al. (2016). Remote sensing of coral reefs for monitoring and management : a review. *Remote Sensing* 8(2), 118.
- Horn, R. A. et C. R. Johnson (2012). *Matrix analysis*. Cambridge university press.
- Hühn, J. et E. Hüllermeier (2009). Furia : an algorithm for unordered fuzzy rule induction. *Data Mining and Knowledge Discovery* 19(3), 293–319.
- Jiménez, F., G. Sánchez, et J. M. Juárez (2014). Multi-objective evolutionary algorithms for fuzzy classification in survival prediction. *Artificial intelligence in medicine* 60(3), 197–219.
- Mustière, S. (2014). *Intégration de données géographiques*. Ph. D. thesis, Université de Grenoble.
- Pontius Jr, R. G. et M. L. Cheuk (2006). A generalized cross-tabulation matrix to compare soft-classified maps at multiple resolutions. *International Journal of Geographical Information Science* 20(1), 1–30.
- Read, J., B. Pfahringer, G. Holmes, et E. Frank (2011). Classifier chains for multi-label classification. *Machine learning* 85(3), 333–359.
- Read, J., P. Reutemann, B. Pfahringer, et G. Holmes (2016). MEKA : A multi-label/multi-target extension to WEKA. *Journal of Machine Learning Research* 17, 21 :1–21 :5.
- Rifqi, M. (2010). *Mesures de similarité, raisonnement et modélisation de l'utilisateur*. Ph. D. thesis, Université Pierre et Marie-Curie.
- Tsoumakas, G., I. Katakis, et I. Vlahavas (2011). Random k-labelsets for multilabel classification. *IEEE Transactions on Knowledge and Data Engineering* 23(7), 1079–1089.

## Summary

In a supervised classification problem, training data often comes from field inventories acquired by domain experts. However, location of these inventories may be approximate (due to intrinsic precision of portable GPS used). This spatial inaccuracy is particularly problematic when these data are used to train a classifier on very high resolution (VHR) satellite images. Indeed, in some cases, spatial accuracy of inventories may be much lower than the one of images. In this paper, we propose three preprocessing methods to correct this spatial inaccuracy of field inventories. The principle of these approaches is to exploit available VHR satellite images to spatially correct field data. Our experiments highlight the interest of these pretreatments and compare the proposed approaches on a dataset consisting of 24 habitat inventories of coral reefs and a VHR satellite image (WorldView-2).<sup>2</sup>

---

2. This work was supported by the “CORAIL” laboratory of excellence and the DigitalGlobe Foundation (*Satellite image courtesy of the DigitalGlobe Foundation*)