

République Algérienne Démocratique et Populaire

Ministère de l'enseignement supérieur et de la recherche scientifique

Université Mouloud Mammeri de Tizi-Ouzou

Faculté de génie électrique et d'informatique

Département d'informatique

Mémoire

De fin d'études

En vue de l'obtention du diplôme de Master en Informatique

Spécialité : Conduite de projets Informatiques

Thème

Classification de structures arborescentes. Cas d'étude : Documents XML.

Proposé et encadré par : M. Aïtelhadj Ali

Réalisé par : M. Otmani Nassim Abdeldjallal

Membres du jury :

Président :	M. Hamadouche Djamel
Examineurs :	M. Daoui Mohamed Mme Aïtelhadj Fatiha
Rapporteur :	M. Aïtelhadj Ali

Date de soutenance : 08/09/2011

Année universitaire : 2010/2011

Index

Résumé	06
Remerciements	06
Introduction Générale	07
 Chapitre 1 : XML	 09
1. Introduction	09
2. Présentation de la technologie XML	10
2.1. Besoin de structuration	10
2.2. Origine et buts de l'XML	10
2.3. Définition de l'XML	11
3. Document XML	11
3.1. Présentation d'un document XML	11
3.2. Documents XML bien formés	12
3.3. Spécification d'un document XML	12
3.3.1. L'arborescence - les nœuds	12
3.3.2. Les nœuds de type élément	13
3.3.3. Les nœuds de type attribut	14
3.3.4. Le prologue – Les instructions de traitement	14
3.3.5. Les nœuds commentaires	14
3.3.6. Les entités	15
3.3.7. Les sections PCDATA et CDATA	15
3.3.8. Encodage des documents	16
3.3.9. Règles d'écriture : document bien formé	16
3.4. Documents XML valides	17
3.4.1. Les DTDs	17
3.4.2. XML Schema	17
3.4.3. Relax NG	17
3.5. Transformation d'un document	18
3.5.1. XPath	18
3.5.2. XSLT	18
3.5.3. XSL-FO	22
4. Rôle du document XML	22
5. Le document XML : orienté document ou données ?	22
6. Conclusion	23

Chapitre2 : Classification	24
1. Introduction	24
2. Notion de classe	24
3. Définition formelle de la classification	26
4. Le ranking	27
5. Le clustering	28
6. Contextes de classification	29
7. Paramètres d'évaluation des performances d'un système de classification	31
8. Évaluation d'un classifieur bi-classe	33
9. Évaluation de classifieurs multi-classes	35
10. Conclusion	36

Chapitre3 : Classification des documents XML	37
1. Introduction	37
2. Travaux faits	39
3. Exemples d'utilisation	40
3.1. Extraction automatique des DTDs	40
3.2. Regroupement structurel	40
3.3. Gestion des données spatiales	40
3.4. La mise à jour des documents	41
3.5. La détection des modifications	41
3.6. Optimisation des échanges	41
3.7. Recherche d'informations	41
3.8. Vérificateurs d'orthographe (faute humaine)	42
4. Formulation du problème	42
5. Modélisation des données semi structurées	43
5.1. Modélisation des documents XML	43
5.2. Edition des arbres	44
5.2.1. Insertion d'un nœud	46
5.2.2. Suppression d'un nœud	47
5.2.3. Remplacement (substitution) d'un nœud	48
6. Les algorithmes de distances d'édérations	49
6.1. Algorithme de Selkow	50
6.2. Algorithme de Dalamagas	53
6.3. L'algorithme de Chawathe	55
6.3.1. Graphe d'édition	56
6.3.2. Algorithme	59
6.4. L'algorithme de Levenshtein	62
7. L'algorithme de Aïtelhadj	64
7.1. Définition préliminaire	64
7.2. Mesure de similarité structurelle	64
8. Clustering des documents XML	67
8.1. La distance et la similarité structurelle	68
8.2. Les algorithmes de clustering	69
8.2.1. Single link	70
8.2.2. Matrice Associée	71
8.2.3. Algorithme de Prim	72
9. Résumé du chapitre	73
10. Conclusion	74

Chapitre 4 : Conception et réalisation d'un logiciel de classification	75
1. Introduction	75
2. Conception	76
2.1. Présentation du logiciel	76
2.2. Schéma général du logiciel	77
2.3. Explication des étapes du logiciel	78
2.3.1. Pré traitement	78
2.3.2. Comparaison des documents	80
2.3.3. Phase classification	82
2.4. Conception détaillée	85
2.4.1. Diagramme de classes	85
2.4.2. Diagramme de séquence	91
3. Réalisation	92
3.1. Outils de développement utilisés	92
3.2. Présentation de l'application (Démonstration)	93
4. Conclusion	102
 Chapitre 5 : Expérimentation des méthodes de classification	 103
1. Introduction	103
2. Présentation	104
3. Corpus de documents XML	105
4. Expérimentation des complexités théoriques	106
5. Evaluation des méthodes	109
5.1. Cas1 : Test des Méthodes de Selkow, Dalamagas et Chawathe	110
5.2. Cas2 : Test de la méthode de Levenshtein	111
5.3. Cas3 : Test de la méthode de Aitelhadj	112
6. Conclusion	116
 Conclusion générale	 117
Les références	118

Résumé

Dans ce mémoire, cinq méthodes (Méthodes de : Selkow, Dalamagas, Chawathe, Levensthein et Aïtelhadj) de classifications des documents XML seront présentées. Leur but est de regrouper les documents qui partagent les mêmes structures. Un système de clustering commun est développé, incorporant les cinq méthodes pour tester la qualité de leur classification et leur complexité temporelle afin de faire une étude comparative et déterminer les meilleurs algorithmes.

Mots clés: Classification, Clustering, Chawathe, Selkow, Dalamagas, Levenshtein, Aïtelhadj, Prim, XML, structure arborescente, Arbre, distance d'édition, distance structurelle, similarité structurelle.

M.Otmani Nassim Abdeldjallal. UMMTO

E-mail : otmani.nassim.lmd@gmail.com

Remerciements

Je tiens à remercier toute personne ayant contribué, de près ou de loin, à l'accomplissement de ce travail de recherche. En particulier mon promoteur M. Aïtelhadj, de m'avoir proposé ce sujet et de m'avoir encadré en m'apportant son aide précieuse durant toute la durée du travail. Je tiens aussi à remercier tous les membres de ma famille qui m'ont toujours soutenu et apporté un environnement d'études agréable. Je remercierai aussi les membres du jury d'avoir accepté de juger mon travail.

Merci. O.N.

Introduction générale

Le Web est en perpétuelle croissance, de plus en plus de documents s'échangent à travers les réseaux. Ces documents sont généralement représentés à l'aide du langage XML qui permet de structurer les données d'un document en une structure hiérarchique arborescente. Cette structure est un moyen d'ajouter aux données des significations afin de les rendre plus perceptibles et plus faciles à gérer. Pour les besoins d'organisation et de meilleure gestion, ces documents sont souvent regroupés en fonction de leur ressemblance afin d'améliorer par exemple la recherche d'information. Or la plupart de ces documents sont issus de sources hétérogènes dont on ne connaît pas leur DTD source, c'est-à-dire qui ne sont pas conformes à une seule DTD spécifique ce qui rend la tâche du regroupement difficile sans connaître les DTD sources. Il y a donc une nécessité de comparer ces documents entre eux en utilisant des méthodes particulières qui permettent de déterminer la ressemblance entre les documents, sans connaître leurs DTDs, afin de les affecter par la suite à des classes spécifiques, c'est ce qu'on appelle la classification.

Actuellement, plusieurs méthodes existent et elles peuvent être classées en 3 groupes, le premier calcule la similarité en utilisant seulement les données contenues dans le document, sans prendre en considération la structure. Le deuxième groupe, qui est notre cas d'étude, classe les données en utilisant seulement la structure arborescente du document, sans prendre en considération les données qu'il contient. Le troisième groupe quant à lui utilise les données et la structure arborescente pour mieux classer les documents.

Le premier groupe contient des méthodes classiques qui existaient bien avant l'apparition du langage XML en 1996. Leur principe est généralement basé sur l'indexation, c.-à-d. calculer la répétition des mots dans un document et déduire la ressemblance entre les documents en comparant les mots existants en fonction de leur occurrence dans les documents.

Avec l'apparition du XML, ces méthodes sont devenues inadaptées à ce genre de documents, et il était nécessaire de concevoir de nouvelles méthodes appropriées. Ces nouvelles méthodes utilisent en plus des données contenues dans le document, la structure de celui-ci. Le deuxième groupe contient ce genre de méthodes. Une partie des méthodes du deuxième groupe existait déjà avant l'apparition de l'XML, vu qu'elles opèrent sur des structures arborescentes similaires à celle des documents XML. La plupart d'entre eux utilisent la technique du calcul de la distance d'édition, une valeur qui sera utilisée pour déterminer la similarité structurelle. Or, réellement, les similarités structurelles ne sont pas suffisamment prises en considération avec ces méthodes lors de la comparaison des documents XML, alors d'autres méthodes ont été conçues, permettant de déterminer d'une manière plus efficace et précise la similarité structurelle de ces documents arborescents.

Le troisième groupe de méthodes vise à augmenter le degré de l'exactitude de la comparaison des documents en comparant non seulement la structure hiérarchique arborescente mais aussi les données contenues dans les documents. Ce sont généralement des méthodes hybrides fusionnant des méthodes du premier et du deuxième groupe. Mais cette exactitude a un coût qui pourrait ralentir la vitesse de la classification.

Le but de ce mémoire est de présenter une étude expérimentale de comparaison de quelques méthodes existantes qui permettent de classifier les structures arborescentes des documents XML. Cinq méthodes seront étudiées et évaluées ; quatre d'entre elles (Méthode de Selkow, Méthode de Dalamagas, Méthode de Chawathe et la méthode de Levenshtein) utilisent la notion de distance d'édition pour calculer la similarité structurelle, une valeur jugée insuffisante pour déterminer avec exactitude la similarité, parce que l'aspect ontologique des nœuds de la structure arborescente n'est pris en considération lors de la comparaison des documents. La cinquième méthode (Méthode de Aïtelhadj) palie à ce problème et permet de déterminer efficacement des valeurs de similarité plus exactes. En revanche cette exactitude a un coût, qui se reflète sur la durée d'exécution du processus de la classification, qui est légèrement supérieure aux méthodes précédentes.

Nous allons faire une étude comparative afin de déterminer la qualité de classification offerte par chaque méthode en mesurant la complexité temporelle et l'exactitude de la classification et essayer de choisir les meilleures méthodes qui présentent le meilleur compromis entre durée d'exécution et qualité de classification.

Ce mémoire est organisé de la manière suivante : Outre le volet introductif, le chapitre 1 est consacré pour donner un aperçu général sur le langage XML en particulier : les documents XML. Le chapitre 2 est dédié pour donner une idée sur la classification des documents. Dans le chapitre 3, nous décrivons les cinq méthodes que nous avons choisi de tester, le chapitre 4 est réservé pour la conception et l'implémentation de notre logiciel de classification. Le chapitre 5 présentera une évaluation expérimentale des cinq méthodes. Enfin, nous terminons par une conclusion générale.

Chapitre 1

XML

7. Introduction

XML (*Extensible Markup Language*, « langage de balisage extensible ») est un langage informatique de balisage *générique* qui dérive du SGML. Sa syntaxe est dite *extensible* car elle permet de définir différents espaces de noms, c'est-à-dire des langages avec chacun leur vocabulaire et leur grammaire, comme XHTML, XSLT, MathML, VoiceXML... Cette syntaxe est reconnaissable par son usage des *chevrons* (< >) encadrant les *balises* [¹]. En plus, XML n'est pas seulement un langage de balisage, mais plutôt, une boîte à outils qui sert à créer, structurer et utiliser des langages de balisage. Il est donc vu comme une famille de langages partageant des caractéristiques communes, et dédiés à une multitude d'usages divers. Les facilités d'écriture de ce format, les possibilités de traitement des données qu'il offre, sa souplesse d'utilisation en font de l'XML un langage extrêmement bien adapté aux échanges de données entre applications aussi bien qu'à leur simple stockage.

Actuellement le XML est devenu un standard largement utilisé pour le stockage, création et l'échange des données. Il est utilisé carrément par toutes les applications actuelles, surtout sur le Web, ce qui a rendu la quantité des documents XML existante très importante. Le grand nombre de documents XML nécessite une organisation pour rendre l'accès et la recherche des informations dans ces entrepôts de fichiers plus pertinente et efficace. La classification des documents XML est une manière d'organiser ces documents en les mettant dans des ensembles en commun selon un facteur de ressemblance des documents.

Avant d'aborder les méthodes existantes pour classer ou organiser ces documents XML, nous allons présenter dans ce chapitre, la technologie XML en abordant principalement la structure des documents XML et les règles d'écriture des documents XML bien formés. Et nous présenterons aussi les outils offerts par cette technologie pour écrire et vérifier la validité des documents XML.

¹ http://fr.wikipedia.org/wiki/Extensible_Markup_Language

8. Présentation de la technologie XML

8.1. Besoin de structuration

Voici un exemple d'un fichier texte contenant des données :

```
Hacking : the art of exploitation  
Jon Erickson
```

Le XML est né d'un besoin universel ; savoir faire cohabiter dans un même document de l'information et de la signification. Ainsi les données deviennent perceptibles par l'humain et la recherche d'information devient plus efficace. Voici l'exemple précédent comme document XML:

```
<bilio>  
  <livre>  
    <titre>Hacking : the art of exploitation</titre>  
    <auteur> Jon Erickson </auteur>  
  </livre>  
</biblio>
```

8.2. Origine et buts de l'XML

A la fin des années 1970 et au début des années 1980, Charles Goldfarb de la société IBM a été l'origine de la création d'un langage de description de données, fondé sur des balises. Ce langage nommé SGML (Standard Generalized Markup Language) fut publié comme norme ISO 8879 :1986 en 1986. Ce langage permettait de décrire la structure d'un document indépendamment de sa visualisation, ou plus généralement de son interprétation par les applications tierces. SGML est un puissant langage de description, mais il est complexe. Il a fallu concevoir un langage qui offre les avantages de SGML et qui soit moins complexe. Ce fut chose faite avec la publication en février 1998 de la première version de la recommandation XML 1.0.

Le XML a été développé par un groupe nommé « XML Working Groupe » formé sous l'égide de la World Wide Web Consortium (W3C) en 1996 [2]. Les objectifs visés lors de la conception du XML sont :

1. XML doit être directement utilisable sur l'Internet.
2. XML doit supporter une grande variété d'applications.
3. XML doit être compatible avec SGML.
4. Il doit être facile d'écrire des programmes qui traitent des documents XML.
5. Le nombre de caractéristiques optionnelles dans XML doit être maintenu au minimum absolu, idéalement à zéro.
6. Les documents XML doivent être lisibles par l'humain et avec une perception claire.
7. La conception XML devrait se faire aisément et rapidement.
8. La conception XML doit être formelle et concise avec une syntaxe simple et sans ambiguïté.
9. Les documents XML doivent être faciles à créer.
10. La concision dans le balisage XML est d'une importance minime.

² <http://www.w3.org/TR/2008/REC-xml-20081126/>

8.3. Définition de l'XML

Extensible Markup Language, abrégé XML, décrit une classe d'objets de données nommée un document XML et décrit aussi, partiellement, le comportement du logiciel qui le traite^[3]. Le XML est un sous ensemble de SGML (Standard Generalized Markup Language) ^[4]. XML est une application ou une restriction de la norme SGML [ISO 8879]^[5]. Ceci dit qu'un document XML est conforme aux documents SGML. Un document XML est un fichier dont les données sont structurées. Son but est de permettre à des documents SGML génériques d'être servis, reçus et traités sur le Web.

9. Document XML

9.1. Présentation d'un document XML

Un document texte est un document XML s'il est **bien formé**, c.-à-d. il correspond à la syntaxe définie par la spécification XML. En plus, le document XML est **valide** s'il vérifie certaines contraintes avancées telle qu'une grammaire DTD ou un Schema XML. Voici un exemple d'un document XML, On utilise souvent le terme « Collection XML » pour désigner un ensemble de documents XML.

```
< ?xml version = "1.0" encoding = "UTF-8" ?>
<!--Element Racine-->
<biblio>
  <!--Premier enfant-->
  <livre>
    <!--Premier élément enfant de livre-->
    <titre>Les Misérables</titre>
    <auteur>Victor Hugo</auteur>
    <nb_tomes>3</nb_tomes>
    <localisation travée="1" armoire="4" étagère="2"/>
  </livre>
  <!--Deuxième enfant-->
  <livre>
    <titre>L'Assommoir</titre>
    <auteur>Emile Zola</auteur>
    <localisation travée="4" armoire="2" étagère="5"/>
  </livre>
</biblio>
```

FIGURE 1.1 : Exemple d'un document XML.

³ <http://www.w3.org/TR/2008/REC-xml-20081126/>

⁴ est un langage de description à balises, de norme ISO (ISO 8879:1986).

⁵ ISO (International Organization for Standardization). *ISO 8879:1986(E). Information processing — Text and Office Systems — Standard Generalized Markup Language (SGML)*. First edition — 1986-10-15. [Geneva]: International Organization for Standardization, 1986.

9.2. Documents XML bien formés

Un fichier texte est un document XML bien formé si :

1. Le fichier texte correspond à la grammaire XML⁶.
2. Il vérifie toutes les contraintes imposées par la spécification.
3. Chaque entité analysée qui est référencée directement ou indirectement par le document est bien formée aussi.

9.3. Spécification d'un document XML

Un document XML est un fichier qui obéit à un format et à une certaines règles de structure que nous allons détailler.

9.3.1. L'arborescence - les nœuds

Un fichier XML peut être considéré comme une base de données hiérarchique. Un document XML possède une racine, de laquelle dépend toute une arborescence d'éléments. Cette arborescence est constituée de nœuds. L'arborescence est la structure hiérarchique des nœuds. La plupart des nœuds d'un document ont un nœud parent, parfois des nœuds frères parfois aussi des nœuds enfants.

La FIGURE 1.2 donne un exemple d'arborescence de nœuds.

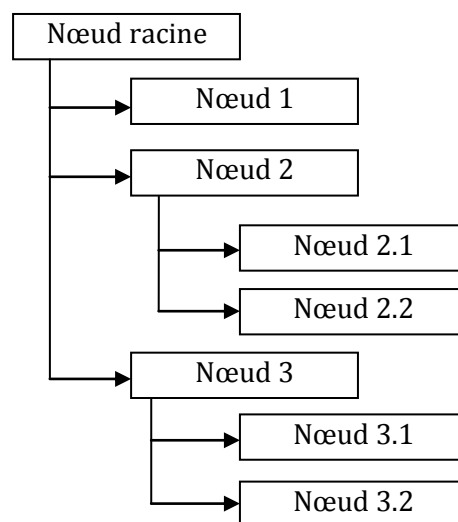


FIGURE 1.2 : Arborescence des nœuds dans un document XML.

⁶ La grammaire complète est décrite sur cette page web : <http://www.w3.org/TR/2008/REC-xml-20081126/>

Pour plus de clarté, nous allons présenter la structure arborescente de l'exemple précédent (FIGURE1.1).

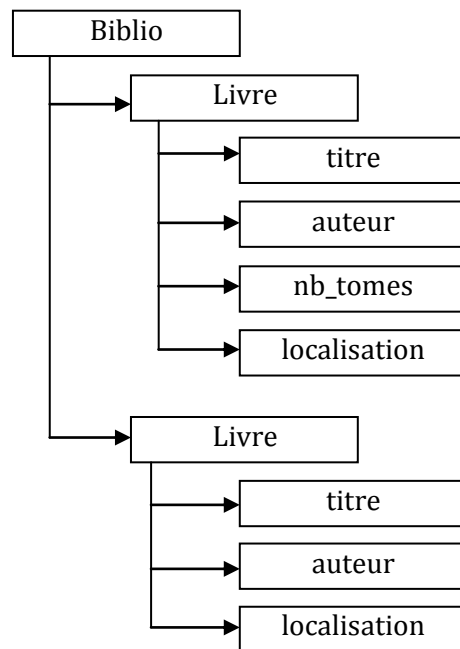


FIGURE 1.3 : représentation arborescente de l'exemple précédent

9.3.2. Les nœuds de type élément

Ce sont les éléments qui composent la structure d'un document XML. Les éléments sont les « étiquettes » associées aux données afin de les caractérisées. Dans l'exemple précédent, les éléments portent les noms *biblio*, *livre*, *titre*, *auteur*, *nb_tomes*, *localisation*. Les éléments possèdent un certain nombre de caractéristiques :

1. Un éléments s'ouvre et se ferme par une balise. Par exemple, <titre> est une balise ouvrante et </titre> est une balise fermante.
2. Certains éléments ne contiennent pas d'éléments fils. La balise fermante est ainsi intégrée à la balise ouvrante qui est alors auto fermante. C'est le cas par exemple de l'élément *localisation*. Notez que la balise qui signale le début de l'élément se termine avec un caractère /. On parle alors de l'élément vide.
3. Dans la balise ouvrante de certains éléments, notez la présence de chaîne de caractères précédés de signe = et d'un mot-clé. Ce sont des attributs.

9.3.3. Les nœuds de type attribut

Dans l'exemple précédent, les balises ouvrantes des éléments *localisation* et *livre* possèdent des attributs. Ces attributs modifient le sens par défaut de la balise, ou bien le précisent. Il est ici implicite que les livres sont écrits en français ; mais un des ouvrages est en Anglais, et il a semblé judicieux de l'indiquer en utilisant un attribut.

Un attribut commence par une chaîne de caractères (c'est le nom de l'attribut) ; il possède une valeur indiquée entre guillemets. Il s'écrit ainsi sous forme générale *nomAttribut="valeur de l'attribut"*. Un élément donné peut posséder plusieurs attributs, mais un même attribut ne peut être présent qu'une seule fois par élément. L'ordre des attributs n'a pas d'importance au sein d'un élément. L'attribut n'est pas repris dans la balise fermante.

9.3.4. Le prologue – Les instructions de traitement

L'exemple précédent débute par la ligne :

```
<?xml version="1.0" encoding="UTF-8" ?>
```

Cette ligne est le prologue du document XML et est facultative. Elle se place toujours en première ligne du fichier. Elle précise le numéro de la version du format XML, ainsi que l'encodage des caractères. Il s'agit d'une information qui est transmise à l'application chargée du traitement et qui lui indique dans quel type de codage le fichier a été enregistré. L'encodage par défaut est UTF-8.

En fait, le prologue est ce que l'on appelle une « instruction de traitement » particulière. Un autre cas très utilisé est l'appel à une feuille de style CSS ou XSL. Une instruction de traitement est destinée à l'application qui devra traiter le document XML. Sa syntaxe est la suivante :

```
<?nomApplication (...) ?>
```

9.3.5. Les nœuds commentaires

Les nœuds de commentaires sont conformes à la norme SGML. Un nœud commentaire commence par la chaîne de caractères `<!--` et se termine par les caractères `-->`

Un nœud commentaire sert à ajouter des commentaires au document XML. Ces commentaires n'ont pas vocation à être interprétés par l'application chargée du traitement du document ; ils peuvent par conséquent servir à insérer une documentation succincte dans le corps du document XML. Dans l'exemple précédent, c'est le cas du commentaire.

```
<!--Elément racine-->
```

9.3.6. Les entités

Une entité est une chaîne de caractères commençant par & et se termine par un point-virgule. Elle est déclarée dans une DTD (Document type Definition). Les entités sont des composants un peu particuliers des documents XML. Elles s'apparentent à des macros et, lors du traitement du document XML par son application, elles sont remplacées par la chaîne de caractères qu'elles représentent. Comme elles doivent être définies dans une DTD et que l'exemple précédent ne dépend pas d'une DTD, il n'en contient pas.

Exemple : L'entité &ADN ; est remplacée par la chaîne de caractères « acide désoxyri-bonucléique » lorsqu'elle est rencontrée, pour peu qu'elle soit définie ainsi dans la DTD dont dépend le document.

Entités prédéfinies

Certains caractères ont un sens particulier en XML. C'est la raison pour laquelle des entités ont été prédéfinies afin de pouvoir les utiliser. On appelle également ces entités des « séquences d'échappement » Il n'est pas nécessaire de déclarer ces entités dans une DTD.

Caractère	Entité
&	& ;
<	< ;
>	> ;
"	" ;
'	&apos ;

FIGURE1.4 : Tableau des entités prédéfinies en XML.

9.3.7. Les sections PCDATA et CDATA

Une section CDATA est une partie du document XML pouvant contenir toute sorte de chaîne de caractères. Une section CDATA permet de définir un bloc de caractères ne devant pas être analysés par le processeur XML. Il est ainsi possible entre autres de garder dans un bloc de texte un exemple de code à afficher tel quel. Il n'est alors pas nécessaire de recourir à des entités pour afficher les caractères réservés de XML. Par exemple, pour afficher les caractères < et > , on peut utiliser une section CDATA, comme dans l'exemple suivant :

```
< ![CDATA[Une balise commence par un < et se termine par un > .]]>
```

A contrario, une section PCDATA (Parsed Character Data) contient du texte destiné à être analysé par le processeur : par exemple, une entité qui s'y trouve sera traduite, et tout caractère < sera interprété comme indiquant le début d'une balise.

9.3.8. Encodage des documents

Le prologue du fichier XML donne des informations sur l'encodage de celui-ci. L'encodage désigne la représentation des caractères sur le disque dur. En raison du grand nombre d'alphabets et de systèmes idéographiques en usage dans le monde, les encodages initialement développés dans le pays occidentaux et fondés sur l'alphabet latin se sont révélés insuffisants. Il a donc fallu développer d'autres encodages donnant accès à des glyphes différents.

Lors de la création d'un nouveau fichier XML, il est important de se passer à priori la question de savoir quels seront les caractères susceptible d'être utilisés, et de choisir en conséquence l'encodage. Il ne suffit pas d'écrire dans le prologue `encoding="ISO-8859-1"` pour que ce soit le codage effectif du fichier ; encore faut-il que le fichier soit sauvegardé sous ce format. Pour peu que cette précaution soit prise, il devient possible de produire des documents XML dont les éléments et attributs contiennent des accents, ou bien écrits dans des alphabets non latins.

9.3.9. Règles d'écriture : document bien formé

Nous n'avons pas encore abordé les contraintes qui s'imposent à la forme d'un document XML. Ces contraintes sont assez simples, mais obligatoires :

1. XML est sensible à la case : les majuscules et les minuscules sont différenciées et Nom désigne un autre élément que nom.
2. Un nom d'élément ne peut pas commencer par un chiffre.
3. Si le nom ne comprend qu'un seul caractère, ce doit être une lettre.
4. Si le nom contient au moins deux caractères, le premier peut être un tiret « - » ou un tiret bas « _ ». le nom peut ensuite être composé de lettres, chiffres, tiret « - », tiret bas « _ » ou deux-points « : ».
5. Tous les éléments doivent être fermés : les éléments non vides par une balise fermante, les éléments vides par une balise auto fermante. Les éléments se ferment dans l'ordre inverse de leur ouverture (le dernier ouvert est le premier à fermer).
6. Les valeurs des attributs doivent être entre guillemets.
7. Il y a exactement un seul élément appelé « Racine », dont aucune partie de celui-ci n'apparaît dans le contenu de n'importe quel autre élément.
8. Pour tous les autres éléments, si la balise de début est dans le contenu d'un autre élément, la balise de fin est dans le contenu du même élément. Autrement dit, les éléments délimités par des balises de début et des balises de fin s'imbriquent correctement les uns dans les autres. Par conséquent, pour chaque élément C du document, autre que la racine, il y a un autre élément P dans le document tel que C est dans le contenu de P, mais il n'est pas dans le contenu de n'importe quel autre élément qui est dans le contenu de P. P se nomme comme parent de C et C comme fils de P.

Un document XML qui respecte ces règles d'écriture est dit bien formé.

9.4. Documents XML valides

Les documents XML sont souvent adaptés à un contexte particulier ; c.-à-d. des documents XML auront un format bien défini pour un cas précis. XML propose des outils qui permettent de définir leur structure conformément à notre contexte. Ces outils sont : DTD, Schemas XML et Relax NG. Ils permettent de définir une syntaxe du document XML. Lors de la réalisation d'un document XML, la liberté donnée au rédacteur est très grande car le seul impératif est que le document soit bien formé. Peu de règles sont définies sur le nom des balises ou l'existence d'attribut. Cette liberté est problématique lors de l'utilisation de documents XML dans le cadre d'échanges d'informations entre applications. Si des balises attendues par le système d'information sont manquantes, un dysfonctionnement pourra avoir lieu. Les DTD, SchemaXML ou Relax NG en Français permettent de résoudre ce type de problème.

Un document XML est valide s'il correspond à une DTD, à un Schema XML ou à et Relax NG.

9.4.1. Les DTDs

Pour s'assurer qu'un document XML n'a pas d'erreurs syntaxiques. On utilise une grammaire qui définira la structure du document. Cette grammaire s'appelle DTD (Document Type Definition ou Définition de type de document).

Exemple : Voici une grammaire DTD (en gras) incluse dans un document XML

```
<?xml version="1.0" encoding="UTF-8"?>

<!DOCTYPE livre [
<!ELEMENT livre (titre? , auteur+ , section+) >
<!ELEMENT titre (#PCDATA) >
<!ELEMENT auteur (nom , prenom+) >
<!ELEMENT nom (#PCDATA)>
<!ELEMENT prenom (#PCDATA)>
<!ELEMENT section (titre? , chapitre+) >
<!ELEMENT chapitre (titre? , paragraphe+) >
<!ELEMENT paragraphe (#PCDATA) >
]>

<livre>
  <titre>
    Titre du livre
  </titre>

  <auteur>
    <nom>
      nomAuteur1
    </nom>
    <prenom>
      prenomAuteur1
    </prenom>
```

```

    <prenom>
      prenomAuteur1
    </prenom>
  </auteur>

  <section>
    <titre>
      section1
    </titre>
    <chapitre>
      <titre>
        chapitre1
      </titre>
      <paragraphe>
        paragraphe
      </paragraphe>
    </chapitre>
  </section>

  <section>
    <titre>
      section1
    </titre>
    <chapitre>
      <titre>
        chapitre1
      </titre>
      <paragraphe>
        paragraphe
      </paragraphe>
    </chapitre>
  </section>
</livre>

```

FIGURE1.5 : Exemple d'un document XML avec une DTD.

Lors de l'échange d'un document XML, les applications analysent intégralement le document reçu pour en vérifier sa validité. De nombreux programmeurs codent alors des fonctions spécifiques dans leurs applications pour mener à bien cette réalisation. Dans cette partie, nous avons vu que les DTD permettent de faire ce travail. Ainsi, il est possible de vérifier la cohérence structurelle d'un document XML sans développer de code spécifique.

9.4.2. XML Schema

Les DTD sont un moyen simple de définir les éléments et attributs autorisés lors de la rédaction d'un document XML. Néanmoins, leur usage est parfois frustrant car ce langage ne permet pas d'exploiter à sa pleine mesure les possibilités de XML. Le W3C a donc développé un autre langage de définition de contenu, le XML schema. Un fichier écrit selon ce format est lui-même un document XML valide.

Le format XML schéma offre plus de possibilités que les DTD pour contraindre les éléments et attributs dans un format XML donné. Ces possibilités supplémentaires sont présentes au prix d'un format un peu plus difficile à analyser à l'œil, car fortement modularisé.

XML Schema permet d'affecter un type aux données d'un document XML, il facilite l'appel à des ensembles d'éléments et d'attributs définis dans plusieurs fichiers différents et parfois épars sur le réseau ; le format rend possible la mise au point d'une documentation détaillée incluse « naturellement » dans le corps du fichier, à fil des déclarations ; il offre une grande souplesse dans la définition des combinaisons d'éléments possibles. Enfin, en tant que document XML lui-même, un fichier de schema XML est un fichier qui peut être traité et manipulé automatiquement par les mêmes outils que ceux qui sont utilisés pour le traitement et la manipulation des fichiers qu'il définit : c'est une simplification du parc des outils, en même temps qu'un élargissement de l'éventail des applications possibles.

9.4.3. Relax NG

Il a été reproché au format XML schema une certaine complexité, rendant le saut d'apprentissage depuis les DTD ardu pour les développeurs qui ont l'habitude de ce langage. En parallèle des travaux du W3C, un consortium a travaillé sur un format alternatif, alliant une syntaxe allégée comme celles des DTD aux raffinements que permet XML Schema : Relax NG.

Relax NG est un format de définition de langage XML alternatif à XML Schema. Un de ses principaux avantages est l'existence d'une syntaxe compacte, plus facile à analyser pour un regard humain. Cette syntaxe reprend des caractéristiques des DTD, mais y ajoute des fonctionnalités qui étaient jusque-là réservées à XML Schema, comme la possibilité d'utiliser des espaces de noms ou celle de donner un type aux données. Ses avantages ont fait de lui une alternative à ce point fiable à XML schema qu'il a été utilisé, et continue à l'être pour la définition de nombreux langages XML.

9.5. Transformation d'un document

9.5.1. XPath

XPath est un langage d'expression qui sert à identifier des groupes de nœuds d'un document XML. Une expression XPath est comparable à la syntaxe d'un chemin d'accès qui décrit l'emplacement d'un fichier dans une arborescence d'un système de fichiers, d'où son nom de XPath, qui pourrait se traduire par : « chemin XML ».

XPath est un langage dont la version 1.0 est devenue une recommandation W3C le 16 novembre 1999 et la version 2.0 le 27 janvier 2007. Il représente l'organisation d'un document sous la forme d'une arborescence de nœuds. Les principaux nœuds qui peuvent composer une arborescence XPath sont les suivants :

- nœud racine ;
- nœud élément ;
- nœud texte ;
- nœud attribut ;
- nœud d'espace de noms ;
- nœud d'instruction de traitement ;
- nœuds commentaires.

9.5.2. XSLT

Un des problèmes majeurs de ces dernières années est l'interopérabilité des applications, c'est-à-dire la possibilité de faire communiquer entre elles des applications diverses et variées. La majorité d'entre elles utilise des formats de données propriétaires. Les échanges ne peuvent donc être réalisés que s'ils ont lieu dans un format commun. XSLT, abréviation de Extensible Stylesheet Language Transformations, est un langage qui apporte une solution à ce problème, comme à celui de la dissociation du contenu d'un document de sa forme. Il permet en effet de transformer un document XML en tout autre document XML ou au format de type texte, comme HTML.

Voici un exemple d'un document XML incluant une feuille de style XSLT (FIGURE 1.6), elle déclarée dans le document XML (FIGURE 1.7) comme suite :

```
<?xml-stylesheet href="exemple2_5.xsl" type="text/xsl"?>
```

L'exemple qui va suivre présente la transformation d'un document XML en un document HTML.

```

<?xml version="1.0" encoding="ISO-8859-1" ?>
<?xml-stylesheet href="exemple.xsl" type="text/xsl">
<liste>
  <livre>
    <titre genre="jeu">Texas Hols'Em Poker online</titre>
    <auteur>Mark Stohan</auteur>
    <auteur>Robert Bluman</auteur>
    <parution>2006</parution>
  </livre>
  <livre>
    <titre genre="jeu">Sudoku Manga</titre>
    <auteur>Sudoku factory</auteur>
    <parution>2007</parution>
  </livre>
  <livre>
    <titre genre="photo">Manuel de la photo</titre>
    <auteur>Jackie Contiboeuf</auteur>
    <auteur>Alain Mocney</auteur>
    <parution>2006</parution>
  </livre>
</liste>

```

FIGURE1.6 : Exemple d'un document (exemple.xml)

```

<xml version="1.0" encoding="ISO-8859-1">
<xsl:stylesheet version="1.0"
xmlns :xsl="http://www.w3.org/1999/XSL/Transform">
<xsl :template match="/liste">
  <html>
    <head>
      <title>Première transformation XSTL</title>
    </head>
    <body>
      <table summary="">
        <xsl :for-each select="livre">
          <tr>
            <td><xsl :value-of select="titre"/></td>
            <td><xsl :value-of select="auteur"/></td>
          </tr>
        </xsl:for-each>
      </table>
    </body>
  </html>
</xsl:template>
</xsl:stylesheet>

```

FIGURE1.7 : Exemple d'une feuille de style (exemple.xsl)

Dans cet exemple, pour chaque élément *livre* rencontré dans l'élément *liste*, les mots *titre* et *auteurs* sont écrits dans une nouvelle ligne du tableau. Pour extraire le contenu texte des éléments *titre* et *auteur*, il existe l'élément `xsl:value-of`. L'instruction `<xsl:value-of select="titre"/>` a pour effet d'extraire le contenu de l'élément *titre*.

9.5.3. XSL-FO

La transformation d'un document XML dans un autre format prise en charge par le langage XSL. Mais ce langage, par sa conception, ne peut produire des fichiers qu'au format texte. Le W3C a donc développé un langage de transformation, nommé XSL-FO (pour XSL-Formatting Object), qui permet de produire un fichier au format imprimable comme le PDF : Portable Document Format.

XSL-FO est un ensemble d'instructions qui offre de multiples possibilités. Pour la réalisation de fichier imprimable, c'est vraisemblablement une des rares solutions qui permet une automatisation complète de la mise en page. C'est donc le complément indispensable à XSLT et XML pour offrir un outil complet d'interopérabilité entre applications. Il ne reste plus qu'à automatiser la vérification syntaxique et structurelle d'un document XML.

10. Rôle du document XML

Un document XML sert à communiquer des informations dans une entreprise. Des informations complémentées par des significations précises. Un document XML peut être analysé pour en extraire les informations nécessaires comme il peut être utilisé pour faciliter la recherche d'information.

11. Le document XML : orienté document ou données ?

Lorsque les données sont élaborées par des êtres humains, on dit que les fichiers XML produits sont orientés document. Lorsque les données sont construites automatiquement par des programmes, on dit que les fichiers XML sont orientés données.

- Un fichier XML orienté document peut être, par exemple, un livre, un article, un message...
- Un fichier XML orienté donnée est, par exemple, un sous-ensemble d'une base de données.

12. Conclusion

La simplicité du codage et l'utilité de son utilisation, ont fait aujourd'hui de l'XML un outil largement utilisé pour le stockage et l'échange des données. Le XML est un langage en perpétuelle expansion et amélioration et le nombre de documents XML augmentent, de plus en plus, et leur transaction via le réseau ne cesse pas d'accroître.

L'aspect semi-structuré des documents XML facilite l'organisation, l'interprétation et la recherche d'informations. L'organisation en une structure hiérarchie arborescente où les nœuds représentent des balises ayant une certaine sémantique facilitant l'interprétation du contenu du document, ainsi la recherche pourra s'effectuer sur le balisage afin d'assurer une meilleure pertinence et diminution du temps de réponse des requêtes.

Les collections XML sont de plus en plus grandes, et les applications qui les utilisent subiront une baisse de performance lorsqu'ils font appel à ces collections XML. Pour absorber cette baisse, et faciliter le traitement et l'accès aux collections XML, il est nécessaire d'organiser les documents XML et de proposer des outils qui permettent une meilleure gestion.

Des travaux ont été faits pour améliorer l'organisation des documents XML. Dans le prochain chapitre nous allons présenter la classification des documents XML, un moyen d'organiser les documents XML.

Chapitre2

Classification

1. Introduction

Apparue vers le début des années 1960 mais qui s'est largement développée durant ces 20 dernières années, la problématique de classification automatique est considérée comme une tâche ancienne de la RI (Recherche d'information). Les méthodes de classification (Taxonomies, Réseaux de Neurones, etc.) étaient déjà connues avant même leur mise en forme algorithmique.

L'objectif de la classification documentaire est d'attribuer à un document une ou plusieurs classes parmi un ensemble prédéfini. Cette problématique a par ailleurs été rencontrée récemment dans de nouvelles applications dans des domaines tels que le filtrage sur Internet, la veille technologique, etc.

2. Notion de classe

Traditionnellement la notion de classe dans la classification documentaire a souvent été synonyme de thème, on parle alors de classification thématique. La tâche revient à partitionner un ensemble de documents autour de thématiques que l'on désigne communément par classes. Dans ce cas, il s'agit par exemple, de savoir si un document concerne le thème botanique, le thème philosophie ou bien le thème physique. Aujourd'hui la tâche de classification a évolué en même temps que les besoins.

Cette évolution a engendré de nouvelles problématiques pour lesquelles les classes ne correspondent plus nécessairement à des thématiques. A ce titre, la tâche de filtrage (text filtering) proposée par la campagne TREC (Text Retrieval Conference) permet de bien comprendre la problématique de classification.

Un système de filtrage est défini par exemple comme étant un processus qui permet d'extraire à partir d'un flot d'informations (News, emails, actualités journalières), celles qui sont susceptibles d'intéresser un ou plusieurs utilisateurs ayant des besoins en information stables [7].

⁷ Tebri Hamid « *Formalisation et spécification d'un système de filtrage incrémental information* » Thèse de Doctorat, Université Paul Sabatier, Toulouse, (2004)

TREC-4 propose d'évaluer un ensemble de classifieurs binaires indiquant si un nouveau document qui se présente est accepté ou pas. On cite ci-après quelques exemples concrets de classification binaire [8].

- Un logiciel qui détecte si un message est un Spam ou non.
- Un système qui sélectionne des informations pertinentes dans des flux de dépêches les envoie à différents organismes concernés.
- Un système d'agent s'intéresse aux flux de type Newsgroup et avertit un utilisateur dès qu'une nouvelle peut potentiellement l'intéresser.

Dans de tels systèmes une classe pertinente est assimilée à un besoin en information qui pourrait intéresser un utilisateur particulier ou une entreprise. D'autres problématiques émergentes telles que la répartition automatique de courriers de clients à la personne compétente dans une entreprise ou bien le tri de courriers électroniques dans différentes boîtes aux lettres personnelles ne peuvent être qualifiées de problématiques thématiques. Si l'on désigne chaque classe par une étiquette associée à des documents.

Le principe du système de classification est simple. En effet, si l'on suppose que l'utilisateur connaît effectivement la sémantique véhiculée par une étiquette, ce n'est à priori pas le cas du système à proprement parler. Par conséquent, le système identifie et classe de manière automatique les documents en se basant uniquement sur leurs labels sans toutefois connaître quelle est leur signification. L'utilisateur par contre, lui qui est censé (re)connaître à priori la signification des étiquettes, sera en mesure de sélectionner directement les documents qui l'intéressent via ces étiquettes.

La FIGURE 2.1 présente un exemple où l'on imagine un système fictif de classification de courriers électroniques basé sur la notion de classes étiquetées. Nous remarquons que les labels des classes possèdent effectivement des sémantiques de différentes natures (messages d'un certain type, messages émanant d'une certaine source, thème, etc.) que l'utilisateur peut facilement interpréter. Ce système organise des emails dans des boîtes aux lettres qui correspondent chacune à une classe du problème de classification.

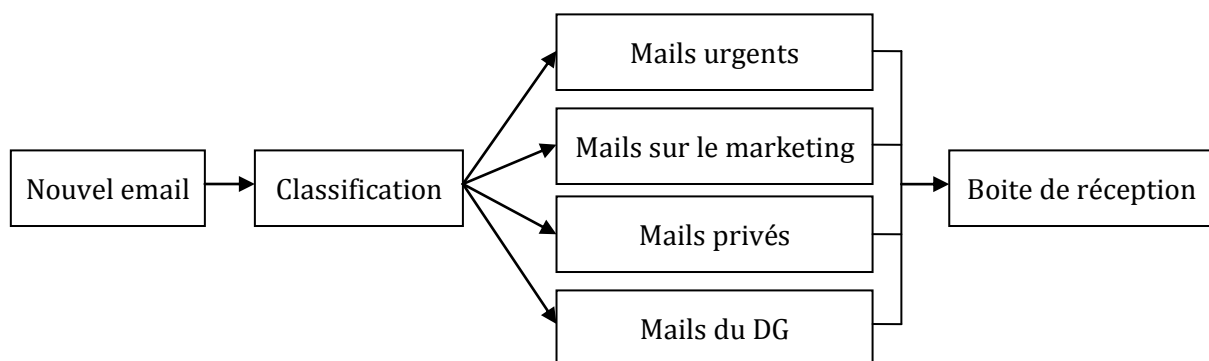


FIGURE 2.1 : Exemple de système de classification d'emails

⁸ Lewis D. *The TREC-4, Filtering Track*. In Proceedings of the 4th Text Retrieval Conference, TREC-4, National Institute of Standards and Technology NIST, Special Publication SP 500-263, Gaithersburg, MD Maryland, USA, 1–3 November, pages 165–180, (1996)

3. Définition formelle de la classification

Soient $C = \{c_1, c_2, \dots, c_n\}$ un ensemble de classes et $D = \{d_1, d_2, \dots, d_m\}$ un ensemble de documents à classer. On dira qu'un système de classification attribue automatiquement à chaque document un ensemble de classes : 0, 1 ou plusieurs. Lorsqu'un document d_i appartient à une classe quelconque c_j , on dira que le document d_i est pertinent pour la classe c_j . Plusieurs formalisations du problème de classification ont été proposées dans la littérature par différents auteurs. On cite à titre indicatif la formalisation de Sebastiani [9] et reprise par Yang [10] qui semble très pratique. Pour cette formalisation on définit deux fonctions [11] :

- Une *fonction de décision* qui consiste à attribuer à chaque document un ensemble de classes
- Une *fonction cible* indiquant la véritable appartenance d'un document à un ensemble de classes.

La fonction de décision permet d'estimer la fonction cible. Plus cette estimation est pertinente (correcte), meilleur est le système de classification. Ces deux fonctions associent à chaque couple $(d_i, c_j) \in D \times C$ une valeur booléenne indiquant si le document $d_i \in$ ou $\notin c_j$.

- La fonction de décision : $\varphi \rightarrow D \times C$ est telle que :

$$\varphi(d_i, c_j) = \begin{cases} \text{vrai si } d_i \text{ est attribué à la classe } c_j \\ \text{faux sinon} \end{cases}$$

- La fonction cible : $\phi \rightarrow D \times C$ est telle que :

$$\phi(d_i, c_j) = \begin{cases} \text{vrai si } d_i \text{ appartient à la classe } c_j \\ \text{faux sinon} \end{cases}$$

Habituellement avec les systèmes de classification basés sur les méthodes d'apprentissage, on évalue la fonction de décision en utilisant un corpus d'entraînement. Un corpus d'entraînement est une collection d'objets répartis en classes dont on connaît à priori les noms (étiquettes) et le nombre. Un modèle d'apprentissage est caractérisé par trois phases principales dont :

- L'apprentissage (phase d'induction) qui consiste en l'élaboration du modèle sur un corpus d'apprentissage dont on connaît la classification. Cela correspond à la fonction cible (training data).
- Le test et validation du modèle sur un nouvel échantillon test dont on connaît aussi la classification pour choisir le meilleur modèle (test data).
- l'application du classifieur (phase de déduction) aux documents à classer. Application de la fonction de décision par estimation de la fonction cible (application data).

⁹ Sebastiani Fabrizio. «*Machine Learning in Automated Text Categorization*». ACM Computing Surveys, 34(1) : 1–47, (2002)

¹⁰ Yang Y, Slattey S, Ghani R. «*A Study of Approaches to Hypertext Categorization*». Journal of Intelligent Information Systems IIIS, 18(2-3) : 219–241, (2002)

¹¹ Denoyer L. «*Apprentissage et inférence statistique dans les bases de documents structurés : Application aux corpus de documents textuels*». Thèse de Doctorat, Université Paris 6, (2004)

4. Le ranking

Contrairement à la classification définie dans la section précédente, qui émet une décision dure sur l'appartenance ou non d'un objet à une classe, le ranking (classement) ordonne (ou trie) ces objets par ordre de leur pertinence pour une classe donnée. On peut aussi dire que le ranking permet d'ordonner les classes par ordre de leur pertinence pour un objet donné. En d'autres termes, la tâche de ranking se base sur le calcul de la valeur d'une fonction de score définie par :

$$\text{Score} : O \times C \rightarrow [0,1]$$

La valeur de Score (O, C) appartenant à l'intervalle [0,1] représente la pertinence de l'objet O pour la classe C. Voici deux exemples typiques d'application du ranking [7] :

- Le filtrage de documents en fixant un seuil de tolérance par rapport au score de ranking.
- Le classement (ranking) de pages Web par rapport à une thématique fixée par l'utilisateur.

Ainsi, la fonction Score (O, C) nous informe sur le degré de pertinence de l'objet O pour la classe C. Plus ce degré est proche de 1, très pertinent est l'objet pour cette classe. Le calcul de la fonction de score permet alors de classer les objets par ordre de pertinence afin de savoir si un objet est plus pertinent ou moins pertinent qu'un autre pour une classe donnée.

Une bonne partie des systèmes probabilistes utilisent des algorithmes de classification automatique basée sur le calcul d'une fonction de score comme celle qu'on vient de définir ci dessus. Par exemple pour chaque document et chaque classe, on calcule la probabilité $P(d|c)$ que le document d soit pertinent pour la classe c. Le principe de ranking (classement) par probabilité [12] s'énonce ainsi :

Effectuer un ranking des documents en fonction de leur probabilité à posteriori est optimal dans le sens où les documents sont indépendants les uns des autres.

Plusieurs systèmes de ranking actuels suivent le principe de l'indépendance des documents. Par conséquent, pour les modèles probabilistes la fonction de score est :

$$\text{Score}(d, c) = P(c|d)$$

Les mêmes systèmes, peuvent cependant être aussi utilisés dans le cadre d'une classification binaire (dure). Il suffit à cet effet, de transformer la fonction de score en une fonction de décision comme suit :

Si λ représente un seuil (ou degré) de pertinence d'un document d vis-à-vis d'une classe donnée c, la fonction de décision recherchée $\varphi(d, c)$ s'exprime comme suit [13] :

$$\varphi(d, c) = \begin{cases} \text{vrai si score}(d, c) > \lambda \\ \text{faux sinon} \end{cases}$$

¹² Robertson S. *The probability ranking principle in IR*. Journal of documentation, 33(4) : 294–304, (1997)

¹³ Denoyer L. «Apprentissage et inférence statistique dans les bases de documents structurés : Application aux corpus de documents textuels». Thèse de Doctorat, Université Paris 6, (2004)

5. Le clustering

Le clustering (ou regroupement) est la tâche qui s'intéresse à trouver de manière automatique une organisation cohérente à un groupe d'objets. Elle correspond à la tâche de classification non-supervisée. La problématique du clustering a été largement étudiée en apprentissage, en reconnaissance des formes et en analyse de documents.

Une autre application importante de la tâche de clustering est en rapport avec la tâche de RD (Recherche Documentaire). Elle consiste à créer à priori sur un ensemble de documents un ensemble de clusters. Pour une requête donnée, le système de RD va tout d'abord chercher le cluster le plus pertinent pour la requête, puis trier les documents par ordre de pertinence dans ce cluster (ranking). Cette méthode permet de limiter la complexité de la recherche dans les grandes bases de données et de renvoyer un sous-ensemble de documents tous pertinents pour la requête puisque le choix du cluster est évalué à partir des caractéristiques communes des documents qui le composent. C'est une technique qui permet habituellement d'augmenter la précision d'un système de RD. En effet, d'une part, elle réduit, de manière plus que significative, la complexité temporelle de la recherche sur de grandes bases documentaires, d'autre part, les documents retournés par le moteur de recherche, en réponse à une requête, sont tous pertinents, car ils présentent logiquement des caractéristiques similaires, voire identiques pour cette requête (le cluster a été construit sur la base de caractéristiques communes aux documents qui s'y trouvent).

On rencontre deux grandes familles de méthodes de clustering :

- Les systèmes fondés sur un calcul de similarité qui rassemblent un sous-ensemble de documents similaires dans un même cluster. Cette similarité peut par exemple ne concerner que la structure du document, tout comme elle peut impliquer un mélange de modèles (structure + contenu).

- Les systèmes probabilistes fondés sur un mélange de modèles (structure+contenu) [14].

¹⁴ Denoyer L. «*Apprentissage et inférence statistique dans les bases de documents structurés : Application aux corpus de documents textuels*». Thèse de Doctorat, Université Paris 6, (2004)

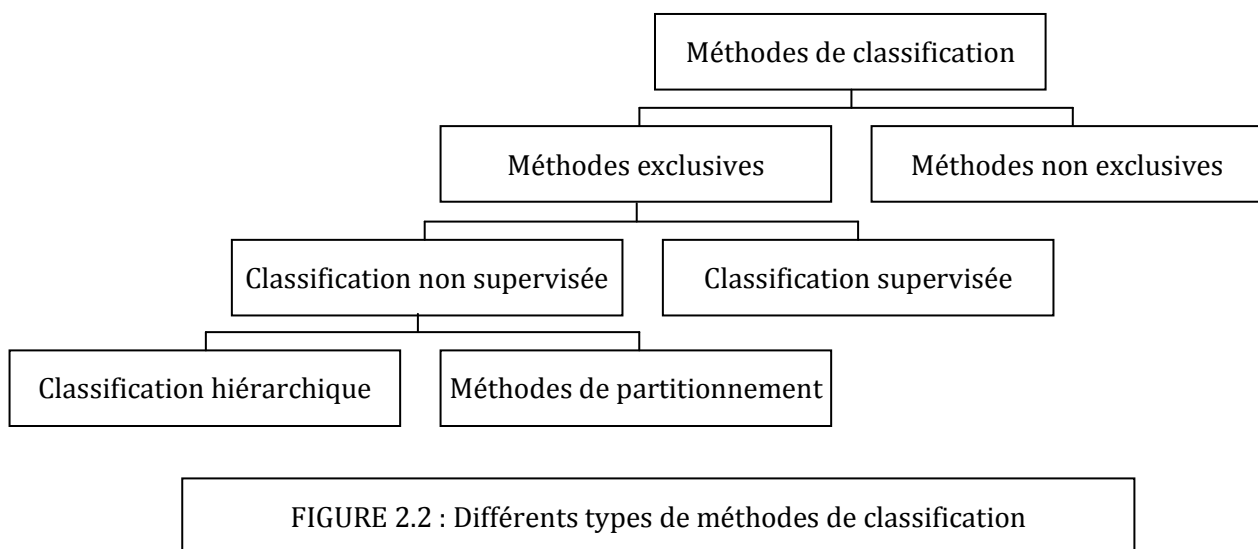
6. Contextes de classification

Le problème de la classification en général est de construire une procédure permettant d'associer une ou plusieurs classes à un objet. Les méthodes de classification peuvent être tout d'abord exclusives, c'est-à-dire qu'un objet n'appartient qu'à une et une seule classe, ou bien non exclusives. Dans ce dernier cas, un objet peut appartenir à plusieurs classes simultanément, éventuellement avec des degrés d'appartenance. On parle alors de recouvrement (overlapping clustering method). On distingue ensuite deux types de traitements : l'approche supervisée et l'approche non supervisée (ou clustering).

Dans la première, on connaît les classes possibles, et on dispose d'un ensemble d'objets déjà classés, servant d'ensemble d'apprentissage. Le problème est alors d'être en mesure d'associer à tout nouvel objet, sa classe la plus appropriée, en se servant des exemples déjà étiquetés.

Dans la seconde, par contre (classification non supervisée), les classes possibles ne sont pas connues à l'avance, et les exemples disponibles sont non étiquetés. Le but est donc de regrouper dans un même cluster (ou groupe), les objets considérés comme similaires, pour constituer les classes.

La FIGURE 2.2 présente les différents types de méthodes, regroupés sous forme d'une hiérarchie par Jain et Dubes dans [15].



¹⁵ Jain A.K, Dubes R.C. *Algorithms for Clustering Data*. Prentice-Hall advanced reference series : Computer Science, Prentice-Hall, Inc, Upper Saddle River, NJ, New Jersey, (1988)

Il existe différents contextes de classification :

Classification bi-classe :

La classification bi-classe est la problématique pour laquelle le système de classification répond à la question .Le document est-il pertinent ou non pour une classe donnée ? Par exemple, un message est-il un Spam ou non ? Cette problématique est le contexte de classification que l'on désigne d'ordinaire sur l'Internet par filtrage.

Classification multi-classe disjointe (exclusive) :

La classification multi-classe disjointe correspond au contexte où nous sommes en présence de plusieurs classes (>1) et pour lequel un document appartient exactement à une seule classe. La classification à classes disjointes est la problématique pour laquelle le système de classification répond à la question .A quelle classe (au singulier) appartient tel document ? [16].

Classification multi-classe (non-exclusive)

C'est le cas le plus général de la classification. Dans cette problématique le système de classification associe zéro (0), une (1) ou plusieurs (>1) classes à un document. Autrement dit, le système doit répondre à la

question .Quelles sont les classes (au pluriel) auxquelles appartient le document ? On est dans un cas de classification non-exclusive. Ce type de classification correspond par exemple à la problématique de classification du corpus Reuters [17]. Dans ce cas, il n'est pas difficile de constater qu'un problème de n classes se résout comme n problèmes biclasses [18]. Cependant, récemment les auteurs dans [92] ont proposé une méthode permettant de résoudre directement un problème multi-classe sans décomposition préalable en plusieurs problèmes bi-classes.

¹⁶ Denoyer L. «*Apprentissage et inférence statistique dans les bases de documents structurés : Application aux corpus de documents textuels*». Thèse de Doctorat, Université Paris 6, (2004)

¹⁷ Le corpus Reuters est un corpus .célèbre. de documents. Il a été conçu par l'agence de presse Reuters afin que des systèmes soient développés pour la classification automatique de dépêches de presse. Il correspond à une problématique de classification en plusieurs classes (un document appartient à une ou plusieurs classes)

¹⁸ Denoyer L. «*Apprentissage et inférence statistique dans les bases de documents structurés : Application aux corpus de documents textuels*». Thèse de Doctorat, Université Paris 6, (2004)

7. Paramètres d'évaluation des performances d'un système de classification

On ne mesure pas la performance absolue d'un système, d'une technique ou d'une approche. La mesure de la performance d'un système apparaît comme une notion subjective dépendant de l'utilisation finale de ce système. De même, l'évaluation de la performance d'un système de classification n'est pas toujours une tâche triviale. Pour évaluer par exemple la performance d'un clustering, d'un ranking ou d'une classification, même si dans le fond on cherche à juger le système sur la base des mêmes critères, les procédés diffèrent d'un système à l'autre. La performance de la classification dépend notamment de l'efficacité de la description de l'objet à classer. De plus, si l'on veut obtenir un système d'apprentissage, la procédure de classification doit permettre de classer efficacement tout nouvel objet (pouvoir prédictif).

Plusieurs critères pour évaluer la performance d'un SRI, ont été proposés :

- Facilité d'utilisation
- Coût d'accès /stockage
- Présentation des résultats
- Capacité d'un système à sélectionner des documents pertinents.

Le dernier critère est vraisemblablement le plus déterminant pour juger de la qualité d'un système de classification. Mais, ceci demande selon les contraintes du problème étudié, une certaine rigueur pour réduire dans la mesure du possible,

- le taux d'erreurs de classification,
- une certaine robustesse à s'accommoder de l'information superflue ou du manque d'information,
- la possibilité de manipuler diverses informations.

Habituellement on fait intervenir conjointement deux mesures dénommées rappel et précision pour mesurer les performances des SRI ou des systèmes de classification. Informellement le rappel est défini par le nombre de documents pertinents au regard du nombre de documents pertinents retournés par le système. Cela signifie que lorsque l'utilisateur interroge la base documentaire il souhaite voir retournés tous les documents susceptibles de répondre à son besoin en information. Si cette adéquation entre le questionnement de l'utilisateur et le nombre de documents présentés est importante alors le taux de rappel est élevé. A l'inverse, si le système possède de nombreux documents intéressants mais que ceux-ci n'apparaissent pas, on parle alors de silence. Le rappel s'oppose au silence.

La précision, quant à elle, représente le nombre de documents pertinents par rapport au nombre total de documents retournés par le système. Le principe est que lorsqu'un utilisateur interroge une base documentaire, il souhaite que les documents retournés en réponse à son interrogation correspondent à son attente. Tous les documents retournés superflus ou non pertinents, constituent du bruit. La précision s'oppose à ce bruit documentaire.

Si elle est élevée, cela signifie que peu de documents inutiles sont proposés par le système et que, par conséquent, ce dernier peut être considéré comme précis.

En bref, on définit dans un système de classification :

- Le rappel [19], comme étant la proportion de documents pertinents correctement classés.
- la précision [13] correspond à la capacité d'un système de classification à ne pas juger comme pertinent, un document qui ne l'est pas.

Formellement, si X_d est le nombre de documents bien classés, N_d le nombre de documents appartenant à la classe c , et N_c le nombre de documents attribués à la classe c , alors les mesures du rappel R et de la précision P sont exprimés respectivement comme suit:

$$R = \frac{X_d}{N_d}$$

$$P = \frac{X_d}{N_c}$$

La précision mesurée indépendamment du rappel et inversement est peu significative.

Il va donc falloir employer conjointement ces deux métriques afin de pouvoir évaluer significativement la performance d'un système. Il existe en ce sens des mesures alternatives combinant ces deux métriques. Une mesure qui combine la précision et le rappel est leur pondération, nommée F -mesure ou F_β introduite par [20] est donnée par la formule suivante :

$$F_\beta = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

On voit très bien d'après cette expression qu'il est possible, grâce au paramètre, d'accorder un poids plus ou moins élevé à la précision du classifieur. On fixe couramment la valeur de β à 1, ce qui donne :

$$F_1 = \frac{(2)PR}{P + R}$$

où la précision et le rappel sont pondérés de manière équitable.

¹⁹ Zhao Y, Karypis G. *Criterion Functions for Document Clustering : Experiments and Analysis*. Technical Report, 01–40, Department of Computer Science, University of Minnesota, Minneapolis, MN 55455, (2001)

²⁰ Van Rijsbergen J.C. *Information retrieval*. Butterworths, London 2ème édition, (1979)

8. Évaluation d'un classifieur bi-classe

Notre problème consiste à évaluer un classifieur binaire qui nous renseigne si un document d est pertinent ou non pour une classe donnée c . Nous nous plaçons alors dans le contexte simple de la problématique de filtrage où c représente une classe unique, celle des documents pertinents, et \bar{c} celle des documents non pertinents. Nous pouvons en fait, appliquer aussi ce principe en ce qui concerne un système de classification multiclasse; il suffit tout simplement de considérer qu'un classifieur multi-classe correspond à un ensemble de classifieurs bi-classes, c'est-à-dire qu'il faut juste combiner les mesures employées avec les différents classifieurs bi-classes.

Pour tester l'efficacité du classifieur, on se base sur un corpus étiqueté de documents, c'est-à-dire que l'on connaît à priori à quelle catégorie appartient un document, pertinent ou non pertinent. Pour ce faire, on construit une matrice de contingence comme celle de la FIGURE 2.3 où sont répertoriées les informations nécessaires au calcul des paramètres d'évaluation des performances du classifieur.

Cette matrice consiste en un tableau dont les cases sont des entiers définis dans un même système de mesure et d'unité homogène [21].

Classe trouvée \ Classe réelle	Pertinent $\phi(d, C) = \text{vrai}$	Non pertinent $\phi(d, C) = \text{faux}$
Pertinent $\phi(d, C) = \text{vrai}$	VP	FN
Non pertinent $\phi(d, C) = \text{faux}$	FP	VN

FIGURE 2.2 : Matrice de contingence [30]

En fait, ces 4 informations sont capitales et sont interprétées comme suit :

- V P : les vrais documents pertinents, c'est-à-dire ceux qui sont pertinents et considérés par le classifieur comme bien classés, on les appelle les « vrais positifs »

- V N : les vrais documents non pertinents, c'est-à-dire ceux qui sont non-pertinents

et considérés par le classifieur comme bien classés, on les appelle les « vrais négatifs »

- FP : les faux documents pertinents, c'est-à-dire ceux qui sont non pertinents et mal classés par le classifieur, on les nomme les « faux positifs »

- FN : les faux documents non pertinents, c'est-à-dire ceux qui sont pertinents et mal classés par le classifieur, nommés « faux négatifs ».

²¹ Denoyer L. « *Apprentissage et inférence statistique dans les bases de documents structurés : Application aux corpus de documents textuels* ». Thèse de Doctorat, Université Paris 6, (2004)

Pour le calcul du rappel R et de la précision P, il suffit d'inférer les probabilités $P(\phi|\varphi)$ et $P(\varphi|\phi)$, d'un document d par rapport à la classe c. A cet effet, on utilise la fonction cible ϕ et la fonction de décision φ . Donc, d'après la matrice de contingence et compte tenu de la définition de la mesure R, on a :

$$R = p((\varphi(d, c) = \text{vrai} | \phi(d, c)) = \text{vrai}) = \frac{VP}{VP + FN}$$

Comme le rappel étant la capacité d'un système à sélectionner tous les documents pertinents d'une collection, alors dans le cas où $FN = 0$, le rappel est = 1. Cependant, que cette estimation soit proche de 0 ou de 1, le rappel demeure insuffisant pour évaluer l'efficacité d'un système de classification. C'est pourquoi il va falloir également calculer la précision. Alors, compte tenu de la définition de la mesure de P et des informations recueillies sur la matrice de contingence, on a :

$$P = p((\phi(d, c) = \text{vrai} | \varphi(d, c)) = \text{vrai}) = \frac{VP}{VP + FP}$$

On voit bien que la précision correspond à la proportion de documents réellement pertinents correctement classés.

Les points suivants illustrent quatre situations « caricaturales » de mesures des performances (R, P et F1) de classifieurs testés sur un corpus de 500 documents [30]. Ce corpus est réparti à priori par la fonction cible ϕ sur deux classes c (300 documents pertinents) et \bar{c} (200 documents non pertinents).

– Si tous les documents sont considérés par le classifieur comme étant tous pertinents, c'est-à-dire $VP = 300, FP = 200, FN = 0$ et $VN = 0$. Par conséquent, $R = 100\%$; $P = 60\%$ et $F1 = 75\%$

– Si tous les documents sont considérés par le classifieur comme étant tous non pertinents, c'est-à-dire $VP = 0, FP = 0, FN = 300, VN = 200$. Par conséquent, $R = 0\%$, $P = 0\%$ et $F1 = 0\%$.

– Si le classifieur ne considère pertinents que les documents qui le sont réellement, il sera désigné par « classifieur parfait », alors $VP = 300, FP = 0, FN = 0, VN = 200$, et par conséquent, $R = 100\%$, $P = 100\%$ et $F1 = 100\%$.

– Si le classifieur prend tous les documents qui sont effectivement pertinents pour des non pertinents, et tous les documents non pertinents pour des pertinents, il sera qualifié de « classifieur le pire », alors $VP = 0, FP = 200, FN = 300, VN = 0$, et par conséquent, $R = 0\%$, $P = 0\%$ et $F1 = 0\%$.

Dans la configuration « classifieur le pire » ainsi que celle où tous les documents, quel que soit leur statut, sont non-pertinents pour le classifieur, la précision et le rappel sont nuls (0%). Par conséquent, il est inutile d'appliquer la formule pour F1. Mais, on considère que $F1 = 0$, lorsqu'on est dans ce cas de figure.

9. Évaluation de classifieurs multi-classes

Comme préconisé, l'évaluation d'un classifieur multi-classe, revient à se placer dans le contexte de plusieurs classifieurs bi-classes. La solution consiste donc à reconsidérer les mêmes métriques (R, P et F1), et à les utiliser à bon escient, suivant le contexte multiclasse. De nouvelles considérations peuvent alors entrer en jeu.

Pour cela, soit $C = (c_1, \dots, c_{|C|})$, où $|C|$ représente le nombre de classes, c'est-à-dire le nombre de classifieurs bi-classes. Sachant que t_j représente le nombre de documents de la classe c_j et que, R_j , P_j et $F1_j$ sont les mesures de la performance d'un classifieur j , on peut alors calculer les mesures globales pour le classifieur à $|C|$ classes connaissant ces mesures. Il suffit pour cela de calculer les moyennes de ces mesures. Cependant, en cas d'inégalité de tailles des classes, ces moyennes risquent d'être incohérentes et ne pas refléter la réalité du classifieur, notamment en ce qui concerne les classes à effectifs élevés. Pour contourner cet inconvénient et améliorer les paramètres d'évaluation, il convient d'introduire deux nouvelles mesures : la *micro_mean* et la *macro_mean* (respectivement la micro-moyenne et la macro-moyenne) [22].

- La *micro_mean* est une moyenne qui prend en compte la taille des classes.
- La *macro_mean* est une moyenne qui ignore totalement la taille des classes.

Les mesures globales (rappel, précision et F-mesure) de type *micro_mean* et *macro_mean* sont alors explicitées comme suit :

La <i>micro_mean</i>	La <i>macro_mean</i>
<p>Le rappel (R)</p> $R = \frac{\sum_{j=1}^{ C } t_j * R_j}{\sum_{j=1}^{ C } t_j}$ <p>La précision (P)</p> $P = \frac{\sum_{j=1}^{ C } t_j * P_j}{\sum_{j=1}^{ C } t_j}$ <p>La F-mesure (F1)</p> $F1 = \frac{\sum_{j=1}^{ C } t_j * F1_j}{\sum_{j=1}^{ C } t_j}$	<p>Le rappel (R)</p> $R = \frac{\sum_{j=1}^{ C } R_j}{ C }$ <p>La précision (P)</p> $P = \frac{\sum_{j=1}^{ C } P_j}{ C }$ <p>La F-mesure (F1)</p> $F1 = \frac{\sum_{j=1}^{ C } F1_j}{ C }$

Tandis que la *micro_mean* permet d'estimer la performance du classifieur en présence de classes à effectifs élevés, la *macro_mean* elle, par contre, permet d'estimer cette performance sur des classes à effectifs réduits.

²² Zhao Y, Karypis G. *Criterion Functions for Document Clustering : Experiments and Analysis*. Technical Report, 01-40, Department of Computer Science, University of Minnesota, Minneapolis, MN 55455, (2001)

10. Conclusion

La classification a conquis aujourd'hui le centre d'intérêt des chercheurs. Plusieurs travaux ont été faits, dans lesquels ils visent à inventer et améliorer des méthodes de classification de plus en plus performantes ; Toujours pour améliorer la classification et diminuer la complexité temporelle de ces méthodes. Avec le progrès rapide sur les matériels d'informatique, le deuxième avantage semble beaucoup moins important maintenant et la qualité et la justesse de la classification devient le but ultime.

Dans ce chapitre le principe de la classification a été présenté, ainsi que les méthodes utilisées pour évaluer la qualité de la classification. Dans le prochain chapitre, quelques méthodes de classification seront décrites en détails afin de les implémenter et les évaluer par la suite.

Chapitre3

Classification des documents XML

1. Introduction

XML est devenu un standard pour le développement des applications. Il s'occupe du stockage et de la récupération des données. Il y a plusieurs articles qui ont mis l'accent sur le problème du traitement, le stockage et la gestion des documents XML. Cependant, dans le domaine du traitement des documents XML, l'organisation des documents XML est devenue un sujet d'actualité.

La classification des documents XML représente l'une des plus implorantes tendances du domaine de la recherche relative à l'organisation des Docs XML. Le but du traitement est de regrouper des documents XML similaires en un seul ensemble, c'est ce qu'on appelle la classification. Le fait que le problème de la classification a été abordé dans de différents contextes à travers des années, ça a mené à développer plusieurs méthodes qui ont été revues et adaptées à ce contexte.

Des méthodes de classification des documents textuelles classiques existaient auparavant, elles sont basées principalement sur l'indexation. A ce stade, on peut dire que les documents XML diffèrent des données textuelles traditionnelles. Un document XML est composé d'une structure et d'un contenu. Ceci dit donc que les algorithmes de classification existants sont considérés inappropriés pour le problème de la classification XML. Ainsi, il y a un besoin de développer des nouveaux algorithmes de classification dédiés aux documents XML. Ces méthodes peuvent être catégorisées en deux groupes. Le premier groupe de méthodes est basé sur la structure et le contenu du document XML [23] tant dis que le second est exclusivement lié à la structure. C'est à ce second groupe que nous allons nous intéresser dans ce chapitre.

La classification des documents XML, selon leur structure, a plusieurs applications. Identification des documents XML avec des structures similaires peut être utile lorsqu'on se focalise sur la détection de

²³ Thierry Despeyroux, Yves Lechevallier, Brigitte Trousse, Anne-Marie Vercoustre, Expériences de classification d'une collection de documents XML de structure homogène 5ème Journées d'Extraction et de Gestion des Connaissances (EGC 2005), 2005.

la similarité structurelle des documents. Par conséquent, ceci aide à résoudre le problème de reconnaissance des différentes sources qui offre le même type d'informations, car les documents qui représentent la même information peuvent être différents, cependant la structure des documents peut être similaire.

Les documents XML peuvent être représentés de différentes manières. Ils peuvent être représentés sous forme de graphe, arbre, ensemble de chemins, ensemble d'étiquettes, série chronologiques, etc... La représentation des documents XML a un impact crucial sur la qualité et la performance du processus de la classification.

Le but de ce chapitre est de présenter cinq méthodes différentes qui peuvent être utilisées dans le processus de la classification des documents XML. Les cinq méthodes sont : la méthode de Selkow^[24], la méthode de Dalamagas^[25], la méthode de Chawathe^[26], la méthode de Levenshtein^[27] et la méthode de Aïtelhadj^[28]. Toutes ces méthodes ont un facteur en commun, elles utilisent toutes la même représentation des documents XML qui est le résumé d'arbre.

Le résumé d'arbre est arbre qui représente la structure hiérarchique du document XML dépourvu des informations qu'il contient. C'est un arbre enraciné, ordonné et étiqueté.

Pour comparer les documents XML entre eux, les quatre premières méthodes utilisent une mesure appelée la distance d'édition pour déterminer la valeur de distance structurelle entre les documents. En revanche, La cinquième méthode utilise une mesure qui détermine la valeur de similarité entre les documents XML. La distance et la similarité sont deux valeurs complémentaires.

La structure du chapitre peut être présentée comme suit : dans un premier temps, on présente le travail relevant du problème de la classification des Doc XML ainsi que des exemples d'utilisation. Puis on déterminera une formulation du problème. Après on présentera la représentation des documents XML qui a été adoptée par les cinq méthodes qui est le résumé d'arbre. Puis on présentera en détails les cinq algorithmes des méthodes choisies. Dans un dernier temps, on présentera un algorithme commun de classification des documents XML pour tester, évaluer et expérimenter les cinq méthodes. Et enfin on terminera par une conclusion.

²⁴ Theodore Dalamagas, Tao Cheng, Klaas-Jan Winkel, Timos Sellis, Clustering XML documents by structure, Inf. Syst., vol 31, 2006, 187--228

²⁵ Theodore Dalamagas, Tao Cheng, Klaas-Jan Winkel, Timos Sellis, Clustering XML documents by structure, Inf. Syst., vol 31, 2006, 187--228

²⁶ Sudarshan Chawathe, Comparing Hierarchical Data in External Memory, in Proceedings of the Twenty-fifth International Conference on Very Large Data Bases, 1999, 90—101.

²⁷ Binary codes capable of correcting deletions, insertions, and reversals, SOVIET PHYSICS-DOKLADY, Vol10, N°8 (1966)

²⁸ Ali Aïtelhadj, Mohamed Mezghiche et Fatiha Souam, Classification de Structures Arborescentes: Cas de Documents XML, CORIA 2009 - Conférence en Recherche d'Information et Applications.

2. Travaux faits

Au début des années 90, le problème de la classification de structures arborescentes des documents XML n'a pas eu une grande attention, donc il y avait peu de solutions proposées [29]. Actuellement, il y a plusieurs domaines de recherche en relation avec soit la représentation des documents XML soit avec le calcul des mesures de similarité soit avec leur classification.

Pour monter une proposition d'une méthode de classification de documents XML selon leur structure arborescente, il est nécessaire d'inclure chaque domaine de recherche. La forme du document peut varier et pourrait être présentée comme un arbre, un graphe, un ensemble de chemins, séries chronologiques, vecteurs et autres comme le SuperTree proposé dans [30]. La plupart des modèles existant pour la représentation des documents XML sont basés sur les arbres étiquetés [31][32][33], parce que c'est une illustration naturelle indiquant la structure hiérarchique des documents XML. Un document XML transformé en un arbre enraciné et étiqueté. Chaque élément du document représente un nœud dans l'arbre. Chaque nœud a son propre nom et est représenté par une étiquette. Un arc décrit une relation existante entre les nœuds. Le calcul des similarités entre des structures de documents XML varie par conséquent. Si la solution est basée sur l'arbre, les auteurs ont utilisés « la distance d'édition entre les arbres » comme mesure de similarité entre documents XML [34]. Si la représentation est basée sur un graphe, la similarité sera calculée à partir de l'ensemble des nœuds et des arcs [35] selon leur fréquence et leur organisation.

Dans [36] les auteurs ont calculé la similarité structurelle entre des sources DTD pour déterminer une DTD générale (XClust [37]) pour lui attribuer par la suite les documents XML qui ont une même structure. Dans [38] les auteurs ont calculé la similarité structurelle entre des Schémas XML pour déterminer des Schémas XML globaux et construire des clusters ayant les mêmes schémas XML. Dans [39] les auteurs ont proposé une solution basée sur les chemins. Ils ont utilisé un « patron séquentiel d'exploitation » (Sequential pattern mining) pour extraire les chemins fréquents du document XML et

²⁹ Anna Lesniewska, Clustering XML Documents by Structure, *Advances in Databases and Information Systems*, 2009, 238-246.

³⁰ Anna Lesniewska, Clustering XML Documents by Structure, *Advances in Databases and Information Systems*, 2009, 238-246.

³¹ S. M. Selkow, The tree-to-tree editing problem, *Information Processing Letters* 6 (1977) 184--186

³² Theodore Dalamagas, Tao Cheng, Klaas-Jan Winkel, Timos Sellis, Clustering XML documents by structure, *Inf. Syst.*, vol 31, 2006, 187--228

³³ Sudarshan Chawathe, Comparing Hierarchical Data in External Memory, in *Proceedings of the Twenty-fifth International Conference on Very Large Data Bases*, 1999, 90--101.

³⁴ Ali Aïtelhadj, Mohamed Mezghiche et Fatiha Souam, Classification de Structures Arborescentes: Cas de Documents XML, CORIA 2009 - Conférence en Recherche d'Information et Applications

³⁵ Charu C. Aggarwal, Haixun Wang, *Managing and Mining Graph Data*, Springer, 2010.

³⁶ Mong Li Lee, Liang Huai Yang, Wynne Hsu, Xia Yang, XClust: Clustering XML Schemas for Effective Integration, 2002

³⁷ XCLUST est une implémentation de la méthode classification par DTD.

³⁸ Alsayed Algergawy, Eike Schallehn, Gunter Saake, A Schema Matching-based Approach to XML Schema Clustering, *Proceedings of the 10th International Conference on Information Integration and Web-based Applications & Services*, 131--136, 2008

³⁹ Calin Garboni, Florent Masseglia, et Brigitte Trousse, Sequential Pattern Mining for Structure-Based XML Document Classification, *Workshop of the Initiative for the Evaluation of XML Retrieval*, 2005.

l'utilisé par la suite pour le clustering. La mesure de similarité peut être calculée en utilisant des normes (Euclidienne par exemple) telle qu'elle a été calculée dans [40].

3. Exemples d'utilisation

Pour illustré l'intérêt de la classification voyant ces exemples. Il ya plusieurs cas où la classification structurelle peut assister des applications dans leurs tâches :

3.1. Extraction automatique des DTDs

Beaucoup de documents XML sont construits à partir des sources de données telles que RDBMSs, flat file, etc, sans DTD. XTRACT et IBM AlphaWorks DDbE4 sont des systèmes qui extraient des DTDs à partir des documents XML. En identifiant des groupes de documents XML qui partagent une structure similaire pourrait être utile pour un tel système, où la collection de documents XML devrait d'abord être regroupée en ensembles qui contiennent des documents XML similaires, ainsi une DTD peut être assignée à chaque ensemble individuellement.

3.2. Regroupement structurel

Depuis que le langage XML encode et représente différents types de données hiérarchiques, la classification des documents XML structurelle peut être exploitée dans des applications qui nécessitent une gestion et un traitement de données hiérarchiques.

Quelques exemples :

- Bioinformatique: la découverte des patrons d'arbres macromoléculaires structurels encodés en documents XML.
- La détection et l'homogénéisation des structures de protéines encodées en documents XML (i.e. un ensemble de structure de protéines partageant une structure similaire)

3.3. Gestion des données spatiales

Les données spatiales sont souvent organisées en catalogues de modèles de données exprimés en un format XML hiérarchique. Par exemple, les espaces qui incluent les forets avec des laques et des firmes peuvent être représentés comme des structures arborescentes en utilisant des documents XML. La classification structurelle peut identifier les entités spatiales avec des structures similaires, exp : les entités avec les espaces incluant les forets et les laques. Plus d'exemple sur la représentation géographique des données sont présentés ici [41]

⁴⁰ Ali Aïtelhadj, Mohamed Mezghiche et Fatiha Souam, Classification de Structures Arborescentes: Cas de Documents XML, CORIA 2009 - Conférence en Recherche d'Information et Applications

⁴¹ R. Wilson, M. Cobb, F. McCreedy, R. Ladner, D. Olivier, T. Lovitt, K. Shaw, F. Petry, M. Abdelguer, Geographical data interchange using xml-enabled technology within the GIDB system, in: A. B. Chaudhri, A. Rashid, R. Zicari (Eds.), XML Data Management, 2003, Addison Wesley.

3.4. La mise à jour des documents

Les algorithmes de comparaison sont utilisés pour calculer et sauvegarder uniquement les différences entre l'ancienne et les nouvelles versions de documents qui ont été modifiés.

3.5. La détection des modifications

L'algorithme permet de comparer deux versions d'un fichier pour détecter où et comment ils se différent. Ces différences sont par la suite présentées en utilisant une interface graphique permettant à l'utilisateur de déterminer laquelle des deux variantes garder.

Exemple, pour des manuels didacticiels, plusieurs versions sont généralement établies, et rarement qu'un manuel soit stable. Il convient d'offrir aux lecteurs, possédant les anciens manuels, les nouveaux manuels. En leur évitant de tout relire, on applique l'algorithme pour qu'il détecte seulement les parties qui ont été modifiée et les mettre en évidence.

3.6. Optimisation des échanges

Les algorithmes de comparaison sont aussi utilisés pour réduire la quantité des données qui transitent via le réseau entre les applications miroirs ou les BDD dupliquées. Les grands serveurs Web et FTP ont souvent des dizaines de sites miroirs répartis dans le monde. Les modifications apportées sur un serveur maître doivent être propagées aux sites miroirs. Idéalement, la personne ou le programme qui effectue la modification devrait marquer exactement quelles données ont été mises à jour pour notifier les autres acteurs. Cependant, en pratique, l'anonymat et la nature désorganisée de tels sites, rendent la relation d'enregistrement des modifications entre eux carrément impossible. Dû à ces difficultés, des programmes de comparaison sont utilisés pour calculer et faire propager seulement les différences entre la version du serveur maître et les sites miroirs.

3.7. Recherche d'informations

La classification (clustering) des documents vise à mettre les documents similaires ensemble. En le faisant, on veut atteindre un des buts suivants:

- Le nombre de clusters, par rapport au nombre de documents, est beaucoup plus petit. Ainsi, on peut accélérer le processus de recherche.
- Si un document est pertinent à une requête, alors les documents similaires ont plus de chance à être pertinents aussi. Ainsi, la classification peut être aussi vue comme un moyen d'expansion.
- Finalement, les réponses du système peuvent être regroupées, plutôt qu'être mises dans une liste individuellement. L'avantage de cette présentation de résultats est que l'utilisateur peut avoir une idée globale des résultats que le système a trouvés assez rapidement.

3.8. Vérificateurs d'orthographe (faute humaine)

La comparaison peut être utilisée pour vérifier la véracité syntaxique du texte introduit par l'utilisateur et de proposer les corrections les plus proches du texte introduit en calculant la similarité entre les mots introduits et les mots du dictionnaire.

4. Formulation du problème

Le problème de la classification des documents XML selon leur structure peut être déclaré comme suit :

Soit un ensemble de documents XML (par exemple Documents Web) le but de la classification XML est de grouper des documents XML similaires en un seul ensemble. Chaque méthode de classification utilise une représentation spécifique des documents et une mesure de similarité adaptée à cette représentation.

Les deux, la représentation des documents XML et les mesures de similarité choisie pour la représentation, jouent un rôle crucial dans l'exploitation des documents XML. En effet, la représentation des documents XML pourrait varier, le document XML pourrait être présenté comme un arbre, un graphe, ...

La mesure de similarité telle qu'elle est appliquée dans l'algorithme de clustering doit être appropriée à la représentation des documents XML choisie. Dans le cas d'un chemin, la mesure de similarité est généralement basée sur la distance métrique (Euclidienne, Norme, etc).

Dans le paragraphe suivant, on va présenter la représentation (résumé d'arbre) utilisée par toutes les méthodes que nous allons tester. Après on présentera les méthodes qui permettent de calculer la mesure de similarité, par conséquent la représentation appliquée aux méthodes.

5. Modélisation des données semi structurées

Toutes les méthodes que nous avons testées utilisent la même représentation des documents XML. Les données XML sont semi structurées. Ils ont une structure hiérarchique qui auto-décrit les informations du document. L'extraction de cette hiérarchie constitue un résumé du document qui peut être représenté avec arbre enraciné, ordonné et étiqueté. Enraciné parce qu'un document XML a un élément unique qui s'appelle l'élément racine, depuis lequel tous les autres éléments se dérivent. Ordonné parce que chaque élément peut être identifié d'une manière unique par un numéro d'ordre. Étiqueté par ce que chaque nœud de l'arbre possède une étiquette qui est la valeur textuelle de la balise.

5.1. Modélisation des documents XML

Les modèles de représentation des données semi structurées sont généralement basés sur des arbres ou plus généralement sur des graphes. Ils sont des modèles simples et flexibles et permettent une vision claire sur le schéma et sur les données.

L'Object Exchange Model (OEM) : modèle d'échange d'objets est un graphe composé d'une collection d'objets. OEM a été introduit dans le projet TSIMMIS [42][43]. Chaque objet OEM a un identifiant (oid) et une valeur, simple (élémentaire) ou complexe. La valeur simple est un entier, un real, une chaîne de caractères, ou n'importe quel autre type donnée simple, tandis que la valeur complexe est un ensemble de oids, qui sont tous liés à un nœud parent en utilisant une étiquette textuelle. Les objets avec des valeurs simples (atomiques) sont appelés des objets atomiques (simples) et les objets avec des valeurs complexes sont appelés des objets complexes. La FIGURE 3.1 présente un exemple d'une base de données OEM. Dans cet exemple il y a quatre objets complexes (la racine, deux objets « SLT Camera » et un objet « Point & Shoot camera ») et sept objets atomiques (trois « brand » trois « price » et un « color »)

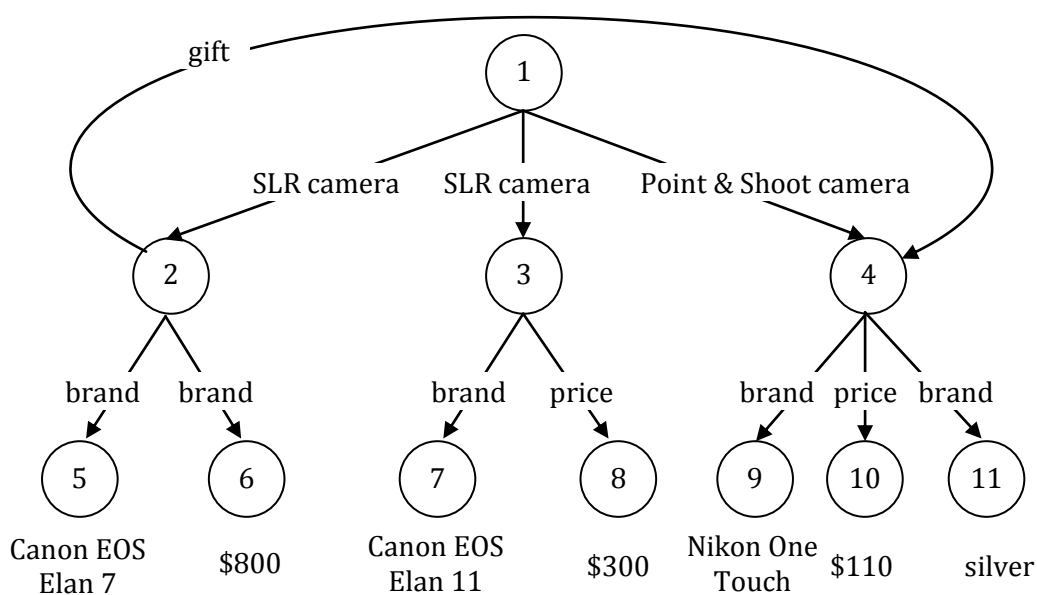


FIGURE 3.1 :
Exemple d'un
OEM
(Object
Exchange
Model)

⁴² Y. Papakonstantinou, H. Garcia-Molina, J. Widom, Object exchange across heterogenous information sources, in: Proceedings of IEEE International Conference on Data Engineering, 1995, pp. 251--260.

⁴³ H. Garcia-Molina, Y. Papakonstantinou, D. Quass, A. Rajaraman, Y. Sagiv, J. D. Ullman, V. Vassalos, J. Widom, The TSIMMIS approach to mediation: Data models and languages, Journal of Intelligent Information Systems 8 (2) (1997) 117--132.

Le XML data model: modèle de données XML est un autre graphe qui permet de représenter une collection d'objets simples (atomiques) et complexes. Cependant, tant que le modèle OEM dénote des graphes avec des arcs étiquetés le XML data model dénote un graphe avec des nœuds étiquetés. La FIGURE 3.2 illustre l'exemple de la FIGURE 3.1 en un XML data model.

Le XML data model offre un mécanisme (IDREFs) pour définir des références, c'est un unique identificateur des éléments. En utilisant des références, on peut référer à un autre élément en utilisant son identificateur. Regarder par exemple la flèche en pointillés dans le FIGURE 3.2, qui réfère de l'élément 7 vers l'élément 4

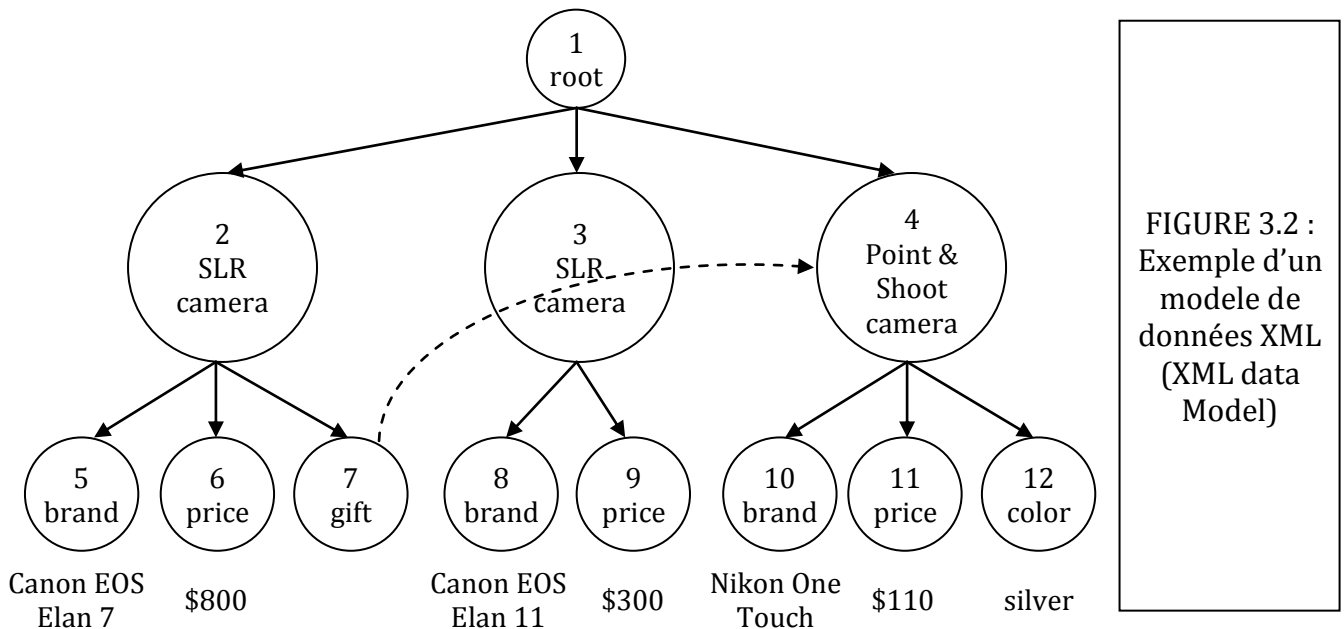


FIGURE 3.2 :
Exemple d'un
modele de
donnees XML
(XML data
Model)

Sans IDREFs, le XML data model devient un **arbre enraciné ordonné et étiqueté**. Nous allons utiliser cette représentation d'arbre enraciné ordonné et étiqueté pour représenter des données XML, en exploitant les idées originales des problèmes d'édition des arbres enracinés, étiquetés et ordonnés.

5.2. Edition des arbres

Un arbre enraciné, ordonné et étiqueté T est un ensemble de $(n + 1)$ nœuds $\{r, n_i\}$ avec $i = 1 \dots n$. Les fils de chaque nœud sont ordonnés. Une étiquette est associée à chaque nœud. R est la racine de T et les nœuds restants $n_1 \dots n_n$ sont partitionnés en m ensembles $T_1 \dots T_m$, chacun est un arbre. Ces arbres $(T_1 \dots T_m)$ sont appelés des sous-arbres de la racine de T . La racine de $T_i, i = 1 \dots m$ est le i^{ieme} fils de la racine de T et la racine de T est le parent de la racine de T_i .

Généralement, si $t_1 \dots t_k$ sont des sous arbres de la racine de l'arbre t , avec t est un sous arbre de T , alors la racine de t est le **parent** de chaque racine des $t_1 \dots t_k$ et les racines de $t_1 \dots t_k$ sont des **fils** de la racine de t . Un nœud x est un **ancêtre** de y et y est un **descendant** de x s'il existe un chemin de nœuds n_0, n_1, \dots, n_k tel que $x = n_0, y = n_k$ et $n_i = \text{parent}(n_{i+1})$ pour chaque $i = 0 \dots k$. Une **feuille** est un nœud sans descendants.

Une **opération atomique** d'édition sur un arbre enraciné ordonné et étiqueté est soit la **suppression** d'un nœud, soit l'**insertion** d'un nœud ou soit le **remplacement** d'un nœud par un autre. Une **opération complexe** d'édition sur arbre est un ensemble d'opérations atomiques d'édition d'arbres, traitées comme une seule opération. Un exemple d'une opération complexe d'édition est l'insertion d'un arbre complet comme un sous arbre dans un autre arbre, qui en fait, une séquence d'opérations atomiques d'insertion de nœuds.

Définition 1 : soit T_1 et T_2 deux arbres enracinés ordonnés et étiquetés. La séquence d'édition d'arbre est une séquence d'opérations d'édition qui transforme T_1 en T_2 .

Définition 2 : soit T_1 et T_2 deux arbres enracinés ordonnés et étiquetés. Etant donné qu'à chaque opération d'édition un coût. Donc à la séquence d'édition un coût. La distance d'édition entre T_1 et T_2 est le coût minimum entre tous les coûts des séquences d'édition, qui peuvent exister, qui transforment T_1 en T_2 . FIGURE 3.3 illustre un exemple d'une séquence d'édition d'insertion, suppression et de remplacement de nœuds pour transformer l'arbre T_1 en T_2 . Etant affecté à chaque opération un coût, la séquence d'édition de cet exemple n'a pas le coût minimum.

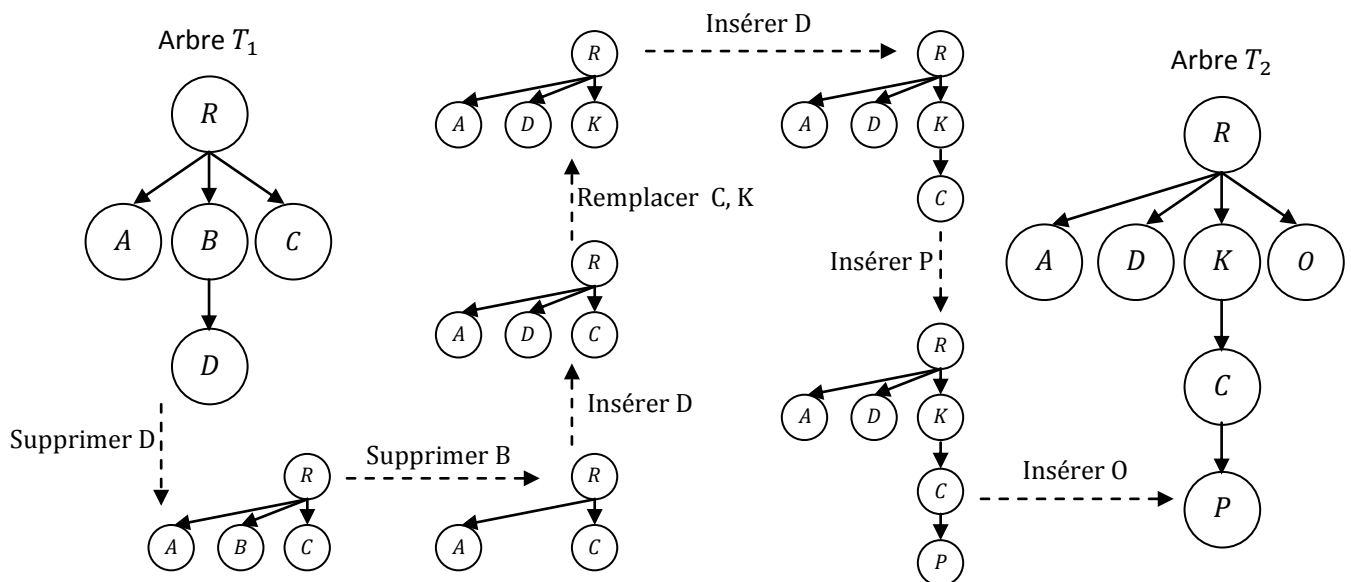


FIGURE 3.3 Exemple d'une séquence d'édition pour transformer T_1 en T_2

Trois Algorithmes parmi les cinq, à savoir Selkow^[44] Dalamagas^[45] Chawathe^[46], qu'on a choisi de tester déterminent la distance d'édition pour comparer les résumés des documents XML. Avant de décrire chaque algorithme en détail, nous présenterons une forme générale de ces opérations d'édition d'arbre que les trois algorithmes utilisent.

⁴⁴ S. M. Selkow, The tree-to-tree editing problem, Information Processing Letters 6 (1977) 184--186

⁴⁵ Theodore Dalamagas, Tao Cheng, Klaas-Jan Winkel, Timos Sellis, Clustering XML documents by structure, Inf. Syst., vol 31, 2006, 187--228

⁴⁶ Sudarshan Chawathe, Comparing Hierarchical Data in External Memory, in Proceedings of the Twenty-fifth International Conference on Very Large Data Bases, 1999, 90--101

5.2.1. Insertion d'un nœud

(a) Variation I ($\text{Ins}^I(x, y, i)$): dans cette variation chaque nœud est inséré seulement comme feuille. Soit x le nœud à insérer comme le $i^{\text{ème}}$ fils du nœud y dans l'arbre t_1 et $y_1 \dots y_n$ sont les fils de y . Dans le nouvel arbre t_2 produit après l'insertion du nœud x , le nœud y aura $y_1 \dots y_{i-1}, x, y_i, y_{i+1} \dots y_n$ comme fils.

(b) Variation II ($\text{Ins}^II(x, y, i)$): dans cette variation, la restriction qui stipule qu'un nœud ne doit être inséré que comme feuille est abrogée. Soit x un nœud à insérer comme le $i^{\text{ème}}$ fils du nœud y dans l'arbre t_1 et $y_1 \dots y_n$ sont des fils de y . Dans le nouvel arbre t_2 produit après l'insertion du nœud x , x possédera une sous-séquence de fils de y comme ses propre fils. Ainsi, étant donné p , un nœud y aura $y_1 \dots y_i, x, y_{p+1} \dots y_n$ comme fils et x aura $y_{i+1}, y_{i+2} \dots y_p$, comme fils.

Nous associons un coût unitaire $c_i(x)$ à l'opération d'insertion d'un nœud x .

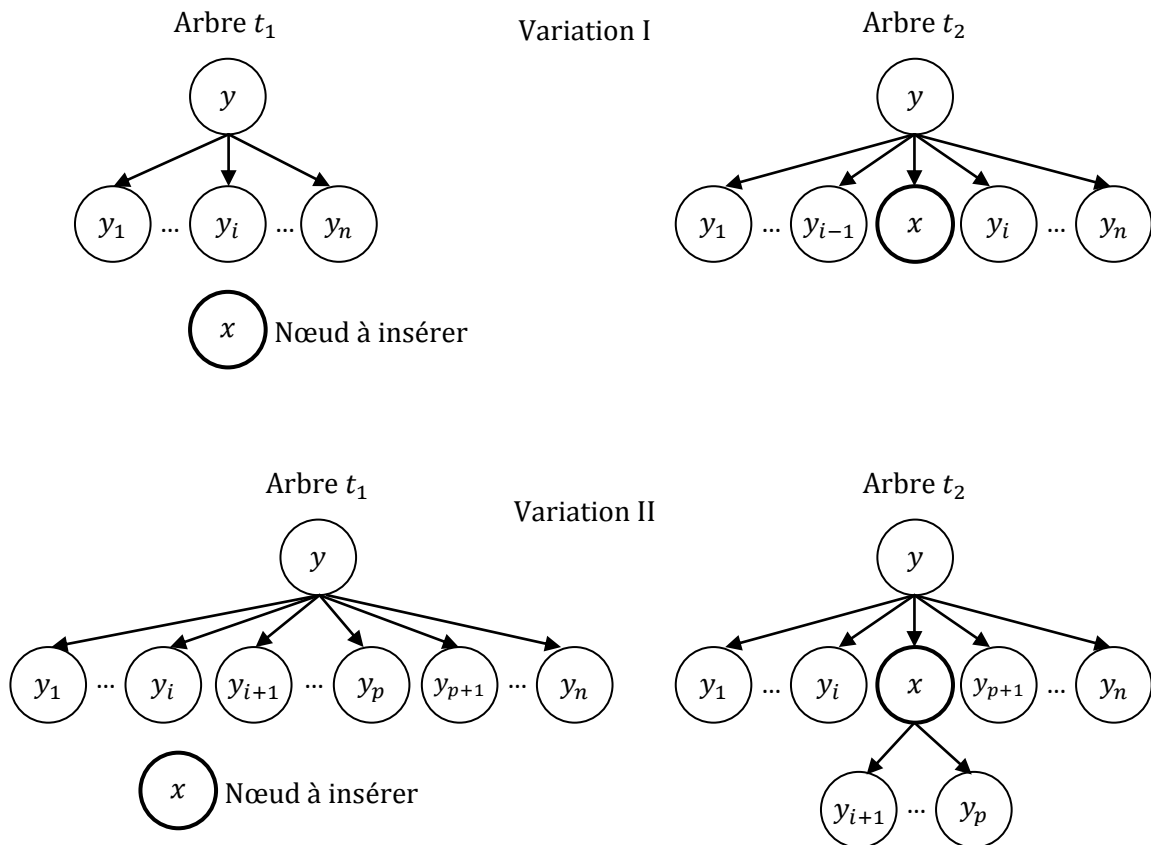


FIGURE 3.4 : Schéma illustratif des opérations d'insertion (Variation I et Variation II)

5.2.2. Suppression d'un nœud

(a) Variation I ($Del(y)$): dans cette variation, les fils du nœud supprimés deviennent les fils de son parent. Soit x un nœud de l'arbre t_1 et $x_1 \dots x_n$ sont les fils de x . Soit $y = x_i$ un des nœuds de $x_1 \dots x_n$ avec $y_1 \dots y_m$ comme nœud fils. Dans le nouvel arbre t_2 produit après la suppression du nœud y , le nœud x aura $x_1 \dots x_{i-1}, y_1 \dots y_m, x_{i+1} \dots x_n$ comme nœuds fils.

(b) variation II ($Del^l(y)$): dans cette variation, seulement les nœuds feuilles peuvent être supprimés. Soit x un nœud dans l'arbre t_1 et $x_1 \dots x_n$ des fils de x . soit $y = x_i$ un des nœuds de $x_1 \dots x_n$ (une feuille). Dans le nouvel arbre t_2 produit après la suppression du nœud y , le nœud x aura $x_1 \dots x_{i-1}, x_{i+1} \dots x_n$ comme fils.

Nous associons un coût unitaire $c_d(y)$ à l'opération de suppression d'un nœud y .

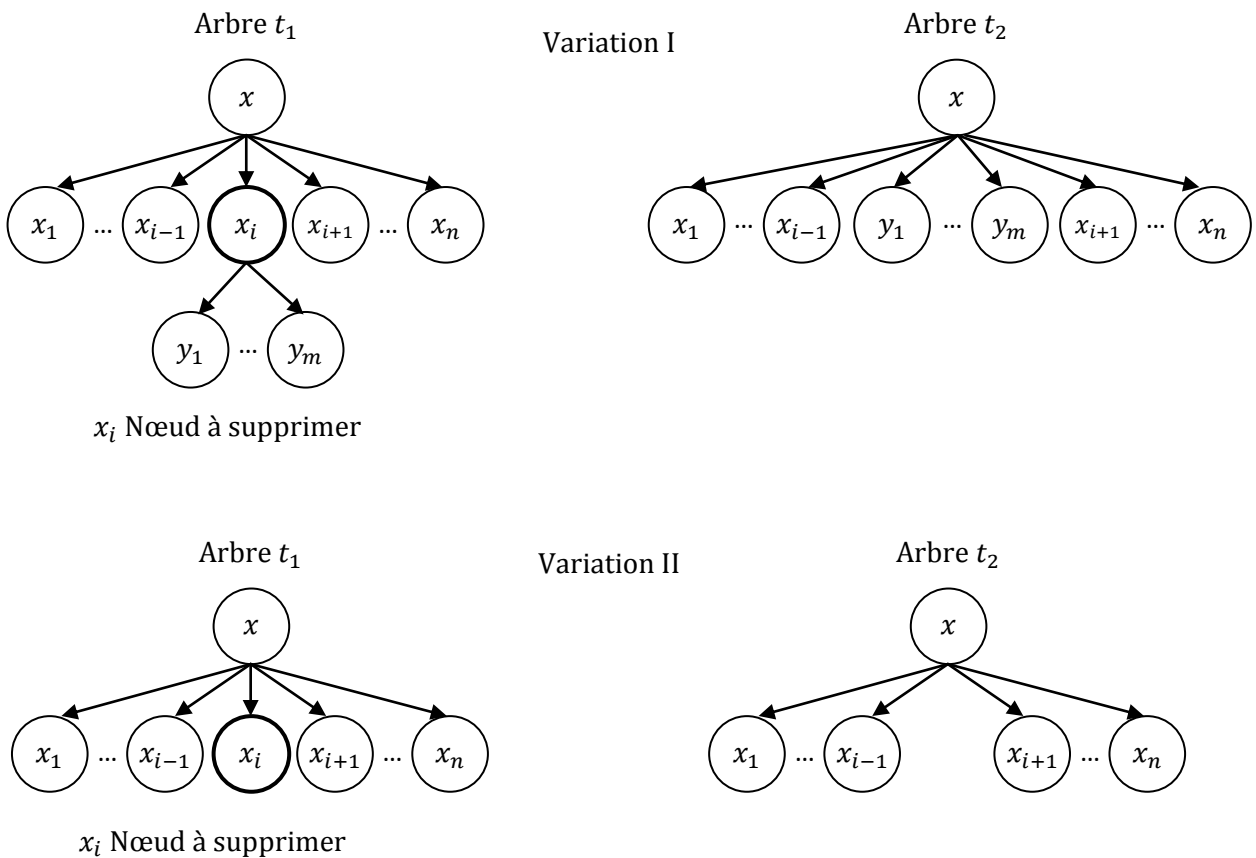


FIGURE 3.5 : Schéma illustratif des opérations de suppression (Variation I et Variation II)

5.2.3. Remplacement (substitution) d'un nœud (Rep(x,y))

Soit y un nouveau nœud et x un nœud dans l'arbre t_1 a remplacé par y . Dans le nouvel arbre t_2 produit après le remplacement d'un nœud x par y , le nœud y aura le même parent et les même fils qu'a x dans t_1 .

Nous associons un coût $c_r(x,y)$ pour l'opération de remplacement d'un nœud x par y . Ce coût peut être variable (par exemple 1 si le nœud qui sera remplacé aura une étiquette différente, 0 sinon) ou une unité constante.

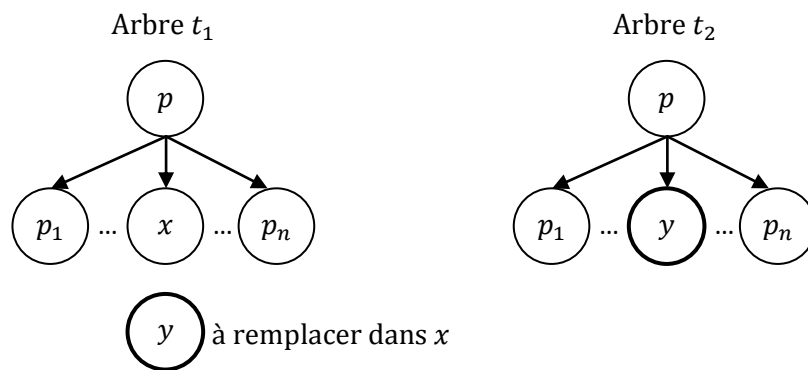


FIGURE 3.6 : Schéma illustratif de l'opération de remplacement

6. Les algorithmes de distances d'édérations

Trois parmi les cinq algorithmes que nous allons décrire, utilisent le même principe pour comparer les documents XML. Ils calculent la distance d'édition entre les résumés d'arbres. Avec cette valeur, on calculera la similarité entre les documents XML et on les classe dans les clusters appropriés. Les algorithmes qui utilisent la distance d'édition sont : Selkow, Chawathe et Dalamagas, c'est trois algorithmes retournent exactement les même résultats parce qu'ils utilisent une mesure commune mais calculée de manières différentes, donc avec des temps d'exécution différents. Nous allons les présenter ces trois algorithmes dans cette partie.

La quatrième méthode, l'algorithme de Leveinshtein utilise le même principe qui est la distance d'édition, mais au lieu de prendre le document XML comme un arbre et le prend comme une chaîne de caractères, dont les étiquètes des nœuds forment ces chaînes de caractères, il sera décrit en détails dans la prochaine partie.

L'algorithme de AïtElhadj, quant à lui, n'utilise pas la distance d'édition, mais une autre mesure de similarité. Il sera décrit en détail ultérieurement.

6.1. Algorithme de Selkow

Selkow [47] propose un algorithme récursif pour calculer la distance d'édition entre deux arbres enracinés, ordonnés étiquetés. L'ensemble des opérations d'édition autorisées par cet algorithme sont :

$\{Ins^l(x, y, i), Del^l(y), Rep(x, y)\}$, avec des coûts $c_i(x)$, $c_d(y)$ et $c_r(x, y)$ respectivement

($c_r(x, y) = 1$ si le nœud a remplacé aura une étiquette différente, $c_r(x, y) = 0$ sinon).

Le coût $W_i(x)$ pour insérer un sous arbre t_2 complet ayant pour racine x , n'importe où dans l'arbre t_1 est le suivant :

$$W_i(x) = \sum_{j=0}^k c_i(x_j) \text{ où } x_0 = x \text{ et } x_1 \dots x_k \text{ sont tous des descendants de } x.$$

Le coût $W_d(y)$ pour supprimer un sous arbre t_2 complet ayant pour racine y , n'importe où dans l'arbre t_1 est le suivant :

$$W_d(y) = \sum_{j=0}^k c_d(y_j) \text{ où } y_0 = y \text{ et } y_1 \dots y_k \text{ sont tous des descendants de } y.$$

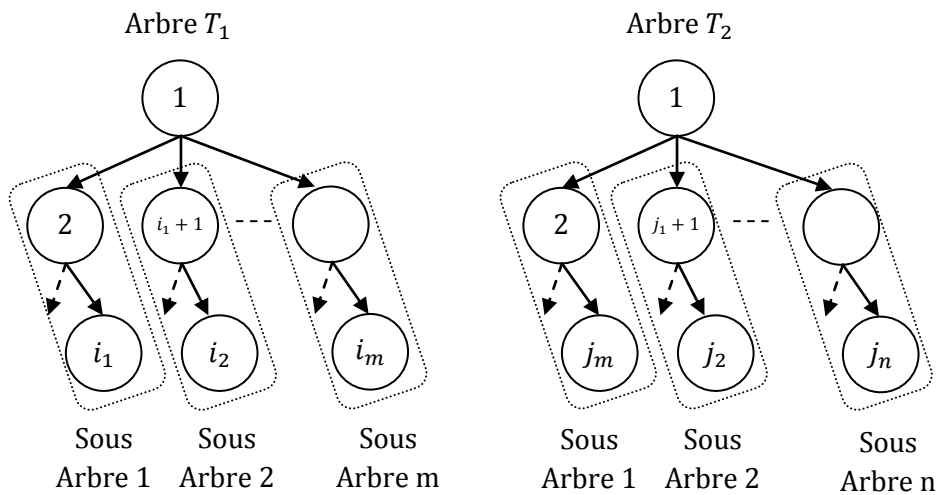


FIGURE 3.7 : Exemple d'arbres

Un arbre T est dénoté $T(1, n_k)$, où 1 est l'étiquette de sa racine, k le nombre des sous arbres connectés à la racine, et n_k le dernier nœud du $k^{\text{ième}}$ sous arbre de T . Tous les nœuds sont étiquetés en respectant un ordre préfixé. FIGURE 3.7 illustre la représentation des deux arbres $T_1(1, i_m)$ et $T_2(1, j_n)$.

⁴⁷ S. M. Selkow, The tree-to-tree editing problem, Information Processing Letters 6 (1977) 184--186.

L'algorithme pour calculer la distance D entre les deux arbres procède récursivement en calculant la distance entre leurs sous arbres. L'idée principale de la récursivité est de calculer la distance d'édition entre deux (sous) arbres t_1 et t_2 requiert le calcul des 4 distances :

1. Distance entre t_1 sans son dernier sous arbre et t_2 .
2. Distance entre t_1 et t_2 sans son dernier sous arbre.
3. Distance entre t_1 sans son dernier sous arbre et t_2 sans son dernier sous arbre
4. Distance entre la dernier sous arbre de t_1 et le dernier sous arbre de t_2 .

Soit r la racine du sous arbre courant t_1 de T_1 , k le nombre de sous arbres dans r , et i le dernier nœud du dernier sous arbre de t_1 ($i = i_k$). Similairement, soit s la racine du sous arbre actuel t_2 de T_2 , l le nombre de sous arbres dans s , et j le dernier nœud du sous arbre de t_2 ($j = j_l$). $D(r, i : s, j)$ dénote la distance structurel entre t_1 et t_2 .

Analogiquement, l'algorithme procède comme suit :

1. *if* ($r == i$) *and* ($s == j$) *then* $D = 0$:

Si t_1 et t_2 consistent en un seul nœud (leur racine), alors le coût pour transformer t_1 vers t_2 est égale à 0 (les racines sont identique parce que les arbres ont été augmentés)

2. *if* ($s == j$) *then* $D = W_d(i_{k-1} + 1) + D(r; i_{k-1} : s; j)$:

Si t_2 consiste en un seul nœud alors le coût pour transformer t_1 en t_2 est égal au coût pour supprimer le $k^{\text{ième}}$ sous arbre de t_1 (qui est le dernier sous arbre de la racine de t_1), plus le coût pour transformer t'_1 (qui est t_1 sans son $k^{\text{ième}}$ sous arbre) en t_2 .

3. *if* ($r == i$) *then* $D = W_i(j_{l-1} + 1) + D(r; i : s; j_{l-1})$:

Si t_1 consiste en un seul nœud, alors le coût pour transformer t_1 en t_2 est égal au coût de l'insertion du $l^{\text{ième}}$ sous arbre de t^2 (qui est le dernier sous arbre de la racine de t^2) dans t^1 plus le coût pour transformer t^1 en t'^2 (qui est t^2 sans son $l^{\text{ième}}$ sous arbre)

4. dans tout autre cas trouver le minimum entre les trois coûts suivants : $D = \min(d_1, d_2, d_3)$

- (a) $d_1 = W_d(i_{k-1} + 1) + D(r; i_{k-1} : s; j)$:

Le coût pour supprimer le $k^{\text{ième}}$ sous arbre de t_1 (qui est la dernier sous arbre de la racine de t_1) plus le coût pour transformer t'_1 (qui est t_1 sans son $k^{\text{ième}}$ sous arbre) en t_2 .

- (b) $d_2 = W_i(j_{l-1} + 1) + D(r; i : s; j_{l-1})$:

Le coût pour insérer le $l^{\text{ième}}$ sous arbre de t_2 (qui est le dernier sous arbre de la racine de t_2) dans t_1 plus le coût pour transformer t_1 en t'_2 (qui est le t_2 sans son $l^{\text{ième}}$ sous arbre)

- (c) $d_3 = D(r; i_{k-1} : s; j_{l-1}) + cr(i_{k-1} + 1; j_{l-1} + 1) + D(i_{k-1} + 1; i_k : j_{l-1} + 1; j_l)$:

Le coût pour transformer t'_1 (qui est t_1 sans son $k^{\text{ième}}$ sous arbre) en t'_2 (qui est t_2 sans son $l^{\text{ième}}$ sous arbre) plus le coût pour remplacer la racine du $k^{\text{ième}}$ sous arbre de t_1 avec la racine du $l^{\text{ième}}$ sous arbre de t_2 plus le coût pour transformer le $k^{\text{ième}}$ sous arbre de t_1 en $l^{\text{ième}}$ sous arbre de t_2 .

Algorithme :

```

D(r; i : s; j)
D(r; i : s; j)
begin
  if ((r == i) and (s == j)) then D = 0 else
    if (s == j) then D = Wd(ik-1 + 1) + D(r; ik-1 : s; j) else
      if (r == i) then D = Wi(jl-1 + 1) + D(r; i : s; jl-1) else
        D = Min {
          {Wd(ik-1 + 1) + D(r; ik-1 : s; j)};
          {Wi(jl-1 + 1) + D(r; i : s; jl-1)};
          {D(r; ik-1 : s; jl-1) + cr(ik-1 + 1; jl-1 + 1) + D(ik-1 + 1; ik : jl-1 + 1; jl)}
        }
  return D
end

```

FIGURE 3.8 : Pseudo-code de l'algorithme de Selkow

La méthode devrait être appelée par : $D(i_0; i_k : j_0; j_l)$, où i_0 est la racine de T_1 , i_k le dernier nœud du $k^{\text{ième}}$ sous arbre (le dernier) de T_1 , j_0 la racine de T_2 , j_l le dernier nœuds de $l^{\text{ième}}$ sous arbre (le dernier) de T_2 . T_1 et T_2 doivent avoir la même racine. Sinon, on créera un nouvel nœud et on le met comme racine pour les deux arbres.

L'algorithme de Selkow était parmi les premiers algorithmes qui permettent de calculer la distance d'édition, et il présente une durée d'exécution, ou une complexité temporelle relativement élevée $O(4^{\min(N,M)})$ où N et M sont le nombre de nœuds de chaque arbre. Depuis son apparition et la constatations de sa lourdeur d'autres algorithmes ont été développés citant, l'algorithme de Chawathe, et l'algorithme de Dalamagas. Ce dernier, est l'algorithme qui présente la meilleure complexité temporelle que nous allons décrire dans le paragraphe suivant. Son principe est proche de celui de Selkow.

6.2. Algorithme de Dalamagas

Cet algorithme a été présenté dans l'article [48], il permet de calculer la distance d'édition entre des arbres enracinés, ordonnés et étiquetés. L'ensemble des opérations d'édition autorisées par cette approche sont : $\{Ins^l(x, y, i), Del^l(y), Rep(x, y)\}$, avec les coûts $c_i(x) = 1$, $c_d(y) = 1$ et $c_r(x, y) = 1$ si le nœud qui sera remplacé aura une étiquette différente, ($c_r(x, y) = 0$ sinon) respectivement. L'algorithme présente une complexité quadratique $O(MN)$ où N et M et le nombre de nœuds de chaque arbre.

Le coût $W_i(x)$ pour insérer un sous arbre complet t_2 , ayant pour racine x , n'importe où dans l'arbre t_1 , est en fait le nombre de nœuds dans t_2 :

$$W_i(x) = \sum_{j=0}^k c_i(x_j) = k + 1 \text{ Où } x_0 = x \text{ et } x_1 \dots x_k \text{ sont tous des descendants de } x.$$

Le coût $W_d(x)$ pour supprimer un sous arbre complet t_2 , ayant pour racine y , n'importe où dans l'arbre t_1 , est en fait le nombre de nœuds dans t_2 :

$$W_d(y) = \sum_{j=0}^k c_d(y_j) = k + 1 \text{ Où } y_0 = y \text{ et } y_1 \dots y_k \text{ sont tous des descendants de } y.$$

Etant donnés deux arbres T_1 et T_2 avec r_1 et r_2 comme racine respectivement, la méthode suivante calcule leur distance d'édition ($CalculateDistance(r_1, r_2)$)

Algorithme

```
int CalculateDistance(TreeNode s, TreeNode t) {
    int[][] D = new int[numOfChildren(s) + 1][numOfChildren(t) + 1];
    D[0][0] = UpdateCost(LabelOf(s), LabelOf(t));
    for (int i = 1; i <= numOfChildren(s); i++)
        D[i][0] = D[i - 1][0] + numOfNodes(si);
    for (int j = 1; j <= numOfChildren(t); j++)
        D[0][j] = D[0][j - 1] + numOfNodes(tj);
    for (int i = 1; i <= numOfChildren(s); i++)
        for (int j = 1; j <= numOfChildren(t); j++)
            D[i][j] = Min(D[i][j - 1] + numOfNodes(tj),
                          D[i - 1][j] + numOfNodes(si),
                          D[i - 1][j - 1] + CalculateDistance(si, tj)
            );
    return D[numOfChildren(s)][numOfChildren(t)];
}
```

FIGURE 3.9 : Pseudo-code de l'algorithme de Dalamagas

⁴⁸ Theodore Dalamagas, Tao Cheng, Klaas-Jan Winkel, Timos Sellis, Clustering XML documents by structure, Inf. Syst., vol 31, 2006, 187--228

Où:

1. s_i est le i^{ieme} fils du nœud s et t_j et j^{ieme} fils du nœud t .
2. $numOfChildren(s)$ retourne le nombre de nœuds fils du nœud s .
3. $numOfNodes(s)$ retourne le nombre de nœuds du sous arbre ayant pour racine s (y compris s).
4. $LabelOf(s)$ retourne l'étiquette du nœud s .
5. $UpdateCost(LabelOf(s); LabelOf(t))$ retourne le coût c_r qui rend l'étiquette de s la même que celle de t : si $LabelOf(s) = LabelOf(t)$ return 1 sinon 0.

6.3. L'algorithme de Chawathe

L'algorithme de Chawathe [49] utilise la même représentation des documents XML que celle utilisée par Selkow et Dalamagas, l'arbre enraciné, ordonné et étiqueté avec une légère différence. Cette différence réside dans l'adjonction des valeurs de profondeurs au model.

L'ensemble des opérations d'édition autorisées par cette approche sont: $\{Ins^l(x,y,i)\}, Del^l(y), Rep(x,y)\}$, avec les coûts $c_i(x) = 1$, $c_d(y) = 1$ et $c_r(x,y) = 1$ si le nœud qui sera remplacé aura une étiquette différente, ($c_r(x,y) = 0$ sinon) respectivement.

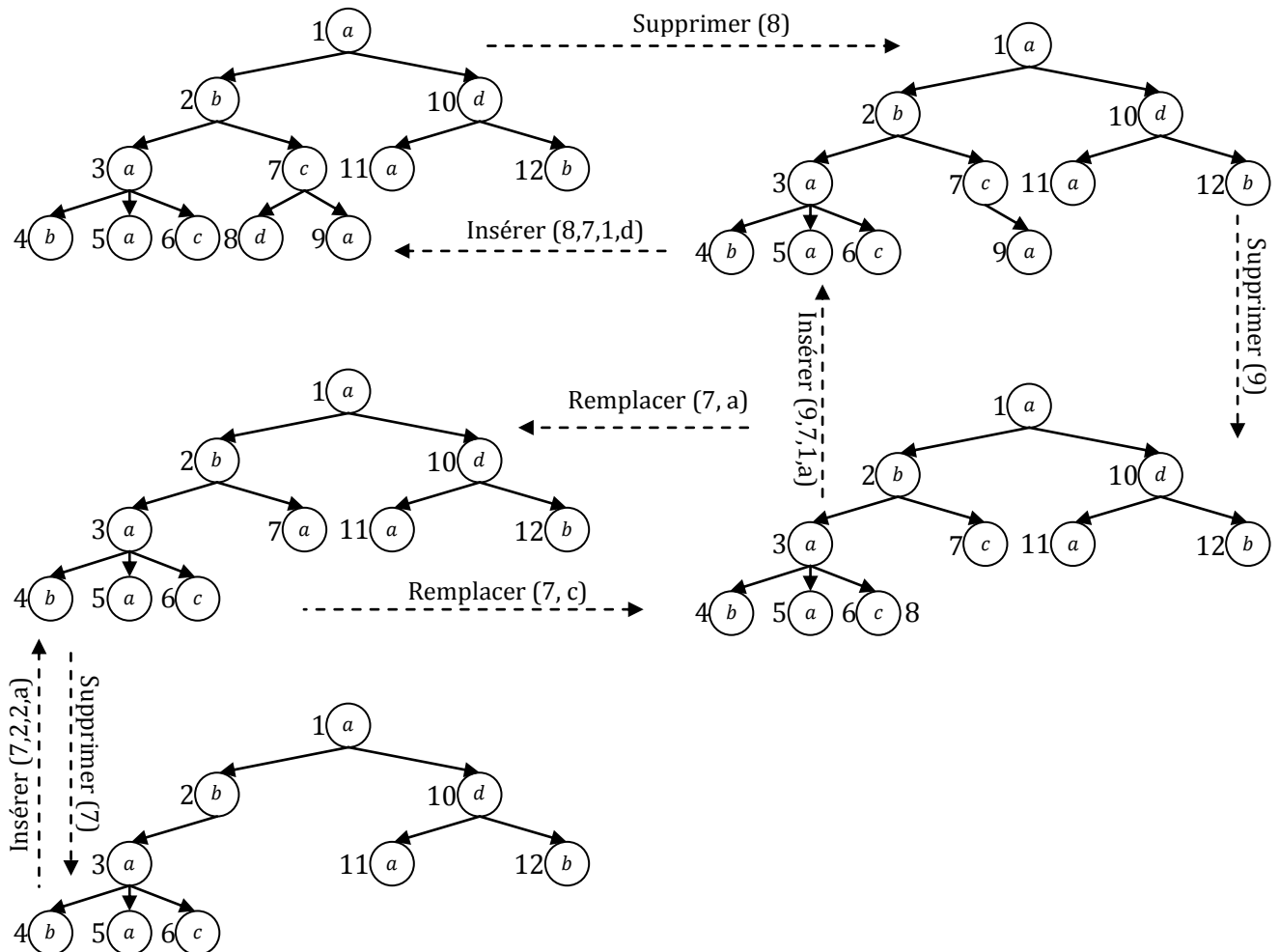


FIGURE 3.10 Des opérations d'édition sur un arbre

Une séquence d'édition (appelé aussi Edit script) est l'ensemble des opérations qui peuvent être appliquées à un arbre pour qu'il se transforme en un autre. L'algorithme de Chawathe diffère des deux autres méthodes précédentes en faisant introduire une structure appelée le **graphe d'édition** (Edit Graph) pour trouver le coût minimum des opérations qu'il faut effectuer pour transformer un arbre A en un Arbre B.

⁴⁹ Sudarshan Chawathe, Comparing Hierarchical Data in External Memory, in Proceedings of the Twenty-fifth International Conference on Very Large Data Bases, 1999, 90--101.

6.3.1. Graphe d'édition

Le graphe d'édition est une structure auxiliaire souvent utilisée dans les algorithmes de comparaison de séquences de données. En fait, le problème de trouver une séquence d'édérations de coût minimum se réduit au problème de trouver le chemin le plus court entre les deux extrémités du graphe. Comme ce graphe a une structure très simple et très régulière, le chemin le plus court peut être trouvé de manière efficace.

Au-dessous, nous présenterons comment le graphe d'édition est utilisé pour les séquences de comparaison et après on introduira les modifications apportées sur ce graphe qui permettront de l'utiliser pour comparer les arbres.

Présentation générale des graphes d'édition

Le graphe d'édition a deux séquences : $A = (A[1] A[2] \dots A[m])$ et $B = (B[1] B[2] \dots B[n])$ et une matrice $(m + 1) \times (n + 1)$.

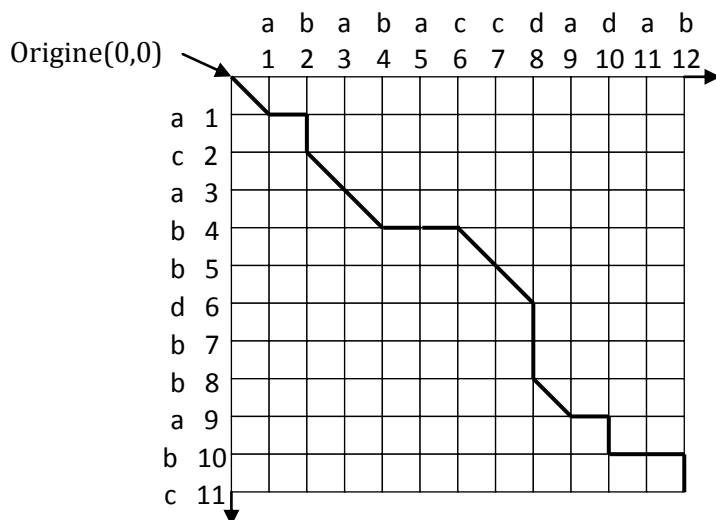


FIGURE 3.11
Exemple d'un
graphe
d'édition pour
séquence de
comparaisons

Chaque croisement d'une ligne verticale avec une ligne horizontale dans le graphe d'édition est un nœud. Un point (x, y) correspond à la paire $(A[x], B[y])$, pour $x \in [1, m]$ et $y \in [1, n]$. Il y'a un arc orienté entre chaque deux nœuds, le cas échéant sur la droite. Il y'a un arc orienté entre chaque deux nœuds, le cas échéant vers le bas. Pour plus de clarté ces arcs orientés sont présentés sans la pointe de flèche dans la FIGURE 3.11. Tous les arcs horizontaux sont dirigés vers la droite et tous les arcs verticaux sont dirigés vers le bas. En plus, il ya un arc diagonal de $(x - 1, y - 1)$ vers (x, y) pour tous $x, y > 0$. Pour des besoins de clarté, ces arcs sont omis dans la FIGURE 3.11.

Le graphe d'édition représenté dans la FIGURE 3.11 correspond aux séquences (chaînes) $A = ababaccdadab$ et $B = acabbdbbabc$. Un arc horizontal $((x - 1, y), (x, y))$ dans le graphe d'édition correspond à une suppression $A[x]$. Un arc vertical $((x, y - 1), (x, y))$ correspond à une insertion $B[y]$. Un arc diagonal $((x - 1, y - 1), (x, y))$ correspond au remplacement de $A[x]$ par $B[y]$.

Les arcs dans le graphe d'édition ont un poids égal aux coûts des opérations d'édition qu'ils représentent. Ainsi un arc horizontal $((x - 1, y), (x, y))$ a un poids égal $c_d(A[x])$. Un arc vertical

$((x, y - 1)(x, y))$ a un poids égal $c_i(B[y])$. Et un arc diagonal $((x - 1, y - 1)(x, y))$ a un poids $c_r(A[x], B[y])$

Il peut être montré que chaque coût-min de la séquence d'édition qui transforme A en B peut être mappé en un chemin de $(0,0)$ vers (M, N) dans le graphe d'édition. Inversement, chaque chemin de $(0,0)$ vers (M, N) correspond à une séquence d'édition qui transforme A en B .

Pour l'exemple suggéré dans la FIGURE 3.11, le chemin surligné correspond à cette séquence d'édition : Supprimer($A[2]$), Insérer($B[2]$), Supprimer ($A[5]$), Supprimer($A[6]$), Remplacer($A[7], b$), Supprimer($A[7]$), Supprimer($A[8]$), Supprimer($A[10]$), Insérer ($B[10]$), Supprimer($A[12]$), Insérer($B[11]$). Il est facile de vérifier, en appliquant les opérations précédentes à A on obtiendra B .

Adaptation des graphes d'édition aux arbres

Dans le but d'utiliser les graphes d'édition pour comparer des arbres, nous avons besoin de modifier leur définition pour incorporer les contraintes imposées par les structures arborescentes. Par exemple, nous devons modéliser la contrainte que si un nœud ancêtre est supprimé, alors tous ses descendants doivent être supprimés.

Pour cela, Chawathe a défini une nouvelle représentation appelée la représentation *id - pair*. Nous définissons un *id - pair* d'un nœud d'un arbre par la paire (l, d) , où l est l'étiquette du nœud et d est sa profondeur dans l'arbre. Nous utilisons $p.l$ et $p.d$ pour référencier respectivement son étiquette et sa profondeur. La représentation *id - pair* d'un arbre en une liste de *id - pair*(s), Cette liste est constituée en parcourant l'arbre en préfixé et en marquant les id-pairs des nœuds rencontrés. Dans le reste de cette partie, les arbres sont supposés être en représentation id-pair. (Un arbre peut être converti à cette représentation en utilisant un parcours préfixé).

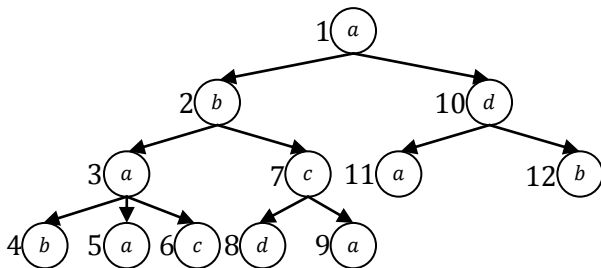


FIGURE 3.12 Arbre A

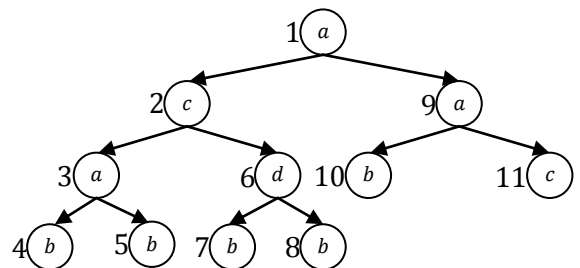


FIGURE 3.13 Arbre B

Soient les arbres A et B représentés dans FIGURE 3.12 et FIGURE 3.13. La lettre de chaque nœud est son étiquette et le numéro à côté et son rang dans l'ordre préfixé. Qui aussi sert d'identificateur. Les représentations *id - pair* de A et B sont comme suit :

$$A = ((a, 0), (b, 1), (a, 2), (b, 3), (a, 3), (c, 3), (c, 2), (d, 3), (a, 3), (d, 1), (a, 2), (b, 2))$$

$$B = ((a, 0), (c, 1), (a, 2), (b, 3), (b, 3), (d, 1), (b, 2), (b, 2), (a, 1), (b, 2), (c, 2))$$

Etant donné un arbre (en représentation *id – pair*) $A = (a_1 a_2 \dots a_M)$, nous utilisons la notation $A[i]$ pour désigner la $i^{\text{ème}}$ nœud a_i de l'arbre A . Ainsi $A[i].l$ et $A[i].d$ dénotent, respectivement, l'étiquette et la profondeur du $i^{\text{ème}}$ nœud de A . Nous définissons le graphe d'édition entre deux arbres A et B comme une matrice $(M + 1) \times (N + 1)$ constituée de nœuds telle qu'elle est présentée dans la FIGURE 3.14.

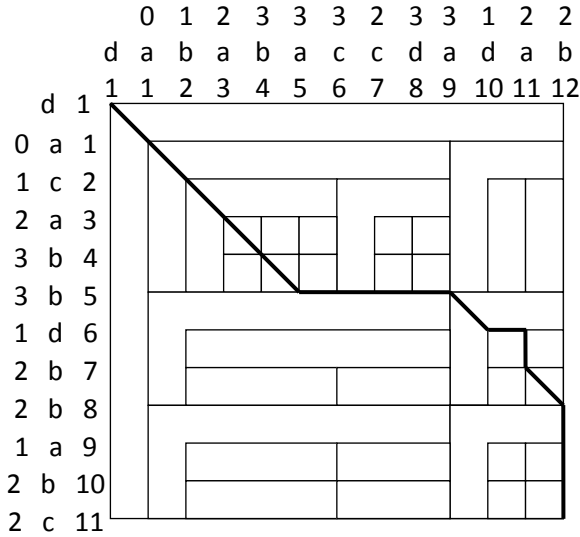


FIGURE 3.14
graphe
d'édition pour
une
comparaison
d'arbre

Il y a un nœud à chaque point (x, y) pour $x \in [0 \dots (M + 1)]$ et $y \in [0 \dots (N + 1)]$. Ces nœuds sont connectés par des arcs orientés comme suit :

- Pour $x \in [0, M - 1]$ et $y \in [0, N - 1]$, il y a un arc diagonal $((x, y), (x + 1, y + 1))$ si et si seulement si $A[x + 1].d = B[y + 1].d$. (pour plus de clarté ces arcs sont omis dans la FIGURE 3.14)
- Pour $x \in [0, M - 1]$ et $y \in [0, N]$, il y a arc horizontal $((x, y), (x + 1, y))$ à moins que $y < N$ et $B[y + 1].d > A[x + 1].d$.
- Pour $x \in [0, M]$ et $y \in [0, N - 1]$, il y a un arc vertical $((x, y), (x, y + 1))$ à moins que $x < M$ et $A[x + 1].d > B[y + 1].d$

Pareil avec le graphe d'édition de la FIGURE 3.11, les arcs horizontaux représentent les suppressions, les arcs verticaux représentent les insertions et les arcs diagonaux représentent l'opération de remplacement. De même, chaque arc a un poids égal au coût de l'opération correspondante.

Cependant, contrairement au graphe d'édition de la FIGURE 3.11, plusieurs arcs horizontaux et verticaux sont absents dans les graphes d'édition de la FIGURE 3.12. Intuitivement, les arcs verticaux manquants assurent la suppression de tous les descendants d'un nœud n lorsque celui-ci si est supprimé; en continuant à traverser les arcs jusqu'à la prochaine opération différente. Similairement, les arcs horizontaux manquants assurent l'insertion, en plus du nœud n , de tous ses nœuds descendants en continuant à parcourir les arcs jusqu'à la prochaine opération différente.

La séquence d'édition à coût-minimum qui transforme A en B peut être mappée en un chemin de $(0,0)$ vers (M,N) dans le graphe d'édition; Réciproquement, chaque chemin de $(0,0)$ vers (M,N) correspond à une séquence d'édition qui transforme A en B .

Par exemple, le chemin indiqué en ligne épaisse dans la FIGURE 3.14 correspond à la séquence d'édition suivante, où n_1, \dots, n_4 sont arbitrairement les identificateurs des nouveaux nœuds insérés :

Remplacer(2,c), Remplacer(5,b), Supprimer(6), Supprimer(8), Supprimer(9), Supprimer(11), Insérer($n_1, 1, 10, b$), Insérer ($n_2, 3, 1, a$), Insérer($n_3, 1, n_2, b$), Insérer($n_4, 2, n_2, c$)

6.3.2. Algorithme

Dans le paragraphe précédent, le problème de calcul du coût-minimum des séquences d'édition entre deux arbres est réduit au problème de trouver le chemin le plus court dans un graphe d'édition.

Etant donné un graphe d'édition G , $(M + 1) \times (N + 1)$, soit D une matrice de $(M + 1) \times (N + 1)$ tel que $D[x, y]$ est la longueur du chemin le plus court de $(0,0)$ vers (x, y) dans le graphe d'édition. Nous appelons D la matrice de distance de G .

Pour des conventions de notation, le poids d'un arc absent dans le graphe d'édition aura une valeur Infini.

Considérons n'importe quel chemin qui connecte l'origine $(0,0)$ à $n = (x, y)$ dans le graphe d'édition. La structure du graphe étant donnée, le nœud précédent dans le chemin est soit le nœud à gauche de n , le nœud au-dessus de n ou le nœud diagonal à gauche et au dessus de n . Donc, la distance de n de l'origine ne peut pas être supérieure à la distance du nœud voisin gauche de l'origine plus le poids de l'arc qui relie le nœud n au nœud voisin gauche. De même pour les relations entre les nœuds voisins au-dessus et diagonal de n , produisant la récurrence suivante pour $D[x, y]$, où $0 < x \leq M$ et $0 < y \leq N$:

$$\begin{aligned}
 D[x, y] &= \min\{m_1, m_2, m_3\} \text{ où} \\
 m_1 &= D[x - 1, y - 1] + c_r(A[x], B[y]), \\
 &\quad \text{si } ((x - 1, y - 1), (x, y)) \in G \\
 &\quad \text{sinon } \infty, \\
 m_2 &= D[x - 1, y] + c_d(A[x]), \\
 &\quad \text{si } ((x - 1, y), (x, y)) \in G \\
 &\quad \text{sinon } \infty, \\
 m_3 &= D[x, y - 1] + c_i(B[y]), \\
 &\quad \text{si } ((x, y - 1), (x, y)) \in G \\
 &\quad \text{sinon } \infty,
 \end{aligned}$$

La récurrence mène à une l'algorithme de programmation dynamique suivant pour calculer la matrice D .

L'algorithme de Chawathe est appelé mmdiff (main-memory differencing):

Input: Deux vecteurs A et B , qui représentent deux arbres en représentation *id – pair*. Ainsi $A[i].l$ et $A[i].d$ dénotent, respectivement, l'étiquette et la profondeur du $i^{ième}$ nœud (en ordre préfixé) de A (et analogiquement pour B). Le nombre d'éléments dans A et B est M et N , respectivement.

Output : La matrice de distance D , où $D[i, j]$ égale à la longueur du chemin le plus court entre $(0,0)$ et (i, j) dans le graphe d'édition de A et B .

Méthode : la FIGURE 3.15 présente le pseudo-code de l'algorithme mmdiff. Nous initialisons la matrice D à l'origine $(0,0)$, suivie par le calcul des distances tout au long des arcs supérieurs et gauches de la matrice. La boucle imbriquée est l'implémentation directe de la récurrence pour $D[i, j]$. Supposons que c_r, c_d, c_i s'exécutent en temps constants α, β, γ respectivement, le temps d'exécution de mmdiff est proportionnel à $1 + \alpha M + \beta N + (\alpha + \beta + \gamma)MN$, ou $O(MN)$

```
D[0,0] = 0;
For i = 1 to M do
    D[i, 0] = D[i - 1, 0] + c_d(A[i]);
For j = 1 to N do
    D[0, j] = D[0, j - 1] + c_i(B[j]);
For i = 1 to M do
    For j = 1 to N do
        begin
            m1 = ∞; m2 = ∞; m3 = ∞;
            if (A[i].d == B[j].d) then m1 = D[i - 1, j - 1] + c_u(A[i], B[j]);
            if (j == N or B[j + 1].d ≤ A[i].d) then m2 = D[i - 1, j] + c_d(A[i]);
            if (i == M or A[j + 1].d ≤ B[j].d) then m3 = D[i, j - 1] + c_i(B[j]);
            D[i, j] = min{m1, m2, m3};
        end;
```

FIGURE 3.15 : Pseudo-code de l'algorithme de Chawathe (mmdiff)

Une fois que nous avons calculé la matrice de distance D , il est facile de récupérer de cette matrice le chemin le plus court entre $(0,0)$ et (M, N) . Nous commencerons à (M, N) et nous traçons la relation de récurrence en arrière pour D . En fonction de l'arc traversé : horizontal, Vertical ou diagonal nous emmètrons l'opération d'édition appropriée : suppression, insertion ou remplacement, respectivement. Ainsi, nous avons l'algorithme suivant pour récupérer la séquence d'édition depuis la matrice d'édition. Ce programme est appelé mmdiff-r (pour mmdiff-recovery)

Algorithme mmdiff-r

Input : deux vecteurs A et B représentant les arbres (comme dans l'algo. mmdiff) et une matrice de distance D calculée par mmdiff.

Output : le coût minimum de la séquence d'édition qui transforme l'arbre A en B .

Méthode : La FIGURE 3.16 présente le pseudo-code de l'algorithme mmdiff-r. Pour chaque itération de la boucle (*while*) i et/ou j sont décrémentés. Ainsi il y a entre $\max(M, N)$ et $M + N$ itérations de la boucle (*while*), avec chaque itération un lot de travail constant est accompli. Ainsi le temps d'exécution de mmdiff-r est $O(M + N)$

```
i = M; j := N;
while (i > 0 and j > 0)do
    if ( $D[i, j] == D[i - 1, j] + c_d(A[i])$  and ( $j == N$  or  $B[j + 1].d \leq A[i].d$ ))then
        begin
            print("supprimer" i);
            i = i - 1
        end
    else
        if ( $D[i, j] == D[i, j - 1] + c_i(B[j])$  and ( $i == M$  or  $A[i + 1].d \leq B[j].d$ ))then
            begin
                print("insérer" j);
                j = j - 1
            end
        else
            begin
                if ( $A[i].l \neq B[j].l$ )then print("remplacer" i B[j].l);
                i = i - 1; j = j - 1;
            end;
    while (i > 0)do
        begin
            print("supprimer" i);
            i = i - 1;
        end;
    while (j > 0)do
        begin
            print("insérer" j);
            j = j - 1;
        end;
```

FIGURE 3.16 : Pseudo-code de l'algorithme de Chawathe (mmdiff-r)

6.4. L'algorithme de Levenshtein

La distance de Levenshtein ^[50] ou distance d'édition, elle mesure la similarité entre deux chaînes de caractères. Elle est égale au nombre minimal de caractères qu'il faut supprimer, insérer ou remplacer pour passer d'une chaîne à une autre.

Cet algorithme diffère des autres dans la mesure où il opère sur des chaînes de caractère. Il a été conçu par Vladimir Levenshtein en 1965. Et il a pour objectif de comparer entre deux chaînes de caractères. Il permet d'effectuer le minimum d'opérations possibles sur une chaîne A pour qu'il la transforme en une chaîne B .

La distance de Levenshtein permet, à partir d'un alphabet et de deux séquences, de déterminer quelle est la longueur d'une séquence minimale d'opérations pour transformer la première séquence en la seconde. Ces opérations sont au nombre de trois et portent sur les lettres de l'alphabet considéré : Remplacement d'une lettre par une autre (R), Suppression d'une lettre (O), Insertion d'une lettre (I)

On a utilisé cet algorithme pour comparer entre deux arbres. Puisque L'algo de Leveï. Opère sur des chaîne de caractères, et non pas sur des arbres, il a fallu transformer l'arbre en une seule grande chaîne de caractères, où les caractères ici représentent les étiquètes des nœuds.

Cette chaîne est construite en faisant le parcours préfixé de l'arbre, et en concaténant les étiquètes à chaque fois qu'on traverse une.

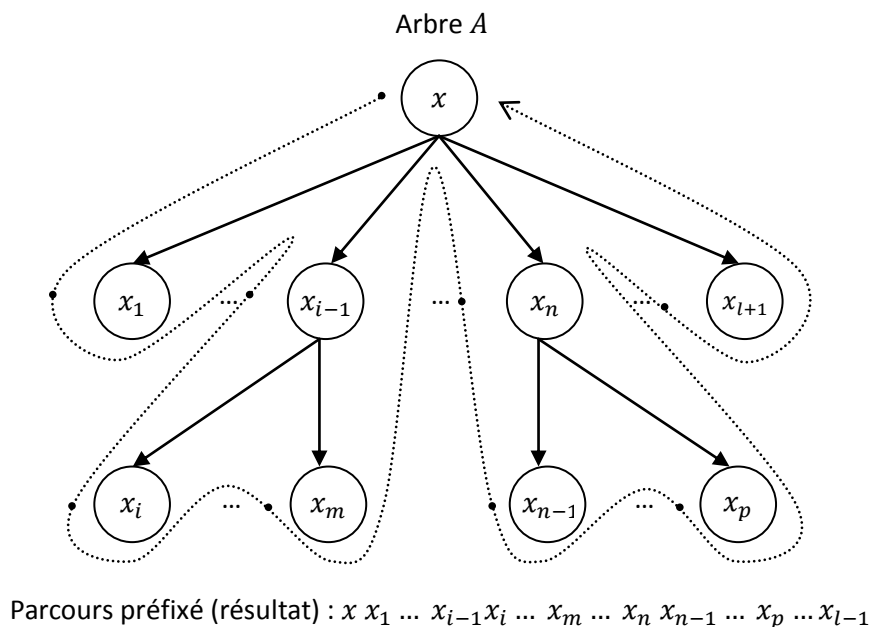


FIGURE 3.17 Le parcours préfixé d'un Arbre A.

⁵⁰ Levenshtein V.I., Binary codes capable of correcting deletions, insertions and reversals, Soviet Physics Doklady 10(8), pp. 707-710, février 1966.

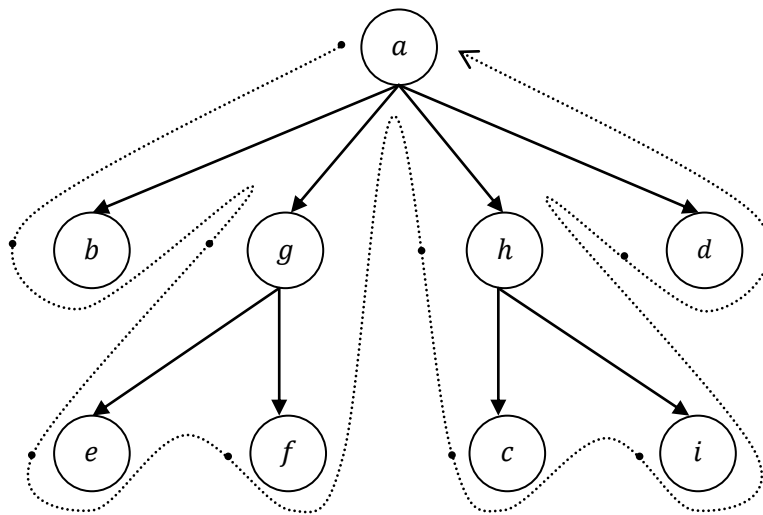


FIGURE 3.18
Exemple
démonstratif
d'un parcours
préfixé

Parcours préfixé (résultat) : *a b g e f h c i d*

L'agorithme de Levenshtein ^[51] :

```
entier DistanceDeLevenshtein(caractere chaine1[1..longueurChaine1],
                             caractere chaine2[1..longueurChaine2])
// d est un tableau de longueurChaine1 + 1 rangées et longueurChaine2 + 1 colonnes
declarer entier d[0..longueurChaine1, 0..longueurChaine2]
// i et j itèrent sur chaine1 et chaine2
declarer entier i, j, coût
pour i de 0 à longueurChaine1
    d[i, 0] := i
pour j de 0 à longueurChaine2
    d[0, j] := j
pour i de 1 à longueurChaine1
    pour j de 1 à longueurChaine2
        si chaine1[i] = chaine2[j] alors coût := 0
        sinon coût := 1
        d[i, j] := minimum(
            d[i - 1, j] + 1, // effacement
            d[i, j - 1] + 1, // insertion
            d[i - 1, j - 1] + coût // substitution
        )
retourner d[longueurChaine1, longueurChaine2]
Sim (c1, c2) = 1 - d[longueurChaine1, longueurChaine2] / dmax
dmax = longueurChaine1 + longueurChaine2
```

FIGURE 3.19 : Pseudo-code de l'algorithme de Levenshteine.

⁵¹ http://fr.wikipedia.org/wiki/Distance_de_Levenshtein

Malgré que l'algorithme de Levenshtein utilise la distance d'édition, qui est la même notion utilisée dans l'algorithme de Selkow, Dalamagas et Chawathe, cette valeur ne donne pas le même résultat. En effet, en transformant l'arbre en une chaîne de caractères, les nœuds perdent la notion de profondeur. Cet algorithme à une complexité de l'ordre $O(MN)$

7. L'algorithme de Aïtelhadj

Le dernier algorithme que nous avons testé est celui de Aïtelhadj [52]. Cet algorithme utilise la même structure que celle utilisée par Chawathe, Selkow et Damagas, c'est l'arbre enraciné, ordonné et étiqueté.

En général, pour comparer deux mots on utilise un dictionnaire. Mais lorsque ces mots correspondent à des nœuds d'un arbre, il est nécessaire de prendre en compte leur contexte hiérarchique, ce qu'il n'était pas le cas pour l'algorithme de Levenshtein, par exemple. L'idée est malgré la représentation de deux nœuds par la même étiquette (même nom), ceci ne veut pas dire qu'ils sont nécessairement similaires dans le contexte où leurs ancêtres ne sont pas complètement identiques. Ainsi la similarité entre deux nœuds ne dépend pas uniquement de leur similarité ontologique (termes devrait être similaire parce qu'ils ont la même chaîne de caractères ou ils sont sémantiquement équivalents en utilisant un dictionnaire), mais aussi de leur ancêtres. Il propose de prendre en compte pour chaque nœud le contexte de ses ancêtres.

7.1. Définition préliminaire

Soit $V = [x_1, x_2, \dots, x_n]$ un vecteur $\in \mathbb{R}^n$, sa norme (Distance Euclidienne) est $\|V\| = \sqrt{\sum_{i=1}^n x_i^2}$

L'usage de la norme permet d'exploiter efficacement le concept du poids. Nous pouvons même étendre son utilisation aux objets qui ne sont pas nécessairement des vecteurs \mathbb{R}^n . en effet, par exemple, si $S = (a, b, b, a, c, b, c)$ est une séquence de nœuds, alors les poids (ou fréquences) de a , b et c sont 2,3,2, respectivement. Cependant, si leur poids sont sauvegardés dans un vecteur tel que $N = [2,3,2]$ alors la norme associée à S est $\|N\| = \sqrt{2^2 + 3^2 + 2^2} = \sqrt{17}$. La norme va servir donc pour la normalisation des valeurs de similarité.

7.2. Mesure de similarité structurelle

Soient T_1 et T_2 deux arbres représentant deux documents XML. Leur similarité est calculée comme suit :

$$\text{Sim}(T_1, T_2) = \frac{\sum_{i=1}^n \sum_{j=1}^m \text{Sim}_{\text{node}}(e_{1i}, e_{2j})}{\|M_1\| * \|M_2\|} \quad (1)$$

⁵² Ali Aïtelhadj, Mohamed Mezghiche et Fatiha Souam, Classification de Structures Arborescentes: Cas de Documents XML, CORIA 2009 - Conférence en Recherche d'Information et Applications.

Où $\lambda \leq \text{Sim}_{\text{node}}(e_{1i}, e_{2j}) < 1$ est la similarité entre les nœuds e_{1i} et e_{2j} . Les nœuds e_{1i} et e_{2j} appartiennent respectivement à T_1 et T_2 . Les nombres $n = |T_1|$ et $m = |T_2|$ sont les tailles (nombre de nœuds) de T_1 et T_2 , respectivement. Le produit $\|M_1\| * \|M_2\|$ permet de normaliser la somme $\sum_{i=1}^n \sum_{j=1}^m \text{Sim}_{\text{node}}(e_{1i}, e_{2j})$. M_1 et M_2 sont deux vecteurs dont les éléments sont les poids des nœuds appartenant respectivement à T_1 et T_2 .

La similarité entre deux nœuds e_{1i} et e_{2j} est exprimée comme suit :

$$\text{Sim}_{\text{node}}(e_{1i}, e_{2j}) = \text{PCC}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j}) * S_0(e_{1i}, e_{2j}) \quad (2)$$

Où $S_0(e_{1i}, e_{2j})$ est la similarité ontologique des nœuds e_{1i} et e_{2j} (obtenu en utilisant un dictionnaire). En d'autres termes, $S_0(e_{1i}, e_{2j}) = 1$ si $(e_{1i} = e_{2j})$, $S_0(e_{1i}, e_{2j}) \cong 1$ si $(e_{1i}$ et $e_{2j})$ sont des synonymes, sinon $S_0(e_{1i}, e_{2j}) = 0$.

$\text{PCC}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j})$ (Path Context Coefficient) Coefficient du chemin contextuel, représente la similarité entre les chemins $r_1 \cdot e_{1i}$ et $r_2 \cdot e_{2j}$ commençant à partir des racines r_1 et r_2 aux nœuds e_{1i} et e_{2j} , respectivement. En effet, $\text{PCC}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j})$ est la similarité entre les ancêtres de e_{1i} et e_{2j} . Les ancêtres jouent un rôle primordial dans le calcul de la similarité. Ceci est basé sur l'idée que même si deux nœuds sont similaires, ils n'ont pas nécessairement un contexte d'ancêtres similaire, ce dernier peut être complètement différent. $\text{PCC}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j})$ est calculé comme suit :

$$\text{PCC}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j}) = \frac{\sum_{k=1}^p \sum_{l=1}^q S_1(n_{1k}, n_{2l})}{\|P_1\| * \|P_2\|} \quad (3)$$

Où $S_1(n_{1k}, n_{2l})$ est la similarité ontologique des nœuds n_{1k} et n_{2l} . Les nœuds n_{1k} et n_{2l} appartiennent respectivement aux chemins $r_1 e_{1i}$ et $r_2 e_{2j}$. Les nombres $p = |r_1 e_{1i}|$ et $q = |r_2 e_{2j}|$ sont les tailles de (longueurs) $r_1 e_{1i}$ et $r_2 e_{2j}$ respectivement. P_1 et P_2 sont deux vecteurs dont les éléments sont les poids des nœuds appartenant respectivement à $r_1 e_{1i}$ et $r_2 e_{2j}$.

Les équations précédentes (1) (2) et (3) doivent en principe être appliquées à tous types de nœuds. Cependant, si au moins un des nœuds e_{1i} et e_{2j} est une racine, il est évident qu'il n'a pas d'ancêtres. Donc, le calcul de la similarité est systématiquement subjective. Pour remédier à cette situation, la similarité $\text{Sim}_{\text{node}}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j})$ est remplacée par la similarité ontologique $S_0(e_{1i}, e_{2j})$. En d'autre terme, $\text{PCC}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j}) = 1$.

Basé sur les équations (1), (2) et (3), le pseudo-code de l'algorithme pour calculer la similarité entre les arbres T_1 et T_2 , est donné dans la FIGURE 3.20.

-Sim est la similarité entre deux arbres T_1 et T_2 , dénotée $\text{Sim}(T_1, T_2)$ dans l'équation (1), alors que PCC représente la similarité entre les ancêtres de e_{1i} et e_{2j} , désigné par $\text{PCC}(r_1 \cdot e_{1i}, r_2 \cdot e_{2j})$ dans l'équation (3).

- $S_0[e_{1i}, e_{2j}]$ est la similarité ontologique des nœuds e_{1i} et e_{2j} telle qu'elle est indiquée dans l'équation (2). La même chose pour $S_1[n_{1k}, n_{2l}]$ avec les nœuds n_{1k} et n_{2l} dans l'équation (3). Donc

for $i \leftarrow 1$ to n et $j \leftarrow 1$ to m , $S_0[e_{1i}, e_{2j}]$ est la matrice de similarité ontologique de (entre) T_1 et T_2 .
 Similairement, for $k \leftarrow 1$ to p et $l \leftarrow 1$ to q , $S_1[n_{1k}, n_{2l}]$ est la matrice de similarité ontologique entre les ancêtres de e_{1i} et e_{2j} .

-les termes $\|P_1\|$, $\|P_2\|$, $\|M_1\|$ et $\|M_2\|$ sont calculés en suivant le norme définie dans le paragraphe « Définition préliminaire » au-dessus.

Algortihme

```

Sim ← 0
for i ← 1 to n do
  for j ← 1 to n do
    if ( $e_{1i}$  est une racine) ou ( $e_{2j}$  est une racine) alors
      PCC ← 1
    sinon
      PCC ← 0
      for k ← 1 to p do
        for l ← 1 to q do
          PCC ← PCC +  $S_1[n_{1k}, n_{2l}]$ 
        end for
      end for
      PCC ←  $\frac{PCC}{\|P_1\| * \|P_2\|}$ 
    end if
    if  $S_0[e_{1i}, e_{2j}] \neq 0$  then
      Sim ← Sim + PCC *  $S_0[e_{1i}, e_{2j}]$ 
    end if
  end for
end for
Sim ←  $\frac{Sim}{\|M_1\| * \|M_2\|}$ 
return Sim

```

FIGURE 3.20 : Pseudo-code de l'algorithme de Aïtelhadj.

L'algorithme à au pire une complexité temporelle de $O(N^2h^2)$. N et h sont respectivement la taille maximale et la longueur maximale des arbres T_1 et T_2 .

8. Clustering des documents XML

Le clustering est le regroupement des documents XML en classes, ou clusters, en fonction d'une valeur de similarité entre les documents XML. Après avoir présenté les algorithmes qui permettent de calculer la similarité entre les documents XML, nous allons présenter la méthode que nous avons utilisée pour classer les documents XML. L'étape de clustering va utiliser :

1. Des résumés de documents XML représentés sous forme d'un arbre enraciné, ordonné et étiqueté.
2. Les distances d'édition entre ces résumés.
3. La distance structurale calculée d'après les distances d'édition d'arbres.
4. L'algorithme de clustering qui utilise les structures de distances pour regrouper les documents.

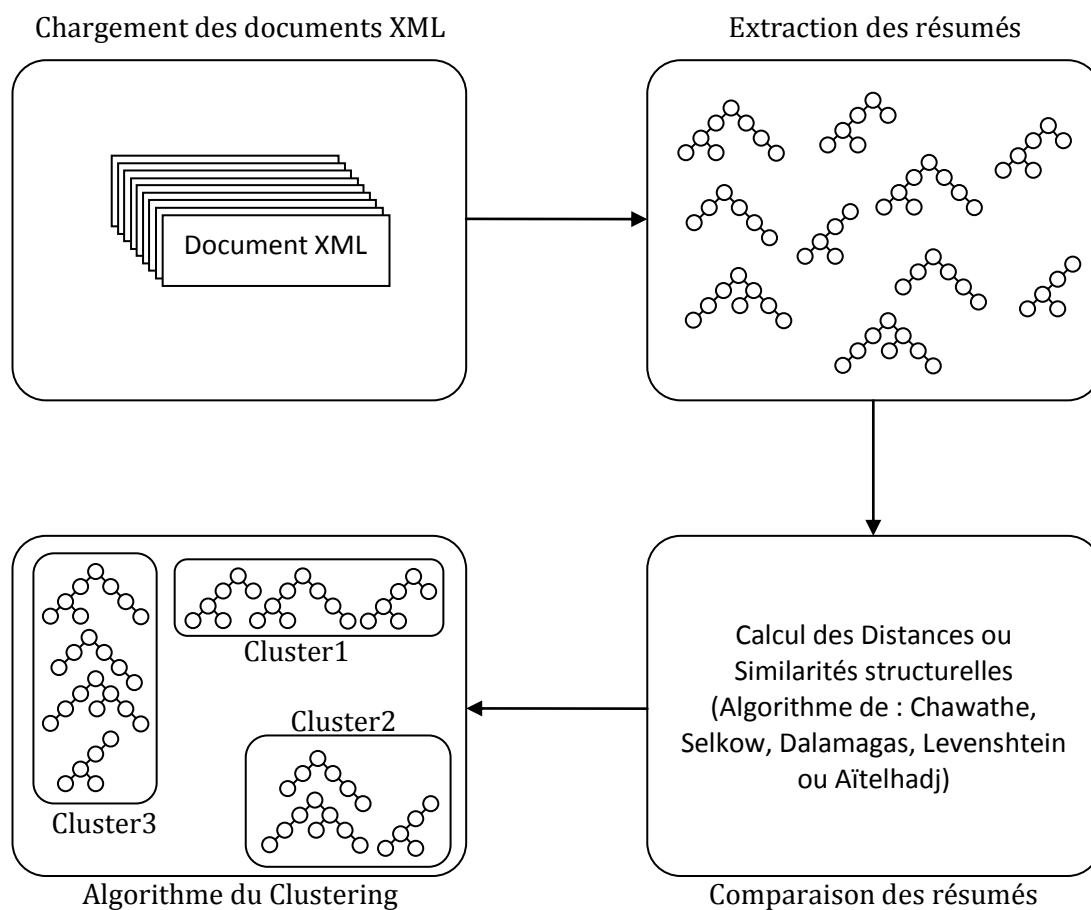


FIGURE 3.21 : Classification des documents XML

8.1. La distance et la similarité structurelle

La distance (similarité) structurelle est une valeur comprise en 0 et 1. Elle détermine la distance (similarité) entre les documents XML. Les 4 premiers algorithmes permettent de calculer la distance d'édition, une mesure utilisée pour calculer la distance structurelle comme suit :

$$\text{DistanceStructurelle}(\text{Arbre1}, \text{Arbre2}) = \frac{\text{distanceD'edition}(\text{Arbre1}, \text{Arbre2})}{d_{\max}}$$

$$\text{avec } d_{\max} = \text{NombreDeNoeuds}(\text{Arbre1}) + \text{NombreDeNoeuds}(\text{Arbre2})$$

Si deux documents sont identiques (hétérogènes), alors la distance (similarité) structurelle entre eux est nulle.

$$\text{DistanceStructurelle} + \text{SimilaritéStructurelle} = 1$$

L'algorithme de AïtElhadj, déduit directement la mesure de similarité structurelle.

Arbre 1	Arbre 2	Distance Structurelle	Similarité Structurelle
A	A	0	1
A	A	1	0
A	B	[0..1]	[0..1]

FIGURE 3.22 : Tableau résumant les valeurs des mesures structurelles.

8.2. Les algorithmes de clustering

Les méthodes de classification sont généralement classées en deux grandes catégories. Les méthodes non hiérarchiques, ils regroupent un ensemble de données en un ensemble de clusters. Les méthodes hiérarchiques, ils créent des fusions entre les ensembles de données, dont lesquels des paires d'éléments (ou clusters) sont successivement fusionner jusqu'à ce qu'il ne y ait qu'un seul grand ensemble contenant tous les éléments.

Les méthodes non hiérarchiques ont de faibles exigences de calcul, $O(k \times n)$, si par exemple n documents à grouper en k clusters, mais certains paramètres tel que le nombre de clusters doivent être connus à priori.

Les méthodes hiérarchiques sont coûteuses en calcul, avec une complexité temporelle de $O(n^2)$, si n documents devraient être classifiés.

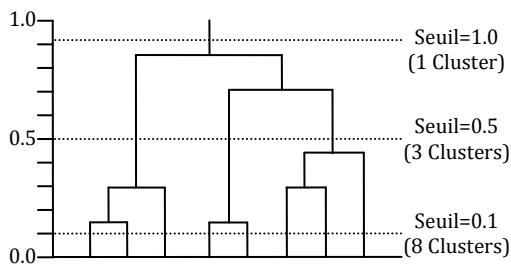


FIGURE 3.23 : Classification hiérarchique.

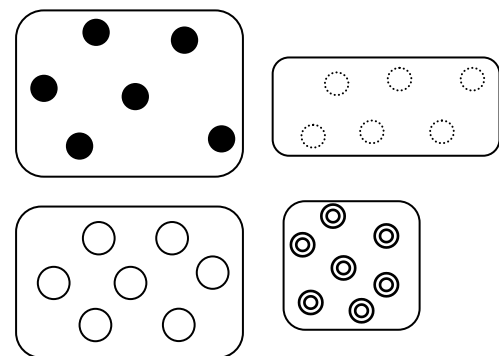


FIGURE 3.24 : Classification non hiérarchique.

Single link, *complete link* et *group link* sont connu comme des méthodes de classification hiérarchiques. Toutes sont basées sur l'idée suivante :

1. Chaque élément de l'ensemble de données à être classifié est considéré comme un cluster.
2. Les clusters avec la distance minimale (i.e similarité maximale) sont regroupés et les distances entre les clusters restant et le nouveau cluster créé, sont recalculés.
3. Tant qu'il y a plus d'un cluster, allez à l'étape2.

Dans le *single link* (*complete link*), la distance entre deux clusters différents est définie comme la distance minimale (maximale) entre tous les paires d'éléments (x, y) , x est dans un cluster et y est dans un autre cluster. Dans *group average link*, la distance entre deux clusters différents est définie comme la distance moyenne entre tous les paires d'éléments (x, y) , x est dans un cluster et y est dans un autre cluster.

Nous avons choisi le *single link* pour comme méthode de clustering pour tester les 5 méthodes décrites précédemment.

8.2.1. Single link

Nous avons implémenté l'algorithme de clustering single link en utilisant l'algorithme de Prim [53] pour calculer le (MST minimum spanning tree ou Le chemin le plus court) d'un graphe.

Etant donné un graphe G avec un ensemble d'arcs pondérés A et un ensemble de nœuds N , un MST est un sous ensemble acyclique $T \subseteq A$ qui chaîne tous les nœuds dont le poids (coût, distance, valeur, ...) total $W(T)$ (la somme des poids des arcs de T) est minimisé. Il a été montré dans [54] que le MST contient toutes les informations requises pour appliquer le single link clustering.

Etant données un ensemble d'arbres enracinés, ordonnés et étiqueté, représentant des documents XML, nous formons un graphe complet G avec n nœuds $\in N$ et $\frac{n(n-1)}{2}$ arcs pondérés $\in A$. Le poids d'un arc correspond à la distance structurelle entre les nœuds (arbres) que cet arc relie.

Le single link clustering pour un seuil de clustering l_1 peut être effectuée en supprimant tous les arcs avec un poids $\geq l_1$ du MST du graphe G . Les nœuds connectés du graphe restant sont des single link clusters.

FIGURE 3.25(a) montre un graphe avec 7 nœuds qui correspondent à 7 documents XML, et 10 arcs. Le poids d'un arc est la distance structurelle entre les documents XML. Par exemple la distance structurelle entre l'arbre 1 et l'arbre 2 est 0.2. Les arcs manquants, se sont les arcs supplémentaires qui rendent le graphe complet, se sont ceux qui ont un poids égal à 1. FIGURE 3.25(b) montre le chemin le plus court (MST) de (a). FIGURE 3.25(c) présente le graphe restant après la suppression de tous les arcs ayant un poids ≥ 0.4 .

Il y a deux composants qui résultent, incluant les nœuds (1,2,3,6) et les nœuds (5,7), respectivement. Ceci indique la présence de deux clusters : cluster1 avec (1,2,3,6) comme membres et cluster 2 avec (5,7) comme membres. Les nœuds qui ne sont pas connectés aux autres nœuds sont considérés des clusters à un seul nœud.

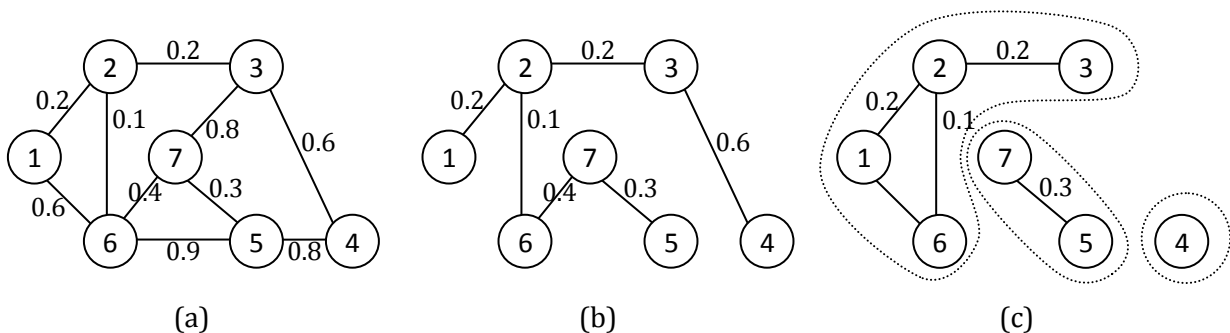


FIGURE 3.25 Le chemin le plus court (MST) et un clustering avec un seuil de 0.4

⁵³ T. Cormen, C. Leiserson, R. Rivest, Introduction to algorithms, MIT Press, 1990

⁵⁴ J. C. Gower, G. J. S. Ross, Minimum spanning trees and single linkage cluster analysis, Applied Statistics 18 (1969) 54--64.

8.2.2. Matrice Associée

Le graphe sera représenté à l'aide d'une matrice appelée la matrice associée. A un graphe $G(N, A)$ d'ordre n on associe une matrice carrée d'ordre n

$$M(n \times n) = \{b_{ij} \mid i = 1..n, j = 1..n\} \text{ tel que } b_{ij} = \{\text{poids si } a(x_i, x_j) \in A\}$$

La représentation du graphe de la FIGURE 3.25(a) en Matrice:

$M(7 \times 7)$	1	2	3	4	5	6	7
1							
2	0.2						
3		0.2					
4			0.6				
5				0.8			
6	0.6	0.1			0.9		
7			0.8		0.3	0.4	

La représentation du graphe de la FIGURE 3.25(b) en Matrice :

$M(7 \times 7)$	1	2	3	4	5	6	7
1							
2	0.2						
3		0.2					
4			0.6				
5							
6		0.1					
7					0.3	0.4	

La représentation du graphe de la FIGURE 3.25(c) en matrice après l'application du seuil 0.4 :

$M(7 \times 7)$	1	2	3	4	5	6	7
1							
2	0.2						
3		0.2					
4							
5							
6		0.1					
7					0.3		

Après il suffit de parcourir la table pour déduire les liens restants entre les documents XML. Pour construire les clusters.

8.2.3. Algorithme de Prim

L'algorithme de Prim ^[55] permet de calculer le plus court chemin dans un graphe pondéré. Etant donné un graphe $G = (N, A)$ avec N = ensemble des nœuds, A = ensemble des arcs. L'algorithme incrémente continuellement la taille de l'arbre, un seul arc à la fois, commençant par un arbre constitué d'un seul nœud, jusqu'à ce qu'un chemin relie tous les nœuds de l'arbre initial.

Input : un graphe non vide et pondéré avec un ensemble de nœuds N et un ensemble d'arcs A .

Initialisation :

$N_{\text{new}} = \{x\}$, où x est un nœud choisi arbitrairement (nœud du départ) de N , $A_{\text{new}} = \{\}$

Répéter jusqu'à $N_{\text{new}} = N$

Choisir un arc $a(u, v)$ avec le poids minimal tel que $u \in N_{\text{new}}$ et $v \notin N_{\text{new}}$.

(s'il ya plusieurs arcs avec le même poids minimal, l'un d'eux est pris)

Ajouter v à N_{new} et (u, v) à A_{new} .

Output : N_{new} et A_{new} représentent le MST (chemin le plus court)

FIGURE 3.26 : Pseudo-code de l'algorithme de Prim.

⁵⁵ R. C. Prim: *Shortest connection networks and some generalizations*. In: *Bell System Technical Journal*, 36 (1957), pp. 1389–1401

9. Résumé du chapitre

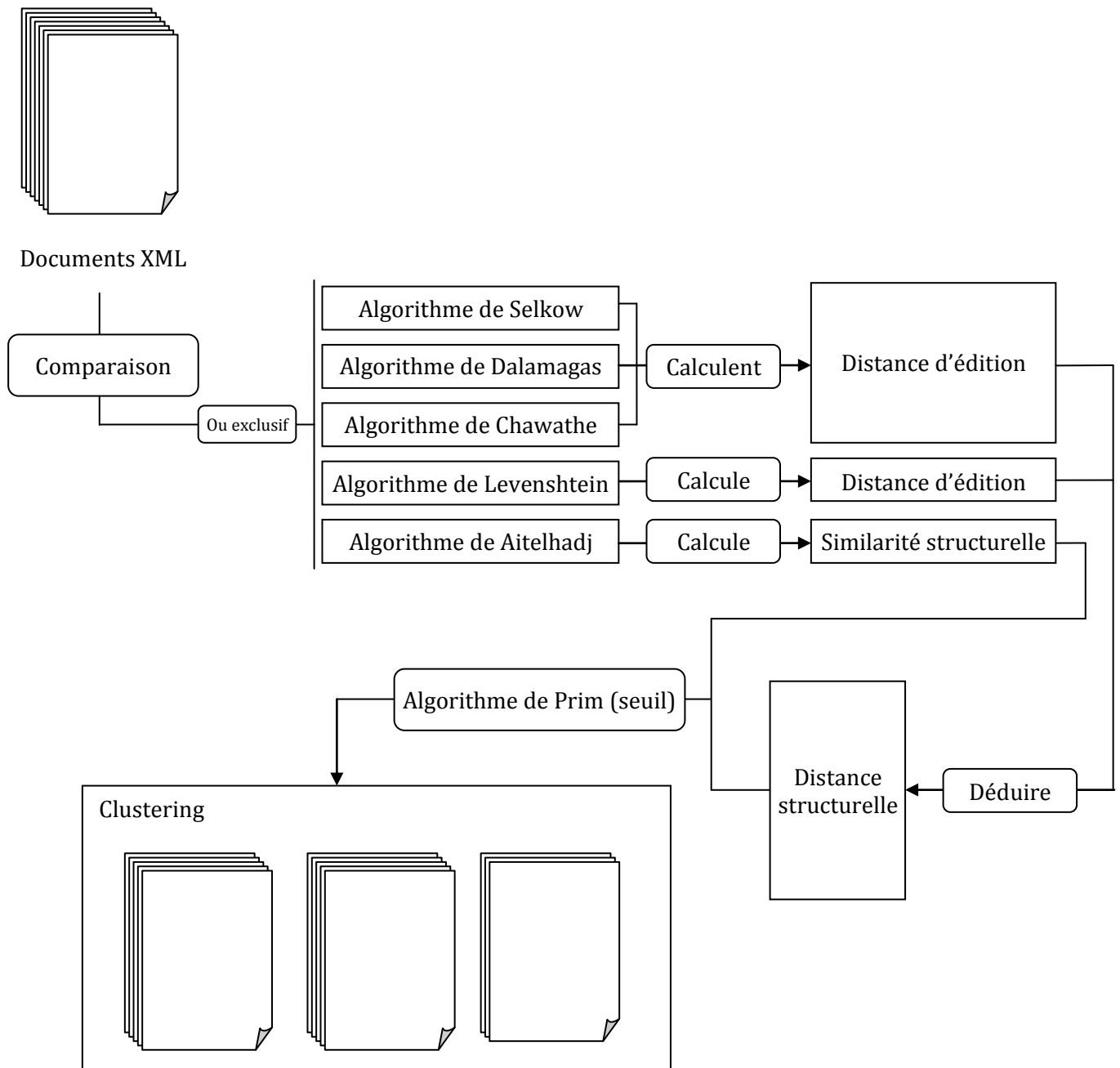


FIGURE 3.27 : Les étapes de la classification.

10. Conclusion

Nous avons présenté dans ce chapitre la classification des documents XML. Plusieurs travaux de recherche relatifs à la classification ont été faits et nous avons présenté cinq algorithmes utilisés pour comparer les documents XML. Chaque algorithme a ses propres performances et qualités de classification.

La performance de la classification dépend de la complexité temporelle de l'algorithme, en revanche la qualité de la classification dépend de la méthode employée pour calculer la similarité structurelle entre les documents XML.

Nous allons dans le prochain chapitre concevoir un logiciel qui incorpore toutes les méthodes de classification vue dans ce chapitre, pour réaliser un logiciel qui fait la classification des documents XML afin d'expérimenter les cinq algorithmes et en déduire les meilleurs algorithmes qui permettent une meilleure classification en qualité et durée d'exécution.

Chapitre 4

Conception et réalisation d'un logiciel de classification

11. Introduction

Dans ce chapitre, nous allons présenter la conception d'un logiciel qui permet de classifier des résumés de documents XML. Un logiciel dans lequel on a implémenté la méthode de Selkow, la méthode de Damalagas, méthode de Chawathe, méthode de Levenshtein et la méthode de AitElhadj. Toutes ces méthodes permettent de calculer la similarité structurelle entre des documents XML pour pouvoir les classer par la suite. Nous présenterons aussi la réalisation du logiciel, le langage de programmation utilisé, la méthode suivie pour le concevoir et aussi l'outil de développement utilisés pour le développer.

12. Conception

12.1. Présentation du logiciel

Le but de la classification est de regrouper des documents XML similaires dans un même dossier (cluster). La similarité entre documents XML peut être calculée en utilisant seulement le contenu du document XML, c.-à-d. seulement les informations qu'il contient, sans prendre en compte le balisage du document. Elle peut aussi être calculée en utilisant seulement la structure hiérarchique du document XML, c'est notre cas d'étude. Comme on peut avoir des méthodes hybrides, plus fiables mais trop coûteuses qui permettent de calculer la similarité en prenant en compte les informations et la structure hiérarchique du document.

Dans notre cas d'étude, on calcule la similarité entre documents XML en utilisant seulement leurs structures hiérarchiques. Donc il faut en extraire cette structure hiérarchique et la représenter sous un format qui pourra être traité par le logiciel que nous allons réaliser.

Étant donné un ensemble de documents XML à classer, il faut extraire leurs représentations hiérarchiques et les représenter sous un format ou une structure de données arborescente qui pourra être passée par la suite au logiciel pour qu'il la traite. C'est la phase de l'Extraction des résumés des documents XML.

Les résumés des documents XML sont des arbres enracinés, ordonnés et étiquetés. Les arbres seront stockés par la suite dans des fichiers textes (Résumés XML).

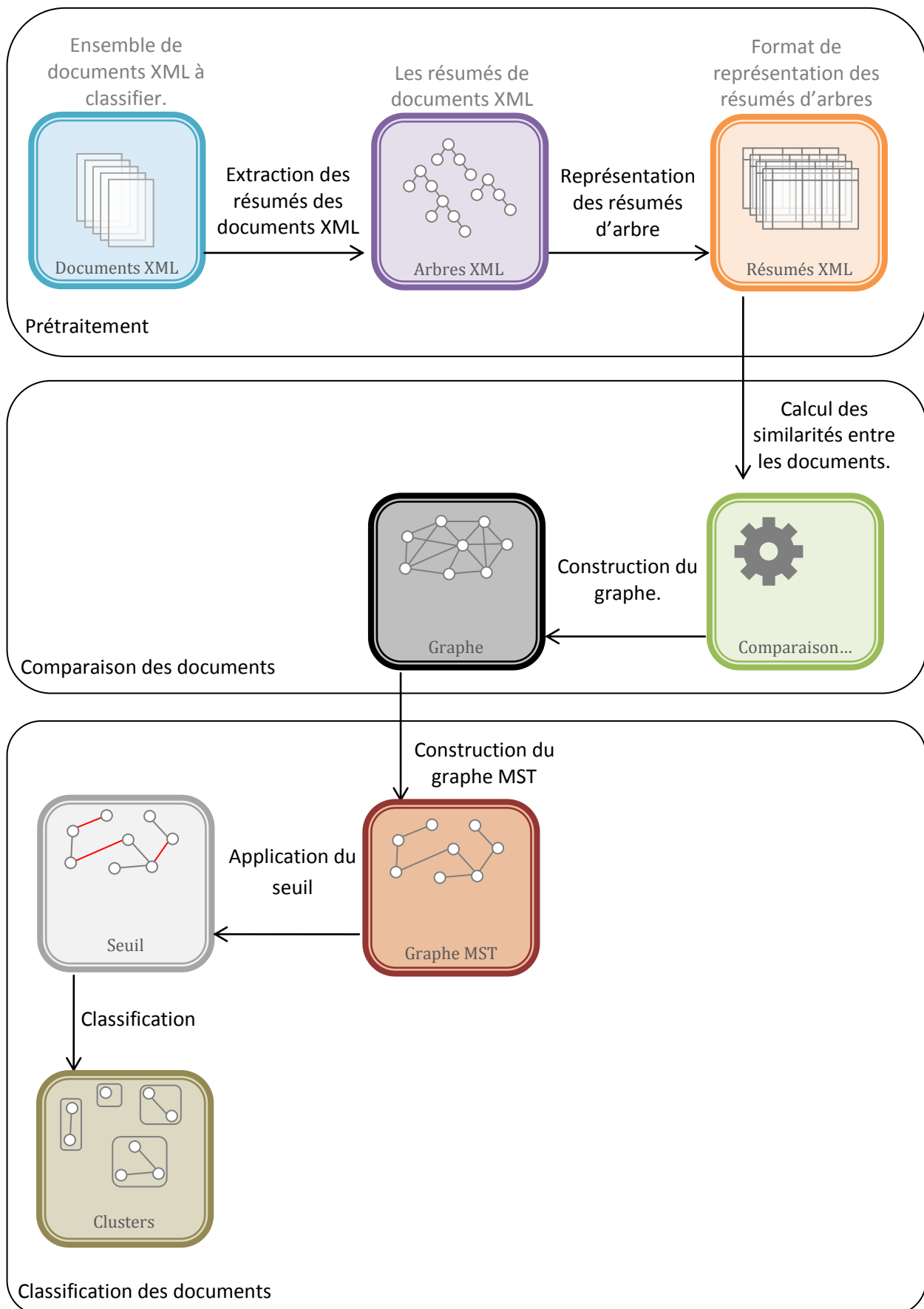
D'après ces représentations hiérarchiques (Résumés XML) qui représentent les documents XML, on calcule la similarité entre eux en utilisant un des algorithmes décrits dans le chapitre 3. Une fois que la similarité entre tous les documents XML, deux à deux, est calculée on construit un graphe complet à arcs pondérés dont les nœuds représentent les documents XML. Un arc représente la similarité entre les deux nœuds qu'il relie et sa pondération (poids) et leur mesure de similarité structurelle (ou la distance structurelle).

Une fois le graphe construit, on génère le graphe MST (Minimum Spanning Tree, le chemin le plus court), en utilisant l'algorithme de Prim, qui relie entre tous les documents XML, avec des arcs ayant les poids les moins élevés. Après il faudra déterminer un seuil pour la classification des documents.

Le seuil permet de déterminer la valeur de précision de la classification, si on veut construire des clusters qui ne contiennent que des documents fortement similaires (identiques) on fixe le seuil à 0, c.-à-d. pas de distance entre les documents, ils sont identiques. Plus le seuil augmente plus la similarité dans les clusters diminue. Le seuil a une valeur comprise entre 0 et 1. $\text{seuil} \in [0,1]$

Depuis le graphe MST, on supprime tous les arcs qui ont un poids \geq seuil. Le poids est la distance structurelle, une valeur comprise entre 0 et 1. En supprimant les arcs qui ont un poids \geq seuil on aura des groupes qui vont se constituer. Tous les éléments qui sont reliés entre eux constituent un cluster. Un nœud qui n'est relié à aucun autre nœud représente un cluster à un seul document.

12.2. Schéma général du logiciel



12.3. Explication des étapes du logiciel

12.3.1. Pré traitement

D'après un document XML, on extrait sa structure hiérarchique. La structure hiérarchique est l'ensemble des balises que contient le document XML. On extrait cette hiérarchie en gardant seulement les titres des balises, pour représenter le document XML sous forme d'un arbre. Dans la structure hiérarchique générée, les informations du document ne sont pas prises en considération.

```
< ?xml version = "1.0" encoding = "UTF-8" ?>
<biblio> <!--Element Racine-->
  <livre> <!--Premier enfant-->
    <!--Premier élément enfant de livre-->
    <titre>Les Misérables</titre>
    <auteur>Victor Hugo</auteur>
    <nb_tomes>3</nb_tomes>
    <localisation travée="1" armoire="4" étagère="2"/>
  </livre>
  <livre> <!--Deuxième enfant-->
    <titre>L'Assommoir</titre>
    <auteur>Emile Zola</auteur>
    <localisation travée="4" armoire="2" étagère="5"/>
  </livre>
  <livre> <!--troisième enfant-->
    <titre>David Copperfield</titre>
    <auteur>Charles Dickens</auteur>
    <nb_tomes>3</nb_tomes>
    <localisation travée="1" armoire="2" étagère="3"/>
  </livre>
</biblio>
```

FIGURE 4.1 : Exemple d'un document XML

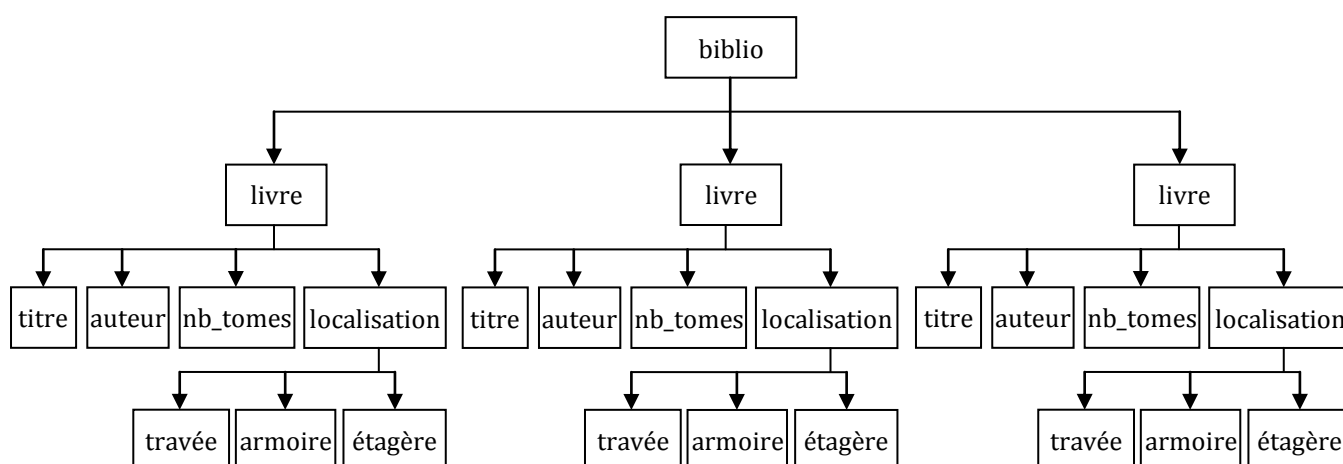


FIGURE 4.2 : La représentation hiérarchique d'un document XML.

Une fois l'arbre construit, on le représente sous forme d'une structure de données enregistrée dans fichiers txt. La structure d'un document XML a la forme suivante :

Numéro du nœud	Noms du nœud	Fils	Frère	Père
Le numéro qui sera attribué au nœud	L'étiquette du nœud (nom de la balise)	Le numéro du nœud fils (le fils le plus à gauche). Les nœuds feuilles auront la valeur -1.	Le numéro du nœud frère voisin droit. Le nœud qui n'a pas de voisin droit aura la valeur -1.	Le numéro du nœud père. Seule la racine a la valeur -1

La présentation de l'arbre précédent donnera le résultat suivant :

Numéro du nœud	Noms du nœud	Fils	Frère	Père
1	biblio	2	-1	-1
2	livre	3	10	1
3	titre	-1	4	2
4	auteur	-1	5	2
5	nb_tomes	-1	6	2
6	localisation	7	-1	2
7	travée	-1	8	6
8	armoire	-1	9	6
9	étagère	-1	-1	6
10	livre	11	18	1
11	titre	-1	12	10
12	auteur	-1	13	10
13	nb_tomes	-1	14	10
14	localisation	15	-1	10
15	travée	-1	16	14
16	armoire	-1	17	14
17	étagère	-1	-1	14
18	livre	19	-1	1
19	titre	-1	20	18
20	auteur	-1	21	18
21	nb_tomes	-1	22	18
22	localisation	23	-1	18
23	travée	-1	24	22
24	armoire	-1	25	22
25	étagère	-1	-1	22

La table précédente sera enregistrée dans un fichier texte qui sera passé comme paramètre au logiciel. Le fichier représente le document XML. Une table pareille sera créée pour chaque document XML qui sera classifié.

12.3.2. Comparaison des documents

La matrice de comparaisons est une matrice contenant les numéros des documents XML avec les valeurs de leur comparaison. Chaque document XML, représenté par le fichier précédant (Résumé XML) aura un numéro qui l'identifiera durant le processus de classification. (Numéros successifs de 1 à T).

Pour calculer la similarité entre les documents, on crée une matrice de (T X T) où T est le nombre de documents à classer. Seulement la moitié de la matrice est nécessaire, car la matrice est symétrique.

Deux documents XML ayant Num1 et Num2 où Num1 l'identificateur du fichier 1 et Num2 est l'identifiant du document 2. La matrice de comparaisons aura la valeur du résultat rendu par un des algorithmes, de calcul de comparaisons des documents, implémentés dans le logiciel qui sera choisi par l'utilisateur (SELKOW, DALAMAGS, CHAWATHE, LEVENSHTAIN, AITELHADJ). La case de la matrice $MAT[num1][num2]$ = similarité entre D1 et D2.

Matrice des distances :

	NumDoc1	NumDoc2	NumDoc3	...	NumDocT
NumDoc1	1.1.val			...	
NumDoc2	2.1.val	2.2.val		...	
NumDoc3	3.1.val	3.2.val	3.3.val	...	
...
NumDocT	T.1.val	T.2.val	T.3.val	...	T.T.val

L'étape de la comparaison s'achève en générant la matrice de comparaison remplie.

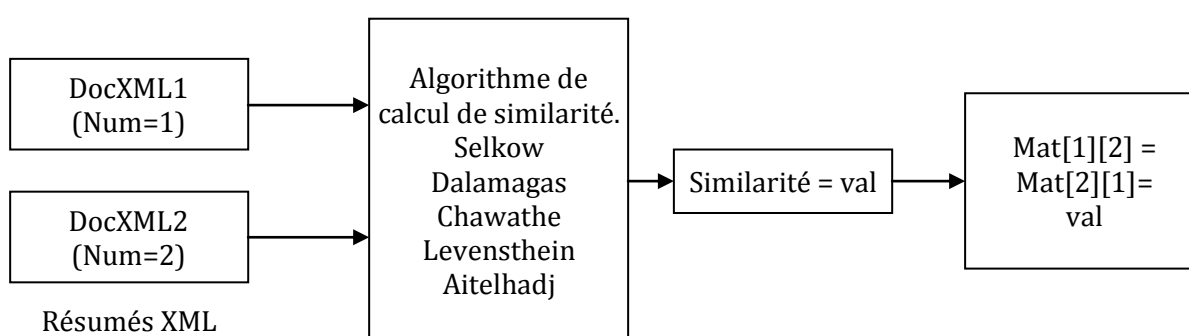


FIGURE 4.3 : Schéma de calcul de la similarité

Construction du graphe

La matrice générée est le graphe de comparaisons des documents XML. En effet les graphes ont plusieurs représentations et parmi ces présentations, on a la matrice d'adjacence. La matrice des distances est une matrice d'adjacence.

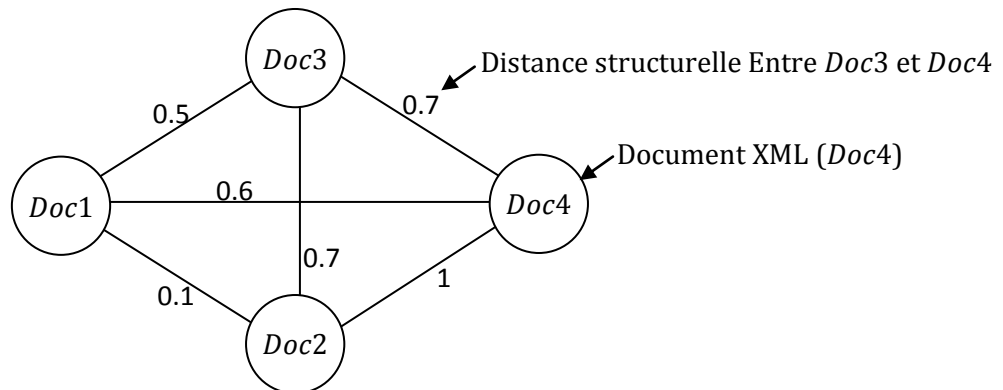


FIGURE 4.4 : Exemple d'un graphe de comparaison de 4 documents XML.

La matrice d'adjacence associée au graphe de la FIGURE 4.4

	<i>Doc1</i>	<i>Doc2</i>	<i>Doc3</i>	<i>Doc4</i>
<i>Doc1</i>	1	0.1	0.5	0.6
<i>Doc2</i>	0.1	1	0.7	1
<i>Doc3</i>	0.5	0.7	1	0.7
<i>Doc4</i>	0.6	1	0.7	1

Soit l'arc $a(\text{Doc3}, \text{Doc4})$, son poids (distance structurelle) = 0,7. Dans la matrice d'adjacence $M[\text{Doc3}][\text{Doc4}] = M[\text{Doc4}][\text{Doc3}] = 0,7$

12.3.3. Phase classification

Le chemin le plus court (MST) : D'après le graphe de la phase précédente, on crée un nouveau graphe qui contient le chemin le plus court qui relie entre tous les documents XML. Le chemin le plus court est calculé en utilisant l'algorithme de Prim. En appliquant l'algorithme de Prim à l'exemple précédent, en commençant par le nom Doc1, on obtient :

Itération0 : (Début)

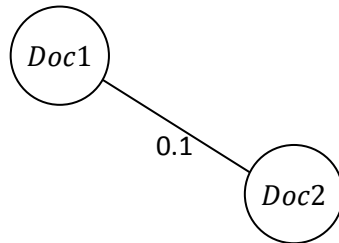


Itération1 :

$$N_{new} = \{Doc1\} A_{new} = \{\}$$

Le minimum entre les arcs $a(Docx, Docy)$ tel que $Docx \in N_{new}$ et $Docy \notin N_{new}$ est $a(Doc1, Doc2) = 0.1$

$$N_{new} = N_{new} \cup \{Doc2\} A_{new} = A_{new} \cup \{a(Doc1, Doc2)\}$$

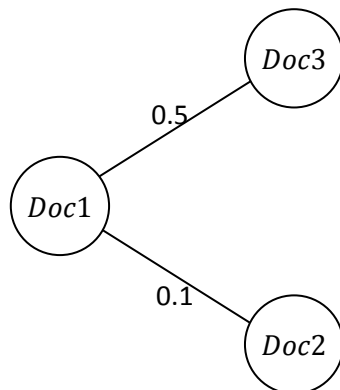


Itération2 :

$$N_{new} = \{Doc1, Doc2\} A_{new} = \{a(Doc1, Doc2)\}$$

Le minimum entre les arcs $b(Docx, Docy)$ tel que $Docx \in N_{new}$ et $Docy \notin N_{new}$ est $b(Doc1, Doc3) = 0.5$

$$N_{new} = N_{new} \cup \{Doc3\} A_{new} = A_{new} \cup \{b(Doc1, Doc3)\}$$



Itération3 :

$$N_{new} = \{Doc1, Doc2, Doc3\} \quad A_{new} = \{a(Doc1, Doc2), b(Doc1, Doc3)\}$$

Le minimum entre les arcs $c(Docx, Docy)$ tel que $Docx \in N_{new}$ et $Docy \notin N_{new}$ est $c(Doc1, Doc4) = 0.6$

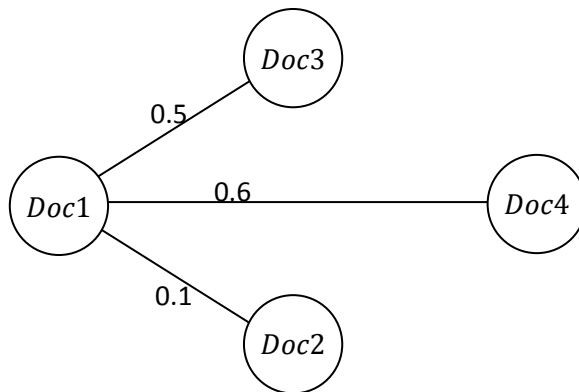
$$N_{new} = N_{new} \cup \{Doc4\} \quad A_{new} = A_{new} \cup \{c(Doc1, Doc4)\}$$

Fin.

Nouveau graphe (MST):

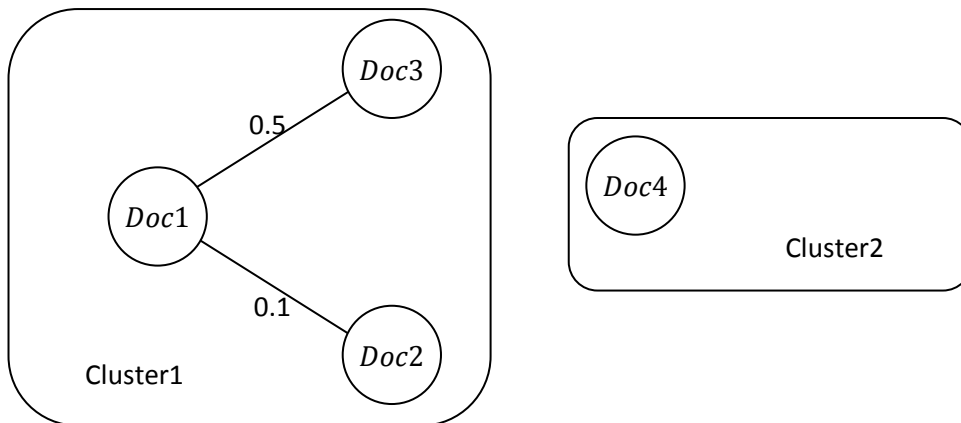
$$N_{new} = \{Doc1, Doc2, Doc3, Doc4\}$$

$$A_{new} = \{a(Doc1, Doc2), b(Doc1, Doc3), c(Doc1, Doc4)\}$$



Détermination du seuil (clustering) : Le seuil est une valeur, comprise en 0 et 1, introduite par l'utilisateur. Elle permet de déterminer le facteur de la précision de la classification. Si deux documents sont identiques cette valeur aura la valeur zéro (0). Plus deux documents divergent plus cette valeur tend vers un (1). Le résultat de la classification dépendra de la valeur du seuil introduite.

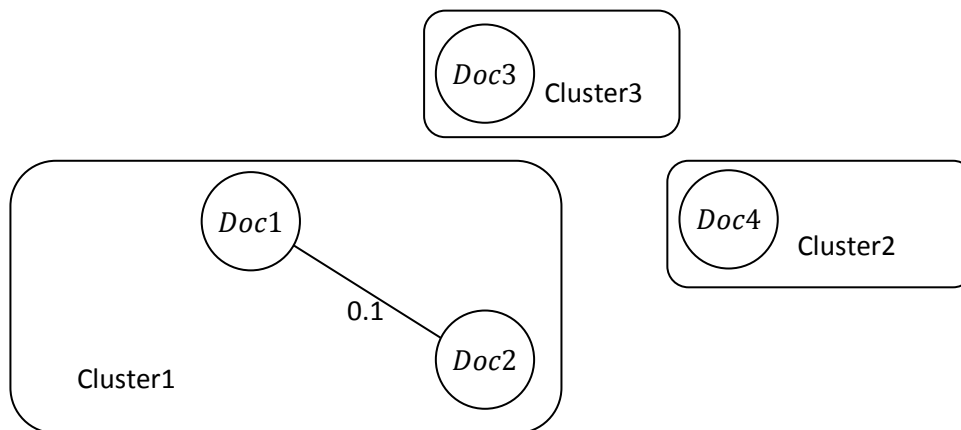
Dans l'exemple précédent, si on fixe le seuil à 0.55 on aura le résultat de la classification suivant : (Tous les arcs qui ont un seuil ≥ 0.55 seront supprimés du MST)



Nous avons eu deux clusters :

$$Cluster1 = \{Doc1, Doc2, Doc3\} \quad Cluster2 = \{Doc4\}$$

Si on fixe le seuil à 0.4 on aura un autre résultat : (Tous les arcs qui ont un seuil ≥ 0.4 seront supprimés du MST)



Nous avons eu trois clusters :

Cluster1 = {Doc1, Doc2} Cluster2 = {Doc4} Cluster3 = {Doc3}

Si on fixe le seuil à zéro on aura 4 clusters à un seul nœud.

Effectuer la classification (classification) : Une fois que l'ensemble des clusters sont créés, il ne reste qu'à procéder à leur classification, en créant des nouveaux dossiers ayant pour nom le « Cluster + numéroDuCluster » puis copier chaque fichier dans le cluster approprié.

12.4. Conception détaillée

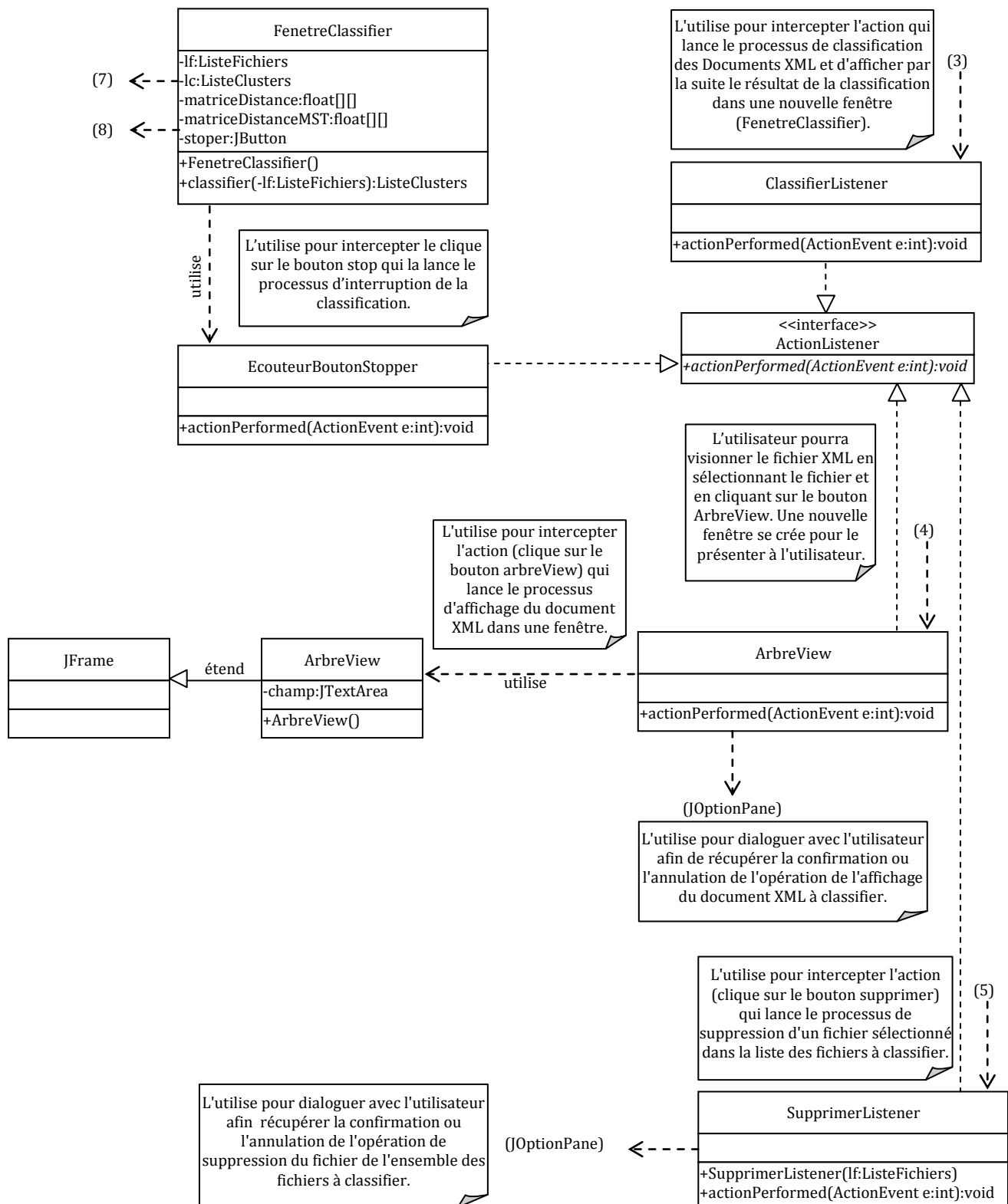
Après avoir présenté le fonctionnement général du logiciel il y a lieu de passer à la présentation de sa conception d'une manière plus détaillée, en utilisant le langage UML. On va présenter les deux principaux diagrammes du logiciel. Un diagramme structurel (ou statique) qui permet de définir la structure du logiciel, qui est le diagramme de classes. Et un diagramme de comportement qui permet de définir le comportement ou la réaction du logiciel aux différentes actions effectuées par l'utilisateur, qui est le diagramme de séquence.

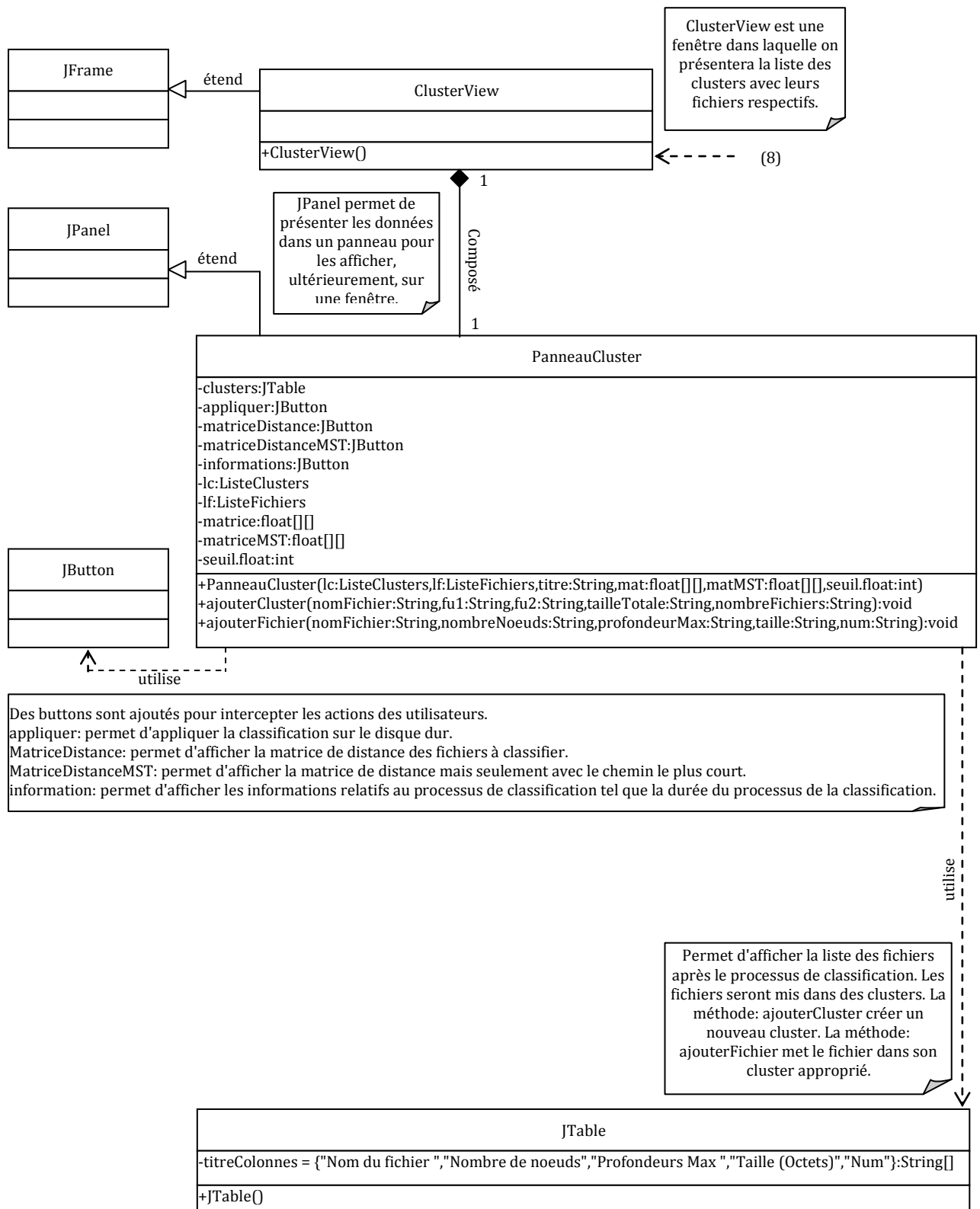
12.4.1. Diagramme de classes

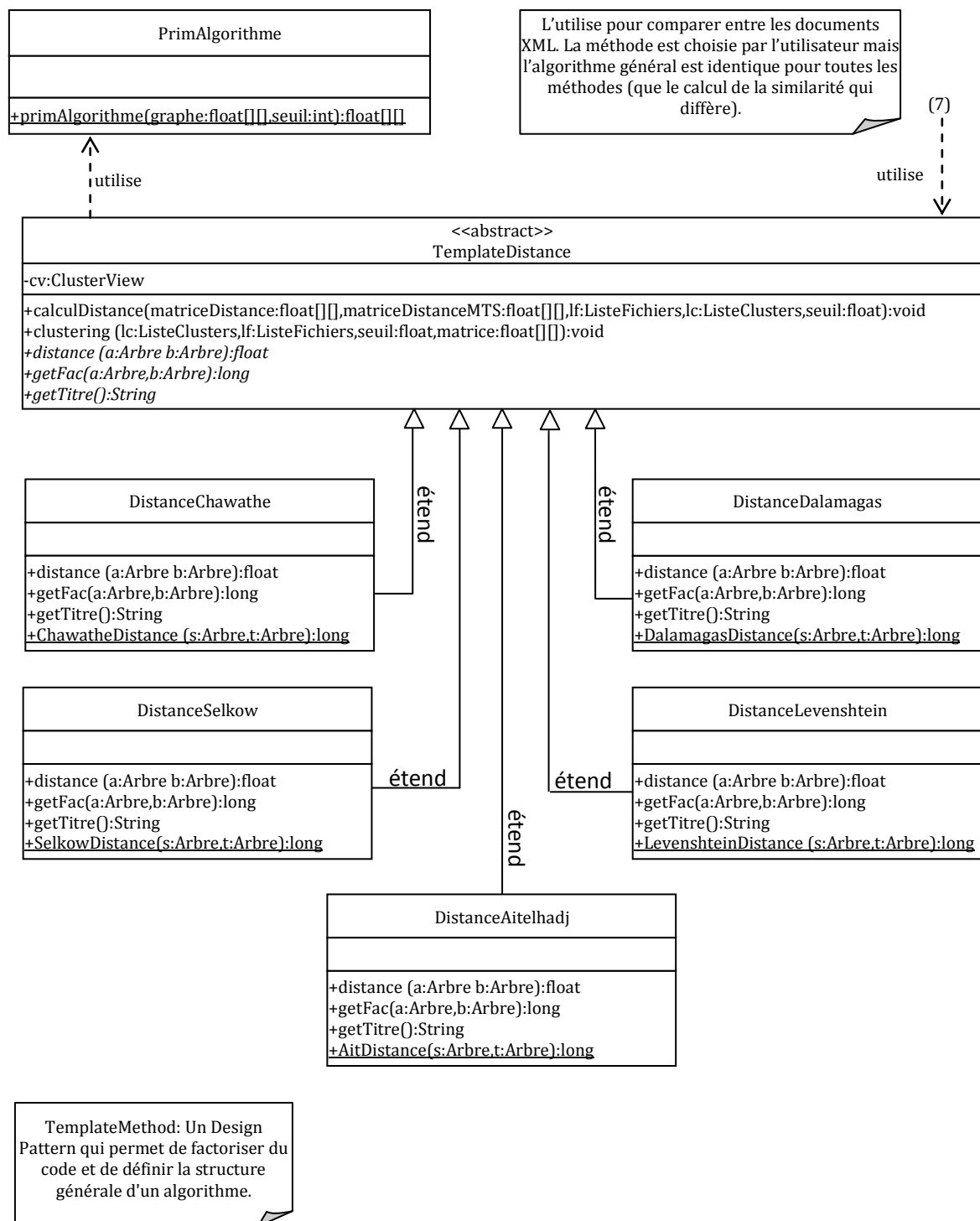
Le diagramme de classes constitue l'un des pivots essentiels de la modélisation avec UML. En effet, ce diagramme permet de donner la représentation statique du système à développer. Cette représentation est centrée sur les concepts de classe et d'association. Chaque classe se décrit par les données et traitements dont elle est responsable pour elle-même et vis-à-vis des autres classes. Les traitements sont matérialisés par des opérations. Le détail des traitements n'est pas représenté directement dans le diagramme de classe ; seul l'algorithme général et le pseudo-code correspondant peuvent être associés à la modélisation.

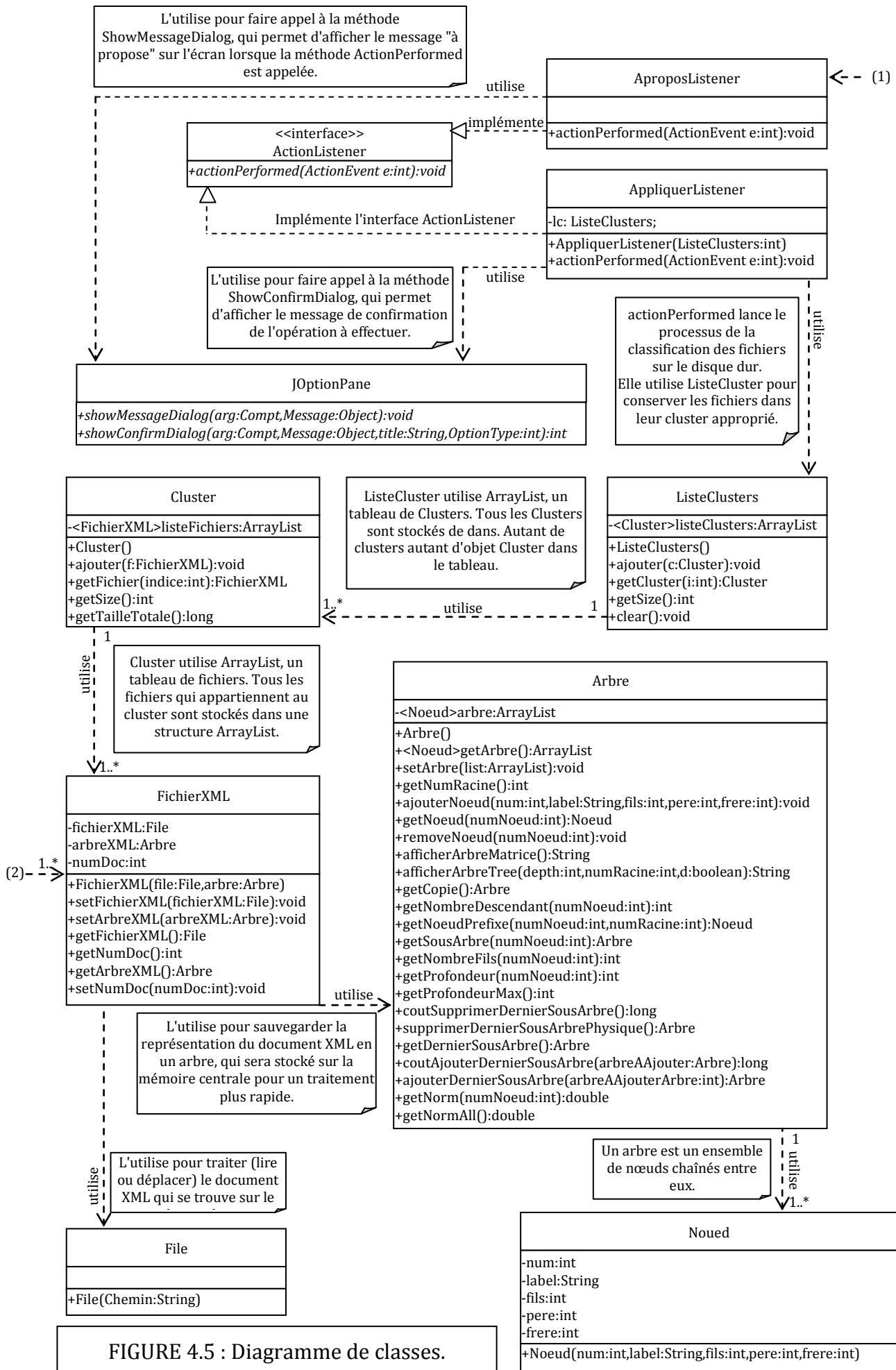
La description du diagramme de classes est fondée sur :

- Le concept d'objet,
- Le concept de classe comprenant les attributs et les opérations,
- Les différents types d'association entre classes.









12.4.2. Diagramme de séquence

L'objectif du diagramme de séquence est de représenter les interactions entre objets en indiquant la chronologie des échanges.

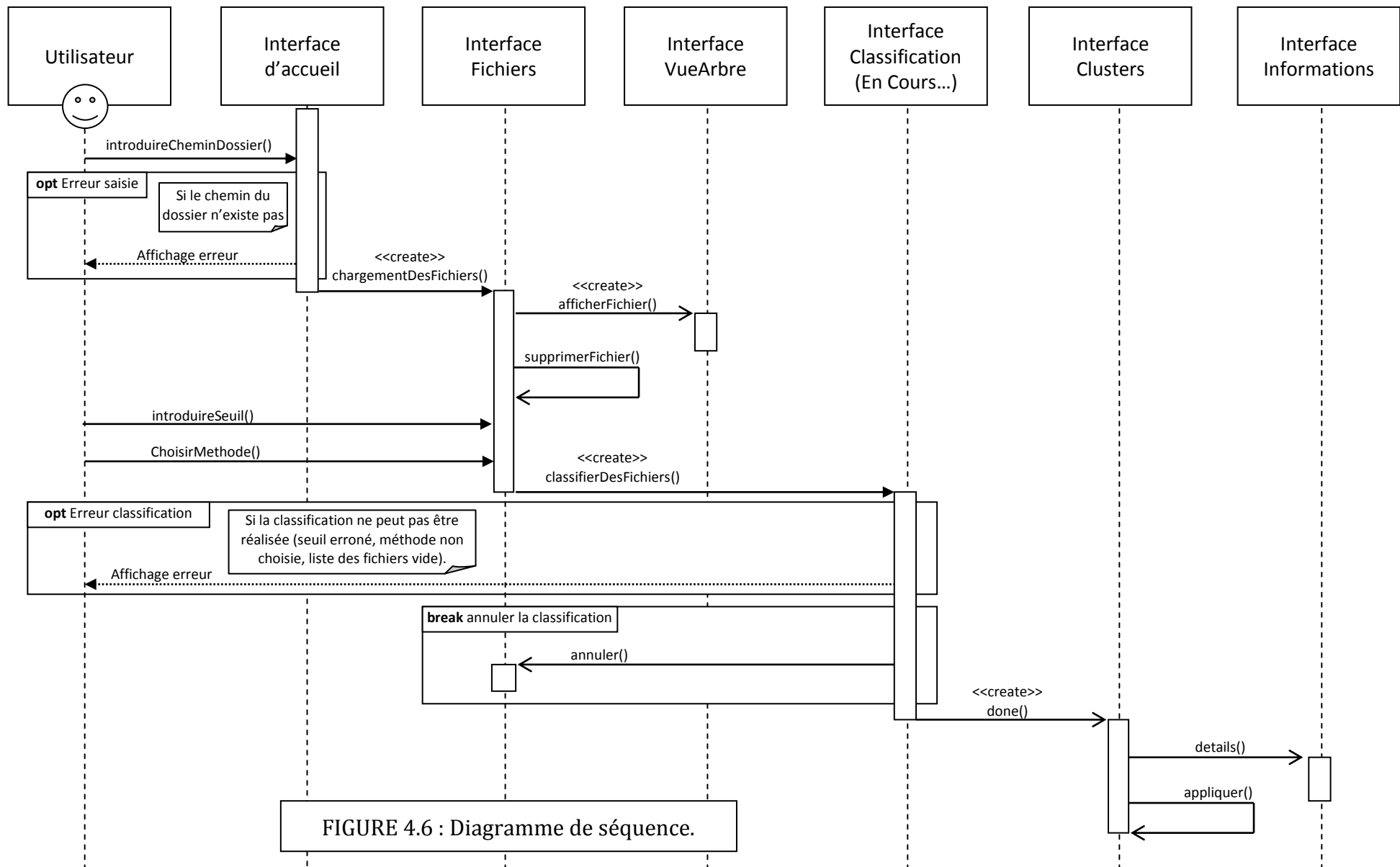


FIGURE 4.6 : Diagramme de séquence.

13. Réalisation

13.1. Outils de développement utilisés

La conception est orientée objet, donc il est nécessaire de choisir un langage orienté objet pour faciliter la réalisation de l'application. Java est un langage orienté objet. Un langage très utilisé actuellement. La compilation d'un code Java génère une application qui peut s'exécuter sur plusieurs systèmes d'exploitation hétérogènes (Unix, Windows, MacOS) sans recompiler le code pour le système d'exploitation cible.

Pour que l'application Java puisse s'exécuter sur un système d'exploitation particulier, il faut juste télécharger et installer une JVM (Java Virtual Machine). JVM est un logiciel, une machine virtuelle, qui permet d'interpréter l'application Java sur le système d'exploitation approprié.

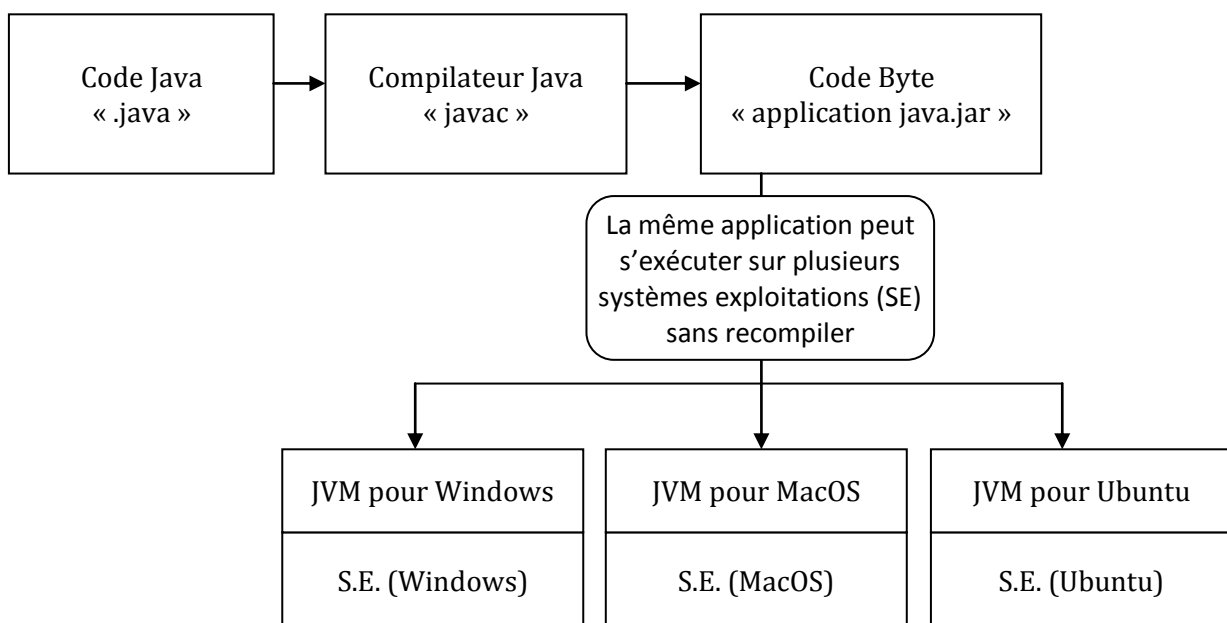


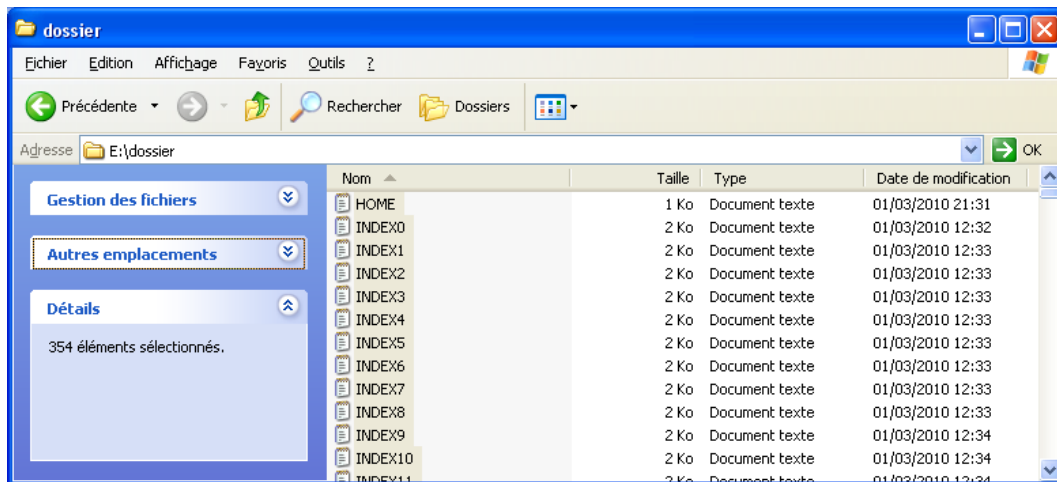
FIGURE 4.7 : Les étapes de compilation et d'interprétation du logiciel

Pour compiler un code Java, il faut un compilateur de code Java. Ce compilateur est téléchargeable gratuitement depuis le site officiel de Java Sun Microsystems. Télécharger le JDK (Java Development Kit), un kit de développement java qui contient plusieurs outils pour réaliser des applications Java. Parmi ces outils il y'a le compilateur (javac comme Java Compiler ou compilateur Java).

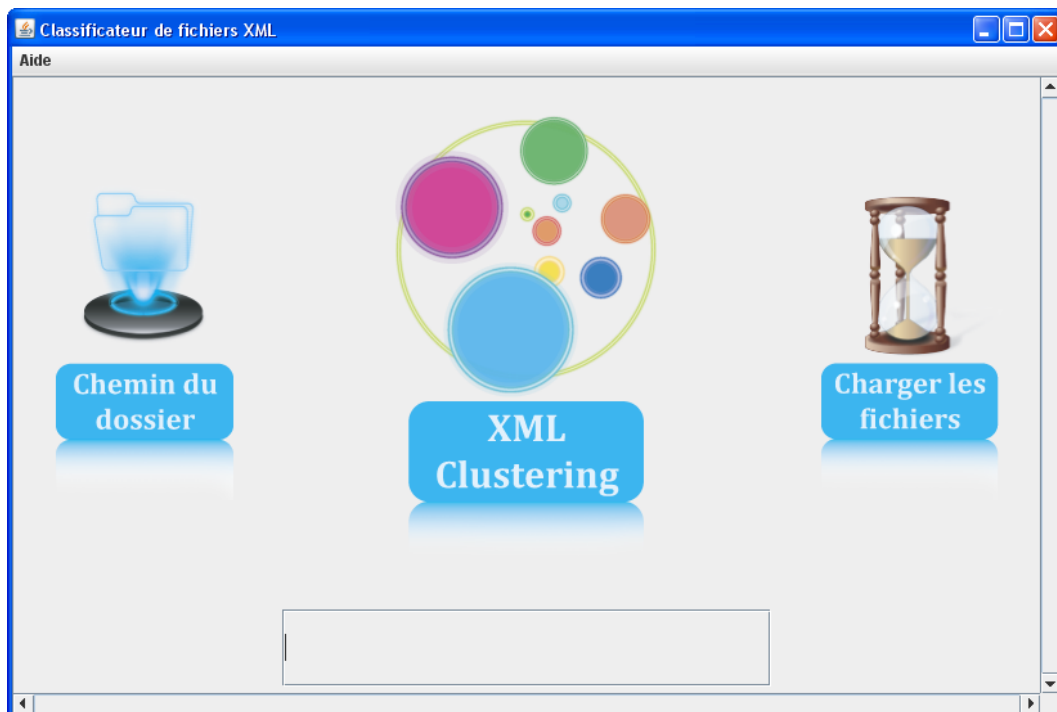
Un code java peut être écrit directement sur un fichier texte avec de simples éditeurs de texte tel que Bloc-Notes (sous Windows). Sauf que l'utilisation d'un éditeur de texte rudimentaire ralentit grandement le développement de l'application. C'est pour ça que l'IDE Eclipse a été utilisé. Eclipse est un logiciel gratuit, c'est un Environnement de Développement Intégré dédié au langage Java. L'outil contient un compilateur intégré et un éditeur de texte avec coloration syntaxique et assistance pour l'écriture du code. Un outil d'une très grande utilité qu'il est vraiment indispensable pour le développement d'applications qui sont relativement grandes telle que la notre.

13.2. Présentation de l'application (Démonstration)

Pour tester le logiciel implémenté nous allons prendre un dossier qui contient un ensemble de documents XML à classer et faire une démonstration. Le chemin du dossier de la Collection XML est : « E:/dossier » il contient 354 documents XML. Pour les classer on lance le logiciel réalisé «ClassificateurXML.jar»



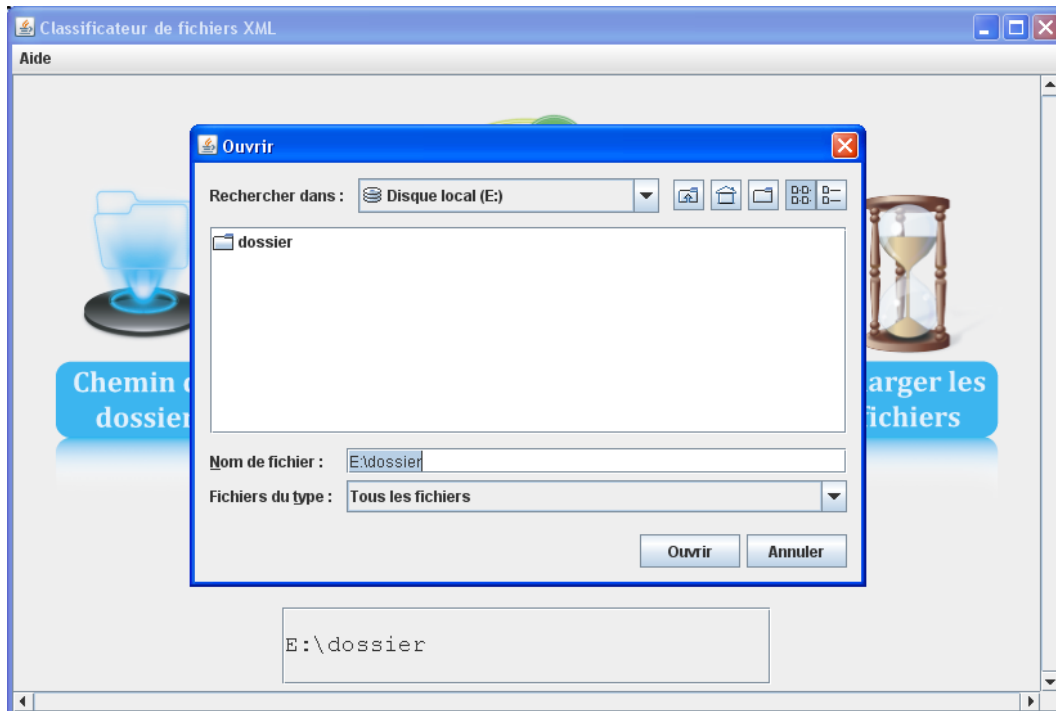
Interface Principale de « ClassificateurXML.jar »



L'interface principale a trois boutons :

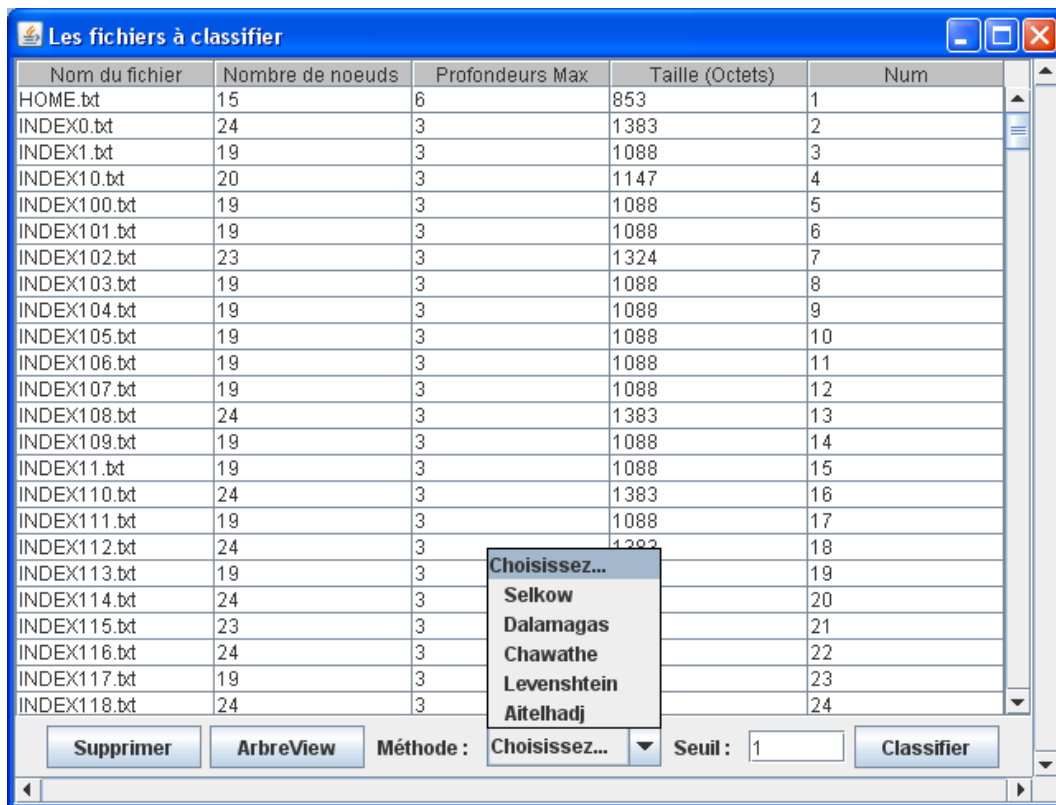
- **Chemin du dossier** : Permet de lancer un explorateur de dossiers.
- **XMLClustering** : Permet d'afficher le « à propos » du logiciel.
- **Charger les fichiers** : Permet de charger depuis le disque vers la mémoire centrale les fichiers à classer.

Pour lancer la classification il faut charger les fichiers depuis le disque dur. Pour dire au logiciel d'où est ce qu'il va récupérer les fichiers à classifier, il faut introduire le chemin du dossier qui contient ces fichiers dans le champ en bas. Si vous ne connaissez pas exactement le chemin complet, vous pouvez vous en servir de l'explorateur de fichiers (Bouton : **Chemin du dossier**) pour récupérer le chemin du dossier.

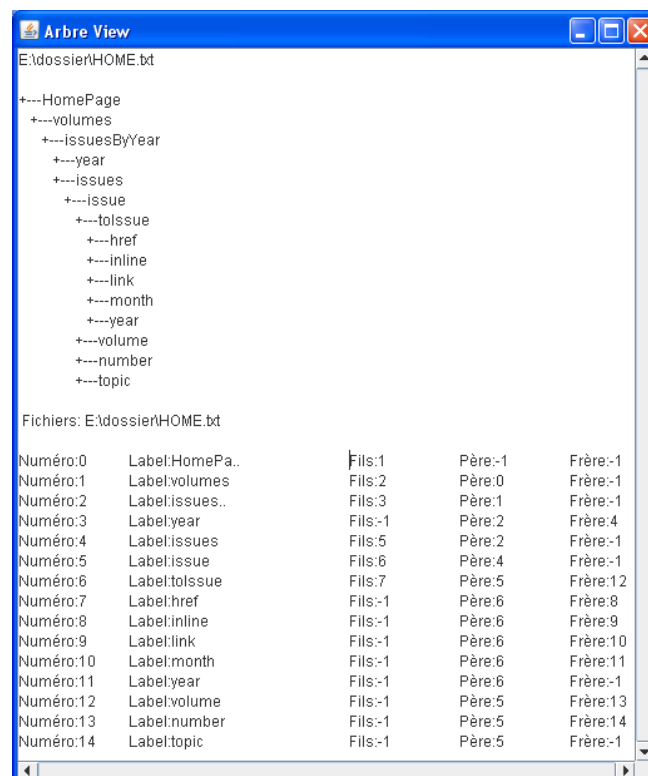


Une fois le chemin du dossier introduit. Il faut charger les fichiers à classifier en cliquant sur le bouton « **Charger les fichiers** ».

Attention, si le chemin est erroné un message d'avertissement sera généré et le processus de chargement s'interrompt. Si le chemin est valide, le logiciel procède au chargement des fichiers et les affichera dans une nouvelle fenêtre (Regardez en bas).

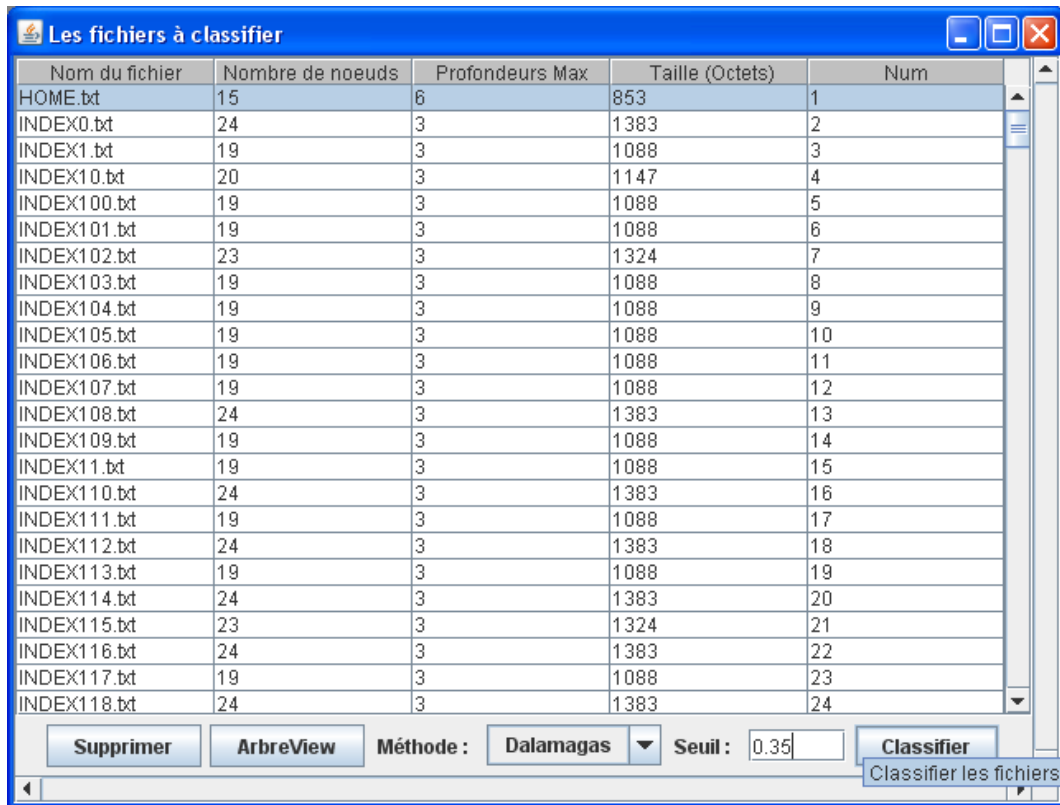


Tous les fichiers chargés seront potentiellement classifiés. Si vous ne voulez pas que quelques fichiers soient classifiés vous n'avez qu'à les sélectionner et les supprimer un par un en cliquant sur le bouton « **Supprimer** ». Vous pouvez aussi visionner le contenu du fichier en le sélectionnant et en cliquant sur le bouton « **ArbreView** ». Exemple pour voir le fichier (HOME.txt) on le sélectionne et on clique sur « **ArbreView** » et on aura le résultat :

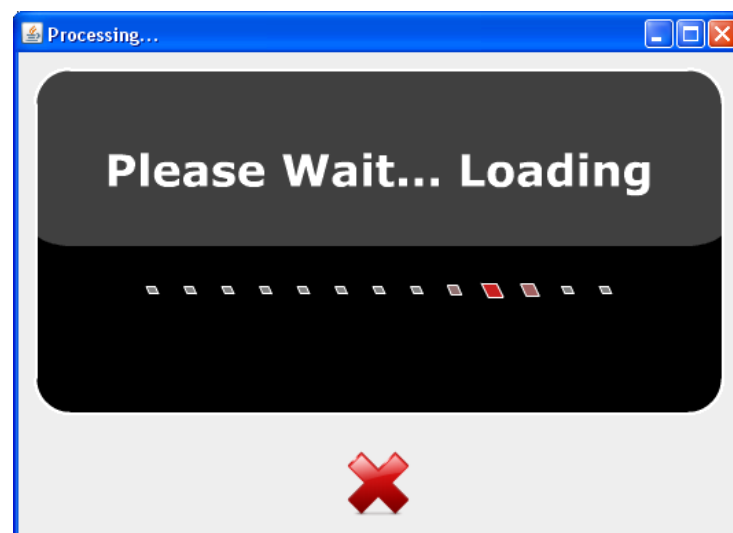


Après avoir chargé les fichiers à classer. On procédera à leur classification. Pour cela on doit choisir une méthode qui sera utilisée pour calculer la similarité structurale entre les fichiers, une des méthodes qu'on a présentées dans le chapitre précédant. Et aussi il faut fixer un seuil (une valeur entre 0 et 1), après il faut cliquer sur le bouton « **Classifier** » pour lancer le processus de la classification.

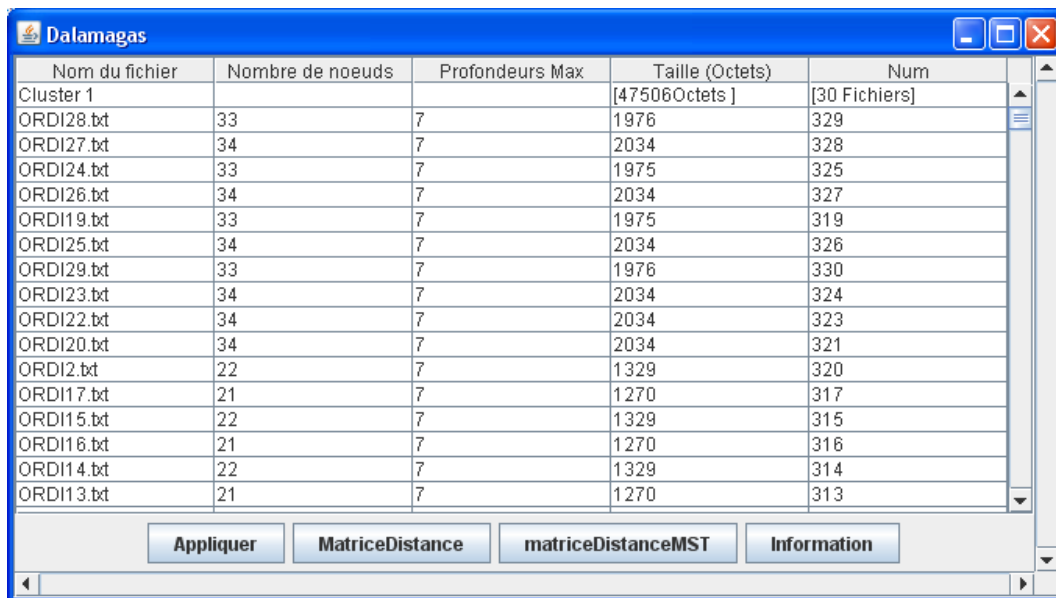
Dans notre démonstration nous allons choisir la méthode de Dalamatag avec un seuil = 0.35



Après le clique sur le bouton « **Classifier** » le processus se lance. La durée du traitement peut être relativement longue en fonction de la méthode choisie et du nombre et de la taille des fichiers à classer. Pour cela une fenêtre de chargement sera générée pour faire patienter l'utilisateur et lui donner la possibilité d'annuler la classification s'il le désire, en cliquant sur le bouton « **X** » en rouge.



Dès que la classification se termine, une nouvelle fenêtre se créera, qui présentera la liste des fichiers dans leurs clusters respectifs. Une fenêtre ayant pour titre le nom de la méthode utilisée.



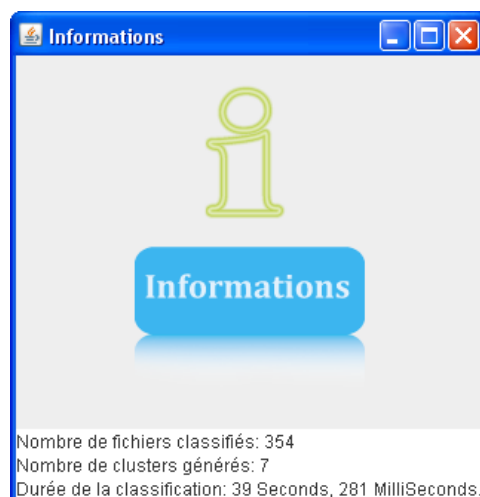
Nom du fichier	Nombre de noeuds	Profondeurs Max	Taille (Octets)	Num
Cluster 1			[47506Octets]	[30 Fichiers]
ORDI28.bt	33	7	1976	329
ORDI27.bt	34	7	2034	328
ORDI24.bt	33	7	1975	325
ORDI26.bt	34	7	2034	327
ORDI19.bt	33	7	1975	319
ORDI25.bt	34	7	2034	326
ORDI29.bt	33	7	1976	330
ORDI23.bt	34	7	2034	324
ORDI22.bt	34	7	2034	323
ORDI20.bt	34	7	2034	321
ORDI2.bt	22	7	1329	320
ORDI17.bt	21	7	1270	317
ORDI15.bt	22	7	1329	315
ORDI16.bt	21	7	1270	316
ORDI14.bt	22	7	1329	314
ORDI13.bt	21	7	1270	313

Buttons: Appliquer, MatriceDistance, matriceDistanceMST, Information

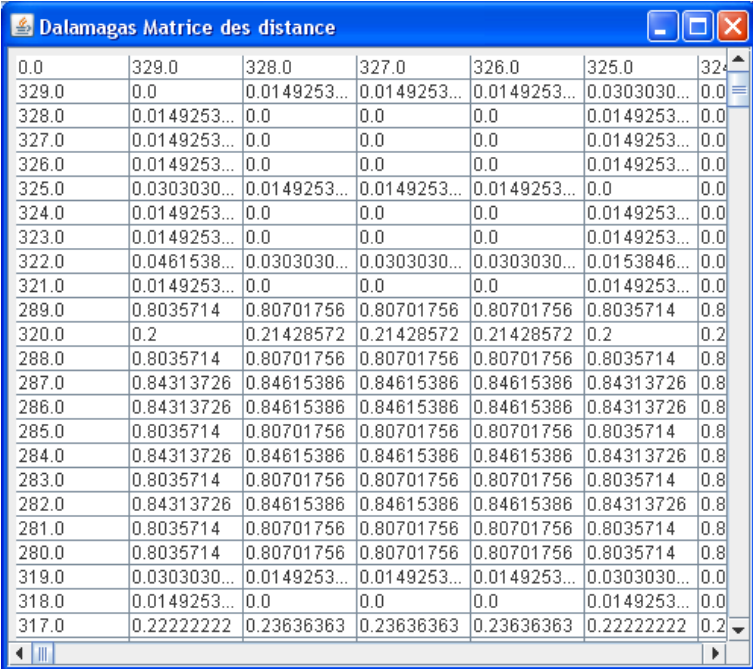
La fenêtre à 4 boutons :

- **Appliquer:** permet de procéder à la classification des fichiers sur le disque dur. Car à la fin du traitement le résultat de la classification est gardé uniquement sur la mémoire centrale, sans déplacer les fichiers dans le disque dur.
- **MatriceDistance:** Permet d'afficher la matrice des distances entre les fichiers XML à classifier.
- **MatriceDistanceMST:** Permet d'afficher la matrice des distances mais avec le chemin le plus court.
- **Information:** permet de générer une fenêtre qui contiendra les informations élémentaires de la classification.

Clique sur **Information** :



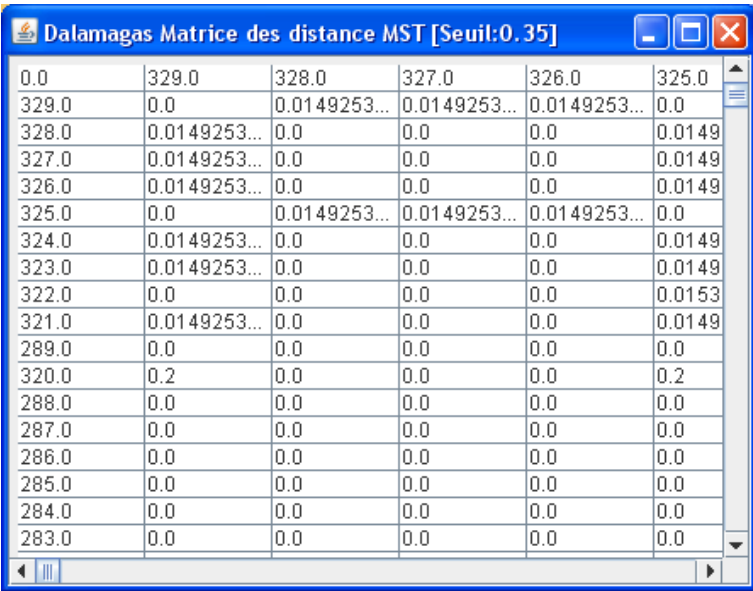
Clique sur le bouton **MatriceDistance** :



0.0	329.0	328.0	327.0	326.0	325.0	324.0
329.0	0.0	0.0149253...	0.0149253...	0.0149253...	0.0303030...	0.0
328.0	0.0149253...	0.0	0.0	0.0	0.0149253...	0.0
327.0	0.0149253...	0.0	0.0	0.0	0.0149253...	0.0
326.0	0.0149253...	0.0	0.0	0.0	0.0149253...	0.0
325.0	0.0303030...	0.0149253...	0.0149253...	0.0149253...	0.0	0.0
324.0	0.0149253...	0.0	0.0	0.0	0.0149253...	0.0
323.0	0.0149253...	0.0	0.0	0.0	0.0149253...	0.0
322.0	0.0461538...	0.0303030...	0.0303030...	0.0303030...	0.0153846...	0.0
321.0	0.0149253...	0.0	0.0	0.0	0.0149253...	0.0
289.0	0.8035714	0.80701756	0.80701756	0.80701756	0.8035714	0.8
320.0	0.2	0.21428572	0.21428572	0.21428572	0.2	0.2
288.0	0.8035714	0.80701756	0.80701756	0.80701756	0.8035714	0.8
287.0	0.84313726	0.84615386	0.84615386	0.84615386	0.84313726	0.8
286.0	0.84313726	0.84615386	0.84615386	0.84615386	0.84313726	0.8
285.0	0.8035714	0.80701756	0.80701756	0.80701756	0.8035714	0.8
284.0	0.84313726	0.84615386	0.84615386	0.84615386	0.84313726	0.8
283.0	0.8035714	0.80701756	0.80701756	0.80701756	0.8035714	0.8
282.0	0.84313726	0.84615386	0.84615386	0.84615386	0.84313726	0.8
281.0	0.8035714	0.80701756	0.80701756	0.80701756	0.8035714	0.8
280.0	0.8035714	0.80701756	0.80701756	0.80701756	0.8035714	0.8
319.0	0.0303030...	0.0149253...	0.0149253...	0.0149253...	0.0303030...	0.0
318.0	0.0149253...	0.0	0.0	0.0	0.0149253...	0.0
317.0	0.22222222	0.23636363	0.23636363	0.23636363	0.22222222	0.2

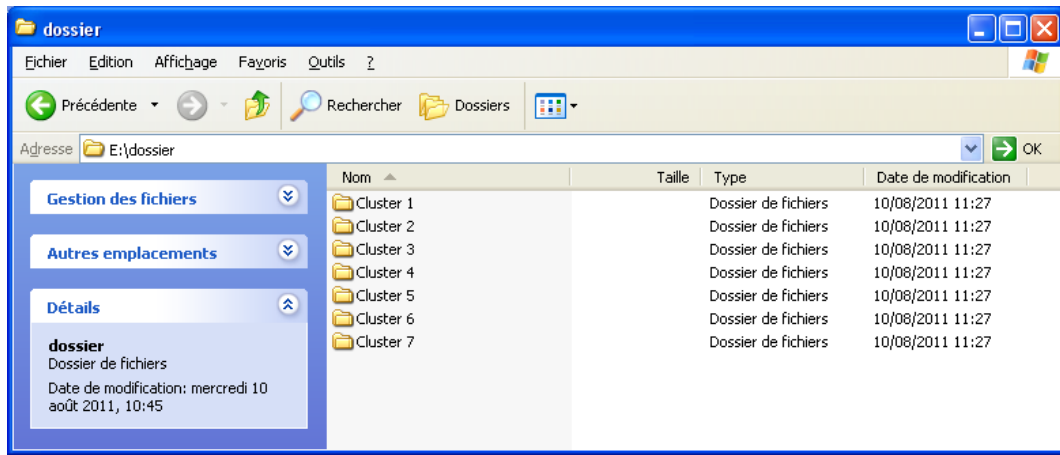
Dans la matrice les fichiers sont représentés par les numéros qui l'auront été attribués.

Clique sur le bouton **MatriceDistanceMST** :



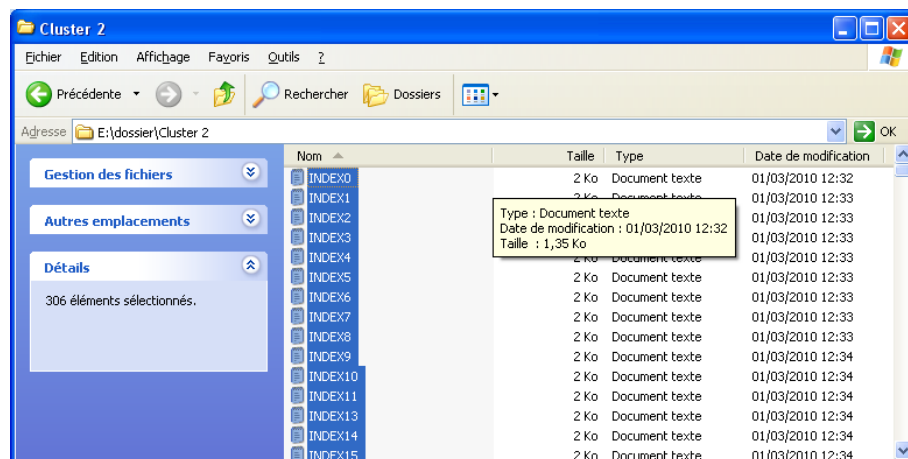
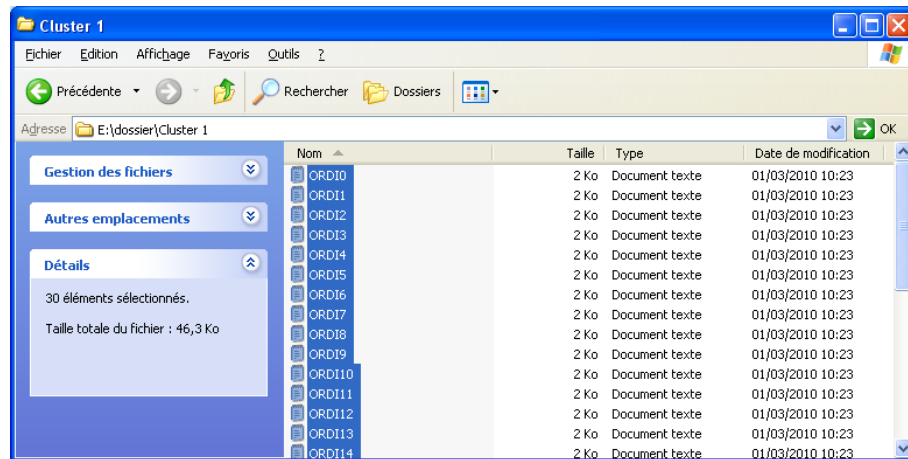
0.0	329.0	328.0	327.0	326.0	325.0	324.0
329.0	0.0	0.0149253...	0.0149253...	0.0149253...	0.0	0.0
328.0	0.0149253...	0.0	0.0	0.0	0.0149	0.0
327.0	0.0149253...	0.0	0.0	0.0	0.0149	0.0
326.0	0.0149253...	0.0	0.0	0.0	0.0149	0.0
325.0	0.0	0.0149253...	0.0149253...	0.0149253...	0.0	0.0
324.0	0.0149253...	0.0	0.0	0.0	0.0149	0.0
323.0	0.0149253...	0.0	0.0	0.0	0.0149	0.0
322.0	0.0	0.0	0.0	0.0	0.0153	0.0
321.0	0.0149253...	0.0	0.0	0.0	0.0149	0.0
289.0	0.0	0.0	0.0	0.0	0.0	0.0
320.0	0.2	0.0	0.0	0.0	0.2	0.0
288.0	0.0	0.0	0.0	0.0	0.0	0.0
287.0	0.0	0.0	0.0	0.0	0.0	0.0
286.0	0.0	0.0	0.0	0.0	0.0	0.0
285.0	0.0	0.0	0.0	0.0	0.0	0.0
284.0	0.0	0.0	0.0	0.0	0.0	0.0
283.0	0.0	0.0	0.0	0.0	0.0	0.0

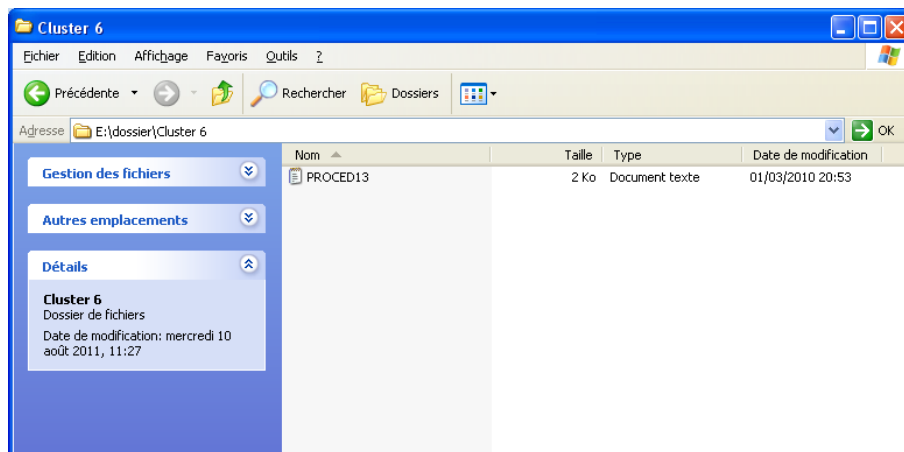
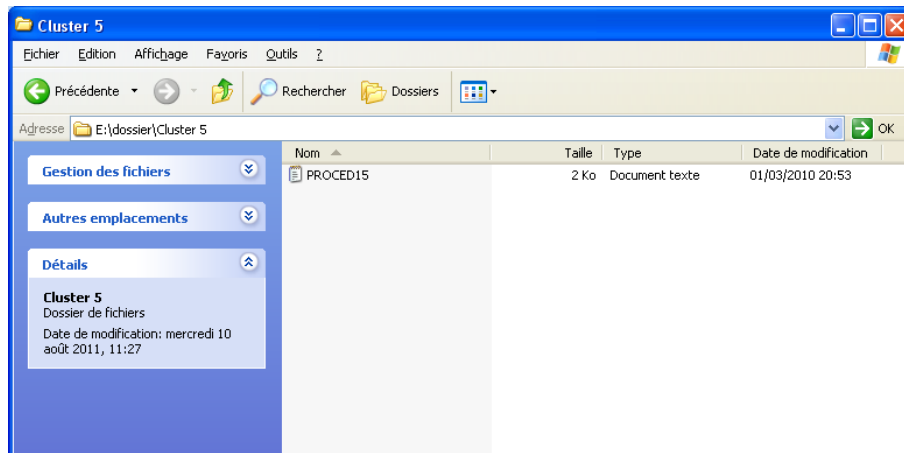
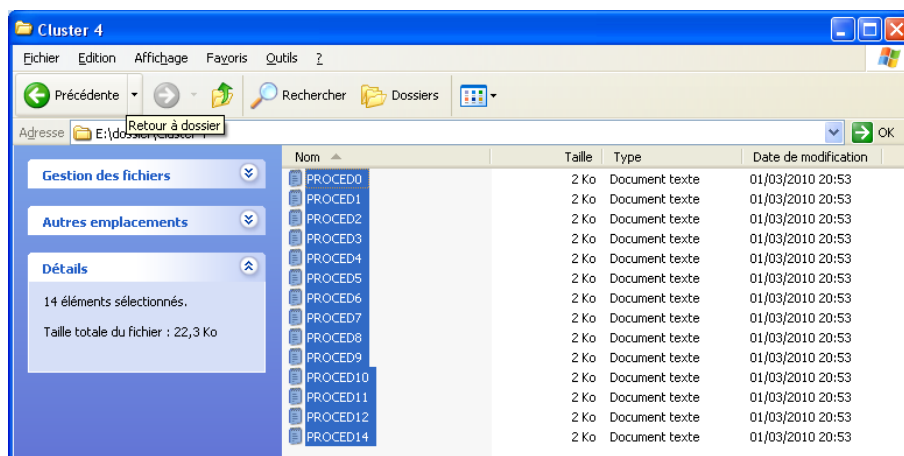
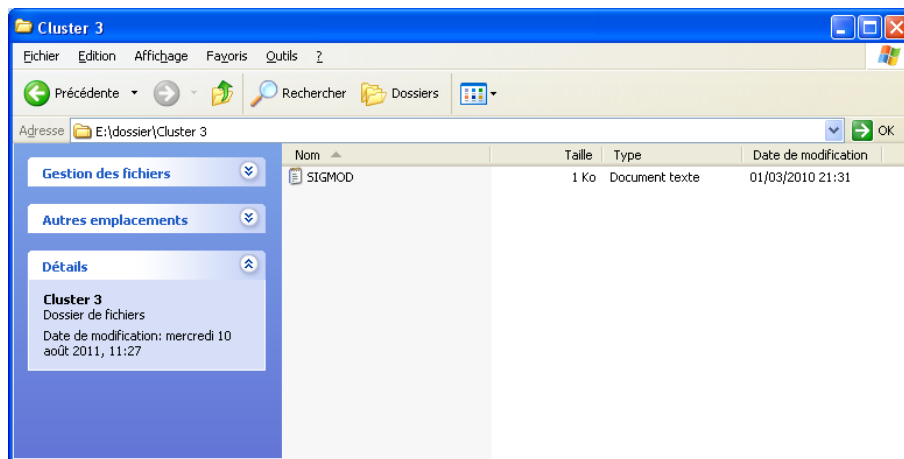
Pour appliquer la classification sur le disque dur, il faut cliquer sur le bouton « **Appliquer** », qui lancera aussi un chargement pour classer les fichiers sur le disque dur. Voici le résultat final de la classification après le clique sur le bouton « **Appliquer** » :

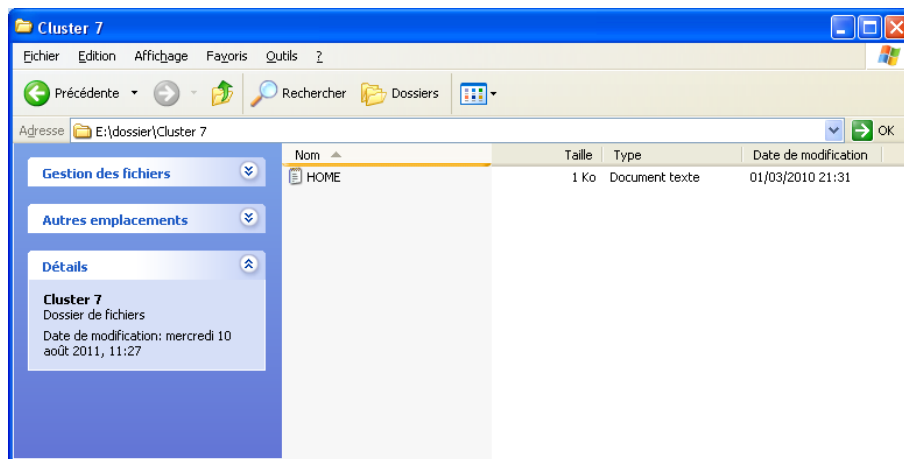


Chaque cluster contiendra des fichiers en fonction de leur ressemblance.

- Cluster1 : 30 fichiers
- Cluster2 : 306 fichiers
- Cluster3 : 1 fichier
- Cluster4 : 14 fichiers
- Cluster5 : 1 fichier
- Cluster6 : 1 fichier
- Cluster7 : 1 fichier







Fin de la démo.

14. Conclusion

Nous avons présenté dans ce chapitre la conception d'un outil qui fait la classification des fichiers contenant la représentation structurelle des documents XML. Nous avons aussi présenté la réalisation du logiciel qui a été programmé avec une démonstration d'utilisation.

Un logiciel qui implémente cinq méthodes différentes de comparaison des documents XML, décrites dans le chapitre précédent. Chacune des méthodes a été expérimentée en lui faisant subir des tests sur des corpus de documents.

Le logiciel sert à regrouper des documents XML qui sont similaires. Une manière efficace d'organiser les documents XML.

Chapitre 5

Expérimentation des méthodes de classification

15. Introduction

Le logiciel du chapitre précédent a été développé dans le but d'expérimenter les cinq méthodes de classifications vues dans le chapitre3. En effet, chaque méthode est caractérisée par une complexité temporelle et par une qualité de classification.

Dans ce chapitre, nous présenterons une expérimentation effectuée sur des échantillons de documents XML. Afin de répondre à la question cruciale, quels sont les meilleurs algorithmes parmi ceux qu'on a testés, qui présentent les meilleurs résultats, vis-à-vis la véracité de la classification et la durée d'exécution ?

16. Présentation

Le résultat de la démonstration présentée dans le chapitre 4 partie réalisation, aurait été le même si on avait utilisé la méthode de Selkow ou la méthode Chawathe car les trois méthodes utilisent la même mesure pour calculer la similarité structurelle entre les fichiers, ils utilisent la distance d'édition. Mais chacune des ces trois méthodes a une complexité temporelle différente des autres. La méthode de Levenshtein et de AitElhadj sont un peu différentes des autres méthodes parce qu'ils utilisent une mesure de similarité calculée différemment. Dans cette partie nous allons expérimenter les 5 méthodes qu'on a implémentées.

Soient A et B deux arbre à comparer, le tableau suivant résume la complexité théorique de chaque méthode :

Méthode	Complexité théorique
Méthode de Selkow	$O(4^{\min(M,N)})$
Méthode de Dalamagas	$O(MN)$
Méthode de Chawathe	$O(MN) + O(M + N)$
Méthode de Levenshtein	$O(MN)$
Méthode de AïtElhadj	$O(\max(M, N)^2 \max(h1, h2)^2)$

Où

- M est le nombre de nœuds de l'arbre A
- N est le nombre de nœuds de l'arbre B
- h1 : la profondeur de l'arbre A
- h2 : la profondeur de l'arbre B

17. Corpus de documents XML

Un corpus est une collection de documents XML dont on connaît à priori leur classification. Pour tester la performance d'un système, on utilise un corpus de documents XML. Pour notre étude, on a choisi un corpus de 147 documents XML. Le corpus contient une collection de fichiers contenant des arbres enracinés ordonnés et étiquetés stockés dans des fichiers textes. Chaque fichier texte représente la structure hiérarchique d'un document XML.

Les caractéristiques (nombre de nœuds et la profondeur des arbres) du corpus sont décrites dans la table suivante :

Table : Caractéristiques des résumés d'arbre de la collection XML à classifier.

Nom de la DTD	Nombre de documents	Nombre de nœuds	Profondeur
Index (IndexTermsPage)	99	[19,24]	3
Ordi(OrdinaryIssuePage)	30	[22,35]	7
Prodic(ProceedingPage)	16	[25,29]	5
Sigmod(Sigmod)	1	12	7
Home(HomePage)	1	15	6
	147 Documents		

Le test va s'effectuer sur une collection de 147 documents XML correspondant à 5 DTDs, nommées HomePage, IndexTermsPage, OrdinaryIssuePage, ProceedingsPage et SigmodRecord. Ces DTD peuvent en réalité être considérées comme les clusters cible de classification. Chaque fichier qui vérifie une DTD doit être classé dans le cluster représentant cette DTD.

Voici la classification cible, le résultat prévu qui devrait être retourné par les méthodes de classification:

Nombre de clusters	C1	C2	C3	C4	C5
5	99	30	16	1	1

18. Expérimentation des complexités théoriques

Pour expérimenter la durée d'exécution. Les cinq méthodes ont été testées sur ce corpus de 147 Documents XML. Les résultats sont comme suit :

Méthode	seuil	Nombre de fichiers classifiés	Nombre de clusters générés	Durée d'exécution
Selkow	1	2	1	Durée de la classification: 1289 Seconds, 719 MilliSeconds (21 mins pour deux fichiers) (147 fichiers : $\sum_{x=1}^{147} x * 21 \text{ min} = 10878 * 21 = 228438 \text{min}$ = 3807,3 heurs = 158,6375 jours) Donc il est carrément impossible de simuler toutes la collection XML avec une machine ayant de telles capacités, et espérer avoir un résultat dans les temps. Donc on va se contenter de tester les autres quartes méthodes seulement.
Dalamagas	1	147	1	Durée de la classification: 8 Seconds, 125 MilliSeconds.
Chawathe	1	147	1	Durée de la classification: 84 Seconds, 94 MilliSeconds.
Leveinstein	1	147	1	Durée de la classification: 27 Seconds, 672 MilliSeconds.
AitElhadj	1	147	1	Durée de la classification: 154 Seconds, 234 MilliSeconds

Caractéristiques de la machine

Le test a été effectué sur une machine ayant la configuration suivante :

Système d'exploitation :	Windows XP SP3.
Processeur :	Pentium(R) Dual-Core CPU E5300 2.60GHz (2CPU)
Mémoire :	2048MO RAM
Disque dur :	500GO SATA
Carte mère :	MSI P31 Neo V2
Carte graphique :	NVIDIA GeForce 9400GT (1024Mb)

Pour tester la complexité temporelle de chaque méthode, nous allons prendre des corpus à tailles variables (2, 5, 10, 20, 50, 100, 200, 500 documents) et appliquer les 5 méthodes aux corpus et calculer le temps d'exécution pour chacune d'elle.

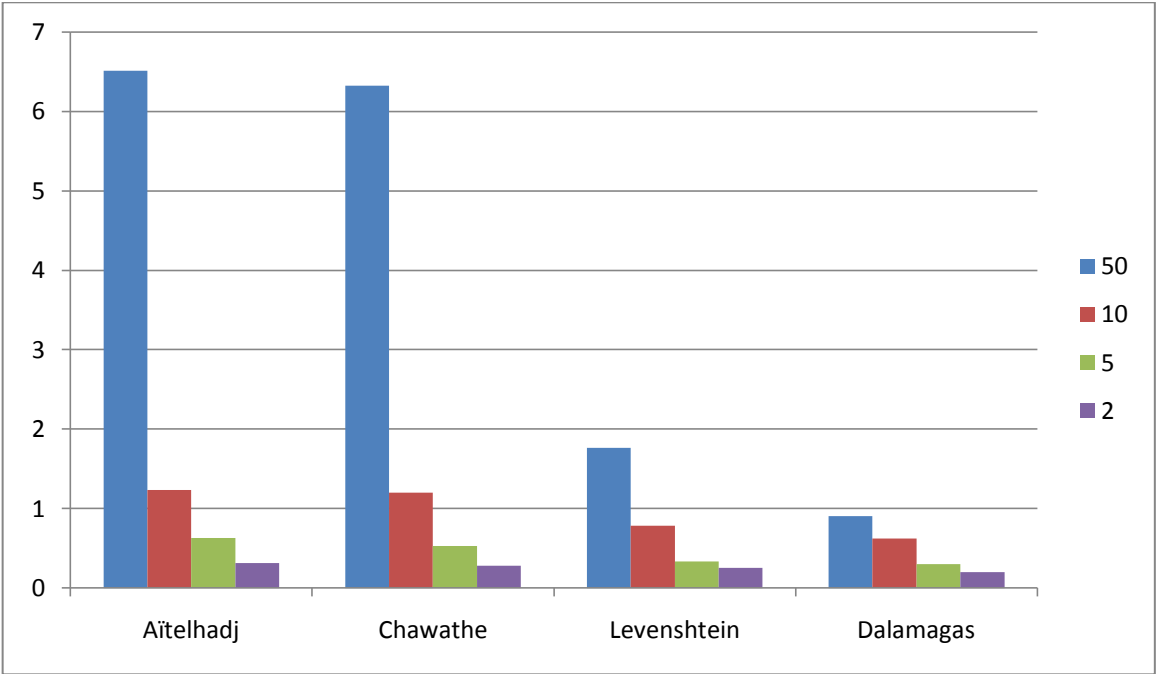
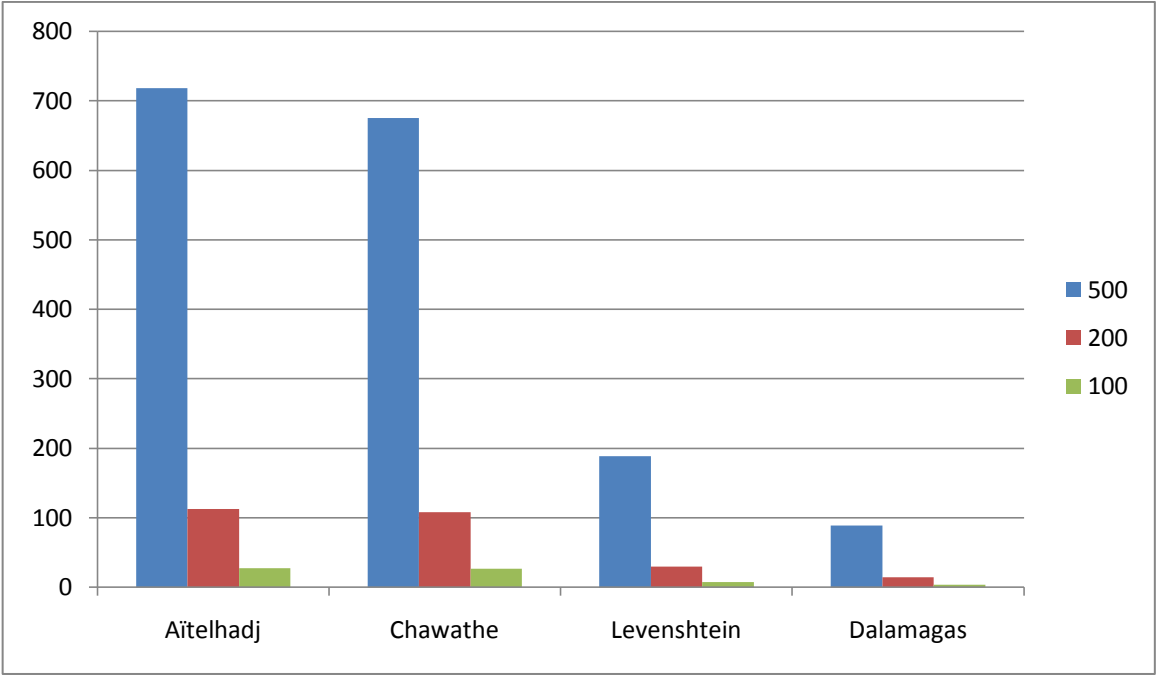
Vue la complexité élevée de l'algorithme de Selkow, nous allons nous contenter de comparer seulement le premier corpus (2 documents) avec cette méthodes.

Méthode Corpus	Selkow (Secondes)	Aïtelhadj (Secondes)	Chawathe (Secondes)	Levenshtein (Secondes)	Dalamagas (Secondes)
500	-	718,469	674,937	188,625	88,891
200	-	112,375	108,219	29,938	14,219
100	-	27,312	26,219	7,203	3,407
50	-	6,516	6,328	1,765	0,906
10	-	1,235	1,203	0,78	0,62
5	-	0,63	0,53	0,33	0,30
2	1435,500	0,31	0,28	0,25	0,20

Nous allons représenter le tableau avec un diagramme à bâtons afin de mieux visualiser et comparer les temps d'exécution pris par chaque méthode. D'après le tableau, il est clair que la méthode de Selkow a une complexité temporelle très élevée. Elle est de l'ordre des minutes pour deux fichiers, alors que, pour deux fichiers, elle est de l'ordre des millièmes de secondes pour les autres méthodes. Pour avoir une meilleure visualisation du diagramme nous allons omettre la méthode de Selkow et nous focaliser sur les quatre autres méthodes.

D'après les résultats des durées d'exécution. Les complexités temporelles sont comme suit, par ordre descendant (celle qui la plus lente à la moins lente): Méthode de Selkow, Méthode de Aïtelhadj, Méthode de Chawathe, Méthode de Levenshtein et la Méthode de Dalamagas. Le diagramme suivant résume ces résultats.

Diagrammes de comparaison des temps d'exécution (Dalamagas, Chawathe Levenshtein Aïtelhadj)



19. Evaluation des méthodes

Le seuil, est ce qu'on appelle la distance structurelle (λ) entre les fichiers, si deux fichiers ont une structure identiques la distance est nulle, si les structures des fichiers n'ont aucune ressemblance la distance est égale à 1. Pour tester les performances et la qualité de classification de chaque méthode, nous allons varier la valeur (λ) de 0 à 1 et appliquer la classification avec chacune des méthodes.

Si on veut détecter les clusters avec des documents trop similaires, la valeur λ doit être fixée à une petite valeur proche de zéro. Inversement, si on affecte une grande valeur à λ (au seuil), nous obtiendront des clusters hétérogènes.

Pour les expérimentés le principe est simple, nous allons exécuter la classification avec chacune des méthodes en fixant le seuil (distance $\in [0,1]$) à des valeurs comprise entre 0.1 et 0.9 avec des écarts de 0.1 entre chaque deux valeurs, et analysons le résultat de la classification. Nous allons calculer pour chaque étape (0.1) les valeurs de la micro-mean et de la macro-mean : pureté, entropie et le F-Mesure vues dans le chapitre 2.

La méthode de Selkow, de Chawathe et de Dalamagas donnent exactement le même résultat de classification parce qu'elles utilisent la distance d'édition qui est une valeur unique. En revanche la méthode de Levenshtein et la méthode de Aïtelhadj donnent des résultats différents donc on a trois cas d'expérimentation :

5.2. Cas1 : Test des Méthodes de Selkow, Dalamagas et Chawathe

Seuil	Nombre de clusters	C1	C2	C3	C4	C5	C6	C7	C8
$\lambda = 0.9$	1	147							
$\lambda = 0.8$	1	147							
$\lambda = 0.7$	3	99	47	1					
$\lambda = 0.6$	6	99	31	14	1	1	1		
$\lambda = 0.5$	6	99	31	14	1	1	1		
$\lambda = 0.4$	7	99	30	14	1	1	1	1	
$\lambda = 0.3$	7	99	30	14	1	1	1	1	
$\lambda = 0.2$	7	99	30	14	1	1	1	1	
$\lambda = 0.1$	8	99	18	14	12	1	1	1	1

Nombre de clusters avec le nombre de document s dans les clusters pour chaque seuil (λ)

Seuil	Macro-mean			Micro-mean		
	Rappel	Précision	F-mesure	Rappel	Précision	F-mesure
$\lambda = 0.9$	1.0000	0.6734	0.8084	1.0000	0.6734	0.8084
$\lambda = 0.8$	1.0000	0.6734	0.8084	1.0000	0.6734	0.8084
$\lambda = 0.7$	1.0000	0.8794	0.9263	1.0000	0.8843	0.9293
$\lambda = 0.6$	0.6666	0.9946	0.9623	0.9753	0.9931	0.9781
$\lambda = 0.5$	0.6666	0.9946	0.9623	0.9753	0.9931	0.9781
$\lambda = 0.4$	0.7142	1.0000	0.7386	0.9753	1.0000	0.9816
$\lambda = 0.3$	0.7142	1.0000	0.7386	0.9753	1.0000	0.9816
$\lambda = 0.2$	0.7142	1.0000	0.7386	0.9753	1.0000	0.9816
$\lambda = 0.1$	0.6250	1.0000	0.6865	0.8773	1.0000	0.9160

Qualité de la classification pour chaque seuil (λ)

5.3. Cas2 : Test de la méthode de Levenshtein

Seuil	Nombre de clusters	C1	C2	C3	C4	C5	C6	C7	C8
$\lambda = 0.9$	1	147							
$\lambda = 0.8$	1	147							
$\lambda = 0.7$	1	147							
$\lambda = 0.6$	1	147							
$\lambda = 0.5$	2	146	1						
$\lambda = 0.4$	4	99	46	1	1				
$\lambda = 0.3$	4	99	46	1	1				
$\lambda = 0.2$	5	99	30	16	1	1			
$\lambda = 0.1$	6	99	18	16	12	1	1		

Nombre de clusters avec le nombre de documents dans les clusters pour chaque seuil (λ)

Seuil	Macro-mean			Micro-mean		
	Rappel	Précision	F-mesure	Rappel	Précision	F-mesure
$\lambda = 0.9$	1.0000	0.6734	0.8084	1.0000	0.6734	0.8084
$\lambda = 0.8$	1.0000	0.6734	0.8084	1.0000	0.6734	0.8084
$\lambda = 0.7$	1.0000	0.6734	0.8084	1.0000	0.6734	0.8084
$\lambda = 0.6$	1.0000	0.6734	0.8084	1.0000	0.6734	0.8084
$\lambda = 0.5$	1.0000	0.8390	0.9040	1.0000	0.8601	0.8094
$\lambda = 0.4$	1.0000	0.9130	0.9473	1.0000	0.8911	0.9340
$\lambda = 0.3$	1.0000	0.9130	0.9473	1.0000	0.8911	0.9340
$\lambda = 0.2$	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
$\lambda = 0.1$	0.8333	1.0000	0.8869	0.9020	1.0000	0.9344

Qualité de la classification pour chaque seuil (λ)

5.4. Cas3 : Test de la méthode de Aitelhadj

Seuil	Nombre de clusters	C1	C2	C3	C4	C5	C6	C7	C8	C9
$\lambda = 0.9$	1	147								
$\lambda = 0.8$	3	99	47	1						
$\lambda = 0.7$	4	99	46	1	1					
$\lambda = 0.6$	5	99	30	16	1	1				
$\lambda = 0.5$	5	99	30	16	1	1				
$\lambda = 0.4$	5	99	30	16	1	1				
$\lambda = 0.3$	5	99	30	16	1	1				
$\lambda = 0.2$	17	99	18	16	1	1	1	1	1	1
$\lambda = 0.1$	18	59	40	18	16	1	1	1	1	1

Seuil	C10	C11	C12	C13	C14	C15	C16	C17	C18
$\lambda = 0.9$									
$\lambda = 0.8$									
$\lambda = 0.7$									
$\lambda = 0.6$									
$\lambda = 0.5$									
$\lambda = 0.4$									
$\lambda = 0.3$									
$\lambda = 0.2$	1	1	1	1	1	1	1	1	
$\lambda = 0.1$	1	1	1	1	1	1	1	1	1

Nombre de clusters avec le nombre de documents dans les clusters pour chaque seuil (λ)

Seuil	Macro-mean			Micro-mean		
	Rappel	Précision	F-mesure	Rappel	Précision	F-mesure
$\lambda = 0.9$	1.0000	0.6734	0.8084	1.0000	0.6734	0.8084
$\lambda = 0.8$	1.0000	0.8794	0.9263	1.0000	0.8843	0.9293
$\lambda = 0.7$	1.0000	0.9130	0.9473	1.0000	0.8911	0.9340
$\lambda = 0.6$	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
$\lambda = 0.5$	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
$\lambda = 0.4$	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
$\lambda = 0.3$	1.0000	1.0000	1.0000	1.0000	1.0000	1.0000
$\lambda = 0.2$	0.2940	1.0000	0.3248	0.8721	1.0000	0.8930
$\lambda = 0.1$	0.2777	1.0000	0.3247	0.5477	1.0000	0.6758

Qualité de la classification pour chaque seuil (λ)

Diagramme de Comparaison des Rappels (cas1 cas2 cas3) : (Macro-Mean)

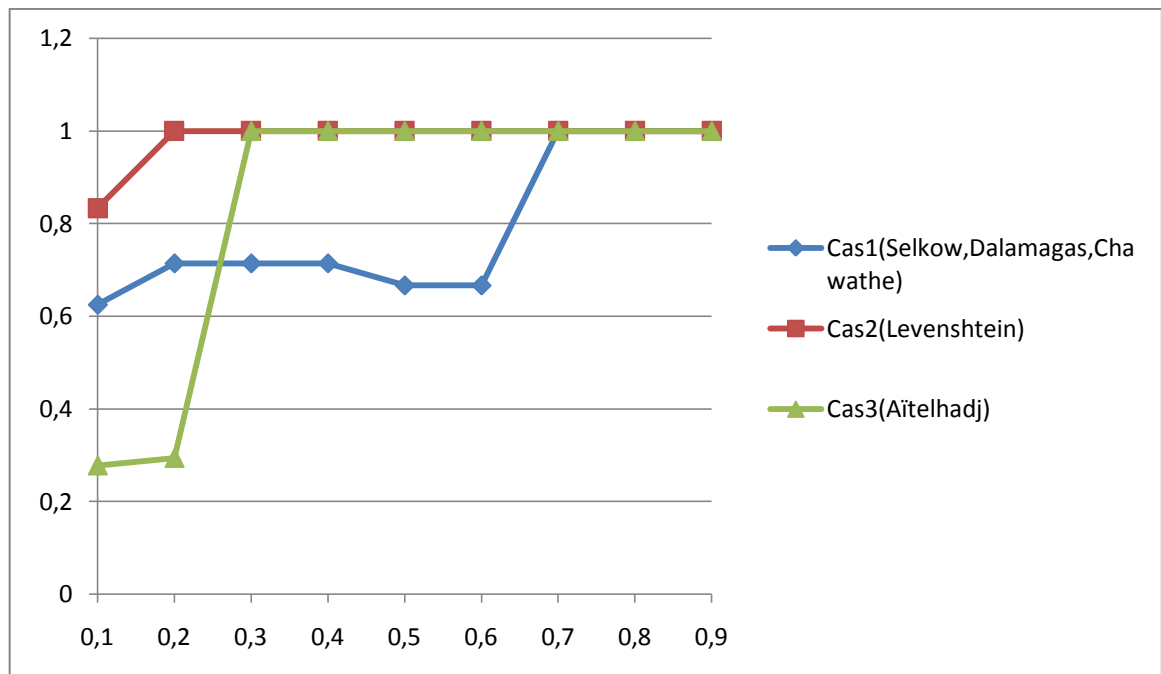


Diagramme de Comparaison des Rappels (cas1 cas2 cas3) : (Micro-Mean)

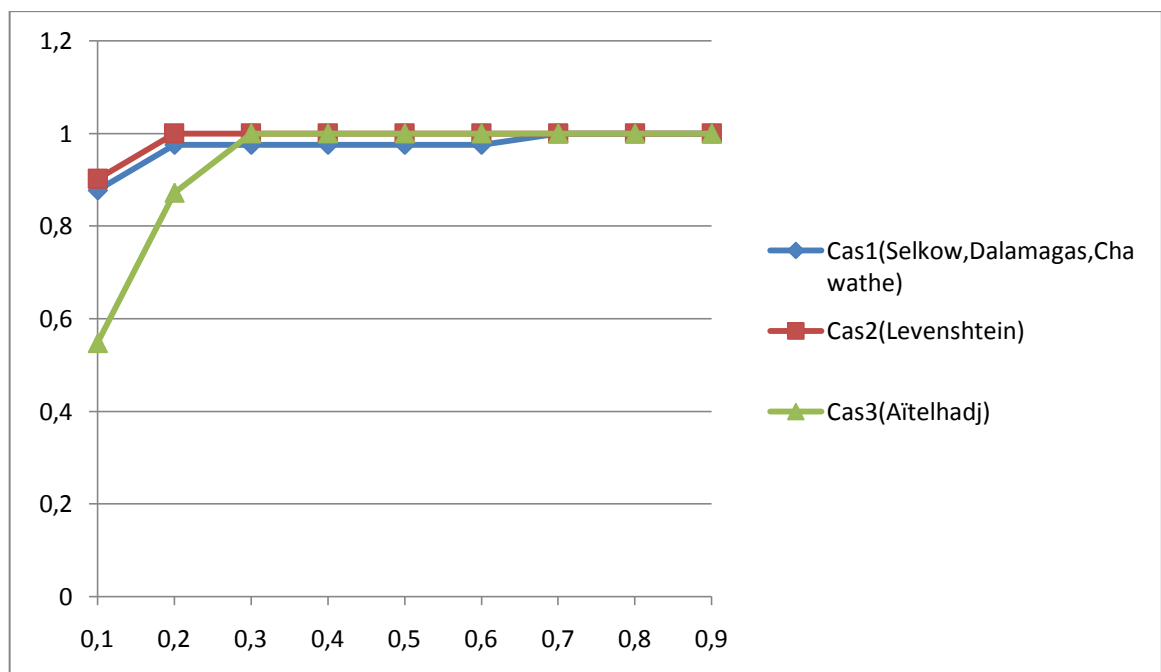


Diagramme de Comparaison des Précisions (cas1 cas2 cas3) (MacroMean)

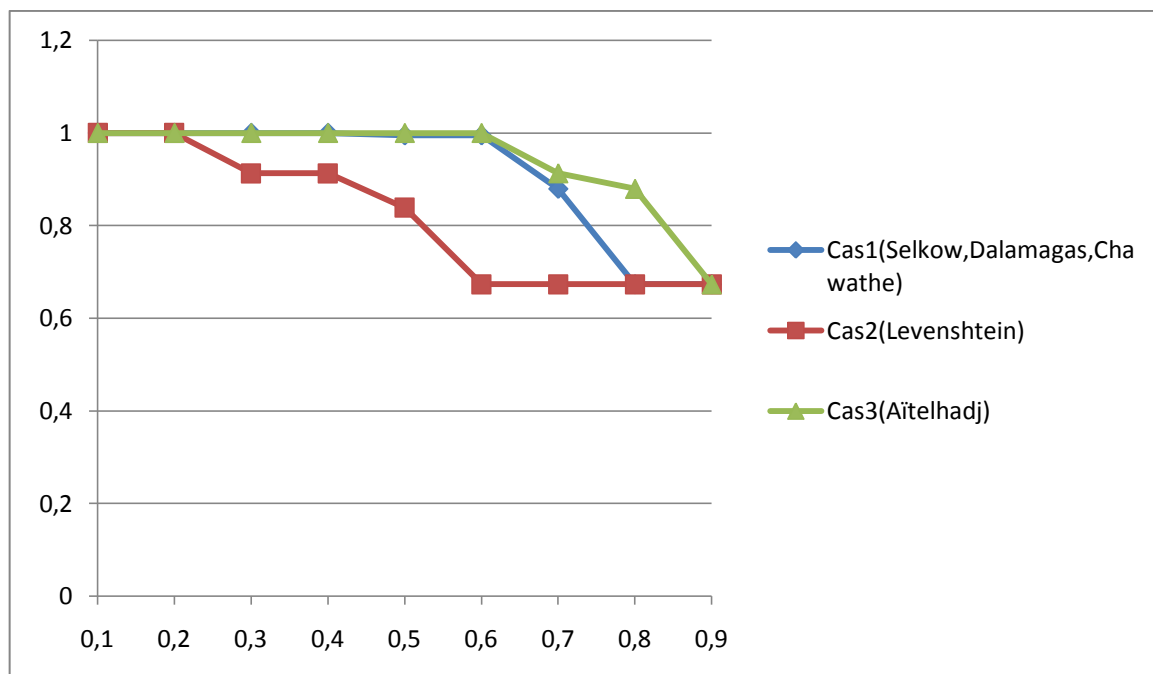


Diagramme de Comparaison des Précisions (cas1 cas2 cas3) (MicroMean)

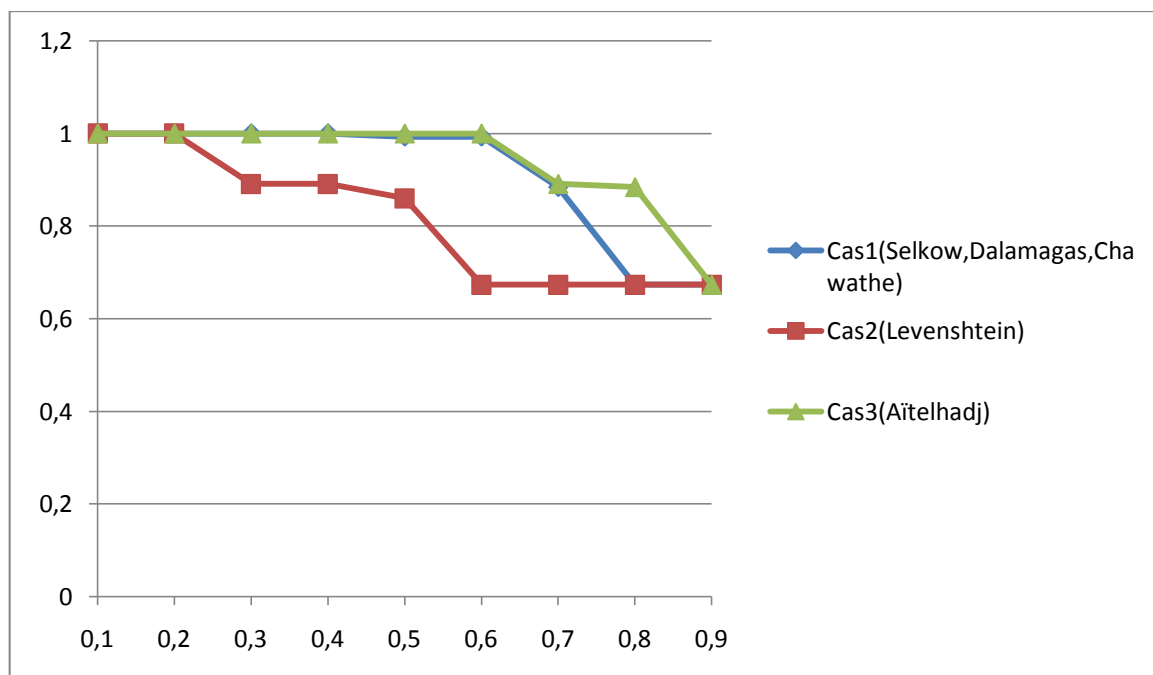


Diagramme de Comparaison des F-mesures (cas1 cas2 cas3) (Macro-Mean)

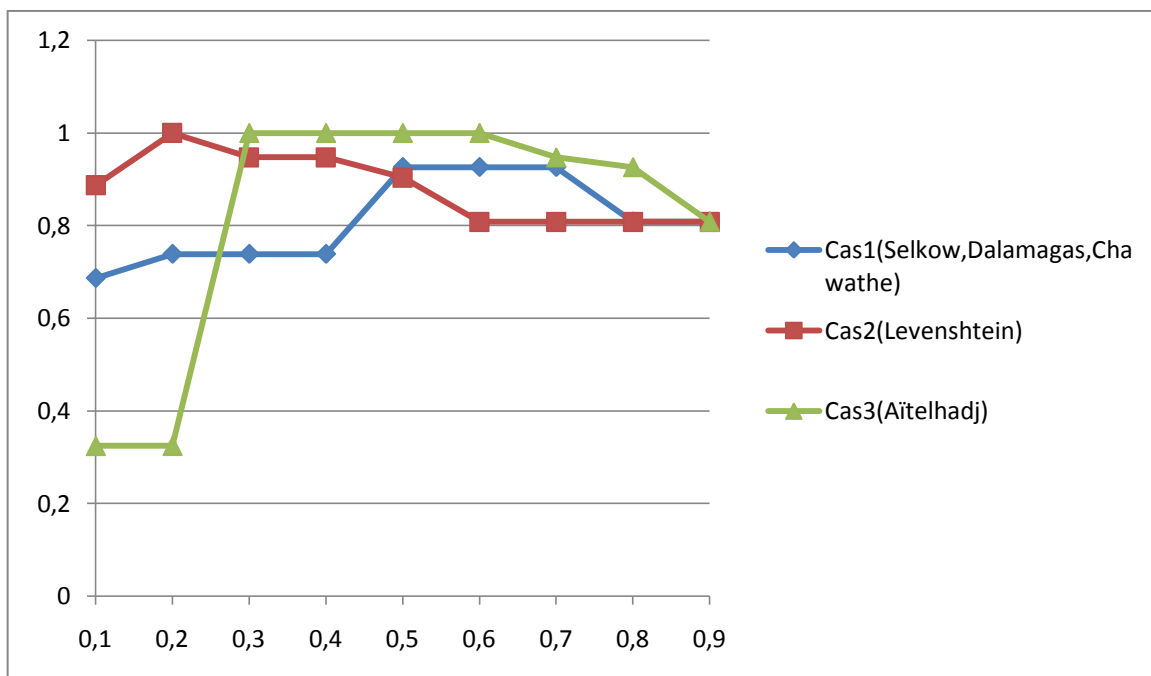
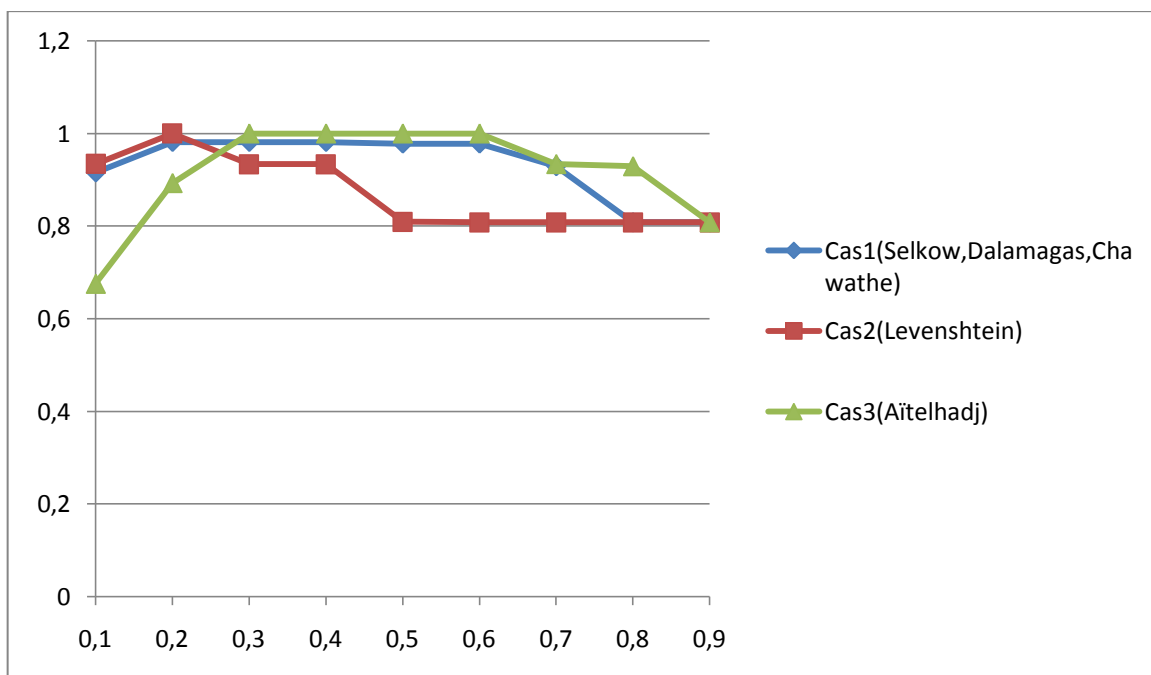


Diagramme de Comparaison des F-mesures (cas1 cas2 cas3) (MicroMean)



6. Conclusion

Nous avons présenté dans le chapitre précédent la conception et la réalisation d'un logiciel qui fait la classification des fichiers contenant la représentation structurelle de documents XML. Un logiciel qui implémente cinq méthodes différentes de comparaison des documents XML, décrites dans le chapitre 3. Chacune des méthodes a été expérimentée dans ce chapitre en lui faisant subir des tests sur des corpus de documents.

Chaque test a présenté une durée d'exécution et une qualité de classification différentes les uns des autres. L'algorithme qui a présenté le meilleur temps d'exécution est l'algorithme de Dalamagas. Et l'algorithme qui a présenté la longue durée d'exécution est l'algorithme de Aïtelhadj (sans compter celui de Selkow). En revanche l'Algorithme qui a retourné la meilleure qualité de classification est l'algorithme de Aïtelhadj. Donc pour une meilleure classification, il est conseillé d'utiliser l'algorithme de Aïtelhadj comme méthode de comparaison malgré que sa complexité temporelle est un peu élevée lorsque la profondeur des arbres est grande.

Les Algorithmes de Selkow, Chawathe et Dalamagas retournent le même résultat de classification donc ils ont la même qualité de classification qui reste la moins juste parmi les méthodes testées. En ce qui concerne l'Algorithme de Levenshtein, il a présenté une bonne durée d'exécution ; la meilleure durée après celle de l'algorithme de Dalamagas. Et il a retourné une véracité de classification meilleure que celle retournée par les algorithmes de Dalamagas, Chawathe et Selkow.

Conclusion générale

A travers ce mémoire, j'ai pu mettre en œuvre les concepts et les techniques vus durant le deux cursus universitaire, et appliquer avec profit les idées et les techniques de conception et de réalisation des logiciels orientés objet. Un mémoire qui m'a appris à bien concevoir des logiciels, et à mieux mener un travail de recherche.

Dans ce mémoire, un logiciel a été développé. Il permet de classifier les documents XML. Le but du logiciel est de regrouper des documents XML en fonction de leur similarité structurelle. En effet, le document XML a une structure arborescente qui le décrit, cette structure constitue un facteur efficace et fiable pour la comparaison et la classification des documents XML. L'extraction de cette structure et faite par un module externe, un module qui analyse les documents XML et extrait leur structure qui sera passée à notre logiciel pour la comparée avec les autres structures représentant les autres documents afin de déterminer les valeurs de similitude entre eux et effectuer par la suite cette classification.

Le but du mémoire était de comparer entre quelques méthodes existantes qui permettent de calculer la similitude entres les structures arborescentes des documents XML. Afin de les expérimenter et déterminer les meilleurs d'entre eux, ceux qui présentent les meilleurs résultats en qualité et en durée de la classification.

Nous avons choisi 5 méthodes à tester, qui sont : la méthode de Selkow, la méthode de Chawathe et la méthode de Dalamagas, méthode de Levenshtein et la méthode de Aïtelhadj. Parmi ces méthodes qu'on a testé, on a la méthode de Dalamagas qui a présenté le meilleur temps d'exécution mais avec une qualité de classification imprécise. Et on a la méthode de Aitelhadj qui a présenté le meilleur résultat pour la qualité de la qualification mais avec une durée d'exécution relativement supérieur à celle de la méthode Dalamagas.

Le logiciel dans sa première version propose cinq méthodes de classifications. Il est doté d'une interface ergonomique mais qui reste adapté à notre cas d'étude. Alors on aimerait améliorer le logiciel en lui implémentant d'autres méthodes de classification qui prennent de prendre en considération d'autres facteurs de comparaison tel que le contenu (données) des documents XML afin de rendre la qualité de la classification plus précise. Aussi, comme on l'a dit, l'extraction des structures des documents XML est faite par un module externe. On aimerait bien comme perspective améliorer le logiciel en lui intégrant un module qui permet d'extraire automatiquement les structures des documents XML et rendre le logiciel autonome.

Ainsi notre travail s'achève, mais le grand travail est à venir. En effet, le test n'est qu'un moyen d'analyser l'existant en matière des méthodes de classification des documents XML. Actuellement, les méthodes de classification ont toujours une complexité temporelle qui reste un peu importante. Et il serait judicieux de proposer un algorithme performant qui combine entre qualité et durée d'exécution.

Les références

Les livres

- Gilles Chagnon, Florent NOLOT, « XML », Marsat France : Pearson Education, 2007
- Alexandre Brillant, « XML cours et exercices », Paris France : Eyrolles, 2007
- Joseph Gabay, David Gabay, UML2 Analyse et Conception, DUNOD, Belgique, 2008.
- Emmanuel Puybaret, Les Cahiers du programmeur Java 1.5 et 5.0, 3^{ème} édition, EYROLLES,

Les sites Web

- **[Norme XML (site officiel En)]** <http://www.w3.org/TR/2008/REC-xml-20081126/>
- **[XML]** http://fr.wikipedia.org/wiki/Extensible_Markup_Language
- **[Prim]** http://en.wikipedia.org/wiki/Prim%27s_algorithm
- **[Levenshtein]** http://fr.wikipedia.org/wiki/Distance_de_Levenshtein

Les articles

- Tebri Hamid « *Formalisation et spécification d'un système de filtrage incrémental information* » Thèse de Doctorat, Université Paul Sabatier, Toulouse, (2004)
- Lewis D. *The TREC-4, Filtering Track*. In Proceedings of the 4th Text Retrieval Conference, TREC-4, National Institute of Standards and Technology NIST, Special Publication SP 500-263, Gaithersburg, MD Maryland, USA, 1–3 November, pages 165–180, (1996)
- Sebastiani Fabrizio. « *Machine Learning in Automated Text Categorization* ». ACM Computing Surveys, 34(1) : 1–47, (2002)
- Yang Y, Slattery S, Ghani R. « *A Study of Approaches to Hypertext Categorization* ». Journal of Intelligent Information Systems IIIS, 18(2-3) : 219–241, (2002)
- Denoyer L. « *Apprentissage et inférence statistique dans les bases de documents structurés : Application aux corpus de documents textuels* ». Thèse de Doctorat, Université Paris 6, (2004)
- Robertson S. *The probability ranking principle in IR*. Journal of documentation, 33(4) : 294–304, (1997)
- Jain A.K, Dubes R.C. *Algorithms for Clustering Data*. Prentice-Hall advanced reference series : Computer Science, Prentice-Hall, Inc, Upper Saddle River, NJ, New Jersey, (1988)

- Zhao Y, Karypis G. *Criterion Functions for Document Clustering : Experiments and Analysis*. Technical Report, 01–40, Department of Computer Science, University of Minnesota, Minneapolis, MN 55455, (2001)
- Ali Aïtelhadj, Mohamed Mezghiche et Fatiha Souam, Classification de Structures Arborescentes : Cas de Documents XML, Conférence en Recherche d'Information et Applications, CORIA 2009
- Thierry Despeyroux, Yves Lechevallier, Brigitte Trousse, Anne-Marie Vercoustre, Expériences de classification d'une collection de documents XML de structure homogène 5ème Journées d'Extraction et de Gestion des Connaissances (EGC 2005), 2005.
- Sudarshan Chawathe, Comparing Hierarchical Data in External Memory,in Proceedings of the Twenty-fifth International Conference on Very Large Data Bases, 1999,90—101.
- Anna Lesniewska, Clustering XML Documents by Structure, Advances in Databases and Information Systems, 2009, 238-246.
- S. M. Selkow, The tree-to-tree editing problem, Information Processing Letters 6 (1977) 184—186
- Sudarshan Chawathe, Comparing Hierarchical Data in External Memory,in Proceedings of the Twenty-fifth International Conference on Very Large Data Bases, 1999,90--101.
- Charu C. Aggarwal,Haixun Wang, Managing and Mining Graph Data, Springer, 2010.
- Mong Li Lee, Liang Huai Yang, Wynne Hsu, Xia Yang, XClust:Clustering XML Schemas for Effective Integration, 2002
- Alsayed Algergawy, Eike Schallehn, Gunter Saake, A Schema Matching-based Approach to XML Schema Clustering, Proceedings of the 10th International Conference on Information Integration and Web-based Applications & Services,131--136, 2008
- Calin Garboni, Florent Masseglia, et Brigitte Trousse, Sequential Pattern Mining for Structure-Based XML Document Classification, Workshop of the Initiative for the Evaluation of XML Retrieval, 2005.
- R. Wilson, M. Cobb, F. McCreedy, R. Ladner, D. Olivier, T. Lovitt, K. Shaw, F. Petry, M. Abdelguer, Geographical data interchange using xml-enabled technology within the GIDB system, in: A. B. Chaudhri, A. Rashid, R. Zicari (Eds.), XML Data Management, 2003, Addison Wesley.
- Y. Papakonstantinou, H. Garcia-Molina, J. Widom, Object exchange across heterogenous information sources, in: Proceedings of IEEE International Conference on Data Engineering, 1995, pp. 251--260.
- H. Garcia-Molina, Y. Papakonstantinou, D. Quass, A. Rajaraman, Y. Sagiv, J. D. Ullman, V. Vassalos, J. Widom, The TSIMMIS approach to mediation: Data models and languages, Journal of Intelligent Information Systems 8 (2) (1997) 117--132.
- S. M. Selkow, The tree-to-tree editing problem, Information Processing Letters 6 (1977) 184--186
- Theodore Dalamagas, Tao Cheng, Klaas-Jan Winkel, Timos Sellis, Clustering XML documents by structure, Inf. Syst., vol 31, 2006, 187--228

- Sudarshan Chawathe, Comparing Hierarchical Data in External Memory,in Proceedings of the Twenty-fifth International Conference on Very Large Data Bases, 1999,90—101
- Levenshtein V.I., Binary codes capable of correcting deletions, insertions and reversals, Soviet Physics Doklady 10(8), pp. 707-710, février 1966.
- T. Cormen, C. Leiserson, R. Rivest, Introduction to algorithms, MIT Press, 1990
- J. C. Gower, G. J. S. Ross, Minimum spanning trees and single linkage cluster analysis, Applied Statistics 18 (1969) 54--64.
- R. C. Prim: *Shortest connection networks and some generalizations*. In: *Bell System Technical Journal*, 36 (1957), pp. 1389–1401
- Sudarshan Chawathe, Comparing Hierarchical Data in External Memory,in Proceedings of the Twenty-fifth International Conference on Very Large Data Bases, 1999,90--101.
- Binary codes capable of correcting deletions, insertions, and reversals, SOVIER PHYSICS-DOKLADY, Vol10, N°8 (1966)
- Ali Aïtelhadj, Mohand Boughanem, Mohamed Mezghiche, Fatiha Souam , Using structural similarity for clustering XML documents, Springer, 2011
- Van Rijsbergen J.C. Information retrieval. Butterworths, London 2ème edition, (1979)