

l'analyse

AUTOMATIQUE des DOCUMENTS

BIBLIOTHEQUE DU CERIST

PAR maurice coyaud
nelly siot-decauville

BIBLIOTHEQUE DU CERIST

INFORMATIQUE

I

Maurice COYAUD
Nelly SIOT-DECAUVILLE

874

L'analyse automatique des documents

mouton
paris
la haye

IST 217

BIBLIOTHEQUE DU CERIST

PLAN

AVANT-PROPOS

PREMIÈRE PARTIE

ÉTUDE THÉORIQUE : LES DIFFÉRENTES MÉTHODES D'ANALYSE AUTOMATIQUE DES DOCUMENTS

0. TYPOLOGIE GÉNÉRALE.

1. LES MÉTHODES DE SÉLECTION.

1.0 PRINCIPE ET CRITÈRES.

1.1 LA TABULATION.

- 11 *Choix des titres.*
- 12 *L'antidictionnaire.*
- 12 *Les séparateurs.*

1.2 L'EXTRACTION DE PHRASES.

1.3 RÉCAPITULATION.

- 31 *L'antidictionnaire.*
- 32 *Les synonymies.*
- 33 *Les polysémies.*
- 34 *Le critère statistique.*
- 35 *Valeur sémantique du critère de co-occurrence.*

1.4 CONCLUSION.

2. LES MÉTHODES LINGUISTIQUES.

2.0 DÉFINITION DU « LANGAGE DOCUMENTAIRE ».

2.1 L'ANALYSE GRAMMATICALE AUTOMATIQUE.

11 *L'analyse par constituants.*

- 111 *Définition.*
- 112 *La méthode des stemmas.*
- 113 *L'analyse prédictive.*
 - A — *Exposé des travaux de Lemmon.*
 - B — *Le travail de Plath.*

12 *L'analyse par « chaînes » (string analysis).*

13 *L'analyse transformationnelle.*

130 *Définition des Transformations par Harris et Chomsky.*

131 *Le projet de Harris.*

132 *Exemples d'analyses transformationnelles (AGT).*

133 *Examen du projet de Harris d'un point de vue linguistique.*

- A — *Définition des éléments terminaux.*
- B — *Problèmes de l'AGT automatique.*

134 *Examen du projet d'un point de vue documentaire.*

- 14 *Récapitulation.*
- 141 Structures formelles et contenus sémantiques.
- 142 Limites des métalangages d'AG.
- 2.2 L'INDEXATION.
 - 21 *Questions préjudicielles.*
 - 211 Le moment de la réalisation du LD.
 - 212 Faut-il définir un « langage intermédiaire » ?
 - A — Passage par une langue standardisée.
 - B — Passage par un langage intermédiaire.
 - 213 Niveau de généralité de l'indexation.
 - 22 *Phases de l'indexation.*
 - 221 Découpage.
 - 222 Analyse morphologique.
 - 223 Traduction lexicale.
 - 224 Traduction syntaxique.
 - 23 *Les outils sémantiques.*
 - 231 Les classifications « manuelles ».
 - 232 La classification automatique.
 - A — Classification automatique à l'aide d'informations grammaticales.
 - B — Classification à l'aide de critères non grammaticaux.
- 2.3 LE RÉSUMÉ AUTOMATIQUE

3. LES MÉTHODES MIXTES.

- 3.1 MÉTHODES DE SÉLECTION ASSOCIÉES A DES MÉTHODES LINGUISTIQUES.
 - 11 *Méthodes de sélection et analyse grammaticale (AG).*
 - 111 Méthode de Storm.
 - 112 Projets de Climenson.
 - 12 *Méthodes de sélection et d'indexation.*
- 3.2 L'ANALYSE GRAMMATICALE COMBINÉE A L'INDEXATION.
 - 21 *La méthode de Stjažkin.*
 - 22 *Les travaux de Jessica Melton.*
- 3.3 SYNOPSIS.

SECONDE PARTIE

ÉTUDE EXPÉRIMENTALE

4. MISE EN ŒUVRE D'UNE EXPÉRIENCE D'INDEXATION AUTOMATIQUE

- 4.1 PRINCIPES DE LA MÉTHODE EXPÉRIMENTÉE.
- 4.2 RÉALISATION PRATIQUE : RÈGLES LINGUISTIQUES.
 - 21 *Moyens utilisés pour la traduction lexicale.*
 - 211 Le lexique documentaire.
 - 212 Le dictionnaire.
 - A — Codage morphologique.
 - B — Règles pour les polysémies.
 - C — « Lexies » et groupes de mots.
 - 22 *Moyens utilisés pour la traduction syntaxique.*
 - 221 Le réseau notionnel.
 - 222 Classes d'outils syntaxiques.
 - 223 Règles de construction syntaxique.
- 4.3 PROGRAMME.
 - 31 *Organisation générale.*
 - 32 *Introduction des données permanentes.*

- 321 Les désinences.
 - A — Les désinences de noms ou d'adjectifs.
 - B — Les désinences de verbes.
 - 322 Le lexique.
 - 323 Le dictionnaire.
 - 324 Le réseau notionnel.
 - 325 Les constructions.
 - 33 *Traitement des résumés.*
 - 331 Introduction des résumés.
 - 332 Analyse lexicale.
 - 333 Analyse syntaxique.
 - 34 *Exploitation.*
 - 341 Conditions
 - 342 Temps.
 - 343 Extensions possibles du programme actuel.
 - 35 *Conclusions.*
5. ANALYSE DES RÉSULTATS.
- 5.0 POUVOIR DE RECONNAISSANCE DU DICTIONNAIRE.
 - 01 *Les formes non reconnues.*
 - 02 *Groupes de formes identiques.*
 - A — Accidents morphologiques.
 - B — Variantes syntaxiques.
 - 5.1 RÉSULTATS DE LA TRADUCTION LEXICALE.
 - 11 *Valeur des traductions obtenues.*
 - 111 Règles de résolution des polysémies.
 - A — Contenu des règles.
 - B — Choix d'une unité de contexte.
 - 112 Mots ambigus non affectés de règles de polysémie.
 - 113 Identification des groupes de mots.
 - 12 *Comparaison entre analyses manuelles et automatiques.*
 - 121 Les temps d'indexation lexicale.
 - 122 Caractéristiques internes des deux types d'indexation.
 - 5.2 RÉSULTATS DE LA TRADUCTION SYNTAXIQUE.
 - 21 *Méthode d'analyse des résultats.*
 - 211 Exemple d'analyse automatique.
 - 212 Critères d'appréciation des résultats.
 - 22 *Évaluation des résultats syntaxiques.*
 - 221 Sondages.
 - 222 Causes d'erreurs et remèdes.
 - A — Absence d'AG.
 - B — Défaillances du système des « outils syntaxiques ».
- 6 CONCLUSION.
- 6.1 INDEXATION LEXICALE.
 - 6.2 INDEXATION SYNTAXIQUE.

ANNEXES.

- 1. Analyse morphologique.
- 2. Idiomatismes.
- 3. Listes d'outils syntaxiques.
- 4. Catalogue des règles syntaxiques.
- 5. Catalogue des règles de résolution de polysémies, avec leurs corrections.
- 6. Exemples d'indexations automatiques (avec évaluations).

ABRÉVIATIONS

I.C.C.L.	International Congress for Computational Linguistics.
S.D.C.	System Development Corporation.
A.M.T.C.L.	Association for Machine Translation and Computational Linguistics.
Amer. Doc.	American Documentation.
Congrès de Moscou'	Conference on Information Processing, Machine Translation and Automatic Reading of Text, 1961.
Congrès de Cleveland	International Conference for Standards on a Common Language for Machine Searching and Translation, New York, Interscience Publishers, 1961.
J.P.R.S.	Joint Publication Research Service, Washington.
Lg.	Language.
I.F.I.P.	International Federation of Information Processing.
C.E.A.	Commissariat à l'Energie Atomique.
T.D.A.P.	Transformation and Discourse Analysis Papers, University of Pennsylvania.

A.G.	Analyse grammaticale.
T.	Transformations.
M.A.G.	Métalangage d'analyse grammaticale.
L.N.	Langue naturelle.
L.D.	Langage documentaire.
L.I.	Langage intermédiaire.

AVANT-PROPOS

1. *Définition*

L'analyse documentaire automatique est un ensemble d'opérations, relevant des mathématiques, de la linguistique, de la programmation, destinées soit à sélectionner certains éléments d'un document sans les changer (notre chapitre 1), soit à transformer la forme et le contenu du document (notre chapitre 2), soit à combiner les deux types de méthodes précédents (notre chapitre 3). Les liens entre la synthèse et l'analyse automatique de données linguistiques sont étroits. Ce n'est pas un hasard si un des chercheurs en analyse syntaxique des textes a pu nommer sa méthode « analyse par synthèse ». Les travaux ne manquent pas sur le thème de la synthèse automatique : « génération » de phrases (Yngve), synthèse de la parole (vocoder, etc.). Nous n'aborderons néanmoins pas le problème des rapports entre analyse et synthèse, qui pourrait faire l'objet d'une monographie séparée.

2. *Intérêt théorique de l'étude*

Il est inutile de souligner l'importance pratique des travaux sur l'analyse automatique des documents. On sait qu'il s'agit là d'un des goulots d'étranglement de la documentation automatique, et que c'est un domaine où il y a le plus de progrès à faire. Par contre, son intérêt théorique n'apparaît pas toujours. Bornons-nous à noter que la définition d'unités minimales de traitement de l'information documentaire est une pierre de touche pour le formalisme linguistique (Ecole de Harris, etc.). C'est un problème théorique important que la définition des moyens de traiter automatiquement les textes linguistiques : pour la segmentation de la chaîne écrite en unités minimales, pour le calcul de la structure de la chaîne (fonctions des morphèmes, rapports implicites entre lexèmes, « analyse du discours »), etc.

3. *Limites de la présente étude*

Nous circonscrivons notre tâche à l'examen des recherches en analyse automatique d'un point de vue surtout linguistique ; en outre, nous nous limiterons à l'analyse de la chaîne écrite.

Un point de vue linguistique : Les questions de programmation et les questions logiques de théorie des langages formels ne seront pas abordées directement. Naturellement, dans la deuxième partie de cet ouvrage, consacrée à une étude expérimentale, on exposera en

détail le programme utilisé ; mais au cours de la première partie, où nous ferons un examen général des méthodes d'analyse automatique, il ne sera pas question de programmes, même d'un point de vue général. Notons néanmoins que c'est là un thème fort souvent abordé, notamment dans la bibliographie américaine. Citons simplement, à titre d'exemple, la communication de P. Garvin à la Conférence Internationale de 1965 sur la linguistique appliquée aux automates¹. Le problème qu'il aborde a son importance ; il s'agit de savoir s'il faut créer un algorithme général, capable de traiter (d'absorber) toutes les grammaires, ou s'il faut au contraire -- par nécessité de fait adapter l'algorithme à la grammaire. A l'encontre des chercheurs qui semblent pencher pour la première solution, Garvin choisit la seconde, avec des arguments fort probants². L'essentiel de ces arguments est qu'il n'est pas possible d'imaginer un algorithme assez général pour traiter une grammaire sensible au contexte.

Une autre implication du point de vue prédominant choisi ici est que notre étude se bornera aux problèmes d'analyse interne³ des documents. Nous ne ferons donc que mentionner les nombreuses recherches sur l'analyse externe, notamment la constitution automatique d'index de références bibliographiques citées par les différents auteurs dans leurs articles (« citation indexes »). Le lecteur pourra se reporter à ce sujet aux travaux cités dans la bibliographie sous les noms de Lesk, Thompson et B.-A. Lipetz.

D'autre part, nous n'aborderons pas les problèmes de transcription ou transformation des notations symboliques artificielles (notation chimique, etc.) ; une bibliographie très abondante leur est consacrée : ces études constituent un tout en elles-mêmes, et méritent une monographie séparée, du fait qu'elles sont indépendantes dans une large mesure, des études concernées par le présent ouvrage.

L'analyse de la chaîne écrite : Nous nous limiterons enfin à l'étude des travaux concernant l'analyse de la chaîne écrite. Malgré des liens évidents, les problèmes posés par l'analyse de la chaîne parlée sont fort différents (linguistiquement et techniquement), et nous les laisserons entièrement de côté.

1. Plan général

Le présent ouvrage est divisé en deux parties. La première doit servir à dresser un tableau général des recherches effectuées sur le problème de l'analyse automatique documentaire. Ce tableau a des lacunes (par exemple en ce qui concerne les recherches en Union Soviétique)

1. « Some comments on algorithm and grammar in the automatic parsing of natural language », International Conference on Computational Linguistics, New York, 1965.

2. La même position est adoptée par Jane Robinson ; cf. « Endocentric constructions and the Cocke parsing logic », p. 1 (*ibid.*).

3. C'est-à-dire fondée sur les éléments proprement linguistiques du document : le texte de l'auteur, débarrassé d'éléments périphériques comme les références bibliographiques, les schémas, etc.

mais nous n'avons pas la prétention d'être complets. Plus qu'un inventaire des recherches effectuées¹, c'est un inventaire des problèmes fondamentaux que nous avons en vue.

La seconde partie de ce livre est consacrée à l'exposé d'une expérience d'indexation automatique de documents, effectuée à la Section d'Automatique Documentaire du C.N.R.S. en 1962-65.

L'expérience d'indexation automatique qui fait l'objet de la deuxième partie de ce livre a été menée sous la direction de M. J.-C. Gardin, à qui nous exprimons nos vifs remerciements. Nous remercions également Mlle Radmila Zygouris, qui a constitué le « réseau notionnel » expérimental, en psychophysiologie, Mlles M. P. Ferry et C. Charmot, qui ont participé à la compilation du dictionnaire expérimental².

La partie linguistique de la présente étude est due à M. Coyaud ; N. Siot-Decauville est l'auteur du programme expérimental d'indexation automatique vers le langage documentaire Syntol, et a rédigé la section 4.3 de ce livre.

M. C. et N. S.-D.

1. On trouvera sur ce thème des bibliographies quasi exhaustives dans *Current research and development in scientific documentation*, publié périodiquement par la National Science Foundation, et dans un rapport récent de Mary Stevens : « Automatic Indexing : A State-of-the-Art-Report », NBS, Washington, 1965.

2. La compilation de ce dictionnaire a été effectuée grâce à un contrat d'un an accordé par l'Euratom, qui s'est également chargé de la perforation du corpus expérimental de cinq cents résumés. A la suite de ce contrat (1961-1962), un rapport final a été adressé à l'Euratom ; il est maintenant publié par les Presses Académiques Européennes à Bruxelles : le *Syntol*, t. 4, « Analyse automatique » (1964).

BIBLIOTHEQUE DU CERIST

O. TYPOLOGIE GÉNÉRALE

Les méthodes d'analyse automatique des documents peuvent être classées en deux grandes catégories, suivant la nature des opérations mises en jeu : ou bien celles-ci consistent en un pur *choix* de certains éléments compris dans les documents (suivi d'un réarrangement ou d'une extraction éventuels de ces éléments); ou bien il s'agit d'une *transformation* du contenu du document. Nous donnerons ci-dessous des exemples détaillés et commentés de méthodes appartenant à ces deux catégories. Indiquons d'abord de façon générale un point de vue qui permet d'apprécier pratiquement la différence entre ces deux grandes catégories : le point de vue des résultats. La première catégorie de méthodes a produit des résultats passés dans l'usage courant depuis plusieurs années : index de permutations, index KWIC (Key Word in Context), entre autres. Le second groupe de méthodes, plus ambitieuses celles-là, n'a pas encore porté à terme ses fruits. A cela, rien d'étonnant, dès qu'on pense à la complexité des opérations engagées dès qu'il ne s'agit plus de *choisir*, sans les modifier, des éléments dans les textes, en vertu de critères statistiques, ou autres, mais de *transformer* le contenu de ces textes dans un langage (artificiel ou naturel) plus condensé et plus précis.

Les produits obtenus peuvent alors revêtir des formes bien diverses : soit le document est soumis à une *analyse grammaticale* tendant à manifester la structure des phrases — et dans ce cas, les graphes obtenus constitueraient des représentations adéquates des documents ; soit le document est traduit dans un langage documentaire¹ et cette traduction est appelée *indexation*, parce qu'elle consiste essentiellement à dégager les notions les plus importantes du document (et parfois leurs relations sémantiques) à l'aide de termes empruntés à un lexique documentaire (*i.e.* un index de descripteurs ou de mots-clés, ou *indexing language*) ; dans un troisième cas le document pourrait être transformé en un *résumé* en langue naturelle, analogue aux résumés faits par des analystes humains. Ce dernier type d'analyse documentaire est évidemment le plus ambitieux, et ne se verra pas illustré par des résultats avant longtemps ; mais puisqu'il y a des chercheurs qui travaillent dans cette voie, il serait injuste de ne pas évoquer leurs travaux, ne serait-ce que pour mémoire.

1. Le langage documentaire sera défini au § 1.20.