

# Extraction de Textes et Analyse de Données Textuelles

## Chapitre 4 : Méthodes et techniques pour la constitution et l'analyse de corpus textuels

### Introduction

Ce chapitre vous guidera à travers les méthodes et techniques essentielles pour constituer, extraire et analyser des corpus textuels en sciences humaines et sociales. L'analyse de données textuelles est devenue incontournable pour comprendre les phénomènes sociaux, culturels et communicationnels contemporains. Vous découvrirez comment collecter des données pertinentes, les organiser en corpus cohérents et les analyser avec les outils appropriés.

## 1. Construction d'un corpus de textes

### 1.1 Principes fondamentaux de constitution d'un corpus

#### 1.1.1 Définition et objectifs d'un corpus textuel

Un corpus textuel est un ensemble structuré de textes assemblés selon des critères explicites pour servir de base à une analyse. Sa constitution repose sur plusieurs principes fondamentaux :

- **Représentativité** : capacité du corpus à refléter fidèlement le phénomène étudié
- **Homogénéité** : cohérence interne des textes sélectionnés
- **Exhaustivité** : couverture complète du domaine ou de la période étudiée
- **Actualité** : pertinence temporelle des données par rapport à la question de recherche
- **Exploitabilité** : format et structure adaptés aux outils d'analyse

#### 1.1.2 Étapes de la constitution d'un corpus

1. Définition précise de l'objet de recherche
2. Identification des sources pertinentes
3. Élaboration des critères de sélection
4. Collecte des données
5. Nettoyage et prétraitement
6. Organisation et structuration
7. Documentation des métadonnées

### 1.2 Types de corpus textuels

#### 1.2.1 Corpus de discours

- **Discours politiques** : allocutions officielles, débats parlementaires, discours de campagne
- **Discours institutionnels** : communications d'organisations, rapports officiels
- **Discours médiatiques** : éditoriaux, tribunes, chroniques

- **Sources** : sites officiels, archives médiatiques, bases de données spécialisées
- **Particularités méthodologiques** : analyse des stratégies rhétoriques, étude des cadres interprétatifs, analyse de l'argumentation

### 1.2.2 Corpus d'articles de presse

- **Presse quotidienne** : nationale, régionale, internationale
- **Presse magazine** : généraliste, spécialisée
- **Presse en ligne** : sites d'information, pure players
- **Sources** : bases de données comme Europresse, Factiva, archives numériques des journaux
- **Particularités méthodologiques** : analyse du cadrage médiatique, étude des lignes éditoriales, analyse de la couverture médiatique d'événements

### 1.2.3 Corpus de médias sociaux

- **Twitter/X** : tweets, threads, conversations
- **Facebook** : posts, commentaires, groupes
- **Forums et blogs** : discussions, billets
- **Sources** : API officielles, outils de scraping, archives spécialisées
- **Particularités méthodologiques** : analyse de sentiment, détection de communautés, étude de la viralité, cartographie des interactions

### 1.2.4 Corpus littéraires et académiques

- **Œuvres littéraires** : romans, poésie, théâtre
- **Publications scientifiques** : articles, thèses, ouvrages
- **Sources** : bibliothèques numériques, bases de données académiques (JSTOR, Cairn, etc.)
- **Particularités méthodologiques** : analyse stylistique, étude des influences, analyse des réseaux de citations

## 1.3 Applications de l'analyse de corpus textuels

### 1.3.1 Extraction d'informations

- **Reconnaissance d'entités nommées** : identification de personnes, lieux, organisations
- **Extraction de relations** : détection des liens entre entités
- **Extraction d'événements** : identification d'actions et de leur chronologie
- **Extraction d'opinions** : détection des jugements et des positionnements
- **Applications concrètes** : veille médiatique automatisée, cartographie des acteurs d'un domaine, suivi de l'évolution d'une controverse

### 1.3.2 Clustering et classification de textes

- **Principes du clustering textuel** : regroupement de textes similaires sans catégories prédéfinies
- **Méthodes de clustering** : k-means, clustering hiérarchique, DBSCAN
- **Classification supervisée** : attribution de textes à des catégories prédéfinies
- **Méthodes de classification** : SVM, forêts aléatoires, réseaux de neurones

- **Applications concrètes** : découverte de thématiques émergentes, segmentation d'un corpus hétérogène, détection automatique de genres textuels

### 1.3.3 Analyse thématique et modélisation de sujets

- **Méthodes d'analyse thématique** : fréquences lexicales, cooccurrences, réseaux de termes
- **Modélisation de sujets (topic modeling)** : LDA, LSI
- **Applications concrètes** : cartographie des préoccupations dans un débat public, évolution temporelle des thématiques dans un corpus médiatique

## 2. Méthodes d'extraction de données textuelles

### 2.1 Approche par scraping (extraction de contenu web)

#### 2.1.1 Principes et enjeux du web scraping

- **Définition** : extraction automatisée de données à partir de sites web
- **Aspects légaux et éthiques** : respect des conditions d'utilisation, droit d'auteur, protection des données personnelles
- **Limites techniques** : défenses anti-scraping, captchas, blocage d'IP

#### 2.1.2 Techniques et outils de scraping

- **Outils sans programmation** : Octoparse, Import.io, Web Scraper
- **Bibliothèques Python** : BeautifulSoup, Scrapy, Selenium
- **Stratégies d'extraction** :
  - Sélecteurs CSS et XPath
  - Navigation automatisée
  - Gestion des cookies et sessions
  - Contournement des limitations

#### 2.1.3 Workflow complet d'un projet de scraping

1. Analyse de la structure du site cible
2. Identification des données pertinentes
3. Conception du script d'extraction
4. Tests et ajustements
5. Planification (scheduling) des extractions
6. Stockage et organisation des données recueillies
7. Documentation du processus

#### 2.1.4 Travaux pratiques

- Extraction d'articles de presse d'un média en ligne
- Constitution d'un corpus de commentaires sur un forum thématique
- Collecte automatisée de publications d'un site institutionnel

### 2.2 Extraction via API (Interface de Programmation d'Application)

### 2.2.1 Principes et avantages des API

- **Définition** : interfaces officielles permettant d'accéder aux données d'un service
- **Avantages** : légalité, fiabilité, structuration des données, documentation
- **Limitations** : quotas d'utilisation, restrictions d'accès, données filtrées

### 2.2.2 API majeures pour les sciences sociales

- **API de réseaux sociaux** : Twitter/X API, Facebook Graph API, Reddit API
- **API d'actualités** : The Guardian, New York Times, Le Monde
- **API académiques** : Scopus API, CrossRef, Google Scholar (non officielle)
- **API gouvernementales et institutionnelles** : data.gouv.fr, Eurostat, Banque mondiale

### 2.2.3 Techniques d'utilisation des API

- **Authentication** : clés API, OAuth
- **Construction de requêtes** : paramètres, filtres, pagination
- **Gestion des réponses** : JSON, XML
- **Bibliothèques Python facilitant l'accès** : tweepy (Twitter), PRAW (Reddit), etc.

## 2.3 Nettoyage et prétraitement des corpus

### 2.3.1 Techniques de nettoyage textuel

- **Normalisation** : mise en minuscules, suppression des accents
- **Filtrage** : élimination des balises HTML, des caractères spéciaux
- **Gestion des doublons** : détection et suppression
- **Correction orthographique** : détection et correction automatique d'erreurs

### 2.3.2 Prétraitement linguistique

- **Tokenisation** : segmentation en unités lexicales
- **Suppression des mots vides (stop words)**
- **Lemmatisation/stemming** : réduction des formes fléchies
- **Annotation morphosyntaxique** : attribution de catégories grammaticales
- **Analyse syntaxique** : identification des relations entre mots

### 2.3.3 Organisation et stockage

- **Formats de stockage** : texte brut, CSV, JSON, XML, bases de données
- **Structuration des métadonnées** : date, source, auteur, contexte
- **Versionning** : gestion des modifications et des mises à jour du corpus

## 3. Les créateurs de contenus et leur rôle dans la production et la diffusion de l'information

### 3.1 Typologie des créateurs de contenus

#### 3.1.1 Acteurs médiatiques traditionnels

- **Journalistes professionnels** : pratiques, contraintes, évolution du métier
- **Rédactions et ligne éditoriale** : processus de sélection et de cadrage de l'information
- **Experts et chroniqueurs** : légitimité et influence dans le débat public

### 3.1.2 Nouveaux acteurs numériques

- **Blogueurs et influenceurs** : émergence, modèles économiques, impact
- **Créateurs de contenu sur plateformes** : YouTubers, streamers, podcasters
- **Lanceurs d'alerte et médias citoyens** : contre-pouvoir et journalisme participatif

### 3.1.3 Institutions et organisations

- **Communication institutionnelle** : stratégies discursives des organisations
- **Think tanks et groupes d'intérêt** : production d'expertise orientée
- **Partis politiques et mouvements sociaux** : communication militante

## 3.2 Dynamiques de production et de circulation de l'information

### 3.2.1 Logiques de production

- **Contraintes professionnelles** : temps, format, audience
- **Modèles économiques** : influence sur le contenu produit
- **Routines et pratiques** : sources, vérification, narration

### 3.2.2 Circuits de diffusion

- **Hierarchisation de l'information** : critères de newsworthiness
- **Amplification et viralité** : mécanismes de propagation
- **Bulles informationnelles** : fragmentation des publics et personnalisation

### 3.2.3 Méthodes d'analyse des producteurs de contenu

- **Analyse de réseau** : cartographie des relations entre acteurs
- **Analyse de discours** : étude des stratégies énonciatives
- **Ethnographie numérique** : observation des pratiques de production

## 3.3 Étude de cas et travaux pratiques

- Analyse comparative des cadrages médiatiques d'un événement entre médias traditionnels et alternatifs
- Cartographie des influenceurs sur une thématique spécifique
- Étude de la circulation d'une information entre différentes sphères médiatiques

## 4. Méthodes qualitatives : questionnaires, questions ouvertes et entretiens

### 4.1 Les questionnaires avec questions ouvertes

#### 4.1.1 Conception de questions ouvertes efficaces

- **Types de questions ouvertes** : descriptives, explicatives, projectives
- **Formulation** : neutralité, clarté, précision
- **Placement stratégique** dans le questionnaire
- **Complémentarité avec les questions fermées**

#### 4.1.2 Administration des questionnaires

- **Modes d'administration** : en ligne, face-à-face, téléphonique
- **Échantillonnage** : représentativité, taille, méthodes de sélection
- **Outils numériques** : LimeSurvey, Google Forms, Qualtrics, SurveyMonkey

#### 4.1.3 Traitement des réponses ouvertes

- **Codification manuelle** : élaboration de grilles d'analyse
- **Analyse textuelle assistée par ordinateur** : lexicométrie, classification
- **Méthodes mixtes** : articulation quanti-quali

### 4.2 Les entretiens de recherche

#### 4.2.1 Types d'entretiens

- **Entretien directif** : questions précises et ordonnées
- **Entretien semi-directif** : guide d'entretien flexible
- **Entretien non-directif** : liberté maximale de l'interviewé
- **Entretien d'expertise** : spécificités méthodologiques
- **Focus groups** : dynamique collective

#### 4.2.2 Préparation et conduite d'entretiens

- **Élaboration du guide d'entretien** : thématiques, questions, relances
- **Recrutement des participants** : critères, prise de contact
- **Techniques d'interview** : écoute active, reformulation, gestion du silence
- **Aspects pratiques** : enregistrement, prise de notes, cadre spatiotemporel

#### 4.2.3 Traitement et analyse des entretiens

- **Transcription** : conventions, outils d'aide à la transcription
- **Analyse thématique** : codage manuel ou assisté par logiciel
- **Analyse de discours** : étude des stratégies rhétoriques et énonciatives
- **Logiciels spécialisés** : NVivo, ATLAS.ti, MAXQDA

### 4.3 De la collecte à l'analyse : intégration des données qualitatives dans un projet de recherche

#### 4.3.1 Triangulation méthodologique

- **Complémentarité des approches** : questionnaires, entretiens, analyse de corpus
- **Validation croisée** des résultats
- **Résolution des contradictions** entre sources de données

### 4.3.2 Analyse intégrée des données textuelles

- **Création d'un corpus unifié** : entretiens, questions ouvertes, textes existants
- **Codification commune** : harmonisation des catégories d'analyse
- **Interprétation contextuelle** : mise en perspective des résultats

### 4.3.3 Restitution et valorisation des résultats

- **Écriture scientifique** : intégration des verbatims et extraits
- **Visualisation** : représentations graphiques des données qualitatives
- **Communication adaptée aux différents publics**

## 5. Travaux pratiques intégrateurs

### 5.1 Projet 1 : Analyse d'un débat public sur les réseaux sociaux

1. Constitution d'un corpus via l'API Twitter/X et techniques de scraping
2. Analyse des acteurs et de leur influence
3. Analyse thématique et identification des cadrages dominants
4. Complément par entretiens avec des participants clés
5. Mise en perspective des résultats

### 5.2 Projet 2 : Étude de la couverture médiatique d'un enjeu social

1. Construction d'un corpus d'articles de presse via scraping et bases de données
2. Analyse comparative entre différents types de médias
3. Entretiens avec des journalistes et créateurs de contenu
4. Questionnaire sur la réception auprès du public
5. Synthèse multidimensionnelle

### 5.3 Projet 3 : Analyse d'un corpus de témoignages

1. Collecte de témoignages via questionnaires à questions ouvertes
2. Approfondissement par entretiens semi-directifs
3. Constitution d'un corpus textuel unifié
4. Analyse thématique et lexicométrique
5. Restitution narrative et visuelle des résultats

## Conclusion

L'extraction et l'analyse de données textuelles constituent un domaine en constante évolution, à l'intersection des méthodes qualitatives traditionnelles et des approches computationnelles. La maîtrise de ces techniques offre aux chercheurs en sciences humaines et sociales des perspectives inédites pour aborder des corpus volumineux et diversifiés, tout en maintenant la profondeur interprétative caractéristique de ces disciplines.

## Bibliographie et ressources

## Ouvrages de référence

- Bardin, L. (2013). *L'analyse de contenu*. PUF.
- Brin, C., Charron, J., & de Bonville, J. (2004). *Nature et transformation du journalisme : théorie et recherches empiriques*. Presses de l'Université Laval.
- Kaufmann, J. C. (2016). *L'entretien compréhensif*. Armand Colin.
- Mitchell, R. (2018). *Web Scraping with Python: Collecting Data from the Modern Web*. O'Reilly Media.
- Paillé, P., & Mucchielli, A. (2016). *L'analyse qualitative en sciences humaines et sociales*. Armand Colin.
- Ratinaud, P., & Marchand, P. (2015). *Des mondes lexicaux aux représentations sociales. Une approche à l'aide de logiciels libres*. *Revue d'Anthropologie des Connaissances*, 9(1), 3-35.

## Ressources en ligne

- Tutoriels Python pour le text mining : <https://www.nltk.org/book/>
- Documentation des API pour la recherche : <https://developer.twitter.com/en/docs/twitter-api/academic-research>
- Guides méthodologiques sur l'entretien : <https://www.cairn.info/revue-recherche-en-soins-infirmiers-2005-3-page-60.htm>
- Communauté d'entraide pour le scraping : <https://www.scrapingbee.com/blog/>

## Logiciels et outils recommandés

- **Analyse qualitative** : NVivo, ATLAS.ti, MAXQDA, Iramuteq
- **Scraping** : BeautifulSoup, Scrapy, Selenium
- **Analyse textuelle** : Voyant Tools, TXM, IRaMuTeQ, Lexico
- **Questionnaires** : LimeSurvey, SurveyMonkey, Sphinx
- **Visualisation** : Gephi, VOSviewer, Tableau