



Université de Rouen Normandie - UFR Sciences et Techniques
Master 2 mention Bioinformatique – Parcours BIMS
2023 - 2024

Rapport de stage

Analyse textuelle d'articles scientifiques évaluant l'impact des vers de terre sur l'environnement

Présenté et soutenu par

Antoine Malet

Campus Agro Paris Saclay, Unité MIA Paris-Saclay
Equipe SOLsTIS

Encadrants :

David Makowski
Sophie Donnet





Université de Rouen Normandie - UFR Sciences et Techniques
Master 2 mention Bioinformatique – Parcours BIMS
2023 - 2024

Rapport de stage

Analyse textuelle d'articles scientifiques évaluant l'impact des vers de terre sur l'environnement

Présenté et soutenu par

Antoine Malet

Campus Agro Paris Saclay, Unité MIA Paris-Saclay
Equipe SOLsTIS

Encadrant :

David Makowski
Sophie Donnet



Remerciements

En premier lieu, j'aimerais remercier mes encadrants pour ce stage, M. David MAKOWSKI et Mme Sophie DONNET, pour m'avoir chaleureusement accueilli au sein de l'équipe et conseillé avec beaucoup de sagesse et de patience. Merci aussi pour leurs conseils et leurs efforts d'accompagnement et de relecture de mes travaux! Je tiens aussi à adresser un mot particulier à mes vaillants collègues de bureau Emré ANAKOK et Caroline COGNOT, pour leur compagnie perpétuelle et leurs très bons conseils.

Je remercie aussi Louis LACOSTE, pour ses excellents conseils en R et en cinématographie, ainsi que François VICTOR, pour ses généreuses explications en statistiques théoriques auxquelles je n'ai souvent pas compris grand-chose. Merci aussi à Armand FAVROT, pour son accueil et ses invitations aux Eventos des repas organisés par la cafet (les fraises maison étaient légendaires!). Enfin, chaleureuses salutations à tous mes collègues de pause café, qui sont toujours restés sympathiques et accueillants même si je n'ai jamais bu la moindre goutte de leur breuvage sacré. Courage, peut-être qu'un jour vous me convertirez à votre religion!

Merci enfin à tous ceux que je n'ai pas nommés, particulièrement aux personnels qui prennent soin des locaux en silence, sans qui l'infrastructure de travail ne pourrait pas fonctionner correctement.

Table des matières

Remerciements	I
Table des matières	III
Liste des Abréviations	VII
1 Introduction	1
1.1 Unité MIA, Campus Agro Paris Saclay	1
1.2 Le vers de terre dans la littérature scientifique	2
1.3 Objectifs de mon travail	3
2 Ressources	5
2.1 Environnement informatique	5
2.2 Pratique Professionnelle	5
2.2.1 Veille bibliographique et technologique	5
2.2.2 Bonnes pratiques	5
2.2.3 Communication des travaux	6
2.3 Outils informatiques et statistiques	6
2.3.1 Récupération des abstracts et des métadonnées avec Python	6
2.3.2 Text Mining avec R	8
2.4 Données	9
3 Résultats	11
3.1 <i>Web-scraping</i> et obtention d'une base de métadonnées en format CSV	11
3.2 Analyse globale des textes bruts	12
3.3 Rang, approche tf-idf et loi de Zipf	14
3.4 Analyses de bigrammes seuls et réseaux de bigrammes	16
3.5 Analyse de sentiment	19
4 Discussion	21
5 Conclusion	23

Table des figures

1.1 Organigramme de l'UMR MIA Paris-Saclay	1
1.2 Les services écosystémiques rendus par les vers de terre.	3
1.3 Distribution régionale de <i>Lumbricus rubellus</i> (espèce exotique) aux USA. . . .	3
3.1 Fréquence des racines de mots dans le corpus entier	12
3.2 Log-log scatter plot montrant les corrélations du choix des mots entre chaque métaanalyse, à l'échelle logarithmique.	13
3.3 Liste des mots les plus représentées pour chaque métaanalyse	14
3.4 Mesures de tf-idf pour chacune des métaanalyses du corpus.	15
3.5 Loi de Zipf pour les abstracts du corpus	16
3.6 Mesures de tf-idf par bigramme pour chacune des métaanalyses du corpus . .	16
3.7 Réseau de bigrammes orienté montrant les connexions les plus fréquentes reliant les mots du corpus.	17
3.8 Réseau de bigrammes orienté montrant les connexions les plus fréquentes reliant les mots de la MA1.	18
3.9 Réseau de bigrammes orienté montrant les connexions les plus fréquentes reliant les mots de la MA2.	18
3.10 Réseau de bigrammes orienté montrant les connexions les plus fréquentes reliant les mots de la MA3.	19
3.11 Réseau de bigrammes orienté montrant les connexions les plus fréquentes reliant les mots de la MA4.	19
3.12 Nuage de mots de l'ensemble des racines de mots positives/négatives retrouvées dans la base de données, colorées par sentiment	20
3.13 Contributions de chaque racine au sentiment global pour chaque métaanalyse	20

Liste des Abréviations

ASCII American Standard Code for Information Interchange

API Application Programming Interface

CSS Cascade Style Sheet

DOI Digital Object Identifier

HTML Hypertext Markup Language

IDE De l'anglais, Environnement de Développement Intégré

MA Métaanalyse

MIA Mathématiques et Informatique Appliquée

Rmd R Markdown

RG ResearchGate

RGS2 ResearchGateScraper2.py

UMR Unité Mixte de Recherche

Glossaire

JSON : JSON (JavaScript Object Notation) est un format de fichier textuel conçu pour la structuration et l'échange de données^[2].

log-log scatter plot : Type de graphique sous forme de nuage de points montrant les relations entre deux variables. Le passage à l'échelle logarithmique permet une meilleure visualisation des résultats.

Markdown : Markdown est un langage de balisage léger qui permet de formater du texte de manière simple et rapide. Il utilise des caractères spéciaux pour indiquer les éléments de mise en forme, tels que les titres, les listes, les liens, etc. Les fichiers Markdown peuvent être convertis en HTML pour être affichés sur un site web ou dans un logiciel de traitement de texte ^[1].

Métaanalyse : Article scientifique présentant la combinaison des résultats statistiques d'une série d'études indépendantes sur un problème donné.

N-gram : Suite de mots consécutifs de taille n (une des possibilités de tokenisation). Utile pour comprendre les relations logiques entre les mots. Les bigrammes sont un cas particulier de n-gram (n-gram de longueur 2).

Racinement (linguistique) : Obtention du radical, par exemple par dépréfixation ou désuffixation (*Exemple* : "enhance", "enhances" et "enhancement" deviennent tous "enhanc").

Réseau de bigrammes : Figure permettant de visualiser les relations entre les différents *tokens* simultanément, plutôt que deux par deux. Cela permet d'aller plus loin que l'analyse de bigrammes séparés les uns des autres. Les noeuds avec plus de deux connexions sont nommés des **centres**.

Token : Unité textuelle souvent réduite, voire ne comprenant qu'un seul mot, issue du processus de **tokenisation**.

Tokenisation : Processus consistant à décomposer un texte ou un corpus de textes en unités textuelles plus réduites, comme des mots, des n-grams ou des phrases.

Text-mining : Processus d'analyse textuelle consistant à transformer un texte non structuré en données structurées pour ensuite procéder à l'analyse statistique. Elle peut être utilisée pour analyser rapidement un corpus de texte massifs, plutôt que d'examiner tous les textes uns à uns.

Web scraping : Technique permettant d'extraire automatiquement de grandes quantités d'informations d'un site Web, sans intervention humaine directe, via un script informatique^[5].

Chapitre 1

Introduction

1.1 Unité MIA, Campus Agro Paris Saclay

Mon stage s'est déroulé au sein de l'Unité MIA (Mathématique et Informatique Appliqués), à l'INRAE du Campus Agro Paris Saclay, sur le plateau de Saclay. L'UMR MIA Paris-Saclay, associée aux tutelles AgroParisTech, INRAE et Université Paris Saclay, regroupe des statisticiens et des informaticiens spécialisés dans la modélisation et l'apprentissage statistique et informatique pour la biologie, l'écologie, l'environnement, l'agronomie et l'agro-alimentaire.

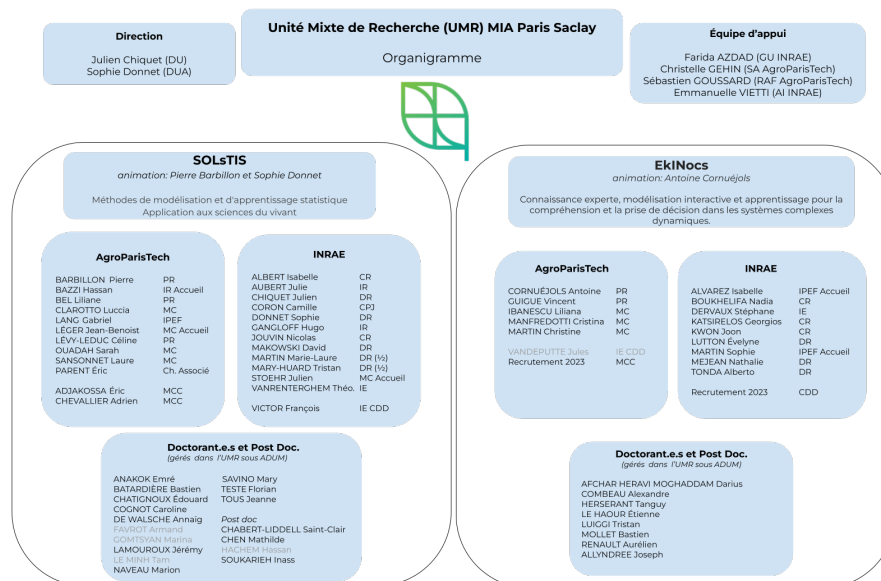


FIGURE 1.1 — Organigramme de l'UMR MIA Paris-Saclay

Les compétences mises en oeuvre portent sur les méthodes d'inférences statistiques et algorithmiques. L'unité développe des méthodes statistiques et informatiques originales, génériques ou motivées par des problèmes précis en science du vivant. Ses activités s'appuient sur une bonne culture dans les disciplines destinataires : écologie, environnement, agro-alimentaire, biologie moléculaire et biologie des systèmes. L'unité MIA est dirigée par Julien CHIQUET et Sophie DONNET, et comprends deux équipes distinctes : l'équipe SOLstIS (Statistical mOdeling and Learning for environnemenT and Life Science) dirigée par Sophie

DONNET et Pierre BARBILLON, et l'équipe EkiNocs (Expert Knowledge, INteractive modelING for understandING and decisiOn makING in dINamic Complex Systems), dirigée par Antoine CORNUÉJOLS.

En tant que stagiaire, j'ai ainsi pu intégrer SOLsTIS (figure 1.1) pour mettre au point des méthodes informatiques et statistiques pour l'analyse textuelle d'abstracts d'articles scientifiques. Dans ce cadre, j'ai aussi pu participer aux interventions de divers spécialistes du domaine, lors des séminaires hebdomadaires usuellement présentés le jeudi.

1.2 Le vers de terre dans la littérature scientifique

Selon la vision des scientifiques Européens, les vers de terre contribuent à l'aération du sol, au recyclage des éléments nutritifs, du carbone, et du phosphore. Ils jouent un rôle important dans le recyclage de la matière organique du sol, participent activement à leur décomposition en rendant d'importants éléments nutritifs accessibles à d'autres organismes vivants du sol, notamment aux végétaux. A ce titre, ils fournissent de nombreux services écosystémiques, notamment en jouant un rôle clé dans la production, la structuration, l'entretien et la productivité des sols, forestiers, prairiaux et agricoles figure 1.2 Kumar u. a. (2023) :

- Ils permettent une meilleure aération des sols (Kim u. a. (2017)).
- Ils favorisent le recyclage des éléments nutritifs, du carbone, et du phosphore (Lemtiri u. a. (2014)).
- Ils jouent un rôle important dans le recyclage de la matière organique des sols (Edwards und Arancon (2022)).
- Leur activité de décomposeurs (minéralisation de la matière organique) favorise la croissance d'autres organismes de l'environnement, augmentant de ce fait la productivité des plantes (Bertrand u. a. (2015)).
- Ils sont essentiels pour la structuration, l'entretien et la productivité des sols, forestiers, prairiaux et agricoles (Sharma u. a. (2017)).

Pourtant, la perception qu'ont les scientifiques des vers de terre est très variable selon la région géographique considérée. Car si ces animaux sont perçus de manière très positive en Europe, à l'inverse, en Amérique du nord :

- Les vers de terre exotiques (venus d'Europe et d'Asie) sont des espèces invasives susceptibles de perturber les écosystèmes natifs (Loss u. a. (2013)). figure 1.3.
- Leur présence dans les sols modifie durablement les propriétés physico-chimiques des écosystèmes souterrains (Bohlen u. a. (2004)).
- Les modifications physico-chimiques induites par les espèces exotiques causent une baisse globale de la biodiversité chez les espèces natives du milieu (Ferlian u. a. (2017)).
- L'impact de ces espèces exotiques sur les émissions de gaz à effet de serre est également sujet à controverse (Lubbers u. a. (2013); Forey u. a. (2023)).

La littérature scientifique à ce sujet se divise donc en deux camps opposés : certains articles soulignent le rôle positif du vers de terre tandis que d'autres le considèrent comme un organisme nuisible dont il faudrait si possible se débarrasser. Afin de mettre en évidence cette divergence d'opinion, je me suis appuyé sur les abstracts (courts résumés **standardisés**) de 116 articles scientifiques issus de 4 métaanalyses distinctes. Une métaanalyse est une démarche statistique qui permet de synthétiser quantitativement les résultats d'études indépendantes ayant trait à une question de recherche bien précise. Dans mon cas, les 4 méta-analyses se répartissaient dans les deux camps.

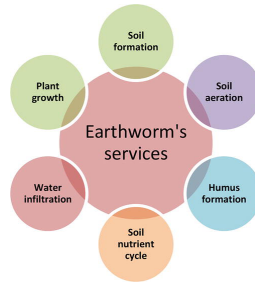


FIGURE 1.2 — Les services écosystémiques rendus par les vers de terre.

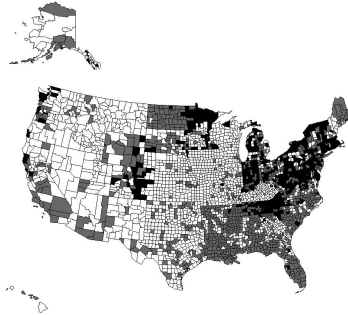


FIGURE 1.3 — Distribution régionale de *Lumbricus rubellus* (espèce exotique) aux USA. Les régions représentent les stations où : *L. rubellus* a été retrouvé (**noir**), une autre espèce que *L. rubellus* a été retrouvée (**gris**), aucune donnée n'est disponible (**blanc**)^[6].

1.3 Objectifs de mon travail

L'analyse informatique des abstracts des articles issus des 4 métaanalyses a été conduite pour répondre à la question de recherche suivante :

Les méthodes statistiques de *Text-mining* permettent-elles d'identifier des groupes distincts d'articles scientifiques défendant une conception opposée du rôle écologique du vers de terre au sein du corpus de texte fourni pour l'analyse ?

Pour répondre à cette question, nous procéderons en deux étapes :

1. Création d'une base de métadonnées (comprenant notamment les **abstracts** de chaque étude indépendante formant les 4 métaanalyses) via un script Python de *Web-scraping*, sous forme d'un fichier CSV brut destiné à être analysé.
2. Analyse textuelle des abstracts d'articles scientifiques présents dans le fichier CSV issu de l'étape précédente via un script R.

Une fois ces étapes achevées, nous discuterons les résultats obtenus et les méthodes employées en les comparant avec ceux d'autres articles de la littérature scientifique spécialisés dans ce domaine.

Chapitre 2

Ressources : pratiques professionnelles, environnement informatique, outils informatiques et statistiques, données

2.1 Environnement informatique

Le matériel qui m'a été fourni par le laboratoire est un ordinateur fixe HP Elite SFF 800 G9 Desktop PC, numéro de série CZC2479ZXS - 4G087AV, avec 31 Go de mémoire vive et un processeur 12th Gen Intel® Core™ i7-12700 × 20, ainsi qu'un processeur graphique Mesa Intel® UHD Graphics 770. Il est géré par un système d'exploitation Ubuntu 22.04 64-bits (distribution Linux), sous la version 42.9 de GNOME (GNU Network Object Model Environment) fournissant une interface utilisateur ergonomique pour interagir avec le système d'exploitation GNU (*GNU's Not Unix*). Il est ainsi composé d'un noyau Linux et de GNU, qui ensemble forment le système d'exploitation communément connu sous le nom "Linux". Le système de fenêtrage est "Wayland", un système chargé de l'affichage et du placement des fenêtres durant l'usage du système d'exploitation par l'utilisateur.

2.2 Pratique Professionnelle

2.2.1 Veille bibliographique et technologique

Pour la veille bibliographique, je me suis aidé du script de *Web-scraping* développé dans le cadre de mon stage, qui, à partir d'un fichier comprenant les titres de nombreux articles cible, m'a permis de récupérer les métadonnées associées de manière semi-automatique. Les journaux consultés (*Applied Soil Ecology*, *Ecosystems*, *Soil Biology and Biochemistry*, etc.) étaient majoritairement orientés sur les études écologiques. Un travail par recherche MESH a aussi été réalisé, en retenant les mots-clés **Oligochaetas**, **Earthworm**, **Lumbricus terrestris**, **Ecosystems**, **Introduced species**, **Soil**, **Data mining** et **Meta-Analysis**. Pour gérer cette bibliographie, j'ai recherché les fichiers BibTeX de chaque référence à l'aide du site spécialisé doi2bib.

2.2.2 Bonne pratique de programmation informatique et de développement logiciel

Pour l'écriture du code Python, j'ai utilisé l'éditeur de code Visual Studio Code (VScode). Pour la conception du script Rmd, j'ai travaillé sous l'IDE RStudio. Conformément aux

bonnes pratiques, l'intégralité des scripts (Python et R) ont été commentés en anglais, pour faciliter la réutilisabilité du code. Pour tester le fonctionnement de chaque script, deux approches différentes ont été mises en place :

Pour le script de *Web-scraping* en Python, les tests de fonctionnement effectués durant le développement ont été réalisées sur des jeux de données réduits (comprenant le plus souvent seulement les dix premiers articles) afin de pouvoir détecter rapidement si la sortie renvoyée était pertinente par rapport à la tâche initiale. Pour la récupération de données via Crossref, des tests ont aussi été réalisées, notamment pour l'exploration des structures JSON renvoyées par la méthode `works()` de l'API. Une fois les données recherchées obtenues en sortie de tests ponctuels, il a été possible de généraliser la méthode, en l'appliquant dans l'ensemble du script principal. Enfin, j'ai souvent travaillé sur deux scripts identiques mais distincts, le premier servant de script principal, tandis que le deuxième servait de support pour développer de nouvelles fonctionnalités, afin d'éviter la régression du code (perte accidentelle de fonctionnalité). De cette façon, le script principal n'était mis à jour que lorsque la sortie du script secondaire correspondait aux attentes. Un dépôt GitHub personnel aurait aussi pu jouer ce rôle, mais comme je travaillais seul sur cette partie, je n'en ai pas ressenti l'utilité.

Pour le script de *Text-mining* en Rmd sous RStudio, la plupart des tests de fonctionnement au cours du développement ont été réalisés dans la console R directement, afin d'éviter d'ajouter de nouveaux objets inutiles de "test" à l'environnement R. Les processus n'ont été intégrés au script Rmd proprement dit qu'une fois leur fonctionnement testé et validé dans la console. Dans cette démarche, le fonctionnement de l'environnement R a été très utile, car cela a permis de réemployer certains objets stockés en mémoire sans avoir à les redéfinir seulement en vue d'effectuer les tests.

2.2.3 Communication des travaux

En concertation avec mes encadrants, j'ai aussi utilisé un dépôt GitHub spécialement mis en place pour le projet entre eux et moi, où mon travail de chaque jour a pu être sauvegardé grâce à un mécanisme de Push/Pull. De cette façon, mes encadrants ont pu facilement suivre l'évolution de mon travail et évaluer la qualité des solutions proposées. Pour faciliter l'utilisation de cet outil, le logiciel GitHub Desktop m'a été présenté, une interface utilisateur graphique facilitant grandement la visualisation et l'usage du dépôt GitHub mis en place pour le projet. Grâce à la fonctionnalité Knitr de Rmd, j'ai pu rendre compte de ma progression quotidienne en produisant automatiquement un fichier rapport au format HTML, directement issu de mon code (figures produites sous R, titres et interprétation rédigées en Markdown ou HTML). Les réunions avec mes encadrants, le plus souvent hebdomadaires, ont été fixées par échange de mails ou bien organisées sur site. Elles m'ont permis de rester bien focalisé sur la mission en me fournissant des objectifs hebdomadaires clairs et précis, tout en permettant à mes encadrants d'être régulièrement informés de ma progression par rapport aux objectifs fixés. Vers la fin de mon stage, je présenterai aussi mes travaux au reste de l'équipe, selon la pratique en vigueur dans la structure d'accueil.

2.3 Outils informatiques et statistiques pour les différentes phases de vos travaux

2.3.1 Récupération des abstracts et des métadonnées avec Python

Pour l'élaboration du script Python, j'ai surtout utilisé le package Python Habanero Crossref pour requêter (en utilisant une API) la base de données Crossref, qui contient une grande

quantité de métadonnées (donnés sur les articles en tant que tel, comme l’abstract, les auteurs, etc.), afin de récupérer les informations voulue à partir d’un fichier texte contenant les titres des articles scientifiques cibles au sujet des vers de terre. Pour compléter les données récupérées (la base de Crossref comprenant de nombreuses données manquantes), j’ai aussi développé en parallèle un module pour récupérer les données d’intérêt dans le code source de ResearchGate (*web scraping*), comme le DOI (*Digital Object Identifier*) de chaque publication, la date de chargement sur la base de RG ou encore le lien vers la page RG correspondante. Pour parvenir à cette solution, voici la liste des modules qui ont été utilisés :

Pandas : Pandas est un module qui permet de manipuler facilement des tableaux de données avec des étiquettes de variables (colonnes) et d’individus (lignes). Il est notamment utilisé dans le script pour exporter les résultats issus du code Python vers un fichier CSV (comma separated values), lisible et modifiable à l’aide d’outils de bureautique courants comme LibreOffice Calc ou Excel.

Numpy : Package conçu pour le calcul scientifique avec Python. Il est très utile pour l’algèbre (comme par exemple pour la manipulation de matrices), et son implémentation en C, C++ et Fortran en fait un outil de calcul rapide et efficace pour l’analyse de données et le calcul scientifique. Dans le code, cela dit, il sert simplement à l’indexage lors de la création du DataFrame de résultats.

Itertools : Module implémentant des outils Python pour maîtriser plus subtilement les itérations. La méthode employée dans le code est *zip_longest*, qui permet de créer un DataFrame à partir d’une liste de listes de tailles potentiellement différentes. La plus longue sera employée en référence (longest), et toutes les autres seront ajustées à cette longueur par l’ajout d’une valeur de remplissage notée *"null"*, dans le code). Dans mon travail, elle a servi à créer le DataFrame requis à partir de listes de données de tailles pas forcément égales (à cause des valeurs manquantes).

Habanero : Module client de bas niveau pour interroger l’API Crossref, une base de données contenant les métadonnées des articles de tous les membres (des informations comme le titre, le nom d’auteur, le DOI etc.). Dans le code Python, elle est employée surtout pour rechercher les noms d’auteurs, les autres champs testés n’étant pas assez fiables pour automatiser complètement la récupération d’informations. Crossref est codé comme une **classe** du module Habanero, comprenant les méthodes *works()*, *members()*, *prefixes()*, *funders()*, *journal()*, *type()* et *licence()*. Dans le code Python, seule la méthode *works()* a été employée pour envoyer une requête à partir du titre de chaque article.

Unicodecode : Module contenant entre autres la fonction éponyme *unicodecode* (employée dans le script) conçue pour transformer les chaînes de caractères contenant des caractères non-ASCII (comme par exemple des idéogrammes chinois) pour les traduire en chaînes de caractères contenant uniquement des caractères ASCII. Dans le code, la fonction *unicodecode* est employée pour rendre l’affichage des noms d’auteur contenant des caractères non-ASCII. Certaines corrections sont imparfaites et mène à des erreurs d’encodage dans le fichier de résultats.

RGS2 : Sous-module codé localement à partir d’un exemple trouvé en ligne^[7], par la suite adapté pour récupérer directement les informations voulues dans le code source du site scientifique ResearchGate, les autres options potentielles (comme Google Scholar) ayant souvent un système de détection et de blocage des bots. Seule la deuxième version du sous-module a été retenue dans le projet final. Il dépend des modules suivants :

1. Module **Parsel**, fonction *Selector* : Module facilitant l’extraction des données pour les formats HTML, JSON et XML. Dans le code, il est utilisé pour trouver les

- données recherchées directement dans le code source de la page (*web scraping*) en s'appuyant sur des sélecteurs CSS.
2. Module **playwright.sync_api**, fonction *sync_playwright* : Module permettant de lancer une session navigateur depuis un script Python. Dans le code, il est utilisé pour se rendre sur le site de ResearchGate via une session Chromium, un navigateur libre développé par Google.
 3. Module **re** : Module fournissant des opérations sur les expressions rationnelles utilisable dans un code Python. Dans le code Python, il est utilisé pour filtrer les résultats HTML bruts issus du *Web scraping*.
 4. Module **time**, fonction *sleep* : Module fournissant différentes fonctions liées au temps. Dans le script, la fonction "sleep(t)" est utilisée pour forcer le système à ne rien faire pendant t secondes, évitant de cette façon de surcharger le serveur cible de requêtes trop rapides et trop nombreuses.

2.3.2 Text Mining avec R

Pour réaliser le *Text-mining*, j'ai utilisé un script R (développé sous Rstudio en Rmd) pour traiter les textes et l'analyser sous forme de figures. Mon travail a donc consisté à adapter les codes R montrés en exemple sur des livres à mes propres données (abstracts d'articles scientifiques), structurées différemment. Les analyses portaient par exemple sur la fréquence des mots, au global et pour chaque MA, ou encore sur l'analyse de sentiment reposant sur l'attribution d'un score positif (+1) ou négatif (-1) à chaque mot du corpus. Cette attribution de sentiment a pu être réalisée grâce à un dictionnaire R conçu pour relier un *token* donné à la valence (positive ou négative) qui lui correspond. Afin de mieux visualiser les résultats, différentes figures ont été réalisées, comme des diagrammes en barre ou des nuages de mots. Le script R développé repose sur les librairies suivantes :

Librairie dplyr : Librairie R conçue pour faciliter la manipulation de larges jeux de données (DataFrame et Tibble), avec des fonctions spcialisées comme *mutate* (ajout de variables), *select* (sélectionner les variables à partir de leurs noms), *filter* (filtrer des cellules selon leur valeur), *summarise* (résumer les informations d'un tibble dans un format très synthétique) et *arrange* (pour réordonner les lignes selon l'ordre / la variable voulue.)

Librairie ggplot2 : Librairie R pour créer déclarativement des graphiques divers et variés (barplots, histogrammes, scatter plots, etc.)

Librairie tidytext : Librairie R conçue pour faciliter l'analyse de texte^[4], se fondant sur le paradigme d'analyse de données "tidy", où chaque variable est une colonne, chaque observation une ligne et chaque ensemble d'observations est un tableau. La fonction la plus utilisée dans le cadre de cette analyse est *unnest_tokens*, qui permet de transformer un texte donné en une tables d'unités textuelles plus réduites (*tokens*), comme des phrases ou des mots.

Librairie Knitr : Librairie R conçue pour récupérer automatiquement l'output d'un code R (par exemple, pour produire une figure) afin de l'inclure dans un autre document (par exemple, format Word, HTML ou PDF) qui contiendra aussi la prose écrite par l'auteur du document, souvent pour interpréter ou commenter des résultats. Cette librairie permet notamment de moduler plus finement l'affichage des résultats, en permettant par exemple de ne pas afficher certaines figures dans un format de sortie donné (pour masquer une figure interactive au format HTML que l'on ne souhaite pas forcément voir apparaître dans un document PDF, par exemple).

Librairie SnowballC : Librairie R implémentant Snowball, un langage conçu pour gérer les chaînes de caractères, les nombres entiers et les booléens. Dans le script R de Text Mining, il a servi à transformer les *tokens* pour ne garder que la racine de chaque mot (processus de *racination*), afin d'éviter les erreurs de comptage (si, pour un humain, "enhance" et "enhancement" sont deux mots ayant approximativement le même sens, informatiquement ce sont deux chaînes de caractères distinctes).

Librairie grid : Librairie R implémentant les fonctions graphiques primitives qui sous-tendent le package ggplot2. Elles permettent de modifier certains détails des graphes produits grâce à ggplot2, offrant ainsi un meilleur contrôle du rendu visuel du résultat obtenu. Dans le script, il est employé, par exemple, pour spécifier exactement l'aspect des flèches composant le réseau de bigrammes.

Librairie ggraph : Librairie R formant une extension de ggplot2, conçue pour permettre de supporter les structures de données relationnelles comme les réseaux, les graphes et les arbres. Dans le script, elle est notamment employée pour produire les réseaux de bigrammes.

Librairie igraph : fonction *graph_from_data_frame()* : Cette fonction appartient à la librairie igraph. Elle crée un objet graphe à partir d'un data frame, où le data frame représente les arêtes entre les nœuds.

Librairie egg : fonction *ggarrange()* : Librairie servant à organiser différents graphes sur une seule et même figure. Dans le script, c'est l'une des librairies employées pour comparer les quatre MA entre elles.

Librairie tidyr : Librairie R implémentant des méthodes utiles pour manipuler des objets de type "tidy". Elle comprend des fonction telles que *pivot_wider* et *pivot_longer* (pour convertir le dataFrame d'un format à un autre), ou encore *bind_rows()*, pour combiner des DataFrames par lignes.

2.4 Données

L'approche "tidy" choisie repose sur la *tokenisation* du corpus de texte en unités textuelles (appelées "*tokens*") plus petites, comme des mots, des phrases ou des *n-grams*. J'ai pour cela pu m'inspirer du livre rédigé par Julia Slige (data scientist) et David Robinson (Directeur de Data Scientist de la plateforme Heap)^[4]. Afin de filtrer les mots d'intérêt seulement, deux stratégies ont été employées : Premièrement, un filtrage brut de tous les mots de liaisons sans rapport direct avec le sujet (nommés "stop words" en *Text-mining*, des mots tels que "the", "and", "is", "of" etc. en anglais) pour ne conserver que les mots sur lesquels les analyses pourront donner des résultats scientifiques significatifs (savoir que le mot "le" est le plus fréquent dans un corpus de texte français ne signifie rien sur le plan biologique). Deuxièmement, une autre approche a été de considérer la fréquence de chaque mot par rapport au nombre total de mots présents dans le corpus ($n \text{ mot} / N \text{ mots}$). De cette façon, les mots les plus fréquents perdent une partie de leur poids statistique, tandis que les mots plus rares en gagnent. On peut ainsi visualiser facilement les mots les plus importants d'un texte, sans même avoir besoin de modifier les données au préalable avec une liste de stop words. Cela permet de conserver le texte dans son ensemble, évitant ainsi un potentiel biais pouvant perturber l'analyse.

Chapitre 3

Résultats

3.1 *Web-scraping* et obtention d'une base de métadonnées en format CSV

En sortie du code de *Web-scraping* développé en Python (et après recherche manuelle des données manquantes), un fichier de métadonnées au format CSV, contenant la métaanalyse d'origine, le titre, les auteurs, l'abstract, la date de publication sur ResearchGate, le DOI et l'URL vers la page RG correspondante a été obtenue (tableau 3.1).

MA	Title	First author	Last author	Abstract	Date	DOI	URL
MA1	Influence of exotic earthworm invasion on soil organic matter, microbial biomass and denitrification potential in forest soils of the northeastern United States.	Amy E Burtelow	Peter M Groffman	Formerly glaciated regions of the northeastern United States have few native earthworm species and the region is dominated by exotic earthworms from Europe and Asia [...].	Sep 1998	10.1016/S0929-1393(98)00075-4	URL

TABLE 3.1 — Tableau présentant un exemple de la structure type du fichier CSV issu du Web-scraping et employé pour l'analyse. Le fichier originel comprend 168 lignes. Les données manquantes (non représentées ici), sont notées "NA".

3.2 Analyse globale des textes bruts

Pour détecter informatiquement les racines les plus fréquentes du corpus, chaque *token* distinct (dans cette section, des mots uniques) présent dans les quatre jeux de données (abstracts d'articles) a été compté. Le graphique ci-dessous présente les racines les plus fréquemment retrouvées dans l'ensemble du corpus. Les trois plus fréquentes sont **earthworm**, **soil** et **plant** (figure 3.1).

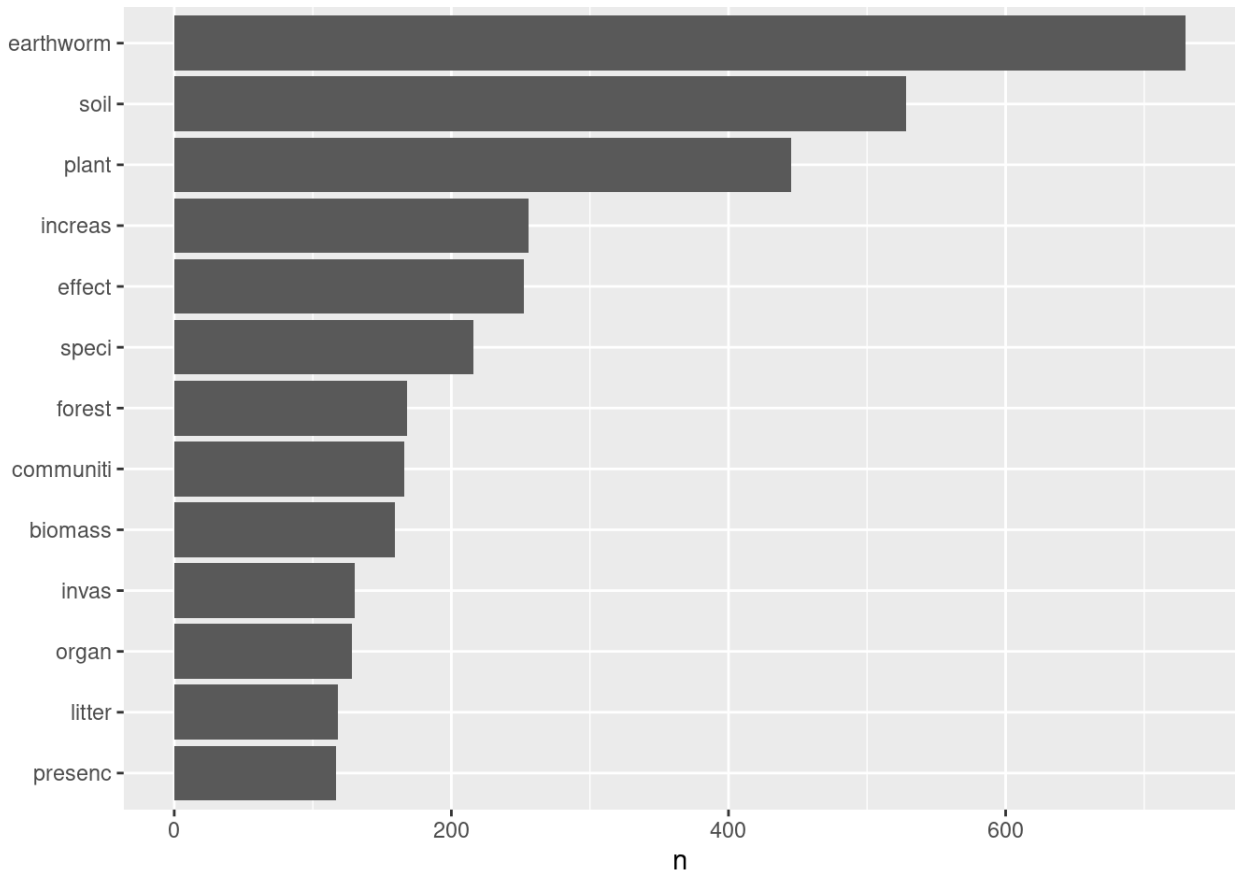


FIGURE 3.1 — Sur l'axe des abscisses, les valeurs numériques représentent le nombre total d'occurrences des mots (après *racination*) écrits en ordonnées. On remarque que les trois mots les plus fréquents dans les articles des quatre métaanalyses confondues sont **earthworm**, **soil** et **plant**.

Dans un second temps, un log-log scatter plot a été réalisé pour analyser les corrélations du choix des mots entre différents textes (figure 3.2). Il en ressort que le terme "earthworm" est très fréquent dans l'ensemble des métaanalyses, car présent dans le quart supérieur droit de tous les graphes. Le terme "species" est lui aussi assez souvent rencontré. Pour la MA1, 2 et 3, on peut remarquer l'importance du terme "soil", présent dans le quart supérieur droit. D'après ce graphique, les termes "carbon", "invasive" et "mineral" sont plus fréquents dans la MA1 que dans la MA2, là où les termes "plant", "community" et "growth" sont plus fréquents dans la MA2. De la même façon, les termes "carbon" et "mineral" sont plus fréquents dans la MA1 que dans la MA3, là où les termes "disturb", "density" ou "plant" sont plus fréquents dans la MA3. Enfin, les termes "forest", "invasive" et "exotic" sont plus fréquents dans la MA1 que dans la MA4, alors que les termes "plant", "growth" et "fertile" sont plus fréquents dans la MA4.



FIGURE 3.2 — Log-log scatter plot montrant les corrélations du choix des mots entre chaque métaanalyse, à l'échelle logarithmique. La ligne grise représente les mots dont la fréquence est similaire dans les deux textes comparés. Les mots au-dessus de la ligne sont plus fréquents dans la MA1, tandis que les mots en-dessous de la ligne sont plus fréquents dans la deuxième MA. Les termes situés dans le quart supérieur droit sont très fréquents dans les deux corpus.

3.3 Rang, approche tf-idf et loi de Zipf

En analyse comparative de texte, il est possible de calculer le **rang** de chaque mot. Dans les tableaux ci-dessous (figure 3.3), on remarque que trois premiers rangs sont toujours occupés par des mots-outils, c'est à dire non sémantiquement pleins. Le mot "earthworm" fait partie des dix premiers de toutes les MA. Le mot "soil" est dans les dix premiers pour les MA1, 3 et 4, là où "plant" est dans les dix premiers pour les MA2, 3 et 4. Pour la MA1, le champ lexical des végétaux est représenté par le mot "forest".

```
> freq_and_rank %>% filter(MA=='MA1')
# A tibble: 1,348 × 6
  MA word      n Total word_frequency rank
  <chr> <chr>    <int> <int>          <dbl> <int>
1 MA1 the      488 5616          0.0862 1
2 MA1 and      456 5616          0.0812 2
3 MA1 of       425 5616          0.0757 3
4 MA1 in       310 5616          0.0552 4
5 MA1 earthworm 302 5616          0.0538 5
6 MA1 soil     267 5616          0.0475 6
7 MA1 to       129 5616          0.0230 7
8 MA1 a        121 5616          0.0215 8
9 MA1 forest   119 5616          0.0212 9
10 MA1 effect   98 5616          0.0175 10
# i 1,338 more rows

> freq_and_rank %>% filter(MA=='MA2')
# A tibble: 911 × 6
  MA word      n Total word_frequency rank
  <chr> <chr>    <int> <int>          <dbl> <int>
1 MA2 the      261 2959          0.0882 1
2 MA2 of       248 2959          0.0838 2
3 MA2 and      225 2959          0.0760 3
4 MA2 plant    170 2959          0.0575 4
5 MA2 in       158 2959          0.0534 5
6 MA2 earthworm 144 2959          0.0487 6
7 MA2 by        84 2959          0.0284 7
8 MA2 a         78 2959          0.0264 8
9 MA2 on        63 2959          0.0213 9
10 MA2 increas   63 2959          0.0213 10
# i 901 more rows

> freq_and_rank %>% filter(MA=='MA3')
# A tibble: 474 × 6
  MA word      n Total word_frequency rank
  <chr> <chr>    <int> <int>          <dbl> <int>
1 MA3 the      495 1860          0.266 1
2 MA3 of       442 1860          0.238 2
3 MA3 and      362 1860          0.195 3
4 MA3 in       282 1860          0.152 4
5 MA3 earthworm 200 1860          0.108 5
6 MA3 plant    185 1860          0.0995 6
7 MA3 a        171 1860          0.0919 7
8 MA3 soil     161 1860          0.0866 8
9 MA3 to       133 1860          0.0715 9
10 MA3 increas 109 1860          0.0586 10
# i 464 more rows

> freq_and_rank %>% filter(MA=='MA4')
# A tibble: 1,447 × 6
  MA word      n Total word_frequency rank
  <chr> <chr>    <int> <int>          <dbl> <int>
1 MA4 the      495 4985          0.0993 1
2 MA4 of       442 4985          0.0887 2
3 MA4 and      362 4985          0.0726 3
4 MA4 in       282 4985          0.0566 4
5 MA4 earthworm 200 4985          0.0401 5
6 MA4 plant    185 4985          0.0371 6
7 MA4 a        171 4985          0.0343 7
8 MA4 soil     161 4985          0.0323 8
9 MA4 to       133 4985          0.0267 9
10 MA4 increas 109 4985          0.0219 10
# i 1,437 more rows
```

FIGURE 3.3 — Impression écran issue de Rstudio, montrant les dix mots les plus représentés dans chaque MA. Ces résultats ont été obtenu par filtration d'un tibble R plus large nommé "freq_and_rank". On remarque que les premiers rangs sont souvent occupés par des mots fonctionnels, comme "the", "of", "and" et "in".

L'objectif de l'approche tf-idf est d'évaluer l'importance d'un terme contenu dans un document, relativement à l'ensemble du document. La fréquence brute du terme est nommée tf, et la fréquence inverse de document (une mesure de l'importance du terme dans l'ensemble du corpus) est nommée idf. Le poids ajusté de chaque terme s'obtient en multipliant ces deux mesures, et il augmente proportionnellement au nombre d'occurrences du mot dans le document. Grâce à l'approche tf-idf (figure 3.4), on peut voir que chaque mot présent sur l'axe Oy est un mot utilisé plus souvent que les autres au sein d'une métaanalyse. Les racines **earthworm**, **soil**, **forest** sont plus importantes dans la MA1, **aphid**, **nematod**, **herbivor** plus importantes dans la MA2, **eastern**, **canopy**, **arthropod** plus importantes dans la MA3 et **grain**, **pb**, **cu** plus importantes dans la MA4.

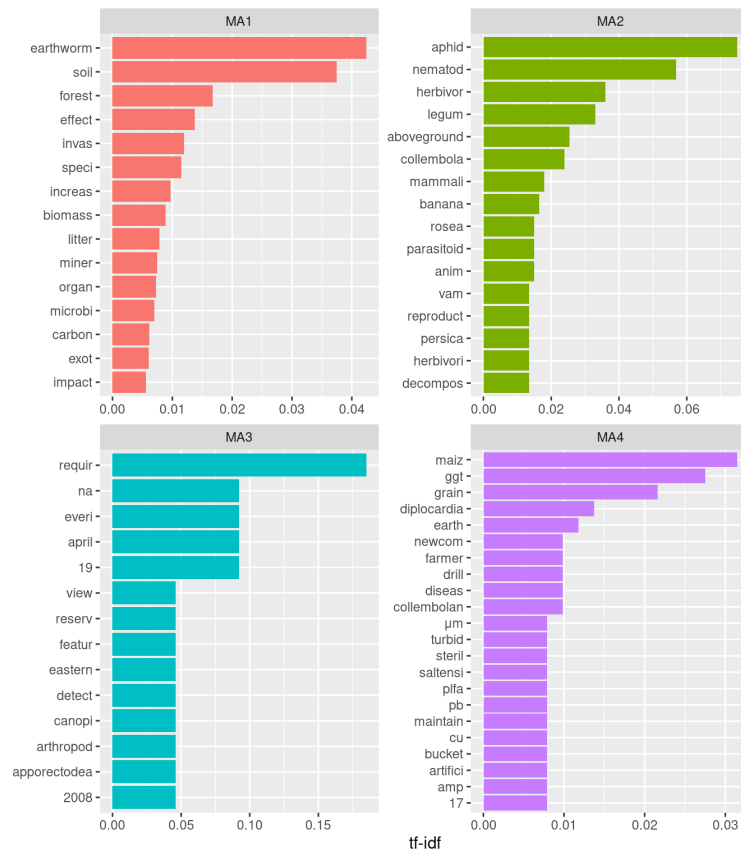


FIGURE 3.4 — Mesures de tf-idf pour chacune des métaanalyses du corpus. Les mots écrits sur l’axe Oy sont plus importants dans leurs MA respectives que dans les autres.

En linguistique, la loi de Zipf (figure 3.5), stipule que $tf \propto \frac{1}{rang}$. La MA3 (courbe grise) contient davantage de mots “rares” que la valeur prédite par le modèle linéaire (quart supérieur gauche), alors que toutes les autres en contiennent moins. La MA3 est donc celle qui contient la plus haute fréquence de mots rares. Pour ce qui est des mots très fréquents (quart inférieur droit sous 0.001), on peut constater que toutes les métaanalyses en contiennent moins que la valeur prédite.

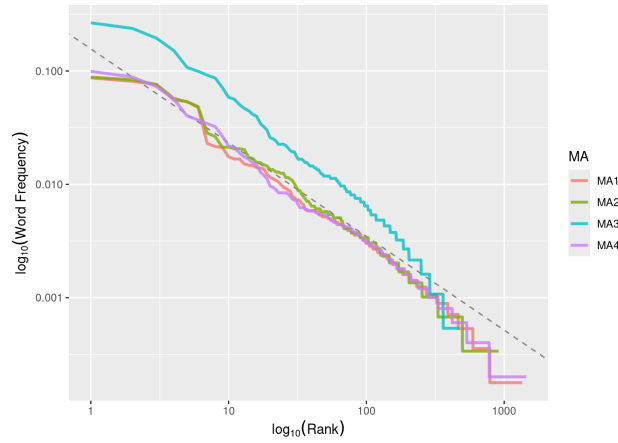


FIGURE 3.5 — La courbe grise en pointillés représente la valeur attendue selon un modèle linéaire sur l'intervalle [11,99]. Sa pente est de -0.8272.

3.4 Analyses de bigrammes seuls et réseaux de bigrammes

Grâce à l'approche tf-idf (figure 3.6), on peut voir que chaque *bigramme* présent sur l'axe Oy est particulièrement important pour la métaanalyse cconcernée. Ainsi, **mineral soil**, **earthworm invasion**, **exotic earthworms** sont plus importants que les autres dans la MA1, **soil organisms**, **plant responses**, **plant mediated** plus importants que les autres dans la MA2, **earthworm invasions**, **native earthworm**, **species richness** plus importants que les autres dans la MA3 et **soil organisms**, **soil fertility**, **dry weight** plus importants que les autres dans la MA4.

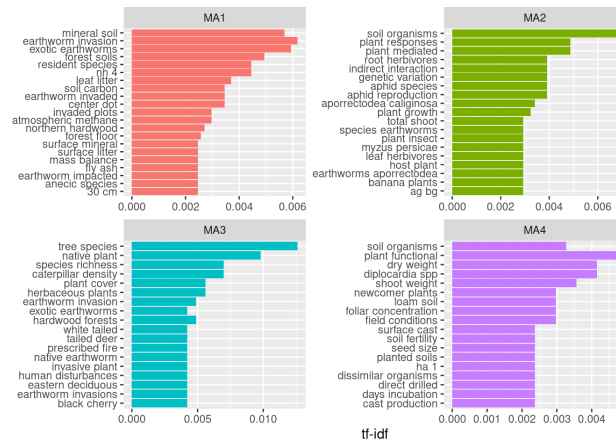


FIGURE 3.6 — Mesures de tf-idf pour chacune des métaanalyses du corpus. Les **bigrammes** écrits sur l'axe Oy sont plus importants dans leurs MA respectives que dans les autres.

Grâce à l'approche en réseau de bigrammes, il est possible de relever des associations de mots au-delà de deux mots seulement. Le réseau ci-dessous (figure 3.7) montre que les trois mots les plus importants du corpus sont "**plant**", "**soil**" et "**earthworm**". Pour le premier, les bigrammes rencontrés sont "*plant community*", "*plant growth*", "*plant communities*", "*native plant*", pour le second "*mineral soil*", "*soil microbial*" / "*microbial biomass*" (*forte association*), "*soil organisms*" et pour le troisième "*earthworm invasion*", "*earthworm species*", "*earthworm activity*". D'autres bigrammes notables sont : "*forest floor*", "*exotic earthworm*", "*organic matter*", "*leaf litter*", "*northern hardwood*" et "*lumbricus terrestris*".

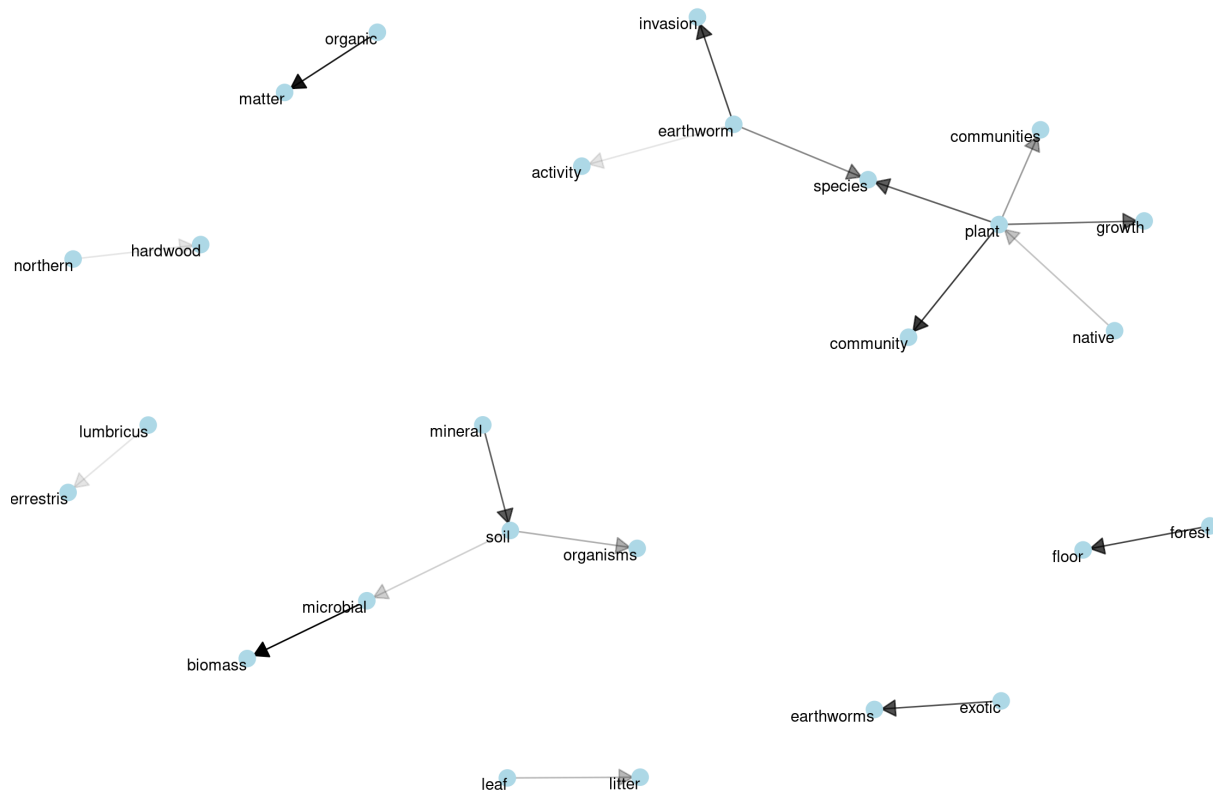


FIGURE 3.7 — Réseau de bigrammes montrant les connexions les plus fréquentes reliant les mots du corpus. Chaque mot est représenté par un nœud du réseau (**bleu clair**) et les connexions entre les mots sont orientées. L'opacité des flèches est proportionnelle à la fréquence d'occurrence du bigramme.

Il peut aussi être intéressant d’analyser l’aspect des réseaux obtenus pour chacune des métaanalyses individuellement (figures 3.8 à 3.11). Pour la MA1 (figure 3.8), le seul centre retrouvé est **"earthworm"**, impliqué dans les connexions *"earthworm species"* et *"earthworm invasion"*. D’autres bigrammes notables sont : *"mineral soil"* → *"soil microbial"* → *"microbial biomass"*, ou encore *"exotic earthworms"*, *"organic matter"*, *"forest floor"*, *"leaf litter"*, *"northern hardwood"*, *"lumbricus terrestris"*

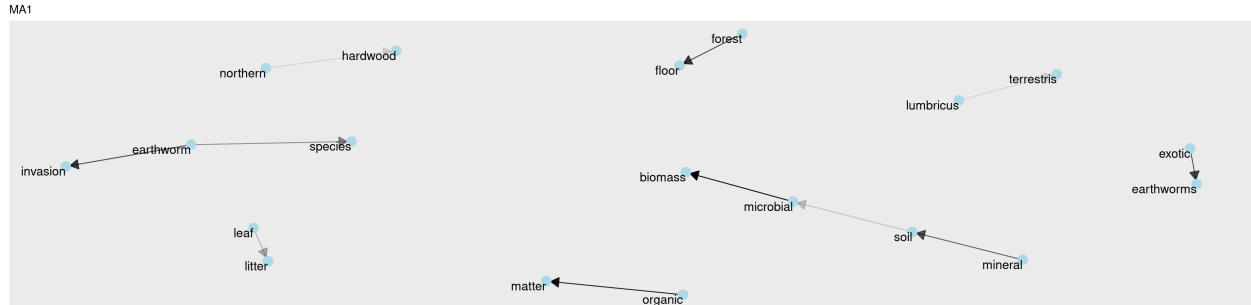


FIGURE 3.8 — Réseau de bigrammes montrant les connexions les plus fréquentes reliant les mots de la MA1. Chaque mot est représenté par un noeud du réseau (**bleu clair**) et les connexions entre les mots sont orientées. L’opacité des flèches est proportionnelle à la fréquence d’occurrence du bigramme.

Pour la MA2 (figure 3.9), le seul centre retrouvé est **"plant"**, impliqué dans les connexions *"plant species"*, *"plant growth"* et *"plant communities"*. "Soil organisms" est aussi un bigramme notable.

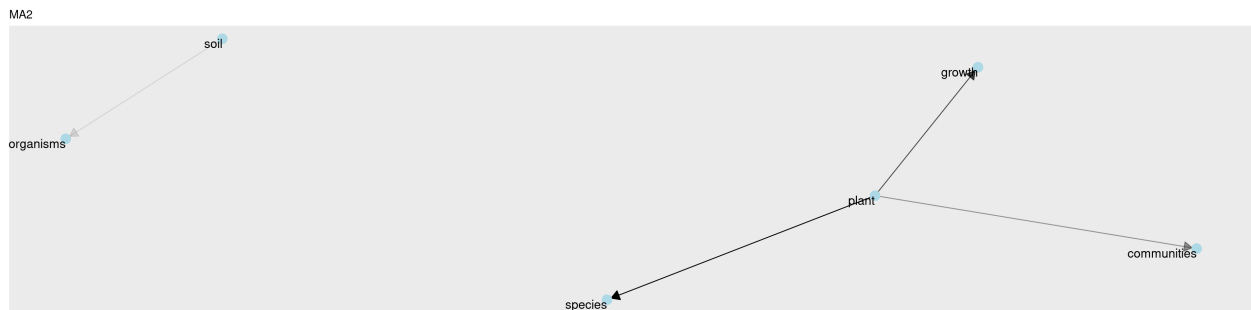


FIGURE 3.9 — Réseau de bigrammes montrant les connexions les plus fréquentes reliant les mots de la MA2. Chaque mot est représenté par un noeud du réseau (**bleu clair**) et les connexions entre les mots sont orientées. L’opacité des flèches est proportionnelle à la fréquence d’occurrence du bigramme.

Pour la MA3 (figure 3.10), aucun centre n'est retrouvé. Par contre, ce réseau est constitué d'une chaîne de mots qui se suivent : "native plants" → "plant species" → "earthworm species".

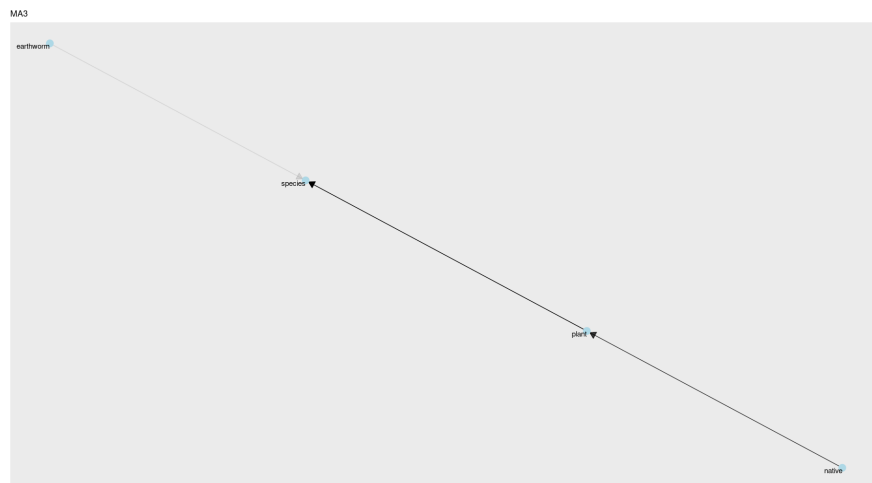


FIGURE 3.10 — Réseau de bigrammes montrant les connexions les plus fréquentes reliant les mots de la MA3. Chaque mot est représenté par un nœud du réseau (**bleu clair**) et les connexions entre les mots sont orientées. L'opacité des flèches est proportionnelle à la fréquence d'occurrence du bigramme.

Pour la MA4 (figure 3.11), le seul centre retrouvé est "**plant**", impliqué dans les connexions "*plant species*", "*plant growth*" et "*plant communities*". "Soil organisms" et "earthworm activity" sont aussi des bigrammes notables.

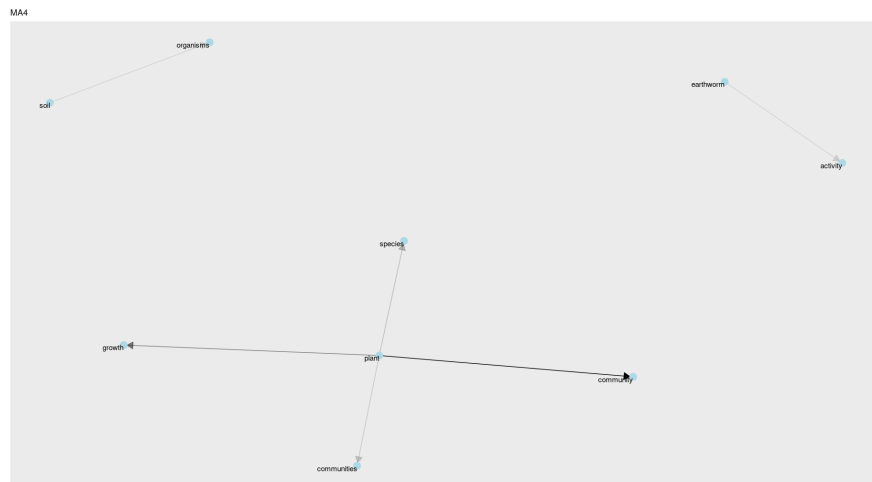


FIGURE 3.11 — Réseau de bigrammes montrant les connexions les plus fréquentes reliant les mots de la MA4. Chaque mot est représenté par un nœud du réseau (**bleu clair**) et les connexions entre les mots sont orientées. L'opacité des flèches est proportionnelle à la fréquence d'occurrence du bigramme.

3.5 Analyse de sentiment

L'objectif de cette série d'analyses est de détecter informatiquement l'avis émis par les abstracts composant le corpus de texte à propos des vers de terre. Pour cela, on s'appuie sur des dictionnaires R prédéfinis faisant correspondre à chaque mot un score chiffré, +1

Chapitre 4

Discussion

Pour la constitution de la base de données de travail (le fichier CSV contenant les abstracts des articles à analyser), une approche semi-automatique a pu être implémentée, mais il serait encore nécessaire d'améliorer l'outil pour pouvoir complètement automatiser le processus. Cependant, de nombreux sites internet combattant activement le *Web-scraping*, il pourrait s'avérer difficile d'obtenir un script à la fois précis, fonctionnel et stable dans le temps. Les abstracts d'articles se sont avérés étonnamment difficiles à récupérer automatiquement, c'est pourquoi une intervention manuelle a souvent été nécessaire.

Une fois les abstracts récupérés, une première analyse de la fréquence brute des termes à montré que les mots les plus retrouvés dans le corpus entier étaient "earthworm", "soil" et "plant", ce qui semble confirmer que tous ces textes traitaient bel et bien de l'influence des vers de terre sur le sol et les plantes.

Individuellement, ils ressort que la MA1 étudie davantage l'écologie des sols que les MA2 et 3, là où les MA2 et 3 semblent davantage étudier les plantes en tant que telles. De la même façon, la MA1 parle davantage d'invasion et d'espèces exotiques que la MA4, qui elle se focalise plus sur la croissance des plantes et la fertilité des sols.

L'étude du thème de chaque corpus par l'approche tf-idf sur des mots uniques confirme que toutes les métaanalyses mentionnent très souvent les vers de terre, les sols et les plantes. La MA2 semblent mentionner les sols moins souvent, cependant, se focalisant plus, comme noté précédemment, sur les réactions physiologiques des plantes en tant que telles face aux pucerons et aux herbivores. Thématiquement, elle semble cependant se rapprocher de la MA4, qui étudie l'effet des vers de terre sur la production végétale et la pollution des sols. Les MA1 et 3 semblent davantage orientées vers la caractérisation de l'influence des vers de terre sur les écosystèmes natifs, notamment les écosystèmes forestiers.

L'application de la loi de Zipf aux différents textes du corpus montre que la métaanalyse comportant le plus de mots "rares" est la MA3, et cela pourrait s'expliquer par la taille plus petite de son corpus d'abstracts par rapport aux autres (avec 13 articles seulement contre 40 pour la MA1, 21 pour la MA2 et 42 pour la MA4, c'est le sous-groupe le plus réduit de la base de données), ce qui fait artificiellement monter la fréquence observée. Cependant, aucun effet notable n'est relevé pour la MA2, qui pourtant s'appuie aussi sur un corpus plutôt réduit.

Pour aller plus loin dans l'étude thématique, la méthode tf-idf a ensuite été appliquée à un ensemble de bigrammes, et non plus uniquement des mots simples. Il en ressort que la MA1 étudie l'influence des vers de terre invasifs sur les sols forestiers, la MA2 l'influence des organismes du sol sur les plantes, la MA3 l'influence des espèces invasives sur la richesse spécifiques des écosystèmes, et la MA4 l'influence des vers de terre sur la pollution des sols et la productivité végétale. L'analyse des réseaux de bigrammes tend aussi à rapprocher thématiquement la MA1 de la MA3, et la MA2 de la MA4.

Enfin, l'analyse de sentiment conduite sur les mots simples a semblé montrer les MA1 et 3 comme étant globalement négatives, et les MA2 et 4 comme étant globalement positives, contribuant à renforcer les deux groupes dégagés au terme des étapes précédentes.

Il semblerait donc que les vers de terre invasifs ont un impact plutôt négatif sur les écosystèmes natifs (MA1 et 3), mais un impact plutôt positif sur les sols et les plantes (MA2 et 4). Ces résultats semblent en accord avec les recherches conduites précédemment par McLean u. a. (2006) et Hodson u. a. (2023). Cependant, cette analyse portant sur de nombreux textes théoriquement très différents les uns des autres, il convient de noter que les groupes ainsi formés ne sont pas absolus, les résultats obtenus montrant simplement l'existence de tendances globales au sein des différents textes analysés.

Chapitre 5

Conclusion

Au cours de ce stage, l'objectif principal fixé par mes encadrants était de savoir s'il était possible d'identifier des groupes distincts d'articles scientifiques défendant une conception opposée du rôle écologique du vers de terre au sein du corpus d'articles scientifiques divers rien qu'à partir des méthodes de *Text-mining* sous R, c'est à dire sans avoir à lire et analyser chacun des nombreux abstracts individuellement. Les résultats obtenus à l'issue de ce stage semblent montrer qu'il est effectivement possible de détecter des tendances générales et des avis au sein de ces données textuelles en quantité, mais les deux groupes ainsi formés (favorables au vers de terres ou bien défavorables aux vers de terre), ne paraissent pas indiscutablement séparés l'un de l'autre. Cependant, la détection, même imprécise, de tels groupes par des méthodes statistiques prouve que la polémique entourant le rôle écologique du vers de terre au sein de la littérature scientifique existe bel et bien, et que des analyses informatiques relativement simples à mettre en place permettent de la mettre en évidence. Pour améliorer la qualité des résultats obtenus, il pourrait être nécessaire de poursuivre les analyses à l'aide de méthodes statistiques plus raffinées, ou bien tenter de confronter ces conclusions avec d'autres résultats comparables, idéalement obtenus grâce à des protocoles différents.

Bibliographie

- [Bertrand u. a. 2015] BERTRAND, Michel ; BAROT, Sébastien ; BLOUIN, Manuel ; WHALEN, Joann ; OLIVEIRA, Tatiana de ; ROGER-ESTRADE, Jean : Earthworm services for cropping systems. A review. In : *Agronomy for Sustainable Development* 35 (2015), Januar, Nr. 2, S. 553–567. – URL <http://dx.doi.org/10.1007/s13593-014-0269-7>. – ISSN 1773-0155
- [Bohlen u. a. 2004] BOHLEN, Patrick J. ; SCHEU, Stefan ; HALE, Cindy M. ; MCLEAN, Mary A. ; MIGGE, Sonja ; GROFFMAN, Peter M. ; PARKINSON, Dennis : Non-native invasive earthworms as agents of change in northern temperate forests. In : *Frontiers in Ecology and the Environment* 2 (2004), Oktober, Nr. 8, S. 427–435. – URL [http://dx.doi.org/10.1890/1540-9295\(2004\)002\[0427:NIEAAO\]2.0.CO;2](http://dx.doi.org/10.1890/1540-9295(2004)002[0427:NIEAAO]2.0.CO;2). – ISSN 1540-9295
- [Edwards und Arancon 2022] EDWARDS, Clive A. ; ARANCON, Norman Q. : *The Role of Earthworms in Organic Matter and Nutrient Cycles*. S. 233–274. In : *Biology and Ecology of Earthworms*, Springer US, 2022. – URL http://dx.doi.org/10.1007/978-0-387-74943-3_8. – ISBN 9780387749433
- [Ferlian u. a. 2017] FERLIAN, Olga ; EISENHAUER, Nico ; AGUIRREBENGOA, Martin ; CAMARA, Mariama ; RAMIREZ-ROJAS, Irene ; SANTOS, Fábio ; TANALGO, Krizler ; THAKUR, Madhav P. : Invasive earthworms erode soil biodiversity : A meta-analysis. In : *Journal of Animal Ecology* 87 (2017), September, Nr. 1, S. 162–172. – URL <http://dx.doi.org/10.1111/1365-2656.12746>. – ISSN 1365-2656
- [Forey u. a. 2023] FOREY, Oswaldo ; SAUZE, Joana ; PIEL, Clément ; GRITTI, Emmanuel S. ; DEVIDAL, Sébastien ; FAEZ, Abdelaziz ; RAVEL, Olivier ; NAHMANI, Johanne ; ROUCH, Laly ; BLOUIN, Manuel ; PÉRÈS, Guénola ; CAPOWIEZ, Yvan ; ROY, Jacques ; MILCU, Alexandru : Earthworms do not increase greenhouse gas emissions (CO₂ and N₂O) in an ecotron experiment simulating a three-crop rotation system. In : *Scientific Reports* 13 (2023), Dezember, Nr. 1. – URL <http://dx.doi.org/10.1038/s41598-023-48765-3>. – ISSN 2045-2322
- [Hodson u. a. 2023] HODSON, M.E. ; BRAILEY-JONES, P. ; BURN, W.L. ; HARPER, A.L. ; HARTLEY, S.E. ; HELGASON, T. ; WALKER, H.F. : Enhanced plant growth in the presence of earthworms correlates with changes in soil microbiota but not nutrient availability. In : *Geoderma* 433 (2023), Mai, S. 116426. – URL <http://dx.doi.org/10.1016/j.geoderma.2023.116426>. – ISSN 0016-7061
- [Kim u. a. 2017] KIM, Young-Nam ; ROBINSON, Brett ; LEE, Keum-Ah ; BOYER, Stephane ; DICKINSON, Nicholas : Interactions between earthworm burrowing, growth of a leguminous shrub and nitrogen cycling in a former agricultural soil. In : *Applied Soil Ecology* 110 (2017), Februar, S. 79–87. – URL <http://dx.doi.org/10.1016/j.apsoil.2016.10.011>. – ISSN 0929-1393

- [Kumar u. a. 2023] KUMAR, Rahul; YADAV, Renu; GUPTA, Rajender K.; YODHA, Kiran; KATARIA, Sudhir K.; KADYAN, Pooja; SHARMA, Pooja; KAUR, Simran : The Earthworms : Charles Darwin's Ecosystem Engineer. In : HAKEEM, Khalid R. (Hrsg.) : *Organic Fertilizers*. Rijeka : IntechOpen, 2023, Kap. 13. – URL <https://doi.org/10.5772/intechopen.1001339>
- [Lemtiri u. a. 2014] LEMTIRI, Aboulkacem; COLINET, Gilles; ALABI, Taofic; CLUZEAU, Daniel; ZIRBES, Lara; HAUBRUGE, Eric; FRANCIS, Frédéric : Impacts of earthworms on soil components and dynamics. A review. In : *Biotechnologie, Agronomie, Société et Environnement / Biotechnology, Agronomy, Society and Environment* (2014), November. – URL <https://hal-bioemco.ccsd.cnrs.fr/ECOBIO-RBPE/hal-01084235v1>
- [Loss u. a. 2013] LOSS, Scott R.; HUEFFMEIER, Ryan M.; HALE, Cindy M.; HOST, George E.; SJERVEN, Gerald; FRELICH, Lee E. : Earthworm Invasions in Northern Hardwood Forests : a Rapid Assessment Method. In : *Natural Areas Journal* 33 (2013), Januar, Nr. 1, S. 21–30. – URL <http://dx.doi.org/10.3375/043.033.0103>. – ISSN 0885-8608
- [Lubbers u. a. 2013] LUBBERS, Ingrid M.; VAN GROENIGEN, Kees J.; FONTE, Steven J.; SIX, Johan; BRUSSAARD, Lijbert; VAN GROENIGEN, Jan W. : Greenhouse-gas emissions from soils increased by earthworms. In : *Nature Climate Change* 3 (2013), Nr. 3, S. 187–194
- [McLean u. a. 2006] MCLEAN, M. A.; MIGGE-KLEIAN, S.; PARKINSON, D. : Earthworm invasions of ecosystems devoid of earthworms : effects on soil microbes. In : *Biological Invasions* 8 (2006), Juli, Nr. 6, S. 1257–1273. – URL <http://dx.doi.org/10.1007/s10530-006-9020-x>. – ISSN 1573-1464
- [Sharma u. a. 2017] SHARMA, D. K.; TOMAR, S.; CHAKRABORTY, D. : Role of earthworm in improving soil structure and functioning. In : *Current Science* 113 (2017), Nr. 6, S. 1064–1071. – URL <http://www.jstor.org/stable/26494167>. – Zugriffsdatum : 2024-06-13. – ISSN 00113891

Sitographie

- [1] BILITY. (s. d.). *Définition Markdown - Bility - Agence de développement web sur-mesure*. Consulté le 16 avril 2024 à l'adresse : <https://bility.fr/definition-markdown/>
- [2] doi2bib. (s. d.). Emploi fréquent à l'adresse : <https://www.doi2bib.org/>
- [3] Krimi, R. (2023, 20 décembre). *Comprendre JSON : syntaxe, stockage et exemples. Hostinger Tutoriels*. Consulté le 3 avril 2024 à l'adresse : <https://www.hostinger.fr/tutoriels/quest-ce-que-json>
- [4] Robinson, J. S. A. D. (s. d.). *Welcome to Text Mining with R | Text Mining with R*. Emploi fréquent à l'adresse : <https://www.tidyttextmining.com/>
- [5] Rod. (2024, 24 avril). *Le Web Scraping en pratique*. MonCoachData. Consulté le 20 mars 2024 à l'adresse : <https://moncoachdata.com/blog/Web-scraping-pratique/>
- [6] Silletti, A. (s. d.). *Invasive earthworms – Drake Lab*. Consulté le 29 mai 2024 à l'adresse : <https://daphnia.ecology.uga.edu/drakelab/?p=318>
- [7] Zdmit (Développeur Python anonyme) (s. d.). *[Script] Scraping ResearchGate, toutes les publications*. Consulté le 23 mars 2024 à l'adresse : https://www.reddit.com/r/Python/comments/v2bxyl/script_scraping_researchgate_all_publications/

Résumé

La perception des vers de terre dans la littérature scientifique mondiale est divisée entre deux positions contradictoires : une première thèse, défendue surtout par les spécialistes Européens, est de considérer les vers de terre comme des agents écosystémiques essentiels pour la productivité et la santé des sols, là où d'autres spécialistes, plutôt du côté Américain, défendent la thèse opposée, considérant les vers de terre Asiatiques et Européens comme des espèces invasives mettant gravement en danger la stabilité des écosystèmes natifs, notamment dans les forêts d'Amérique du Nord. L'objectif de cette étude est d'essayer d'apporter une réponse à cette controverse via l'analyse d'un jeu de données textuelles de 168 abstracts d'articles scientifiques du domaine, grâce à des méthodes informatiques et statistiques de *Web-scraping* et de *Text-mining*, développées respectivement dans les langages informatiques Python et R. Les résultats obtenus ont montré qu'il était effectivement possible de détecter informatiquement l'existence de cette controverse en observant les principaux sujets et opinions développés dans le corpus de textes fourni, même s'il n'a pas été possible d'affirmer catégoriquement l'existence de deux groupes thématiques clairement distincts. Cette étude suggère cependant que les vers de terre auraient un effet bénéfique sur les sols et les plantes, mais un effet néfaste sur les écosystèmes natifs où ils sont importés.

Mots-clés de référencement type MESH : *Lumbricus terrestris*, Soil, Ecosystem, Introduced species, Meta-Analysis.

Mots-clés des acquis techniques : Web-scraping, Text-mining, dplyr, R Markdown, GitHub.