
TD1 - G01

Guillaume Bonenfant

Pierre de Saxce

Nathanaël Lefèvre ***** Scrum Master

Sébastien Rocca

Raphaël Rey

Sprint 5 - Génie Logiciel

01-04 Décembre 2020

Message du scrum Master

“E-VA-LU-A-TION !!! Pas de panique : on s'évalue nous même, alors ça va ! Bravo à tous, je trouve qu'on a vraiment réussi à avoir de bon résultat, petit à petit. Il nous a peut-être fallu un peu de temps pour se "mettre dans le bain" (surtout moi) mais au final on est parvenu à quelque chose.. Ensemble ! Et c'est pour faire les choses ensemble que je suis venu en CMI, c'était dans ma lettre de motivation.”

VUE D'ENSEMBLE

Noter le parsing

faire le rapport final

Répartition des tâches:

- Nathanaël :
 - A Benders Decomposition Approach to Correlation Clustering
 - A memetic algorithm for community detection in signed networks
- Raphaël :
 - An Improved Branch-and-Cut Code for the Maximum Balanced Subgraph of a Signed Graph
 - Cabrera_RESUMES_2019
- Pierre :
 - Conversational Networks for Automatic Online Moderation
 - Dynamical Models Explaining Social Balance and Evolution of Cooperation
- Sébastien :
 - Exact Clustering via Integer Programming and Maximum Satisfiability

-
- LDA_resume
 - Guillaume :
 - Partitioning
 - Polibits_42_02 :

Resultats

PST : Précision stricte / PSS : Précision souple

- A Benders Decomposition Approach to Correlation Clustering : PST=11/15, PSS=13/15
- A memetic algorithm for community detection in signed networks : PST=11/16, PSS=11/16
- An Improved Branch-and-Cut Code for the Maximum Balanced Subgraph of a Signed Graph
- Cabrera_RESUMES_2019
- Conversational Networks for Automatic Online Moderation
- Dynamical Models Explaining Social Balance and Evolution of Cooperation
- Exact Clustering via Integer Programming and Maximum Satisfiability: PSS: 14/16, PST: 14/16
- LDA_resume: PSS= 11/18, PST= 14/18
- Partitioning
- Polibits_42_02 :
-

PRÉPARATION DU RAPPORT

TITRE

- Analyse des performances d'un parseurs de pdf réalisé par des étudiants de niveau L3 Informatique.

AUTEURS

- Guillaume Bonenfant (guillaume.bonenfant@alumni.univ-avignon.fr)
- Nathanaël Lefèvre (nathanael.lefevre@alumni.univ-avignon.fr)
- Raphaël Rey (raphael.rey@alumni.univ-avignon.fr)
- Sébastien Rocca (sebastien.rocca@alumni.univ-avignon.fr)
- Pierre de Saxce (pierre.de-saxce-nahon@alumni.univ-avignon.fr)

ABSTRACT

Dans cette étude, nous nous intéresserons aux performances d'un parseur de pdf que nous avons réalisé dans le cadre de l'Unité d'Enseignement "génie logiciel" dispensée par Avignon Université. Pour ce qui est du cadre de notre expérience, notre but était de récupérer des fichiers de sortie en .txt ou .xml contenant les informations principales d'articles scientifiques au format pdf. Pour ce faire, nous avons comparé deux systèmes d'OCR (pdf2txt et pdftotext) pour sélectionner le plus fiable. Par la suite, nous nous sommes basé sur les fichiers .txt générés par le système d'OCR pour les parser et récupérer le titre, l'abstract, les auteurs, l'introduction, le corps de l'article, la discussion, la conclusion et les références bibliographiques. Nous nous intéresserons dans un premier temps à l'implémentation de notre parseur puis nous analyserons les performances de notre système.

METHODE

Pour parser des articles scientifiques en pdf vers une sortie en .txt ou .xml, nous nous appuyons sur la transcription des fichiers d'entrée en pure texte par le programme d'OCR pdftotext avec les options -npgbrk et -raw, la première évitant les problèmes liés aux sauts de page et la seconde permettant de reconnaître les textes sur plusieurs colonnes.

Notre structure de récupération des données est basée sur un dictionnaire data associant les noms des différentes parties à récupérer à leur contenu détecté dans le fichier d'entrée. Pour ce qui est de la récupération des données à proprement parler, notre système se base sur une méthode itérative, ligne par ligne, et fonctionne grâce à des expressions Regex. Lorsqu'une partie est détectée, la variable "copy" prend la valeur de la clé correspondante dans le dictionnaire "data". Tant que "copy" est égal au nom de la partie en cours et n'est pas vide, les lignes lues sont ajoutées à la valeur liée à la clé correspondant à "copy" dans data. Par exemple, si on détecte l'abstract, alors copy="abstract" et les lignes lues sont enregistrées dans data["abstract"].

Par ailleurs, nous nous assurons de ne récupérer chaque partie qu'une seule fois. Ainsi, si plusieurs lignes matchent la règle de détection d'une partie, seule la première correspondance est prise en compte. Pour ce faire, on ajoute la clé de la partie en cours dans la liste "copied". On ne détecte le début d'une partie que si celle-ci n'est pas inscrite dans "copied".

Le nom original du document est directement récupéré lors du choix du pdf à parser de façon triviale. Pour le titre du document, le système tente de le récupérer depuis les données au format doc du pdf parsé par le module pdfminer. Si cela ne fonctionne pas, on récupère la première ligne du document et nous considérons qu'il s'agit du titre de l'article.

Les autres récupérations de données sont purement basées sur les regex et celles-ci ont un certain degré d'indulgence.

Voici, pour chaque partie, la règle appliquée pour détecter le début de la partie en question:

Auteurs : du début du document jusqu'au début de l'abstract.

Abstract : le mot "abstract" en début de ligne précédé de 0 à 5 caractères et suivi de 0 à 100 caractères puis d'un retour à la ligne. La casse est ignorée.

Introduction : le mot "introduction" en début de ligne précédé de 0 à 5 caractères et suivi de 0 à 3 caractères puis d'un retour à la ligne. La casse est ignorée.

Corps : de la fin de l'introduction jusqu'à ce qu'une autre partie soit détectée.

Discussion : le mot "Discussion" avec une majuscule ou tout en majuscule, précédé par 0 à 5 caractères et suivi par 0 à 40 caractères et un retour à la ligne.

Conclusion : Le mot "Conclusion" avec ou sans "s" ou "Summary" avec une majuscule ou tout en majuscule, précédé de 0 à 10 caractères et suivi de 0 à 40 caractères pour "Conclusion" et de 0 à 10 caractères pour "Summary". Le tout doit être suivi d'un retour à la ligne.

Références : Le mot "Reference" avec ou sans "s" avec une majuscule ou tout en majuscule, précédé par 0 à 5 caractères et suivi par 0 à 3 caractères et un retour à la ligne.

En règle générale, une partie se termine lors de la détection du début d'une partie suivante. Dans certains cas, une partie peut être terminée prématurément lors de la détection du numéro d'une nouvelle partie, notamment la fin de l'introduction lorsque l'on détecte une partie numérotée par un "2".

RESULTATS

!!! Ces valeurs ont été corrigées dans le rapport final en enlevant les nom des pdf de l'évaluation. (on enlève donc 2 points au numérateur et au dénominateur)

	Précision souple	Précision stricte
A Benders Decomposition Approach to Correlation Clustering	13/15	11/15
A memetic algorithm for community detection in signed networks	11/16	11/16
An Improved Branch-and-Cut Code for the Maximum Balanced Subgraph of a Signed Graph	13/14	11/14
Cabrera_RESUMES_2019	14/18	13/18
Conversational Networks for Automatic Online Moderation	14/16	13/16
Dynamical Models Explaining Social Balance and Evolution of Cooperation	9/16	7/16
Exact Clustering via Integer Programming and Maximum Satisfiability	14/16	14/16
LDA_resume	14/18	11/18
Partitioning large signed two-mode networks/ Problems and prospects	13/18	13/18
Polibits 42 02	16/18	13/18
MOYENNES	418/525 = 79,6%	71,1%

CONCLUSION

En précision souple (+/-2 phrases de précision sur les frontières des parties), notre système réalise une performance de 79,6% de réussite sur un corpus de 10 pdf mis à notre

disposition pour évaluation. En précision stricte (à la phrase près), les performances baissent à 71,1% ce qui est très correct selon nous.

Finalement nous pouvons considérer que le parsing d'articles scientifiques n'est pas une tâche très difficile dans le monde de l'informatique. Si des étudiants de L3 sont capables d'obtenir un pourcentage de réussite supérieur à 70% en peu de temps, nous pouvons imaginer que des professionnels sont capables de faire beaucoup mieux. De plus l'appel à une IA permettrait d'obtenir des résultats encore meilleurs. Notons que sans IA, le parsing des articles ne doit pas pouvoir se faire très facilement avec un pourcentage de réussite vraiment élevé. Sans IA, si les articles scientifique suivaient une norme plus stricte, le parsing serait trivial.

DEBRIEFING

- Le projet s'est déroulé sans encombre.
- Notons que certains d'entre nous doivent encore s'améliorer dans l'usage de github.
- On peut être fier de nos résultats, félicitation à tous !!

- Le scrum master -