

M1 Informatique –UE Projet

Carnet de bord : les coulisses de la recherche documentaire

Les éléments que vous indiquez dans ce carnet donneront lieu à une notation

Noms, prénoms et spécialité :

Amine Djeghri
Idles Mamou
Master 1 DAC

Sujet :

Sparse Embeddings pour la Recherche d'Information

Consigne :

1. **Introduction (5- 10 lignes max)** : Décrivez rapidement votre sujet de recherche, ses différents aspects et enjeux, ainsi que l'angle sous lequel vous avez décidé de le traiter.
2. **Les mots clés retenus (5- 10 lignes max)** : Listez les mots clés que vous avez utilisés pour votre recherche bibliographique. Organisez-les sous forme de carte heuristique.
3. **Descriptif de la recherche documentaire (10-15 lignes)** : Décrivez votre utilisation des différents outils de recherche (moteurs de recherche, base de donnée, catalogues, recherche par rebond etc.) et comparez les outils entre eux ? A quelles sources vous ont-ils permis d'accéder ? Quelles sont leurs spécificités ? Leur niveau de spécialisation ?
4. **Bibliographie produite dans le cadre du projet** : Utilisez la norme ACM ou IEEE.
5. **Evaluation des sources (5 lignes minimum par sources)** : Choisissez 3 sources parmi votre bibliographie, décrivez la manière dont vous les avez trouvées et faites-en une évaluation critique en utilisant les critères vus en TD.

Votre carnet de bord doit être remis en mains propres au formateur LE JOUR DU TUTORAT. Une copie numérique devra être envoyée à l'adresse suivante : Adrien.Demilly@scd.upmc.fr

Rappel : les supports de TD sont disponibles à l'adresse suivante:
<http://www.pearltrees.com/formationbsu/master-info/id23514400>

Introduction :

La recherche d'information est le domaine qui étudie la manière de retrouver des informations dans un corpus. Notre projet consiste à implémenter un modèle de ranking de recherche d'information, qui consiste à retrouver les documents tout en les triant pour l'utilisateur dans un ordre du plus pertinent au moins pertinent sur ce qu'il recherche.

Dans notre projet, nous aborderons les modèles de recherche d'information, à base de réseaux ainsi que des modèles de représentation de textes. Une des manières de représenter du texte dans une machine dans le domaine du Traitement Automatique de la Langue, c'est, d'utiliser une représentation continue, appelée « embedding ».

L'un des objectifs du projet consiste à implémenter les découvertes récentes en recherche d'information.

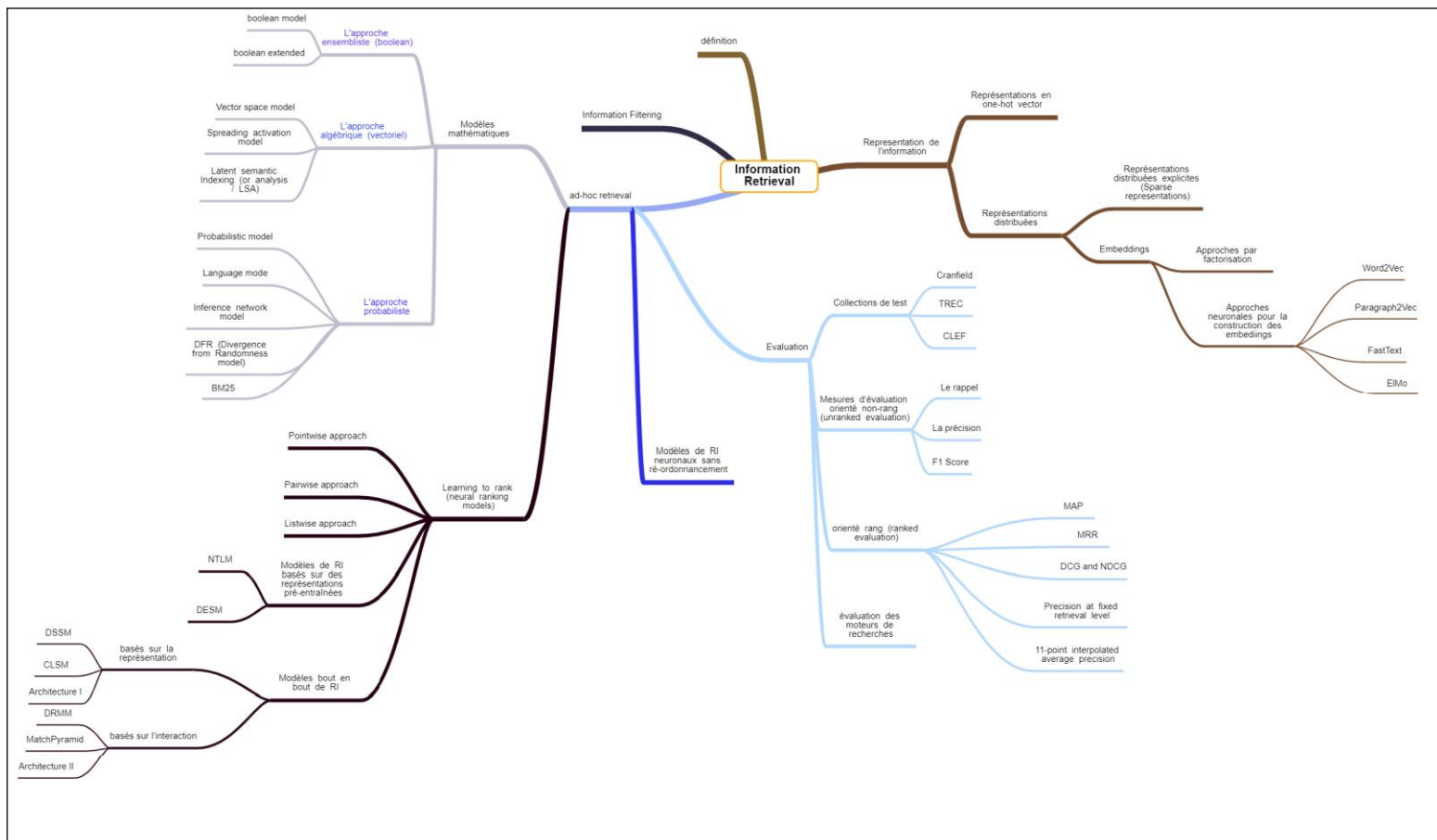
Les mots clés :

Information Retrieval, ad-hoc retrieval, Information Filtering, boolean model, Vector space model, Latent semantic Indexing, Probabilistic model, Language model, Inference network model, BM25, Learning to rank, neural ranking models, Pointwise approach, Pairwise approach, Listwise approach, NTLM, DESM, DSSM, CLSM, Architecture I, DRMM, MatchPyramid, Architecture II, one-hot vector, Sparse representations, Embeddings, Word2Vec, Paragraph2Vec, FastText, EIMo, Evaluation metrics information retrieval, Collections de test, unranked evaluation, ranked evaluation

(Généralement, nous combinons les mots clés cités au-dessus, avec le mot clé 'information retrieval' par exemple : 'Evaluation information retrieval' afin d'obtenir les bons résultats qui tournent autour du sujet notre projet).

Carte heuristique :

Pour dessiner la carte heuristique, nous avons utilisé Framindmap.org



En partant des articles cités dans les références du sujet de notre projet et des sources primaires :

Références

- [1] ZAMANI, Hamed, DEGHANI, Mostafa, CROFT, W. Bruce, et al. From neural re-ranking to neural ranking : Learning a sparse representation for inverted indexing. In : Proceedings of the 27th ACM International Conference on Information and Knowledge Management. ACM, 2018 [link](#)
- [2] DEGHANI, Mostafa, ZAMANI, Hamed, SEVERYN, Aliaksei, et al. Neural ranking models with weak supervision. In : Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval. ACM, 2017. [link](#)
- [3] MIKOLOV, Tomas, SUTSKEVER, Ilya, CHEN, Kai, et al. Distributed representations of words and phrases and their compositionality. In : Advances in neural information processing systems. 2013. [link](#)
- [4] Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L. (2018). Deep contextualized word representations. [link](#)
- [5] Huang, P. S., He, X., Gao, J., Deng, L., Acero, A., Heck, L. (2013, October). Learning deep structured semantic models for web search using clickthrough data. In Proceedings of the 22nd ACM international conference on Information Knowledge Management ACM. [link](#)

Nous avons déterminé quelques mots clés qui nous aideront à débiter notre recherche documentaire.

D'abord, après avoir recherché les mots clés sur Google, ce dernier nous affiche généralement en premier lieu wikipédia, que nous utilisons uniquement pour cadrer le sujet et avoir le bon vocabulaire. Après avoir construit le début de la carte heuristique, cette dernière va nous permettre de se concentrer sur les branches une par une et de s'approfondir grâce à Arxiv, Google Scholar et Web of science pour retrouver les articles en lien avec nos recherches. Par exemple lors de notre utilisation de **Google Scholar**, nous avons remarqué qu'il contient des millions de résultats et permet d'avoir une grande variété de documents, mais propose peu de filtres, par conséquent, nous l'avons utilisé que pour une recherche précise d'un article.

D'autre part, le **Portail documentaire Sorbonne / web of science** propose beaucoup de filtres pour la recherche de livres, d'articles et de reviews et nous fournit l'accès complet grâce à l'abonnement de l'université. Quant à **Arxiv/ACM**, nous l'avons utilisé principalement pour la recherche d'articles vu qu'il est spécialisé dans le domaine de l'informatique, cependant, il contient beaucoup de pré-publications, et d'articles qui n'ont pas été encore publiés.

Bibliographie produite dans le cadre du projet

Et enfin, pour ce qui est des sources, grâce à Zotero et son extension sur le navigateur, ce logiciel nous a permis d'insérer et de citer facilement les références et les bibliographies des différents articles que nous avons utilisés pour la rédaction de notre rapport.

Ci-dessous, quelques sources parmi les 60 sources que nous avons citées, ainsi que le fichier en extension « .bib » généré par Zoteron norme : IEEE

- [3]J. Guo, Y. Fan, Q. Ai, et W. B. Croft, « A Deep Relevance Matching Model for Ad-hoc Retrieval » ,in *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management - CIKM '16*, p. 55–64, 2016, doi: 10.1145/2983323.2983769.
- [4]Y. Yue, T. Finley, F. Radlinski, et T. Joachims, « A support vector method for optimizing average precision », in *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval - SIGIR '07*, Amsterdam, The Netherlands, 2007, p. 271, doi: 10.1145/1277741.1277790.
- [5]Q. Wu, C. J. C. Burges, K. M. Svore, et J. Gao, « Adapting boosting for information retrieval measures », *Information Retrieval*, vol. 13, n° 3, p. 254–270, juin 2010, doi: 10.1007/s10791-009-9112-1.
- [12]Q. Le et T. Mikolov, « Distributed Representations of Sentences and Documents », in *Proceedings of the 31st International Conference on Machine Learning*, PMLR 32(2):1188-1196, p. 9.
- [15]T. Mikolov, K. Chen, G. Corrado, et J. Dean, « Efficient Estimation of Word Representations in Vector Space », In : *Advances in neural information processing systems*, sept. 2013, arXiv:1301.3781

- [16]Y. Nakamoto, « FOREWORD », *IEICE Transactions on Information and Systems*, vol. E94-D, n° 1, p. 1–2, 2011, doi: 10.1587/transinf.E94.D.1.
- [17]P. D. Turney et P. Pantel, « From Frequency to Meaning: Vector Space Models of Semantics », in *Journal of Artificial Intelligence Research*, vol. 37, p. 141–188, févr. 2010, doi: 10.1613/jair.2934.
- [18]J. Pennington, R. Socher, et C. Manning, « Glove: Global Vectors for Word Representation », in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Doha, Qatar, 2014, p. 1532–1543, doi: 10.3115/v1/D14-1162.
- [19]S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, et R. Harshman, « Indexing by latent semantic analysis », in *Journal of the American Society for Information Science*, vol. 41, n° 6, p. 391–407, sept. 1990, doi: 10.1002/(SICI)1097-4571(199009)41:6<391::AID-ASI1>3.0.CO;2-9.
- [20]A. Berger, « Information Retrieval as Statistical Translation », *ACM SIGIR Forum*, vol. 51, n° 2, p. 8, 2017, doi: 10.1145/3130348.3130371
- [21]G. Zuccon, B. Koopman, P. Bruza, et L. Azzopardi, « Integrating and Evaluating Neural Word Embeddings in Information Retrieval », in *Proceedings of the 20th Australasian Document Computing Symposium on ZZZ - ADCS '15*, Parramatta, NSW, Australia, 2015, p. 1–8, doi: 10.1145/2838931.2838936.
- [22] T. Francesiaz, R. Graille, et B. Metahri, « Introduction aux modèles probabilistes utilisés en Fouille de Données », Laboratoire Jean Kuntzmann, Grenoble, 2015, p. 27

[Evaluation des sources](#)

[Source 1 :](#)

[15]T. Mikolov, K. Chen, G. Corrado, et J. Dean, « Efficient Estimation of Word Representations in Vector Space », In : *Advances in neural information processing systems*, sept. 2013, arXiv:1301.3781

Cet article est l'une des sources principales fournies par notre encadrant. Il a été publié en septembre 2016 dans « Neural Information Processing Systems » par T. Mikolov, K. Chen, G. Corrado, et J. Dean et l'article a été cité plus de 15000 fois (source : Semantic Scholar). Les quatre auteurs sont des chercheurs chez google, précisément, le premier est le créateur de la méthode word2vec, le second est un professeur à Hong Kong Université, le dernier quand à lui est le directeur de google AI Division. Les auteurs débutent leur article avec une revue de l'état de l'art en citant chaque source, ils proposent par la suite leur modèle et extensions, les résultats de leurs tests et à la fin une comparaison avec des résultats publiés par d'autres auteurs du modèle. Le but de ce papier est de présenter plusieurs extensions qui améliorent à la fois la qualité des vecteurs et la vitesse de l'apprentissage.

[Source 2 :](#)

[12]Q. Le et T. Mikolov, « Distributed Representations of Sentences and Documents », in *Proceedings of the 31st International Conference on Machine Learning*, PMLR 32(2):1188-1196, 2014

Cet article écrit par Q. Le et T. Mikolov qui sont des chercheurs membre de l'équipe de recherche Google Brain. Cet article a été publié dans « Proceedings of the 31st International Conference on Machine Learning » en 2014.

Nous avons trouvé l'article par une simple recherche sur arxiv, et nous avons également constaté qu'il a été cité 4000 fois (source : Semantic Scholar)

Les auteurs débutent leur article avec une revue de l'état de l'art en citant chaque source, puis par une présentation des algorithmes dont ils se sont inspirés, ensuite leur framework avec les expérimentations, et enfin, ils donnent les résultats de leurs tests et comparent avec des résultats publiés par d'autres auteurs.

Le but de ce papier est de proposer « Paragraph Vector », un framework non supervisé qui apprend des vecteurs de représentations continues et distribuées de morceaux de texte.

[Source 3 :](#)

[21]G. Zuccon, B. Koopman, P. Bruza, et L. Azzopardi, « Integrating and Evaluating Neural Word Embeddings in Information Retrieval », in *Proceedings of the 20th Australasian Document Computing Symposium on ZZZ - ADCS '15*, Parramatta, NSW, Australia, 2015,

Cet article est l'une des sources trouvées comme référence l'un des autres articles. Il a été publié en septembre 2015 dans « Proceedings of the 20th Australasian Document Computing Symposium on ZZZ - ADCS '15, Parramatta, NSW, Australia » par G. Zuccon, B. Koopman, P. Bruza, et L. Azzopard. Les trois premiers sont des chercheurs à Queensland University of Technology, le

deuxième particulièrement est également chercheur à Australian e-Health Research, enfin, le dernier est chercheur à University of Glasgow.

Les auteurs débutent leur article avec une revue de l'état de l'art tout en essayant de montrer les améliorations qu'ils peuvent apporter, puis ils présentent le « language modelling framework » et quelques autres modèles, Enfin ils décrivent les expérimentations et les tests qu'ils ont faits, et donnent les résultats obtenus.

Le but de ce papier est de déterminer comment les 'words embeddings' peuvent être utilisées dans un modèle de recherche d'information et quels avantages pourraient-ils apporter.