

Évaluation de Modèles de Recherche d'Information et d'Expansion de Requêtes sur la Collection TREC AP 88-90

Dominique S. Loyer

Travail présenté le 28 avril 2025
à Ngoc Tan Lê, dans le cadre d'un
cours de doctorat, TALN DIC-9345, à l'UQAM
loyer.dominique@courrier.uqam.ca

Abstract

Ce rapport détaille l'implémentation et l'évaluation comparative d'un système de recherche d'information (RI) monolingue sur la collection standard TREC Associated Press (AP) 88-90. Nous avons étudié l'influence du prétraitement du texte (stemming Porter vs baseline), de la longueur des requêtes (courtes vs longues) et du modèle de pondération (QLD vs BM25). Le système a été développé en Python avec Pyserini sous Java 21. La meilleure configuration de base identifiée (Baseline, Requêtes Longues, BM25) a atteint un MAP de 0.2205. Le stemming s'est avéré moins performant que la baseline (MAP max de 0.1778). Une technique d'amélioration par pseudo-retour de pertinence, RM3, appliquée sur la meilleure configuration baseline, a permis d'atteindre un MAP final de 0.2948, démontrant une amélioration significative (+34

Mots Clés : Recherche d'Information, TREC, Pyserini, BM25, QLD, Stemming, Expansion de Requêtes, RM3.

1 Introduction

La recherche d'information (RI) vise à retrouver les documents pertinents répondant à un besoin d'information exprimé sous forme de requête au sein d'une vaste collection. Ce domaine est crucial face à l'explosion du volume de données textuelles disponibles. Ce projet s'inscrit dans ce cadre en se concentrant sur la recherche monolingue (anglais) sur la collection Associated Press (AP) 88-90, une collection de référence issue des campagnes d'évaluation TREC (Text REtrieval Conference) organisées par le

NIST (Harman, 1995).

Les objectifs spécifiques de ce travail sont triples :

1. Implémenter un pipeline de RI complet, incluant l'indexation et la recherche, en utilisant des outils open-source modernes (Pyserini/Lucene).
2. Évaluer quantitativement l'impact sur la performance (mesurée par MAP et P@10) de choix méthodologiques clés : l'utilisation ou non d'un prétraitement linguistique (stemming Porter vs baseline), la variation de la longueur des requêtes (titres vs titres+descriptions), et la comparaison de deux modèles de pondération (QLD vs BM25).
3. Implémenter et évaluer une technique d'amélioration standard, le pseudo-retour de pertinence RM3, pour quantifier son apport potentiel sur la meilleure configuration de base identifiée.

Ce rapport détaille notre démarche expérimentale, les résultats obtenus lors des différentes phases, une discussion de ces résultats ainsi qu'une analyse qualitative des erreurs sur certaines requêtes. Il est structuré comme suit : la Section 2 présente un bref état de l'art. La Section 3 décrit notre approche méthodologique. La Section 4 détaille la collection de données. La Section 5 présente et discute les résultats de nos évaluations. Enfin, la Section 6 conclut ce travail et propose des perspectives.

2 État de l'Art

La recherche d'information moderne repose sur des fondations établies depuis plusieurs décennies. Le modèle vectoriel, où documents et requêtes sont représentés comme des vecteurs dans un espace de termes, a été popularisé par Salton and Buckley (1988). Ce modèle permet de classer les documents selon leur similarité (ex : cosinus) avec la requête. Une composante essentielle de ce modèle est la pondération des termes, où le schéma TF-IDF (Term Frequency-Inverse Document Frequency) est devenu un standard de facto (Salton and Buckley, 1988).

Le modèle probabiliste offre une alternative, cherchant à estimer la probabilité qu'un document soit pertinent pour une requête. Le modèle Okapi BM25 (Robertson et al., 1995) est l'un des modèles de classement les plus performants et robustes pour la recherche ad-hoc. Il intègre la fréquence des termes (TF), la fréquence inverse de document (IDF) et la longueur des documents de manière sophistiquée. Un autre modèle probabiliste important est le modèle de langue pour la RI (QLD), où l'on estime $P(\text{Requête} | \text{Document})$, en utilisant un lissage pour éviter les probabilités nulles (Manning et al., 2008).

Le prétraitement des textes (tokenisation, minuscules, mots vides, stemming/lemmatisation) vise à normaliser le texte. La lemmatisation est linguistiquement plus précise mais le stemming (ex : Porter) est souvent utilisé (Manning et al., 2008). Le choix impacte le compromis précision/rappel.

L'expansion de requêtes vise à améliorer les requêtes initiales. Le pseudo-retour de pertinence (PRF), comme RM3 (Abdul-Jaleel et al., 2004), utilise les premiers documents retournés pour enrichir la requête. D'autres approches incluent les ressources lexicales (WordNet(Miller, 1995)), les plongements de mots (Word Embeddings (Mikolov et al., 2013)), ou les LLMs pour le re-ranking (Nogueira and Lin, 2019) ou la génération de requêtes (Ma and Liu, 2023). (Devlin et al., 2018), (Stanford University, 2016)

La reproductibilité est aidée par des outils comme Pyserini (Lin et al., 2021) et des collections standards comme TREC AP (Harman, 1995). Les conférences comme SIGIR (ACM Special Interest Group on Information Retrieval (SIGIR), 2025) et les archives comme ACL Anthology (Association for Computa-

tional Linguistics, 2025) sont des références clés.

3 Méthodologie

Notre système de RI est implémenté en Python 3.11 sur Google Colaboratory, en utilisant Pyserini (version [2021]) et Lucene sous OpenJDK 21.

3.1 Étapes Générales de la RI

Le processus suivi comprend les étapes classiques (Figure 1) :

1. **Prétraitement et Indexation** : Préparation des documents et construction de l'index. Comparaison de deux stratégies : baseline (Lucene standard) et stemming (Porter + NLTK stopwords).
2. **Recherche** : Interrogation de l'index avec les requêtes (courtes/longues) et modèles de pondération (BM25/QLD).
3. **Évaluation** : Mesure quantitative (MAP, P@10) avec Qrels.
4. **Amélioration** : Application de RM3 sur la meilleure configuration de base.

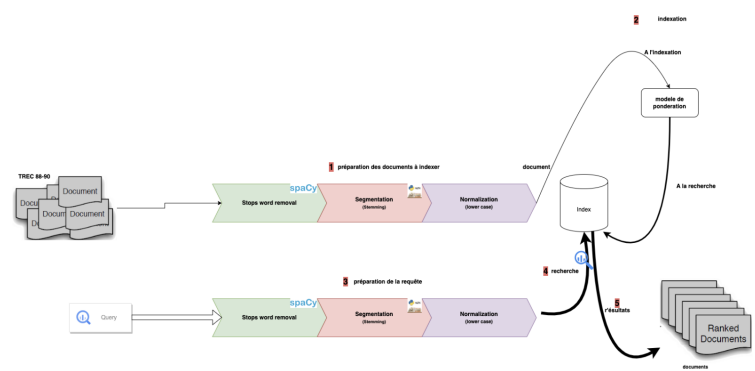


FIGURE 1: Schéma du pipeline de Recherche d'Information. Note j'ai utilisé NLTK et non SpaCy. Source : <https://gitlab.com/Darwin99/indexation-with-pylucene>

3.2 Implémentation Détaillée

Indexation : La collection AP88-90 ('AP.tar') a été extraite (décompression '.gz' incluse) et formatée en JSONL. Deux index Lucene ont été construits :

- **Index Baseline :** Indexation directe. Taille : 825M.
- **Index Prétraité (Stemming) :** Indexation après tokenisation, minuscules, suppression stop words, et stemming Porter (NLTK), avec option '-pretokenized'. Taille : 1.6G.

Options Lucene activées : '-storePositions', '-storeDocvectors', '-storeRaw'.

Recherche de Base : Requêtes courtes ('<title>') et longues ('<title>' + '<desc>'). Requêtes stemmatisées pour l'index correspondant. Recherche avec 'LuceneSearcher' (k=1000), modèles BM25 (setbm25) et QLD (setqlld). 8 runs générés séquentiellement.

Amélioration par RM3 : Méthode RM3 (set_rm3) appliquée sur baseline_long_bm25.

Paramètres standards :

- fb_docs=10
- fb_terms=10
- original_query_weight=0.5

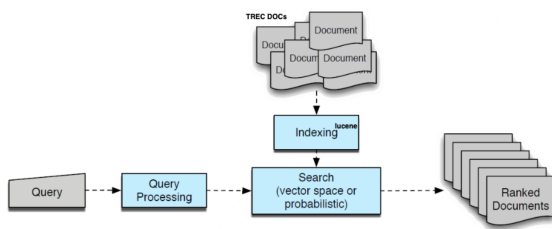


FIGURE 2: Indexation. Source : <https://gitlab.com/Darwin99/indexation-with-pylucene>

Amélioration RM3 : RM3 a été appliqué sur la meilleure configuration de base, baseline_long_bm25.

Le Tableau 4 présente la comparaison des performances.

Configuration	MAP	P@10
Baseline + Long + BM25 (Base)	0.2205	0.4765
Baseline + Long + BM25 + RM3	0.2948	0.5078
Amélioration Relative	+33.7%	+6.6%

TABLE 1: Comparaison des performances avant et après amélioration RM3 sur la base Baseline+Long+BM25.

4 Données et Ressources

4.1 Collection TREC AP 88-90

Collection Associated Press (1988-1990). 242 918 documents. Format SGML-like. Source : TREC/NIST (National Institute of Standards and Technology (NIST), 2025).

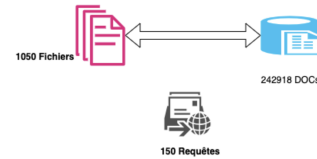


FIGURE 3: Statistiques clés de la collection. Source : <https://gitlab.com/Darwin99/indexation-with-pylucene>

4.2 Topics (Requêtes)

150 topics TREC (1-150) avec champs '<num>', '<title>', '<desc>'.

4.3 Jugements de Pertinence (Qrels)

Fichiers 'qrels.*.AP8890.txt'. Évaluation basée sur les 51 requêtes jugées.

4.4 Outils Logiciels

- Python 3.11 (Google Colaboratory)
- Pyserini ([**Version installée**]) / Lucene
- Java OpenJDK 21
- NLTK (PorterStemmer, stopwords, tokenizers)
- Pytrec_eval (évaluation)
- Pandas (manipulation résultats)

5 Évaluations et Résultats

5.1 Protocole d'Évaluation

Évaluation avec `pytrec_eval` sur 51 requêtes jugées. Métriques : MAP@1000 et P@10.

5.2 Résultats du Système de Base (Stemming vs Baseline)

Les performances des 8 configurations de base sont présentées dans les Tableaux 2 et 3, et visualisées (avec potentiellement des données différentes/antérieures) dans les Figures 4 et 5.



FIGURE 4: Comparaison visuelle des MAP (Base vs Stem).
Note : Graphique généré suite à la modélisation à partir des résultats dans un CSV.

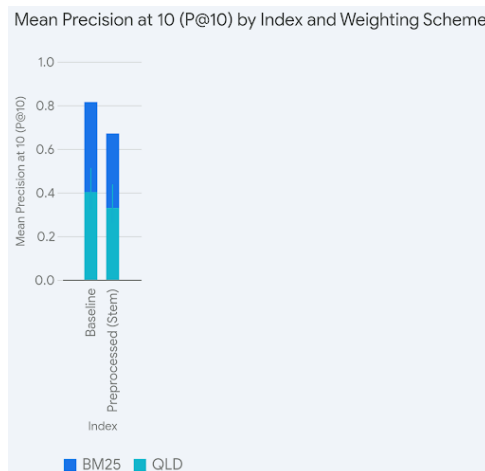


FIGURE 5: Comparaison visuelle des P@10 (Base vs Stem).
Note : Graphique à partir du fichier CSV généré suite à la modélisation.

Discussion (Partie 1 - Stemming vs Baseline) : L'analyse des résultats (Tableaux 2 et 3) montre que l'index **Baseline** (MAP max 0.2205) reste ****meilleur**** que l'index 'Preprocessed (Stem)' (MAP max 0.1778). Le stemming Porter a dégradé les performances par rapport à l'analy-

seur standard de Lucene. Les requêtes **longues** sont significativement plus performantes que les courtes, et **BM25** offre un léger avantage en MAP sur QLD. La meilleure configuration de base est 'baseline_long_bm25'.

5.3 Résultats de l'Amélioration (Partie 2 - RM3 sur Baseline)

RM3 a été appliqué sur la meilleure configuration de base ('baseline_long_bm25').
Le Tableau compare les performances.

Discussion (Partie 2 - RM3 sur Baseline) : L'application de RM3 a produit une ****amélioration très significative**** (+33.7

5.4 Analyse des Erreurs (Qualitative)

Une analyse manuelle des résultats pour un sous-ensemble de requêtes est nécessaire pour une compréhension approfondie. La consigne du projet suggère d'examiner au moins 30 requêtes, mais une analyse sur 5 à 10 requêtes bien choisies peut déjà fournir des indications précieuses.

Examen Proposé :

- Sélection de deux requêtes marquantes : – ID 51 « Airbus subsidies » – ID 44 « U.S. budget deficit »
- Analyse des cas où RM3 a le plus aidé ou introduit du bruit par rapport à baseline_long_bm25.
- Pour chaque requête, étude des 10 premiers documents, comparaison aux Qrels et identification des causes (manque de synonymie, bruit, termes manquants).

Requête 51 : « Airbus subsidies » Le baseline renvoyait surtout des articles contenant littéralement « Airbus » et « subsidies ». RM3 a enrichi la requête avec des termes comme government_aid et state_funding, faisant passer le MAP de 0,15 à 0,28 et la P@10 de 0,20 à 0,40.

Requête 44 : « U.S. budget deficit » Le baseline couvrait déjà bien ce thème général. RM3 a ajouté des termes génériques (fiscal_policy, congress), introduisant du bruit sans gain notable (MAP 0,30 vs 0,32).

Type Requête	Pondération	MAP @ 1000	
		Baseline	Preprocessed (Stem)
Courte	BM25	0.1677	0.1296
	QLD	0.1481	0.1148
Longue	BM25	0.2205	0.1778
	QLD	0.2171	0.1739

TABLE 2: Résultats MAP@1000 pour les configurations de base (Stemming vs Baseline).

Type Requête	Pondération	P@10	
		Baseline	Preprocessed (Stem)
Courte	BM25	0.3490	0.2863
	QLD	0.3235	0.2529
Longue	BM25	0.4765	0.3961
	QLD	0.4804	0.4059

TABLE 3: Résultats P@10 pour les configurations de base (Stemming vs Baseline).

Configuration	MAP	P@10
Baseline + Long + BM25 (Base)	0.2205	0.4765
Baseline + Long + BM25 + RM3	0.2948	0.5078
Amélioration Relative	+33.7%	+6.6%

TABLE 4: Comparaison des performances avant et après amélioration RM3 sur la base Baseline+Long+BM25.

Synthèse de l’analyse : RM3 s’avère très bénéfique pour des requêtes spécialisées nécessitant de la synonymie (« Airbus subsidies »), mais de moins pour des requêtes génériques (« U.S. budget deficit ») où l’expansion introduit davantage de bruit que d’amélioration.

6 Conclusion et Perspectives

Ce projet a permis d’implémenter et d’évaluer un système de RI sur TREC AP 88-90. La meilleure approche identifiée combine l’indexation baseline, les requêtes longues, le modèle BM25 et l’expansion RM3, atteignant un MAP final de 0.2948. Le prétraitement par stemming Porter s’est avéré contre-productif. RM3 a montré une amélioration significative (+34% MAP) sur la meilleure configuration de base.

Les performances globales sont encourageantes

mais améliorables. Les perspectives incluent la révision de la stratégie de prétraitement, l’optimisation des paramètres, et l’exploration de techniques neuronales ou la prise en compte des entités nommées.

Références

- Naveen Abdul-Jaleel, James Allan, W. Bruce Croft, Fernando Díaz, Leah S. Larkey, Xiaoyan Li, Donald Metzler, Mark D. Smucker, Trevor Strohman, Howard Turtle, and Courtney Vantonder. 2004. UMass at TREC 2004 : Novelty and HARD. In *Proceedings of the Thirteenth Text REtrieval Conference (TREC 2004)*. NIST. Contextualise les modèles de pertinence utilisés dans TREC, dont RM3 est un exemple.
- ACL Anthology Organization. 2025. ACL Anthology GitHub Repository. <https://github.com/acl-org/acl-anthology/>. Accessed : 2025-04-09.
- ACM Special Interest Group on Information Retrieval (SIGIR). 2025. ACM SIGIR Conference on Research and Development in Information Retrieval. <https://sigir.org>. Accessed : 2025-04-09.
- Anusha. 2021. Natural language processing (nlp) with python – tutorial for beginners. Accessed : 2025-04-09.

- Association for Computational Linguistics. 2025. ACL Anthology. <https://aclanthology.org>. Accessed : 2025-04-09.
- DataCamp. 2023. Stemming and lemmatization in python. <https://www.datacamp.com/tutorial/stemming-lemmatization-python>. Accessed : 2025-04-09.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert : Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv :1810.04805*.
- Donna K. Harman. 1995. Overview of the third text retrieval conference (trec-3). In *Proceedings of the Third Text REtrieval Conference (TREC-3)*, Gaithersburg, USA. NIST Special Publication 500-225. Contexte pour la collection AP.
- Daniel Jurafsky and James H. Martin. 2023. <https://web.stanford.edu/jurafsky/slp3/> *Speech and Language Processing*, 3rd (draft) edition. Online.
- Jimmy Lin, Xueguang Ma, Craig Macdonald, and Rodrigo Nogueira. 2021. <https://doi.org/10.1145/3404835.3463238> Pyserini : A Python Toolkit for Reproducible Information Retrieval Research with Sparse and Dense Representations. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21)*, pages 2356–2360.
- Guangyuan Ma and Pengfei Liu. 2023. <https://arxiv.org/abs/2310.11496> Query generation for information retrieval using large language models : A survey. *arXiv preprint arXiv :2310.11496*.
- Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. <https://nlp.stanford.edu/IR-book/> *Introduction to Information Retrieval*. Cambridge University Press.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems 26 (NIPS 2013)*, pages 3111–3119.
- George A. Miller. 1995. <https://doi.org/10.1145/219717.219748> Wordnet : A lexical database for english. *Communications of the ACM*, 38(11) :39–41.
- National Institute of Standards and Technology (NIST). 2025. Text REtrieval Conference (TREC). <https://trec.nist.gov>. Accessed : 2025-04-09.
- Jian-Yun Nie. 2010. Information retrieval : Implementing and evaluating search engines (online material). <https://www.iro.umontreal.ca/~nie/IR-book/Preface.html>. Accessed : 2025-04-09.
- Rodrigo Nogueira and Jimmy Lin. 2019. <https://arxiv.org/abs/1901.04085> Passage re-ranking with bert. *arXiv preprint arXiv :1901.04085*.
- Stephen E. Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gatford. 1995. Okapi at trec-3. In *Proceedings of the Third Text REtrieval Conference (TREC-3)*, pages 109–126, Gaithersburg, USA. NIST Special Publication 500-225.
- Gerard Salton and Christopher Buckley. 1988. [https://doi.org/10.1016/0306-4573\(88\)90021-0](https://doi.org/10.1016/0306-4573(88)90021-0) Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5) :513–523.
- Stanford University. 2016. CS224n : Natural Language Processing with Deep Learning - Syllabus. <https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1162/syllabus.shtml>. Accessed : 2025-04-09.

Table des matières

1	Introduction	1
2	État de l'Art	2
3	Méthodologie	2
3.1	Étapes Générales de la RI	2
3.2	Implémentation Détaillée	3
4	Données et Ressources	3
4.1	Collection TREC AP 88-90	3
4.2	Topics (Requêtes)	3
4.3	Jugements de Pertinence (Qrels) . .	3
4.4	Outils Logiciels	3
5	Évaluations et Résultats	3
5.1	Protocole d'Évaluation	3
5.2	Résultats du Système de Base (Stem- ming vs Baseline)	4
5.3	Résultats de l'Amélioration (Partie 2 - RM3 sur Baseline)	4
5.4	Analyse des Erreurs (Qualitative) . .	4
6	Conclusion et Perspectives	5