

# Da mettere alla fine degli esperimenti

Alessandro Stefani<sup>1</sup>, Caterina Buranelli<sup>2</sup>, and Cristi Gutu<sup>3</sup>

<sup>1</sup> Corso di laurea in Statistica per le tecnologie e le scienze, matricola 1148387  
`alessandro.stefani.6@studenti.unipd.it`

<sup>2</sup> Corso di laurea in Statistica per le tecnologie e le scienze, matricola 1234567  
`caterina.buranelli@studenti.unipd.it`

<sup>3</sup> Corso di laurea in Statistica per le tecnologie e le scienze, matricola 1147351  
`gheorghecristi.gutu@studenti.unipd.it`

**Sommario** DA FINIRE QUANDO ABBIAMO MESSO A POSTO LA SEZIONE ESPERIMENTI. Questo progetto tratta la realizzazione attraverso il pacchetto *Whoosh*, di un motore di ricerca volto al reperimento di documenti della collezione sperimentale *OHSUMED* indicizzata opportunamente. Il progetto è anche corredato di un webserver che permette all'utente di interrogare il motore di ricerca in forma interattiva attraverso un browser a scelta.

**Keywords:** *Information Retrieval · IR · statistica · reperimento · indicizzazione · Whoosh · Python · webserver · web · interrogazione web*

## 1 Introduzione

Lo scopo finale di un Sistema di Information Retrieval (IRS) è quello di reperire documenti rilevanti relativi a una certa esigenza informativa<sup>4</sup>; dunque i documenti sono il primo ingresso del sistema, mentre il secondo è costituito dalle interrogazioni; i documenti devono essere indicizzati e nell'indice creato, si andrà a effettuare la ricerca per reperire documenti rilevanti<sup>5</sup>. Questa seconda parte è detta reperimento e non si occupa solo di ricercare tra i documenti, ma anche di riordinare secondo un certo ordine di rilevanza. Ciò che descrive i documenti e ciò che descrive le interrogazioni deve essere confrontabile, infatti nei programmi di indicizzazione e di reperimento si usa uno schema, che deve essere uguale in entrambi i casi. L'indicizzazione è un trade-off tra il miglioramento della rappresentazione del contenuto informativo dei documenti (efficacia) e la gestione degli indici (efficienza). Benché ci sia sempre tensione tra queste due caratteristiche, ad oggi la questione più studiata è quella dell'efficienza.

---

<sup>4</sup> insieme delle circostanze in cui una persona ha un problema da risolvere o un compito da svolgere e richiede informazioni importanti, utili o necessarie per la soluzione del problema o lo svolgimento del compito

<sup>5</sup> la rilevanza è la proprietà che rende l'informazione importante, utile o necessaria a soddisfare l'esigenza informativa dell'utente

Esistono diversi modi per risolvere questo problema, nel nostro caso abbiamo ricercato la configurazione migliore tra un sistema di reperimento con o senza uso di stopword e il tuning dei parametri dello schema di pesatura<sup>6</sup> BM25F .

L'obiettivo principale della relazione è da una parte la documentazione del progetto di un servizio di IR e dall'altra una misura del grado in cui si sia riusciti a mettere in pratica i contenuti della disciplina illustrati durante le lezioni.

A tal scopo la relazione dovrà illustrare nelle sezioni successive:

- i metodi di indicizzazione,
- i modelli di reperimento,
- l'interfaccia basata su un *browser* per il WWW
- i risultati della *valutazione* condotta con la collezione sperimentale OHSU-MED.

Il lettore della relazione è lo studente medio di un corso di laurea in statistica al quale la relazione deve dare tutti gli strumenti per comprendere il contenuto. Ci si metta nei suoi panni e si scriva tutto ciò e solo ciò che serve. Chiedersi qual è il messaggio che lo studente deve “portarsi a casa”, esplicitarlo in questo paragrafo e concentrarsi su quello nel resto della relazione.

L'introduzione della relazione deve servire al lettore a capire se vale la pena continuare a leggere il resto. Si possono riassumere i contenuti delle sezioni successive e metterne in evidenza i punti principali. La relazione consiste di tre paragrafi principali dopo questa introduzione e prima della bibliografia, per la quale si suggerisce Bib<sub>TEX</sub> se si scrive con L<sub>A</sub>T<sub>E</sub>X.

## 2 Base di partenza

La base di partenza è formata dai metodi documentati nei libri di testo. Si eviti di trascrivere pari pari, si cerchi piuttosto di rielaborare i contenuti in modo da renderli *coerenti* col resto della relazione; in particolare, si descrivano tutti e solo i metodi usati negli esperimenti e si eviti di parlare di quei metodi che poi non sono stati usati; ad esempio, se si conducono degli esperimenti con BM25F, si deve descrivere questo schema di pesatura in questa sezione.

”Best Match 25 Model with Extension to Multiple Weighted Fields”

Le stopwords[2] sono parole che non portano informazione significativa al contenuto informativo come congiunzioni, articoli, avverbi..

## 3 Esperimenti e Risultati

Per gli esperimenti si è utilizzata parte della già citata collezione sperimentale chiamata OHSUMED<sup>7</sup>.

<sup>6</sup> funzione che assegna per ogni documento diversi livelli d'importanza dei termini mediante dei pesi, che possono variare con l'interrogazione

<sup>7</sup> <https://bit.ly/2wpOynZ>

Si tratta di una collezione di oltre 300'000 articoli provenienti dal database bibliografico MEDLINE pubblicati tra il 1987 ed il 1991, di questi documenti ne sono stati utilizzati 54'711.

I programmi utilizzati per effettuare indicizzazione e reperimento sono stati scritti in python(versione 2.7), sfruttando oggetti e funzioni del modulo whoosh<sup>8</sup>.

Per la valutazione dei risultati ottenuti nei vari esperimenti si è invece utilizzato lo strumento standard trec\_eval<sup>9</sup>. Questo programma permette di ottenere alcune misure della qualità di un sistema di information retrieval se si conoscono i documenti rilevanti per le query sperimentali utilizzate.

Come misura da usare nel confronto tra risultati di diverse run si è scelta la Mean Average Precision(MAP)[3] , quantità calcolabile sia a livello di singola query che a livello complessivo su un insieme di query. Inoltre, per verificare se i valori di MAP ottenuti in run diverse si discostano in modo statisticamente significativo gli uni dagli altri si effettuano test di Wilcoxon[] su coppie di risultati.

### 3.1 Scelte iniziali

Prima di iniziare gli esperimenti sono state fatte delle scelte riguardo certi aspetti dei metodi di reperimento che rimanessero costanti durante tutti gli esperimenti. Queste scelte sono state fatte in seguito ad alcune run preliminari.

Per la funzione di reperimento e ranking come schema di pesatura si è deciso di utilizzare il BM25F, descritto nella sezione precedente, inquanto "state-of-the-art" dal punto di vista dei modelli per information retrieval su documenti strutturati, e come operatore logico per raggruppare le parole delle query si è deciso di utilizzare l'operatore OR, ovvero si cercano documenti in cui è presente almeno una parola della query.

Sempre a seguito delle run preliminari si è deciso di utilizzare il campo "desc" delle query anzichè il campo "title" inquanto dà risultati migliori, inoltre, si è notato che in alcune query sperimentali sono presenti parole scritte in modo errato e che quindi non risultano presenti nell'indice. Per questo dopo aver verificato, leggendo i documenti rilevanti per quelle query, quali fossero gli errori si è ritenuto opportuno effettuare una correzione al momento del reperimento. A grandi linee, la correzione viene fatta su parole che non risultano presenti nell'indice cercando in quest'ultimo delle parole alternative che si discostano di al più una lettera dalla parola originale; queste parole vengono quindi aggiunte alla query di partenza e si procede con la ricerca.

Un'ultima scelta è stata fatta riguardo il numero massimo di documenti reperimenti. Si è optato per 100 documenti poiche un numero maggiore non porta grandi miglioramenti dal punto di vista del MAP.

<sup>8</sup> <https://whoosh.readthedocs.io/en/latest/index.html>

<sup>9</sup> [https://trec.nist.gov/trec\\_eval/](https://trec.nist.gov/trec_eval/)

Affinchè Whoosh possa indicizzare una collezione di documenti, necessita la specificazione di uno schema che include, per ogni possibile campo dei documenti della collezione, il nome del campo ed il tipo del campo.

Lo schema di base per la collezione qui utilizzata è il seguente:

```
schema = Schema(docid      = ID(stored=True),
                title      = TEXT(stored=True),
                identifier  = ID(stored=True),
                terms      = NGRAM(stored=True),
                authors    = NGRAM(stored=True),
                abstract    = TEXT(stored=True),
                publication = TEXT(stored=True),
                source      = TEXT(stored=True))
```

**Figura 1.** Schema necessario all'indicizzazione dei documenti. "stored=True" indica che il campo viene salvato ed è quindi successivamente accessibile.

### 3.2 Baseline

I valori di MAP scelti come baseline sono quelli ottenuti effettuando ricerche su un indice con lo schema di base, riportato sopra, senza l'utilizzo di una lista di stopwords e cercando in due campi dei documenti: "title" e "abstract". Cercando su un campo ("title") o tre campi ("title", "abstract" e "terms") risulta solo in un peggioramento dei risultati.

Si possono vedere i risultati ottenuti usando la configurazione "baseline" riassunti in Tabella 1.

un campo	due campi	tre campi
numq all 63	numq all 63	numq all 63
numret all 6228	numret all 6244	numret all 6271
numrel all 670	numrel all 670	numrel all 670
numrelret all 349	numrelret all 419	numrelret all 367
map all 0.2130	map all <b>0.2744</b>	map all 0.2155

**Tabella 1.** Risultati treceval complessivi per tutte le query, nessuna manipolazione del testo, numero risultati restituiti per ogni query = 100, pesatura BM25F.

Come si nota dalla tabella il valore baseline di *M.A.P* è 0.2744.

Il processo ed il codice per creare l'indice sono facilmente comprensibili visionando il file *indicizzatore\_batch\_baseline.py*.

Per eseguire l'indicizzazione baseline è sufficiente lanciare lo script python con il seguente comando:

```
indicizzatore_batch_baseline.py \
cartella_indice file_documenti.xml
```

Per eseguire il reperimento, che poi produce il file con i risultati in formato compatibile con `trec_eval`, è sufficiente lanciare lo script python con il comando:

```
python search_tk.py cartella_indice \
file_query.xml cartella_risultati nome_file_risultati
```

### 3.3 Esperimenti

Gli esperimenti sono costituiti in tre prove in cui si cambia l'insieme di stopwords usato nell'indicizzazione. In particolare ad ogni prova si utilizza un insieme di stopwords con parole in più rispetto a quello delle prove precedenti.

Come con la baseline, per ciascuna prova si effettua la ricerca per uno, due e tre campi.

**Prima prova:** Per la prima prova sono state utilizzate le stopwords generali della lingua inglese. Queste comprendono termini come "the", "that", "is"; sono poi state aggiunte anche le lettere dell'alfabeto ed i numeri.

Come si vede in Tabella 2, non sembra si sia ottenuto un miglioramento significativo rispetto alla baseline. Il MAP ottenuto utilizzando due campi "migliora" solo dello 0.0001.

un campo	due campi	tre campi
<i>numq all 63</i>	<i>numq all 63</i>	<i>numq all 63</i>
<i>numret all 6228</i>	<i>numret all 6244</i>	<i>numret all 6271</i>
<i>numrel all 670</i>	<i>numrel all 670</i>	<i>numrel all 670</i>
<i>numrelret all 350</i>	<i>numrelret all 418</i>	<i>numrelret all 367</i>
<i>map all 0.2137</i>	<i>map all <b>0.2745</b></i>	<i>map all 0.2152</i>

**Tabella 2.** Risultati treceval, rimozione delle stopwords generali, numero risultati restituiti per ogni query = 100, pesatura BM25F.

**Seconda prova:** Nella seconda prova si è allora cercato di migliorare i risultati aggiungendo alle stopwords generali, alcune parole che non sono normalmente considerate stopwords ma che nel contesto medico, come nel caso della collezione OHSUMED, diventano tali.

Per questo abbiamo utilizzato un insieme di stopwords denominate `stopword_cliniche` che oltre ad avere le stopwords generali ha anche parole come "medical", "condition" o "family"<sup>10</sup>

I risultati sono stati accettabili, ma siamo riusciti a migliorarli togliendo anche stopwords cliniche, cioè strettamente inerenti al contesto medico e abbiamo ottenuto i seguenti risultati:

<sup>10</sup> <https://github.com/kavgan/clinical-concepts/blob/master/clinical-stopwords.txt>

un campo	due campi	tre campi
<i>numq all 63</i>	<i>numq all 63</i>	<i>numq all 63</i>
<i>numret all 6228</i>	<i>numret all 6244</i>	<i>numret all 6271</i>
<i>numrel all 670</i>	<i>numrel all 670</i>	<i>numrel all 670</i>
<i>numrelret all 350</i>	<i>numrelret all 419</i>	<i>numrelret all 376</i>
<i>map all 0.2156</i>	<i>map all <b>0.2837</b></i>	<i>map all 0.2166</i>

**Tabella 3.** Risultati treceval, rimozione delle stopword cliniche, numero risultati restituiti per ogni query = 100, pesatura BM25F.

**Terza prova:** Nell’ultima prova sono state aggiunte altre parole alla lista delle stopword, la scelta delle parole da aggiungere è stata fatta, prima prendendo in considerazione le parole delle query e la frequenza con cui queste si trovano nei documenti rilevanti e poi cercando la presenza di possibili stopword nelle query che davano risultati peggiori.

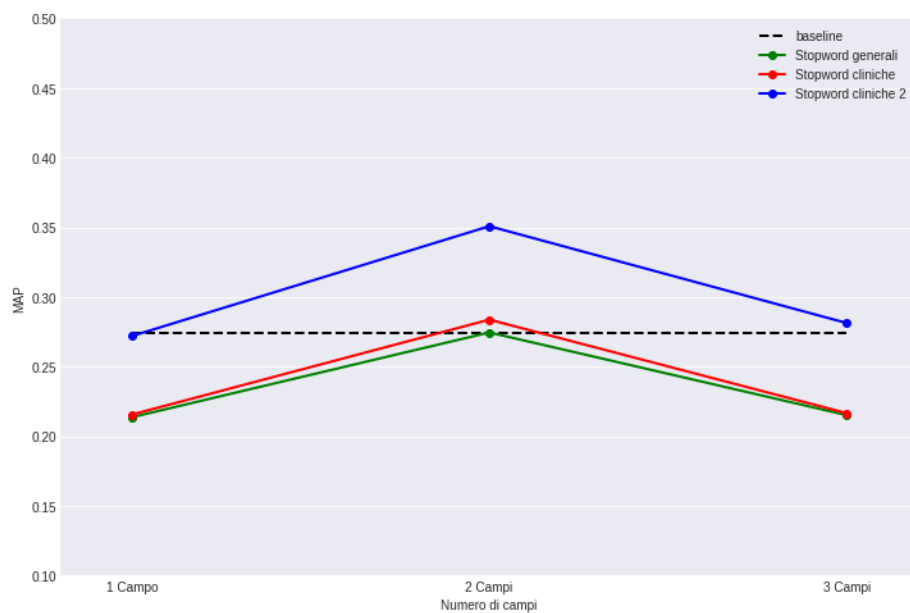
I risultati sono stati accettabili, ma siamo riusciti a migliorarli togliendo anche stopword cliniche, cioè strettamente inerenti al contesto medico e abbiamo ottenuto i seguenti risultati:

un campo	due campi	tre campi
<i>numq all 63</i>	<i>numq all 63</i>	<i>numq all 63</i>
<i>numret all 5282</i>	<i>numret all 5690</i>	<i>numret all 5788</i>
<i>numrel all 670</i>	<i>numrel all 670</i>	<i>numrel all 670</i>
<i>numrelret all 386</i>	<i>numrelret all 463</i>	<i>numrelret all 435</i>
<i>map all 0.2720</i>	<i>map all <b>0.3508</b></i>	<i>map all 0.2813</i>

**Tabella 4.** Risultati treceval, rimozione delle stopword cliniche migliorate, numero risultati restituiti per ogni query = 100, pesatura BM25F.

## Riferimenti bibliografici

1. W. Bruce Croft and Donald Metzler and Trevor Strohman. Search Engines: Information Retrieval in Practice. Addison Wesley, (2009), pp. 250-252
2. W. Bruce Croft and Donald Metzler and Trevor Strohman. Search Engines: Information Retrieval in Practice. Addison Wesley, (2009), pp. 90
3. W. Bruce Croft and Donald Metzler and Trevor Strohman. Search Engines: Information Retrieval in Practice. Addison Wesley, (2009), pp. 313
4. Discovering Related Clinical Concepts Using Large Amounts of Clinical Notes. Ganesan, Kavita and Lloyd, Shane and Sarkar, Vikren, (2016), pp. 27-33



**Figura 2.** Performace modelli di classificazione sul test set in funzione della dimensione del training set.



(a) Alessandro Stefani



(b) Caterina Buranelli



(c) Cristi Gutu