

ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΑΤΡΩΝ · ΤΜΗΜΑ ΜΗΧΑΝΙΚΩΝ Η/Υ ΚΑΙ ΠΛΗΡΟΦΟΡΙΚΗΣ

# ΑΝΑΚΤΗΣΗ ΠΛΗΡΟΦΟΡΙΑΣ

ΕΡΓΑΣΤΗΡΙΑΚΗ ΑΣΚΗΣΗ · 2023 – 2024

**ΠΕΡΙΕΧΟΜΕΝΑ**

<b>1</b>	<b>ΕΙΣΑΓΩΓΗ</b>	<b>2</b>
1.1	ΕΠΙΛΟΓΗ ΣΥΣΤΗΜΑΤΩΝ ΣΤΑΘΜΙΣΗΣ TF-IDF . . . . .	2
1.1.1	ΠΡΩΤΟ ΣΥΣΤΗΜΑ ΣΤΑΘΜΙΣΗΣ . . . . .	2
1.1.2	ΔΕΥΤΕΡΟ ΣΥΣΤΗΜΑ ΣΤΑΘΜΙΣΗΣ . . . . .	2
<b>2</b>	<b>ΥΛΟΠΟΙΗΣΗ ΜΟΝΤΕΛΩΝ</b>	<b>3</b>
2.1	ΠΡΟΕΠΕΞΕΡΓΑΣΙΑ ΕΓΓΡΑΦΩΝ & ΒΟΗΘΗΤΙΚΕΣ ΣΥΝΑΡΤΗΣΕΙΣ . . . . .	3
2.2	ΑΝΕΣΤΡΑΜΜΕΝΟ ΕΥΡΕΤΗΡΙΟ . . . . .	4
2.3	ΥΛΟΠΟΙΗΣΗ VECTOR SPACE ΜΟΝΤΕΛΟΥ . . . . .	4
2.4	ΥΛΟΠΟΙΗΣΗ coIBERT . . . . .	5
<b>3</b>	<b>ΜΕΤΡΙΚΕΣ ΑΞΙΟΛΟΓΗΣΗΣ</b>	<b>5</b>
3.1	ΑΠΟΤΕΛΕΣΜΑΤΑ . . . . .	6

## 1 ΕΙΣΑΓΩΓΗ

### 1.1 ΕΠΙΛΟΓΗ ΣΥΣΤΗΜΑΤΩΝ ΣΤΑΘΜΙΣΗΣ TF-IDF

Καταρχάς πρέπει να επιλέξουμε δύο συστήματα στάθμισης των βαρών για τους όρους των εγγράφων και των ερωτημάτων.

#### 1.1.1 ΠΡΩΤΟ ΣΥΣΤΗΜΑ ΣΤΑΘΜΙΣΗΣ

Το πρώτο σύστημα στάθμισης είναι μια παραλλαγή<sup>1</sup> του προτεινόμενου ως καλύτερου πλήρως σταθμισμένου συστήματος σύμφωνα με τους Salton-Buckley<sup>2</sup> (best fully weighted system). Θα χρησιμοποιήσουμε την **απλή συχνότητα εμφάνισης** (raw term frequency) για το TF βάρος των εγγράφων,

$$\text{Σύστημα \#1: } TF_{\text{εγγράφων}} = f_{i,j}$$

όπου  $f_{i,j}$  οι φορές που ο όρος εμφανίζεται σε ένα έγγραφο, τη **διπλή 0,5 κανονικοποίηση** για το TF βάρος των ερωτημάτων (augmented normalized TF),

$$\text{Σύστημα \#1: } TF_{\text{ερωτημάτων}} = 0.5 + 0.5 \frac{f_{i,j}}{\max_i f_{i,j}}$$

και τέλος την **απλή ανάστροφη συχνότητα εμφάνισης** για το IDF βάρος και των εγγράφων και των ερωτημάτων:

$$\text{Σύστημα \#1: } IDF_{\text{ερωτημάτων}} = \log \frac{N}{n_i}$$

όπου  $N$  το πλήθος των εγγράφων και  $n_i$  ο αριθμός των εγγράφων στα οποία εμπεριέχεται ο όρος.

#### 1.1.2 ΔΕΥΤΕΡΟ ΣΥΣΤΗΜΑ ΣΤΑΘΜΙΣΗΣ

Ως δεύτερο σύστημα στάθμισης θα χρησιμοποιήσουμε το καλύτερα σταθμισμένο πιθανολογικό σύστημα σύμφωνα με τους Salton-Buckley<sup>1</sup> (best weighted probabilistic weight) με

$$\text{Σύστημα \#2: } \text{βάρος όρου}_{\text{εγγράφων}} = 0.5 + 0.5 \frac{f_{i,j}}{\max_i f_{i,j}}$$

$$\text{Σύστημα \#2: } \text{βάρος όρου}_{\text{ερωτημάτων}} = \log \frac{N - n_i}{n_i}$$

Και τα δύο αυτά συστήματα στάθμισης έχουν επιφέρει τα ακριβέστερα αποτελέσματα στο σύνολο των συλλογών στα οποία έχουν εξεταστεί αλληλά και ειδικότερα για ιατρικές (MED) συλλογές. Επομένως, συνολικά έχουμε:

<sup>1</sup>Το σύστημα αναφέρεται ως παραλλαγή των Salton-Buckley για το λόγο ότι δεν έχει συμπεριληφθεί κάποιος παράγοντας κανονικοποίησης, μιας και τα έγγραφα είναι περίπου ισομεγέθη (μέσος όρος 350 λέξεις).

<sup>2</sup>Gerard Salton, Christopher Buckley, Term-weighting approaches in automatic text retrieval, Information Processing & Management, Volume 24, Issue 5, 1988, Pages 513-523, ISSN 0306-4573

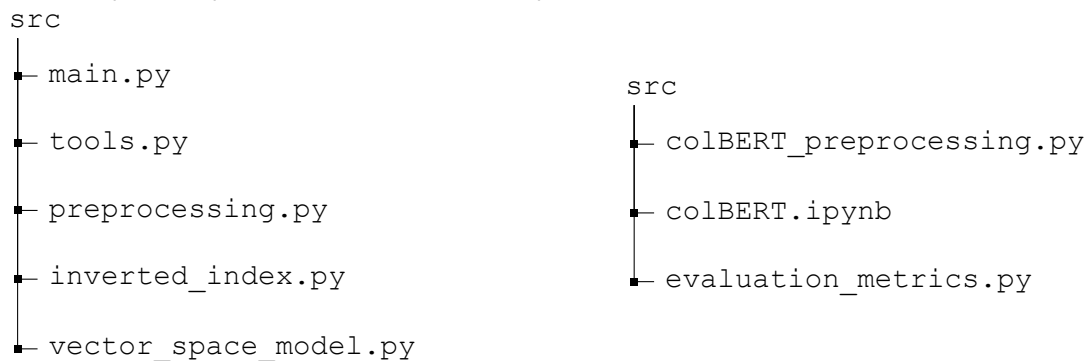
Σύστημα στάθμισης	Βάρος όρου εγγράφου	Βάρος όρου ερωτήματος
1	$f_{i,j} \times \log \frac{N}{n_i}$	$(0.5 + 0.5 \frac{f_{i,j}}{\max_i f_{i,j}}) \times \log \frac{N}{n_i}$
2	$0.5 + 0.5 \frac{f_{i,j}}{\max_i f_{i,j}}$	$\log \frac{N - n_i}{n_i}$

## 2 ΥΛΟΠΟΙΗΣΗ ΜΟΝΤΕΛΩΝ

Η εργασία υλοποιήθηκε σε Python χρησιμοποιώντας τις βιβλιοθήκες:

Βιβλιοθήκη	Περιγραφή
<b>os</b>	σύνδεση με λειτουργικό σύστημα
<b>nltk</b>	αφαίρεση stopwords, stemming
<b>numpy</b>	υπολογισμός ομοιότητας συνημιτόνου
<b>math</b>	υπολογισμός λογαρίθμων
<b>matplotlib</b>	γραφικές παραστάσεις
<b>json</b>	αποθήκευση-ανάγνωση JSON αρχείων
<b>pickle</b>	αποθήκευση αντικειμένων Python

Αυτή είναι η δομή των αρχείων κώδικα όπου έχει χωριστεί η υλοποίηση:



Η υλοποίηση του colBERT μοντέλου έχει πραγματοποιηθεί στο Google Colab ως Jupyter Notebook.

### 2.1 ΠΡΟΕΠΕΞΕΡΓΑΣΙΑ ΕΓΓΡΑΦΩΝ & ΒΟΗΘΗΤΙΚΕΣ ΣΥΝΑΡΤΗΣΕΙΣ

Το αρχείο `tools.py` περιλαμβάνει βοηθητικές συναρτήσεις για κάποιες επαναλαμβανόμενες διαδικασίες της υλοποίησης. Περιλαμβάνονται οι συναρτήσεις `get_docs()` και `get_queries()`.

Η συνάρτηση `get_docs()`, χρησιμοποιώντας την `os` βιβλιοθήκη διαβάζει το πλήθος των αρχείων της βιβλιοθήκης.<sup>3</sup> Η συνάρτηση, αφού αφαιρέσει το escape character `'\n'`, ο οποίος προκύπτει από την μορφολογία των εγγράφων (η κάθε λέξη είναι τοποθετημένη σε διαφορετική γραμμή), δημιουργεί

<sup>3</sup>Να σημειωθεί ότι το πλήθος των εγγράφων διαφέρει από την αύξουσα αριθμότητά τους. Συγκεκριμένα έχουμε 1209 έγγραφα αριθμημένα από το 000001 ως 01239. Με άλλα λόγια υπάρχουν αριθμοί στη συλλογή που δεν αντιστοιχούν σε έγγραφα. Συνεπώς δεν θα μπορούσαμε να χρησιμοποιήσουμε κάποια αριθμητική επανάληψη, για παράδειγμα, για την εισαγωγή των εγγράφων.

και επιστρέφει μια λίστα από tuples, με κάθε tuple να αντιστοιχεί σε κάθε αρχείο-έγγραφο. Τα tuples έχουν την δομή:

```
doc_tuple = ('docID', ['λήμμα_1', 'λήμμα_2' ...])
```

όπου docID η αρίθμηση του κάθε εγγράφου και doc\_term\_n η κάθε λέξη-λήμμα του εγγράφου. Η συνάρτηση strip() είναι απαραίτητη για την αφαίρεση των \n χαρακτήρων που προέκυψαν από την μορφολογία των εγγράφων (κάθε λέξη είναι σε νέα γραμμή). Με παρόμοιο τρόπο η συνάρτηση **get\_queries()** επιστρέφει τη λίστα με τα ερωτήματα της συλλογής.

Στις συναρτήσεις **preprocess\_collection()** και **preprocess\_queries()** πραγματοποιείται η προεπεξεργασία των εγγράφων, συγκεκριμένα η αφαίρεση των stopwords και το stemming.

Χρησιμοποιούμε τη nltk βιβλιοθήκη. Τα doc\_tuples της get\_docs() αφού περάσουν από τον PorterStemmer της nltk αποθηκεύονται σε μια λίστα, η οποία στη συνέχεια επιστρέφεται. Αντίστοιχη διαδικασία πραγματοποιείται και για την προεπεξεργασία των ερωτημάτων, στην preprocess\_queries().

## 2.2 ΑΝΕΣΤΡΑΜΜΕΝΟ ΕΥΡΕΤΗΡΙΟ

Το ανεστραμμένο ευρετήριο δημιουργείται στη συνάρτηση **create\_inverted\_index()** του αρχείου inverted\_index.py. Στη συνάρτηση εισάγονται οι λίστες που δημιουργήθηκαν στις προηγούμενες συναρτήσεις. Τα tuples από τα οποία αποτελούνται, αποθηκεύονται σε ένα dictionary που θα αποτελέσει το ανεστραμμένο ευρετήριο με την εξής δομή:

```
inverted_index['λήμμα'] =
{ ('docID στο οποίο εμφανίζεται' = <φορές εμφάνισης>), (...), ... }
```

Κάθε value του dictionary είναι ένα set<sup>4</sup> το οποίο περιλαμβάνει ένα ή περισσότερα tuples με το docID και τη συχνότητα εμφάνισης του λήμματος στο συγκεκριμένο έγγραφο. Η συχνότητα υπολογίζεται μέσω της count() σε όλο το έγγραφο ανά λήμμα.

Αυτό είναι ένα παράδειγμα του τελικού ανεστραμμένου ευρετηρίου<sup>5</sup>:

```
inverted_index = { ... 'coronari': {('01217', 2), ('00779', 1)},
                  'mobil': {('00673', 2)}, 'strain': {('00179', 7)}, ... }
```

## 2.3 ΥΛΟΠΟΙΗΣΗ VECTOR SPACE ΜΟΝΤΕΛΟΥ

Η υλοποίηση του vector space μοντέλου γίνεται στο αρχείο vector\_space\_model.py. Η συνάρτηση **run\_vsm(doc\_collection, queries, inverted\_index)**, μέσω μιας επανάληψης ώστε να καλυφθούν όλα τα queries, καλεί την συνάρτηση **vsm1(doc\_collection, query, inverted\_index)** ή **vsm2()**, ανάλογα με το σύστημα στάθμισης που θέλουμε να χρησιμοποιήσουμε.

Πριν προχωρήσουμε στις συναρτήσεις vsmX(), από το τρόπο που έχει δομηθεί το ανεστραμμένο ευρετήριο, είναι σαφές πως μπορούμε να υπολογίσουμε αμέσως το πλήθος των εγγράφων όπου υπάρχει ένα συγκεκριμένο λήμμα ως len(inverted\_index[term]), με άλλα λόγια δηλαδή το  $n_i$ . Είναι επίσης προφανές ότι  $N = \text{len}(\text{inverted\_index})$ .

Στην συνάρτηση vsm1() αρχικά υπολογίζονται στο λεξικό query\_tfs οι συχνότητες  $f_{i,j}$  κάθε όρου του εκάστοτε ερωτήματος:

<sup>4</sup>Έχει επιλεχθεί set για εξοικονόμηση μνήμης, μας και δεν μας ενδιαφέρει η σειρά των tuples.

<sup>5</sup>Τα λήμματα έχουν τη stemming μορφή τους.

```
query_tfs len(query) = { ... 'calcium': 1, 'effect': 2 ... }
```

Στα λεξικά `query_tfidf` και `doc_tfidf` θα αποθηκευτούν οι τελικές τιμές-βάρη που θα αποτελέσουν τα διανύσματα ερωτήματος και εγγράφων. Το λεξικό `doc_tfidf` αρχικοποιείται με άδειες λίστες για κάθε έγγραφο.

Διατρέχουμε κάθε όρο του ερωτήματος `query` και αν ο όρος υπάρχει στο ανεστραμμένο ευρετήριο υπολογίζουμε την IDF τιμή του στην συνάρτηση `idfX()`, και στη συνέχεια το TF-IDF βάρος του ερωτήματος στο `query_tfidf`.

```
query_ifidf len(query) = { ... 'calcium': 3.6318, 'effect': 1.76483 ... }
```

Στη συνέχεια διατρέχουμε όλα τα έγγραφα της συλλογής. Αν το έγγραφο περιλαμβάνει τον όρο του ερωτήματος που εξετάζουμε (κάτι που ελέγχουμε από τη λίστα `docs_containing_term`), τότε μέσω του ανεστραμμένου ευρετηρίου, αποθηκεύεται η TF τιμή του, υπολογίζεται το TF-IDF βάρος του εγγράφου και αποθηκεύεται στο `doc_tfidf[docID]`. Αν το έγγραφο δεν περιλαμβάνει τον όρο, αποθηκεύουμε 0.

```
doc_ifidf len(doc_collection) = {'000001': [0, 0, 0.0, 1.76, 1.5] len(query) ...}
```

Εν τέλει έχουμε δημιουργήσει το διάνυσμα ερωτήματος και 1209 διανύσματα εγγράφων, τα οποία έχουν μήκος όσο και το διάνυσμα ερωτήματος. Στην λίστα `similarity` υπολογίζουμε την ομοιότητα συνημιτόνου μεταξύ του διανύσματος ερωτήματος και κάθε διανύσματος εγγράφου. Τέλος, ταξινομούμε την λίστα και επιστρέφουμε τα 100 κείμενα με την μεγαλύτερη ομοιότητα.

Η ίδια διαδικασία ακολουθείται και στην συνάρτηση `vsm2()`, με τις απαραίτητες αλλαγές για το διαφορετικό σύστημα στάθμισης.

## 2.4 ΥΛΟΠΟΙΗΣΗ colBERT

Στο αρχείο `colBERT_preprocessing.py` γίνεται η προεπεξεργασία των εγγράφων και ερωτημάτων ώστε να δωθούν ως `inputs` στο colBERT μοντέλο. Συγκεκριμένα έγινε μετατροπή των ερωτημάτων σε κεφαλαία και αποθηκεύτηκαν σε JSON αρχεία τα έγγραφα και ερωτήματα σε λεξικά ως εξής:

```
colBERTdocs = {"docIDi": "<doci>", ... }
colBERTqueries = {"queryIDi": "<queryi>", ... }
```

Η εκτέλεση του colBERT μοντέλου πραγματοποιείται μέσω του αρχείου `colBERT.ipynb` στο Google Colab. Τα λεξικά των JSON αρχείων μετατρέπονται σε 2 ξεχωριστές λίστες (`id` και `περιεχόμενο`). Μετά τη δημιουργία του `Indexer` και του `Searcher`, επιστρέφονται τα 100 κείμενα με την μεγαλύτερη ομοιότητα.

## 3 ΜΕΤΡΙΚΕΣ ΑΞΙΟΛΟΓΗΣΗΣ

Για την σύγκριση των μοντέλων χρησιμοποιήθηκαν:

- **Διάγραμμα ανάκλησης-ακρίβειας (Precision-Recall curve).** Το διάγραμμα παρουσιάζει ενδιαφέρον γιατί -συνδυάζοντας τις μετρικές ανάκλησης και ακρίβειας- παρουσιάζει με ακρίβεια τα σημεία που το μοντέλο παρουσιάζει μεγαλύτερη ακρίβεια και λιγότερη ανάκληση και αντίστροφα.

$$\text{Ακρίβεια} = \frac{\text{αριθμός σχετικών ανακτηθέντων κειμένων}}{\text{αριθμός εγγράφων που ανακτήθηκαν}}$$

$$\text{Ανάκληση} = \frac{\text{αριθμός σχετικών ανακτηθέντων κειμένων}}{\text{αριθμός σχετικών εγγράφων στη συλλογή}}$$

- **??????? (Mean Average Precision – mAP).** Χρησιμοποιούμε την συγκεκριμένη μετρική λόγω των πολλών ερωτημάτων όπου εφαρμόζουμε τα μοντέλα σε αυτά, για την αξιολόγηση των αποτελεσμάτων.

$$\sum_{i=1}^{i=k} (\text{ανάκληση}[i] - \text{ανάκληση}[i - 1]) \times \text{ακρίβεια}[i]$$

### 3.1 ΑΠΟΤΕΛΕΣΜΑΤΑ