

ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΑΤΡΩΝ · ΤΜΗΜΑ ΜΗΧΑΝΙΚΩΝ Η/Υ ΚΑΙ ΠΛΗΡΟΦΟΡΙΚΗΣ

ΑΝΑΚΤΗΣΗ ΠΛΗΡΟΦΟΡΙΑΣ

ΕΡΓΑΣΤΗΡΙΑΚΗ ΑΣΚΗΣΗ · 2023 – 2024

ΠΕΡΙΕΧΟΜΕΝΑ

1	ΕΙΣΑΓΩΓΗ	2
1.1	ΕΠΙΛΟΓΗ ΣΥΣΤΗΜΑΤΩΝ ΣΤΑΘΜΙΣΗΣ TF-IDF	2
1.1.1	ΠΡΩΤΟ ΣΥΣΤΗΜΑ ΣΤΑΘΜΙΣΗΣ	2
1.1.2	ΔΕΥΤΕΡΟ ΣΥΣΤΗΜΑ ΣΤΑΘΜΙΣΗΣ	2
2	ΥΛΟΠΟΙΗΣΗ ΜΟΝΤΕΛΩΝ	3
2.1	ΠΡΟΕΠΕΞΕΡΓΑΣΙΑ ΕΓΓΡΑΦΩΝ & ΒΟΗΘΗΤΙΚΕΣ ΣΥΝΑΡΤΗΣΕΙΣ	3
2.1.1	tools.py	3
2.1.2	preprocessing.py	4
2.2	ΑΝΕΣΤΡΑΜΜΕΝΟ ΕΥΡΕΤΗΡΙΟ	4
2.3	ΥΛΟΠΟΙΗΣΗ VECTOR SPACE ΜΟΝΤΕΛΟΥ	4
2.4	ΥΛΟΠΟΙΗΣΗ colBERT	5
3	ΜΕΤΡΙΚΕΣ ΑΞΙΟΛΟΓΗΣΗΣ	5
3.1	ΑΠΟΤΕΛΕΣΜΑΤΑ	6
4	ΠΑΡΑΡΤΗΜΑ	9
4.0.1	tools.py	9
4.0.2	preprocessing.py	10
4.0.3	inverted_index.py	10
4.0.4	vector_space_model.py	10
4.0.5	colBERT_preprocessing.py	12
4.0.6	colBERT_.ipynb	12
4.0.7	evaluation_metrics.py	14
4.0.8	main.py	15

1 ΕΙΣΑΓΩΓΗ

1.1 ΕΠΙΛΟΓΗ ΣΥΣΤΗΜΑΤΩΝ ΣΤΑΘΜΙΣΗΣ TF-IDF

Καταρχάς πρέπει να επιλέξουμε τα δύο συστήματα στάθμισης των βαρών για τα διανύσματα που θα χρησιμοποιήσουμε.

1.1.1 ΠΡΩΤΟ ΣΥΣΤΗΜΑ ΣΤΑΘΜΙΣΗΣ

Το πρώτο σύστημα στάθμισης θα είναι μια παραλλαγή¹ του προτεινόμενου ως καλύτερου πλήρως σταθμισμένου συστήματος σύμφωνα με τους Salton-Buckley² (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-Buckley¹ (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 και στο σύνολο των συλλογών στα οποία έχουν εξεταστεί αλλά και ειδικότερα σε ιατρικές (MED) συλλογές. Επομένως, συνολικά έχουμε:

¹Το 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).

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

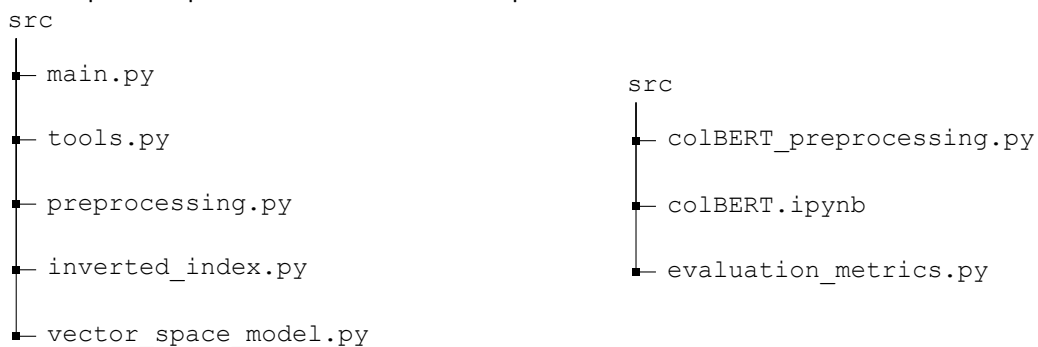
Σύστημα στάθμισης	Βάρος όρου εγγράφου	Βάρος όρου ερωτήματος
1	$f_{i,j} \times \log \frac{N}{n_i}$	$(0.5 + 0.5 \frac{f_{i,j}}{\max_i f_{i,j}}) \times \log \frac{N}{n_i}$
2	$0.5 + 0.5 \frac{f_{i,j}}{\max_i f_{i,j}}$	$\log \frac{N - n_i}{n_i}$

2 ΥΛΟΠΟΙΗΣΗ ΜΟΝΤΕΛΩΝ

Η εργασία υλοποιήθηκε σε Python χρησιμοποιώντας τις βιβλιοθήκες:

Βιβλιοθήκη	Περιγραφή
os	σύνδεση με λειτουργικό σύστημα
json	αποθήκευση-ανάγνωση JSON αρχείων
nltk	αφαίρεση stopwords, stemming
numpy	υπολογισμός ομοιότητας συνημιτόνου
math	υπολογισμός λογαρίθμων
matplotlib	γραφικές παραστάσεις

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



Τέλος, η υλοποίηση του colBERT μοντέλου έχει πραγματοποιηθεί στο Google Colab ως Jupyter Notebook.

2.1 ΠΡΟΕΠΕΞΕΡΓΑΣΙΑ ΕΓΓΡΑΦΩΝ & ΒΟΗΘΗΤΙΚΕΣ ΣΥΝΑΡΤΗΣΕΙΣ

2.1.1 tools.py

Το αρχείο `tools.py` περιλαμβάνει βοηθητικές συναρτήσεις για κάποιες επαναλαμβανόμενες διαδικασίες της υλοποίησης.

Η συνάρτηση `get_docs()`, χρησιμοποιώντας την `os` βιβλιοθήκη, διαβάζει το πλήθος των αρχείων της συλλογής.³ Η συνάρτηση, αφού αφαιρέσει το escape character `'\n'`, το οποίο προκύπτει από την μορφολογία

³Να σημειωθεί ότι το πλήθος των εγγράφων διαφέρει από την αύξουσα αρίθμησή 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.

των εγγράφων (η κάθε λέξη είναι τοποθετημένη σε διαφορετική γραμμή), δημιουργεί και επιστρέφει μια λίστα από tuples, με κάθε tuple να αντιστοιχεί σε κάθε αρχείο-έγγραφο. Έτσι τα έγγραφα έχουν τη δομή:

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

όπου docID η αρίθμηση του κάθε εγγράφου και λήμμα_n η κάθε λέξη-λήμμα του εγγράφου. Αυτός θα είναι ο τρόπος αποθήκευσης και προσπέλασης των εγγράφων σε όλη την υλοποίηση.

Η συνάρτηση `strip()` της `get_docs()` είναι απαραίτητη για την αφαίρεση των `\n` χαρακτήρων που προκύπτουν από τη μορφολογία των εγγράφων (μιας και κάθε λέξη είναι σε νέα γραμμή). Με παρόμοιο τρόπο η συνάρτηση `get_queries()` επιστρέφει τη λίστα με τα ερωτήματα της συλλογής, η `get_relevant()` τη λίστα με τα σχετικά έγγραφα ανά ερώτημα και η `get_json_file()` τα περιεχόμενα ενός JSON αρχείου.

2.1.2 preprocessing.py

Στο αρχείο `preprocessing.py`, και συγκεκριμένα στις συναρτήσεις `preprocess_collection()` και `preprocess_queries()` πραγματοποιείται η προεπεξεργασία των εγγράφων, συγκεκριμένα η αφαίρεση των stopwords και το stemming.

Χρησιμοποιούμε τη `nltk` βιβλιοθήκη. Κάθε λέξη από τη συλλογή των εγγράφων αφού περάσουν από τον `PorterStemmer` της `nltk` και ελεγχθούν ότι δεν ανήκουν στη λίστα των stopwords, επιστρέφονται σε παρόμοια μορφή όπως δημιουργήθηκαν στη `get_docs()`. Αντίστοιχη διαδικασία πραγματοποιείται και για την προεπεξεργασία των ερωτημάτων, στη `preprocess_queries()`.

2.2 ΑΝΕΣΤΡΑΜΜΕΝΟ ΕΥΡΕΤΗΡΙΟ

Το ανεστραμμένο ευρετήριο δημιουργείται στη συνάρτηση `create_inverted_index()` του αρχείου `inverted_index.py`. Τα έγγραφα, μετά την προεπεξεργασία τους, εισέρχονται σε μια δομή επανάληψης, η οποία δημιουργεί το ανεστραμμένο ευρετήριο ως μορφή λεξικού ως εξής:

```
inverted_index['λήμμα'] =
{('docID1': <φορές εμφάνισης>), ('docID2': <φορές εμφάνισης>), ... }
```

Κάθε value του dictionary είναι ένα σύνολο⁴ 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:

```
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 `vsm(doc_collection, query, inverted_index, model_type)`, όπου `model_type = '1' ή '2'`, ανάλογα με το σύστημα στάθμισης που θέλουμε να χρησιμοποιήσουμε.

Η συνάρτηση `vsm()` αρχικά υπολογίζει τις συχνότητες $f_{i,j}$ κάθε όρου του εκάστοτε ερωτήματος:

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

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

Στη συνέχεια διατρέχεται κάθε όρος από το ερώτημα `query`. Αν ο όρος του ερωτήματος υπάρχει στο ανεστραμμένο ευρετήριο, υπολογίζεται η IDF τιμή του. και στη συνέχεια το TF-IDF βάρος του ερωτήματος, όπου θα αποθηκευτεί στο λεξικό `query_tfidf`.

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

Στη συνέχεια, για κάθε όρο του ερωτήματος, διατρέχουμε όλα τα έγγραφα της συλλογής. Αν το έγγραφο περιλαμβάνει τον όρο του ερωτήματος που εξετάζουμε (κάτι που ελέγχουμε μέσω της `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_idflen(doc_collection) = {'000001': [0, 0, 0.0, 1.76, 1.5]len(query) ...}
```

Η διαδικασία επαναλαμβάνεται για κάθε όρο του ερωτήματος. Επειδή τα δύο μοντέλα έχουν διαφορετικό σύστημα στάθμισης, ο υπολογισμός των βαρών βρίσκεται μέσα σε συνθήκες, ανάλογα με το `model_type` που χρησιμοποιούμε.

Το αποτέλεσμα είναι να έχουμε δημιουργήσει το διάνυσμα ερωτήματος και 1209 διανύσματα εγγράφων. Στην λίστα `similarity` υπολογίζουμε την ομοιότητα συνημιτόνου μεταξύ του διανύσματος ερωτήματος και κάθε διανύσματος εγγράφου. Επειδή κάποια διανύσματα αποτελούνται από μηδενικά, στους υπολογισμούς οδηγούμαστε σε `nan` τιμές, που μετατρέπουμε σε 0. Η συνάρτηση εν τέλει επιστρέφει τα 100 κείμενα με την μεγαλύτερη ομοιότητα.

2.4 ΥΛΟΠΟΙΗΣΗ colBERT

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

```
colBERTdocs = {"docIDi": "<docii": "<queryi

```

Η εκτέλεση του colBERT μοντέλου πραγματοποιείται μέσω του αρχείου `colBERT.ipynb` στο Google Colab. Τα λεξικά των JSON αρχείων αφού διαβαστούν, διασπώνται σε 2 ξεχωριστές λίστες: `doc_list.keys()` για τα IDs, `doc_list.values()` για το περιεχόμενο και αντίστοιχα για τα ερωτήματα.

Μετά τη δημιουργία του `Indexer` και του `Searcher` με `doc_maxlen = 300`, `kmeans_niters = 20` και `k = 100`, δημιουργείται μια λίστα τα 100 πιο σχετικά κείμενα, ο οποία γίνεται `download` ως JSON αρχείο.

3 ΜΕΤΡΙΚΕΣ ΑΞΙΟΛΟΓΗΣΗΣ

Μιας και ως δεδομένα έχουμε τα πραγματικά σχετικά έγγραφα (`Relevant_20`), 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 προφανή τρόπο τα σημεία που το μοντέλο παρουσιάζει μεγαλύτερη ακρίβεια και ταυτόχρονα λιγότερη ανάκληση και αντίστροφα.

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

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

- **Μέσος Όρος Μέσης Ακρίβειας (Mean Average Precision – MAP).** Χρησιμοποιούμε την συγκεκριμένη μετρική για τη δημιουργία μιας συνοπτικής ενιαίας τιμής της κατάταξης, χρησιμοποιώντας τα πολλαπλά ερωτήματα που έχουμε.

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

$$\text{MAP} = \frac{\text{Μέση Ακρίβεια για κάθε ερώτημα}}{\text{Σύνολο ερωτημάτων}}$$

Η συνάρτηση `recall_precision()` με εισόδους τα αποτελέσματα των μοντέλων και τα πραγματικά σχετικά έγγραφα, υπολογίζει τις τιμές ανάκλησης και ακρίβειας για κάθε ερώτημα και τις επιστρέφει σε μία λίστα. Η λίστα αυτή εισάγεται στην `mean_average_precision()` όπου υπολογίζεται ο Μέσος Όρος Μέσης Ακρίβειας.

3.1 ΑΠΟΤΕΛΕΣΜΑΤΑ

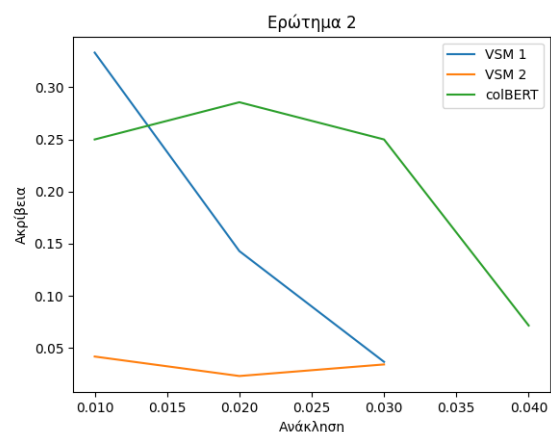
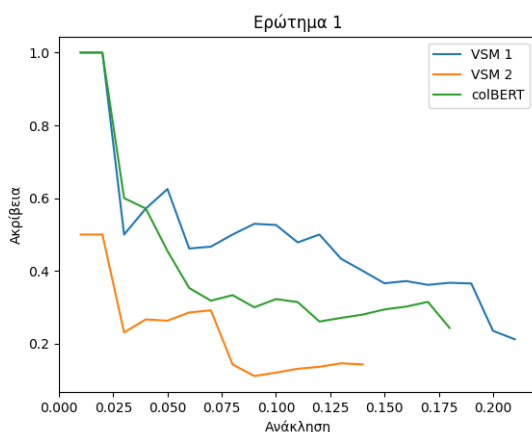
Οι τιμές του Μέσου Όρου Μέσης Ακρίβειας (MAP) για τα τρία μοντέλα είναι:

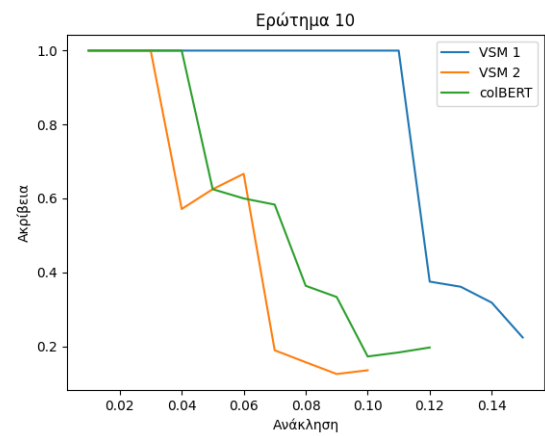
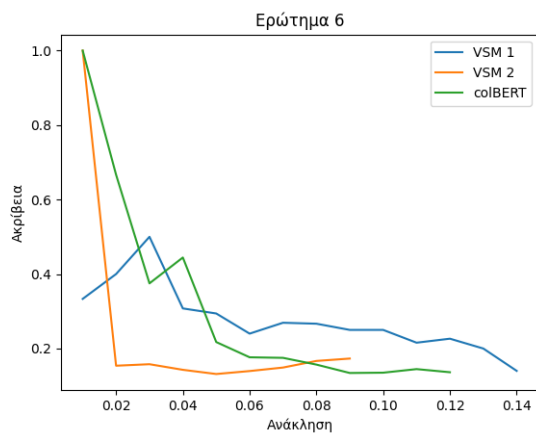
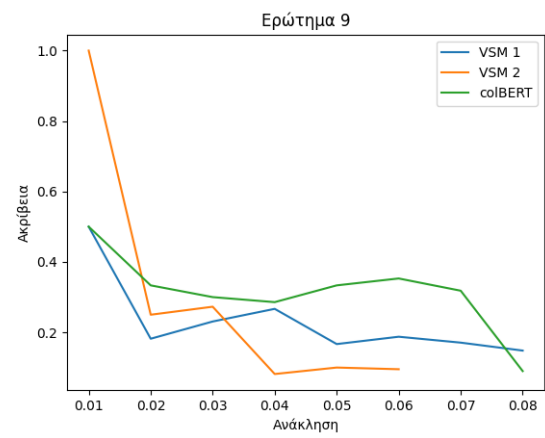
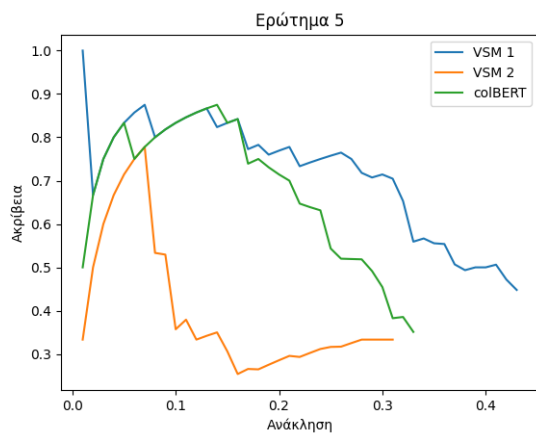
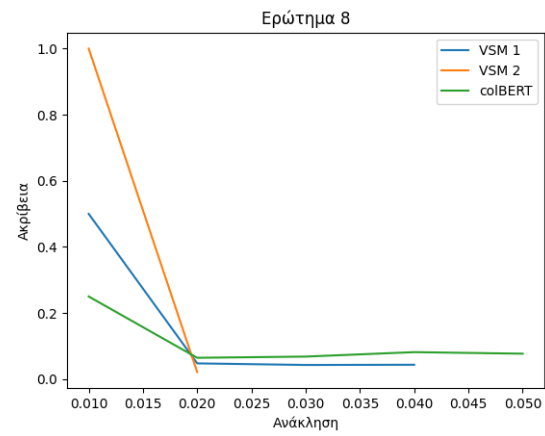
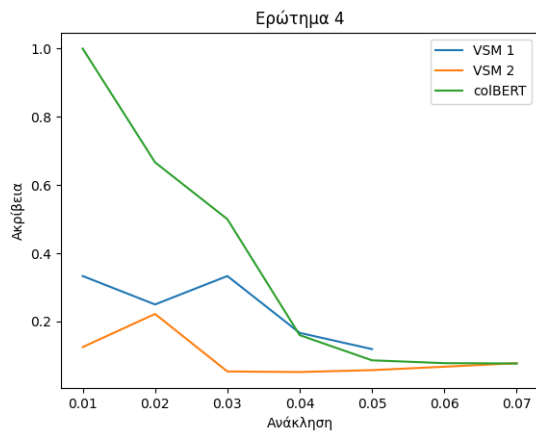
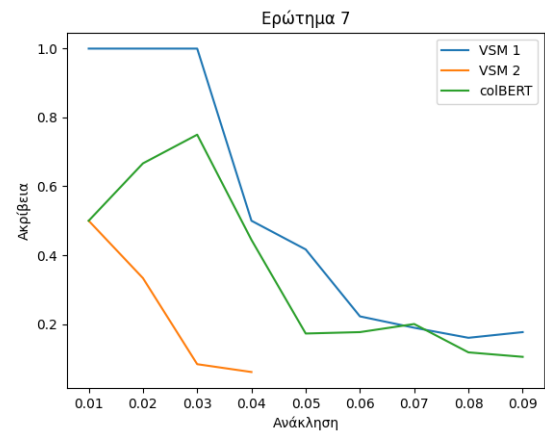
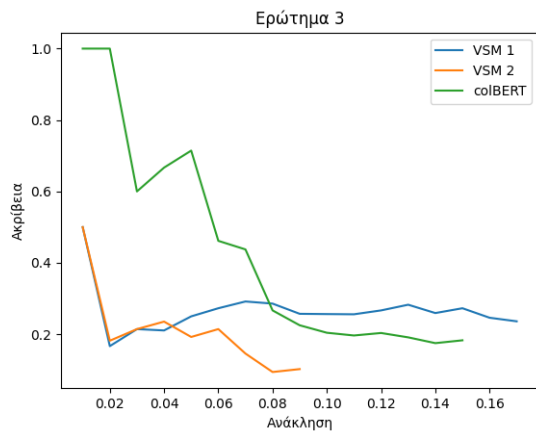
VSM 1	0.016352
VSM 2	0.008598
colBERT	0.007089

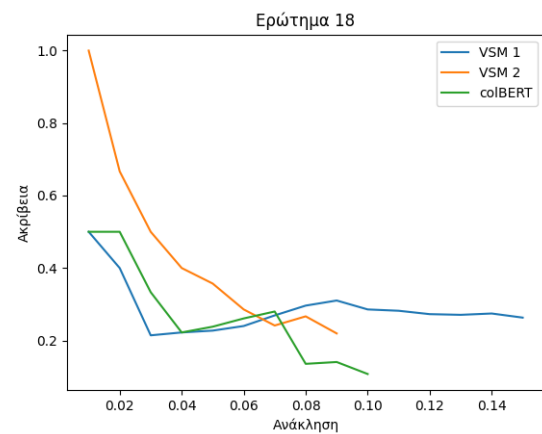
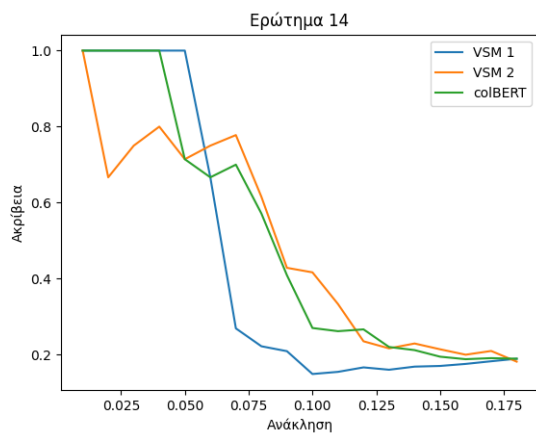
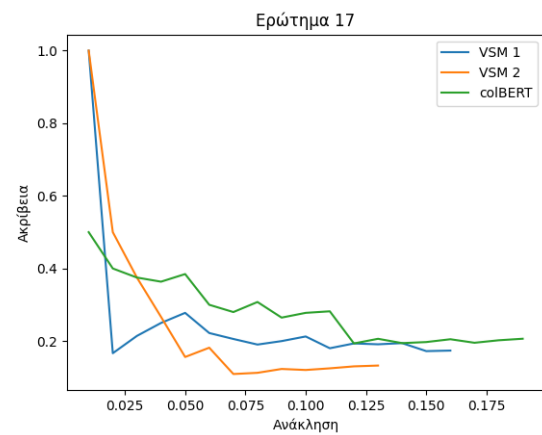
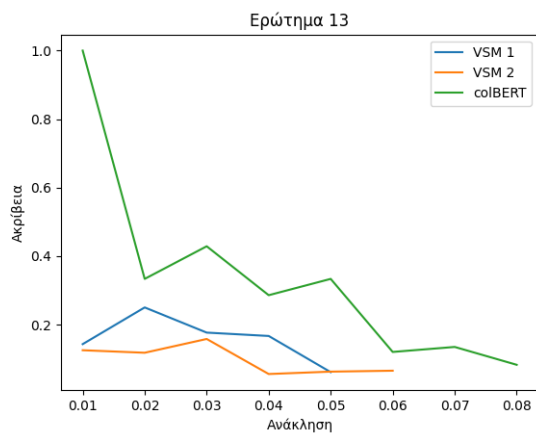
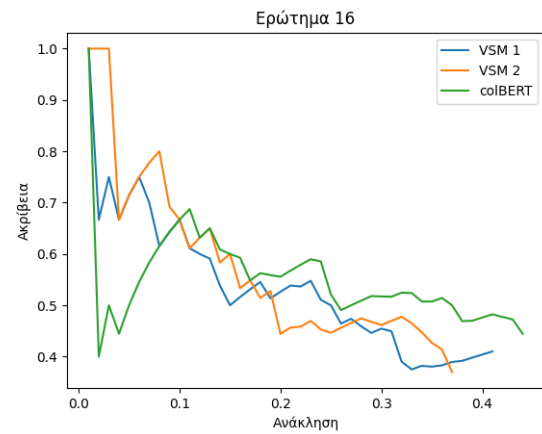
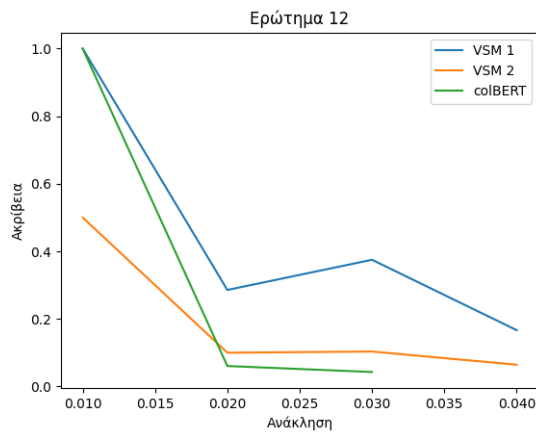
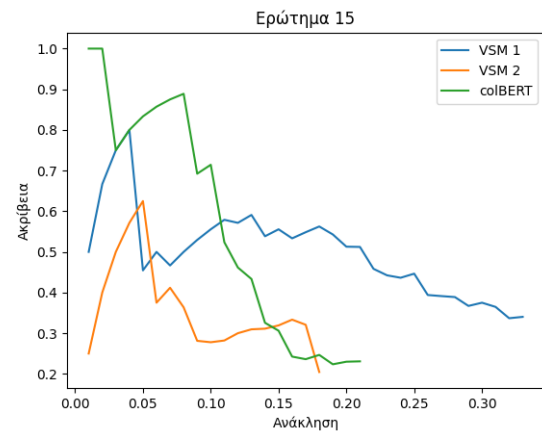
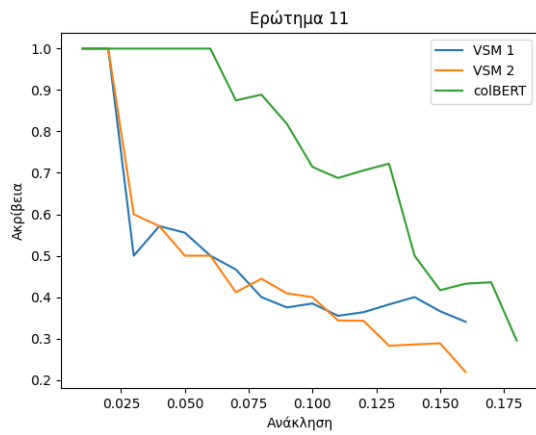
Στα διαγράμματα ανάκλησης-ακρίβειας που παρουσιάζονται παρακάτω, παρατηρούμε πως σχεδόν πάντα το **VSM #2** μοντέλο τερματίζει τελευταίο. Αυτό είναι κάτι που συνάδει με τα αποτελέσματα των Salton-Buckley, αλλιώς είναι και λογικό αν λάβουμε υπόψη τον τρόπο υπολογισμού των βαρών του, στα οποία δεν εμπεριέχεται ο παράγοντας IDF στα βάρη των εγγράφων.

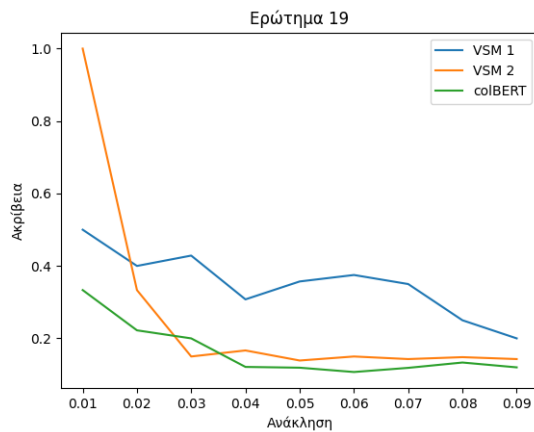
Από την άλλη, ξεκάθαρη υπεροχή παρουσιάζει το **VSM #1** μοντέλο και σε ακρίβεια αλλιώς και σε ανάκληση, κάτι που είναι εξίσου σημαντικό όταν έχουμε να κάνουμε με επιστημονικές/ιατρικές αναζητήσεις. Η υπεροχή του **VSM #1** αποτυπώνεται και από τη MAP τιμή του που είναι διπλάσια από τα υπόλοιπα μοντέλα, δείχνοντας μια σταθερότητα εν' όψει όλων των ερωτημάτων στο οποίο εξετάστηκε.

Θεωρώ πως το **colBERT** θα έπρεπε να τα είχε πάει καλύτερα, δεδομένης της σημασιολογικής του ικανότητας, αν και σε μεμονωμένες περιπτώσεις φαίνεται να υπερτερεί, ειδικά ως προς την ακρίβεια. Ίσως η συλλογή των κειμένων που χρησιμοποιήθηκε να μην ήταν η ιδανική για το συγκεκριμένο μοντέλο, ή ίσως το μοντέλο να επιδέχεται διαφορετικές παραμετροποιήσεις που θα μπορούσαν να βελτιώσουν την απόδοσή του.









4 ΠΑΡΑΡΤΗΜΑ

4.0.1 tools.py

```

1 import json
2 import os
3
4 def get_docs():
5     docs_directory = 'Collection/docs'
6     filename_list = [file for file in os.listdir(docs_directory)]
7
8     doc_tuples_list = []
9
10    for filename in filename_list:
11        with open(os.path.join(docs_directory, filename), 'r') as doc_file:
12            doc = doc_file.readlines()
13            doc = [word.strip() for word in doc]
14            doc_tuple = (filename, doc) # (DocID, <doc>)
15            doc_tuples_list.append(doc_tuple)
16
17    return doc_tuples_list
18
19 def get_queries():
20     filename = 'Collection/Queries_20'
21
22     with open(filename, 'r') as queries_file:
23         queries = queries_file.readlines()
24
25     return queries
26
27
28 def get_json_file(json_file):
29     with open(json_file, 'r') as file:
30         json_data = json.load(file)
31
32     return json_data
33
34
35 def get_relevant():
36     relevant_docs = {}
37     with open('Collection/Relevant_20', 'r') as file:
38         for i, line in enumerate(file):
39             relevant_docs[i] = [int(item) for item in line.split()]
40
41     return relevant_docs

```

4.0.2 preprocessing.py

```

1 from nltk.corpus import stopwords
2 from nltk.stem.porter import PorterStemmer
3 import tools
4
5 def preprocess_collection():
6     doc_tuples_list = tools.get_docs()
7
8     stop_words = set(stopwords.words('english'))
9     stemmer = PorterStemmer()
10
11     stripped_docs = []
12
13     for doc_tuple in doc_tuples_list:
14         stripped_doc = [stemmer.stem(word.lower()) for word in doc_tuple[1] if word.lower() not in
15             stop_words]
16         stripped_doc_tuple = (doc_tuple[0], stripped_doc) # (DocID, <stripped_doc>)
17         stripped_docs.append(stripped_doc_tuple)
18
19     return stripped_docs
20
21 def preprocess_queries():
22     queries = tools.get_queries()
23
24     stop_words = set(stopwords.words('english'))
25     stemmer = PorterStemmer()
26
27     stripped_queries = []
28
29     for query in queries:
30         stripped_query = [stemmer.stem(word.lower()) for word in query.split() if word.lower() not
31             in stop_words]
32         stripped_queries.append(stripped_query)
33
34     stripped_queries_new = []
35     for query in stripped_queries:
36         query = list(set(query))
37         stripped_queries_new.append(query)
38
39     return stripped_queries_new

```

4.0.3 inverted_index.py

```

1 def create_inverted_index(stripped_docs_tuples):
2     inverted_index = {}
3
4     for doc_tuple in stripped_docs_tuples:
5         for term in doc_tuple[1]:
6             if term not in inverted_index:
7                 inverted_index[term] = set()
8             inverted_index[term].add((doc_tuple[0], doc_tuple[1].count(term)))
9
10    return inverted_index

```

4.0.4 vector_space_model.py

```

1 import math

```

```

2 import numpy as np
3 from numpy.linalg import norm
4
5 def idf1(N, n):
6     return math.log(N / n)
7
8 def idf2(N, n):
9     return math.log(N - n / n)
10
11 def get_value(item):
12     return item[1]
13
14
15 def vsm(doc_collection, query, inverted_index, model_type):
16     query_tfidf = {}
17     doc_tfidf = {}
18     for doc in doc_collection:
19         doc_tfidf[doc[0]] = [] # Αρχικοποίηση λεξικού tf-idf τιμών με άδειες λίστες για κάθε
20                                 # έγγραφο
21
22     # Υπολογισμός TF όρων για το ερώτημα:
23     query_tfs = {}
24     for term in query:
25         query_tfs[term] = query.count(term)
26
27     # Υπολογισμός TF-IDF τιμών. Διατρέχουμε κάθε όρο από το ερώτημα...
28     for term in query:
29         docs_containing_term = set()
30
31         # Αν ο όρος του ερωτήματος υπάρχει στο ανεστραμμένο ευρετήριο, υπολογίζουμε την IDF τιμή
32         # του.
33         if term in inverted_index:
34             if model_type == "1":
35                 idf = idf1(len(doc_collection), len(inverted_index[term]))
36             elif model_type == "2":
37                 idf = idf2(len(doc_collection), len(inverted_index[term]))
38
39             # και την TF-IDF τιμή, δηλαδή το ΒΑΡΟΣ ΟΡΟΥ ΕΡΩΤΗΜΑΤΟΣ
40             if model_type == "1":
41                 query_tfidf[term] = (0.5 + 0.5 * (query_tfs[term] / max(query_tfs.values())) * idf
42             elif model_type == "2":
43                 query_tfidf[term] = idf
44
45     # Υπολογισμός TF-IDF τιμών για τα έγγραφα:
46     for item in inverted_index[term]:
47         docs_containing_term.add(item[0])
48         # το docs_containing_term περιέχει τα έγγραφα που περιέχουν τον συγκεκριμένο όρο.
49
50     # διατρέχουμε τα έγγραφα του docs_containing_term...
51     for doc in doc_collection:
52         if doc[0] in docs_containing_term:
53             # έγγραφο βρέθηκε - υπολογίζουμε TF
54             for item in inverted_index[term]:
55                 if item[0] == doc[0]:
56                     doc_tf = item[1]
57                     break
58             # υπολογισμός TFIDF
59             if model_type == "1":
60                 doc_tfidf[doc[0]].append(doc_tf * idf1(len(doc_collection), len(
61                     inverted_index[term])))
62             elif model_type == "2":

```

```

60         doc_tfidf[doc[0]].append(0.5 + 0.5 * (query_tfs[term] / max(query_tfs.
        values()))))
61     else:
62         # δε βρέθηκε έγγραφο, 0 στο διάνυσμα
63         doc_tfidf[doc[0]].append(0)
64
65     similarity = {}
66     for doc in doc_tfidf:
67         similarity[doc] = np.dot(list(query_tfidf.values()), doc_tfidf[doc]) / (
68             norm(list(query_tfidf.values())) * norm(doc_tfidf[doc]))
69
70     # μετατροπή nan τιμών σε 0
71     similarity = {k: 0 if np.isnan(v) else v for k, v in similarity.items()}
72     sort_similarity = sorted(similarity.items(), key=get_value)
73
74     return sort_similarity[-100:][::-1]
75
76
77 def run_vsm(doc_collection, queries, inverted_index):
78     results1 = []
79     results2 = []
80
81     for query in queries:
82         results1.append(vsm(doc_collection, query, inverted_index, "1"))
83         results2.append(vsm(doc_collection, query, inverted_index, "2"))
84
85     return results1, results2

```

4.05 colBERT_preprocessing.py

```

1 import json
2 import tools
3
4 docs = tools.get_docs()
5 queries = tools.get_queries()
6
7 ColBERTqueries = {}
8 for i, query in enumerate(queries):
9     ColBERTqueries[str(i)] = query.strip().upper()
10
11 colBERTdocs = {}
12 for i, doc in enumerate(docs):
13     colBERTdocs[doc[0]] = doc[1]
14
15 with open("colBERT_input/colBERT_docs", "w") as filename:
16     json.dump(colBERTdocs, filename)
17     filename.close()
18
19 with open("colBERT_input/colBERT_queries", "w") as filename:
20     json.dump(ColBERTqueries, filename)
21     filename.close()

```

4.06 colBERT_.ipynb

Εισαγωγή Βιβλιοθηκών/ColBERT

```

1 !git -C ColBERT/ pull || git clone https://github.com/stanford-futuredata/ColBERT.git
2 import sys; sys.path.insert(0, 'ColBERT/')
3 try: # When on google Colab, let's install all dependencies with pip.

```

```

4     import google.colab
5     !pip install -U pip
6     !pip install -e ColBERT/['faiss-gpu','torch']
7 except Exception:
8     import sys; sys.path.insert(0, 'ColBERT/')
9     try:
10         from colbert import Indexer, Searcher
11     except Exception:
12         print("If you're running outside Colab, please make sure you install ColBERT in conda following
13             the instructions in our README. You can also install (as above) with pip but it may install
14             slower or less stable faiss or torch dependencies. Conda is recommended.")
15     assert False
16 import colbert
17 from colbert import Indexer, Searcher
18 from colbert.infra import Run, RunConfig, ColBERTConfig
19 from colbert.data import Queries, Collection

```

Εισαγωγή αρχείων

```

1 import json
2
3 with open('colBERT_docs', "r") as file:
4     doc_list = json.load(file)
5
6 with open('colBERT_queries', "r") as file:
7     query_list = json.load(file)
8
9 doc_ids = list(doc_list.keys())
10 doc_content = list(doc_list.values())
11
12 query_ids = list(query_list.keys())
13 query_content = list(query_list.values())
14
15 nbits = 2
16 doc_maxlen = 300

```

Δημιουργία Indexer - Searcher

```

1 checkpoint = 'colbert-ir/colbertv2.0'
2
3 with Run().context(RunConfig(nranks=1, experiment='notebook')): # nranks specifies the number of
4     config = ColBERTConfig(doc_maxlen=doc_maxlen, nbits=nbits, kmeans_niters=20) # kmeans_niters
5     # specifies the number of iterations of k-means clustering; 4 is a good and fast default.
6     # Consider larger
7     # numbers for small datasets.
8
9     indexer = Indexer(checkpoint=checkpoint, config=config)
10    indexer.index(name='index_name', collection=doc_content, overwrite=True)
11
12 with Run().context(RunConfig(experiment='notebook')):
13    searcher = Searcher(index='index_name', collection=doc_content)

```

Αποτελέσματα:

```

1 results = []
2 for query in query_content:
3     result = searcher.search(query, k=100)
4
5     passages = []
6     for passage_id, passage_rank, passage_score in zip(*result):
7         passages.append(int(doc_ids[passage_id]))
8     results.append(passages)

```

```

9 print(results)
10
11 from google.colab import files
12 import json
13
14 with open("colBERT_output.json", "w") as file:
15     json.dump(results, file)
16
17 files.download("colBERT_output.json")

```

4.0.7 evaluation_metrics.py

```

1 import tools
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 def recall_precision(results, model):
6     # results: λίστα όλων( των queries) που περιλαμβάνει λίστες με τα πιο σχετικά έγγραφα, όπως
7     υπολογίστηκαν από το cosine similarity
8     relevant_docs = tools.get_relevant()
9     # relevant_docs: dictionary με τις τιμές των πραγματικά σχετικών εγγράφων του Relevant_20
10
11     recall_precision_values = []
12
13     if model == "vsm":
14         # Το results περιλαμβάνει συγκεντρωτικά όλα τα αποτελέσματα για κάθε query. Διατρέχουμε το
15         # κάθε query μεμονομένα...
16         for i, results_query in enumerate(results):
17             precision = []
18             recall = []
19             truly_relevant_docs = 0
20             # διατρέχουμε τα πιο σχετικά έγγραφα για το συγκεκριμένο query...
21             for j, doc in enumerate(results[i]):
22                 if type(doc) == tuple:
23                     current_doc = int(doc[0])
24                     if current_doc in relevant_docs[i]:
25                         truly_relevant_docs += 1
26
27             recall.append(truly_relevant_docs / len(results_query))
28             precision.append(truly_relevant_docs / (j + 1))
29
30             recall_precision_values.append((recall, precision))
31
32     if model == "colBERT":
33         for i, results_query in enumerate(results):
34             precision = []
35             recall = []
36             truly_relevant_docs = 0
37             for j, doc in enumerate(results[i]):
38                 if doc in relevant_docs[i]:
39                     truly_relevant_docs += 1
40
41             recall.append(truly_relevant_docs / len(results_query))
42             precision.append(truly_relevant_docs / (j + 1))
43
44             recall_precision_values.append((recall, precision))
45
46     return recall_precision_values
47
48 def mean_average_precision(recall_precision_values):

```

```

47 average_precision_list = []
48 for i in range(len(recall_precision_values)):
49     recalls, precisions = recall_precision_values[i]
50     average_precision = 0
51
52     for j in range(1, len(recalls)):
53         average_precision += (recalls[j] - recalls[j - 1]) * precisions[j]
54     average_precision_list.append(average_precision)
55
56     return np.mean(average_precision_list)
57
58
59 def run_metrics(vsm1_results, vsm2_results):
60     colBERT_results = tools.get_json_file("colBERT_output.json")
61
62     recall_precision_vsm1 = recall_precision(vsm1_results, "vsm")
63     recall_precision_vsm2 = recall_precision(vsm2_results, "vsm")
64     recall_precision_colBERT = recall_precision(colBERT_results, "colBERT")
65
66     map_vsm1 = mean_average_precision(recall_precision_vsm1)
67     map_vsm2 = mean_average_precision(recall_precision_vsm2)
68     map_colBERT = mean_average_precision(recall_precision_colBERT)
69
70     print("vsm1", map_vsm1)
71     print("vsm2", map_vsm2)
72     print("colBERT", map_colBERT)
73
74     for i in range(len(recall_precision_vsm1)):
75         plt.figure()
76         plt.plot(recall_precision_vsm1[i][0], recall_precision_vsm1[i][1], label="VSM 1")
77         plt.plot(recall_precision_vsm2[i][0], recall_precision_vsm2[i][1], label="VSM 2")
78         plt.plot(recall_precision_colBERT[i][0], recall_precision_colBERT[i][1], label="colBERT")
79         plt.xlabel('Ανάκληση')
80         plt.ylabel('Ακρίβεια')
81         plt.title(f'Ερώτημα {i + 1}')
82         plt.legend(loc='upper right')
83         plt.savefig('Precision_Recall_Curve/' + str(i + 1) + '.png')

```

4.0.8 main.py

```

1 import inverted_index as ii
2 import vector_space_model as vsm
3 import preprocessing as prep
4 import evaluation_metrics as metrics
5
6 def main_app():
7     doc_collection = prep.preprocess_collection() # stripped_docs
8     queries = prep.preprocess_queries() # stripped_queries
9
10    inverted_index = ii.create_inverted_index(doc_collection)
11    results1, results2 = vsm.run_vsm(doc_collection, queries, inverted_index)
12
13    metrics.run_metrics(results1, results2)
14
15
16 if __name__ == '__main__':
17     main_app()

```