

# Προγραμματιστική εργασία: Σημασιολογική ανάκτηση

## Φάση 1

Καραγιαννάκος Άριστείδης AM 3220066

Κοσμίδη Έλλη AM 3220086

Elastic serach: bin\elasticsearch.bat , curl <http://localhost:9200/>, notepad

.\config\elasticsearch.yml , xpck.security.enabled: false

xpack.security.http.ssl.enabled: false

```
[2026-01-10T21:13:51,101][INFO ][o.e.c.h.NettyAllocator ] [DESKTOP-0VV32SM] creating NettyAllocator with the following configs: [name=elasticsearch_comfigured, chunk_size=1mb, suggested_max_allocation_size=1mb, factors=[es.unsafe.use.netty.default_chunk_and_page_size=false, glibc_enabled=true, glibc_region_size=4mb]]  
[2026-01-10T21:13:49,230][INFO ][o.e.d.DiscoveryModule ] [DESKTOP-0VV32SM] using discovery type [multi-node] and seed hosts providers [settings]  
[2026-01-10T21:13:50,358][INFO ][o.e.n.Node ] [DESKTOP-0VV32SM] initialized  
[2026-01-10T21:13:50,359][INFO ][o.e.n.Node ] [DESKTOP-0VV32SM] starting ...  
[2026-01-10T21:13:51,227][INFO ][o.e.x.c.PersistentCache] [DESKTOP-0VV32SM] persistent cache index loaded  
[2026-01-10T21:13:51,228][INFO ][o.e.x.d.DeprecationIndexingComponent] [DESKTOP-0VV32SM] deprecation component started  
[2026-01-10T21:13:51,292][INFO ][o.e.t.TransportService ] [DESKTOP-0VV32SM] publish_address {127.0.0.1:9300}, bound_addresses {[::1]:9300}, {127.0.0.1:9300}  
[2026-01-10T21:13:51,661][WARN ][o.e.c.ClusterBootstrapService] [DESKTOP-0VV32SM] this node is locked into cluster UUID [RY0_dGJQiuCsAGxb19mcA] but [cluster.initial_master_nodes] is set to [DESKTOP-0VV32SM]; remove this setting to avoid possible data loss caused by subsequent cluster bootstrap attempts; for further information see https://www.elastic.co/docs/deploy-manage/deploy-self-managed/importtant-settings-configuration?version=9.%24initial_master_nodes  
[2026-01-10T21:13:51,757][INFO ][o.e.s.MasterService ] [DESKTOP-0VV32SM] elected-as-master ([1] nodes joined in term 8) [FINISH_ELECTION, {DESKTOP-0VV32SM}{3yZyPEoSSP0JoelXAochBQ}{e8XFkuJAQlCeaNz64Q2e1A}{DESKTOP-0VV32SM}{127.0.0.1}{cdfh1mrstw}{9.2.1}{8000099-9039001} completing election], term: 8, version: 200, delta: master node changed [previous [], current [{DESKTOP-0VV32SM}{3yZyPEoSSP0JoelXAochBQ}{e8XFkuJAQlCeaNz64Q2e1A}{DESKTOP-0VV32SM}{127.0.0.1}{9300}{cdfh1mrstw}{9.2.1}{8000099-9039001}]]  
[2026-01-10T21:13:51,871][INFO ][o.e.c.s.ClusterApplierService] [DESKTOP-0VV32SM] master node changed [previous [], current [{DESKTOP-0VV32SM}{3yZyPEoSSP0JoelXAochBQ}{e8XFkuJAQlCeaNz64Q2e1A}{DESKTOP-0VV32SM}{127.0.0.1}{cdfh1mrstw}{9.2.1}{8000099-9039001}]], term: 8, version: 200, reason: Publication{term=8, version=200}  
[2026-01-10T21:13:51,912][INFO ][o.e.c.NodeJoinExecutor] [DESKTOP-0VV32SM] node-join: [{DESKTOP-0VV32SM}{3yZyPEoSSP0JoelXAochBQ}{e8XFkuJAQlCeaNz64Q2e1A}{DESKTOP-0VV32SM}{127.0.0.1}{127.0.0.1:9300}{cdfh1mrstw}{9.2.1}{8000099-9039001}] with reason [completing election]  
[2026-01-10T21:13:51,915][INFO ][o.e.h.AbstractHttpServerTransport] [DESKTOP-0VV32SM] publish_address {10.5.0.2:9200}, bound_addresses {[::]:9200}  
[2026-01-10T21:13:51,925][INFO ][o.e.w.LicensedWriteLoadForecaster] [DESKTOP-0VV32SM] license state changed, now [valid]  
[2026-01-10T21:13:51,926][INFO ][o.e.n.Node ] [DESKTOP-0VV32SM] started [DESKTOP-0VV32SM]{3yZyPEoSSP0JoelXAochBQ}{e8XFkuJAQlCeaNz64Q2e1A}{DESKTOP-0VV32SM}{127.0.0.1}{127.0.0.1:9300}{cdfh1mrstw}{9.2.1}{8000099-9039001}[ml.allocated_processors=4, ml.machine_memory=17111879680, transform.config_version=10.0.0, xpack.installed=true, ml.config_version=12.0.0, ml.max_jm_size=8556380160, ml.allocated_processors.double=4.0}  
[2026-01-10T21:13:51,958][WARN ][o.e.x.i.s.e.a.ElasticInferenceServiceAuthorizationHandler] [DESKTOP-0VV32SM] Failed to revoke access to default inference endpoint IDs: [elser_mode]_2, jina-embeddings-v3, rainbow-sprinkles, elastic-rerank-v1, error: org.elasticsearch.cluster.block.ClusterBlockException: blocked by: [SERVICE_UNAVAILABLE/1/state not recovered / initialized];  
[2026-01-10T21:13:52,007][INFO ][o.e.m.MIndexRollover ] [DESKTOP-0VV32SM] ML legacy indices rolled over  
[2026-01-10T21:13:52,210][INFO ][o.e.x.m.AnomaliesIndexUpdate] [DESKTOP-0VV32SM] legacy ml anomalies indices rolled over and aliases updated  
[2026-01-10T21:13:52,330][INFO ][o.e.l.ClusterStateLicenseService] [DESKTOP-0VV32SM] license [e3b442b2-/cda-4d24-81fb-67395470fd46] mode [basic] - valid  
[2026-01-10T21:13:52,333][INFO ][o.e.c.f.AbstractFileWatchingService] [DESKTOP-0VV32SM] starting file watcher ...  
[2026-01-10T21:13:52,338][INFO ][o.e.c.f.AbstractFileWatchingService] [DESKTOP-0VV32SM] file settings service up and running [tid:72]  
[2026-01-10T21:13:52,340][INFO ][o.e.r.s.FileSettingsService] [DESKTOP-0VV32SM] setting file [C:\elasticsearch-9.2.1\config\operator\settings.json] not found, initializing [file_settings] as empty  
[2026-01-10T21:13:52,345][INFO ][o.e.g.GatewayService ] [DESKTOP-0VV32SM] recovered [7] indices into cluster_state  
[2026-01-10T21:13:52,357][INFO ][o.e.w.LicensedWriteLoadForecaster] [DESKTOP-0VV32SM] license state changed, now [not valid]  
[2026-01-10T21:13:53,271][INFO ][o.e.h.n.HealthNodeTaskExecutor] [DESKTOP-0VV32SM] Node [{DESKTOP-0VV32SM}{3yZyPEoSSP0JoelXAochBQ}] is selected as the current health node.  
[2026-01-10T21:13:53,274][INFO ][o.e.c.a.a.AllocationService] [DESKTOP-0VV32SM] current.health="YELLOW" message="Cluster health status changed from [RED] to [YELLOW] (reason: [shards started [[documents][0], [.ds-lm-history/-2025.11.27-00001][0]]]. previous.health="RED" reason="shards started [[documents][0], [.ds-lm-history/-2025.11.27-00001][0]]]"
```

## A1 και A2

Αρχικά δημιουργήσαμε την συνάρτηση **process\_text** που δέχεται παράμετρο **text** και το επεξεργάζεται, δηλαδή αφαιρεί τα πολλά κενά και σύμβολα αλλαγής γραμμής, tabs κλπ (\r,\n,\t) και τα αντικαθιστά με ένα κενό.

Υστερα για το index της Elasticsearch:

- Δημιουργήσαμε BM25\_optimized\_mapping που ορίζει τον τρόπο που θα επεξεργάζεται και θα αντιστοιχεί το index τα δεδομένα του. Συγκεκριμένα:

Analyzer υπεύθυνος για να σπάει το text σε tokens, να αφαιρεί τα stopwords σύμφωνα με τη λίστα “\_english\_”.

Similarity καθορίσαμε τον BM25 αλγόριθμο βαθμολόγησης που θα χρησιμοποιήσει η Elasticsearch (αλγόριθμος βαθμολόγησης δίνει υψηλότερη βαθμολογία σε έγγραφα που είναι πιο σχετικά με το query και χαμηλότερη σε αυτά που είναι λιγότερο σχετικά).

Mapping στο ευρετήριο ορίσαμε σαν keyword το id και text (περιεχόμενο ευρετηρίου) το text που θα τροποποιηθεί με τον "english\_analyzer" που ορίσαμε πριν.

- Δημιουργήσαμε την συνάρτηση **generate\_actions**, η οποία ορίζει τον τρόπο που θα εισαχθούν τα κείμενα μέσα στο index όταν θα χρησιμοποιήσουμε την συνάρτηση bulk. Δηλαδή η generate\_actions δέχεται το documents.csv, διαβάζει μια μια τις γραμμές και ουσιαστικά για κάθε γραμμή/ κείμενο φτιάχνει ένα dict που ορίζει σε ποιο index θα μπει, τι id θα έχει και τι περιεχόμενο, με δομή κατάλληλη για την συνάρτηση bulk (κλειδιά \_index, \_id, \_source και οι αντίστοιχες τιμές τους).
- \* το \_source περιέχει τα πραγματικά δεδομένα του κάθε κειμένου σε μορφή dictionary. Για το text καλούμε την συνάρτηση process\_text για προεπεξεργασία.
- Χρησιμοποιήσαμε την έτοιμη συνάρτηση bulk για το γέμισμα του ευρετηρίου με παραμέτρους την σύνδεση με την Elasticsearch και την generate\_actions

```
[INFO] DATA_DIR: C:\Users\ArisK\Desktop\IR20252026
[INFO] OUT_DIR : C:\Users\ArisK\Desktop\dawdawdawd
```

```
[DOCS] Sample after preprocess:
```

	doc_id	text
0	193157	Support towards the Europe PMC initiative-Cont...
1	193158	Support to the Vice-Presidents of the ERC Scie...
2	193159	Implementation of activities described in the ...

```
[DOCS] Count: 18316
```

```
[QUERIES] Sample after preprocess:
```

	qid	query
0	Q01	EUTRAVEL Optimodal European Travel Ecosystem E...
1	Q02	Track And Know Big Data for Mobility Tracking ...
2	Q03	SELIS, Towards a Shared European Logistics Int...
3	Q04	TYPHON Polyglot and Hybrid Persistence Archite...
4	Q05	CHARIOT Cognitive Heterogeneous Architecture f...

```
[QUERIES] Count: 10
```

## A2: Δημιουργία ευρετηρίου ElasticSearch (BM25) + Bulk indexing

To screenshot δεν φαίνοταν καλά αυτό είναι το print φαίνεται και στον κώδικα.

```
[OK] Elasticsearch: elasticsearch | 9.2.1
```

```
[i] Deleted existing index: ir2025_phase1_bm25
```

```
[OK] Created index 'ir2025_phase1_bm25' with BM25(k1=1.2, b=0.75)
```

```
[BULK] Indexing documents...
```

C:\Users\ArisK\AppData\Local\Temp\ipykernel\_3308\1684811186.py:61: DeprecationWarning: Passing transport options in the API method is deprecated. Use 'Elasticsearch.options()' instead.

```
    helpers.bulk(es, doc_actions(df_docs, INDEX_NAME), chunk_size=1000,  
    request_timeout=120)
```

[OK] Indexed docs in ES: 18316

### A3

Για το A3 φορτώσαμε τα ερωτήματα από το αρχείο queries.csv χρησιμοποιώντας τη βιβλιοθήκη pandas. Το αρχείο περιέχει για κάθε ερώτημα ένα μοναδικό αναγνωριστικό (ID) και το αντίστοιχο κείμενο (Text). Οι στήλες μετονομάστηκαν σε qid και query αντίστοιχα και εφαρμόστηκε βασική προεπεξεργασία στο κείμενο των ερωτημάτων με τη συνάρτηση clean\_text, ώστε να είναι συμβατή με την επεξεργασία των εγγράφων.

Στη συνέχεια, για κάθε ερώτημα εκτελέσαμε αναζήτηση στο ευρετήριο της Elasticsearch, συγκρίνοντας το κείμενο του ερωτήματος με τα κείμενα των εγγράφων του index. Η κατάταξη των εγγράφων πραγματοποιήθηκε με βάση τη βαθμολογία σχετικότητας που υπολογίζεται από τον αλγόριθμο BM25, ο οποίος ορίστηκε στο στάδιο A2.

Για κάθε ερώτημα συλλέξαμε τα top-k αποτελέσματα και δημιουργήσαμε τρία αρχεία run για k = 20, 30 και 50 (run\_bm25\_k20.txt, run\_bm25\_k30.txt, run\_bm25\_k50.txt).

To scresshot δεν φαίνοταν καλά αυτό είναι το print φαίνεται και στον κώδικα.

Queries k=20: 100% [██████████] 10/10 [00:00<00:00, 13.06it/s]

[OK] Wrote run\_bm25\_k20.txt (200 lines)

Queries k=30: 100% [██████████] 10/10 [00:00<00:00, 34.77it/s]

[OK] Wrote run\_bm25\_k30.txt (300 lines)

Queries k=50: 100% [██████████] 10/10 [00:00<00:00, 37.98it/s]

[OK] Wrote run\_bm25\_k50.txt (500 lines)

### A4

Για την αξιολόγηση χρησιμοποιήσαμε το εργαλείο trec\_eval, συγκρίνοντας τα run files με το ground truth αρχείο qrels.txt. Τρέξαμε το trec\_eval για τα μέτρα:

MAP (Mean Average Precision), που συνοψίζει την ποιότητα κατάταξης σε όλο το ranking,

Precision @k για k = 5, 10, 15, 20, που μετρά πόσα από τα πρώτα k αποτελέσματα είναι συναφή.

Για κάθε run file (k=20,30,50) εκτελέσαμε:

-m map για MAP

-m P.5, -m P.10, -m P.15, -m P.20 για precision στα αντίστοιχα cutoff.

Στη συνέχεια κάναμε parsing της εξόδου σε pandas και δημιουργήσαμε:

πίνακες ανά query (per-query results),

καθώς και συγκεντρωτικό πίνακα “all” που περιέχει τα συνολικά MAP και P@k για κάθε run.

Τέλος, αποθηκεύσαμε τα αποτελέσματα σε CSV ώστε να τα ενσωματώσουμε στην αναφορά και να κάνουμε σύγκριση μεταξύ των διαφορετικών τιμών k.

Συμπέρασμα:

Από τα αποτελέσματα παρατηρούμε ότι το MAP αυξάνεται όσο μεγαλώνει το k, γεγονός που δείχνει ότι το BM25 κατατάσσει σωστά περισσότερα συναφή έγγραφα όταν εξετάζουμε μεγαλύτερο μέρος της λίστας αποτελεσμάτων. Αντίθετα, οι τιμές Precision@5, Precision@10, Precision@15 και Precision@20 παραμένουν ίδιες για όλα τα k, κάτι που σημαίνει ότι τα πιο σχετικά έγγραφα εμφανίζονται σταθερά στις πρώτες θέσεις.

Συνολικά, το BM25 παρουσιάζει σταθερή και καλή απόδοση ως baseline μοντέλο, αποτελώντας μια αξιόπιστη βάση για σύγκριση με πιο προχωρημένες μεθόδους ανάκτησης στις επόμενες φάσεις της εργασίας.

Εκτελείται αξιολόγηση για: run\_bm25\_k20.txt

Metric	map	P@5	P@10	P@15	P@20
Query					
Q01	0.62280	0.80000	0.80000	0.73330	0.65000
Q02	0.32610	0.60000	0.30000	0.26670	0.30000
Q03	0.70970	0.80000	0.70000	0.66670	0.60000
Q04	0.29750	0.60000	0.40000	0.26670	0.30000
Q05	0.70660	1.00000	0.80000	0.73330	0.65000
Q06	0.69260	1.00000	1.00000	0.80000	0.70000
Q07	0.64130	1.00000	0.80000	0.73330	0.60000
Q08	0.46940	0.80000	0.70000	0.46670	0.45000
Q09	0.67690	1.00000	0.90000	0.86670	0.75000
Q10	0.31980	0.60000	0.30000	0.26670	0.25000
all	0.54630	0.82000	0.67000	0.58000	0.52500

■ Saved: [C:\Users\ArisK\Desktop\dawdawdawd\run\\_bm25\\_k20\\_per\\_query.csv](C:\Users\ArisK\Desktop\dawdawdawd\run_bm25_k20_per_query.csv)

Εκτελείται αξιολόγηση για: run\_bm25\_k30.txt

Metric	map	P@5	P@10	P@15	P@20
Query					
Q01	0.69690	0.80000	0.80000	0.73330	0.65000
Q02	0.32610	0.60000	0.30000	0.26670	0.30000
Q03	0.74070	0.80000	0.70000	0.66670	0.60000
Q04	0.31840	0.60000	0.40000	0.26670	0.30000
Q05	0.73780	1.00000	0.80000	0.73330	0.65000
Q06	0.76510	1.00000	1.00000	0.80000	0.70000
Q07	0.71360	1.00000	0.80000	0.73330	0.60000
Q08	0.57030	0.80000	0.70000	0.46670	0.45000
Q09	0.74260	1.00000	0.90000	0.86670	0.75000
Q10	0.39930	0.60000	0.30000	0.26670	0.25000
all	0.60110	0.82000	0.67000	0.58000	0.52500

■ Saved: [C:\Users\ArisK\Desktop\dawdawdawd\run\\_bm25\\_k30\\_per\\_query.csv](C:\Users\ArisK\Desktop\dawdawdawd\run_bm25_k30_per_query.csv)

Εκτελείται αξιολόγηση για: run\_bm25\_k50.txt

Εκτελείται αξιολόγηση για: run\_bm25\_k50.txt

Metric	map	P@5	P@10	P@15	P@20
<b>Query</b>					
Q01	0.72910	0.80000	0.80000	0.73330	0.65000
Q02	0.38270	0.60000	0.30000	0.26670	0.30000
Q03	0.74070	0.80000	0.70000	0.66670	0.60000
Q04	0.34880	0.60000	0.40000	0.26670	0.30000
Q05	0.79040	1.00000	0.80000	0.73330	0.65000
Q06	0.78490	1.00000	1.00000	0.80000	0.70000
Q07	0.74290	1.00000	0.80000	0.73330	0.60000
Q08	0.59000	0.80000	0.70000	0.46670	0.45000
Q09	0.84870	1.00000	0.90000	0.86670	0.75000
Q10	0.42660	0.60000	0.30000	0.26670	0.25000
all	0.63850	0.82000	0.67000	0.58000	0.52500

■ Saved: [C:\Users\ArisK\Desktop\dawdawdawd\run\\_bm25\\_k50\\_per\\_query.csv](C:\Users\ArisK\Desktop\dawdawdawd\run_bm25_k50_per_query.csv)

==== ΣΥΓΚΕΝΤΡΩΤΙΚΑ (all) για όλα τα runs ===

	map	P@5	P@10	P@15	P@20
<b>run</b>					
run_bm25_k20.txt	0.54630	0.82000	0.67000	0.58000	0.52500
run_bm25_k30.txt	0.60110	0.82000	0.67000	0.58000	0.52500
run_bm25_k50.txt	0.63850	0.82000	0.67000	0.58000	0.52500

■ Saved: [C:\Users\ArisK\Desktop\dawdawdawd\phase1\\_trec\\_eval\\_summary\\_all.csv](C:\Users\ArisK\Desktop\dawdawdawd\phase1_trec_eval_summary_all.csv)

