

Πανεπιστήμιο Κρήτης
Τμήμα Επιστήμης Υπολογιστών
ΗΥ463 Συστήματα Ανάκτησης Πληροφοριών
Εξάμηνο: Άνοιξη 2019

Γραπτή Αναφορά Έργου

Στοιχεία Φοιτητών

Μέλος	1 ^ο
Ονοματwpώνυμο	Κατσιφαράκης Εμμανουήλ
ΑΜ	3195
Email	csd3195@csd.uoc.gr

Μέλος	2 ^ο
Ονοματwpώνυμο	Κανάρια Ιωάννα
ΑΜ	3609
Email	csd3609@csd.uoc.gr

Πίνακας Περιεχομένων

1	Εισαγωγή	3
2	Τύποι και μέτρα που χρησιμοποιήθηκαν	3
2.1	Τύποι για την δημιουργία του ανεστραμμένου ευρετηρίου – Φάση Α	
2.1.1	Κανονικοποίηση	
2.1.2	Βάρος όρου i στο έγγραφο j	
2.1.3	Νόρμα εγγράφου	
2.2	Μέτρα για την αξιολόγηση του συστήματος – Φάση Β	
2.2.1	bPref (binary preference-based measure)	
2.2.2	AveP (Average Precision)	
2.2.3	nDCG(Normalized Discounted Cumulative Gain)	
3	Μετρήσεις	5
4	Απαιτούμενα αρχεία για την λειτουργία του συστήματος	8
5	Επίλογος	8
6	Αναφορές	8

- Εισαγωγή

Στο project του μαθήματος, ασχοληθήκαμε με την δημιουργία και την αξιολόγηση, ενός δικού μας Συστήματος Ανάκτησης Πληροφοριών. Αναλυτικότερα, στην Α' Φάση της εργασίας, κατασκευάσαμε ένα ανεστραμμένο ευρετήριο, πάνω στο οποίο κάναμε επερωτήσεις σε μια πραγματική συλλογή βιοϊατρικών άρθρων. Επιπροσθέτως, στην Β' Φάση της εργασίας, αξιολογήσαμε την αποτελεσματικότητα του συστήματος μας, βασιζόμενοι σε 30 ιατρικά θέματα και καταγράφοντας μετρήσεις, βάσει των μέτρων bpref (μέτρηση με βάση το άθροισμα του αριθμού των σχετικών εγγράφων που κατατάσσονται πριν από μη συναφή έγγραφα), AveP (υπολογισμός της τιμής της ακρίβειας σε κάθε δυνατό recall level) και nDCG (κανονικοποιημένο μειωμένο σωρευτικό κέρδος).

Ολοκληρώσαμε επιτυχώς όλα τα ερωτήματα της εκφώνησης για την Β' Φάση, ενώ διορθώσαμε την διαδικασία παραγωγής της νόρμας κάθε εγγράφου, την οποία δεν είχαμε επιτύχει στην Α' Φάση.

Τα αδύνατο σημείο της εργασίας μας, είναι ότι η διαδικασία της ευρετηρίασης δεν είναι ιδιαίτερα γρήγορη όπως επίσης ότι κατά την αξιολόγηση του συστήματος με βάση τα μέτρα bpref, AveP, nDCG, είχαμε πολλά μηδενικά αποτελέσματα.

Όσον αφορά τις χρονικές επιδόσεις του συστήματος, η διαδικασία επερωτήσεων διήρκεσε 25 λεπτά και οι μετρήσεις έγιναν σε σύστημα:

CPU: AMD FX-8320 3.5Ghz

Ram: 8GB DDR3-1600MHz

SSD: 120Gb Samsung 840 Evo

Η εργασία έγινε συνεργατικά από τα μέλη της ομάδας, οπότε και τα δύο μέλη ασχολήθηκαν με όλα τα κομμάτια της εργασίας.

- Τύποι και μέτρα που χρησιμοποιήθηκαν

1. Τύποι για την δημιουργία του ανεστραμμένου ευρετηρίου – Φάση Α

Κανονικοποίηση

$tf_{ij} = freq_{ij} / \max_k \{freq_{kj}\}$,

όπου $\max_k \{freq_{kj}\}$ το μεγαλύτερο πλήθος εμφανίσεων ενός όρου στο έγγραφο j.

Βάρος όρου i στο έγγραφο j

$w_{ij} = tf_{ij} \cdot df_i$,

όπου df_i είναι το πλήθος εγγράφων που περιέχουν τον όρο i.

Νόρμα εγγράφου

$|W_{ij}| = \sqrt{(\sum_{1 \leq k \leq N} (w_{ij}^k)^2)}$

2. Μέτρα για την αξιολόγηση του συστήματος – Φάση Β

bPref (binary preference-based measure) [1]

$$bpref = \frac{1}{R} \sum_r 1 - \frac{|n \text{ ranked higher than } r|}{R}$$

όπου R είναι ο αριθμός των συναφών εγγράφων και r είναι το πρώτο συναφές έγγραφο που εντοπίζουμε μετά από η μη συναφή έγγραφα.

AveP (Average Precision)

$$AP = \frac{1}{\Sigma} \sum_{i=1}^{|E|} rel(i)P @ i$$

$$P@k = \frac{|E_k \cap \Sigma|}{k}$$

όπου E είναι ο αριθμός των ευρεθέντων εγγράφων, E_k είναι το ευρεθέν έγγραφο στην θέση k, Σ είναι ο αριθμός των συναφών εγγράφων και rel(i) είναι μία συνάρτηση η οποία μας επιστρέφει 1 αν το έγγραφο i είναι συναφές, διαφορετικά 0.

nDCG(Normalized Discounted Cumulative Gain)

$$\begin{array}{ll} \text{Για } i < 2 & \text{Για } i > 2 \\ DCG = \sum_{i=1}^E (rel(i)) & DCG = \sum_{i=1}^E \left(\frac{rel(i)}{\log_2(i)} \right) \end{array}$$

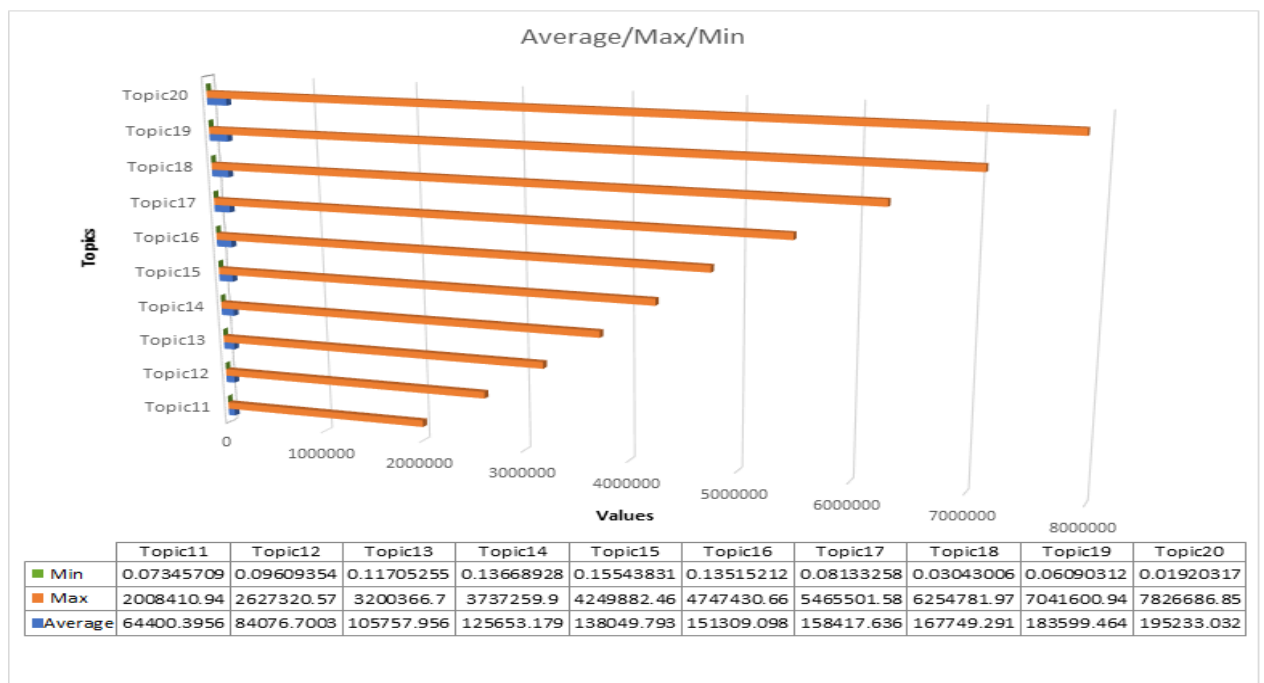
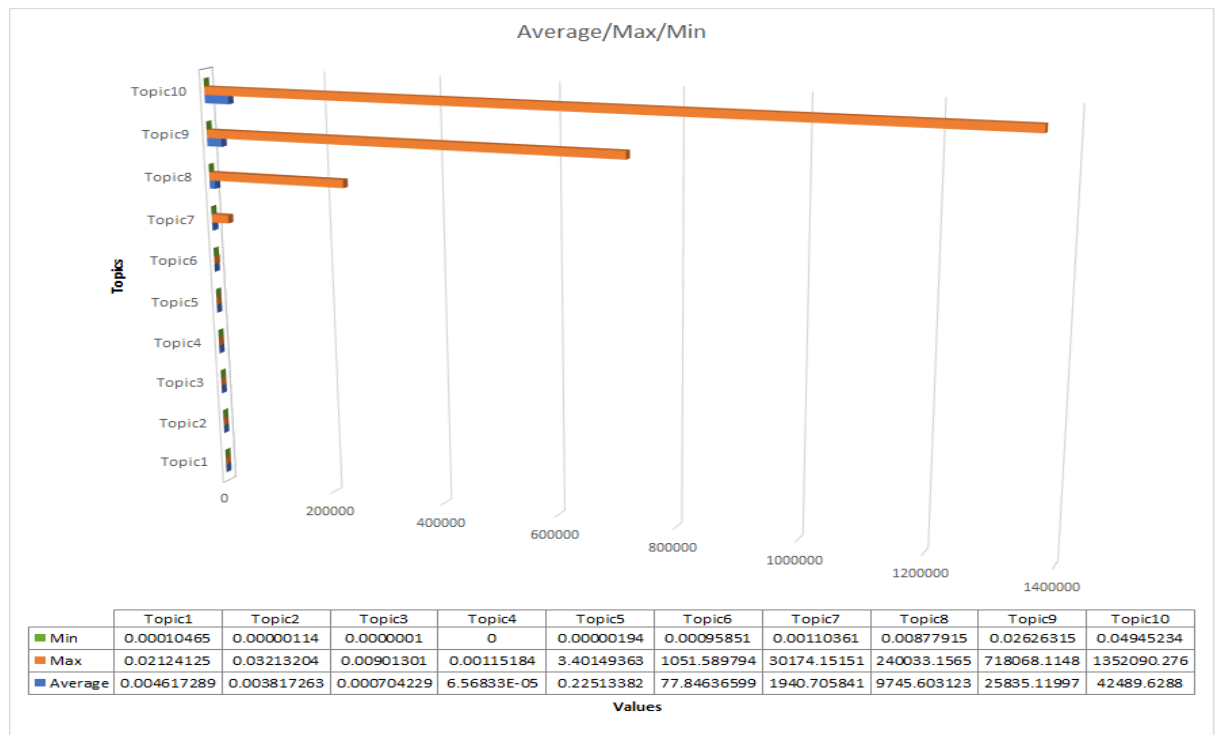
$$\begin{array}{ll} \text{Για } i < 2 & \text{Για } i > 2 \\ IDCG = \sum_{i=1}^{\Sigma} (rel(i)) & IDCG = \sum_{i=1}^{\Sigma} \left(\frac{rel(i)}{\log_2(i)} \right) \end{array}$$

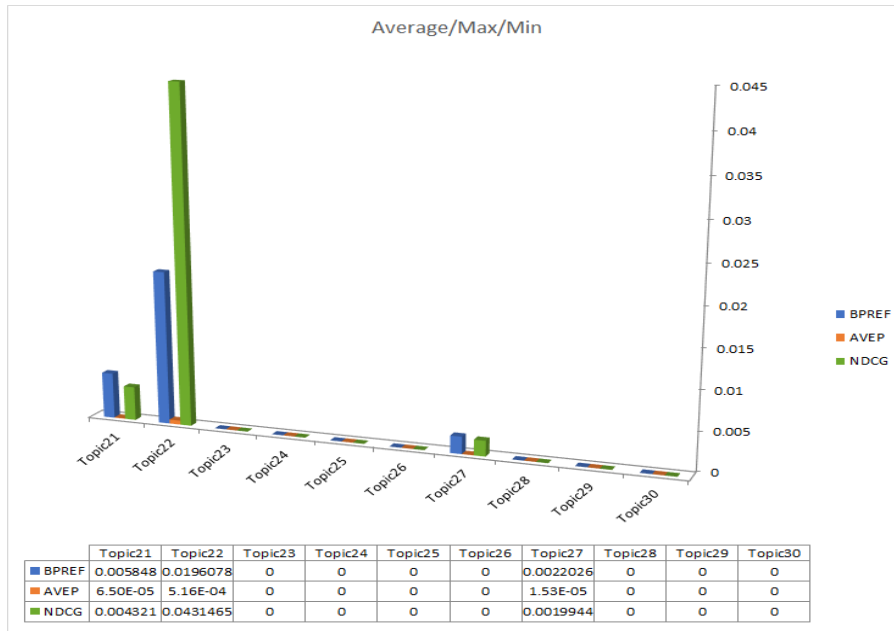
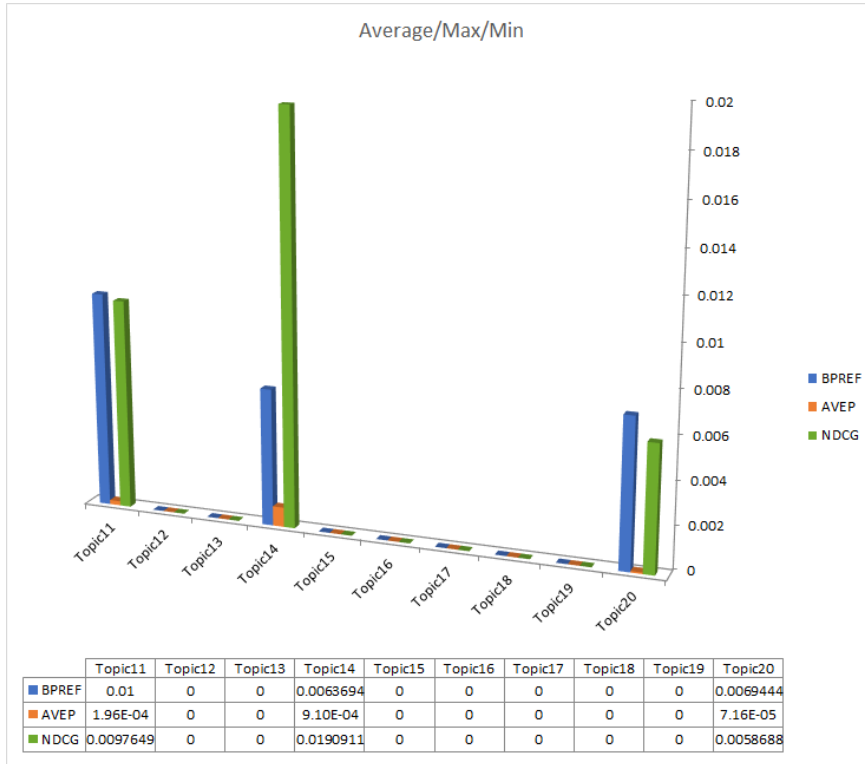
$$NDCG = \frac{DCG}{IDCG}$$

όπου E είναι τα ευρεθέντα έγγραφα, Σ είναι τα συναφή έγγραφα, και rel(i) είναι μία συνάρτηση η οποία μας επιστρέφει 1 αν το έγγραφο i είναι συναφές, διαφορετικά 0.

- Μετρήσεις

Υπολογίσαμε τις min, max και average τιμές για κάθε topic από το results.txt αρχείο. Παρακάτω παρουσιάζονται τα δεδομένα που εξάγαμε διαγραμματικά:





- *Απαιτούμενα αρχεία για την λειτουργία του συστήματος*

Σχετικά με το indexing και το queryevaluator , θα πρέπει να υπάρχει ένας φάκελος με το όνομα stoplists , ο οποίος θα περιέχει τα αρχεία stopwordsEn.txt , stopwordsGr.txt και θα βρίσκεται στον ίδιο φάκελο με τα indexer.jar queryevaluator.jar καθώς χρειάζονται για την αφαίρεση των stopwords. Για το queryevaluator.jar χρειαζόμαστε επίσης το αρχείο grels.txt και τον έναν φάκελο με όνομα topics ο οποίος θα περιέχει το αρχείο topics.xml καθώς και τον φάκελο CollectionIndex ο οποίος περιέχει τα αρχεία VocabularyFile.txt, PostingFile.txt και DocumentsFile.txt τα οποία παράχθηκαν από τον Indexer. Όλα τα παραπάνω αρχεία θα πρέπει να βρίσκονται στο ίδιο path με τα jar αρχεία.

- *Επίλογος*

Συνοψίζοντας, όπως μπορεί κανείς να δει και στα διαγράμματα, το σύστημα μας είχε αρκετά καλή απόδοση, στις περιπτώσεις όπου η τομή των ευρεθέντων εγγράφων με τα συναφή δεν είναι το κενό σύνολο ή ο βαθμός συνάφειας των ευρεθέντων εγγράφων δεν είναι παντού μηδενικός. Η αποτυχία του συστήματος στις περιπτώσεις μηδενισμού, είναι πιθανόν να οφείλεται στον λάθος υπολογισμό της νόρμας και αυτό να έχει ως συνέπεια το σύστημα μας να επιστρέφει λάθος έγγραφα ως ευρεθέντα για κάποια επερώτηση.

- *Αναφορές*

1. Chris Buckley and Ellen M. Voorhees, "Retrieval evaluation with incomplete information.", Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval. ACM, 2004.
2. Tetsuya Sakai, "Alternatives to bpref.", Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval. ACM, 2007.