

ΛΥΣΕΙΣ 2^{ης} ΣΕΙΡΑΣ ΑΣΚΗΣΕΩΝ

Άσκηση 1

Θεωρείστε μια συλλογή κειμένων που περιέχει τα ακόλουθα 5 έγγραφα:

Έγγραφο 1: «Computer Games»

Έγγραφο 2: «Computer Games Computer Games»

Έγγραφο 3: « Games Theory and Computer »

Έγγραφο 4: «Computer for Computer »

Έγγραφο 5: «Cheap Games Computer Games»

- 1) Δώστε τη διανυσματική παράσταση του κάθε εγγράφου με βάρη TF-IDF. Θεωρείστε ότι η θέση της κάθε λέξης στα διανύσματα γίνεται κατά αλφαβητική σειρά.
- 2) Θεωρείστε την επερώτηση q_1 = «Computer Games». Υπολογίστε το TF-IDF διάνυσμα αυτής της επερώτησης και δώστε την διάταξη των εγγράφων που θα επιστρέψει ένα σύστημα που βασίζεται στο διανυσματικό μοντέλο. Σχεδιάστε το ανεστραμμένο ευρετήριο για αυτή τη συλλογή.

Λύση

1)

Έγγραφο 1: «Computer Games»

Έγγραφο 2: «Computer Games Computer Games»

Έγγραφο 3: « Games Theory and Computer »

Έγγραφο 4: «Computer for Computer »

Έγγραφο 5: «Cheap Games Computer Games»

	And	Cheap	Computer	for	Games	Theory	$\text{MAX}_k \{ \text{FREQ}_{ij} \}$
D_1	0	0	1	0	1	0	1
D_2	0	0	2	0	2	0	2
D_3	1	0	1	0	1	1	1
D_4	0	0	2	1	0	0	2
D_5	0	1	1	0	2	0	2
DF	1	1	5	1	4	1	
IDF	5/1	5/1	5/5	5/1	5/4	5/1	

- FREQ_{ij} = το πλήθος των εμφανίσεων του όρου i στο έγγραφο j
- $\text{IDF} = N / \text{DF}$
- $\text{MAX}_k \{ \text{FREQ}_{ij} \}$ = συχνότητα της λέξης με τη μέγιστη συχνότητα στο κείμενο

TF-IDF							
	And	Cheap	Computer	for	Games	Theory	MAX _k {FREQ _{ij} }
D ₁	0	0	1/1*5/5	0	1/1*5/4	0	1
D ₂	0	0	2/2*5/5	0	2/2*5/4	0	2
D ₃	1/1*5/1	0	1/1*5/5	0	1/1*5/4	1/1*5/1	1
D ₄	0	0	2/2*5/5	1/2*5/1	0	0	2
D ₅	0	1/2*5/1	1/2*5/5	0	2/2*5/4	0	2
DF	1	1	5	1	4	1	
IDF	5/1	5/1	5/5	5/1	5/4	5/1	

- $TF_{ij} = FREQ_{ij} / MAX_k \{FREQ_{ij}\}$
- $V_{ij} = TF_{ij} * IDF_i$

Οι διανυσματικές παραστάσεις των κειμένων είναι :

$$\begin{aligned}
 V_1 &= \{0, 0, 1, 0, 1.25, 0\}, & |V_1| &= 2,5625 \\
 V_2 &= \{0, 0, 1, 0, 1.25, 0\}, & |V_2| &= 2,5625 \\
 V_3 &= \{5, 0, 1, 0, 1.25, 5\}, & |V_3| &= 52,5625 \\
 V_4 &= \{0, 0, 1, 2.5, 0, 0\}, & |V_4| &= 7,25 \\
 V_5 &= \{0, 2.5, 0.5, 0, 1.25, 0\}, & |V_5| &= 8,0625
 \end{aligned}$$

2)

	And	Cheap	Computer	For	Games	Theory
q1= Computer Games	0	0	1/1*5/5	0	1/1*5/4	0
IDF	5/1	5/1	5/5	5/1	5/4	5/1

$$q_1 = \{0, 0, 1, 0, 1.25, 0\}, |q_1| = 2,5625$$

$$V_1 * q_1 = 1 * 1 + 1,25 * 1,25 = 1 + 1,5625 = 2,5625$$

$$V_2 * q_1 = 2,5625$$

$$V_3 * q_1 = 5 * 0 + 1 * 1 + 1,25 * 1,25 + 5 * 0 = 2,5625$$

$$V_4 * q_1 = 1$$

$$V_5 * q_1 = 0,5 * 1 + 1,25 * 1,25 = 2,0625$$

$$\text{CosSim}(d_j, q) = \frac{\vec{d}_j \cdot \vec{q}}{|\vec{d}_j| \cdot |\vec{q}|} = \frac{\sum_{i=1}^t (w_{ij} \cdot w_{iq})}{\sqrt{\sum_{i=1}^t w_{ij}^2 \cdot \sum_{i=1}^t w_{iq}^2}}$$

Με βάση τον παραπάνω τύπο, υπολογίζουμε το μέτρο ομοιότητας συνημίτονου για κάθε έγγραφο D_j

$$R(D_1, q_1) = 2,5625 / (2,5625 * 2,5625)^{1/2} \Rightarrow R(D_1, q_1) = 1$$

$$R(D_2, q_1) = 2,5625 / (2,5625 * 2,5625)^{1/2} \Rightarrow R(D_2, q_1) = 1$$

$$R(D_3, q_1) = 2,5625 / (52,5625 * 2,5625)^{1/2} = 2,5625 / (134,69140625)^{1/2} = 2,5625 / 11,6056627 \Rightarrow R(D_3, q_1) = \mathbf{0,21928589}$$

$$R(D_4, q_1) = 1 / (7,25 * 2,5625)^{1/2} = 1 / (18,578125)^{1/2} = 1 / 4,31023491 \Rightarrow R(D_4, q_1) = \mathbf{0,23200592}$$

$$R(D_5, q_1) = 2,0625 / (8,0625 * 2,5625)^{1/2} = 2,0625 / (20,660156)^{1/2} = 2,0625 / 4,54534443 \Rightarrow R(D_5, q_1) = \mathbf{0,45376099}$$

Με βάση το διανυσματικό μοντέλο η διάταξη των εγγράφων είναι :
 $\langle \{D_1, D_2\}, D_5, D_4, D_3 \rangle$

Θα περιμέναμε το έγγραφο D_3 να έρθει στη σειρά πριν το D_4 , επειδή περιέχει όλους τους όρους της επερώτησης, όμως περιέχει και τον όρο Theory ο οποίος εμφανίζεται μόνο σε αυτό το έγγραφο και αυτό επηρέασε το βάρος του D_3 .

Ανεστραμμένο ευρετήριο

Μία μορφή του ανεστραμμένου ευρετηρίου στο οποίο εμφανίζονται μόνο οι θέσεις των όρων είναι :

Term	< Document Frequency, (Document; Position) >
Computer	< 5 (D ₁ ;1), (D ₂ ;1), (D ₂ ;3), (D ₃ ;4), (D ₄ ;1), (D ₄ ;3), (D ₅ ;3)>
Games	< 4 (D ₁ ;2), (D ₂ ;2), (D ₂ ;4), (D ₃ ;1), (D ₅ ;2), (D ₅ ;4) >
Theory	< 1 (D ₃ ;2) >
Cheap	< 1 (D ₅ ;1) >
and	< 1 (D ₃ ;3)>
for	< 1 (D ₄ ;2)>

Μία άλλη μορφή του ανεστραμμένου ευρετηρίου στο οποίο εμφανίζεται το TF του κάθε όρου σε κάθε έγγραφο είναι :

Term	< Document : Term Frequency : { Position } >
Computer	< D ₁ : 1 : { 1 } >
	< D ₂ : 1 : { 1 } >
	< D ₃ : 1 : { 4 } >
	< D ₄ : 1 : { 1, 3 } >
	< D ₅ : 0.5 : { 3 } >
Games	< D ₁ : 1.25 : { 2 } >
	< D ₂ : 1.25 : { 2, 4 } >
	< D ₃ : 1.25 : { 1 } >
	< D ₅ : 1.25 : { 2 } >
Theory	< D ₃ : 5 : { 2 } >
Cheap	< D ₅ : 2.5 : { 1 } >
and	< D ₃ : 5 : { 3 } >
for	< D ₄ : 2.5 : { 2 } >

Άσκηση 2

Έστω μια συλλογή από κείμενα D , και έστω A ένα διατεταγμένο υποσύνολο αυτής. Έστω ότι μας δίνουν το A και μας ζητούν να βρούμε αν υπάρχει επερώτηση q τ.ω. η απάντηση της να έχει στην αρχή της το διατεταγμένο σύνολο A . Για παράδειγμα, αν $A = \langle d_1, d_2, d_3 \rangle$ και βρούμε μια επερώτηση q τ.ω. $\text{Answer}(q) = \langle d_1, d_2, d_3, d_8, \dots \rangle$ τότε αυτή είναι μια λύση του προβλήματος μας. Θεωρώντας ότι το σύστημα σας βασίζεται στο διανυσματικό μοντέλο, απαντήστε τα παρακάτω ερωτήματα.

(α) Πως μπορούμε να βρούμε αν υπάρχει τέτοια επερώτηση;

(β) Αν υπάρχει ποια είναι;

(γ) Αν δεν υπάρχει τέτοια επερώτηση, πως θα χαλαρώνατε το πρόβλημα και τι θα μπορούσατε να επιστρέψετε; Μπορείτε να αναπτύξετε τις σκέψεις σας όσο θέλετε.

Σημείωση: Προσέξτε ώστε το υπολογιστικό κόστος των λύσεων που θα προτείνετε για τα (α) και (β) να μην είναι απαγορευτικό.

Λύση

(α) Για να υπάρχει μία τέτοια επερώτηση q θα πρέπει να ισχύουν οι δύο παρακάτω συνθήκες:

(α1) Έστω q η επερώτηση που ψάχνουμε. Για να επιστρέφει η q όλα τα έγγραφα του A και μάλιστα με την σχετική διάταξη που έχουν στο A , θα πρέπει το μέτρο ομοιότητας του συνημίτονου μεταξύ της q και του πρώτου εγγράφου στο A να είναι μεγαλύτερο από το μέτρο του δευτέρου εγγράφου και εκείνο μεγαλύτερο από του τρίτου εγγράφου κ.ο.κ., και όλα να είναι μεγαλύτερα του μηδενός

Δηλαδή αν $A = \langle d_1, d_2, d_3 \rangle$ τότε θα έπρεπε

$$\text{Sim}(d_1, q) > \text{Sim}(d_2, q) > \text{Sim}(d_3, q) > 0 \quad (1)$$

(α2) Για να μας επιστρέφει η q όλα τα έγγραφα που ανήκουν στο A , με την διάταξη που έχουν σε αυτό και πριν από οποιοδήποτε άλλο έγγραφο θα πρέπει όλα τα έγγραφα που μας επιστρέφει η επερώτηση q και δεν ανήκουν στο A , να έχουν μέτρο ομοιότητας του συνημίτονου μικρότερο από την τιμή του τελευταίου εγγράφου που ανήκει στο A .

Δηλαδή αν $A = \langle d_1, d_2, d_3 \rangle$ τότε θα έπρεπε

$$\text{Sim}(d_1, q) > \text{Sim}(d_2, q) > \text{Sim}(d_3, q) > \text{Sim}(d_x, q) \quad \text{για κάθε } d_x \text{ που δεν ανήκει στο } A.$$

(β)

Εύκολα μπορούμε να δείξουμε ότι αν ικανοποιούνται οι παραπάνω συνθήκες τότε το $q=d1$ είναι μια επιθυμητή απάντηση (από τις ενδεχομένως πολλές). Ο έλεγχος της συνθήκης (α1) είναι απλός και όχι ιδιαίτερα ακριβός. Συγκεκριμένα απαιτεί $|A|-1$ υπολογισμούς βαθμού ομοιότητας. Ένας εύκολος τρόπος για να δούμε αν ισχύει η συνθήκη (α2) είναι να υπολογίσουμε το $\text{Answer}(d1)$ και να δούμε εάν τα πρώτα $|A|$ στοιχεία του είναι τα στοιχεία του A .

(γ)

Μια χαλάρωση του προβλήματος είναι η εξής: Η συνθήκη (α1) ικανοποιείται αλλά δεν ικανοποιείται η συνθήκη (α2). Και σε αυτήν την περίπτωση το $q=d1$ θα ήταν μια πιθανή λύση του προβλήματος. Απλά η απάντηση του q θα μπορούσε να είχε τη μορφή: $A(q) = \langle \mathbf{d1}, d5, d6, \mathbf{d2}, d4, d8, d9, \mathbf{d3}, d10, d7 \rangle$

Μια άλλη χαλάρωση του προβλήματος θα ήταν να μειώσουμε το σύνολο των εγγράφων του συνόλου A αρχίζοντας από το τέλος. Δηλαδή αντί για $A = \langle d_1, d_2, d_3 \rangle$ να δούμε αν υπάρχει λύση για το σύνολο $A' = \langle \mathbf{d1}, \mathbf{d2} \rangle$. Αν δεν υπάρχει ούτε για το A' να δούμε αν υπάρχει για το $A'' = \langle \mathbf{d1} \rangle$.

Για περισσότερα δείτε τις διαφάνειες του Μαθήματος 11 καθώς και το άρθρο:

http://www.ics.forth.gr/~tzitzik/publications/2007_TzitzikasTheoharisNaming.pdf

Άσκηση 3

Στο μάθημα είδαμε δύο μοντέλα ανάκτησης του βασίζονται στη Θεωρία Ασαφών Συνόλων. Το πρώτο θεωρεί βάρυνση $TF*IDF$, ενώ το δεύτερο είναι εκείνο που προτάθηκε από τους [Ogawa, Morita, Kobayashi, 1991]. Θεωρείστε έναν όρο t_i ενός εγγράφου d_j . Συγκρίνετε την συμπεριφορά των δύο αυτών μοντέλων για διάφορες περιπτώσεις, π.χ.: για μικρές και μεγάλες τιμές του tf_{ij} , για μικρές και μεγάλες τιμές του idf_i , για μικρές και μεγάλες τιμές του w_{ij} αν προκύπτει από $tf*idf$.

Λύση

Για το πρώτο μοντέλο που βασίζεται σε βάρυνση $TF-IDF$ ξέρουμε ότι:

$d_j = (w_{1,j}, w_{2,j}, \dots, w_{t_j})$ όπου $w_{i,j} \in [0,1]$

$R(d_j, t_i) = \mu_{t_i}(d_j) = w_{i,j} = tf_{ij} * idf_i$ όπου $tf_{ij} = \text{freq}_{ij} / \text{MAX}_k \{ \text{freq}_{kj} \}$, $idf_i = \log_2 (N / df_i)$

και N ο αριθμός των εγγράφων.

Για το μοντέλο που προτάθηκε από τους [Ogawa, Morita, and Kobayashi,1991]

ξέρουμε ότι: $d_j = (w_{1,j}, w_{2,j}, \dots, w_{t_j})$ όπου $w_{i,j} \in \{0,1\}$ και

$w_{i,j} = 1$ όταν ο όρος t_i εμφανίζεται στο κείμενο d_j . (αλλιώς $w_{i,j} = 0$)

$R(d_j, t_i) = \mu_{ti}(d_j)$ το οποίο ορίζεται ως εξής:

Αρχικά ορίζεται η εγγύτητα μεταξύ των όρων με τον εξής τύπο:

$c(i, j) = n(i, j) / n_i + n_j - n(i, j)$ όπου

$n(i, j)$: το πλήθος των εγγράφων που περιέχουν τον όρο k_i και τον k_j .

n_i : το πλήθος των εγγράφων που περιέχουν τον όρο k_i .

n_j : το πλήθος των εγγράφων που περιέχουν τον όρο k_j .

Κατόπιν θέτουμε $\mu_i(j) = \sum_{t_w \in d_j} c(i, w)$,

$\mu_i(j) = 1 - \prod(1 - c(i, w))$, $t_w \in d_j$

Έστω ότι θέλουμε να κάνουμε μία επερώτηση σε μία σύλλογη εγγράφων και η επερώτηση αποτελείται από ένα όρο k , δηλαδή $q = "k"$.

Μικρές τιμές tf_{ij}

Έστω ένα έγγραφο j . Αν το tf_{ij} του όρου k_i είναι μικρό, αυτό σημαίνει ότι ο όρος k_i εμφανίζεται λίγες φορές στο έγγραφο αυτό άρα η κατάταξη του εγγράφου θα είναι χαμηλή σύμφωνα με το πρώτο μοντέλο ανάκτησης. Το δεύτερο μοντέλο ([Ogawa, Morita, Kobayashi, 1991]) αγνοεί το πλήθος των εμφανίσεων ενός όρου σε ένα έγγραφο. Λαμβάνει όμως υπόψη τον βαθμό συνεμφάνισης των όρων στη συλλογή. Αυτό σημαίνει ότι αν το έγγραφο j περιέχει πολλούς όρους οι οποίοι έχουν μεγάλη εγγύτητα (συνεμφάνιση) με τον όρο k_i , τότε το έγγραφο αυτό μπορεί να σταθμιστεί υψηλά.

Μεγάλες τιμές tf_{ij}

Έστω ένα έγγραφο j . Αν το tf_{ij} του όρου k_i είναι μεγάλο, αυτό σημαίνει ότι ο όρος k_i εμφανίζεται πολλές φορές στο έγγραφο αυτό άρα η κατάταξη του εγγράφου θα είναι υψηλή σύμφωνα με το πρώτο μοντέλο ανάκτησης. Στο δεύτερο μοντέλο (που αγνοεί το tf) το έγγραφο αυτό θα μπορούσε να σταθμιστεί χαμηλά αν περιέχει λίγους όρους που έχουν εγγύτητα (συνεμφάνιση) με τον όρο k_i .

Μικρές τιμές idf_i

Αν το idf του όρου k_i είναι μικρό αυτό σημαίνει ότι το πλήθος των εγγράφων που περιέχουν τον όρο k_i είναι μεγάλο (αφού $idf = N/df$).

Έστω ένα έγγραφο j που περιέχει τον όρο k_i .

Αν λαμβάναμε υπόψη μόνο το idf τότε το 1^ο μοντέλο, θα έδινε μικρό βαθμό συνάφειας στο έγγραφο j (καθώς και σε όλα τα υπόλοιπα έγγραφα). Αντίθετα το 2^ο μοντέλο ενδεχομένως να έδινε μεγαλύτερο βαθμό συνάφειας στο έγγραφο j διότι από τη στιγμή που ο όρος k_i εμφανίζεται σε πολλά κείμενα μπορεί να έχει υψηλό βαθμό συνεμφάνισης με άλλους όρους που περιέχει το έγγραφο j .

Μεγάλες τιμές idf_j

Αν το idf του όρου **ki** είναι μεγάλο αυτό σημαίνει ότι το πλήθος των εγγράφων που περιέχουν τον όρο **ki** είναι μικρό (αφού $idf = N/df$).

Έστω ένα έγγραφο j που περιέχει τον όρο ki.

Αν λαμβάναμε υπόψη μόνο το idf τότε το 1^ο μοντέλο, θα έδινε μεγάλο βαθμό συνάφειας στο έγγραφο j (καθώς και σε όλα τα υπόλοιπα έγγραφα). Αντίθετα το 2^ο μοντέλο ενδεχομένως να έδινε μικρότερο βαθμό συνάφειας στο έγγραφο j διότι από τη στιγμή που ο όρος ki εμφανίζεται σε πολλά κείμενα μπορεί να έχει υψηλό βαθμό συνεμφάνισης με άλλους όρους που περιέχει το έγγραφο j.

Μικρές και μεγάλες τιμές του w_{ij}

Θα μπορούσαμε να διακρίνουμε τις παρακάτω περιπτώσεις:

1) Μικρό TF*IDF λόγω πολύ μικρού TF

Έστω ένα έγγραφο j. Αν το **tf_{ij}** του όρου **ki** είναι πολύ μικρό, αυτό σημαίνει ότι ο όρος **ki** εμφανίζεται πολύ λίγες φορές στο έγγραφο αυτό άρα η κατάταξη του εγγράφου θα είναι χαμηλή σύμφωνα με το πρώτο μοντέλο ανάκτησης. Το δεύτερο μοντέλο αγνοεί το πλήθος των εμφανίσεων ενός όρου σε ένα έγγραφο. Λαμβάνει όμως υπόψη τον βαθμό συνεμφάνισης των όρων στη συλλογή. Αυτό σημαίνει ότι αν το έγγραφο j περιέχει πολλούς όρους οι οποίοι έχουν μεγάλη εγγύτητα (συνεμφάνιση) με τον όρο ki, τότε το έγγραφο αυτό μπορεί να σταθμιστεί υψηλά.

2) Μικρό TF*IDF λόγω πολύ μικρού IDF

Αν το idf του όρου **ki** είναι πολύ μικρό αυτό σημαίνει ότι το πλήθος των εγγράφων που περιέχουν τον όρο **ki** είναι πολύ μεγάλο (αφού $idf = N/df$).

Έστω ένα έγγραφο j που περιέχει τον όρο ki.

Αν λαμβάναμε υπόψη μόνο το idf τότε το 1^ο μοντέλο, θα έδινε πολύ μικρό βαθμό συνάφειας στο έγγραφο j (καθώς και σε όλα τα υπόλοιπα έγγραφα). Αντίθετα το 2^ο μοντέλο ενδεχομένως να έδινε πολύ μεγαλύτερο βαθμό συνάφειας στο έγγραφο j διότι από τη στιγμή που ο όρος ki εμφανίζεται σε πολλά κείμενα μπορεί να έχει υψηλό βαθμό συνεμφάνισης με άλλους όρους που περιέχει το έγγραφο j.

3) Μικρό TF*IDF λόγω μικρού TF και μικρού IDF

Έστω ένα έγγραφο j . Αν το tf_{ij} του όρου ki είναι μικρό, αυτό σημαίνει ότι ο όρος ki εμφανίζεται λίγες φορές στο έγγραφο αυτό και αν το idf του όρου ki είναι μικρό αυτό σημαίνει ότι το πλήθος των εγγράφων που περιέχουν τον όρο ki είναι μεγάλο (αφού $idf = N/df$). Επομένως, το πρώτο μοντέλο θα έδινε μικρό βαθμό συνάφειας στο έγγραφο j . Αντίθετα το 2^ο μοντέλο ενδεχομένως να έδινε πολύ μεγαλύτερο βαθμό συνάφειας στο έγγραφο j διότι από τη στιγμή που ο όρος ki εμφανίζεται σε πολλά κείμενα μπορεί να έχει υψηλό βαθμό συνεμφάνισης με άλλους όρους που περιέχει το έγγραφο j .

4) Μεγάλο TF*IDF λόγω πολύ μεγάλου TF

Έστω ένα έγγραφο j . Αν το tf_{ij} του όρου ki είναι πολύ μεγάλο, αυτό σημαίνει ότι ο όρος ki εμφανίζεται πολλές φορές στο έγγραφο αυτό άρα η κατάταξη του εγγράφου θα είναι πολύ υψηλή σύμφωνα με το πρώτο μοντέλο ανάκτησης. Στο δεύτερο μοντέλο (που αγνοεί το tf) το έγγραφο αυτό θα μπορούσε να σταθμιστεί χαμηλά αν περιέχει λίγους όρους που έχουν εγγύτητα (συνεμφάνιση) με τον όρο ki .

5) Μεγάλο TF*IDF λόγω πολύ μεγάλου IDF

Αν το idf του όρου ki είναι πολύ μεγάλο αυτό σημαίνει ότι το πλήθος των εγγράφων που περιέχουν τον όρο ki είναι πολύ μικρό (αφού $idf = N/df$).

Έστω ένα έγγραφο j που περιέχει τον όρο ki .

Αν λαμβάναμε υπόψη μόνο το idf τότε το 1^ο μοντέλο, θα έδινε πολύ μεγάλο βαθμό συνάφειας στο έγγραφο j (καθώς και σε όλα τα υπόλοιπα έγγραφα). Αντίθετα το 2^ο μοντέλο ενδεχομένως να έδινε πολύ μικρότερο βαθμό συνάφειας στο έγγραφο j διότι από τη στιγμή που ο όρος ki εμφανίζεται σε πολλά κείμενα μπορεί να έχει υψηλό βαθμό συνεμφάνισης με άλλους όρους που περιέχει το έγγραφο j .

6) Μεγάλο TF*IDF λόγω μεγάλου TF και μεγάλου IDF

Έστω ένα έγγραφο j . Αν το tf_{ij} του όρου ki είναι μεγάλο, αυτό σημαίνει ότι ο όρος ki εμφανίζεται πολλές φορές στο έγγραφο αυτό και αν το idf του όρου ki είναι μεγάλο αυτό σημαίνει ότι το πλήθος των εγγράφων που περιέχουν τον όρο ki είναι μικρό (αφού $idf = N/df$). Επομένως, το πρώτο μοντέλο θα έδινε μεγάλο βαθμό συνάφειας στο έγγραφο j .

Αντίθετα το 2^ο μοντέλο ενδεχομένως να έδινε πολύ μικρότερο βαθμό συνάφειας στο έγγραφο j διότι από τη στιγμή που ο όρος ki εμφανίζεται σε λίγα κείμενα μπορεί να έχει μικρό βαθμό συνεμφάνισης με άλλους όρους που περιέχει το έγγραφο j .