



**Ανάκληση Πληροφορίας**

**Information Retrieval**

Διδάσκων –  
Δημήτριος Κατσαρός



# Το μειονέκτημα της Boolean ανάκτησης

- Τα Boolean ερωτήματα δίνουν είτε πολύ λίγα αποτελέσματα (=0) είτε πάρα πολλά (1000s)
  - Query 1: “*standard user dlink 650*” → 200,000 αποτελέσματα
  - Query 2: “*standard user dlink 650 no card found*”: 0 αποτελέσματα
- Απαιτείται μεγάλη δεξιότητα και εμπειρία για να φτιάξεις ένα ερώτημα που παράγει διαχειρίσιμο αριθμό αποτελεσμάτων
  - Το AND δίνει λίγα και το OR δίνει πολλά αποτελέσματα



# Μοντέλα διαβαθμισμένης ανάκτησης

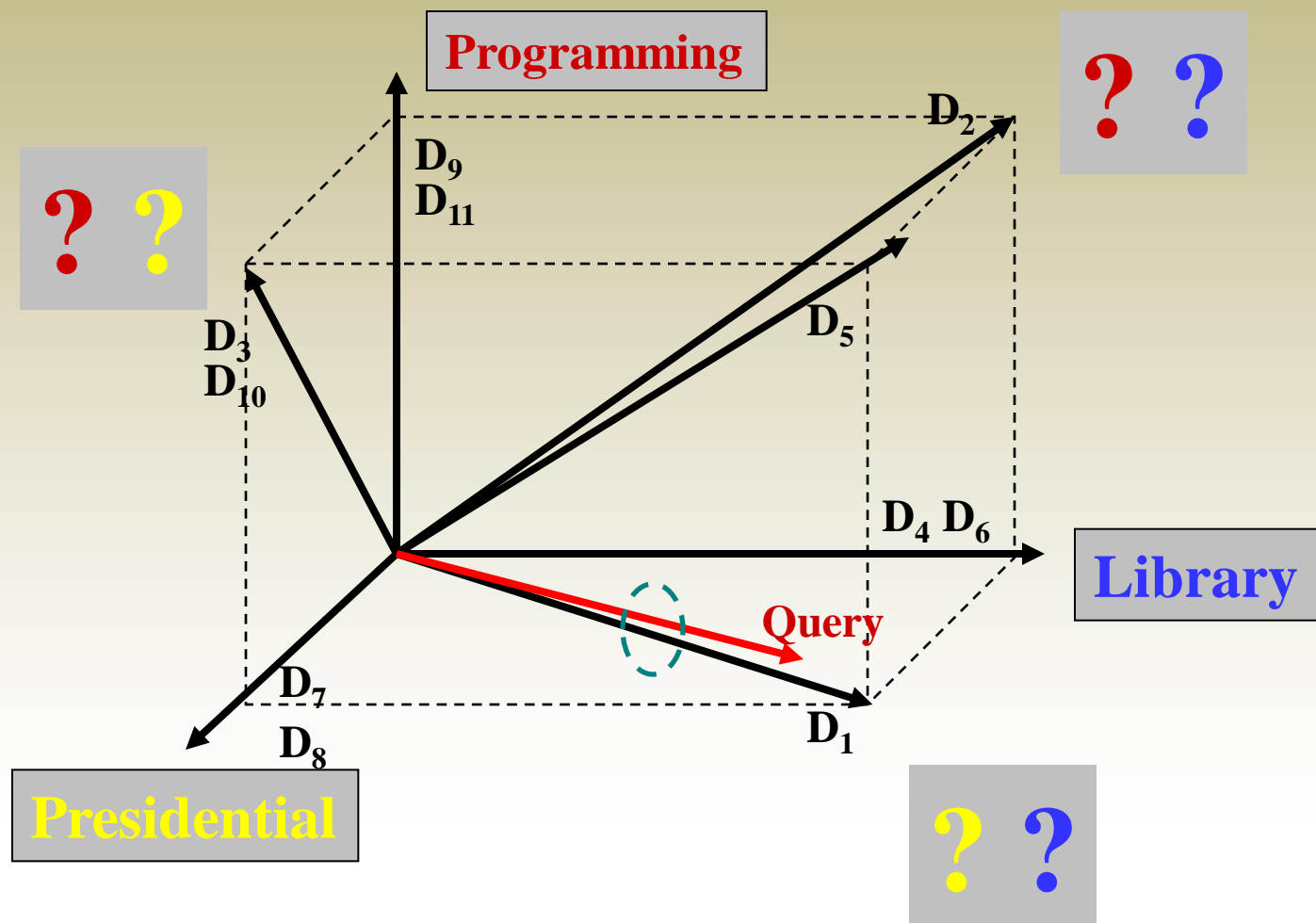
- Αντί για ένα σύνολο εγγράφων που ικανοποιεί ακριβώς ένα ερώτημα, στα **ranked retrieval models**, το σύστημα επιστρέφει μια διάταξη (διαβάθμιση) των (top) εγγράφων της συλλογής σε σχέση με το ερώτημα
- **Free text queries**
  - Αντί για μια γλώσσα ερωτημάτων με τελεστές και εκφράσεις,
  - το ερώτημα του χρήστη αποτελείται απλά μια ή περισσότερες λέξεις κάποιας ανθρώπινης γλώσσας
- Στην πράξη τα ranked retrieval models έχουν συσχετιστεί με τα free text ερωτήματα, και αντιστρόφως



# Το σκορ ως βάση της διαβαθμισμένης ανάκτησης

- Θα επιθυμούσαμε να επιστρέψουμε τα έγγραφα σε σειρά που να είναι πιο χρήσιμα στον χρήστη
- Πώς θα μπορούσαμε να το επιτύχουμε;
- Με την ανάθεση ενός score – π.χ., στο διάστημα τιμών  $[0, 1]$  – για κάθε έγγραφο
- Αυτό το score μετρά πόσο καλά “ταιριάζει” το έγγραφο με το ερώτημα

# Μοντέλο Vector Space Retrieval





# Υιοθέτηση του Vector Space μοντέλου

- Το VSM είναι στην πραγματικότητα ένα πλαίσιο (framework) και απαιτείται:
  - Ορισμός των διαστάσεων
  - Απόφαση για το πού θα τοποθετήσουμε τα έγγραφα και τα ερωτήματα ως διανύσματα στον διανυσματικό χώρο (δηλαδή το μήκος τους σε κάθε διάσταση)
  - Ορισμός της ομοιότητας μεταξύ δυο διανυσμάτων



# Αναπαράσταση με bit vectors & ομοιότητα

$$q = (x_1, \dots, x_N)$$

$$x_i, y_i \in \{0, 1\}$$

1: word  $W_i$  is **present**

$$d = (y_1, \dots, y_N)$$

0: word  $W_i$  is **absent**

$$\text{Sim}(q, d) = q \cdot d = x_1 y_1 + \dots + x_N y_N = \sum_{i=1}^N x_i y_i$$

# Παράδειγμα εφαρμογής των bit vectors

Query = “news about presidential campaign”

Ideal ranking?

$d_1$  ... news about ...

$d_4 +$

$d_2$  ... news about organic food campaign ...

$d_3 +$

$d_3$  ... news of presidential campaign ...

$d_4$  ... news of presidential campaign ...  
... presidential candidate ...

$d_1 -$

$d_2 -$

$d_5$  ... news of organic food campaign ...  
campaign ... campaign ... campaign ...

$d_5 -$



# Υπολογισμός της ομοιότητας

Query = "news about presidential campaign"

$d_1$  ... news about ...

$d_3$  ... news of presidential campaign ...

$V = \{\text{news, about, presidential, campaign, food ...}\}$

$q = (1, 1, 1, 1, 0, \dots)$

$d_1 = (1, 1, 0, 0, 0, \dots)$

$$f(q, d_1) = 1 * 1 + 1 * 1 + 1 * 0 + 1 * 0 + 0 * 0 + \dots = 2$$


$d_3 = (1, 0, 1, 1, 0, \dots)$

$$f(q, d_3) = 1 * 1 + 1 * 0 + 1 * 1 + 1 * 1 + 0 * 0 + \dots = 3$$

# Διαβάθμιση (ranking) των εγγράφων

Query = “news about presidential campaign”

$d_1$	... news about ...	$f(q, d_1) = 2$
$d_2$	... news about organic food campaign ...	$f(q, d_2) = 3$
$d_3$	... news of presidential campaign ...	$f(q, d_3) = 3$
$d_4$	... news of presidential campaign ... ... presidential candidate ...	$f(q, d_4) = 3$
$d_5$	... news of organic food campaign ... campaign ... campaign ... campaign ...	$f(q, d_5) = 2$



# Αναπαράσταση με διάνυσμα συχνότητας

$$q = (x_1, \dots, x_N) \quad x_i = \text{count of word } W_i \text{ in query}$$

$$d = (y_1, \dots, y_N) \quad y_i = \text{count of word } W_i \text{ in doc}$$

$$\text{Sim}(q, d) = q \cdot d = x_1 y_1 + \dots + x_N y_N = \sum_{i=1}^N x_i y_i$$

# Υπολογισμός της ομοιότητας

$d_2$  ... news about organic food campaign ...  $f(q, d_2) = 3$

$q = (1, 1, 1, 1, 0, \dots)$   
 $d_2 = (1, 1, 0, 1, 1, \dots)$

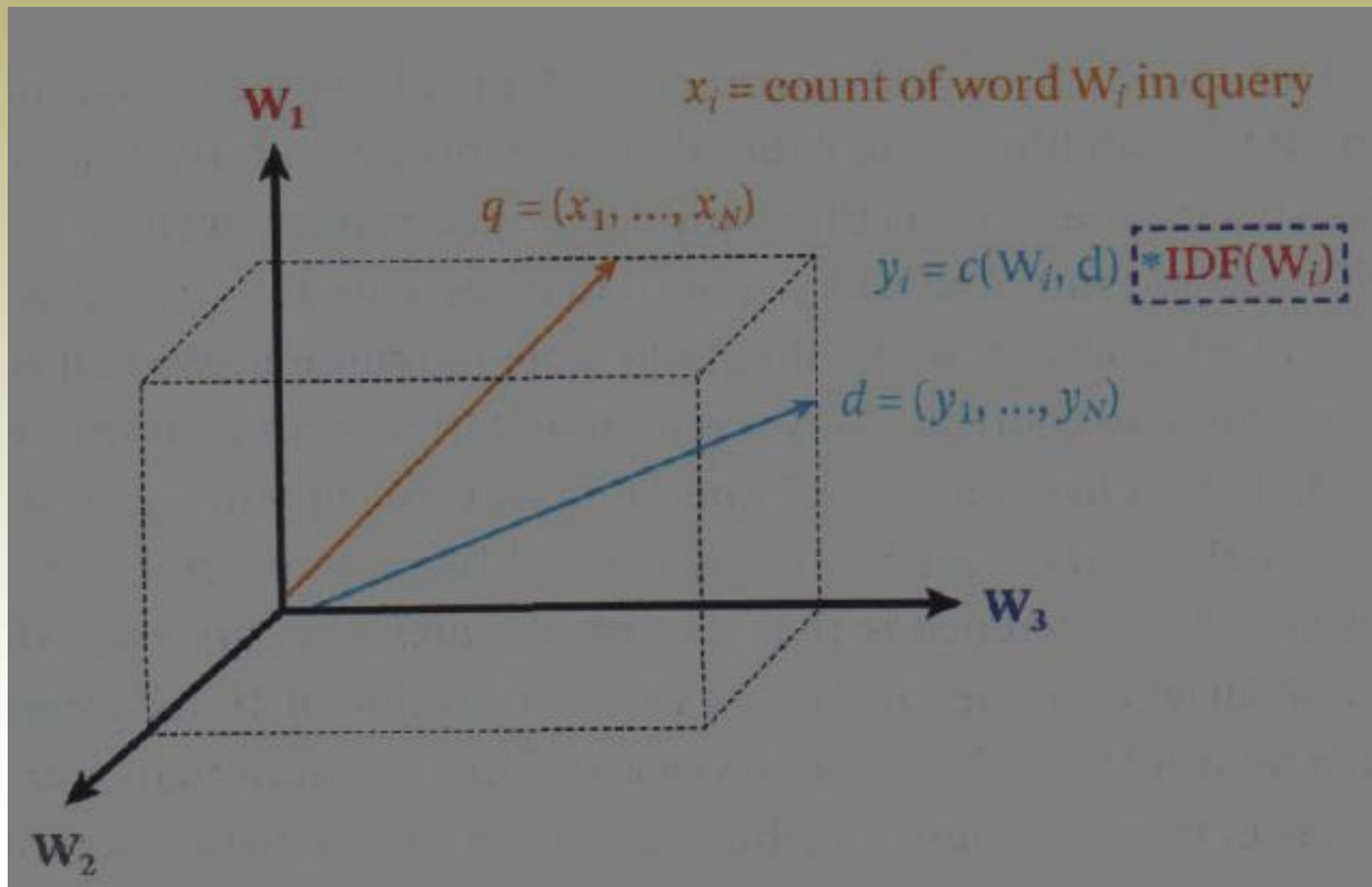
$d_3$  ... news of presidential campaign ...  $f(q, d_3) = 3$

$q = (1, 1, 1, 1, 0, \dots)$   
 $d_3 = (1, 0, 1, 1, 0, \dots)$

$d_4$  ... news of presidential campaign ...  
... presidential candidate ...  $f(q, d_4) = 4!$

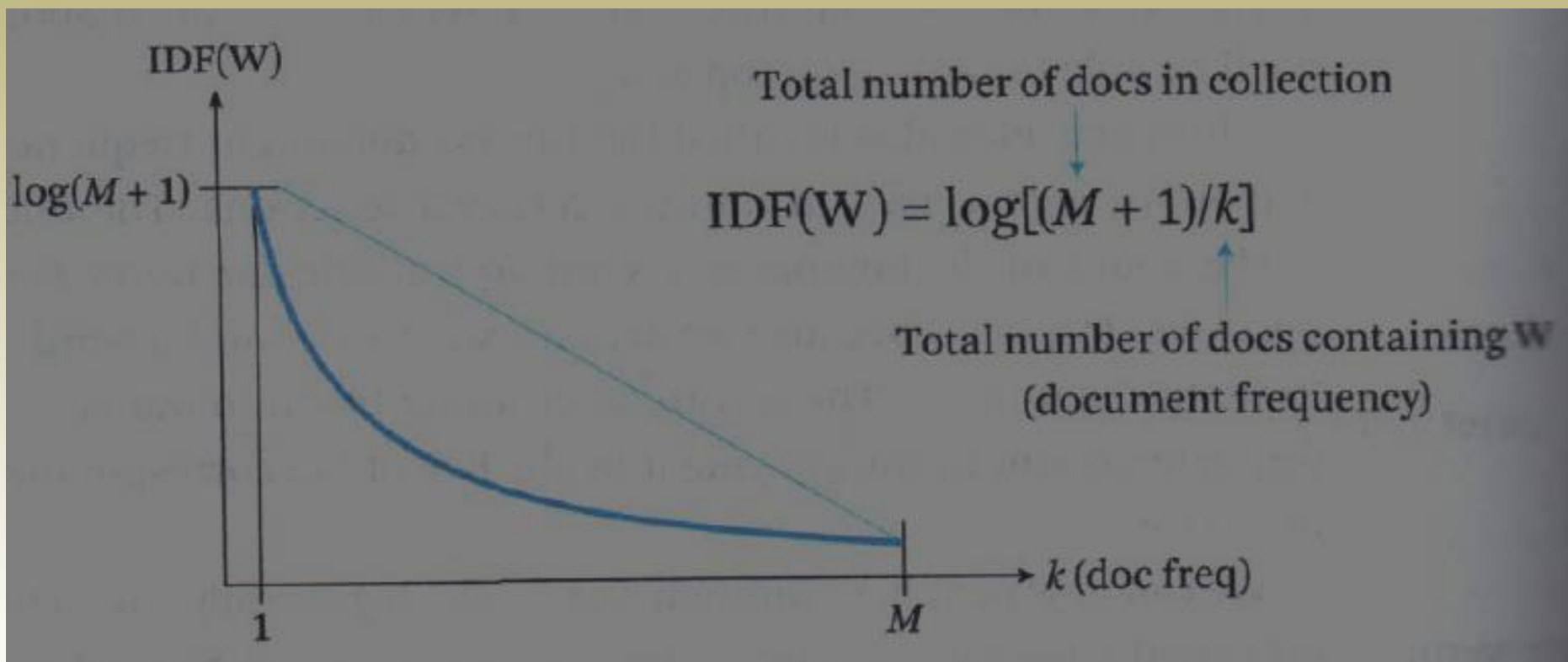
$q = (1, 1, 1, 1, 0, \dots)$   
 $d_4 = (1, 0, 2, 1, 0, \dots)$

# Αναπαράσταση με TF-IDF ζύγισμα





# IDF ως συνάρτηση της συχνότητας στα έγγραφα



# Η επίδραση του IDF στην διαβάθμιση

$d_2$  ... news about organic food campaign ...

$d_3$  ... news of presidential campaign ...

$V =$	{news,	about,	presidential,	campaign,	food	...}
IDF(W) =	1.5	1.0	2.5	3.1	1.8	

$q =$	(1,	1,	1,	1,	0,	...)
$d_2 =$	(1 * 1.5,	1 * 1.0,	0,	1 * 3.1,	0,	...)
$q =$	(1,	1,	1,	1,	0,	...)
$d_3 =$	(1 * 1.5,	0,	1 * 2.5,	1 * 3.1,	0,	...)

$$f(q, d_2) = 5.6 < f(q, d_3) = 7.1$$





# Scores με το TF-IDF ζύγισμα

Query = “news about presidential campaign”

$d_1$	... news about ...	$f(q, d_1) = 2.5$
$d_2$	... news about organic food campaign ...	$f(q, d_2) = 5.6$
$d_3$	... news of presidential campaign ...	$f(q, d_3) = 7.1$
$d_4$	... news of presidential campaign ... ... presidential candidate ...	$f(q, d_4) = 9.6$
$d_5$	... news of organic food campaign ... campaign ... campaign ... campaign ...	$f(q, d_5) = 13.9!$

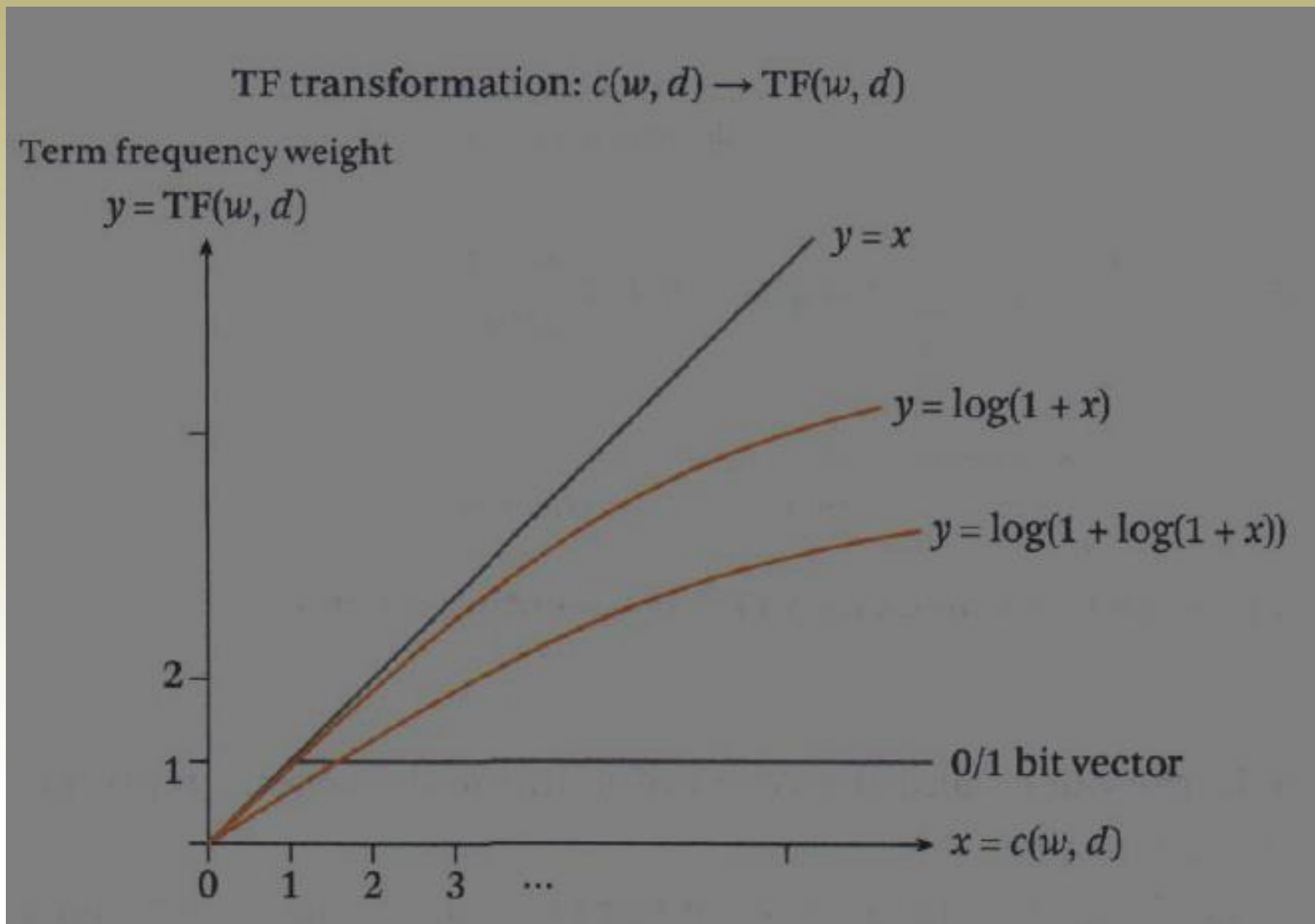


# Συνάρτηση διαβάθμισης (ranking) μέχρι στιγμής

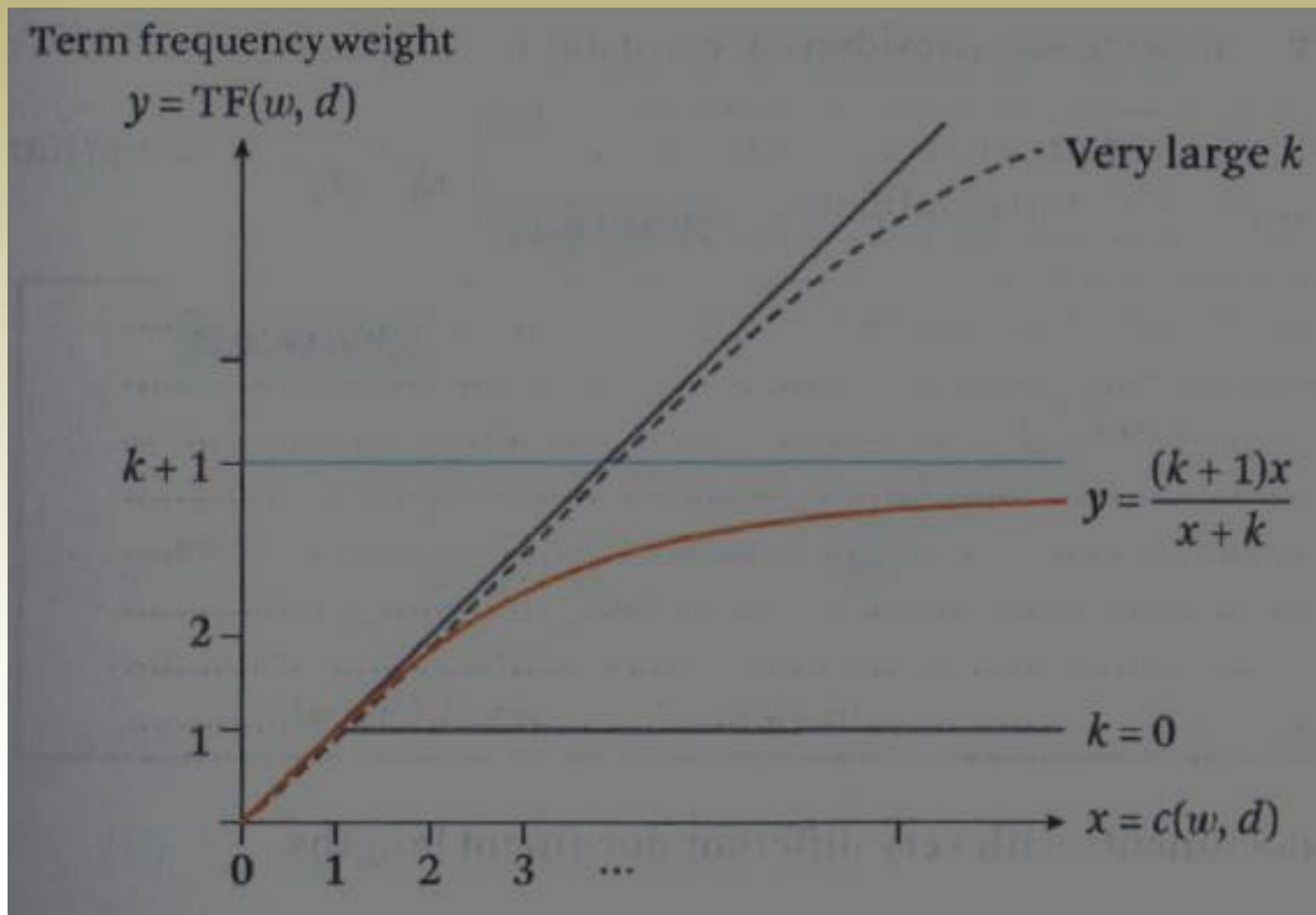
The diagram illustrates the ranking function  $f(q, d)$  with the following components and annotations:

- Equation:** 
$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum_{w \in q \cap d} c(w, q) c(w, d) \log \frac{M + 1}{df(w)}$$
- Annotations:**
  - Total number of documents in collection:** Points to  $M$  in the denominator of the log term.
  - Document frequency:** Points to  $df(w)$  in the denominator of the log term.
  - All matched query words in document:** Points to the summation index  $w \in q \cap d$ .

# Διαφορετικοί τρόποι ορισμού του TF



# Ορισμός του BM25 ως TF





# Έγγραφα με διαφορετικά μήκη

Query = "news about presidential campaign"

$d_4$

... news of presidential campaign ...  
... presidential candidate ...

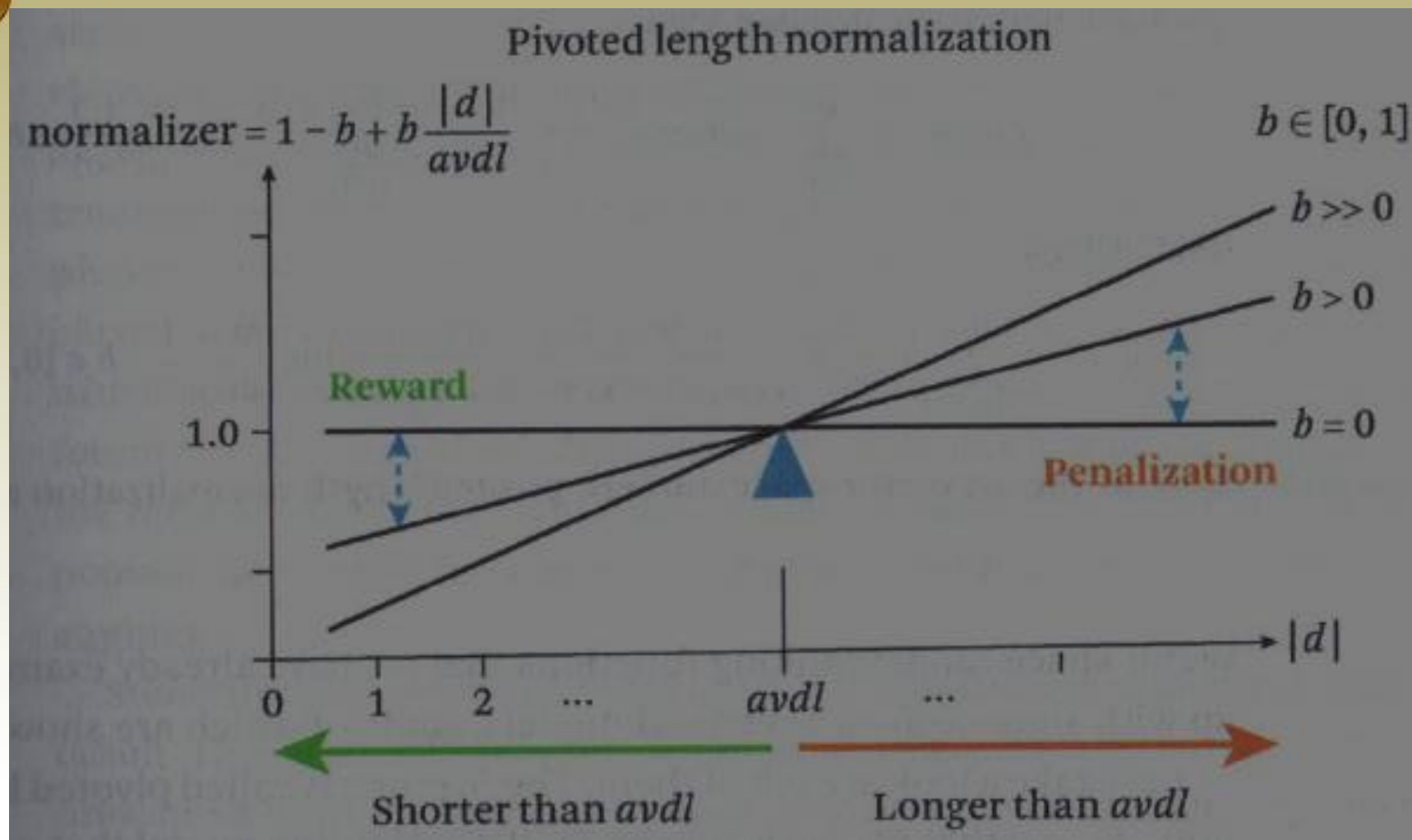
100 words

$d_6 > d_4?$

$d_6$

... campaign.....campaign ..... 5000 words ....  
.....  
..... news.....  
.....  
..... news.....  
.....  
..... presidential.....presidential.....

# Pivoted document length κανονικοποίηση





# State-of-the-art VSR μοντέλα

Pivoted length normalization VSM

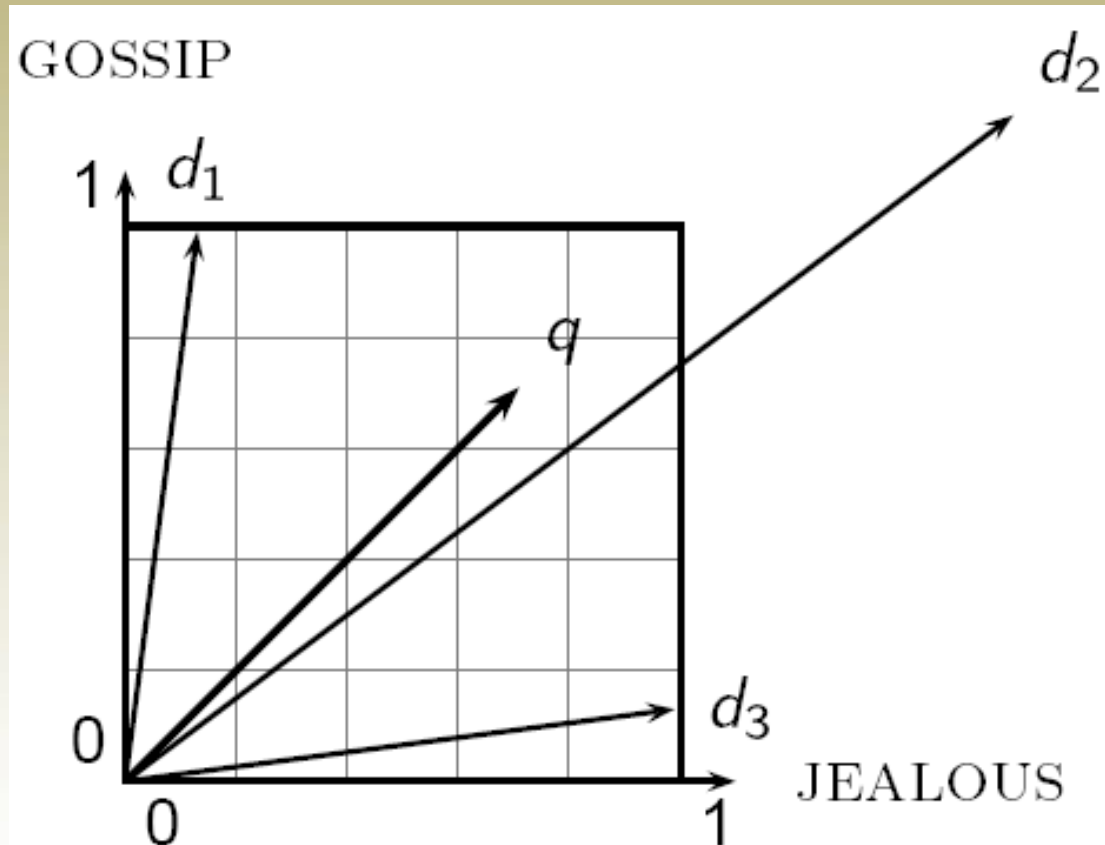
$$f(q, d) = \sum_{w \in q \cap d} c(w, q) \frac{\ln(1 + \ln(1 + c(w, d)))}{1 - b + b \frac{|d|}{\text{avdl}}} \log \frac{M + 1}{\text{df}(w)} \quad b \in [0, 1]$$

BM25/Okapi

$$f(q, d) = \sum_{w \in q \cap d} c(w, q) \frac{(k + 1)c(w, d)}{c(w, d) + k(1 - b + b \frac{|d|}{\text{avdl}})} \log \frac{M + 1}{\text{df}(w)} \quad b \in [0, 1], k \in [0, +\infty)$$

# Γιατί η απόσταση είναι κακή ιδέα

Η Ευκλείδεια απόσταση μεταξύ του  $\vec{q}$  και του  $\vec{d}_2$  είναι μεγάλη, παρόλο που η κατανομή των όρων στο ερώτημα  $\vec{q}$  και η κατανομή των όρων στο έγγραφο  $\vec{d}_2$  είναι πολύ όμοια





## Χρήση της γωνίας αντί της απόστασης

- Πείραμα σκέψης: πάρτε ένα έγγραφο  $d$  και append στον εαυτό του. Καλέστε αυτό το νέο έγγραφο ως  $d'$
- “Σημασιολογικά” τα  $d$  και  $d'$  έχουν το ίδιο περιεχόμενο
- Η Ευκλείδεια απόσταση μεταξύ των δυο εγγράφων (μπορεί) να είναι αρκετά μεγάλη
- Η γωνία όμως μεταξύ των δυο εγγράφων είναι 0, δηλαδή αντιστοιχεί στην μέγιστη ομοιότητα





# Από τις γωνίες στα συνημίτονα

- Οι επόμενες έννοιες είναι ισοδύναμες:
  - Διάταξη των εγγράφων σε φθίνουσα γωνία μεταξύ του ερωτήματος και του εγγράφου
  - Διάταξη των εγγράφων σε αύξουσα τιμή του συνημιτόνου  $\text{cosine}(\text{query}, \text{document})$
- Θυμηθείτε ότι το συνημίτονα είναι γνησίως φθίνουσα συνάρτηση στο διάστημα  $[0^\circ, 180^\circ]$



# Κανονικοποίηση μήκους

- Ένα διάνυσμα μπορεί να κανονικοποιηθεί (ως προς το μήκος) διαιρώντας κάθε μια από τις συνιστώσες σου με το μήκος του (ως μήκος τους χρησιμοποιούμε την  $L_2$  νόρμα):

$$\|\vec{x}\|_2 = \sqrt{\sum_i x_i^2}$$

- Αυτή η διαίρεση οδηγεί σε μοναδιαίο διάνυσμα



# cosine(query, document)

$$\cos(\vec{q}, \vec{d}) = \frac{\text{Dot product}}{|\vec{q}| |\vec{d}|} = \frac{\text{Unit vectors}}{|\vec{q}|} \bullet \frac{\vec{d}}{|\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

$q_i$  είναι το tf-idf βάρος του όρου  $i$  στο ερώτημα

$d_i$  είναι το tf-idf βάρος του όρου  $i$  στο έγγραφο

$\cos(\vec{q}, \vec{d})$  είναι η cosine ομοιότητα των  $\vec{q}$  και  $\vec{d}$  ... ή, ισοδύναμα, το cosine της γωνίας μεταξύ των  $\vec{q}$  και  $\vec{d}$



# Cosine για length-normalized διανύσματα

- Για length-normalized διανύσματα, η ομοιότητα cosine είναι απλά το εσωτερικό γινόμενο (dot product) ή βαθμωτό γινόμενο (scalar product):

$$\cos(\vec{q}, \vec{d}) = \vec{q} \bullet \vec{d} = \sum_{i=1}^{|V|} q_i d_i$$

για  $q, d$  length-normalized.



## Υπολογισμός των cosine scores

COSINESCORE( $q$ )

```
1  float Scores[ $N$ ] = 0
2  float Length[ $N$ ]
3  for each query term  $t$ 
4  do calculate  $w_{t,q}$  and fetch postings list for  $t$ 
5      for each pair( $d, tf_{t,d}$ ) in postings list
6      do  $Scores[d] + = w_{t,d} \times w_{t,q}$ 
7  Read the array Length
8  for each  $d$ 
9  do  $Scores[d] = Scores[d] / Length[d]$ 
10 return Top  $K$  components of Scores[]
```

# Το σχήμα tf-idf έχει πολλές παραλλαγές

Term frequency		Document frequency		Normalization	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_t(tf_{t,d})}$	p (prob idf)	$\max\{0, \log \frac{N - df_t}{df_t}\}$	u (pivoted unique)	$1/u$
b (boolean)	$\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			b (byte size)	$1/CharLength^\alpha$ , $\alpha < 1$
L (log ave)	$\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$				

Οι στήλες με τίτλο ‘n’ είναι ακρωνύμια για σχήματα βάρους



# Στενωποί

- Η κύρια υπολογιστική στενωπός στον υπολογισμό των scores: cosine computation
- Μπορούμε ν' αποφύγουμε αυτόν τον υπολογισμό;
- Ναι, αλλά μερικές φορές σφάλουμε:
  - Ένα doc που *δεν είναι* στα top  $K$  μπορεί να εισδύσει στην λίστα των  $K$  output docs
  - Προκαλεί αυτό δυσχέρειες;



# Η ομοιότητα cosine είναι απλά μια προσέγγιση

- Ο κάθε χρήστης έχει μια δραστηριότητα και μια διατύπωση ερωτήματος
- Το cosine ταιριάζει docs στο ερώτημα
- Συνεπώς, το cosine είναι απλά μια προσέγγιση στην ικανοποίηση που λαμβάνει ο χρήστης
- Εάν τού επιστραφεί (του χρήστη) μια λίστα με  $K$  docs “κοντά” στα πραγματικά top  $K$  (τα οποία θα επέστρεφε η μετρική cosine), τότε δεν θα υπάρχει σημαντικό πρόβλημα





# Η γενική προσέγγιση

- Εύρεση ενός συνόλου  $A$  αποτελούμενο από *contenders*, με  $K < |A| \ll N$ 
  - Το  $A$  δεν περιέχει αναγκαστικά τα top  $K$ , αλλά περιέχει πολλά docs από αυτά του top  $K$
  - Επιστροφή των top  $K$  docs του  $A$
- Σκεφτείτε το  $A$  ως pruning των non-contenders
- Η ίδια προσέγγιση μπορεί να χρησιμοποιηθεί και για άλλες (non-cosine) συναρτήσεις score
- Θα εξετάσουμε διάφορα σχήματα τα οποία ακολουθούν αυτήν την προσέγγιση



## Μόνο οι query terms με υψηλό idf

- Για ένα ερώτημα όπως: *catcher in the rye*
- Συσσωρεύουμε μόνο τα score από τα *catcher* και *rye*
- Διαίσθηση: οι *in* και *the* συνεισφέρουν λίγο στα scores και έτσι δεν μεταβάλλουν την διάταξη σημαντικά
- Ώφελος:
  - Οι postings των of low-idf όρων έχουν πολλά → αυτά τα (πολλά) docs απαλείφονται από το σύνολο  $A$  των contenders



## Docs που περιέχουν πολλά query terms

- Κάθε doc με τουλάχιστον ένα query term είναι υποψήφιο για την top  $K$  output list
- Σε ερωτήματα με πολλαπλούς όρους, υπολογίζουμε τα scores των docs τα οποία περιέχουν “αρκετούς” από τους query terms
  - Π.χ., τουλάχιστον 3 από τους 4
  - Επιβάλλει ένα “soft conjunction” σε ερωτήματα που υποβάλλονται στις Web search engines
- Εύκολο στην υλοποίηση κατά την διάσχιση των postings



# Champion lists

- Προϋπολογίζουμε για κάθε dictionary term  $t$ , τα  $r$  docs με το υψηλότερο στην postings list του  $t$ 
  - Αποκαλούμε αυτό το σύνολο των  $r$  docs ως η champion list για τον  $t$
  - (ή και fancy list ή top docs για τον  $t$ )
- Παρατηρήστε ότι το  $r$  πρέπει να επιλεγεί κατά την διάρκεια κατασκευής του inverted index
  - Συνεπώς, είναι πιθανόν να ισχύει ότι:  $r < K$
- Όταν υποβληθεί ένα ερώτημα, υπολογίζουμε τα scores για τα docs στην champion list του query term
  - Επιλέγουμε τα  $K$  top-scoring docs μεταξύ αυτών



# Score στατικής ποιότητας

- Ζητούμε τα top-ranking documents να είναι ταυτόχρονα *relevant* και *authoritative*
- Η *relevance* μοντελοποιείται με τα cosine scores
- Η *authority* είναι τυπικά μια ιδιότητα που δεν εξαρτάται από το ερώτημα, δηλ., query-independent
- Παραδείγματα ποσοτικοποίησης της authority
  - Wikipedia μεταξύ των websites
  - Άρθρα σε κάποιες εφημερίδες
  - Ένα άρθρο με πολλά citations
  - Πολλά diggs, Y!buzzes ή del.icio.us marks
  - (Pagerank)

Ποσοτικά





# Μοντελοποίηση της αυθεντίας (authority)

- Ανάθεση σε κάθε έγγραφο  $d$  ένα *query-independent quality score* με τιμές στο διάστημα  $[0,1]$ 
  - Ας το συμβολίσουμε με  $g(d)$
- Επομένως, μια ποσότητα, όπως για παράδειγμα ο αριθμός των incoming links, την κανονικοποιούμε στο διάστημα  $[0,1]$ 
  - Άσκηση: Προτείνετε έναν μαθηματικό τύπο γι' αυτό



## Η έννοια του net score

- Θεωρήστε ένα απλό συνολικό score το οποίο συνδυάζει την cosine relevance και την αυθεντία (authority), δηλαδή:
- **$\text{net-score}(q,d) = g(d) + \text{cosine}(q,d)$** 
  - Φυσικά, μπορούμε να χρησιμοποιήσουμε κάποιον άλλον γραμμικό συνδυασμό και όχι την ίση συνεισφορά
  - Στην πραγματικότητα, κάθε συνάρτηση των δυο “σημάτων” που ποσοτικοποιούν την ικανοποίηση του χρήστη
- Συνεπώς, αναζητούμε τα top  $K$  docs με βάση το by net score



# Top $K$ με το net score – γρήγορες μέθοδοι

- Πρώτη ιδέα: Διάταξη όλων των postings με βάση το  $g(d)$
- Το κλειδί: αυτό είναι κοινή διάταξη για όλα τα postings
- Συνεπώς, μπορούμε ταυτόχρονα να διασχίσουμε τις postings λίστες των query terms για
  - Intersection των postings
  - Υπολογισμό του cosine score





## Γιατί διάταξη των postings ως προς $g(d)$ ;

- Με τη  $g(d)$ -ordering, τα top-scoring docs είναι πιθανό να εμφανίζονται νωρίς κατά την διάσχιση των postings λιστών
- Σε εφαρμογές με αυστηρές απαιτήσεις σε χρονική απόκριση (π.χ., απαιτείται να επιστρέψουμε οποιαδήποτε αποτελέσματα αναζήτηση βρήκαμε μέσα σε 50 ms), μάς επιτρέπει να σταματήσουμε την διάσχιση των postings σχετικά νωρίς



## Champion lists στο $g(d)$ -ordering

- Μπορούμε να συνδυάσουμε τις champion lists με το  $g(d)$ -ordering
- Διατηρούμε για κάθε όρο μια champion list των  $r$  docs με το υψηλότερο  $g(d) + tf-idf_{td}$
- Αναζητούμε τα top- $K$  αποτελέσματα μόνο από αυτές τις champion lists

# Η πλήρης εικόνα

