



# Ανάκτηση Πληροφορίας

## Ενότητα 8: Λανθάνουσα Σημασιολογική Ανάλυση (Latent Semantic Analysis)

Απόστολος Παπαδόπουλος  
Τμήμα Πληροφορικής



Ευρωπαϊκή Ένωση  
Ευρωπαϊκό Κοινωνικό Ταμείο



ΕΠΙΧΕΙΡΗΣΙΑΚΟ ΠΡΟΓΡΑΜΜΑ  
ΕΚΠΑΙΔΕΥΣΗ ΚΑΙ ΔΙΑ ΒΙΟΥ ΜΑΘΗΣΗ  
*επένδυση στην κοινωνία της γνώσης*

ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ  
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ

# Άδειες Χρήσης

- Το παρόν εκπαιδευτικό υλικό υπόκειται σε άδειες χρήσης Creative Commons.
- Για εκπαιδευτικό υλικό, όπως εικόνες, που υπόκειται σε άλλου τύπου άδειας χρήσης, η άδεια χρήσης αναφέρεται ρητώς.



# Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «Ανοικτά Ακαδημαϊκά Μαθήματα στο Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης» έχει χρηματοδοτήσει μόνο την αναδιαμόρφωση του εκπαιδευτικού υλικού.
- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.





# Λανθάνουσα Σημασιολογική Ανάλυση (Latent Semantic Analysis)



Ευρωπαϊκή Ένωση  
Ευρωπαϊκό Κοινωνικό Ταμείο



ΕΠΙΧΕΙΡΗΣΙΑΚΟ ΠΡΟΓΡΑΜΜΑ  
ΕΚΠΑΙΔΕΥΣΗ ΚΑΙ ΔΙΑ ΒΙΟΥ ΜΑΘΗΣΗ  
*επένδυση στην κοινωνία της γνώσης*

ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ  
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΣΠΑ  
2007-2013  
πρόγραμμα για την ανάπτυξη  
ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ

# Περιεχόμενα ενότητας

---

## 1. Λανθάνουσα Σημασιολογική Ανάλυση





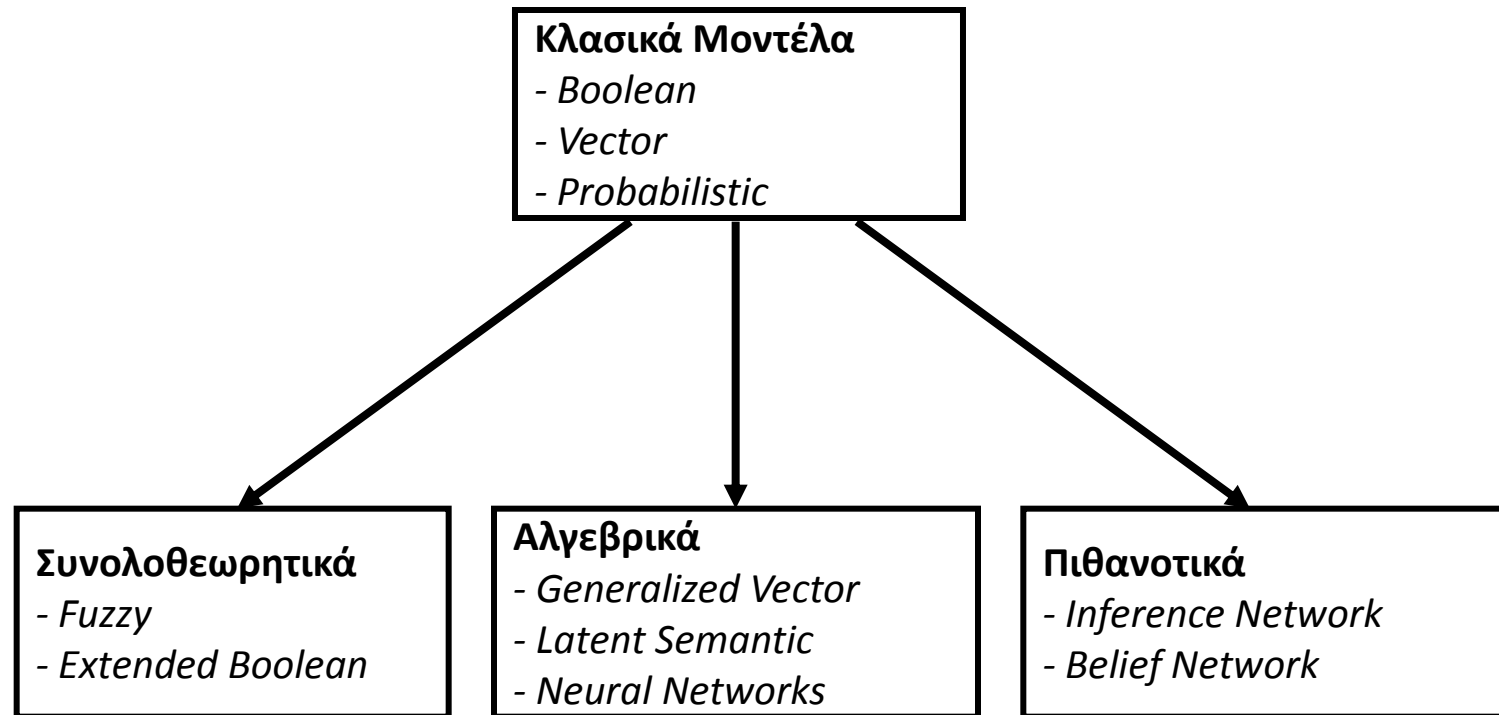
ΑΡΙΣΤΟΤΕΛΕΙΟ  
ΠΑΝΕΠΙΣΤΗΜΙΟ  
ΘΕΣΣΑΛΟΝΙΚΗΣ

---

Latent Semantic Analysis

# Λανθάνουσα Σημασιολογική Ανάλυση

# Μοντέλα IR



# Κίνητρο-1

Βασικό μειονέκτημα των μέχρι τώρα μοντέλων είναι ότι υποθέτουν την ανεξαρτησία των όρων (η πιθανότητα ένας όρος να ανήκει σε ένα έγγραφο δεν εξαρτάται από την ύπαρξη κάποιου άλλου όρου στο έγγραφο).

Επίσης, έχουμε δύο βασικά προβλήματα με τους όρους: **πολυσημία** (ένας όρος έχει διαφορετικό νόημα σε διαφορετικές περιοχές, π.χ. ποντίκι) και **συνωνυμία** (δύο όροι έχουν ακριβώς το ίδιο ή πολύ κοντινό νόημα, άστρο, αστέρας).





# Κίνητρο-2

Ποια η επίδραση της **πολυσημίας** και της **συνωνυμίας** στην **ανάκληση** και την **ακρίβεια**;

Η **πολυσημία** μειώνει την **ακρίβεια** και  
η **συνωνυμία** μειώνει την **ανάκληση**

Γιατί;



# Στόχος

- Να μπορέσουμε να βρούμε συσχετίσεις μεταξύ των όρων έτσι ώστε να μεταβούμε από τους όρους στις θεματικές περιοχές (concepts).
- Να μειώσουμε την επίδραση της πολυσημίας και της συνωνυμίας.
- Να μειώσουμε τον αριθμό των διαστάσεων.



# Singular Value Decomposition-1

- Η τεχνική LSI (Latent Semantic Indexing) βασίζεται στο μαθηματικό εργαλείο SVD (Singular Value Decomposition).



# Singular Value Decomposition-2

- Problem:
  - #1: Find concepts in text
  - #2: Reduce dimensionality

<b>term</b> <b>document</b>	<b>data</b>	<b>information</b>	<b>retrieval</b>	<b>brain</b>	<b>lung</b>
<b>CS-TR1</b>	1	1	1	0	0
<b>CS-TR2</b>	2	2	2	0	0
<b>CS-TR3</b>	1	1	1	0	0
<b>CS-TR4</b>	5	5	5	0	0
<b>MED-TR1</b>	0	0	0	2	2
<b>MED-TR2</b>	0	0	0	3	3
<b>MED-TR3</b>	0	0	0	1	1



# SVD - Definition

$$\mathbf{A}_{[n \times m]} = \mathbf{U}_{[n \times r]} \mathbf{\Lambda}_{[r \times r]} (\mathbf{V}_{[m \times r]})^T$$

- **A**:  $n \times m$  matrix (e.g.,  $n$  documents,  $m$  terms)
- **U**:  $n \times r$  matrix ( $n$  documents,  $r$  concepts)
- **$\Lambda$** :  $r \times r$  diagonal matrix (strength of each 'concept') ( $r$ : rank of the matrix)
- **V**:  $m \times r$  matrix ( $m$  terms,  $r$  concepts)



# SVD – Properties-1

**THEOREM** [Press+92]: always possible to decompose matrix  $\mathbf{A}$  into  $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T$ , where

- $\mathbf{U}, \mathbf{\Lambda}, \mathbf{V}$ : unique (\*)
- $\mathbf{U}, \mathbf{V}$ : column orthonormal (ie., columns are unit vectors, orthogonal to each other)
  - $\mathbf{U}^T \mathbf{U} = \mathbf{I}; \mathbf{V}^T \mathbf{V} = \mathbf{I}$  ( $\mathbf{I}$ : identity matrix)
- $\mathbf{\Lambda}$ : singular value are positive, and sorted in decreasing order



# SVD – Properties-2

‘spectral decomposition’ of the matrix:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} | & | \\ u_1 & u_2 \\ | & | \end{bmatrix} \times \begin{bmatrix} \lambda_1 & \emptyset \\ \emptyset & \lambda_2 \end{bmatrix} \times \begin{bmatrix} \text{---} & v_1 & \text{---} \\ \text{---} & v_2 & \text{---} \end{bmatrix}$$



# SVD - Interpretation

‘documents’, ‘terms’ and ‘concepts’:

- **U**: document-to-concept similarity matrix
- **V**: term-to-concept similarity matrix
- **$\Lambda$** : its diagonal elements: ‘strength’ of each concept

Projection:

- best axis to project on: (‘best’ = min sum of squares of projection errors)





# SVD – Example-1

- $A = U \Lambda V^T$  - example:

retrieval

inf. ↓

data    brain    lung

$$\begin{array}{c} \uparrow \\ \text{CS} \\ \downarrow \\ \uparrow \\ \text{MD} \\ \downarrow \end{array}
 \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}
 =
 \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}
 \times
 \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix}
 \times
 \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$



# SVD – Example-2

- $A = U \Lambda V^T$  - example:

**doc-to-concept similarity matrix**

↑ CS

↓

↑ MD

↓

		retrieval					
	inf.	↓					
	data		brain	lung			

**CS-concept**

↓

**MD-concept**

↙

X


X




# SVD – Example-3

- $A = U \Lambda V^T$  - example:

retrieval

		inf. ↓	brain	lung															
data																			

↑

CS

↓

↑

MD

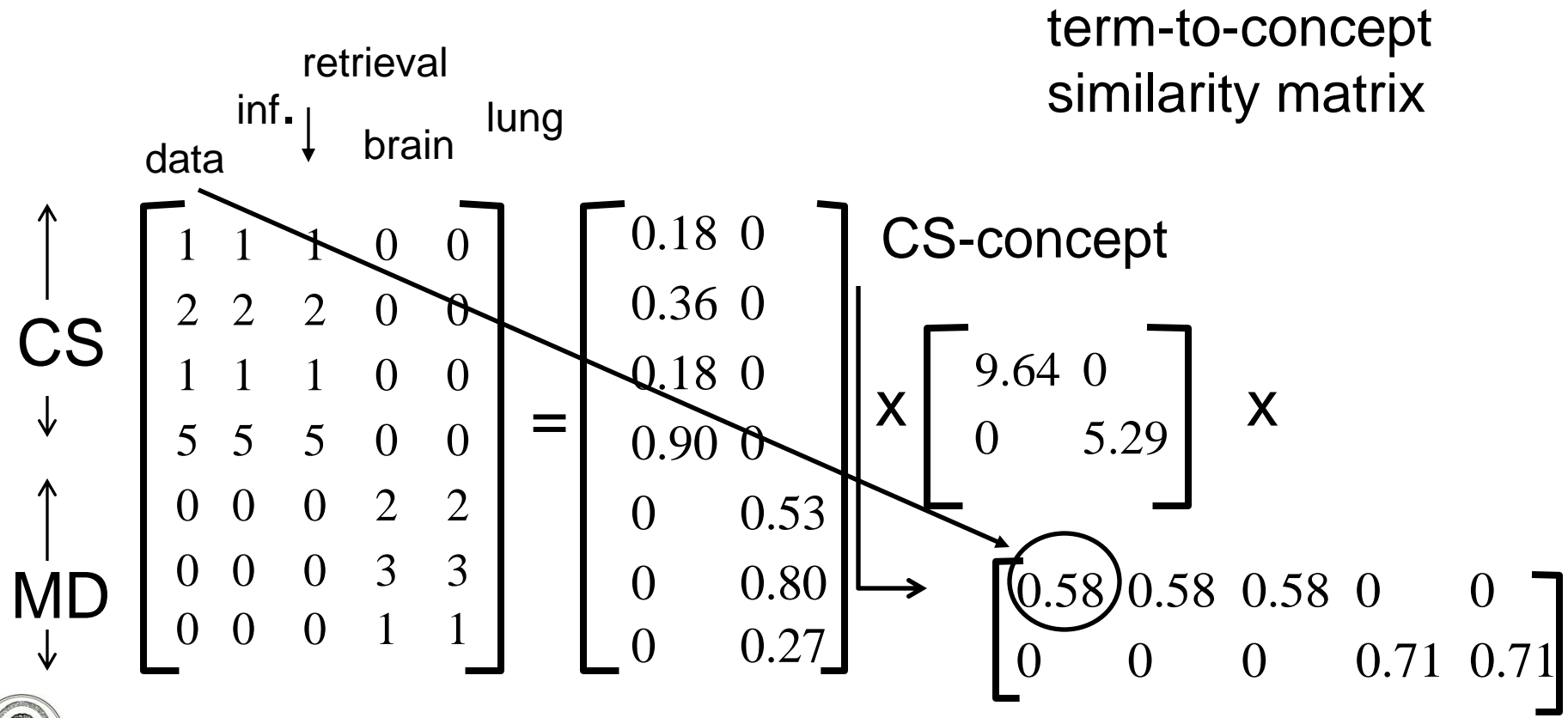
↓

$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}$	=	$\begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix}$	x	$\begin{bmatrix} 0 & 5.29 \\ 0 & 5.29 \end{bmatrix}$	x	$\begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$
---	---	--	---	--	---	---

‘strength’ of CS-concept

# SVD – Example-4

- $A = U \Lambda V^T$  - example:



# SVD – Dimensionality reduction-1

- Q: how exactly is dim. reduction done?
- A: set the smallest singular values to zero:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

The image shows the SVD decomposition of a matrix. The first matrix is a 7x5 matrix. The second matrix is a 7x2 matrix of singular values, with the second column crossed out. The third matrix is a 2x2 matrix of singular values, with the second value (5.29) crossed out. The fourth matrix is a 7x5 matrix of singular vectors, with the second and third columns crossed out.



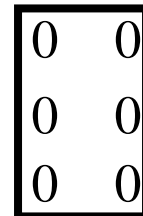
# SVD - Dimensionality reduction-2

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} \sim \begin{bmatrix} 0.18 \\ 0.36 \\ 0.18 \\ 0.90 \\ 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 9.64 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \end{bmatrix}$$



# SVD - Dimensionality reduction-3

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} \sim \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$



# LSI-1

Q1: How to do queries with LSI?

A: map query vectors into 'concept space' – how?

$$\begin{array}{c}
 \uparrow \\
 \text{CS} \\
 \downarrow \\
 \uparrow \\
 \text{MD} \\
 \downarrow
 \end{array}
 \begin{array}{c}
 \text{data} \\
 \text{inf.} \\
 \text{retrieval} \\
 \downarrow \\
 \text{brain} \\
 \text{lung}
 \end{array}
 \begin{bmatrix}
 1 & 1 & 1 & 0 & 0 \\
 2 & 2 & 2 & 0 & 0 \\
 1 & 1 & 1 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 0 & 0 & 2 & 2 \\
 0 & 0 & 0 & 3 & 3 \\
 0 & 0 & 0 & 1 & 1
 \end{bmatrix}
 =
 \begin{bmatrix}
 0.18 & 0 \\
 0.36 & 0 \\
 0.18 & 0 \\
 0.90 & 0 \\
 0 & 0.53 \\
 0 & 0.80 \\
 0 & 0.27
 \end{bmatrix}
 \times
 \begin{bmatrix}
 9.64 & 0 \\
 0 & 5.29
 \end{bmatrix}
 \times
 \begin{bmatrix}
 0.58 & 0.58 & 0.58 & 0 & 0 \\
 0 & 0 & 0 & 0.71 & 0.71
 \end{bmatrix}$$





# LSI-2

Q: How to do queries with LSI?

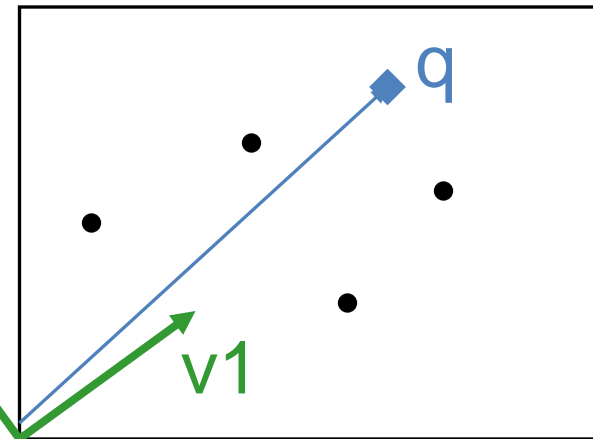
A: map query vectors into 'concept space' – how?

$$q = \begin{matrix} & \text{data} & \text{inf.} & \text{retrieval} & \text{brain} & \text{lung} \\ & & \downarrow & & & \\ \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

A: inner product  
(cosine similarity)  
with each 'concept' vector  $v_i$

term2

$v_2$



$v_1$

term1



# LSI-3

compactly, we have:

$$q_{\text{concept}} = q \mathbf{V}$$

e.g.:

$$q = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{matrix} \text{data} & \text{inf.} & \text{retrieval} & \text{brain} & \text{lung} \\ \downarrow & & & & \end{matrix} \begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix} \begin{matrix} \text{CS-concept} \\ \downarrow \\ \end{matrix} = \begin{bmatrix} 0.58 & 0 \end{bmatrix}$$

term-to-concept similarities



# Little example

How would the document ('information', 'retrieval') handled by LSI? A: SAME:

$$d_{\text{concept}} = d \mathbf{V}$$

Eg:

$$d = \begin{matrix} & \text{data} & \text{inf.} & \text{retrieval} & \text{brain} & \text{lung} \\ \begin{bmatrix} 0 & 1 & 1 & 0 & 0 \end{bmatrix} & & \downarrow & & & \end{matrix} \begin{matrix} \begin{bmatrix} 0.58 & 0 \\ 0.58 & 0 \\ 0.58 & 0 \\ 0 & 0.71 \\ 0 & 0.71 \end{bmatrix} \\ \text{term-to-concept} \\ \text{similarities} \end{matrix} \begin{matrix} \text{CS-concept} \\ \downarrow \\ = \begin{bmatrix} 1.16 & 0 \end{bmatrix} \end{matrix}$$



# Προσοχή

- Εάν στον αρχικό πίνακα  $A$  τα διανύσματα των εγγράφων είναι στις γραμμές του πίνακα τότε για τη μείωση των διαστάσεων χρησιμοποιείται ο μετασχηματισμός:

$$A' = A * V_k$$

- Εάν στον αρχικό πίνακα  $A$  τα διανύσματα των εγγράφων είναι στις στήλες του  $A$  τότε:

$$A' = U_k^T * A$$



# Σημείωμα Αναφοράς

Copyright Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης, Απόστολος Παπαδόπουλος. «Ανάκτηση πληροφορίας. Λανθάνουσα Σημασιολογική Ανάλυση (Latent Semantic Analysis)». Έκδοση: 1.0. Θεσσαλονίκη 2014.  
Διαθέσιμο από τη δικτυακή διεύθυνση:  
<http://eclass.auth.gr/courses/OCRS388/>



# Σημείωμα Αδειοδότησης

Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά - Μη Εμπορική Χρήση - Όχι Παράγωγα Έργα 4.0 [1] ή μεταγενέστερη, Διεθνής Έκδοση. Εξαιρούνται τα αυτοτελή έργα τρίτων π.χ. φωτογραφίες, διαγράμματα κ.λ.π., τα οποία εμπεριέχονται σε αυτό και τα οποία αναφέρονται μαζί με τους όρους χρήσης τους στο «Σημείωμα Χρήσης Έργων Τρίτων».



Ο δικαιούχος μπορεί να παρέχει στον αδειοδόχο ξεχωριστή άδεια να χρησιμοποιεί το έργο για εμπορική χρήση, εφόσον αυτό του ζητηθεί.

Ως **Μη Εμπορική** ορίζεται η χρήση:

- που δεν περιλαμβάνει άμεσο ή έμμεσο οικονομικό όφελος από την χρήση του έργου, για το διανομέα του έργου και αδειοδόχο
- που δεν περιλαμβάνει οικονομική συναλλαγή ως προϋπόθεση για τη χρήση ή πρόσβαση στο έργο
- που δεν προσπορίζει στο διανομέα του έργου και αδειοδόχο έμμεσο οικονομικό όφελος (π.χ. διαφημίσεις) από την προβολή του έργου σε διαδικτυακό τόπο

[1] <http://creativecommons.org/licenses/by-nc-nd/4.0/>





# Τέλος ενότητας

Επεξεργασία: <Μαυρίδης Απόστολος>  
Θεσσαλονίκη, <Εαρινό εξάμηνο 2013-2014>



Ευρωπαϊκή Ένωση  
Ευρωπαϊκό Κοινωνικό Ταμείο



ΕΠΙΧΕΙΡΗΣΙΑΚΟ ΠΡΟΓΡΑΜΜΑ  
ΕΚΠΑΙΔΕΥΣΗ ΚΑΙ ΔΙΑ ΒΙΟΥ ΜΑΘΗΣΗ  
*επένδυση στην κοινωνία της γνώσης*  
ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ  
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΣΠΑ  
2007-2013  
πρόγραμμα για την ανάπτυξη  
ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ



ΑΡΙΣΤΟΤΕΛΕΙΟ  
ΠΑΝΕΠΙΣΤΗΜΙΟ  
ΘΕΣΣΑΛΟΝΙΚΗΣ

---

# Σημειώματα



# Διατήρηση Σημειωμάτων

Οποιαδήποτε αναπαραγωγή ή διασκευή του υλικού θα πρέπει να συμπεριλαμβάνει:

- το Σημείωμα Αναφοράς
- το Σημείωμα Αδειοδότησης
- τη δήλωση Διατήρησης Σημειωμάτων
- το Σημείωμα Χρήσης Έργων Τρίτων (εφόσον υπάρχει)

μαζί με τους συνοδευόμενους υπερσυνδέσμους.

