

ΜΥΕ003: Ανάκτηση Πληροφορίας

Διδάσκουσα: Ευαγγελία Πιτουρά

Εισαγωγή

Ακαδημαϊκό Έτος 2023-2024

Ανάκτηση Πληροφορίας (Information Retrieval) - (IR)

είναι η εύρεση αντικειμένων (κυρίως εγγράφων) από μεγάλες συλλογές τα οποία ικανοποιούν μια ανάγκη για πληροφόρηση

Συνήθως αναζήτηση με λέξεις κλειδιά

Εφαρμογές;

Γιατί να μας ενδιαφέρει;

Παλιότερα,
Βιβλιοθηκονόμους, βοηθούς νομικών
επαγγελματιών κλπ;

ISBN: 0-201-12227-8

Author: Salton, Gerard

Title: Automatic text processing: the transformation, analysis,
and retrieval of information by computer

Editor: Addison-Wesley

Date: 1989

Content: <Text>

external attributes (metadata) and **internal attribute** (content)

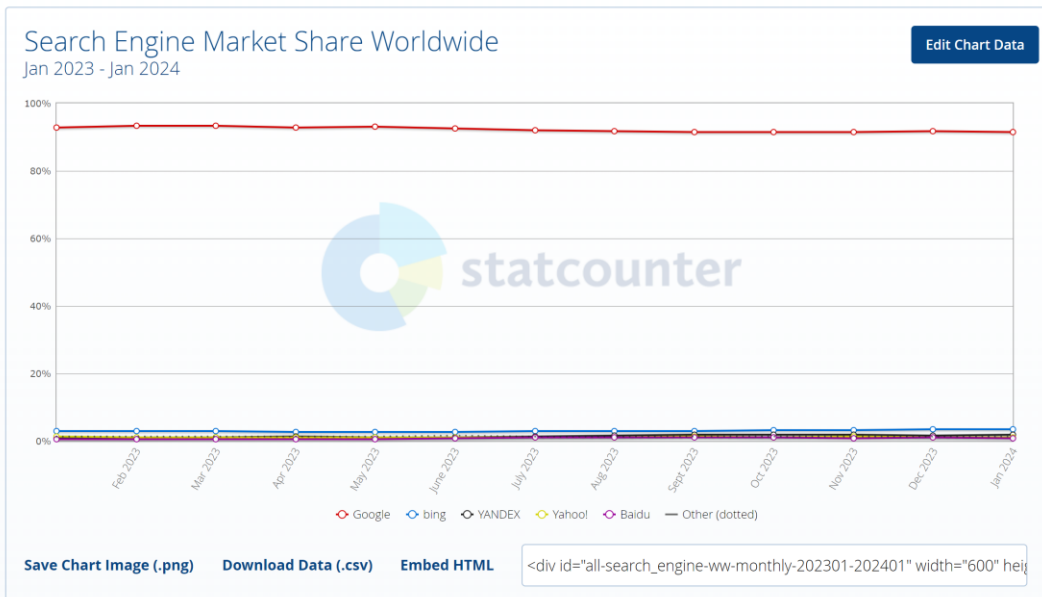
Search by external attributes = Search in DB

IR: search by content



Εφαρμογές

Search Engines



Last 12 months, worldwide, all platforms

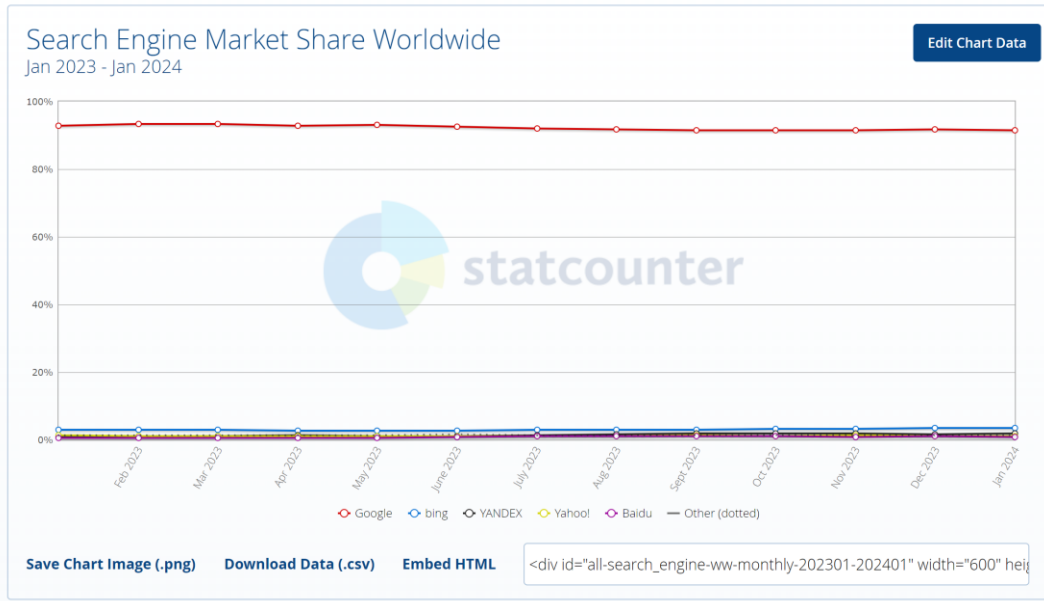
<https://gs.statcounter.com/search-engine-market-share>

Google: 3.5 billion searches per day

- **Bing:** Microsoft
Build on previous search engines (MSN Search, Windows Live Search, Live Search),
In 2023, faster than Google to release an AI Chatbot service (Bing Chat) based on GPT4, but despite initial acceptance, market share remained at low levels.
- **Yandex:** ρωσική πολυεθνική τεχνολογική εταιρεία που ειδικεύεται σε υπηρεσίες και προϊόντα σχετικά με το Διαδίκτυο και θεωρείται η μεγαλύτερη εταιρεία τεχνολογίας της *Ρωσίας*. Λειτουργεί την μεγαλύτερη μηχανή αναζήτησης στη Ρωσία με μερίδιο αγοράς περίπου 65%.

Πηγή: Wikipedia

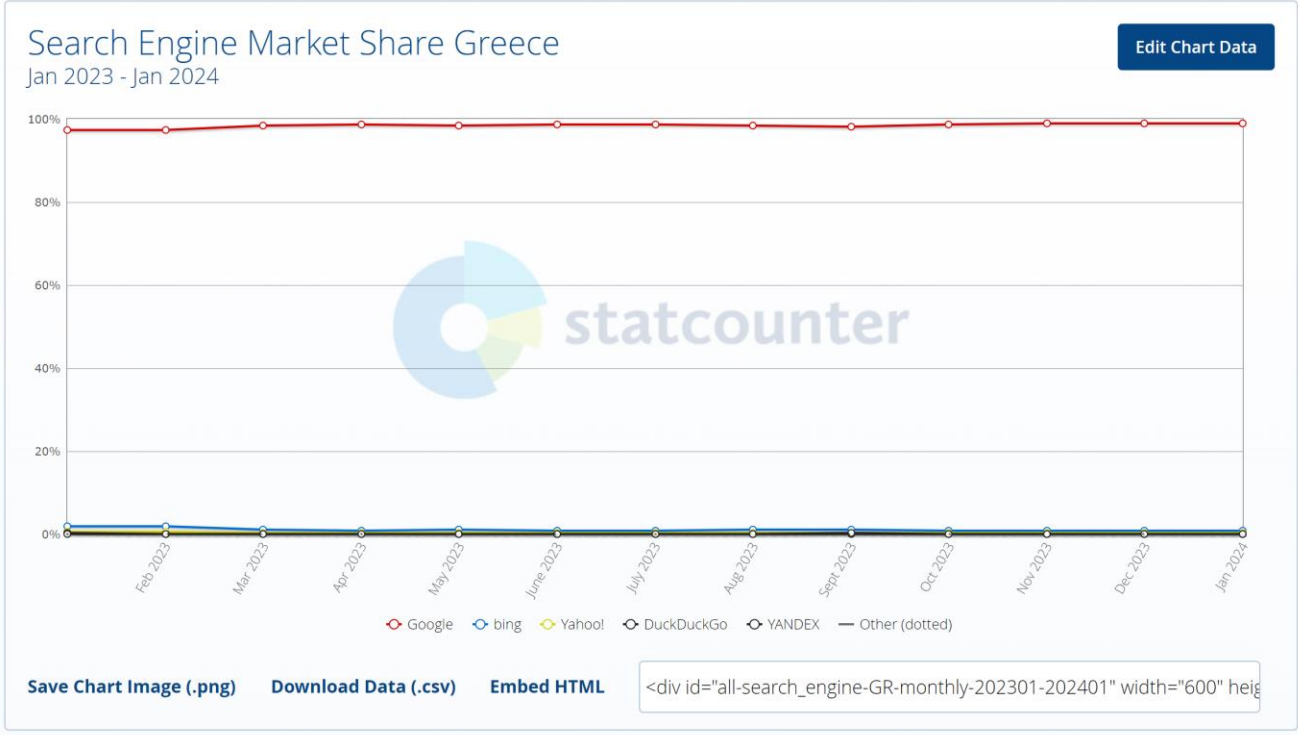
Search Engines



- **Baidu:** 2η μεγαλύτερη μηχανή αναζήτησης στον κόσμο, και κατέχει το 76.05% του μεριδίου αγοράς στην αγορά μηχανών αναζήτησης της *Κίνας*.
- **DuckDuckGo** είναι διαδικτυακή μηχανή αναζήτησης που δίνει έμφαση στην *προστασία της ιδιωτικής ζωής* των χρηστών της και στην αποφυγή του “*φίλτρου φυσαλίδας*” των εξατομικευμένων αποτελεσμάτων αναζήτησης.
 - Πηγή: *Wikipedia*

Last 12 months, worldwide, all platforms

<https://gs.statcounter.com/search-engine-market-share>



Petal search:
Developed by Huawei

Desktop search



IR tools are designed to find information on a PC, including web browser history, e-mail archives, text documents, images, video, etc

Εφαρμογές

Email search

Social search

Enterprise search

helps people in an organization find the information they need to perform their jobs-- data extracted from inside the business, along with external data sources like document management systems, databases, etc

Domain specific search: Legal information retrieval,
Digital libraries

Διαφορετικές απαιτήσεις ανάλογα με την εφαρμογή

Κατηγορίες εφαρμογών

- Στο web/διαδίκτυο (search engines)
Δισεκατομμύρια έγγραφα σε εκατομμύρια υπολογιστές.
Συλλογή εγγράφων, κλίμακα, διάταξη αποτελεσμάτων, ..
- Προσωπική ανάκτηση πληροφορίας
(στον προσωπικό υπολογιστή, email, κλπ)
Διαφορετικά είδη αρχείων, light-weight, maintenance-free, ...
- Σε επίπεδο επιχείρησης, οργανισμού
(enterprise, institutional)
- Αναζήτηση ειδικού σκοπού (domain-specific search) – πχ ερευνητικά άρθρα σε βιοχημεία

Ορισμός

Ανάκτηση Πληροφορίας (**Information Retrieval**) - (IR)

- είναι η εύρεση αντικειμένων κυρίως εγγράφων (**documents**) αδόμητης φύσης (*) (**unstructured**) που συνήθως έχουν τη μορφή κειμένου (**text**)
- από μεγάλες συλλογές (συνήθως αποθηκευμένες σε υπολογιστές)
- τα οποία ικανοποιούν μια ανάγκη πληροφόρησης (**information need**)

() όχι ακριβώς!*

Αδόμητα δεδομένα

- Τυπικά αναφέρεται σε *ελεύθερο κείμενο*
- Επιτρέπει
 - Ερωτήματα με *λέξεις κλειδιά* (keyword) με πιθανούς τελεστές
 - Ποιο περίπλοκες ερωτήσεις για *έννοιες*: π.χ.,
 - Βρες όλες τις web σελίδες για την απελευθέρωση των Ιωαννίνων
- Κλασσικό μοντέλο για αναζήτηση σε έγγραφα κειμένου

structured

	A	B	C	D	E	F	G
1	Purchase ID	Last name	First name	Birth day	Country	Date of purchase	Amount of purchase
2	1	Davidson	Michael	04/03/1986	United States	10/12/2016	37
3	2	Vito	Jim	09/01/1994	United Kingdom	02/02/2016	85
4	3	Johnson	Tom	23/08/1972	France	02/11/2016	83
5	4	Lewis	Peter	18/10/1979	Germany	22/11/2016	27
6	5	Koenig	Edward	13/05/1983	Argentina	26/03/2015	43
7	6	Preston	Jack	16/06/1991	United States	06/11/2016	77
8	7	Smith	David	11/03/1965	Canada	15/11/2016	23
9	8	Brown	Luis	03/09/1997	Australia	03/07/2015	74
10	9	Miller	Thomas	07/01/1980	Germany	07/11/2016	13
11	10	Williams	Bill	26/07/1960	United States	20/11/2015	80
12	11	Gemini	Alexia	12/09/1995	Canada	11/03/2017	35
13	12	Bond	James	25/02/1975	United Kingdom	12/08/2017	40
14	13	Burgle	Patricia	01/12/1990	United States	18/01/2015	55
15	14	Reding	Michelle	07/04/1985	Canada	23/02/2017	28
16	15	Harvey	Billy	14/07/1971	United Kingdom	12/01/2016	41
17							

unstructured

Introducing one of Australia's greatest treks and one of our newest trips! 🏔️

The Cradle Mountain Overland track offers some of Tasmania's most stunning scenery – dramatic valleys, temperate rainforests, beautiful lakes and more. And this 6-day camping trip is the perfect way to experience it, alongside like-minded adventurers and experienced guides.

Why travel there with us? ... See More

View Similar Products

116

11 Comments 12 Shares

Jessica Chapman
sigh I miss travel. With COVID travel restrictions here it's going to be a while.
Like · Reply · 1w
↳ 1 Reply

Abderrahman Chafiq
snow in the Atlas mountains from Morocco 🇲🇴🇲🇴🇲🇴🇲🇴

Dian Clayton
Loved it, but snow on Mt Osser in November! 😊
Like · Reply · 1w

Hany Sayed
Beautiful 🌟
Like · Reply · 1w

Kelly McCarthy
Cradle Mountain is one of my favorite places. So many wild wombats to observe!!
Like · Reply · 1w

Most Relevant is selected, so some comments may have been filtered out.

Ημιδομημένα δεδομένα

- Στην πραγματικότητα, δεν υπάρχουν αμιγώς μη δομημένα δεδομένα
 - π.χ., αυτή η διαφάνεια έχει διακριτές ζώνες όπως *Title* και *Bullets*
 - *Web pages?*
 - *Emails?*
- «Ημιδομημένη» αναζήτηση όπως:
 - *Title* contains ημιδομημένα AND *Bullets* contain αναζήτηση

... και βέβαια υπάρχει πάντα η γλωσσική δομή

Τι είναι η Ανάκτηση Πληροφορίας (Information Retrieval);

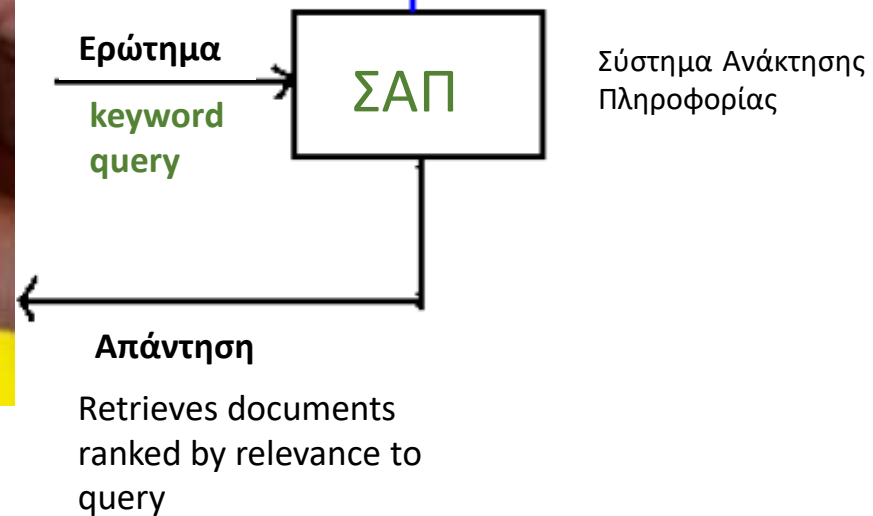
information need

Ανάγκη
πληροφόρησης



corpus/collection

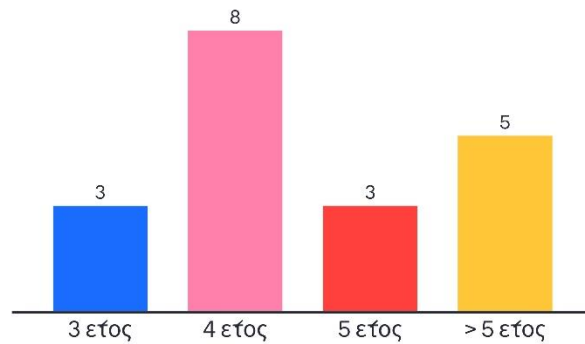
Συλλογή
Εγγράφων



Ρολλ 1: Ποιοι είστε;

Multiple Choice

Mentimeter



6 13 6



LLMs (ChatGPT) και Ανάκτηση Πληροφορία

ChatGPT

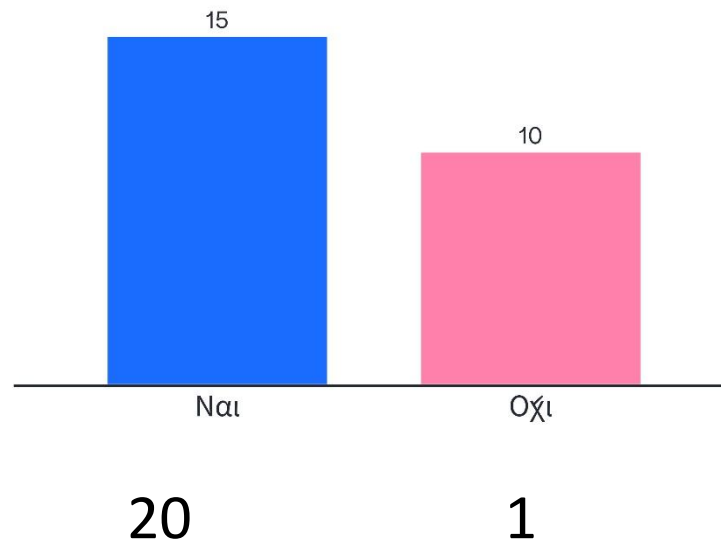
Poll 2

Το έχετε χρησιμοποιήσει;

ChatGPT

Έχετε χρησιμοποιήσει το ChatGPT;

Mentimeter



ChatGPT

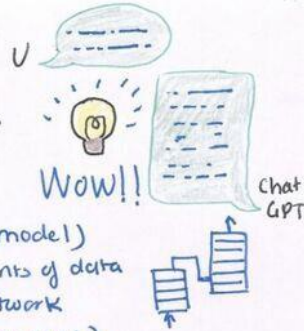
Τι είναι σε μια πρόταση



Openai.com

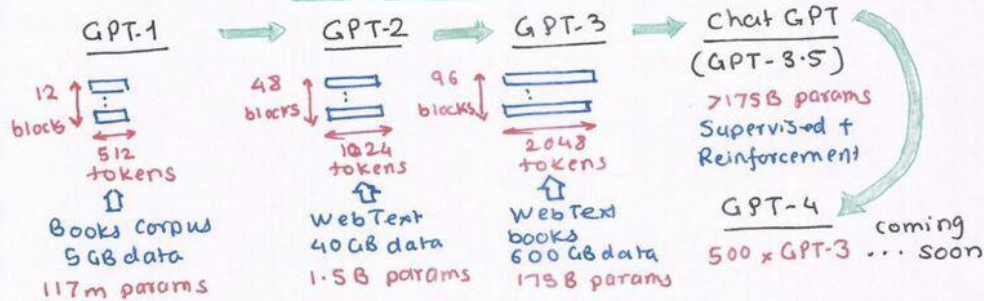
CHATGPT

* Generative Pre-trained Transformers

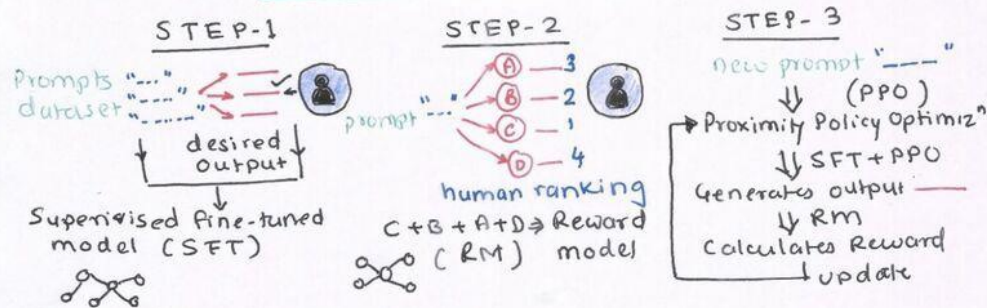


- Generative**: predicting next word (Language model)
- Pre-trained**: previously trained on large amounts of data
- Transformer**: Encoder-Decoder based neural network
- Chat GPT**: GPT fine tuned for conversations (chatbot)

GPT Progression



How ChatGPT was trained?



Advantages

- human like writing
- complex instructions
- mix-match ideas

Dis-advantages

- hallucinate at times
- not truly creative
- expensive \$\$\$

Applications

- writing blogs, emails etc
- search descriptive info
- idea generation

Contact
 📞 91-9890251406
 ✉️ yogeshkulkarni@yahoo.com

Question
 Will chatGPT pass the Turing Test?

ChatGPT

Poll 2

Μέχρι 3 πιο θετικά

Μέχρι 3 πιο αρνητικά

ChatGPT

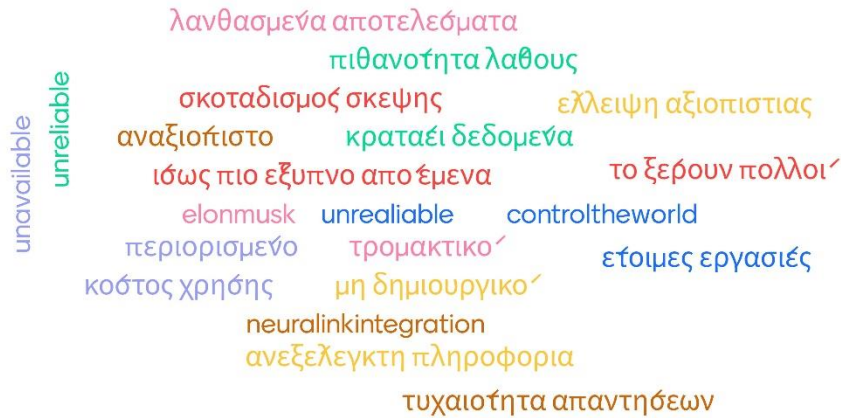
Μέχρι 3 πιο θετικά

Mentimeter



ChatGPT

Μέχρι 3 πιο αρνητικά



Συζήτηση

Λάθη

Λογοκλοπή (οικειοποιείται κείμενο άλλων)

ChatGPT και ΜΥΕ003

Δύο οπτικές γωνίες/“disruptive” σε δύο άξονες

- Εκπαίδευση γενικά
- Στο αντικείμενο του μαθήματος: σχέση με αναζήτηση

ChatGPT και εκπαίδευση

Ως πηγή γνώσης (educational resource)

Ιδέες:

ChatGPT

13/2/23, 12:43 μ.μ.

Information Retrieval (IR)



New chat



What is information retrieval?



Information retrieval (IR) is the process of retrieving relevant and useful information from a collection of data or a database, in response to a query or search request made by a user. The goal of IR is to provide the most relevant and accurate information to the user in a timely manner, out of the potentially large and complex collection of information available. IR systems typically use a combination of algorithms, data structures, and indexing techniques to search and rank the information in the database, based on factors such as relevance, popularity, or user preference. IR is widely used in many applications, such as web search engines, digital libraries, enterprise search systems, and recommendation systems, among others.



ChatGPT και εκπαίδευση

Επίδραση στην αξιολόγηση των μαθητών/φοιτητών

Προσεγγίσεις:

1. Απαγόρευση

Πως:

Άλλου τύπου αξιολόγηση, αποκλεισμός πρόσβασης

2. Όχι απαγόρευση

- Χρήση με αναφορά
- Χρήση χωρίς αναφορά

The three models of satisfying information needs

	Structures queries	Search Engines	Large Language Models
Systems	DBMS	Google search	ChatGPT
Corpus	Relational tables	Web+	Web texts, books, Wikipedia
Input	SQL	Keywords	Natural language
Output	Data in tables	Ranked sources	New content

- Ο “κόσμος” παράγει περισσότερο από **2 exabytes** νέας πληροφορίας το χρόνο, 90% της οποίας είναι σε ψηφιακή μορφή και με 50% ετήσια αύξηση
- Τι θα γίνει με το κείμενο που παράγει το ChatGPT και γενικά από AI tools
 - Θα γίνει index?

Retrieval augmented generation (RAG)

- Rather than relying on a fixed LLM to deliver the answer to a query, if we **first find relevant documents** (online or elsewhere) and **then use an LLM to process the query and the documents** into an answer, this could provide an alternative to current web search.
- Executing this efficiently and at scale would be complex, but the effect would be akin to having an LLM do a web search and summarize the results.

Πίσω στο ΜΥΕ003: τι θα μάθουμε;

Αναζήτηση

- Μοντελοποίηση
- Επεξεργασία (φυσικές γλώσσες)
- Βασικές δομές
- Αξιολόγηση

Μηχανές Αναζήτησης

Δημιουργία της μηχανής αναζήτησης

- Δημιουργία συλλογής (αν δεν υπάρχει):
 - crawl, scrap, use specific APIs (social media)
- Ανάλυση του κειμένου
 - Ποιοι θα είναι οι όροι, περιλαμβάνει και γλωσσολογική επεξεργασία
- Κατασκευή ευρετηρίων
 - Η πιο απλή μορφή: για κάθε όρο σε ποια έγγραφα εμφανίζεται
(επεκτάσεις πχ με συχνότητα εμφάνισης, πληροφορία θέσης, συμπίεση, κλπ)

Μηχανές Αναζήτησης

Λειτουργία

Ο χρήστης υποβάλει ένα ερώτημα

Επεξεργασία του ερωτήματος

Χρήση του ευρετηρίου για να βρούμε (retrieve) τα συναφή έγγραφα και να τα **διατάξουμε** με βάση τη **συνάφεια**

Πως ορίζουμε τη συνάφεια:

- Boolean model: υπάρχει, δεν υπάρχει ο όρος
- tf-idf (βασισμένη σε κείμενο)
- Source importance: Pagerank (οι συνδέσεις), clickthrough (γενικά traffic), document age, authority (wikipedia)
- Personalization, contextualization
- User intent
- Learning to rank

Μηχανές Αναζήτησης

Αξιολόγηση

Δε μας ενδιαφέρει η διάταξη:

Recall/Precision/AUC

Μας ενδιαφέρει η διάταξη:

MAP, NDCG

Word embeddings + ML

Διανυσματικές αναπαραστάσεις όρων

word2vec (σε βάθος)

Transformers, LLMs (σε υψηλό επίπεδο)

Lucene

Open-source software για τη δημιουργία μηχανών αναζήτησης

Μηχανές Αναζήτησης

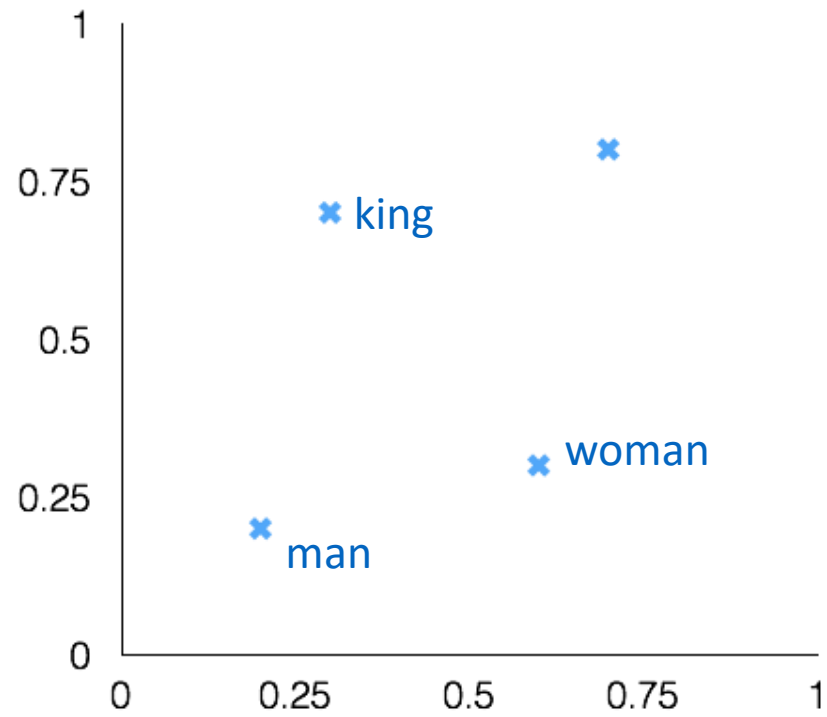
Lucene

Open-source software για τη δημιουργία μηχανών αναζήτησης

Μερικά Στοιχεία Μηχανικής Μάθησης

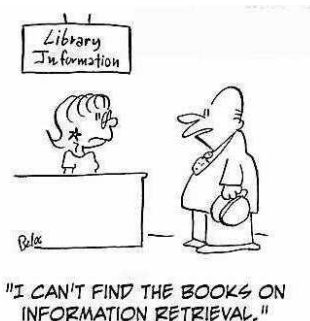
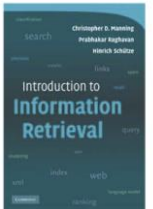
Word embeddings

+	king	[0.30 0.70]
-	man	[0.20 0.20]
+	woman	[0.60 0.30]
<hr/>		
	queen	[0.70 0.80]



Διαδικαστικά

- Ιστοσελίδα
- Βιβλίο
 - Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze. *Εισαγωγή στην Ανάκτηση Πληροφοριών*, Εκδόσεις Κλειδάριθμος
 - Η αγγλική έκδοση διαθέσιμη δωρεάν
- Ricardo Baeza-Yates and Berthier Ribeiro-Neto, *Ανάκτηση Πληροφορίας*, 2^η Έκδοση, Εκδόσεις Τζιόλα



Διαδικαστικά

- Βαθμολογία (μπορεί να αλλάξει):
 - Εργασία (έως 2 άτομα) – σε φάσεις: 50%
 - Προαιρετική εργασία – ανεξάρτητη μελέτη ML μοντέλων (20%) ή σύνδεση τους με τη μηχανή αναζήτησης
 - Τελικό Διαγώνισμα: 50% (αν όχι την προαιρετική εργασία)
30% (αν την προαιρετική εργασία)
- Η εργασία δεν «κρατιέται»
- Για να περάσετε το μάθημα, βαθμός διαγωνίσματος ≥ 4

Λίγα λόγια για την εργασία

Μηχανή αναζήτησης

Lucene (solar)

2 καταληκτικές ημερομηνίες

Θεματικό (πέρυσι αναζήτηση τραγουδιών (στίχοι)

ΜΥΕ003: Ανάκτηση Πληροφορίας

Διδάσκουσα: Ευαγγελία Πιτουρά

Εισαγωγή

Ακαδημαϊκό Έτος 2023-2024