

VISUAL INFORMATION RETRIEVAL

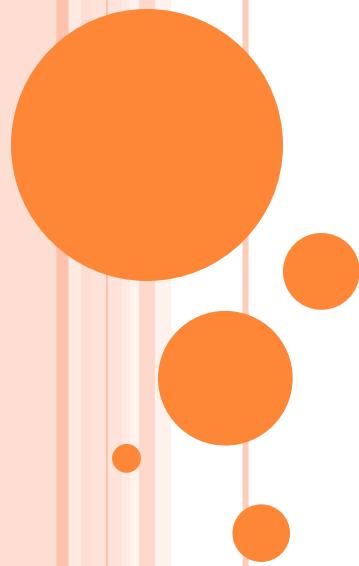
Gang Hua

Department of Computer Science
Stevens Institute of Technology

01/22/2014

Acknowledgement:

Thanks Cees Snoek for sharing his course slides!
A portion of the slides in this course is adapted from his slides.



CS598 VISUAL INFORMATION RETRIEVAL

Lecture I: Part I: Introduction to the course

WHO AM I?



□ Prof. Gang Hua

华冈

- Associate Professor of Computer Science
- Stevens Institute of Technology
- Research Staff Member (07/2010—08/2011)
● IBM T J. Watson Research Center
- Senior Researcher (08/2009—07/2010)
● Nokia Research Center Hollywood
- Scientist (07/2006—08/2009)
● Microsoft Live Labs Research
- Ph.D. in ECE, Northwestern University, 06/2006

Now, WHO ARE YOU?

- Why do you choose this class?
-
- What do you expect to learn?

- Previous experience?



COURSE PROMISE

- You will learn the theory and practice of visual information retrieval.
- You will be able to recall the major scientific problems in visual information retrieval.
- You will be able to understand and explain how state-of-the-art visual search systems work.



COURSE INFORMATION (1)

- **CS598** Visual Information Retrieval
- **Term:** Spring 2014
- **Instructor:** Prof. Gang Hua
- **Class time:** Monday 6:15pm—8:40pm
- **Location:** McLean Chemical Sciences Building 414
- **Office Hour:** Wednesday 4:00pm—5:00pm by appointment
- **Office:** Lieb/Room 305
- **Course Assistant:** Haoxiang Li
- **Course Website:**
<http://www.cs.stevens.edu/~ghua/ghweb/Teaching/CS598Spring2014.htm>

COURSE INFORMATION (2)

- **Text Book:**

- No required text book

- Recommended Reading

- See course website for recommended

- **Grading:**

- Class participation: 10%

- One written homework: 10%

- 4 course project:

- Project #1(10%), Project #2(10%), Project #3 (15%), Project #4 (15%)

- Final Project & Presentation: (15% System demo and presentation, 15% final report)

- Final grade: A-- 90% to 100%; B--80% to 89%; C-- 65% to 79%; F -- < 65% .



COURSE INFORMATION (3)

- Projects are all group project
 - You will be grouped into a team with 5-6 members
- Your final project will be integrating your first four project into an end-to-end image search system
- So it is important that you don't lag behind!!
- Program language is not specified, pick up your favorite one with the team!



RULES

- **Need to be absent from class?**
 - 1 point per class: please send notification and justification at least 2 days before the class
- **Late submission of homework?**
 - The maximum grade you can get from your late homework decreases 50% per day
- **Zero tolerance on plagiarism!!**
 - The first time you receive zero grade for the assignment
 - The second time you get “F” in your final grade
 - Refer to Stevens honor system for your behavior



PREREQUISITES

- No definite pre-requisite
- However, it would be helpful if you possess
 - Basic knowledge of computer vision, pattern recognition and information retrieval
 - e.g.,
 - CS558 Computer Vision
 - CS559 Machine Learning
- We will learning something about information retrieval today



WHAT IS A VISUAL INFORMATION RETRIEVAL?

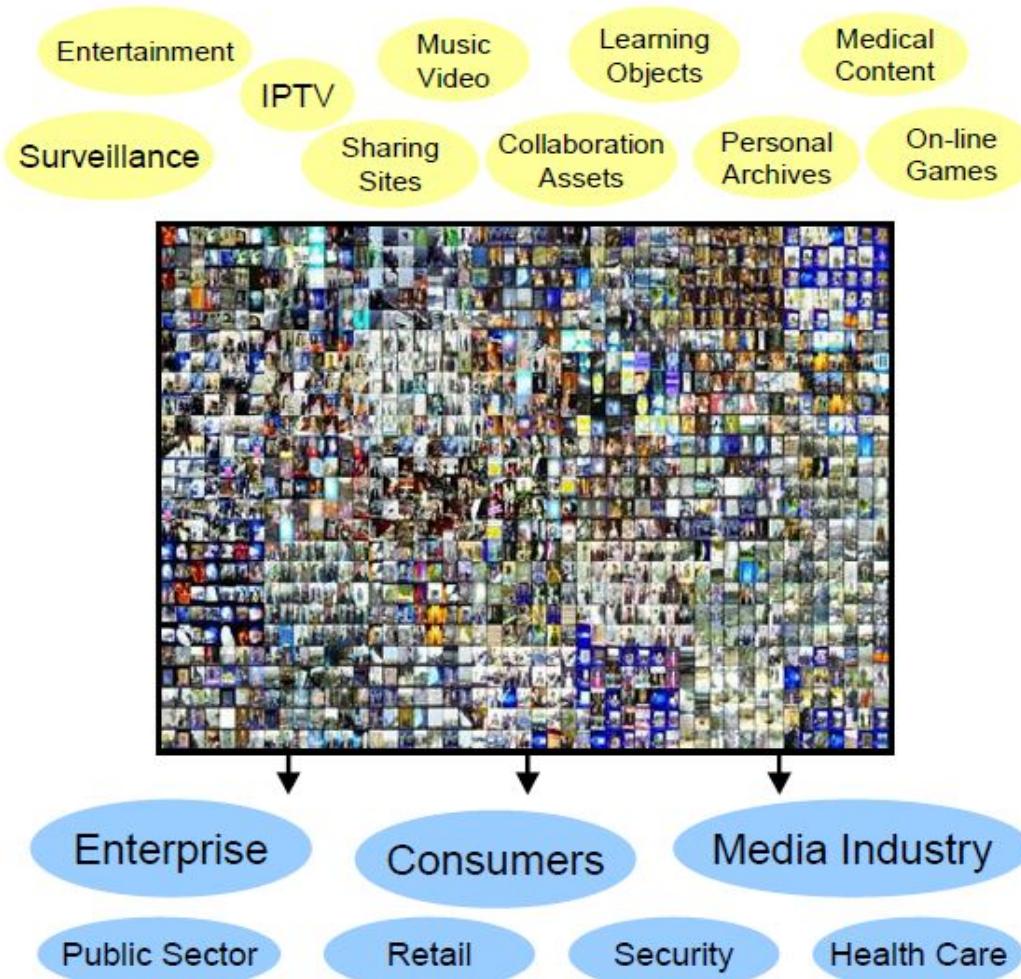
□ Definition

- Visual information retrieval is the activity of obtaining (*i.e.*, **search and retrieve**) visual information resources relevant to an information need (*i.e.*, **a query**) from a collection of visual information resources (*e.g.*, **image and video database**).



WHY NEEDED?

- Growing deluge requires more effective solutions for organizing, managing & searching video content
- Manual indexing is costly, time-consuming and inadequate
- New technologies are needed to automate processing and unlock value of large repositories
- Metadata standards are needed to support interoperable search



WHY BOTHER?

> \$100 million

Women's Handbags - patent At Like.com - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://www.like.com/womenshandbags-search--patent--style-580002

★ Women's Handbags - patent At Like... +

like.com shop visually

patent search

Women's Handbags

ALL CATEGORIES ▾

Refine By Price Brand Color leather hobo X Material Site

On Sale Free Shipping View

ZYNC FROM AMERICAN EXPRESSSM
CUSTOMIZE THE CARD TO FIT YOUR LIFE
LEARN MORE

Women's Bags Designer Bags Leather Bags Totes Hobo Clutch Wristlet Shoulder Bags Evening Bags Backpacks Wallets LIKEMAG

Michael By Michael K. Derek Alexander Leal Hobo Michael By Michael K. Hobo International Julia Cocco Hobo Coach Jimmy Choo Hobo

SHOP \$449.00 \$358.40 SHOP \$149.00 \$79.99 eBags SHOP \$98.00 Lordandtaylor SHOP \$448.00 Saks Fifth Avenue SHOP \$159.00 \$148.00 Endless SHOP \$425.00 Forzieri SHOP \$99.00 Lordandtaylor SHOP \$298.00 coach SHOP \$1,295 Nordstrom

FREE SHIPPING 48% OFF FREE SHIPPING FREE SHIPPING FREE SHIPPING

Visual Search Visual Search

Guess Michael By Michael K. Hobo International Bravo Coach Hobo International Coach Jimmy Choo Michael By Michael K.

SHOP \$101.25 Forzieri SHOP \$448.00 Zappos SHOP \$228.00 Endless SHOP \$430.00 \$238.95 Shoebuy SHOP \$298.00 coach SHOP \$228.00 Endless SHOP \$298.00 coach SHOP \$1,350 Nordstrom SHOP \$248.00 Bloomingdales

FREE SHIPPING FREE SHIPPING FREE SHIPPING

Visual Search Visual Search

Complete the Look

Hobo Hobo

Shop \$295.00 Lordandtaylor \$94.99

Complete the Look

Hobo Hobo

Shop \$149.00 Lordandtaylor \$44.00

Done

One active download (2 minutes, 30 seconds remaining)

Slide credit: C. Snoek (Univ. of Amsterdam)

\$1.65 billion

STILL NOT CONVINCED?

The screenshot shows a Mozilla Firefox browser window with the title bar "YouTube - A Message From Chad & Steve...The Real Story of YouTube - Mozilla Firefox". The menu bar includes "Bestand", "Bewerken", "Beeld", "Ga", "Bladwijzers", "Extra", and "Help". The address bar shows the URL "http://www.youtube.com/watch?v=xtr1UX3Ygc". The main content area displays the YouTube homepage with the "Broadcast Yourself" slogan. Navigation tabs include "Home", "Videos", "Channels", "Groups", "Categories", and "Upload". Below these are links for "Most Recent", "Most Viewed", "Top Rated", "Most Discussed", "Top Favorites", "Most Linked", and "Recently Featured". Promotional banners for "Gratis cabaret DVD?", "1001 Speakers Secrets", "Knee Pads, P&P free (UK)", and "More than 50.000 pages" are visible. A video player on the left shows a video of two men, Chad and Steve, with the title "A Message From Chad & Steve...The Real Story of YouTube". The video player interface includes a play button, a progress bar at 00:18 / 00:46, and a view count of 444. To the right of the video player is the video's metadata: "Added October 11, 2006", "From fff123", "SUBSCRIBE", and a link to "to fff123". The video summary reads: "How do they really feel about selling ou ... (more)". The category is "Entertainment", and the tags are "Urban", "Rush", "Michael", "Eckford" (more). The URL is "http://www.youtube.com/watch?v=xtr1UX3Ygc". Below the video player, there are "Related", "More from this user", and "Playlists" buttons. A sidebar on the right lists "Director Videos" and other videos from the same user, such as "Gizmodo Video Review: Honda Civic Hybrid -- iPod On Board 03:59" and "Where the Hell is Matt? 03:42". The bottom of the page shows a footer with the text "ash-v26.ash.youtube.com gelezen" and a large orange circular graphic.

Slide credit: C. Snoek (Univ. of Amsterdam)

QUIZ: WHAT IS THE SIZE OF THIS PROBLEM?

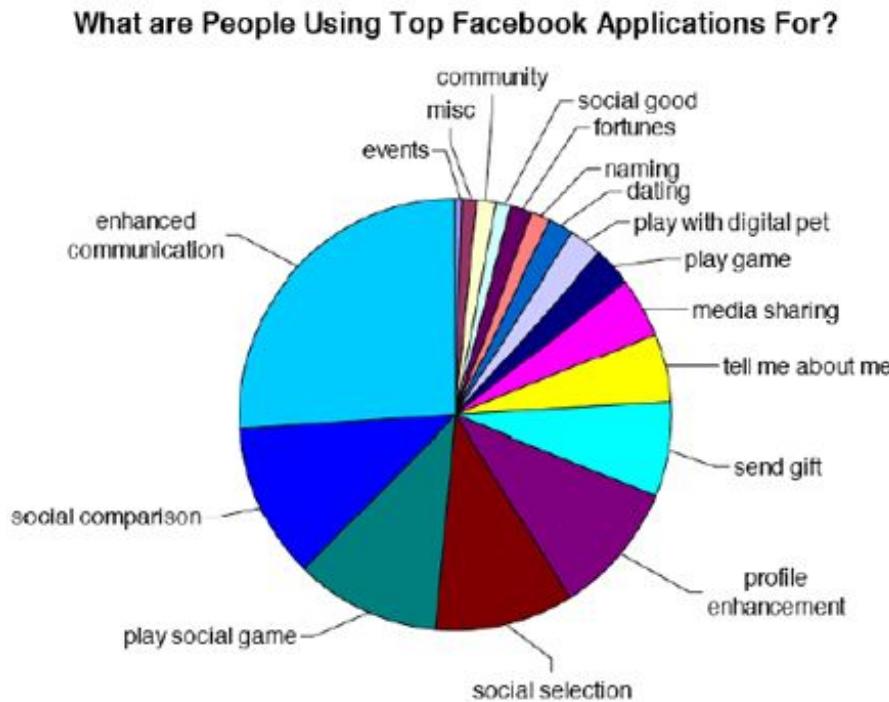
- How many images uploads to Facebook?
- How many uploads on YouTube per day?
- What is the size of a typical national broadcast video archive?





STATISTICS FROM FACEBOOK

- 350 millions of photos uploaded to Facebook daily
- Photo application was NO.1 utilized application for a long time since the inception of Facebook





STATISTICS FROM YOUTUBE

- More than 1 billion unique users visit YouTube each month
- Over 6 billion hours of video are watched each month on
- 100 hours of video are uploaded to YouTube every minute
- 80% of YouTube traffic comes from outside the US
- YouTube is localized in 61 countries and across 61 languages
- According to Nielsen, YouTube reaches more US adults ages 18-34 than any cable network
- Millions of subscriptions happen each day. The number of people subscribing daily is up more than 3x since last year, and the number of daily subscriptions is up more than 4x since last year



ANSWER FROM THE NETHERLANDS

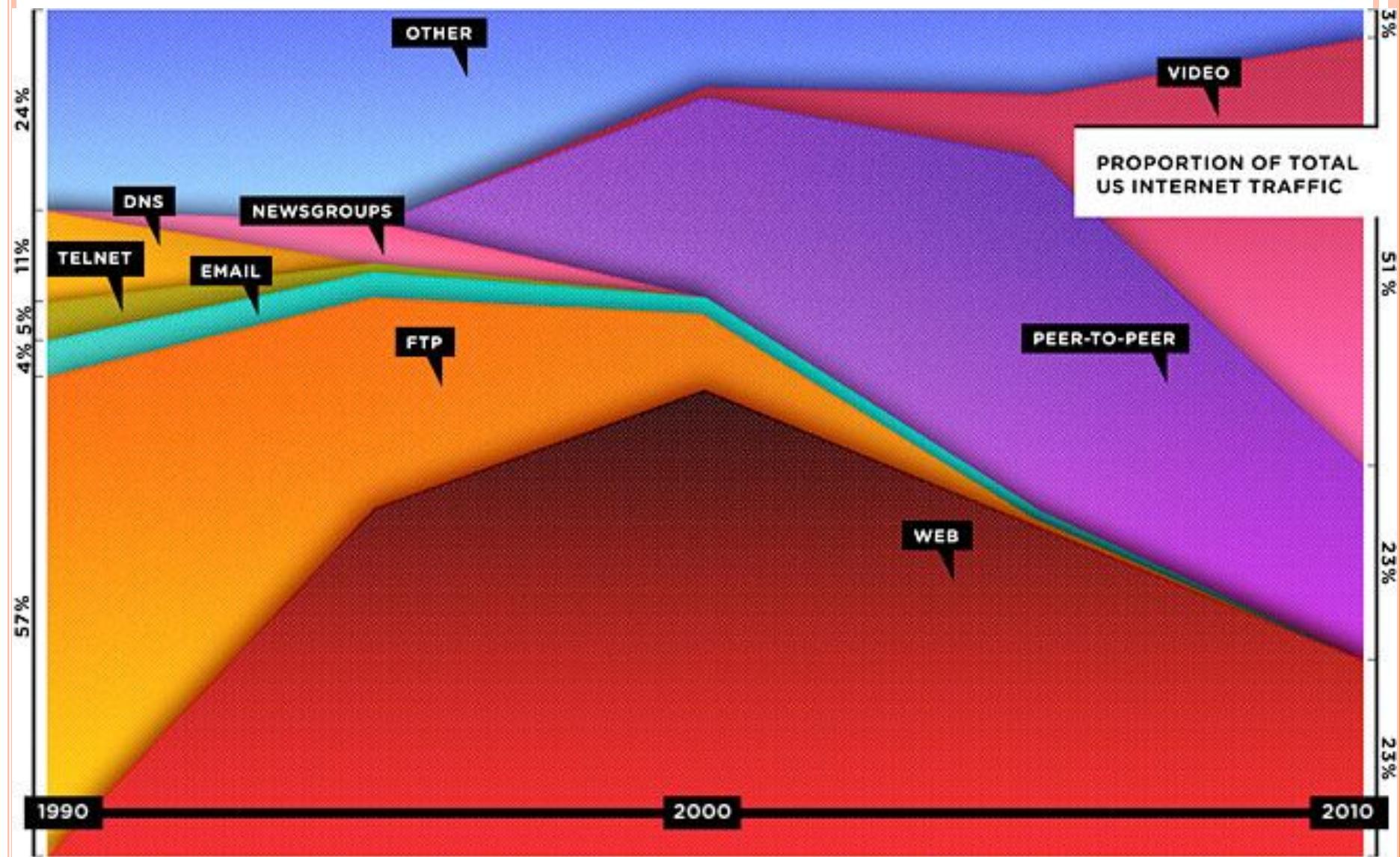
- Europe's largest digitization project
- BEELDEN VOOR DE TOEKOMST
- Yearly ingest
 - 15,000 hours of video
 - 40,000 hours of radio
 - >1 peta-byte per year
 - Next 6 years
 - 137,200 hours of video
 - 22,510 hours of film
 - 123,900 hours of audio
 - 2,900,000 photo's



Lack of metadata



LATEST UPDATE FROM WIRED (THE WEB IS DYING!)



VISUAL INFORMATION RETRIEVAL: THREE CUES

- Expert-driven search
 - Exploit professional annotations
- Crowd-driven search
 - Exploit what others are saying about the image
- Content-driven search
 - Exploit the content of the images and video



EXPERT DRIVEN SEARCH



CROWD-DRIVEN SEARCH

Amsterdam red light zone - Bing Images - Windows Internet Explorer
http://www.bing.com/images/search?q=Amsterdam+red+light+zone&form=QBIR&qs=n&sk=&sc=8-12#

File Edit View Favorites Tools Help

Favorites Amsterdam red light zone - Bing Images

Web Images Videos Shopping News Maps More | MSN Hotmail

Sign in Berkeley, California Preferences

bing Images

Amsterdam red light zone

Web Videos Images

SafeSearch moderate Change

1-24 of 187 results

SIZE
Small
Medium
Large
Wallpaper

LAYOUT
Square
Wide
Tall

COLOR
Full color
Black & white
Specific color

STYLE
Photograph
Illustration

PEOPLE
Just faces
Head & shoulders

SEARCH HISTORY
Amsterdam dam square
Amsterdam
Batu caves

See all
Clear all - Turn off

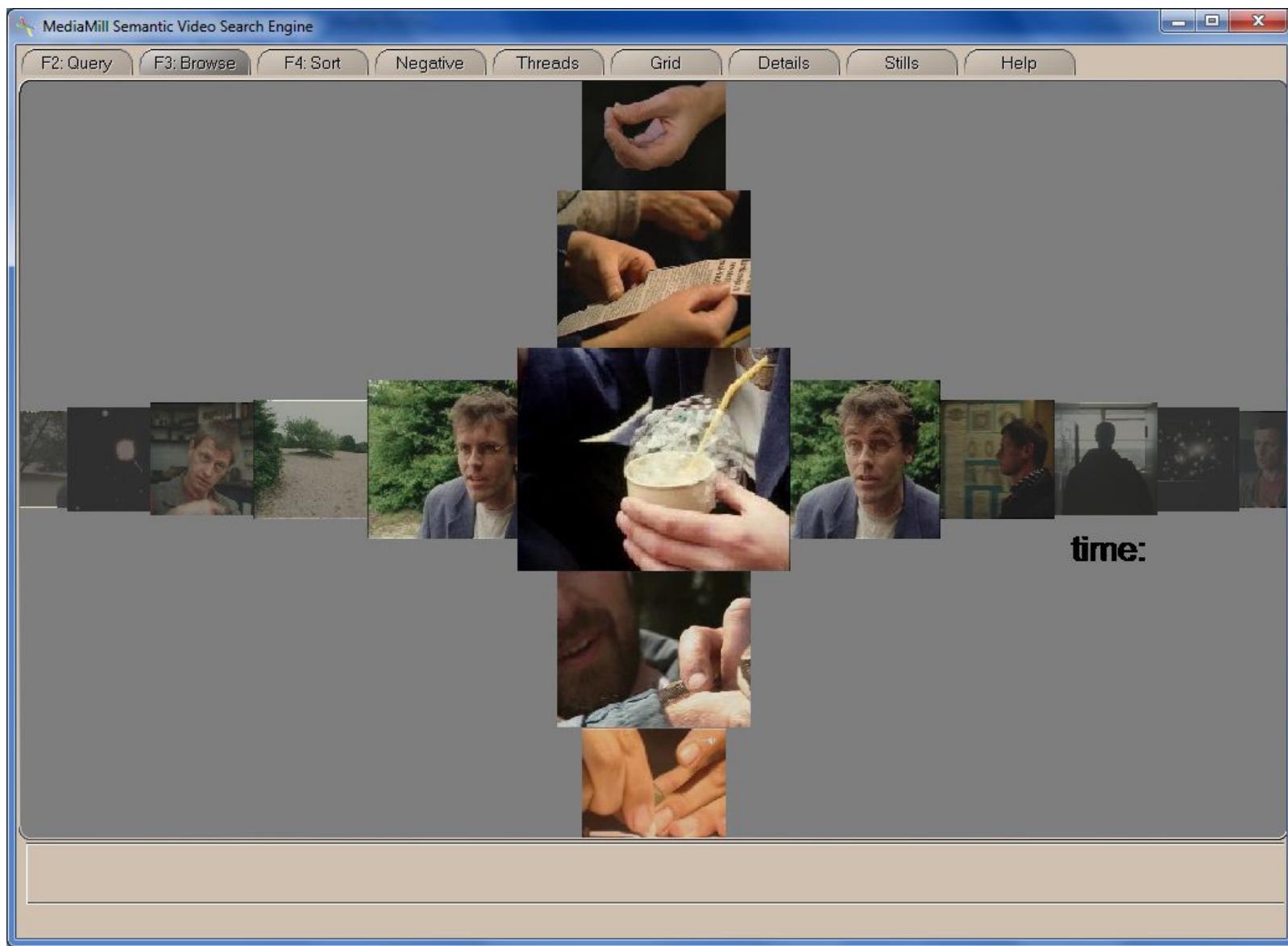
Amsterdam Red Light District
1280 x 960 · 598 kB · jpeg
travel.smart-guide.net

More sizes · Similar images

QUIZ: WHAT'S THE LIMITATIONS OF META DATA BASED SEARCH?

Issue	What's wrong
Too sparse	Few video objects have any metadata
Inadequate	Mainly tags or few keywords, program-guide info for broadcast video, speech available in few cases
Coarse-grain	At level of digital objects only
Not visual	Does not describe what is visually depicted
Ambiguous	Taxonomies not widely used; folksonomies creating new problems
Inconsistent	Vocabularies and taxonomies not standardized
Subjective	Limited verification across users
Not trustworthy	Professional metadata mixed-in with noise

CONTENT-DRIVEN SEARCH



WHAT ARE THE CHALLENGES?



PROBLEM 1: VISUAL VARIATION

- Many images of one thing, due to minor differences in
 - Illumination, Background, Occlusion, Viewpoint, ...



This is also called the “sensory gap.”



PROBLEM 2: VISUAL REPRESENTATION & DESCRIPTION



PROBLEM 3: MANY THINGS IN THE WORLD

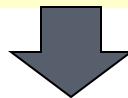


□ This is the model gap.

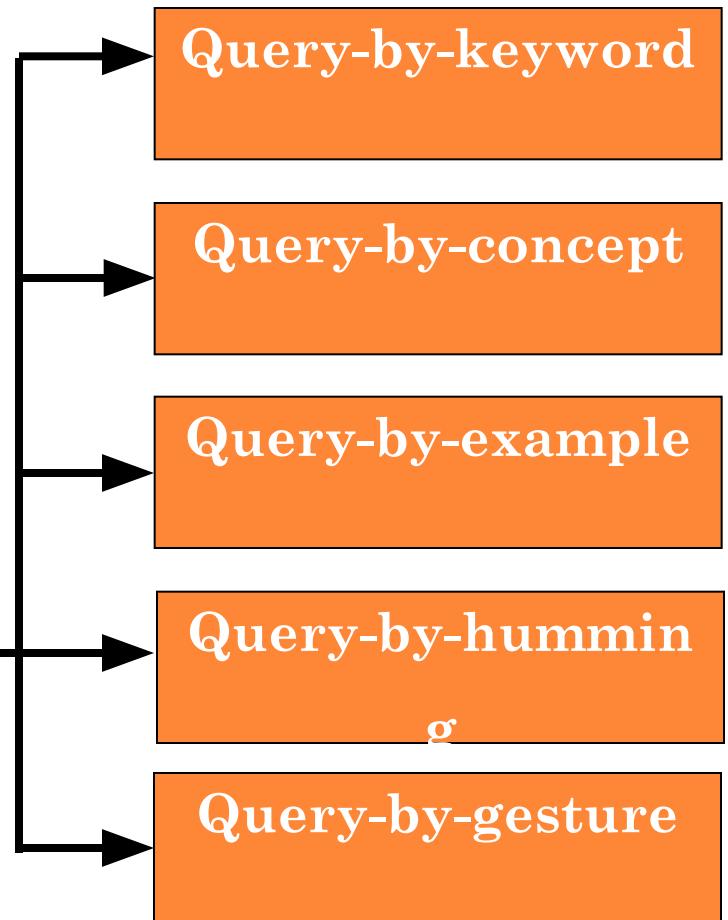
PROBLEM 4: QUERY VARIATION



Find shots of people
shaking hands



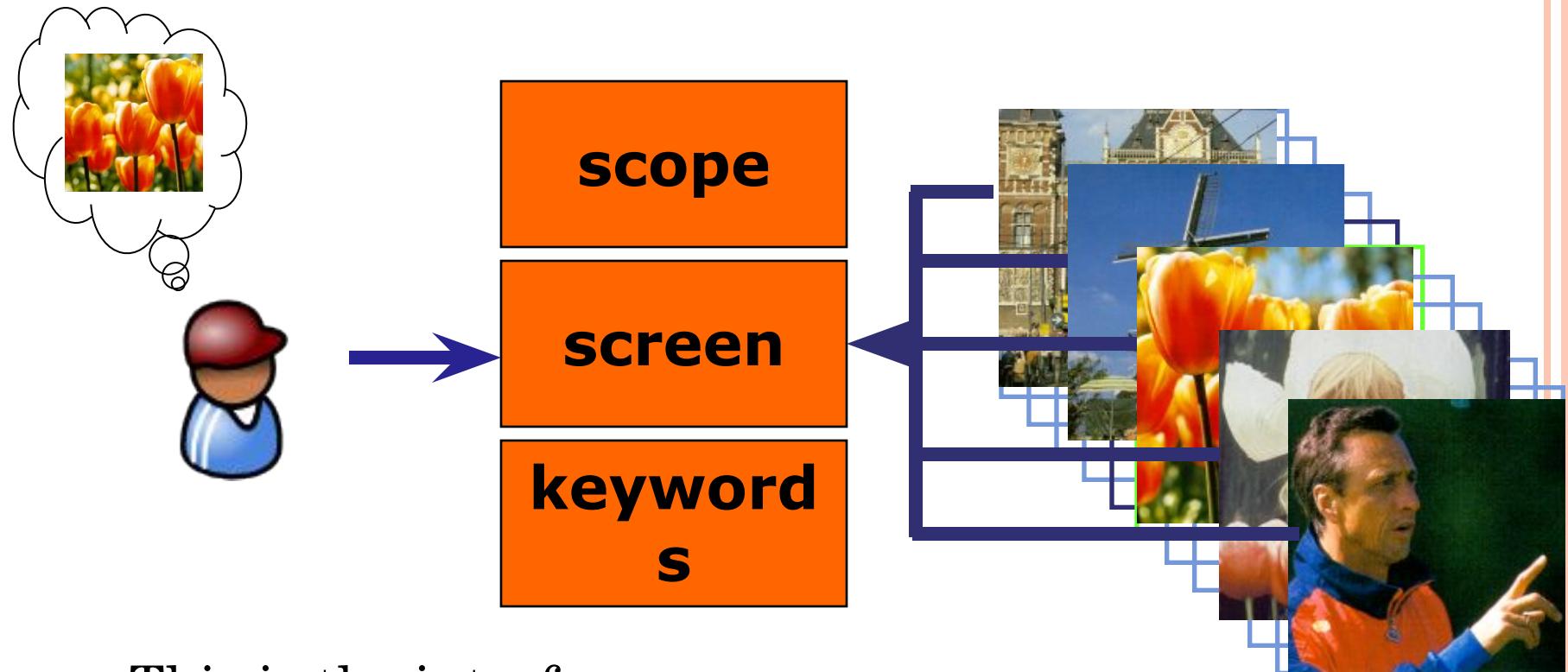
Query
Prediction



This is the query-context gap

Any combination, any sequence?

PROBLEM 5: USE IS OPEN-ENDED



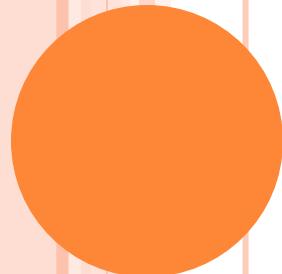
SUMMARY ON PROBLEMS

- Visual search is a diverse and challenging research topic
 - Sensory gap
 - Semantic gap
 - Model gap
 - Query-context gap
 - Interface gap



Q&A





CS598 VISUAL INFORMATION RETRIEVAL

Lecture I: Part II: Basics of Information Retrieval

INFORMATION RETRIEVAL (IR)

- Definition
 - Information retrieval is the activity of obtaining (**i.e., search and retrieve**) information resources relevant to an information need (**i.e., a query**) from a collection of information resources (**e.g., a document database**).



BASIC ASSUMPTIONS OF IR

- Collection: Fixed set of documents
- Goal: Retrieve documents with information that is relevant to the user's information need and helps the user complete a task



HOW GOOD ARE THE RETRIEVED DOCS?

- *Precision* : The fraction of retrieved docs, i.e., those are claimed by the retrieval algorithm to be relevant to the query, that are ***actually relevant*** to user's information need
- *Recall* : The fraction of relevant docs in the collection that are retrieved (i.e., the collection of document that are claimed by the retrieval algorithm to be relevant to the query).
- More measurements to follow in later lectures



INFORMATION RETRIEVAL IN 1680

- Which plays of Shakespeare contain the words ***Brutus*** AND ***Caesar*** but *NOT Calpurnia*?
- One could grep all of Shakespeare's plays for ***Brutus*** and ***Caesar***, then strip out lines containing ***Calpurnia***, but...
 - Slow (for large corpora)
 - *NOT* ***Calpurnia*** is non-trivial
 - Other operations (e.g., find the word ***Romans*** near ***countrymen***) not feasible



TERM-DOCUMENT INCIDENCE

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	1	1	0	0	0	1
Brutus	1	1	0	1	0	0
Caesar	1	1	0	1	1	1
Calpurnia	0	1	0	0	0	0
Cleopatra	1	0	0	0	0	0
mercy	1	0	1	1	1	1
worser	1	0	1	1	1	0

1 if play contains word, 0 otherwise

BOOLEAN QUERY SOLVING FOR DUMMIES

- Computing the answer to the Boolean query

Brutus AND Caesar BUT NOT Calpurnia

- Take the vectors for ***Brutus***, ***Caesar*** and ***Calpurnia*** (complemented) □ bitwise *AND*.
 - 110100 *AND* 110111 *AND* 101111 = 100100.



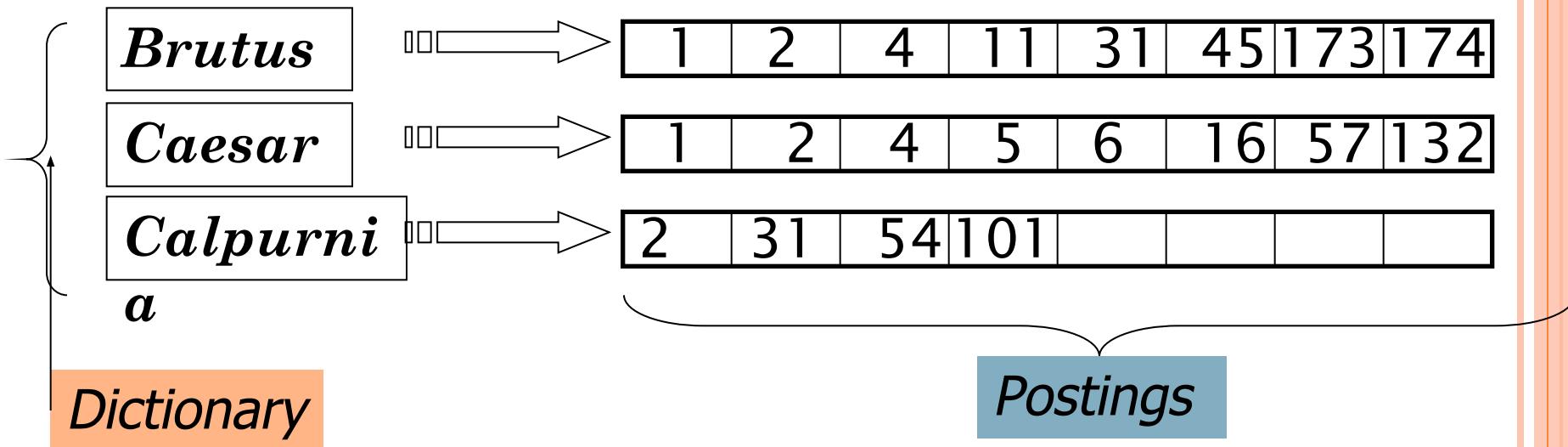
QUIZ

- Bitwise AND is slow and memory intensive for large text collections, how to speed this up?



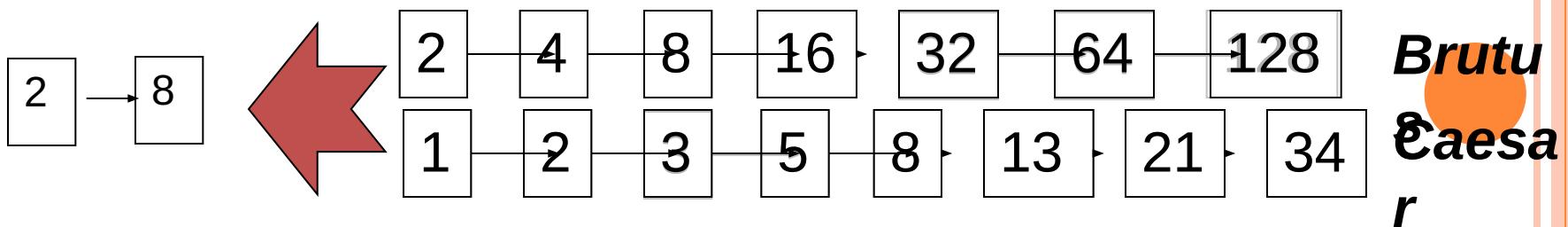
INVERTED INDEX

- For each term t , we must store a list of all documents that contain t .
 - Identify each by a **docID**, a document serial number



QUERY PROCESSING: AND

- Consider processing the query: ***Brutus AND Caesar***
 - Locate ***Brutus*** in the Dictionary;
 - Retrieve its postings.
 - Locate ***Caesar*** in the Dictionary;
 - Retrieve its postings.
 - “Merge” the two postings:



BOOLEAN RETRIEVAL MODEL

- Thus far, our queries have all been Boolean.
 - Documents either match or don't
- Good for expert users with precise understanding of their needs and the collection
- Not good for the majority of users
 - Most users incapable or too lazy to write Boolean queries
 - Most users don't want to wade through 1000s of results



RANKED RETRIEVAL MODELS

- The system returns an ordering over the (top) documents in the collection with respect to a query
- Free text queries
 - The user's query is just one or more words in a human language



SCORING AS THE BASIS OF RANKED RETRIEVAL

- We wish to return an ordered set of documents most likely to be useful to the searcher
- How can we rank-order the documents in the collection with respect to a query?
 - Assign a score – say in [0, 1] – to each document
 - This score measures how well document and query “match”.



TERM-DOCUMENT FREQUENCY MATRICES

- Consider the number of occurrences of a term in a document:
 - Each document is a count vector in \mathbb{N}^v : a column

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0



BAG OF WORDS MODEL

- Vector representation doesn't consider the ordering of words in a document
 - “*John is quicker than Mary*” and “*Mary is quicker than John*” have the same vectors
- This is called the bag of words model



TERM FREQUENCY (TF)

- The term frequency $\text{tf}_{t,d}$ of term t in document d is defined as the number of times that t occurs in d .
- Raw term frequency is not what we want:
 - A document with 10 occurrences of the term is more relevant than a document with 1 occurrence of the term.
 - But not 10 times more relevant.
- Relevance does not increase proportionally with term frequency



LOG FREQUENCY WEIGHTING

- The log frequency weight of term t in d defined as follows:
- Score for a document t in both q and d :
$$w_{t,d} = \begin{cases} 1 + \log_{10} \text{tf}_{t,d} & \text{if } \text{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$
- The score is 0 if none of the query terms is present in the document

$$\text{matching-score} = \sum_{t \in q \cap d} (1 + \log \text{tf}_{t,d})$$



DOCUMENT FREQUENCY (DF)

- Rare terms are more informative than frequent terms
- Consider a term in the query that is rare in the collection (e.g., *arachnocentric*)
 - A document containing this term is very likely to be relevant to the query *arachnocentric*
→ We want a high weight for rare terms like *arachnocentric*



DOCUMENT FREQUENCY (CONT'D)

- Frequent terms are less informative than rare terms
- Consider a query term that is frequent in the collection (e.g., *high*, *increase*, *line*)
 - A document containing such a term is more likely to be relevant than a document that doesn't
 - But it's not a sure indicator of relevance.
- For frequent terms, we want high positive weights but lower weights than for rare terms.
 - We will use document frequency (df) to capture this.



IDF WEIGHT

- df_t is the document frequency of t : the number of documents that contain t
 - df_t is an inverse measure of the informativeness of t
 - $\text{df}_t \leq N$
- We define the idf (inverse document frequency) of t by
 - We use $\log(N/\text{df}_t)$ instead of N/df_t to “dampen” the effect of idf.

$$\text{idf}_t = \log_{10}(N/\text{df}_t)$$



IDF EXAMPLE (N=1 MILLION)

term	df_t	idf_t
calpurnia		1
animal		100
sunday		1,000
fly		10,000
under		100,000
the		1,000,000

$$idf_t = \log_{10} (N/df_t)$$

There is one idf value for each term t in a collection



TF.IDF WEIGHTING

- The tf.idf weight of a term is the product of its tf weight and its idf weight.

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \cdot \log \frac{N}{\text{df}_t}$$

- Best known weighting scheme in information retrieval
- Increases with
 - the number of occurrences within a document
 - the rarity of the term in the collection



FINAL RANKING OF DOCUMENTS FOR A QUERY

$$\text{Score}(q, d) = \sum_{t \in q \cap d} \text{tf.idf}_{t,d}$$



BINARY → COUNT → WEIGHT MATRIX

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.21	6.1	0	1	0	0
Caesar	8.59	2.54	0	1.51	0.25	0
Calpurnia	0	1.54	0	0	0	0
Cleopatra	2.85	0	0	0	0	0
mercy	1.51	0	1.9	0.12	5.25	0.88
worser	1.37	0	0.11	4.15	0.25	1.95

- Each document is now represented by a real-valued vector of tf.idf weights $\in \mathbb{R}^{|V|}$

DOCUMENTS AS VECTORS

- So we have a $|V|$ -dimensional vector space
 - Terms are axes of the space
 - Documents are points or vectors in this space
- Very high-dimensional:
 - Tens of millions of dimensions when you apply this to a web search engine
 - These are very sparse vectors - most entries are zero



QUERIES AS VECTORS

- **Key idea 1:** Do the same for queries: represent them as vectors in the space
- **Key idea 2:** Rank documents according to their proximity to the query in this space
 - proximity = similarity of vectors



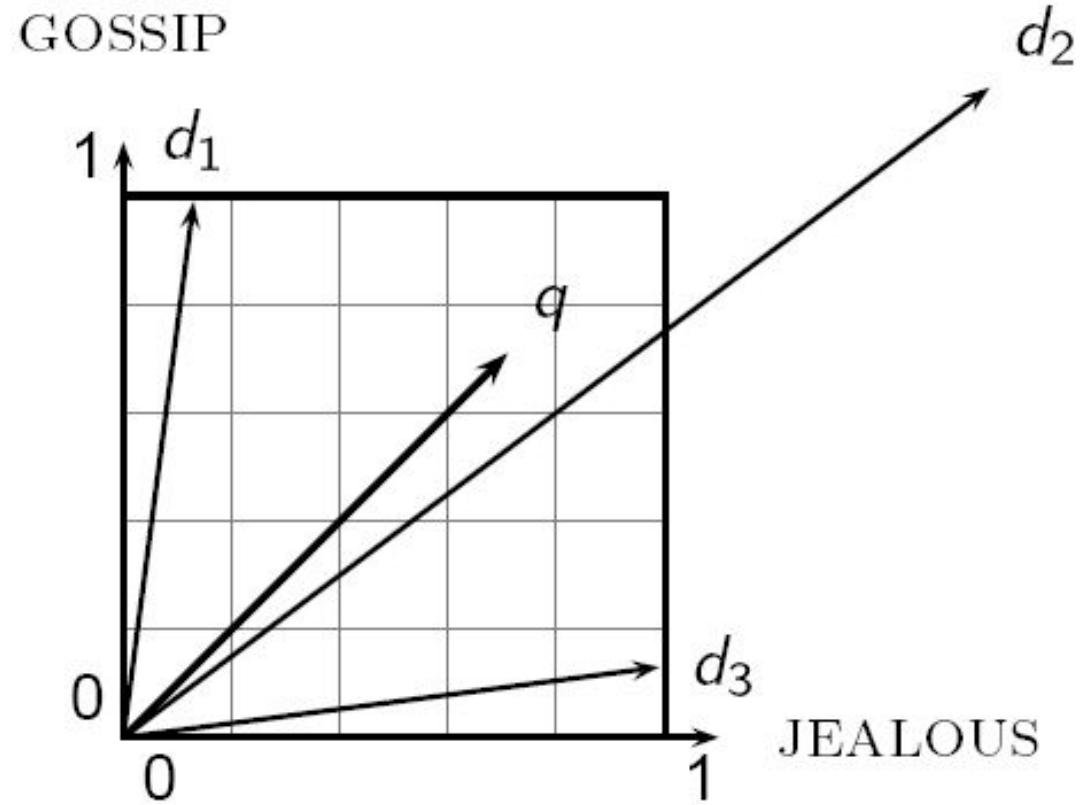
QUIZ: WHY IS EUCLIDEAN DISTANCE A BAD ESTIMATOR FOR PROXIMITY?

- . . . because Euclidean distance is **large** for vectors of **different lengths**.



WHY DISTANCE IS A BAD IDEA

The Euclidean distance between \vec{q} and \vec{d}_2 is large even though the distribution of terms in the query \vec{q} and the distribution of terms in the document \vec{d}_2 are very similar.



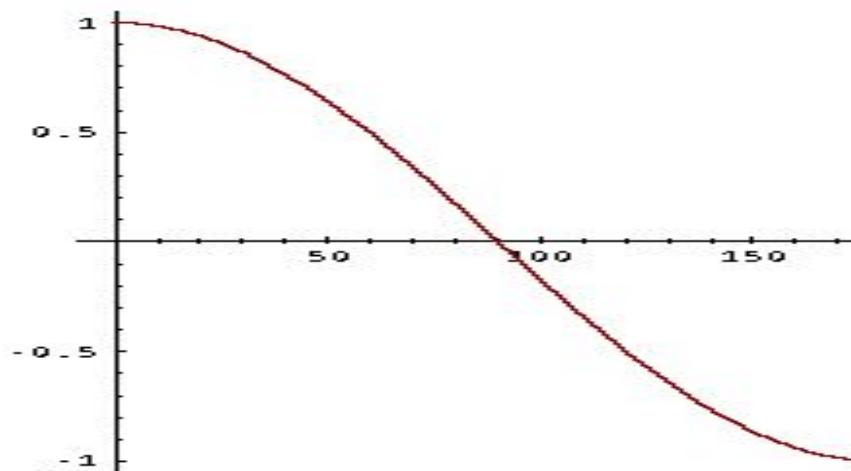
USE ANGLE INSTEAD OF DISTANCE

- Thought experiment: take a document d and append it to itself. Call this document d' .
 - “Semantically” d and d' have the same content
- The Euclidean distance between the two documents can be quite large
 - The angle between the two documents is 0, corresponding to maximal similarity.
- **Key idea:** Rank documents according to angle with query



FROM ANGLES TO COSINES

- The following two notions are equivalent.
 - Rank documents in decreasing order of the angle between query and document
 - Rank documents in increasing order of cosine(query,document)
- Cosine is a monotonically decreasing function for the interval $[0^\circ, 180^\circ]$



LENGTH NORMALIZATION

- A vector can be (length-) normalized by dividing each of its components by its length – for this we use the L_2 norm:

$$\left\| \mathbf{x} \right\|_2 = \sqrt{\sum_i x_i^2}$$

- Dividing a vector by its L_2 norm makes it a unit (length) vector
 - on surface of unit hypersphere



LENGTH NORMALIZATION (CONT'D)

- Effect on the two documents d and d' (d appended to itself) from earlier slide:
 - they have identical vectors after length-normalization
- Long and short documents now have comparable weights



COSINE SIMILARITY BETWEEN QUERY AND DOCUMENT

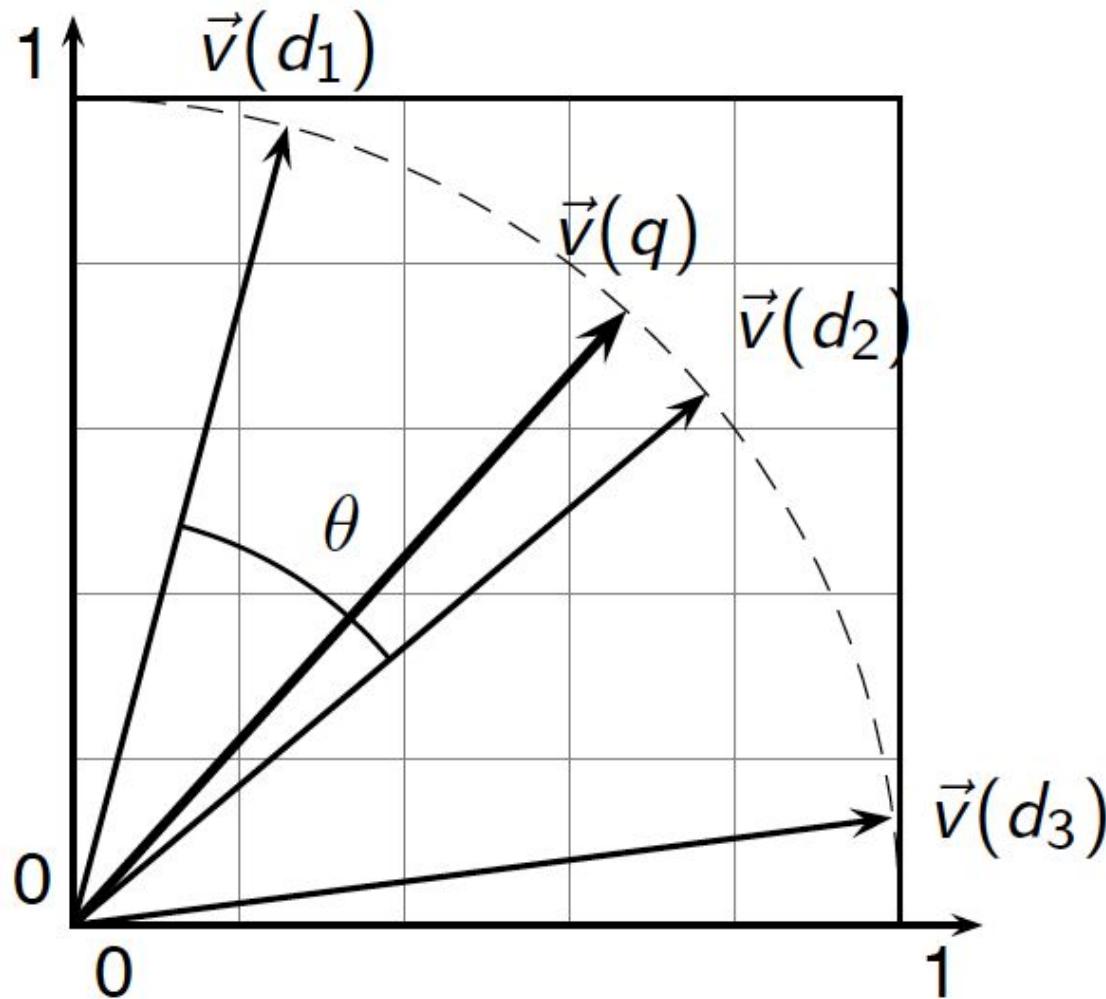
$$\cos(\vec{q}, \vec{d}) = \text{SIM}(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| |\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

- q_i is the tf-idf weight of term i in the query
- d_i is the tf-idf weight of term i in the document
- This is the cosine similarity of q and d



COSINE SIMILARITY ILLUSTRATED

POOR



TF.IDF WEIGHTING HAS MANY VARIANTS

Term frequency	Document frequency	Normalization
n (natural) $tf_{t,d}$	n (no) 1	n (none) 1
I (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) $\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$		



SUMMARY – VECTOR SPACE MODEL AND RANKING

- Represent the query as a weighted tf.idf vector
- Represent each document as a weighted tf.idf vector
- Compute the cosine similarity score for the query vector and each document vector
- Rank documents with respect to the query by score
- Return the top K (e.g., $K = 10$) to the user



AN ALTERNATIVE: OKAPI/BM25

- Not the full theory
- From idf weighting to BM25 in a few steps
- Simplest score for doc d is idf weighting of the query terms present:

$$RSV_d = \sum_{t \in q} \log \frac{N}{df_t}$$

- (RSV: retrieval status value, the “score”)
- Next step, factoring in frequency of each term and document length:



OKAPI/BM25

$$RSV_d = \sum_{t \in q} \log \left[\frac{N}{df_t} \right] \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d / L_{ave})) + tf_{td}}$$

- Here, tf_{td} is the tf of t in d and L_d and L_{ave} are the length of d and the average doc length for the whole collection
- Parameters
 - k_1 is a positive tuning parameter that calibrates the doc term frequency scaling
 - b calibrates the doc length scaling



OKAPI/BM25

$$RSV_d = \sum_{t \in q} \left[\log \frac{N}{df_t} \right] \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d / L_{ave})) + tf_{td}} \cdot \frac{(k_3 + 1)tf_{tq}}{k_3 + tf_{tq}}$$

- If the query is long use weighting for query terms
 - tf_{tq} : frequency of t in q
 - k_3 : tuning frequency scaling of query



QUIZ: THREE KEY NOTIONS?

- Many “document retrieval models” share three key ingredients:
- **Term frequency**
 - A term is a better indicator of the contents of a document if it occurs more frequently in that document
- **Document frequency**
 - Terms that occur in many documents are not good discriminators
- **Document length normalization**
 - Longer documents are more likely to contain a term (without necessarily being “about” the term)



QUERY INDEPENDENT RANKING

- PageRank (named after Larry Page)
 - <http://ilpubs.stanford.edu:8090/422/1/1999-66.pdf>
- Citation analysis: analysis of citations in the scientific literature
 - Example citation: “Miller (2001) has shown that physical activity alters the metabolism of estrogens.”
 - “Miller (2001)” is a hyperlink linking two scientific articles.



QUERY INDEPENDENT RANKING

- One application of these “hyperlinks” in the scientific literature:
 - Measure the similarity of two articles by the overlap of other articles citing them
 - This is called cocitation similarity
- **Key idea:** Cocitation similarity on the web?



LINK-BASED RANKING FOR WEB SEARCH

- Simple version of using links for ranking on the web
 - First retrieve all pages satisfying the query (say *venture capital*)
 - Order these by the number of in-links
- Simple link popularity (= number of in-links) is easy to spam.
Why?



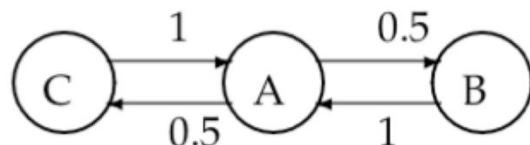
PAGERANK BASIS: RANDOM WALK

- Imagine a web surfer doing a random walk on the web
 - Start at a random page
 - At each step, go out of the current page along one of the links on that page, equiprobably
- In the steady state, each page has a **long-term visit rate**.
- This long-term visit rate is the page's **PageRank**.
- **PageRank = steady state probability = long-term visit rate**
- Concept of long-term visit rate clear?



MARKOV CHAINS

- A Markov chain consists of N states, plus an $N \times N$ transition probability matrix P .
- state = page
- At each step, we are on exactly one of the pages.
- For $1 \leq i, j \leq N$, the matrix entry P_{ij} tells us the probability of j being the next page, given we are currently on page i .

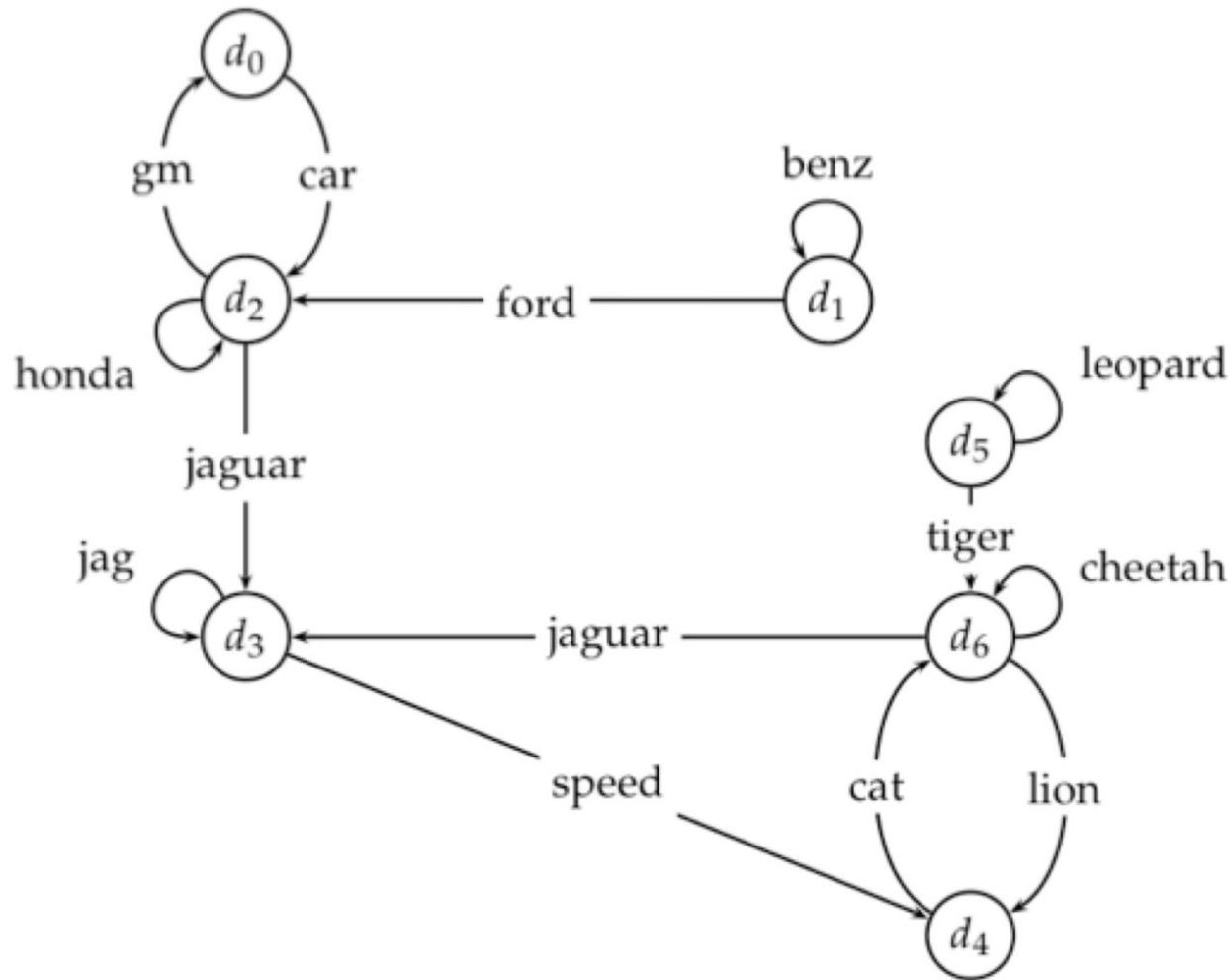


MARKOV CHAINS

- Clearly, for all i , $\sum_{j=1}^N P_{ij} = 1$
- Markov chains are abstractions of random walks.



EXAMPLE WEB GRAPH



LINK MATRIX FOR EXAMPLE

	d_0	d_1	d_2	d_3	d_4	d_5	d_6
d_0	0	0	1	0	0	0	0
d_1	0	1	1	0	0	0	0
d_2	1	0	1	1	0	0	0
d_3	0	0	0	1	1	0	0
d_4	0	0	0	0	0	0	1
d_5	0	0	0	0	0	1	1
d_6	0	0	0	1	1	0	1



TRANSITION MATRIC FOR EXAMPLE

	d_0	d_1	d_2	d_3	d_4	d_5	d_6
d_0	0.00	0.00	1.00	0.00	0.00	0.00	0.00
d_1	0.00	0.50	0.50	0.00	0.00	0.00	0.00
d_2	0.33	0.00	0.33	0.33	0.00	0.00	0.00
d_3	0.00	0.00	0.00	0.50	0.50	0.00	0.00
d_4	0.00	0.00	0.00	0.00	0.00	0.00	1.00
d_5	0.00	0.00	0.00	0.00	0.00	0.50	0.50
d_6	0.00	0.00	0.00	0.33	0.33	0.00	0.33



LONG-TERM VISIT RATE

- Recall: PageRank = long-term visit rate
- Long-term visit rate of page d is the probability that a web surfer is at page d at a given point in time.
- Next: what properties must hold of the web graph for the long-term visit rate to be well defined?
- The web graph must correspond to an **ergodic** Markov chain.
- First a special case: The web graph must not contain **dead ends**.



DEAD ENDS

- The web is full of dead ends.
- Random walk can get stuck in dead ends.
- If there are dead ends, long-term visit rates are not well-defined (or non-sensical).



QUIZ: HOW TO COPE WITH DEAD ENDS?



TELEPORTING

- At a dead end, jump to a random web page
- At any non-dead end, with probability 10%, jump to a random web page
- With remaining probability (90%), go out on a random hyperlink (e.g., randomly choose with probability $(1-0.1)/4=0.225$ one of the four hyperlinks of the page)
- 10% is a parameter.



RESULT OF TELEPORTING

- With teleporting, we cannot get stuck in a dead end.
- Concept of teleporting clear?
- Even without dead-ends, a graph may not have well-defined long-term visit rates.
- More generally, we require that the Markov chain be ergodic.



ERGODIC MARKOV CHAINS

- A Markov chain is ergodic iff it is irreducible and aperiodic.
- **Irreducibility.** Roughly: there is a path from any page to any other page.
- **Aperiodicity.** Roughly: The pages cannot be partitioned such that the random walker visits the partitions sequentially.



ERGODIC MARKOV CHAINS

- Theorem: For any ergodic Markov chain, there is a unique long-term visit rate for each state.
- This is the steady-state probability distribution.
- Over a long time period, we visit each state in proportion to this rate.
- It doesn't matter where we start.



FORMALIZATION OF “VISIT”

- A probability (row) vector $\vec{x} = (x_1, \dots, x_N)$ tells us where the random walk is at any point.
- Example:
$$\begin{pmatrix} 0 & 0 & 0 & \dots & 1 & \dots & 0 & 0 & 0 \\ 1 & 2 & 3 & \dots & i & \dots & N-2 & N-1 & N \end{pmatrix}$$
- More generally: the random walk is on page i with probability x_i .
- Example:
$$\begin{pmatrix} 0.05 & 0.01 & 0.0 & \dots & 0.2 & \dots & 0.01 & 0.05 & 0.03 \\ 1 & 2 & 3 & \dots & i & \dots & N-2 & N-1 & N \end{pmatrix}$$
- $\sum x_i = 1$



CHANGE IN PROBABILITY VECTOR

- If the probability vector is $\vec{x} = (x_1, \dots, x_N)$ at this step, what is it at the next step?
- Recall that row i of the transition probability matrix P tells us where we go next from state i .
- Equivalently: column j of P tells us “where we came from” (and with which probability).
- So from \vec{x} , our next state is distributed as $\vec{x}P$.



STEADY STATE IN VECTOR NOTATION

- The steady state in vector notation is simply a vector $\vec{\pi} = (\pi_1, \pi_2, \dots, \pi_N)$ of probabilities.
- (We use $\vec{\pi}$ to distinguish it from the generic notation \vec{x} for a probability vector.)
- π_i is the long-term visit rate (or PageRank) of page i .
- So we can think of PageRank as a very long vector – one entry per page.



HOW DO WE COMPUTE THE STEADY STATE VECTOR?

- In other words: how do we compute PageRank?
- Recall: $\vec{\pi} = (\pi_1, \pi_2, \dots, \pi_N)$ is the PageRank vector, the vector of steady-state probabilities ...
- ... and if the distribution in this step is described by \vec{x} , then the distribution in the next step is distributed as $\vec{x}P$.
- But $\vec{\pi}$ is the steady state! So: $\vec{\pi} = \vec{\pi}P$
- Solving this matrix equation gives us $\vec{\pi}$.
- $\vec{\pi}$ is the principal left eigenvector for P ...
- ... that is, $\vec{\pi}$ is the left eigenvector with the largest eigenvalue.
- Transition probability matrices always have largest eigenvalue 1.



ONE WAY OF COMPUTING THE PAGERANK

- Recall: regardless of where we start (except for pathological cases), we eventually reach the steady state $\vec{\pi}$.
- Start with (almost) any distribution \vec{x} , e.g., uniform distribution
- After one step, we're at $\vec{x}P$.
- After two steps, we're at $\vec{x}P^2$.
- After k steps, we're at $\vec{x}P^k$.
- Algorithm: multiply \vec{x} by increasing powers of P until the product looks stable.
- This is called the **power method**.



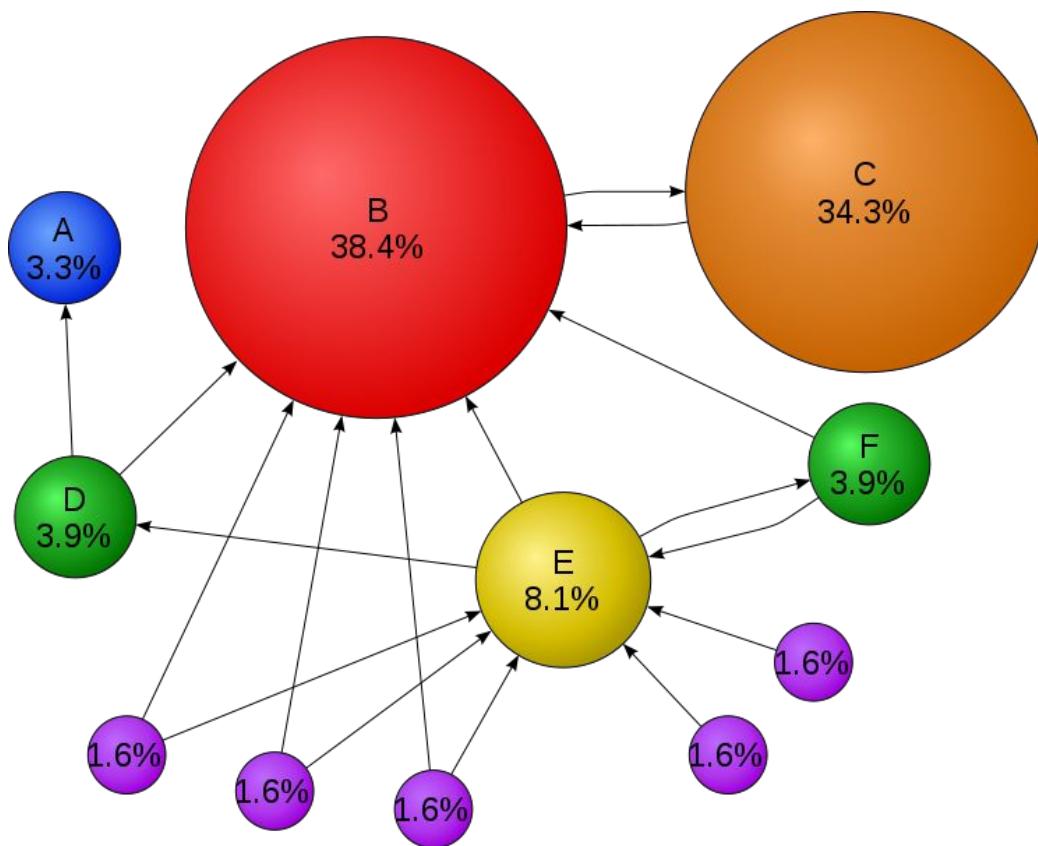
POWER METHOD: EXAMPLE

- Two-node example: $\vec{x} = (0.5, 0.5)$, $P = \begin{pmatrix} 0.25 & 0.75 \\ 0.25 & 0.75 \end{pmatrix}$
- $\vec{x}P = (0.25, 0.75)$
- $\vec{x}P^2 = (0.25, 0.75)$
- Convergence in one iteration!



QUIZ: PAGE RANK

- Why is website C's page rank higher than E's?



PAGERANK SUMMARY

- Preprocessing
 - Given graph of links, build matrix P
 - Apply teleportation
 - From modified matrix, compute $\vec{\pi}$
 - $\vec{\pi}_i$ is the PageRank of page i .
- Query processing
 - Retrieve pages satisfying the query
 - Rank them by their PageRank
 - Return reranked list to the user



PAGERANK ISSUES

- Real surfers are not random surfers – Markov model is not a good model of surfing.
 - Issues: back button, short vs. long paths, bookmarks, directories – and search!
- Simple PageRank ranking (as described on previous slide) produces bad results for many pages.
 - Consider the query *video service*
 - The Yahoo home page (i) has a very high PageRank and (ii) contains both words.
 - If we rank all Boolean hits according to PageRank, then the Yahoo home page would be top-ranked.
 - Clearly not desirable
- In practice: rank according to weighted combination of many factors, including raw text match, anchor text match, PageRank and many other factors

HOW IMPORTANT IS PAGERANK?

- Frequent claim: PageRank is the most important component of web ranking.
- The reality:
 - There are several components that are at least as important: e.g., anchor text, phrases, proximity, tiered indexes ...
 - Rumor has it that PageRank in its original form (as presented here) has a negligible impact on ranking!
 - However, variants of a page's PageRank are still an essential part of ranking.
 - Addressing link spam is difficult and crucial.



NEXT WEEK

- We will start from simple visual representations of an image



ACKNOWLEDGEMENTS

- Information Retrieval slide credit:
 - C. Manning and P. Raghavan, Stanford University
 - M. de Rijke, University of Amsterdam

