

LO1: Search Engine (SE)

Names for different retrieval information (IR) | Document (DK) | Text (TR)

Specific SE

Vertical	One type of data (job search, news search)
Site	Just one site
Custom	narrow search to small set of websites
Enterprise	Corporate intranet w/ collections, metadata, roles, security

Federated Search system: IR search 'glue' system

Federated search (union)	IR search (filtering rules)
Search with full-time SE	It itself is not a SE
Agreement obtained	No agreement: many block
Standard query/res. top.	novel and parsing
etc. collaborate	SEs do their own thing

Difficulties of IR

Unstructured text/media	Standard	Semantic	Dynamic, Relevance
-------------------------	----------	----------	--------------------

Cloud computing advantages (central IR same)

Virtualization, scaling, many replicas and updates, redundancy for loss

High cost (Inefficient Investment)

Timeline of IR goes 0 to 3
 1990: first web (1993: cat, crawler 1997: hit)
 2000: Business model (ad)

Other things being

search service
 cell the software for internet search
 customized content portals, webpage by ad (e.g. Yahoo)

Search engine

Do the crawling/indexing robot for rights
 Page-based algorithm: res. so pop up & banner

Change end user

Refuse to pay unless unique valuable info
 by queries, for indexing, for ranking

Change website

Refuse to pay unless unique valuable info
 by queries, for indexing, for ranking

Rules

Know what Advanced (key, low), Address (page content)
 (content) Portal: SE used to need portals + now supported, need ad channels

Test (Marketing): SEIT company help get information, find

low rank high
 Pay-by-impression (low per click = low)
 Pay-per-click
 make website rank high in algorithm

Performance

relevance legal problems such as lawsuits

retrieval needs

performance retrieval
 precision: % relevant in
 model (noise in second)

model features details

doc repr. | query repr. | retrieval function
 current class of models
 best model / vector space (linear) probabilistic

Cons: the model, no complex requests, bad IR/interpretation, no control of docs, no ranking, no auto relevance feedback

Latent Semantic model (LSM) (LSA, LSI)

Doc. repr. Set of term w/ stat info
 D = doc. 1, 2, 3, 4, 5, 6, 7, 8, 9
 weight matrix, not ratios
 imp of word
 weights for each vector, 0 for none

Query repr. | low with optimal weights: no balance
 1. if present, else 0

retrieval for similarity

Challenges: word importance (not doc, also words)
 degree of sim b/w doc & q

Language is context: words, context, frequency
 point info
 result with error

LSA: Theory: word importance

(within doc & entire colln. IR)

word importance (term - doc)
 T = doc. 1, 2, 3, 4, 5, 6, 7, 8, 9

Principle: repeat = important
 value: when doc
 fix: frequency
 normalized by dimension
 T = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

word importance (user - doc)
 U = row, d of (unique) words
 load: 0.5 ~ 1

LSA: matrix math
 doc & query space, no bad change, not 1-st
 index: query space
 inv. file (doc matrix)

Data structure: host file, P, H, T, S, etc. by index
 Doc. | Query | Index | Value | Location | Frequency

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

Index: query space

any more evidence

There is no further evidence of the 4th or 5th

There is no further evidence of the 4th or 5th

There is no further evidence of the 4th or 5th