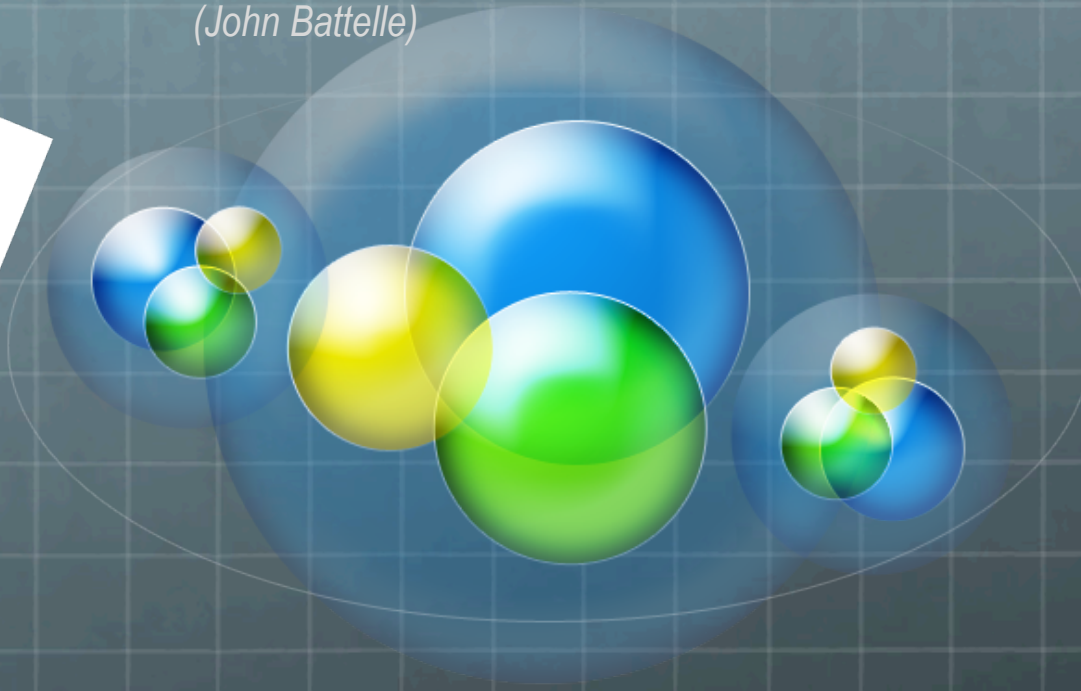
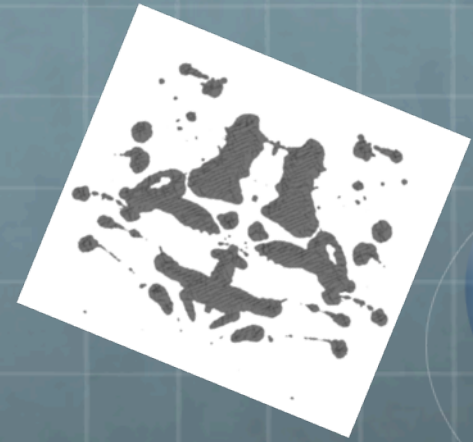


«Google è un globale test di Rorschach. Vediamo
in esso ciò che vogliamo vedere.»

(John Battelle)



TECNOLOGIE INFORMATICHE MULTIMEDIALI


Corso di Laurea “Scienze e Tecnologie della Comunicazione”

Prof. Giorgio Poletti (giorgio.poletti@unife.it)

a.a. 2013-2014

Sviluppo della lezione

Contenuti

-  Algoritmi di ricerca delle informazioni

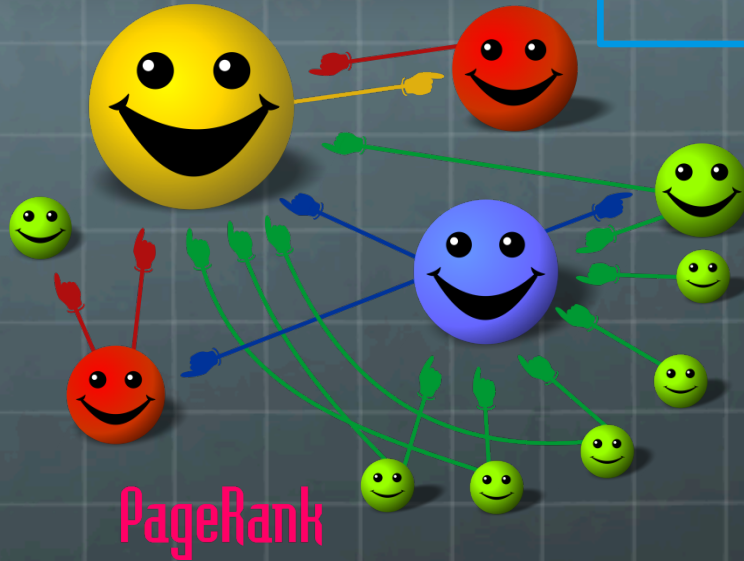
Attività

-  Analisi dell'algoritmo di page rank

Algoritmi di Ricerca delle Informazioni

ALGORITMO DI ANALISI

Assegnazione di un peso ad ogni elemento di collegamento ipertestuale di un insieme di documenti



PageRank è un marchio Google

Algoritmo brevettato dall'università di Stanford

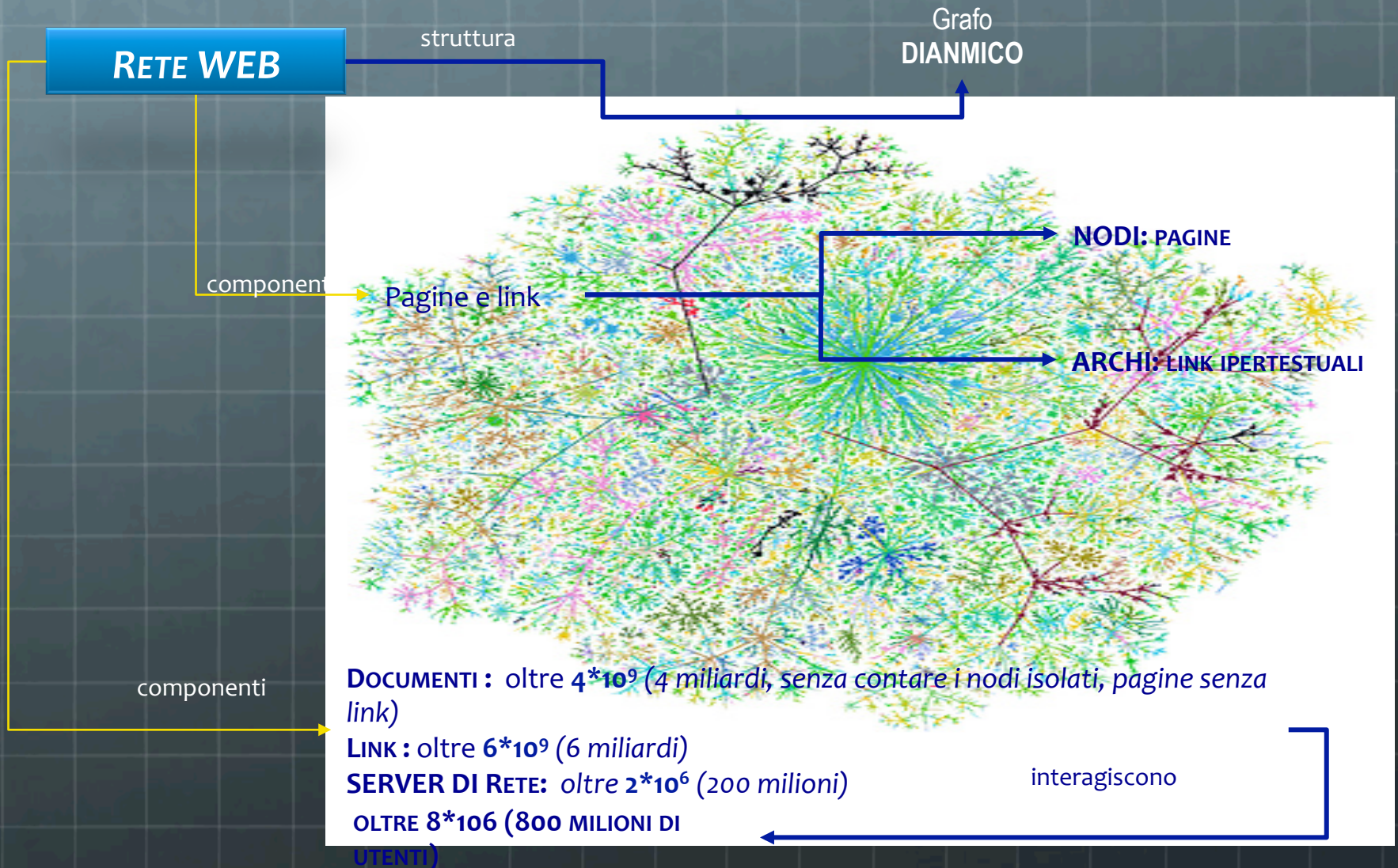
$PR(E)$ = PageRank di E è il valore, peso numerico, associato ad un elemento E

- Una pagina è tanto più importante quanto più numerose sono le pagine che la puntano

Nomen omen...

Algoritmo ideato, tra gli altri da LAWRENCE EDWARD "LARRY" PAGE fondatore, con SERGEY BRIN, di Google.

Algoritmo PageRank



Algoritmo PageRank

GRAFO SEMPLIFICATO PER IL WEB

finito

Grafo $G(P,L)$, insieme di pagine e

link

Se orientato

Relazione \leftrightarrow

$p_1 \leftrightarrow p_2$ se

ESISTE UN CAMMINO DA p_1 A

p_2

e

ESISTE UN CAMMINO DA p_2 A

p_1

è

RELAZIONE DI EQUIVALENZA

DEFINIRE UN GRAFO RIDOTTO G^* , i nodi sono le classi e due classi C_1 e C_2 sono connesse se esiste un nodo in C_1 collegato a un nodo in C_2 , esiste un arco da C_1 a C_2

Classi

composto da

dette

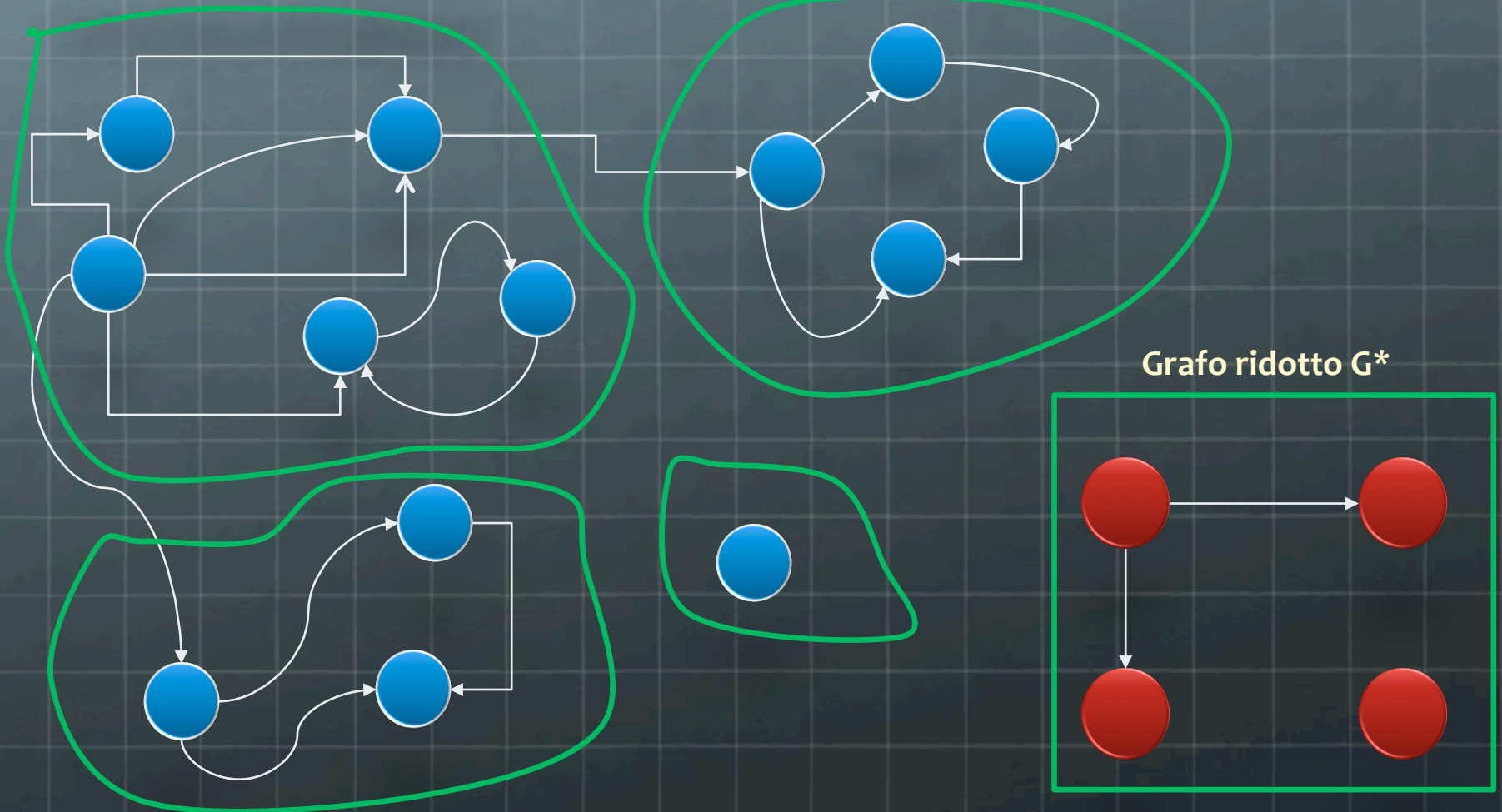
COMPONENTI (FORTEMENTE) CONNESSE DEL GRAFO

Permette di

Algoritmo PageRank

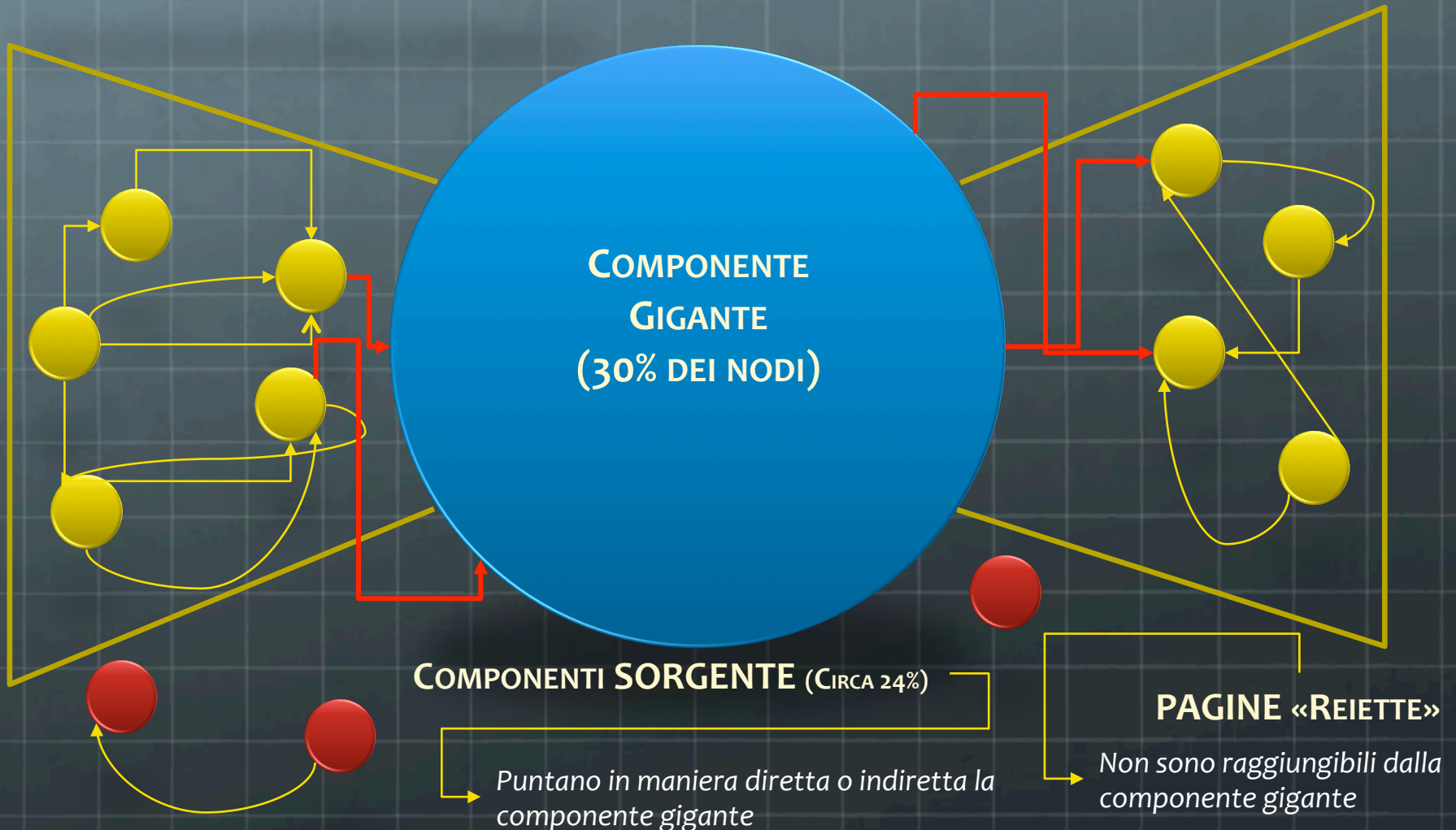
GRAFO RIDOTTO G^*

esempio



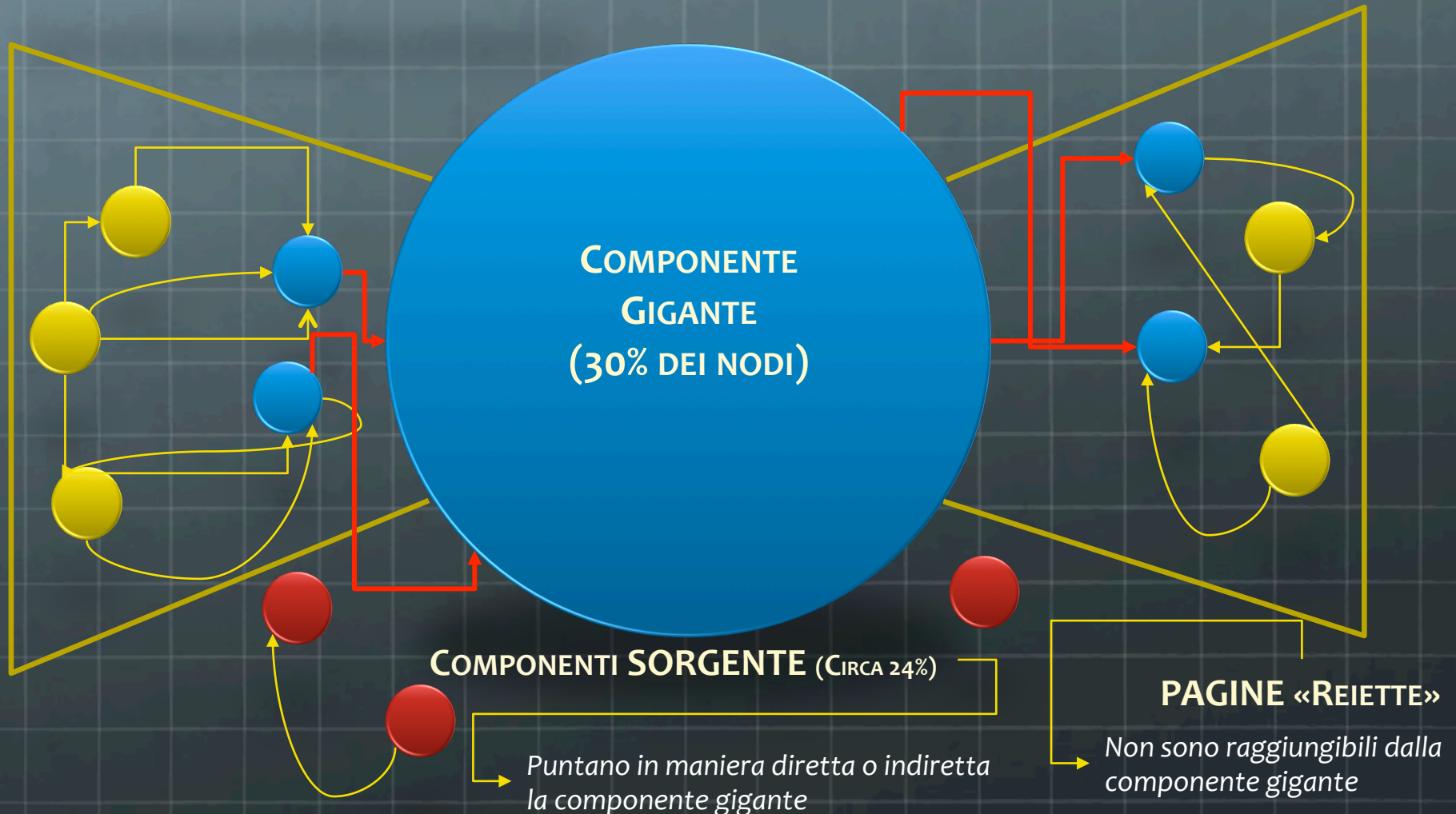
Algoritmo PageRank

STRUTTURA A CAMELLA DEL WEB



Algoritmo PageRank

STRUTTURA A CARAMELLA DEL WEB



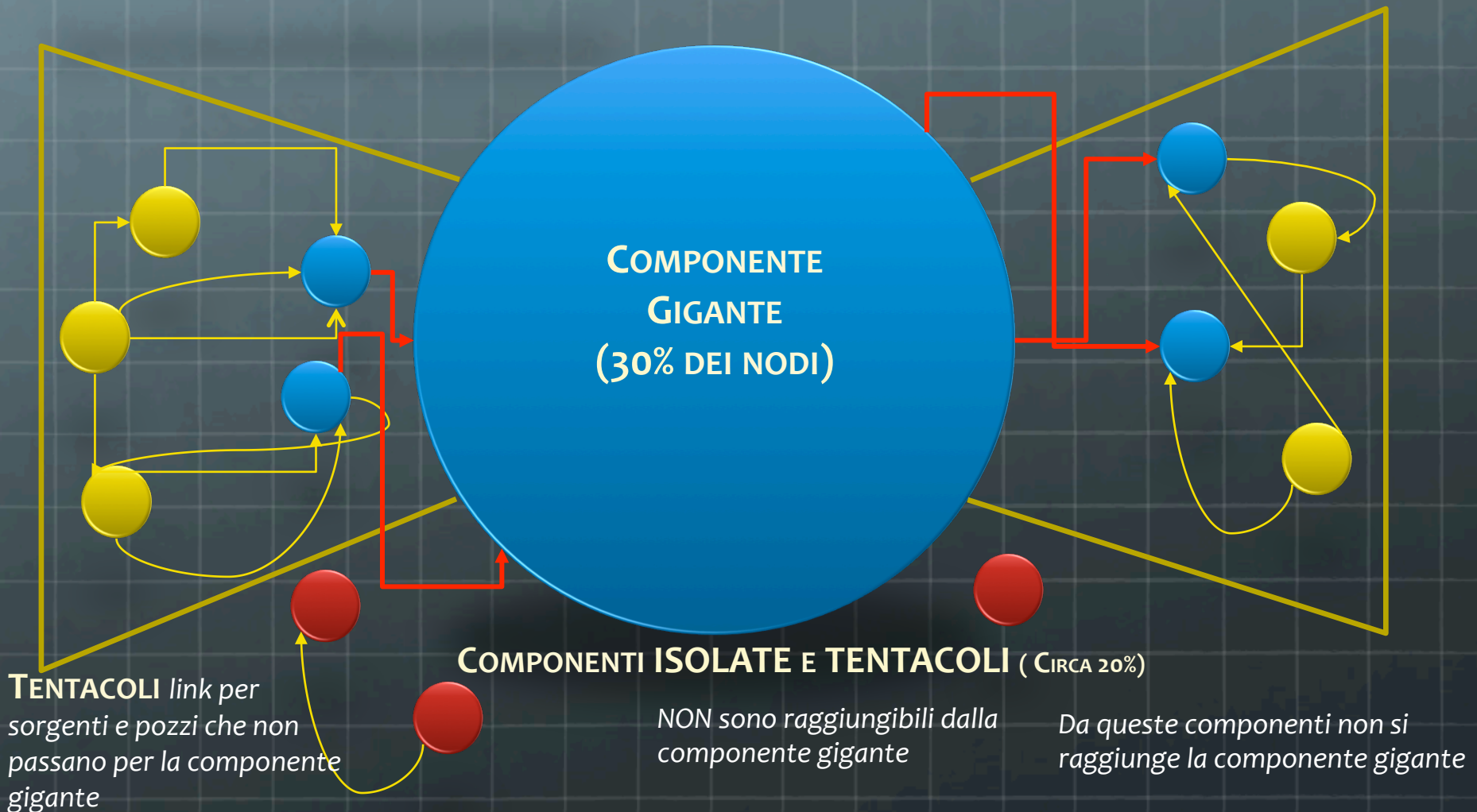
Algoritmo PageRank

STRUTTURA A CAMELLA DEL WEB



Algoritmo PageRank

STRUTTURA A CAMELLA DEL WEB



Algoritmo PageRank

DIFFICOLTÀ DI RICERCA

80% DEGLI UTENTI USA MOTORI DI RICERCA

Quantità di informazioni
troppo elevata

Eterogeneità della
qualità e formato delle
informazioni

Rapida modifica delle
informazioni

Assenza di SEMANTICA e
STRUTTURA

Algoritmo PageRank

MOTORI DI RICERCA

AZIONI DEI MOTORI DI RICERCA

→ Raccolta dati

Elaborazione e catalogazione
dei dati raccolti

→ Elaborazione e risposta alle
interrogazioni (Query) degli utenti



Algoritmo PageRank

MOTORI DI RICERCA

RACCOLTA DATI

Raccolta del contenuto delle pagine Web (informazioni di tipo testuale ma anche immagini ([google immagini](#) e [documentazione](#)))

robots.txt

Si usa uno SPIDER O CRAWLER O ROBOT

SPIDER

googlebot

fast

scooter

mercator

Ask Jeeves

teoma_agent

ia_archiver

Slurp

Romilda

MOTORE DI RICERCA

Google

Fast - Alltheweb

Altavista

Altavista

Ask Jeeves

Teoma

!Alexa - Internet Archive

Yahoo

Facebook

Simulatore di Spider

Spider traps

PROBLEMI

Quantità di dati e larghezza di banda

Aggiornamento frequente delle pagine

Pagine nascoste (*pagine isolate della struttura a caramella*)

Mancanza di standard condivisi e rispettati

Algoritmo PageRank

MOTORI DI RICERCA

ELABORAZIONE E CATALOGAZIONE DEI DATI RACCOLTI

PARSING, (analisi: estrazione di informazioni)

Rilevazione delle ridondanze (presenza di MIRRORING)

INDICIZZAZIONE dei dati

Reperimento e analisi delle informazioni per il calcolo del **RANKING**



Rilevazione di presenza di SPAMMING

In sketch comico del [Monty Python's Flying Circus](#) che ha come luogo un locale nel quale ogni pietanza proposta dalla cameriera era a base di Spam (un tipo di carne in scatola).



Algoritmo PageRank

MOTORI DI RICERCA

ELABORAZIONE E RISPOSTA ALLE INTERROGAZIONI (QUERY) DEGLI UTENTI

Ricerche testuali raffinate **AND** (Pozzo **NEAR** Pizza),il forse cercavi...

→ Suggerimenti ontologici

Ontologia fondamentale o primitiva per ha come obiettivo quello di descrivere "ciò che esiste" secondo un insieme di entità ritenuto non ulteriormente definibile (vocabolari)

→ Sistemi di catalogazione automatica

→ Analisi linguistiche (frequenze, relazioni...)

→ Analisi dei profili utente (ad esempio i bookmarks)

I contenuti indicizzati sono classificati rispetto ai concetti definiti all'interno dell'ontologia



ALGORITMO PAGERANK

MOTORI DI RICERCA

ELABORAZIONE DEI DATI

obiettivo

→ INDICIZZAZIONE DEI DOCUMENTI RACCOLTI

esigenze

→ INDICIZZAZIONE → RENDE EFFICIENTE E VELOCE LA RISPOSTA ALLE QUERY

→ INDICIZZAZIONE → RENDE POSSIBILE IL **RANKING** DEI DOCUMENTI



Algoritmo PageRank

RANKING

RANKING

definizione

DATO UN INSIEME DI PAGINE P E UNA QUERY Q IL RANKING È DEFINITO DA UNA FUNZIONE:

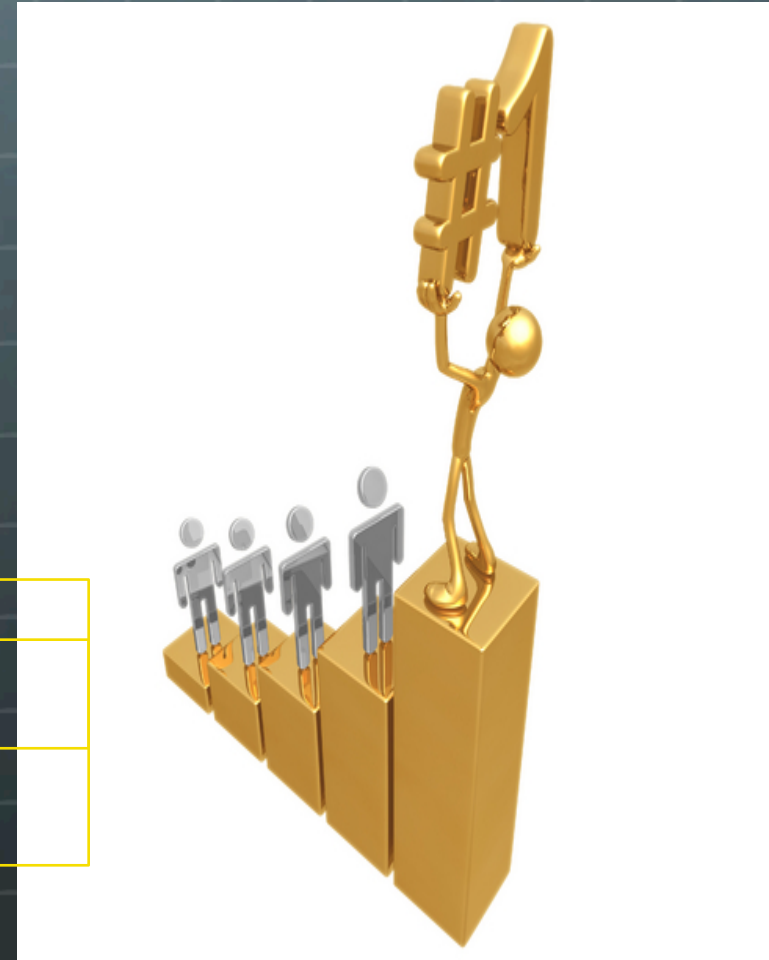
$$R_q: P \rightarrow \mathbf{R} \text{ (INSIEME DEI NUMERI REALI)}$$

CHE ASSOCIA AD OGNI PAGINA UN NUMERO REALE CHE INDICA LA «**RILEVANZA**» DI QUELLA PAGINA NEL CONTESTO DI QUELLA QUERY.

tecniche

Analisi del contenuto testuale
(ALTAVISTA)

Analisi della struttura dei link
(GOOGLE)



Algoritmo PageRank

RANKING

RANKING

Tecniche

LATENT SEMANTIC INDEX
Analisi del contenuto testuale
(ALTAVISTA)

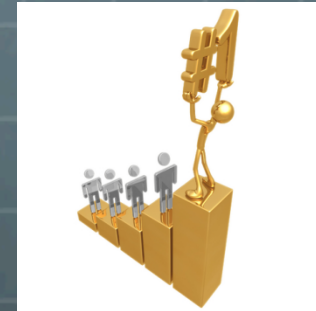
procedimento

t numero di termini presi in considerazione, appartenenti ad un vocabolario o individuati durante la raccolta delle pagine

Ad ogni pagina P è associato un vettore

con

$(d_p)_j = \text{numero di occorrenze del termine } j \text{ in } P$



Roma

Pioggia

t=2

Oggi la pioggia è stata
abbondante in tutta Italia.
Roma con la pioggia è piacevole
da visitare

Pagina P

$(d_p)_1 = 1$ numero di occorrenze di Roma in P

$(d_p)_2 = 2$ numero di occorrenze di Pioggia in P

Algoritmo PageRank

RANKING

RANKING

Tecniche

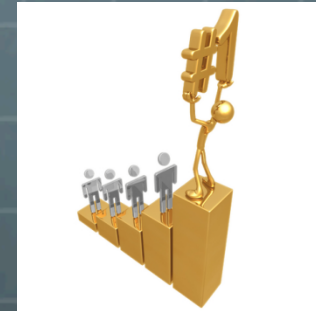
LATENT SEMANTIC INDEX
Analisi del contenuto testuale
(ALTAVISTA)

procedimento

Ad ogni query Q è associato un vettore

con

$(d_p)_j = 1$ se il termine j compare in P
 $(d_p)_j = 0$ se il termine j non compare in P



Roma

Pioggia

Oggi la pioggia è stata
abbondante in tutta Italia.
Roma con la pioggia è piacevole
da visitare

Pagina P

Pagine su Roma

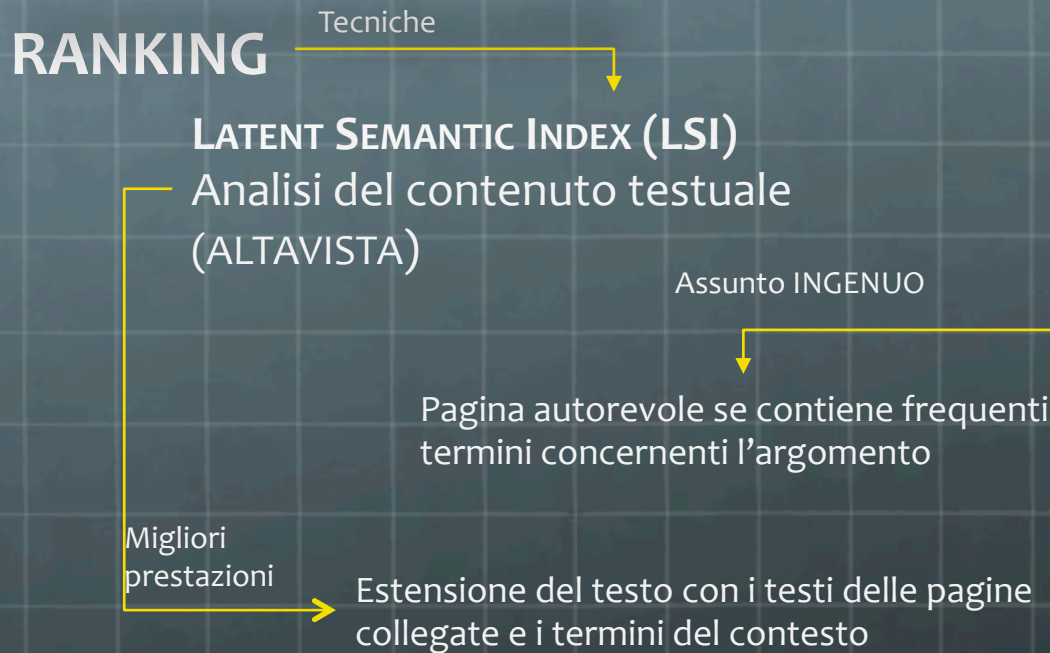
Query Q

$(d_p)_1 = 1$ Roma è nella query Q

$(d_p)_2 = 0$ Pioggia non è nella query Q

Algoritmo PageRank

RANKING



$$\cos d_P d_Q = \frac{\hat{d}_P \times \hat{d}_Q}{\|d_P\| \cdot \|d_Q\|}$$

CONSIDERAZIONE: LSI funziona bene su query multiple, le normali query sono semplici, 2 o 3 elementi al massimo.

Algoritmo PageRank

RANKING

RANKING

Tecniche

PAGERANK

Analisi dei link (GOOGLE)

procedimento

Ad ogni pagina j viene assegnato un valore reale, un **rank** R_j statico, indipendente cioè dalla query.

Data la query Q si ordinano i risultati in base al rank delle pagine individuate

L'importanza delle pagine è determinata **ESCLUSIVAMENTE** in relazione ai link che presenta o di cui è target. L'assunto è che il contenuto **NON** è **AUTODESCRITTIVO** e l'importanza di una pagina è il risultato di un processo **ESOGENO** (esterno all'ambito di riferimento)



Algoritmo PageRank

RANKING

RANKING

Tecniche

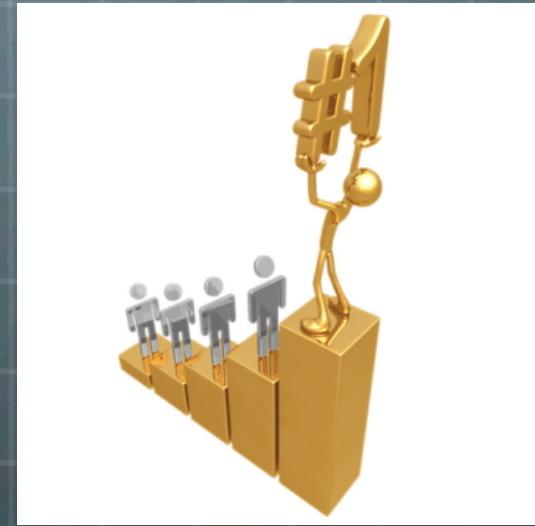
PAGERANK

Analisi dei link (GOOGLE)

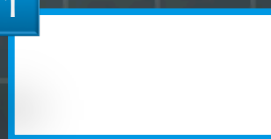
principio

L'importanza di una pagina è direttamente proporzionale al numero di pagine che la puntano

Se R_j è il rank (RANGO) di una pagina, la sua IMPORTANZA, la pagina trasmette l'importanza alle pagine che punta, distribuisce il suo rango in maniera uniforme



1



2



3



Algoritmo PageRank

RANKING

RANKING

Tecniche

PAGERANK

Analisi dei link (GOOGLE)

problema

Principio ispiratore

STOCASTICO (greco *stochastikós*, "congetturale", derivato di *stocházesthai*, "fare congetture«)

PROCESSO STOCASTICO una famiglia di variabili aleatorie dipendenti dal tempo, definite su un unico spazio campione finito e che assumono valori in un insieme definito spazio degli stati del processo

Processo **STOCASTICO MARKOVIANO** o (processo di Markov): *processo stocastico nel quale la probabilità di transizione che determina il passaggio ad uno stato di sistema dipende unicamente dallo stato del sistema immediatamente precedente (PROPRIETÀ DI MARKOV) e non dal come si è giunti a tale stato.*

DEFINIZIONE INTUITIVA: *un processo stocastico é un insieme ordinato di variabili casuali, indicizzate dal parametro t , spesso detto tempo. (Quantità di pioggia)*

La funzione ha un risultato unico solo se il grafo è connesso

soluzione

Si introduce un fattore che equivale a inserire link random (casuali) al grafo

d fattore di **dumping** (deciso da Google), fattore di spargimento, passa da una pagina all'altra ed é valore di PageRank minimo attribuito ad ogni pagina in archivio.

Nella documentazione originale $d=0,85$



Algoritmo PageRank

RANKING

- Caffeine, algoritmo 2010 (10 giugno)
- Hummingbird, algoritmo 2013 (27 settembre)

VANTAGGI

- *valore del RANK calcolato in maniera accurata*
- *processo iterativo converge molto rapidamente (tempi rapidi di esecuzione)*
- *RANK calcolato indipendentemente dalle query*

SVANTAGGI

- *si possono pubblicare insiemi di pagine «TRAPPOLA» «ARTEFATTE» che influiscono sul ranking delle pagine*
- *è un limite l'indipendenza dalla query*

Googlebombing

L'algoritmo di PageRank può essere corretto usando *aggiustato usando un secondo ranking basato sul contenuto, come ad esempio LSI (LATENT SEMANTIC INDEX).*