

UNIDAD 3: ANÁLISIS DE VÍNCULOS

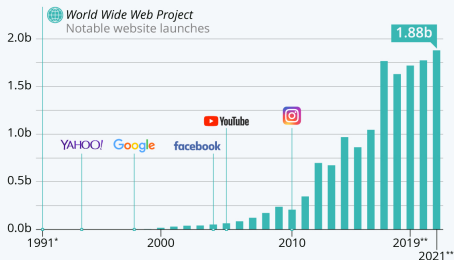
PAGE RANK SENSIBLE AL TÓPICO

Blanca Vázquez y Gibran Fuentes-Pineda

11 de octubre de 2022

How Many Websites Are There?

Number of websites online from 1991 to 2021



* As of August 1, 1991.

** Latest available data for 2019: October 28, for 2020: June 2, for 2021: August 6.

Source: Internet Live Stats



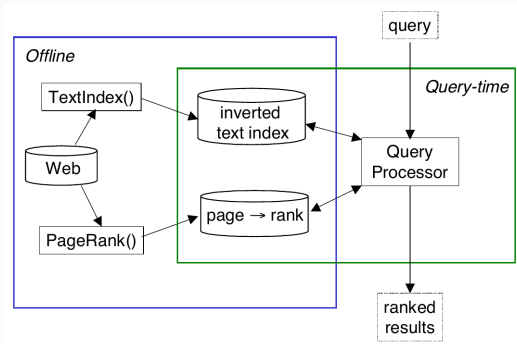
statista

Casi el 90 % del tráfico se encuentra en los buscadores.

CRÍTICAS A LAS SOLUCIONES EXISTENTES

	HITS	PageRank
Ventajas	<ul style="list-style-type: none">- Simple e iterativo- Puntuación específica de la consulta	<ul style="list-style-type: none">- Poco costoso (en tiempo de ejecución)- Las puntuaciones se calculan utilizando el grafo completo
Desventajas	<ul style="list-style-type: none">- Costoso (tiempo de ejecución)- Las puntuaciones se calculan utilizando un subgrafo a partir de todo el grafo.	<ul style="list-style-type: none">- La puntuación es independiente de la consulta- El algoritmo es propenso a manipulaciones (granjas de enlaces)

ALGORITMO PAGERANK

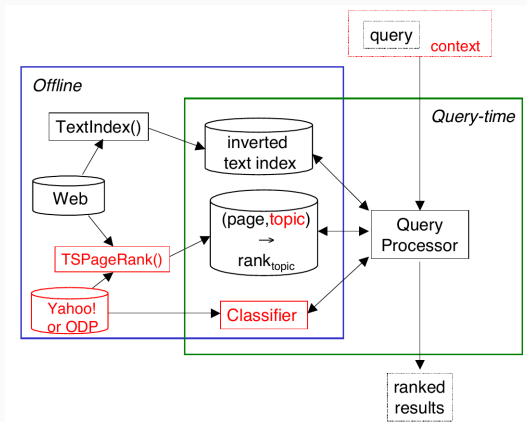


Motor de búsqueda usando el algoritmo de PageRank

Imagen tomada de Haveliwala, 2003.

- TSPR son las siglas de Topic-Sensitive PageRank
- Propuesto por Taher H. Haveliwala de la Universidad de Stanford en el 2003.
- Es la versión personalizada de *Page Rank*.
- En lugar de calcular un solo vector de rango, ¿por qué no calcular un conjunto de vectores de rango (uno por cada tópico)?

ALGORITMO PAGE RANK SENSIBLE AL TÓPICO



Motor de búsqueda usando el algoritmo de PageRank sensible al topico
Imagen tomada de Haveliwala, 2003.

- Supongamos que creamos un vector único para cada tópico usando PageRank.
- Si se pudiera determinar ¿cuál de estos tópicos son de interés para el usuario?, entonces :
 - Se podría usar el vector de Page Rank de ese tópico cuando se clasifiquen las páginas por relevancia.

- Usa el proyecto *Open Directory Project* como fuente de selección de tópicos [*https://dmoz-odp.org/*](https://dmoz-odp.org/)
 - También conocido como *DMoz* por *directory.mozilla.org*
- Es una colección de páginas web clasificadas por humanos.
- Consta de 16 tópicos (deportes, medicina,...)

Su formulación es similar a la de Page Rank:

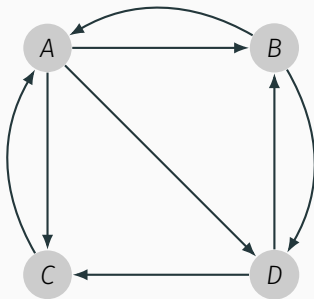
$$v' = \beta Mv + (1 - \beta)e_S/|S|$$

dónde:

- β es la probabilidad de elegir un vínculo de forma aleatoria
- M es la matriz de adyacencia
- v es el vector de Page Rank
- S indica la páginas que pertenecen a cierto tópico
- e_S es un vector que tiene 1s en los componentes S y 0s en el resto.
- $|S|$ es el tamaño del conjunto S .

EJEMPLO

Calcular el Page Rank sensible al t3pico, donde $\beta = 0.8$ y $S = \{B, D\}$

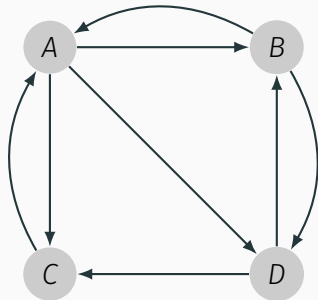


EJEMPLO

$$v' = \beta Mv + (1 - \beta)e_s/|S|$$

Paso 1: matriz de adyacencia

$$M = \begin{bmatrix} 0 & 1/2 & 1 & 0 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{bmatrix}$$



$$v' = \beta Mv + (1 - \beta)e_s/|S|$$

Paso 2: matriz de adyacencia * β

$$\begin{bmatrix} 0 & 1/2 & 1 & 0 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 0 & 0 & 1/2 \\ 1/3 & 1/2 & 0 & 0 \end{bmatrix} * 0.8 = \begin{bmatrix} 0 & 2/5 & 4/5 & 0 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 2/5 & 0 & 0 \end{bmatrix}$$

$$v' = \beta Mv + (1 - \beta)e_S/|S|$$

Paso 3: resolver $(1 - \beta)e_S/|S|$

$$(1 - 0.8) \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} / 2 = \frac{1}{5} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} / 2 = \begin{bmatrix} 0 \\ 1/10 \\ 0 \\ 1/10 \end{bmatrix}$$

$$v' = \beta Mv + (1 - \beta)e_S/|S|$$

Unimos los resultados previos

$$v' = \begin{bmatrix} 0 & 2/5 & 4/5 & 0 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 2/5 & 0 & 0 \end{bmatrix} v + \begin{bmatrix} 0 \\ 1/10 \\ 0 \\ 1/10 \end{bmatrix}$$

Cálculo de las primeras iteraciones: t_0

$$t_0 = \begin{bmatrix} 0 \\ 1/2 \\ 0 \\ 1/2 \end{bmatrix}$$

Recordemos, solo aplica en los nodos del conjunto S

EJEMPLO

Cálculo de las primeras iteraciones: t_1

$$t_1 = \begin{bmatrix} 0 & 2/5 & 4/5 & 0 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 2/5 & 0 & 0 \end{bmatrix} v + \begin{bmatrix} 0 \\ 1/10 \\ 0 \\ 1/10 \end{bmatrix}$$

$$t_1 = \begin{bmatrix} 0 & 2/5 & 4/5 & 0 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 0 & 0 & 2/5 \\ 4/15 & 2/5 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1/2 \\ 0 \\ 1/2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1/10 \\ 0 \\ 1/10 \end{bmatrix} = \begin{bmatrix} 1/5 \\ 3/10 \\ 1/5 \\ 3/10 \end{bmatrix}$$

Iteraciones

$$t_0 = \begin{bmatrix} 0/2 \\ 1/2 \\ 0/2 \\ 1/2 \end{bmatrix}, t_1 = \begin{bmatrix} 1/5 \\ 3/10 \\ 1/5 \\ 3/10 \end{bmatrix}, t_2 = \begin{bmatrix} 42/150 \\ 41/150 \\ 25/150 \\ 41/150 \end{bmatrix}, t_3 = \begin{bmatrix} 62/250 \\ 71/250 \\ 46/250 \\ 71/250 \end{bmatrix} \dots \begin{bmatrix} 54/210 \\ 59/210 \\ 38/210 \\ 59/210 \end{bmatrix}$$

¿CÓMO INTEGRAR TSPR AL BUSCADOR?

1. Decidir sobre los tópicos para crear vectores de Page Rank especializados
2. Encontrar una manera de determinar el tópico o los tópicos que sean más relevantes
3. Usar los vectores de Page Rank de esos tópicos para responder la consulta del usuario.

¿CÓMO IDENTIFICAR LOS TÓPICOS?

- Permitir que el usuario los seleccione (usando un menú)
- Inferir los tópicos usando:
 - Las búsquedas previas del usuario.
 - La información del usuario (marcadores, Facebook).

- Ejemplo: las palabras **sarampión** y **gol** aparecen frecuentemente en las páginas web:
 - **Sarampión** — — — $> T_{medicina}$
 - **Gol** — — — $> T_{deportes}$
- Supongamos que identificamos las palabras más frecuentes de cada página.
- Supongamos que tomamos un conjunto de páginas especializadas de un cierto tópico, y extraemos las palabras más frecuentes.

- Sea $S_1, S_2 \dots S_k$ son el conjunto de palabras que definen cada tópico.
- Sea P el conjunto de palabras que aparecen en una página p .
- Calcular la medida de similitud de Jaccard entre P y en cada uno de S_i .
- Clasificar la página al tópico con mayor similitud.

INFERIR TÓPICOS BASADO EN PALABRAS

computer vision	
COMPUTERS	0.24
BUSINESS	0.14
REFERENCE	0.09

gardening	
HOME	0.63
SHOPPING	0.14
REGIONAL	0.04

java	
COMPUTERS	0.53
GAMES	0.10
KIDS & TEENS	0.06

national parks	
REGIONAL	0.42
RECREATION	0.16
KIDS & TEENS	0.09

cruises	
RECREATION	0.65
REGIONAL	0.18
SPORTS	0.04

graphic design	
COMPUTERS	0.36
BUSINESS	0.23
SHOPPING	0.09

lipari	
HOME	0.19
KIDS & TEENS	0.17
NEWS	0.13

parallel architecture	
COMPUTERS	0.70
SCIENCE	0.10
REFERENCE	0.07

death valley	
REGIONAL	0.28
SOCIETY	0.14
NEWS	0.10

gulf war	
SOCIETY	0.21
KIDS & TEENS	0.18
REGIONAL	0.17

lyme disease	
HEALTH	0.96
REGIONAL	0.01
RECREATION	0.01

recycling cans	
HOME	0.42
BUSINESS	0.38
KIDS & TEENS	0.06

Imagen tomada de Haveliwala, 2003.