

# Análisis de Vínculos

Alondra Berzunza



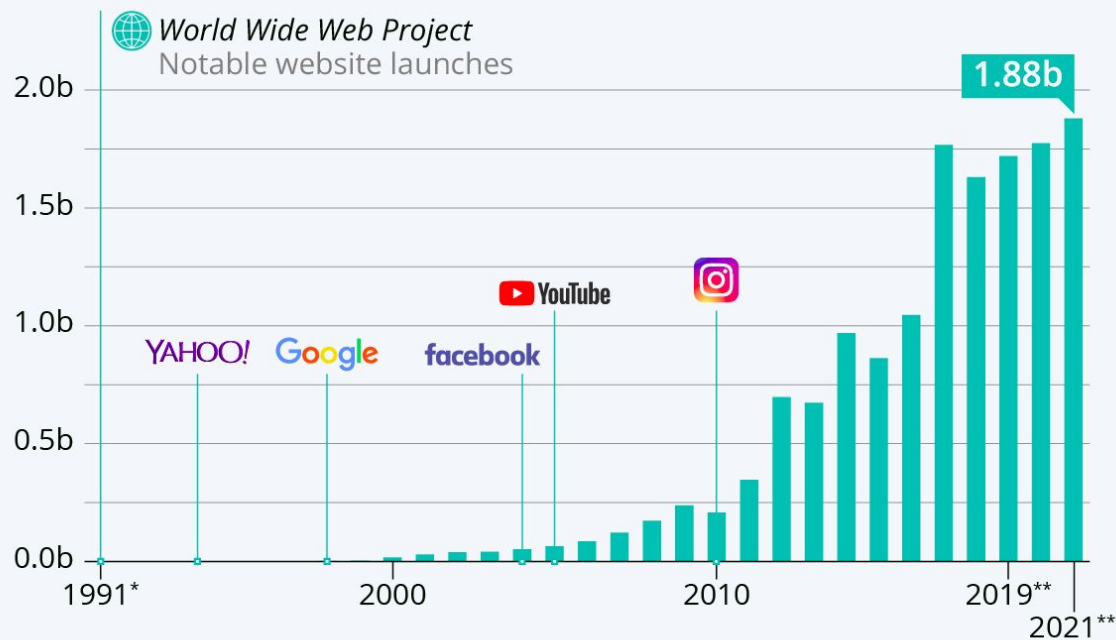


# Page Rank Sensible al Tópico



# How Many Websites Are There?

Number of websites online from 1991 to 2021



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

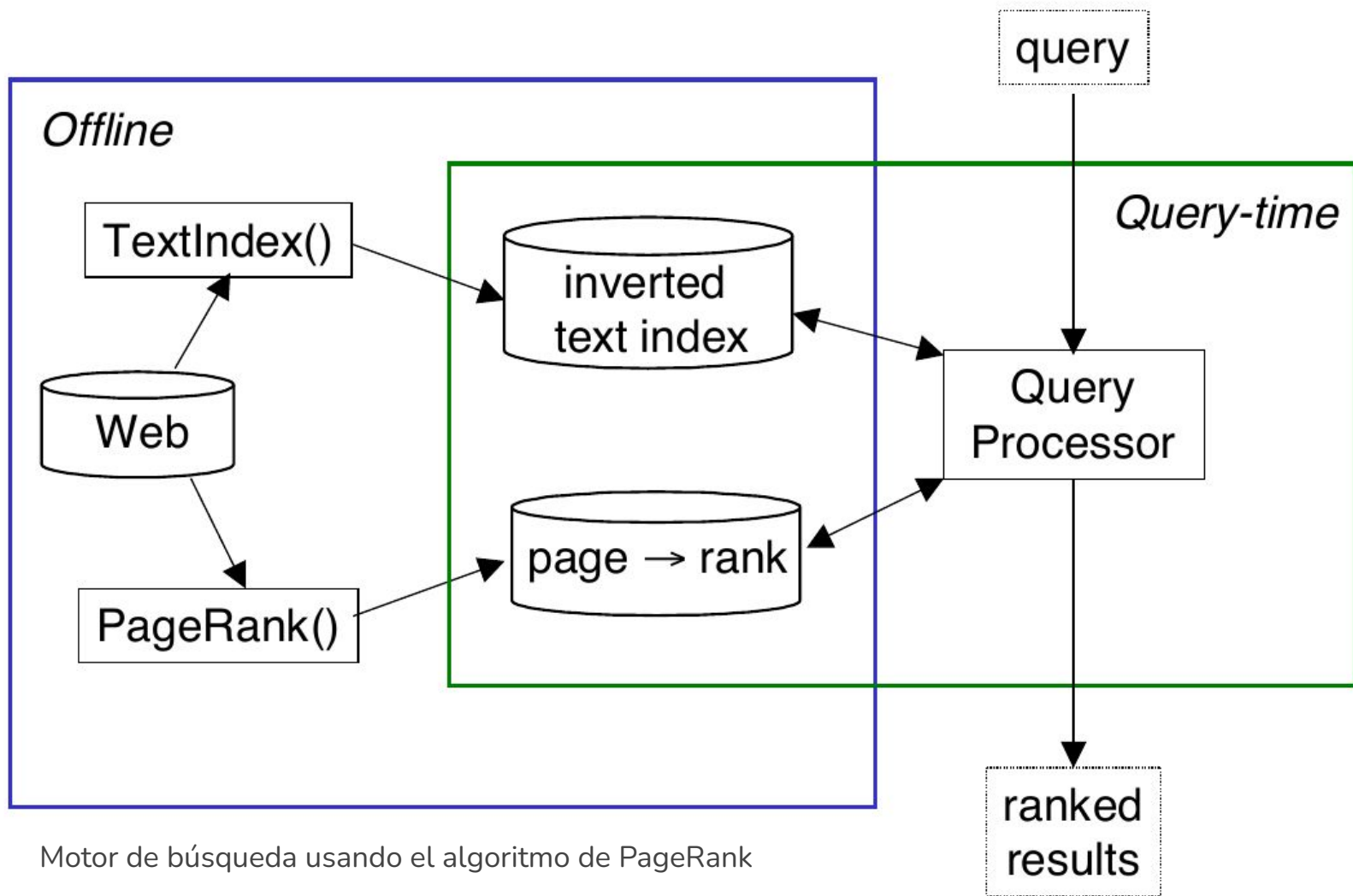
\* As of August 1, 1991.

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

Source: Internet Live Stats



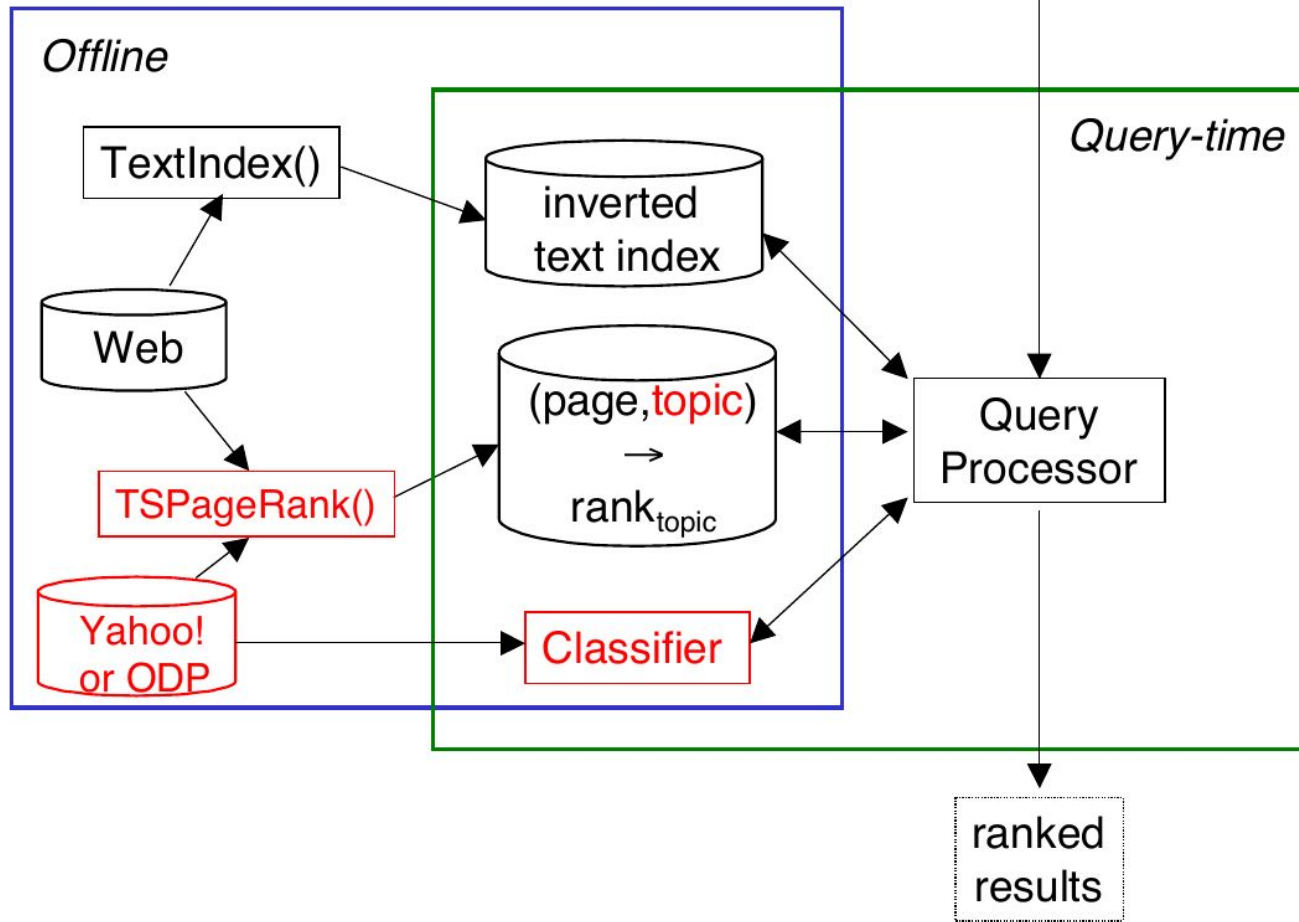
	HITS	Page Rank
<b>Ventajas</b>	<ul style="list-style-type: none"><li>★ Simple e iterativo.</li><li>★ Puntuación específica de la consulta.</li></ul>	<ul style="list-style-type: none"><li>★ Poco costoso (en tiempo de ejecución).</li><li>★ Las puntuaciones se calculan utilizando el grafo completo.</li>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





# Page Rank Sensible al Tópico

- 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)?



Motor de búsqueda usando el algoritmo de PageRank sensible al tópico.



# TSPR

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.

Una variante de PageRank donde la importancia está sesgada por los vectores de temas, mejorando las clasificaciones contextuales.





# Caminatas Aleatorias Sesgadas

Su formulación es similar a la de Page Rank:

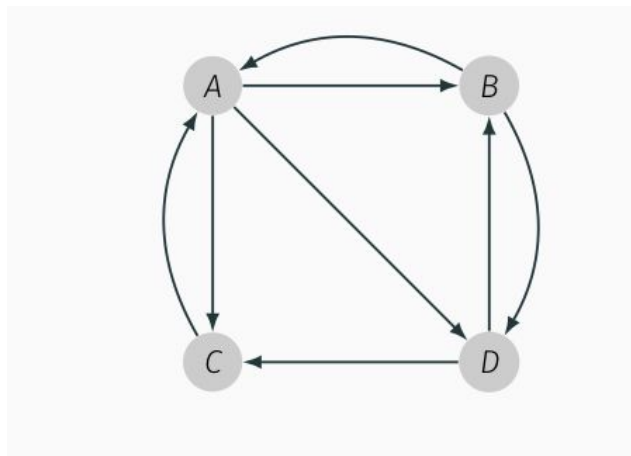
dónde:

- $\beta$  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, d3nde  $\beta = 0.8$  y  $S = \{B, D\}$





# 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.



# 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).



# Inferir Tópicos Basado en Palabras

Ejemplo: las palabras sarampión y gol aparecen frecuentemente en las páginas web:

- Sarampión – – – > T medicina
- Gol – – – > T deportes

Identificamos las palabras más frecuentes de cada página.

Tomamos un conjunto de páginas especializadas de un cierto tópico, y extraemos las palabras más frecuentes.



# Inferir Tópicos Basado en Palabras

- Sea  $S_1, S_2 \dots S_k$  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 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