



DEPARTAMENTO  
DE COMPUTACION

Facultad de Ciencias Exactas y Naturales - UBA

## Trabajo Práctico 2

Tirate un qué, tirate un *ranking*...

---

Métodos Numéricos  
Segundo Cuatrimestre de 2014

Integrante	LU	Correo electrónico
Fosco, Martin Esteban	449/13	mfosco65@gmail.com
Minces Müller, Javier Nicolás	231/13	javijavi1994@gmail.com
Chibana, Christian Ezequiel	586/13	christian.chiba93@gmail.com



Facultad de Ciencias Exactas y Naturales  
Universidad de Buenos Aires

Ciudad Universitaria - (Pabellón I/Planta Baja)

Intendente Güiraldes 2160 - C1428EGA

Ciudad Autónoma de Buenos Aires - Rep. Argentina

Tel/Fax: (54 11) 4576-3359

<http://www.fcen.uba.ar>

## Resumen

En este trabajo se analizan tres métodos para crear un *ranking* de páginas web: PageRank, HITS e In-Deg. Estos algoritmos utilizan matrices para modelar la Web y utilizan métodos numéricos, como el método de la potencia, para ordenar las páginas a partir de los links que las relacionan.

## Índice

<b>1. Introducción Teórica</b>	<b>3</b>
<b>2. Desarrollo</b>	<b>5</b>
2.1. Criterio de almacenamiento . . . . .	5
2.2. PageRank . . . . .	6
2.2.1. Método de la Potencia . . . . .	7
2.3. Hits . . . . .	9
2.4. In-Deg . . . . .	10
2.5. Experimentos realizados . . . . .	10
<b>3. Resultados y Análisis</b>	<b>11</b>
3.1. Convergencia de Normas . . . . .	11
3.2. Tiempos de Cómputo . . . . .	13
3.3. Rankings obtenidos . . . . .	15
3.4. Redes pequeñas . . . . .	18
<b>4. Conclusiones</b>	<b>22</b>
<b>5. Apéndice A</b>	<b>24</b>

## 1. Introducción Teórica

El problema a resolver es encontrar una manera confiable de definir en qué páginas es conveniente comprar espacios de propaganda para garantizar la mayor publicidad posible a un grupo de música.

El objetivo de este trabajo, por tanto, es resolver este problema analizando diferentes métodos que nos permitan definir un ranking de importancia de páginas web.

Para trabajar con HITS, se debe comenzar por acordar un método de selección de las páginas web sobre las cuales se trabaja, es decir que es necesario adquirir una subred de páginas relevantes al tema en cuestión.

Se recurre en primer lugar a un buscador textual, pero este no asegura un resultado confiable en cuanto a la relevancia de las páginas recibidas. En el paper escrito por Kleinberg <sup>1</sup> se sugiere un conjunto inicial de páginas  $S_\sigma$  que cumple las siguientes condiciones:

- $S_\sigma$  es relativamente pequeña.
- $S_\sigma$  contiene una gran cantidad de páginas relevantes.
- $S_\sigma$  contiene muchas de las autoridades de más peso.

El autor propone una forma de ampliar este conjunto inicial para asegurar que se cumpla la tercera condición: buscar autoridades de gran peso a las que haya links en la subred obtenida del buscador e incluir éstas también en  $S_\sigma$ .

Se construye luego un mapa de este conjunto de páginas en forma de grafo dirigido, observando los links entre ellas. Hasta este punto, el proceso escapa al objetivo del trabajo, que considera dado este mapa.

Terminado este paso, se procede a construir la matriz de conectividad  $W \in n \times n$ , tal que  $w_{ij} = 1$  si existe un link de la página  $j$  a la página  $i$  y  $w_{ij} = 0$  en caso contrario.

Usando esta matriz (con la misma notación), se pueden aplicar distintos métodos para procesar y definir un cierto *ranking* de las páginas. Esta matriz es esparsa, porque incluso en redes de varios cientos de miles de páginas, en promedio cada una tiene 7 u 8 links salientes. <sup>2</sup> Esto debe ser tenido en cuenta a la hora de guardar la matriz en memoria.

Los métodos a analizar son:

**PageRank**, el más conocido de los utilizados por Google. En este buscador se aplica el algoritmo al conjunto de todas las páginas de la Web, con un tiempo de cómputo del orden de un mes. Al buscar, se retornan las páginas con mejor puntaje que coincidan con la búsqueda. En este trabajo se limitará el conjunto de páginas, utilizando el mismo que para HITS.

PageRank asigna un puntaje a cada página según el puntaje de las que apuntan a ella. De esta forma se evita que una página con muchos links de páginas sin valor obtenga un puntaje demasiado alto. Además, modifica  $W$  para tener en cuenta que un *navegante aleatorio* puede con una determinada probabilidad saltar a cualquier página de la web, algo que también hace si encuentra una página sin links salientes.

Esto se traduce como buscar un vector  $x$  tal que  $Wx = x$ , es decir, encontrar el autovector para el autovalor 1. Si se sabe que 1 es el autovalor de mayor módulo, se puede encontrar su autovector con el método de las potencias. Para esto se requiere modificar  $W$  para que cumpla ciertas condiciones, que se explican en la próxima sección.

**HITS** (Hyperlink-Induced Topic Search), este se basa en la existencia de páginas de gran valor como 'autoridad' y como "hub". En cada iteración, se actualiza el puntaje de cada página como autoridad a partir del puntaje hub de las páginas que la citan. Luego, se actualiza el puntaje como hub de cada

---

<sup>1</sup>Kleinberg - 1999 - *Authoritative sources in a hyperlinked environment*

<sup>2</sup>Kamvar, Haveliwala - 2003 - *Extrapolation methods for accelerating PageRank computations*

página a partir del puntaje como autoridad de las páginas a las que apunta. Esto se puede reducir a una operación matricial.

**In-Deg**, El tercer método, el más intuitivo, se incluye principalmente como control. In-Deg consiste simplemente en contar los links que apuntan a cada página.

## 2. Desarrollo

### 2.1. Criterio de almacenamiento

Debimos primero encontrar una forma de almacenar los datos de los links entre páginas provistos (independientemente del formato en que fueron dados). Estos se guardan en la matriz  $W$ , ya mencionada.

Entre las 3 opciones sugeridas en el enunciado se optó finalmente por el almacenamiento en forma de Compressed Sparse Column (CSC)<sup>3</sup> porque notamos que esto facilita aquellas operaciones que pueden ser necesarias para los algoritmos de mayor complejidad temporal (Pagerank y HITS), algunos ejemplos de estas son:

- La suma de los elementos de una columna (ya que se trata de sumar una sección de elementos contiguos del array en el que se guardan los valores de la matriz distintos de 0).
- La división de una columna por una constante, debido a que al momento de definir en pagerank la probabilidad de salir de una página a otra, debemos dividir las columnas de forma que sus elementos (las probabilidades de salir hacia alguna página) sumen 1.
- El producto de una matriz transpuesta por un vector, ya que puede recorrer una sección del vector de valores y multiplicarlos por una posición del vector. Para el producto de la matriz sin transponer el algoritmo es más complejo pero tiende a ser igual de eficiente: alcanza con recorrer los valores e ir acumulando en las distintas posiciones del vector resultado los valores parciales que se obtienen.

Además, tiene una eficiencia espacial óptima para el almacenamiento de matrices que cuentan con muchos ceros (como es el caso de las matrices con las que trabajamos). Más específicamente : la memoria usada es  $2 \times O(m) + O(n)$ .<sup>4</sup>

En comparación con las otras dos opciones de almacenamiento:

**Dictionary of Keys** cuenta con sólo la última ventaja del CSC, que permite obviar a las posiciones en las que se guardan ceros, pero no es capaz de mejorar en ninguna manera la eficiencia de las operaciones matriciales, que se efectúan en los métodos de ranking más complicados. Una de las desventajas más sobresalientes del Dictionary of Keys es su poca óptima complejidad espacial, ya que al guardar varios elementos de una columna en particular debemos guardar esta columna una vez por cada uno de estos elementos, es decir  $m$  veces mas. La memoria usada por el *DoK* es  $3 \times O(m) + O(n)$ .

**Compressed Sparse Row (CSR)** cuenta con la última ventaja del CSC también. Optamos por la *Compressed Sparse Column* debido a que facilita la optimización de las operaciones por columnas que intervienen en los algoritmos mas complejos (es posible modificar, sin embargo, las operaciones entre columnas del CSR para que se puedan realizar con una eficiencia análoga a las efectuadas por el CSC).

---

<sup>3</sup>Una matriz esparsa que se recorre por columnas.

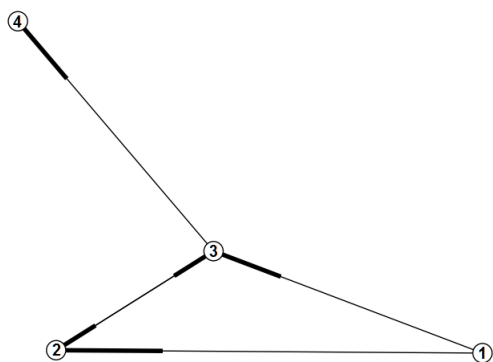
<sup>4</sup> $m$  es la cantidad de elementos diferentes de 0 y  $n$  las filas de la matriz.

## 2.2. PageRank

El algoritmo PageRank asigna un puntaje de importancia a cada página entre 0 y 1. Este no depende únicamente de la cantidad de links entrantes, sino que pondera cada link según el puntaje de la página al que pertenece. Sea  $x$  el vector en  $\mathbb{R}^n$  donde  $x_j$  es el puntaje de la página  $j$  y  $n_j$  el *grado* de  $j$ , es decir, la cantidad de links salientes, entonces se quiere ver que:

$$x_k = \sum_{j=1, w_{kj}=1}^n \frac{x_j}{n_j}, \quad k = 1, \dots, n.$$

Supongamos que se tienen cuatro páginas para hacer el ranking (1, 2, 3, 4) y que están conectadas de esta forma:



Entonces el problema se reduciría a buscar la solución a este sistema de ecuaciones:

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ \frac{1}{2} & 1 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

Llamamos  $P$  a la matriz, que se construye dividiendo cada columna de  $W$  por su suma. Se puede ver que la solución de este sistema de ecuaciones es un vector  $x$  tal que  $Px = x$ .

Se define un autovalor  $\lambda$  de una matriz  $A$  a un valor tal que  $Ax = \lambda x$ , y a  $x$  como su autovector. En este caso, entonces, el problema se reduce a buscar un autovector para el autovalor 1. Dado que este autovector representa probabilidades, debe tener norma 1 igual a 1.

El método de la potencia permite encontrar el autovalor de mayor módulo de una matriz. Para poder aplicarlo en este caso, deben cumplirse las siguientes condiciones <sup>5</sup>:

- $P$  tiene un autovector asociado al autovalor 1
- Los demás autovalores de la matriz cumplen  $1 = \lambda_1 > |\lambda_2| \geq \dots \geq |\lambda_n|$
- La dimensión del autoespacio asociado al autovalor  $\lambda_1$  es 1: esto garantiza la unicidad del autovector de norma 1 igual a 1 asociado al autovalor 1.

Si la matriz  $P$  es estocástica por columnas, se puede ver que se cumple las primeras dos condiciones: Sea  $A$  una matriz en  $\mathbb{R}^{n \times n}$  estocástica por columnas y sea  $e$  un vector en  $\mathbb{R}^n$  tal que  $e_i = 1 \forall i$ , sabemos que  $A$  y su transpuesta  $A^T$  tienen los mismos autovalores. Como  $A$  es estocástica por columnas, se puede ver que  $A^T e = e$ , por lo que 1 es un autovalor de  $A^T$  y, por tanto, de  $A$ . Además al estar dividida cada columna por su norma 1, la matriz es efectivamente estocástica por columnas, salvo que haya columnas cuya suma es 0. Esto se da en el caso de los llamados "dangling nodes", páginas sin links salientes.

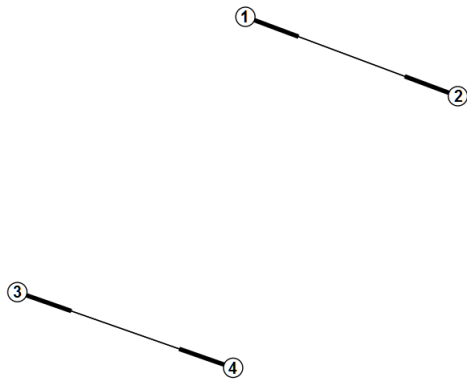
Se puede interpretar la navegación en la Web de un *navegante aleatorio* como un proceso de Markov, a  $P$  como su matriz de transición y a la componente  $x_j$  del vector solución de norma 1 del sistema  $Px = x$  como la proporción del tiempo que el navegante pasa en la página  $j$ . Entonces, se puede suponer que al encontrarse en una página sin links salientes (*dangling node*), irá a cualquiera con probabilidad

<sup>5</sup>Propuestas en Bryan, Leise - 2006 - *The Linear Algebra behind Google* \* y Kamvar, Haveliwala - 2003 - *Extrapolation methods for accelerating PageRank computations*

$1/n$ . Se define entonces la matriz  $P_1$ , en la que se asigna  $1/n$ , entonces, a cada fila de cada columna que representa un dangling node. En el ejemplo anterior:

$$P_1 = \begin{bmatrix} 0 & 0 & 0 & \frac{1}{4} \\ \frac{1}{2} & 0 & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{2} & 1 & 0 & \frac{1}{4} \\ 0 & 0 & \frac{1}{2} & \frac{1}{4} \end{bmatrix}$$

Sin embargo, la última condición no se cumple necesariamente. Si se tienen dos subredes aisladas, entonces la dimensión del autoespacio será 2. Por ejemplo, sean las páginas (1,2,3,4), con los siguientes links:



Su matriz de conectividad es:

$$W = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

Se puede ver que existen dos vectores linealmente independientes de norma 1 tales que  $Wx = x$ ,  $x = (\frac{1}{2}, \frac{1}{2}, 0, 0)$  y  $x = (0, 0, \frac{1}{2}, \frac{1}{2})$

Para evitar esto, se agrega un parámetro  $c > 0$ , donde  $1 - c$  representa la probabilidad de que un navegante aleatorio vaya a una página cualquiera. Se define entonces  $P_2$  tal que  $P_{2ij} = cP_{1ij} + \frac{1-c}{n} \forall i$

$P_2$  cumple las tres condiciones, luego se puede buscar la solución de  $P_2x = x$  con el método de la potencia. Sin embargo, no es esparsa; de hecho, no tiene ceros.

En el último ejemplo, con  $c=4/5$ :

$$W = \begin{bmatrix} \frac{1}{20} & \frac{17}{20} & \frac{1}{20} & \frac{1}{20} \\ \frac{17}{20} & \frac{1}{20} & \frac{1}{20} & \frac{1}{20} \\ \frac{1}{20} & \frac{1}{20} & \frac{1}{20} & \frac{17}{20} \\ \frac{1}{20} & \frac{1}{20} & \frac{17}{20} & \frac{1}{20} \end{bmatrix}$$

Y el PageRank resultante sería:

- 1: 0,25
- 2: 0,25
- 3: 0,25
- 4: 0,25

### 2.2.1. Método de la Potencia

Este método permite encontrar al autovalor de mayor módulo de una matriz.

Sea  $A \in \mathbb{R}^{n \times n}$ ,  $\lambda_1 \dots \lambda_n$  sus autovalores tales que  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$ ,  $v_1 \dots v_n$  una base de autovectores y  $x$  un vector cualquiera.  $x$  se puede escribir como combinación lineal de los autovectores de  $A$ :

$$\begin{aligned} x &= \sum_{j=1}^n \alpha_j v_j \\ Ax &= \sum_{j=1}^n \alpha_j A v_j \end{aligned}$$

Como  $v_j$  es autovector de  $A$ , entonces:

$$\begin{aligned} Ax &= \sum_{j=1}^n \alpha_j \lambda_j v_j \\ A^2 x &= \sum_{j=1}^n \alpha_j \lambda_j^2 v_j \\ A^2 x &= \sum_{j=1}^n \alpha_j \lambda_j^2 v_j \\ &\vdots \\ A^k x &= \lambda_1^k \cdot \left( \sum_{j=1}^n \frac{\lambda_j^k}{\lambda_1^k} v_j \right) \\ A^k x &= \lambda_1^k \cdot \left( \alpha_1 v_1 \sum_{j=2}^n \frac{\lambda_j^k}{\lambda_1^k} v_j \right) \\ \frac{\|A^k x\|}{\|A^{k-1} x\|} &= \frac{\|\lambda_1^k \cdot (\alpha_1 v_1 \sum_{j=2}^n \frac{\lambda_j^k}{\lambda_1^k} v_j)\|}{\|\lambda_1^{k-1} \cdot (\alpha_1 v_1 \sum_{j=2}^n \frac{\lambda_j^{k-1}}{\lambda_1^{k-1}} v_j)\|} \end{aligned}$$

cuando  $k \rightarrow \infty$ ,  $\frac{\lambda_j}{\lambda_1} \rightarrow 0$  porque  $\lambda_1 > \lambda_j \forall j$ , luego

$$\frac{\|A^k x\|}{\|A^{k-1} x\|} \rightarrow \frac{\|\lambda_1^k \cdot (\alpha_1 v_1)\|}{\|\lambda_1^{k-1} \cdot (\alpha_1 v_1)\|} = \lambda_1$$

$$A^k x = ((\dots((AA)A)\dots A)A)x = A(A(A(\dots(A(Ax))))$$

Es decir que se puede implementar un algoritmo que calcule  $Ax$  y asigne el resultado a  $x$ , siendo este el autovector buscado.

Para el contexto particular del PageRank, utilizamos el algoritmo 1 propuesto por Kamvar<sup>6</sup>, que permite calcular  $P_2 x$  a partir de  $Px$ . Esto permite ahorrar tiempo y memoria, ya que  $P$  es una matriz esparsa. El algoritmo consiste de tres pasos:

$$\begin{aligned} y &= c(Px) \\ w &= \|x\|_1 - \|y\|_1 \\ y &= y + wv \end{aligned}$$

Donde  $v$  es un vector en  $\mathbf{R}^n$  tal que  $v_i = \frac{1}{n} \forall i$

Es decir,  $y_i = (cPx)_i + \frac{(\|x\|_1 - c\|Px\|_1)}{n}$

<sup>6</sup>Kamvar, Haveliwala - 2003 - Extrapolation methods for accelerating PageRank computations



recordemos que

$$(P_2x)_i = \sum_j P_{2ij}x_j$$

y que

$$P_{2ij} = cP_{ij} + \frac{(1-c)}{n}$$

$$(P_2x)_i = \sum_{j=1}^n (cP_{ij} + \frac{1-c}{n})x_j$$

$$(P_2x)_i = \sum_{j=1}^n cP_{ij}x_j + \sum_{j=1}^n \frac{(1-c)x_j}{n}$$

$$(P_2x)_i = c(Px)_i + \frac{1-c}{n} \sum_{j=1}^n x_j$$

$$\text{Como } x_i > 0 \forall i, \sum_{j=0}^n x_i = \|x\|_1$$

$$(P_2x)_i = c(Px)_i + \frac{(1-c)\|x\|_1}{n}$$

Como  $P$  es estocástica por columnas si no tiene dangling nodes:

$$\|Px\|_1 = \|x\|_1$$

$$(1-c)\|x\|_1 = \|x\|_1 - c\|Px\|_1$$

$$(P_2x)_i = c(Px)_i + \frac{(1-c)\|x\|_1}{n} = c(Px)_i + \frac{\|x\|_1 - c\|Px\|_1}{n} = y_i$$

## 2.3. Hits

La idea detrás de este método es, como se ha explicado brevemente antes, realizar un análisis de la red y sus páginas viéndolas como hubs y autoridades, nociones que son definidas por las relaciones (links) que tienen con el resto de la red.

Dada una búsqueda, el objetivo es retornar un subconjunto de páginas relevantes. Se asigna entonces sobre las páginas el rol de *autoridad* en el tema que se busca. Su valor como *autoridad* está definido por el valor de *hub* de las páginas que apuntan a ella mientras que el valor o peso de *hub* está definido, a su vez, por el valor de *autoridad* de las páginas a las que apunta.

Recordemos que los links están almacenados en la matriz de conectividad  $W$  donde  $w_{ij} = 1$  si existe un link de la página  $j$  a la página  $i$ ,  $w_{ij} = 0$  en caso contrario. Para cada página  $i \in Web$  se define como peso de autoridad  $x_i$  y como peso de hub  $y_i$ . A partir de esto, podemos definir los vectores  $x, y \in \mathbb{R}^n$  como los vectores de pesos de autoridad y hubs, y se puede suponer que se encuentran normalizados. Evidentemente, se desprende de esto que las páginas con mayores  $x_i$  o  $y_i$  son consideradas las mejores *autoridades* o *hubs*.

Se puede expresar la transferencia de pesos de hubs y autoridad, respectivamente, de la siguiente manera:

$$\text{Peso de Autoridad } x_j = \sum_i w_{ij}y_i$$

$$\text{Peso de Hub } y_i = \sum_j w_{ij}x_j$$

Consecuentemente, si expresamos con operaciones matriciales los vectores de hub y autoridad:

$$x = Wy$$

$$y = W^t x$$

Se aplica luego el paso de normalización  $y \leftarrow \frac{y}{\|y\|_2}$  y se vuelve a comenzar hasta llegar a la convergencia. Para la primera iteración es posible empezar con un  $y_0$  inicial cualquiera, y el método convergerá si se cumplen las siguientes condiciones:

- $W$  tiene un autovalor  $\lambda_1$  tal que  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$
- $y_0$  no es ortogonal al autovector asociado a  $\lambda_1$

El autor recomienda tomar  $y_{0i} = 1 \forall i$ . En este trabajo se tomó  $y_{0i} = \frac{1}{n} \forall i$ .

## 2.4. In-Deg

Este método es considerablemente más simple y sigue una idea más intuitiva de ranking de páginas sin llegar a entrar en complejizaciones para tener en cuenta distintos casos, lo que consideramos que lleva a unos resultados menos fiables. Esto se debería notar especialmente cuando son contrastados con aquellos que obtenemos de analizar redes con métodos que dan más peso a detalles que potencialmente cambiarían los resultados (por ejemplo: una página apuntada por pocas páginas de gran valor caería muy por debajo en el ranking, mientras que probablemente PageRank y HITS le asignarían una posición más alta).

La idea básica es definir el valor de una página según la cantidad de páginas que la apuntan.

En otras palabras: el método consiste en sumar los links a cada página. Teniendo en cuenta que ya contamos con la matriz  $W$ , este método se reduce a sumar la fila correspondiente a cada página y obtener así un vector que nos muestra el valor de cada una (la  $i$ -ésima posición corresponde a la  $i$ -ésima página, según su posición asignada en la matriz  $W$ ).

## 2.5. Experimentos realizados

Para los dos primeros métodos, se realizaron experimentos para ver la convergencia de la norma del autovector buscado en cada iteración. En PageRank, se utilizó la norma 1 del autovector  $x$ . En el caso de HITS, se analizaron los vectores  $x$  e  $y$  tomando norma 2. Luego, se comparó el tiempo de ejecución de cada uno utilizando la librería `time.h` de C++. No se tomó en cuenta el tiempo que toma cargar la matriz. Como una iteración de HITS hace dos productos de una matriz por un vector, mientras que PageRank hace uno, y que esta operación es la que domina el costo temporal de la ejecución, se esperaba que el tiempo de ejecución de HITS fuera alrededor del doble del de PageRank, independientemente del valor de  $c$  usado. Este influye en el tiempo de ejecución<sup>7</sup> pero no esperábamos que modificara sustancialmente esta relación.

Dado que indexamos la matriz desde 1, reemplazamos cada página numerada como 0 por  $n + 1$  para poder realizar los tests.

Las instancias de test usadas fueron las siguientes:

$\sigma$	nodos	links	$\frac{\text{link}}{\text{nodo}}$
Movie	5.757	24.451	4,24
Censorship	2.947	9.555	3,24
Genetic	3.468	12.689	3,65
NotreDame	325.729	1.497.134	4,59
Stanford	281.903	2.312.497	8,20

Los valores de  $c$  usados para el pageRank fueron los siguientes:

- 0.5
- 0.7
- 0.85
- 0.95

Se toma 0.85 por ser el valor recomendado por los autores de *Brin, Page - The anatomy of a large-scale hypertextual Web search engine*. También se tomaron valores mayores y menores que 0.85 para compararlos con este. La tolerancia para asumir convergencia la fijamos en  $10^{-5}$ .

A continuación se experimentó con una muestra de links entre varias páginas obtenida mediante el script provisto por la cátedra. De esta forma esperábamos analizar los resultados obtenidos en función de qué páginas se esperaban con mayor ranking. Por último, se realizaron experimentos con redes pequeñas, para mostrar el comportamiento esperado de PageRank y HITS.

---

<sup>7</sup>Kamvar, Haveliwala - 2003 - Extrapolation methods for accelerating PageRank computations

### 3. Resultados y Análisis

#### 3.1. Convergencia de Normas

A continuación se presentan los resultados obtenidos al experimentar con Pagerank (con  $c = 0.50, 0.70, 0.85$  y  $0.95$ ), observando la convergencia de la Norma Manhattan hacia un valor que cumpla con la tolerancia definida, y comparando la cantidad de iteraciones necesarias para llegar a dicho valor.

En los gráficos de PageRank se compara la evolución de  $\|x^{(k)} - x^{(k-1)}\|_1$  para distintos valores de  $c$ . Se decidió graficar la red Censorship al considerarla representativa de los otros casos probados para redes medianas, y Notre Dame al considerarlo representativo de los casos de redes grandes.

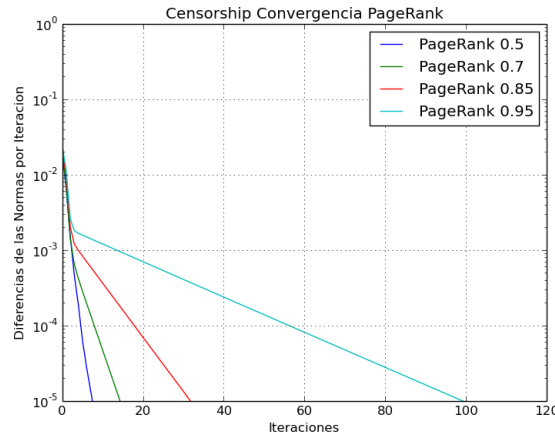


Figura 1: Normas (PageRank): escala logarítmica

Se observa en la figura 1 que la cantidad de iteraciones necesarias para llegar a un ranking aceptable parece tener una clara relación con el  $c$  que se define: cuanto mayor es el  $c$ , mayor cantidad de iteraciones debe hacer antes de llegar a una norma que el criterio de parada acepte.

De esto se puede deducir que cuanto menor importancia se le dé a la estructura de la red, menor cantidad de cálculos habrá que hacer. De allí se puede ver que un  $c$  pequeño implica menor probabilidad de que el navegante aleatorio siga algún camino predecible, y mayor de que salte a una página cualquiera de la web; esto a su vez llevará a que los caminos que nosotros tratamos de establecer tengan un menor peso a la hora de determinar sus acciones.

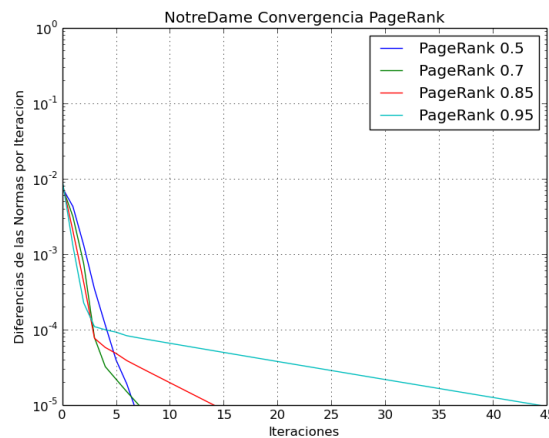
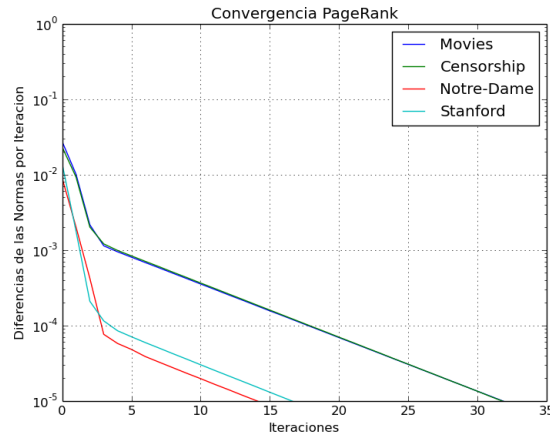


Figura 2: Normas (PageRank): escala logarítmica

Las normas, al igual que en el caso anterior, convergen en forma exponencial.

La relación entre las iteraciones necesarias según el  $c$  que se elija, además, son análogas a aquel presentado en redes medianas.

Se decidió luego analizar estas convergencias de normas de *Pagerank* de diferentes redes entre sí, tomando como  $c$  aquel recomendado en *Brin, Page - The anatomy of a large-scale hypertextual Web search engine*, esto es, tomando  $c = 0,85$ .



Contrario a lo que se esperaba, las normas de las matrices más grandes (Stanford y Notre-Dame) presentaron una convergencia pronunciadamente más rápida que las de aquellas más pequeñas (Movies, Censorship).

En PageRank, las normas convergen a 0 en forma exponencial luego de las primeras iteraciones. Cuanto mayor sea el valor de  $c$ , más tarda en alcanzarse el criterio de parada. Con  $c = 1$ , la norma podría no converger, ya que la matriz  $P_2$  es igual a  $P$ , que puede no ser estocástica por columnas.

Luego de experimentar con el algoritmo de PageRank se decidió observar la convergencia de la diferencia de la norma 2 en HITS:  $\|x^{(k)} - x^{(k-1)}\|_2$  y  $\|y^{(k)} - y^{(k-1)}\|_2$ . Se presentan en los siguientes gráficos dos casos que consideramos representativos de los resultados de experimentos sobre redes medianas y grandes, respectivamente.

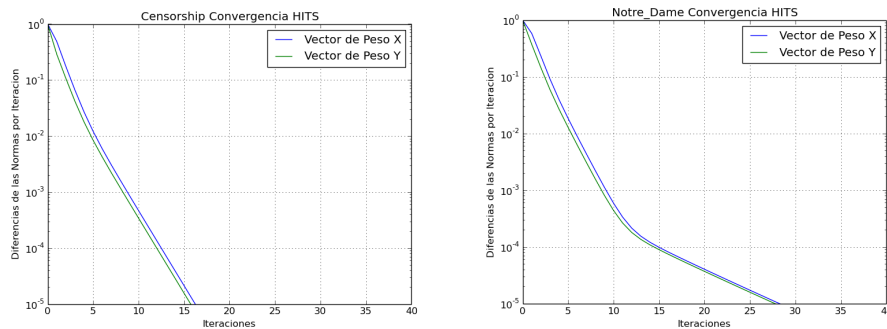


Figura 3: Normas (HITS): escala logarítmica

En la figura 3 se puede observar que en HITS las normas también convergen a 0 en forma exponencial. Sin embargo, en las muestras de mayor tamaño se observa que tarda más iteraciones en converger.

Se procede a continuación a comparar una red que consideramos representativa de la mayoría de las redes medianas tomadas y *Genetic*, que mostró una diferencia importante en comparación con las demás.

Métodos	Iteraciones de Llegada	
	Censorship	Genetic
HITS	18	37
Pagerank 0.5	9	10
Pagerank 0.7	16	19
Pagerank 0.85	33	39
Pagerank 0.95	101	121

Sorprendentemente, se encontró que al correr HITS sobre las redes medianas, *Genetic* precisó aproximadamente el doble de iteraciones que fueron necesarias para *Censorship* (y para la mayoría de las redes medianas tomadas).

Por otro lado, en Pagerank existe una cierta diferencia entre *Genetic* y *Censorship* también, pero está evidentemente muy lejos de ser tan grande como la presentada en HITS. Se puede ver, sin embargo, que ésta escala de forma proporcional al  $c$  que se tome.

Se consideró en un principio la idea de que la cantidad de nodos, links o la relación entre estos dos valores tuvieran alguna influencia en esta diferencia, sin embargo no parece ser ésta la principal razón ya que los valores de *Genetic* no muestran ninguna discrepancia especial con la mayoría de los casos de tests tomados (ver sección 2.5.: Experimentos realizados)

No consideramos que sea posible determinar de manera fehaciente la causa de esta diferencia entre redes de tamaño tan grande sin entrar en experimentos más grandes y que requerirían un análisis mucho más profundo.

### 3.2. Tiempos de Cómputo

A continuación se presentan los resultados de tiempo de cómputo obtenidos, en ciclos de clock:

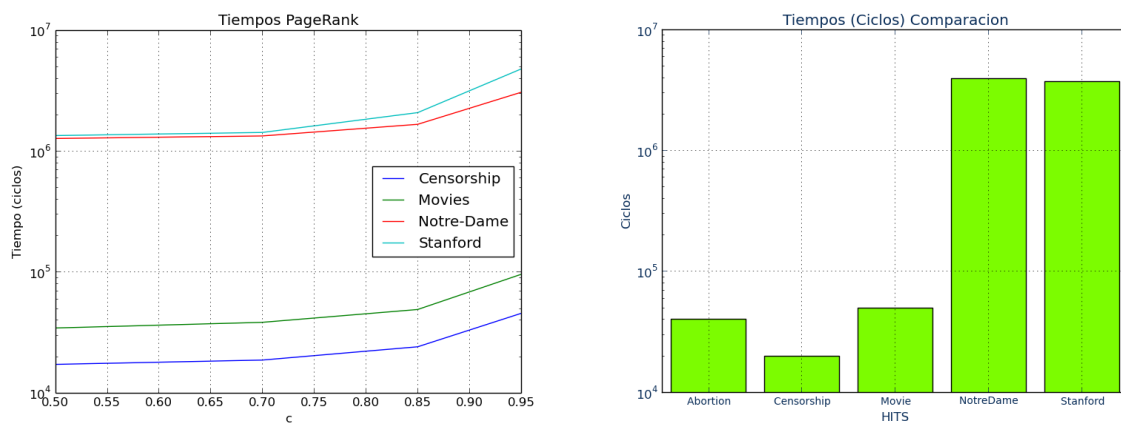


Figura 4: Tiempo de cómputo

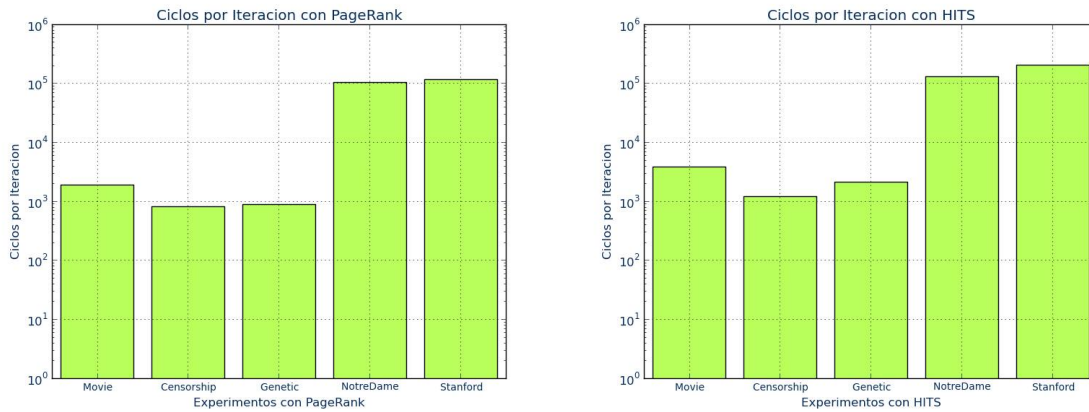


Figura 5: Ciclos por iteración

En todos los casos analizados se obtuvo que el tiempo de cómputo de PageRank aumentó con el valor del  $c$  elegido.

Se puede notar en ambos gráficos que dicho tiempo de cómputo es proporcional al tamaño de la red para ambos métodos.

Los tests de mayor tamaño, aunque no requerían una mayor cantidad de iteraciones, sí tienen mayor tiempo de cómputo. Esto se debe a que la multiplicación entre la matriz y el vector que es necesaria en cada iteración tiene una mayor cantidad de elementos. Esto se puede observar en la figura de Ciclos por iteración.

El algoritmo HITS resultó más lento que PageRank para todos los valores de  $c < 0,85$  que utilizamos, en todas las redes para las que medimos los tiempos. Sin embargo, consideramos que esta diferencia se debe principalmente a que el algoritmo HITS debe realizar dos productos entre una matriz y un vector, mientras que PageRank solo realiza una.

### 3.3. Rankings obtenidos

Se presenta aquí una experimentación sobre una red *real*; se cuenta en esta red con las siguientes páginas:

- ole.com.ar
- clasificados.clarin.com
- ciudad.com.ar
- lanacion.com.ar
- canchallena.lanacion.com.ar
- clarin.com
- clarin.com/deportes
- rollingstone.com.ar
- zonaprop.com.ar
- google.com
- infobae.com
- mamapuntocero.com.ar
- pagina12.com.ar
- yahoo.com

Utilizando el *webParser* provisto por la cátedra, se ha buscado obtener el ranking de estas páginas utilizando los distintos métodos.

En cuanto al PageRank obtenido, se esperaba que las páginas que consideramos más importantes como *google.com* y *yahoo.com* estuvieran encabezando la lista. Se esperaba además que por debajo de las más importantes estuvieran las páginas relacionadas con diarios y revistas, tales como *lanacion.com.ar* y *clarin.com*. Finalmente considerabamos razonable encontrar en último lugar secciones de los diarios (tales como *clarin.com/deportes*) y páginas menos relevantes como *mamapuntocero.com.ar*.

Para el puntaje de hub de HITS, se esperaba algo similar: un gran puntaje en los buscadores (Google, Yahoo) y un puntaje algo menor en los diarios y sus secciones, que deben referenciarse entre sí, y últimas, páginas como *mamapuntocero.com.ar*.

Para el puntaje de autoridad, se esperaba un puntaje bajo para los buscadores, que no suelen tener más links entrantes que los generados por otros buscadores. El puntaje para los diarios y revistas y sus secciones debería ser alto, ya que se espera que se referencien entre sí.

Sin embargo, es necesario notar que la estructura de la red no es representativa de la *World Wide Web*: hay cinco sitios del diario Clarín (considerando Olé y Ciudad) y dos de La Nación. Es esperable que estos sitios estén relacionados entre sí, generando un puntaje que puede superar al de los buscadores.

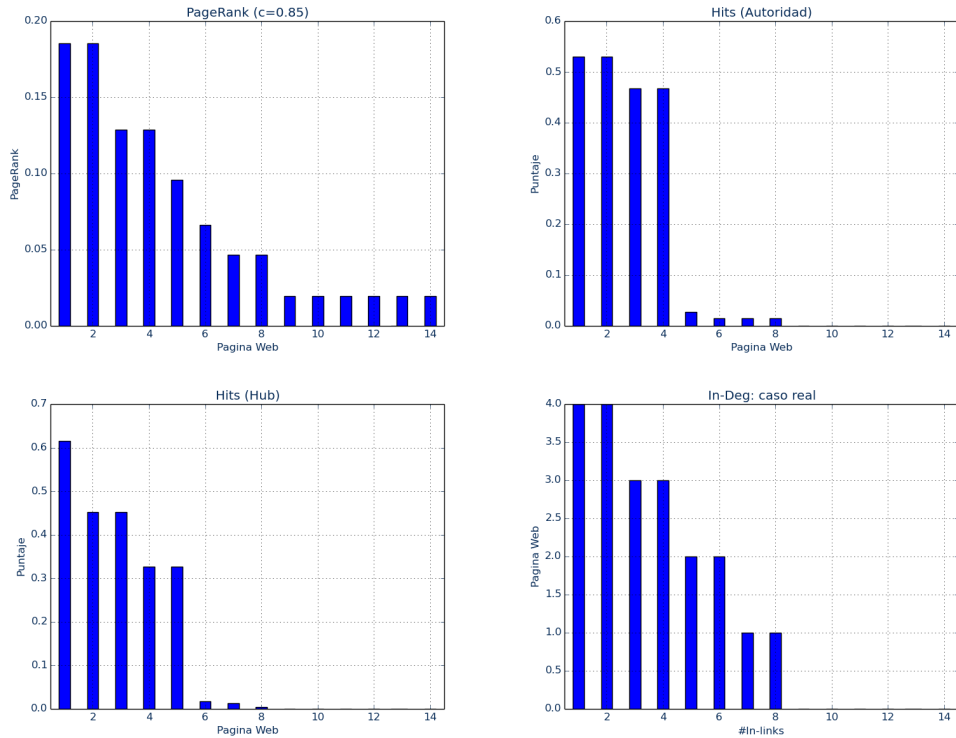


Figura 6: Puntajes obtenidos con los distintos métodos por distintas páginas

Ranking	PageRank		Hits (Autoridad)		Hits(Hub)
1	clarin.com	1	clarin.com	1	clarin.com/deportes
	ole.com.ar		ole.com.ar	2	ole.com.ar
2	clasificados.clarin.com	2	clasificados.clarin.com		clarin.com
	ciudad.com.ar		ciudad.com.ar	3	clasificados.clarin.com
3	lanacion.com.ar	3	canchallena.lanacion.com.ar		ciudad.com.ar
4	canchallena.lanacion.com.ar	4	lanacion.com.ar	4	lanacion.com.ar
5	rollingstone.com.ar	5	rollingstone.com.ar	5	rollingstone.com.ar
	zonaprop.com.ar		zonaprop.com.ar	6	canchallena.lanacion.com.ar
6	google.com	6	google.com	7	zonaprop.com.ar
	clarin.com/deportes		clarin.com/deportes		yahoo.com
	infobae.com		infobae.com		pagina12.com.ar
	mamapuntocero.com.ar		mamapuntocero.com.ar		mamapuntocero.com.ar
	pagina12.com.ar		pagina12.com.ar		infobae.com
	yahoo.com		yahoo.com		google.com



Ranking	InDeg
1	clarin.com ole.com.ar
2	clasificados.clarin.com ciudad.com.ar
3	lanacion.com.ar canchallena.lanacion.com.ar
4	rollingstone.com.ar zonaprop.com.ar
5	google.com clarin.com/deportes infobae.com mamapuntocero.com.ar pagina.com.ar yahoo.com

Al hacer un análisis cualitativo de los puntajes obtenidos se puede observar que los rankings en general no fueron los esperados. Por ejemplo, se esperaba que *Google* y *Yahoo* obtuvieran los mayores puntajes de hub, hecho que no tuvo lugar.

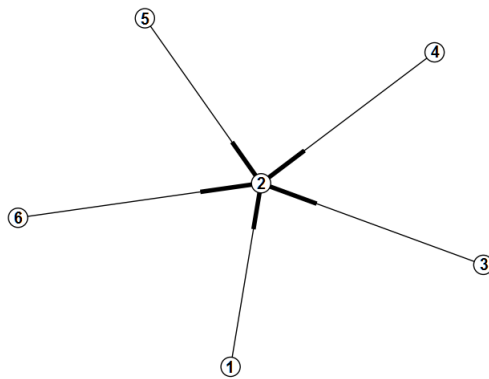
Esto se debe a que la porción de la red elegida no es representativa. Se ve que *clarin.com* obtuvo el mayor puntaje de PageRank, pero hay que tener en cuenta que en la lista hay sitios que pertenecen al mismo dominio, que también obtuvieron un puntaje alto. En general, diarios y revistas se ubicaron en primer lugar. Páginas que hemos considerado de menor importancia se han ubicado por debajo, como era esperado.

Además, se tomaron las páginas principales de *Google* y *Yahoo* que no tienen links salientes. Para obtenerlos, es necesario realizar una búsqueda.

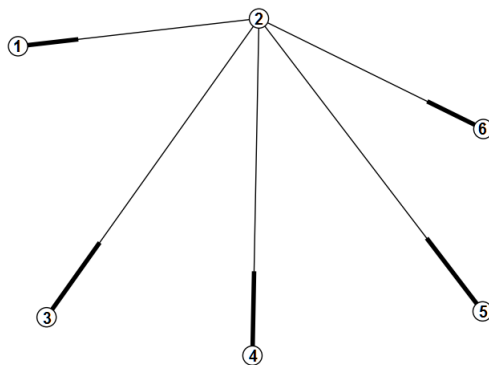
### 3.4. Redes pequeñas

Para observar el comportamiento de PageRank y HITS, consideramos los siguientes casos de seis nodos, tomando para PageRank  $c = 0,85$ :

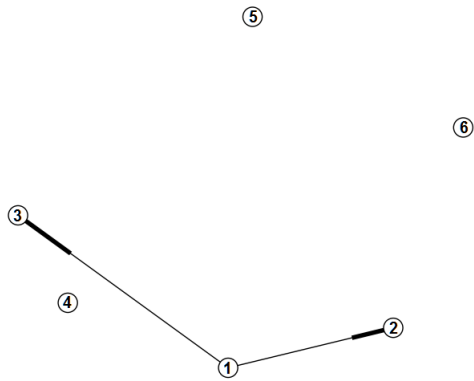
- Una red en la que todas las páginas apuntan a una. El PageRank de una página puede considerarse como la probabilidad del navegante aleatorio de encontrarse en ella. En este caso, el navegante aleatorio comienza en una página, sigue el link a la única a la que esta apunta y, como ésta es un dangling node, salta a cualquiera y vuelve a seguir el único link. Esto se traduciría en que la mitad del tiempo el navegante se encuentra en la página a la que apuntan todas, y la otra mitad, en cualquier otra. Esto no es exacto por la probabilidad de un salto entre páginas sin seguir links, pero para  $c = 0,85$  esta probabilidad es baja. Además, se podría saltar directamente a la página central. Para HITS, esperamos ver un puntaje de autoridad de 1 para esta página y un puntaje de hub repartido entre el resto.



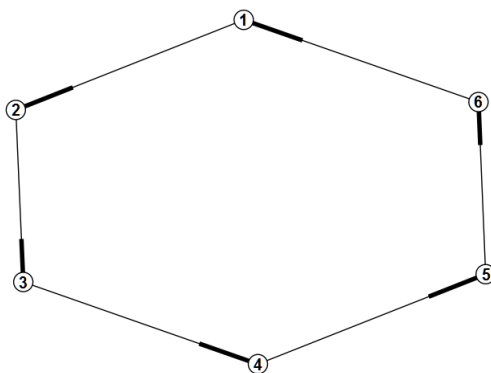
- Una red en la que una página apunta a todas. La intención es mostrar un caso en el que el PageRank no consigue representar la estructura de la red, ya que debería ser igual para todas. Esto es porque la página que apunta a todas se comporta exactamente como un dangling node, y el resto lo son. En cambio, en HITS, el puntaje de hub de esta página debería ser 1, mientras que el puntaje de autoridad debe estar repartido entre el resto de las páginas.



- Una red con muy pocos links. Con esto se puede ver la diferencia entre los dos métodos causada por los dangling nodes: en HITS el puntaje está concentrado en los pocos links, mientras que en PageRank se "fabrican" links salientes de cada dangling node, generando un puntaje similar para todas las páginas.



- Una red en la que se forma un ciclo. El navegante aleatorio tiene en todo momento la misma probabilidad de estar en cualquier página, ya que, o bien sigue el ciclo, o bien salta con igual probabilidad a cualquier página, por lo que el PageRank debería ser igual para todas. El puntaje de HITS también, porque todas las páginas tienen un link saliente y un link entrante a páginas que no son distinguibles de sí misma. Lo que se espera ver en este ejemplo es un caso en que las páginas tienen un grado de entrada (In-Deg) muy similar al del ejemplo anterior, pero su estructura es muy diferente.



Los resultados obtenidos se presentan a continuación:

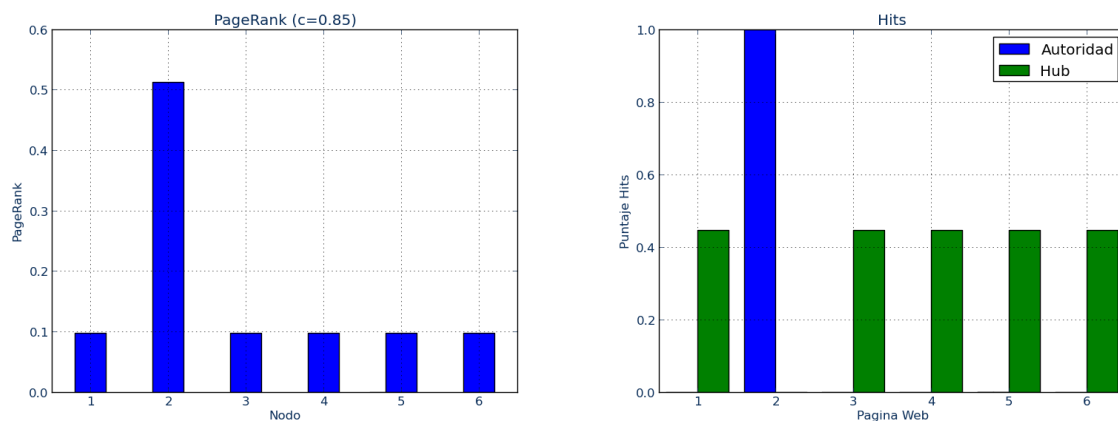


Figura 7: Puntajes obtenidos para una red en la que todas las páginas apuntan a la n° 2

En la figura 7 se ve que la página n° 2 es la que, como se esperaba, concentra todo el puntaje de autoridad; y es también la que tiene mayor PageRank. Este es muy cercano a 0.5, pero mayor.

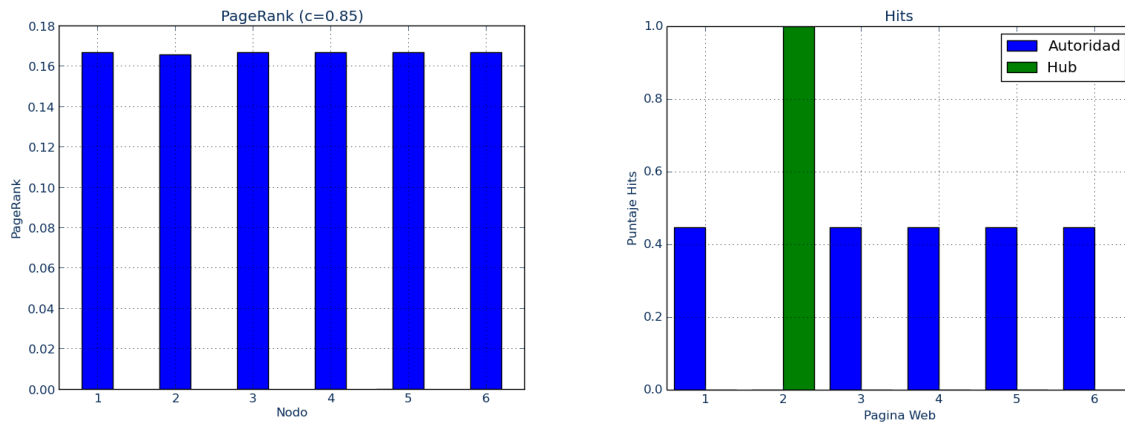


Figura 8: Puntajes obtenidos para una red en la que la pagina n° 2 apunta a todas

Por otro lado en la Figura 8, como era esperable, los puntajes de HITS se han invertido completamente, aquellos que quedaban en segundo lugar ahora ocupan el primero. En cambio, el puntaje de PageRank fue igual para todas las páginas.

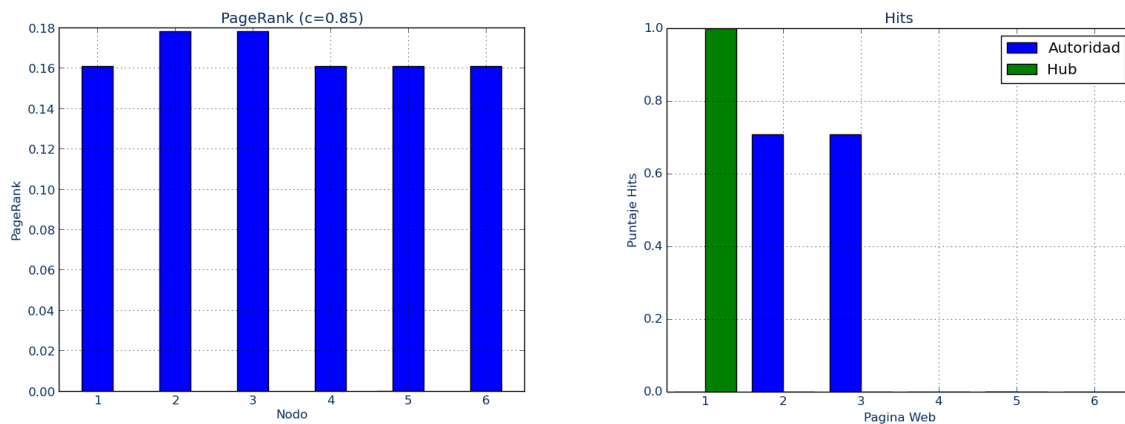


Figura 9: Puntajes obtenidos para una red con pocos links

En HITS, tanto el puntajes de autoridad como el de hub están concentrados en las páginas que tienen links. En cambio, en PageRank los puntajes fueron similares para todas las páginas, con un puntaje mayor en las que tienen links entrantes. En este caso, como la *página hub* no tiene ningún link entrante, tiene el mismo puntaje que los nodos aislados. De hecho, en este caso, PageRank no refleja ninguna información sobre los links salientes.

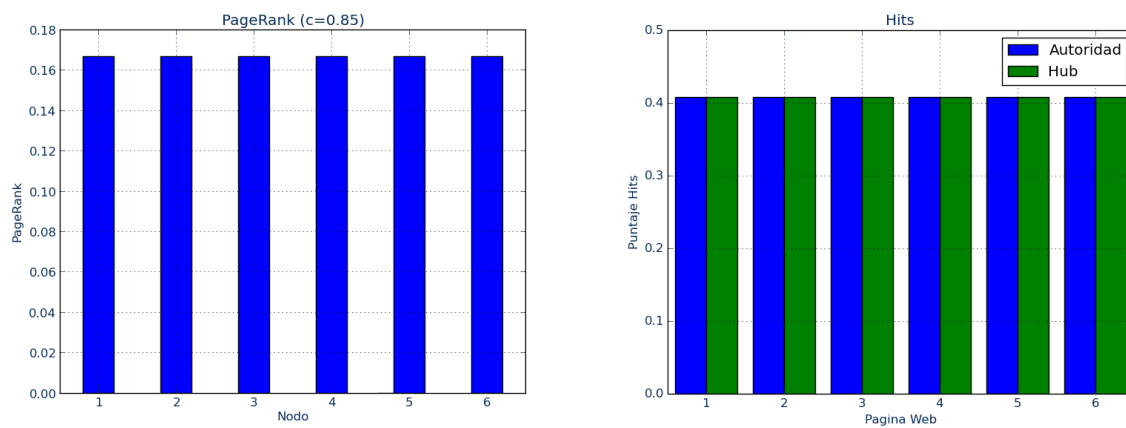


Figura 10: Puntajes obtenidos para una red cíclica

Para la red que forma un ciclo, el puntaje fue exactamente igual para todos, tanto en PageRank como en Hits. Esto es lo que se esperaba.

## 4. Conclusiones

En el paper escrito por Brin y Page <sup>8</sup> se recomienda tomar un  $c$  (probabilidad de un navegante aleatorio de seguir un link de la página en la que está) de 0.85. Consideramos sin embargo que, de contar con los recursos necesarios para realizar mediciones de la magnitud necesaria para obtener resultados representativos, la mejor manera de definir un valor de  $c$  fiable en la mayor cantidad de casos posibles sería acercarnos a la probabilidad real con la que el navegante aleatorio *salta* a otra página sin seguir un link, para esto haría falta hacer experimentos con navegantes reales y realizar un promedio de los  $c$  que obtuvimos de estos navegantes; es decir que el mejor  $c$  a tomar sería la esperanza de los  $c$  obtenidos en estas mediciones.

El hecho de que la norma converja exponencialmente permite que la cantidad de iteraciones sea lo suficientemente acotada, teniendo en cuenta que cada iteración es una multiplicación de una matriz por un vector que puede tener varios millones de posiciones.

Luego de haber contrastado estos distintos métodos de ranking podemos decir con seguridad que:

**In-Deg** es completamente descartable como método para decidir la relevancia de una página, se limita a realizar el análisis más previsible y propenso a caer en errores. Buscadores puramente textuales que siguieron este método, como AltaVista, fueron reemplazados eventualmente por otros más útiles, si bien fueron un paso necesario para la creación de buscadores mejores.

**HITS** es mucho más efectivo y confiable ya que realiza un análisis más inteligente de la red que le es proporcionada, confiriéndole distintos roles a las páginas y contrastándolas a partir de las relaciones entre estas páginas según el rol que estén cumpliendo.

Hemos comprobado que la cantidad de links a una página no son relevantes por sí solos para considerarla una buena autoridad sino que lo más importante es la relevancia como hubs de aquellas páginas que la apuntan, esto revela un análisis mucho más profundo que el de *Indeg*.

Por la naturaleza del análisis que realiza es, sin embargo, más susceptible de caer en engaños tales como que páginas acuerden apuntarse entre sí para ganar peso en el ranking que *HITS* les asigne, mientras que con *PageRank* este artilugio no daría frutos debido a que cuantos más links salientes una página tenga, menos valor le estará transfiriendo con cada link a las páginas apuntadas.

Por ejemplo, si en el caso considerado en la sección 3.3 se agrega un link de *Clarín*, una página que tiene varios links salientes, a una que está en el último lugar de los rankings, como *mamapuntocero*, la página apuntada asciende al quinto puesto en el ranking de Hubs de Hits, mientras que solo llega al séptimo puesto en el ranking de *PageRank* (con  $c=0.85$ ).

Mirando el concepto detrás de *HITS* (a grandes rasgos), es posible definir una clara relación entre el concepto de Autoridad de una página en *HITS* y el concepto de importancia de una en *Pagerank*, la relación, sin embargo, no va más allá del concepto ya que los cálculos y la forma de definir el valor de Autoridad o la importancia de una página son evidentemente muy distintos.

**PageRank** es el método que consideramos mejor, superando tanto a *HITS* como a *Indeg* debido a que realiza un análisis profundo de las relaciones entre páginas, distinto de aquel proporcionado por *HITS*. En ciertos casos parece que el procedimiento es opuesto. Por ejemplo, una página que apunta a todas le da poco valor a cada una en *PageRank*, pero al apuntar a muchas autoridades es un buen hub según *HITS*, luego le da valor a los que apunta. Un inconveniente que tiene este algoritmo es que trata una página con links a toda la web igual que un dangling node, mientras que *HITS* le da un gran puntaje como hub.

Por ultimo, si el cliente quiere tener un mejor *PageRank*, debe conseguir que lo apunten páginas que tengan la mejor proporción entre *PageRank* y cantidad de links salientes, ya que el puntaje se distribuye entre todas las páginas. En cambio, si quiere conseguir un mejor puntaje como autoridad en *HITS*, debe buscar que la mayor cantidad posible de las mejores páginas lo apunten. Si desea tener un mejor puntaje de hub tendría que apuntar a las páginas con mayor autoridad en *HITS*.

---

<sup>8</sup>Brin, Page - 1998 - The anatomy of a large-scale hypertextual Web search engine

Podemos finalmente concluir que la mejor estrategia a sugerir a clientes sería darle prioridad al ranking proporcionado por Pagerank, está además la prueba de la efectividad del Pagerank en el éxito que tuvo Google al usarlo. Sugeriríamos entonces apoyarse en el orden de resultados que proporcione este buscador (que es además el más usado) para decidir en qué páginas comprar espacio de publicidad.

## 5. Apéndice A

### Tirate un qué, tirate un *ranking*...

#### Motivación

Luego de su repentina y efímera irrupción durante el año 2011, un grupo de la movida tropical<sup>9</sup> está buscando recuperar la notoriedad y los niveles de popularidad otrora alcanzados. El retorno incluye, entre otras cosas, un mega recital gratuito, giras por las principales *bailantas* y por el interior del país.<sup>10</sup>

Para que toda esta movida sea exitosa, los miembros del grupo han acordado con su *community manager* que, además de tener una participación destacada en Pasión de Sábado, es necesario que la llegada a través de los medios electrónicos y las redes sociales sea muy efectiva, al igual que en 2011, alcanzando a la mayor cantidad posible de gente y poder, nuevamente, sentarse en el living de *la diva de los teléfonos*. La conclusión a la que llegaron es que necesitan que cada vez que realiza una búsqueda relacionada con la movida tropical, su página se encuentre entre las primeras que muestran los buscadores.

Con ese motivo, se han contactado con el equipo de R+D de Métodos Numéricos, donde en la primera reunión el cliente propuso *comprar clicks en publicidades*. Esta, si bien es una alternativa viable, representa un gasto importante para la escala de inversión con la que se dispone. Luego de una reunión del equipo técnico, se les hizo una contrapropuesta: estudiar el comportamiento de los buscadores y, a cambio de shows libres de costo y presentaciones privadas, buscar en qué páginas conviene figurar para mejorar el posicionamiento virtual del grupo.

#### Contexto

A partir de la evolución de Internet durante la década de 1990, el desarrollo de motores de búsqueda se ha convertido en uno de los aspectos centrales para su efectiva utilización. Hoy en día, sitios como Yahoo, Google y Bing ofrecen distintas alternativas para realizar búsquedas complejas dentro de un red que contiene miles de millones de páginas web.

En sus comienzos, una de las características que distinguió a Google respecto de los motores de búsqueda de la época fue la calidad de los resultados obtenidos, mostrando al usuario páginas relevantes a la búsqueda realizada. El esquema general de los orígenes de este motor de búsqueda es brevemente explicando en Brin y Page [?], donde se mencionan aspectos técnicos que van desde la etapa de obtención de información de las páginas disponibles en la red, su almacenamiento e indexado y su posterior procesamiento, buscando ordenar cada página de acuerdo a su importancia relativa dentro de la red. El algoritmo utilizado para esta última etapa es denominado PageRank y es uno (no el único) de los criterios utilizados para ponderar la importancia de los resultados de una búsqueda. En este trabajo nos concentraremos en el estudio y desarrollo del algoritmo PageRank.

#### Los métodos, Parte I: PageRank

El algoritmo PageRank se basa en la construcción del siguiente modelo. Supongamos que tenemos una red con  $n$  páginas web  $Web = \{1, \dots, n\}$  donde el objetivo es asignar a cada una de ellas un puntaje que determine la importancia relativa de la misma respecto de las demás. Para modelar las relaciones entre ellas, definimos la *matriz de conectividad*  $W \in \{0, 1\}^{n \times n}$  de forma tal que  $w_{ij} = 1$  si la página  $j$  tiene un link a la página  $i$ , y  $w_{ij} = 0$  en caso contrario. Además, ignoramos los *autolinks*, es decir, links de una página a sí misma, definiendo  $w_{ii} = 0$ . Tomando esta matriz, definimos el grado de la página  $j$ ,  $n_j$ , como la cantidad de links salientes hacia otras páginas de la red, donde  $n_j = \sum_{i=1}^n w_{ij}$ . Además, notamos con  $x_j$  al puntaje asignado a la página  $j \in Web$ , que es lo que buscamos calcular.

La importancia de una página puede ser modelada de diferentes formas. Un link de la página  $u \in Web$  a la página  $v \in Web$  puede ser visto como que  $v$  es una página importante. Sin embargo, no queremos que una página obtenga mayor importancia simplemente porque es apuntada desde muchas páginas. Una forma de limitar esto es ponderar los links utilizando la importancia de la página de origen. En

<sup>9</sup>Por cuestiones de privacidad, no haremos público de qué grupo se trata.

<sup>10</sup>A riesgo de exponer su edad, los miembros de la cátedra quieren destacar a aquellos próceres que llevaron a este género musical a las primeras planas, como Alcides, Sebastián, Miguel Conejito Alejandro, Ráfaga, La Nueva Luna, Comanche y, como dejar fuera, al MAESTRO Antonio Ríos.



otras palabras, pocos links de páginas importantes pueden valer más que muchos links de páginas poco importantes. En particular, consideramos que la importancia de la página  $v$  obtenida mediante el link de la página  $u$  es proporcional a la importancia de la página  $u$  e inversamente proporcional al grado de  $u$ . Si la página  $u$  contiene  $n_u$  links, uno de los cuales apunta a la página  $v$ , entonces el aporte de ese link a la página  $v$  será  $x_u/n_u$ . Luego, sea  $L_k \subseteq Web$  el conjunto de páginas que tienen un link a la página  $k$ . Para cada página pedimos que

$$x_k = \sum_{j \in L_k} \frac{x_j}{n_j}, \quad k = 1, \dots, n. \quad (1)$$

Definimos  $P \in \mathbb{R}^{n \times n}$  tal que  $p_{ij} = 1/n_j$  si  $w_{ij} = 1$ , y  $p_{ij} = 0$  en caso contrario. Luego, el modelo planteado en (1) es equivalente a encontrar un  $x \in \mathbb{R}^n$  tal que  $Px = x$ , es decir, encontrar (suponiendo que existe) un autovector asociado al autovalor 1 de una matriz cuadrada, tal que  $x_i \geq 0$  y  $\sum_{i=1}^n x_i = 1$ . En Bryan y Leise [?] y Kamvar et al. [?, Sección 1] se analizan ciertas condiciones que debe cumplir la red de páginas para garantizar la existencia de este autovector.

Una interpretación equivalente para el problema es considerar al *navegante aleatorio*. Éste empieza en una página cualquiera del conjunto, y luego en cada página  $j$  que visita sigue navegando a través de sus links, eligiendo el mismo con probabilidad  $1/n_j$ . Una situación particular se da cuando la página no tiene links salientes. En ese caso, consideramos que el navegante aleatorio pasa a cualquiera de las páginas de la red con probabilidad  $1/n$ . Para representar esta situación, definimos  $v \in \mathbb{R}^{n \times n}$ , con  $v_i = 1/n$  y  $d \in \{0, 1\}^n$  donde  $d_i = 1$  si  $n_i = 0$ , y  $d_i = 0$  en caso contrario. La nueva matriz de transición es

$$\begin{aligned} D &= v d^t \\ P_1 &= P + D. \end{aligned}$$

Además, consideraremos el caso de que el navegante aleatorio, dado que se encuentra en la página  $j$ , decida visitar una página cualquiera del conjunto, independientemente de si esta se encuentra o no referenciada por  $j$  (fenómeno conocido como *teletransportación*). Para ello, consideramos que esta decisión se toma con una probabilidad  $c \geq 0$ , y podemos incluirlo al modelo de la siguiente forma:

$$\begin{aligned} E &= v \bar{1}^t \\ P_2 &= cP_1 + (1 - c)E, \end{aligned}$$

donde  $\bar{1} \in \mathbb{R}^n$  es un vector tal que todas sus componentes valen 1. La matriz resultante  $P_2$  corresponde a un enriquecimiento del modelo formulado en (1). Probabilísticamente, la componente  $x_j$  del vector solución (normalizado) del sistema  $P_2x = x$  representa la proporción del tiempo que, en el largo plazo, el navegante aleatorio pasa en la página  $j \in Web$ .

En particular,  $P_2$  corresponde a una matriz *estocástica por columnas* que cumple las hipótesis planteadas en Bryan y Leise [?] y Kamvar et al. [?], tal que  $P_2$  tiene un autovector asociado al autovalor 1, los demás autovalores de la matriz cumplen  $1 = \lambda_1 > |\lambda_2| \geq \dots \geq |\lambda_n|$  y, además, la dimensión del autoespacio asociado al autovalor  $\lambda_1$  es 1. Luego, la solución al sistema  $P_2x = x$  puede ser calculada de forma estándar utilizando el método de la potencia.

Una vez calculado el ranking, se retorna al usuario las  $t$  páginas con mayor ranking.

## Los métodos, Parte II: Hyperlink-Induced Topic Search

Un método alternativo es propuesto en Kleinberg [?], denominado *Hyperlink-Induced Topic Search* (HITS). La intuición del método se basa en el análisis intrínseco de la red, donde una noción de *autoridad* se transfiere de una página a otra mediante los links que las relacionan. El objetivo es, dada una búsqueda concreta, retornar un subconjunto acotado de páginas relevantes. Con este fin, se considera que existen páginas que cumplen un rol de *autoridad* sobre un tema específico y se busca modelar la relación entre estas páginas y aquellas que apuntan a varias de estas autoridades, denominadas *hubs*. En la práctica, los autores observan que suele existir una especie de equilibrio en la relación entre hubs y autoridades, y se busca aprovechar esta relación para el desarrollo del algoritmo. Intuitivamente, un buen *hub* es una página que apunta a muchas autoridades, y una buena *autoridad* es una página que es apuntada por muchos *hubs*.

El procedimiento consiste en los siguientes pasos. Dada una búsqueda concreta, se utiliza en primer lugar un *buscador* simple (por ejemplo, basado en texto) para obtener un conjunto acotado de páginas (digamos, 200), llamado *root set*. Luego, asumiendo que la estructura de la red es conocida, se busca extender este conjunto agregando páginas que son apuntadas y que apuntan a las páginas de *root set*, hasta llegar a una sub-red de un tamaño determinado. En el contexto del trabajo práctico, asumiremos que este paso ha sido realizado y que contamos con el grafo que considera la sub-red.

Formalmente, y retomando la notación introducida en la sección anterior, consideramos que las páginas de nuestra sub-red se encuentran en el conjunto  $Web = \{1, \dots, n\}$ . Para modelar las relaciones entre las páginas, adoptamos una definición similar: consideramos la matriz de adyacencia  $A \in \{0, 1\}^{n \times n}$  donde  $a_{ij} = 1$  si existe un link de la página  $i$  a la página  $j$ .<sup>11</sup> Para cada página  $i \in Web$  se considera el *peso de autoridad*  $x_i$  y el *peso de hub*  $y_i$ . Consecuentemente, se definen los vectores  $x, y \in \mathbb{R}^n$  los vectores de pesos de autoridad y hubs, respectivamente, y supondremos además que se encuentran normalizados. Las páginas con mayores valores de  $x_i$  e  $y_i$  son consideradas mejores *autoridades* y *hubs*, respectivamente.

La relación mencionada entre los distintos tipos de páginas se expresan numéricamente de la siguiente forma. Dados los vectores  $x, y$ , la operación de transferencia de los *hubs* a la autoridad  $j \in Web$  puede expresarse de la siguiente forma:

$$x_j = \sum_{i:i \rightarrow j} y_i. \quad (2)$$

Análogamente, el peso de un hub está dado por la siguiente ecuación

$$y_i = \sum_{j:i \rightarrow j} x_j. \quad (3)$$

Las ecuaciones (2) y (3) podemos expresarlas matricialmente de la siguiente manera:

$$x = A^t y \quad (4)$$

$$y = Ax, \quad (5)$$

aplicando luego el paso de normalización correspondiente. Los autores proponen comenzar con un  $y_0$  inicial, aplicar estas ecuaciones iterativamente y demuestran que, bajo ciertas condiciones, el método converge. Finalmente, en base a los rankings obtenidos, se retorna al usuario las mejores  $t$  *autoridades* y los mejores  $t$  *hubs*.

### Enunciado

El objetivo del trabajo es experimentar en el contexto planteado utilizando los algoritmos de ranking propuestos. Para ello, se considera un entorno que, dentro de nuestras posibilidades, simule el contexto real de aplicación donde se abordan instancias de gran escala (es decir,  $n$ , el número total de páginas, es grande). El archivo tomará como entrada un archivo que especifique el algoritmo, los parámetros del mismo y un puntero al grafo de la red y retorne como resultado el ranking obtenido para cada página. Los detalles sobre el input/output del programa son especificados en la siguiente sección.

El trabajo consistirá en estudiar distintos aspectos de los siguientes métodos: PageRank, HITS, e IN-DEG, éste último consiste en definir el ranking de las páginas utilizando solamente la cantidad de ejes entrantes a cada una de ellas, ordenándolos en forma decreciente. Para tener una descripción más completa de los dos primeros métodos, se propone:

1. Considerar el trabajo de Kleinberg [?] con los detalles sobre HITS, en particular las secciones 1, 2 y 3.
2. Considerar el trabajo de Bryan y Leise [?] donde se explica la intuición y algunos detalles técnicos respecto a PageRank. Además, en Kamvar et al. [?] se propone una mejora del mismo. Si bien esta mejora queda fuera de los alcances del trabajo, en la Sección 1 se presenta una buena formulación del algoritmo. En base a su definición,  $P_2$  no es una matriz esparsa. Sin embargo, en Kamvar et al. [?, Algoritmo 1] se propone una forma alternativa para computar  $x^{(k+1)} = P_2 x^{(k)}$ . Este resultado puede ser utilizado para mejorar el almacenamiento de los datos.

<sup>11</sup>Notar que  $A = W^t$ .

3. (Opcional) Completar la demostración del Teorema 3.1 de Kleinberg [?], incluyendo el detalle de los puntos que el autor asume como triviales.

En la práctica, el grafo que representa la red de páginas suele ser esparso, es decir, una página posee relativamente pocos links de salida comparada con el número total de páginas. A su vez, dado que  $n$  tiende a ser un número muy grande, es importante tener en cuenta este hecho a la hora de definir las estructuras de datos a utilizar. Luego, desde el punto de vista de implementación se pide utilizar alguna de las siguientes estructuras de datos para la representación de las matrices esparsas: *Dictionary of Keys* (dok), *Compressed Sparse Row* (CSR) o *Compressed Sparse Column* (CSC). Se deberá incluir una justificación respecto a la elección que consdiere el contexto de aplicación. Una vez definida la estructura a utilizar, se deberá implementar el algoritmo HITS utilizando las ecuaciones (4) y (5). Para el caso de PageRank, se debe implementar el método de la potencia para calcular el autovector principal.

En función de la experimentación, se deberá realizar un estudio particular para cada algoritmo (tanto en términos de comportamiento del mismo, como una evaluación de los resultados obtenidos) y luego se procederá a comparar cualitativamente los rankings generados. La experimentación deberá incluir como mínimo los siguientes experimentos:

1. Estudiar la convergencia de PageRank, analizando la evolución de la norma Manhattan (norma  $L_1$ ) entre dos iteraciones sucesivas. Comparar los resultados obtenidos para al menos dos instancias de tamaño mediano-grande, variando el valor de  $c$ . Opcional: Establecer una relación con la proporción entre  $\lambda_1 = 1$  y  $|\lambda_2|$ .
2. Estudiar la convergencia de los vectores de peso  $x$  e  $y$  para HITS de forma similar al punto anterior.
3. Estudiar el tiempo de cómputo requerido por PageRank y HITS. Si bien ambos pueden se aplicados sobre una red genérica, cada algoritmo tiene un contexto particular de aplicación. Estudiar como impacta el factor temporal en este sentido.
4. Estudiar cualitativamente los rankings obtenidos por los tres métodos. Para ello, se sugiere considerar distintos ejemplos de búsquedas de páginas web<sup>12</sup>. Analizar los resultados individualmente en una primera etapa, y luego realizar un análisis comparativo entre los tres rankings obtenidos.
5. Para cada algoritmo, proponer ejemplos de tamaño pequeño que ilustren el comportamiento esperado (puede ser utilizando las instancias provistas por la cátedra o generadas por el grupo).

Finalmente, y en base a la experimentación realizada, buscamos resolver el problema planteado originalmente: dada una foto de la red, con sus interconexiones entre páginas, supongamos que tenemos los pesos (ranking) asignados por uno de los algoritmos estudiados. ¿Cuál sería la estrategia que le sugiere al cliente para mejorar su correspondiente ranking? Para este último punto, suponer que es posible *negociar* que una página apunte a nuestro sitio, y que la cantidad de estas negociaciones que podemos tener es acotada.

### Parámetros y formato de archivos

El programa deberá tomar por línea de comandos dos parámetros. El primero de ellos contendrá la información del experimento, incluyendo el método a ejecutar (`alg`, 0 para PageRank, 1 para HITS, 2 para IN-DEG), la probabilidad de teletransportación  $c$  en el caso de PageRank (que valdrá -1 si `alg` no es 0), el tipo de instancia, el `path` al archivo/directorio conteniendo la definición de la red (que debe ser relativa al ejecutable, o el path absoluto al archivo) y el valor de tolerancia utilizado en el criterio de parada impuesto a cada método. El siguiente ejemplo muestra un caso donde se pide ejecutar PageRank, con una probabilidad de teletransportación de 0.85, sobre la red descrita en `red-1.txt` (que se encuentra en el directorio `tests/`) y con una tolerancia de corte de 0,0001.

```
0 0.85 0 tests/red-1.txt 0.0001
```

Para la definición del grafo que representa la red, se consideran dos bases de datos de instancias con sus correspondientes formatos. La primera de ellas es el conjunto provisto en SNAP [?] (el tipo de

---

<sup>12</sup>La cátedra adjunta casos de *benchmark* que representan sub-redes obtenidas en base a búsquedas temáticas

instancia es 0), con redes de tamaño grande obtenidos a partir de datos reales. Además, se consideran las instancias propuestas en [?]. Estas instancias son de tamaño mediano, obtenidas también en base a datos reales, y corresponden a redes temáticas obtenidas a partir de una búsqueda particular. Para cada nodo de la red se tiene: la dirección URL, una breve descripción, y las páginas a las cuales apunta. Si bien algunas de las URL ya no son válidas, la descripción permite tener algo más de información para realizar un análisis cualitativo.

En el caso de la base de SNAP, los archivos contienen primero cuatro líneas con información sobre la instancia (entre ellas,  $n$  y la cantidad total de links,  $m$ ) y luego  $m$  líneas con los pares  $i, j$  indicando que  $i$  apunta a  $j$ . A modo de ejemplo, a continuación se muestra el archivo de entrada correspondiente a la red propuesta en Bryan y Leise [?, Figura 1]:

```
# Directed graph (each unordered pair of nodes is saved once):
# Example shown in Bryan and Leise.
# Nodes: 4 Edges: 8
# FromNodeId      ToNodeId
1      2
1      3
1      4
2      3
2      4
3      1
4      1
4      3
```

Para la otras instancias, en [?] puede encontrarse una descripción del formato propuesto (el tipo de instancia será 1 en este caso).

Una vez ejecutado el algoritmo, el programa deberá generar un archivo de salida que contenga una línea por cada página ( $n$  líneas en total), acompañada del puntaje obtenido por el algoritmo Page-Rank/IN-DEG. En el caso de HITS, el archivo contendrá  $2n$  líneas, las primeras  $n$  con el *peso de autoridad* y las segundas  $n$  con el *peso de hub* para los vértices  $1, \dots, n$ .

Para generar instancias, es posible utilizar el código Python provisto por la cátedra. La utilización del mismo se encuentra descripta en el archivo README. Es importante mencionar que, para que el mismo funcione, es necesario tener acceso a Internet. En caso de encontrar un bug en el mismo, por favor contactar a los docentes de la materia a través de la lista. Desde ya, el código puede ser modificado por los respectivos grupos agregando todas aquellas funcionalidades que consideren necesarias.

---

### **Fechas de entrega**

- **Formato Electrónico:** Sábado 11 de Octubre de 2014, hasta las 23:59 hs, enviando el trabajo (informe + código) a la dirección `metnum.lab@gmail.com`. El subject del email debe comenzar con el texto [TP2] seguido de la lista de apellidos de los integrantes del grupo.
- **Formato físico:** Miércoles 15 de Octubre de 2014, a las 17 hs. en la clase teórica.

**Importante:** El horario es estricto. Los correos recibidos después de la hora indicada serán considerados re-entrega.