

# Métodos Numéricos

## TP2

6 de octubre de 2015

*Years Later for Guillermo Vilas, He's Still Not the One*



Integrante	LU	Correo electrónico
Martin Baigorria	575/14	<a href="mailto:martinbaigorria@gmail.com">martinbaigorria@gmail.com</a>
Federico Beuter	827/13	<a href="mailto:federicobeuter@gmail.com">federicobeuter@gmail.com</a>
Mauro Cherubini	835/13	<a href="mailto:cheru.mf@gmail.com">cheru.mf@gmail.com</a>
Rodrigo Kapobel	695/12	<a href="mailto:rok_35@live.com.ar">rok_35@live.com.ar</a>

**Reservado para la cátedra**

Instancia	Docente	Nota
Primera entrega		
Segunda entrega		

**Resumen: TODO**  
**Keywords: TODO**

# Índice

<b>1. Introduccion</b>	<b>3</b>
<b>2. PageRank</b>	<b>5</b>
2.1. Modelado para paginas web . . . . .	5
2.1.1. Propiedades . . . . .	5
2.1.2. Existencia y Unicidad . . . . .	6
2.2. Modelado para Tenis . . . . .	7
2.2.1. Sistema de puntos y consideraciones . . . . .	8
2.3. Eliminacion Gausiana . . . . .	8
2.4. Representacion del grafo . . . . .	8
2.5. Computo: Método de la Potencia . . . . .	10
2.5.1. Correctitud . . . . .	10
2.5.2. Complejidad . . . . .	10
<b>3. Experimentación</b>	<b>11</b>
3.1. PageRank . . . . .	11
3.1.1. Complejidad . . . . .	11
3.1.2. Casos Patologicos . . . . .	11
3.2. Paginas Web . . . . .	11
3.2.1. Comparacion PageRank vs In-Deg (RODRI) . . . . .	11
3.2.2. Manipulacion . . . . .	11
3.3. Ranking ATP . . . . .	11
3.3.1. Ranking ATP oficial vs. Ranking PageRank vs. Sort por diferencia de victorias/derrotas . . . . .	12
3.3.2. Eleccion del factor de 'teletransportacion' $c$ (RODRI) . . . . .	13
3.4. Metodo de la Potencia . . . . .	13
3.4.1. Representacion de la Matriz de Transicion (RODRI/FEDE) . . . . .	13
3.4.2. Evolucion de la norma entre iteraciones . . . . .	14
3.4.3. Convergencia . . . . .	14
3.4.4. Eleccion del $x_0$ . . . . .	14
<b>4. Conclusiones</b>	<b>15</b>
<b>5. Apéndice A: Enunciado</b>	<b>16</b>
<b>6. Apéndice B: Código</b>	<b>21</b>
6.1. system.cpp . . . . .	21
6.2. matrix.h . . . . .	29
6.3. sparseMatrix.h . . . . .	35

## 1. Introduccion

El 25 de Mayo de 2015 el diario The New York Times publico un articulo titulado "Years Later for Guillermo Vilas, He's Still Not the One", donde se repasa el rendimiento del tenista argentino durante los años 1975/1976 y se discute el calculo del ranking de la ATP en ese momento. Aunque hoy en día Vilas es un icono del tenis argentino, nunca logró estar en la cima del ranking de la ATP.



Figura 1: Guillermo Vilas after winning a tournament in Stockholm in 1975. A journalist has asserted that Vilas deserved to be ranked No. 1 during that year.

En 2016, un grupo de investigadores argentinos decidió analizar el ranking de la ATP en 1975 y 1976 para determinar si Vilas debió haber sido número 1. Dado que los rankings no se actualizaban constantemente en ese momento, los investigadores mostraron que de haberse actualizado de forma periódica, Vilas hubiese sido número 1 por durante 7 semanas en 1975 y 1976.

Existen precedentes donde se actualizó un ranking de tenis de forma retroactiva. Este es el caso de la WTA, que determinó que Evonne Goolagong Cawley debió haber sido número 1 por dos semanas en 1976. Por esta razón el grupo de investigación argentino considera que revisar estos rankings no es un esfuerzo en vano. Cuando buscábamos los datos de la ATP en el 1975/1976, uno de los investigadores de este equipo que contactamos nos comento: "Es interesante tu decisión de indagar sobre el tema. Tal vez no estás al tanto del trabajo y lucha que estamos realizando contra la ATP, por el ranking de los 70 en el que perjudicaron a Vilas y muchísimos otros jugadores."

En ese momento, el calculo del ranking de la ATP era bastante rudimentario: "It was a system based on an average of a player's results, and it often rewarded top players who played fewer tournaments. Vilas was a workhorse, which is how he managed not to reach No. 1 in the ATP rankings in 1977, when he won the French Open, the United States Open and 14 other tournaments." [5].

Los métodos para calcular rankings no solo son relevantes para definir las posiciones de equipos y jugadores en eventos deportivos, sino que aparecen constantemente en todo tipo de situaciones donde se debe imponer algún tipo de orden. Este es el caso por ejemplo de los concursos docentes, donde se ponderan los diferentes antecedentes para decidir cual es el candidato *idoneo* para el puesto.

Otro caso sumamente relevante en cuanto algoritmos de ranqueo es el de los motores de búsqueda. Los motores de búsqueda deben encontrar alguna forma de ordenar de forma relevante los sitios web que están relacionados con una consulta. El caso icónico es el de Google con su algoritmo PageRank. Los buscadores antes de 1990 eran sumamente rudimentarios, utilizaban algoritmos de ranqueo vulnerables en el sentido que podían ser manipulados y no se explotaba gran parte de la estructura de la web. Esta fue una de las razones por las cuales una consulta no siempre devolvía resultados relevantes. Este fue el caso por ejemplo de algunos buscadores en ese momento como Yahoo! Search o AltaVista.

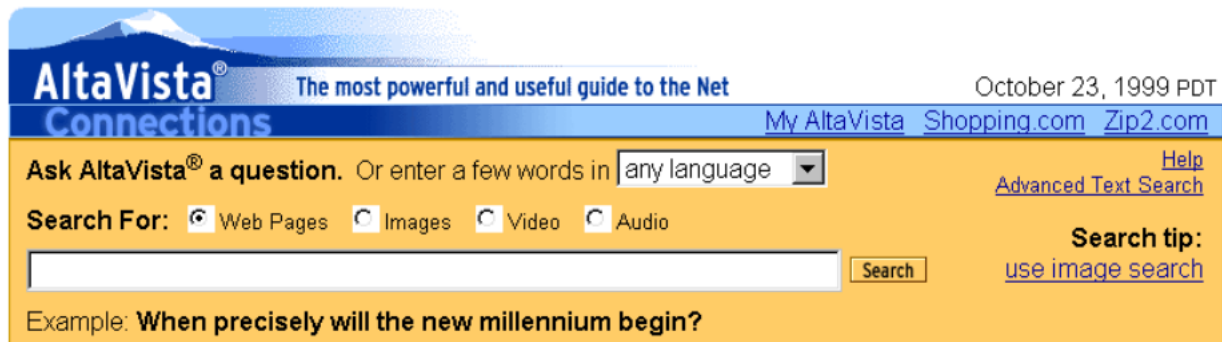


Figura 2: Sitio Web de Altavista, ano 1999.

El clásico paper de Brin y Page, “The anatomy of large-scale hypertextual Web search engine.” [3] explica brevemente el origen del motor de búsqueda de Google y del algoritmo PageRank. La idea es básicamente la siguiente, en primer lugar se implementa un crawler distribuido para poder solicitar y armar el grafo de la web. Las palabras de cada sitio son indexadas y guardadas en una base de datos. Al llegar una consulta al buscador, un programa busca la consulta en los índices de paginas. De esta forma llegamos a un conjunto de paginas que están relacionadas con la consulta. Luego, antes de devolverle al usuario los resultados, estas paginas son ordenadas utilizando el famoso algoritmo PageRank. Este algoritmo se basa en la idea de que para medir la relevancia de un sitio se puede usar como proxy la cantidad de sitios que tienen un link al mismo. Para evitar que un usuario malintencionado manipule los resultados del mismo, la relevancia otorgada por un sitio web que linkea a otro es proporcional a su propia relevancia e inversamente proporcional a la cantidad de links (o grado de salida) del mismo.

El presente trabajo práctico tendrá como objetivo implementar el algoritmo PageRank para luego utilizarlo para generar rankings de todo tipo, ya sea para ordenar la relevancia de paginas webs o generar rankings deportivos. PageRank es un algoritmo que basa su ranking en encontrar el autovector de una matriz de transiciones. A priori esto puede sonar complicado, pero luego mostraremos que en realidad es bastante simple y elegante. Dado que ordenar la relevancia de millones de sitios web no es un problema trivial, en la practica este problema se resuelve utilizando álgebra lineal y métodos numéricos. Una muy buena introducción teórica se puede encontrar en el trabajo de Bryan y Leise [4]. Otros autores como Kamvar et al. [7] han buscado otros enfoques y métodos para poder acelerar este algoritmo. La idea es encontrar una forma eficiente de poder computar este modelo, calibrando sus diferentes parámetros de modelado y convergencia para lograr un orden relevante.

Una vez planteado el procedimiento, experimentaremos con la complejidad temporal de los métodos implementados y evaluaremos los diferentes parámetros a calibrar. Finalmente concluiremos si según el algoritmo PageRank y nuestra matriz de transición Vilas efectivamente debió haber estado en la punta del ATP en 1975/1976. En caso afirmativo, sin dudas nos comunicaremos con la ATP.



Figura 3: Guillermo Vilas apoya este TP.

## 2. PageRank

### 2.1. Modelado para paginas web

El algoritmo PageRank fue ideado en un principio para buscar de darle alguna medida de relevancia a los sitios web en internet. El mismo tiene dos interpretaciones equivalentes, que serán expuestas a continuación.

El problema se modela a partir de un grafo  $G(Web, Links)$  donde  $Web$  es el conjunto de sitios web y  $Links$  es la cantidad de conexiones entre sitios. Consideremos que toda pagina web  $u \in Web$  esta representada por un vértice y la relación entre paginas por un link con una arista. Una representación posible del grafo es mediante matrices de adyacencia. Definimos la matriz de adyacencia o conectividad  $W \in \{0,1\}^{n \times n}$  de forma tal que  $w_{ij} = 1$  si la pagina  $j$  tiene un link a la pagina  $i$  y  $w_{ij} = 0$  en caso contrario. Por lo tanto, la cantidad de paginas a las que la pagina  $u$  apunta ( $d_{out}(u)$ ) se puede calcular como  $n_j = \sum_{i=1}^n w_{ij}$ .

#### 2.1.1. Propiedades

Sea  $x_j$  el puntaje asignado a la pagina o vértice  $j \in Web$  y otra pagina  $u \in Web$ . La idea es buscar una medida que cumpla con las siguientes propiedades:

- La relevancia de todo sitio web es positiva.
- La relevancia de un sitio web debe aumentar a medida que mas sitios unicos lo apuntan.
- La relevancia derivada de otro sitio web debe depender de su propia relevancia. Es decir, es mas valioso que me linkee un sitio relevante que uno no relevante. En caso de no cumplirse esta propiedad, el ranking seria fácilmente manipulable al permitir que un usuario cree muchos sitios que linkeen a uno para darle relevancia.
- La relevancia de todos los sitios web debe sumar uno. De esta manera estamos ante una distribución de probabilidad de los sitios. Mas adelante veremos que al interpretar esto mediante Cadenas de Markov existe una interpretación directa: la relevancia se puede ver como la proporción del tiempo total que un usuario pasa en ese sitio.

Por lo tanto, estamos buscando una medida de relevancia tal que la importancia obtenida por la pagina  $u$  obtenida por el link de la pagina  $v$  sea proporcional a la relevancia de  $v$  e inversamente proporcional al grado de  $v$ . El aporte del link de  $v$  a  $u$  entonces es  $x_u = x_v/n_v$ . Luego, sea  $L_k \subseteq Web$  el conjunto de paginas que tienen un link a la pagina  $k$ . Por lo tanto, la relevancia total de un sitio sera:

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

Notar que esta es de cierta manera una definición recursiva. La relevancia de un sitio  $u$  puede depender de la relevancia de un sitio  $v$ , y luego la de  $v$  puede depender de la de  $u$ . A priori calcular la relevancia de un sitio puede parecer sumamente complicado, pero luego veremos que al plantearlo como un sistema de ecuaciones esta dificultad per se ya no se presenta.

Definimos entonces una matriz de transición o adyacencia con pesos en las aristas  $P \in \mathbb{R}^{n \times n}$  tal que  $p_{ij} = 1/n_j$  si  $w_{ij} = 1$  y  $p_{ij} = 0$  en caso contrario. Luego, el modelo planteado en (1) para toda pagina web se puede expresar  $Px = x$  donde  $x \in \mathbb{R}^n$ . Notar que esto es equivalente a encontrar el autovector de autovalor 1 tal que  $x_i > 0$  y  $\sum_{i=1}^n x_i = 1$ . Notar que si logramos probar que bajo ciertas condiciones nuestra matriz de transición tiene autovalor 1, el signo de todos los elementos de un autovector es el mismo y la dimension del autoespacio es 1 ya tenemos un ranking valido. Esto se debe a que cualquier autovector puede ser reescalado a uno de norma unitaria con  $x_i \geq 0$ .

### 2.1.2. Existencia y Unicidad

Bryan y Leise [4] analiza y prueba las condiciones bajo las que podemos garantizar que:

- La matriz de transición tiene autovalor 1.
- La dimension del autoespacio asociado al autovalor 1 es 1. Es deseable que el ranking asociado a una matriz de transición sea único.
- El signo de todos los elementos del autovector asociado al autovalor 1 es el mismo.

Veamos bajo que condiciones nuestra matriz de transición cumple con estas propiedades:

**Definición** Una matriz cuadrada se llama estocástica por columnas si todos sus elementos son positivos y la suma de cada columna es igual a 1.

A partir de esta definición se puede probar la siguiente proposición:

**Proposición 2.1** *Toda matriz estocástica por columnas tiene a 1 como autovalor.*

Esto significa que si no existen **dangling nodes**, es decir, vértices con  $d_{out} = 0$ , podemos garantizar que nuestra matriz de transición es estocástica por columnas.

Notar que bajo las condiciones actuales no podemos garantizar que si existe el autoespacio asociado al autovalor 1, el mismo tenga dimension 1. Intuitivamente, esto se debe a que el grafo de la web puede tener varias componentes conexas. ¿Como comparamos sitios web que no están relacionados? Justamente la relación, ya sea directa o indirecta mediante transitividad me da algún tipo de relación de orden. Al no tener una relación de orden entre dos sitios web bien definida, es razonable que existan múltiples autovectores, es decir, rankings. Esto se puede ver claramente en la pagina 4 del paper de Bryan y Leise [4].

Por lo tanto, la idea es básicamente buscar algún tipo de transformación relevante de mi matriz de transición que me permita garantizar que no voy a tener **dangling nodes** y ademas que solo tenga una componente conexa, es decir, que el grafo sea conexo. Definimos la siguiente matriz de transición, donde  $v \in \mathbb{R}^{n \times n}$ , con  $v_i = 1/n$  y  $d \in \{0, 1\}^n$ ,  $d_i = 1$  si  $n_i = 0$  y  $d_i = 0$  en caso contrario, como:

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

De esta manera, en caso de tener una pagina web que es un **dangling node**, le asignamos un link uniforme a todos los sitios web  $u \in Web$ . Una interpretación equivalente es tomar a la matriz de transiciones como la matriz que describe una Cadena de Markov, donde el link pesado representa la probabilidad de dirigirse de una pagina a la otra. Por lo tanto, esta transformación se puede interpretar como que existe una probabilidad uniforme de ir de uno de estos sitios a cualquiera de la web. Esto normalmente se conoce como el **navegante aleatorio**.

Tambien podemos considerar la posibilidad de que el navegante aleatorio se dirija a una pagina web que no esta linkeada a la pagina a la que esta actualmente. Este fenómeno se conoce como teletransportación. Para incluirlo al modelo, tomemos un numero  $c \in [0, 1]$  y transformemos la matriz de transiciones de la siguiente manera, donde  $\bar{1} \in \mathbb{R}^n$  es un vector tal que todos sus componentes valen 1:

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

Notar que en caso de tener  $c = 1$ , estamos en la matriz de transición sin teletransportación. Por otro lado, si  $c = 0$  estamos en el caso donde solo hay teletransportación y no importa la estructura del grafo de la web.

Esta nueva matriz de transición, dado que es estocástica por columnas y no tiene **dangling nodes**, nos garantiza que la dimension del autoespacio generado por el autovector de autovalor 1 es unitaria. Solo nos falta mostrar que todo autovector tiene todos sus elementos del mismo signo. Es facil probar la siguiente proposición:

**Proposición 2.2** *Si la matriz  $M$  es positiva y estocástica por columnas, entonces todo autovector en  $V_1(M)$  tiene todos sus elementos positivos o negativos.*

Por lo tanto, ya probamos la existencia del autovector de norma 1 asociado al autovalor 1 de la matriz de transición transformada. El siguiente lema nos garantiza su unicidad. Su respectiva demostración se encuentra nuevamente en la pagina 7 del paper de Bryan y Leise [4].

### Lemma 2.3

Si  $M$  es positiva y estocástica por columnas, entonces  $V_1(M)$  tiene dimension 1.

**Proposición 2.4** Sea  $M$  una matriz real de  $n \times n$  positiva y estocástica por columnas, y  $V$  el subespacio de  $\mathbb{R}^n$  que consiste de aquellos vectores  $v$  tales que la suma de sus componentes sea 0, entonces  $Mv \in V$  y  $\|Mv\|_1 \leq \alpha \|v\|_1$  donde  $\alpha = \max_{1 \leq j \leq n} |1 - 2\min_{1 \leq i \leq n}| < 1$ .

## 2.2. Modelado para Tenis

El modelo GeM Ranking Method presente en el paper de Govan et al. plantea un método alternativo para rankear ligas deportivas. Utilizaremos el mismo para analizar los rankings del ATP entre 1975 y 1977. El algoritmo está motivado en el PageRank de Page y Brin, pero con unos ligeros retoques para adaptarlo a torneos y competencias. En lo que sigue, explicaremos como se construye la matriz de transiciones que utiliza el algoritmo, pero antes de comenzar con las definiciones, vale aclarar que para mantener la homogeneidad haremos esta definición sobre la transpuesta de la matriz que se usa en el paper de Govan et al.

Se representa una temporada como un grafo directo con  $n$  nodos. Cada nodo representa un participante o equipo y cada partido entre dos participantes representa una arista desde el perdedor al ganador igual a la diferencia positiva de los puntos obtenidos por cada participante en el partido.

La matriz  $A^t$  de adyacencias queda definida de la siguiente manera:

$$A^t = \begin{cases} w_{ij} & \text{si el participante } j \text{ pierde con el participante } i \\ 0 & \text{en cualquier otro caso} \end{cases}$$

Donde cada  $w_{ij}$  representa la diferencia positiva del puntaje de cada participante en ese enfrentamiento. En caso de que un equipo pierda mas de una vez en el mismo torneo, será la suma de las diferencias de cada partido entre  $i$  y  $j$ . Esta es la mayor generalización en cuanto a PageRank.

Luego se define la matriz  $H$  de la siguiente forma:

$$H = \begin{cases} w_{ij} / \sum_{k=1}^n A_{kj}^t & \text{si hay un link de } j \text{ a } k \\ 0 & \text{en cualquier otro caso} \end{cases}$$

(Cabe destacar que en el paper de Govan et al. hay un error de tipeo y la division se define como  $1 / \sum_{k=1}^n A_{kj}^t$ )

Luego definimos  $G$  a la matriz resultante de:

$$G = \alpha(A^t + uv^t) + (1 - \alpha)ev^t$$

Donde  $0 < \alpha < 1$ ,  $v$  es un vector de probabilidades,  $a$  es tal que  $a_i$  es 1 si la fila  $j$  de  $H$  es 0 y  $a_i$  es 0 en cualquier otro caso y  $u$  sea un vector de probabilidad de  $n \times 1$ . e se define como un vector fila con todas sus entradas igual a 1.

El vector que contendrá los puntajes de cada equipo será un  $\pi$  tal que:

$$G\pi = \pi$$

Cada entrada  $H_{ij}$  de  $H$  se puede interpretar como la probabilidad de que el participante  $j$  pierda contra el participante  $i$ . Para los participantes invictos un simple ajuste que se propone es elegir un vector  $u$  donde todas las entradas son  $1/n$ . Esto significa cambiar la probabilidad de los invictos a que puedan perder contra cualquiera de los otros participantes (incluido si mismo) con probabilidad uniforme. El modelo básico utiliza  $\alpha$  de la misma forma que en PageRank y  $v^t$  como vector de personalización del sistema. Una simple elección puede ser  $v = (1/n)e$ . Aunque  $v$  ofrece mucha mas flexibilidad. Podría usarse con el resultado de un ranking previo, aumentando así las probabilidades de que cierto participante gane el presente torneo. La elección del  $\alpha$  determina la importancia de la matriz de personalización  $ev^t$ . Una buena elección en el calculo de  $\alpha$  podría ser la que se menciona en la siguiente investigación de microsoft research: Tracking the random surfer: Empirically measured teleportation parameters in PageRank del cual mencionaremos la siguiente abstracción:

PageRank computes the importance of each page in a directed graph under a random surfer model governed by a teleportation parameter. Commonly denoted alpha, this parameter models the probability of following an edge inside the graph or,

when the graph comes from a network of web pages and links, clicking a link on a web page. We empirically measure the teleportation parameter based on browser toolbar logs and a click trail analysis. For a particular user or machine, such analysis produces a value of  $\alpha$ . We find that these values nicely fit a Beta distribution with mean edge-following probability between 0.3 and 0.7, depending on the site. Using these distributions, we compute PageRank scores where PageRank is computed with respect to a distribution as the teleportation parameter, rather than a constant teleportation parameter. These new metrics are evaluated on the graph of pages in Wikipedia.

Igualmente en nuestro caso computaremos los rankings con un  $\alpha$  de 0.85 que es el utilizado en Govan et al. y que a su vez fue el standard de Google en sus comienzos.

### 2.2.1. Sistema de puntos y consideraciones

Para procesar los rankings, obtuvimos todos los partidos realizados entre 1975 y 1977 de todas las competencias a nivel internacional de tenis que sirven para clasificar al ATP. (Australian Open, Indianapolis, Roma, Roland Garros, US Open, Wimbledon, Davis Cup, y muchos más) junto con los rankings de cada año. Uno a mitad de año y otro a finales de año.

Los datos están completos a nivel enfrentamientos y participantes, pero carecen de los puntajes que cada jugador obtuvo por cada partido ganado, por avanzar de ronda y/o ganar un torneo. Esta es una parte importante a aclarar del sistema de rankings de la ATP. Se obtienen puntos por cada partido ganado, ronda superada y torneo ganado, los cuales además varían por torneo en importancia, siendo los Grand Slam los de mayor jerarquía. Utiliza un sistema de defensa de puntos que fué variando con los años. En aquella época, la defensa se hacía por año. Es decir, los puntos realizados durante un año, se defienden al año siguiente.

Igualmente esto no supone una limitación grave para computar los rankings y comprobar si Guillermo Vilas fue realmente o no 1ro al menos uno de esos años.

Existe una relación directa entre sumar puntos y ganar partidos, y es que naturalmente el que más puntos tiene, más partidos ganó, naturalmente. Además, dada la naturaleza del comportamiento de page rank, no solo importa a cuantos contrincantes derrotó un participante dado, si no a quienes derrotó. Ésta es la primera consideración a tener en cuenta y una carta a favor que será utilizada por el algoritmo. Dado que no disponemos del sistema de puntuación real y no sabemos a priori cuantos partidos jugó exactamente cada participante, es esencial que además de la puntuación que definamos por cada partido, un jugador obtenga un extra por el rango del rival derrotado. Esto mismo se puede interpretar como que ese jugador ganó un partido importante, dado que los partidos importantes son los de instancias decisivas (cuartos, semifinales o una final), y generalmente, al que solo llegan los mejores jugadores.

Dadas estas consideraciones podemos definir un sistema de puntos con el que poder computar partidos. El mismo es muy sencillo. 3 puntos al ganador y 1 punto al perdedor. Solo utilizaremos el ranking de fin de año, totalizando por los puntos acumulados en el mismo, obteniendo así 3 rankings por cada año.

Veremos luego en la experimentación, como este sistema resuelve el cálculo de los rankings de manera apropiada y explicaremos que es lo que sucede en cada año computado observando algunos detalles importantes y comparandolos con los rankings oficiales de la ATP para luego concluir si Vilas fue o no, nuevamente, 1ro entre 1975 y 1977.

### 2.3. Eliminacion Gausiana

Esta seccion solo la pongo para que la consideres. Se podra hacer eliminacion gausiana con pivoteo para  $(P-I)x = 0$ ? Igual si es posible es de orden cubico, con la web de millones de paginas se te va al carajo. Es solo para enriquecer la discusion.

### 2.4. Representacion del grafo

Ya hemos demostrado las condiciones necesarias para poder obtener el autovector asociado al autovalor dominante de una matriz de Markov. Ahora debemos proceder a calcular el mismo. Para esto, tenemos que tener en cuenta las cualidades del sistema y el método de resolución del algoritmo. Recordemos que en general, el grafo que representa la web tenderá a ser desconexo y muy grande, es decir, que podrán existir dos o mas rankings diferentes. Por lo tanto la matriz de transiciones puede ser muy esparsa e inclusive puede suceder que una página no tenga links de salida, dando lugar a dangling nodes. Para solucionar estos inconvenientes, con lo visto anteriormente disponemos de dos soluciones. Para los dangling nodes, la solución consiste en sumar una columna con probabilidad  $1/n$  a la columna de ceros, esto en si, se puede interpretar como la probabilidad de navegación aleatoria que previamente describimos. Aunque con esto no solucionamos el problema de la esparsidad de la matriz en si y el de poder tener mas de un ranking diferente. Para esto último, se agregó la matriz de probabilidad de teletransportación.



Dada esta definición, la matriz de transiciones resultante no es esparsa. Para sistemas muy grandes, esto puede resultar contraproducente a la hora de obtener el autovector asociado, dado que la complejidad espacial y temporal aumenta considerablemente con la cantidad de información representada en la matriz. Sin embargo existe un resultado que podremos utilizar para mejorar la eficiencia del algoritmo en términos de complejidad temporal y espacial. El mismo se basa en la idea de Kamvar et al. [7, Algoritmo 1] para el calculo del autovector. Este resultado nos permite utilizar la matriz original de transiciones sin modificar en lo absoluto, pero si cambiando su representación, valiendonos de una buena estructura para almacenar las entradas de la misma.

Las cualidades de la matriz hacen que sea razonable intentar pensar en una forma de representar solo las entradas que no sean ceros, y dado que la matriz suele ser esparsa, la misma contendrá muchos ceros que podrían no ser representados. Para esto optamos por una de entre las 3 siguientes estructuras de representación:

- Dictionary of Keys (*DOK*)
- Compressed Sparse Row (*CSR*)
- Compressed Sparse Column (*CSC*)

De todas estas representaciones posibles, para este t.p optamos por *CSR*. Aún así no haremos una elección sin una justificación apropiada del porque consideramos que es la mejor para nuestro trabajo, dado que como en toda estructura de datos, siempre existen pros y contras. Nos encargaremos en lo que sigue de exponer estos detalles para dejar en claro nuestro punto de vista.

- Dictionary of Keys (*DOK*)

Consiste en un diccionario que mapea pares de fila-columna a la entrada. No se representan las entradas nulas. El formato es bueno para gradualmente construir una matriz esparsa en orden aleatorio, pero pobre para iterar sobre valores distintos de cero en orden lexicográfico. Uno construye típicamente una matriz en este formato y luego se convierte en otro formato más eficiente para su procesamiento.

- Compressed Sparse Row (*CSR*)

Pone las entradas no nulas de las filas de la matriz en posiciones de memoria contiguas. Suponiendo que tenemos una matriz dispersa no simétrica, creamos vectores: uno para los números de punto flotante (*val*), y los otros dos para enteros (*col\_ind*, *row\_ptr*). El vector *val* almacena los valores de los elementos distintos de cero de la matriz, de izquierda a derecha y de arriba hacia abajo. El vector *col\_ind* almacena los índices de columna de los elementos en el vector *val*. Es decir, si  $val(k) = a_{ij}$  entonces  $col\_ind(k) = j$ . El vector *row\_ptr* almacena los lugares en el vector *val* que comienza y termina una fila, es decir, si  $val(k) = a_{ij}$  entonces  $row\_ptr(i) \leq k \leq row\_ptr(i + 1)$ . Por convención, se define  $row\_ptr(n + 1) = nnz$ , en donde *nnz* es el número de entradas no nulas en la matriz. Los ahorros de almacenamiento de este enfoque es significativo. En lugar de almacenar elementos  $n^2$ , solamente necesitamos  $2nnz + n$  lugares de almacenamiento.

Veamos con un ejemplo como seria la representacion:

$$\begin{pmatrix} 0 & 0 & 0 & 0 \\ 5 & 8 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 6 & 0 & 0 \end{pmatrix}$$

Es una matrix de 4x4 con 4 entradas no nulas. Luego:

$$\begin{aligned} val &= [ 5 \ 8 \ 3 \ 6 ] \\ row\_ptr &= [ 0 \ 0 \ 2 \ 3 \ 4 ] \\ col\_ind &= [ 0 \ 1 \ 2 \ 1 ] \end{aligned}$$

- Compressed Sparse Column (*CSC*) La idea es analoga a *CSR*, pero la compresion se hace por columnas es decir, si *CSR* comprime *A*, *CSC* comprime  $A^t$

Sobre la matriz definida para *CSR*, con *CSC* obtenemos lo siguiente:

$$\begin{aligned} val &= [ 5 \ 8 \ 6 \ 3 ] \\ col\_ptr &= [ 0 \ 1 \ 3 \ 4 \ 4 ] \\ row\_ind &= [ 1 \ 1 \ 2 \ 3 ] \end{aligned}$$

Todos los resultados anteriores permiten evitar representar valores nulos. El motivo de nuestra elección se debe a que *CSR* ofrece una buena representación de las filas de la matriz y es más eficiente a la hora de hacer operaciones del tipo  $A \cdot x$  (matriz-vector) que es lo que nos interesa en el método de la potencia que realiza pageRank. *CSC* en cambio, es efectiva para el producto  $x \cdot A$  (vector-matriz) dado que la misma ofrece una mejor representación de las columnas. En contra partida, tanto *CSR* como *CSC*, no permiten construcción incremental aleatoria, que si ofrece *DOK*, es decir, que cambios a la esparsidad de la matriz son costosos. En general están pensadas para ser estáticas, pero esto no es un inconveniente en nuestro caso, dado que no se realizaran cambios en la esparsidad de la matriz durante el proceso.

En el presente trabajo utilizaremos la idea de Kamvar et al. [7, Algoritmo 1] para el calculo del autovector valiendonos de nuestra estructura de representación elegida y compararemos los resultados con el algoritmo standard para mostrar que al final de cuentas, si el sistema es muy grande y esparso, puede resultar muy beneficioso en terminos de complejidad espacial y temporal.

## 2.5. Computo: Método de la Potencia

Habiendo definido la matriz de transiciones como la matriz  $P_2$  de la sección anterior, debemos hallar el vector de relevancias  $x$  que define en  $x_j$  la relevancia de la  $j$ -ésima pagina  $\forall 0 \leq j \leq n$ . Debido a la definición recursiva de este vector,  $x$  debe satisfacer que  $P_2 x = x$  o lo que es equivalente, que  $x$  sea un autovector asociado al autovalor 1 (i.e  $x \in V_1(P_2)$ ). Para este tipo de problemas es que contamos con el método de la potencia, un algoritmo iterativo que nos devuelve una aproximación lineal al  $x$  deseado.

### 2.5.1. Correctitud

La idea detras de este algoritmo consiste en generar a partir de un vector inicial  $x_0$ , una secuencia  $x_k = \frac{P_2 x_{k-1}}{\|P_2 x_{k-1}\|}$  que aproxime al autovector  $q$  (en nuestro caso el vector  $x$  de relevancias) asociado al autovalor 1. Siguiendo esta propuesta elegimos nuestro  $x_0$  de manera que satisfaga que todas sus componetes sean positivas y  $\|x_0\|_1 = 1$ , definiendolo además como  $x_0 = q + v$ , en donde  $q$  es el único vector unitario que satisface que  $q \in V_1(P_2)$  y que posee todas sus componentes positivas; y  $v$  siendo un vector tal que la suma de sus componentes sea 0. De esta manera, dado que  $P_2$  es una matriz positiva y estocástica por columnas, y  $x_0$  un vector unitario con todas sus componentes también positivas, resulta que  $P_2 x_0$  también es un vector unitario con todas sus componentes positivas. Esto nos facilita la secuencia de la siguiente forma  $x_k = P_2 x_{k-1} = P_2^k x_0$ , pues desde un principio  $x_1 = P_2 x_0$  ( $\|P_2 x_0\| = 1$ ); luego  $x_k = P_2^k x_0 = P_2^k q + P_2^k v = q + P_2^k v \implies P_2^k v = P_2^k x_0 - q$ . Ahora bien, por la desigualdad de la proposición 2.4,  $\|P_2^k x_0 - q\|_1 = \|P_2^k v\|_1 \leq \alpha^k \|v\|_1$ , haciendo tender  $k \rightarrow \infty$ ,  $\alpha \rightarrow 0$ , de lo que se deduce que  $P_2^k x_0 - q \rightarrow 0$  o equivalentemente  $P_2^k x_0 \rightarrow q$ , es decir al  $x$  buscado.

### 2.5.2. Complejidad

En cuanto a complejidad temporal es evidente que el costo del algoritmo radica en cuantas operaciones elementales nos conlleva calcular  $P_2^k x$ , en especial si se requiere realizar una gran cantidad de iteraciones. Recordando como la definimos,  $P_2 = cP_1 + (1 - c)E$  donde  $E$  es la matriz con todos sus elementos iguales a  $\frac{1}{n}$  y  $P_1$  es una matriz positiva y estocástica por columnas, de lo que podemos concluir que todos los elementos de  $P_2$  son estrictamente positivos (si  $c \in (0,1)$ ). Esta conclusión no es muy alentadora, pues sólo tener que hacer  $P_2 x$  tiene una complejidad espacial y temporal de  $\Theta(n^2)$ . No obstante, sabiendo que  $\|x\|_1 = 1$ , resulta que  $P_2 x = (cP_1 + (1 - c)E)x = cP_1 x + (1 - c)Ex = cP_1 x + (1 - c)e$  en donde  $e$  ahora es un vector con todos sus elementos iguales a  $\frac{1}{n}$ , luego la atención se concentra en el costo de  $P_1 x$ . Esta última matriz ya no sólo no tiene todos sus elementos estrictamente positivos, sino que suele tratarse de una matriz esparsa, ahorrandonos varias multiplicaciones (las que corresponderian a los elementos iguales a cero) aventajando en velocidad a la anterior situación; a su vez podremos aprovechar la redundancia de ceros para utilizar alguna estructura más eficiente a la hora de almacenar la matriz de la actual iteración. En conclusión dandonos una complejidad temporal total de  $\mathcal{O}(k * n^2)$ , pero un caso promedio notablemente menor.

## 3. Experimentación

### 3.1. PageRank

#### 3.1.1. Complejidad

tiempo de computo en funcion de size del grafo, eje x, cantidad de sitios web, eje y, tiempo en ms a convergencia.

#### 3.1.2. Casos Patologicos

Caso particular chiquito, pagina 3. Fijate el parrafo que arranca en A simple approach..... y despues This approach ignores that... La idea es armar el mismo grafo y mostrar el mismo ejemplo jaja

### 3.2. Paginas Web

#### 3.2.1. Comparacion PageRank vs In-Deg (RODRI)

Comparar solo los rankings, nada de complejidad. Podes mencionar que In-Deg usa un algoritmo  $\mathcal{O}(n \times \log(n))$ , pero nada mas. Comparar top 10 con los dos y discutir diferencias.

#### 3.2.2. Manipulacion

Pagina 5, ejercicio 1. La idea es que plantees un caso de un tipo que quiere manipular el ranking, muestra que aunque agregues miles de nuevas paginas apuntando no puedes hacer demasiado, hacelo en funcion de la cantidad de paginas que agregas?

Se puede manipular entonces o no? Agarra, en el eje x pone cantidad de sitios web que apuntan solamente al sitio u que le quiero subir el ranking, y en el eje y el ranking de ese sitio. Fijate que aumenta, y fijate si puedes hacer algun tipo de curva de nivel con c (cuanto mayor c, mas manipulable es la cosa). Citar el paper de Sergei y Brin, que dicen que hacen promedios de muchas cosas en la practica para evitar este problema. Usan muchos criterios promediados.

### 3.3. Ranking ATP

Empezemos con el apartado que seguro el lector más esperaba de todo el t.p.. Vilas fue o no 1ro entre 1975 y 1977? Esta pregunta podemos contestarla. Pero previo a esto, necesitamos poner al lector al tanto de la situación. Empezamos viendo el ranking oficial de la época.

En 1975, Vilas llegó a la primera final de un Grand Slam, el Roland Garros, en donde fue derrotado por el sueco Björn Borg por 2-6 3-6 4-6 y cuartos de final de Wimbledon.

Este es el top 10 para 1975 segun la ATP, donde Guillermo Vilas se hubica 2do:

173	<i>Jimmy</i>	<i>Connors</i>
127	<i>Guillermo</i>	<i>Vilas</i>
34	<i>Bjorn</i>	<i>Borg</i>
19	<i>Arthur</i>	<i>Ashe</i>
225	<i>Manuel</i>	<i>Orantes</i>
210	<i>Ken</i>	<i>Rosewall</i>
144	<i>Ilie</i>	<i>Nastase</i>
180	<i>John</i>	<i>Alexander</i>
318	<i>Roscoe</i>	<i>Tanner</i>
309	<i>Rod</i>	<i>Laver</i>

En 1976 la ATP hubica a Vilas 6to dentro del top 10:

173	<i>Jimmy</i>	<i>Connors</i>
34	<i>Bjorn</i>	<i>Borg</i>
144	<i>Ilie</i>	<i>Nastase</i>
225	<i>Manuel</i>	<i>Orantes</i>
288	<i>Raul</i>	<i>Ramirez</i>
127	<i>Guillermo</i>	<i>Vilas</i>
1	<i>Adriano</i>	<i>Panatta</i>
134	<i>Harold</i>	<i>Solomon</i>
87	<i>Eddie</i>	<i>Dibbs</i>
41	<i>Brian</i>	<i>Gottfried</i>

En 1977 Vilas conquistó 17 torneos (récord todavía vigente), se consagró en Roland Garros frente a Brian Gottfried por 6-0 6-3 6-0 y el US Open frente a Jimmy Connors por 2-6 6-3 7-6 6-0, fue finalista de Australia donde cayó frente a Roscoe Tanner por 3-6 3-6 3-6, logró una seguidilla de 50 partidos consecutivos sin conocer la derrota y, durante esos doce meses, ganó 145 de los 159 encuentros que jugó (91,1 %) ubicándose como 2do en el ranking de la ATP de finales de ese año. Según el ranking de la ATP, el mejor del año fue el norteamericano Jimmy Connors, batido por Vilas en la final del US Open (detalle a tener en cuenta), que en ese 1977 no obtuvo ningún torneo de Grand Slam y ganó ocho torneos, menos de la mitad de los logrados por Vilas.

Veamos el top 10 oficial para este año:

173	<i>Jimmy</i>	<i>Connors</i>
127	<i>Guillermo</i>	<i>Vilas</i>
34	<i>Bjorn</i>	<i>Borg</i>
371	<i>Vitas</i>	<i>Gerulaitis</i>
41	<i>Brian</i>	<i>Gottfried</i>
87	<i>Eddie</i>	<i>Dibbs</i>
225	<i>Manuel</i>	<i>Orantes</i>
288	<i>Raul</i>	<i>Ramirez</i>
144	<i>Ilie</i>	<i>Nastase</i>
81	<i>Dick</i>	<i>Stockton</i>

Todos los años tienen como líder indiscutido al estadounidense Jimmy Connors.

### 3.3.1. Ranking ATP oficial vs. Ranking PageRank vs. Sort por diferencia de victorias/derrotas

Introducimos aquí un algoritmo de ranqueo con un criterio de ordenamiento basado en victorias/derrotas. Además, si hay empate, define por diferencia de puntos, pero para el caso de tenis y como está definido el sistema de puntos del mismo no tiene ninguna relevancia. Solo nos servirá para poder hacer un análisis cualitativo de las virtudes de pagerank sobre algoritmos más básicos.

Veamos que obtuvimos en cada año con este algoritmo, empezando por 1975:

19	96	18
225	90	20
127	86	18
173	79	11
144	88	22
34	82	17
288	70	28
180	65	23
318	64	24
155	60	21

Podemos ver que Vilas no aparece segundo, si no tercero. Connors fue desplazado al 4to lugar y en primer lugar aparece Arthur Ashe quien ocupaba el lugar que Connors ocupa ahora. Podemos ver varios cambios relacionados al top 10 oficial. Todos ellos particularmente relacionados al hecho de que como se habia anticipado la cantidad de victorias sobre derrotas conformaria otro ranking diferente en el cual no importa exactamente que clase de victorias a conseguido o derrotas sufrido un determinado participante. Podemos ver que Vilas con más victorias sobre derrotas que Connors, se encontraba por debajo de él en el ranking oficial. Esto se debe al sistema de puntos manejados por el ATP donde se suma más puntaje cuanto mas avancemos en un torneo y cuantas más finales ganemos y por la jerarquía de ese torneo.

Observemos el ranking de 1976:

173	86	9
127	83	20
144	73	15
288	91	33
225	76	19
87	83	28
34	57	12
318	71	27
380	75	32
134	70	27

Vemos que Vilas avanzó del 6to lugar al 2do puesto!! que era ocupado por una bestia como Bjorn Borg y que fue desplazado al 7mo. Si analizamos un poco los partidos veremos que Vilas fué derrotado por Bjorn Borg 3 veces (Winbledon, Dallas WCT, Sao Paulo WCT) mientras que Vilas nunca pudo derrotarlo, Contra Connors se enfrentó una vez y fue derrotado (US Open) al igual que Bjorn Borg que se enfrentó en 3 ocasiones y cayó en la misma cantidad...esto nuevamente es un indicador de que determinadas victorias son más importantes que otras.

Veamos que sucede en 1977:

12712614411052234718173691587772837160152256525134642728861243806630

...humo humo humo. mostrar que vilas en 1975 y 1976 no le gano a nadie groso y por eso no esta primero, eso en consecuencia es que perdio en instancias mas o menos decisivas contra esos rivales grosos y por eso mantiene un lugar alto en el ranking..ademas todos ellos le ganaron a muchos que a su vez le ganaron muchos. Mostrar que esto no sucede para 1977 donde ampliamente vilas se los viola porque gano todo y ademas y lo mas importante le gano a jimmy neutron y por eso rankeo mucho mas puntaje. que a su vez ellos dos son los unicos con 0.30 y pico en el ranking y que los anteriores estan mas parejos porque todos estaban mas parejos (?)

### 3.3.2. Eleccion del factor de 'teletransportacion' $c$ (RODRI)

probar relevancia a medida que cambias ese valor  $= 0.85$ , creo que  $c$ .

Citar paper de google, que usan 0.85. Discutir que si  $c$  es uno, ignoras la estructura del grafo al hacer el ranking, todos rankean igual.

Pagina 6.... This is the ultimately egalitarian case: the only... blah. La idea es jugar con  $c$  aca, como dije arriba. Es un buen exp, hay que pensar bien como graficarlo y que quede lindo, creo que es facil.

## 3.4. Metodo de la Potencia

### 3.4.1. Representacion de la Matriz de Transicion (RODRI/FEDE)

Este experimento lo pueden hacer directo o usando al PageRank. Si pueden, implementen todas las representaciones de matrices y luego comparen el tiempo de computo del producto  $N$  veces. Comparen la matriz normal vs el resto. Discutan que en paginas web la cantidad de vertices del grafo se va al carajo, pero para deportes es super acotada, asi que la eleccion de estructura no afecta tanto.

Aca podes argumentar que lo que domina al metodo de la potencia es la cantidad de productos, asi que no hace falta probar PageRank directo. Igual si queres metelo con pagerank de una, a fin de cuentas es lo mismo.

### 3.4.2. Evolucion de la norma entre iteraciones

Como va evolucionando la norma manhattan entre dos iteraciones sucesivas. Eje x, iteraciones, eje y, norma manhattan.

### 3.4.3. Convergencia

Aca tienen que calcular el vector posta, y luego tomar algun tipo de norma. En el eje x van a tener la cantidad de iteraciones, y en el eje y van a tener la norma de  $x^* - x_{actual}$ .

### 3.4.4. Eleccion del $x_0$

Aca pongan que te conviene arrancar con una buena 'adivinanza' de la solucion, asi se acerca mas rapido. Muestren la cantidad de iteraciones a la convergencia (norma manhattan |epsilon) dependiendo de la distancia de la solucion inicial a la solucion posta. Si arranco con la posta de una, converge de una. Si arranco con una sol asquerosa inicial, tarda mas iteraciones en cumplir nuestro epsilon.

Mostrar dos instancias, una donde arranco desde el valor inicial donde todos tienen  $1/n$  y otra donde una tiene 1 y el resto 0, mostrar la cantidad de pasos y como evoluciona la norma.

## 4. Conclusiones

Una vez que ya este todo lo leo y escribo esto bien a los pedos, incluyendo la caratula.

## 5. Apéndice A: Enunciado

**Métodos Numéricos**  
Segundo Cuatrimestre 2015  
**Trabajo Práctico 2**



*Departamento de Computación*  
*Facultad de Ciencias Exactas y Naturales*  
*Universidad de Buenos Aires*

*Ohhh solo tiran  $\pi$ -edras...*

---

### Contexto y motivación

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 explicado en Brin y Page [3], 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.

Por otro lado, las competencias deportivas, en todas sus variantes y disciplinas, requieren casi inevitablemente la comparación entre competidores mediante la confección de *Tablas de Posiciones* y *Rankings* en base a resultados obtenidos en un período de tiempo determinado. Estos ordenamientos de equipos están generalmente (aunque no siempre) basados en reglas relativamente claras y simples, como proporción de victorias sobre partidos jugados o el clásico sistema de puntajes por partidos ganados, empatados y perdidos. Sin embargo, estos métodos simples y conocidos por todos muchas veces no logran capturar la complejidad de la competencia y la comparación. Esto es particularmente evidente en ligas donde, por ejemplo, todos los equipos no juegan la misma cantidad de veces entre sí.

A modo de ejemplo, la NBA y NFL representan dos ligas con fixtures de temporadas regulares con estas características. Recientemente, el Torneo de Primera División de AFA se suma a este tipo de competencias, ya que la incorporación de la *Fecha de Clásicos* parece ser una interesante idea comercial, pero no tanto desde el punto de vista deportivo ya que cada equipo juega contra su *clásico* más veces que el resto. Como contraparte, éstos rankings son utilizados muchas veces como criterio de decisión, como por ejemplo para determinar la participación en alguna competencia de nivel internacional, con lo cual la confección de los mismos constituye un elemento sensible, afectando intereses deportivos y económicos de gran relevancia.

### El problema, Parte I: PageRank y páginas web

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 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 [4] y Kamvar et al. [7, 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 página 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

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

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  $j \in \text{Web}$ . Denotaremos con  $\pi$  al vector solución de la ecuación  $P_2x = x$ , que es comúnmente denominado *estado estacionario*.

En particular,  $P_2$  corresponde a una matriz *estocástica por columnas* que cumple las hipótesis planteadas en Bryan y Leise [4] y Kamvar et al. [7], 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,  $\pi$  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 puntaje.

## El problema, Parte II: PageRank y ligas deportivas

Existen en la literatura distintos enfoques para abordar el problema de determinar el *ranking* de equipos de una competencia en base a los resultados de un conjunto de partidos. En Govan et al. [6] se hace una breve reseña de dos ellos, y los autores proponen un nuevo método basado en el algoritmo PageRank que denominan GeM<sup>1</sup>. Conceptualmente, el método GeM representa la temporada como un red (grafo) donde las páginas web representan a los equipos, y existe un link (que tiene un valor, llamado peso, asociado) entre dos equipos que los relaciona modelando los resultados de los posibles enfrentamientos entre ellos. En base a este modelo, Govan et al. [6] proponen calcular el ranking de la misma forma que en el caso de las páginas web.

En su versión básica, que es la que consideraremos en el presente trabajo, el método GeM (ver, e.g., [6, Sección GeM Ranking Method]) es el siguiente<sup>2</sup>:

1. La temporada se representa mediante un grafo donde cada equipo representa un nodo y existe un link de  $i$  a  $j$  si el equipo  $i$  perdió al menos una vez con el equipo  $j$ .
2. Se define la matriz  $A^t \in \mathbb{R}^{n \times n}$

$$A_{ji}^t = \begin{cases} w_{ji} & \text{si el equipo } i \text{ perdió con el equipo } j, \\ 0 & \text{en caso contrario,} \end{cases}$$

donde  $w_{ji}$  es la diferencia absoluta en el marcador. En caso de que  $i$  pierda más de una vez con  $j$ ,  $w_{ji}$  representa la suma acumulada de diferencias. Notar que  $A^t$  es una generalización de la matriz de conectividad  $W$  definida en la sección anterior.

3. Definir la matriz  $H_{ji}^t \in \mathbb{R}^{n \times n}$  como

$$H_{ji}^t = \begin{cases} A_{ji}^t / \sum_{k=1}^n A_{ki}^t & \text{si hay un link } i \text{ a } j, \\ 0 & \text{en caso contrario.} \end{cases}$$

<sup>1</sup>Aunque no se especifica, asumimos que el nombre se debe a las iniciales de los autores.

<sup>2</sup>Notar que en artículo, Govan et al. [6] lo definen sobre la traspuesta. La definición y las cuentas son equivalentes, simplemente se modifica para mantener la consistencia a lo largo del enunciado.

4. Tomar  $P = H^t$ , y aplicar el método PageRank como fue definido previamente, siendo  $\pi$  la solución a la ecuación  $P_2x = x$ . Notar que los páginas sin links salientes, en este contexto se corresponden con aquellos equipos que se encuentran invictos.
5. Utilizar los puntajes obtenidos en  $\pi$  para ordenar los equipos.

En función del contexto planteado previamente, el método GeM define una estructura que relaciona equipos dependiendo de los resultados parciales y obtener un ranking utilizando solamente esta información.

### Enunciado

El objetivo del trabajo es experimentar en el contexto planteado utilizando el algoritmo PageRank con las variantes propuestas. A su vez, se busca comparar los resultados obtenidos cualitativa y cuantitativamente con los algoritmos tradicionales utilizados en cada uno de los contextos planteados. Los métodos a implementar (como mínimo) en ambos contextos planteados por el trabajo son los siguientes:

1. *Búsqueda de páginas web*: PageRank 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.
2. *Rankings en competencias deportivas*: GeM y al menos un método estándar propuesto por el grupo (ordenar por victorias/derrotas, puntaje por ganado/empatado/perdido, etc.) en función del deporte(s) considerado(s).

El contexto considerado en 1., en la búsqueda de páginas web, representa un desafío no sólo desde el modelado, si no también desde el punto de vista computacional considerando la dimensión de la información y los datos a procesar. Luego, dentro de nuestras posibilidades, consideramos un entorno que simule el contexto real de aplicación donde se abordan instancias de gran escala (es decir,  $n$ , el número total de páginas, es grande). Para el desarrollo de PageRank, se pide entonces considerar el trabajo de Bryan y Leise [4] donde se explica la intuición y algunos detalles técnicos respecto a PageRank. Además, en Kamvar et al. [7] 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. [7, Algoritmo 1] se propone una forma alternativa para computar  $x^{(k+1)} = P_2x^{(k)}$ . Este resultado debe ser utilizado para mejorar el almacenamiento de los datos.

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 considere el contexto de aplicación. Además, para PageRank se debe implementar el método de la potencia para calcular el autovector principal. Esta implementación debe ser realizada íntegramente en C++.

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  $c$ .
2. Estudiar el tiempo de cómputo requerido por PageRank.
3. Para cada algoritmo, proponer ejemplos de tamaño pequeño que ilustren el comportamiento esperado (puede ser utilizando las herramientas provistas por la cátedra o bien generadas por el grupo).

Puntos opcionales:

1. Demostrar que los pasos del Algoritmo 1 propuesto en Kamvar et al. [7] son correctos y computan  $P_2x$ .
2. Establecer una relación con la proporción entre  $\lambda_1 = 1$  y  $|\lambda_2|$  para la convergencia de PageRank.

El segundo contexto de aplicación no presenta mayores desafíos desde la perspectiva computacional, ya que en el peor de los casos una liga no suele tener mas que unas pocas decenas de equipos. Más aún, es de esperar que en general la matriz que se obtiene no sea esparsa, ya que probablemente un equipo juegue contra un número significativo de contrincantes. Sin embargo, la popularidad y sensibilidad del problema planteado requieren de un estudio detallado y pormenorizado de la calidad de los resultados obtenidos. El objetivo en este segundo caso de estudio es puramente experimental.

En función de la implementación, aún cuando no represente la mejor opción, es posible reutilizar y adaptar el desarrollo realizado para páginas web. También es posible realizar una nueva implementación desde cero, simplificando la operatoria y las estructuras, en C++, MATLAB o PYTHON.

La experimentación debe ser realizada con cuidado, analizando (y, eventualmente, modificando) el modelo de GeM:

1. Considerar al menos un conjunto de datos reales, con los resultados de cada fecha para alguna liga de algún deporte.
2. Notar que el método GeM asume que no se producen empates entre los equipos (o que si se producen, son poco frecuentes). En caso de considerar un deporte donde el empate se da con cierta frecuencia no despreciable (por ejemplo, fútbol), es fundamental aclarar como se refleja esto en el modelo y analizar su eventual impacto.
3. Realizar experimentos variando el parámetro  $c$ , indicando como impacta en los resultados. Analizar la evolución del ranking de los equipos a través del tiempo, evaluando también la evolución de los rankings e identificar características/hechos particulares que puedan ser determinantes para el modelo, si es que existe alguno.
4. Comparar los resultados obtenidos con los reales de la liga utilizando el sistema estándar para la misma.

Puntos opcionales:

1. Proponer (al menos) dos formas alternativas de modelar el empate entre equipos en GeM.

### 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 el método alternativo), la probabilidad de teletransportación  $c$ , el tipo de instancia (0 páginas web, 1 deportes), 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 del método de la potencia.

El siguiente ejemplo muestra un caso donde se pide ejecutar PageRank, con una probabilidad de teletransportación de 0.85, sobre la red descrita en **test1.txt** (que se encuentra en el directorio **tests/**), correspondiente a una instancia de ranking aplicado a deportes y con una tolerancia de corte de 0,0001.

```
0 0.85 1 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 [2] (el tipo de instancia es 0), con redes de tamaño grande obtenidos a partir de datos reales. Además, se consideran las instancias que se forman a partir de resultados de partidos entre equipos, para algún deporte elegido por el grupo.

En el caso de la base de SNAP, los archivos contiene 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 [4, 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 el caso de rankings en ligas deportivas, el archivo contiene primero una línea con información sobre la cantidad de equipos ( $n$ ), y la cantidad de partidos totales a considerar ( $k$ ). Luego, siguen  $k$  líneas donde cada una de ellas representa un partido y contiene la siguiente información: número de fecha (es un dato opcional al problema, pero que puede ayudar a la hora de experimentar), equipo  $i$ , goles equipo  $i$ , equipo  $j$ , goles equipo  $j$ . A continuación se muestra el archivo de entrada con la información del ejemplo utilizado en Govan et al. [6]:

```
6 10
1 1 16 4 13
1 2 38 5 17
1 2 28 6 23
1 3 34 1 21
1 3 23 4 10
1 4 31 1 6
1 5 33 6 25
1 5 38 4 23
1 6 27 2 6
1 6 20 5 12
```

Es importante destacar que, en este último caso, los equipos son identificados mediante un número. Opcionalmente podrá considerarse un archivo que contenga, para cada equipo, cuál es el código con el que se lo identifica.

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 PageRank/IN-DEG/método alternativo.

Para generar instancias de páginas web, 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.

Para instancias correspondientes a resultados entre equipos, la cátedra provee un conjunto de archivos con los resultados del Torneo de Primera División del Fútbol Argentino hasta la Fecha 23. Es importante aclarar que los dos partidos suspendidos, River - Defensa y Justicia y Racing - Godoy Cruz han sido arbitrariamente completados con un resultado inventado, para simplificar la instancia. En función de datos reales, una alternativa es considerar el repositorio DataHub [1], que contiene información estadística y resultados para distintas ligas y deportes de todo el mundo.

---

### Fechas de entrega

- *Formato Electrónico*: Martes 6 de Octubre de 2015, 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 7 de Octubre de 2015, a las 18 hs. en la clase práctica.

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

## 6. Apéndice B: Código

### 6.1. system.cpp

---

```
1  #include <iostream>
2  #include <math.h>
3  #include <fstream>
4  #include <sstream>
5  #include <stdio.h>
6  #include <string.h>
7  #include <time.h>
8  #include <new>
9  #include <regex>
10 #include <iterator>
11 #include "sparseMatrix.h"
12
13 using namespace std;
14
15 struct dataNode {
16     //dataNode(int n) : node(n), edgesCount(0) {} // no lo pude compilar con listas
17     //de inicializacion...
18     int node;
19     int edgesCount;
20 };
21
22 struct matchesStats {
23     //matchesStats(int t) : team(t), matchesWin(0), matchesDefeat(0), pointsScored(0)
24     // , pointsReceived(0) {}
25     int team;
26     int matchesWin;
27     int matchesDefeat;
28     int pointsScored;
29     int pointsReceived;
30 };
31
32 //webs / sports
33 Matrix<double> pageRank(Matrix<double>& M, double c, double d, vector<int>&
34     nodesCount);
35
36 Matrix<double> enhancementPageRank(Matrix<double>& M, double c, double d, vector<int>
37     & nodesCount);
38
39 //webs
40 void in_deg(vector<dataNode>& nodesCount);
41
42 //sports
43 void basic_sort(vector<matchesStats>& stats);
44
45 //out data
46 void saveResultPageRank(FILE * pFile, Matrix<double>& data);
47 void saveResultInDeg(FILE * pFile, vector<dataNode>& data);
48 void saveResultBasicSort(FILE * pFile, vector<matchesStats>& data);
49
50 //utils
51 double uniform_rand(double a, double b);
52
53 int main(int argc, char** argv) {
```

```

52     if (argc < 3) {
53         printf("Usage %s: parametros.in salida.out \n", argv[0]);
54         return 0;
55     }
56
57     ifstream inputFile(argv[1]);
58
59     if (!inputFile.good()) {
60         printf("can't open input file.\n");
61         return 0;
62     }
63
64     FILE * outputFile = NULL;
65
66     outputFile = fopen(argv[2], "w");
67     if (outputFile == NULL) {
68         printf("can't open output file.\n");
69         return 0;
70     }
71
72     int alg = 0;
73     double c = 0;
74     int inst = 0;
75     double e = 0;
76     char testFileName[100];
77
78     string line;
79     getline(inputFile, line);
80
81     sscanf(line.c_str(), "%d %d %d %s %d", &alg, &c, &inst, testFileName, &e);
82
83     ifstream testFile(testFileName);
84     getline(testFile, line);
85
86     if (inst == 0) {
87         int nodes = 0;
88         int edges = 0;
89
90         sscanf(line.c_str(), "%d %d", &nodes, &edges);
91         cout << "nodes: " << nodes << " edges: " << edges << endl;
92
93         if (alg == 0) {
94             Matrix<double> M(nodes, nodes);
95
96             vector<int> nodesCount(nodes);
97
98             for (int i = 0; i < nodes; i++) {
99                 nodesCount[i] = 0;
100             }
101
102             int i = 0;
103             while (i < edges) {
104                 int node_from = 0;
105                 int node_to = 0;
106
107                 getline(testFile, line);
108                 sscanf(line.c_str(), "%d %d", &node_from, &node_to);
109                 //cout << "node_from: " << node_from << " node_to: " << node_to <<
                     endl;

```

```

110         nodesCount[node_from-1] += 1;
111         M(node_to-1, node_from-1) = 1;
112         i++;
113     }
114
115     Matrix<double> res = pageRank(M, c, e, nodesCount);
116     // Matrix<double> res = enhancementPageRank(M, c, e, nodesCount);
117
118     cout << "page rank result: \n" << endl;
119     res.printMatrix();
120
121     saveResultPageRank(outputFile, res);
122 } else {
123     // group algorithm webs
124     vector<dataNode> nodesCount(nodes);
125
126     for(int i = 0; i < nodes; i++){
127         dataNode nod;
128         nod.node = i+1;
129         nod.edgesCount = 0;
130         nodesCount[i] = nod;
131     }
132
133     int i = 0;
134     while(i < edges){
135         int node_from = 0;
136         int node_to = 0;
137
138         getline(testFile, line);
139         sscanf(line.c_str(), "%d %d", &node_from, &node_to);
140         //cout << "node_from: " << node_from << " node_to: " << node_to <<
141             endl;
142
143         dataNode nod = nodesCount[node_from-1];
144
145         nod.edgesCount += 1;
146
147         nodesCount[node_from-1] = nod;
148
149         i++;
150     }
151
152     in_deg(nodesCount);
153
154     saveResultInDeg(outputFile, nodesCount);
155 }
156 } else {
157     int teams = 0;
158     int matches = 0;
159
160     sscanf(line.c_str(), "%d %d", &teams, &matches);
161     cout << "teams: " << teams << " matches: " << matches << endl;
162
163     int day = 0;
164     int local = 0;
165     int visitor = 0;
166     int local_score = 0;
167     int visitor_score = 0;

```

```

168
169     if (alg == 0){
170         // page rank sports
171         Matrix<double> M(teams, teams);
172
173         vector<int> totalAbs(teams);
174
175         for(int i = 0; i < teams; i++) {
176             totalAbs[i] = 0;
177         }
178
179         int i = 0;
180         while(i < matches){
181             getline(testFile, line);
182             sscanf(line.c_str(), "%d %d %d %d %d", &day, &local, &local_score, &
                visitor, &visitor_score);
183
184             int abs_score = abs(local_score - visitor_score);
185             if (local_score > visitor_score) {
186                 M(local-1, visitor-1) += abs_score;
187                 totalAbs[visitor-1] += abs_score;
188             } else if (visitor_score > local_score) {
189                 M(visitor-1, local-1) += abs_score;
190                 totalAbs[local-1] += abs_score;
191             }
192
193             i++;
194         }
195
196         Matrix<double> res = pageRank(M, c, e, totalAbs);
197
198         cout << "gem result: \n" << endl;
199         res.printMatrix();
200
201         saveResultPageRank(outputFile, res);
202     } else {
203         // group algorithm sports
204         vector<matchesStats> stats(teams);
205
206         for(int i = 0; i < teams; i++){
207             matchesStats teamStats;
208             teamStats.team = i+1;
209             teamStats.matchesWin = 0;
210             teamStats.matchesDefeat = 0;
211             teamStats.pointsScored = 0;
212             teamStats.pointsReceived = 0;
213             stats[i] = teamStats;
214         }
215
216         int i = 0;
217         while(i < matches){
218             getline(testFile, line);
219             sscanf(line.c_str(), "%d %d %d %d %d", &day, &local, &local_score, &
                visitor, &visitor_score);
220
221             matchesStats localStats = stats[local-1];
222             matchesStats visitorStats = stats[visitor-1];
223
224             if (local_score > visitor_score) {

```



```

225         localStats.matchesWin += 1;
226         visitorStats.matchesDefeat += 1;
227     }else if(visitor_score > local_score) {
228         localStats.matchesDefeat += 1;
229         visitorStats.matchesWin += 1;
230     }
231
232     localStats.pointsScored += local_score;
233     localStats.pointsReceived += visitor_score;
234     visitorStats.pointsScored += visitor_score;
235     visitorStats.pointsReceived += local_score;
236
237     stats[local-1] = localStats;
238     stats[visitor-1] = visitorStats;
239
240     i++;
241 }
242
243     basic_sort(stats);
244
245     saveResultBasicSort(outputFile, stats);
246 }
247 }
248
249     testFile.close();
250
251     inputFile.close();
252
253     if (outputFile != NULL) fclose(outputFile);
254
255     //M.printMatrix();
256
257     //Depending on input data, create a matrix with the input file and call rank with
        matrix a values
258
259     return 0;
260 }
261
262 Matrix<double> pageRank(Matrix<double>& M, double c, double d, vector<int>&
    nodesCount) {
263     srand(45);
264
265     int n = M.rows();
266     double dbl_n = M.rows();
267
268     int j = 0;
269     while(j < n){
270         int i = 0;
271         while(i < n){
272             if(M(i, j) != 0){
273                 M(i, j) = M(i, j) / (double)nodesCount[j];
274             }else if(nodesCount[j] == 0){
275                 M(i, j) = 1/dbl_n; // dangling node / undefeated team
276             }
277             i++;
278         }
279         j++;
280     }
281

```

```

282 Matrix<double> v(n, 1/dbl_n);
283
284 Matrix<double> E(n, n, (1 - c)*1/dbl_n); // PRE: rows == columns
285
286 //Salvo que sea c = 1 no tiene sentido usar Sparse Matrix
287 Matrix<double> A = M*c + E;
288
289 Matrix<double> x(n, 1, 1/dbl_n);
290
291 //for (int i = 0; i < M.rows(); i++) {
292 //    x(i) = uniform_rand(0, 1);
293 //}
294
295 Matrix<double> last_x(n);
296
297 double delta = 0;
298
299 do {
300     last_x = x;
301     x = A*x;
302     delta = x.L1(last_x);
303 } while (delta > d);
304
305 printf("delta is %f\r\n", delta);
306
307 return x;
308 }
309
310 Matrix<double> enhancementPageRank(Matrix<double>& M, double c, double d, vector<int
>& nodesCount) {
311     srand(45);
312
313     int n = M.rows();
314     double dbl_n = M.rows();
315
316     int j = 0;
317     while(j < n){
318         int i = 0;
319         while(i < n){
320             if(M(i, j) != 0){
321                 M(i, j) = M(i, j) / (double)nodesCount[j];
322             }
323             i++;
324         }
325         j++;
326     }
327
328     SparseMatrix<double> A(M);
329
330     SparseMatrix<double> x(n, 1/dbl_n);
331
332     //for (int i = 0; i < M.rows(); i++) {
333     //    x(i) = uniform_rand(0, 1);
334     //}
335
336     SparseMatrix<double> last_x(n);
337
338     SparseMatrix<double> v(n, 1/dbl_n);
339

```

```

340     double delta = 0;
341
342     do {
343         last_x = x;
344
345         x = A*x;
346         x = x*c;
347
348         double w = last_x.norm1() - x.norm1();
349
350         x = x + v*w;
351
352         delta = x.L1(last_x);
353     } while (delta > d);
354
355     printf("delta is %f\r\n", delta); //Deberia devolverse.
356
357     return x.descompress();
358 }
359
360 void in_deg(vector<dataNode>& nodesCount) {
361     sort(nodesCount.begin(), nodesCount.end(), [](dataNode a, dataNode b) {
362         return b.edgesCount < a.edgesCount;
363     });
364
365     cout << "IN-DEG result \n" << endl;
366     for (dataNode a : nodesCount) {
367         cout << "node: " << a.node << " points: " << a.edgesCount << "\n";
368     }
369 }
370
371 void basic_sort(vector<matchesStats>& stats) {
372     sort(stats.begin(), stats.end(), [](matchesStats a, matchesStats b) {
373         if (a.matchesWin - a.matchesDefeat != b.matchesWin - b.matchesDefeat) {
374             return b.matchesWin - b.matchesDefeat < a.matchesWin - a.matchesDefeat;
375         } else {
376             return b.pointsScored - b.pointsReceived < a.pointsScored - a.
                 pointsReceived;
377         }
378     });
379
380     cout << "basic sort result \n" << endl;
381     for (matchesStats a : stats) {
382         cout << "team: " << a.team
383             << " matches win: " << a.matchesWin << " matches defeat: " << a.matchesDefeat
384             << " points scored: " << a.pointsScored << " points received: " << a.
                 pointsReceived << "\n";
385     }
386 }
387
388 void saveResultPageRank(FILE * pFile, Matrix<double>& data) {
389     int n = data.rows();
390     int i = 0;
391
392     while(i < n){
393         fprintf(pFile, "%f\r\n", data(i));
394         i++;
395     }
396 }

```

```

397
398 void saveResultInDeg(FILE * pFile , vector<dataNode>& data) {
399     for (dataNode a : data){
400         fprintf(pFile , "%d %d\r\n" , a.node , a.edgesCount);
401     }
402 }
403
404 void saveResultBasicSort(FILE * pFile , vector<matchesStats>& data) {
405     for (matchesStats a : data){
406         fprintf(pFile , "%d %d %d %d %d\r\n" , a.team , a.matchesWin , a.matchesDefeat , a
            .pointsScored , a.pointsReceived);
407     }
408 }
409
410 double uniform_rand(double a, double b) {
411     return ((b-a)*((double)rand()/RAND_MAX))+a;
412 }

```

---

## 6.2. matrix.h

---

```
1  /*
2  * File:    matrix.h
3  * Author:  Federico
4  *
5  * Created on August 16, 2015, 9:54 PM
6  */
7
8  #ifndef MATRIX_H
9  #define MATRIX_H
10
11  #include <algorithm>
12  #include <math.h>
13  #include <vector>
14  #include <stdio.h>
15
16  using namespace std;
17
18  // La matriz respeta la notacion de la catedra, es decir, el primer subindice
19  // es la fila y el segundo es la columna
20
21  template<class T>
22  class Matrix {
23  public:
24      Matrix();
25      Matrix(int rows); // Columnas impllicitas (col = 1)
26      Matrix(int rows, int col);
27      Matrix(int rows, int col, const T& init);
28      Matrix(const Matrix<T>& other);
29      ~Matrix();
30
31      Matrix<T>& operator=(const Matrix<T>& other);
32      Matrix<T> operator*(const Matrix<T>& other);
33      Matrix<T>& operator*=(const Matrix<T>& other);
34      Matrix<T> operator+(const Matrix<T>& other);
35      Matrix<T>& operator+=(const Matrix<T>& other);
36      Matrix<T> operator-(const Matrix<T>& other);
37      Matrix<T>& operator-=(const Matrix<T>& other);
38
39      Matrix<T> operator*(const T& scalar);
40      Matrix<T> operator/(const T& scalar);
41
42      T& operator()(int a, int b);
43      const T& operator()(const int a, const int b) const;
44      T& operator()(int a);
45      const T& operator()(const int a) const;
46
47      //PRE are vectors
48      double L1(const Matrix<T>& other);
49      double norm1();
50
51      int rows();
52      int columns();
53
54      int rows() const;
55      int columns() const;
56
57      void printMatrix();
```

```

58
59     private:
60         vector<vector<T> > _values;
61         int _rows;
62         int _columns;
63
64 };
65
66 template<class T>
67 Matrix<T>::Matrix()
68     : _values(1), _rows(1), _columns(1)
69 {
70     _values[0].resize(1);
71 }
72
73 template<class T>
74 Matrix<T>::Matrix(int rows)
75     : _values(rows), _rows(rows), _columns(1)
76 {
77     for(int i = 0; i < rows; i++) {
78         _values[i].resize(1);
79     }
80 }
81
82 template<class T>
83 Matrix<T>::Matrix(int rows, int col)
84     : _values(rows), _rows(rows), _columns(col)
85 {
86     for(int i = 0; i < rows; i++) {
87         _values[i].resize(col);
88     }
89 }
90
91 template<class T>
92 Matrix<T>::Matrix(int rows, int col, const T& init)
93     : _values(rows), _rows(rows), _columns(col)
94 {
95     for(int i = 0; i < rows; i++) {
96         _values[i].resize(col, init);
97     }
98 }
99
100 template<class T>
101 Matrix<T>::Matrix(const Matrix<T>& other)
102     : _values(other._values), _rows(other._rows), _columns(other._columns)
103 {}
104
105 template<class T>
106 Matrix<T>::~~Matrix() {}
107
108 template<class T>
109 Matrix<T>& Matrix<T>::operator=(const Matrix<T>& other) {
110     if (&other == this)
111         return *this;
112
113     int new_rows = other._rows;
114     int new_columns = other._columns;
115
116     _rows = new_rows;

```

```

117     _columns = new_columns;
118
119     _values.resize(new_rows);
120     for (int i = 0; i < new_columns; i++) {
121         _values[i].resize(new_columns);
122     }
123
124     for(int i = 0; i < new_rows; i++) {
125         for(int j = 0; j < new_columns; j++) {
126             _values[i][j] = other(i, j);
127         }
128     }
129
130     return *this;
131 }
132
133 template<class T>
134 Matrix<T> Matrix<T>::operator*(const Matrix<T>& other) {
135     // ASUME QUE LAS DIMENSIONES DAN
136     Matrix<T> result(_rows, other._columns);
137
138     int innerDim = _columns; // Tambien podria ser other._rows
139
140     for(int i = 0; i < result._rows; i++) {
141         for(int j = 0; j < result._columns; j++) {
142             result(i,j) = 0;
143             for(int k = 0; k < innerDim; k++) {
144                 result(i,j) += _values[i][k] * other(k,j);
145             }
146         }
147     }
148
149     return result;
150 }
151
152 template<class T>
153 Matrix<T>& Matrix<T>::operator*=(const Matrix<T>& other) {
154     Matrix<T> result = (*this) * other;
155     (*this) = result;
156     return (*this);
157 }
158
159 template<class T>
160 Matrix<T> Matrix<T>::operator+(const Matrix<T>& other) {
161     // ASUME QUE LAS DIMENSIONES DAN
162     Matrix<T> result(_rows, other._columns);
163
164     for(int i = 0; i < result._rows; i++) {
165         for(int j = 0; j < result._columns; j++) {
166             result(i,j) = _values[i][j] + other(i,j);
167         }
168     }
169
170     return result;
171 }
172
173 template<class T>
174 Matrix<T>& Matrix<T>::operator+=(const Matrix<T>& other) {
175     Matrix<T> result = (*this) + other;

```

```

176     (*this) = result;
177     return (*this);
178 }
179
180 template<class T>
181 Matrix<T> Matrix<T>::operator-(const Matrix<T>& other) {
182     // ASUME QUE LAS DIMENSIONES DAN
183     Matrix<T> result(_rows, other._columns);
184
185     for(int i = 0; i < result._rows; i++) {
186         for(int j = 0; j < result._columns; j++) {
187             result(i,j) = _values[i][j] - other(i,j);
188         }
189     }
190
191     return result;
192 }
193
194 template<class T>
195 Matrix<T>& Matrix<T>::operator-=(const Matrix<T>& other) {
196     Matrix<T> result = (*this) - other;
197     (*this) = result;
198     return (*this);
199 }
200
201 template<class T>
202 Matrix<T> Matrix<T>::operator*(const T& scalar) {
203     Matrix<T> result(_rows, _columns);
204
205     for(int i = 0; i < result._rows; i++) {
206         for(int j = 0; j < result._columns; j++) {
207             result(i,j) = _values[i][j] * scalar;
208         }
209     }
210
211     return result;
212 }
213
214 template<class T>
215 Matrix<T> Matrix<T>::operator/(const T& scalar) {
216     Matrix<T> result(_rows, _columns);
217
218     for(int i = 0; i < result._rows; i++) {
219         for(int j = 0; j < result._columns; j++) {
220             result(i,j) = _values[i][j] / scalar;
221         }
222     }
223
224     return result;
225 }
226
227 template<class T>
228 T& Matrix<T>::operator()(int a, int b) {
229     return _values[a][b];
230 }
231
232 template<class T>
233 const T& Matrix<T>::operator()(const int a, const int b) const {
234     return _values[a][b];

```



```

235 }
236
237 template<class T>
238 T& Matrix<T>::operator()(int a) {
239     return _values[a][0];
240 }
241
242 template<class T>
243 const T& Matrix<T>::operator()(const int a) const {
244     return _values[a][0];
245 }
246
247 template<class T>
248 double Matrix<T>::L1(const Matrix<T>& other) {
249     Matrix<T> vectorSubs = *this - other;
250
251     double res = 0;
252     for(int i = 0; i < rows(); i++){
253         res = res + abs(vectorSubs(i));
254     }
255
256     return res;
257 }
258
259 template<class T>
260 double Matrix<T>::norm1() {
261     double res = 0;
262     for(int i = 0; i < rows(); i++){
263         res = res + abs(_values[i][0]);
264     }
265
266     return res;
267 }
268
269
270 template<class T>
271 int Matrix<T>::rows() {
272     return _rows;
273 }
274
275 template<class T>
276 int Matrix<T>::columns() {
277     return _columns;
278 }
279
280 template<class T>
281 int Matrix<T>::rows() const{
282     return _rows;
283 }
284
285 template<class T>
286 int Matrix<T>::columns() const{
287     return _columns;
288 }
289
290 template<class T>
291 void Matrix<T>::printMatrix() {
292     for(int i = 0; i < _rows; i++) {;
293         for(int j = 0; j < _columns; j++) {

```

```
294         printf("%f", _values[i][j]);
295         // cout << _values[i][j] << " ";
296     }
297     cout << endl;
298 }
299 cout << endl;
300 }
301
302 #endif /* MATRIX.H */
```

---

### 6.3. sparseMatrix.h

---

```
1  /*
2  * File:    SparseMatrix.h
3  * Author:  Rodrigo Kapobel
4  *
5  * Created on August 21, 2015, 23:05 PM
6  */
7
8  #ifndef SparseMatrix_H
9  #define SparseMatrix_H
10
11  #include <algorithm>
12  #include <math.h>
13  #include <vector>
14  #include <stdio.h>
15  #include <cmath>
16  #include "matrix.h"
17
18  using namespace std;
19
20  //CSR implementation for sparse matrix. For more information check the next link:
21  // https://en.wikipedia.org/wiki/Sparse\_matrix#Compressed\_row\_Storage\_.28CRS\_or\_CSR.29
22
23  template<class T>
24  class SparseMatrix {
25  public:
26      SparseMatrix();
27      SparseMatrix(int rows); // column vector with 0's
28      SparseMatrix(int rows, T value);
29      SparseMatrix(vector<T>& values, vector<int>& iValues, vector<int>& jValues, int
30          columns); // PRE: jValues::size == values::size & values of iValues[0..iValues
31          ::size-2] are indices of values & iValue[iValues::size-1] == values::size
32      SparseMatrix(const SparseMatrix<T>& other); // compress matrix
33      SparseMatrix(const Matrix<T>& other);
34      ~SparseMatrix();
35
36      SparseMatrix<T>& operator=(const SparseMatrix<T>& other);
37      SparseMatrix<T> operator*(const SparseMatrix<T>& other); // performs a matrix-
38          vector multiplication
39      SparseMatrix<T>& operator*=(const SparseMatrix<T>& other); // performs a matrix-
40          vector multiplication
41      SparseMatrix<T> operator+(const SparseMatrix<T>& other); // performs a vector-
42          vector sum
43      SparseMatrix<T>& operator+=(const SparseMatrix<T>& other); // performs a vector-
44          vector sum
45      SparseMatrix<T> operator-(const SparseMatrix<T>& other); // performs a vector-
46          vector subs
47      SparseMatrix<T>& operator-=(const SparseMatrix<T>& other); // performs a vector-
48          vector subs
49
50      SparseMatrix<T> operator*(const T& scalar);
51      SparseMatrix<T> operator/(const T& scalar);
52
53      T& operator()(int a, int b);
54      const T& operator()(const int a, const int b) const;
55      T& operator()(int a);
56      const T& operator()(const int a) const;
```

```

48
49     Matrix<T> descompress();
50     //PRE are vectors
51     double L1(const SparseMatrix<T>& other);
52     double norm1();
53
54     int rows();
55     int columns();
56     int rows() const;
57     int columns() const;
58     void printSparseMatrix();
59
60 private:
61     vector<T> _values;
62     vector<int> _iValues;
63     vector<int> _jValues;
64
65     int _columns;
66 };
67
68 template<class T>
69 SparseMatrix<T>::SparseMatrix()
70 : _values(1), _iValues(1), _jValues(1), _columns(1)
71 {}
72
73 template<class T>
74 SparseMatrix<T>::SparseMatrix(int rows)
75 : _values(rows), _iValues(1), _jValues(rows), _columns(1)
76 {
77     for(int i = 0; i <= rows; i++){
78         _iValues.resize(i+1, i);
79     }
80 }
81
82 template<class T>
83 SparseMatrix<T>::SparseMatrix(int rows, T value)
84 : _values(rows), _iValues(1), _jValues(rows), _columns(1)
85 {
86     for(int i = 0; i <= rows; i++){
87         _iValues.resize(i+1, i);
88         if(i < rows){
89             _values[i] = value;
90         }
91     }
92 }
93
94 template<class T>
95 SparseMatrix<T>::SparseMatrix(vector<T>& values, vector<int>& iValues, vector<int>&
96     jValues, int columns)
97 : _values(values), _iValues(iValues), _jValues(jValues), _columns(columns)
98 {}
99
100 template<class T>
101 SparseMatrix<T>::SparseMatrix(const SparseMatrix<T>& other)
102 : _values(other._values), _iValues(other._iValues), _jValues(other._jValues),
103     _columns(other._columns)
104 {}
105
106 template<class T>

```

```

105 SparseMatrix<T>::SparseMatrix(const Matrix<T>& other)
106 : _values(0), _iValues(0), _jValues(0), _columns(other.columns())
107 {
108     int new_rows = other.rows();
109     int new_columns = other.columns();
110
111     for(int i = 0; i < new_rows; i++) {
112         _iValues.resize(_iValues.size()+1, _values.size());
113         for(int j = 0; j < new_columns; j++) {
114             if (other(i, j) != 0) {
115                 _values.resize(_values.size()+1, other(i, j));
116                 _jValues.resize(_jValues.size()+1, j);
117             }
118         }
119     }
120     _iValues.resize(_iValues.size()+1, _values.size());
121 }
122
123 template<class T>
124 SparseMatrix<T>::~SparseMatrix() {}
125
126 template<class T>
127 SparseMatrix<T>& SparseMatrix<T>::operator=(const SparseMatrix<T>& other) {
128     if (&other == this)
129         return *this;
130
131     _values = other._values;
132     _iValues = other._iValues;
133     _jValues = other._jValues;
134     _columns = other.columns();
135
136     return *this;
137 }
138
139 template<class T>
140 SparseMatrix<T> SparseMatrix<T>::operator*(const SparseMatrix<T>& other) {
141     // ASUME QUE LAS DIMENSIONES DAN
142     SparseMatrix<T> result(rows());
143
144     for (int i = 0; i < rows(); i++) {
145         for (int j = _iValues[i]; j < _iValues[i+1]; j++) {
146             result(i) += _values[j]*other(_jValues[j]);
147         }
148     }
149
150     return result;
151 }
152
153 template<class T>
154 SparseMatrix<T>& SparseMatrix<T>::operator*=(const SparseMatrix<T>& other) {
155     SparseMatrix<T> result = (*this) * other;
156     (*this) = result;
157     return (*this);
158 }
159
160 template<class T>
161 SparseMatrix<T> SparseMatrix<T>::operator+(const SparseMatrix<T>& other) {
162     SparseMatrix<T> result(rows());
163

```

```

164     for(int i = 0; i < rows(); i++) {
165         result(i) = _values[i] + other(i);
166     }
167
168     return result;
169 }
170
171 template<class T>
172 SparseMatrix<T>& SparseMatrix<T>::operator+=(const SparseMatrix<T>& other) {
173     SparseMatrix<T> result = (*this) + other;
174     (*this) = result;
175     return (*this);
176 }
177
178 template<class T>
179 SparseMatrix<T> SparseMatrix<T>::operator-(const SparseMatrix<T>& other) {
180     SparseMatrix<T> result(rows());
181
182     for(int i = 0; i < rows(); i++) {
183         result(i) = _values[i] - other(i);
184     }
185
186     return result;
187 }
188
189 template<class T>
190 SparseMatrix<T>& SparseMatrix<T>::operator-=(const SparseMatrix<T>& other) {
191     SparseMatrix<T> result = (*this) - other;
192     (*this) = result;
193     return (*this);
194 }
195
196 template<class T>
197 SparseMatrix<T> SparseMatrix<T>::operator*(const T& scalar) {
198
199     vector<T> resValues(_values.size());
200
201     for(int i = 0; i < _values.size(); i++) {
202         resValues[i] = _values[i] * scalar;
203     }
204
205     SparseMatrix<T> result(resValues, _iValues, _jValues, _columns);
206
207     return result;
208 }
209
210 template<class T>
211 SparseMatrix<T> SparseMatrix<T>::operator/(const T& scalar) {
212     vector<T> resValues(_values.size());
213
214     for(int i = 0; i < _values.size(); i++) {
215         resValues[i] = _values[i] / scalar;
216     }
217
218     SparseMatrix<T> result(resValues, _iValues, _jValues, _columns);
219
220     return result;
221 }
222

```

```

223 template<class T>
224 T& SparseMatrix<T>::operator()(int a) {
225     return _values[a];
226 }
227
228 template<class T>
229 const T& SparseMatrix<T>::operator()(const int a) const {
230     return _values[a];
231 }
232
233 template<class T>
234 Matrix<T> SparseMatrix<T>::descompress() {
235     Matrix<T> result(rows(), _columns, 0);
236     for (int i = 0; i < rows(); i++) {
237         for (int j = _iValues[i]; j < _iValues[i+1]; j++) {
238             result(i, _jValues[j]) = _values[j];
239         }
240     }
241
242     return result;
243 }
244
245 template<class T>
246 double SparseMatrix<T>::L1(const SparseMatrix<T>& other) {
247     SparseMatrix<T> vectorSubs = *this - other;
248     double res = 0;
249     for (int i = 0; i < rows(); i++){
250         res = res + abs(vectorSubs(i));
251     }
252
253     return res;
254 }
255
256 template<class T>
257 double SparseMatrix<T>::norm1() {
258     double res = 0;
259     for (int i = 0; i < rows(); i++){
260         res = res + abs(_values[i]);
261     }
262
263     return res;
264 }
265
266 template<class T>
267 int SparseMatrix<T>::rows() {
268     return _iValues.size() - 1;
269 }
270
271 template<class T>
272 int SparseMatrix<T>::columns() {
273     return _columns;
274 }
275
276 template<class T>
277 int SparseMatrix<T>::rows() const {
278     return _iValues.size() - 1;
279 }
280
281 template<class T>

```

```
282 int SparseMatrix<T>::columns() const{
283     return _columns;
284 }
285
286 template<class T>
287 void SparseMatrix<T>::printSparseMatrix() {
288     Matrix<T> descompressedMatrix = descompress();
289     descompressedMatrix.printMatrix();
290 }
291
292 #endif /* SparseMatrix_H */
```

---



## Referencias

- [1] Datahub. <http://datahub.io>.
- [2] Stanford large network dataset collection. <http://snap.stanford.edu/data/#web>.
- [3] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems*, 30(1-7):107–117, April 1998.
- [4] Kurt Bryan and Tanya Leise. The linear algebra behind google. *SIAM Review*, 48(3):569–581, 2006.
- [5] Christopher Clarey. Years later for guillermo vilas, he’s still not the one. *The New York Times*.
- [6] Angela Y. Govan, Carl D. Meyer, and Rusell Albright. Generalizing google’s pagerank to rank national football league teams. In *Proceedings of SAS Global Forum 2008*, 2008.
- [7] Sepandar D. Kamvar, Taher H. Haveliwala, Christopher D. Manning, and Gene H. Golub. Extrapolation methods for accelerating pagerank computations. In *Proceedings of the 12th international conference on World Wide Web*, WWW ’03, pages 261–270, New York, NY, USA, 2003. ACM.