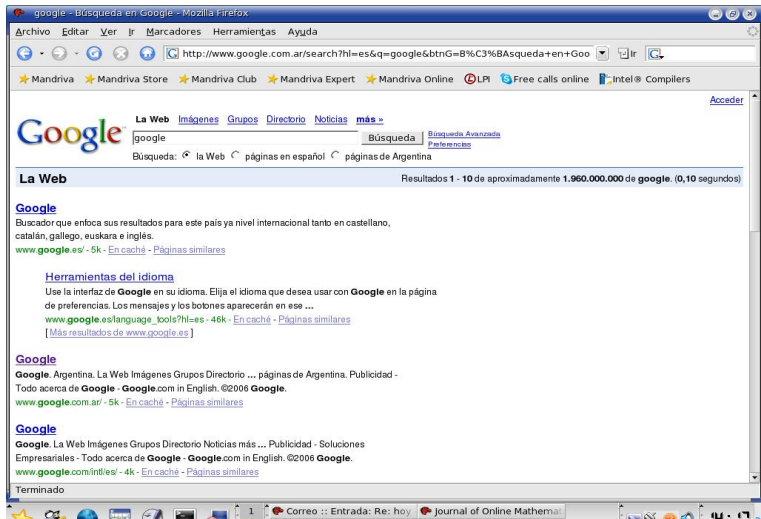


Buscando "Google" en Google

Gabriel Acosta

April 27, 2007

- 1 Introducción
- 2 Ordenando Contenidos
- 3 Ordenando Las Páginas
- 4 Historia...



Algunos detalles

- Numero de paginas! = 1.900.000.000.

Algunos detalles

- Numero de paginas! = 1.900.000.000.
- Tiempo de la busqueda = 0.10 segundos!.

Algunos detalles

- Numero de paginas! = 1.900.000.000.
- Tiempo de la busqueda = 0.10 segundos!.
- El paginado es cada 10 paginas: visitando 1 por segundo ...
aproximadamente 6,3 años.

Algunos detalles

- Numero de paginas! = 1.900.000.000.
- Tiempo de la busqueda = 0.10 segundos!.
- El paginado es cada 10 paginas: visitando 1 por segundo ...
aproximadamente 6,3 años.
- Parece que solo muestran los 1000 primeros (quien tiene tanta paciencia!)

Algunos detalles

- Numero de paginas! = 1.900.000.000.
- Tiempo de la busqueda = 0.10 segundos!.
- El paginado es cada 10 paginas: visitando 1 por segundo ...
aproximadamente 6,3 años.
- Parece que solo muestran los 1000 primeros (quien tiene tanta
paciencia!)
- Cuantos "Juanes" hay????

Algunos detalles

- Numero de paginas! = 1.900.000.000.
- Tiempo de la busqueda = 0.10 segundos!.
- El paginado es cada 10 paginas: visitando 1 por segundo ...
aproximadamente 6,3 años.
- Parece que solo muestran los 1000 primeros (quien tiene tanta
paciencia!)
- Cuantos "Juanes" hay????
- Criterios de presentación de la información.

La Importancia de las Cosas ...

- 1 Importancia de los contenidos.

La Importancia de las Cosas ...

- 1 Importancia de los contenidos.
- 2 Importancia de las páginas

La Importancia de Los Contenidos

- Pagina de la Biblioteca Digital de la Mathematical Association of America
<http://mathdl.maa.org/>

La Importancia de Los Contenidos

- Pagina de la Biblioteca Digital de la Mathematical Association of America
<http://mathdl.maa.org/>
- Supongamos que tenemos siete páginas P_1, P_2, \dots, P_7 con contenidos sobre "Postres", "Panes" y "Vegetales".

La Importancia de Los Contenidos

- Pagina de la Biblioteca Digital de la Mathematical Association of America
<http://mathdl.maa.org/>
- Supongamos que tenemos siete páginas P_1, P_2, \dots, P_7 con contenidos sobre "Postres", "Panes" y "Vegetales".
- Aunque no todas tienen la misma cantidad de información sobre cada item. Por ejemplo P_1 se dedica solo a "Postres", mientras que por ejemplo P_3 dedica el 70% a "Panes" y el 30% a "Vegetales" sin dar información sobre "Postres"

La Importancia de Los Contenidos

- Pagina de la Biblioteca Digital de la Mathematical Association of America
<http://mathdl.maa.org/>
- Supongamos que tenemos siete páginas P_1, P_2, \dots, P_7 con contenidos sobre "Postres", "Panes" y "Vegetales".
- Aunque no todas tienen la misma cantidad de información sobre cada item. Por ejemplo P_1 se dedica solo a "Postres", mientras que por ejemplo P_3 dedica el 70% a "Panes" y el 30% a "Vegetales" sin dar información sobre "Postres"
- Cómo podemos cuantificar esto?.

Vectorizamos los datos

- Podemos tomar para P_1 y P_3

$$\begin{array}{l} \text{Postres} \\ \text{Panes} \\ \text{Vegetales} \end{array} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{array}{l} \text{Postres} \\ \text{Panes} \\ \text{Vegetales} \end{array} \begin{pmatrix} 0 \\ 0.7 \\ 0.3 \end{pmatrix}$$

Vectorizamos los datos

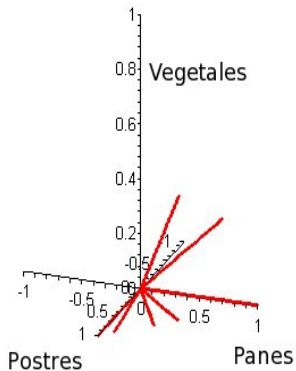
- Podemos tomar para P_1 y P_3

Postres	($\begin{matrix} 1 \\ 0 \\ 0 \end{matrix}$)	Postres	($\begin{matrix} 0 \\ 0.7 \\ 0.3 \end{matrix}$)
Panes		Panes	
Vegetales		Vegetales	

- Si hacemos lo propio con las restantes páginas podemos juntar la información en una "matriz"

Tipo de Comida	P_1	P_2	...	P_7	
Postres	($\begin{matrix} 1 & 0 & \dots & 0.9 \\ 0 & 1 & \dots & 0.1 \\ 0 & 0 & \dots & 0 \end{matrix}$)	1	0	...	0.9
Panes		0	1	...	0.1
Vegetales		0	0	...	0

Interpretación Geometrica

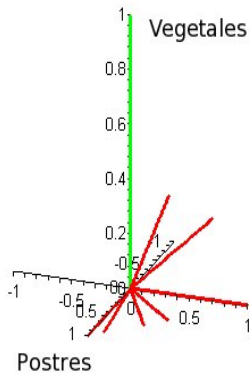


Busqueda "Geometrica"

Imaginemos que buscamos información sobre "Vegetales". Nuestro criterio vectorizado será

$$\begin{array}{l} \text{Postres} \\ \text{Panes} \\ \text{Vegetales} \end{array} \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

Graficamente



Busqueda con un vector general

- Si tenemos un vector general V con el criterio que nos interesa

Busqueda con un vector general

- Si tenemos un vector general V con el criterio que nos interesa
-

$$\begin{array}{l} \text{Postres} \\ \text{Panes} \\ \text{Vegetales} \end{array} \begin{pmatrix} a \\ b \\ c \end{pmatrix}$$

Busqueda con un vector general

- Si tenemos un vector general V con el criterio que nos interesa



$$\begin{array}{l} \text{Postres} \\ \text{Panes} \\ \text{Vegetales} \end{array} \begin{pmatrix} a \\ b \\ c \end{pmatrix}$$

- Como elegimos las páginas relevantes?

Cuestión de ángulos

- Fijamos un cierto ángulo α

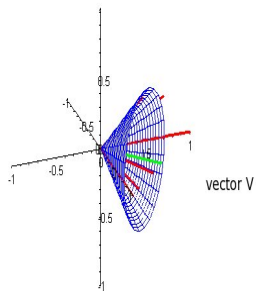
Cuestión de ángulos

- Fijamos un cierto ángulo α
- Y elegimos las páginas cuyo vector forme un ángulo menor a α respecto de V .

Cuestión de ángulos

- Fijamos un cierto ángulo α
- Y elegimos las páginas cuyo vector forme un ángulo menor a α respecto de V .
- Dicho de otro modo: construimos un "cono" con eje en V y elegimos las paginas que caen dentro del cono

Graficamente



Orden????

Damos las páginas (vectores) dentro del cono ordenados por ángulos crecientes!.

En la realidad ...

- Elección cuidadosa del α

En la realidad ...

- Elección cuidadosa del α
- Un α chico nos deja con una búsqueda muy estricta (tal vez sin resultados!), un α grande con demasiados resultados!!!

En la realidad ...

- Elección cuidadosa del α
- Un α chico nos deja con una búsqueda muy estricta (tal vez sin resultados!), un α grande con demasiados resultados!!!
- Consideremos **todas** las palabras del diccionario (en inglés hay cerca de 300.000).

En la realidad ...

- Elección cuidadosa del α
- Un α chico nos deja con una búsqueda muy estricta (tal vez sin resultados!), un α grande con demasiados resultados!!!
- Consideremos **todas** las palabras del diccionario (en inglés hay cerca de 300.000).
- Ordenadas alfabéticamente de "aarónico" a "zuzón". Una búsqueda por las palabras "deportes" y "playa" puede vectorizarse como hicimos antes ...

Vector de búsqueda

aarónico	0
.	.
.	.
deportes	1
.	.
.	.
playa	1
.	.
.	.
zuzón	0

Muchas cuentas para hacer!!!

Paginas en Internet $\sim 3.000.000.000$

$$G = \begin{pmatrix} \text{Palabras} & \text{Pag1} & \text{Pag2} & \dots & \text{Pag } 3.000.000.000 \\ \text{palabra 1} & 0 & 1 & \dots & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & 1 & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & 1 \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \text{palabra } 300.000 & 0 & 0 & \dots & \cdot \end{pmatrix}$$

Demasiada Información!

- Bueno: se puede paralelizar!

Demasiada Información!

- Bueno: se puede paralelizar!
- Malo: hay que almacenar y mover mucha información!

Demasiada Información!

- Bueno: se puede paralelizar!
- Malo: hay que almacenar y mover mucha información!
- Bit 0 o 1 mínima información posible.

Demasiada Información!

- Bueno: se puede paralelizar!
- Malo: hay que almacenar y mover mucha información!
- Bit 0 o 1 mínima información posible.
- Byte son 8 bits: p.ej. 10101110

Demasiada Información!

- Bueno: se puede paralelizar!
- Malo: hay que almacenar y mover mucha información!
- Bit 0 o 1 mínima información posible.
- Byte son 8 bits: p.ej. 10101110
- Almacenar un número en una PC ocupa 4 bytes (precisión simple)

Demasiada Información!

- Bueno: se puede paralelizar!
- Malo: hay que almacenar y mover mucha información!
- Bit 0 o 1 mínima información posible.
- Byte son 8 bits: p.ej. 10101110
- Almacenar un número en una PC ocupa 4 bytes (precisión simple)
- Información en la matriz:

$$300.000 \times 3.000.000.000 = 90.000.000.000.000 \sim 360.000 \text{Gigas}$$

Demasiada Información!

- Bueno: se puede paralelizar!
- Malo: hay que almacenar y mover mucha información!
- Bit 0 o 1 mínima información posible.
- Byte son 8 bits: p.ej. 10101110
- Almacenar un número en una PC ocupa 4 bytes (precisión simple)
- Información en la matriz:

$$300.000 \times 3.000.000.000 = 90.000.000.000.000 \sim 360.000 \text{ Gigas}$$

- Con discos de 100 Gigas ... 3600 discos rígidos!!!

La matematica vuelve a ayudar!

Se usan técnicas para "achicar" la información.

LSI (Indexado por Semantica Latente)

Descomposicion en valores singulares:

$$G = UDV$$

No hace falta tomar toda la matriz D para "aproximar" G . Para saber mas de este tema:

La matematica vuelve a ayudar!

Se usan técnicas para "achicar" la información.

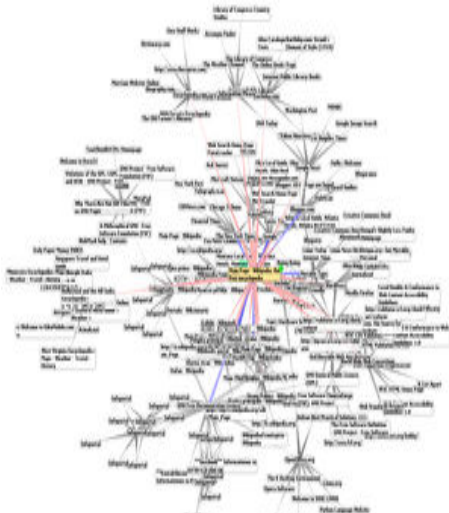
LSI (Indexado por Semantica Latente)

Descomposicion en valores singulares:

$$G = UDV$$

No hace falta tomar toda la matriz D para "aproximar" G . Para saber mas de este tema: Pueden hacer la Licenciatura en Matemática!

El Problema



Criterio de Importancia de Páginas: el PageRank

- www.uam.es/personal_pdi/ciencias/gallardo/upm_google.pdf

Criterio de Importancia de Páginas: el PageRank

- www.uam.es/personal_pdi/ciencias/gallardo/upm_google.pdf
- Como cuantificar la importancia de una página?

Criterio de Importancia de Páginas: el PageRank

- www.uam.es/personal_pdi/ciencias/gallardo/upm_google.pdf
- Como cuantificar la importancia de una página?
- Una vez más consideremos todas las páginas de Internet:

$$P_1, P_2, \dots, P_{3.000.000.000}$$

y llamemos

$$x_1, x_2, \dots, x_{3.000.000.000}$$

a la importancia que le atribuimos a cada una.

Criterio de Importancia de Páginas: el PageRank

- www.uam.es/personal_pdi/ciencias/gallardo/upm_google.pdf
- Como cuantificar la importancia de una página?
- Una vez más consideremos todas las páginas de Internet:

$$P_1, P_2, \dots, P_{3.000.000.000}$$

y llamemos

$$x_1, x_2, \dots, x_{3.000.000.000}$$

a la importancia que le atribuimos a cada una.

- Google utiliza un algoritmo llamado PageRank

Mirando "Links"

El "truco" para cuantificar la importancia de cada pagina esta en relacionarla con las demás: esto se hace viendo los links que recibe.

- Si una pagina recibe muchos links debe ser importante

Mirando "Links"

El "truco" para cuantificar la importancia de cada pagina esta en relacionarla con las demás: esto se hace viendo los links que recibe.

- Si una pagina recibe muchos links debe ser importante
- Pero si yo recibo un solo link proveniente de Microsoft?????

Mirando "Links"

El "truco" para cuantificar la importancia de cada pagina esta en relacionarla con las demás: esto se hace viendo los links que recibe.

- Si una pagina recibe muchos links debe ser importante
- Pero si yo recibo un solo link proveniente de Microsoft?????
- La importancia de una pagina es proporcional a la suma de las importancias de las paginas que linkean a ella.

Mas Matemática!

- $x_1 = K(x_{27} + x_{1235})$

Mas Matemática!

- $x_1 = K(x_{27} + x_{1235})$
- Viejo conocido: Sistema de ecuaciones!!!

$$\left\{ \begin{array}{lcl} x_1 & = & K(x_{27} + x_{1235}) \\ x_2 & = & K(x_{132} + x_{1256} + x_{5689}) \\ \cdot & = & \dots\dots\dots \\ \cdot & = & \dots\dots\dots \\ \cdot & = & \dots\dots\dots \\ x_{3.000.000.000} & = & K(x_{45} + x_{67}) \end{array} \right.$$

Mas Matemática!

- $x_1 = K(x_{27} + x_{1235})$
- Viejo conocido: Sistema de ecuaciones!!!

$$\left\{ \begin{array}{lcl} x_1 & = & K(x_{27} + x_{1235}) \\ x_2 & = & K(x_{132} + x_{1256} + x_{5689}) \\ \cdot & = & \dots\dots\dots \\ \cdot & = & \dots\dots\dots \\ \cdot & = & \dots\dots\dots \\ x_{3.000.000.000} & = & K(x_{45} + x_{67}) \end{array} \right.$$

- 3.000.000.000 de ecuaciones y 3.000.000.001 incognitas!!!

Problemas de autovalores

- Se llama un problema de autovalores

Problemas de autovalores

- Se llama un problema de autovalores
- Aplicaciones en Ingeniería, Física, Biología, Economía

Problemas de autovalores

- Se llama un problema de autovalores
- Aplicaciones en Ingeniería, Física, Biología, Economía
- Mucha teoría sobre esto: y no es nueva ...

Problemas de autovalores

- Se llama un problema de autovalores
- Aplicaciones en Ingeniería, Física, Biología, Economía
- Mucha teoría sobre esto: y no es nueva ...
- ~ 1910 Perron y Frobenius

Problemas de autovalores

- Se llama un problema de autovalores
- Aplicaciones en Ingeniería, Física, Biología, Economía
- Mucha teoría sobre esto: y no es nueva ...
- ~ 1910 Perron y Frobenius
- Para saber mas ...

Problemas de autovalores

- Se llama un problema de autovalores
- Aplicaciones en Ingeniería, Física, Biología, Economía
- Mucha teoría sobre esto: y no es nueva ...
- ~ 1910 Perron y Frobenius
- Para saber mas ...
- Licenciatura en Matemática!!!

Terminando!

- En 1939 Edward Kasner 10^{100}

Terminando!

- En 1939 Edward Kasner 10^{100}



$$10^{100} = 1 \underbrace{0000000000000000...0}_{100 \text{ ceros}}$$

Terminando!

- En 1939 Edward Kasner 10^{100}



$$10^{100} = 1 \underbrace{0000000000000000...0}_{100 \text{ ceros}}$$

- Googol

Terminando!

- En 1939 Edward Kasner 10^{100}



$$10^{100} = 1 \underbrace{0000000000000000\dots0}_{100 \text{ ceros}}$$

- Googol
- De Arena y de Atomos.

Terminando!

- En 1939 Edward Kasner 10^{100}



$$10^{100} = 1 \underbrace{0000000000000000...0}_{100 \text{ ceros}}$$

- Googol
- De Arena y de Atomos.
- En 1998 Sergei Brin y Lawrence Page crearon Google .