

Visualización de Datos con Python 3.7 y D3.js



Alfonso Neil Jiménez Casallas
Ingeniero de sistemas
Pontificia Universidad Javeriana

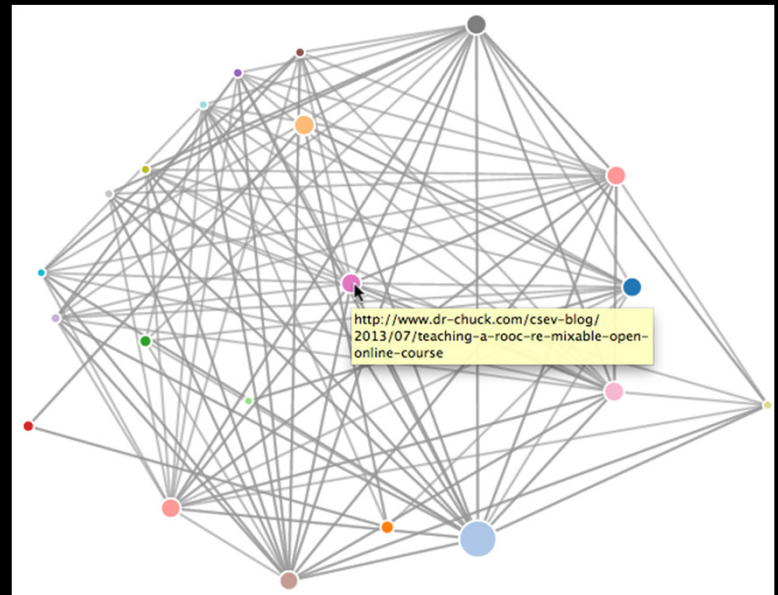
[https://github.com/pochito427/
alfonso.jimenez@javeriana.edu.co](https://github.com/pochito427/alfonso.jimenez@javeriana.edu.co)

<https://www.linkedin.com/in/alfonso-jimenez-3208a120/>



Caso de estudio: Page Rank

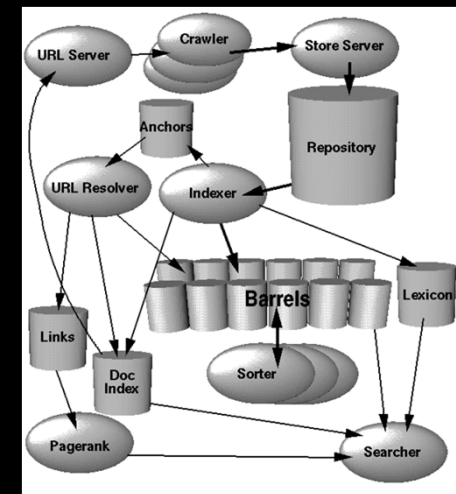
- Escribir un rastreador (crawler) simple de página Web
- Calcular una versión simple del algoritmo Page Rank de Google
- Visualizar la red resultante



<http://www.py4e.com/code3/pagerank.zip>

Arquitectura de Motores de Búsqueda

- Rastreo Web (Web Crawling)
- Indexación (Index Building)
- Búsqueda (Searching)



<http://infolab.stanford.edu/~backrub/google.html>

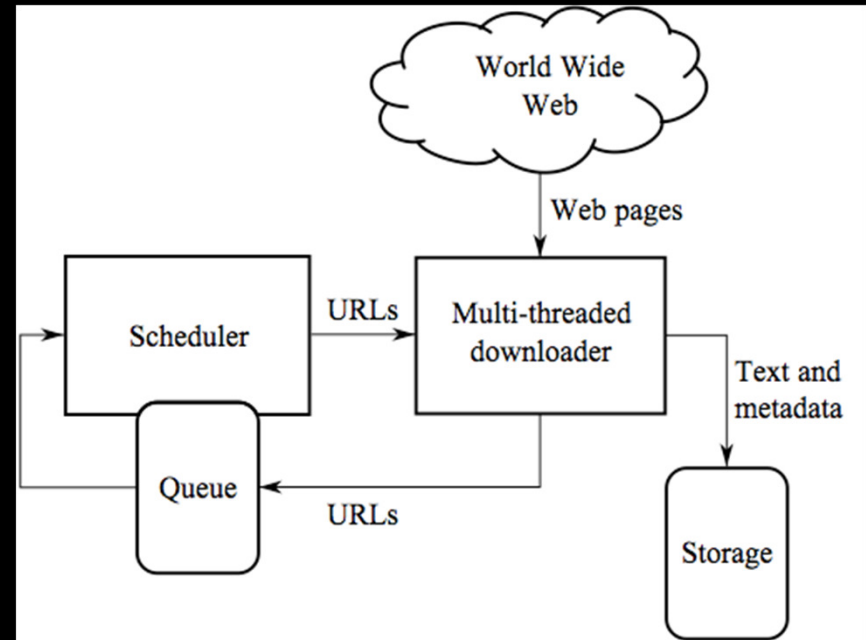
Rastreador Web (Web Crawler)

Programa en computador que comprende la Internet (World Wide Web) de una manera metódica y automatizada. Es usado principalmente para crear una copia de todas las páginas visitadas para procesamiento posterior por un motor de búsquedas que indexará las páginas descargadas para proveer búsquedas rápidas.

http://en.wikipedia.org/wiki/Web_crawler

Web Crawler

- Recupera una página
- Examina la página para encontrar enlaces
- Agrega los enlaces a una lista de sitios “a ser recuperados”
- Repite el proceso...



http://en.wikipedia.org/wiki/Web_crawler

Políticas de Rastreo Web

- Una **política de selección** que establece cuáles páginas a descargar,
- Una **política de re-visitas** que establece cuándo verificar cambios a las páginas,
- Una **política de cortesía** que establece cómo evitar la sobrecarga de sitios Web, y
- Una **política de paralelización** que establece cómo coordinar rastreadores Web distribuidos

robots.txt

- Una vía para un sitio web de comunicar con los rastreadores web
- Un estándar informal y voluntario

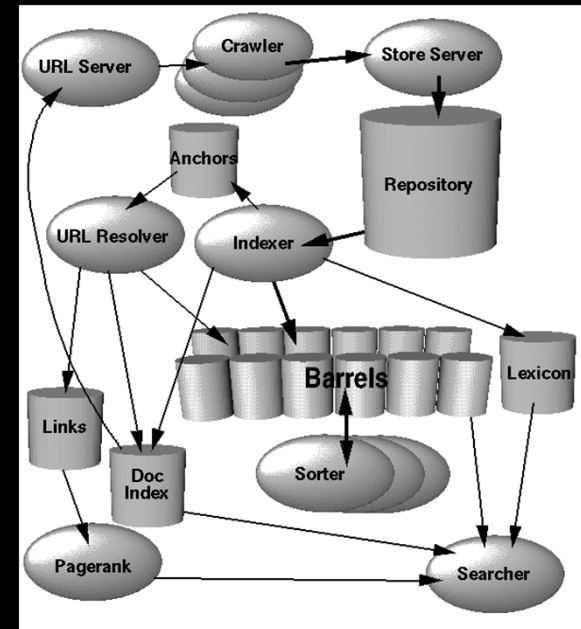
User-agent: *
Disallow: /cgi-bin/
Disallow: /images/
Disallow: /tmp/
Disallow: /private/



http://en.wikipedia.org/wiki/Robots_Exclusion_Standard
http://en.wikipedia.org/wiki/Spider_trap

Arquitectura Google

- Web Crawling
- Index Building
- Searching



<http://infolab.stanford.edu/~backrub/google.html>

Indexación de Búsquedas

La indexación en motores de búsqueda recopila, analiza, y almacena datos para facilitar la recuperación rápida y precisa de información. El propósito de almacenar un índice es optimizar velocidad y rendimiento en encontrar documentos relevantes para una consulta de búsqueda. Sin un índice, el motor de búsquedas exploraría cada documento en el cuerpo, lo cual requeriría un tiempo y poder de cómputo considerables.

[http://en.wikipedia.org/wiki/Index_\(search_engine\)](http://en.wikipedia.org/wiki/Index_(search_engine))



D3.js

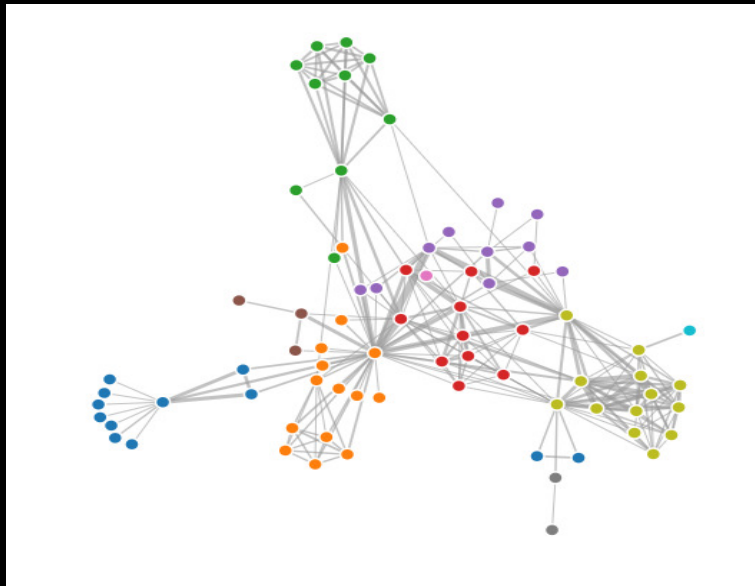
“D3.js es una biblioteca de JavaScript para manipular documentos basados en datos. D3 ayuda a dar vida a los datos usando HTML, SVG y CSS. El énfasis de D3 en los estándares web le brinda todas las capacidades de los navegadores modernos sin atarse a un marco propietario, que combina componentes de visualización potentes y un enfoque basado en datos para la manipulación DOM.”

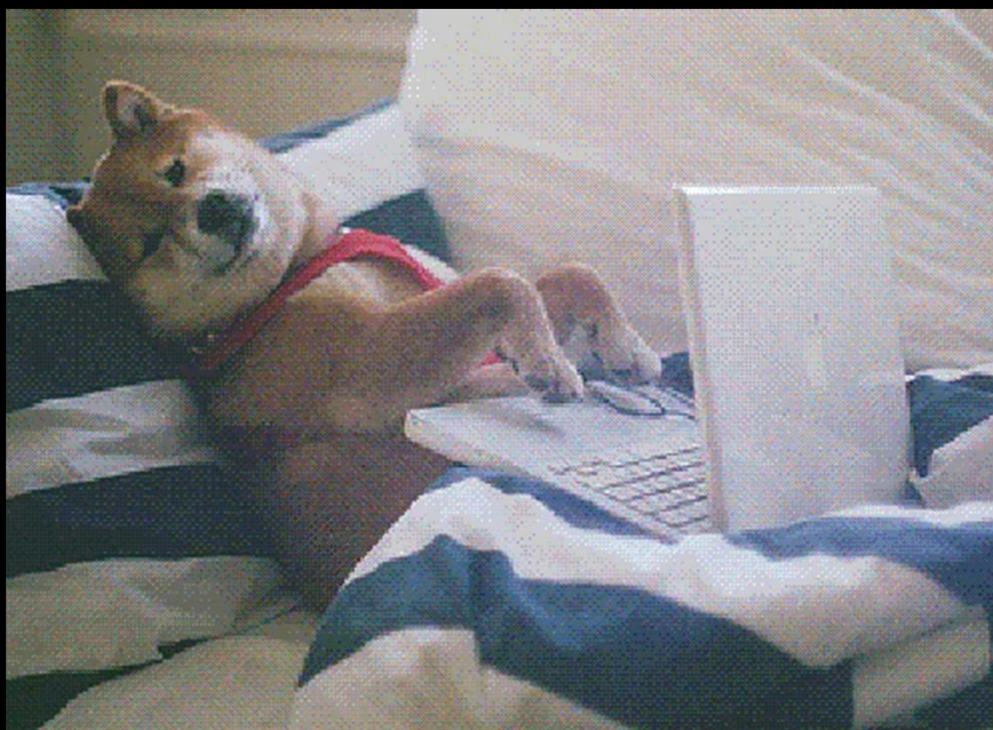
<https://d3js.org/>

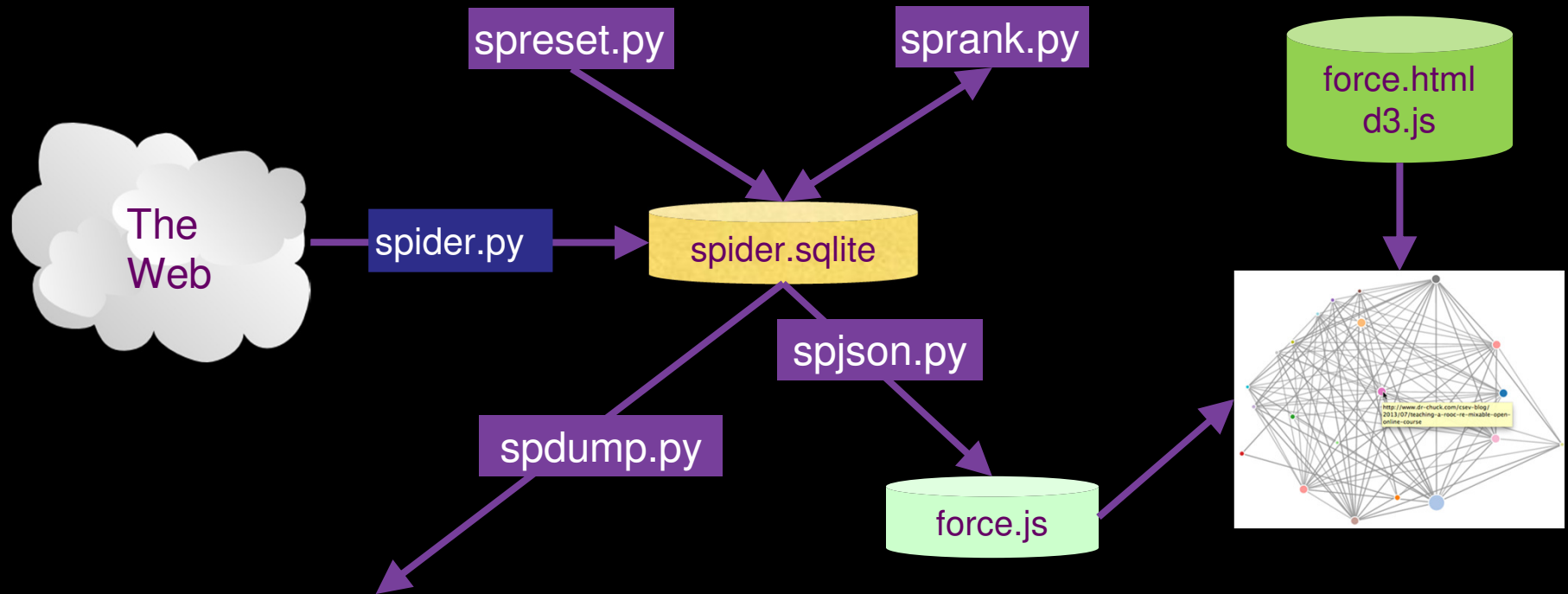
Aplicaciones en grafos dirigidos

Red de co-ocurrencias entre caracteres en la novela “Los Miserables”

<https://observablehq.com/@d3/force-directed-graph>







(5, None, 1.0, 3, u'http://www.dr-chuck.com/csev-blog')
(3, None, 1.0, 4, u'http://www.dr-chuck.com/dr-chuck/resume/speaking.htm')
(1, None, 1.0, 2, u'http://www.dr-chuck.com/csev-blog')
(1, None, 1.0, 5, u'http://www.dr-chuck.com/dr-chuck/resume/index.htm')
4 rows.

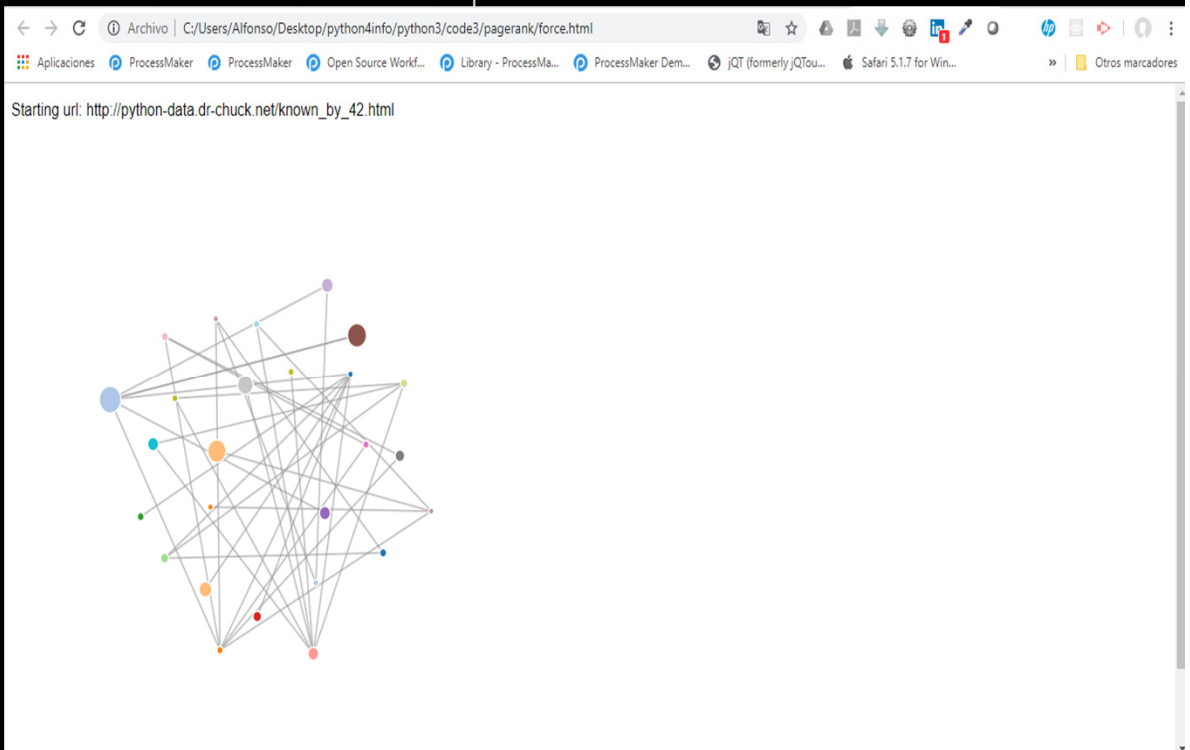
<http://www.py4e.com/code3/pagerank.zip>

Prueba del PageRank para la URL

<http://python-data.dr-chuck.net>

```
99 1.6076249242715756e-17
100 9.488864938011857e-18
[(1, 0.07935445138366898), (7, 0.15870890276733796), (20, 0.8659724559084117), (137, 0.20363623131344208), (187, 0.20363623131344208)]
```

```
C:\Users\Alfonso\Desktop\python4info\python3\code3\pagerank\python3_spdup.py
(6, 2.497974849443514, 2.497974849443512, 167, 'http://python-data.dr-chuck.net/known_by_Caitaidh.html')
(6, 2.0063914094917608, 2.006391409491774, 217, 'http://python-data.dr-chuck.net/known_by_Peebles.html')
(5, 1.4480960457637058, 1.4480960457636758, 8, 'http://python-data.dr-chuck.net/known_by_Miranne.html')
(5, 3.1582031724081374, 3.15820317240812, 545, 'http://python-data.dr-chuck.net/known_by_Joojo.html')
(4, 1.2335786043101776, 1.2335786043101805, 424, 'http://python-data.dr-chuck.net/known_by_Data.html')
(4, 1.648026818689624, 1.648026818689495, 730, 'http://python-data.dr-chuck.net/known_by_Charyl.html')
(4, 2.835276303045442, 2.835276303045417, 987, 'http://python-data.dr-chuck.net/known_by_Josese.html')
(4, 2.2881799417389597, 2.2881799417389543, 1282, 'http://python-data.dr-chuck.net/known_by_Constance.html')
(4, 4.339563444773513, 4.339563444773419, 1342, 'http://python-data.dr-chuck.net/known_by_Ziya.html')
(4, 1.992283503264014, 1.9922835032639936, 2093, 'http://python-data.dr-chuck.net/known_by_Denise.html')
(4, 1.8624241535669017, 1.8624241535669072, 2206, 'http://python-data.dr-chuck.net/known_by_Iseaseel.html')
(3, 0.2444103408649326, 0.24441034086493212, 149, 'http://python-data.dr-chuck.net/known_by_Abbeygail.html')
(3, 1.6811212196759728, 1.6811212196760015, 164, 'http://python-data.dr-chuck.net/known_by_Khyla.html')
(3, 3.085221427230209, 3.085221427230229, 307, 'http://python-data.dr-chuck.net/known_by_Terry.html')
(3, 1.5113028237495718, 1.5113028237495914, 321, 'http://python-data.dr-chuck.net/known_by_Jayvi.html')
(3, 0.6091310439689915, 0.6091310439689912, 354, 'http://python-data.dr-chuck.net/known_by_Patricia.html')
(3, 1.4981198898079757, 1.4981198898080001, 357, 'http://python-data.dr-chuck.net/known_by_Miller.html')
(3, 1.5095226580053214, 1.509522658005301, 489, 'http://python-data.dr-chuck.net/known_by_Jing.html')
(3, 1.9006265377924811, 1.900626537792471, 507, 'http://python-data.dr-chuck.net/known_by_Saranta.html')
(3, 0.8687066123785291, 0.8687066123785212, 538, 'http://python-data.dr-chuck.net/known_by_Prinsrose.html')
(3, 0.7384613274092983, 0.7384613274092806, 595, 'http://python-data.dr-chuck.net/known_by_Shawnpaul.html')
(3, 0.7502507040641068, 0.7502507040641102, 598, 'http://python-data.dr-chuck.net/known_by_Aliiza.html')
(3, 1.6342186170851623, 1.6342186170851405, 695, 'http://python-data.dr-chuck.net/known_by_Ioanna.html')
(3, 1.293583132084372, 1.293583132084266, 1106, 'http://python-data.dr-chuck.net/known_by_Clio.html')
(3, 1.3828509371082147, 1.3828509371082198, 1236, 'http://python-data.dr-chuck.net/known_by_Kristoffer.html')
(3, 3.503973288820543, 3.5039732888205166, 1426, 'http://python-data.dr-chuck.net/known_by_Chelsey.html')
(3, 0.925661434404678, 0.9256614344046692, 1529, 'http://python-data.dr-chuck.net/known_by_Saiba.html')
(3, 1.4194292180343904, 1.419429218034385, 1567, 'http://python-data.dr-chuck.net/known_by_Kadin.html')
(3, 1.910817490097222, 1.910817490097201, 1644, 'http://python-data.dr-chuck.net/known_by_Dilano.html')
(3, 0.9951473408301005, 0.9951473409301028, 1842, 'http://python-data.dr-chuck.net/known_by_Iristain.html')
(3, 0.9409792280921202, 0.9409792280921127, 2449, 'http://python-data.dr-chuck.net/known_by_Kodi.html')
(3, 2.589854837527781, 2.5898548375278367, 2609, 'http://python-data.dr-chuck.net/known_by_Luella.html')
(3, 1.3289246970699113, 1.3289246970699053, 3165, 'http://python-data.dr-chuck.net/known_by_Ander.html')
(3, 2.855752325228588, 2.855752325228586, 3296, 'http://python-data.dr-chuck.net/known_by_Kofi.html')
(3, 2.6457840142542155, 2.6457840142542572, 3524, 'http://python-data.dr-chuck.net/known_by_Iyllor.html')
(3, 1.5719400805373084, 1.5719400805372968, 4015, 'http://python-data.dr-chuck.net/known_by_Daisy.html')
(2, 0.8699724559084274, 0.8699724559084117, 20, 'http://python-data.dr-chuck.net/known_by_London.html')
(2, 0.6534634660675016, 0.6534634660675178, 125, 'http://python-data.dr-chuck.net/known_by_Harris.html')
(2, 0.3147506142819535, 0.31475061428195, 159, 'http://python-data.dr-chuck.net/known_by_Cale.html')
(2, 0.14161727908655408, 0.1416172790865533, 273, 'http://python-data.dr-chuck.net/known_by_Maillie.html')
(2, 0.15825819673809394, 0.15825819673809266, 280, 'http://python-data.dr-chuck.net/known_by_Millie.html')
(2, 0.282398132706626, 0.282398132706666, 454, 'http://python-data.dr-chuck.net/known_by_Suvi.html')
(2, 0.5385834005152125, 0.53858340051523, 492, 'http://python-data.dr-chuck.net/known_by_Harry.html')
(2, 0.638354027240914, 0.63835402724089, 542, 'http://python-data.dr-chuck.net/known_by_Mehmet.html')
(2, 0.8530121861453319, 0.8530121861453168, 841, 'http://python-data.dr-chuck.net/known_by_Forbes.html')
(2, 0.7377150405356286, 0.7377150405356301, 870, 'http://python-data.dr-chuck.net/known_by_Charlie.html')
(2, 0.618371924941718, 0.6183719249417226, 887, 'http://python-data.dr-chuck.net/known_by_Khalen.html')
(2, 1.7183690761066418, 1.7183690761066265, 1804, 'http://python-data.dr-chuck.net/known_by_Kalen.html')
(2, 1.467526126929276, 1.4675261269292803, 1024, 'http://python-data.dr-chuck.net/known_by_Avinash.html')
(2, 0.2339474326191034, 0.23394743261909962, 1104, 'http://python-data.dr-chuck.net/known_by_Fatiha.html')
100 rows.
```



Prueba del PageRank para la URL <https://www.javeriana.edu.co>

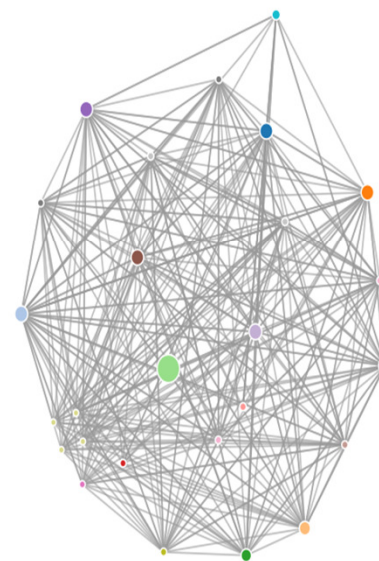
```
(1, 0.0013671286293981483), (24, 0.00011703442315089238), (2, 0.001408339412141812), (101, 0.0016018235083929142), (48, 0.001341412368514083)
```

```
C:\Users\Alfonso\Desktop\python4info\python\code3\pagerank\python3\spdump.py
146, 0.000514140741757064, 0.001408335412141812, 2, 'https://www.javeriana.edu.co/educon'
108, 23.886879451437462, 23.90138366370368, 246, 'https://www.javeriana.edu.co/pesquisa/perfiles/huellas'
108, 21.23279224482974, 21.2456764001068, 247, 'https://www.javeriana.edu.co/pesquisa/perfiles/entrevistas'
104, 0.00041614077933761, 0.001307186230918415, 1, 'https://www.javeriana.edu.co/'
187, 0.000321199151211728, 0.001366230804180302, 46, 'https://www.javeriana.edu.co/programas/posgrados'
174, 0.004277408870995728, 0.0009247416726240122, 15, 'https://www.javeriana.edu.co/extension'
66, 0.006285473717999223, 0.001414062619348786, 4, 'https://www.javeriana.edu.co/editorial'
68, 0.004869242742863759, 0.0018527139736375471, 13, 'https://www.javeriana.edu.co/institucional/financiacin-becas'
68, 0.00626497333517837, 0.001343417236851408, 48, 'https://www.javeriana.edu.co/programas/asignaturas-de-libre-escoencia'
68, 0.00626497333517837, 0.001343417236851408, 48, 'https://www.javeriana.edu.co/programas/listado-de-curso-preuniversitarios'
68, 0.000307067427452106, 0.001320420740954516, 51, 'https://www.javeriana.edu.co/institucional/centros-institutos-y-observatorios'
53, 4.540975863111828, 4.552648484342692, 273, 'https://www.javeriana.edu.co/pesquisa/hipopotamos-en-colombia-un-problema-de-enormes-dimensiones'
151, 23.886879451437466, 23.90138366370368, 244, 'https://www.javeriana.edu.co/pesquisa/perfiles'
48, 0.0007441554417619359, 0.0001688839342907155, 104, 'https://www.javeriana.edu.co/educon/politica-y-sociedad'
35, 0.0007489104134584624, 0.00016818235083929142, 108, 'https://www.javeriana.edu.co/educon/comunicacion-y-lenguaje'
34, 0.0007489104134584624, 0.00016818235083929142, 101, 'https://www.javeriana.edu.co/educon/artes'
32, 0.0007369824975918429, 0.00015933147782745, 89, 'https://www.javeriana.edu.co/educon/salud'
32, 0.0007369824975918429, 0.00015933147782745, 98, 'https://www.javeriana.edu.co/educon/estudios-ambientales-y-rurales'
108, 0.0007321307839593358, 0.00015828422154506564, 88, 'https://www.javeriana.edu.co/educon/ciencias-juridicas'
108, 0.0007321307839593357, 0.00015828422154506564, 111, 'https://www.javeriana.edu.co/educon/region-oriente'
23, 0.0007346913510673567, 0.00015883781001072822, 70, 'https://www.javeriana.edu.co/web/educon/programacion'
26, 0.00258378365165238, 0.00054229234442421, 138, 'https://www.javeriana.edu.co/web/programas/especializaciones'
25, 0.00258378365165238, 0.00054229234442421, 135, 'https://www.javeriana.edu.co/web/programas/mestrrias'
25, 0.00019732514081375846, 0.0000973212052135, 159, 'https://www.javeriana.edu.co/programas/especializaciones'
23, 0.00232474776626376, 0.0005430438739435207, 227, 'https://www.javeriana.edu.co/web/programas/home'
23, 0.04251740217494831, 0.042470747822123485, 261, 'https://www.javeriana.edu.co/pesquisa/category/multimedia'
18, 0.021908040848632724, 0.02190637492418162, 640, 'https://www.javeriana.edu.co/pesquisa/tag/medio-ambiente'
16, 0.0005790719696538814, 0.00012593428049474528, 165, 'https://www.javeriana.edu.co/admisiones/admisiones-paso-a-paso-2'
15, 0.11740802634386783, 0.11747088088203761, 547, 'https://www.javeriana.edu.co/pesquisa/tag/peces'
15, 0.09744541080512065, 0.0975054851789339, 898, 'https://www.javeriana.edu.co/pesquisa/tag/ecosistema'
14, 0.000371844110014225, 0.35840404514634384-48, 489, 'https://www.javeriana.edu.co/web/institucional/pos-icetes'
14, 0.07384646199015406, 0.0177801491973754793, 711, 'https://www.javeriana.edu.co/pesquisa/tag/neorria'
13, 0.07384646029463804, 0.073867835880208, 275, 'https://www.javeriana.edu.co/pesquisa/tag/ingenieria'
12, 15.92598886658797, 15.93484121875219, 381, 'https://www.javeriana.edu.co/pesquisa/nas-alla-de-lo-evidente'
9, 0.0016245718549746247, 0.001611877787121188, 18, 'https://www.javeriana.edu.co/pesquisa/antes-de-votar-revise-a-que-se-pueden-comprometer-los-candidatos'
9, 0.0004267973751910924, 9.2893320234739664-48, 78, 'https://www.javeriana.edu.co/educon/apocaliptica-judeo-cristiana-y-de-san-juan'
9, 0.00109720897351513, 0.000231451204942327, 147, 'https://www.javeriana.edu.co/vice-rectoria-academica/home'
9, 7.21515136530959598-48, 1.54623453578455593-48, 284, 'https://www.javeriana.edu.co/web/educon/google-markettag-plataforma'
9, 7.151611955209598-48, 1.5461453578455593-48, 212, 'https://www.javeriana.edu.co/web/educon/proyectos-de-ficcion-y-no-ficcion'
9, 7.151611955209598-48, 1.5461453578455593-48, 214, 'https://www.javeriana.edu.co/web/educon/aplicacion-de-las-rda-en-bibliotecas'
8, 0.00037604938473214, 0.000374313835126233, 643, 'https://www.javeriana.edu.co/pesquisa/la-era-de-la-hegemonia-cuantica'
7, 0.000465178137920471, 0.00018129158224808438, 190, 'https://www.javeriana.edu.co/admisiones/registro'
7, 3.5993425068221285-48, 7.78166392339648-48, 492, 'https://www.javeriana.edu.co/web/educon/acopanamiento-en-el-duelo'
7, 3.40052508635382-48, 4.091211833651546-48, 496, 'https://www.javeriana.edu.co/web/educon/politica-y-sociedad/p_id=101_INSTANCE_641NVQhX8Kp_p_lifecycle=0p_p_state=normalp_p_mode=view0p_p_col_id=column-10p_p_col_pos=10p_p_col_count=10p_p_354232324_category=57879351'
7, 0.00303698367263736, 0.003034273470140866, 569, 'https://www.javeriana.edu.co/pesquisa/los-inborrables-anos-70'
7, 0.00063152198077985, 0.00063524893803697, 582, 'https://www.javeriana.edu.co/pesquisa/tag/presupuesto'
7, 0.00037604938473214, 0.000374313835126233, 652, 'https://www.javeriana.edu.co/pesquisa/medicina-e-ingenieria-se-unen-para-salvar-vidas'
7, 0.002528338143861826, 0.00252874217876598, 835, 'https://www.javeriana.edu.co/pesquisa/tag/caspeinos'
8, 3.74808498251144-48, 8.092421874826744-48, 210, 'https://www.javeriana.edu.co/web/educon/comunicacion-y-lenguaje/p_id=101_INSTANCE_641NVQhX8Kp_p_lifecycle=0p_p_state=normalp_p_mode=view0p_p_col_id=column-10p_p_col_pos=10p_p_col_count=10p_p_354232324_category=57879351'
6, 0.4354547764786114, 0.43582871842482155, 388, 'https://www.javeriana.edu.co/pesquisa/tag/agrocdema'
808 rows.
```

Archivo | C:\Users\Alfonso\Desktop\python4info\python3\code3\pagerank\force.html

Aplicaciones ProcessMaker ProcessMaker Open Source Work... Library - ProcessMa... ProcessMaker Dem... jQT (formerly jQTou... Safari 5.1.7

Starting url: <https://www.javeriana.edu.co>



Referencias y recursos adicionales

- <https://www.py4e.com/html3/16-viz>
- <http://www.dr-chuck.com/>
- <https://www.youtube.com/playlist?list=PLIRFEj9H3Oj7Bp8-DfGpfAfDBibIRfl5p>
- <http://infolab.stanford.edu/~backrub/google.html>
- <https://d3js.org/>
- <https://github.com/d3/d3-force>
- <https://observablehq.com/@d3/force-directed-graph>
- <https://bl.ocks.org/mbostock/ad70335eeef6d167bc36fd3c04378048>
- <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>
- <https://www.statisticshowto.datasciencecentral.com/absolute-error/>



Acknowledgements / Contributions



These slides are Copyright 2010- Charles R. Severance (www.dr-chuck.com) of the University of Michigan School of Information and open.umich.edu and made available under a Creative Commons Attribution 4.0 License. Please maintain this last slide in all copies of the document to comply with the attribution requirements of the license. If you make a change, feel free to add your name and organization to the list of contributors on this page as you republish the materials.

Initial Development: Charles Severance, University of Michigan School of Information

