

# گوگل چگونه کار می کند؟

«مقدمه‌ای بر پیوندکاوی»

نویسنده

سید جمال میرکمالی

گروه دانش آماری - مهر ۱۳۸۸



## ایده‌ی گوگل

وب‌گردها در هنگام گشت و گذار در اینترنت معمولاً از موتورهای جستجو برای یافتن صفحات وب مرتبط با موضوع مورد علاقه‌ی خود استفاده می‌کنند. متأسفانه (!) اغلب هزاران صفحه‌ی وب در ارتباط با این پرس و جوها<sup>۱</sup> وجود دارد. بنابراین وجود یک فهرست از صفحات وب که بر اساس اهمیت مرتب شده‌اند ضروری به نظر می‌رسد. این فهرست بایستی به طور منظم و مکرراً روزآمد<sup>۲</sup> شود. بنابراین یافتن الگوریتم‌های سریع برای محاسبه‌ی رتبه‌ی صفحه بگونه‌ای که تاخیر در روزآمدسازی را کاهش دهد از اهمیت بسزایی برخوردار است. به نظر می‌رسد که این مساله‌ی دشواری است. نه فقط به دلیل اینکه تعداد صفحات وب در اینترنت بسیار زیاد است بلکه تعداد آن‌ها به سرعت در حال افزایش است.

رتبه‌ی صفحه<sup>۳</sup> توسط پیج و همکاران (۱۹۹۸) برای بیان اهمیت هر صفحه‌ی وب پیشنهاد شد (اینجا را ببینید). لری پیج و سرگی برین پایه‌گذاران گوگل هستند. در واقع، می‌توان عبارتهای زیر را در سایت گوگل (در اینجا) مشاهده کرد: "قلب نرم‌افزار ما رتبه‌ی صفحه است؛ سامانه‌ای که برای رتبه‌بندی صفحات وب توسط موسسین ما یعنی لری پیج و سرگی برین در دانشگاه استنفورد تولید شده است. در حالی که مهندسان بسیاری به طور روزمره جنبه‌های مختلف گوگل را بهبود می‌بخشند با این وجود رتبه‌ی صفحه به عنوان مبنای همه‌ی ابزارهای جستجوی ما همچنان به کار خود ادامه می‌دهد."

Query<sup>1</sup>  
Update<sup>2</sup>  
PageRank<sup>3</sup>



ایده‌ی مشابهی در خصوص رتبه‌بندی مجله‌ها توسط گارفیلد (۱۹۵۵ و ۱۹۷۲) به عنوان درجه‌ی اعتبار مجله‌ها تحت عنوان ضریب تاثیرگذاری<sup>۱</sup> معرفی شد. ضریب تاثیرگذاری یک مجله از تقسیم تعداد ارجاع‌ها به آن مجله بر تعداد مقالات آن مجله به دست می‌آید. در واقع ضریب تاثیرگذاری متوسط تعداد ارجاع‌ها به یک مقاله از یک مجله است. با در نظر گرفتن هر صفحه‌ی وب به عنوان یک مجله، این ایده برای اندازه‌گیری اهمیت صفحات وب در الگوریتم رتبه‌ی صفحه گسترش داده شده است.

رتبه‌ی صفحه به صورت زیر تعریف می‌شود. فرض کنید  $N$  تعداد کل صفحات وب موجود در اینترنت باشد و ماتریس  $Q$  را به نام ماتریس فرایوند<sup>۲</sup> به صورت زیر تعریف کنید:

$$Q_{ij} = \begin{cases} 1/k_i & \text{اگر صفحه‌ی } i \text{ یک پیوند بیرونی صفحه‌ی } j \text{ باشد} \\ 0 & \text{در غیر این صورت} \end{cases}$$

و  $k_i$  تعداد کل پیوندهای بیرونی صفحه‌ی  $i$  است. برای سادگی در اینجا فرض می‌کنیم که برای هر  $i$  داشته باشیم  $Q_{ii} > 0$ . این یعنی هر صفحه‌ی وب یک پیوند به خودش دارد. بنابراین  $Q$  را می‌توان به عنوان یک ماتریس احتمال انتقال یک زنجیر مارکوف در یک قدم تصادفی در نظر گرفت. مساله شبیه به این است که وب‌گرد را به عنوان شخصی که به طور تصادفی قدم برمی‌دارد و صفحات وب را به عنوان وضعیت‌های یک زنجیر مارکوف در نظر بگیریم. با فرض اینکه زنجیر مارکوف ناشکننده<sup>۳</sup> باشد توزیع احتمال ایستای

$$(p_1, p_2, \dots, p_N)^T$$

Impact Factor<sup>1</sup>  
Hyperlink<sup>2</sup>  
Irreducible<sup>3</sup>



برای وضعیت‌ها (صفحات وب) وجود دارد. در اینجا  $p_i$  نسبتی از زمان است که شخصی که به طور تصادفی قدم برمی دارد (وب گرد) وضعیت (صفحه‌ی وب)  $i$  را ملاقات می کند. هر چه مقدار  $p_i$  بیشتر باشد به معنی اهمیت صفحه‌ی وب  $i$  است. بنابراین رتبه‌ی صفحه‌ی  $i$  همان  $p_i$  تعریف می شود.

### یک مثال

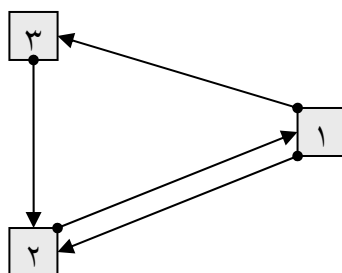
یک شبکه‌ی اینترنت شامل ۳ صفحه‌ی وب در نظر بگیرید: ۱، ۲، ۳ بگونه‌ای که

$$1 \rightarrow 1, 1 \rightarrow 2, 1 \rightarrow 3$$

$$2 \rightarrow 1, 2 \rightarrow 2,$$

$$3 \rightarrow 2, 3 \rightarrow 3.$$

این روابط را می توان به صورت زنجیر مارکوف زیر نمایش داد.



شکل ۱: یک مثال با سه صفحه‌ی وب

ماتریس احتمال انتقال زنجیر مارکوف به صورت زیر به دست می آید:

$$Q = \begin{matrix} 1 & \begin{pmatrix} 1/3 & 1/2 & 0 \\ 1/3 & 1/2 & 1/2 \\ 1/3 & 0 & 1/2 \end{pmatrix} \\ 2 \\ 3 \end{matrix}$$

توزیع احتمال ایستای زنجیر مارکوف  $\mathbf{p} = (p_1, p_2, p_3)$  در شرط  $\mathbf{p} = Q\mathbf{p}$  و

$$p_1 + p_2 + p_3 = 1$$

صدق می کند. با حل این معادلات داریم:



$$p = \left(\frac{3}{9}, \frac{4}{9}, \frac{2}{9}\right).$$

بنابراین ترتیب اهمیت صفحات وب به صورت زیر است:

$$\text{صفحه‌ی وب ۳} > \text{صفحه‌ی وب ۱} > \text{صفحه‌ی وب ۲}$$

از آنجا که زنجیر مارکوف بسیار بزرگ است و مدت زمان محاسبه‌ی رتبه‌ی صفحه توسط گوگل فقط چند روز است، روش مستقیم محاسبه‌ی توزیع احتمال ایستا مطلوب نیست. به همین دلیل روش‌های متنوعی بر اساس تجزیه و تکرار ارائه شده‌اند. مطالعات گسترده‌ای در این خصوص تحت عنوان *بهینه‌سازی موتورهای جستجو*<sup>۱</sup> (SEO) شکل گرفته است که نتیجه‌ی آن معرفی شیوه‌های تطبیقی، الگوریتم‌های موازی و روش هیبرید بوده است.

از آنجا که ماتریس  $Q$  ممکن است شکننده باشد محاسبه‌ی توزیع احتمال ایستا با مشکل

مواجه است. در این صورت می‌توان از ماتریس  $P$

$$P = \alpha Q_{N \times N} + \frac{(1 - \alpha)}{N} \mathbf{1}_{N \times N}$$

استفاده کرد که در آن  $0 < \alpha < 1$  و  $\mathbf{1}_{N \times N}$  یک ماتریس  $N \times N$  با مقادیر ۱ است. در این صورت ماتریس  $P$  ناشکننده و نادره‌ای<sup>۲</sup> است و توزیع احتمال ایستا وجود دارد و یکتاست (راس ۲۰۰۰). مقادیری که برای  $\alpha$  در نظر گرفته می‌شود معمولاً ۰/۸۵ یا  $1 - 1/N$  است.

تعبیری که برای ماتریس  $P$  می‌توان ارائه کرد این است که در یک شبکه شامل  $N$

صفحه‌ی وب، هر صفحه دارای اهمیت ذاتی به اندازه‌ی  $(1 - \alpha)/N$  است. اگر صفحه‌ی وب  $i$

Search Engine Optimization<sup>1</sup>  
Aperiodic<sup>2</sup>

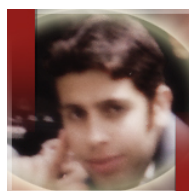


دارای اهمیت  $p_i$  باشد آنگاه به اندازه‌ی  $\alpha p_i$  در اهمیت صفحات وبی که به آن‌ها پیوند داده است شریک خواهد بود. اهمیت صفحه‌ی وب  $i$  از حل معادله‌ی زیر به دست می‌آید:

$$\mathbf{p} = \alpha Q\mathbf{p} + \frac{(1-\alpha)}{N} \mathbf{1}_{N \times 1}$$

برای توضیحات بیشتر در خصوص الگوریتم‌های محاسباتی به وایکی و میخائیل (۲۰۰۶) یا لنگویل و می‌یر (۲۰۰۶) مراجعه کنید.

در آخر به این نکته توجه کنید که در سامانه‌ی اندازه‌گیری اهمیت صفحات وب گوگل هیچ‌گونه هوشمندی از نوع تفسیر محتوا صورت نگرفته است و آنچه که موجب محبوبیت موتور جستجوی گوگل شده است در واقع استفاده از هوش نهفته‌ای است که کاربران اینترنت در پیوندها جاسازی می‌کنند. قاعدتا این پایان راه گوگل نخواهد بود و ممکن است روزی برسد که بتوان در گوگل یک قطعه‌ی آواز، یک فیلم یا یک عکس شبیه



را جستجو کرد و گوگل شما را مثلاً به اینجا راهنمایی کند.

### مرجع‌ها:

<p>Garfield E (1955) <i>Citation Indexes for Science: A New Dimension in Documentation Through Association of Ideas</i>, Science, 122:108–111.</p> <p>Garfield E (1972) <i>Citation Analysis as a Tool in Journal Evaluation</i>, Science, 178:471–479.</p> <p>Langville A N, Meyer C D (2006) <i>Google's PageRank and Beyond: The Science of Search Engine Rankings</i>, Princeton University.</p> <p>Page L, Brin S, Motwani R and Winograd T (1998) <i>The PageRank Citation Ranking: Bring Order to the Web</i>, Technical Report, Stanford University.</p> <p>Ross S (2000) <i>Introduction to Probability Models, 7th Edition</i>, Academic Press.</p> <p>Wai-Ki C, Michael K (2006) <i>Markov Chains Models, Algorithms and Applications</i>, Springer.</p>
---