

خزنده و ساختواره وب

شعله ارسطوپور¹

چکیده

وب به عنوان بستر فعالیت موتورهای جستجو، ساختاری نموداری دارد. این ساختار حرکت خزنده‌ها در موتورهای جستجو را به روشهایی منطبق بر خود محدود می‌سازد. مقاله حاضر، به بررسی تأثیر ساختار وب بر چگونگی حرکت خزنده‌ها و فعالیت نمایه‌سازها در موتورهای جستجو می‌پردازد. پس از بحثی مقدماتی در باب نمودارهای جهت‌دار و فرایند کار خزنده، عمده‌ترین روشهای حرکت خزنده در سطح وب شامل حرکت «عمق - شروع»، «توزیع - شروع» و «بهترین - شروع» مطرح شده و سپس واحد سازه‌یابی و چگونگی تشکیل درختهای سازه‌یابی از قالب HTML مورد بررسی قرار خواهد گرفت.

کلیدواژه‌ها: ساختار نموداری وب، خزنده‌ها، حرکت عمق - شروع، حرکت توزیع - شروع، حرکت بهترین - شروع، نمایه‌سازی وب

مقدمه

وب، مجموعه‌ای عظیم از مدارک است که هر یک برای پاسخگویی به یک نیاز بالقوه ایجاد و در بستر اینترنت منتشر گردیده است. حجم عظیم صفحات و اطلاعات موجود در وب و لزوم وجود ابزارهایی جهت سازماندهی دست کم گوشه‌ای از این مجموعه تاکنون به حد کفایت مورد بحث قرار گرفته است و اینکه موتورهای جستجو با تمام کاستیهای خود از عمده‌ترین ابزارهای دسترسی به محتوای مدارک پیش گفته هستند، اصلی پذیرفته شده میان بسیاری از کاربران اینترنتی است. وب، میلیاردها صفحه الکترونیکی را از طریق شبکه‌ای از اتصالها با یکدیگر مرتبط می‌سازد و از آنجا که ساختار آن هیچ محدودیتی بر چگونگی نشر صفحات، قالب انتشار، و یا چگونگی برقراری اتصالها و تعداد صفحات وب اعمال نمی‌نماید، گسترش چشمگیری یافته است. اما این مسئله هرگز بدین معنا نیست که وب بدون ساختار است. بسیاری از پژوهشها بیانگر وجود روابطی مشخص و منطقی در ساختار وب هستند.

(Barabasi & Albert, 1999; Kleinberg et al., 1999; Albert et al., 1999)

هر یک از این تحقیقات به نوبه خود به نتایج متفاوتی رسیده اند. فصل مشترک آنها مسئله ناهماهنگی توزیع صفحات وب است به گونه‌ای که آلبرت و دیگران² (1999) از وب به عنوان یک شبکه جهانی کوچک یاد می‌کنند در حالی که برادر و همکاران³ (2000) توزیع وب را پاپیونی شکل توصیف می‌کنند. بر

۱. دانشجوی کارشناسی ارشد کتابداری و اطلاع‌رسانی دانشگاه فردوسی مشهد

1. Albert et al.
2. Broder et al.

این اساس، وب از یک هسته که بیشترین اتصالات به آن برقرار شده و دو بال که با توجه به جهت اتصالات، توزیع متفاوتی از صفحات را در بر می‌گیرد، تشکیل شده است. کلیه پژوهش‌های انجام شده و تمام نتایج یک مسئله را اساس کار خود قرار داده‌اند و آن نموداری بودن ساختار وب⁴ است. با توجه به دستاوردهای پژوهشها، ساختار وب را به یقین می‌توان به صورت یک نمودار عظیم ترسیم نمود؛ نموداری که در برگیرنده گره‌ها و اتصالات متعدد است. در این نمودار، صفحات وب همان گره‌ها و اتصالات همان لینکهای برقرار شده میان صفحات مختلف است

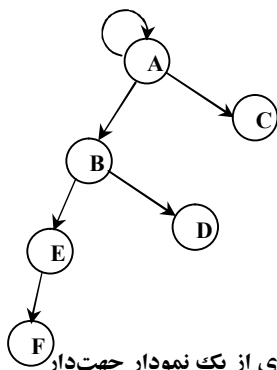
(Albert et al., 1999; Thelwall, 2002; Evans & Walker, 2004; Yu & Meng, 2004; Coathy, 2004).

بنابراین وب نه تنها بدون ساختار نیست بلکه می‌توان آن را به صورت یک ساختار نموداری نیز ترسیم نمود. این ساختار خواسته یا ناخواسته بر چگونگی گردآوری صفحات و نمایه‌سازی آنها توسط موتورهای جستجو تأثیر گذار بوده و در نهایت نتایج جستجو و میزان ربط آنها را رقم می‌زند.

نمودارهای جهت‌دار⁵

نمودار در واقع مجموعه‌ای از گره‌ها و خطوط است که آن را به صورت ریاضی $G(V, E)$ نشان می‌دهند. هر E حتماً دو گره را به یکدیگر متصل می‌کند. نمودار جهت‌دار، نموداری است که بتوان در آن جهت حرکت از هر گره به گره دیگر را به راحتی تشخیص داد.

در یک نمودار جهت‌دار، چنانچه تعداد اتصالات را با i نمایش دهیم، برای هر گره که به عنوان گره اصلی در نظر گرفته شود، رابطه $i \in [1, n]$ برقرار خواهد بود. در این شرایط، هر یک از اتصالات تنها به طرف یک گره حرکت خواهد کرد. از دیگر ویژگیهای عمده یک نمودار جهت‌دار آن است که هیچ کدام از گره‌ها منفرد و بدون اتصال نخواهد بود. (تصویر 1).



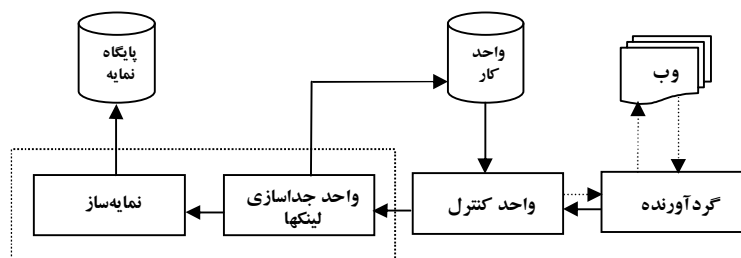
تصویر 1. نمونه‌ای از یک نمودار جهت‌دار

بیشتر بیان شد که هر صفحه از وب دارای تعدادی اتصال است. این اتصالات یا از صفحه هسته⁶ به سایر صفحات و یا از سایر صفحات به صفحه مورد نظر برقرار شده است. بنابراین، تمام ویژگیهای پیش گفته، در باب صفحات وب و لینکهای آنها نیز صدق می‌کند. با توجه به جهت برقراری لینکها، می‌توان وب را نیز به صورت یک نمودار جهت‌دار ترسیم نمود.

خزنده

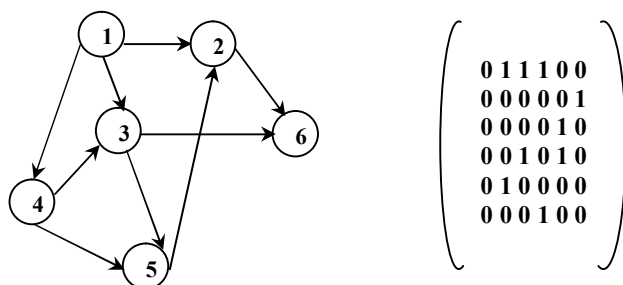
تقریباً تمام خزنده‌ها در موتورهای جستجو دارای چهار بخش گردآورنده⁷، واحد کنترل⁸، واحد سازه‌یابی⁹ و واحد کار¹⁰ هستند. گردآورنده تحت نظارت واحد کنترل، به صفحات هسته (گره‌های مختلف) رفته و مدارک را به واحد جداسازی لینکها می‌فرستد. پس از جداسازی، لینکهای مناسب به واحد کار ارسال شده و در فهرست دستور کار بعدی گردآورنده قرار می‌گیرند. واحد سازه‌یابی در واقع از دو بخش جداسازی لینکها¹¹ و نمایه‌سازی تشکیل شده است.

آنچه نهایتاً در پایگاه ذخیره می‌شود در واقع حاصل فرایند نمایه‌سازی است که تحت قالب تعریف شده در الگوریتم موتور جستجو به صورت واژگان و عبارات مختلف در آمده است (Cothey, 2004). تصویر 2 ساختار یک خزنده را نشان می‌دهد.



تصویر 2. ساختار خزنده و چهار جزء اصلی آن

بیشتر اشاره شد که ساختار وب همچون نموداری جهت‌دار است، لذا خزنده نیز برای حرکت روی این ساختواره، چاره‌ای جز تبعیت از ویژگیها و جهت‌های از پیش تعیین شده ندارد. از نظر جبری، حرکت روی نمودارهای جهت‌دار را می‌توان به صورت ماتریسهای حرکتی¹² ساده کرد. در موتورهای جستجو نیز، با توجه به شباهت ساختاری، استفاده از قواعد ریاضی حاکم بر این ماتریسها، تحلیل و مقایسه را آسان می‌نماید. (تصویر 3)



تصویر 3. ساختواره نموداری وب و ماتریس حرکتی متناظر با جهت‌های مشخص شده

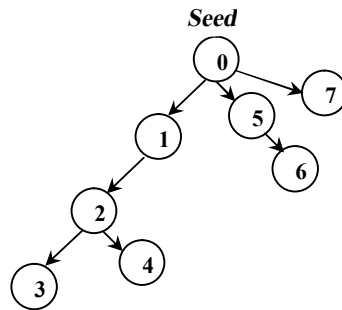
به طور کلی، سه روش برای حرکت خزنده در شبکه لینکهای وب وجود دارد. این سه عبارتند از

حرکت عمق - شروع¹³، توزیع - شروع¹⁴، و بهترین - شروع¹⁵.

1. Fetcher
2. Controller
3. Parsing Unit
4. Workload Unit
5. Link Extracting
6. Traversal Matrix
1. Depth - First
2. Breadth - First

حرکت عمق - شروع

در این حرکت، واحد کنترل خزنده یک صفحه را به عنوان صفحه هسته برای گردآورنده مشخص می‌سازد. پس از جداسازی لینکها، واحد کنترل یکی از لینکهای خارجی صفحه را انتخاب و گره مقصد را به گردآورنده معرفی می‌کند. این فرایند تا زمانی که برای واحد کنترل تعریف شده باشد، ادامه پیدا می‌کند. به عنوان نمونه، ترتیب حرکت گردآورنده به صفحه‌های مختلف با الگوریتم عمق - شروع، مانند تصویر 4 خواهد بود.

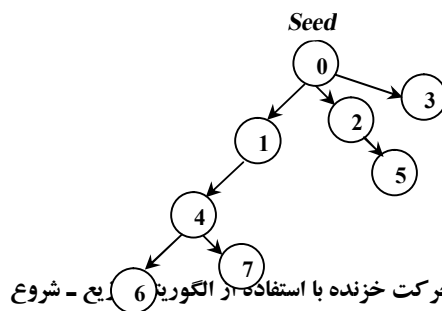


تصویر 4. صفحه هسته و ترتیب حرکت خزنده با استفاده از الگوریتم عمق - شروع

از آنجا که تقریباً تمام صفحه‌های وب لینکهایی به سایر صفحات برقرار می‌کنند، چنانچه سطح عمق برای واحد کنترل تعریف شده نباشد، حرکت به عمق آن قدر ادامه خواهد یافت که به مرور زمان، عملاً پایگاه نمایه موتور جستجو را از مطالب بی کیفیت خواهد انباشت. به همین دلیل، در بیشتر موتورهای جستجو، سطح عمق برای واحد کنترل تعریف می‌شود. در تصویر 4، چنانچه سطح عمق تعریف شده 2 باشد، ترتیب حرکت گردآورنده 1، 2، 5، 6 و 7 بوده و عملاً صفحه‌های 3 و 4 مورد بررسی قرار نخواهند گرفت.

حرکت توزیع - شروع

در این حرکت، واحد کنترل پس از تعیین صفحه هسته، کلیه گره‌های هم عمق با یکدیگر را تعیین و به ترتیب به گردآورنده معرفی می‌کند. پس از رجوع به کلیه صفحات مشخص شده در آن سطح، واحد کنترل سطح دوم را مورد بررسی قرار می‌دهد. به عنوان نمونه، ترتیب حرکت گردآورنده تحت نظارت واحد کنترل و با استفاده از الگوریتم توزیع - شروع در صفحه‌های مختلف مانند تصویر شماره 5 خواهد بود.



تصویر 5. صفحه هسته و ترتیب حرکت خزنده با استفاده از الگوریتم توزیع - شروع

روش حرکت توزیع - شروع مورد علاقه بسیاری از طراحان برنامه‌های خزنده در موتورهای جستجو است، زیرا طراحی و اجرای آن به صورت رایانه‌ای بسیار ساده‌تر از روش حرکت عمق - شروع بوده و در صورت تعیین سیاست دقیق، به لحاظ محدود بودن دامنه لینکهای هر صفحه به عنوان صفحه هسته، حجم پایگاه موتور جستجو بهبود یافته خواهد یافت (Chakrabarti et al., 2002).

حال، چنانچه صفحه هسته به یک مطلب خاص پردازد، با توجه به آنکه گردآورنده تمام لینکهای موجود در صفحه و یا صفحات بعد را دنبال نمی‌کند، حرکت خزنده تأثیر بسیار زیادی بر نمایه‌سازی و در نهایت بازیابی اطلاعات خواهد داشت.

در حرکت عمق - شروع، با انتخاب هر لینک و رفتن به صفحه بعدی و ادامه این کار، یک مطلب خاص (حوزه موضوعی مربوط به سطح عمق اول حرکت) به صورت اختصاصی دنبال شده و از آنجا که گرایش واحد کنترل نسبت به حرکت عمقی گردآورنده بیشتر از حرکت در سطح است، در نهایت صفحاتی که برای نمایه‌سازی فرستاده می‌شوند به احتمال، اغلب حول یک مطلب یا موضوع خواهند بود. در حالی که در حرکت توزیع - شروع گرایش واحد کنترل به حرکت در سطح است و لذا گردآورنده ابتدا به گره‌های تعیین شده در سطح سرکشی خواهد کرد. در چنین شرایطی صفحاتی برای واحد نمایه‌سازی ارسال می‌شوند که در کل دیدی عامتر دارند (مانند آنچه در صفحه هسته آمده است). این مسئله، ناشی از آن است که معمولاً لینکهایی که از هر صفحه برقرار می‌شوند، به بخشی از مطلب مطرح شده در صفحه هسته مربوط می‌شوند.

بنابراین، با حرکت از سطح به عمق، دید نمایه‌سازی خواه نا خواه جزء نگر بوده و با حرکت در سطح دید نمایه‌سازی در حول یک مطلب با گستره‌ای وسیع‌تر و جامع‌تر متمرکز خواهد بود (Herrman, 2003). با توجه به مطالعات دهه 1990، استفاده از هر یک از این روشها برای حرکت روی ساختار وب با توجه به حجم عظیم صفحات و مطالب هر یک، در نهایت تفاوت چندانی در بازیابی بهتر موتورهای جستجوی مختلف در سطح وب نداشت.

تصمیم‌گیری در باب انتخاب هر یک از این دو روش و تردیدهای موجود، روش سوم را پیش پای طراحان و برنامه‌نویسان قرار داد و آن، حرکت «بهترین - شروع» بود.

حرکت بهترین - شروع

بهترین در حوزه حرکت خزنده روی ساختار وب، در واقع معانی متفاوتی دارد. الگوریتمهای مختلفی برای حرکت بهترین - شروع وجود دارند که بر اساس فرمول محاسبه بهترین صفحه بعدی، اسامی متفاوت دارند. از این بین می‌توان به خزنده متمرکز¹⁶، جستجوی کوسه‌ای¹⁷، عنکبوتهای اطلاعاتی¹⁸ و... اشاره کرد. در ساده‌ترین حالت، از سیاستهای رتبه‌بندی همچون "رتبه‌بندی صفحات"¹⁹ به عنوان معیار بهترین بودن استفاده می‌شود. در این روش واحد کنترل با توجه به رتبه هر صفحه میان سایر صفحات، گردآورنده را به صفحه بعدی می‌فرستد.

1. Focused Crawler
2. Shark Search
3. Info Spiders
4. Page Rank

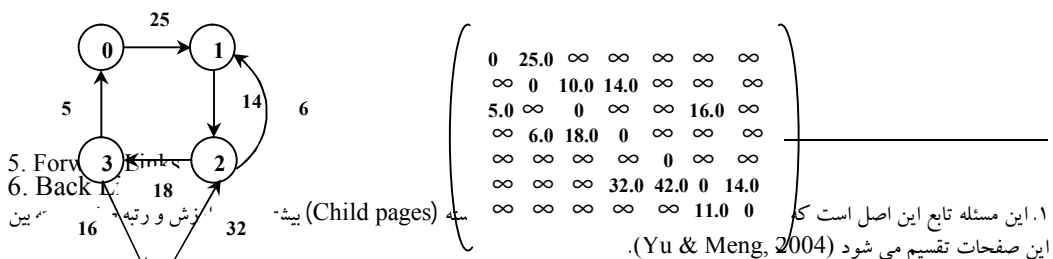
حرکت بهترین - شروع، بر این اصل مبتنی است که پدیدآورنده هر صفحه زمانی از صفحه خود (A) به صفحه دیگری (B) لینک برقرار می کند که B از نظر پدیدآور A ارزشمند باشد. در بحث رتبه بندی صفحات چنین عملی برای B یک امتیاز مثبت محسوب می شود. بنابراین، در ساختار خزنده ها الگوریتم رتبه بندی صفحات در واقع برنامه ای است که اهمیت نسبی هر صفحه را بر اساس لینکهای برقرار شده به آن مشخص می سازد. بر این پایه در خزنده هایی که برای حرکت خود از روش بهترین - شروع استفاده می کنند در واقع از سه اصل پیروی می شود:

- صفحاتی که لینکهای بیشتری به آنها برقرار شده است، اهمیت بیشتری دارند. تعداد بیشتر لینکها به نوعی نشانگر شهرت و یا محبوبیت صفحه مذکور در سطح وب است.
 - چنانچه این لینکها از صفحات معتبرتری برقرار شده باشند، اعتبار صفحه مورد مطالعه افزایش خواهد یافت.
 - از طرفی، هرچه تعداد لینکهایی که از صفحه مورد مطالعه به سایر صفحات برقرار می شود بیشتر باشد، ارزش آن صفحات کمتر خواهد بود.
- بر این اساس، در اکثر خزنده ها چنانچه U صفحه اصلی، $F_{U|}$ نشانگر صفحاتی که از U به آنها لینک برقرار شده ²⁰ و $B_{U|}$ نشانگر صفحاتی باشد که به U لینک برقرار کرده اند ²¹ رتبه صفحه از طریق فرمول حاضر و یا فرمولهایی با عناصر اصلی مشابه، قابل سنجش خواهد بود.

$$R(U) = \sum_{V \in B_u} \frac{R(V)}{|F_u|}$$

در این فرمول، علامت جمع به معنای وجود رابطه مثبت میان تعداد لینکهای برقرار شده به U و رتبه U است. $R(V)$ یا رتبه صفحات برقرار کننده لینک به U نیز چون در صورت کسر قرار گرفته اند، رابطه مثبت با رتبه U دارند، در حالی که وجود $|F_u|$ در مخرج، نشانگر وجود رابطه معکوس رتبه U و تعداد لینکهایی است که به صفحات دیگر برقرار کرده است ²².

استفاده از نتیجه این فرمول در گردآوری صفحات ارزشمندتر، مؤثر خواهد بود. با کمک این روش و استفاده از رتبه بندی محتوای مدارک در واحد نمایه سازی به احتمال پاسخهای بهتری به نیاز کاربر موتور جستجو داده خواهد شد. با رتبه بندی صفحات، در واقع رتبه لینکهای منتهی به آنها مشخص شده و خود به خود ترتیب حرکت گردآورنده بر روی ساختواره نموداری عظیم وب روشن می گردد. تصویر 6 نشانگر این ترتیب بر اساس رتبه های مشخص شده در ماتریس حرکتی خواهد بود.



تصویر 6. لیستها و ارزش عددی هر یک و نمایش آنها در قالب ماتریس حرکتی

این مسئله در نظر کاملاً منطقی است، اما ناجورک و وینر²³ (2001) در عمل ثابت می‌کنند که با توجه به هزینه بالا و زمان بر بودن فرایند رتبه‌بندی صفحات، استفاده از روش توزیع - شروع به نسبت توجیه پذیرتر می‌نماید. آنها بر مبنای دستاوردهای خود بیان می‌دارند که با توجه به امکانات فعلی فناوری، روش رتبه‌بندی صفحات بسیار هزینه بر بوده، زمان زیادی را می‌طلبد و از طرف دیگر با توجه به عدم ثبات رتبه صفحات در طول زمان بایگانی نگهداری رتبه‌های صفحات را به سرعت باید روزآمدسازی نمود. این در حالی است که توجه به حرکت عمق - شروع به عنوان یک گزینه مطرح، کمتر صورت می‌پذیرد.

واحد سازه‌یابی

واحد نمایه ساز موتورهای جستجو باید صفحات حاوی اطلاعات را از گردآورنده دریافت کند و عبارات و واژگان آنها را استخراج و در پایگاه خود ذخیره سازی نمایند. بنابراین، چنانچه هر صفحه در مقام خود یک واحد کلی باشد، نمایه ساز آن را به اجزای کوچکتر از قبیل واژه و یا عبارت تبدیل کرده و در پایگاه خود ذخیره می‌سازد. نرم‌افزاری که توانایی انجام این عمل را داشته باشد، «سازه یاب» نام دارد. در فرایند سازه‌یابی، اولین کار تشخیص زبان رشته نشانه‌های ورودی²⁴ است. پس از آن، بر اساس دستور آن زبان خاص، سازه‌یاب به تعیین ساختار ترکیبی آن رشته می‌پردازد.²⁵

با این توصیف برنامه سازه‌یاب در ساختار یک موتور جستجو کار جداسازی و یکسان‌سازی آدرسهای اینترنتی موجود در مدرک، نگهداری فهرست واژگان غیر مجاز و تهیه درخت سازه‌یابی را انجام می‌دهد. از آنجا که سازه‌یاب بر اساس دستور زبان از قبل تعریف شده به هدف دستیابی به محتوای مشخصی عمل می‌کند، تقسیم‌بندی واژگان استخراج شده و وزن دهی به آنها کار ساده‌ای خواهد بود (Fischer, 2005)

زبان HTML زبان غالب در سطح وب به شمار می‌آید، لذا کلیه موتورهای جستجو دارای نرم افزارهای سازه‌یابی سازگار با HTML برای زبانهای مختلف هستند. به واسطه این نرم افزارها، برچسبهای HTML و ارزش آنها به سرعت شناسایی می‌شوند.

برای جداسازی و یکسان نمودن آدرسهای موجود در یک صفحه، از سازه‌یابها به منظور شناسایی برچسبهای مختلف و ارزش آنها استفاده می‌شود. این کار معمولاً به منظور کمک به واحد کنترل جهت هدایت گردآورنده انجام می‌شود. اما کار معرفی آدرسها به گردآورنده، اغلب مشکل‌تر از این است. گاه لازم است بسیاری از آدرسهای ذکر شده در صفحه یکدست و تصحیح شوند. به منظور یکسان سازی آدرسهای اینترنتی

1. Najork & Wiener

2. Input Symbols' String

۳. پاکروان، امیرحسین (۱۳۷۶)، فرهنگ کامپیوتر یادواره انگلیسی - فارسی، (تهران: یادواره اسدی؛ فرهنگستان یادواره).

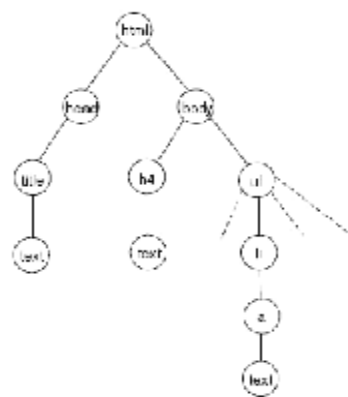
دستورالعملهایی برای تبدیل حروف بزرگ آدرسها به حروف کوچک، برداشتن انشعابهای اضافی از دنباله آدرس، تصحیح و یا تکمیل برخی از آدرسها و ...، برای سازه یابهای مختلف تعریف می شود (Pant, Sirinivan & Meczer, 2004).

سازه یابهای مختلف ممکن است سیاهه واژگان غیر مجاز متفاوتی داشته باشند و یا اصولاً فاقد ویژگی حذف واژگان بدون بار معنایی در طول فرایند نمایه سازی باشند. در سطوح بالاتر، برخی از سازه یابها با توجه به دستور زبان از پیش تعریف شده، توانایی تشخیص ریشه کلمات و ذخیره کلیه واژگان هم ریشه را در یک محل دارند.

در نهایت، وظیفه هر سازه یاب ایجاد درخت سازه یابی²⁶ است. در این مرحله، واحد سازه یابی آدرس و یا واژه موجود در صفحه را با کمک محتوا و محل برچسب ارزیابی کرده، درختواره ای از ساختار صفحه تشکیل می دهد. نمونه ای از این درخت و قالب HTML متناظر با آن، در تصویر 7 نشان داده شده است. گره اول در این نمودار نشانگر قالب مدرک، گره های میانی برچسبهای مختلف مدرک و آخرین گره ها نماینده محتوای میان دو برچسب ابتدایی و انتهایی است.

تشکیل این درخت، کار وزن دهی به هر متن و واژه و عبارت استخراج شده از آن را ساده می نماید. رتبه متون در قسمتهای مختلف صفحه با توجه به الگوریتم رتبه بندی خاص هر موتور متفاوت است. واژگان وزن دهی شده را می توان به راحتی در قالب یک مقیاس عددی ریخته و برای الگوریتم رتبه بندی موتور جستجو، امکان سنجش و مقایسه سؤال کاربر و واژگان موجود پایگاه را فراهم آورد.

```
<html>
<head>
<title>Projects</title>
</head>
<body>
<h1>Projects</h1>
<ul>
<li><a href="blink.html">LAMP</a> Linkage analysis with multiple processors.</li>
<li><a href="nice.html">NICE</a> The network infrastructure for combinatorial exploration.</li>
<li><a href="amass.html">AMASS</a> A DNA sequence assembly algorithm.</li>
<li><a href="dal.html">DAL</a> A distributed, adaptive, first-order logic theorem prover.</li>
</ul>
</body>
</html>
```



تصویر 7. درخت سازه یابی و قالب HTML متناظر با آن

جمع بندی

با وجود ماهیت متغیر وب، باز هم ساختار وب بر نحوه سازماندهی آن تأثیرگذار خواهد بود. با توجه به اینکه هم اکنون موتورهای جستجو از مهمترین سازمان دهندگان به شمار می روند و از طرفی با در نظر گرفتن اینکه ساختار وب روشهای متفاوت و نه متعدد گردآوری اطلاعات را به طراحان نرم افزارهای خزنده دیکته

می‌کند، می‌توان بیان داشت که ساختواره جهت‌دار وب بی‌تأثیر بر بازیابی‌های مفید و یا بی‌تأثیر خواهد داشت. آنچه در این زمان مورد توجه بیشتر محققان حوزه قرار گرفته، چگونگی بهینه‌سازی استفاده از امکاناتی است که وب در اختیار طراحان قرار می‌دهد. تصمیم‌گیری در باب انتخاب شیوه حرکت خزننده - اعم از حرکت به عمق یا حرکت در سطح و یا انتخاب هر صفحه بسته به کیفیت آن - یکی از مباحث مورد توجه علاقه‌مندان به این حوزه است. مطالعه در باب بازدهی هر روش در طول زمان و یا امکان‌سنجی استفاده از یک روش در حال حاضر از مطالعات مطرح در این حوزه به شمار می‌آید.

این در حالی است که از زاویه‌ای دیگر، اعمال از پیش تعریف شده برای هر سازه‌یاب - چه کوتاه و مختصر (سازه‌یاب‌های ساده) و چه پیچیده (سازه‌یابهای سطح بالا) - نیز به بهینه‌سازی استفاده از امکانات وب توجه می‌کند. با تعریف جزئیات بیشتر، نمایه‌سازی دقیق‌تر شده و در نهایت بازیابی بهتری حاصل خواهد شد. مدارک به واسطه‌ی شیوه حرکت در سطح وب گردآوری شده‌اند و نمایه‌سازی روی ساختار ساختمند وب به اجرا درآمده است؛ ساختاری که حتی زبان نگارش آن را می‌توان در قالب نمودار ترسیم نمود. الگوریتم‌های مختلف رتبه‌بندی بر اساس این ساختار و اجزای آن کار خود را انجام می‌دهند، پس ساختار وب به صورتی غیر مستقیم اما با قدرتی بسیار بر آنچه بازیابی می‌شود تأثیر خواهد داشت.

منابع

- Albert, R., Jeung, H. & Barabasi, A. (1999). "The Diameter of the World Wide Web". Nature . Vol. 401, P. 130 Available online: www10.org/cdrom/papers/208 [Accessed on Oct. 2005]

- Barabasi A.L. & Albert, R. (1999). "Emergence of Scaling in random Networks". Science. Vol. 286, P 509 - 512. Available online: <http://www.nd.edu/~networks> [Accessed on Oct. 2005]

- Broder, A., Kumar, R., Maghoul, F., Raghavan, P., Rajagopalan, s., Stata, R., (2000). "Graph Structure in the Web". Computer Networks. Vol. 33, P. 309 - 320. Available online: <http://www9.org/w9cdrom/160/160.html> [Accessed on Oct. 2005]

- Chakrabarti, s. , Joshi, M.M., Punera, K. & Pennock, D.M. (2002). "The Structure of Broad Topics On the Web". Proceedings of the 11th World Wide Web Conference (p.508 – 510). Honolulu, Hawaii, May 7- 11 . New York: ACM. Available online: <http://http.cs.berkeley.edu/~soumen/doc/www2002t/p338-chakrabarti.pdf> [Accessed on Oct. 2005]

- Cothey, Viv (2004). "Web Crawling reliability". Journal of the American Society for Information Science and Technology. 55(14). P. 1228 – 1238.

- Evans, Michael P. & Walker, Andrew (2004). "Using The Web Graph to Influence Application Behavior". Internet Research. 14(5). P. 372 – 378.

- Fischer, Hendrik (2005). Decisions To Go: An Intelligent Mobile Decision Support System[Dissertation]. Georgia: The University of Georgia. Available online: <http://graduate.gradsch.uga.edu/etdarchive/summer2005/fischer-hendrik-200508-ms.pdf> [Accessed on Oct. 2005]

- Herrmann, Frank (2003). Web search engines. Available online: <http://graduate.gradsch.uga.edu/etdarchive/summer2005/fischer-hendrik-200508-ms.pdf> [Accessed on Oct. 2005]

- Kleinberg, J., Kumar, R., Raghava, P., Rajagopalan, S., & Tomkins, A. (1999). "The Web as a Graph: Measurements, Models, and Methods". Proceedings of the International Conference on Combinatorics and Computing , Tokyo , Japan, July 26 – 28. London: Springer, P. 1 - 17. Available online: <http://www.tomkinshome.com/papers/archive/cocoon99.pdf> [Accessed on Oct. 2005]
- Najork,M. & Wiener, J.L. (2001). "Breadth – First Crawling Yields High Quality Pages". Proceedings of the 10th World Wide Web Conference (p.114 - 118). Hongkong. May 1 - 5 . New York: ACM. Available online: <http://www10.org/cdrom/papers/208/> [Accessed on Oct. 2005]
- Pant, G., Srinivasan, p. Menczer,F. (2004). "Crawling the Web". Web Dynamics. Springer. Availableonline: <http://mia.ece.uic.edu/~papers/MediaBot/pdf00001.pdf> [Accessed on Oct. 2005]
- Thelwal, Mike (2002). "Methodologies for Crawler Based Web Surveys". Internet Research. 12(2), P. 124 – 138. Available online: www.scms.rgu.ac.uk/staff/fh/CM1008/documents/lecture3.pdf [Accessed on Oct. 2005]
- Yu, Clement & Meng, Weiyi (2004). "Web Search Technology". The Internet Encyclopedia. Hoboken, NJ: Wiley. P 738 – 753.