

تعیین Bias موتورهای جستجوگر از Robots.txt

Yang Sun, Ziming Zhuang, Isaac G. Councill, and C. Lee Giles

Information Sciences and Technology
The Pennsylvania State University
University Park, PA 16801, USA
{ysun, zzhuang, icouncill, giles}@ist.psu.edu

مترجم : مسعود مجلسی

دانشگاه غیر انتفاعی اشراق بجنورد
Masoud.Majlesi@Gmail.com

خلاصه

موتورهای جستجوگر برای جمع آوری اطلاعات از وب بسیار به ربات‌ها متکی هستند. این چنین فعالیت‌های خزیدن (در وب) می‌تواند در طرف سرور با به کارگیری پروتکل‌های ممانعت از ربات‌ها در فایل‌ی به نام Robots.txt کنترل گردد. ربات‌های متعهد از قواعد مشخص شده در Robots.txt پیروی می‌کنند. وب سایت‌ها می‌توانند اولویت‌های دسترسی برای هر ربات را بوسیله نام تعیین کنند. این bias ممکن است منجر به قدرت‌مندتر شدن قدرتمندان گردد، به طوری که تعداد محدودی موتورهای جستجوگر محبوب در نهایت بر وب حکمفرما شوند به دلیل اینکه آن‌ها تقدم دسترسی به منابعی را دارند که دیگران به آن دسترسی ندارند. این موضوع به ندرت مورد توجه قرار می‌گیرد، اگرچه قرارداد Robots.txt به یک استاندارد غیر رسمی برای تنظیم ربات‌ها تبدیل شده و موتورهای جستجوگر یک ابزار واجب برای دسترسی به اطلاعات شده است. ما یک معیار برای ارزیابی مقدار Bias برای ربات‌های ویژه‌ای که گفته شد پیشنهاد می‌کنیم. ما ۷۵۹۳ وب سایت که شامل دومین‌های آموزشی، دولتی، خبری و تجاری را مورد بررسی قرار داده و ۲۹۲۵ فایل Robots.txt مجزا را جمع‌آوری کرده‌ایم. نتایج و تحلیل‌های آماری از اطلاعات ثبت شده این است که ربات‌های موتورهای جستجوگر محبوب و پرتال‌های اطلاعاتی توسط اکثر وب سایت‌هایی که ما به عنوان نمونه نام برده‌ایم مورد همکاری قرار گرفته است. همچنین نتایج نشان می‌دهد که یک رابط قوی میان Bias و سهم موتورهای جستجوگر از بازار بخصوص در مورد ربات‌های موتور جستجوگر وجود دارد.

۱. مقدمه :

شاید بدون ربات‌ها موتورهای جستجوگری نبود. موتورهای جستجو وب، کتابخانه‌های دیجیتال، و بسیاری دیگر از برنامه‌های کاربردی وب همچون مرورگرهای آفلاین، نرم‌افزارهای خرید اینترنتی و عاملین هوشمند جستجو برای بدست آوردن مستندات بسیار به ربات‌ها وابسته‌اند. ربات‌هایی که به نام های spider, crawler, bot, harvesters شناخته می‌شوند، عامل‌های خودکاری هستند برای هدایت هایپرلینک‌های وب، کنکاش در منابع موضوعی برای مدیریت انسانی مبدأ هزینه است. به دلیل ذات خودکار ربات‌ها، قواعد باید طوری تنظیم گردند که چنین فعالیت‌های کنکاشی مانع از اثر نامطلوب بر تراکم کاری سرورها و نیز دسترسی به اطلاعات غیر عمومی و شخصی گردد.

پروتکل‌های ممانعت از ربات‌ها قواعدی را برای پیروی ربات‌ها پیشنهاد کرده است. فایل‌ی به نام Robots.txt که شامل سیاست‌های دسترسی ربات‌ها می‌باشد در فهرست ریشه یک وب سایت به کار گرفته می‌شود و برای تمام ربات‌ها قابل دسترسی است. ربات‌های متعهد این فایل را می‌خوانند و در طول سرکشی به وب سایت قوانین آن را اجرا می‌کنند. قرارداد Robots.txt از اواخر دهه ۹۰ در مجامع پذیرفته شده و همچنان به عنوان ابزاری فراگیر و غالب برای تنظیمات ربات‌ها به کار می‌رود. هر چند با وجود مشخص بودن قرارداد Robots.txt برای هر دو تولید کنندگان محتوا و نیز کنکاشگران، کارهای کمی برای بررسی فواید آن با جزئیات انجام شده است بخصوص در مقیاس وب.

وب سایت‌ها ممکن است قواعد دسترسی مختلفی را به ربات‌ها اختصاص دهند، این Bias می‌تواند منجر به قدرتمند شدن قدرتمندان گردد که بوسیله آن بعضی از موتورهای جستجوگر محبوب دارای امتیاز انحصاری دسترسی به منابع معینی می‌شوند که در نتیجه می‌تواند آن‌ها را باز هم محبوب تر کند. نظر به این واقعیت که کاربران اغلب موتورهای جستجوگر با پوشش اطلاعات گسترده را ترجیح می‌دهند این پدیده (قدرتمندان قدرتمندتر می‌شوند) نشانگر تاثیر زیاد انتخاب موتور جستجوگر توسط کاربران است، که سرانجام به سهم بیشتر موتورهای جستجوگر در بازار می‌انجامد، به عبارت دیگر (اگر چه شاید اغراق باشد) "چیزی که قابل جستجو نیست پس وجود ندارد" این پدیده دید جهت‌گیرانه اطلاعات روی وب را نشان می‌دهد.

مطالعه‌ی استفاده از Robots.txt در دانشگاه‌های انگلیس و ۱۶۳ وب سایت و ۵۳ Robots.txt دانشگاهی بررسی شده در ۱۹۹۹. فایل‌های Robots.txt در زمینه حجم فایل و نیز پروتکل‌های ممانعت روبات در داخل دومین‌های دانشگاهی انگلیس مورد بررسی قرار گرفتند. در سال ۲۰۰۲ Drott استفاده‌ی Robots.txt را در کمک به شاخص‌گذاری (Indexing) اطلاعات محافظت شده را مورد مطالعه قرار داد. ۶۰ نمونه از وب‌سایت‌های ۵۰۰ شرکت جهانی که به صورت دستی مورد آزمون قرار گرفت که در این آزمایش این نتیجه حاصل شد: Robots.txt به طور زیاد در گروه مورد آزمایش استفاده نمی‌شود و در اکثر آن‌هایی که دیده می‌شود زائد است... آن‌ها روبات‌ها را از فهرست‌هایی محروم کرده‌اند که در هر صورت قفل شده است). بررسی ما نتایج متفاوتی را نشان می‌دهد که ممکن است ناشی از نمونه حجم، دومین و زمان متفاوت باشد. نمونه دیگر اشاره به جنبه حقوقی اجرای Robots.txt و نظر کلی در مورد روبات‌های وب و استفاده از Robots.txt که ضمیمه شده. هیچ یک از تحقیقات مذکور محتویات Robots.txt را بر حسب Bias نسبت به روبات‌های مختلف بررسی نکرده است. به علاوه نمونه مطالعات گذشته با تمایل نسبتاً کمی به حجم وب انجام شده است. در این مقاله ما برای اولین بار یک مطالعه کمی در مورد این Bias و یک تحقیق جامع‌تر در زمینه استفاده از Robots.txt در وب را ارائه می‌دهیم. با به کار گرفتن کنکاشگر مخصوص Robots.txt خودمان، مقدار قابل توجهی داده در دنیای واقعی از وب سایت‌های مختلف مجزا شامل دومین‌های آموزشی، دولتی، خبری و تجاری جمع‌آوری کرده‌ایم. ما این پرسش‌ها را مورد بررسی قرار داده‌ایم:

- آیا Bias روبات وجود دارد؟
 - چگونه چنین Bias به طور کمی اندازه‌گیری می‌شود؟
 - معنی و استنباط از این Bias چیست؟
- پیشنهادات ما:
- پیشنهاد یک معیار کمی برای اندازه‌گیری خودکار Bias روبات.
 - با به بکارگیری معیار در نمونه‌ی بزرگی از وب‌سایت‌ها یافته‌های خود را در مورد روبات‌های پرترفدار و بالعکس ارائه می‌دهیم.

ادامه مقاله به این صورت سازمان یافته: در بخش ۲ به طور مختصر پروتکل‌های ممانعت از روبات‌ها را معرفی می‌کنیم. در بخش ۳ جمع‌آوری اطلاعات این مطالعات را ارائه می‌دهیم. در بخش ۴ معیار Bias را ارائه می‌دهیم و شرح می‌دهیم که چگونه آن را برای اندازه‌گیری bias روبات‌ها استفاده کرد. در بخش ۵ مشاهداتمان را روی Robots.txt و در مورد مفاهیم آن بحث می‌کنیم. در بخش ۶ نتیجه‌گیری مقاله و برای برای کارهای آینده.

۲. پروتکل ممانعت از روبات‌ها:

پروتکل‌های ممانعت از روبات‌ها یک قرارداد است که به مدیران وب سایت اجازه می‌دهد که به روبات‌ها نشان دهد که کدام بخش از سایت آن‌ها نباید مورد بازدید قرار گیرد. اگر در وب سایت فایل Robots.txt نباشد روبات‌ها برای دسترسی به تمام محتویات آزاد هستند.

قالب پروتکل‌های ممانعت از روبات‌ها در [12] توضیح داده شده است. فایلی که Robots.txt نام دارد تحت فهرست ریشه وب سرور قرار می‌گیرد. هر خط فایل Robots.txt به صورت زیر می‌باشد:

<field>:<optionalspace><value><optionalspace>

سه نوع برجسب حساس به کوچک و بزرگ بودن حروف برای <field> وجود دارد که دستورات را معین می‌کند: User-Agent, Disallow, Allow. و نیز دستور غیر رسمی Crawl-Delay توسط بسیاری از وب سایت‌ها برای محدود کردن تعدد بازدیدهای روبات استفاده می‌گردد.

فایل Robots.txt با یک یا چند فیلد User-Agent برای مشخص کردن اینکه کدام روبات‌ها را شامل می‌شود آغاز می‌گردد و در ادامه تعدادی فیلدهای Allow یا Disallow که نمایان‌گر تنظیمات روبات است می‌باشد. توضیحات در همه جای فال مجاز هستند و می‌توانند شامل فضای خالی اختیاری باشند. توضیحات با حرف توضیح '#!' آغاز و Linebreak خاتمه می‌یابد. یک نمونه از Robots.txt را در زیر می‌بینید.

```

User-Agent: *
Disallow: /robots/
Disallow: /src/
Disallow: /botseer
Disallow: /uastiring
Disallow: /srcseer
Disallow: /robotstxtanalysis
Disallow: /whois

User-Agent: googlebot
Disallow: /robotstxtanalysis
Disallow: /uastiring

User-Agent: botseer
Disallow:

```

این فایل نشان می‌دهد که Googlebot نمی‌تواند "/robotstxtanalysis" و "/uastiring" را مشاهده کند. Botseer می‌تواند تمام فهرست‌ها و فایل‌های روی سرور را بازدید کند. تمام ربات‌های دیگر باید از قوانین * : User-agent پیروی کنند که بر طبق آن نمی‌توانند فایل‌ها و فهرست‌های زیر را بازدید کنند:

"/robots/" , "/src/" , "/botseer" , "/uastiring" , "/srcseer" , "/robotstxtanalysis" , "/whois"

۳. bias ربات :

ما $\Delta P(r)$ یک مقیاس برای میزان مساعدت ربات‌ها در یک فایل نمونه Robots.txt را پیشنهاد می‌کنیم، که برای اندازه‌گیری میزان اعتبار یک ربات توسط مجموعه‌ای از وب سایت‌ها می‌باشد. تعریف رسمی Bias ربات در زیر توضیح داده شده است.

۳-۱ الگوریتم بدست آوردن Bias :

تعریف ما از یک ربات معتبر یا همکار رباتی است که بر اساس فایل Robots.txt اجازه دسترسی به فهرست‌های بیشتری نسبت به ربات‌های همه منظوره را دارد. ربات‌های همه منظوره به هر رباتی گویند که با هیچ کدام از اسامی User-Agent در فایل Robots.txt مطابقت ندارد. به عبارت دیگر نام ربات‌های همه منظوره در فایل‌های Robots.txt ظاهر نمی‌شود.

مجموعه فایل‌های Robots.txt در مجموعه داده‌هایمان را F می‌گذاریم. یک فایل Robots.txt معین را $f \in F$ می‌گذاریم، و مجموعه از ربات‌های نام برده شده در Robots.txt فایل f را R می‌گذاریم. برای هر ربات $r \in R$ الگوریتم $\text{GetBias}(r, f)$ را که در الگوریتم ۱ آمده است را تعریف می‌کنیم. Getbias درجه اعتبار یک ربات r را که در Robots.txt فایل f آمده را بیان می‌کند.

Algorithm 1 $\text{GetBias}(r, f)$

```

1: if  $r$  is * then
2:   return 0
3: end if
4: Construct  $DIR$  for  $f$ ;
5:  $bias = 0$ 
6: for all  $d \in DIR$  do
7:   if  $d$  is allowed for * then
8:      $D_u \leftarrow d$ 
9:   end if
10: end for
11: for all  $d \in DIR$  do
12:   if  $d$  is allowed for  $r$  then
13:      $D_r \leftarrow d$ 
14:   end if
15: end for
16:  $bias = |D_r| - |D_u|$ 
17: return  $bias$ 

```

مجموعه فهرست‌هایی که در Robots.txt فایل f در یک وب سایت مشخص ظاهر می‌شود را DIR می‌گذاریم. DIR برای برآورد ساختار واقعی یک فهرست به استفاده می‌شود به این دلیل که پروتکل ممانعت از روبات‌ها به طور پیش فرض فهرست‌هایی که با Robots.txt مطابقت ندارد را به عنوان فهرست‌های مجاز در نظر می‌گیرد.

$D_u \in DIR$ مجموعه‌ای از فهرست‌هایی است یک روبات همه منظوره برای بازدید از آن مجاز است. اگر هیچ شرطی برای User-Agent* وجود نداشته باشد روبات به طور پیش فرض می‌تواند به همه چیز دسترسی داشته باشد. $D_r \in DIR$ مجموعه‌ی فهرست‌هایی است که روبات اجازه بازدید از آن را دارد. $|D_u|$ و $|D_r|$ تعداد فهرست‌ها در D_u و D_r است.

برای روبات داده شده r الگوریتم ابتدا محاسبه می‌کند چه تعداد از فهرست‌ها در DIR در r مجاز هستند. سپس مقدار Bias را که تفاوت بین تعداد فهرست‌هایی موجود در DIR که برای روبات r مجاز است و تعداد فهرست‌هایی که برای روبات‌های همه منظوره مجاز هستند را بدست می‌آورد. در الگوریتم Getbias برای روبات‌های همه منظوره مقدار مرجع صفر را برمیگرداند (GetBias Returns 0). مقدار Bias برای روبات‌های معتبر مقداری مثبت را برمی‌گرداند. مقدار بیشتر برای یک روبات به معنای اعتبار بیشتر روبات است. بالعکس مقدار bias برای روبات‌های نامعتبر منفی می‌باشد که با ما برای Bias سازگار است. بنابراین Bias یک روبات در Robots.txt می‌تواند توسط متغیرهای طبقه‌بندی شده نمایش داده شود که سه نوع طبقه‌بندی معتبر، نامعتبر و بدون Bias را شامل می‌شود. به عنوان مثال Robots.txt را در <http://BotSeer.ist.psu.edu> مورد بررسی قرار می‌دهیم.

$DIR = \{"/robots/", "/src/", "/botseer/", "/uastring/", "/srcseer/", "/robotstxtanalysis/", "/whois/"\}$.

بر طبق الگوریتم داریم:

$$D_u = \{\text{null}\}$$

$$D_{\text{botseer}} = \{"/robots/", "/src/", "/botseer/", "/uastring/", "/srcseer/", "/robotstxtanalysis/", "/whois/"\}$$

$$D_{\text{Google}} = \{"/robots/", "/src/", "/botseer/", "/srcseer/", "/whois/"\}.$$

بنابراین

$$|D_u|=0, |D_{\text{botseer}}|=7, |D_{\text{Google}}|=5$$

بنابر الگوریتم ۱:

$$|D_u|-|D_u|=0, \text{bias}_{\text{botseer}} = |D_{\text{botseer}}|-|D_u|=7, \text{bias}_{\text{Google}} = |D_{\text{Google}}|-|D_u|=5$$

بنابراین روبات‌های googlebet و botseer توسط این وب سایت مورد اعتماد هستند و در طبقه معتبر قرار می‌گیرند و بقیه روبات‌ها در طبقه بدون Bias.

۳-۲ اندازه گیری bias کلی:

بر اساس مقدار Bias برای هر فایل، $\Delta P(r)$ را به منظور محاسبه برای مجموعه‌ای از فایل‌های Robots.txt ارائه می‌کنیم. $N=|F|$ تعداد کل فایل‌های Robots.txt در مجموعه داده‌ها است. $\Delta P(r)$ می‌تواند به صورت زیر تعریف گردد:

$$\Delta P(r) = P_{\text{favor}}(r) - P_{\text{disfavor}}(r) = \frac{N_{\text{favor}}(r) - N_{\text{disfavor}}(r)}{N} \quad (1)$$

که $N_{\text{favor}}(r)$ و $N_{\text{disfavor}}(r)$ به ترتیب برابر تعداد دفعاتی است که روبات مورد تأیید و عدم تأیید قرار می‌گیرد.

$P_{\text{favor}}(r)$ نسبتی از فایل Robots.txt r که مورد تأیید قرار می‌گیرد و $P_{\text{disfavor}}(r)$ نسبتی از فایل Robots.txt r است که مورد تأیید قرار نمی‌گیرد.

نسبت فایل Robots.txt که مورد تأیید و یا عدم تأیید قرار می‌گیرد برای بررسی آماری محاسبه ساده‌ای دارد. از آنجائیکه بیشتر از دو مقدار داریم (favor, disfavor, no bias) جداسازی دو مقدار برای محاسبه Bias کلی منتج به نتیجه دقیق نمی‌شود. به عبارت

$$P_{\text{favor}}(r) + P_{\text{disfavor}}(r) < 1$$

هر مقدار فقط یک وجه از bias را نتیجه می‌دهد. برای مثال روباتی به نام ia_archiver 0.24% مورد تأیید وبسایت‌های مجموعه اطلاعات ما قرار می‌گیرد و در مورد momspider 0.21%. و به ترتیب عدم تأیید ia_archiver و momspider برابر است با 1.9% و 0%. اگر ما فقط نسبت تأیید را مدنظر قرار دهیم ia_archiver نسبت به momspider درجه تأیید بیشتری دارد.

$\Delta P(r)$ اختلاف بین نسبت سایت‌های که یک روبات خاص را مورد تایید و یا عدم تایید قرار داده‌اند می‌باشد بنابراین هر دو مورد هماهنگ هستند. برای مثال بالا $\Delta p(\text{ia_archiver})=1.66\%$ و $\Delta p(\text{momspider})=0.21\%$ بنابراین momspider درجه تایید بیشتری از ia_archiver دارد. برای فایل r بدون bias داریم $\Delta p(r)=0$. محاسبه bias می‌تواند موارد گمراه کننده را از بین ببرد.

۳-۳ آزمون تایید پذیری

تایید پذیری در واقع یک تابع رتبه بندی برای روبات‌ها است. برای محاسبه دقت این تابع رتبه‌بندی ما تست کارایی رتبه‌بندی را بر اساس امتد ضریب رتبه بندی Kendall [11] اجرا کردیم. ضریب رتبه‌بندی به طور خلاصه در ادامه آمده است. جزئیات ارزیابی کارایی رتبه‌بندی می‌توانید در [9] ببینید. برای فایل Robots.txt به نام f ، m_a را برابر اندازه تابع bias برای تمام روبات‌هایی که در f دیده می‌شود قرار می‌دهیم، r_i و r_j عضو R را دو روبات در f فرض می‌کنیم. نشان می‌دهیم $r_i < m_a$ اگر r_i در رتبه بالاتری از R_j با اندازه‌گیری m_a باشد. بنابراین برای هر دو مقدار m_a و m_b Kendall τ را تعریف کرده است براساس P_f که مقدار دوتایی‌های هماهنگ و Q_f که مقدار دوتایی‌های ناهماهنگ باشند. دو مقدار $r_i \neq r_j$ هماهنگ هستند اگر m_a و m_b با یکدیگر مطابقت داشته باشد و ناهماهنگ هستند اگر مطابق نباشند. در این مورد Kendall τ را به این صورت بیان کرده است:

$$\tau_f(m_a, m_b) = \frac{P_f - Q_f}{P_f + Q_f} \quad (2)$$

برای هر مقداری m_a و m_b داشته باشند $\tau_f(m_a, m_b)$ نشان می‌دهد که دو مقدار چگونه در فایل f با هم مطابقت دارند. m_a را نشان‌گر مقدار واقعی تابع رتبه بندی روبات می‌دانیم. در صورتیکه که ما تابع رتبه‌بندی واقعی را نمی‌دانیم و فقط رتبه‌بندی جزئی روبات را برای هر Robots.txt داریم که بر اساس مقدار bias که قبلاً توضیح دادیم بدست می‌آید. بنابراین $\tau_f(m_a, m_b)$ برای هر فایل Robots.txt در مجموعه داده را محاسبه می‌کنیم. اگر برای یک فایل داده شده Robots.txt $\tau_f(m_a, m_b)=1$ به این نتیجه می‌رسیم که فایل مذکور فایلی متوازن است و در غیر این صورت فایلی غیر متوازن. با محاسبه فایل‌های متوازن و نامتوازن در مجموعه داده‌ها داریم: $(P+Q=N)$

$$\tau(m_a, m_b) = \frac{P - Q}{P + Q} = 1 - \frac{2Q}{N} \quad (3)$$

ما با استفاده از δP روبات‌ها را رتبه‌بندی کردیم. لسیت رتبه‌بندی شده با رتبه‌بندی واقعی مقایسه شد. مقدار میانگین τ برابر 0.957 که معتقدیم به اندازه کافی برای اندازه‌گیری مقدار bias دقیق می‌باشد.

۴- جمع آوری اطلاعات:

به منظور پتانسیل bias روبات‌ها در وب فعالیت مطالعاتی ما در طیف گسترده‌ای از وب سایت‌ها از دومین‌های مختلف و نیز مکان‌های مختلفی انجام شده است که با توضیحات در زیر آمده است.

۴-۱ منابع اطلاعاتی:

[6] Open Directory Project بزرگترین و گسترده‌ترین فهرست وب است که با نیروی انسانی جمع‌آوری شده است. منبع اصلی ما برای جمع‌آوری URL‌های اولیه DMOZ بود به دلیل اینکه Open Directory Project تعداد گسترده‌ای از URL‌ها را در طبقات مختلفی طبقه‌بندی کرده است. که ما را قادر ساخت تا داده‌هایمان را از دومین‌های مختلف و از مکان‌های مختلف جمع‌آوری کنیم. اطلاعات ما از این منبع سه نوع دومین آموزشی، خبری و دولتی بود. دومین‌های دانشگاهی بعداً به دومین دانشگاهی آمریکایی، اروپایی و آسیایی تقسیم شد.

از آنجائیکه ساختار فهرست دومین‌های تجاری در DMOZ پیچیده بود و بسیاری از آن‌ها مشترک بودند ما از لیست ۱۰۰۰ شرکت ثروتمند در ۲۰۰۵ استفاده کردیم. قطعاً در داده‌های جمع‌آوری شده ما محدودیت‌هایی وجود دارد. اولاً بدلیل اینکه مجموعه وبسایت‌های DMOZ برای سایر کشورها بخصوص وبسایت‌های غیر انگلیسی زبان محدود بود و اکثر وبسایت‌ها از آمریکا بودند. دوماً بدلیل اینکه ورورهای DMOZ توسط نیروی انسانی ویرایش و سازماندهی می‌شده احتمال وجود خطا می‌باشد. و نهایتاً ما

وبسایت‌های تجاری را از میان لیست ۱۰۰۰ شرکت ثروتمند جمع‌آوری کردیم که محتوی اطلاعاتی در مورد شرکت‌های بزرگ بود که این دومین‌ها نمی‌تواند نمایند شرکت‌های تجاری کوچک باشد و ما قصد داریم در تحقیقات آینده به این محدودیت‌ها اشاره کنیم.

۲-۴ کاوشگری برای Robots.txt

ما یک کاوشگر مخصوص برای این مطالعات پیاده‌سازی کردیم. کاوشگر با کاوش در وب سایت‌های که در DMOZ با طبقه‌بندی موضوعی، نام وبسایت و سازمان‌های وابسته در مکان‌های مختلف شروع به کار کرد. سپس کاوشگر وجود Robots.txt را در دومین‌های مذکور بررسی و در صورت وجود آن‌را برای تحلیل‌های آفلاین دانلود می‌کرد. یک مازول جدا کننده و فیلتر کننده نیز کاوشگر اغام شده بود تا تکراری‌ها را از بین ببرد و از اینکه صفحات بدست آمده در دومین‌های مقصد هستند مطمئن گردد. علاوه بر فهرست سطح ریشه کاوشگر ما سایر مکان‌های ممکن را برای Robots.txt بررسی می‌کرد. زیر فهرست‌های وب سایت تا سطح ۳ سرکشی شد. نتایج نشان داد که موارد کمی بود که Robots.txt در فهرست ریشه نباشد که بر طبق پروتکل ممانعت از روبات است. فایل‌های با نام غلط نیز توسط کاوشگر ما مورد بررسی قرار گرفت که دذ موارد نادری Robot.txt به جای Robots.txt استفاده شده بود.

برای مشاهده کردن مشخصات زودگذر کاوش ۵ بار برای همان مجموعه وبسایت از دسامبر ۲۰۰۵ تا اکتبر ۲۰۰۶. به منظور تحلیل مشخصات موقتی فایل‌های Robots.txt دانلود شده بر اساس تاریخ کاوش آرشیو شد.

۳-۴ آمار:

ما ۷۵۹۳ وب سایت مجزا را که شامل ۶۰۰ وب سایت دولتی، ۲۰۴۷ وب سایت روزنامه، ۱۴۸۷ وب سایت دانشگاهی آمریکایی، ۱۴۲۰ وب سایت دانشگاهی اروپایی، ۱۰۳۹ وب سایت دانشگاهی آسیایی و وب سایت ۱۰۰۰ شرکت را مورد کاوش و بررسی قرار دادیم. تعداد وب سایت‌هایی که Robots.txt داشتند در ۵ مرحله کاوش در جدول ۱ آمده است.

	Websites	Collected robots.txt files				
		Dec. 2005	Jan. 2006	May. 2006	Sep. 2006	Oct. 2006
Government	600	248	257	263	262	264
Newspaper	2047	859	868	876	937	942
USA Univ.	1487	615	634	650	678	683
European Univ.	1420	497	510	508	524	537
Company	1000	303	306	319	341	339
Asian Univ.	1039	140	248	149	165	160
Total	7593	2662	2823	2765	2907	2925

Table 1. Number of robots.txt found in each domain for each crawl.

توضیح بهتر استفاده از Robots.txt در وبسایت‌های دومین‌های مختلف در شکل ۱ که نشانگر نسبت وبسایت‌هایی که دارای این فایل هستند مشخص شده است.

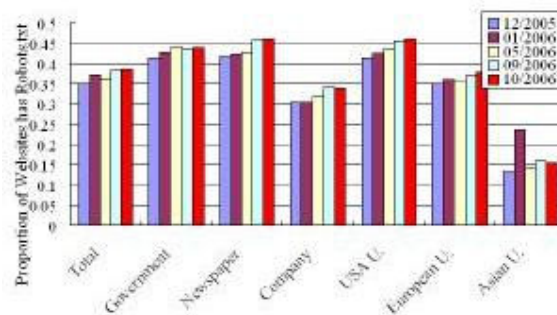
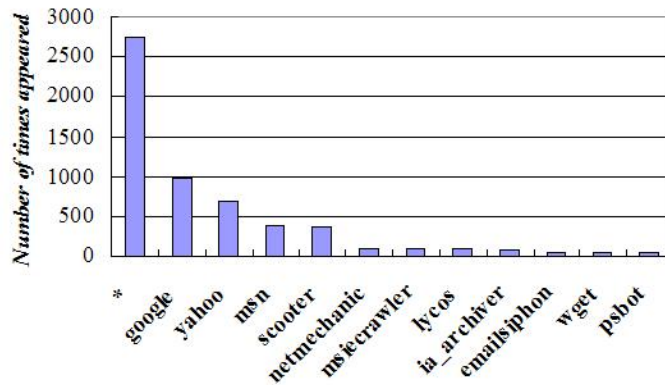


Figure 1. Probability of a website that has robots.txt in each domain.

به طور کلی به جز در مورد وب سایت‌های دانشگاه‌های آسیایی استفاده از Robots.txt افزایش یافته است. 46.02% از وب سایت‌های روزنامه‌ها اکنون Robots.txt را اجرا کرده‌اند و دومین‌های روزنامه که بیشترین پذیرش پروتکل ممنوعت از روبات را داشته‌اند. 45.93% از وب سایت‌های دانشگاهی آمریکا پروتکل ممنوعت از روبات را پذیرفته‌اند که بیشتر از سایت‌های اروپا (37.8%) و آسیا (15.4%) است. دلیل اینکه موتورهای جستجوگر و عاملین جستجوی هوشمند بسیار برای دسترسی به اطلاعات وب با اهمیت شده‌اند این نتایج مورد انتظار بود. پروتکل ممنوعت از روبات بیشتر توسط وب سایت‌های دولتی، روزنامه‌ها و دانشگاهی در آمریکا مورد پذیرش واقع شده است که به طور گسترده‌ای از آن برای حفاظت از اطلاعاتی که نباید در اختیار عموم قرار بگیرد و توازن تراکم این وب سایت‌ها استفاده می‌شود. یک گزارش تفصیلی از استفاده از Robots.txt را می‌توانید در [6] ببینید.

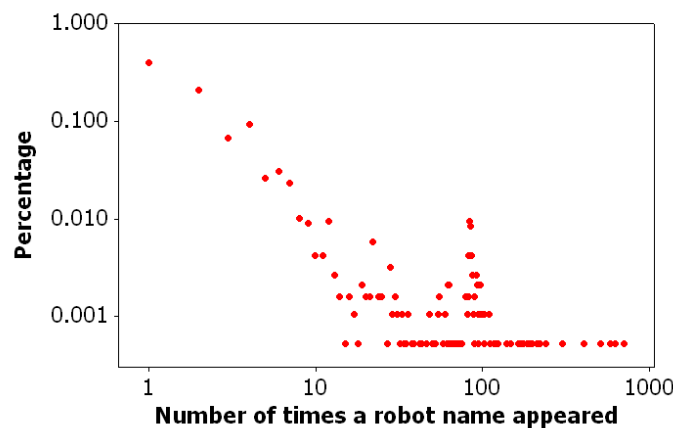
۵- نتایج :

در مجموعه داده‌های ما ۱۰۵۶ روبات یافت می‌شود. روبات‌های عمومی بیشترین استفاده را در فیلد User-Agent با ۲۷۴۴ بار داشته است، که به معنی این است که 93.8% از فایل‌های Robots.txt برای روبات‌های عمومی اجرا می‌شود. 72.4% از روبات‌ها یک یا دو بار ظاهر شده‌اند. روبات‌هایی که بیشتر از همه در مجموعه داده‌های ما پدیدار شده در شکل ۲ می‌بینید.



۱-۵ تاریخچه Bias

توزیع اینکه چند بار یک روبات مورد استفاده قرار گرفته است (شکل ۳) طی ۱۱ ماه گذشته تغییری نکرده است. بنابراین ما نتایج بدست آمده برای Bias را در آخرین کاوش که تغییری نکرده نشان می‌دهیم



از آنجائیکه اغلب روبات‌ها یک یا دو بار در مجموعه داده‌ها ظاهر شدند، امتیاز رده‌بندی آن‌ها در رده میانی قرار گرفته و غیر قابل تشخیص بود. با فقط روبات‌های با امتیازات تاییدپذیری بالا و پایین را مورد بررسی قرار دادیم. ۱۰ روبات با بیشترین امتیاز و ۱۰ روبات با کمترین امتیاز را در جدول ۲ می‌بینید که N مقدار نمونه است و N_{favor} تعداد دفعاتی است که روبات مورد تایید قرار گرفته و $N_{disfavor}$ تعداد دفعاتی است که روبات مورد تایید نبوده است و انحراف استاندارد مطلق [3] برای $\Delta P(r)$ است. انحراف استاندارد مطلق زمانیکه

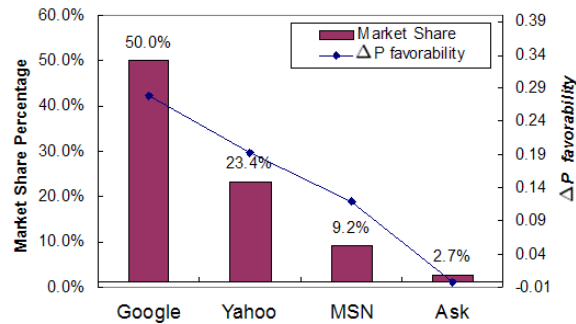
از Δp برای محاسبه درجه تایید پذیری روبات در وب استفاده می‌شود واریانس را می‌دهد . اندازه‌گیری Bias ما نشان می‌دهد که روبات‌هایی که بسیار مورد تایید قرار گرفته‌اند از موتورهای جستجوگر و سازمان‌های شناخته شده مانند Google و Yahoo و MSN بوده که بیشتر از سایر روبات‌ها مورد تایید قرار گرفته است ، دقت داشته باشید که در رده‌بندی روبات‌های تایید نشده Δp آن‌ها نشان دهنده این نیست که آن‌ها به ندرت در Robots.txt مورد آزمایش ظاهر شده‌اند . به عبارت دیگر اغلب روبات‌های تایید نشده گردآورنده‌های ایمیل ("CherryPicker" , "emailsiphon") و مرورگرهای آفلاین ("Wget" , Webzip") بودند . از لحاظ شخصی معقول به نظر می‌رسد که مدیران وب سایت‌ها از روبات‌هایی که هدف اصلی آن‌ها جمع‌آوری اطلاعات شخصی است ممانعت بعمل آورند . همچنین معمولاً نمی‌خواهند وب سایت آن‌ها به طور کلی توسط دیگران کپی شود . اگر چه روبات‌هایی از شرکت‌های شناخته شده نیز هستند که مورد تایید قرار نگرفته‌اند مانند "MSIEcrawler" و "ia_archiver" . روباتی است که در IE جاسازی شده و زمانیکه کاربران IE یک صفحه را Bookmark می‌کنند MSIEcrawler تمام صفحه و لینک‌های مرتب با آن را دانلود می‌کند (لینک ها ، تصاویر ، جاوااسکریپت و ...) . "ia_archiver" یک کاوشگر از archive.org و Alexa.com است . یک لیست با توضیحات تفصیلی از روبات‌های شناخته شده را در <http://botseer.ist.psu.edu/namedrobots.html> می‌توانید ببینید .

Favored Robots (Sample size $N = 2925$)				
robot name	N_{favor}	N_{disf}	$\Delta P(r)$	σ
google	877	25	0.2913	0.0084
yahoo	631	34	0.2041	0.0075
msn	349	9	0.1162	0.0059
scooter	341	15	0.1104	0.0058
lycos	91	5	0.0294	0.0031
netmechanic	84	10	0.0253	0.0029
htdig	15	3	0.0041	0.0012
teoma	13	3	0.0034	0.0011
oodlebot*	8	0	0.0027	0.0010
momspider	6	0	0.0021	0.0008
Disfavored Robots (Sample size $N = 2925$)				
robot name	N_{favor}	N_{disf}	$\Delta P(r)$	σ
msiecrawler	0	85	-0.0291	0.0031
ia_archiver	7	55	-0.0164	0.0023
cherry picker	0	37	-0.0126	0.0021
emailsiphon	3	34	-0.0106	0.0019
roverbot	2	27	-0.0085	0.0017
psbot	0	23	-0.0079	0.0016
webzip	0	21	-0.0072	0.0016
wget	1	22	-0.0072	0.0016
linkwalker	2	20	-0.0062	0.0015
asterias	0	18	-0.0062	0.0015

Table 2. Top 10 favored and disfavored robots. σ is the standard deviation of $\Delta P(r)$.

ما همچنین به این نتیجه رسیدیم که bias روبات‌ها در دومین‌های مختلف تغییر می‌کند . Google همیشه تایید شده ترین روبات است . روبات‌های بسیار تایید شده دیگر در دومین مختلف تغییر می‌کند . Yahoo (روبات Slurt روبات Yahoo است) و MSN نیز توسط اکثر دومین‌ها مورد تایید قرار گرفته‌اند اما آن سایر روبات‌ها را مورد تایید قرار نداده‌اند . سایر روبات‌های رده بالا تایید شده مربوط به کاوشگران کد باز و سازمان‌های معروف است . روبات‌های مورد تایید قرار نگرفته برای دومین‌های مختلف تغییر می‌کند . اغلب آن‌ها گردآورندگان ایمیل و مرورگرهای آفلاین هستند . تفاوت‌ها ناشی از رفتار مختلف روبات‌ها در دومین‌های مختلف است .

بر اساس مطالعات اهمیت تاثیر قدرتمندتر شدن قدرتمندان، ما برای شرکت‌های خاصی ارتباط بین Bias روبات و سهم از بازار موتورهای جستجوگر را محاسبه نمودیم. سهم از بازار Google، Yahoo، MSN و Ask در ۱۱ ماه گذشته و نیز ΔP تاییدپذیری که با دو متغیر مستقل برای روبات‌ها مطرح می‌شود. ضریب همبستگی (PMCC)[15] Pearson Product-moment برای دو متغیر اندازه‌ای است برای میل دو متغیر X و Y که در یک مبحث یا ترکیب برای افزایش یا کاهش نسبت به هم اندازه‌گیری می‌شود. در مجموعه داده‌ها ما ضریب Pearson از سهم بازار ۴ شرکت و $\Delta P(r)$ روبات آنها 0.930 با مقدار $P < 0.001$ است. سهم از بازار موتورهای جستجوگر و Bias روبات در سپتامبر ۲۰۰۶ در شکل ۴ نشان داده شده است.



۶- نتیجه

ما یک تحقیق گسترده بر روی bias روبات‌ها در وب بر روی اطلاعات موثق و تحلیل‌های آماری تعداد زیادی از فایل‌های Robots.txt ارائه دادیم. نتایج نشان‌دهنده این است که روبات‌های موتورهای جستجوگر محبوب و درگاه‌های اطلاعاتی همچون Google، Yahoo و MSN اغلب توسط وبسایت‌هایی که ما مورد بررسی قرار دادیم مورد تایید قرار گرفته است. این نکته اشاره به موضوع قدرتمند شدن قدرتمندان در موتورهای جستجوگر محبوب نیز دارد. ما همچنین یک تناسب مطمئن بین سهم از بازار موتورهای جستجوگر و Bias روبات‌ها بیان کردیم. مطالعه ما نشان می‌دهد استفاده از Robots.txt در ۱۱ ماه افزایش داشته است به طوری در کاوش در مرحله اول ۲۶۶۲ و در مرحله آخر ۲۹۲۵ فایل را شامل می‌شود. ما مشاهده کردیم که 46.02% از وبسایت‌های روزنامه‌ها اکنون فایل Robots.txt را استفاده می‌کنند و دومین‌های روزنامه‌ها دومین‌هایی است که بیشترین پذیرش پروتکل ممانعت از روبات را دارد. 45.93% از وبسایت دانشگاه‌های آمریکایی پروتکل ممانعت از وبسایت را پذیرفته‌اند که به طور مشخص از دانشگاه‌های اروپایی (37.8%) و آسیایی (15.4%) بیشتر است. فعالیت‌های آینده ما در آینده بیشتر تقسیم می‌شود به تحلیل مکان‌های جغرافیایی برای بررسی تایید پذیری روبات در هر کشور.

فعالیت‌های آینده ما بیشتر در مورد تحلیل عمیق‌تر و در مقیاس بزرگتر تنظیمات و رفتارهای روبات‌ها می‌باشد. ما سایر مقیاس‌های مربوط به bias روبات را مورد بررسی قرار خواهیم داد. بررسی‌های تجربی ما در زمینه چگونگی تفسیر پروتکل ممانعت از روبات‌ها در روبات‌ها خواهد بود.

- [1] Fortune magazine .<http://money.cnn.com/magazines/fortune/-fortune500> , 2005 .
- [2] M.L.Boonk , D . R . A . d .Groot , F . M . T . Brazier , and A . Oskamp . Agent exclusion on websites . In proceeding of The 4th Workshop on the Law and Eletronic Agent , 2005
- [3] S.Chakrabarti , M.van den berg , and B. Dom . Focused crawling : A new approach to topic-specific web resource discovery . In Proc. of the 8th WWW conference,pages 545-562 ,1999
- [4] J. Cho , H. Garcia-Molina , and L.page . efficient crawling through url ordering . In proceeding of the 7th International WWW conference , 1998
- [5] T.Y Chum . World wide Web robots : An overview . online information review , 23(3): 135-142 , 1999
- [6] Dmoz The open directory project . <http://dmoz.org> , 2005
- [7] M.Drott . Indexing aids at corporate Websites: the use of robots.txt and meta tags . information processing and management , 38(2):209-219 , 2002
- [8] D.Eichmann . Ethical web agent . Computer Networks and ISDN Systems , 28(1-20):127-136 ,1995
- [9] t.joachims .optimizing search engine using clickthrough data . In KDD 02 : proceeding of the eight ACM SIGKDD international Conference on Knowledge discovery and data mining , pages 133-142 New York , NY , USA , 2002 . ACM Press
- [10] B.Kelly and I.Peacock . Webwaching Uk web Communities : Final report for the webWatch project . British library research and innovation report , 1999
- [11] M.Kendall . Rank Correlation Method .hafner , 1955
- [12] M. Koster A method for web robots control . In the internet draft , The internet Engineering Task Force (IETF) , 1996
- [13] R.L.Ott and M.T . Longnecker . an introduction to statistical method and data analysis . Duxbery press ; 5edition 2000
- [14] G.Pant , P. Srinivasan , and F.menczer . crawling the web,chapter web dynamic . springer-verlag , 2004
- [15] R.R Sokal and F.J.Rohlf . biometry . Freeman new york , 2001
- [16] Y.sun , Z.Zhuang and C.L.Giles . A large Scale study for robots.txt In WWW '07 : proceeding of the 16th international conference on world wide web , page 1123-1124 New York , NY USA ,2007 .ACM Press .