

مستندات خزنده‌ی وب^۱

علی علی محمدی

۹۶۱۳۰۲۷

- در فایل فشرده‌ی zip، فایل Web_Scraping.ipynb که یک Jupyter Notebook شامل اسکریپت‌هایی به زبان Python وجود دارد.
- این کُد بر روی بستر colab توسعه داده شده است و با آن سازگار است.
- از کتابخانه‌های BeautifulSoup برای عملیات اشتقاق^۲ و استخراج داده‌های موجود در سند^۳ HTML، requests برای دریافت محتویات HTML یک صفحه با استفاده از URL آن صفحه، json برای سهولت در استفاده‌ی فرمت^۴ JSON، Pillow برای نمایش تصاویر و عکس‌ها و base64 برای کدگذاری^۵ عکس‌ها به نحوی که در قالب JSON Object قابل ذخیره باشد (زیرا نمی‌توان آرایه‌ی بایت^۶ را بدون کدگذاری و به صورت خام در JSON Object ذخیره کرد)، استفاده شده است.
- نام سلول‌های کُد به همراه کامنت‌های برنامه به وضوح روند کار را تشریح می‌کنند. همچنین نام توابع به صراحت بیان‌گر کاربرد و منظور آن تابع است. همچنین نحوه‌ی ذخیره‌سازی اطلاعات در JSON Object نیز با توجه کلیدها^۷، گویا و واضح است.
- برنامه از قابلیت اطمینان بالایی برخوردار است و داده‌ها را به طور پیوسته در طی عملیات خزش ذخیره می‌کند. بدین ترتیب، در صورت به وجود آمدن وقفه^۸ در اجرای برنامه، مثل قطعی اینترنت یا توقف عملیات توسط کاربر یا سیستم‌عامل، داده‌هایی که تا آن لحظه خزش شده‌اند، از دست نمی‌رود و درون فایل خروجی موجود خواهند بود.
- در هنگام اجرای برنامه، می‌توان تعداد کمینه‌ی محصولاتی لازم است مورد خزش قرار گیرند، در ورودی تابع scrape_category وارد نمود؛ بدین ترتیب، برنامه پس از اتمام خزش یک صفحه از فهرست محصولات، به صفحه‌ی بعدی می‌رود و آن را مورد خزش قرار می‌دهد و این چرخه را تا زمانی که حداقل item_min عدد محصول مورد خزش قرار گرفته باشند، ادامه می‌دهد.
- به ازای خزش هر دسته (هر بار فراخوانی تابع scrape_category)، یک فایل با پسوند json و با نام آن دسته ساخته می‌شود و در هر سطر آن، یک JSON Object حاوی اطلاعات یک محصول وجود دارد؛ به عنوان مثال، برای دسته‌ی «تخت‌ها»، اطلاعات در فایل beds.json ذخیره می‌شوند.
- در رابطه با فایل robots.txt می‌توان گفت که این فایل به صورت محترمانه (!) از ربات‌های جست‌وجوگر و خزنده تقاضا می‌کند که محدودیت‌هایی را در نظر بگیرند و آن‌ها را رعایت کنند؛ به عنوان مثال، اجازه‌ی جست‌وجوی محتویات وبسایت را برای دسته‌ی خاصی از bot‌های جست‌وجوگر صادر می‌کند و دسته‌ای دیگر را ممنوع می‌کند. در برخی از موارد نیز اجرای درخواست‌ها^۹ را ممنوع اعلام می‌کند و از bot‌ها می‌خواهد که دستورات و درخواست‌های این‌چنینی را به سمت سرور ارسال نکنند یا از آن‌ها می‌خواهد که برای جست‌وجو در آن وبسایت، استفاده از چه فیلترهایی مجاز است. البته شایان به ذکر است

^۱ Web Scraper

^۲ Parsing

^۳ Document

^۴ JavaScript Object Notation

^۵ Encoding

^۶ Byte-Array

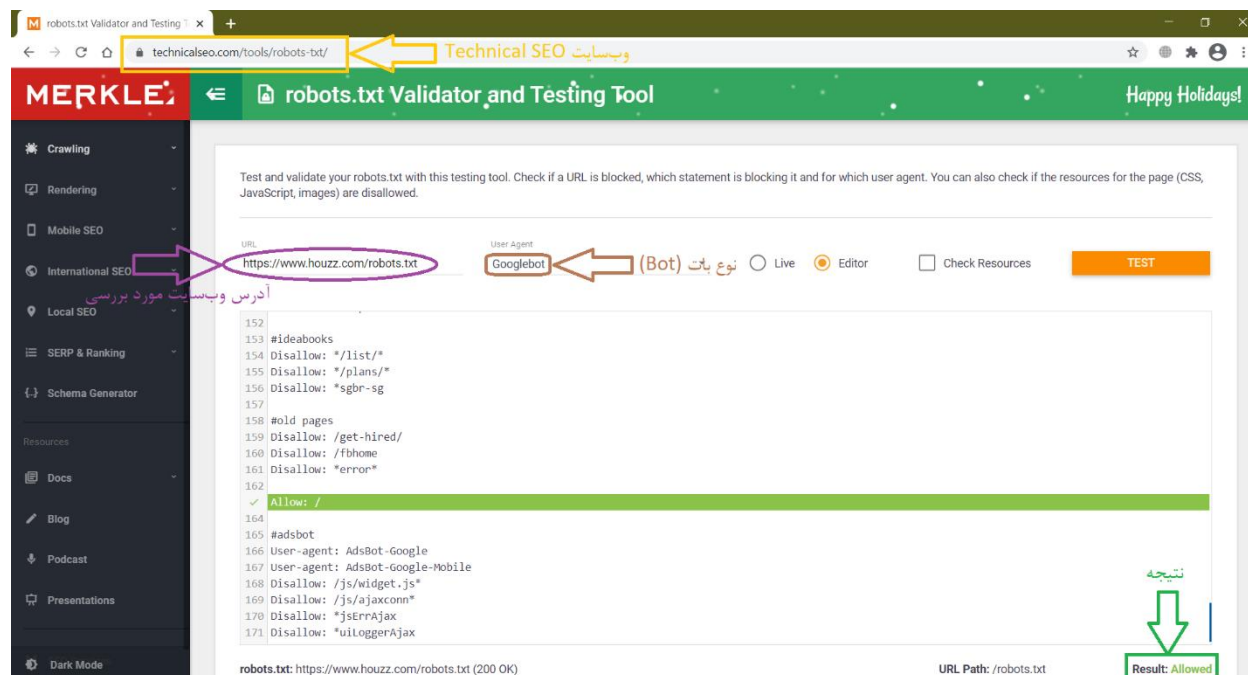
^۷ Keys

^۸ Interruption

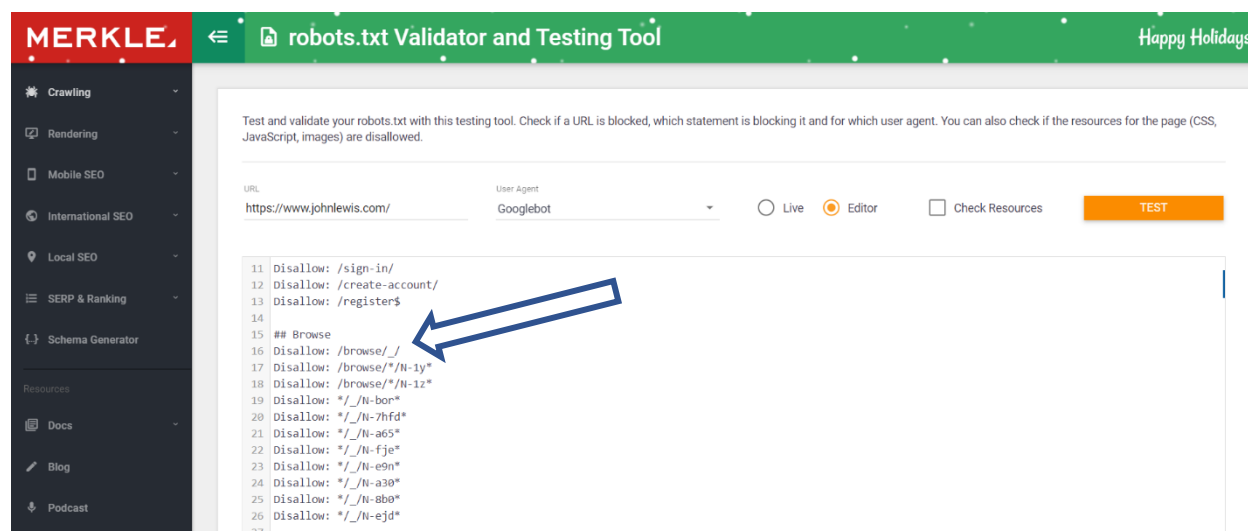
^۹ Queries

که ربات‌های جست‌وجوگر لزوماً این قوانین را رعایت نمی‌کنند؛ به عنوان مثال، موتور جست‌وجوی گوگل^{۱۰} اعلام کرده است که از تاریخ یکم سپتامبر سال ۲۰۱۹، دیگر از قوانین **robots.txt** پیروی نمی‌کند.

- برای درک بهتر این فایل و بررسی دقیق‌تر آن می‌توان از ابزار مربوط به آن در وبسایت «تکنیکال سئو»^{۱۱} استفاده کرد. این وبسایت امکان استفاده از انواع bot خزنده را فراهم می‌سازد و در پایان نتیجه‌ی کلی را نیز اعلام می‌کند.
- ✓ برای بررسی و مقایسه‌ی وبسایت‌های **Houzz** و **John Lewis** از ابزار فوق استفاده شده است؛ نتیجه در تصاویر زیر مشاهده می‌شود:



نتیجه‌ی بررسی وبسایت Houzz



نتیجه‌ی بررسی وبسایت John Lewis

¹⁰ Google Search Engine

¹¹ Technical SEO

✓ همان‌طور که مشاهده می‌شود، وب‌سایت Houzz اجازه‌ی جست‌وجوی محصولات را در فایل robots.txt صادر کرده است اما درباره‌ی وب‌سایت John Lewis چنین نیست؛ بلکه اجازه‌ی دسترسی را مسدود کرده است. بنابراین به نظر می‌رسد که وب‌سایت Houzz برای خزش مناسب‌تر از وب‌سایت John Lewis است.

✓ همچنین نحوه‌ی آدرس‌دهی در URL صفحه‌های وب‌سایت Houzz دارای ساختاری بسیار مشخص هستند و به راحتی می‌توان آدرس صفحه‌های مورد نیاز را از روی الگوی آن ساخت و مورد استفاده قرار داد اما درباره‌ی وب‌سایت John Lewis چنین نیست و الگوها^{۱۲} به راحتی قابل تشخیص نیستند و الگوهای پیچیده‌ای دارد که کار را برای دسترسی خودکار خزنده به صفحه‌های مختلف آن وب‌سایت دشوار می‌سازد.

➤ درباره‌ی Selenium و Scrapy می‌توان گفت که هر دو فریم‌ورک‌های پایتون^{۱۳} هستند که برای خزش وب استفاده می‌شوند. این در حالی است که Selenium برای خودکارسازی تعاملات مرورگر وب^{۱۴} و تست خودکار Web Application ها استفاده می‌شود اما Scrapy برای دانلود HTML، پردازش و ذخیره‌سازی آن مورد استفاده قرار می‌گیرد.

➤ از ویژگی‌های مهم Selenium می‌توان به این نکته اشاره کرد که به راحتی می‌توان با هسته‌ی مفاهیم JavaScript^{۱۵} (DOM) کار کرد. همچنین مدیریت و به کار بردن AJAX و PJAX در آن آسان است.

➤ یکی از مزایای اصلی Scrapy این است که بر روی بستر Twisted که یک فریم‌ورک شبکه‌ی ناهمزمان^{۱۶} است، ساخته شده است که یعنی Scrapy هنگام ارسال درخواست‌ها به کاربران از مکانیسم غیر مسدودکننده^{۱۷} استفاده می‌کند.

➤ به طور کلی، Scrapy در دسته‌ی Web Scraping API قرار می‌گیرد، در حالی که Selenium در دسته‌ی Browser Testing Tools قرار دارد.

در صورت وجود هرگونه سؤال یا ابهام، با ایمیل alialimohammadi@ce.aut.ac.ir در تماس باشید.

¹² Patterns

¹³ Python Frameworks

¹⁴ Automate Web-browser Interaction

¹⁵ Core JavaScript Concepts

¹⁶ Asynchronous Networking Framework

¹⁷ Non-blocking Mechanism