

بسم الله الرحمن الرحيم

راهنمای Web ماژول crawler

تیر ۹۷

۱-۱ نگارش ها	۲
۲-۱ توضیحات اولیه	۲
۳-۱ نحوه کار ماژول	۲
۴-۱ محتوای موجود در socket request	۴
۵-۱ جدول خطاها	۵

فصل ۱

۱-۱ نگارش ها

در جدول زیر تاریخچه نگارش های مختلف API آورده شده است.

تاریخ	نسخه	توضیحات
اسفند ۹۶	۱,۰	نسخه اولیه document آماده شد.
تیر ۹۷	۱,۱	خروجی تصحیح شد. خطایابی ماژول انجام شد. Progressbar اضافه شد.

۲-۱ توضیحات اولیه

این ماژول جهت بازدید از یک url و ذخیره محتوای آن نوشته شده است. با توجه به عمق در نظر گرفته شده و قابل تنظیم ، این ماژول، لینک های موجود در url اولیه را هم بازدید کرده و محتوای آنها را نیز ذخیره می کند. در کد مربوط به این ماژول از سرویس Hrobot استفاده شده است. محل ذخیره محتوای url ها مسیر و پوشه مشخصی است و در نهایت ماژول لسیته از ادرس محل ذخیره url های بازدید شده را در قالب entity_property برای کاربر ارسال می کند.

۳-۱ نحوه کار ماژول

این ماژول دو آرگومان به نام های url و depth را به عنوان ورودی دریافت می کند و لیستی از فایل های crawl شده را به صورت entity_property به عنوان خروجی بر می گرداند. جداول زیر ورودی ها و پارامترهای خروجی این ماژول را نشان می دهد. همچنین تصویر یک نمونه خروجی این ماژول در این [لینک](#) قابل مشاهده است.

ورودی	توضیحات
url	ادرس url اولیه که کاربر قصد crawl کردن آن را دارد
depth	عمق درخواستی برای crawl

جدول مربوط به ورودی ماژول

برای بهبود عملکرد ماژول crawler، دو فاکتور optional به ورودی crawler اضافه شده است، برای استفاده از این موارد کافی است به آن را صورت key-value در ورودی ارسال کنید که توضیحات هر یک به شرح زیر است:

کلید	مقدار	نوضیحات
base_url_constraint	True, False	ارسال این مقدار باعث می شود که crawler فقط صفحاتی را ذخیره کند که دارای url پایه مشابه با url پایه مورد درخواست کاربر باشد. در صورتی که این مقدار ارسال نشود مقدار آن false خواهد بود.
link_limit	2000	این عدد مشخص کننده ی تعداد صفحات مجاز برای ذخیره شدن توسط crawler می باشد. اگر تعداد صفحات ذخیره شده از این عدد بگذرد خطا رخ می دهد. در صورتی که این مقدار ارسال نشود مقدار آن برابر ۲۰۰۰ خواهد بود.

جدول مربوط به ورودی های اختیاری ماژول

توجه !

دقت شود که در جدول بالا پارامتر depth که عمق درخواستی را نشان میدهد به صورت اختیاری بوده و در صورت تنظیم نکردن عدد ۱ برای عمق crawl به صورت پیش فرض در نظر گرفته می شود.

entity	Output_parameters
File_content	<pre> { "results": [{ "data": "/home/dpe/Documents/inprogress_project", "type": 9, "properties": [{ "url": "http://www.mybabyname.com/", "type": 1 }], "ref": { "task": "crawl", "depth": 1 } }] }</pre>

جدول مربوط به پارامترهای خروجی ماژول

۴-۱ محتوای موجود در socket request

برای استفاده از ماژول crawler باید یک درخواست از طریق socket به فریم ورک ارسال شود. بدنه درخواست ارسالی شامل پارامترهای data و type است. پارامتر type مشخص کننده نوع ورودی ماژول (unstructue,file, domain,path,...) است که id مربوط به انواع مختلف ورودی ها و توضیحات مربوطه که در این [لینک](#) موجود است، در این فیلد قرار می گیرد. پارامتر data هم همان داده ارسالی است که شامل پارامترهای ورودی ماژول می باشد.

```

{
  "data": {
    "url": "http://www.mybabyname.com/",
    "depth": 2,
  },
  "type": 9
}
```

نمونه body یک socket request

در تصویر نمونه بالا داریم:

data : داده ارسالی برای استفاده در ماژول

url : آدرس url درخواستی جهت crawl

depth : عمق درخواستی برای crawl

type : نوع ورودی ماژول را مشخص می کند که می تواند از نوع unstructured,path,file و یا انواع دیگر باشد و باید id ورودی موردنظر در این قسمت قرار گیرد.

توجه !

دقت شود که در جدول بالا پارامتر type شماره id مربوط ورودی ماژول را نشان می دهد. و چون ورودی ماژول url است ، مقدار این پارامتر همیشه 9 می باشد.

۵-۱ جدول خطاها

جدول زیر شامل خطاهایی است که ماژول crawler در پاسخ به کلاینت بر می گرداند. خطا ها با شماره مشخصی نشان داده می شوند.

کد	نام	توضیحات
۱۰۲	InvalidInputError	فرمت ورودی نادرست است
۱۱۹	CrawlLimit	صفحات ذخیره شده از تعداد مجاز عبور کرده است.
۱۲۱	TimeoutError	زمانی که hrobot نتواند در مدتی مشخص (۶۰ ثانیه) وارد صفحه ای شود
۱۲۲	ResultNotSetError	زمانی که نتیجه ای برای خروجی set نشده باشد.
۱۱۴	ResultNotFoundError	زمانی که نتیجه ای یافت نشود.

۵-۱ progressbar

برای نمایش progressbar به صورت real time. محل به روز رسانی progressbar درون ماژول است. بدین منظور هر قسمت از ماژول که قسمتی از نتیجه را محاسبه می کند. قبل یا بعد از آن باید update_progressbar صدا زده شود.

در ماژول crawler از آن جایی که نمی توان با اطمینان از مقدار کار باقی مانده صحبت کرد. در قسمت مربوط به پیام این ماژول عمق فعلی و تعداد صفحات ذخیره شده و تعداد صفحات باقی مانده در این عمق چاپ می شود. درصد پیشرفت نیز حاصل تقسیم تعداد صفحات ذخیره شده بر تعداد لینک های یافت شده تا این لحظه می باشد. بدیهی است که با زیاد شدن تعداد لینک های پیدا شده در یک لحظه ممکن است درصد پیشرفت کاهش یابد.