

geturls.py

概要

第一引数に指定した URL の Web ページにアクセスし、ページ内に含まれているリンクをたどって行く。

アクセスした URL を標準出力に出力する。

第二引数は最大で何階層までのページリンクをたどるかを指定する。スクリプトに指定した URL は 0 階層目。そのページ内に含まれるリンクが 1 階層目。1 階層目のページに含まれていたリンクは 2 階層目、と数えて何階層までのページリンクをたどるかを指定する。デフォルトは 5 が指定されている（5 階層目までのページリンクをたどる、6 階層目はたどらない）。

スクリプトの引数で指定できるパラメーターは上記の 2 つだけ。あと、いくつかスクリプト内でユーザーが調整できるパラメーターが存在する。詳細はソース参照。

このスクリプトは robots.txt に対応している。デフォルトで 10000 個の robots.txt 設定をキャッシュし、禁止されている URL は GET せずにスキップする。

拡張子が正規表現 gif|png|bmp|jpg|jpeg|img|mov|mpg|mpeg|avi|exe|gz|zip|css|pdf|rdf にマッチするコンテンツにはアクセスしない。これはスクリプト内のカスタムパラメーター (EXC_EXTENSIONS_RESTRICT) で指定しているので、ユーザーが追加できる。