

```

1 import sqlite3
2 import urllib.error
3 import ssl
4 from urllib.parse import urljoin
5 from urllib.parse import urlparse
6 from urllib.request import urlopen
7 from bs4 import BeautifulSoup
8
9 # Ignore SSL certificate errors
10 ctx = ssl.create_default_context()
11 ||
12 ctx.check_hostname = False
13 ctx.verify_mode = ssl.CERT_NONE
14
15 conn = sqlite3.connect('spider.sqlite')
16 cur = conn.cursor()
17
18 cur.execute('''CREATE TABLE IF NOT EXISTS Pages
19     (id INTEGER PRIMARY KEY, url TEXT UNIQUE, html TEXT,
20      error INTEGER, old_rank REAL, new_rank REAL)'''')
21
22 cur.execute('''CREATE TABLE IF NOT EXISTS Links
23     (from_id INTEGER, to_id INTEGER, UNIQUE(from_id, to_id))'''')
24
25 cur.execute('''CREATE TABLE IF NOT EXISTS Webs (url TEXT UNIQUE)'''')
26
27 # Check to see if we are already in progress...
28 cur.execute('SELECT id,url FROM Pages WHERE html is NULL and error is NULL ORDER BY
29 RANDOM() LIMIT 1')
30 row = cur.fetchone()
31 if row is not None:
32     print("Restarting existing crawl. Remove spider.sqlite to start a fresh crawl.")
33 else :
34     starturl = input('Enter web url or enter: ')
35     if ( len(starturl) < 1 ) : starturl = 'http://www.dr-chuck.com/'
36     if ( starturl.endswith('/') ) : starturl = starturl[:-1]
37     web = starturl
38     if ( starturl.endswith('.htm') or starturl.endswith('.html') ) :
39         pos = starturl.rfind('/')
40         web = starturl[:pos]
41
42     if ( len(web) > 1 ) :
43         cur.execute('INSERT OR IGNORE INTO Webs (url) VALUES ( ? )', ( web, ) )
44         cur.execute('INSERT OR IGNORE INTO Pages (url, html, new_rank) VALUES ( ?, ?, 1.0 )', ( starturl, ) )
45         conn.commit()
46
47 # Get the current webs
48 cur.execute('''SELECT url FROM Webs''')
49 webs = list()
50 for row in cur:
51     webs.append(str(row[0]))
52
53 print(webs)
54 many = 0
55 while True:
56     if ( many < 1 ) :
57         sval = input('How many pages: ')

```

```

58         if ( len(sval) < 1 ) : break
59     many = int(sval)
60     many = many - 1
61
62     cur.execute('SELECT id,url FROM Pages WHERE html is NULL and error is NULL ORDER
63 BY RANDOM() LIMIT 1')
64     try:
65         row = cur.fetchone()
66         # print row
67         fromid = row[0]
68         url = row[1]
69     except:
70         print('No unretrieved HTML pages found')
71         many = 0
72         break
73
74     print(fromid, url, end=' ')
75
76     # If we are retrieving this page, there should be no links from it
77     cur.execute('DELETE from Links WHERE from_id=?', (fromid, ) )
78     try:
79         document = urlopen(url, context=ctx)
80
81         html = document.read()
82         if document.getcode() != 200 :
83             print("Error on page: ",document.getcode())
84             cur.execute('UPDATE Pages SET error=? WHERE url=?', (document.getcode(),
85 url) )
86
87         if 'text/html' != document.info().get_content_type() :
88             print("Ignore non text/html page")
89             cur.execute('DELETE FROM Pages WHERE url=?', ( url, ) )
90             conn.commit()
91             continue
92
93             print('('+str(len(html))+')', end=' ')
94
95             soup = BeautifulSoup(html, "html.parser")
96     except KeyboardInterrupt:
97         print('')
98         print('Program interrupted by user...')
99         break
100    except:
101        print("Unable to retrieve or parse page")
102        cur.execute('UPDATE Pages SET error=-1 WHERE url=?', (url, ) )
103        conn.commit()
104        continue
105
106        cur.execute('INSERT OR IGNORE INTO Pages (url, html, new_rank) VALUES ( ?, NULL,
1.0 )', ( url, ) )
107        cur.execute('UPDATE Pages SET html=? WHERE url=?', (memoryview(html), url ) )
108        conn.commit()
109
110        # Retrieve all of the anchor tags
111        tags = soup('a')
112        count = 0
113        for tag in tags:
114            href = tag.get('href', None)
115            if ( href is None ) : continue
116            # Resolve relative references like href="/contact"

```

```
115     up = urlparse(href)
116     if ( len(up.scheme) < 1 ) :
117         href = urljoin(url, href)
118     ipos = href.find('#')
119     if ( ipos > 1 ) : href = href[:ipos]
120     if ( href.endswith('.png') or href.endswith('.jpg') or href.endswith('.gif') )
121         ) : continue
122         if ( href.endswith('/') ) : href = href[:-1]
123         # print href
124         if ( len(href) < 1 ) : continue
125
126     # Check if the URL is in any of the webs
127     found = False
128     for web in webs:
129         if ( href.startswith(web) ) :
130             found = True
131             break
132     if not found : continue
133
134     cur.execute('INSERT OR IGNORE INTO Pages (url, html, new_rank) VALUES ( ?, NULL, 1.0 )', ( href, ) )
135     count = count + 1
136     conn.commit()
137
138     cur.execute('SELECT id FROM Pages WHERE url=? LIMIT 1', ( href, ) )
139     try:
140         row = cur.fetchone()
141         toid = row[0]
142     except:
143         print('Could not retrieve id')
144         continue
145     # print fromid, toid
146     cur.execute('INSERT OR IGNORE INTO Links (from_id, to_id) VALUES ( ?, ? )', ( fromid, toid ) )
147
148     print(count)
149
150 cur.close()
151
```