

'''

evolving-search-engine-using-genetic-algorithm

<https://code.google.com/p/evolving-search-engine-using-genetic-algorithm/>

Our venture presents a new method of developing a search engine using Genetic Algorithm. It takes a word, a phrase etc as input and lists the URLs containing the mentioned content.

The project is being implemented in three steps namely web crawling, web indexing and searching. A Web crawler is an Internet bot that systematically browses the World Wide Web, typically for the purpose of Web indexing. Web indexing refers to various methods for indexing the contents of a website or of the Internet as a whole. A genetic algorithm is a search heuristic that mimics the process of natural evolution. Genetic algorithm can solve every optimization problem, it can also solve problems with multiple solutions. The project is applicable in all fields related to searching and finding the optimal solution.

'''

from StringIO import StringIO

import pycurl

import BeautifulSoup

import re

import MySQLdb as mdb

import urllib2

def visible(element):

if element.parent.name in ['style', 'script', '[document]', 'head', 'title']:

return False

elif re.match('<!--.*-->', element):

return False

return True

url=raw_input("url ? <http://url>: ")

page=urllib2.urlopen(url)

content=page.read()

soup = BeautifulSoup.BeautifulSoup(content)

texts = soup.findAll(text=True)

visible_texts = filter(visible, texts)

print visible_texts

save=""

for word in visible_texts:

save+=word

print save

save=save.encode('ascii','ignore')

print "Saving to Database Now..."

con = mdb.connect('localhost', 'root','foobar', 'testdb')

with con:

print "Connection Established."

cur=con.cursor()

query="INSERT INTO urllist (url,content) VALUES ('"+url+"','"+save+"')"

cur.execute(query)

cur.execute("SELECT * FROM urllist")

rows=cur.fetchall()

for row in rows:

print row