# Program No:15

**Aim:**Implement a simple web crawler(csv).

**Code:**

```python
from bs4 import BeautifulSoup
import requests
pages_crawled = []
def crawler(url):
    page = requests.get(url)
    soup = BeautifulSoup(page.text, 'html.parser')
    links = soup.find_all('a')
    for link in links:
        if 'href' in link.attrs:
            if link['href'].startswith('/wiki') and ':' not in link['href']:
                if link['href'] not in pages_crawled:
                    new_link = f"https://en.wikipedia.org{link['href']}"
                    pages_crawled.append(link['href'])
                    try:
                        with open('data.csv', 'a') as file:
                            file.write(f'{soup.title.text}; {soup.h1.text}; {link["href"]}\n')
                        crawler(new_link)
                    except:
                        continue
crawler('https://en.wikipedia.org')
```

**Output:**

| | |
|---|---|
| Wikipedia | the fr |
| Wikipedia - Wikipedia; Wikipedia; /wiki/Main_Page | |
| Wikipedia | the fr |
| Free content - Wikipedia; Free content; /wiki/Definition_of_Free_Cultural_Works | |
| Definition of Free Cultural Works - Wikipedia; Definition of Free Cultural Works; /wiki/Free_content_movement | |
| Free-culture movement - Wikipedia; Free-culture movement; /wiki/Free_culture_(disambiguation) | |
| Free Culture - Wikipedia; Free Culture; /wiki/Free_Culture_(book) | |
| Free Culture (book) - Wikipedia; Free Culture (book); /wiki/Lawrence_Lessig | |
| Lawrence Lessig - Wikipedia; Lawrence Lessig; /wiki/Lawrence_Lessing | |
| Lawrence Lessing - Wikipedia; Lawrence Lessing; /wiki/Science_writer | |
| Science journalism - Wikipedia; Science journalism; /wiki/Scientific_journalism | |
| Scientific journalism - Wikipedia; Scientific journalism; /wiki/Science_journalism | |
| Science journalism - Wikipedia; Science journalism; /wiki/Scientific_writing | |
| Scientific writing - Wikipedia; Scientific writing; /wiki/Science_writing | |
| Science journalism - Wikipedia; Science journalism; /wiki/Science_communication | |
| Science communication - Wikipedia; Science communication; /wiki/Science_publishing | |
| Scientific literature - Wikipedia; Scientific literature; /wiki/Medical_literature | |
| Medical literature - Wikipedia; Medical literature; /wiki/Edwin_Smith_Papyrus | |
| Edwin Smith Papyrus - Wikipedia; Edwin Smith Papyrus; /wiki/New_York_Academy_of_Medicine | |
| New York Academy of Medicine - Wikipedia; New York Academy of Medicine; /wiki/Eclecticism_in_architecture | |
| Eclecticism in architecture - Wikipedia; Eclecticism in architecture; /wiki/Basilica | |
| Basilica - Wikipedia; Basilica; /wiki/Basilicas_in_the_Catholic_Church | |
| Basilicas in the Catholic Church - Wikipedia; Basilicas in the Catholic Church; /wiki/List_of_Catholic_basilicas | |
| List of Catholic basilicas - Wikipedia; List of Catholic basilicas; /wiki/Catholic_Church | |
| Catholic Church - Wikipedia; Catholic Church; /wiki/Catholic_Church_(disambiguation) | |
| Catholic Church (disambiguation) - Wikipedia; Catholic Church (disambiguation); /wiki/Catholic_(disambiguation) | |

| | |
|---|---|
| Wikipedia | t |
| Encyclopedia - Wikipedia; Encyclopedia; /wiki/Online_encyclopedia | |
| Wikipedia | t |
| Wikipedia - Wikipedia; Wikipedia; /wiki/Main_Page | |
| Wikipedia | t |
| Free content - Wikipedia; Free content; /wiki/Definition_of_Free_Cultural_Works | |
| Definition of Free Cultural Works - Wikipedia; Definition of Free Cultural Works; /wiki/Free_content_movement | |
| Free-culture movement - Wikipedia; Free-culture movement; /wiki/Free_culture_(disambiguation) | |
| Free Culture - Wikipedia; Free Culture; /wiki/Free_Culture_(book) | |
| Free Culture (book) - Wikipedia; Free Culture (book); /wiki/Lawrence_Lessig | |
| Lawrence Lessig - Wikipedia; Lawrence Lessig; /wiki/Lawrence_Lessing | |
| Lawrence Lessing - Wikipedia; Lawrence Lessing; /wiki/Science_writer | |
| Science journalism - Wikipedia; Science journalism; /wiki/Scientific_journalism | |
| Scientific journalism - Wikipedia; Scientific journalism; /wiki/Science_journalism | |
| Science journalism - Wikipedia; Science journalism; /wiki/Scientific_writing | |
| Scientific writing - Wikipedia; Scientific writing; /wiki/Science_writing | |
| Science journalism - Wikipedia; Science journalism; /wiki/Science_communication | |
| Science communication - Wikipedia; Science communication; /wiki/Science_publishing | |
| Scientific literature - Wikipedia; Scientific literature; /wiki/Medical_literature | |
| Medical literature - Wikipedia; Medical literature; /wiki/Edwin_Smith_Papyrus | |
| Edwin Smith Papyrus - Wikipedia; Edwin Smith Papyrus; /wiki/New_York_Academy_of_Medicine | |
| New York Academy of Medicine - Wikipedia; New York Academy of Medicine; /wiki/Eclecticism_in_architecture | |
| Eclecticism in architecture - Wikipedia; Eclecticism in architecture; /wiki/Basilica | |
| Basilica - Wikipedia; Basilica; /wiki/Basilicas_in_the_Catholic_Church | |
| Basilicas in the Catholic Church - Wikipedia; Basilicas in the Catholic Church; /wiki/List_of_Catholic_basilicas | |
| List of Catholic basilicas - Wikipedia; List of Catholic basilicas; /wiki/Catholic_Church | |