

Abstract geometric lines in the top-left corner of the page, consisting of several thin black lines forming various polygons and intersecting at different points.

WWII WIKIPEDIA SCRAPER

Berjer Ding
Toronto, Canada

PREMISE

- Scrapes information from Wikipedia
- Written in Python
- Centered around World War 2 pages as boundary
- Testable

```
project.py X test_project.py
1 import wikipediaapi, re, random, sys
2
3
4 wiki = wikipediaapi.Wikipedia('en')
5
6
7 def main():
8     x = input("Would you like to know about a [specific] event, a [random] event, or an event in a [general] year? ")
9     if x.lower() == "specific":
10         print(specific(input("What would you like to know about? ")))
11     elif x.lower() == "random":
12         randomlist = []
13         battles = wiki.page("List of World War II battles").links
14         for link in battles.keys():
15             randomlist.append(link)
16         i = random.randint(0, len(randomlist))
17         print(random_event(i))
18     elif x.lower() == "general":
19         print(general_event())
20     else:
21         print("Please enter [specific] or [random] or [general] to continue")
22
23
24 def specific_event(i):
25     page = wiki.page(i)
26     if page.exists():
27         if matches := re.search(r"1939|194[0-5]", page.summary):
28             return page.summary
29         else:
30             return "Page is not about the Second World War or WW2-adjacent."
31     else:
32         return "Page does not exist"
33
34
35 def random_event(x):
36     randomlist = []
37     battles = wiki.page("List of World War II battles").links
38     for link in battles.keys():
39         randomlist.append(link)
40     battle_page = wiki.page(randomlist[x])
41     return battle_page.summary
42
43
44 def general_event():
45     generallist = []
46     portal = wiki.page("Category:World War II").categorymembers
47     for category in portal.keys():
48         generallist.append(category.replace("Category:", ""))
49     print(generallist)
50     while True:
51         inputted = input("")
52         if inputted in generallist:
53             generallist = []
54             portal = wiki.page(f"Category:{inputted}").categorymembers
55             for category in portal.keys():
56                 generallist.append(category.replace("Category:", ""))
57             if generallist == []:
58                 return wiki.page(inputted).summary
59             else:
60                 print(generallist)
61         else:
62             sys.exit("Incorrect selection")
63
64
65 if __name__ == "__main__":
66     main()
```

PRELIMINARIES

Importing 4 libraries:

1. Wikipedia API
2. Regular Expressions
3. Random Module
4. Sys Module

Shorten the Wikipedia API to the wiki keyword and in English

```
1 import wikipediaapi, re, random, sys
2
3
4 wiki = wikipediaapi.Wikipedia('en')
```

```

7 def main():
8     x = input("Would you like to know about a |specific| event, a |random| event, or an event in a |general| year? ")
9     if x.lower() == "specific":
10        print(specific(input("What would you like to know about? ")))
11    elif x.lower() == "random":
12        randomlist = []
13        battles = wiki.page("List of World War II battles").links
14        for link in battles.keys():
15            randomlist.append(link)
16        i = random.randint(0, len(randomlist))
17        print(random_event(i))
18    elif x.lower() == "general":
19        print(general_event())
20    else:
21        print("Please enter |specific| or |random| or |general| to continue")
22

```

MAIN() FUNCTION

- For the user to get to the individual functions
- Three options: specific, random, and general
- Hardcoded to input the options

```
24 def specific_event(i):
25     page = wiki.page(i)
26     if page.exists():
27         if matches := re.search(r"1939|194[0-5]", page.summary):
28             return page.summary
29         else:
30             return "Page is not about the Second World War or WW2-adjacent."
31     else:
32         return "Page does not exist"
33
```

SPECIFIC_EVENT(I)

- Easiest to make and test
- Receives as input “i” which is a topic the user inputted in main()
- Wikipedia API returns the page if it exists
- Re.search confirms if it is about WW2
- Returns summary of the page

RANDOM_EVENT(X)

- Takes as input a randomly generated number from main()
- Selects from the wiki page of WWII battles
- Uses the wiki page links as a list
- Returns a randomly selected page summary

```
34
35 def random_event(x):
36     randomlist = []
37     battles = wiki.page("List of World War II battles").links
38     for link in battles.keys():
39         randomlist.append(link)
40     battle_page = wiki.page(randomlist[x])
41     return battle_page.summary
42
```

```

43
44 def general_event():
45     generallist = []
46     portal = wiki.page("Category:World War II").categorymembers
47     for category in portal.keys():
48         generallist.append(category.replace("Category:", ""))
49     print(generallist)
50     while True:
51         inputted = input("")
52         if inputted in generallist:
53             generallist = []
54             portal = wiki.page(f"Category:{inputted}").categorymembers
55             for category in portal.keys():
56                 generallist.append(category.replace("Category:", ""))
57             if generallist == []:
58                 return wiki.page(inputted).summary
59             else:
60                 print(generallist)
61         else:
62             sys.exit("Incorrect selection")
63

```

GENERAL_EVENT()

- Selects from the WWII category portal
- First captures the available categories on the page to output
- In a while-loop, allows user input to delve deeper into the categories
- Ends when a page is found, which is then outputted

```

1 import pytest, project, mock, builtins
2
3 def test_specific_event():
4     assert project.specific_event("List of French divisions in World War II") == "This is a listing of French divisions that served between 1939 and 1945."
5     assert project.specific_event("Cat") == "Page is not about the Second World War or WW2-adjacent."
6     assert project.specific_event("BerjerDing") == "Page does not exist"
7
8 def test_random_event():
9     assert project.random_event(5) == "The Action at Bir el Gubi (November 1941) (Arabic: بئر الغبي, romanized: Bīr al-Ġubbiyy, lit. 'Well of the Depressed Terrain',
10
11
12 def test_general_event():
13     with mock.patch.object(builtins, 'input', lambda _: 'USCGC Northland (WPG-49)'):
14         assert project.general_event() == 'USCGC Northland (WPG-49) was a United States Coast Guard cruising class of gunboat especially designed for Arctic operations
15     with mock.patch.object(builtins, 'input', lambda _: 'Cat'):
16         with pytest.raises(SystemExit):
17             project.general_event()
18     with mock.patch.object(builtins, 'input', lambda _: 'BerjerDing'):
19         with pytest.raises(SystemExit):
20             project.general_event()

```

TESTING

- Testing `specific_event()` was confirming the inputted page matched the output summary
- Testing `random_event()` was through the parameter, which matched the argument page to the out page
- Testing `general_event()` was harder:
 - Required mocking the input of the user
 - Asserting that the mocked input returned the correct output