

Projekt zaliczeniowy

Proces ETL

Grupa projektowa: 25

Imię	Nazwisko	Numer albumu	Grupa działańska	Wkład w prace nad projektem ¹	Udział procentowy
Ewa	Bąk	192207	WZISS2-1111	Stworzenie skryptu ETL, Frontend, komunikacja między podstronami	
Gabriela	Lenard	193670	WZISS2-1111	PostgreSQL	
Karol	Skoczyk	194677	WZISS2-1111	Python, Flask	

___/70 pkt

¹ proszę wymienić konkretne zadania

Nazwy i wersje użytych technologii:

- Język programowania wykorzystany w skrypcie - Python 3.7.1
- Web framework - Flask 1.0.2
- Framework - Bootstrap
- System zarządzania bazą danych - PostgreSQL

Informacje na temat środowiska:

- Wykorzystane biblioteki:
 - bs4
 - requests
 - pandas
 - psycopg2
- Wymagania sprzętowe: procesor taktowany częstotliwością co najmniej 1 GHz, 0,5 GB pamięci na dysku twardym, 1 GB pamięci operacyjnej.

Linki do oprogramowania tworzącego środowisko:

- <https://www.python.org/downloads/release/python-371/>
- <https://pypi.org/project/Flask/1.0.2/>
- <https://pypi.org/project/beautifulsoup4/>
- <https://pypi.org/project/pandas/>
- <https://pypi.org/project/requests/>
- <https://getbootstrap.com/>
- <https://www.postgresql.org/download/>

Struktura aplikacji:

- scrape.py - plik wywołujący funkcje zadane w poleceniu projektu - pełny proces ETL oraz osobne procesy E, T oraz L.
- index.html - plik zawierający strukturę strony głównej aplikacji
- wynik.html - plik zawierający strukturę podstrony po wybraniu opcji pełnego procesu ETL
- extract.html - plik zawierający strukturę podstrony po wybraniu opcji Extract, która zawiera wyodrębnione dane w nieprzetworzonej postaci
- transform.html - plik zawierający strukturę podstrony po wybraniu opcji Transform, która zawiera dane w postaci tabeli

- clean.html - plik zawierający strukturę podstrony po wyczyszczeniu bazy danych
- style.css - plik zawierający szczegółowe opcje dotyczące wyglądu aplikacji.

Nazwa pliku	Nazwa zmiennej	Opis zmiennej
scrape.py	cena_p	cena początkowa oferty, która jest pobierana od użytkownika
	cena_k	cena końcowa oferty, która jest pobierana od użytkownika
	nazwy nazwa	przechowuje nazwę/tytuł ogłoszenia
	dzielnice podpis	przechowuje dokładniejszą lokalizację miejsca z ogłoszenia
	pokoj pokoje	przechowuje ilość pokoi, które oferuje miejsce z ogłoszenia
	metry metryy	przechowuje ilość metrów kwadratowych mieszkania w ofercie
	cena_metr cena_meryy	przechowuje cenę jednego metra kwadratowego mieszkania w ofercie
	cena ceny	przechowuje cenę mieszkania zamieszczonego na stronie
	page	przechowuje adres strony, z której pobierane są dane
	html	pobiera adres strony
	soup	definiuje podstawowy interfejs strony
	mieszkanie	????????????????????
	df	wyświetla dane pobrane ze strony internetowej
	count	przechowuje numer nowo dodanych do bazy rekordów
	datat	przechowuje wiersze pobrane z tabeli
	column_names	Przechowuje nagłówki tabeli
	dt	Wyświetla tabelę zawierającą pobrane rekordy

```

15     if request.method == 'POST':
16         if request.form['submit_button'] == 'etl':
17             cena_p = request.form['od']
18             cena_k = request.form['do']
19
20             csv_file = open('cms_scrape.csv', 'w', encoding='utf-8', errors = 'ignore')
21             csv_writer = csv.writer(csv_file)
22             csv_writer.writerow(['nazwa', 'dzielnica', 'pokoj', 'metry', 'cena_metr', 'cena'])
23
24             nazwy=[]
25             dzielnice=[]
26             pokoje=[]
27             metry=[]
28             cena_metry=[]
29             ceny=[]
30             count = 0
31
32             connection = psycopg2.connect(user = "postgres", password = "haslo", database = "postgres")
33             cursor = connection.cursor()
34             sql_add = """INSERT INTO otodom ("Tytul", "Dzielnica", "Liczba pokoi", "Metraz", "Cena za metr", "Cena") VALUES (%s,%s,%s,%s,%s,%s) ON CONFLICT DO NOTHING"""
35             sql_display = """SELECT * FROM otodom"""
36
37             for i in range(1, 10):
38                 page = "https://www.otodom.pl/sprzedaz/mieszkanie/krakow/?search%5Bfilter_float_price%3Afrom%5D="+ cena_p + "&search%5Bfilter_float_price%3Ato%5D="+ cena_k + "&page={}".format(i)
39                 html = requests.get(page)
40                 soup = BeautifulSoup(html.text, 'lxml')
41
42                 for mieszkanie in soup.find_all('div', class_='offer-item-details'):
43                     nazwa = mieszkanie.find('span', class_='offer-item-title').text
44                     nazwy.append(nazwa)
45
46                     podpis = mieszkanie.find('p', class_='text-nowrap hidden-xs').text
47                     dzielnice = podpis.split(':')[1]
48                     dzielnice.append(dzielnice)
49
50                     pokoj = mieszkanie.find('li', class_='offer-item-rooms hidden-xs').text
51                     pokoje.append(pokoj)
52
53                     metry = mieszkanie.find('li', class_='hidden-xs offer-item-area').text
54                     metry.append(metry)
55
56                     cena_metr = mieszkanie.find('li', class_='hidden-xs offer-item-price-per-m').text
57                     cena_metry.append(cena_metr)
58
59                     cena = mieszkanie.find('li', class_='offer-item-price').text.replace(' ', '').replace('\n', '')
60                     ceny.append(cena)
61
62                     to_insert = (nazwa, dzielnice, pokoj, metry, cena_metr, cena)
63                     cursor.execute(sql_add, to_insert)
64                     connection.commit()
65

```

```

67         if cursor.rowcount==1:
68             count += 1
69
70             csv_writer.writerow([nazwa, dzielnica, pokoj, metry, cena_metr, cena])
71
72             df = pd.DataFrame({'Nazwa':nazwy, 'Dzielnica':dzielnice, 'Pokoj':pokoje, 'Metry':metry, 'Cena za metr':cena_metry, 'Cena':ceny})
73
74             csv_file.close()
75
76             cursor.execute(sql_display)
77             datat = cursor.fetchall()
78             cursor.close()
79             connection.close()
80
81             return render_template('wynik.html', c1=cena_p, c2=cena_k, data=df, rec=count, datat=datat)
82

```

Funkcja odpowiedzialna za wykonanie pełnego procesu ETL. Zwraca tabelę zawierającą ogłoszenia zgodne z kryteriami użytkownika wpisanymi na stronie startowej aplikacji. Po wyciągnięciu danych ze strony internetowej istnieje również możliwość pobrania ich do pliku CSV. Wygląd strony końcowej warunkuje plik wynik.html.

```
File Edit Selection View Go Debug Terminal Help
scrape.py - Projekt-HD - Visual Studio Code

69 csv_file = open('cms_scrape.csv', 'w', encoding='utf-8', errors = 'ignore')
70 csv_writer = csv.writer(csv_file)
71 csv_writer.writerow(['nazwa', 'dzielnica', 'pokoj', 'metry', 'cena_metr', 'cena'])
72
73 nazwy=[]
74 dzielnice=[]
75 pokoje=[]
76 metry=[]
77 cena_metry=[]
78 ceny=[]
79
80
81 for i in range(1, 10):
82     page = "https://www.otodom.pl/sprzedaz/mieszkanie/krakow/?search%5Bfilter_float_price%3Afrom%5D="+ cena_p +"&search%5Bfilter_float_price%3Ato%5D="+ cena_k +"&page="
83     html = requests.get(page)
84     soup = BeautifulSoup(html.text, 'lxml')
85
86     for mieszkanie in soup.find_all('div', class_='offer-item-details'):
87
88         nazwa = mieszkanie.find('span', class_='offer-item-title')
89         nazwy.append(nazwa)
90
91         podpis = mieszkanie.find('p', class_='text-nowrap hidden-xs')
92         dzielnice.append(podpis)
93
94         pokoj = mieszkanie.find('li', class_='offer-item-rooms hidden-xs')
95         pokoje.append(pokoj)
96
97         metry = mieszkanie.find('li', class_='hidden-xs offer-item-area')
98         metry.append(metry)
99
100         cena_metr = mieszkanie.find('li', class_='hidden-xs offer-item-price-per-m')
101         cena_metry.append(cena_metr)
102
103         cena = mieszkanie.find('li', class_='offer-item-price')
104         ceny.append(cena)
105
106         #print(nazwa, podpis, pokoj, metry, cena_metr, cena)
107
108         df = (nazwy, dzielnice, pokoje, metry, cena_metry, ceny)
109
110     return render_template('extract.html', data=df)
```

Funkcja odpowiedzialna za wykonanie procesu Extract. Struktura kodu jest niemal identyczna jak w poprzednim przypadku z różnicą na wyświetlenie danych, które są przedstawione w formie nieprzetworzonej. Wygląd strony warunkuje plik extract.html.

```
@app.route('/csv_file')
def csv_file():
    return send_file('cms_scrape.csv', attachment_filename='cms_scrape.csv', as_attachment=True)
```

Skrypt odpowiedzialny za określenie ścieżki dostępu, która umożliwia pobranie pliku z danymi ogłoszeń w formacie CSV.

```
140 @app.route('/clean')
141 def clean():
142     connection = psycopg2.connect(user = "postgres", password = "haslo", database = "postgres")
143     cursor = connection.cursor()
144     sql_del = """DELETE FROM otodom"""
145
146     cursor.execute(sql_del)
147     connection.commit()
148     cursor.close()
149     connection.close()
150     return render_template('clean.html')
151
```

Skrypt odpowiedzialny za usunięcie wszystkich rekordów z bazy danych

Instrukcja obsługi aplikacji:

Instrukcja uruchomienia aplikacji:

1. Pobranie wszystkich elementów środowiska aplikacji
2. Pobranie skompresowanego pliku ZIP

3. Wyodrębnienie plików z pliku ZIP
4. Uruchomienie aplikacji poprzez:
 - 4.1. Uruchomienie skryptu do tworzenia tabeli w programie PostgreSQL
 - 4.2. Otworzenie folderu aplikacji w edytorze kodu takim jak na przykład Visual Studio Code
 - 4.2.1. Uruchomienie pliku scrape.py
 - 4.3. Lokalizacja folderu aplikacji za pomocą Windows Command Line/ konsoli Linux
 - 4.3.1. Uruchomienie pliku scrape.py
5. ...

Opis funkcjonalności aplikacji:

Po wykonaniu powyższych kroków, edytor kodu lub konsola dowolnego systemu skieruje użytkownika do strony głównej aplikacji w przeglądarce internetowej.

SCREEN STRONY GŁÓWNEJ

Wyświetlone są na niej pola do uzupełnienia kryteriów cenowych przy poszukiwaniu mieszkania, a także przycisk wykonania procesu ETL oraz Extract. W górnym rogu znajdują się również odnośniki do dokumentacji oraz do repozytorium z kodem aplikacji.

SCREEN KLIKNIECIA ETL

Wybór wykonania całego procesu ETL skutkuje przejściem do podstrony, w której znajduje się tabela z wyszczególnionymi ofertami wedle (lub bez podanego) kryterium wyboru. Oprócz ceny, w tabeli widnieje nazwa ogłoszenia, dzielnica, w jakiej znajduje się oferowana nieruchomość, liczba pokoi, metraż mieszkania oraz jego cena za metr kwadratowy.

ZBLIŻENIE NA PRZYCISK CSV W OSOBNYM SCREENIE XD

Przycisk CSV służy do pobrania zebranych danych w tym kroku i zapisu ich na dysk komputera. Możliwe jest przeglądanie tych danych i edytowanie według własnych dowolnych potrzeb w programie typu Excel.

ZBLIŻENIE NA PRZYCISK CSV W OSOBNYM SCREENIE MOZE TEŻ???

POUZUPEŁNIAJCIE JAK WYGLĄDA TO Z PRZYCISKIEM BAZY DANYCH I CO ROBI FAJNIUTKIEGO ;) JAK PRZECHODZI GDZIEŚ DALEJ TO PISZTA

SCREEN KLIKNIECIA EXTRACT

Wracając do ekranu głównego, wybór opcji Extract przeniesie użytkownika do podstrony, w której znajdują się pobrane ze strony internetowej dane bez przetworzenia.

Dopiero w kolejnym etapie aplikacji zebrane informacje ulegają segregacji i przetworzeniu.
XDDD :)

Opis scenariuszy użycia aplikacji:

Opis widoków okna aplikacji:

Model danych użytych w projekcie: