

# Projekt zaliczeniowy

## Proces ETL

### Grupa projektowa: 25

Imię	Nazwisko	Numer albumu	Grupa dziekańska	Wkład w prace nad projektem <sup>1</sup>	Udział procentowy
Ewa	Bąk	192207	WZISS2-1111	Stworzenie skryptu ETL, Frontend, komunikacja między podstronami, stworzenie formularza dla użytkownika, dokumentacja, stworzenie skryptu Extract	40%
Gabriela	Lenard	193670	WZISS2-1111	Stworzenie tabeli i połączenie z bazą danych, skrypt zapisujący rekordy do tabeli, skrypt wyświetlający zawartość bazy danych, operacje na tabeli, dokumentacja	25%
Karol	Skoczyk	194677	WZISS2-1111	Stworzenie skryptu ETL, Stworzenie pętli pobierającej konkretną ilość danych, stworzenie formularza dla użytkownika, dokumentacja	35%

\_\_\_/70 pkt

<sup>1</sup> proszę wymienić konkretne zadania

## *Dokumentacja techniczna*

---

### *Nazwy i wersje użytych technologii:*

- Język programowania wykorzystany w skrypcie - Python 3.7.1
- Web framework - Flask 1.0.2
- Framework - Bootstrap 4.1
- System zarządzania bazą danych - PostgreSQL 11.1

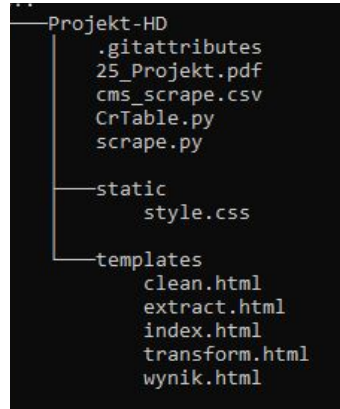
### *Informacje na temat środowiska:*

- 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.
- CrTable.py - plik tworzący tabelę w bazie danych
- 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.

```

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_scraper.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             metryy=[]
28             cena_metryy=[]
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                     dzielnica = podpis.split(':')[1]
48                     dzielnice.append(dzielnica)
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                     metryy.append(metry)
55
56                     cena_metr = mieszkanie.find('li', class_='hidden-xs offer-item-price-per-m').text
57                     cena_metryy.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, dzielnica, pokoj, metry, cena_metr, cena)
63                     cursor.execute(sql_add, to_insert)
64                     connection.commit()
65
66

```

```

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':metryy, 'Cena za metr':cena_metryy, '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:

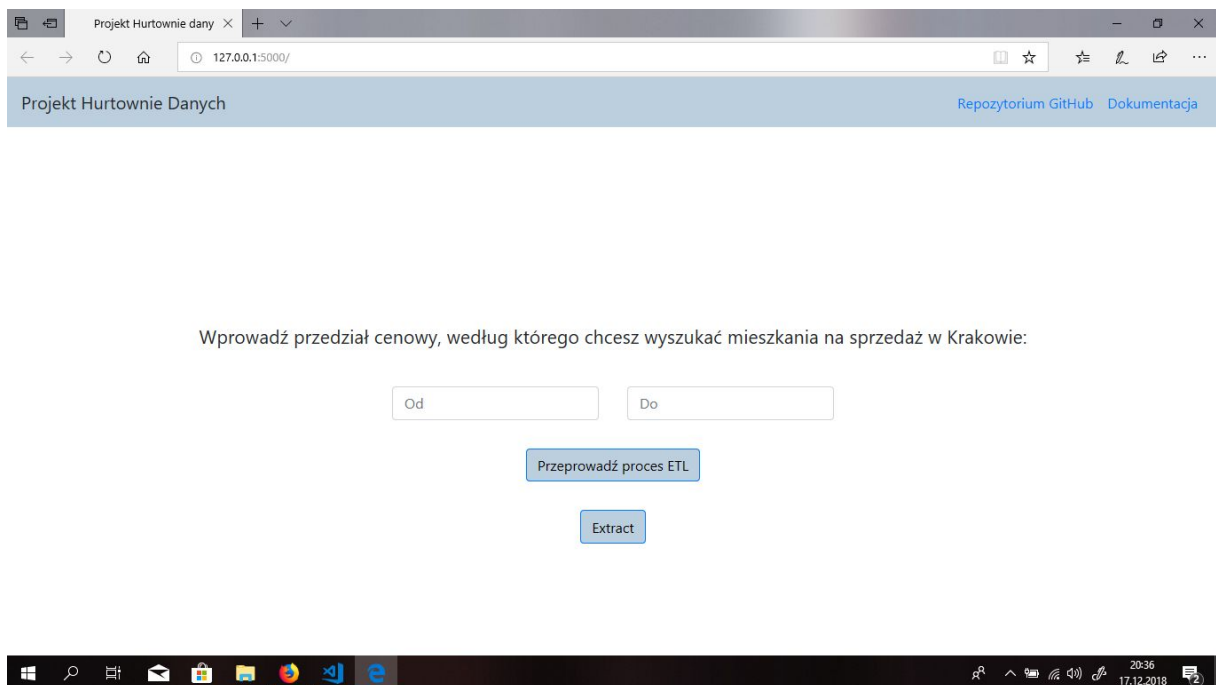
### Instrukcja uruchomienia aplikacji:

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 CrTable.py
    - 4.2.2. Uruchomienie pliku scrape.py
  - 4.3. Lokalizacja folderu aplikacji za pomocą Windows Command Line/ konsoli Linux
    - 4.3.1. Uruchomienie pliku CrTable.py
    - 4.3.2. Uruchomienie pliku scrape.py

### **Opis funkcjonalności aplikacji:**

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



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 pobrania pliku dokumentacji oraz do repozytorium z kodem aplikacji.

Przeprowadź proces ETL

Projekt Hurtownie dany X + -

127.0.0.1:5000/ceny/

Projekt Hurtownie Danych [Repozytorium GitHub](#) [Dokumentacja](#)

Liczba nowych rekordów dodanych do bazy danych: 159

[Pobierz plik CSV](#) [Wyczyść bazę danych](#)

	Tytuł	Dzielnica	Liczba pokoi	Metraz	Cena za metr	Cena
0	2 pokoje - jasna kuchnia - gotowe do wejścia	Kraków, Krzesławice	2 pokoje	41,10 m²	7 056 zł/m²	290000zł
1	Atrakcyjne 100 metrowe poddasze	Kraków, Krowodrza	3 pokoje	107 m²	4 056 zł/m²	434000zł
2	ul.Palacha -Świetne 1pok Parze/Singlom. Krowod...	Kraków, Azory	1 pokój	24,10 m²	8 714 zł/m²	210000zł
3	Stare Miasto - 2 - 3 pokoje - Mega Okazja	Kraków, Grzegórzki	2 pokoje	49 m²	8 469 zł/m²	415000zł

20:38 17.12.2018

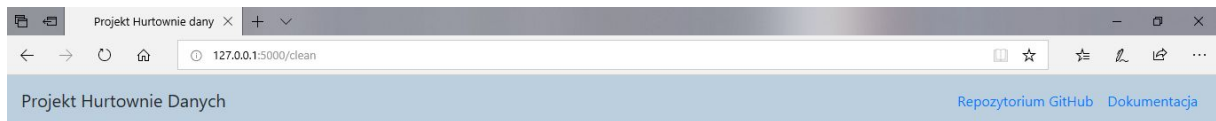
Wybór wykonania całego procesu ETL skutkuje przejściem do podstrony, w której znajduje się tabela z wyszczególnionymi ofertami wedle (lub bez podanego) kryterium wyboru. Oprócz ceny, w tabeli widnieją nazwa ogłoszenia, dzielnica, w jakiej znajduje się oferowana nieruchomość, liczba pokoi, metraż mieszkania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

[Pobierz plik CSV](#)

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

[Wyczyść bazę danych](#)



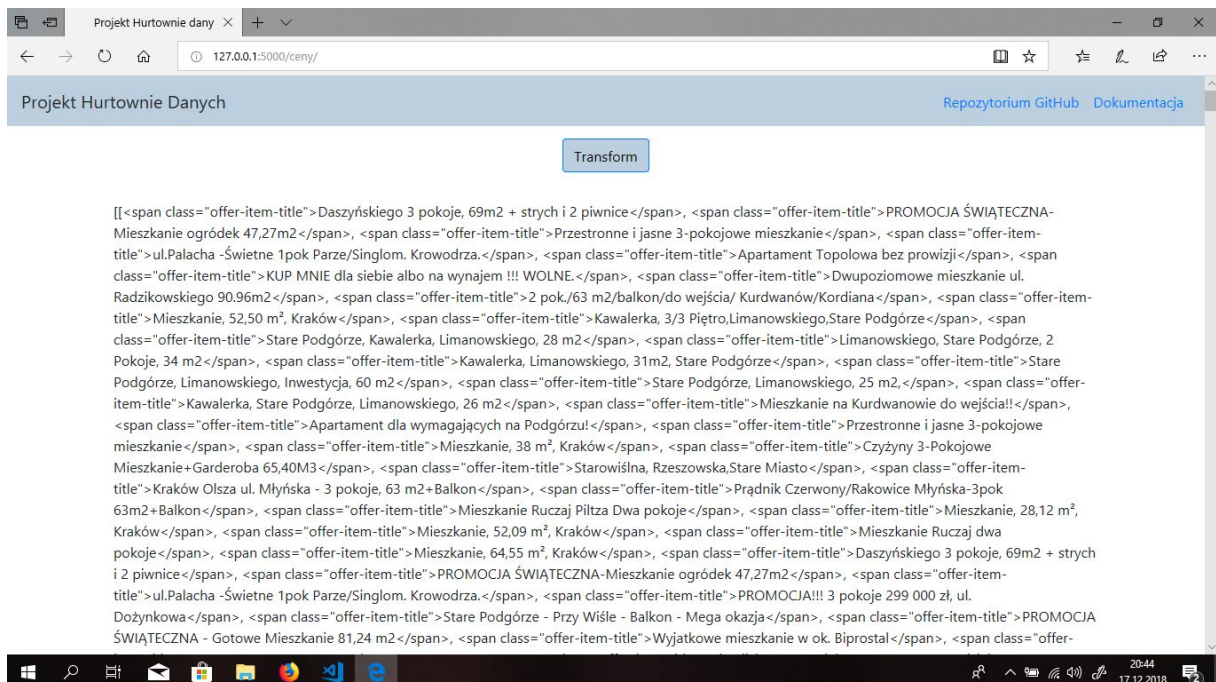


Baza danych została pomyślnie wyczyszczona!



Przycisk “Wyczyść bazę danych” pozwala na usunięcie z bazy danych przechowywanych w niej danych.

Extract



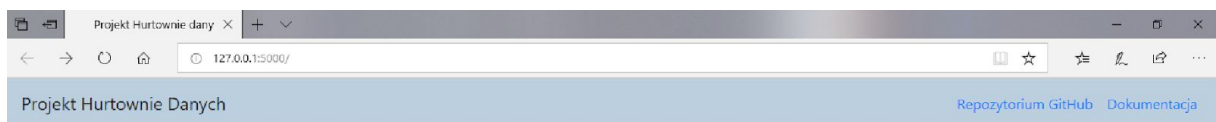
Wracając do ekranu głównego, wybór opcji Extract przeniesie użytkownika do podstrony, w której wykonany zostanie pierwszy etap procesu ETL i wyświetli się jego wynik w formie kodu Html.



### Opis scenariuszy użycia aplikacji:

1. Użytkownik zainteresowany jest wyszukaniem mieszkań na sprzedaż w obrębie miasta Kraków
2. Uzupełnia pola "od" oraz "do" w interesujące go kwoty
3. Naciska przycisk "Przeprowadź proces ETL", strona przekierowuje go do podstrony zawierającej spis ogłoszeń pasujących do kryteriów
4. Użytkownik ma możliwość pobrania pliku CSV z interesującymi go danymi

### Opis widoków okna aplikacji:



5

6

Wprowadź przedział cenowy, według którego chcesz wyszukać mieszkania na sprzedaż w Krakowie:

1   2

3

4



Strona główna

1. Pole formularza ceny początkowej

2. Pole formularza ceny końcowej
3. Przycisk przeprowadzania procesu ETL
4. Przycisk przeprowadzenia procesu Extract
5. Odnośnik do repozytorium projektu
6. Odnośnik do pobrania dokumentacji

Projekt Hurtownie Danych

Repozytorium GitHub Dokumentacja

Liczba nowych rekordów dodanych do bazy danych: 212

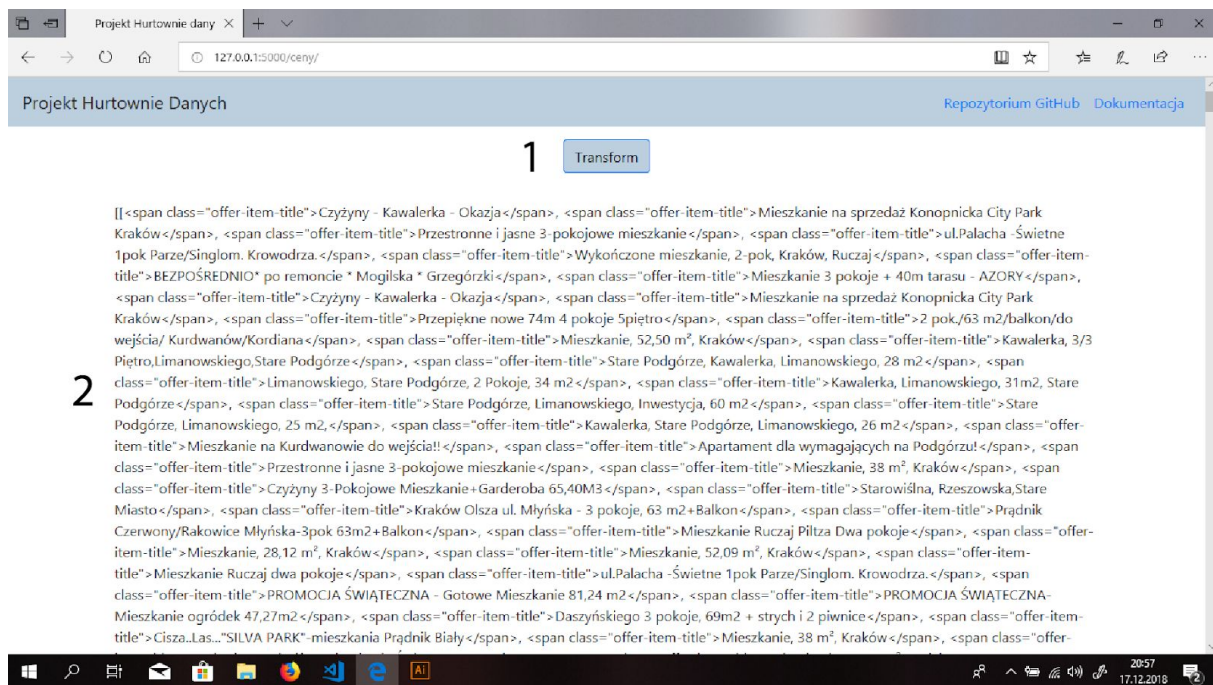
1 Pobierz plik CSV Wyczyść bazę danych 2

3

	Tytuł	Dzielnica	Liczba pokoi	Metraz	Cena za metr	Cena
0	Przestronne i jasne 3-pokojowe mieszkanie	Kraków, Kazimierz	3 pokoje	63,30 m²	10 500 zł/m²	664650zł
1	Stare Miasto - Wysoka stopa zwrotu - Wawel	Kraków, Stare Miasto	1 pokój	55 m²	11 800 zł/m²	649000zł
2	Daszyńskiego 3 pokoje, 69m2 + strych i 2 piwnice	Kraków, Grzegórzki	3 pokoje	69 m²	10 000 zł/m²	690000zł

## Strona procesu ETL

1. Odnośnik do pobrania pliku CSV
2. Przycisk czyszczący bazę danych
3. Tabela zawierająca pobrane dane



## Strona procesu Extract

1. Przycisk do przeprowadzenia procesu Transform
2. Pobrane dane w formie kodu Html

## Model danych użytych w projekcie:

Model danych strukturalny:

Nazwa pliku	Nazwa zmiennej	Opis zmiennej
	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

scrape.py	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	zbiór wszystkich danych pojedynczego ogłoszenia
	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

Nazwa funkcji	Wartość zwracana	Przyjmowane parametry	Opis działania
scrape.py			
getValue()	strona wynik.html wartości c1, c2, data, rec, datat	brak	Funkcja przeprowadza proces ETL
app.route()	plik cms_scrape.csv	/csv_file	Funkcja przekierowuje do pobrania pliku cms_scrape.csv
app.route()	strona clean.html	/clean	Funkcja wykonuje skrypt usuwający wszystkie rekordy z bazy danych
app.route()	plik 25_Projekt.pdf	/dokumentacja_file	Funkcja przekierowuje do pobrania pliku 25_Projekt.pdf
app.route()	strona index.html	/	Przekierowanie podczas uruchamiania skryptu do strony głównej

