

Instrukcja obsługi

1. Aby zainstalować aplikację na serwerze potrzebujemy maszynę z zainstalowanym Apache HTTP Server, PHP oraz bazą danych MySQL. W praktyce posłużyć nam do tego może dowolny hosting internetowy, nawet darmowy.

Wszystkie potrzebne pliki znajdują się w repozytorium GIT'a pod adresem

<https://github.com/klapaucius4/ETL-Process>

Możemy pliki te pobrać w formacie .zip, rozpakować a następnie poprzez klienta FTP (FileZilla, Total Commander) wgrywamy. Jeśli aplikację instalujemy na fizycznej maszynie (komputer, serwer) lub posiadamy dostęp SSH do naszego hostingu, możemy po prostu w wybranych katalogu wywołać komendę:

```
git clone https://github.com/klapaucius4/ETL-Process
```

Oczywiście zakładamy, że Apache, PHP i MySQL mamy na tej maszynie zainstalowane i skonfigurowane.

Następnie tworzymy bazę danych. Na postawionej już stronie którą umieściliśmy na serwerze, znajduje się link do pliku .sql zawierającego strukturę danych tabeli (files/struktura_bazy_danych.sql). Możemy go zaimportować do naszej bazy danych, np. za pomocą phpMyAdmin lub jakiegokolwiek klienta bazodanowego. Możemy też za pomocą tych narzędzi ręcznie stworzyć strukturę naszej bazy danych, gdyż tak naprawdę jest to tylko jedna tabela. Wystarczy użyć jednego zapytania MySQL:

```
-- Zrzut struktury tabela palkora_etl.job_adverts
CREATE TABLE IF NOT EXISTS `job_adverts` (
  `id` int(11) NOT NULL AUTO_INCREMENT,
  `gumtree_id` varchar(255) NOT NULL,
  `title` varchar(255) NOT NULL,
  `date` varchar(255) NOT NULL,
  `location` varchar(255) NOT NULL,
  `by` varchar(255) NOT NULL,
  `type_of_work` varchar(255) NOT NULL,
  `type_of_contact` varchar(255) NOT NULL,
  PRIMARY KEY (`id`)
) ENGINE=InnoDB AUTO_INCREMENT=186 DEFAULT CHARSET=latin2;
```

Posiadając bazę, możemy w pliku config.php ustalić parametry połączeniowe z bazą danych, takie jak: host, nazwa bazy, nazwa użytkownika i hasło.

Gdy to już zrobimy, mamy aplikację poprawnie skonfigurowaną na serwerze.

Jako klient, posłużyć nam już w tym wypadku dowolna przeglądarka internetowa.

2. Interfejs aplikacji jest intuicyjny i nie powinno być raczej większych problemów ze zrozumieniem działania aplikacji.

W pierwszej, tytułowej sekcji mamy możliwość przejścia do repozytorium na Github'ie, pobrania strony tytułowej projektu lub pobrania struktury bazy danych w pliku .sql.

W kolejnych sekcjach mamy możliwość pobrania dokumentacji technicznej i niniejszej instrukcji obsługi.

Główny interfejs aplikacji znajdziemy po kliknięciu w menu na odnośnik „Aplikacja”. Możemy wywoływać odpowiednie operacje za pomocą widocznych przycisków:



APLIKACJA

WYKONANIE PROCESU ETL W CAŁOŚCI

Wykonaj proces ETL

WYKONAJ PROCES ETL ODDZIELNIE

E (extract)

T (transform)

L (load)

ZARZĄDZANIE DANYMI W BAZIE

Zobacz dane w bazie

Pobierz dane w formacie CSV

Wyczyść bazę danych

Oto ich funkcjonalności:

- „Wykonaj proces ETL” – zostaje wykonany pełen proces ETL – trzy etapy – Extract, Transform, Load. Wykonane zostanie to samo, gdybyśmy wywoływali poniżej znajdując się 3 przyciski ręcznie. Przejdźmy więc do nich w kolejnych punktach.
- „E – extract” – na tym etapie wywoływana jest funkcja „extractStep()”. Na samym początku sprawdzamy w pliku etlStatus.json, na jakim ostatnio etapie skończyliśmy proces ETL. Jeśli literka w tym pliku to „L”, oznacza to, że był to proces Load, czyli zakończyliśmy cały proces ETL. W tej sytuacji możemy przejść dalej. W przeciwnym razie pojawi się komunikat o niepowodzeniu. Zadaniem funkcji extractStep() jest połączenie się za pomocą funkcji file_get_html z wybraną stroną (definiowaną wcześniej w pliku config.php) i pobranie całej treści html ze strony. Następnie za pomocą biblioteki „PHP Simple HTML DOM Parser” przekształcamy ten „tekst” do postaci DOM (Obiektowy model dokumentu). Biblioteka dostarcza nam wiele przydatnych funkcji, dzięki którym możemy swobodnie przemieszczać się po strukturze drzewa DOM i swobodnie pobierać wybrane przez nas elementy ze strony. Z jej pomocą tworzymy tablicę z ogłoszeniami z portalu Gumtree, w której dane są walidowane i sprawdzane, czy posiadają odpowiedni format itd. Oprócz tego usuwane są z nich zbędne znaki, takie jak tzw. „white spaces” czy też znaczniki HTML. Ostatecznie dane te zapisywane są do pliku w formacie .csv o nazwie data_from_extract.csv. Na samym końcu w pliku etlStatus.json wstawiamy literkę „E” oraz obecną godzinę. Ma to za zadanie informować aplikację, na jakim etapie procesu ETL skończyliśmy.
- „T – transform” – Na tym etapie uruchamiamy funkcję „transformStep()”, która sprawdzi status procesu, a następnie, jeśli się okaże, że poprzednio wywołano proces Extract – przejdzie dalej. Na tym etapie otwarty zostaje plik data_from_extract.csv, z którego dane są pobrane do tablicy, przekształcone, a następnie zapisane do pliku data_from_transform.csv. Jednocześnie usunięty zostaje plik z poprzedniego procesu o nazwie data_from_extract.csv. Na końcu oczywiście ustalamy status aplikacji w pliku etlStatus.json.

- „L - Load” – ostatni etap procesu ETL, czyli zapisanie danych do bazy. Sprawdzany jest status aplikacji, następnie dane są pobierane z pliku data_from_transform.csv, przekształcane na tablicę i następnie zapisywane do bazy danych MySQL za pomocą zapytania „INSERT”. Zapis jest możliwy wyłącznie wtedy, gdy w bazie danych nie znajduje się jeszcze ogłoszenie o danym ID ogłoszenia. Jeśli jest inaczej, ogłoszenie zostaje pominięte. Oczywiście o każdym takim przypadku użytkownik zostaje poinformowany.
Na końcu następuje usunięcie pliku data_from_transform.csv i ustawienie statusu aplikacji na status „L”. Dzięki temu jesteśmy w stanie wykonać proces ETL od początku.
- „Zobacz dane w bazie” – opcja ta wyświetla nam listę wszystkich rekordów z bazy w tabeli. Dzięki bibliotece JavaScript – DataTable – mamy możliwość swobodnego sortowania danych, paginacji itd.
- „Pobierz dane w formacie CSV” – pobieramy całą zawartość rekordów z bazy danych do pliku w formacie .CSV.

O PROJEKCIE
DOKUMENTACJA
INSTRUKCJA
APLIKACJA

DANE Z BAZY DANYCH

Pokaż 10 pozycji

Szukaj:

Id w bazie	Id z Gumtree	Tytuł	Data dodania ogłoszenia	Lokalizacja	Ogłoszone przez	Rodzaj pracy	Rodzaj umowy
152	1003985839520911544407609	Helpdesk (servicedesk) informatyk z językiem Ukraińskim! WARSZAWA, umowa o pracę + ubezpieczenie!	12/01/2019	Wawer,Warszawa	Firma/Agencja	Pełny etat	Umowa o pracę
153	1002726581400910476863609	IT Researcher/ Bazy Danych	12/01/2019	Żoliborz,Warszawa	Firma/Agencja	Staż	Inne
154	1003980765230911544003509	Informatyk, wsparcie Office Outlook, Strony internetowe, IT ogólnie, Geometria 3D.	11/01/2019	Poznań,Wielkopolskie	Osoba prywatna	Dorywca/Tymczasowa	Inne
155	1002234875610911062594109	System Administrator	11/01/2019	Praga Południe,Warszawa	Firma/Agencja	Pełny etat	Inne
156	1003979691020911543778009	Obsługa sklepu internetowego	11/01/2019	Targówek,Warszawa	Firma/Agencja	Pełny etat	Umowa o pracę
157	1003582463760910471129109	Front-End Developer	11/01/2019	Kraków,Małopolskie	Firma/Agencja	Pełny etat	Umowa o pracę
158	1003979507700910471129109	Front-End Developer	11/01/2019	Kraków,Małopolskie	Firma/Agencja	Pełny etat	Inne
159	1003590299300910471129109	Junior Front-End Developer	11/01/2019	Kraków,Małopolskie	Firma/Agencja	Pełny etat	Umowa o pracę
160	1003978778590910943542309	Osoba do rozliczeń z NFZ	11/01/2019	Kraków,Małopolskie	Firma/Agencja	Pełny etat	Umowa o pracę
161	100397776890911050080109	Wsparcie IT - Obsługa LAN/WLAN z j. angielskim	11/01/2019	Mokotów,Warszawa	Firma/Agencja	Pełny etat	Umowa zlecenie

Pozycje od 1 do 10 z 54 łącznie

Poprzednia 1 2 3 4 5 6 Następna

- „Wyczyść bazę danych” – wyczyści wszystkie istniejące rekordy w bazie danych.

Dodatkowo należy zaznaczyć, że wykonując proces ETL w całości, jak i dzieląc go ręcznie na etapy, nie jesteśmy w stanie tych właśnie etapów wykonać w różnej kolejności. Zawsze jest to schemat Extract -> Transform -> Load. Przy próbie zmienienia kolejności otrzymamy komunikat o niepowodzeniu. Dzieje się tak za sprawą pliku etlStatus.json, który przechowuje informacje o ostatnio wykonywanym etapie procesu ETL.

Wszelkiego rodzaju „produkty uboczne” procesu ETL (takie jak pliki .csv generowane w etapach Extract i Transform) zostają automatycznie usuwane.

3. Przykładowym scenariuszem wykorzystania aplikacji może być chęć analizy ogłoszeń o pracę związanych z branżą IT – aby np. wiedzieć jakie stanowiska, technologie, miejsca zatrudnienia itd. są popularne.

Można by też traktować aplikację jako pewne API, z którego moglibyśmy pobierać ogłoszenia do innych portali, aplikacji itd.