

Implementasi Web Scraping Pada Media Sosial Instagram

by John Doe

Submission date: 27-Nov-2021 10:25PM (UTC+0700)

Submission ID: 1713611194

File name: Makalah_TA.docx (775.01K)

Word count: 1044

Character count: 6614

Implementasi *Web Scraping* Pada Media Sosial Instagram

1 Rio Baskara
Fakultas Teknologi Industri
Universitas Islam Indonesia
Yogyakarta, Indonesia
18523254@students.uii.ac.id

1 Syruz Rahma, S.T., M.Eng.
Fakultas Teknologi Industri
Universitas Islam Indonesia
Yogyakarta, Indonesia
175230101@uui.ac.id

Abstract—Dinas Komunikasi dan Informatika Daerah Istimewa Yogyakarta (Diskominfo DIY) merupakan instansi pemerintah yang memiliki tugas membantu Gubernur melaksanakan urusan pemerintahan bidang komunikasi dan informatika dan urusan pemerintahan bidang persandian. Diskominfo DIY memiliki sebuah sistem analitik berbasis big data, yang difokuskan pada pengembangan data analitik dan pendukung pengambilan keputusan, serta merujuk pada dimensi-dimensi *Jogja Smart Province (JSP)* yang diberi nama *Jogja Center*. Data yang digunakan dalam sistem tersebut dikumpulkan menggunakan teknik *web scraping*. Tujuan dari penelitian ini adalah memenuhi kebutuhan data yang menjadi dasar kebutuhan analisis.

Keywords—*Web Scraping, Media Sosial, Extract, Transform, Load.*

I. PENDAHULUAN

Teknologi adalah penerapan pengetahuan ilmiah dengan tujuan memudahkan kehidupan manusia, menimbulkan perubahan dan manipulasi kehidupan. Seiring dengan berjalannya zaman, perlunya mengembangkan teknologi guna memajukan kehidupan bangsa. Kebutuhan informasi sebagai salah satu alasan, teknologi perlu dikembangkan agar pengolahan data dan informasi lebih mudah dilakukan, mempermudah untuk mendapatkan suatu informasi yang dibutuhkan. *Web Scraping* adalah salah satu teknik pengambilan data, umumnya berupa halaman-halaman web yang terdapat di internet.

Web scraping bukanlah bagian dari data mining, karena *web scraping* terfokus pada cara memperoleh data melalui pengambilan data. Sedangkan data mining berupaya memahami tren atau pola dari data yang telah diperoleh.

Dengan *Web Scraping*, data yang dikumpulkan akan terpusat kepada tujuan yang telah ditentukan, sehingga akan memudahkan dalam proses pencarian. Menciptakan sebuah program yang dapat mempelajari dokumen dengan Bahasa pemrograman HTML pada situs yang dituju, merupakan salah satu cara untuk mengembangkan teknik *web scraping*.

II. KAJIAN PUSTAKA

A. *Web Scraping*

Web scraping merupakan sebuah teknik yang digunakan untuk mengumpulkan data secara manual yaitu dengan melakukan copy paste data yang diinginkan maupun secara otomatis yaitu dengan membuat sebuah program atau kode yang dapat melakukan proses pengambilan data dari sebuah halaman web [1]. Dalam melakukan *web scraping* terdapat beberapa metode yang dapat dilakukan yaitu:

1. Menyalin data secara manual
2. *Regular Expression*

3. *Parsing HTML*
4. Analisa *Document Object Model (DOM)*

Tidak dipungkiri teknik *web scraping* memiliki kekurangan yaitu sampai saat ini belum ada teknik *web scraping* yang 100% efektif. Selain itu hasil yang didapatkan tidak selalu rapih, maka perlu juga untuk memahami struktur halaman website yang dituju. Karena tidak semua data dapat diekstrak dengan mudah, sering kali program harus dijalankan berulang kali yang mengakibatkan akses terhadap halaman web tersebut terblokir. Namun ada pula manfaat dari melakukan *web scraping* yaitu data-data yang didapatkan akan lebih terfokus yang dapat memudahkan dalam pencarian sesuatu.

B. *Extract, Transform, Load (ETL)*

Extract, transform, load atau dapat disingkat menjadi ETL merupakan prosedur umum dalam menyalin data dari suatu sumber ke sistem yang dituju [2]. Proses ini dibagi menjadi tiga tahap yaitu:

1. *Extract*: Ekstraksi merupakan proses paling penting karena dapat menjadi sebuah acuan keberhasilan tahap selanjutnya. Sebagian besar dari proyek database menggabungkan data dari sumber yang berbeda. Setiap sumber juga dapat menggunakan format data yang berbeda. Format data pada umumnya berbentuk JSON, namun ada juga yang memiliki format XML.
2. *Transform*: Pada tahap transformasi data disiapkan untuk dimuat pada target akhir dengan melakukan serangkaian fungsi pada data yang telah diekstrak. Proses *transform* ini berfungsi untuk membersihkan data atau memisahkan data-data yang tepat untuk digunakan. Namun terdapat tantangan ketika sistem yang berbeda melakukan interaksi yaitu antarmuka dan komunikasi sistem yang relevan. Kumpulan karakter yang mungkin tersedia di satu sistem mungkin tidak tersedia di sistem lainnya.
3. *Load*: Tahap *load* merupakan tahap dimana data dimasukan ke target akhir. Proses *load* dapat dibagi menjadi menjadi dua, yaitu *full load* dan *incremental load*. *Full load* merupakan metode untuk memuat seluruh data secara bersamaan untuk menjadi catatan baru pada database. Metode ini berguna untuk menghasilkan data yang tumbuh secara eksponensial namun sulit untuk di atur. Sedangkan *incremental load* merupakan metode untuk memuat data secara *interval* terjadwal. Karena *incremental load* membandingkan data yang masuk dengan data yang sudah ada, metode ini menghasilkan data tambahan jika ditemukan data yang baru dan memudahkan untuk di atur.

V. KESIMPULAN

Berdasarkan implementasi *web scraping* yang telah dilakukan, dapat disimpulkan bahwa:

1. *Web scraping* merupakan suatu teknik yang sangat bermanfaat untuk mengumpulkan data secara cepat.
2. *Web scraping* merupakan suatu hal yang sah untuk dilakukan selama tidak disalah gunakan seperti pencurian data, dan lain-lain.
3. Selanjutnya penelitian ini dapat dikembangkan dengan mengaplikasikannya pada media sosial lain dengan menggunakan bahasa pemrograman lainnya.

VI. DAFTAR PUSTAKA

- [1] Anand V., Kedar G., Schweta A. An Overview On Web Scraping Techniques And Tools. IJFRCSCCE. 2018; 4(4): 363-367.
- [2] Srividya K., Sebastian K. Integrating Big Data: A Semantic Extract-Transform-Load Framework. IEEE. 2015; 48(3): 42-50.

Implementasi Web Scraping Pada Media Sosial Instagram

ORIGINALITY REPORT

9%

SIMILARITY INDEX

9%

INTERNET SOURCES

3%

PUBLICATIONS

0%

STUDENT PAPERS

PRIMARY SOURCES

1	journal.uii.ac.id Internet Source	3%
2	diskominfo.jogjaprov.go.id Internet Source	1%
3	repository.tudelft.nl Internet Source	1%
4	www.techscience.com Internet Source	1%
5	docplayer.info Internet Source	1%
6	pt.scribd.com Internet Source	1%
7	www.scribd.com Internet Source	1%
8	www.smarthipnotis.com Internet Source	1%

Exclude quotes On

Exclude matches Off

Exclude bibliography On