

Implementasi *Web Scraping* Pada Media Sosial Instagram

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

Fayruz 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) 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*. Penelitian ini bertujuan untuk memenuhi kebutuhan data *Jogja Center* yang menjadi dasar kebutuhan analisis. Pengumpulan data yang akan digunakan menggunakan teknik *web scraping* dengan bahasa pemrograman *Python*. Metode yang akan digunakan terdiri dari *analysis*, *coding*, dan *testing*. Hasil dari penelitian ini berupa data dalam bentuk format dokumen *.CSV* yang telah dibersihkan. Dapat disimpulkan bahwa dengan menggunakan teknik *web scraping* dapat memudahkan pengguna dalam mengumpulkan data secara cepat.

Keywords—*Web Scraping, Python*.

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. LANDASAN TEORI

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 rapi, 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 diatur. 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

datetime	timestamp	image_link	caption	link	media_id	comment
05/08/21	1628129892	https://instagr	Hai Lur!\nKab	https://www.	2633266855686853562_3535868226	
04/08/21	1628086393	https://instagr	Hasil laporan	https://www.	2632901954711	Total jumlah keseluruhan
04/08/21	1628063733	https://instagr	ÀuSaya sudah	https://www.	2632711871105217551_3535868226	
04/08/21	1628051653	https://instagr	Berkaitan den	https://www.	2632610541166	Yang sumbangan 2T, te
03/08/21	1627998671	https://instagr	Rekapitulasi P	https://www.	2632166088460	Mekanisme nya lewat
03/08/21	1627996397	https://instagr	Hasil laporan	https://www.	2632147015559	Total jumlah keseluruh
02/08/21	1627916220	https://instagr	Hai Lur! Menji	https://www.	2631474440667690661_3535868226	
02/08/21	1627911833	https://instagr	Hasil laporan	https://www.	2631437642981	Terimakasih atas Infom
02/08/21	1627887588	https://instagr	Hai Lur! \nMa	https://www.	2631234017189875923_3535868226	
01/08/21	1627824884	https://instagr	Hasil laporan	https://www.	2630708262777	Yok Jogja masih tinggi
01/08/21	1627809711	https://instagr	Pemda DIY lah	https://www.	2630580978334312085_3535868226	
01/08/21	1627784612	https://instagr	Hai Lur!\n\nKi	https://www.	2630369108461722592_3535868226	
31/07/21	1627740685	https://instagr	Hasil laporan	https://www.	2630001952091809450_3535868226	
31/07/21	1627723789	https://instagr	Berkaitan den	https://www.	2629860217718142644_3535868226	
30/07/21	1627645103	https://instagr	Hasil laporan	https://www.	2629200148932455629_3535868226	
30/07/21	1627635341	https://instagr	Gubernur DIY	https://www.	2629118259894690767_3535868226	
30/07/21	1627610659	https://instagr	Hai Lur!\n\nM	https://www.	2628910891072300586_3535868226	
29/07/21	1627559657	https://instagr	Hasil laporan	https://www.	2628483375625587372_3535868226	
28/07/21	1627484781	https://instagr	Halo Sahabat	https://www.	2627855267964750978_3535868226	

Gambar 6. Hasil Setelah Parsing

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 disalahgunakan 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. *IJFRCSE*. 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.
- [3] Wijaya Y., Astuti M. Sistem Informasi Penjualan Tiket Wisata Berbasis *Web* Menggunakan Metode *Waterfall*. e-ISSN. 2019; 1(2): 273-276.