

Rajiv Gandhi University of Knowledge Technologies

R.K Valley, Y.S.R Kadapa (Dist) - 516330

A
project report
on

Data Delivery System through Web Scraping

Submitted by

T.Sujay Kumar

R170044



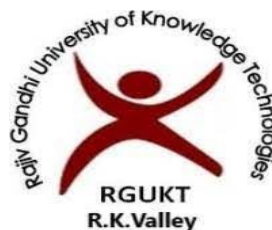
Under the guidance of

MRS.CHALLA RATNAKUMARI
(Assistant Professor, CSE)

Department of Computer Science Engineering

This project report has been submitted in fulfillment of the requirements for the Degree of Bachelor of Technology in software Engineering.

**Rajiv Gandhi University of Knowledge Technologies
IIIT, R.K. Valley, YSR Kadapa (Dist) - 516330**



CERTIFICATE

This is to certify that report entitled “**Data Delivery System through Web scraping**” Submitted by T.Sujay Kumar (R170044) in partial fulfillment of the requirements of the award of bachelor of technology in Computer Science Engineering is a bonafide work carried by the them under the supervision and guidance. The report has been not submitted previously in part or full to this or any other university or institute for the award of any degree or diploma.

GUIDE
MRS.CHALLA RATNAKUAMRI

HEAD OF THE DEPARTMENT
N.SATYANANDARAM

ACKNOWLEDGEMENT

The satisfaction that accompanies the successful completion of any task would be incomplete without the mention of the people who made it possible and who's constant guidance and encouragement crown all the efforts success. I would like to express my sincere gratitude to **Mrs.Challa Ratnakumari**, my project guide for valuable suggestions and keen interest throughout the progress of my project. I am grateful to **Mr.Satyanandaram sir HOD CSE**, for providing excellent computing facilities and congenial atmosphere for progressing my project. At the outset, I would like to thank **Rajiv Gandhi Of University of Knowledge Technologies (RGUKT)**, for providing all the necessary resources and support for the successful completion of my course work.

DECLARATION

We hereby declare that this report entitled “**Data Delivery System through Web scraping**” Submitted by me under the guidance and Supervision of **Mrs.Challa Ratnakumari**, is a bonafide work. We also declare that it has not been of Submitted previously in part or in full to this University or other institution for the award of any degree or diploma.

Date:- 03-05-2022

Place:-RK VALLEY

T.Sujay Kumar (R170044)

INDEX

S.NO	INDEX	PAGE NO
1	Abstract	6
2	Introduction	7
3	Design and Analysis	8 - 9
4	Implementation and Analysis	10 - 18
5	Testing	19
6	Results	20 - 21
7	Conclusion	22
8	References	22

ABSTRACT

Web scraping is the process of extracting data and content from web pages. When we use the traditional approach of web scraping for data delivery it takes more time and some data or details may be excluded. So in order to extract each and every minute details or data from the web page, we implemented the same web scraping process in 6 phases in our project, namely **recrawl**, **feedcrawl**, **extraction**, **dedup**, **normalisation**, **upload**. By following this approach we can extract each and every minute details or data from the web page in a short time.

The primary objective of this project is to provide a complete, accurate, reliable data to the customer in a short time and in an efficient manner.

INTRODUCTION

As we can see in our daily lives, we come across huge volumes of data from web pages through the internet. Most of the websites do not provide a mechanism such as web services, APIs to collect their data, or the provided APIs are either poorly documented or difficult to use. So for extracting these huge volumes of data we are using the process of web scraping. Web scraping is a process of extracting the content and data in large quantities from the web pages. Most of this data is in unstructured format i.e. in HTML format which is needed to be converted into structured data so that it can be used in various applications.

By using this process of web scraping we can extract huge volumes of data from the web page and provide that data in a structured format like xml, json, csv to the customer. Another concern while extracting data is some data may not be included in the process of extraction. So we implemented this web scraping process in a 6 phase approach namely **recrawl**, **feedcrawl**, **extraction**, **dedup**, **normalisation**, **upload**. When we scrape huge volumes of data through this approach, we can extract each and every minute details and data from the web page that too in a short time and in an efficient manner.

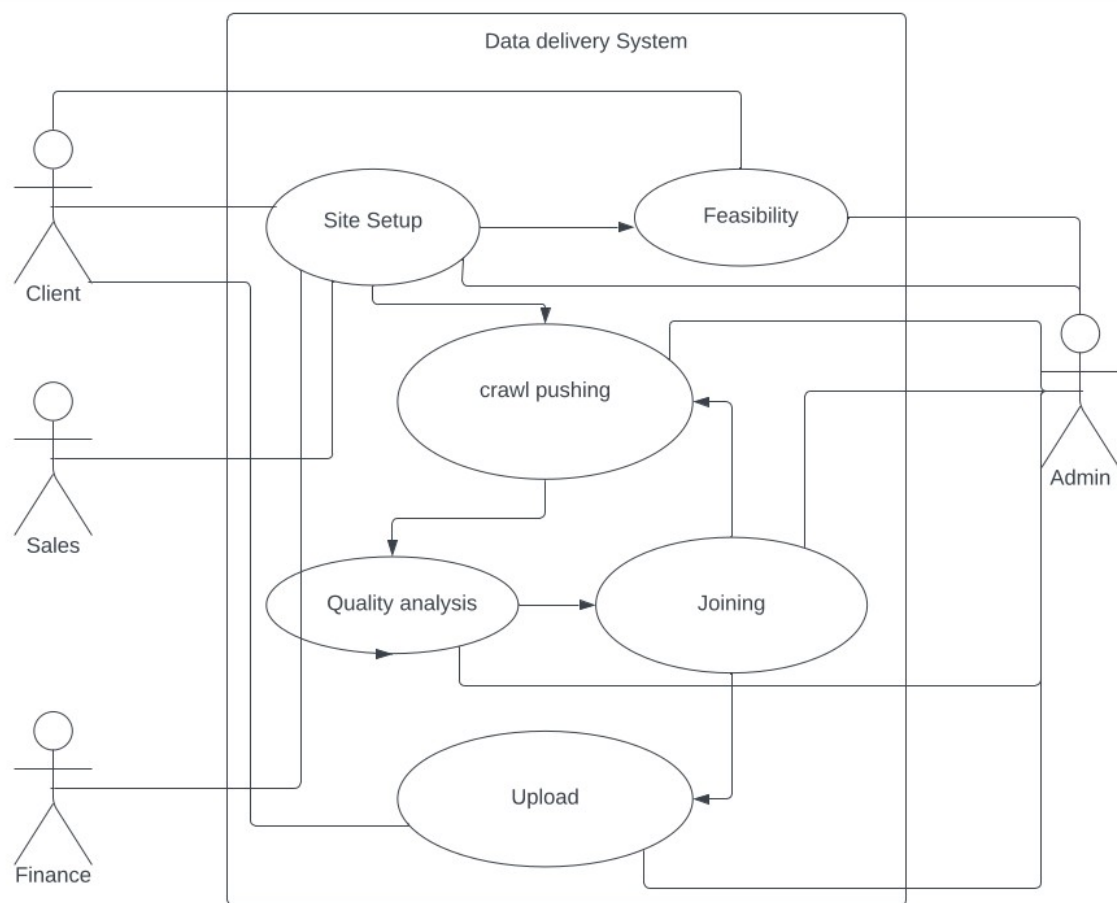
FUNCTIONAL REQUIREMENTS

The functionalities we are implementing in this project are:

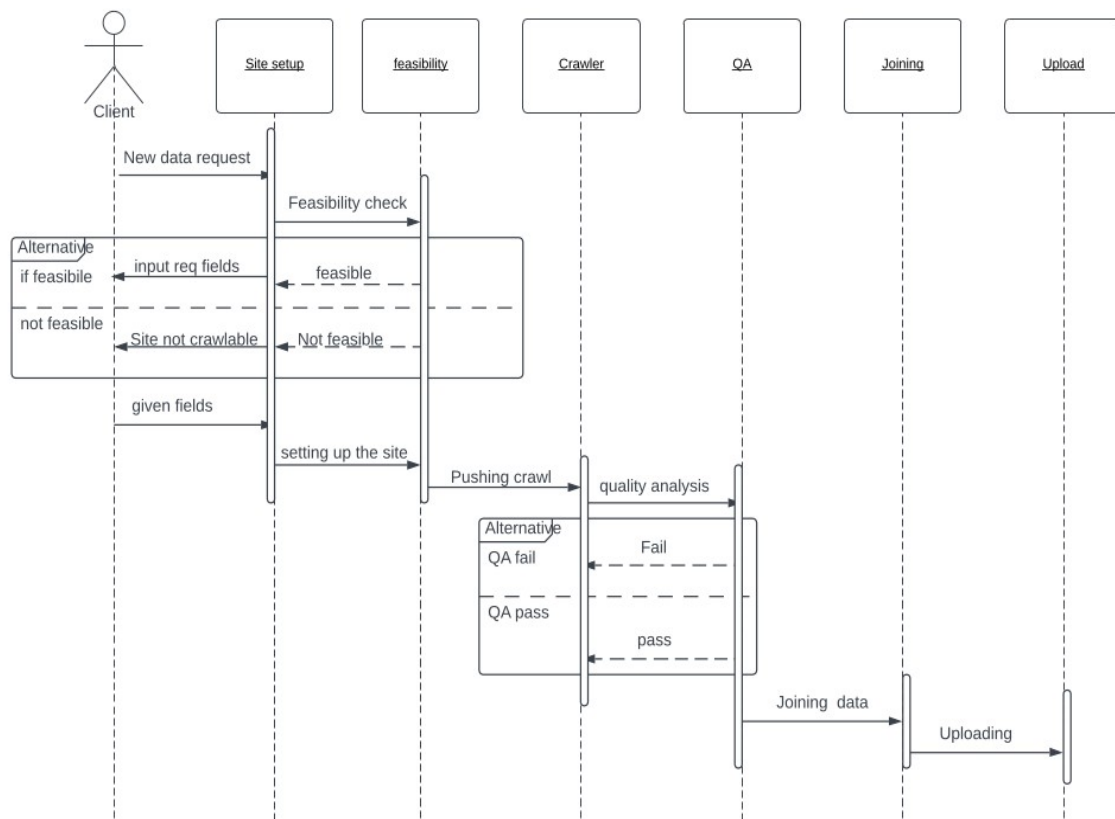
1. **Recrawl**: Here we collect all final product urls from the site
2. **Feedcrawl**: Here we fetch the page source of collected urls and store it.
3. **Extracton**: Here we extract the required data fields from the fetched page source.
4. **Dedup**: Here we remove duplicate records from the extracted data.
5. **Normalisation**: Here we convert the format of data file from xml to required format.
6. **Upload**: Here we upload the data to the clients location.

DESIGN AND ANALYSIS

Use Case Diagram :



Sequence Diagram :



IMPLEMENTATION AND ANALYSIS

TOOLS AND TECHNOLOGIES

RUBY

Ruby is a dynamic, reflective, object-oriented, general-purpose programming language. Ruby is a pure Object-Oriented language developed by Yukihiro Matsumoto. Everything in Ruby is an object except the blocks but there are replacements too for it i.e procs and lambda. The objective of Ruby's development was to make it act as a sensible buffer between human programmers and the underlying computing machinery. In this project we are using this ruby programming language for building rss, dsk, ext logics.

MYSQL

MySQL is an open source relational database management system. It is the acronym for structured query language. A relational database organises data into one more data tables in which data may be related to each other. In our project we are using this to store app.promptcloud data and redmine data.

OPEN URI

Open uri is a ruby library which is used for fetching an url of a web page. With the help of this open uri we can make HTTP requests for fetching the page source of the web page. Open uri is an easy-to-use wrapper for Net::HTTP, Net::HTTPS and Net::FTP.

NOKOGIRI

Nokogiri is another ruby library which is used to for parsing the HTML, XML document. Nokogiri is a ruby library for parsing the HTML or Xml documents. It is not possible to directly extract the data from the HTML page, first it need to be converted into Dom. To do so we are using this Nokogiri.

ELASTIC SEARCH

Elastic Search is a distributed, free and open search and analytics engine for all types of data, including textual, numerical, geospatial, structured and unstructured. ES allows us to store and search and analyze huge volumes of data quickly. It is used to store the data and query the data. In our project it is used in Rss, Dedup phases.

RABBITMQ

RabbitMq is one of the most popular open source message brokers. RabbitMq is the lightweight and easy to deploy on premises in the cloud. It supports multiple messaging protocols. RabbitMQ can be deployed in distributed and federated configurations to meet high-scale, high-availability requirements. Here we store some unnecessary data or the data which is required only for some time. Rss, Drss output will be stored in RabbitMQ.

CRAWLBOARD

Crawlboard is a dedicated management system for managing all the crawls that have been pushed. We are calling this as app.promptcloud and it contains major feature called as internal dashboard which is used to get the crawl data of any sites that we are crawling. And we can collect failed urls, respush failed urls and also lot of functionalities with this crawlboard.

KIBANA

Kibana is a free and open frontend application that sits on top of the Elastic Stack, providing search and data visualization capabilities for the data indexed in Elastic Search. Commonly known as the charting tool for the Elastic Stack. Kibana also acts as the user interface for monitoring, managing and securing an Elastic Stack cluster.

IMPLEMENTATION

We implemented this web scraping in our pipeline in 6 stages, they are:

RECRAWL(RSS):

Here we extract the actual page URL's using a plugin called rss plugin. This plugin creates a dom out of a file and extracts page urls using Xpath and then pushes it to a queue from where feed crawl happens.

In this initial stage we follow depth rules. In the depth0 we take the home page url (seed url) of a website and extract all the category urls from the website. In the depth1 we take the total page count of a category. In this depth2 we extract job urls from each page.

Based on the configuration given in yml file of rss, product urls will be pushed to rss_queue. The RSS plugin will create a rss_queue having final urls. The rss_queue is located in the following path:

```
<crawl_home><pipeline_v2><rss_crawl><site_name>_<yyyy_mm_dd>_<crawl_time stamp>
```

The linux command to test this rss code is :

```
rtest -t r -s <site_name>
```

Our pipeline first begins with this stage. In this stage we will collect all the product urls from the given seed_url by the client. And then we will store these urls in a temporary storage called as rss_queue.

Example:

```
1 1.-
2 depth0:
3   seed_urls:
4   xpath: "//div[contains(@class,'center')]/p"
5   method_name: get_pagination_urls
6   verification_xpath: "//meta[contains(@content,'Hourglass')]"
7   request_type: curl
8   rss_crawl_using_proxies: true
9   proxy_source: webshare_proxies, new_rotating_proxies, rack_proxies
10  fetch_retry_attempt_per_url: 5
11  ignore_cache_for_retries: true
12 depth1:
13  xpath: ".//a[contains(@class,'grid')]/@href"
14  method_name: get_product_urls
15  verification_xpath: ".//a[contains(@class,'grid')]/@href"
16  request_type: curl
17  rss_crawl_using_proxies: true
18  proxy_source: webshare_proxies, new_rotating_proxies, rack_proxies
19  fetch_retry_attempt_per_url: 5
20  ignore_cache_for_retries: true
```

rss.yml file

```
18 class SiteSpecificPlugin::RssPlugin::PromptcloudMaster::RssHourglasscosmeticsUkSujayTrainingPromptcloudMaster < PipelineV2::Stages::Recrawl
19   def get_inputs_for_all_depth()
20     return get_all_depth_from_yaml_file()
21   end
22
23   def get_pagination_urls(url, current_depth_level, current_depth_hash)
24     pagination_urls = []
25
26     begin
27       page_no = 1
28       while true
29         page_url = url + "?page=#{page_no}"
30         dom = get_html_dom(page_url, current_depth_hash)
31         raise Exception.new("Dom not found") if not dom
32         msg = dom.xpath(current_depth_hash["xpath"]).first
33         break if msg
34         pagination_urls << page_url
35         page_no = page_no + 1
36       end
37     rescue Exception=>e
38       $log.error "Error in #{__method__}, for url:#{url}, message:#{e.message}, class:#{e.class}", @log_hash
39     end
40
41     return pagination_urls
42   end
43
44   def get_product_urls(url,current_depth_level, current_depth_hash)
45     product_urls = []
46
47     begin
48       dom = get_html_dom(url, current_depth_hash)
49       raise Exception.new("dom not found") if not dom
50
51       product_nodes = dom.xpath(current_depth_hash["xpath"])
52       product_nodes and product_nodes.each do |product_url|
53         next if not product_url
54         product_urls << "https://www.hourglasscosmetics.co.uk#{product_url.content}|#{url.split("/").last.split("?").first}"
55       end
56     rescue Exception=>e
57       $log.error "Error in #{__method__}, for url:#{url}, message:#{e.message}, class:#{e.class}", @log_hash
58     end
59
60     return product_urls
61   end
62
63 end
```

rss.rb file

DISKFETCHER(FEEDCRAWL):

It is the second phase of our pipeline. Where we fetch the pages of urls which are stored in rss_queue. It takes rss_queue as an input and fetches all the pages in that queue and then it stores those pages.

Diskfetcher takes popped urls from rss queues as Input. All the URLs from rss queue pushed as input and we are hitting that particular Product/ review URLs with different types of request types those are **get, post, curl, open_uri, scraper_api**. In this stage it checks the daily_recrawl_lists(drl) for reading the inputs of frequency time, verification xpath, proxy country and other required fields for a project nothing but all the configuration. By taking all these fields we can fetch the page source of the url and send the page source file to next pipeline stage. In this output will be stored in our local. This page source is send to next pipeline stage i.e; Extraction.

Command:

```
rtest -t f -s <site_name>
```

Example:

```
18 class SiteSpecificPlugin::DiskfetcherPlugin::ReviewsPlusPlusPhase4::EnagRoPhase4ProductsReviewsPlusPlusPhase4DiskfetcherPlugin < SiteSpecificPlugin::DiskfetcherPlugin::ReviewsPlusPlusPhase3::EnagRoPhase3Prod
   uctsReviewsPlusPlusPhase3DiskfetcherPlugin
19
20   def get_fetch_timeout_per_url_from_drl(drl_site_hash)
21     return 700
22   end
23   def get_page_content_hash(url,args_hash={})
24     begin
25       request = {}
26       request["Authority"] = "www.enag.ro"
27       request["Accept"] = "text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3;q=0.7"
28       request["Accept-Language"] = "en-GB,en-US;q=0.9,en;q=0.8"
29       @user_agents_pool = SDF::get_user_agent_pool_from_agent_config
30       request["User-Agent"] = @user_agents_pool.shuffle.sample
31       args_hash[:extended_header] = request
32       page_content_hash = super
33     rescue Exception=>e
34       $log.error "Error in method #{__method__} for url #{url}", @log_hash
35     end
36     return page_content_hash
37   end
38 end
39
```

diskfetcher file

EXTRACTION:

It is the main phase where the actual data is being collected. We apply the xpaths on the DOM of the pages which are fetched in feedcrawl stage. We extract the fields which are required by the client and then store them in xml format.

We will extract the required fields from the source code and parse the text to required format. It consider the output of Diskfetcher as input which is stored in local. Here we write all the methods for each and every field which client required separately. We use Xpaths to extract the field from the source page. We will extract the required fields from the source code and parse the text to required format. For any website uniq_id and crawl_timestamp are the required fields these two fields are generated by md5 hash. If any required fields have not generated for the url then it will drop in extraction stage itself. Those urls data neet not send to next pipeline stage.

Extraction plugin will create an ext segment of having the details of a web page url. Along with that it also contains the Quality Analysis report, start time and end time required to crawl the data.

Command:

Testing a single webpage url : `rtest -t e -s <site_name> -u <url>`

Testing an entire segment : `rtest -t e -s <site_name> -q`

After ruuning above commands, if all the code is proper and error free then it will generate xml file consisting of required fields. Then this xml file will be passed to next phase called as dedup.

Example:

```
94 price:
95   desc_of_xpath: "//meta[contains(@property,'amount')]/@content"
96   standard_nodeset_range: first
97   standard_nodeset_join_char: "|"
98   standard_post_processing_functions_on_text:
99     - remove_newlines_and_whitespace
100 views:
101   desc_of_xpath:
102   standard_nodeset_range: first
103   standard_nodeset_join_char: "|"
104   standard_post_processing_functions_on_text:
105     - remove_newlines_and_whitespace
106 images:
107   desc_of_xpath: "//meta[contains(@property,'image')]/@content"
108   standard_nodeset_range: all
109   standard_nodeset_join_char: "|"
110   standard_post_processing_functions_on_text:
111     - remove_newlines_and_whitespace
112 test_urls:
113   url_given_via_dashboard: www.hourglasscosmetics.UK
```

ext yml file

```
156 def get_contact_person_profile(page_doc, inhash)
157   #uncomment one of these lines as necessary
158   #inhash[:return_type] = "standard_function_value"
159   #inhash[:return_type] = "tld_function_value"
160   #inhash[:return_type] = "field_hash"
161   value = get_automated_value_from_page_doc(page_doc, inhash, field_name='contact_person_profile')
162 end
163
164 def get_location(page_doc, inhash)
165   #uncomment one of these lines as necessary
166   #inhash[:return_type] = "standard_function_value"
167   #inhash[:return_type] = "tld_function_value"
168   #inhash[:return_type] = "field_hash"
169   value = get_automated_value_from_page_doc(page_doc, inhash, field_name='reviewer_location') #Autosuggested
170 end
171
172 def get_zip(page_doc, inhash)
173   #uncomment one of these lines as necessary
174   #inhash[:return_type] = "standard_function_value"
175   #inhash[:return_type] = "tld_function_value"
176   #inhash[:return_type] = "field_hash"
177   value = get_automated_value_from_page_doc(page_doc, inhash, field_name='zip')
178 end
179
180 def get_price(page_doc, inhash)
181   #uncomment one of these lines as necessary
182   #inhash[:return_type] = "standard_function_value"
183   #inhash[:return_type] = "tld_function_value"
184   #inhash[:return_type] = "field_hash"
185   value = get_automated_value_from_page_doc(page_doc, inhash, field_name='price')
186 end
187
188 def get_views(page_doc, inhash)
189   #uncomment one of these lines as necessary
190   #inhash[:return_type] = "standard_function_value"
191   #inhash[:return_type] = "tld_function_value"
192   #inhash[:return_type] = "field_hash"
193   value = get_automated_value_from_page_doc(page_doc, inhash, field_name='views')
194 end
195
196 def get_images(page_doc, inhash)
197   #uncomment one of these lines as necessary
198   #inhash[:return_type] = "standard_function_value"
199   #inhash[:return_type] = "tld_function_value"
200   #inhash[:return_type] = "field_hash"
201   value = get_automated_value_from_page_doc(page_doc, inhash, field_name='images')
202 end
203
204 end
205
```

ext rb file

DEDUP

DEDUP is the stage where all the duplicates are removed from the data. If there are any duplicates present in the data then we eliminate this with the help of elastic search.

In this stage we remove the redundant or duplicate records from the extracted records. The duplication process can be identified by using the `uniq_id` which has been generated from the extraction. Dedup plugin will create a dedup segment of having the details of a web page url. Along with that it also contains the start time and end time required to crawl the data.

Command:

```
rtest -t d -s <site_name> -q
```

This phase will be automatically done by the plugin. No need to write any logic here, Just we need to mention the required format of data to be delivered.

NORMALIZATION

Normalization is the process where the format of data is being converted. By default the data will be in xml format, if client want data to be delivered in other formats like json, csv then we will convert that xml into required format.

In this stage we Normalize the data and change some fields according to the client's requirements. Normalize plugin will create a normalize segment of having the details of a web page url. Along with that it also contains the start time and end time required to crawl the data.

Command:

```
rtest -t n -s <site_name> -q
```

UPLOAD

It is the final stage where the data is being delivered to the clients location. In this stage we upload the extracted data to our API server. If client wanted to deliver to their location the we will upload the data to clients location.

Command:

```
rtest -t u -s <site_name> -q
```

This is the phase where the actual data is being delivered to the client's location. Here we upload the extracted data to two locations. They are:

1. Application Program Interface (API link)
2. Clients Bucket/location

We are having another type upload called as depend data upload which means when the client wants the data which is in the form of image, audio, video we use this concept of dependent data upload.

The data which not in the form of text will be treated as dependent data. We store this dependent data in a temporary storage called as RabbitMq which then uploaded to the client's location.

TESTING

S.NO	ACTION	INPUT	EXPECTED RESULTS	ACTUAL RESULTS	STATUS
1	To check whether seed urls are proper.	Open any seed url.	It should display the web page of that url	Web page displayed	Pass
2	To check whether seed urls are redirecting.	Open any seed url an check whether it is redirecting or not.	It sholud not redirect to any other page.	Url is not reedirecting	Pass
3	To verify whether all the category urls are added in crawl configs	Open crawl configs and check whether all the required categories urls are added.	It should display all the required category urls.	Category urls are displayed	Pass
4	To verify pagination	Open any seed url and check whether all pagination urls are added.	It should display all the pagination urls	Pagination urls are displayed	Pass
5	To verify fetched page source.	Fetch the page source of any product url and check whether it contains all required fields or not.	It should fetch the correct page source.	Correct page source is fetched	Pass
6	To verify data accuracy	Perform extraction on any product url.	It should display all the required data fields according the page source.	All the required data fields are extracted correctly.	Pass
7	To verify upload	Upload the extracted data to tthe api server and clients location.	It should upload the data to clients location.	Data uploaded to the clients location	Pass

RESULTS

INPUT

We take seed_url(site url) and the fields to be extracted as an input from the client.

Site Details

Site URL: <https://www.cycletrader.com/motorcycles-for-sale>

Project name: sujay_training_promptcloud_master

Crawl Type: fullcrawl

Schema Details: Attached as a excel file

Additional Information:

Attachments:

Delivery Format: csv

Delivery Frequency: Daily

Delivery Method: PromptCloud API

..update from app

← → ×

cycletrader.com/motorcycles-for-sale

HR Stop Redmine The Predictive... Mattermost QnA | Questions Slab Talent Lms YouTube Gmail Crawlboard

Motorcycles For Sale

View Makes | View New | View Used | View States | Under \$5000 | Under \$2000

WE'RE HIRING!

JOIN OUR TEAM

Cycle Trader

Home

Buy Online

Complete your Motorcycle purchase 100% online.

Show listings I can buy online (2,598)

Save this search Clear All

Begin searching by adding filters below

Location

Keyword

Price

Make

Model

Trim

Featured

2005 Harley-Davidson® SCREA...
World of Powersports -

\$11,499

38,315 miles

Featured

2016 Harley-Davidson FLHXS - S...
DREAM MACHINES OF TEXAS

Managers Special

20,190 miles

Featured

2014 Harley-Davidson® FLHTK ...
DREAM MACHINES OF TEXAS

Managers Special

24,081 miles

Featured

2017 Harley-Davidson® FLHTK ...
DREAM MACHINES OF TEXAS

\$18,950

31,273 miles

Featured

2020 Harley-Davidson FLHXSE - ...
DREAM MACHINES OF TEXAS

\$37,977

4,910 miles

Featured

2018 H...
DREAM MACHINES OF TEXAS

\$27,9...

...

191,441 Motorcycles For Sale

Sorted By Premium

Available For Purchase Online

Premium

2016 Harley-Davidson® Touring
FLHXS - Street Glide® Special

15,691 miles

\$17,990

Used 2016 Harley-Davidson® Touring
FLHXS - Street Glide® Special

2016 Harley-Davidson® FLHXS - Street Glide® Special,
YOU ARE LOOKING AT A 2016 HARLEY DAVIDSON STREET GLIDE SPECIAL (FLHXS) WITH 15,691 MILES

Email

1-844-823-0736

DREAM MACHINES OF TEXAS - Farmers Branch, TX

Chat

Available For Purchase Online

OUTPUT

We take site urls as an input and produce an xml document containing the required fields of data as an output.

```
<?xml version="1.0" encoding="UTF-8"?>
<root>
  <page>
    <pageurl>https://www.cycletrader.com/listing/2023-Indian+Motorcycle-PursuitKC2NAE+Dark+HorseKC2NAE-5024124701</pageurl>
    <entity>
      <record>
        <uniq_id>Scab3920b854954aecec648699b224bc</uniq_id>
        <crawl_timestamp>2023-01-05 09:56:47 UTC</crawl_timestamp>
        <pageurl>https://www.cycletrader.com/listing/2023-Indian+Motorcycle-PursuitKC2NAE+Dark+HorseKC2NAE-5024124701</pageurl>
        <ad_id>5024124701</ad_id>
        <ad_title>Pursuit Dark Horse</ad_title>
        <category>Touring</category>
        <sub_category>Indian</sub_category>
        <category_tree>Home<B>Browse Motorcycles<B>Touring<B>Indian<B>PURSUIT<B>Wisconsin<B>Racine<B>2023</category_tree>
        <contact_person>Indian Motorcycle of Racine</contact_person>
        <contact_person_profile>https://rideracine.com/dxinventory/e72db419-03c5-4f38-814d-af640101edd4</contact_person_profile>
        <location>Racine, WI</location>
        <zip>53403</zip>
        <prices>$31,749</prices>
        <images>https://cdn1.cycletrader.com/v1/media/639b65ff9a6e046d7513022c.jpg?width=96&height=72&quality=60&bestfit=true&upsized=true&blurBackground=true&blurValue=100</images>
      </record>
    </entity>
  </page>
</root>
```

CONCLUSION

we have implemented the process of extracting the data from the web pages in six phases. This implementation of web scraping process through the six phases namely **recrawl**, **feedcrawl**, **extraction**, **dedup**, **normalisation**, **upload** has enhanced the data delivery system.

By implementing this we are able to deliver the accurate, complete, reliable data to the customer in a short time and in an efficient manner. We have brought some improvements in the data delivery system with the help of this approach of six phases. Improvements include fast data delivery, accurate data delivery, complete data delivery.

REFERENCES

1. <https://redmine.promptcloud.com/issues/328349>
2. <https://www.geeksforgeeks.org/what-is-web-scraping-and-how-to-use-it/>
3. <https://promptcloud.slabs.com/posts/complete-project-setup-in-the-pipeline-step-by-step-d-scraper-api-0464lbrv>
4. https://redmine.promptcloud.com/projects/sdftech/wiki/SE_SSE_Training_Module_2
5. <https://www.tutorialspoint.com/ruby/index.htm>