

Programming assignment 2 - Data extraction

Kim Ana Badovinac, Jakob Petek, Jovan Prodanov
Faculty of Computer and Information Science, University of Ljubljana

I. INTRODUCTION

Data extraction from websites is a critical task for businesses and researchers who need to gather valuable information from the internet. One of the most efficient ways to accomplish this task is through the use of automated tools that can extract relevant data from web pages. In this context, regular expressions, XPath, and the RoadRunner algorithm are three popular methods for extracting structured data from web pages which, as per our programming assignment, we had to implement.

II. IMPLEMENTATION

All three extraction methods were implemented in Python. We have extracted data from 6 websites, 4 of which were already specified (rtvslo.si and overstock.com) and 2 optional websites which we decided to be from the Steam store domain. The data collected from Steam store was the name of the video game, its release date, price, rating and amount of reviews as well as any tags associated with the game. The relevant data extracted is shown in figure 1.

only one regular expression. For this we used Python's `re` module. Each function takes in a html file on which we then perform regular expressions to gather data. Most of the gathered data does not need any additional editing or filtering (except `rtvslo` article's content). The function then packs the data to a JSON file and prints out its content on the standard output. The statements used for `rtvslo`, `overstock` and `steam` can be seen on figures 2, 3 and 4 respectively.

```
title = re.search(r'<h1>(.+?)</h1>', html_content).group(1)
published_time = re.search(r'<div class="publish-meta">(.+?)</div>', html_content).group(1)
author = re.search(r'<div class="author-name">(.+?)</div>', html_content).group(1)
subtitle = re.search(r'<div class="subtitle">(.+?)</div>', html_content).group(1)
lead = re.search(r'<p class="lead">(.+?)</p>', html_content).group(1)
content = re.search(r'<article class="article">(.+?)</article>', html_content).group(1)
```

Fig. 2. Regular expressions on rtvslo

```
titles = re.findall('title=([^&]+)', html_content)
list_prices = re.findall('listPrice=([0-9.]+)', html_content)
links = re.findall('link=([^&]+)', html_content)
savings = re.findall('savings=([0-9.]+)', html_content)
savings_percent = re.findall('savingsPercent=([0-9.]+)', html_content)
contents_data = re.findall('contentsData=([0-9.]+)', html_content)
```

Fig. 3. Regular expressions on overstock

```

page_title = re.search('<title>(.*)shopspass pass(.*)</title>', html_content).group(1)
reply = re.search('<reply>(.*)</reply>', html_content).group(1)
amount_of_reviews = re.search('<div class="review">(.*)</div>', html_content).group(1)
release_date = re.search('<div class="date">(.*)</div>', html_content).group(1)
price = re.search('<div class="discount final price">(.*)</div>', html_content).group(1)
tags = re.findall('<div class="label">(.*)</div>', html_content)

```

Fig. 4. Regular expressions on Steam store

B. XPath implementation

Similarly to the implementation using regular expressions, we have extracted same data using XPath. We have again implemented 3 functions that take in the html document on which XPath expressions were used to gather data. We have again only used one expression per required data, except at certain parts where we had to clean up the output. The XPath statements used for rtvslo, overstock and steam can be seen on figures 5, 6 and 7 respectively.

```
title = tree.xpath("//title/text()")[0]
published_time = tree.xpath("//a[@class='publish-meta']/text()")[0]
published_time = published_time.lstrip('w\n')
author = tree.xpath("//a[@class='author-name']/text()")[0]
subtitle = tree.xpath("//a[@class='subtitle']/text()")[0]
lead = tree.xpath("//a[@class='lead']/text()")[0]
contents = tree.xpath("//a[@class='article-body']/article/p/text()//a[@class='article-body']/article/p/strong/text()")
```

Fig. 5. XPath expressions on rtvslo

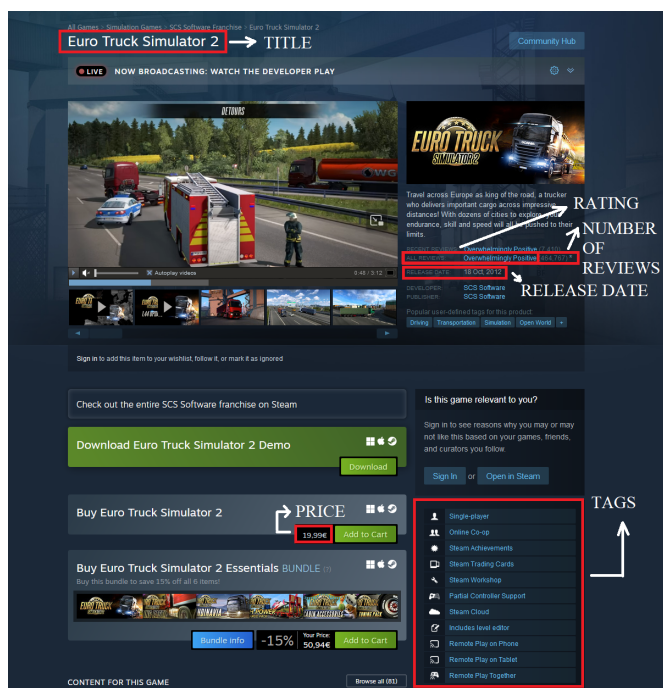


Fig. 1. Data on Steam store

A. Regular expression implementation

As per requirements we have implemented 3 functions (for each type of website) and extracted the required data using

```

titles = tree.xpath(r"//*[@valign='top']/a[@href]/b/text()")
list_prices = tree.xpath(r"//*[@nowrap='nowrap']/s/text()")
prices = tree.xpath(r"//*[@class='bigred']/b/text()")
savings_data = tree.xpath(r"//*[@class='littleorange']/text()")
contents_data = tree.xpath(r"//*[@class='normal']/text()")

```

Fig. 6. XPath expressions on overstock

```

game_title = tree.xpath(r"//*[@id='apphubAppName'][@class='apphub_AppName']/text()")
review = tree.xpath("//span[@class='game_review_summary positive'][@itemprop='description']/text()")
amount_of_reviews = tree.xpath(r"//meta[@itemprop='reviewCount']/@content")
release_date = tree.xpath(r"//div[@class='date']/text()")
price = tree.xpath(r"//div[@class='discount_final_price']/text()")
tags = tree.xpath(r"//div[@class='label']/text()")

```

Fig. 7. XPath expressions on Steam store

return regex

Unwanted tags in our implementation are: ["script", "input", "option", "style", "br", "hr"] and in the `getCommonAttrs` we avoid style and bracket style attributes.

2) *Wrappers*: There are only parts of the wrappers produced in pictures, however they are currently present in the GitHub repository in .html format.

C. RoadRunner implementation

Roadrunner is implemented in a way that it generates a common wrapper in the form of HTML code between two web pages. It generates common tags, common attributes of the tags and common content of tags.

Some details are when the content in the tags differ in the two websites (most commonly - a string), when that happens the `#PCDATA` symbol is inserted. Also when the children of tags differ, the `(<tag.*?>.*?</tag>)+` symbol is inserted.

1) Pseudo code:

```

generateWrapper(webpage1, webpage2):
    removeUnwantedTags(webpage1)
    removeUnwantedTags(webpage2)
    regex += generateRegex(head1, head2)
    regex += generateRegex(body1, body2)
    regex = "<html>" + regex + "</html>"
    return regex

generateRegex(tag1, tag2):
    attrs = getCommonAttrs(tag1, tag2)
    regex = "<" + tag1.name + attrs + ">"
    children = zip(tag1.children,
                  tag2.children)
    for child1, child2 in children:
        if isinstance(child1, Comment):
            break

        if contentMatches(tag1, tag2):
            regex += tag1.content
        else:
            regex += "#PCDATA"

        if tagsMatch(child1, child2):
            regex +=
                generateRegex(child1, child2)
        else:
            regex += "(<" + child1 +
                ".*?>.*?</" + child1 + ">)+"

    regex += "</" + tag1.name + ">"

```