# Implementation of Web Scraper Bot: Web Harvesting

**NandanAS[1], Rahul S Niranjan[2], PB Rahul Choudhary[3], Karthik Srinivas[4], Nilesh Kumar Singh[5], Kavya P Hathwar[6]**

Student, Department of CSE, BNM Institute of Technology, Bangalore, India [1]

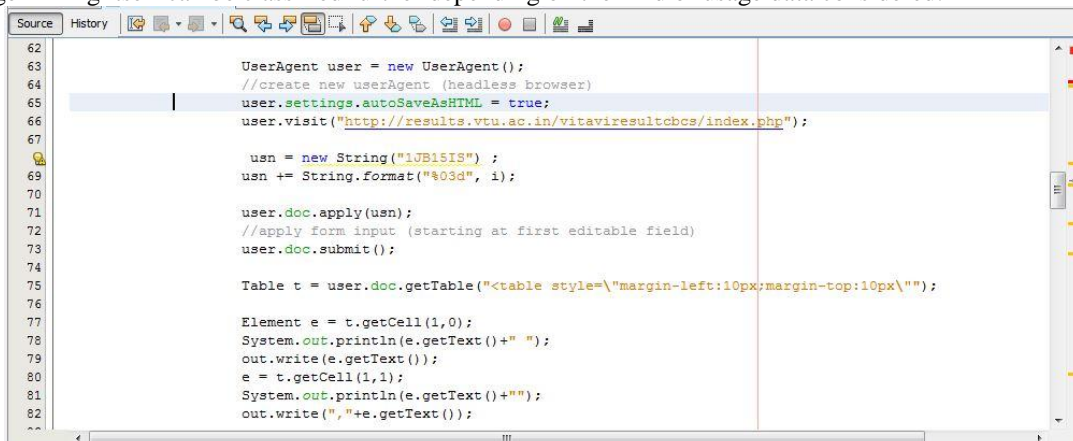Student, Department of CSE, BNMIT, Bangalore, India [2, 3, 6]

Student, Department of ECE, JSSATE, Bangalore, India[4, 5]

**Abstract**: Web Harvesting or Web-Scraping also called web data extraction, are various methods of collecting information from across the internet. It is essentially a form of data mining. Programs are written to mine the data and to convert it into a meaningful and useful structure. In this paper we demonstrate a code/program written to harvest web-data from a particular web-site and to display the same in different file formats. The university results of a class of Students arescraped from the web and are being stored and calculated. This code can be reused several number of times and may also be altered to suit the desired/intended application. Creating a customized score sheet of all students in the college or university is a tedious task. In this paper a web scraper bot is employed to do the same within minutes.

**Keywords**: Web Harvesting or Web-Scraping, essentially a form of data mining, reused several number of times, Creating a customized score sheet, web scraper bot, regular expression (regex).

## I. INTRODUCTION

Web mining is employed to extract data from the web [1]. It is the application of data mining techniques to extract knowledge from web data, including web documents, usage logs of website and also web content which can be further conditioned to provide monetary benefits. Figure1 Shows the Taxonomy of web mining. Web content mining is the process of extracting useful information from the contents of web documents.Content data is the collection of facts a web page is designed to contain [2]. It may consist of text, image, audio, video or structured records such as list and tables. Web structure mining is the process of discovering structure information from the web. It is further divided into hyperlinks and document structure. Web usage mining is the application of data mining techniques to discover interesting usage patterns from web usage data, in order to understand and better serve the needs of web-based applications. Usage data captures the identity or origin of web users along with their browsing behaviours at a website. Web usage mining itself can be classified further depending on the kind of usage data considered.



```
62
63      UserAgent user = new UserAgent();
64      //create new userAgent (headless browser)
65      user.settings.autoSaveAsHTML = true;
66      user.visit("http://results.vtu.ac.in/vitaviresultcbcs/index.php");
67
        usn = new String("1JB15IS") ;
69      usn += String.format("%03d", i);
70
71      user.doc.apply(usn);
72      //apply form input (starting at first editable field)
73      user.doc.submit();
74
75      Table t = user.doc.getTable("<table style=\"margin-left:10px;margin-top:10px\"");
76
77      Element e = t.getCell(1,0);
78      System.out.println(e.getText()+" ");
79      out.write(e.getText());
80      e = t.getCell(1,1);
81      System.out.println(e.getText()+"");
82      out.write(","+e.getText());
```

Figure 2 shows the URL of a university score display website

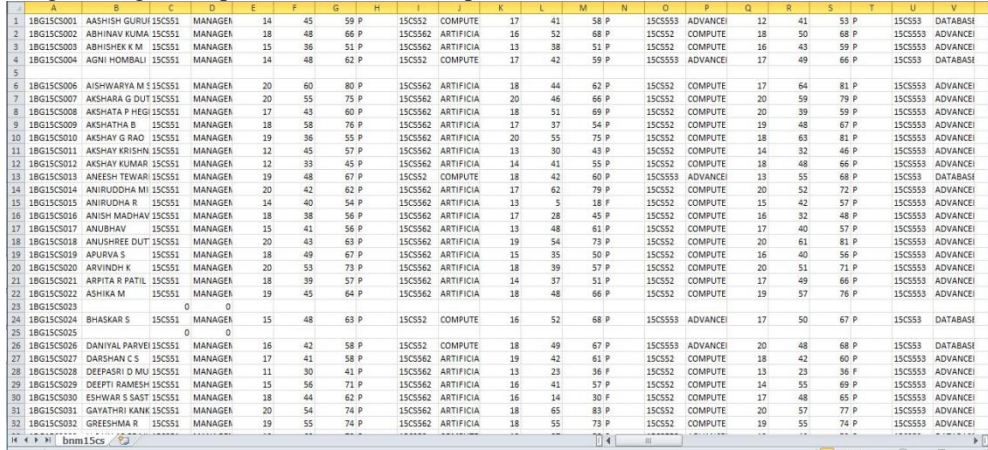## II. PROBLEM STATEMENT AND POSSIBLE SOLUTION

A.    Problem

Creating a customized score sheet of all the students in a university is a tedious task. A  DEO has to fill more than 50,000 rows on an average. CGPA calculation, credit score realization and categorization is another overhead that cannot be ignored.

B.    Solution
- A Web Scraper Bot

1. An intelligent way of handling things and is very much required in this space. An Algorithm is the need of the day.
2. People can expect a lot of accuracy in the results obtained.
3. The work of gathering the required data will be completed within minutes/seconds.



Figure 3 shows the USNs being scraped from the web and are displayed in CSV format.



Figure 4 shows the if-else statement for subject credits

## III.    METHODOLOGY

A java code is written to systematically harvest the data from the web. The URL of interest is selected and is pointed to by the admin to the web scrapper bot.Figure 2 shows the URL of a university score display website. The University Seat Number (USN) is incremented from its initial value to a stoppage value. USNs are scraped from the web and are displayed in CSV format as shown in figure 3.

- SUBJECT CREDITS

Each Subject has its own credit Points. Core subjects carry 4 points. Optional subjects carry 3 points. Lab/Practical sessions carry 2 points. Due to the variation in credit points the CGPA calculation is bit tedious when done manually. The web scraper bot does the job for you. Figure 4 shows the code for the above.

- Case Statements for awarding GRADE

Figure 5 shows the code, where case statement is used to calculate the GRADE of a result of a candidate. Figure 6 shows pattern matching using regex.Everything is parsed using regular expression (regex) [3]. Regex can be used to provide an effective and compact solution to a problem.



Figure 5 shows the code, where case statement is used to calculate the GRADE of a result of a candidate.

Figure 6 shows pattern matching using regex.

## IV.    RESULTS

The flow diagram of a 'VTU Result Scraper Bot' which calculates the CGPA and awards grades to a candidate belonging to certain stream is as shown in figure 7.The code is executed and the results are displayed in the output window. Figure 8 shows the results with CGPA and total credit points [4].
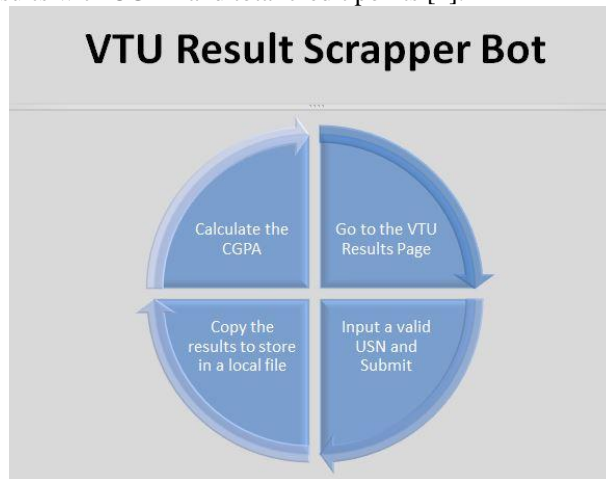


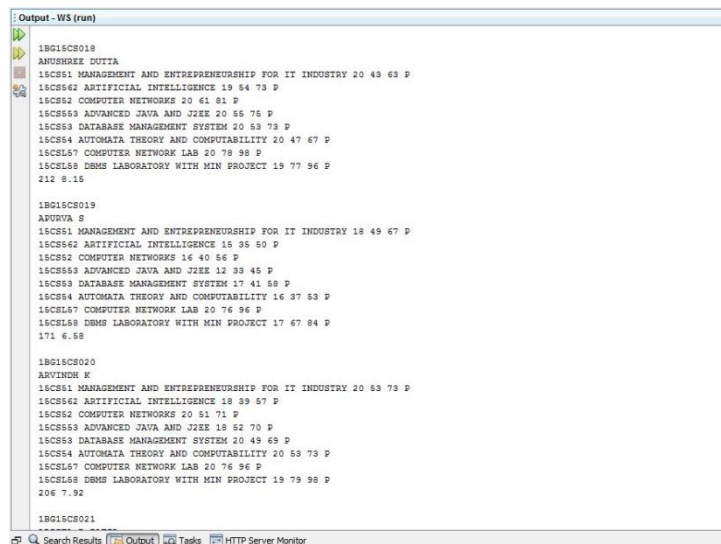Figure 7 shows 'VTU Result Scraper Bot'



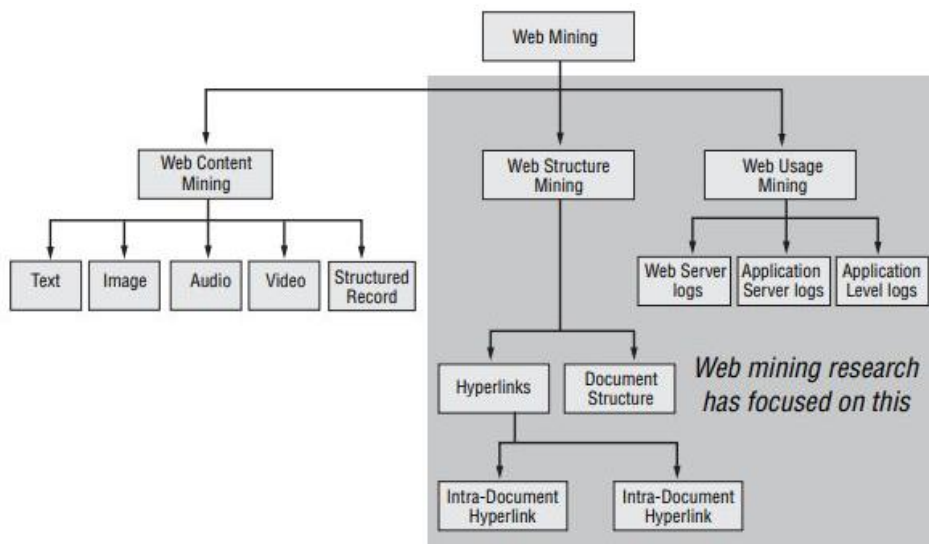Figure 8shows the results with CGPA and total credit points.

Figure 1 shows the taxonomy of web mining

## V.    CONCLUSIONS

In this paper we have developed a 'VTU Web Scraper Bot'. An intelligent way of handling things and is very much required in this space. An Algorithm is written and executed. People can expect a lot of accuracy in the results obtained. The work of gathering the required data will be completed within minutes/seconds. The 'VTU Result Scraper Bot' calculates the CGPA and awards grades to a candidate belonging to certain stream automatically. Creating a customized score sheet of all students in the college or university is a tedious task.This code can be reused several number of times and may also be altered to suit the desired/intended application.

## REFERENCES

[1]    Data Mining, Internet Marketing and Web Mining - ijarcce
[2]    https://www.ijarcce.com/upload/2017/march-17/IJARCCE%20117.pdf
[3]    Web Mining - Data Analysis and Management Research Group dmr.cs.umn.edu/papers/p2004_4.pdf
[4]    RegExr: Learn, Build, & Test RegEx https://regexr.com/
[5]    (cbcs) regulations governing the - VTU-www.vtu.ac.in/pdf/regulation/becbcs.pdf

## OUR GUIDE

**VISHESH S** born on 13th June 1992, hails from Bangalore (Karnataka) and has completed B.E in Telecommunication Engineering from VTU, Belgaum, Karnataka in 2015. He also worked as an intern under Dr. Shivananju BN, former Research Scholar, Department of Instrumentation, IISc, Bangalore. His research interests include Embedded Systems, Wireless Communication, BAN and Medical Electronics. He is also the Founder and Managing Director of the corporate company Konigtronics Private Limited. He has guided over a hundred students/interns/professionals in their research work and projects. He is also the co-author of many International Research Papers. He is currently pursuing his MBA in e-Business and PG Diploma in International Business. Presently Konigtronics Private Limited has extended its services in the field of Software Engineering and Webpage Designing. Konigtronics also conducts technical and non-technical workshops on various topics.