**Web Scarping Toolkit Report**

For this Module 6 assignment I am using 3 Web Scraper Toolkits to extract Super Bowl (Most valuable players') data from url

http://www.espn.com/nfl/superbowl/history/mvps

Information being extracted are as follows:

1)Number(Super Bowl Number)

2)Player Name(Most Valuable Player)

3)Highlight (Highlight of the match)

Web scraper tool kits used are

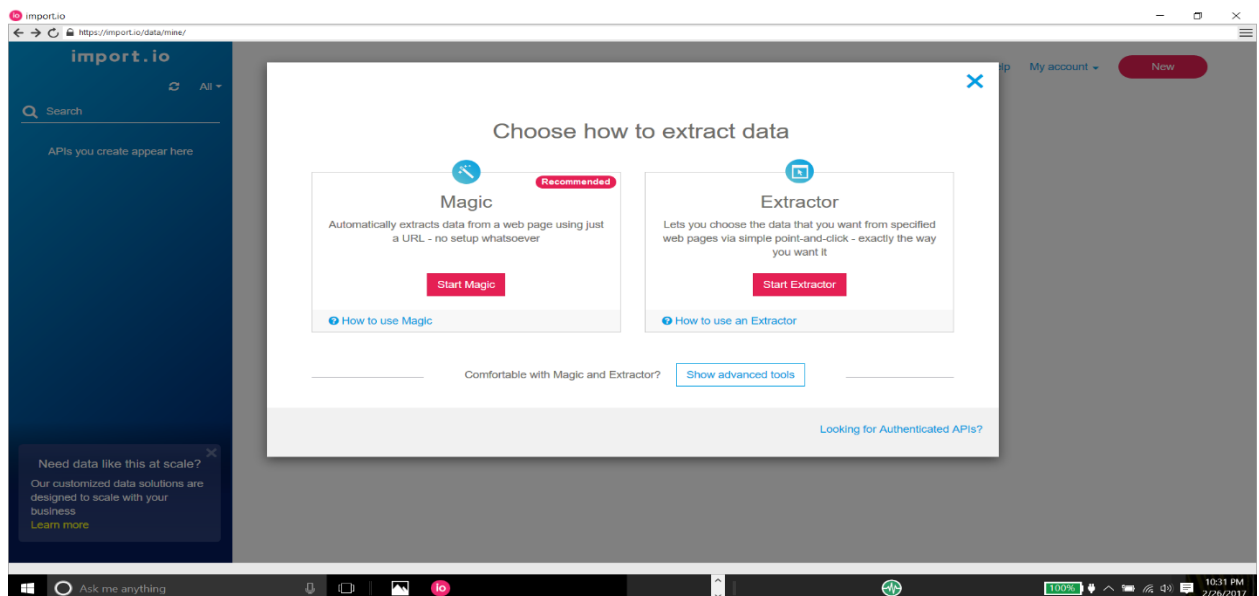**1) Impot.io**

Pros :

1) It automatically can detect tables and simply extract data without choosing it. Smarter algorithm.
2) It has features like "Magic" (automatically extract data) and "Classic Extractor"(user has to choose which data to add)

Cons :

1) It is not free, though free trial  version is available.
2) Slower in extraction sometimes.
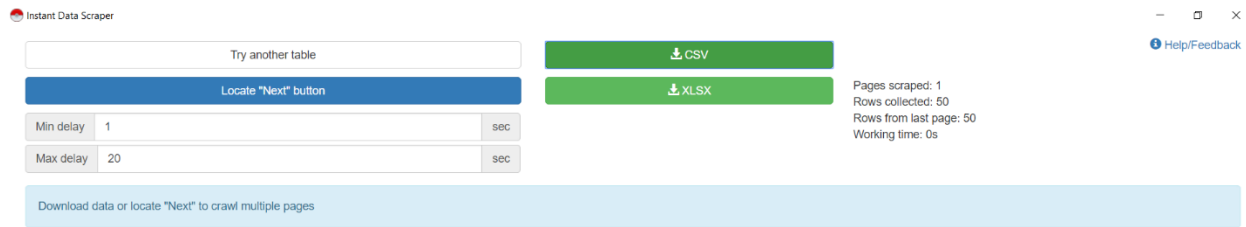3) Not able to modify searching pattern.
4) No access to Xpath.

**2)Instant Data Scraper (Google Extension)**

Pros :

1)Easy to install as being google extension.

2)Easy to use : go to page you want scrape and click on the logo.

3)It is free and fast.

4)Extract tables automatically.

Cons:

1)Not able to choose all vectors and tables of the page together.

2)Data cleaning and tangling should be done manually.

**3)rvest**

Pros:

1)Easily available as rvest package and its free.

2)Coding is easier than RCurl toolkit.

3)Reading data is easier.

4) Less cleaning and editing is required as we can select required data only and directly store into data frame.

Cons:

1)Requires source code of required web page .

2)Knowledge of HTML and nodes is required to extract.

3)Required coding to extract data instead of drag and drop data as import.io

Easiest to use -> import.io

Recommended to use ->rvest