

MEMO

What	Why	How
<p><b>Import.io scraper data without programming.</b></p>	<p>Import.io is a web scraper, which is automated to a high degree. Many tasks can be solved simply by copy and pasting the URL into its web page. You must register as a user to retrieve data as a csv file. It is free.</p> <p>You can also download a free version to your computer, where you earlier have been able to do more advanced scraping tasks. Finally, you can buy the enterprise version, and then import.io solves scraping tasks for you. All of their current and future product development efforts are now focused on the new web version of Import.io. They are working to bring across all of the old capability of the Desktop version, and they promise to include new features.</p>	<p>Basic site: <a href="https://import.io/">https://import.io/</a></p> <p>Free download: <a href="https://www.import.io/download/download-info/">https://www.import.io/download/download-info/</a></p> <p>Tutorial: <a href="http://support.import.io/">http://support.import.io/</a></p>
<p><b>Magic</b></p>	<p>Magic is the most basic function, where you insert an url and then the software extracts the data in a structured way. If the data is on multiple pages, the limit is 20 pages. Magic handles a number of basic tasks - both in the web version and your downloaded version. import.io uses English separator of comma, which causes problems when loading other language data in Excel. However, it can be solved by using search and replace systematically ( " ", " : and " , "to ; ), moving columns and text columns where ";" is used as a separator. If you are importing a table from the web, it is often easier just to use an add-on for Google Chrome (Table Capture), which detects tables on a webpage, and you can then select the relevant table and then copy-paste into Excel. You can also use web import from Excel directly.</p>	
<p><b>Extractor</b></p>	<p>After installing import.io on your computer, you can use the more advanced methods. Extractor is not working on the web yet. Start by clicking New and select Extractor:</p> <div data-bbox="375 1435 509 1503" style="border: 1px solid black; padding: 2px; display: inline-block; background-color: #e91e63; color: white;">New ▾</div> <p>Insert the link into the browser window and press enter :</p> <div data-bbox="371 1574 1238 1648" style="border: 1px solid gray; padding: 2px;"> <span style="font-size: 0.8em;">io import.io Magic   Web Data Platform &amp; Free Web Scraping Tool</span> <span style="font-size: 0.8em;">io My Data: Laeger1</span> <span style="font-size: 0.8em;">io Get started with import.io</span> <span style="float: right; font-size: 0.8em;">New</span> </div> <p>Change the button from OFF to ON:</p> <div data-bbox="375 1704 587 1765" style="border: 1px solid gray; padding: 2px; display: inline-block; background-color: #ccc;">OFF</div> <p>Name the first column and the click on the first cell with the content of the column that you want extracted:</p> <div data-bbox="375 1845 647 1982" style="border: 1px solid gray; padding: 5px; background-color: #fff9c4; width: fit-content;">Aage Christiansen</div> <p>Click on Many rows.</p>	

MEMO



**Many rows**

Give your number a header name. Click the + New column. Select the first cell in the next column. Give it a name until you have selected all the columns. Click Done. Give your API a name and click Publish. You can delete columns and you can later edit the Extractor.

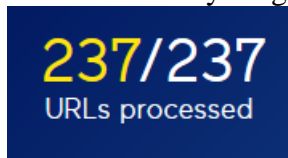
**Bulk Extraction**

Once you have built your Extractor, you can put it to extract data from many URLs, if they have the same structure. You need to understand the method for going to the next page. The easiest way if it is just numbers from 1 to the last page like this: <http://www.fda.gov/Safety/Recalls/ArchiveRecalls/2015/default.htm?Page=1> in the last page: <http://www.fda.gov/Safety/Recalls/ArchiveRecalls/2015/default.htm?Page=15> Here you can make the 15 urls in a spreadsheet. First part of the url is this and it is constant: <http://www.fda.gov/Safety/Recalls/ArchiveRecalls/2015/default.htm?Page=> Second part is the number (1 to 15). You can merge that by the formula =A1&B1. In total we will then build 15 URL's to be pasted into import.io Here is a Google Spreadsheet with the URL's prepared: [kortlink.dk/ks7f](http://kortlink.dk/ks7f) In the example below the webpages are built by three parts.

D3 : X ✓ fx =A3&B3&C3

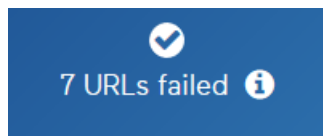
	A	B	C	D	E
1	Url1	No	Url2	Samlet Url	
2	<a href="http://sundhedsstyrels">http://sundhedsstyrels</a>	0 &max=108	<a href="http://sundhedsstyrelsen.dk/da">http://sundhedsstyrelsen.dk/da</a>		
3	<a href="http://sundhedsstyrels">http://sundhedsstyrels</a>	10 &max=108	<a href="http://sundhedsstyrelsen.dk/da">http://sundhedsstyrelsen.dk/da</a>		
4	<a href="http://sundhedsstyrels">http://sundhedsstyrels</a>	20 &max=108	<a href="http://sundhedsstyrelsen.dk/da">http://sundhedsstyrelsen.dk/da</a>		

You collect the three parts of the formula: = A3 & B3 and C3 in cell D3. Then select all the 237 Url 's in Excel and copy them. You go back to import.io, choose Bulk Extract, insert the many urls and press Run Queries waiting. After a while, the screen shows that everything is pulled out :

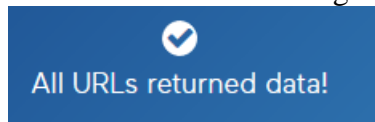


MEMO

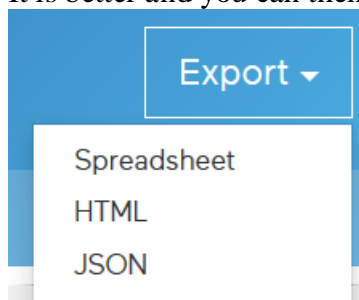
And yet, sometimes. Error detected :



Click on i and run them again. Then it displays :



It is better and you can then download all - select spreadsheet



If you get problems with special characters, you can instead save it as HTML and open the file directly in Excel.

**Extraction of tax data**

Identify the structure of links in this page:

<https://www.revenue.wi.gov/delqlist/nmalla.htm>

For the letter A the url is:

<https://www.revenue.wi.gov/delqlist/NMALLA.htm>

For the letter B the url is:

<https://www.revenue.wi.gov/delqlist/NMALLB.htm>

And for the last url there is a slight change:

<https://www.revenue.wi.gov/delqlist/NUMALL.htm>

You can now construct the 27 urls copying the letters to word, replacing spaces with a new line, copy-pasting it to Excel and then make a formula combining 3 cells.

https://www.revenue.wi.gov/delqlist/NMALL	A	.htm	=a1&b1&c1
---	---	------	-----------

And then it's the same as in the example below, making a new API from the first website and then copy and paste all the constructed urls into a bulk extraction.

**Extraction of many urls**

Often it is not so easy to identify the structure of links and you have to split the scraper in two API's. One for extraction of links. The other for extraction of data. You then use the first API as input to the second.

MEMO

link Link ▾

Show advanced settings (What is this?)

Cancel ✓ Done

You save this scraper and then build a new scraper to extract the data.

Then write the name of the relevant API into :

How would you like to use this API?

URLs from another API ▾

Chain APIs i

- Source API or Data Set i
  - vl114urls|
    - Source URL
      - 
      - link i

vl114urls returned 114 rows, this is a preview of the first 100

```

http://iframe.vl.dk/vlgruppe.php?area=København&id=1
http://iframe.vl.dk/vlgruppe.php?area=København&id=2
http://iframe.vl.dk/vlgruppe.php?area=København&id=3
http://iframe.vl.dk/vlgruppe.php?area=København&id=4
http://iframe.vl.dk/vlgruppe.php?area=København&id=5
    
```

**Plan for training in Kathmandu:**

Sign up and get import.io downloaded. Good if it's done in advance. Then we will scrape these two websites:  
<https://www.revenue.wi.gov/delqlist/nmalla.htm>  
<http://www.fda.gov/Safety/Recalls/ArchiveRecalls/2015/default.htm>  
 The first website we will try in Magic. And when that is done we will do a bulk extract from magic, copy-pasting the links from this Google Spreadsheet:  
[https://docs.google.com/spreadsheets/d/1\\_T7eT3VvR2gPW8JFXX-RCnNEuyY7cCdUmFHYzShYpH8/edit?usp=sharing](https://docs.google.com/spreadsheets/d/1_T7eT3VvR2gPW8JFXX-RCnNEuyY7cCdUmFHYzShYpH8/edit?usp=sharing)  
 And then we will download the data as a spreadsheet and html and open it in Excel. Next step is to use the Extractor on the second website, the fda-data. And then also from the spreadsheet copy-pasting the links for bulk extraction. Run it and then download and clean the data.  
 In the end a demo on how one API can be the input for another scraper.