

A Data Scraper's Toolbox



Alice Corona

Data Journalist &
Researcher

data-journalism-tools.silk.co

[@alice_corona](https://twitter.com/alice_corona)
alice-corona.nl





“process of automatically
collecting information
from the Web”

~ Wikipedia on “Web Scraping”

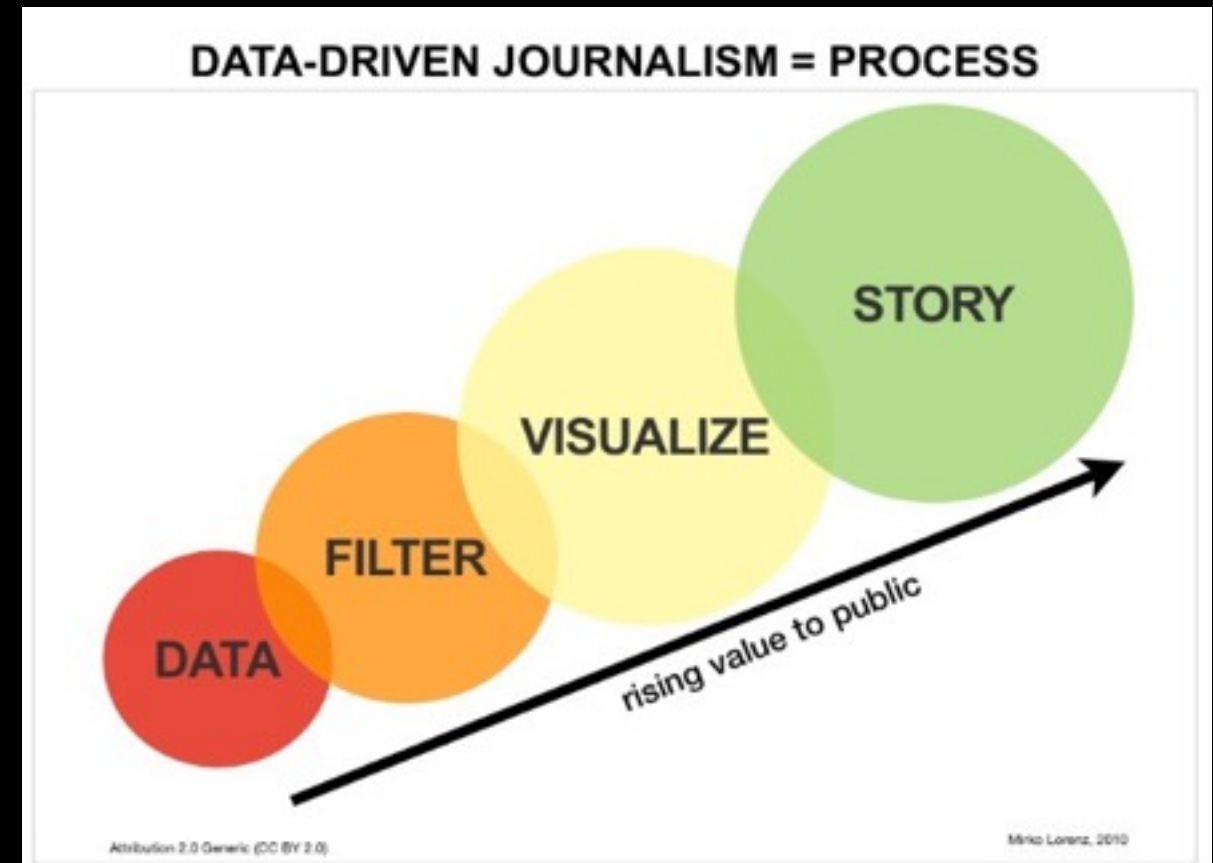


- Manual copy-paste
- Text grepping / RegEx
- HTML parsing (CSS, XPATH...)
- Web Scraping Software



Despite its fame, scraping is
a highly creative process.
It requires...

- Problem solving
- Lateral thinking
- Planning for efficiency
- Prediction and planning for possible outcomes and uses
- Methodological choices affect the whole outcome





There cannot be a course on web scraping teaching you 'solutions'

We can dive into the fine art of scraping by learning tools
and techniques. Then each **individual** personality can
adapt, mix and mash this knowledge according to each
individual problem



Copy-Paste



Data Import in Google Sheets



Chrome's *Scrape-Similar*



X-Path Parsing in Google Sheets



SearchLink

“SearchLink is a System Service for OS X which handles searching multiple sources and automatically generating Markdown links for text”

<http://brettterpstra.com/projects/searchlink/>



APIs



Scripts



Scraping Software



“Hidden Secrets”



is scraping legal?



is scraping ethical?



Consider...

- Brute force attacks and DoS (prefer APIs, when available)
- Your use of the data (copyright, privacy...)

A Data Scraper's Toolbox



Alice Corona

Data Journalist &
Researcher

data-journalism-tools.silk.co

[@alice_corona](https://twitter.com/alice_corona)
alice-corona.nl