

# STORY HUNT

*Workshop #3: Scraping & Cleaning*

# Today's Program:

Introduction to Scraping

Web Scraping

*Demo: Chrome Scraper Plugin*

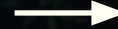
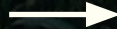
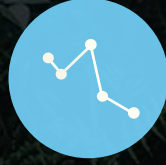
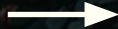
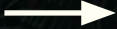
Introduction to Data Cleaning

Basic Data Cleaning Techniques:

*Demo: Open Refine*



# The Data Pipeline: **Theory**



## Finding & Getting

- \* Where to find Data?
- \* Data portals
- \* Freedom of Information Requests
- \* **Webs-scraping**

## **Cleaning**

- \* **transform documents into a structured data**

## Analyzing

- \* How do I spot:
- \* Patterns & Trends ?
- \* Outlier ?
- \* Connections between Columns, Datasets
- \* relation between Data & Context

## Visualizing

- \* How to visualize the essence of a finding:
- \* appealing
- \* informative

## Communicating

- \* What's the right:
- \* Format?
- \* Medium?
- \* Target group?

# The Data Pipeline: **Tools**



## Finding & Getting

- \* Google Sheets
- \* Webscraper
- \* Unbubble
- \* Tabula



## Cleaning

- \* Open Refine
- \* Trifacta
- \* Google Sheets / Excel
- \* Python, R

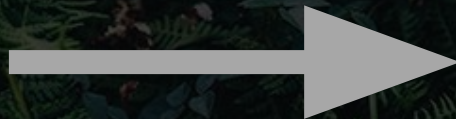


# What is Scraping?

# What is Scrapping:



unstructured  
Documents







Structured  
tabular data

3 Steps:

Fetch the document  
(most likely from the  
Internet)

Follow futher linked  
documents (especially in  
web scraping)

Extract interesting parts  
and store them in a  
structured format



# When is Web-Scraping useful?

- *avoid tedious big copy-pasting jobs*

## What is the challenge?

- *one scraper only works for documents from the same source (e.g. Website, PDF Documents frame the same collection)*

# Chrome Scrapping Plugin

*A web scraper with user interface*



# Chrome Scrapping Plugin

*Download-Link: <http://bit.ly/RmLaUv>*

# Basics: Data Cleaning



# Basics: Clean Data

*Structured, machine-readable Data*

# What makes Data Clean?

Consistent spelling of corresponding values (U.S. ≠ USA)

WatchOut: No trailing Spaces in front/ behind a value

Unique column labels

No hidden rows / columns

No empty rows / columns



# Open Refine

*Free Open Source Tool for working with messy  
Data*

# Open Refine

<http://openrefine.org/download.html>



# NEXT WEEK: Basics of Statistical Analysis