# Four-day workshop
# Automated Content Analysis with Python
Day 4

Damian Trilling

d.c.trilling@uva.nl
@damian0604
www.damiantrilling.net

Afdeling Communicatiewetenschap
Universiteit van Amsterdam

26–9–2017

# Today

**❶** We put on our magic cap, pretend we are Firefox, scrape all comments from GeenStijl, clean up the mess, and put the comments in a neat CSV table

**❷** OK, but this surely can be doe more elegantly? Yes!

**❸** Next steps

We put on our magic cap, pretend we are Firefox, scrape all comments from GeenStijl, clean up the mess, and put the comments in a neat CSV table

http://www.geenstijl.nl/mt/archieven/2014/05/ik_ben_niet_gek_ik_ben_een_ehm_sja_eeeh_wie_het_weet_mag_het...  Reader

Q Google

GS
GEENSTIJL

Vreemde vogel vogelvrij bij dodenherdenking Dam

Vreemde vogel bij dodenherdenking Dam

0:00 / 4:38

### HEADLINES

05-05

Daar is pathologische leugenaar Rehwinkel weer!

Kijk. Zo leren kindjes in Tsjechie over de EU

Gezocht. Dief met tatoeage en motorfiets

Wientjes: 'Hey paupers, de EU is fantastisch hoor'

Anarchie! Wildplassen is ook vrijheid

Vreemde vogel vogelvrij bij dodenherdenking Dam

Aangifte tegen GeenStijl – De Q en A Kloosried

Generaal van Uhm: 'Polderjihadisten zijn

Filmpje uit het rauwrealistische archief van Amsterdamse stadscineast Frank Buis. En met 'rauwrealistisch' bedoelen we: de beelden die u niet ziet op de stream-met-nietszeggend-commentaar-in-stemmige-fluistertoon van de Staatsomroep. Vlak voor de Dodenherdenking dook bovenstaande dodo op de Dam op. Plastic Hannibal Lecter maskertje voor de mond. Rare oranje klomp aan de voet. Kogelhuizen in de kroon. Daar wilde de politie wel even mee babbelen, zo vlak voor de aankomst van de koninklijke nazi-nazaat en zijn Junta-meisje. Want doden herdenken op je eigen manier is 1 ding. Maar jezelf daarbij kòdig maken wel buiten de voorschriften van de vrijheid. Wat moesten die 69 kindjes in het buitenlandse gasten-vak wel niet denken? Daarom. Goed dat de politie deze malle meneer en zijn aluminium verpakkingje eventjes naar een veilige afstand van de vrijheidsviering voerde. (Video: Frank Buis/Royal Press Amsterdam)

Van Rossem | 05-05-14 | 11:11 | Link | 229 reacties

#### REAGUURSELS

Tsja, ik zou m ook niet in mn buurt willen hebben, die vreemde vogel.

Lepo | 05-05-14 | 11:13

Volgend jaar stropdascontrole voor de heren en hoedjescheck voor de dames. De volledige lijst met goedgekeurde kleding kunt u vinden op Postbus51.nl.

DUMPERT

DE BESTE DEALS!
kortingen tot 80%
De Telegraaf    BEKIJK DEALS ►

POWNEWS
POWNEWS – ELKE WERKDAG, 22:45, NED3
Eerten Delft strijdt tegen Volkert
16:40 Jongens melden zich na mishandeling
16:32 Polio gaat wereld slopen
16:24 Verkiezingsposters plakken niet
15:49 Stemfie = celstraf in Zuid-Afrika
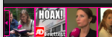15:38 Update: ontvoeringen bevestigd
15:26 Timmermans sipt op rijm over Roda
15:19 Tientallen migranten verdronken
15:13 Natte vrouwen op supergiltbaan

GEENSTIJL TV
HOAX!

DASKAPITAL
16:58 Pensioenpolder jaagt fondsen naar België

# REAGUURSELS

Tsja, ik zou m ook niet in mn buurt willen hebben, die vreemde vogel.

Lepo | 05-05-14 | 11:13

Volgend jaar stropdascontrole voor de heren en hoedjescheck voor de dames. De volledige lijst van goedgekeurde kleding kunt u vinden op Postbus51.nl.

Uiteraard bent u extra welkom als u Abercrombie & Fitch draagt.

rara | 05-05-14 | 11:15

Gewoon even de afdeling psychiatrie bellen, wie ze missen: Klaar!

Mazzeltov | 05-05-14 | 11:16

netjes opgelost toch?

--sql error-- | 05-05-14 | 11:16

```html
rond de Duitse grens. In dat geval moet u niet bellen met 0900-8844, maar het heerschap vriendelijk
rechtmatige eigenaar te vinden. Hierbij willen we u nogmaals wijzen op het feit dat deze man een ta
heeft. En u weet hoe mensen met tatoeages kunnen zijn. Dank voor uw aandacht. GeenStijl RegioNieuws

        <footer>Van Rossem | <time pubdate datetime="2014-05-05T14:14">05-05-14 | 14:14</time>
href="http://www.geenstijl.nl/mt/archieven/2014/05/das_toch_niet_normaal.html">Link</a> | <a
href="http://www.geenstijl.nl/mt/archieven/2014/05/das_toch_niet_normaal.html#comments" class="comm
    </article>
<!-- google_ad_section_end -->

        <section id="comments">
            <header>
                <h2>Reaguursels</h2>
            </header>

        <div class="commentlist">
            <article id="c192607851">
                <p>Kijk volgende week eens in Irak.</p>
                <footer>Mammeloe | 05-05-14 | 14:16 </footer>
            </article>
            <article id="c192607941">
                <p>Tattooboer, kom dr maar in.</p>
                <footer>RegseHippy | 05-05-14 | 14:16 </footer>
            </article>
            <article id="c192607951">
                <p>Door de afkomst is Bennie Jolink dus al uitgesloten.</p>
                <footer>Ing. aclonen | 05-05-14 | 14:16 </footer>
```

# Let's make a plan!

## Which elements from the page do we need?

- What do they mean?
- How are they represented in the source code?

## How should our output look like?

- What *lists* do we want?
- . . .

And how can we achieve this?

## Operation Magic Cap

## Operation Magic Cap

❶ Download the page

## Operation Magic Cap

❶ Download the page
- They might block us, so let's do as if we were a web browser!

## Operation Magic Cap

1. Download the page
   - They might block us, so let's do as if we were a web browser!

2. Remove all line breaks (\n, but maybe also \n\r or \r) and TABs (\t): We want one long string

## Operation Magic Cap

1. Download the page
   - They might block us, so let's do as if we were a web browser!
2. Remove all line breaks (\n, but maybe also \n\r or \r) and TABs (\t): We want one long string
3. Isolate the comment section (it started with <div class="commentlist"> and ended with </div>)

## Operation Magic Cap

1. Download the page
   - They might block us, so let's do as if we were a web browser!

2. Remove all line breaks (\n, but maybe also \n\r or \r) and TABs (\t): We want one long string

3. Isolate the comment section (it started with <div class="commentlist"> and ended with </div>)

4. Within the comment section, identify each comment (<article>)

## Operation Magic Cap

①  Download the page
   - They might block us, so let's do as if we were a web browser!
②  Remove all line breaks (\n, but maybe also \n\r or \r) and TABs (\t): We want one long string
③  Isolate the comment section (it started with <div class="commentlist"> and ended with </div>)
④  Within the comment section, identify each comment (<article>)
⑤  Within each comment, seperate the text (<p>) from the metadata <footer>)

## Operation Magic Cap

**1** Download the page
  - They might block us, so let's do as if we were a web browser!

**2** Remove all line breaks (\n, but maybe also \n\r or \r) and TABs (\t): We want one long string

**3** Isolate the comment section (it started with <div class="commentlist"> and ended with </div>)

**4** Within the comment section, identify each comment (<article>)

**5** Within each comment, seperate the text (<p>) from the metadata <footer>)

**6** Put text and metadata in lists and save them to a csv file

## Operation Magic Cap

**1** Download the page
   - They might block us, so let's do as if we were a web browser!

**2** Remove all line breaks (\n, but maybe also \n\r or \r) and TABs (\t): We want one long string

**3** Isolate the comment section (it started with <div class="commentlist"> and ended with </div>)

**4** Within the comment section, identify each comment (<article>)

**5** Within each comment, seperate the text (<p>) from the metadata <footer>)

**6** Put text and metadata in lists and save them to a csv file

**And how can we achieve this?**

```
1   from urllib import request
2   import re
3   import csv
4
5   onlycommentslist=[]
6   metalist=[]
7
8   req = request.Request("http://www.geenstijl.nl/mt/archieven/2014/05/
        das_toch_niet_normaal.html", headers={"User-Agent" : "Mozilla
        /5.0"})
9   tekst=request.urlopen(req).read()
10  tekst=tekst.decode(encoding="utf-8",errors="ignore").replace("\n"," ").
        replace("\t"," ")
11
12  commentsection=re.findall(r'<div class="commentlist">.*?</div>',tekst)
13  print (commentsection)
14  comments=re.findall(r'<article.*?>(.*?)</article>',commentsection[0])
15  print (comments)
16  print ("There are",len(comments),"comments")
17  for co in comments:
18      metalist.append(re.findall(r'<footer>(.*?)</footer>',co))
19      onlycommentslist.append(re.findall(r'<p>(.*?)</p>',co))
20  writer=csv.writer(open("geenstijlcomments.csv",mode="w",encoding="utf
        -8"))
21  output=zip(onlycommentslist,metalist)
22  writer.writerows(output)
```

| | |
|---|---|
| ['Kijk volgende week eens in Irak.'] | ['Mammeloe \| 05-05-14 \| 14:16 '] |
| ['Tattooboer, kom dr maar in.'] | ['RegseHippy \| 05-05-14 \| 14:16 '] |
| ['Door de afkomst is Bennie Jolink dus al uitgesloten.'] | ['Ing. eslapen \| 05-05-14 \| 14:16 '] |
| ['He psstt, meisje - motor kopen?'] | ['Snackbar van Allah \| 05-05-14 \| 14:17 '] |
| ['"Daarnaast is de vermoedelijk van Marokkaanse afkomst."<br /> Is dit reeds de gekuiste versie, we mogen Marokkanen ook al geen dader meer noemen?'] | ['Causa Sui \| 05-05-14 \| 14:17 '] |
| ['Vast drie puntjes tussen duim en wijsvinger.'] | ['Mark Smith \| 05-05-14 \| 14:17 '] |
| ['3 stippen op de rechterhand, in de buurt van de duim? '] | ['Harreeee \| 05-05-14 \| 14:17 '] |
| ['De dader reed zomaar weg met de machine. Ja dat is dan ook de bedoeling van een proefrit.'] | ['Graaf van Egmont \| 05-05-14 \| 14:18 '] |
| ['Wat een kneus. Een jappenchopper die \x802500 of minder waard is stelen.'] | ['eelendsknook \| 05-05-14 \| 14:18 '] |
| ['Dacht dat er in de Achterhoek niet zoveel Marokkanen zaten maar die paar die er zitten weten het ook al te verkloten. Hulde... wat een prachtvolk.'] | ['drekzooi \| 05-05-14 \| 14:19 '] |
| ['Okk Aalten is niet veilig meer qua cultuurverrijking..'] | ['Space2012 \| 05-05-14 \| 14:20 '] |
| ['Ik wist niet dat Marokkanen van ezel en geit al waren overgestapt op motorfietsen. Gaat dat niet veel te snel voor ze? Een fin op een motorfiets, moet niet gekker worden. Nog even en je ziet ze nog eens in politieuniform...'] | ['FrankVeer \| 05-05-14 \| 14:20 '] |
| ['Ook dus..'] | ['Space2012 \| 05-05-14 \| 14:20 '] |

## Some remarks

### The regexp

- `.*?` instead of `.*` means *lazy* matching. As `.*` matches everything, the part where the regexp should stop would not be analyzed (*greedy* matching) – we would get the whole rest of the document (or the line, but we removed all line breaks).

## Some remarks

### The regexp

- .*? instead of .* means *lazy* matching. As .* matches everything, the part where the regexp should stop would not be analyzed (*greedy* matching) – we would get the whole rest of the document (or the line, but we removed all line breaks).

- The parentheses in (.*?) make sure that the function only returns what's between them and not the surrounding stuff (like <footer> and </footer>)

# Some remarks

## The regexp

- .*? instead of .* means *lazy* matching. As .* matches everything, the part where the regexp should stop would not be analyzed (*greedy* matching) – we would get the whole rest of the document (or the line, but we removed all line breaks).

- The parentheses in (.*?) make sure that the function only returns what's between them and not the surrounding stuff (like <footer> and </footer>)

## Optimization

- Only save the 0th (and only) element of the list
- Seperate the username and interpret date and time

# Further reading

## Doing this with other sites?

- It's basically puzzling with regular expressions.
- Look at the source code of the website to see how well-structured it is.

OK, but this surely can be doe more elegantly? Yes!

# Scraping

## Geenstijl-example

- Worked well (and we could do it with the knowledge we already had)

# Scraping

## Geenstijl-example

- Worked well (and we could do it with the knowledge we already had)
- But we can also use existing parsers (that can interpret the structure of the html page)

# Scraping

## Geenstijl-example

- Worked well (and we could do it with the knowledge we already had)
- But we can also use existing parsers (that can interpret the structure of the html page)
- especially when the structure of the site is more complex

# Scraping

## Geenstijl-example

- Worked well (and we could do it with the knowledge we already had)

- But we can also use existing parsers (that can interpret the structure of the html page)

- especially when the structure of the site is more complex

The following example is based on `http://www.chicagoreader.com/chicago/best-of-chicago-2011-food-drink/BestOf?oid=4106228`. It uses the module `lxml`

# What do we need?

- the URL (of course)
- the XPATH of the element we want to scrape (you'll see in a minute what this is)

Let's install a useful add-on: the Firefox XPath Checker



. . . and after that, inspect the website we're interested in:

**Samsung Galaxy Tab 3 10.1 WiFi 16GB prijzen**

★★★★☆ 8,5 (328 reviews)

**Bezorgprijzen vanaf:**

PDAWERELD.NL  € 215,00

redcoon.nl  € 224,82

Beeldschermgrootte 10,1 inch
Totale opslagcapaciteit 16 GB
Beeldschermresolutie 1280 x 800

❯ Alternatieve producten
❯ Alle tablets van Samsung

⊕ Deel   ⊕ Bewaar   ❯ Volg   ⊕ Print

ASUS

Work easy. Play hard.   ⊞ Windows 8

| Prijzen | Overzicht | Specificaties | Reviews | Vragen en Antwoorden | Lokale winkels |

| Winkel (beoordeling) ◆ | Winkelinformatie | Levertijd ◆ | Afhaalprijs ◆ | Bezorgprijs ◆ | ❯ Toon filters |
|---|---|---|---|---|---|
| PDAWERELD.NL ★★★★☆ 8,6 (84 reviews) | Samsung tab 3 p5210 wi... | 1-2 dagen | € 215,00 | € 215,00 | ❯ Bekijk |
| redcoon.nl ★★★★☆ 8,5 (4669 reviews) | Samsung GALAXY Tab 3 1... ✉ | 1-2 dagen | | € 224,82 | ❯ Bekijk  ⬤ Bestel via Kieskeurig.nl |
| 4LAUNCH❯❯ ★★★★☆ 8,7 (283 reviews) | Samsung Galaxy Tab 3 ... ✉ 🛡 | 24 uur | € 222,99 | € 224,99 | ❯ Bekijk |
| CONRAD ★★★★☆ 8,2 (119 reviews) | SAMSUNG Samsung Galaxy... ✉ | 3-4 dagen | | € 229,00 | ❯ Bekijk |
| YOURTABLET ★★★★☆ 8,7 (7 reviews) | samsung galaxy tab 3 1... ✉ | 3-4 dagen | € 229,00 | € 229,00 | ❯ Bekijk |
| belsimpel.nl ★★★★★ 9,0 (1074 reviews) | Samsung Galaxy Tab 3 1... ✉ | 24 uur | € 225,00 | € 231,95 | ❯ Bekijk |
| Trendy Computers | Samsung Galaxy Tab 3 W... | 24 uur | € 225,00 | € 231,99 | ❯ Bekijk |

Stel ons je vraag

kieskeurig | Prijzen | Overzicht | Specificaties | **Reviews** | Vragen en Antwoorden | Lokale winkels

Reviews Samsung Galaxy Tab 3 10.1 WiFi 16GB | ▶ Plaats een review | ▶ Plaats een vraag | Naar prijzen ▶

Consumentenbond

juli 2013

Deze tablet heeft een schermgrootte van 10 inch. Het geteste besturingssysteem is Android 4.2. Dankzij een dikte van 9 millimeter past hij makkelijk in je tas. Voor een 10 inch tablet is hij best licht (511 gram). Het scherm heeft een gemiddelde resolutie van 800 x 1280. Het aantal pixels per inch is 149: vooral teksten zien er redelijk uit. Het geheugen voor opslag bestaat uit 16GB. Het besturingssysteem neemt zelf ook ruimte in beslag waardoor je 9,6GB overhoudt om je apps en bestanden op te slaan. Je kunt het opslaggeheugen uitbreiden met een SD-kaart. Onderweg internetten via 3G is niet mogelijk. Wel kun je altijd je locatie weergeven met de ingebouwde GPS-ontvanger. Een aparte HDMI-aansluiting ontbreekt. Deze tablet heeft 1 USB-on-the go aansluiting. Hiermee kun je bijvoorbeeld een USB-stick aansluiten. Samsung heeft de tablet ook uitgerust met een camera aan de voor- en de achterkant.

✚ Gebruiksvriendelijke bediening
✚ Start zeer snel op
✚ Veel aansluitingen

━ Kan standaard geen websites met Flash inho...

◉ Lees de uitgebreide testresultaten op Consumentenbond.nl

**Review: Karine**
16-04-2014

Ik heb de tab3 nu enkele maanden en ben er heel tevreden over.
Scherm is prima en de batterij gaat toch best lang mee.
De vergelijking met de Ipad gaat voor mij niet op daar de Ipad enkele maanden geleden 400 euro koste en ik de tab3 gekocht...

◉ Lees de review van Karine over de Samsung Galaxy Tab 3 10.1 WiFi 16GB

▶ Reageer op deze review

| Totaal 8.7 ★★★★☆ | beeldscherm 9 | opstartsnelheid 9 | accuduur 8 |
|---|---|---|---|

**Review: lesley84**
14-04-2014

Wat een mooi ding zeg!
Heb zelf de note 10.1 zag deze van de eek en echt een top ding wil deze dus dan ook graag hebben.

◉ Lees de review van lesley84 over de Samsung Galaxy Tab 3 10.1 WiFi 16GB

▶ Reageer op deze review

| Totaal 10 | beeldscherm | opstartsnelheid | accuduur |
|---|---|---|---|

Stel ons je vraag

# Playing around with the Firefox XPath Checker

# Playing around with the Firefox XPath Checker

# Playing around with the Firefox XPath Checker

Some things to play around with:

- // means 'arbitrary depth' (=may be nested in many higher levels)

# Playing around with the Firefox XPath Checker

Some things to play around with:

- // means 'arbitrary depth' (=may be nested in many higher levels)
- * means 'anything'. (p[2] is the second paragraph, p[*] are all

# Playing around with the Firefox XPath Checker

Some things to play around with:

- // means 'arbitrary depth' (=may be nested in many higher levels)
- * means 'anything'. (p[2] is the second paragraph, p[*] are all
- If you want to refer to a specific attribute of a HTML tag, you can use @. For example, every
  *[@id="reviews-container"] would grap a tag like <div id=reviews-container" class="'user-content'

# Playing around with the Firefox XPath Checker

Some things to play around with:

- // means 'arbitrary depth' (=may be nested in many higher levels)
- \* means 'anything'. (p[2] is the second paragraph, p[\*] are all
- If you want to refer to a specific attribute of a HTML tag, you can use @. For example, every
  \*[@id="reviews-container"] would grap a tag like <div id=reviews-container" class="'user-content'
- Let the XPATH end with /text() to get all text

# Playing around with the Firefox XPath Checker

Some things to play around with:

- // means 'arbitrary depth' (=may be nested in many higher levels)
- * means 'anything'. (p[2] is the second paragraph, p[*] are all
- If you want to refer to a specific attribute of a HTML tag, you can use @. For example, every
  *[@id="reviews-container"] would grap a tag like <div id=reviews-container" class="'user-content'
- Let the XPATH end with /text() to get all text
- Have a look at the source code of the web page to think of other possible XPATHs!

# The XPATH

## You get something like

```
//*[@id="tabbedReviewsDiv"]/dl[1]/dd
//*[@id="tabbedReviewsDiv"]/dl[2]/dd
```

# The XPATH

## You get something like

```
//*[@id="tabbedReviewsDiv"]/dl[1]/dd
//*[@id="tabbedReviewsDiv"]/dl[2]/dd
```

The * means "every".
Also, to get the text of the element, the XPATH should end on
/text().

# The XPATH

## You get something like

```
//*[@id="tabbedReviewsDiv"]/dl[1]/dd
//*[@id="tabbedReviewsDiv"]/dl[2]/dd
```

The * means "every".
Also, to get the text of the element, the XPATH should end on
/text().

## We can infer that we (probably) get all comments with

```
//*[@id="tabbedReviewsDiv"]/dl[*]/dd/text()
```

## Let's scrape them!

```
 1   from lxml import html
 2   from urllib import request
 3
 4   req=request.Request("http://www.kieskeurig.nl/smartphone/product
         /2334455-apple-iphone-6/reviews")
 5   tree =
 6
 7   html.fromstring(request.urlopen(req).read().decode(encoding="utf-8",
         errors="ignore"))
 8   reviews = tree.xpath('//*[@class="reviews-single__text"]/text()')
 9
10   # remove empty reviews and remove leading/trailing whitespace
11   reviews = [r.strip() for r in reviews if r.strip()!=""]
12
13   print (len(reviews),"reviews scraped. Showing the first 60 characters of
           each:")
14   i=0
15   for review in reviews:
16       print("Review",i,":",review[:60])
17       i+=1
```

# The output – perfect!

```
1  63 reviews scraped. Showing the first 60 characters of each:
2  Review 0 : Apple maakt mooie toestellen met hard- en software uitsteken
3  Review 1 : Vanaf de iPhone 4 ben ik erg te spreken over de intuitieve i
4  Review 2 : Helaas ontbreekt het Apple toestellen aan een noodzakelijk i
5  Review 3 : Met een enorme mate van pech hebben wij als (beschaafd!!) ge
```

# Recap

### General idea

1. Identify each element by its XPATH (look it up in your browser)
2. Read the webpage into a (looooooong) string
3. Use the XPATH to extract the relevant text into a list (with a module like lxml)
4. Do something with the list (preprocess, analyze, save)

Alternatives: scrapy, beautifulsoup, regular expressions, . . .

# Last remarks

## There is often more than one way to specify an XPATH

1. You can usually leave away the namespace (the x:)
2. Sometimes, you might want to use a different suggestion to be able to generalize better (e.g., using the *attributes* rather than the *tags*)
3. in that case, it makes sense to look deeper into the structure of the HTML code, for example with "Inspect Element" and use that information to play around with in the XPATH Checker