

Natural Language Processing

Regexes

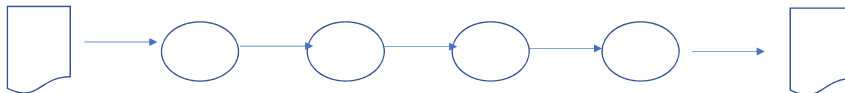
Beautiful Soup

NLTK (Natural Language Toolkit)

Objective:

For some given webpage or blog or tweet:

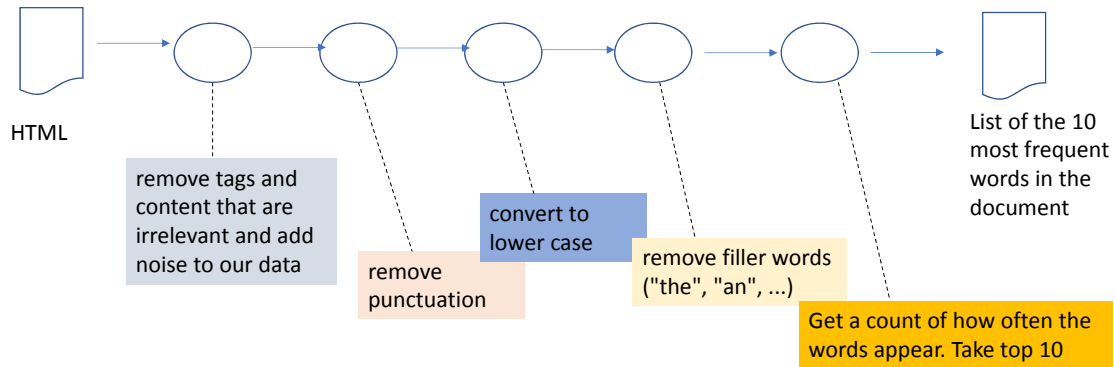
- write code to extract all the text
- analyze the text to determine what the page is about
- build JSON-LD based on the text
- do it using code



Data Flow : data flows through a pipeline of functions, each of which performs a single data transformation task.

Data Flow Thinking is Functional Thinking

Data Flow Pipeline : data flows through a pipeline of functions, each of which performs a single data transformation task.



The Challenge: Removing non-essential text from web pages

Where is the content, where is it not?

```

1 <!doctype html>
2 <html xmlns:og="http://opengraphprotocol.org/schema/" xmlns:fb="http://www.facebook.com/2008/fbml" xmlns:website="http://ogp.me/ns/website" lang="en-US">
3 <head>
4 <meta http-equiv="X-UA-Compatible" content="IE=edge,chrome=1">
5
6 <meta name="viewport" content="width=device-width, initial-scale=1">
7
8 <!-- This is Squarespace. --><!-- artificial-intelligence-blog -->
9 <base href="">
10 <meta charset="utf-8" />
11 <title>2019 Artificial Intelligence News - AI News</title>
12 <link rel="shortcut icon" type="image/x-icon" href="https://static1.squarespace.com/static/585e8f70beba6b65339930b8/t/589cd24eff7c507ac47f8597/favicon.1" />
13 <link rel="canonical" href="https://www.artificial-intelligence.blog/news/">
14 <meta property="og:site_name" content="Artificial Intelligence Blog - AI News"/>
15 <meta property="og:title" content="2019 Artificial Intelligence News - AI News"/>
16 <meta property="og:url" content="https://www.artificial-intelligence.blog/news/">
17 <meta property="og:type" content="website"/>
18 <meta property="og:description" content="Daily fresh news from the field of Artificial Intelligence. There is so much happening in AI ... here's your one-stop-shop for all artificial intelligence news - AI news." />
19 <meta property="og:image" content="http://static1.squarespace.com/static/585e8f70beba6b65339930b8/t/5c61cac6652deal549e1a8bb/1549912776477/artificial-intelligence-news-ai-news-1" />
20 <meta property="og:image:width" content="1500"/>
21 <meta property="og:image:height" content="750"/>
22 <meta itemprop="name" content="2019 Artificial Intelligence News - AI News"/>
23 <meta itemprop="url" content="https://www.artificial-intelligence.blog/news/">
24 <meta itemprop="description" content="Daily fresh news from the field of Artificial Intelligence. There is so much happening in AI ... here's your one-stop-shop for all artificial intelligence news - AI news." />
25 <meta itemprop="thumbnailUrl" content="http://static1.squarespace.com/static/585e8f70beba6b65339930b8/t/5c61cac6652deal549e1a8bb/1549912776477/artificial-intelligence-news-ai-news-1" />
26 <link rel="image_src" href="http://static1.squarespace.com/static/585e8f70beba6b65339930b8/t/5c61cac6652deal549e1a8bb/1549912776477/artificial-intelligence-news-ai-news-1" />
27 <meta itemprop="image" content="http://static1.squarespace.com/static/585e8f70beba6b65339930b8/t/5c61cac6652deal549e1a8bb/1549912776477/artificial-intelligence-news-ai-news-1" />
28 <meta name="twitter:title" content="2019 Artificial Intelligence News - AI News"/>
29 <meta name="twitter:image" content="http://static1.squarespace.com/static/585e8f70beba6b65339930b8/t/5c61cac6652deal549e1a8bb/1549912776477/artificial-intelligence-news-ai-news-1" />
30 <meta name="twitter:url" content="https://www.artificial-intelligence.blog/news/">
31 <meta name="twitter:card" content="summary"/>
32 <meta name="twitter:description" content="Daily fresh news from the field of Artificial Intelligence. There is so much happening in AI ... here's your one-stop-shop for all artificial intelligence news - AI news." />
33 <meta name="description" content="Daily fresh news from the field of Artificial Intelligence. There is so much happening in AI ... here's your one-stop-shop for all artificial intelligence news - AI news." />
34 <script type="text/javascript" src="//use.typekit.net/ik/VcFFX-zDKdf tqeEUEH4H1BSiDrIMLmwPAFAhXMotLCFeChffFHN4UJLFRbh52jhND9tFA9tZRS3jcmKjA2KFR9DwDqowDgI" />
35 <script type="text/javascript">try{Typekit.load()}catch(e){!</script>
36 <link rel="stylesheet" type="text/css" href="//fonts.googleapis.com/css?family=Lato:700,400,900"/>
37 <script type="text/javascript" src="//use.typekit.net/ik/tZgicQd3xyTS-CN-9Wkzturd3191CoV11k JWjIWleqfenGff4e6pUJ6wRMUSQwXFWvu52m8SeJaw48jC8jRjUwDwo5Qv" />
38 <script type="text/javascript">try{Typekit.load()}catch(e){!</script>
39 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
40 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
41 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
42 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
43 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
44 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
45 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
46 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
47 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
48 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
49 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
50 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
51 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
52 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
53 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
54 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
55 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
56 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
57 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
58 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
59 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
60 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
61 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
62 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
63 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
64 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
65 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
66 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
67 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
68 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
69 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
70 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
71 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
72 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
73 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
74 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
75 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
76 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
77 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
78 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
79 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
80 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
81 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
82 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
83 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
84 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
85 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
86 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
87 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
88 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
89 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
90 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
91 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
92 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
93 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
94 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
95 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
96 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
97 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
98 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
99 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>
100 <script type="text/javascript">var x=22;var y=99;print(x,y)</script>

```

Simpler page: Where is the content, where is it not?

```

<html>
<head>
<script> var x = 22; var y = 99; print(x,y) </script>
</head>
<body>
<h1 style="text-align: center;">The Best Artificial Intelligence Blogs</h1>
<h2>OpenAI</h2>
<p>The AI researchers at the non-profit AI research company OpenAI are working hard to help us all understand the power of AI as well as the issues that society must work through on this fascinating topic. An important nuance, they seek to enact the path to safe artificial general intelligence.</p>
<p>For more see: <a href="https://openai.com/" target="_blank" rel="nofollow">OpenAI.com</a></p>
<h2>The a16z AI Playbook</h2>
<p>This is not as frequently updated as a blog but it is such a tremendous resource it belongs high on our list. They bring insights on AI topics with a special focus on the creators who are building AI solutions.</p>
<p>For more see: <a href="http://a16z.com/" target="_blank" rel="nofollow">The a16z AI Playbook</a></p>
<h2>Artificial Intelligence Blog</h2>
<p>They cover AI news, research, books, and thought leaders in the industry. Track for insights into companies and conferences as well.</p>
<p>For more see: <a href="https://www.artificial-intelligence.blog/news/" target="_blank" rel="nofollow">Artificial-Intelligence.blog</a></p>
<h2>Machine Learning Mastery</h2>
<p>Dr. Jason Brownlee is a respected practitioner and master of machine learning and he writes for others seeking to really excel at machine learning. Since his focus is on the people who can really execute the blog gets technical, but it is still understandable for the non-technical person who needs to track the big issues.</p>
...

```

Goal:

Extract relevant text and create a Word Cloud



Beautiful Soup –
a Python Module for
extracting text from web pages

```
1 import bs4
2 from bs4 import BeautifulSoup
3 import numpy as np
4 import pandas as pd
5 from os import path
6
7 from PIL import Image
8 from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
9 import requests
```

Extract text from simple web page about Rollo using BS

```

1 html = '''
2 <html>
3 <head>
4
5 <h1>Rollo WORLD</h1>
6 <p>Rollo ate lunch and then went to the pool to study AI.</p>
7 He knew AI was the way to take over the world.
8 Rollo dreamed of AI as the solution to all his problems.
9 <p>Perhaps the world will someday
10 know about Rollo and his ambitions. Yet now, the world does not know the Rollo.
11 </p>
12 </html>
13 '''
14
15 soup = BeautifulSoup(html, "lxml")
16 # we don't want to strip line breaks
17 htmlwords = soup.get_text(strip=True)
18 print (htmlwords)

```

Rollo WORLD Rollo ate lunch and then went to the pool to study AI. He knew AI was the way to take over the world.
 Rollo dreamed of AI as the solution to all his problems. Perhaps the world will someday
 know about Rollo and his ambitions. Yet now, the world does not know the Rollo.

when using strip="True", we remove extra space
 and linebreaks - but we may merge words ☹️

```

1 html = '''
2 <html>
3 <head>
4
5 <h1>Rollo WORLD</h1>
6 <p>Rollo ate lunch and then went to the pool to study AI.</p>
7 He knew AI was the way to take over the world.
8 Rollo dreamed of AI as the solution to all his problems.
9 <p>Perhaps the world will someday
10 know about Rollo and his ambitions. Yet now, the world does not know the Rollo.
11 </p>
12 </html>
13 '''
14
15 soup = BeautifulSoup(html, "lxml")
16 # we don't want to strip line breaks
17 htmlwords = soup.get_text(strip=False)
18 print (htmlwords)

```

Rollo WORLD
 Rollo ate lunch and then went to the pool to study AI.
 He knew AI was the way to take over the world.
 Rollo dreamed of AI as the solution to all his problems.
 Perhaps the world will someday
 know about Rollo and his ambitions. Yet now, the world does not know the Rollo.

↑
 a string

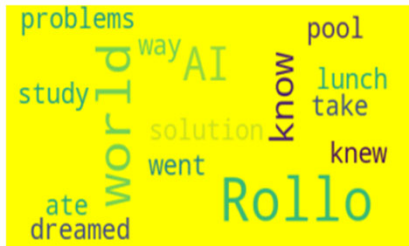
We don't want to strip extra spaces and
 linebreaks. We get the single words!

WordCloud – from a string of words

```

1 wordcloud = WordCloud(max_font_size=50, max_words=15, background_color="yellow").generate(htmlwords)
2
3 # Display the generated image:
4 plt.imshow(wordcloud, interpolation='bilinear')
5 plt.axis("off")
6 plt.show()

```



OK BUT what if our web page is more complex...

```

4 html = '''
5 <html>
6 <head>
7 <script> print("hellow world")</script>
8 </script>
9
10 (function(){
11   print ("hello"); print ("bye"); print ("now");
12   print ("later"); print ("doc"); print ("yes");
13   print ("doctext"); print ("doctext"); print ("doctext");
14 </script>
15 </head>
16 <h1>Rollo WORLD</h1>
17 <p class="fool">pret1</p>
18 <p class="fool">date is now</p>
19 <p class="fool">zipfest</p>
20 <p>code</p>
21 <p>data</p>
22 <p class="fool">aT DAWN Rollo <strong class='xyz'>ate
23 </strong>
24 that the lunch left for him. Oh that Rollo. He wants to take lunch and lunch over the that that world.
25 The world does not know.
26 </p>
27 <html>
28 '''
29
30 soup = BeautifulSoup(html, "lxml")
31 # we don't want to strip line breaks
32 htmlwords = soup.get_text(strip=False)
33 print (htmlwords)

```

We get irrelevant words from the <script> element



We get the TEXT inside every HTML element – even the script tag!

```

print("hellow world")

(function(){
  print ("hello"); print ("bye"); print ("now");
  print ("later"); print ("doc"); print ("yes");
  print ("doctext"); print ("doctext"); print ("doctext");

Rollo WORLD
pret1
date is now
zipfest
code
data
aT DAWN Rollo ate

that the lunch left for him. Oh that Rollo. He wants to take lunch and lunch
over the that that world.
The world does not know.

```

We need the power of regex

Quick Regex Overview....

regular expressions with the Python re module

Syntax

```
import re

result = re.sub(pattern, repl, string, count=0, flags=0);
```

Simple Examples

```
result = re.sub('abc', '', input)           # Delete pattern abc
result = re.sub('abc', 'def', input)         # Replace pattern abc -> def
result = re.sub(r'\s+', ' ', input)          # Eliminate duplicate whitespaces
result = re.sub('abc(def)ghi', r'\1', input) # Replace a string with a part of itself
```

Note: Take care to always prefix patterns containing \ escapes with raw strings (by adding an r in front of the string). Otherwise the \ is used as an escape sequence and the regex won't work.

Experiment with REGEXes at <https://www.regexpal.com/>

Basics

*	Match preceding character 0 or more times
+	Match preceding character 1 or more times
.	Match any single character
x y	Match either 'x' or 'y'
\	Escape a special character
b	The character b
abc	The string abc

Character Classes

[abc]	Match any one of the characters in the set 'abc'
[^abc]	Match anything not in character set 'abc'

\d	Match a digit character
\D	Match a non-digit character
\s	Match a single white space character (space, tab, form feed, or line feed)
\S	Match a single character other than white space
\w	Match any alphanumeric character (including underscore)
\W	Match any non-word character

How to use regex to eliminate a tag (e.g. <script>) and all its content

The secret sauce

the '?' turns OFF greedy matching.
note that the <script> element is gone
in the text!

```
1 # remove the script element and its content
2 blogtext = re.sub(r'<script>.*?</script>', ' ', blogtext)
3 blogtext
```

```
'<html>\n<head>\n \n</head>\n<body>\n<h1 style="text-align: center;">The Best Artificial Intelligence Blogs</h1>\n<h2>OpenAI</h2>\n<p>The AI researchers at the non-profit AI research company OpenAI are working hard to help us all understand the power of AI as well as the issues that society must work through on this fascinating topic. An important nuance, they seek to enact the path to safe artificial general intelligence.</p>\n<p>For more see: <a href="https://openai.com/" target="_blank" rel="nofollow">OpenAI.com</a></p>\n<h2>The a16z AI Playbook</h2>\n<p>This is not as frequently updated as a blog but it is such a tremendous resource it belongs high on our list. They bring insights on AI topics with a special focus on the creators who are building AI solutions.</p>\n<p>For more see: <a href="http://a16z.com/" target="_blank" rel="nofollow">a16z.com</a></p>\n</body>\n</html>'
```

Challenge: Create WordCloud from a web page

The Best Artificial Intelligence Blogs

OpenAI

The AI researchers at the non-profit AI research company OpenAI are working hard to help us all understand the power of AI as well as the issues that society must work through on this fascinating topic. An important nuance, they seek to enact the path to safe artificial general intelligence.

For more see: [OpenAI.com](https://openai.com/)

The a16z AI Playbook

This is not as frequently updated as a blog but it is such a tremendous resource it belongs high on our list. They bring insights on AI topics with a special focus on the creators who are building AI solutions.

For more see: [The a16z AI Playbook](http://aiplaybook.a16z.com/)

Artificial Intelligence Blog

They cover AI news, research, books, and thought leaders in the industry. Track for insights into companies and conferences as well.

For more see: [Artificial-Intelligence.blog](https://www.artificial-intelligence.blog/news/)

Machine Learning Mastery

Dr. Jason Browlee is a respected practitioner and master of machine learning and he writes for others seeking to really excel at machine learning. Since his focus is on the people who can really execute the blog gets technical, but it is still understandable for the non-technical person who needs to track the big issues.

For more see: [Machine Learning Mastery](http://machinelearningmastery.com/blog/)



```
<html>
<head>
<script> var x = 22; var y = 99; print(x,y) </script>
</head>
<body>
<h1 style="text-align: center;">The Best Artificial Intelligence Blogs</h1>
<h2>OpenAI</h2>
<p>The AI researchers at the non-profit AI research company OpenAI are working hard to help us all understand the power of
AI as well as the issues that society must work through on this fascinating topic. An important nuance, they seek to enact the
path to safe artificial general intelligence.</p>
<p>For more see: <a href="https://openai.com/" target="_blank" rel="nofollow">OpenAI.com</a></p>
<h2>The a16z AI Playbook</h2>
<p>This is not as frequently updated as a blog but it is such a tremendous resource it belongs high on our list. They bring
insights on AI topics with a special focus on the creators who are building AI solutions.</p>
<p>For more see: <a href="http://aiplaybook.a16z.com/" target="_blank" rel="nofollow">The a16z AI Playbook</a></p>
<h2>Artificial Intelligence Blog</h2>
<p>They cover AI news, research, books, and thought leaders in the industry. Track for insights into companies and
conferences as well.</p>
<p>For more see: <a href="https://www.artificial-intelligence.blog/news/" target="_blank" rel="nofollow">Artificial-
Intelligence.blog</a></p>
<h2>Machine Learning Mastery</h2>
<p>Dr. Jason Browlee is a respected practitioner and master of machine learning and he writes for others seeking to really
excel at machine learning. Since his focus is on the people who can really execute the blog gets technical, but it is still
understandable for the non-technical person who needs to track the big issues.</p>
<p>For more see: <a href="http://machinelearningmastery.com/blog/" rel="nofollow" target="_blank">Machine Learning
Mastery</a></p>
<h2>The Algorithmia Blog</h2>
...
```

```

1 # read file and save content as text string
2 file = open("blogpage.html", "r")
3 blogtext = file.read()
4 blogtext

```

```

'<html>\n<head>\n<script> var x = 22; var y = 99; print(x,y) </script>\n</head>\n<body>\n<h1 style="text-align: center;">The
e Best Artificial Intelligence Blogs</h1>\n<h2>OpenAI</h2>\n<p>The AI researchers at the non-profit AI research company Ope
nAI are working hard to help us all understand the power of AI as well as the issues that society must work through on this
fascinating topic. An important nuance, they seek to enact the path to safe artificial general intelligence.</p>\n<p>For mo
re see: <a href="https://openai.com/" target="_blank" rel="nofollow">OpenAI.com</a></p>\n<h2>The al6z AI Playbook</h2>\n<p>
This is not as frequently updated as a blog but it is such a tremendous resource it belongs high on our list. They bring in
sights on AI topics with a special focus on the creators who are building AI solutions.</p>\n<p>For more see: <a href="http
://aiplaybook.al6z.com/" target="_blank" rel="nofollow">The al6z AI Playbook</a></p>\n<h2>Artificial Intelligence Blog</h2>
\n<p>They cover AI news, research, books, and thought leaders in the industry. Track for insights into companies and confer
ences as well.</p>\n<p>For more see: <a href="https://www.artificial-intelligence.blog/news/" target="_blank" rel="nofollow

```

Tasks:

- 1) remove tags and content not relevant to the page (e.g. script and others?)
- 2) remove \n
- 3) remove tags
- 4) build Word Cloud for relevant text