FULL TEXT
SEARCH

박 경 규

# Introduction

## Building a full-text search engine in 150 lines of Python code

Mar 24, 2021

`how-to`  `search`  `full-text search`  `python`

Full-text search is everywhere. From finding a book on Scribd, a movie on Netflix, toilet paper on Amazon, or anything else on the web through Google (like how to do your job as a software engineer), you've searched vast amounts of unstructured data multiple times today. What's even more amazing, is that you've even though you searched millions (or billions) of records, you got a response in milliseconds. In this post, we are going to explore the basic components of a full-text search engine, and use them to build one that can search across millions of documents and rank them according to their relevance in milliseconds, in less than 150 lines of Python code!

# Data

## ■ 소스 다운로드

**https://github.com/bartdegoede/python-searchengine/**

# Data

## ■ 파이썬 패키지 설치

pip install -r requirements.txt

## ■ 프로그램 실행

python run.py

**download.py**

```python
import requests


def download_wikipedia_abstracts():
    URL = 'https://dumps.wikimedia.org/enwiki/latest/enwiki-latest-abstract.xml.gz'
    with requests.get(URL, stream=True) as r:
        r.raise_for_status()
        with open('data/enwiki-latest-abstract.xml.gz', 'wb') as f:
            # write every 1mb
            for i, chunk in enumerate(r.iter_content(chunk_size=1024*1024)):
                f.write(chunk)
                if i % 10 == 0:
                    print(f'Downloaded {i} megabytes', end='\r')
```

# Data preparation

```
<doc>
    <title>Wikipedia: London Beer Flood</title>
    <url>https://en.wikipedia.org/wiki/London_Beer_Flood</url>
    <abstract>The London Beer Flood was an accident at Meux & Co's Horse Shoe Brewery, London, on 17 October 1814. It took place when one of the  wooden vats of fermenting porter burst.</abstract>
    ...
</doc>
```

https://bart.degoe.de/building-a-full-text-search-engine-150-lines-of-code/

5

# Data preparation

**documents.py**

```python
from collections import Counter
from dataclasses import dataclass

from .analysis import analyze

@dataclass
class Abstract:
    """Wikipedia abstract"""
    ID: int
    title: str
    abstract: str
    url: str

    @property
    def fulltext(self):
        return ' '.join([self.title, self.abstract])

    def analyze(self):
        self.term_frequencies = Counter(analyze(self.fulltext))

    def term_frequency(self, term):
        return self.term_frequencies.get(term, 0)
```

# Data preparation

**load.py**

```python
import gzip
from lxml import etree
import time

from search.documents import Abstract

def load_documents():
    start = time.time()
    with gzip.open('data/enwiki-latest-abstract.xml.gz', 'rb') as f:
        doc_id = 0
        for _, element in etree.iterparse(f, events=('end',), tag='doc'):
            title = element.findtext('./title')
            url = element.findtext('./url')
            abstract = element.findtext('./abstract')

            yield Abstract(ID=doc_id, title=title, url=url, abstract=abstract)

            doc_id += 1
            element.clear()
    end = time.time()
    print(f'Parsing XML took {end - start} seconds')
```
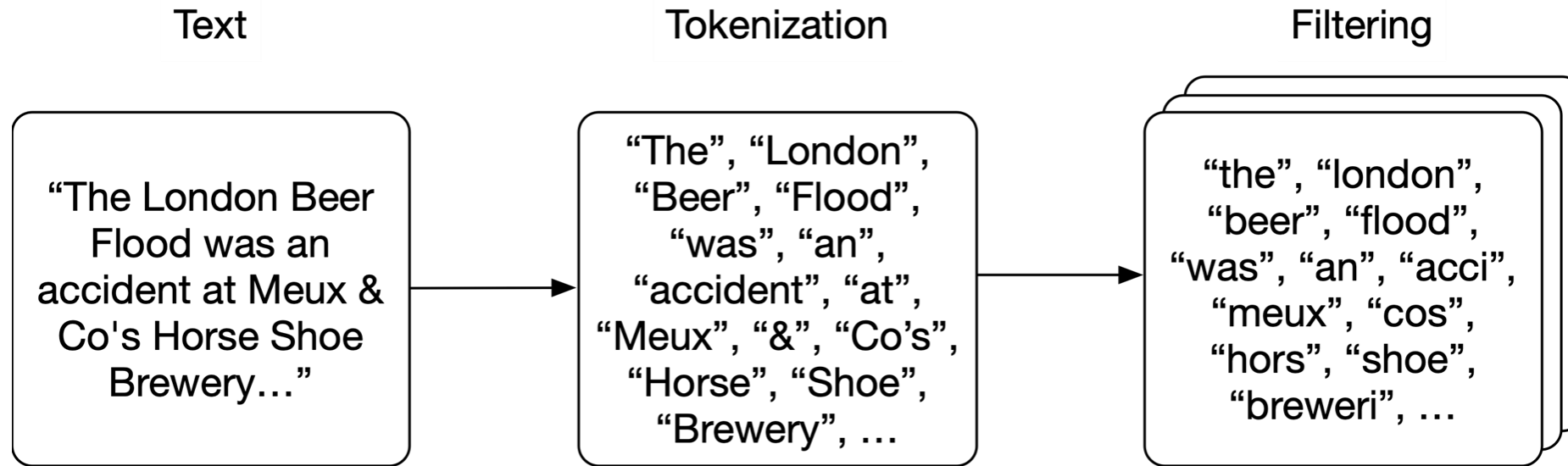
7

# Indexing

https://bart.degoe.de/building-a-full-text-search-engine-150-lines-of-code/

# Indexing

| Text | Tokenization | Filtering |
|---|---|---|

"The London Beer Flood was an accident at Meux & Co's Horse Shoe Brewery…"

→

"The", "London", "Beer", "Flood", "was", "an", "accident", "at", "Meux", "&", "Co's", "Horse", "Shoe", "Brewery", …

→

"the", "london", "beer", "flood", "was", "an", "acci", "meux", "cos", "hors", "shoe", "breweri", …

```
{
    ...
    "london": [5245250, 2623812, 133455, 3672401, ...],
    "beer": [1921376, 4411744, 684389, 2019685, ...],
    "flood": [3772355, 2895814, 3461065, 5132238, ...],
    ...
}
```

https://bart.degoe.de/building-a-full-text-search-engine-150-lines-of-code/

9

# Indexing

**analysis.py**

```python
import re
import string
import Stemmer

# top 25 most common words in English and "wikipedia":
# https://en.wikipedia.org/wiki/Most_common_words_in_English
STOPWORDS = set(['the', 'be', 'to', 'of', 'and', 'a', 'in', 'that', 'have',
                 'i', 'it', 'for', 'not', 'on', 'with', 'he', 'as', 'you',
                 'do', 'at', 'this', 'but', 'his', 'by', 'from', 'wikipedia'])
PUNCTUATION = re.compile('[%s]' % re.escape(string.punctuation))
STEMMER = Stemmer.Stemmer('english')

def tokenize(text):
    return text.split()

def lowercase_filter(tokens):
    return [token.lower() for token in tokens]
```

# Indexing

analysis.py

```python
19    def punctuation_filter(tokens):
20        return [PUNCTUATION.sub('', token) for token in tokens]
21
22    def stopword_filter(tokens):
23        return [token for token in tokens if token not in STOPWORDS]
24
25    def stem_filter(tokens):
26        return STEMMER.stemWords(tokens)
27
28    def analyze(text):
29        tokens = tokenize(text)
30        tokens = lowercase_filter(tokens)
31        tokens = punctuation_filter(tokens)
32        tokens = stopword_filter(tokens)
33        tokens = stem_filter(tokens)
34
35        return [token for token in tokens if token]
```

# Indexing the corpus

index.py

```python
import math

from .timing import timing
from .analysis import analyze

class Index:
    def __init__(self):
        self.index = {}
        self.documents = {}

    def index_document(self, document):
        if document.ID not in self.documents:
            self.documents[document.ID] = document
            document.analyze()

        for token in analyze(document.fulltext):
            if token not in self.index:
                self.index[token] = set()
            self.index[token].add(document.ID)
```

# Indexing the corpus

**index.py**

```python
import math

from .timing import timing
from .analysis import analyze

class Index:
    def __init__(self):
        self.index = {}
        self.documents = {}

    def index_document(self, document):
        if document.ID not in self.documents:
            self.documents[document.ID] = document
            document.analyze()

        for token in analyze(document.fulltext):
            if token not in self.index:
                self.index[token] = set()
            self.index[token].add(document.ID)
```

# Searching

```python
33    @timing
34    def search(self, query, search_type='AND', rank=False):
35        """
36        Search; this will return documents that contain words from the query,
37        and rank them if requested (sets are fast, but unordered).
38
39        Parameters:
40          - query: the query string
41          - search_type: ('AND', 'OR') do all query terms have to match, or just one
42          - score: (True, False) if True, rank results based on TF-IDF score
43        """
44        if search_type not in ('AND', 'OR'):
45            return []
46
47        analyzed_query = analyze(query)
48        results = self._results(analyzed_query)
49        if search_type == 'AND':
50            # all tokens must be in the document
51            documents = [self.documents[doc_id] for doc_id in set.intersection(*results)]
52        if search_type == 'OR':
53            # only one token has to be in the document
54            documents = [self.documents[doc_id] for doc_id in set.union(*results)]
55
56        if rank:
57            return self.rank(analyzed_query, documents)
58        return documents
```

14

# Run

**python 실행**

**run.py 복사 실행**

```
선택 명령 프롬프트 - python                                    —   □   ×

>>>
>>> sr = index.search('Red Flags', search_type='AND', rank=True)
search took 0.0 milliseconds
>>>
>>> sr[0]
(Abstract(ID=87261, title='Wikipedia: Van Dorn battle flag', abstract='The Van
Dorn battle flag is a historical Confederate flag with a red field depicting a
white crescent moon in the canton and thirteen white stars; and trimmed with go
ld cord. In February, 1862, Confederate general Earl Van Dorn ordered that all
units under his command use this flag as their regimental colors.', url='https:
//en.wikipedia.org/wiki/Van_Dorn_battle_flag'), 15.828240750671078)
>>>
>>> sr[1]
(Abstract(ID=141116, title='Wikipedia: Flag of Krasnoyarsk Krai', abstract="The
 flag of Krasnoyarsk Krai, in the Russian Federation, is a red field charged wi
th the krai's coat of arms in the center.  Two fifths of the flag's height, it
displays a golden lion holding a sickle in its left hand and a shovel in its ri
ght hand.", url='https://en.wikipedia.org/wiki/Flag_of_Krasnoyarsk_Krai'), 12.5
06210320685472)
>>>
```

```python
1   import os.path
2   import requests
3
4   from download import download_wikipedia_abstracts
5   from load import load_documents
6   from search.timing import timing
7   from search.index import Index
8
9
10  @timing
11  def index_documents(documents, index):
12      for i, document in enumerate(documents):
13          index.index_document(document)
14          if i % 5000 == 0:
15              print(f'Indexed {i} documents', end='\r')
16      return index
17
```

15

# Thank you

kgpark88@gmail.com