

Popular topics analysis

Stone Fang (Student ID: 19049045)
Computers and Information Sciences
Auckland University of Technology
Auckland, New Zealand
fnk7060@autuni.ac.nz

I. METRICS

xx

APPENDIX

```
1 #!/usr/bin/env python
2 # -*- coding: utf-8 -*-
3
4 import sys
5 import os.path
6 from glob import glob
7 from tqdm import tqdm
8 import pickle
9
10 from collections import namedtuple
11 import pandas as pd
12 from xml.etree import ElementTree
13 from xml.etree.ElementTree import ParseError
14 from bs4 import BeautifulSoup
15 from datetime import datetime
16
17 #####
18 # Codes for data reading &
19 #####
20
21 Record = namedtuple('Record', ['meta', 'posts'])
22 Post = namedtuple('Post', ['date', 'text'])
23 MetaData = namedtuple('MetaData', ['id', 'gender', 'age', 'category', 'zodiac'])
24
25 def parse_meta_data(meta_data_str):
26     arr = meta_data_str.strip().split('.')
27     return MetaData(arr[0], arr[1], int(arr[2]), arr[3], arr[4])
28
29 # _parser = ElementTree.XMLParser(encoding="utf-8")
30 def read_blog_file(fpath):
31     try:
32         # tree = ElementTree.parse(fpath, parser=_parser)
33         with open(fpath, encoding='utf-8', errors='ignore') as f:
34             soup = BeautifulSoup(f.read(), "xml")
35             blog = soup.Blog
36     except ParseError:
37         print('Error: invalid xml file {}'.format(fpath))
38         raise
39     return []
40
41 posts = []
42 state = 'date'
43 for c in blog.find_all(recursive=False):
44     # print(c)
45     # print(c.text)
46     # check the <date> and <post> tags appear alternately
47     if c.name != state:
48         print('Warning: inconsistent format in file {}'.format(fpath))
49     if state == 'date':
50         try:
51             date_str = c.text.strip()
52             # date_str = date_str.replace('janvier', 'january') \
53             #     .replace('mars', 'march') \
54             #     .replace('avril', 'april') \
55             #     .replace('mai', 'may') \
56             #     .replace('juin', 'june') \
57             #     .replace('juillet', 'july')
58             # date = pd.to_datetime(date_str, format='%d,%B,%Y')
59             date = date_str
60         except ValueError:
61             print('Warning: invalid date {} in file {}'.format(c.text, fpath))
62             date = datetime.fromtimestamp(0)
63         state = 'post'
64     else:
65
```

```

66         text = c.text.strip()
67         state = 'date'
68         posts.append(Post(date, text))
69         # print(c, c.text)
70     posts.sort(key=lambda p: p.date)
71     # print(posts)
72     # print(pd.DataFrame(posts))
73     # sys.exit()
74     return posts
75
76 def read_blogs(path, force=False, cache_file='blogs.pkl'):
77     if not force and os.path.exists(cache_file):
78         print('load dataset from cached pickle file ' + cache_file)
79         with open(cache_file, 'rb') as f:
80             dataset = pickle.load(f)
81         return dataset
82
83     dataset = read_blogs_xml(path)
84
85     # save to pickle file for fast loading next time
86     with open(cache_file, 'wb') as f:
87         print('save dataset to pickle file ' + cache_file)
88         pickle.dump(dataset, f)
89
90     return dataset
91
92 def read_blogs_xml(path):
93     print('reading all data files from directory {} ...'.format(path))
94     dataset = []
95     for fpath in tqdm(glob(os.path.join(path, '*'))):
96         # for fpath in list(glob(os.path.join(path, '*'))[:10]):
97             # print(fpath)
98             fname = os.path.basename(fpath)
99             meta_data = parse_meta_data(fname)
100             # print(meta_data)
101             posts = read_blog_file(fpath)
102             rec = Record(meta_data, posts)
103             dataset.append(rec)
104     return dataset
105
106 def show_summary(dataset):
107     df = pd.DataFrame([d.meta for d in dataset])
108     df['blog_count'] = [len(d.posts) for d in dataset]
109     # print(df)
110     print(df.describe(include='all'))
111     print('{} possible values for "gender": {}'.format(
112         len(df.gender.unique()), ', '.join(sorted(df.gender.unique()))))
113     # print('{} possible values for "{}": {}'.format(
114     #     len(df.age.unique()), ', '.join(sorted(df.age.unique()))))
115     print('{} possible values for category: {}'.format(
116         len(df.category.unique()), ', '.join(sorted(df.category.unique()))))
117     print('{} possible values for zodiac: {}'.format(
118         len(df.zodiac.unique()), ', '.join(sorted(df.zodiac.unique()))))
119
120 def main():
121     read_blogs('blogs')
122     return
123
124 if __name__ == '__main__':
125     main()

```