

Using topic modeling to add similar posts links to a Pelican powered blog

Similar Posts

- [Using topic modeling to find related blog posts](#), Score: 0.989
- [Analysis of Shakespeare characters speech topics](#), Score: 0.907
- [Using sed to make specific text lowercase in place](#), Score: 0.593
- [When joins go wrong, check data types](#), Score: 0.478
- [Filter common words from documents](#), Score: 0.400

Frank Cleary

www.datasciencebytes.com (hobby)

Member Technical Staff at Index (day job)

Pelican – site generation in Python

IP[y]: Notebook bart-annotate Last Checkpoint: Feb 10 07:00 (autosaved)

File Edit View Insert Cell Kernel Help

Markdown Cell Toolbar: None

To extend on my [post about plotting and reshaping data](#) from the [BART API](#), I worked a bit with the matplotlib annotation interface to add text and arrows to a plot. The meat of this post is in cell #4 below. [Download notebook](#).

(github.com/danielfrg/pelican-ipynb)

```
In [2]: %matplotlib inline
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import datetime

def prettify_axis(ax, ylabel='', xlabel=
label_format_dict = dict(fontsize=20
tick_format_dict = dict(labelsize=16
length=4, wi
ax.set_xlabel(xlabel, label_format_d:
ax.set_ylabel(ylabel, label_format_d:
ax.tick_params(**tick_format_dict)
```

The file [mill_pivot.csv](#) contains the estimated row is a minute of the day and each column is a s

```
In [3]: def time_parser(time_string, string_form:
return datetime.datetime.strptime(t:
mill_pivot = pd.read_csv('data/mill_pivo:
converters={'tir
index_col=0)
mill_pivot.ix[:, :4]
```

```
Out[3]:
```

	2014-11-28	2014-12-01	2014-1-
time_of_day			
03:29:00	52	52	52
03:30:00	51	51	51
03:31:00	50	50	50

Data Science Bytes

News Tips Tutorials Recommended Books Recommended Videos About

Annotating matplotlib plots

Tweet 0

Published: Mon 29 December 2014

By [Frank Cleary](#)

In [Tips](#).

tags: [pandas](#) [matplotlib](#) [data](#) [python](#)

To extend on my [post about plotting and reshaping data](#) from the [BART API](#), I worked a bit with the matplotlib annotation interface to add text and arrows to a plot. The meat of this post is in cell #4 below. [Download notebook](#).

In [2]:

```
%matplotlib inline
import matplotlib.pyplot as plt
import pandas as pd
```

Adding similar posts

Data Science Bytes

News Tips Tutorials Recommended Books Recommended Videos About

Spark 1.2.0 released

 Tweet 0

Published: Fri 19 December 2014

By [Frank Cleary](#)

In [News](#).

tags: [Spark](#)

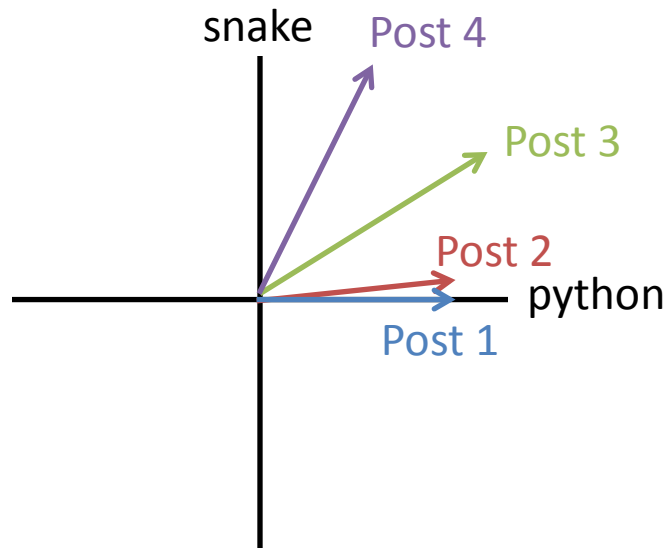
Spark 1.2.0 was released yesterday ([release notes](#)). I'm curious to see how the new machine learning API's in spark.ml evolve.

Similar Posts

- [Scikit-learn machine learning algorithm flowchart](#), Score: 0.967
- [Installing python for data science](#), Score: 0.780
- [Using sed to make specific text lowercase in place](#), Score: 0.719
- [How to Transition from Ph.D. Student to Data Scientist](#), Score: 0.535
- [When joins go wrong, check data types](#), Score: 0.487

Finding similar posts

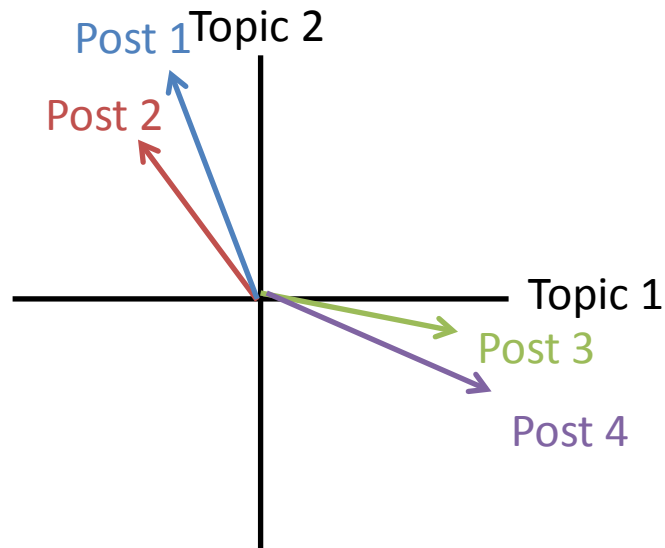
	Word Count			
Word	Post 1	Post 2	Post 3	Post 4
python	10	5	12	6
snake	0	1	5	10
iterator	5	4	0	0
data	8	2	0	1
...



- Here I've already selected words that differentiate topics. Topic models can do this automatically.

Topic model

	Contribution from Topic			
Topic	Post 1	Post 2	Post 3	Post 4
python, snake, scale	-.2	-.3	.5	.6
iterator, data, import	.8	.5	-.2	-.4
sweet, awesome, bug	.2	.2	.2	.2
...



Gensim – topic modeling for humans

Here are two topics of Shakespeare character’s speech.

www.datasciencebytes.com/bytes/2014/12/31/analysis-of-shakespeare-character-speech-topics/

x topic:

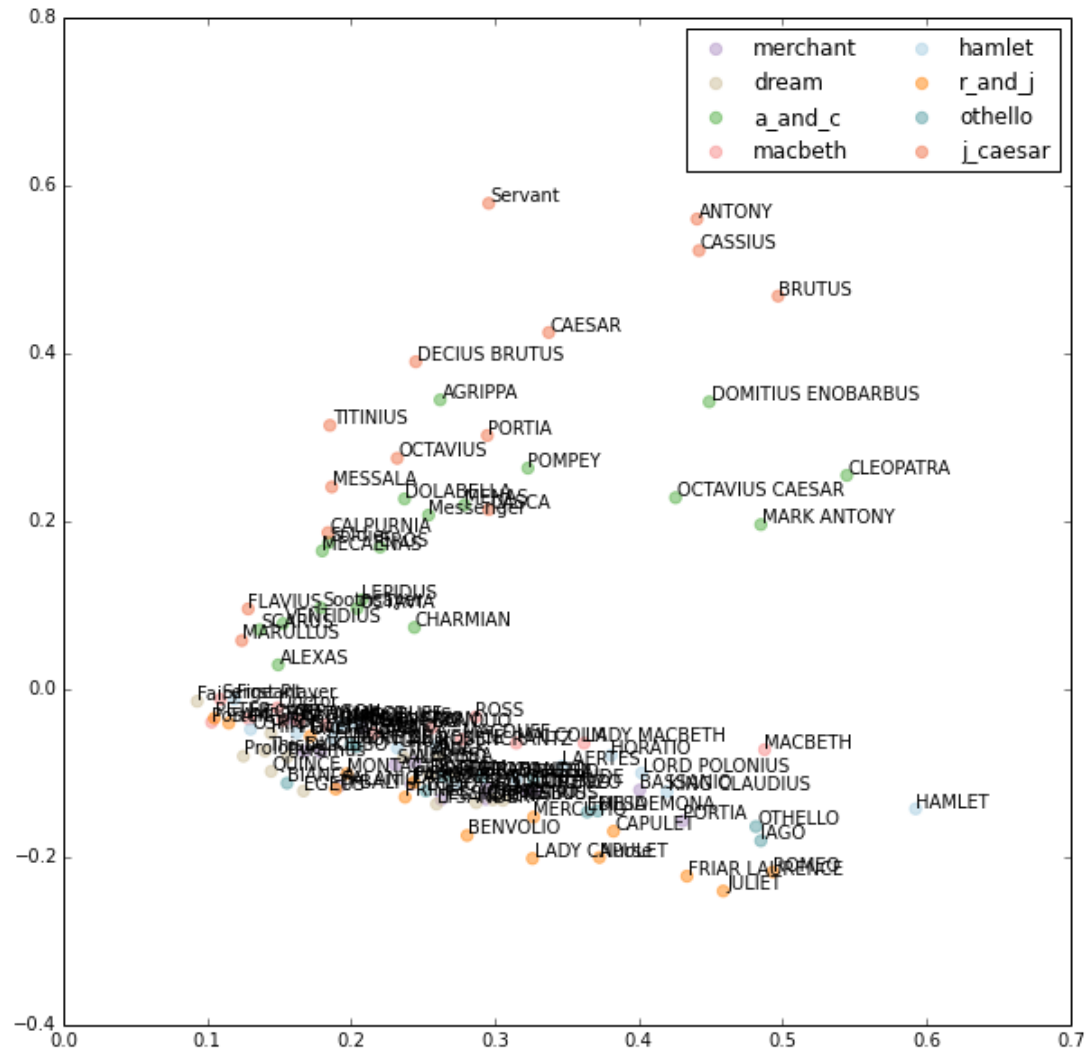
```
( '0.192', u'caesar' ),
( '0.125', u'lord' ),
( '0.121', u'antony' ),
( '0.112', u'brutus' ),
( '0.106', u'thou' ),
( '0.105', u'romeo' ),
( '0.093', u'cassio' ),
( '0.091', u'love' ),
( '0.084', u'thee' ),
( '0.078', u'madam' ),
...

```

y topic:

```
( '0.513', u'caesar' ),
( '0.378', u'brutus' ),
( '0.286', u'antony' ),
( '0.192', u'cassius' ),
( '-0.151', u'romeo' ),
( '0.139', u'rome' ),
( '-0.108', u'cassio' ),
( '0.090', u'octavius' ),
( '0.081', u'lepidus' ),
( '-0.073', u'tybalt' )
...

```



Implementation

- Hook into Pelican at site-generation time:

```
def register():  
    """Entry point for ArticleGenerator from pelican"""  
    signals.article_generator_finalized.connect(add_related_posts)
```

- Process articles:

```
def add_related_posts(generator, default_max_related_posts=5):  
    max_posts = generator.settings.get("MAX_RELATED_POSTS",  
                                       default_max_related_posts)  
    similarity_scores = recommend_articles(generator.articles, max_posts)  
    articles_by_path = {art.source_path: art for art in generator.articles}  
    for article in generator.articles:  
        related_posts = similarity_scores[article.source_path]  
        article.related_posts = []  
        article.score = {}  
        for source_path, similarity in related_posts:  
            related_post = articles_by_path[source_path]  
            article.score[related_post.source_path] = similarity  
            article.related_posts.append(related_post)
```


Implementation

- Add to article template:

```
{% if article.related_posts %}
<h1>Related Posts</h1>
<ul>
  {% for related_post in article.related_posts %}
    <li><a href="{{ SITEURL }}/{{ related_post.url }}">
      {{ related_post.title }}
    </a>,
    Score: {{ '%0.3f' % article.score.get(related_post.source_path) }}</li>
  {% endfor %}
</ul>
{% endif %}
```

Thanks!

- Pelican: docs.getpelican.com
- Gensim: radimrehurek.com/gensim/
- IPython notebook plugin: github.com/danielfrg/pelican-ipy nb
- Pelican tag-based related posts: github.com/getpelican/pelican-plugins/tree/master/related_posts
- SF Python Meetup

- See More: www.datasciencebytes.com (site source: github.com/frankcleary/data-science-bytes)
- Careers at Index: index.com/careers