

Clustering-Based Sentiment Analysis for Media Agenda Setting

Opinion Lab Group 2.3

Wing Sheung Leung, Qiaoxi Liu

May 3, 2020



TUM Uhrenturm

Outline

- Aims
- Dataset
- Expected outputs
- Todo
 - Stage 1
 - Generate word / token / sentence embeddings with our corpus
 - Build a k-mean clustering model for identifying sub-topics in organic dataset
 - Select a sentiment pre-trained model
 - Stage 2
 - Visualize distribution of samples among clusters and corresponding sentiment frequency
 - Stage 3
 - Investigate Media Agenda Setting
- Milestones

Aims

Measure the influence of two online newspapers onto the social media according to Agenda setting theory

Agenda setting theory:

Suggest the news item which is covered more frequently and prominently indicates that the audience will regard the issue as more important.

Dataset

Articles and respective comments on the domain of organic food with search terms *organic food* and *organic farming*

- Articles from two online newspapers, *New York Times* (English) and *Der Spiegel* (German)
- Direct response (bilingual): comments right under those articles
- Indirect response: posts in unrelated discussion forums, *Quora*

	Start	End	No. of articles
<i>New York Times</i>			
With comments	2006	2017	99
Without comments	1970	2017	228
<i>Der Spiegel</i>			
With comments	2007	2017	61
Without comments	2007	2017	91
<i>Quora</i>			
With comments	2009	2017	1304
Without comments	2010	2017	193

Table: Statistics for data crawled from New York Times, Der Spiegel and Quora with 'relevant' labelled as 1.0

Expected outputs (Processes)

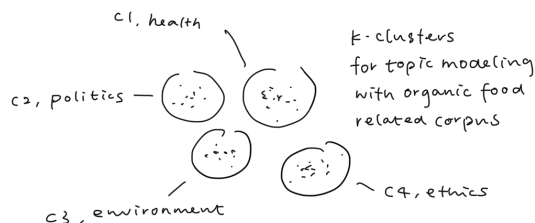


Figure: Stage 1: Cluster, tuple(aspect,sentiment)

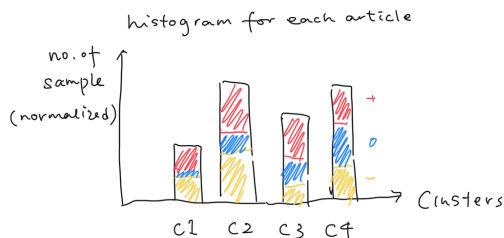


Figure: Stage 2: Distribution of clusters per articles

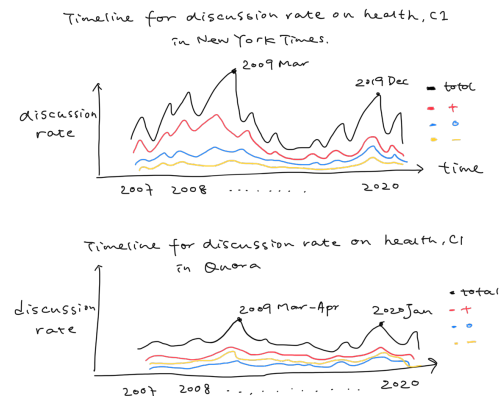


Figure: Stage 3: Timeline analysis, influence on social media

Todo Stage 1.1

Generate embeddings

text → *lines split* → list of sentences → *sentence embedding* → **vectors**

¹Chidambaram et al. 2018.

²dalequark 2019.

Todo Stage 1.1

Generate embeddings

text \rightarrow *lines split* \rightarrow list of sentences \rightarrow *sentence embedding* \rightarrow **vectors**

Example: *sentence* = "You are what you eat."

¹Chidambaram et al. 2018.

²dalequark 2019.

Todo Stage 1.1

Generate embeddings

text → *lines split* → list of sentences → *sentence embedding* → **vectors**

Example: *sentence* = "You are what you eat."

Possible embedding models (on TensorFlow Hub) to apply:

- Universal sentence encoder¹
- Bert²

¹Chidambaram et al. 2018.

²dalequark 2019.

Todo Stage 1.1

Generate embeddings

text → *lines split* → list of sentences → *sentence embedding* → **vectors**

Example: *sentence* = "You are what you eat."

Possible embedding models (on TensorFlow Hub) to apply:

- Universal sentence encoder¹
- Bert²

```
embed = hub.Module("https://tfhub.dev/google/universal-sentence-encoder-xling/en-de/1")
with tf.Session() as session:
    session.run()
    print(session.run(embed(sentence)))
```

→ [0.02, 0.01, ..., -0.12]

¹Chidambaram et al. 2018.

²dalequark 2019.

Todo Stage 1.2

Build a k-mean clustering model for identifying sub-topics

Now, let's start working with **vectors**!

→ k-means cluster, optimize k with human observation (samples: sentences)

Todo Stage 1.2

Build a k-mean clustering model for identifying sub-topics

Now, let's start working with **vectors**!

→ k-means cluster, optimize k with human observation (samples: sentences)

```
s_emb_matrix = [s0, s2, ..., s9] # s are embedding vectors of sentences
nclusters= 3
km = KMeans(nclusters)
km.fit(s_emb_matrix)
clusters = {} # key: label, values: index of sentence
for i, label in enumerate(km.labels_):
    clusters[label].append(i) # append the index to the corresponding label
```

Todo Stage 1.2

Build a k-mean clustering model for identifying sub-topics

Now, let's start working with **vectors**!

→ k-means cluster, optimize k with human observation (samples: sentences)

```
s_emb_matrix = [s0, s2, ..., s9] # s are embedding vectors of sentences
nclusters= 3
km = KMeans(nclusters)
km.fit(s_emb_matrix)
clusters = {} # key: label, values: index of sentence
for i, label in enumerate(km.labels_):
    clusters[label].append(i) # append the index to the corresponding label
```

→ get *clusters*[0] : s_0, s_4, s_5 with top-n words list *topWords*[0] : $w_0, w_1, \dots w_n$

Todo Stage 1.2

Build a k-mean clustering model for identifying sub-topics

Now, let's start working with **vectors**!

→ k-means cluster, optimize k with human observation (samples: sentences)

```
s_emb_matrix = [s0, s2, ..., s9] # s are embedding vectors of sentences
nclusters= 3
km = KMeans(nclusters)
km.fit(s_emb_matrix)
clusters = {} # key: label, values: index of sentence
for i, label in enumerate(km.labels_):
    clusters[label].append(i) # append the index to the corresponding label
```

→ get *clusters*[0] : s_0, s_4, s_5 with top-n words list *topWords*[0] : $w_0, w_1, \dots w_n$

→ under each cluster, check similarities of words (samples: words)

Todo Stage 1.2

Build a k-mean clustering model for identifying sub-topics

Now, let's start working with **vectors**!

→ k-means cluster, optimize k with human observation (**samples: sentences**)

```
s_emb_matrix = [s0, s2, ..., s9] # s are embedding vectors of sentences
nclusters= 3
km = KMeans(nclusters)
km.fit(s_emb_matrix)
clusters = {} # key: label, values: index of sentence
for i, label in enumerate(km.labels_):
    clusters[label].append(i) # append the index to the corresponding label
```

→ get *clusters*[0] : s_0, s_4, s_5 with top-n words list *topWords*[0] : w_0, w_1, \dots, w_n

→ under each cluster, check similarities of words (**samples: words**)

→ interpret/extract *clusters*[0] to an aspect (like environment).

→ sentences s_0, s_4, s_5 are assigned to environment.

Todo Stage 1.3

Select a sentiment pre-trained model

Since we get $s_0 \rightarrow$ "environment", we use pre-trained VADER classifier³ for each sentence, output a **2-tuple**.

³Hutto and Gilbert 2014.

Todo Stage 1.3

Select a sentiment pre-trained model

Since we get $s_0 \rightarrow$ "environment", we use pre-trained VADER classifier³ for each sentence, output a **2-tuple**.

Example

s_0 ="not only improves the fruit quality, but is a lot better on our environment as well."

³Hutto and Gilbert 2014.

Todo Stage 1.3

Select a sentiment pre-trained model

Since we get $s_0 \rightarrow$ "environment", we use pre-trained VADER classifier³ for each sentence, output a **2-tuple**.

Example

$s_0 =$ "not only improves the fruit quality, but is a lot better on our environment as well."

$\rightarrow arr = [0.2, \dots, 0.1] \rightarrow kmean\ clustering \rightarrow$ "Environment" (result from stage 1.2)

³Hutto and Gilbert 2014.

Todo Stage 1.3

Select a sentiment pre-trained model

Since we get $s_0 \rightarrow$ "environment", we use pre-trained VADER classifier³ for each sentence, output a **2-tuple**.

Example

s_0 ="not only improves the fruit quality, but is a lot better on our environment as well."

\rightarrow $arr = [0.2, \dots, 0.1] \rightarrow$ *kmean clustering* \rightarrow "Environment" (result from stage 1.2)

\rightarrow *VADER classifier*

```
analyzer = SentimentIntensityAnalyzer()  
print(analyzer.polarity_scores(arr))
```

\rightarrow {'pos': **0.74**, 'neu': 0.26, 'neg': 0.0}

³Hutto and Gilbert 2014.

Todo Stage 1.3

Select a sentiment pre-trained model

Since we get $s_0 \rightarrow$ "environment", we use pre-trained VADER classifier³ for each sentence, output a **2-tuple**.

Example

s_0 ="not only improves the fruit quality, but is a lot better on our environment as well."

\rightarrow $arr = [0.2, \dots, 0.1] \rightarrow$ *kmean clustering* \rightarrow "Environment" (result from stage 1.2)

\rightarrow *VADER classifier*

```
analyzer = SentimentIntensityAnalyzer()  
print(analyzer.polarity_scores(arr))
```

\rightarrow {'pos': **0.74**, 'neu': 0.26, 'neg': 0.0}

\rightarrow (**Environment**, **pos**)

³Hutto and Gilbert 2014.

Todo Stage 2

Visualize distribution on single article (and its comments)

Now, the sentence s is converted into a tuple $t = (Aspect, +/ - / o)$.

Todo Stage 2

Visualize distribution on single article (and its comments)

Now, the sentence s is converted into a tuple $t = (Aspect, +/ - / o)$.

→ Take all samples $[t_0, t_1, t_2, \dots]$ from one article

→ accumulate and normalize them

Todo Stage 2

Visualize distribution on single article (and its comments)

Now, the sentence s is converted into a tuple $t = (Aspect, + / - / o)$.

→ Take all samples $[t_0, t_1, t_2, \dots]$ from one article

→ accumulate and normalize them

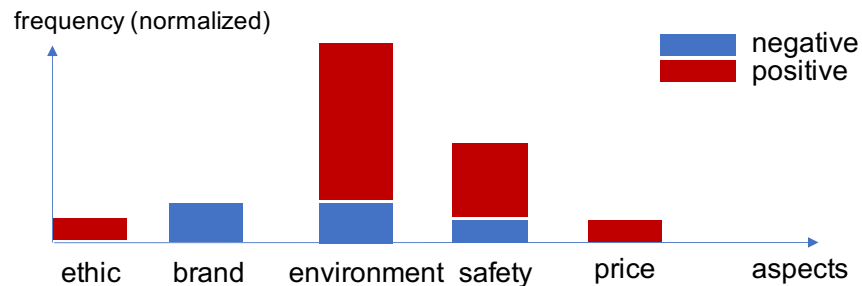


Figure: Aspects distribution on one text

Todo Stage 3

Investigate Media Agenda Setting

Does Y (news) has influence on X (comments/Quora)?

Todo Stage 3

Investigate Media Agenda Setting

Does Y (news) has influence on X (comments/Quora)?

→ discover/measure the relationship between X, Y

- Pearson correlations (time not considered)
- Local Similarity Analysis (LSA) statistic identifies the existence of local and lagged relationships
- Granger causality (Does Y_t help to predict X_{t+1} ?)
- Lagged Correlation (response after a lapse of time, how strong correlation)

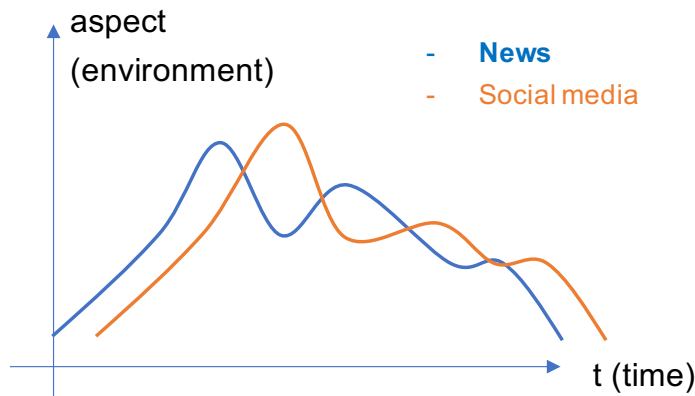
Todo Stage 3

Investigate Media Agenda Setting

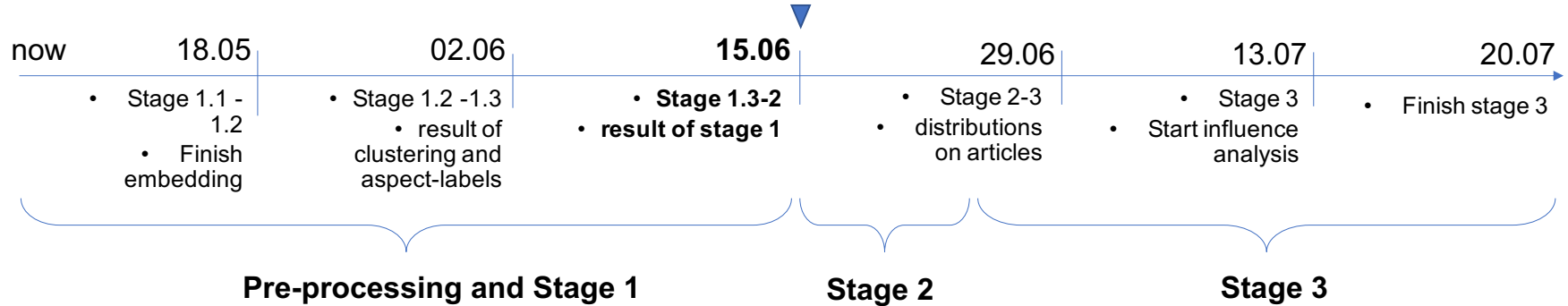
Does Y (news) has influence on X (comments/Quora)?

→ discover/measure the relationship between X, Y




- Pearson correlations (time not considered)
- Local Similarity Analysis (LSA) statistic identifies the existence of local and lagged relationships
- Granger causality (Does Y_t help to predict X_{t+1} ?)
- Lagged Correlation (response after a lapse of time, how strong correlation)



Milestones



References

-  Chidambaram, M. et al. (2018). “Learning cross-lingual sentence representations via a multi-task dual-encoder model”. In: *arXiv preprint arXiv:1810.12836*.
-  dalequark (2019). *predicting movie reviews with bert on tf hub*. URL: https://github.com/google-research/bert/blob/master/predicting_movie_reviews_with_bert_on_tf_hub.ipynb (visited on 2019).
-  Hutto, C. J. and E. Gilbert (2014). “Vader: A parsimonious rule-based model for sentiment analysis of social media text”. In: *Eighth international AAAI conference on weblogs and social media*.