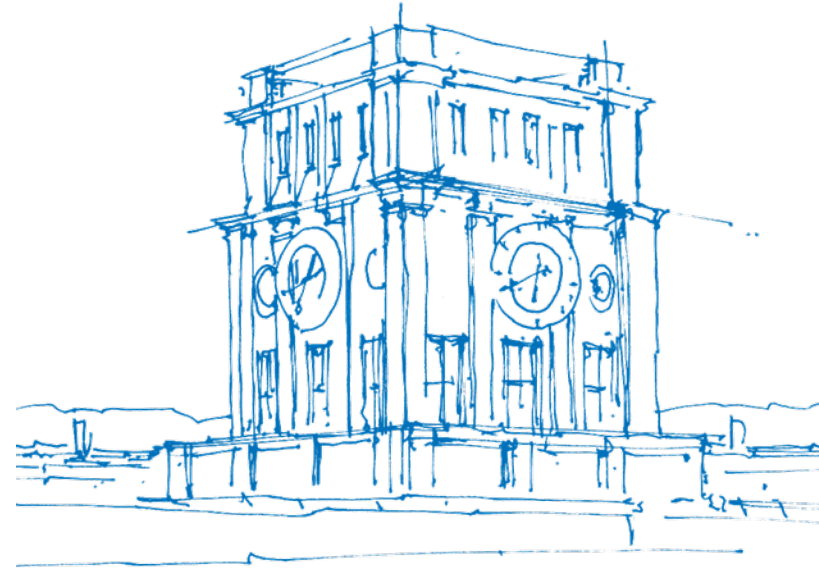


Clustering-Based Sentiment Analysis for Media Agenda Setting

Opinion Lab Group 2.3

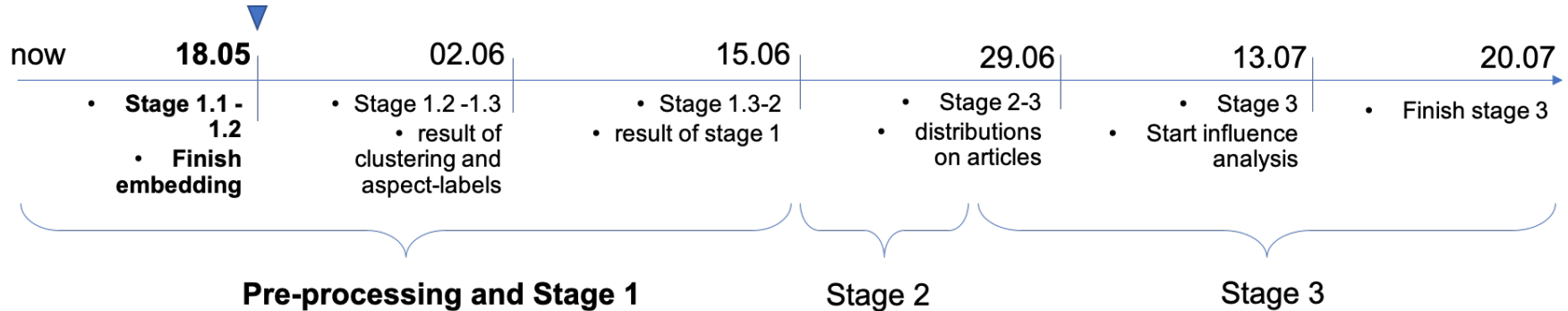
Wing Sheung Leung, Qiaoxi Liu

May 18, 2020



TUM Uhrenturm

Milestones



Overview

Stage 1: Generate sentence embeddings with our corpus

1.1 Embeddings

XLING sentence-level embeddings

Indexing sentences

1.2 Kmeans and Elbow Method

`sklearn.cluster.MinibatchKMeans`

Elbow Method for determining optimal k

Future Plan

Stage 1.1: XLING sentence-level embeddings

Our encoder:

```
en_de_embed = hub.Module("https://tfhub.dev/google/universal-sentence-encoder-xling/en-de/1")
```

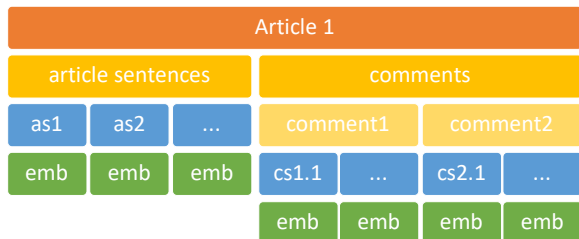


Figure: hierarchy of json

```

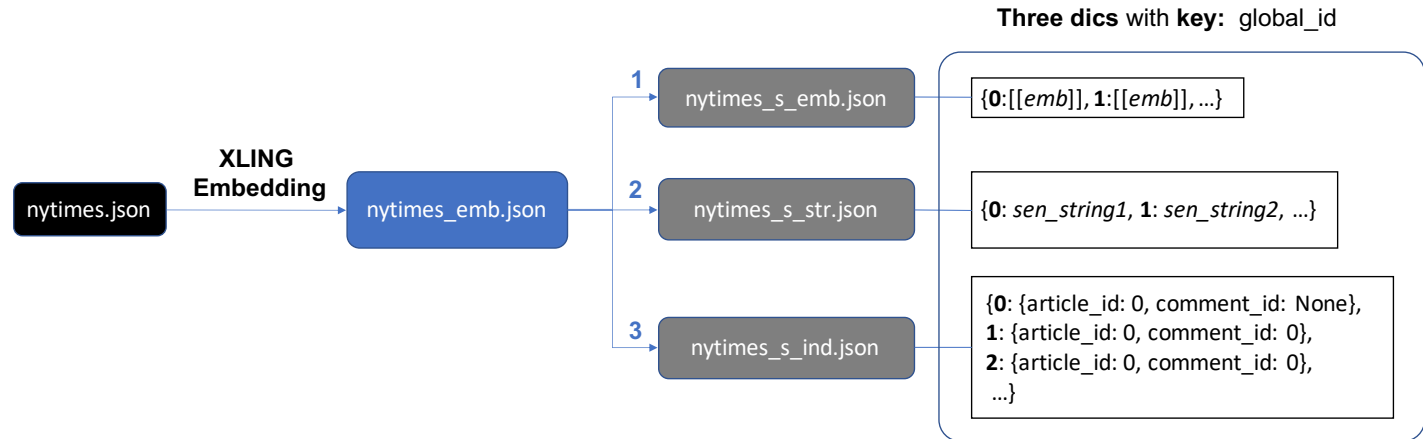
{
  "id": ObjectId("5ebb9ff15edf889717ef47c2"),
  "article_title": "Why are organic eggs typically brown?",
  "article_author": Array,
  "article_time": "2017-10-12 06:23:00",
  "article_text": "",
  "comments": Array
    ~ 0: Object
      "comment_id": "answer_59679995"
      "comment_author": Object
      "comment_time": "2017-10-12 09:22:00"
      "comment_text": "The main reason real organic eggs are brown is that most the heirloom ..."
      "comment_rating": 7
      "processed_comment_text": "The main reason real organic eggs are brown is that most the heirloom ..."
      "comment_sentences": Array
        ~ 0: Object
          "sentence": "The main reason real organic eggs are brown is that most the heirloom ..."
          "embedding": Array
  }
}
  
```

Figure: example from one article in quora

Stage 1.1: Reorganize sentences from articles

Extract embeddings from nested dict to get three separate files:

- (i) embedding vectors → do clustering
- (ii) strings → after clustering to generate word list, helping for define aspect labels
- (iii) indexes → after sentence labeling, to do article (comments) level statistic



Stage 1.2: sklearn.cluster.MinibatchKMeans

```
class KMeansClustering():
def __init__(self, k, X, is_mini_batch = True, plot_bar_chart = True):
    self.k = k
    self.X = np.array(X).reshape(len(X), 512)
    self.km = MiniBatchKMeans(n_clusters=k, init='k-means++', batch_size=3000, compute_labels=True).fit(self.X)
```

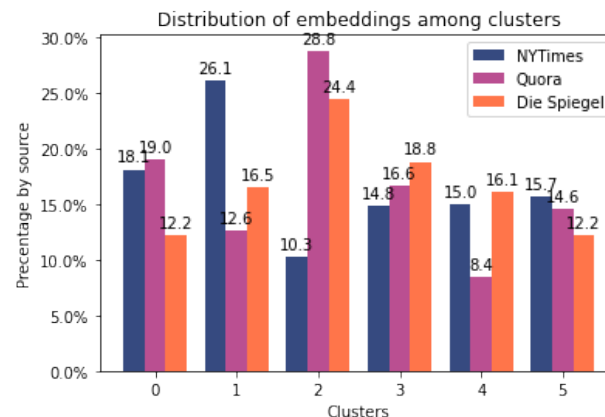
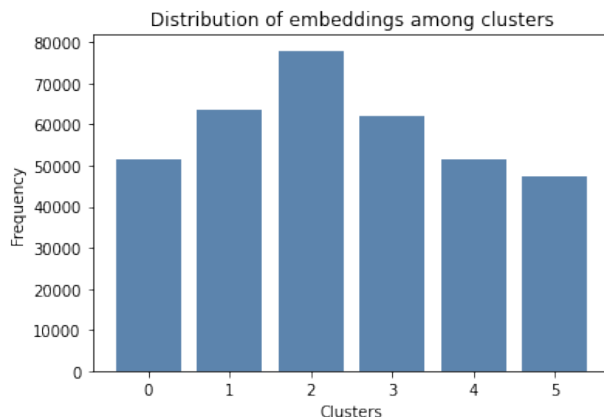


Figure: Example of distribution of embeddings of all tokenized sentences from the three sources among 6 clusters

Stage 1.2: sklearn.cluster.MinibatchKMeans

Why MiniBatchKmeans instead of original sklearn.cluster.KMeans

XLING sentence level embeddings is generated in 512 dimensions for each tokenized sentence by NLTK.

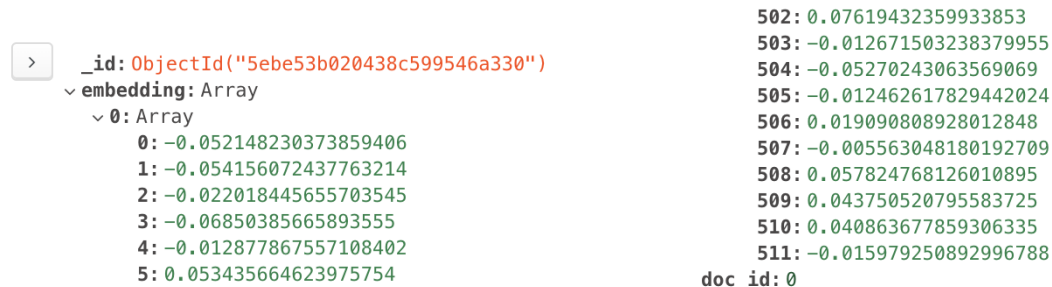


Figure: XLING embedding output for a sample sentence. Left: First 6 dimensions. Right: Last 10 dimensions

Source	Embedding JSON size	Original corpus size
New York Times	827 MB	55.9 MB
Quora	638 MB	15.9 MB
Die Speigel	2.3 GB	131 MB

Table: Embeddings generated are greatly larger then the original corpus size

Stage 1.2: sklearn.cluster.MinibatchKMeans

Why MiniBatchKmeans instead of original sklearn.cluster.KMeans

Just loading all sentence embeddings in Google Colaboratory, 6.36 GB out of the given 12.72 GB RAM had already been used up.

MiniBatchKMeans is faster and helps to prevent the session from crushing, however, gives slightly different results.

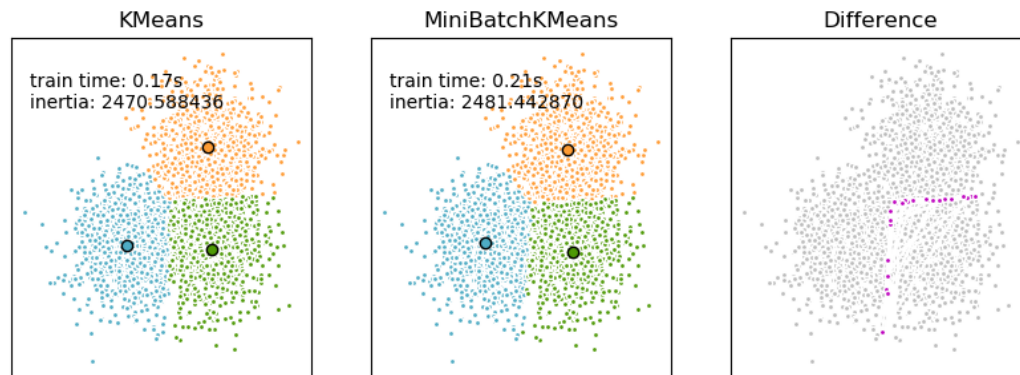


Figure: Extracted from scikit-learn; Data points classified differently are shown as purple points in 'Difference' block

https://scikit-learn.org/stable/auto_examples/cluster/plot_mini_batch_kmeans.html

Stage 1.2: Elbow Method for determining optimal k

```
K = range(2, 21)
for k in K:
    model = KMeansClustering(k, X)
    distortions.append(model.km.inertia_)
```

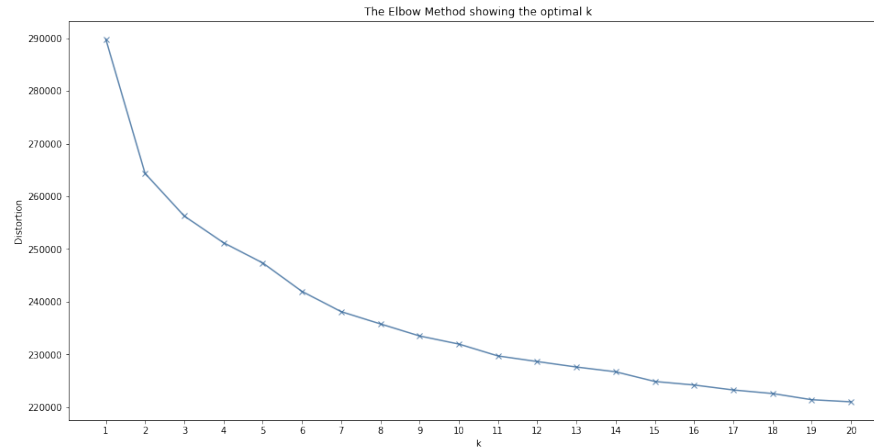


Figure: No distinguishable elbow of the curve for determination of optimal k

Future Plan

