Confronting Sparseness and High Dimensionality in Short Text Clustering via Feature Vector Projections

Leonidas Akritidis¹, Miltiadis Alamaniotis², Athanasios Fevgas³, Panayiotis Bozanis¹

¹School of Science and Technology, International Hellenic University ²Department of Electrical and Computer Engineering, University of Texas at San Antonio

³Department of Electrical and Computer Engineering, University of Thessaly

The 32nd IEEE International Conference on Tools with Artificial Intelligence (ICTAI 2020) November 9-11, 2020, Virtual Conference

Short Text Clustering

- A problem that focuses on the unsupervised grouping of similar short text documents, or entitled entities.
- Specialization of the general text clustering problem which in turn is a specialization of the generic data clustering problem.
- Quite popular. Important applications include cluster creation from:
 - microblogs,
 - entitled entities (e.g. news headlines, product titles, etc.),
 - FAQs,
 - search result snippets,
 - etc.

Short-Texts: 2 difficult problems

- Short texts suffer from 2 challenging issues:
- **Data sparseness:** it concerns the absence of important features (terms) from a portion of the input records.

- It blurs the similarities among the involved entities.

- High dimensionality: originates from the usage of a huge number of features (terms, n-grams, skip grams etc.) in short texts.
 - It is the source of the "Curse of Dimensionality".

Clustering algorithms: Primary goals

- A clustering algorithm must group its input records by fulfilling (at least) two goals:
- Homogeneity: inclusion of only similar elements within a cluster, and
- **Completeness:** *all* similar elements are grouped into the same cluster.

VEPHC

- A two stage short text clustering algorithm.
- Stage 1 (VEP part) projects the original feature vectors onto a lower dimensional space by generating all the feature combinations of the initial text vectors.
- A feature combination corresponds to a projection vector and is treated as a candidate cluster label.
- Each projection vector is assigned a score which favors completeness and the homogeneity.
- All documents whose feature vectors are projected onto the same space are grouped into the same cluster.
- Stage 2 (HC part) consists of a post-processing algorithm that enhances the homogeneity and the completeness of the clusters that were generated by the previous stage.

VEPHC Part 1: Feature combinations as vector projections

- An integer hyper-parameter K is initially set.
- For each input vector x_i we construct a set X'_i containing all the possible projections of x_i with a dimensionality of 1, 2, ... K elements.
- In other words, X[']_i contains all the {1, 2, ... K} combinations of the components of x_i.
- If $\mathbf{x}_i = \{x_1, x_2, x_3\}$, then $\mathbf{X}'_i = \{\{x_1\}, \{x_2\}, \{x_3\}, \{x_1x_2\}, \{x_1x_3\}, \{x_2x_3\}\}$.

Scoring the vector projections

For each input document, each projection vector x'_{ij} in X'_i is assigned a score:

$$S_{\mathbf{x}'_{ij}} = \frac{1}{l_{\mathbf{x}'_{ij}}} \log f_{\mathbf{x}'_{ij}} \sum_{\forall x_{ij} \in \mathbf{x}'_{ij}} x_{ij}^K,$$

- Where $l_{\mathbf{x}_{ij}'}$ is the length of \mathbf{x}_{ij}'
- $f_{\mathbf{x}'_{ij}}$ is the global frequency of \mathbf{x}'_{ij} .
- x_{ij} : a component of \mathbf{x}'_{ij} .

Scoring the vector projections (2)

If the original feature vectors have been created with tf-idf:

$$S_{\mathbf{x}'_{ij}} = \frac{1}{l_{\mathbf{x}'_{ij}}} \log f_{\mathbf{x}'_{ij}} \sum_{\forall x_{ij} \in \mathbf{x}'_{ij}} \left(tf_{ij} \cdot \log \frac{n}{f_j} \right)^K$$

- Where *n* is the number of the input records.
- f_j is the freq of the *j*-th component of \mathbf{x}'_{ij} .
- The highest-scoring projection of x_i is declared as the dominant cluster label x^{*}_i.
- The input is clustered in \mathbf{x}_i^* .

Example of clustering

- $\mathbf{x}_1 = \{x_1, x_2, x_3, x_4, x_5\}$
- $\mathbf{x}_2 = \{0, x_2, x_3, x_4, 0\}$
- If $\mathbf{x}_1^* = \mathbf{x}_2^* = \{x_2, x_4\}$
- Then x₁ and x₂ will be placed inside the same cluster.

Consequences

- The higher the dimensionality of the dominant vector, the more similar documents the corresponding cluster contains.
- The documents included in that cluster will have more words in common.
- High/Low dimensional dominant projection vectors lead to:
 - More/Lesser homogeneous clusters.
 - Lesser/More complete clusters.

VEPHC Part 2: Cluster refinement stage

- At this point, the dominant projections of the input short text documents have been computed.
- The input documents have been clustered in the corresponding cluster labels.
- The second part introduces a refinement algorithm which was designed to further improve homogeneity and completeness.

VEPHC Part 2 – Stage 1

- Split the existing clusters
 & create new clusters.
- In short:
- Compute the clustroid point of each cluster.
 - i.e. the point that has max sim with the rest of the elements of the cluster.

```
Algorithm 1: HC Part, Phase 1: Element Transfer and
 Creation of New Clusters
 1 initialize an empty cluster C;
 2 for each cluster c \in C do
          \mathbf{u}_c \leftarrow \text{clustroid of } c;
 3
 4 end
 5 for each cluster c \in C do
           for each element \mathbf{x} \in c do
                 T_{\mathbf{x},\mathbf{u}_c} \leftarrow \text{similarity between } \mathbf{x} \text{ and } \mathbf{u}_c;
 7
                 if T_{\mathbf{x},\mathbf{u}_c} < T_h then
  8
                       remove x from c;
  Q
                       T_{\max} \leftarrow 0, c^* \leftarrow \text{NULL};
 10
                       for each cluster c' \in C do
 11
                             T_{\mathbf{x},\mathbf{u}'_{c}} \leftarrow \text{similarity between } \mathbf{x} \text{ and } \mathbf{u}'_{c};
 12
                             if T_{\mathbf{x},\mathbf{u}_{\alpha}'} > T_{\max} then
 13
                                   T_{\max} \leftarrow T_{\mathbf{x},\mathbf{u}_{\alpha}'};
 14
                                    c^* \leftarrow c':
 15
                              end
 16
 17
                       end
                       if T_{\rm max} > T_c then
 18
                             c^* \leftarrow c^* \cup \{\mathbf{x}\};
 19
                       else
 20
                             \mathcal{C} \leftarrow \mathcal{C} \cup \{\mathbf{x}\};
 21
22
                       end
23
                 end
24
           end
25 end
26 C_{new} \leftarrow \emptyset;
27 for each element \mathbf{x} \in \mathcal{C} do
          remove x from C;
28
          T_{\max}, c^* \leftarrow perform steps 11–17 on C_{new};
29
          if T_{\max} \geq T_c then
30
                 c^* \leftarrow c^* \cup \{\mathbf{x}\};
31
32
           else
                 create and initialize new cluster c_{new} \leftarrow \emptyset;
33
                 c_{new} \leftarrow c_{new} \cup \{\mathbf{x}\}, \mathbf{u}_{c_{new}} \leftarrow \mathbf{x};
34
                 C \leftarrow C \cup \{c_{new}\}, C_{new} \leftarrow C_{new} \cup \{c_{new}\};
35
36
           end
37 end
```

VEPHC Part 2 – Stage 1

- For each cluster:
- Evict the cluster elements that are too far away from the clustroid.
- Evicted elements may be:
 - inserted into another,
 more similar cluster, OR
 - placed into a new, empty cluster.

```
L. Akritidis, M. Alamaniotis, A. Fevgas, P. Bozanis
```

```
Algorithm 1: HC Part, Phase 1: Element Transfer and
 Creation of New Clusters
 1 initialize an empty cluster C;
 2 for each cluster c \in C do
          \mathbf{u}_c \leftarrow \text{clustroid of } c;
 3
 4 end
 5 for each cluster c \in C do
           for each element \mathbf{x} \in c do
                 T_{\mathbf{x},\mathbf{u}_c} \leftarrow \text{similarity between } \mathbf{x} \text{ and } \mathbf{u}_c;
                if T_{\mathbf{x},\mathbf{u}_c} < T_h then
  8
                       remove x from c;
                       T_{\max} \leftarrow 0, c^* \leftarrow \text{NULL};
 10
                       for each cluster c' \in C do
 11
                             T_{\mathbf{x},\mathbf{u}'_{a}} \leftarrow \text{similarity between x and } \mathbf{u}'_{c};
 12
                             if T_{\mathbf{x},\mathbf{u}_{\alpha}'} > T_{\max} then
 13
                                   T_{\max} \leftarrow T_{\mathbf{x},\mathbf{u}_{\alpha}'};
 14
                                    c^* \leftarrow c':
 15
                             end
 16
 17
                       end
                       if T_{\rm max} > T_c then
 18
                             c^* \leftarrow c^* \cup \{\mathbf{x}\};
 19
                       else
 20
                             \mathcal{C} \leftarrow \mathcal{C} \cup \{\mathbf{x}\};
 21
22
                       end
23
                end
           end
24
25 end
26 C_{new} \leftarrow \emptyset;
27 for each element \mathbf{x} \in \mathcal{C} do
           remove x from C;
28
          T_{\max}, c^* \leftarrow perform steps 11–17 on C_{new};
29
          if T_{\max} \geq T_c then
30
                c^* \leftarrow c^* \cup \{\mathbf{x}\};
31
32
           else
                 create and initialize new cluster c_{new} \leftarrow \emptyset;
33
                 c_{new} \leftarrow c_{new} \cup \{\mathbf{x}\}, \mathbf{u}_{c_{new}} \leftarrow \mathbf{x};
34
                C \leftarrow C \cup \{c_{new}\}, C_{new} \leftarrow C_{new} \cup \{c_{new}\};
35
36
           end
37 end
```

VEPHC Part 2 – Stage 2

- Cluster merging stage.
- For each cluster:
- Search for other, highly similar clusters.
- If the similarity with the most similar cluster exceeds a threshold, then we merge the two clusters.

Algorithm 2: HC Part, Phase 2: Cluster Merging

```
1 for each cluster c \in C do
        \mathbf{u}_c \leftarrow \text{clustroid of } c;
3 end
 4 sim\_array[*] \leftarrow 0, msc\_array[*] \leftarrow -1;
 5 for each cluster c \in C do
        msc\_array[c] \leftarrow most similar cluster (MSC)
        sim\_array[c] \leftarrow similarity with MSC;
 8 end
9 while merge do
        merge \leftarrow FALSE;
10
        c' \leftarrow \text{NULL}, c'' \leftarrow \text{NULL};
11
        for each cluster c \in C do
12
             T_{\text{max}} \leftarrow 0;
13
             if sim\_array[c] > T_{max} then
14
                  T_{\max} \leftarrow sim\_array[c];
15
                  c' \leftarrow c, c'' \leftarrow msc\_array[c];
16
17
             end
18
        end
        if T_{\rm max} > T_c then
19
             merge \leftarrow TRUE;
20
             c' \leftarrow c' \cup c'':
21
             \mathbf{u}_c' \leftarrow \text{clustroid of } c';
22
             msc\_array[c'] \leftarrow MSC;
23
             sim\_array[c'] \leftarrow similarity with MSC;
24
             C \leftarrow C - \{c''\};
25
             update msc_array, sim_array;
26
        end
27
```

Experiments

- We attested VEPHC by using two datasets:
 - A dataset of product titles from Pricerunner.
 - A dataset of news titles (headlines).
- Versus 4 generic data clustering + 3 short text clustering methods.
- In terms of both clustering quality (P, R, F1, NMI) and execution times.

	n	C	$l_{ m max}$	l_{ave}
PriceRunner (PRUN)	35311	13233	41	8.23
UCI News Aggregator (UNA)	50138	823	15	7.05

EXPERIMENTAL DATASETS

Dependence from the hyper-parameters



Fig. 1. Performance fluctuation of VEPHC (without the verification stage) against varying values of the hyper parameter K, on the PriceRunner (left), and UCI News Aggregator (right) datasets.

Clustering performance

PERFORMANCE EVALUATION OF VARIOUS CLUSTERING ALGORITHMS ON THE PRICERUNNER AND UCI NEWS AGGREGATOR DATASETS

Method	Setup	PriceRunner			UCI	UCI News Aggregator		
		F1	NMI	Time	F1	NMI	Time	
VEPHC	$K_1 = 6, K_2 = 2$	0.385	0.946	305.1	0.471	0.851	10.2	
k-Means	k = C , I = 10	0.048	0.399	494.2	0.028	0.697	104.3	
Agglomerative	_	0.314	0.935	166.7	0.454	0.835	2492.8	
Leader Clustering	-	0.306	0.934	17.2	0.284	0.778	17.1	
DBSCAN	minPoints = 2	0.055	0.425	58.4	0.256	0.790	130.5	
GSDMM-1	$\alpha = \beta = 0.1, k = 1.5 \cdot C $	0.002	0.404	3246.2	0.171	0.753	390.4	
GSDMM-2	$\alpha = 0.001, \beta = 0.01, k = 1.5 \cdot C $	0.008	0.631	3317.3	0.424	0.820	453.9	
Spherical k-Means	k = C , I = 10	0.226	0.903	426.1	0.335	0.757	119.6	
vk-Means	k = C , I = 10	0.220	0.900	484.4	0.154	0.614	104.7	
Cosine similarity	tf - idf weights	0.248	-	31.3	0.388	_	75.2	
Jaccard Index	tf - idf weights	0.237	-	31.5	0.388	-	74.9	

Conclusions

- We introduced VEPHC, a two-stage clustering algorithm for short text documents.
- Designed to limit the native sparseness and high dimensionality of short texts.
- Key elements:
 - Projection of the feature vectors onto a lower dimensional space by the construction and scoring of feature combinations.
 - A feature combination can be viewed as a projection vector.
 - Identify the dominant (i.e., the highest scoring) projection vector.
 - Two or inputs having the same projection vector are grouped into the same cluster.

Conclusions

- Key elements (continued):
 - Post processing stage: Split the initial clusters by removing the most distant items.
 - Place the evicted items into other clusters, or create new ones.
 - Merge the most similar clusters.
- Experiments conducted with two test datasets.
- The results showed significant performance improvements over the current state-of-the-art methods.
- In terms of both effectiveness and efficiency.

Thank you Any questions?