

# Xử lý ngôn ngữ tự nhiên

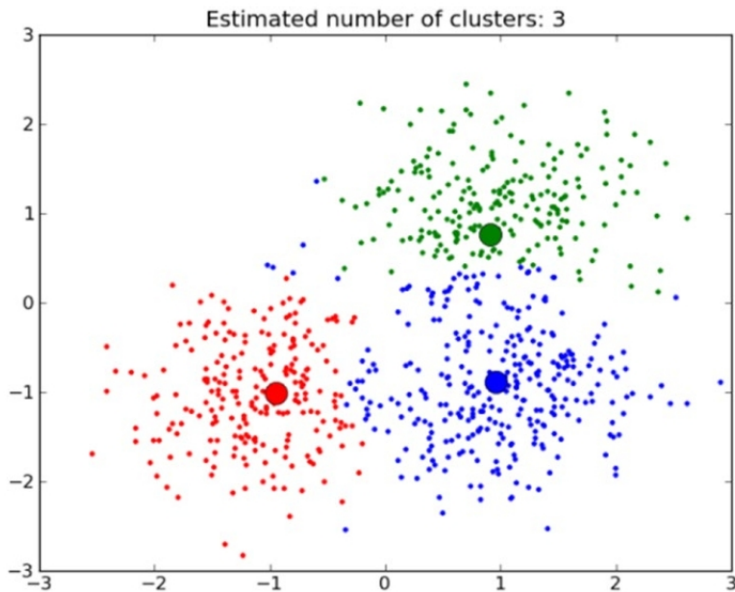
## Chương 8: Text similarity and clustering - part 2

Khoa CNTT, Đại học Kỹ thuật - Công nghệ Cần Thơ  
Lưu hành nội bộ

# Document clustering

- Document clustering, cluster analytics là lĩnh vực nghiên cứu của NLP ứng dụng học không giám sát (unsupervised learning)
- Cơ sở của document clustering là sự bắt đầu với toàn bộ corpus, sau đó phân tách (segregation) thành các nhóm (groups) dựa trên tính chất khác biệt (distinctive properties), thuộc tính (attributes), features của documents

- Không biết trước bao nhiêu group cần phân tách
- Các documents thuộc về 1 group sẽ có sự tương đồng cao hơn so với các documents thuộc group khác
- `http://scikit-learn.org/stable/modules/clustering.html`



# Một số kỹ thuật và phương pháp tìm cluster

- Hierarchical clustering models
- Centroid-based clustering models
- Distribution-based clustering models
- Density-based clustering models

# Clustering greatest movies of all time

- Chúng ta tiến hành clustering 100 bộ phim phổ biến dựa trên internet movie database<sup>1</sup>
- Mỗi bộ phim gồm tên phim và phần tóm tắt (synopses)
- Top 100 Greatest Movies of All Time (The Ultimate List):  
<https://www.imdb.com/list/ls055592025/>
- movie\_data.csv trên hệ thống LearnWeb

---

<sup>1</sup>[www.imdb.com](http://www.imdb.com)

```
1 import numpy as np
2 import pandas as pd
3
4 movie_data = pd.read_csv('movie_data.csv')
5
6 print(movie_data.head())
7
8 movie_titles = movie_data['Title'].tolist()
9 movie_synopses = movie_data['Synopsis'].tolist()
10
11 print('Movie:', movie_titles[0])
12 # print first 1000 characters
13 print('Movie Synopsis:', movie_synopses[0][:1000])
```

#	Title	Synopsis
#0	The Godfather	In late summer 1945, guests are gathered for t...
#1	The Shawshank Redemption	In 1947, Andy Dufresne (Tim Robbins), a banker...
#2	Schindler's List	The relocation of Polish Jews from surrounding...
#3	Raging Bull	The film opens in 1964, where an older and fat...
#4	Casablanca	In the early years of World War II, December 1...

#Movie: The Godfather

#Movie Synopsis: In late summer 1945, guests are gathered for the wedding reception of Don Vito Corleone's daughter Connie (Talia Shire) and Carlo Rizzi (Gianni Russo). Vito (Marlon Brando), the head of the Corleone Mafia family, is known to friends and associates as "Godfather." He and Tom Hagen (Robert Duvall), the Corleone family lawyer, are hearing requests for favors because, according to Italian tradition, "no Sicilian can refuse a request



```

22
23 def remove_characters_after_tokenization(tokens):
24     import re
25     import string
26     pattern = re.compile("[{}]" .format(re.escape(string.punctuation)))
27     filtered_tokens = list(filter(None, [pattern.sub("", token) for token
28         in tokens]))
29     return filtered_tokens
30
31 def remove_stopwords(tokens, language="english"):
32     stopword_list = nltk.corpus.stopwords.words(language)
33     filter_tokens = [token for token in tokens if token not in
34         stopword_list]
35     return filter_tokens

```

34

35 `def remove_repeated_characters(tokens):`36  `from nltk.corpus import wordnet`37  `import re`

38

39  `repeat_pattern = re.compile(r'(\w*)(\w)\2(\w*)')`40  `match_substitution = r'\1\2\3'`41  `def replace(old_word):`42  `if wordnet.synsets(old_word):`43  `return old_word`44  `new_word = repeat_pattern.sub(match_substitution, old_word)`45  `return replace(new_word) if new_word != old_word else new_word`

46

47  `correct_tokens = [replace(word) for word in tokens]`48  `return correct_tokens`

49

```
50 def text_stemming(tokens):
51     from nltk.stem.snowball import SnowballStemmer
52     stemmer = SnowballStemmer("english")
53     def replace(old_word):
54         return stemmer.stem(old_word)
55
56     filter_tokens = [replace(token) for token in tokens]
57     return filter_tokens
58
59 def text_normalization(tokens):
60     tokens = remove_characters_after_tokenization(tokens)
61     tokens = text_stemming(tokens)
62     tokens = remove_stopwords(tokens)
63     tokens = remove_repeated_characters(tokens)
64     return tokens
65
66 def tokens_to_string(tokens):
67     return " ".join(tokens)
```



```
68
69 import nltk
70 default_wt = nltk.word_tokenize
71
72 tokens_synoses = []
73 for syn in movie_synopses:
74     tokens_synoses.append(default_wt(text=syn))
75
76 tokens_synoses_norm = []
77 for tokens in tokens_synoses:
78     tokens_synoses_norm.append(text_normalization(tokens))
79
80 tokens_synoses_norm_string = []
81 for tokens_norm in tokens_synoses_norm:
82     tokens_synoses_norm_string.append(tokens_to_string(tokens_norm))
```

83

```
84 from sklearn.feature_extraction.text import CountVectorizer,  
    TfidfVectorizer
```

85

```
86 def build_feature_matrix(documents, feature_type='frequency',  
    ngram_range=(1, 1), min_df=0.0, max_df=1.0):  
87     feature_type = feature_type.lower().strip()
```

88

```
89     if feature_type == 'binary':
```

```
90         vectorizer = CountVectorizer(binary=True, min_df=min_df, max_df=  
    max_df, ngram_range=ngram_range)
```

```
91     elif feature_type == 'frequency':
```

```
92         vectorizer = CountVectorizer(binary=False, min_df=min_df, max_df=  
    max_df, ngram_range=ngram_range)
```

```
93     elif feature_type == 'tfidf':
```

```
94         vectorizer = TfidfVectorizer(min_df=min_df, max_df=max_df,  
    ngram_range=ngram_range)
```

```

95     else:
96         raise Exception("Wrong feature type entered. Possible values: '
          binary', 'frequency', 'tfidf'")
97
98     feature_matrix = vectorizer.fit_transform(documents).astype(float)
99     return vectorizer, feature_matrix
100
101 vectorizer, feature_matrix = build_feature_matrix(
          tokens_synoses_norm_string, feature_type="tfidf")
102 print(feature_matrix.shape)
103 #(100, 10797)
104
105 #get feature names
106 feature_names = vectorizer.get_feature_names()
107 print(feature_names[:20])
108 #['10', '100', '1000', '10000', '1000000', '101st', '108', '11', '110',
          '11th', '12', '120pound', '1212', '1280', '#1297', '1298', '12fot
          ', '12yearold', '13', '130']

```

# K-means clustering

- k-means clustering algorithm thuộc nhóm centroid-based clustering
- Nhóm data vào các clusters dựa trên sự khác biệt / giống nhau
- Tiêu chuẩn (criteria) hoặc phép đo (measure) mà giải thuật sử dụng là tối thiểu hóa *inertia*, within-cluster sum-of-squares
- Xác định giá trị  $k$  trước khi thực thi giải thuật

- Giải sử ta có dataset  $X$  với  $N$  data points
- Ta muốn nhóm thành  $K$  clusters
- Giải thuật k-means sẽ tuần tự phân chia  $N$  data points vào  $K$  clusters riêng biệt  $C_k$
- Mỗi cluster  $C_k$  được xem như trung bình của các cluster samples (data point thuộc cluster được gọi là cluster samples)
- Giá trị trung bình là cluster centroid  $\mu_k$
- Chọn các giá trị  $\mu_k$  sao cho inertia cực tiểu



$$\min \sum_{i=1}^K \sum_{x_n \in C_i} \|x_n - \mu_i\|^2 \quad (1)$$

- Chọn  $k$  centroids  $\mu_k$  ban đầu bằng ngẫu nhiên
- Cập nhật clusters bằng cách gán data points vào centroid point gần nhất

$$C_k = \{x_n : \|x_n - \mu_k\| \leq \text{all} \|x_n - \mu_l\|\}$$

- Cập nhật giá trị centroids

$$\mu_k = \frac{1}{C_k} \sum_{x_n \in C_k} x_n$$

```
109
110 from sklearn.cluster import KMeans
111
112 def k_means(feature_matrix, num_clusters=5):
113     km = KMeans(n_clusters=num_clusters, max_iter=10000)
114     km.fit(feature_matrix)
115     clusters = km.labels_
116     return km, clusters
117
118 num_clusters = 5
119 km_obj, clusters = k_means(feature_matrix=feature_matrix, num_clusters=
    num_clusters)
120 movie_data['Cluster'] = clusters
```

```
121
122 from collections import Counter
123 # get the total number of movies per cluster
124 c = Counter(clusters)
125 print (c.items())
126 #dict_items([(0, 6), (1, 53), (2, 8), (3, 29), (4, 4)])
```

- Kết quả trả về là danh sách các clusters từ 0 đến 4 cùng với số data points thuộc về clusters đó

127

```
128 def get_cluster_data(clustering_obj, movie_data,
129     feature_names, num_clusters,
130     topn_features=10):
131
132     cluster_details = {}
133     # get cluster centroids
134     ordered_centroids = clustering_obj.cluster_centers_.argsort()[:,
135         ::-1]
136     # get key features for each cluster
137     # get movies belonging to each cluster
138     for cluster_num in range(num_clusters):
139         cluster_details[cluster_num] = {}
140         cluster_details[cluster_num]['cluster_num'] = cluster_num
141         key_features = [feature_names[index]
142             for index
143             in ordered_centroids[cluster_num, :topn_features]]
144         cluster_details[cluster_num]['key_features'] = key_features
```

144

```
movies = movie_data[movie_data['Cluster'] == cluster_num]['Title']  
            .values.tolist()
```

145

```
cluster_details[cluster_num]['movies'] = movies
```

146

147

```
return cluster_details
```

- Kết quả trả về là thông tin cluster:

- key features có tầm quan trọng xác định cluster từ centroids
- movie titles thuộc về từng cluster

```
148
149 def print_cluster_data(cluster_data):
150     # print cluster details
151     for cluster_num, cluster_details in cluster_data.items():
152         print('Cluster {} details:'.format(cluster_num))
153         print('-'*20)
154         print('Key features:', cluster_details['key_features'])
155         print('Movies in this cluster:')
156         print(', '.join(cluster_details['movies']))
157         print('='*40)
158
159 import matplotlib.pyplot as plt
160 from sklearn.manifold import MDS
161 from sklearn.metrics.pairwise import cosine_similarity
162 import random
163 from matplotlib.font_manager import FontProperties
```

164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179

```
def plot_clusters(num_clusters, feature_matrix, cluster_data,
                  movie_data, plot_size=(16,8)):
    # generate random color for clusters
    def generate_random_color():
        color = '#%06x' % random.randint(0, 0xFFFFFF)
        return color

    # define markers for clusters
    markers = ['o', 'v', '^', '<', '>', '8', 's', 'p', '*', 'h', 'H', 'D',
               ', 'd']

    # build cosine distance matrix
    cosine_distance = 1 - cosine_similarity(feature_matrix)

    # dimensionality reduction using MDS
    mds = MDS(n_components=2, dissimilarity="precomputed", random_state
              =1)

    # get coordinates of clusters in new low-dimensional space
    plot_positions = mds.fit_transform(cosine_distance)
    x_pos, y_pos = plot_positions[:, 0], plot_positions[:, 1]
```



```

180 # build cluster plotting data
181 cluster_color_map = {}
182 cluster_name_map = {}
183
184 for cluster_num, cluster_details in cluster_data.items():
185     # assign cluster features to unique label
186     cluster_color_map[cluster_num] = generate_random_color()
187     cluster_name_map[cluster_num] = ', '.join(cluster_details['
key_features'][:5]).strip()
188 # map each unique cluster label with its coordinates and movies
189 cluster_plot_frame = pd.DataFrame({'x': x_pos,
190     'y': y_pos,
191     'label': movie_data['Cluster'].values.tolist(),
192     'title': movie_data['Title'].values.tolist()
193 })
194 grouped_plot_frame = cluster_plot_frame.groupby('label')

```



```

195 # set plot figure size and axes
196 fig, ax = plt.subplots(figsize=plot_size)
197 ax.margins(0.05)
198 # plot each cluster using co-ordinates and movie titles
199 for cluster_num, cluster_frame in grouped_plot_frame:
200     marker = markers[cluster_num] if cluster_num < len(markers) \
201         else np.random.choice(markers, size=1)[0]
202     ax.plot(cluster_frame['x'], cluster_frame['y'],
203             marker=marker, linestyle='', ms=12,
204             label=cluster_name_map[cluster_num],
205             color=cluster_color_map[cluster_num], mec='none')
206     ax.set_aspect('auto')
207     ax.tick_params(axis='x', which='both', bottom='off', top='off',
208 labelbottom='off')
209     ax.tick_params(axis='y', which='both', left='off', top='off',
labelleft='off')

```

```

210 fontP = FontProperties()
211 fontP.set_size('small')
212 ax.legend(loc='upper center', bbox_to_anchor=(0.5, -0.01), fancybox=
213 True, shadow=True, ncol=5, numpoints=1, prop=fontP)
214 #add labels as the film titles
215 for index in range(len(cluster_plot_frame)):
216     ax.text(cluster_plot_frame.ix[index]['x'],
217             cluster_plot_frame.ix[index]['y'],
218             cluster_plot_frame.ix[index]['title'], size=8)
219 # show the plot
220 plt.show()

```

```
221
222 cluster_data = get_cluster_data(clustering_obj=km_obj,
223     movie_data=movie_data,
224     feature_names=feature_names,
225     num_clusters=num_clusters,
226     topn_features=5)
227
228 print_cluster_data(cluster_data)
229
230 plot_clusters(num_clusters=num_clusters,
231     feature_matrix=feature_matrix,
232     cluster_data=cluster_data,
233     movie_data=movie_data,
234     plot_size=(16,8))
```

235

236 #Cluster 0 details:

237 #-----

238 #Key features: ['michael', 'tramp', 'forest', 'eliot', 'corleon']

239 #Movies in this cluster:

240 #The Godfather, The Godfather: Part II, Forrest Gump, E.T. the Extra-  
Terrestrial, City Lights, Tootsie

241 #=====

242

243 #Cluster 1 details:

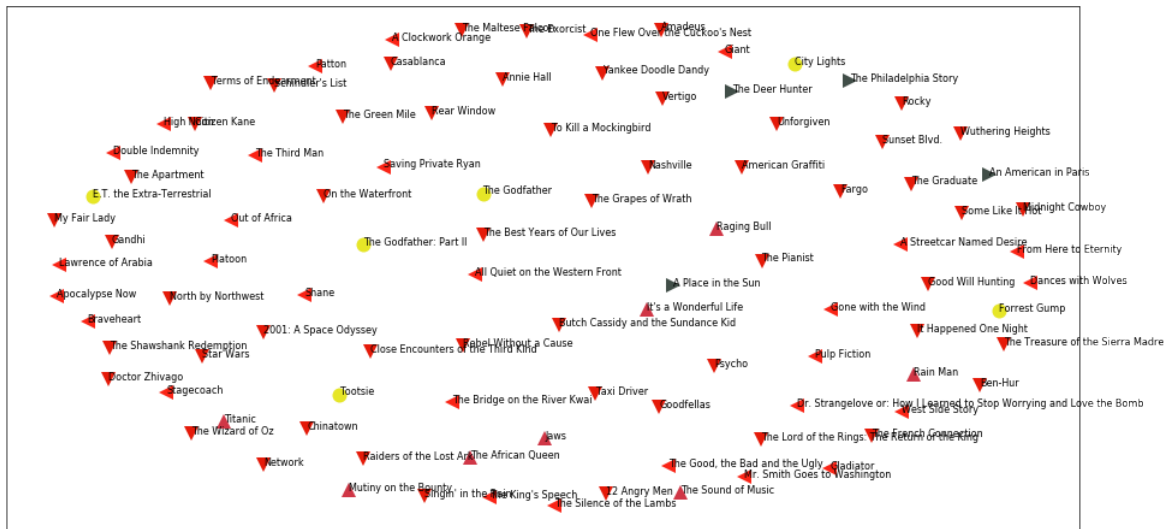
244 #

245 #Key features: ['joe', 'tell', 'teri', 'jeri', 'get']

246 #Movies in this cluster:

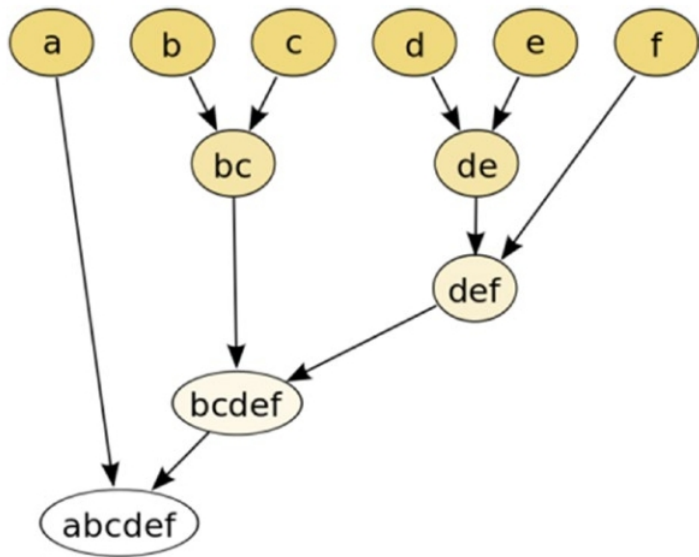
247 #The Shawshank Redemption, Schindler's List, Casablanca, Citizen Kane,  
The Wizard of Oz, Psycho, Sunset Blvd., #Vertigo, On the Waterfront,  
, Star Wars, 2001: A Space Odyssey, Chinatown, Singin' in the Rain,  
Some Like It Hot, #12 Angry Men, Amadeus, Gandhi, The Lord of the  
Rings: The Return of the King, Unforgiven, Raiders of the Lost #Ark  
, Rocky, To Kill a Mockingbird, The Best Years of Our Lives, My  
Fair Lady, Ben-Hur, Doctor Zhivago, Butch #Cassidy and the Sundance  
Kid, The Treasure of the Sierra Madre, The Apartment, The Pianist,  
Goodfellas, The #Exorcist, The French Connection, It Happened One  
Night, Midnight Cowboy, Annie Hall, Good Will Hunting, Terms of #  
Endearment, Fargo, The Grapes of Wrath, The Green Mile, Close  
Encounters of the Third Kind, Network, Nashville, #The Graduate,  
American Graffiti, The Maltese Falcon, Taxi Driver, Wuthering  
Heights, Rebel Without a Cause, Rear #Window, North by Northwest,





# Ward's agglomerative hierarchical clustering

- Cluster thứ bậc (hierarchical clustering) xây dựng cluster lồng vào nhau (nested hierarchy) bởi kết hợp (merging) hoặc chia cách (splitting) cluster kế tiếp (succession)
  - Agglomerative: giải thuật sử dụng tiếp cận bottom-up, bắt đầu từ các data points, kết hợp dần lên cluster lớn
  - Divisive: giải thuật sử dụng tiếp cận top-down, bắt đầu với tất cả data points thuộc về 1 cluster lớn, sau đó chia nhỏ dần đến từng data points





# Ta xét 2 công cụ

- distance metric: dùng để đo độ giống (similarity) nhau hoặc khác nhau giữa các data points; Sử dụng Cosine similarity
- linkage criterion: xác định các mà metric sử dụng để kết hợp clusters; Sử dụng phương pháp Ward's

# The Ward's linkage criterion

- Tối thiểu hóa tổng của bình phương (sum of squared differences) sự khác nhau trong tất cả các clusters

$$d_{ij} = d(\{C_i, C_j\}) = \|C_i - C_j\|^2 \quad (2)$$

248

```
249 from scipy.cluster.hierarchy import ward, dendrogram
```

250

```
251 def ward_hierarchical_clustering(feature_matrix):
```

```
252     cosine_distance = 1 - cosine_similarity(feature_matrix)
```

```
253     linkage_matrix = ward(cosine_distance)
```

```
254     return linkage_matrix
```

255

```
256 def plot_hierarchical_clusters(linkage_matrix, movie_data, figure_size  
    =(8,12)):
```

```
257     # set size
```

```
258     fig, ax = plt.subplots(figsize=figure_size)
```

```
259     movie_titles = movie_data['Title'].values.tolist()
```

```
260     # plot dendrogram
```

```
261     ax = dendrogram(linkage_matrix, orientation="left", labels=  
        movie_titles)
```

```
262 plt.tick_params(axis= 'x',  
263                 which='both',  
264                 bottom='off',  
265                 top='off',  
266                 labelbottom='off')  
267 plt.tight_layout()  
268 plt.savefig('ward_hierachical_clusters.png', dpi=200)  
269 plt.show()  
270
```

```
1  
2 # build ward's linkage matrix  
3 linkage_matrix = ward_hierarchical_clustering(feature_matrix)  
4 # plot the dendrogram  
5 plot_hierarchical_clusters(linkage_matrix=linkage_matrix , movie_data=  
    movie_data , figure_size=(8,10))
```

