

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.feature_extraction.text import TfidfVectorizer
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
9 from sklearn.cluster import KMeans
10 from sklearn.decomposition import PCA
11 import nltk
12 from nltk.corpus import stopwords
13 from nltk.tokenize import word_tokenize
14 from nltk.sentiment import SentimentIntensityAnalyzer
15 from collections import Counter
16 import tensorflow as tf
17 from tensorflow.keras.preprocessing.text import Tokenizer
18 from tensorflow.keras.preprocessing.sequence import pad_sequences
19 from tensorflow.keras.models import Sequential
20 from tensorflow.keras.layers import Embedding, Conv1D, GlobalMaxPooling1D, Dense, LSTM, Dropout
21 import tensorflow_datasets as tfds
22 import plotly.express as px
23 import plotly.graph_objects as go
24 from wordcloud import WordCloud
25
26 # Download required NLTK data
27 nltk.download('stopwords')
28 nltk.download('punkt')
29 nltk.download('vader_lexicon')

```

```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
True

```

```

1 from sklearn.datasets import fetch_20newsgroups
2 import pandas as pd
3
4 news = fetch_20newsgroups(subset="all", remove=("headers", "footers", "quotes"))
5 news_df = pd.DataFrame({"text": news.data, "category": news.target})
6 news_df['category_name'] = news_df['category'].map(lambda x: news.target_names[x])
7 news_df.head()
8

```

```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
True

```

	text	category	category_name
0	\n\nI am sure some bashers of Pens fans are pr...	10	rec.sport.hockey
1	My brother is in the market for a high-perform...	3	comp.sys.ibm.pc.hardware
2	\n\n\n\nFinally you said what you dream abou...	17	talk.politics.mideast
3	\nThink\n\nIt's the SCSI card doing the DMA t...	3	comp.sys.ibm.pc.hardware
4	1) I have an old lasmine drive which I cann...	4	comp.sys.mac.hardware

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4
5 train_df = news_df.copy()

```

✓ VADER score:

VADER uses a sentiment lexicon with words annotated with a sentiment score ranging from -4 to 4,

```

1 from nltk.sentiment import SentimentIntensityAnalyzer
2 sia = SentimentIntensityAnalyzer()
3
4 # Add word count and character count
5 train_df['word_count'] = train_df['text'].apply(lambda x: len(str(x).split()))
6 train_df['char_count'] = train_df['text'].apply(len)
7
8 # Calculate VADER sentiment scores
9 train_df['vader_scores'] = train_df['text'].apply(
10     lambda x: sia.polarity_scores(x)['compound']
11 )
12

```

✓ Data Preprocessing

- Removes generic words (e.g., would, could, well, also, much, many, even, still, always, get, take, thing).
- Lemmatizes words to unify different forms (running, runs, ran → run).
- Removes non-essential POS tags like adverbs (RB), conjunctions (CC), determiners (DT), and auxiliary verbs (MD).
- Keeps important nouns and verbs to preserve meaning.
- Cleans text without over-filtering punctuation (keeps spaces, removes unwanted symbols).

```

1 nltk.download('all')

```

```

1 import nltk
2 import re
3 import unicodedata
4 from nltk.corpus import stopwords
5 from nltk.tokenize import word_tokenize
6 from nltk import pos_tag
7 from nltk.stem import WordNetLemmatizer
8 from nltk.corpus import wordnet
9
10 nltk.download("stopwords")
11 nltk.download("punkt")
12 nltk.download("averaged_perceptron_tagger")
13 nltk.download("wordnet")
14
15 # Initialize lemmatizer
16 lemmatizer = WordNetLemmatizer()
17
18 # Expanded stopword list

```

```

19 custom_stopwords = set(stopwords.words("english")) | {
20     "would", "could", "like", "well", "also", "way", "may",
21     "much", "many", "even", "still", "always", "get", "take", "thing", 'ax', 'one', 'good'
22 }
23
24 # POS tags to remove
25 excluded_pos_tags = {"RB", "DT", "IN", "CC", "MD", "PRP", "WP", "EX"}
26
27 def get_wordnet_pos(tag):
28     """Convert POS tag to WordNet format for lemmatization."""
29     if tag.startswith("J"):
30         return wordnet.ADJ
31     elif tag.startswith("V"):
32         return wordnet.VERB
33     elif tag.startswith("N"):
34         return wordnet.NOUN
35     elif tag.startswith("R"):
36         return wordnet.ADV
37     else:
38         return wordnet.NOUN
39
40 def clean_text(text):
41     """Cleans text by removing HTML, URLs, and unwanted symbols."""
42     text = re.sub('<[>]*>', '', text) # Remove HTML tags
43     text = re.sub('https?://[^\s]*', '', text) # Remove URLs
44     text = unicodedata.normalize("NFKD", text) # Normalize special characters
45     text = re.sub('[^\w\s]', '', text) # Remove non-word characters but keep spaces
46     return text
47
48 def preprocess_text(text):
49     """Full preprocessing pipeline for text cleaning and filtering."""
50     text = clean_text(text) # Initial text cleaning
51     text = text.lower() # Convert to lowercase
52     words = word_tokenize(text) # Tokenize text
53     tagged_words = pos_tag(words) # POS tagging
54
55     filtered_words = [
56         lemmatizer.lemmatize(word, get_wordnet_pos(tag))
57         for word, tag in tagged_words
58         if word not in custom_stopwords and # Remove stopwords
59         len(word) > 1 and # Remove single characters
60         not word.isdigit() and # Remove numbers
61         tag not in excluded_pos_tags # Remove unnecessary POS
62     ]
63
64     return " ".join(filtered_words)
65
66 train_df["cleaned_text"] = train_df["text"].apply(preprocess_text)
67 train_df["word_count"] = train_df["cleaned_text"].str.split().str.len()
68
69 print("Number of empty sentences:", (train_df["word_count"] == 0).sum())
70 print("Average words per sentence:", train_df["word_count"].mean())
71 print("\nSample of very short cleaned texts:")
72 print(train_df[train_df["word_count"] < 3][["text", "cleaned_text"]].head())
73

```

```

Number of empty sentences: 1
Average words per sentence: 85.65664461815996

```

```

Sample of very short cleaned texts:
                                     text \
99          Just opened up the distribution.\n
174          please subscribe me.
200  \nJesus did and so do I.\n\nPeace be with you,
206  Inquiry by address:eri@eridan.chuvashia.su\n
335  \nsubscribe min@stella.sku.ac.kr\n\n

cleaned_text
99          open distribution
174          please subscribe
200          jesus peace
206  inquiry addresser1eridanchuvashiasu
335  subscribe minstellaskkuackr

```

```

1 # Remove rows where cleaned_text is empty or just whitespace
2 train_df = train_df[train_df['cleaned_text'].str.strip() != '']
3 train_df[train_df['cleaned_text']!='']

```

text	category	category	name	word count	char count	vader scores	cleaned text
------	----------	----------	------	------------	------------	--------------	--------------

```

1 train_df[["text", 'cleaned_text']].head(10)

```

	text	cleaned_text
0	\n\nI am sure some bashers of Pens fans are pr...	sure bashers pen fan confuse lack kind post re...
1	My brother is in the market for a high-perform...	brother market highperformance video card supp...
2	\n\n\n\n\nFinally you said what you dream abou...	say dream mediterranean new area great year ho...
3	\nThink\n\nIt's the SCSI card doing the DMA t...	think scsi card dma transfer disk scsi card dm...
4	1) I have an old Jasmine drive which I cann...	old jasmine drive use new system understanding...
5	\n\nBack in high school I worked as a lab assi...	high school work lab assistant bunch experimen...
6	\n\nAE is in Dallas...try 214/241-6060 or 214/...	ae dallastry tech support line start
7	\n[stuff deleted]\n\nOk, here's the solution t...	stuff delete ok heres solution problem move ca...
8	\n\n\nYeah, it's the second one. And I believ...	second believe price ive try look bruinsabre t...
9	\nIf a Christian means smenne who believes in	christian mean smenne believe divinity iesus

Perform EDA

```

1 def create_length_distribution_plot():
2     """Create interactive distribution plot for review lengths"""
3     fig = go.Figure()
4     fig.add_trace(go.Histogram(
5         x=train_df['word_count'],
6         name='Word Count Distribution',
7         nbinsx=100
8     ))
9     fig.update_layout(
10         title='Distribution of Review Lengths',
11         xaxis_title='Word Count',

```

```

12     yaxis_title='Frequency'
13 )
14 fig.show()
15
16 def create_sentiment_distribution_plot():
17     """Create interactive sentiment distribution plot"""
18     fig = go.Figure()
19     fig.add_trace(go.Violin(
20         y=train_df['vader_scores'],
21         box_visible=True,
22         line_color='black',
23         name='Sentiment Distribution'
24     ))
25     fig.update_layout(
26         title='Distribution of Sentiment Scores',
27         yaxis_title='VADER Sentiment Score'
28     )
29     fig.show()
30
31 def create_wordcloud():
32     """Generate and display word cloud"""
33     all_words = ' '.join(train_df['cleaned_text'])
34     wordcloud = WordCloud(width=800, height=400, background_color='white').generate(all_words)
35
36     plt.figure(figsize=(15, 8))
37     plt.imshow(wordcloud, interpolation='bilinear')
38     plt.axis('off')
39     plt.title('Word Cloud of Reviews')
40     plt.show()
41

```

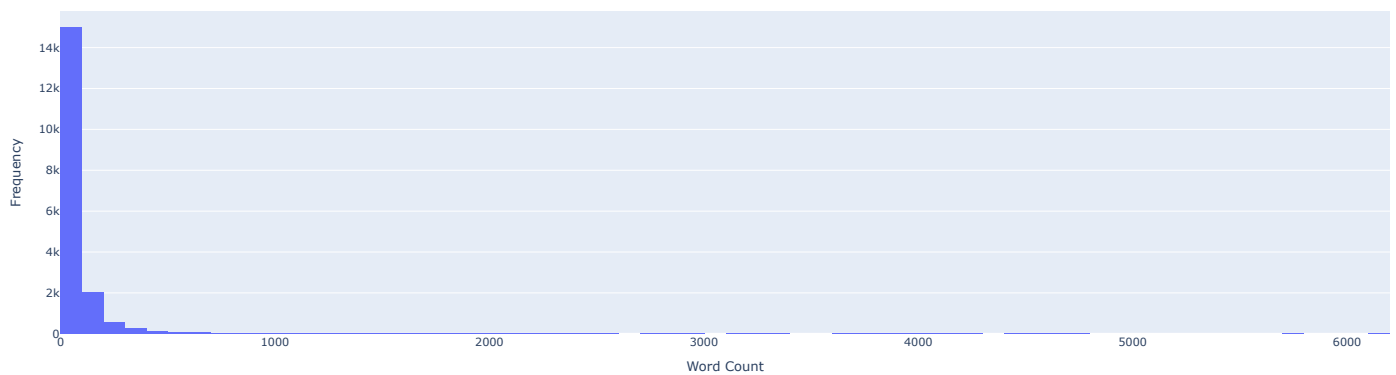
```

1 create_length_distribution_plot()
2

```



Distribution of Review Lengths



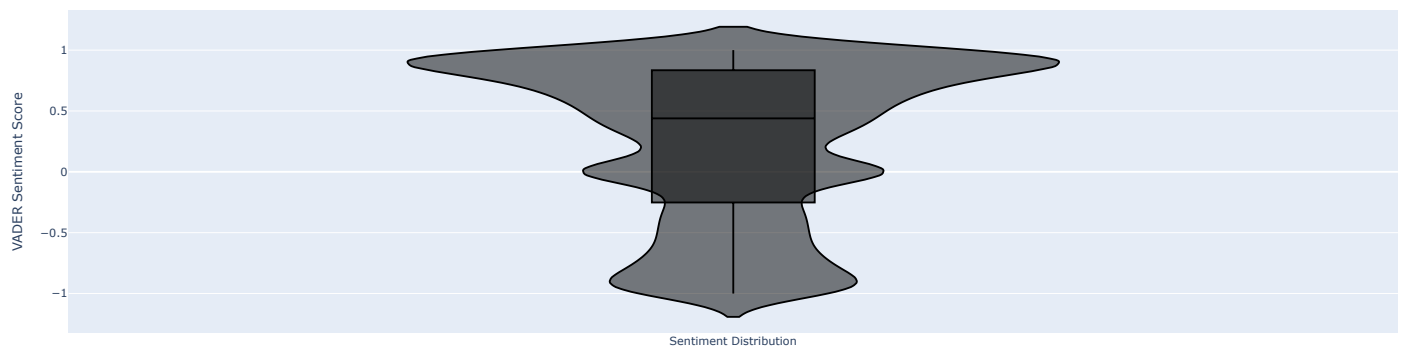
```

1 create_sentiment_distribution_plot()

```



Distribution of Sentiment Scores



```

1 create_wordcloud()

```

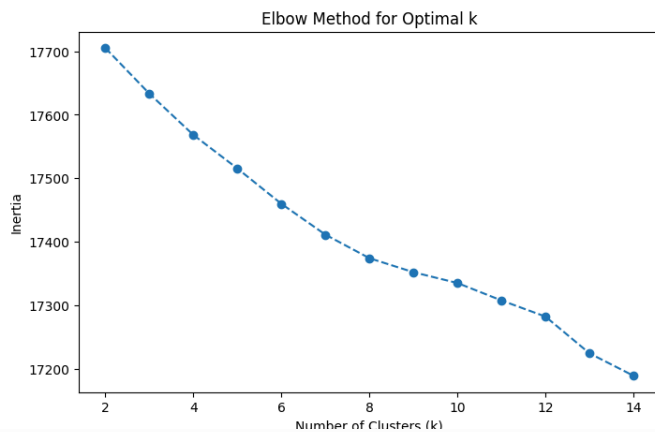


- For unsupervised sentiment analysis, we usually aim for:
 - $ARI \geq 0.3$ (decent for text data)
 - $NMI \geq 0.5$ (clusters contain meaningful info)

```

1 import matplotlib.pyplot as plt
2 from sklearn.cluster import KMeans
3
4 tfidf = TfidfVectorizer(max_features=2000)
5 tfidf_matrix = tfidf.fit_transform(train_df["cleaned_text"])
6
7 inertia = []
8 k_values = range(2, 15)
9
10 for k in k_values:
11     kmeans = KMeans(n_clusters=k, random_state=42)
12     kmeans.fit(tfidf_matrix)
13     inertia.append(kmeans.inertia_)
14
15 # Plot inertia vs. k
16 plt.figure(figsize=(8, 5))
17 plt.plot(k_values, inertia, marker="o", linestyle="--")
18 plt.xlabel("Number of Clusters (k)")
19 plt.ylabel("Inertia")
20 plt.title("Elbow Method for Optimal k")
21 plt.show()
22

```



```

1 pip install kneed

```

```

1 from kneed import KneeLocator
2
3 knee_locator = KneeLocator(k_values, inertia, curve="convex", direction="decreasing")
4 optimal_k = knee_locator.elbow
5
6 print(f"Optimal k found: {optimal_k}")
7

```



Optimal k found: 7

```

1 optimal_k = 7
2
3 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
4 train_df["knn_cluster"] = kmeans.fit_predict(tfidf_matrix)
5

```

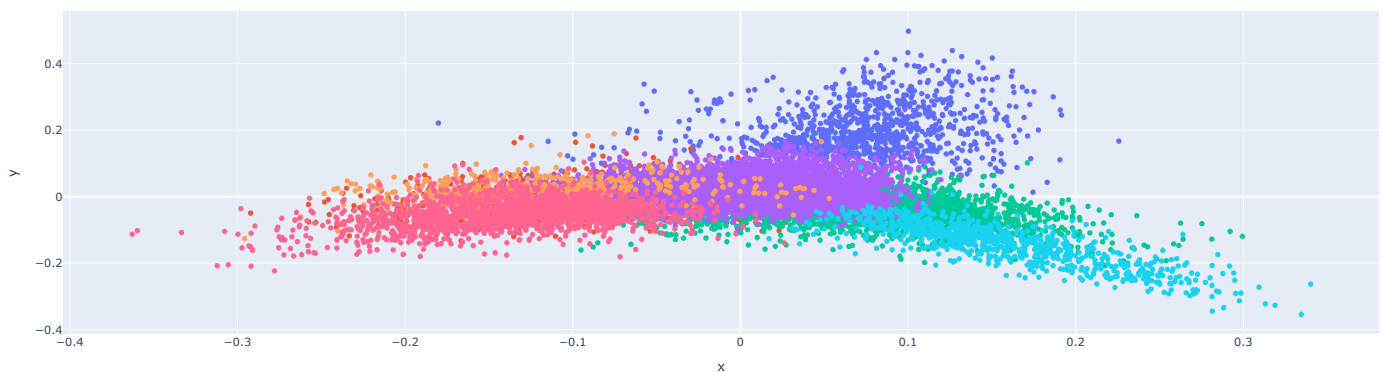
```

1 import plotly.express as px
2 from sklearn.decomposition import PCA
3
4 pca = PCA(n_components=2)
5 reduced_features = pca.fit_transform(tfidf_matrix.toarray())
6
7 fig = px.scatter(
8     x=reduced_features[:, 0],
9     y=reduced_features[:, 1],
10    color=train_df["knn_cluster"].astype(str),
11    title="Review Clusters based on Content Similarity",
12    labels={"color": "Cluster"},
13 )
14 fig.show()
15

```



Review Clusters based on Content Similarity



```

1 import numpy as np

```

```

2
3 # Get top words per cluster
4 def get_top_keywords_per_cluster(tfidf_matrix, clusters, feature_names, top_n=10):
5     cluster_centers = kmeans.cluster_centers_
6     top_words = {}
7
8     for cluster_idx in range(optimal_k):
9         top_word_indices = np.argsort(cluster_centers[cluster_idx])[-top_n:]
10        top_words[cluster_idx] = [feature_names[i] for i in top_word_indices]
11
12    return top_words
13
14 feature_names = tfidf.get_feature_names_out()
15 top_words = get_top_keywords_per_cluster(tfidf_matrix, train_df["knn_cluster"], feature_names)
16
17
18 for cluster, words in top_words.items():
19     print(f"\n🟢 Cluster {cluster} Keywords: {' '.join(words)}")
20

```

```

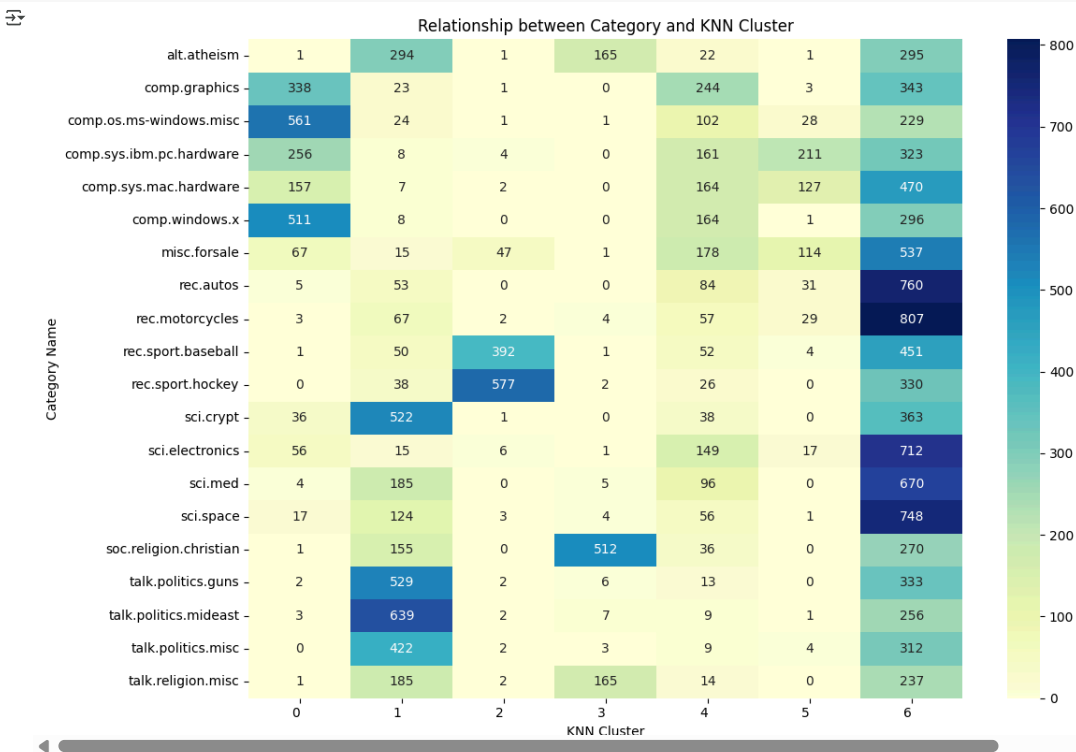
🟢 Cluster 0 Keywords: version, work, driver, problem, run, program, card, use, file, window
🟢 Cluster 1 Keywords: see, go, know, right, make, government, think, dont, say, people
🟢 Cluster 2 Keywords: last, go, season, win, hockey, year, play, player, team, game
🟢 Cluster 3 Keywords: faith, church, bible, sin, christ, say, believe, jesus, christian, god
🟢 Cluster 4 Keywords: reply, im, look, hi, advance, know, anyone, email, please, thanks
🟢 Cluster 5 Keywords: system, problem, use, ide, floppy, controller, hard, disk, scsi, drive
🟢 Cluster 6 Keywords: get, time, dont, know, new, make, car, go, think, use

```

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 cluster_category_ct = pd.crosstab(train_df['category_name'], train_df['knn_cluster'])
6
7 plt.figure(figsize=(12, 8))
8 sns.heatmap(cluster_category_ct, annot=True, cmap="YlGnBu", fmt="d")
9 plt.title("Relationship between Category and KNN Cluster")
10 plt.xlabel("KNN Cluster")
11 plt.ylabel("Category Name")
12 plt.xticks(rotation=0, ha="right")
13 plt.tight_layout()
14 plt.show()

```



Observations:

- Cluster 0 appears to group together all the hardware and electronics related articles.
- Cluster 1 is dominated by politics and religion topics.
- Cluster 2 is all about sports.
- Cluster 3 contains mostly religion related news, overlapping with cluster 1.
- Clusters 4,5 and 6 have mixed topics.
- The clusters are not evenly populated. Most news lie in cluster 6.

DBSCAN clustering with BERT embeddings

DBSCAN Clustering

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a **density-based clustering algorithm** that groups data points based on their density.

- **Finds dense regions** and groups them into clusters.
- **Does not require specifying the number of clusters** (unlike KMeans).
- **Can detect noise (outliers)** instead of forcing all points into clusters.
- Uses two key parameters:
 - ϵ (**epsilon**) → Defines how close points must be to be considered neighbors.
 - **min_samples** → Minimum points needed to form a dense region (cluster).
- Works **well for high-dimensional embeddings** (e.g., BERT).
- Can **discover natural groupings** (topics) instead of forcing a fixed number.
- Can **handle noise** (irrelevant or vague reviews).

BERT Embeddings

BERT (**Bidirectional Encoder Representations from Transformers**) is a deep learning model that creates **contextual word embeddings**.

- They are **vector representations of text**, capturing **meaning and context**.
- Unlike TF-IDF, **BERT understands synonyms & context** (e.g., "bank" as a financial institution vs. a riverbank).
- Converts reviews into **dense numerical vectors**, making clustering more meaningful.
- **Understands context** better than simple word frequency methods (TF-IDF).
- Works well with **DBSCAN** since it finds natural clusters in high-dimensional space.

Why Use DBSCAN + BERT Here?

1. **We don't know the exact number of clusters** → DBSCAN finds natural clusters without predefining them.
2. **Reviews may contain noise** → DBSCAN can ignore irrelevant reviews instead of forcing them into a category.
3. **Sentiment-based clustering is unclear** → BERT embeddings help capture review **semantics**, making clustering more meaningful.
4. **We want a more flexible model** → Unlike KMeans, DBSCAN doesn't assume clusters are spherical (it handles non-uniform shapes).

```
1 !pip install transformers sentence-transformers scikit-learn scipy nltk tensorflow_datasets
```

```
1 from sentence_transformers import SentenceTransformer
2 from sklearn.cluster import DBSCAN
3 from sklearn.metrics import silhouette_score
4 import numpy as np
5
6 model = SentenceTransformer("paraphrase-MiniLM-L6-v2")
7 train_sentences = train_df["cleaned_text"].tolist()
8 train_embeddings = model.encode(train_sentences, convert_to_numpy=True, batch_size=32, show_progress_bar=True)
```



Batches: 100% 572/572 [13:24<00:00 6.14it/s]

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.cluster import DBSCAN
4 from sklearn.neighbors import NearestNeighbors
5 import matplotlib.pyplot as plt
6 import plotly.express as px
7 from sklearn.decomposition import PCA
8 from sklearn.preprocessing import StandardScaler
9 import seaborn as sns
10 from kneed import KneeLocator
11
12 def find_optimal_epsilon(data, min_samples, n_neighbors=5):
13     """
14     Find optimal epsilon using k-nearest neighbors distance plot
15
16     Args:
17         data: feature matrix
18         min_samples: min_samples parameter for DBSCAN
19         n_neighbors: number of neighbors to consider (default: 5)
20
21     Returns:
22         optimal epsilon value
23     """
24     # Calculate distances to nearest neighbors
25     neigh = NearestNeighbors(n_neighbors=n_neighbors)
26     neigh.fit(data)
27     distances, _ = neigh.kneighbors(data)
28
29     # Sort distances in ascending order
30     distances = np.sort(distances[:, -1])
31
32     # Plot k-distance graph
33     plt.figure(figsize=(10, 6))
34     plt.plot(range(len(distances)), distances)
35     plt.xlabel('Points sorted by distance')
36     plt.ylabel(f'Distance to {n_neighbors}th nearest neighbor')
37     plt.title('K-distance Graph')
38     plt.grid(True)
39     plt.show()
40
41     # Find elbow point using KneeLocator
42     kneedle = KneeLocator(
43         range(len(distances)),
44         distances,
45         S=1.0,
46         curve='convex',
47         direction='increasing'
48     )
49
50     optimal_epsilon = distances[kneedle.knee]
51
52     plt.figure(figsize=(10, 6))
53     plt.plot(range(len(distances)), distances)
54     plt.axhline(y=optimal_epsilon, color='r', linestyle='--',
55               label=f'Optimal epsilon: {optimal_epsilon:.2f}')
56     plt.plot(kneedle.knee, distances[kneedle.knee], 'ro')
57     plt.xlabel('Points sorted by distance')
58     plt.ylabel(f'Distance to {n_neighbors}th nearest neighbor')
59     plt.title('K-distance Graph with Optimal Epsilon')
60     plt.legend()
61     plt.grid(True)
62     plt.show()
63
64     return optimal_epsilon
65
```

```

66 def perform_dbscan_clustering(data, epsilon, min_samples):
67     """
68     Perform DBSCAN clustering with given parameters
69
70     Args:
71         data: feature matrix
72         epsilon: epsilon parameter for DBSCAN
73         min_samples: min_samples parameter for DBSCAN
74
75     Returns:
76         cluster labels and DBSCAN model
77     """
78     dbscan = DBSCAN(eps=epsilon, min_samples=min_samples)
79     clusters = dbscan.fit_predict(data)
80     return clusters, dbscan
81
82 def visualize_clusters(data, clusters):
83     """
84     Create visualizations for clustering results
85
86     Args:
87         data: feature matrix
88         clusters: cluster labels
89     """
90     # Reduce dimensionality for visualization
91     pca = PCA(n_components=2)
92     reduced_features = pca.fit_transform(data)
93
94     # Create DataFrame for plotting
95     df_plot = pd.DataFrame({
96         'x': reduced_features[:, 0],
97         'y': reduced_features[:, 1],
98         'cluster': clusters
99     })
100
101     # Interactive scatter plot with Plotly
102     fig = px.scatter(
103         df_plot,
104         x='x',
105         y='y',
106         color='cluster',
107         title='DBSCAN Clustering Results (PCA)',
108         labels={'x': 'First Principal Component',
109                'y': 'Second Principal Component'},
110         color_continuous_scale='viridis'
111     )
112     fig.show()
113
114     # Plot cluster distribution
115     plt.figure(figsize=(10, 6))
116     cluster_counts = pd.Series(clusters).value_counts().sort_index()
117     sns.barplot(x=cluster_counts.index, y=cluster_counts.values)
118     plt.title('Distribution of Clusters')
119     plt.xlabel('Cluster')
120     plt.ylabel('Number of Points')
121     plt.show()
122
123 def analyze_clusters(data, clusters):
124
125     n_clusters = len(set(clusters)) - (1 if -1 in clusters else 0)
126     n_noise = list(clusters).count(-1)
127
128     metrics = {
129         'n_clusters': n_clusters,
130         'n_noise_points': n_noise,
131         'noise_percentage': n_noise / len(clusters) * 100,
132         'cluster_sizes': pd.Series(clusters).value_counts().sort_index().to_dict()
133     }
134
135     return metrics
136
137 # Main clustering pipeline
138 def run_dbscan_analysis(data, min_samples=5):
139
140     scaler = StandardScaler()
141     scaled_data = scaler.fit_transform(data)
142
143     # Find optimal epsilon
144     print("Finding optimal epsilon...")
145     epsilon = find_optimal_epsilon(scaled_data, min_samples)
146     print(f"Optimal epsilon: {epsilon:.4f}")
147
148     # Perform clustering
149     print("\nPerforming DBSCAN clustering...")
150     clusters, dbscan = perform_dbscan_clustering(scaled_data, epsilon, min_samples)
151
152     print("\nCreating visualizations...")
153     visualize_clusters(scaled_data, clusters)
154
155     print("\nAnalyzing clustering results...")
156     metrics = analyze_clusters(scaled_data, clusters)
157
158     print("\nClustering Results:")
159     print(f"Number of clusters: {metrics['n_clusters']}")
160     print(f"Number of noise points: {metrics['n_noise_points']}")
161     print(f"Percentage of noise points: {metrics['noise_percentage']:.2f}%")
162     print("\nCluster sizes:")
163     for cluster, size in metrics['cluster_sizes'].items():
164         print(f"Cluster {cluster}: {size} points")
165
166     return clusters, metrics
167
168 clusters, metrics = run_dbscan_analysis(train_embeddings)

```

```

1 epsilon = 0.02
2 min_samples = 10
3 dbscan = DBSCAN(eps=epsilon, min_samples=min_samples)
4 train_df['DBSCAN_clusters'] = dbscan.fit_predict(train_embeddings)

```

```

1 import plotly.express as px
2 from sklearn.decomposition import PCA
3
4 pca = PCA(n_components=2)
5 reduced_embeddings = pca.fit_transform(train_embeddings)
6
7 fig = px.scatter(
8     x=reduced_embeddings[:, 0],
9     y=reduced_embeddings[:, 1],
10    color=train_df["DBSCAN_clusters"].astype(str),
11    title="DBSCAN Clustering of IMDB Reviews",
12    labels={"color": "Cluster"}
13 )

```



```

14 fig.show()
15

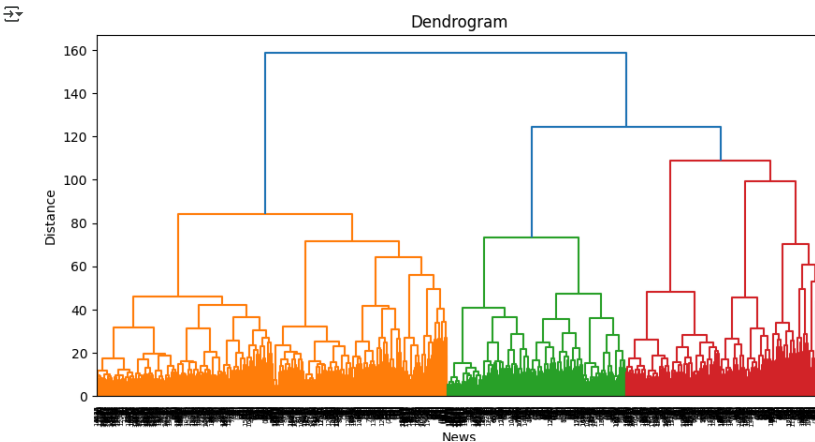
```

Hierarchical clustering

```

1 from scipy.cluster.hierarchy import linkage, dendrogram, fcluster
2 import matplotlib.pyplot as plt
3
4 linkage_matrix = linkage(train_embeddings, method="ward")
5
6 plt.figure(figsize=(10, 5))
7 dendrogram(linkage_matrix, truncate_mode="level", p=10)
8 plt.title("Dendrogram")
9 plt.xlabel("News")
10 plt.ylabel("Distance")
11 plt.show()
12

```



```

1 import numpy as np
2 import pandas as pd
3 from scipy.cluster.hierarchy import linkage, dendrogram, fcluster
4 import matplotlib.pyplot as plt
5 from kneed import KneeLocator
6 from sklearn.metrics import silhouette_score
7 import plotly.graph_objects as go
8
9 def find_optimal_clusters(embeddings, max_clusters=20):
10     """
11     Find optimal number of clusters using multiple methods:
12     1. Elbow method using distortion
13     2. Silhouette analysis
14     3. Dendrogram analysis
15     """
16
17     # Calculate distortion scores (within-cluster sum of squares)
18     distortions = []
19     silhouette_scores = []
20     n_clusters_range = range(2, max_clusters + 1)
21
22     for n_clusters in n_clusters_range:
23         cluster_labels = fcluster(linkage_matrix, n_clusters, criterion='maxclust')
24
25         # Calculate distortion
26         distortion = 0
27         for i in range(1, n_clusters + 1):
28             cluster_points = embeddings[cluster_labels == i]
29             centroid = np.mean(cluster_points, axis=0)
30             distortion += np.sum((cluster_points - centroid) ** 2)
31         distortions.append(distortion)
32
33         # Calculate silhouette score
34         if len(np.unique(cluster_labels)) > 1: # Silhouette requires at least 2 clusters
35             silhouette_scores.append(silhouette_score(embeddings, cluster_labels))
36         else:
37             silhouette_scores.append(0)
38
39     fig = go.Figure()
40     fig.add_trace(go.Scatter(x=list(n_clusters_range),
41                             y=distortions,
42                             mode='lines+markers',
43                             name='Distortion')))
44
45     kn = KneeLocator(list(n_clusters_range),
46                     distortions,
47                     curve='convex',
48                     direction='decreasing')
49     elbow_point = kn.knee
50
51     fig.add_trace(go.Scatter(x=[elbow_point],
52                             y=[distortions[elbow_point-2]],
53                             mode='markers',
54                             marker=dict(size=15, color='red'),
55                             name=f'Elbow Point (k={elbow_point})))
56
57     fig.update_layout(title='Elbow Method for Optimal k',
58                       xaxis_title='Number of Clusters (k)',
59                       yaxis_title='Distortion Score',
60                       showLegend=True)
61     fig.show()
62
63     fig = go.Figure()
64     fig.add_trace(go.Scatter(x=list(n_clusters_range),
65                             y=silhouette_scores,
66                             mode='lines+markers',
67                             name='Silhouette Score')))
68
69     # Find optimal k using silhouette score
70     optimal_k_silhouette = n_clusters_range[np.argmax(silhouette_scores)]
71
72     fig.add_trace(go.Scatter(x=[optimal_k_silhouette],

```

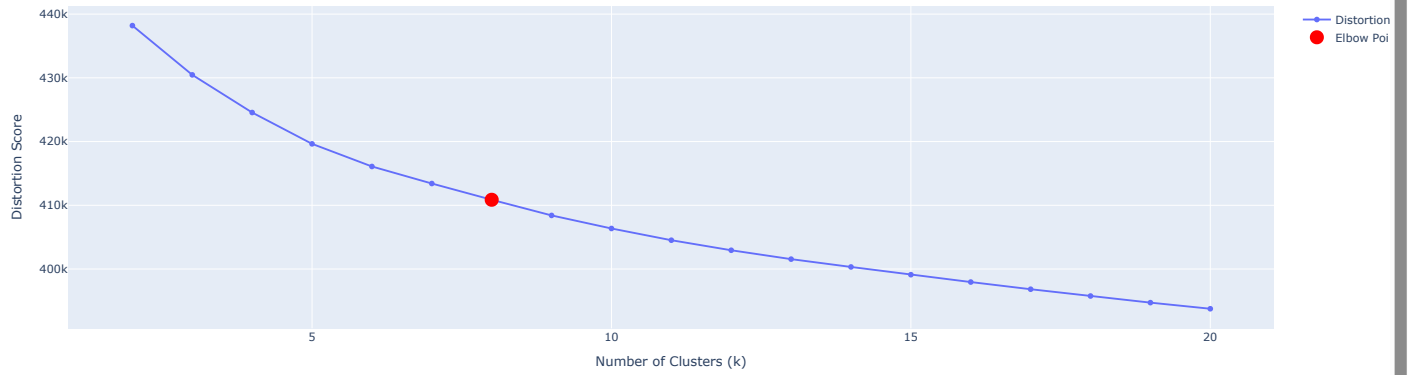
```

72         y=[max(silhouette_scores)],
73         mode='markers',
74         marker=dict(size=15, color='red'),
75         name=f'Optimal k={optimal_k_silhouette}'))
76
77 fig.update_layout(title='Silhouette Scores for Different k',
78                   xaxis_title='Number of Clusters (k)',
79                   yaxis_title='Silhouette Score',
80                   showLegend=True)
81 fig.show()
82
83 # Calculate inconsistency coefficients
84 def get_inconsistency(linkage_matrix):
85     n = len(linkage_matrix) + 1
86     heights = linkage_matrix[:, 2]
87     mean_heights = np.mean(heights)
88     std_heights = np.std(heights)
89     inconsistency = (heights - mean_heights) / std_heights if std_heights > 0 else heights
90     return inconsistency
91
92 inconsistency = get_inconsistency(linkage_matrix)
93
94 plt.figure(figsize=(12, 6))
95 plt.plot(range(len(inconsistency)), sorted(inconsistency, reverse=True))
96 plt.title('Sorted Inconsistency Coefficients')
97 plt.xlabel('Merge Index')
98 plt.ylabel('Inconsistency Coefficient')
99 plt.show()
100
101 print("\nOptimal Cluster Analysis Results:")
102 print(f"1. Elbow Method suggests {elbow_point} clusters")
103 print(f"2. Silhouette Analysis suggests {optimal_k_silhouette} clusters")
104 print(f"3. Dendrogram Analysis: Check the dendrogram plot for major splits")
105
106 return {
107     'elbow_k': elbow_point,
108     'silhouette_k': optimal_k_silhouette,
109     'distortions': distortions,
110     'silhouette_scores': silhouette_scores,
111     'linkage_matrix': linkage_matrix
112 }
113
114 def apply_clustering(embeddings, n_clusters, linkage_matrix=None):
115     """
116     Apply hierarchical clustering with the optimal number of clusters
117     """
118     if linkage_matrix is None:
119         linkage_matrix = linkage(embeddings, method="ward")
120
121     cluster_labels = fcluster(linkage_matrix, n_clusters, criterion='maxclust')
122     cluster_sizes = pd.Series(cluster_labels).value_counts().sort_index()
123
124     print("\nCluster Statistics:")
125     print(f"Number of clusters: {n_clusters}")
126     print("\nCluster sizes:")
127     for cluster, size in cluster_sizes.items():
128         print(f"Cluster {cluster}: {size} samples")
129
130     return cluster_labels
131
132
133 results = find_optimal_clusters(train_embeddings, max_clusters=20)
134 n_clusters = results['elbow_k'] # or results['silhouette_k']
135

```



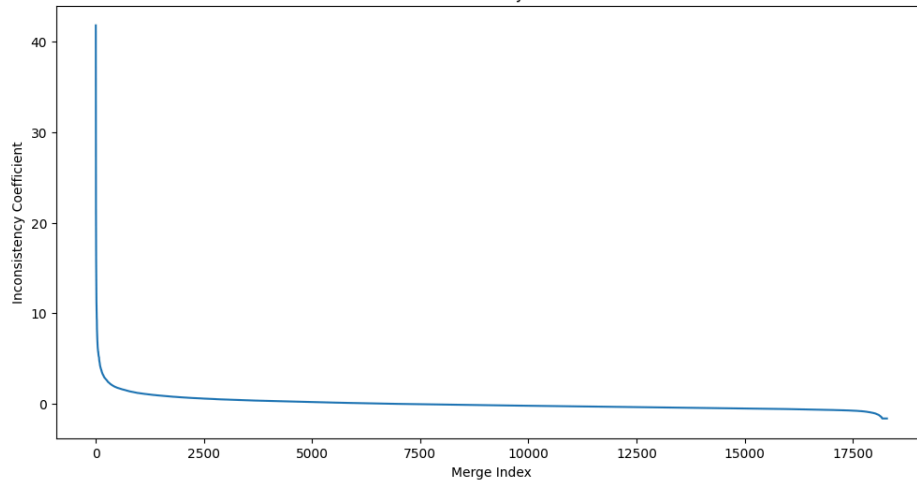
Elbow Method for Optimal k



Silhouette Scores for Different k



Sorted Inconsistency Coefficients

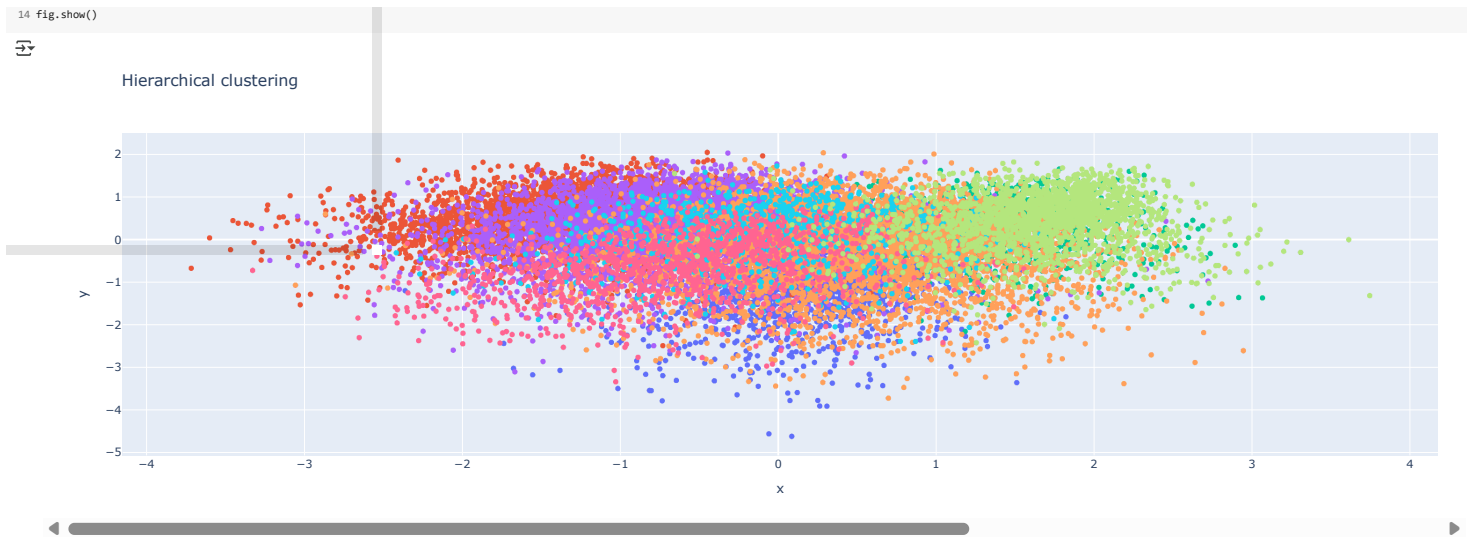


Optimal Cluster Analysis Results:

1. Elbow Method suggests 8 clusters
2. Silhouette Analysis suggests 2 clusters
3. Dendrogram Analysis: Check the dendrogram plot for major splits

```
1
2 num_clusters = 8
3 clusters = fcluster(linkage_matrix, num_clusters, criterion="maxclust")
4
5 train_df["h_cluster"] = clusters
6
```

```
1 import plotly.express as px
2 from sklearn.decomposition import PCA
3
4 pca = PCA(n_components=2)
5 reduced_embeddings = pca.fit_transform(train_embeddings)
6
7 fig = px.scatter(
8     x=reduced_embeddings[:, 0],
9     y=reduced_embeddings[:, 1],
10    color=train_df["h_cluster"].astype(str),
11    title="Hierarchical clustering",
12    labels={"color": "Cluster"}
13 )
```



```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 cluster_category_ct = pd.crosstab(train_df['category_name'], train_df['h_cluster'])
6
7 plt.figure(figsize=(12, 8))
8 sns.heatmap(cluster_category_ct, annot=True, cmap="YlGnBu", fmt="d")
9 plt.title("Relationship between Category and Hierarchical Cluster")
10 plt.xlabel("Hierarchical Cluster")
11 plt.ylabel("Category Name")
12 plt.xticks(rotation=0, ha='right')
13 plt.tight_layout()
14 plt.show()
```

