

```
1 #Setup
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import re
7 from pathlib import Path
```

```
1 # Set style
2 sns.set_style('whitegrid')
3 plt.rcParams['figure.figsize'] = (12, 6)
```

```
1 #Load Data
2 df = pd.read_csv('texts.csv')
3 print(f"Loaded {len(df)} texts")
4 df.head()
```

Loaded 60 texts

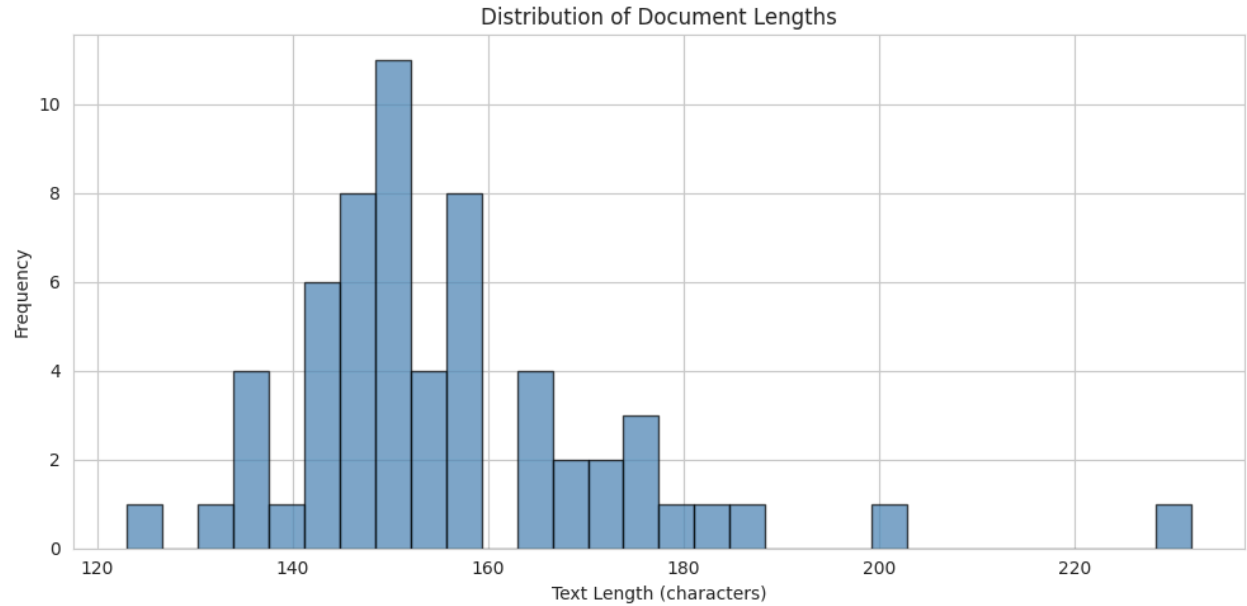
	text_id		title	text	
0	T0001	APOE ε4 allele and Alzheimer disease risk	The rs429358 variant in the APOE gene is stron...		
1	T0002	TREM2 variants in neurodegeneration	Rare coding variants in TREM2, including rs759...		
2	T0003	BIN1 locus and tau pathology	Common variants near BIN1 at rs744373 are asso...		
3	T0004	CLU genetic associations in AD	The clusterin gene CLU contains several varian...		
4	T0005	CR1 complement receptor variants	Variants in the complement receptor 1 gene CR1...		

Next steps: [Generate code with df](#) [New interactive sheet](#)

```
1 #1. Data Exploration
2
3 #Basic Statistics
4 print(f"Total documents: {len(df)}")
5 print(f"Average text length: {df['text'].str.len().mean():.0f} characters")
6 print(f"Unique text IDs: {df['text_id'].nunique()}")
```

Total documents: 60
Average text length: 156 characters
Unique text IDs: 60

```
1 # Text length distribution
2 plt.figure(figsize=(10, 5))
3 plt.hist(df['text'].str.len(), bins=30, edgecolor='black', alpha=0.7, color='steelblue')
4 plt.xlabel('Text Length (characters)')
5 plt.ylabel('Frequency')
6 plt.title('Distribution of Document Lengths')
7 plt.tight_layout()
8 plt.show()
```



```
1 # ## 2. Quick Entity Extraction Test
2
3 #Test on Sample
4 sample_text = df.iloc[0]['text']
5 print("Sample text:")
6 print(sample_text)
7 print("\n" + "="*60 + "\n")
```

Sample text:
The rs429358 variant in the APOE gene is strongly associated with increased Alzheimer disease risk across multiple populations. This v

```

1 # Extract variants
2 variants = re.findall(r'\b[rs\d+\b]', sample_text)
3 print(f"Variants found: {variants}")
4
5 # Extract gene-like tokens
6 genes = re.findall(r'\b[A-Z][A-Z0-9]{2,5}\b', sample_text)
7 genes = [g for g in genes if g not in ['THE', 'AND', 'FOR', 'WITH', 'RNA', 'DNA', 'USA']]
8 print(f"Potential genes: {genes}")

```

Variants found: ['rs429358']
 Potential genes: ['APOE']

```

1 # Look for disease terms
2 disease_keywords = ['alzheimer', 'dementia', 'disease', 'disorder', 'cognitive', 'parkinson']
3 found_diseases = [kw for kw in disease_keywords if kw.lower() in sample_text.lower()]
4 print(f"Disease keywords: {found_diseases}")

```

Disease keywords: ['alzheimer', 'disease']

```

1 # ## 3. Build Entity Extractor
2
3 #Entity Extractor Class
4 class GenomicEntityExtractor:
5     def __init__(self):
6         self.variant_pattern = re.compile(r'\b[rs\d+\b]')
7         self.gene_pattern = re.compile(r'\b[A-Z][A-Z0-9]{2,5}\b')
8
9         self.relations = [
10             'associated with', 'linked to', 'increases risk',
11             'affects', 'modulates', 'influences', 'causes',
12             'contributes to', 'related to', 'implicated in'
13         ]
14
15         self.disease_terms = [
16             'alzheimer', 'dementia', 'cognitive', 'parkinson',
17             'disease', 'disorder', 'syndrome', 'impairment'
18         ]
19
20         self.gene_stopwords = {'THE', 'AND', 'FOR', 'WITH', 'RNA', 'DNA', 'USA', 'BUT', 'NOT'}
21
22     def extract_variants(self, text):
23         return list(set(self.variant_pattern.findall(text)))
24
25     def extract_genes(self, text):
26         candidates = self.gene_pattern.findall(text)
27         return list(set([g for g in candidates if g not in self.gene_stopwords]))
28
29     def extract_diseases(self, text):
30         diseases = []
31         text_lower = text.lower()
32
33         # Find sentences with disease terms
34         sentences = text.split('.')
35         for sent in sentences:
36             sent_lower = sent.lower()
37             for term in self.disease_terms:
38                 if term in sent_lower:
39                     # Extract capitalized phrases
40                     words = sent.strip().split()
41                     for i, word in enumerate(words):
42                         if term in word.lower() and len(word) > 2:
43                             # Get context window
44                             start = max(0, i-2)
45                             end = min(len(words), i+3)
46                             phrase = ' '.join(words[start:end])
47                             if len(phrase) > 5:
48                                 diseases.append(phrase.strip())
49
50         return list(set(diseases))
51
52     def extract_relations(self, text):
53         found = []
54         text_lower = text.lower()
55         for rel in self.relations:
56             if rel in text_lower:
57                 found.append(rel)
58         return found
59
60     def extract_all(self, text, text_id):
61         variants = self.extract_variants(text)
62         genes = self.extract_genes(text)
63         diseases = self.extract_diseases(text)
64         relations = self.extract_relations(text)
65
66         records = []
67
68         # Create combinations

```

```

69         if variants and genes and diseases and relations:
70             for variant in variants:
71                 for gene in genes:
72                     for disease in diseases:
73                         for relation in relations:
74                             # Find evidence
75                             sentences = text.split('.')
76                             evidence = ""
77                             for sent in sentences:
78                                 if variant in sent and gene in sent:
79                                     evidence = sent.strip()
80                                     break
81
82                             records.append({
83                                 'text_id': text_id,
84                                 'variant': variant,
85                                 'gene': gene,
86                                 'phenotype': disease,
87                                 'relation': relation,
88                                 'evidence_span': evidence[:200]
89                             })
90     else:
91         # Create partial records
92         if variants or genes or diseases:
93             records.append({
94                 'text_id': text_id,
95                 'variant': variants[0] if variants else None,
96                 'gene': genes[0] if genes else None,
97                 'phenotype': diseases[0] if diseases else None,
98                 'relation': relations[0] if relations else None,
99                 'evidence_span': text[:200]
100             })
101
102     return records

```

```

1 # ## 4. Extract Entities from All Texts
2
3 # Full Extraction
4 extractor = GenomicEntityExtractor()
5
6 all_entities = []
7 for idx, row in df.iterrows():
8     entities = extractor.extract_all(row['text'], row['text_id'])
9     all_entities.extend(entities)
10
11 entities_df = pd.DataFrame(all_entities)
12 print(f"\nExtracted {len(entities_df)} entity records")
13 print(f"From {len(df)} documents")
14
15 # Remove None values
16 entities_df = entities_df.dropna(subset=['variant', 'gene'])
17
18 print(f"After filtering: {len(entities_df)} records")
19 entities_df.head(10)

```

Extracted 97 entity records
From 60 documents
After filtering: 69 records

	text_id	variant	gene	phenotype	relation	evidence_span	
0	T0001	rs429358	APOE	variant increases disease risk by	associated with	The rs429358 variant in the APOE gene is stron...	
1	T0001	rs429358	APOE	with increased Alzheimer disease risk	associated with	The rs429358 variant in the APOE gene is stron...	
2	T0001	rs429358	APOE	increased Alzheimer disease risk across	associated with	The rs429358 variant in the APOE gene is stron...	
3	T0002	rs75932628	TREM2	late-onset Alzheimer disease	None	Rare coding variants in TREM2, including rs759...	
4	T0003	rs744373	BIN1	with Alzheimer disease risk	associated with	Common variants near BIN1 at rs744373 are asso...	
5	T0003	rs744373	BIN1	associated with Alzheimer disease risk	associated with	Common variants near BIN1 at rs744373 are asso...	
6	T0004	rs11136000	CLU	to Alzheimer disease susceptibility	linked to	The clusterin gene CLU contains several varian...	
7	T0004	rs11136000	CLU	to Alzheimer disease susceptibility	influences	The clusterin gene CLU contains several varian...	
8	T0004	rs11136000	CLU	linked to Alzheimer disease susceptibility	linked to	The clusterin gene CLU contains several varian...	
9	T0004	rs11136000	CLU	linked to Alzheimer disease susceptibility	influences	The clusterin gene CLU contains several varian...	

Next steps: [Generate code with entities_df](#) [New interactive sheet](#)

```

1 # ## 5. Entity Analysis & Visualization
2
3 # Entity Statistics
4 print("="*60)
5 print("ENTITY STATISTICS")
6 print("="*60)
7
8 print("\nTop 10 Variants:")
9 print(entities_df['variant'].value_counts().head(10))
10
11 print("\nTop 10 Genes:")

```

```

12 print(entities_df['gene'].value_counts().head(10))
13
14 print("\nRelation Types:")
15 print(entities_df['relation'].value_counts())
16
17 print("\nTop Phenotypes:")
18 print(entities_df['phenotype'].value_counts().head(10))
19

```

ENTITY STATISTICS

Top 10 Variants:

```

variant
rs3851179      12
rs11136000     4
rs3865444      4
rs17125944     4
rs4420638      4
rs429358       3
rs11767557     3
rs242557       3
rs28834970     3
rs74615166     2
Name: count, dtype: int64

```

Top 10 Genes:

```

gene
APP          7
APOE         6
PICALM       6
CLU          4
CD33         4
FERMT2       4
EPHA1        3
PTK2B        3
MAPT         3
MS4A         2
Name: count, dtype: int64

```

Relation Types:

```

relation
associated with    21
linked to         17
affects           7
influences        6
modulates         5
implicated in     2
contributes to    2
Name: count, dtype: int64

```

Top Phenotypes:

```

phenotype
reduced Alzheimer disease risk          6
with reduced Alzheimer disease risk      6
associated with Alzheimer disease risk    5
with Alzheimer disease risk through      3
linked to Alzheimer disease susceptibility 2
with Alzheimer disease risk              2
contributes to Alzheimer disease susceptibility 2
to Alzheimer disease susceptibility through 2
linked to Alzheimer disease and           2
to Alzheimer disease susceptibility       2
Name: count, dtype: int64

```

```

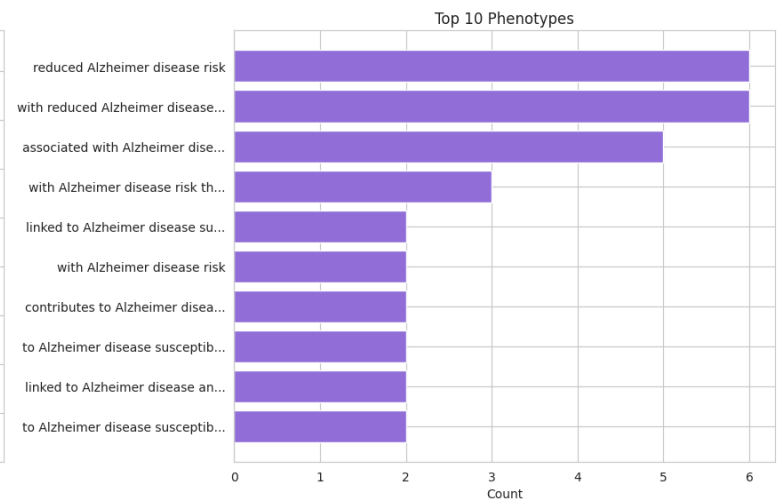
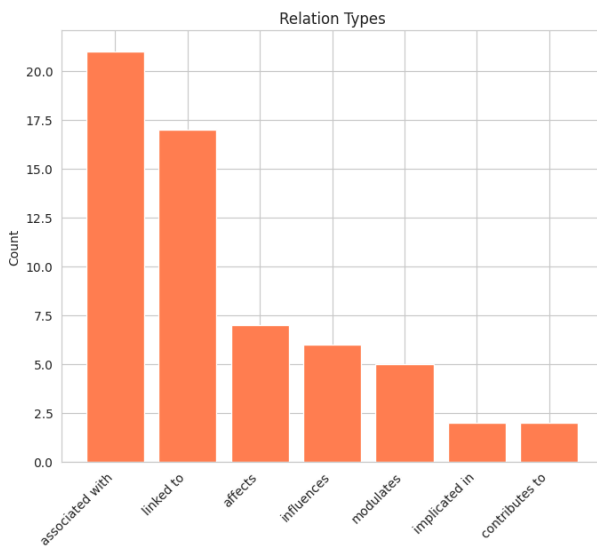
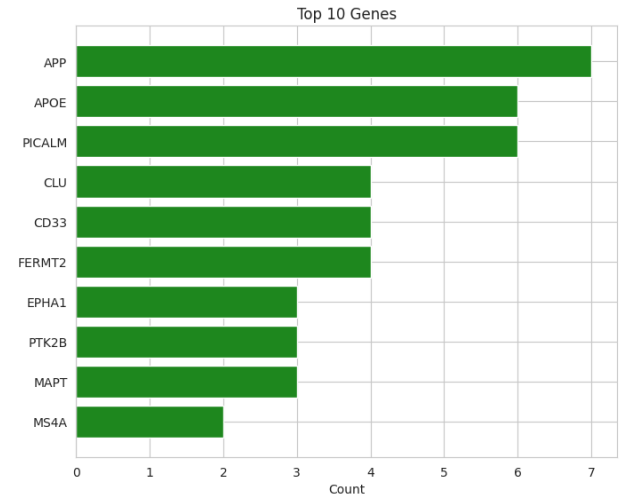
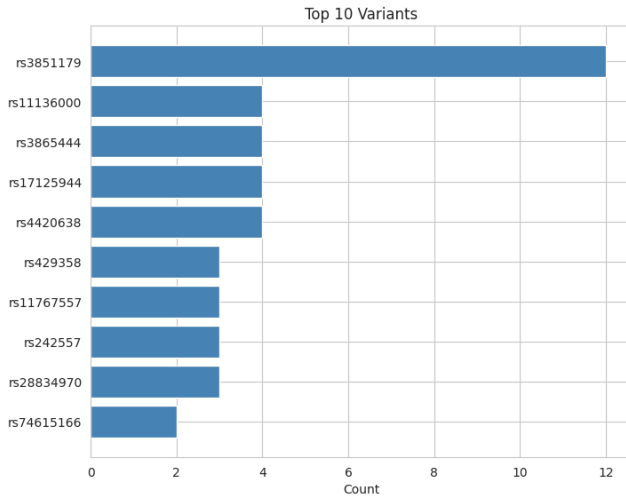
1 #Visualize Entities
2 fig, axes = plt.subplots(2, 2, figsize=(16, 12))
3 # Top variants
4 top_variants = entities_df['variant'].value_counts().head(10)
5 axes[0, 0].barh(range(len(top_variants)), top_variants.values, color='steelblue')
6 axes[0, 0].set_yticks(range(len(top_variants)))
7 axes[0, 0].set_yticklabels(top_variants.index)
8 axes[0, 0].set_xlabel('Count')
9 axes[0, 0].set_title('Top 10 Variants')
10 axes[0, 0].invert_yaxis()
11
12 # Top genes
13 top_genes = entities_df['gene'].value_counts().head(10)
14 axes[0, 1].barh(range(len(top_genes)), top_genes.values, color='forestgreen')
15 axes[0, 1].set_yticks(range(len(top_genes)))
16 axes[0, 1].set_yticklabels(top_genes.index)
17 axes[0, 1].set_xlabel('Count')
18 axes[0, 1].set_title('Top 10 Genes')
19 axes[0, 1].invert_yaxis()
20
21 # Relations
22 relation_counts = entities_df['relation'].value_counts()
23 axes[1, 0].bar(range(len(relation_counts)), relation_counts.values, color='coral')
24 axes[1, 0].set_xticks(range(len(relation_counts)))
25 axes[1, 0].set_xticklabels(relation_counts.index, rotation=45, ha='right')
26 axes[1, 0].set_ylabel('Count')
27 axes[1, 0].set_title('Relation Types')
28
29 # Top phenotypes
30 top_phenotypes = entities_df['phenotype'].value_counts().head(10)
31 axes[1, 1].barh(range(len(top_phenotypes)), top_phenotypes.values, color='mediumpurple')

```

```

32 axes[1, 1].set_yticks(range(len(top_phenotypes)))
33 axes[1, 1].set_yticklabels([p[:30] + '...' if len(p) > 30 else p for p in top_phenotypes.index])
34 axes[1, 1].set_xlabel('Count')
35 axes[1, 1].set_title('Top 10 Phenotypes')
36 axes[1, 1].invert_yaxis()
37
38 plt.tight_layout()
39 plt.show()

```



```

1 # ## 6. Topic Modeling - Find Optimal k
2
3 # Elbow Method
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.cluster import KMeans

```

```

1 # Create TF-IDF features
2 vectorizer = TfidfVectorizer(
3     max_features=200,
4     stop_words='english',
5     ngram_range=(1, 2),
6     min_df=2
7 )
8 X = vectorizer.fit_transform(df['text'])
9
10 print(f"TF-IDF matrix shape: {X.shape}")

```

TF-IDF matrix shape: (60, 200)

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```

1 # Try different k values
2 inertias = []

```

```

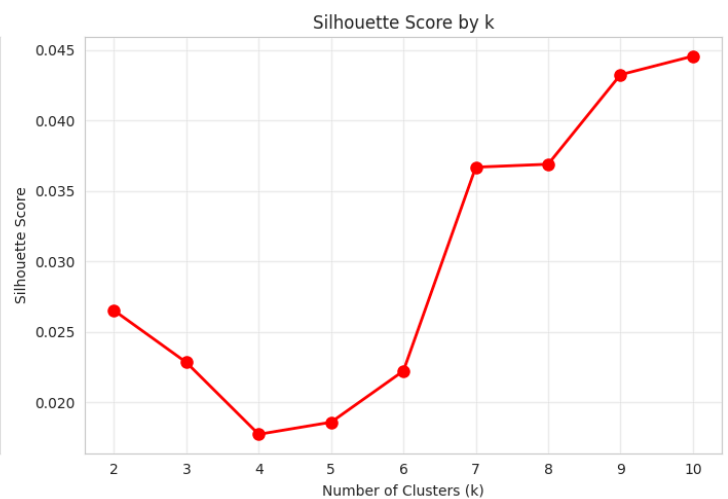
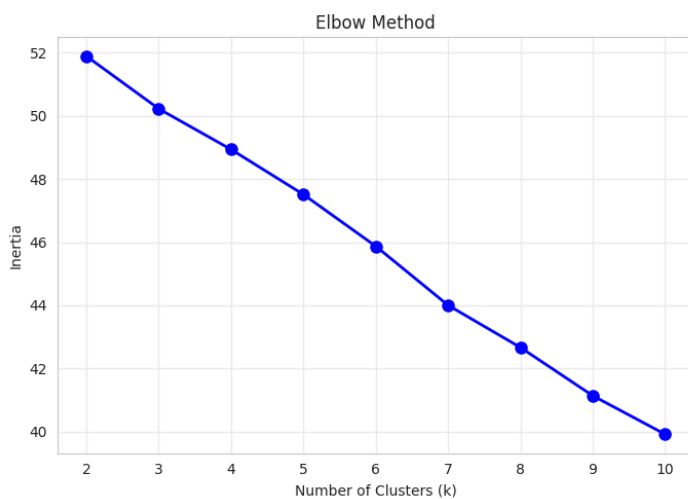
3 silhouette_scores = []
4 K_range = range(2, 11)
5
6 from sklearn.metrics import silhouette_score
7
8 for k in K_range:
9     kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
10    labels = kmeans.fit_predict(X)
11    inertias.append(kmeans.inertia_)
12
13    # Silhouette score (only if we have enough samples)
14    if len(set(labels)) > 1:
15        score = silhouette_score(X, labels)
16        silhouette_scores.append(score)
17    else:
18        silhouette_scores.append(0)

```

```

1 # Plot
2 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(14, 5))
3
4 ax1.plot(K_range, inertias, 'bo-', linewidth=2, markersize=8)
5 ax1.set_xlabel('Number of Clusters (k)')
6 ax1.set_ylabel('Inertia')
7 ax1.set_title('Elbow Method')
8 ax1.grid(True, alpha=0.3)
9
10 ax2.plot(K_range, silhouette_scores, 'ro-', linewidth=2, markersize=8)
11 ax2.set_xlabel('Number of Clusters (k)')
12 ax2.set_ylabel('Silhouette Score')
13 ax2.set_title('Silhouette Score by k')
14 ax2.grid(True, alpha=0.3)
15
16 plt.tight_layout()
17 plt.show()
18
19 print(f"Best k by silhouette score: {K_range[np.argmax(silhouette_scores)]}")

```



Best k by silhouette score: 10

```

1 #7. Final Clustering with k=5
2
3 # Cluster Documents
4 n_clusters = 5
5 kmeans = KMeans(n_clusters=n_clusters, random_state=42, n_init=20)
6 df['topic'] = kmeans.fit_predict(X)
7
8 # Topic distribution
9 topic_dist = df['topic'].value_counts().sort_index()
10 print("\nTopic Distribution:")
11 print(topic_dist)

```

Topic Distribution:

```

topic
0      9
1     22
2      7
3      5
4     17
Name: count, dtype: int64

```

```

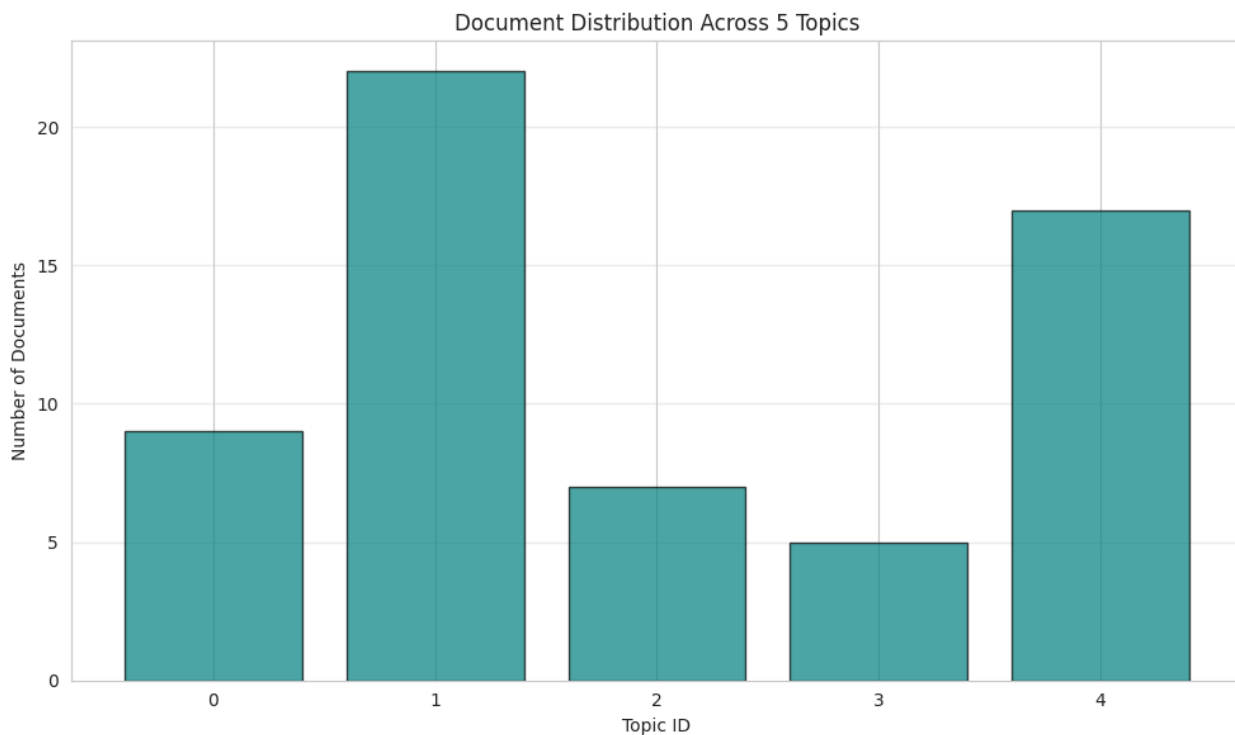
1 # Visualize distribution
2 plt.figure(figsize=(10, 6))
3 plt.bar(topic_dist.index, topic_dist.values, color='teal', alpha=0.7, edgecolor='black')
4 plt.xlabel('Topic ID')
5 plt.ylabel('Number of Documents')

```

```

6 plt.title(f'Document Distribution Across {n_clusters} Topics')
7 plt.xticks(topic_dist.index)
8 plt.grid(axis='y', alpha=0.3)
9 plt.tight_layout()
10 plt.show()

```



```

1 # %% Topic Keywords
2 feature_names = vectorizer.get_feature_names_out()
3
4 print("\n" + "="*60)
5 print("TOPIC KEYWORDS")
6 print("="*60)
7
8 for i in range(n_clusters):
9     center = kmeans.cluster_centers_[i]
10    top_indices = center.argsort()[-15:][::-1]
11    top_words = [feature_names[idx] for idx in top_indices]
12
13    print(f"\nTopic {i} (n={topic_dist[i]} documents):")
14    print(f"Keywords: {'', '.join(top_words[:10])}")
15
16    # Show 2 example texts
17    examples = df[df['topic'] == i].head(2)
18    for idx, row in examples.iterrows():
19        print(f"    Example: {row['text'][:150]}...")

```

=====

TOPIC KEYWORDS

=====

Topic 0 (n=9 documents):

Keywords: linked, linked alzheimer, synaptic, neuronal, disease, alzheimer disease, alzheimer, variants, cognitive decline, decline
 Example: The clusterin gene CLU contains several variants including rs11136000 that are linked to Alzheimer disease susceptibility.
 Example: The protein tyrosine kinase PTK2B contains variant rs28834970 linked to cognitive decline and Alzheimer disease. PTK2B regu

Topic 1 (n=22 documents):

Keywords: genetic, ad, variants, including, loci, abca7, genetic variants, populations, risk, ad risk
 Example: Loss-of-function variants in ABCA7, including rs3764650, increase AD susceptibility. ABCA7 regulates lipid homeostasis and
 Example: Large-scale genome-wide association study meta-analysis identified 29 risk loci for Alzheimer disease including ZCWPW1, ADA

Topic 2 (n=7 documents):

Keywords: contribute, cd2ap, pathogenesis, blood brain, brain barrier, barrier, blood, brain, genes, contributes
 Example: Variants in the complement receptor 1 gene CR1, particularly rs6656401, contribute to AD risk through altered complement ca
 Example: The immune-related phosphatase gene INPP5D harbors rs35349669 which modulates microglial responses and contributes to AD pa

Topic 3 (n=5 documents):

Keywords: app, beta, amyloid beta, production, beta production, amyloid, psen1, psen2, app psen1, psen1 psen2
 Example: The PICALM gene variant rs3851179 affects clathrin-mediated endocytosis and is associated with reduced Alzheimer disease ri
 Example: Variants in SORL1 such as rs11218343 influence APP trafficking and processing. SORL1 dysfunction leads to increased amyloid

Topic 4 (n=17 documents):

Keywords: risk, disease risk, alzheimer disease, alzheimer, disease, variant, associated, function, gene, multiple
 Example: The rs429358 variant in the APOE gene is strongly associated with increased Alzheimer disease risk across multiple populati
 Example: Rare coding variants in TREM2, including rs75932628, significantly affect microglial function and increase risk for late-on

```

1 # ## 8. Visualize Clusters in 2D (PCA)
2
3 #PCA Visualization
4 from sklearn.decomposition import PCA
5
6 # Reduce to 2D

```

```

7 pca = PCA(n_components=2, random_state=42)
8 X_2d = pca.fit_transform(X.toarray())
9
10 print(f"Explained variance: {pca.explained_variance_ratio_.sum():.2%}")

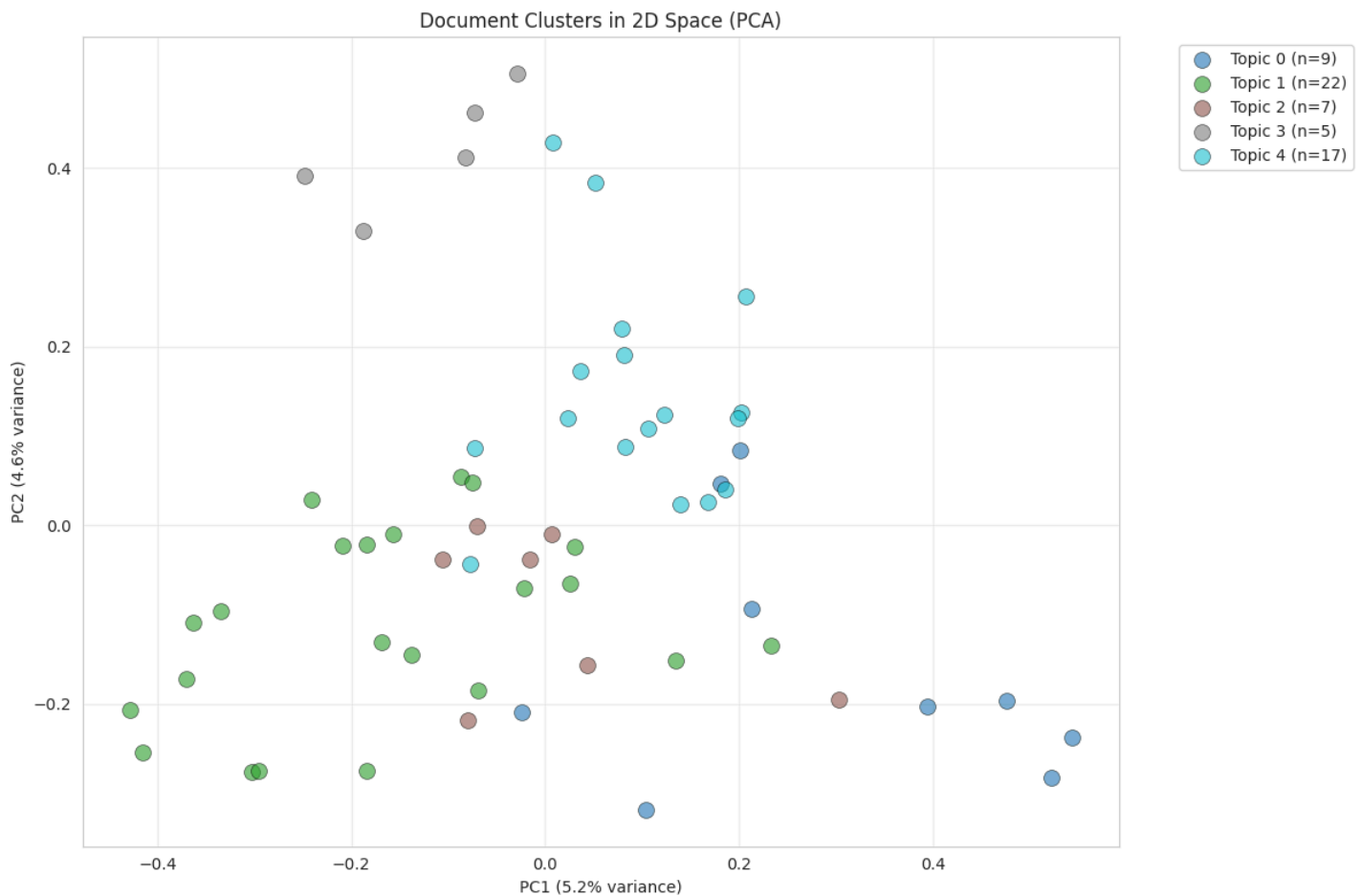
```

Explained variance: 9.88%

```

1 # Plot
2 plt.figure(figsize=(12, 8))
3 colors = plt.cm.tab10(np.linspace(0, 1, n_clusters))
4
5 for i in range(n_clusters):
6     mask = df['topic'] == i
7     plt.scatter(
8         X_2d[mask, 0],
9         X_2d[mask, 1],
10        c=[colors[i]],
11        label=f'Topic {i} (n={sum(mask)}),
12        alpha=0.6,
13        s=100,
14        edgecolors='black',
15        linewidth=0.5
16    )
17
18 plt.xlabel(f'PC1 ( {pca.explained_variance_ratio_[0]:.1%} variance)')
19 plt.ylabel(f'PC2 ( {pca.explained_variance_ratio_[1]:.1%} variance)')
20 plt.title('Document Clusters in 2D Space (PCA)')
21 plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
22 plt.grid(True, alpha=0.3)
23 plt.tight_layout()
24 plt.show()

```



```

1 # ## 9. Advanced: Semantic Embeddings with UMAP
2
3 # Sentence Transformers
4
5 try:
6     from sentence_transformers import SentenceTransformer
7     import umap
8
9     print("Creating semantic embeddings...")
10    model = SentenceTransformer('all-MiniLM-L6-v2')
11    embeddings = model.encode(df['text'].tolist(), show_progress_bar=True)
12
13    print("Reducing dimensions with UMAP...")
14    reducer = umap.UMAP(
15        n_components=2,
16        random_state=42,
17        n_neighbors=15,

```



```

18     min_dist=0.1,
19     metric='cosine'
20 )
21 embeddings_2d = reducer.fit_transform(embeddings)
22
23 # Re-cluster on embeddings
24 kmeans_emb = KMeans(n_clusters=n_clusters, random_state=42)
25 df['topic_semantic'] = kmeans_emb.fit_predict(embeddings)
26
27 # Plot
28 plt.figure(figsize=(12, 8))
29 colors = plt.cm.viridis(np.linspace(0, 1, n_clusters))
30
31 for i in range(n_clusters):
32     mask = df['topic_semantic'] == i
33     plt.scatter(
34         embeddings_2d[mask, 0],
35         embeddings_2d[mask, 1],
36         c=[colors[i]],
37         label=f'Topic {i} (n={sum(mask)}),
38         alpha=0.6,
39         s=100,
40         edgecolors='black',
41         linewidth=0.5
42     )
43
44 plt.xlabel('UMAP Dimension 1')
45 plt.ylabel('UMAP Dimension 2')
46 plt.title('Semantic Clustering with Sentence Transformers + UMAP')
47 plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
48 plt.grid(True, alpha=0.3)
49 plt.tight_layout()
50 plt.show()
51
52 print("\nSemantic clustering complete!")
53
54 except ImportError:
55     print("\nSkipping semantic embeddings (sentence-transformers not installed)")
56     print("To enable: pip install sentence-transformers umap-learn")
57

```

Creating semantic embeddings...

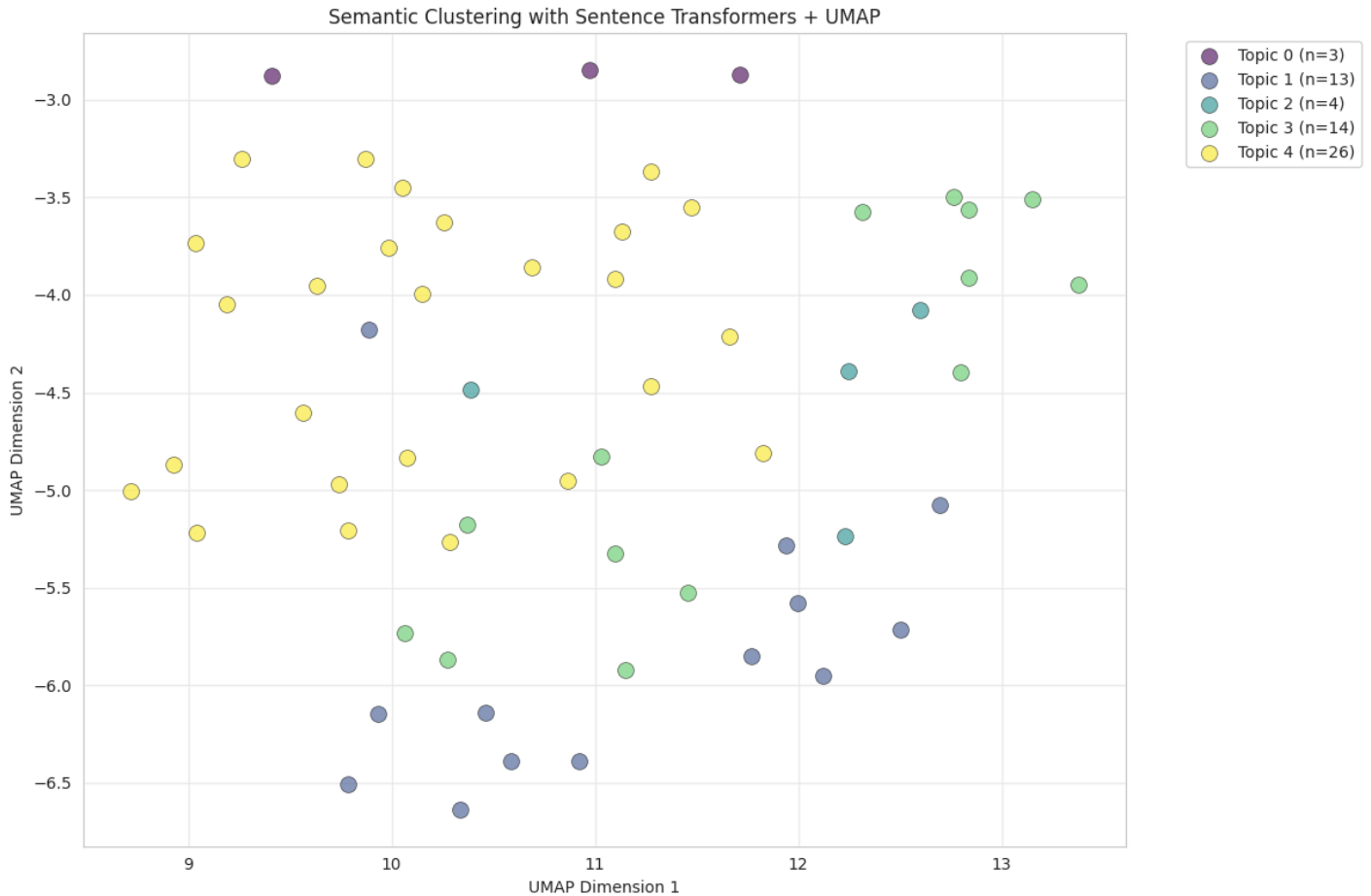
```
/usr/local/lib/python3.12/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as s
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
```

Batches: 100%

2/2 [00:02<00:00, 1.36s/it]

Reducing dimensions with UMAP...

```
/usr/local/lib/python3.12/dist-packages/umap/umap_.py:1952: UserWarning: n_jobs value 1 overridden to 1 by setting random_state. Use
warn(
```



Semantic clustering complete!

```
1 # ## 10. Save Results
2
3 # Export Data
4 # Create output directory
5 Path('../output').mkdir(exist_ok=True)
6
7 # Save entities
8 entities_df.to_csv('../output/entities.csv', index=False)
9 entities_df.to_json('../output/entities.json', orient='records', indent=2)
10
11 # Save clustered documents
12 df.to_csv('../output/texts_with_topics.csv', index=False)
13
14 # Create topic summary
15 topic_summary = []
16 for i in range(n_clusters):
17     center = kmeans.cluster_centers_[i]
18     top_indices = center.argsort()[-10:][::-1]
19     top_words = [feature_names[idx] for idx in top_indices]
20
21     topic_summary.append({
22         'topic_id': i,
23         'size': int(topic_dist[i]),
24         'keywords': top_words,
25         'example_texts': df[df['topic'] == i]['text'].head(3).tolist()
26     })
27
28 import json
29 with open('../output/topics.json', 'w') as f:
30     json.dump(topic_summary, f, indent=2)
31
32 print("\n" + "="*60)
33 print("RESULTS SAVED")
34 print("="*60)
35 print("✓ entities.csv")
```