## ⌄ 🏛 AI-Powered News Classification & Indexing

```
1 import os
2 os.environ["OPENAI_API_KEY"] = "sk-proj-oSqbiX-zZAtEAsFS0PMsLBCnK8AaTBtaQ-lOOlvNqjy9MG1v2v19C8tI8TC7CFkyI1jLP3egMvT3Blbk
```

```
1 # Install Dependencies
2 !pip install --upgrade pip
3 !pip install pandas openpyxl chromadb langchain umap-learn hdbscan tiktoken -q
```

⮂ Show hidden output

## ⌄ 1. Imports & Configuration

```
1 !pip install chromadb langchain-community -q
2 import os
3 import pandas as pd
4 from langchain.embeddings import OpenAIEmbeddings
5 from chromadb import Client
6 from chromadb.config import Settings
7 import umap
8 import hdbscan
```

## ⌄ 2. Load & Prepare CSV

```
 1 # 2.1 Load CSV
 2 df = pd.read_csv('Aggregated News Dataset - Sheet1.csv')
 3
 4 # Combine title + summary/content into one field
 5 df['text'] = (df['Title'].fillna('') + ' \n\n ' +
 6               df['news_summary'].fillna(''))
 7
 8 # 2.3 Keep only needed columns
 9 # Check if 'id' column exists, if not, use a different unique identifier or create one
10 if 'id' not in df.columns:
11   df['id'] = df.index
12
13 df = df[['id', 'Title', 'text']]
14 df
```

⮂ Show hidden output

--------------------------------------------------

Next steps:  ( Explain error )

## ⌄ 3. Define Category Prototypes

```
1 CATEGORIES = {
2     "sports": "Sports news about matches, athletes, teams and sporting events",
3     "finance": "Financial markets, stocks, economy and business news",
4     "politics": "Political news, elections, government policies",
5     "lifestyle": "Lifestyle, health, travel, food and culture",
6     "music": "Music artists, albums, concerts and industry news"
7 }
```

## ⌄ 4. Spin Up ChromaDB – Categories Collection

```
 1 import chromadb
 2 from chromadb.config import Settings
 3 from langchain.embeddings import OpenAIEmbeddings
 4 from chromadb.utils import embedding_functions
 5
 6 # 4.1 Initialize Chroma with persistant storage
 7 # PersistentClient will automatically use DuckDB with Parquet serialization
 8 chroma_client = chromadb.PersistentClient(path="./chroma_db")
 9
10 # 4.2 Create (or get) collection for categories,
11 #     using Chroma's own OpenAI wrapper
```

```
12 openai_ef = embedding_functions.OpenAIEmbeddingFunction(
13     api_key=os.getenv("OPENAI_API_KEY"),
14     model_name="text-embedding-ada-002"
15 )
16
17 cat_coll = chroma_client.get_or_create_collection(
18     name="news_categories",
19     embedding_function=openai_ef,
20     metadata={"description": "Prototype embeddings for news categories"}
21 )
22
23 # 4.3 Upsert one "document" per category — Chroma will call OpenAIEmbeddingFunction internally
24 for cat_name, cat_prompt in CATEGORIES.items():
25     cat_coll.upsert(
26         ids=[cat_name],
27         documents=[cat_prompt],
28         metadatas=[{"category": cat_name}]
29     )
30 # No need to call persist()—PersistentClient writes every upsert to disk automatically.
```

## 5. Classify Articles via Nearest-Neighbor

```
1 # 5.1 Helper to classify one text
2 !pip install tiktoken
3 import tiktoken
4 embeddings = OpenAIEmbeddings()
5 def classify(text: str):
6     emb = embeddings.embed_query(text)
7     results = cat_coll.query(query_embeddings=[emb], n_results=1)
8     return results['ids'][0][0], results['distances'][0][0]
9
10 # 5.2 Run classification
11 df[['predicted_category', 'similarity']] = df['text'].apply(
12     lambda t: pd.Series(classify(t))
13 )
14
15 # Inspect
16 df.head()
17
```

Show hidden output

Next steps:    Generate code with df      View recommended plots      New interactive sheet

## 6. Article Clustering (for RAG-friendly groups)

```
1 # 6.1 Compute all article embeddings
2 embs = embeddings.embed_documents(df['text'].tolist())
3
4 # 6.2 UMAP reduction to 5D (speeds HDBSCAN)
5 umap_model = umap.UMAP(n_components=5, random_state=42)
6 reduced = umap_model.fit_transform(embs)
7 # ------------------------
8 # STEP 6: Deduplication & Prioritization
9 # ------------------------
10 # Clustering for deduplication
11 print("🔄 Detecting duplicates...")
12 # 6.3 HDBSCAN clustering
13 clusterer = hdbscan.HDBSCAN(min_cluster_size=15, prediction_data=True)
14 clusters = clusterer.fit_predict(reduced)
15 df['cluster'] = clusters
16
```

Show hidden output

## 7. Priority Tagging & Highlight Extraction

```
1 import re
2
3 # 5.1 Cluster sizes → proxy for frequency across sources
4 df['cluster_size'] = df.groupby('cluster')['id'].transform('count')
5
6 # 5.2 Compile per-category keyword patterns once
7 PRIORITY_KEYWORDS = {
```

```
 8      "sports": ["breaking", "championship", "cup final", "olympics", "transfer",
 9                  "injury update", "record broken", "doping scandal", "last minute",
10                  "trophy", "victory", "defeat", "comeback"],
11      "finance": ["alert", "breaking", "market crash", "rate hike", "earnings report",
12                  "recession", "inflation", "fed decision", "stock plunge", "merger",
13                  "bankruptcy", "crypto crash", "dividend", "short squeeze", "ipo"],
14      "politics": ["election", "scandal", "resignation", "law passed", "protest",
15                  "breaking", "impeachment", "summit", "sanctions", "diplomatic crisis",
16                  "war", "treaty", "vote result", "corruption", "speech"],
17      "lifestyle": ["viral", "trending", "recipe", "hack", "celebrity",
18                  "wedding", "divorce", "pregnancy", "royal family", "makeover",
19                  "controversial", "banned", "recall", "study finds"],
20      "music": ["tour announced", "new album", "grammy", "controversy", "feud",
21                  "concert disaster", "chart-topping", "streaming record", "comeback",
22                  "breakup", "collaboration", "lyrics decoded", "canceled", "viral hit"]
23 }
24
25 # lower-case titles once
26 df['title_lc'] = df['Title'].str.lower()
27
28 # build regex dict
29 keyword_patterns = {
30     cat: re.compile(r'\b(' + '|'.join(map(re.escape, kws)) + r')\b')
31     for cat, kws in PRIORITY_KEYWORDS.items()
32 }
33
34 # 5.3 Flag priority articles
35 def _is_priority(row):
36     pat = keyword_patterns.get(row['predicted_category'])
37     return bool(pat.search(row['title_lc'])) if pat else False
38
39 df['is_priority'] = df.apply(_is_priority, axis=1)
40
41 # 5.4 Compute a single highlight_score
42 #    → boost keyword hits heavily, then cluster size
43 df['highlight_score'] = df['is_priority'].astype(int) * 1000 + df['cluster_size']
44
45 # 5.5 Extract top-5 highlights per category
46 highlights = (
47     df.sort_values(['predicted_category','highlight_score'], ascending=[True, False])
48       .groupby('predicted_category')
49       .head(5)
50       .reset_index(drop=True)
51 )
52
53 # 5.6 Preview or export to CSV
54 print(highlights[['predicted_category','Title','cluster_size','is_priority','highlight_score']])
55 highlights.to_csv('daily_highlights.csv', index=False)
```

⇥ Show hidden output

## ⌄ 7. Index All Articles for Retrieval

```
 1 # 7.1 Create/get an "articles" collection
 2 art_coll = chroma_client.get_or_create_collection(
 3     name="news_articles",
 4     metadata={"description": "All articles indexed for retrieval"}
 5 )
 6
 7 # 7.1.1 Ensure IDs are strings
 8 if 'id' not in df.columns:
 9     df['id'] = df.index
10 df['id'] = df['id'].astype(str)
11
12 # 7.1.2 Prepare payload
13 docs = df['text'].tolist()
14 ids  = df['id'].tolist()
15 metas = df[['Title', 'predicted_category', 'cluster']].to_dict(orient='records')
16 # assume `embs` is your precomputed list of embeddings for df['text']
17
18 # 7.2 Upsert with retry and logging
19 from tenacity import retry, stop_after_attempt, wait_exponential, retry_if_exception_type
20 import logging
21
22 logger = logging.getLogger(__name__)
23 logging.basicConfig(level=logging.INFO)
24
25 @retry(
26     reraise=True,
27     stop=stop_after_attempt(3),
```

```
28        wait=wait_exponential(multiplier=1, min=2, max=10),
29        retry=retry_if_exception_type(Exception)
30 )
31 def _upsert_all(docs, embs, ids, metas):
32        art_coll.upsert(
33            documents=docs,
34            embeddings=embs,
35            ids=ids,
36            metadatas=metas
37        )
38
39 try:
40        _upsert_all(docs, embs, ids, metas)
41        logger.info(f"Successfully upserted {len(ids)} articles to 'news_articles'.")
42 except Exception as err:
43        logger.error(f"Failed to upsert articles after retries: {err}")
44        # fallback or raise, as needed
45
46 # No explicit persist()—PersistentClient writes every upsert to disk automatically.
```

## ⌄ 8. Export Classified Articles Results

```
1 # 9.1 To CSV
2 df.to_csv('classified_articles.csv', index=False)
```