

# BE-NIG — News Insight Generator

Step-by-step project guide (A → Z) — Ready to follow

Prepared for: █████ AIML Project (2000+ samples)

## 0. Project Overview (Overview)

**Project Name:** BE-NIG — Bangla/English News Insight Generator  
**Project Goal:** Build a machine learning model to analyze news data and generate insights.  
**Dataset:** News data (Bangla/English news articles) from various sources.  
**Tasks:** • **Classification (Category)** — Politics, Sports, Tech, Business, Entertainment, Health, Education • **Summarization** (Summary) • **Similarity** (similar headlines) • **Deployment** (Streamlit app)  
**Tools/Libraries:** Python, pandas, sklearn, nltk, sumy, Streamlit, etc.

## 1. What you'll submit (What you'll submit)

**Submission Requirements:** - notebooks/01\_eda\_clean.ipynb (EDA & cleaning) - notebooks/02\_train\_eval.ipynb (training + evaluation) - app/streamlit\_app.py (Streamlit app) - models/ (vectorizer, classifier, label encoder saved) - data/news\_clean.parquet (cleaned dataset) - README.md (Project report.pdf)

## 2. Step-by-step plan (Step-by-step plan)

- Step 0 — Environment (Environment): • Python 3.8+ • Jupyter Notebook • VS Code
- Step 1 — Decide categories (Decide categories): • Categories: Politics, Sports, Business, Tech, Entertainment, Health, Education — Total 7 categories • Minimum samples per category: 250-400 → Total 2000+
- Step 2 — Collect data (Collect data): • Kaggle/news datasets • Minimum columns: id, title, text, category
- Step 3 — Data cleaning (Data cleaning): • title, text, category • clean\_text() • Remove whitespace, punctuation, digits, etc.
- Step 4 — TF-IDF (TF-IDF): • TfidfVectorizer(ngram\_range=(1,2), min\_df=2, max\_df=0.9, sublinear\_tf=True)
- Step 5 — Train (Train): • Stratified 80/20 split • Model: Linear SVM (SGDClassifier(loss='hinge')) • LogisticRegression • Metrics: Accuracy, Macro F1, Confusion matrix
- Step 6 — Summary & Keywords (Summary & Keywords): • Summary: TextRank (sumy) • Keywords: TF-IDF top terms • Similar: cosine\_similarity with corpus TF-IDF
- Step 7 — Demo (Demo): • Streamlit app: text input → category, summary, keywords, similar headlines
- Step 8 — Report (Report): • Confusion matrix, per-class metrics, few error-analysis examples, limitations, future work

## 3. Recommended folder layout (Recommended folder layout)

project\_root/ ├── data/ | ├── news\_raw.csv (raw collected data) | ├── notebooks/ | ├── app/ | ├── streamlit\_app.py | ├── utils.py | ├── models/ | ├── requirements.txt | └── README.md

## 4. Core Python code (core.py)

```
import re, string, joblib, numpy as np, pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
```

```

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import SGDClassifier
from sklearn.pipeline import Pipeline
from sklearn.metrics import classification_report, confusion_matrix
from joblib import dump
BN_PUNCT = "'\"'_,-..."
PUNCT_TABLE = str.maketrans("", "", string.punctuation + BN_PUNCT)
def clean_text(s):
    s = str(s)
    s = s.replace("\u200c", " ").replace("\u200d", " ")
    s = s.lower()
    s = re.sub(r"\d+", " ", s)
    s = s.translate(PUNCT_TABLE)
    s = re.sub(r"\s+", " ", s).strip()
    return s
df = pd.read_csv("data/news_raw.csv") # id,title,text,category
df["text_clean"] = df["title"].fillna("") + ". " + df["text"].fillna("")
df["text_clean"] = df["text_clean"].apply(clean_text)
le = LabelEncoder()
y = le.fit_transform(df["category"])
X_train, X_test, y_train, y_test = train_test_split(
    df["text_clean"], y, test_size=0.2, stratify=y, random_state=42
)
pipe = Pipeline([
    ("tfidf", TfidfVectorizer(ngram_range=(1,2), min_df=2, max_df=0.9, sublinear_tf=True)),
    ("clf", SGDClassifier(loss="hinge", alpha=1e-4, random_state=42))
])
pipe.fit(X_train, y_train)
pred = pipe.predict(X_test)
print(classification_report(y_test, pred, target_names=le.classes_))
dump(pipe, "models/news_clf.joblib")
dump(le, "models/label_encoder.joblib")

```

## 4. 新闻摘要生成 — 基于TFIDF的摘要生成 (app/utlis.py)

```

import numpy as np, pandas as pd, joblib
from sumy.summarizers.text_rank import TextRankSummarizer
from sumy.parsers.plaintext import PlaintextParser
from sumy.nlp.tokenizers import Tokenizer
pipe = joblib.load("models/news_clf.joblib")
le = joblib.load("models/label_encoder.joblib")
tfidf_only = joblib.load("models/tfidf_only.joblib")
M = joblib.load("models/corpus_sparse.joblib")
meta = pd.read_parquet("data/corpus_meta.parquet")
def predict_info(text: str, topk_sim=3):
    proba = None
    if hasattr(pipe.named_steps["clf"], "predict_proba"):
        proba = pipe.predict_proba([text])[0]
    pred = pipe.predict([text])[0]
    cat = le.inverse_transform([pred])[0]
    # summary
    try:
        parser = PlaintextParser.from_string(text, Tokenizer("english"))
        summ = TextRankSummarizer()
        sents = [str(s) for s in summ(parser.document, 2)]
        summary = " ".join(sents) if sents else text[:180]
    except:
        summary = text[:180]
    # keywords
    v = tfidf_only.transform([text])
    terms = np.array(tfidf_only.get_feature_names_out())
    top_idx = np.asarray(v.tocoo().col)
    top_data = np.asarray(v.tocoo().data)
    order = np.argsort(-top_data)[:5]

```

```

keywords = terms[top_idx[order]].tolist()
# similar docs
from sklearn.metrics.pairwise import cosine_similarity
sims = cosine_similarity(v, M).ravel()
top = sims.argsort()[-topk_sim:][::-1]
recs = meta.iloc[top].to_dict(orient="records")
return cat, summary, keywords, proba, recs

```

## □. Streamlit □□□□ (app/streamlit\_app.py) — □□□ □□□□□□

```

import streamlit as st
from utils import predict_info
st.set_page_config(page_title="News Insight Generator", layout="centered")
st.title("■ BE-NIG – News Insight Generator")
text = st.text_area("Paste a news text/headline:", height=180)
if st.button("Analyze"):
    if not text.strip():
        st.warning("Please paste some text.")
    else:
        cat, summary, keywords, proba, recs = predict_info(text)
        st.success(f"Predicted Category: **{cat}**")
        st.write("***Summary:** ", summary)
        st.write("***Keywords:** ", ", ".join(keywords) if keywords else "-")
        st.subheader("Similar news")
        for r in recs:
            st.write(f"• {r['title']} – _{r['category']}_")

```

## □. □□□□□ □□□ □□□□□□□ □ □□□□□□□□□

```

# Virtual environment
python -m venv venv
source venv/bin/activate # Windows: venv\Scripts\activate
pip install -r requirements.txt
# Run training script / notebook
python notebooks/core.py
# Run demo
streamlit run app/streamlit_app.py

```

□□□□ □□□□ □□□□ □□□□□□□□ □□□□□ (Immediate next steps): 1) □□□□□□□□□□□□□□ □□□□□□ □□□ (□□□□□ □□□□□□□□□□□□□□□ 7□□ □□□ □□□) 2) data/news\_raw.csv □□□□ □□□□ □□□: columns = id,title,text,category 3) □□□ □□□, □□□ □□□□□ □□□□ □□□ sample dataset (100-300 □□□□) □□□□ □□□ □□□□ □□□□□ 4) □□□□ '□□□ □□□□□□□ □□□' □□□□□ □□□□ training □□□ □□□□□□□ sample □□□□ □□□□ □□□ □□□□□□□ □□ □□□□□□□ □□□□ A→Z □□□□ □□□ — □□□□□ □□□ □□□□ □□□□ sample dataset □ □□□□□□ □□□□□□□ □□□ □□□ □□□□ □□□□□□□