

**** Step 1: Exploratory Data Analysis (EDA)****

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Load the dataset
6 df = pd.read_csv("/Users/md.jubayerhossain/Downloads/Bangla_news.csv")
7
8 # Basic info
9 print(df.info())
10 print(df.head())
11
12 print(df['category'].value_counts())
13
14 # Check category distribution
15 plt.figure(figsize=(8,5))
16 sns.countplot(data=df, x='category', palette='Set2')
17 plt.title("category Distribution")
18 plt.xticks(rotation=45)
19 plt.tight_layout()
20 plt.show()
21
22 # Check missing values
23 print(df.isnull().sum())
24
25 # Text length distribution
26 df['text_length'] = df['content'].astype(str).apply(lambda x: len(x.split()))
27 sns.histplot(df['text_length'], bins=40)
28 plt.title("Text Length Distribution")
29 plt.xlabel("Number of Words")
30 plt.show()
31
```



RangeIndex: 11904 entries, 0 to 11903

Data columns (total 6 columns):

#	Column	Non-Null Count	Dtype
0	title	11904 non-null	object
1	published_date	11904 non-null	object
2	reporter	10914 non-null	object
3	category	11904 non-null	object
4	url	11904 non-null	object
5	content	11904 non-null	object

dtypes: object(6)

memory usage: 558.1+ KB

None

```

title \
0 সিন নদীতে ব্যাপক দূষণ, স্থগিত করা হলো ট্রায়াল...
1 এসিসির নতুন সভাপতি হচ্ছেন মহসিন নাকভি!
2 ২০৩০ ও ২০৩৪ বিশ্বকাপের হোস্ট জানা যাবে দুইদিন পর!
3 প্যারিস অলিম্পিক: কোয়ার্টার ফাইনালে ওঠার লক্ষ্য...
4 আজ টিভিতে যা দেখবেন (৩০ জুলাই ২০২৪)

```

```

published_date reporter category url \
0 30th July, 2024 5:59 pm আরআইএম sports https://jamuna.tv/news/552127
1 30th July, 2024 5:25 pm NaN sports https://jamuna.tv/news/552123
2 30th July, 2024 1:22 pm আরআইএম sports https://jamuna.tv/news/552066
3 30th July, 2024 11:54 am এনকে sports https://jamuna.tv/news/552046
4 30th July, 2024 7:01 am এনকে sports https://jamuna.tv/news/552016

```

```

content
0 ছবি: সংগৃহীত সিন নদীতে অলিম্পিকের উদ্বোধনী অনু...
1 ২০২৫ সালের জানুয়ারিতে এশিয়ান ক্রিকেট কাউন্সিল...
2 ফিফা বিশ্বকাপ ২০৩০ ও ৩৪ সালের স্বাগতিক হবার ল...
3 প্যারিস অলিম্পিকে কোয়ার্টার ফাইনাল রেসে টিকে ...
4 আজ ভারত-শ্রীলঙ্কার তৃতীয় টি-টোয়েন্টি অলিম্পিকে...

```

category

sports 2976

international 2976

entertainment 2976

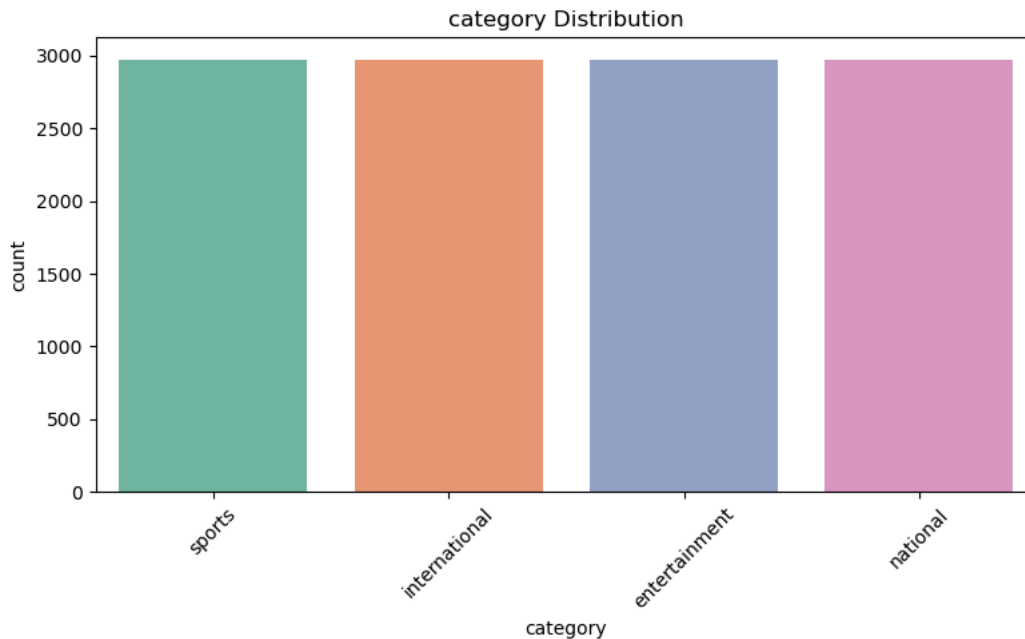
national 2976

Name: count, dtype: int64

/var/folders/nm/l7zhs8816hv5d_c6lvpbzn8m0000gn/T/ipykernel_22789/1749407327.py:16: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `leg`

sns.countplot(data=df, x='category', palette='Set2')



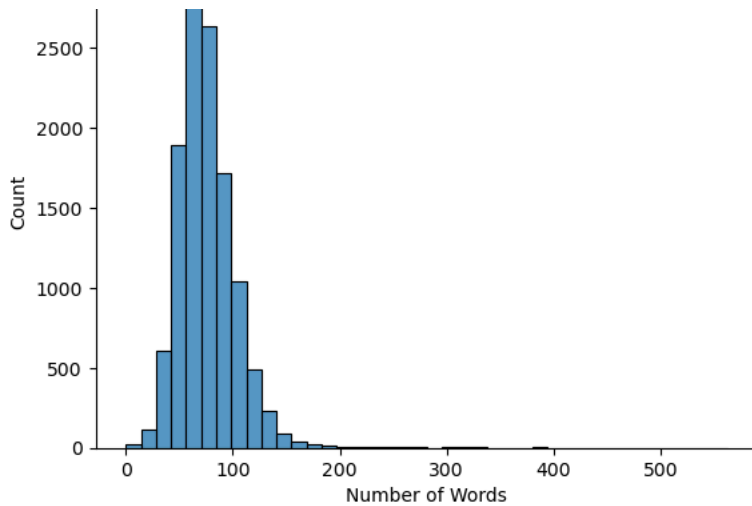
```

title 0
published_date 0
reporter 990
category 0
url 0
content 0
dtype: int64

```

Text Length Distribution





Step 2: Data Preprocessing

```

1 import re
2 import nltk
3 from nltk.corpus import stopwords
4 nltk.download('stopwords')
5
6 # Remove non-Bangla characters & stopwords
7 stop_words = set(stopwords.words('bengali'))
8
9 def clean_text(text):
10     text = str(text).lower()
11     text = re.sub(r'http\S+', '', text)
12     text = re.sub(r'^অ-ঔক-হড়-ঝ\s', '', text)
13     text = " ".join([word for word in text.split() if word not in stop_words])
14     return text.strip()
15
16 # Apply cleaning
17 df['clean_content'] = df['content'].apply(clean_text)
18 df.dropna(subset=['clean_content', 'category'], inplace=True)
19 df = df[df['clean_content'].str.strip() != '']
20

```

```

[nltk_data] Downloading package stopwords to
[nltk_data]    /Users/md.jubayerhossain/nltk_data...
[nltk_data] Package stopwords is already up-to-date!

```

Step 3: Feature Engineering (Label Encoding & Train-Test Split)

```

1 from sklearn.preprocessing import LabelEncoder
2 from sklearn.model_selection import train_test_split
3
4 # Encode labels
5 le = LabelEncoder()
6 df['label'] = le.fit_transform(df['category'])
7
8 # Split
9 X_train, X_test, y_train, y_test = train_test_split(
10     df['clean_content'], df['label'], test_size=0.2, stratify=df['label'], random_state=42
11 )
12

```

Step 4: Model Creation (Transformer - BanglaBERT)

```


1 import torch
2 from transformers import AutoTokenizer, AutoModelForSequenceClassification, TrainingArguments, Trainer
3 from datasets import Dataset
4 from tensorflow import keras
5
6 # Load tokenizer
7 tokenizer = AutoTokenizer.from_pretrained("sagorsarker/bangla-bert-base")
8

```

```

9 # Prepare datasets
10 train_dataset = Dataset.from_dict({'text': X_train.tolist(), 'label': y_train.tolist()})
11 test_dataset = Dataset.from_dict({'text': X_test.tolist(), 'label': y_test.tolist()})
12
13 # Tokenization
14 def tokenize(batch):
15     return tokenizer(batch['text'], padding=True, truncation=True, max_length=256)
16
17 train_dataset = train_dataset.map(tokenize, batched=True)
18 test_dataset = test_dataset.map(tokenize, batched=True)
19
20 # Load model
21 model = AutoModelForSequenceClassification.from_pretrained("sagorsarker/bangla-bert-base", num_labels=len(le.classes_))
22

```

 The OrderedVocab you are attempting to save contains holes for indices [1015, 1016, 1017, 1018, 1053, 1054, 1055, 1056, 1057, 1060, 1061]

 Map: 0% | 0/9521 [00:00<?, ? examples/s]

 The OrderedVocab you are attempting to save contains holes for indices [1015, 1016, 1017, 1018, 1053, 1054, 1055, 1056, 1057, 1060, 1061]

 Map: 0% | 0/2381 [00:00<?, ? examples/s]

 Some weights of BertForSequenceClassification were not initialized from the model checkpoint at sagorsarker/bangla-bert-base and are new

 You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

Step 5: Model Evaluation

```

1 import evaluate
2 import numpy as np
3
4 # Load metrics once
5 accuracy_metric = evaluate.load("accuracy")
6 f1_metric = evaluate.load("f1")
7 precision_metric = evaluate.load("precision")
8 recall_metric = evaluate.load("recall")
9
10 def compute_metrics(eval_pred):
11     logits, labels = eval_pred
12     preds = np.argmax(logits, axis=-1)
13     return {
14         "accuracy": accuracy_metric.compute(predictions=preds, references=labels)["accuracy"],
15         "f1": f1_metric.compute(predictions=preds, references=labels, average="weighted")["f1"],
16         "precision": precision_metric.compute(predictions=preds, references=labels, average="weighted")["precision"],
17         "recall": recall_metric.compute(predictions=preds, references=labels, average="weighted")["recall"]
18     }


```

Step 6: Model Tuning & Training

```

1 print(type(model))
2 print(len(train_dataset))
3 print(len(test_dataset))

```

 <class 'transformers.models.bert.modeling_bert.BertForSequenceClassification'>


 9521

 2381

```

1 import transformers
2 print(transformers.__version__)

```

 4.46.0

```

1 import numpy as np
2 import evaluate
3 from transformers import Trainer, TrainingArguments
4
5 # Load evaluation metrics
6 accuracy = evaluate.load("accuracy")
7 f1 = evaluate.load("f1")
8 precision = evaluate.load("precision")
9 recall = evaluate.load("recall")
10
11 # Define compute_metrics function
12 def compute_metrics(eval_pred):
13     logits, labels = eval_pred
14     preds = np.argmax(logits, axis=-1)
15     return {

```

```
16     "accuracy": accuracy.compute(predictions=preds, references=labels)["accuracy"],
17     "f1": f1.compute(predictions=preds, references=labels, average="weighted")["f1"],
18     "precision": precision.compute(predictions=preds, references=labels, average="weighted")["precision"],
19     "recall": recall.compute(predictions=preds, references=labels, average="weighted")["recall"]
20 }
21
22 # Set training arguments
23 training_args = TrainingArguments(
24     output_dir="./results",
25     evaluation_strategy="epoch",          # ← this works in your version
26     save_strategy="epoch",
27     learning_rate=2e-5,
28     per_device_train_batch_size=32,
29     per_device_eval_batch_size=32,
30     num_train_epochs=1,
31     metric_for_best_model="f1",
32     load_best_model_at_end=True,
33     logging_dir='./logs',
34     logging_steps=100
35 )
36
37 # Define Trainer
38 trainer = Trainer(
39     model=model,                          # ← make sure your model is loaded
40     args=training_args,
41     train_dataset=train_dataset,          # ← make sure train_dataset is ready
42     eval_dataset=test_dataset,           # ← make sure test_dataset is ready
43     tokenizer=tokenizer,                  # ← ensure tokenizer is passed correctly
44     compute_metrics=compute_metrics
45 )
46
47 # Train the model
48 trainer.train()
49
50 # Evaluate on test set
51 results = trainer.evaluate()
52 print("Evaluation results:", results)
```