

5/2/25

Ex: 4

5/2/25

## Text Preprocessing And Analytics Pipeline

Aim:

To perform text cleaning (remove stop words, special characters) and tokenisation on a text dataset.

PROBLEM CODE:

```
import pandas as pd
import re
import spacy

nlp = spacy.load("en-core-web-sm")
df = pd.read_csv('amazon-reviews.csv')
print(df['reviewText'].head())

def clean_text(text):
    if pd.isnull(text):
        return []
    text = text.lower()
    text = re.sub(r'[\n\r\t]', "", text)
    text = text.encode('ascii', 'ignore').
        decode('ascii')

    doc = nlp(text)
    tokens = [token.text for token in doc
               if not token.is_stop and
               token.is_punct]

    return tokens
```

```
df['cleaned_tokens'] = df['reviewText'].
    apply(clean_text)
print(df[['reviewText', 'cleaned_tokens']].
    head(5))
```

2025/10/25

## OUTPUT:

We get this OTR for my husband who is an (OTR) : ....

- 1) I'm professional OTR truck driver, and
- 2) well what can I say. I have had this unit....
- 3) Not going to write a long review even though....
- 4) I've had mine for a year and what it's go ....

N

all\_tokens = [token for token in df['cleaned\_token']]

for token in tokens]

from collections import Counter

word\_freq = Counter(all\_tokens)

print("\nTop 15 frequent words in Amazon reviews:")

print(word\_freq.most\_common(15))

RESULT:

Thus the text cleaning performance has been executed successfully.