

TEST PREPROCESSING AND ANALYTICS PIPELINE.

AIM:

TO perform test cleaning (remove stop words, special characters) and tokenization on a text dataset.

PROGRAM 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_spacy (text):
    if pd.isnull (text):
        return()
    text = text.lower()
    text = re.sub & r'[" \w\w\s]', "", text)
    text = text.encode ('ascii', 'ignore').decode
    ('ascii')
    doc = nlp (text)
    tokens = [token.text for token in doc
              if not token.is_stop and not
              token.is_punct]
    return tokens
df ['cleaned_tokens'] = df ['reviewText'].apply (clean_Text_spacy)
print (df ['reviewText', 'cleaned_tokens'],
      head (5))

```

OUTPUT:

me got this GPS for my husband who is an (OTR):...

- 1) I'm professional OTR truck driver, and I bou...
- 2) Well, what can I say. I've had this unit in m...
- 3) Not going to write a long review, even thought...
- 4) I've had mine for a year and here's what me go..

Name : newtext, alttype: Object.

all_tokens = [token for token in df['cleaned_token']
for token in tokens]

from collections import Counter

word_freq = Counter(all_tokens)

Print("Top 15 frequent words in Amazon
Reviews: ")

Print(word_freq.most_common(15))

RESULT:

Thus the text cleaning performance has
been executed successfully.