

AIM:

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

PROGRAM CODE:

```
import pandas as pd
import re
import spacy
nlp = spacy.load("en_core_web_sm")

df = pd.read_csv('amazon_review.csv')
print(df['review_text'].head())

def clean_text_spacy(text):
    if pd.isnull(text):
        return ''
    else:
        text = text.lower()
        text = re.sub(r'[^a-zA-Z]', " ", text)
        text = text.encode('ascii', 'ignore').decode('utf-8')
        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['review_text'].apply(clean_text_spacy)

print(df[['review_text', 'cleaned_tokens']].head(5))
```



all -

July 27, 2019

we got some friends that night
but we were invited for a wedding
. but

OUTPUT:

we got this GPS for my husband who is
an OTR

- 1) I'm professional OTR truck driver, and I
considered no 1) best price = the best
- 2) well, what can I say, I've had this
(about 1 year) and I put it in
- 3) not going to write a long review,
even though.
- 4) I've had mine for a year and
here, what we go...
name : reinertext, alttype, object.

✓
P



Shot on OnePlus

all - doc = \sum tokens for tokens in df
for token in tokens]

from collections import Counter

word_freq = Counter(all_tokens)

print("In Top 15 frequent words in

Amazon Reviews!")

print(word_freq.most_common(15))

* Result:

Thus the doc cleaning performance

has been executed successfully.



Shot on OnePlus