

Ques:

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

Code

```
import pandas as pd
import spacy
nlp = spacy.load("en-core-web-sm")
df = pd.read_csv("amazon_reviews.csv")
print(df['reviewText'].head(1))
```

```
def clean_text_spacy(text):
```

```
    if pd.isnull(text):
```

```
        return []
```

```
    text = text.lower()
```

```
    text = re.sub(r'[^\w\s]', "", text)
```

```
    doc = nlp(text)
```

```
    tokens = [token.text for token in doc if not
```

```
              token.is_stop and not token.is_punct]
```

```
return tokens
```

REVIEW:

we got this one for my husband who is an
(OTR) ...

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

name : reviewText, dtype : object

df_cleaned_tokenize = df['content'].apply(lambda text:
print(df[[df['content'] == 'cleaned_tokens']].head(5))

all_tokens = [tokens for tokens in df['cleaned_tokens']
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))

Top 15 frequent words in Amazon Reviews
are as follows:

1. the
2. and
3. to
4. of
5. in
6. a
7. for
8. on
9. is
10. that
11. with
12. it
13. for
14. this
15. an

After executing above code, output is as follows:

1. the
2. and
3. to
4. of
5. in
6. a
7. for
8. on
9. is
10. that
11. with
12. it
13. for
14. this
15. an

Results :

The program to perform text
cleaning has been executed successfully