

AIM:

To clean and preprocess amazon review
tools for analysis.

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['reviewText'].head())

def clean(text):
    if pd.isnull(text):
        return ()
    text = text.lower()
    text = re.sub('[^a-zA-Z]', '', text)
    text = text.encode('ascii', ignorecode)

    doc = nlp(text)

    if not doc.is_stop and not doc.is_punct:
        return tokens.
```

cleaned - token:

0 [got, ghs, respond, other, road ...

1 [m, professional, after, much ...

2 [duck, dogs, prior, ...

3 [going, write, thought, ...

4 [use, got, rocks, got, there, ...

Top 15 frequent words in Amazon reviews:

[('1', 3203), ('root', 1947), ('it', 962),

('books', 005), ('kindle', 501), ...

('great', 422), ('use', 120), ...]

```
df['clean'] = df['review_text'].apply  
(clean_text_story)
```

```
print cat ['review_text', 'cleaned'] - head  
(5)
```

```
all_token = [ token for token in df  
(cleaned_token)
```

```
for token in means ]
```

```
from collections import counter
```

```
word_freq = counter all_token
```

```
print ('word_freq', word_freq)
```

✓
RESULT:

Cleaned token contains meaningful word
without noise and frequent words reflect
from the review.