... preparation and analyzing pipeline

**Aim:** To learn and program once have day to a list . the idally the most frequent words-

**Algorithm:**

1- load dataset :-

Import raw data from program.-
raw csv.

2- clean data :-

convert to lower remove punctation
and special character the two
apars

3- remove stopwords :-

frite and common stopwords
and munilalis token

4- Apply cleaning :-

merge the entire raw list
colour to generate class.

5- Analyze freqs :-

plot all item sent word freqg
very law

**Program:**

import pandas as pd

import "

import spacy

txt = spacy.load (" en = web , sci" )

df = pd . read -csv ( "data - raw - csv" )

Print ( df ( raw txt )) . head () )

df colum data_present txt :

df pd- Inull (data

return c]

Output:

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

1.) I am a proffessional OTR toull driver

2.) well, what can I say

3.) not going to write a long review, even thayt

4.) I've had mine for a year name. review Gp iTout

```python
text = text.lower()
text = re.sub(r'[^\w\n\s]', '', text)
text = stopwords.words("english - stop py" - about ('auii'))
doc = nlp(text)
# now is .et and not do is read
words = text

print text ['our text', 'element.doth]. head.(500)
all.tokens = [ token of the in dtype dont .thu ]
          for the .words an token ]

from collection import count

word.free = count (all-itr )

print ('TP is from words in Agron text''')
print ( word.frequent, most. common (15))
```

Result:

cleaned tokens contains meaningful words without
not used this words frequently than output the
terms from the value.