**Program:**

```
from nltk.corpus import stopwords
from nltk.corpus import words

def tokenize(s):
    #split by white text
    space_split = list(s.split())

    #dealing with punctuations
    punctuation = [',', '.', '!', '?']

    tokens=[]
    for i in space_split:
        f=0
        for j in punctuation:
            if j in i:
                tokens.append(i.replace(j, ''))
                tokens.append(j)
                f=1
        if f==0:
            tokens.append(i)

    #dealing with 'm and 's
    tokens2 = []
    for i in tokens:
        if '\'s' in i:
            tokens2.append(i.replace('\'s', ''))
            tokens2.append('\'s')
        elif '\'m' in i:
            tokens2.append(i.replace('\'m', ''))
            tokens2.append('am')
        else:
            tokens2.append(i)
    for i in tokens:
        if '-' in i:
            a,b = i.split('-',1)
            tokens2.append(a)
            tokens2.append('-')
            tokens2.append(b)

    return tokens2
def remove_stop_words(a):
    stop_words = set(stopwords.words('english'))
```

```
    b = [w for w in a if w not in stop_words]
    return b

def filter(s):
    b=[]
    for i in s:
        if i in words.words():
            b.append(i)
    return b

with open('text.txt', 'r+') as f:
    data=f.read()
print("Text:","\n",data)
tokens = tokenize(data)
print("\nTokens: ",tokens)

filtered_words = filter(tokens)
print("\n\nFiltered Words: ", filtered_words)
```

**Output:**
Text:
 This is a sample sentence, showing off the stop words filtration.
I'm rey's friend. I worked hard.
Khana khaya kya
Python is a high-level, interpreted, interactive and object-oriented scripting language. Python is designed to be highly readable. It uses English keywords frequently where as other languages use punctuation.

Tokens:  ['This', 'is', 'a', 'sample', 'sentence', ',', 'showing', 'off', 'the', 'stop', 'words', 'filtration', '.', 'I', 'am', 'rey', "'s", 'friend', '.', 'I', 'worked', 'hard', '.', 'Khana', 'khaya', 'kya', 'Python', 'is', 'a', 'high-level', ',', 'interpreted', ',', 'interactive', 'and', 'object-oriented', 'scripting', 'language', '.', 'Python', 'is', 'designed', 'to', 'be', 'highly', 'readable', '.', 'It', 'uses', 'English', 'keywords', 'frequently', 'where', 'as', 'other', 'languages', 'use', 'punctuation', '.', 'high', '-', 'level', 'object', '-', 'oriented']


Filtered Words:  ['is', 'a', 'sample', 'sentence', 'showing', 'off', 'the', 'stop', 'filtration', 'I', 'am', 'friend', 'I', 'worked', 'hard', 'is', 'a', 'interactive', 'and', 'language', 'is', 'designed', 'to', 'be', 'highly', 'readable', 'English', 'frequently', 'where', 'as', 'other', 'use', 'punctuation', 'high', 'level', 'object']