```
In [1]:
import numpy as np
import pandas as pd
from wordcloud import WordCloud
import nltk
from nltk.corpus import stopwords
import matplotlib.pyplot as plt
from sklearn.model selection import train test split
from sklearn.feature extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.metrics import accuracy score
In [2]:
Data=pd.read csv("C:/Users/nehal/Downloads/spam.csv")
In [3]:
Data.head()
Out[3]:
                                              v2 Unnamed: 2 Unnamed: 3 Unnamed: 4
     v1
    ham
           Go until jurong point, crazy.. Available only ...
                                                        NaN
                                                                   NaN
                                                                               NaN
                           Ok lar... Joking wif u oni...
                                                        NaN
                                                                   NaN
                                                                              NaN
    ham
              Free entry in 2 a wkly comp to win FA Cup
2 spam
                                                        NaN
                                                                   NaN
                                                                               NaN
                                                                               NaN
          U dun say so early hor... U c already then say...
                                                        NaN
                                                                   NaN
3
    ham
                                                        NaN
                                                                   NaN
                                                                               NaN
    ham
           Nah I don't think he goes to usf, he lives aro...
In [4]:
Data=Data[['v1','v2']]
In [5]:
Data.head()
Out[5]:
                                              v2
     v1
    ham
           Go until jurong point, crazy.. Available only ...
0
    ham
                           Ok lar... Joking wif u oni...
              Free entry in 2 a wkly comp to win FA Cup
2 spam
                                           fina...
    ham
          U dun say so early hor... U c already then say...
           Nah I don't think he goes to usf, he lives aro...
    ham
In [6]:
Data=Data.rename(columns={'v1' : 'Label', "v2" : 'Text'})
In [7]:
Data['Label'].value counts(ascending=True)
Out[7]:
           747
spam
         4825
ham
Name: Label, dtype: int64
```

In [8]:

```
nltk.download("punkt")
import warnings
warnings.filterwarnings('ignore')

[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\nehal\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

In [9]:

```
ham_words = ''
spam_words = ''
```

In [10]:

```
for val in Data[Data['Label'] == 'spam'].Text:
    Text = val.lower()
    tokens = nltk.word_tokenize(Text)
    for words in tokens:
        spam_words = spam_words + words + ' '

for val in Data[Data['Label'] == 'ham'].Text:
    Text = Text.lower()
    tokens = nltk.word_tokenize(Text)
    for words in tokens:
        ham_words = ham_words + words + ' '
```

In [11]:

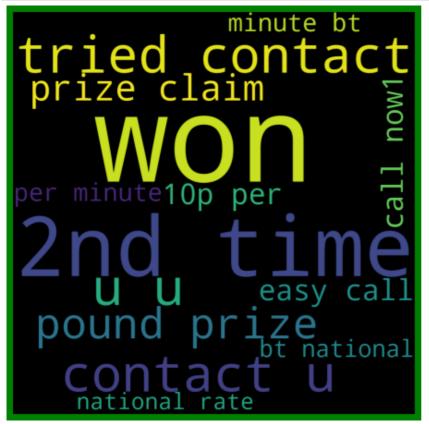
```
spam_wordcloud = WordCloud(width=1000, height=1000).generate(spam_words)
ham_wordcloud = WordCloud(width=1000, height=1000).generate(ham_words)
```

In [12]:

```
plt.figure( figsize=(10,8), facecolor='w')
plt.imshow(spam_wordcloud)
plt.axis("off")
plt.show()
```



```
plt.figure( figsize=(10,8), facecolor='g')
plt.imshow(ham_wordcloud)
plt.axis("off")
plt.show()
```



In [14]:

```
Data = Data.replace(['ham', 'spam'], [0, 1])
Data.head(10)
```

Out[14]:

	Label	Text
0	0	Go until jurong point, crazy Available only
1	0	Ok lar Joking wif u oni
2	1	Free entry in 2 a wkly comp to win FA Cup fina
3	0	U dun say so early hor U c already then say
4	0	Nah I don't think he goes to usf, he lives aro
5	1	FreeMsg Hey there darling it's been 3 week's n
6	0	Even my brother is not like to speak with me
7	0	As per your request 'Melle Melle (Oru Minnamin
8	1	WINNER!! As a valued network customer you have
9	1	Had your mobile 11 months or more? U R entitle

In [15]:

```
nltk.download('stopwords')
import string
def text_process(Text):

    Text = Text.translate(str.maketrans('', '', string.punctuation))
    Text = [word for word in Text.split() if word.lower() not in stopwords.words('englis h')]
    return " ".join(Text)
[nltk_data] Downloading package stopwords to
```

```
[nltk data]
                C:\Users\nehal\AppData\Roaming\nltk data...
               Package stopwords is already up-to-date!
[nltk_data]
In [16]:
Data['Text'] = Data['Text'].apply(text process)
Data.head()
Out[16]:
  Label
                                        Text
0
      0 Go jurong point crazy Available bugis n great ...
1
                           Ok lar Joking wif u oni
2
      1 Free entry 2 wkly comp win FA Cup final tkts 2...
3
      0
                 U dun say early hor U c already say
           Nah dont think goes usf lives around though
In [17]:
Text = pd.DataFrame(Data['Text'])
Label = pd.DataFrame(Data['Label'])
In [18]:
vectorizer = TfidfVectorizer()
vectors = vectorizer.fit transform(Data['Text'])
vectors.shape
Out[18]:
(5572, 9333)
In [19]:
features = vectors
In [20]:
X_train, X_test, y_train, y_test = train_test_split(features, Data['Label'], test_size=0
.2, random state=111)
In [21]:
from sklearn.naive bayes import MultinomialNB
from sklearn.tree import DecisionTreeClassifier
In [22]:
mnb = MultinomialNB(alpha=0.2)
dtc = DecisionTreeClassifier(min_samples_split=7, random_state=111)
In [23]:
clfs = { 'NB': mnb, 'DT': dtc}
In [24]:
def train(clf, features, targets):
    clf.fit(features, targets)
def predict(clf, features):
    return (clf.predict(features))
In [25]:
pred scores word vectors = []
```

```
for k,v in clfs.items():
   train(v, X_train, y_train)
   pred = predict(v, X_test)
    pred_scores_word_vectors.append((k, [accuracy_score(y_test , pred)]))
In [26]:
pred scores word vectors
Out[26]:
[('NB', [0.9856502242152466]), ('DT', [0.9650224215246637])]
In [27]:
def SpamOrNot(x):
    if x == 1:
       print ("Message is SPAM")
    elif x==0:
       print ("Message is NOT Spam")
In [28]:
newtext = ["prize"]
integers = vectorizer.transform(newtext)
In [29]:
x = mnb.predict(integers)
SpamOrNot(x)
Message is SPAM
In [ ]:
```