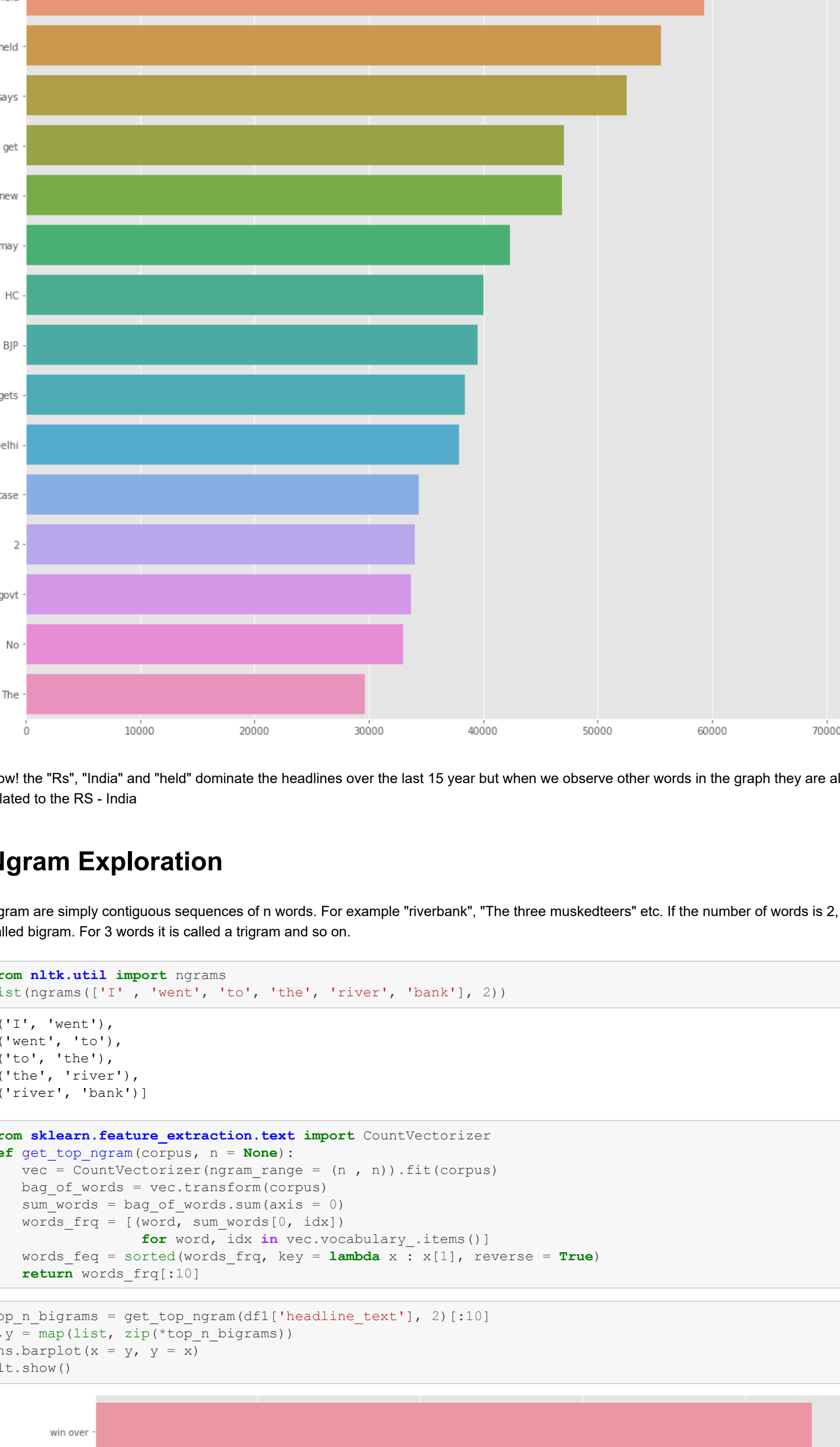



```
[35]: sns.barplot(x = y, y = x)
plt.show()
```



Wow! the "Rs", "India" and "held" dominate the headlines over the last 15 year but when we observe other words in the graph they are all related to the RS - India

Ngram Exploration

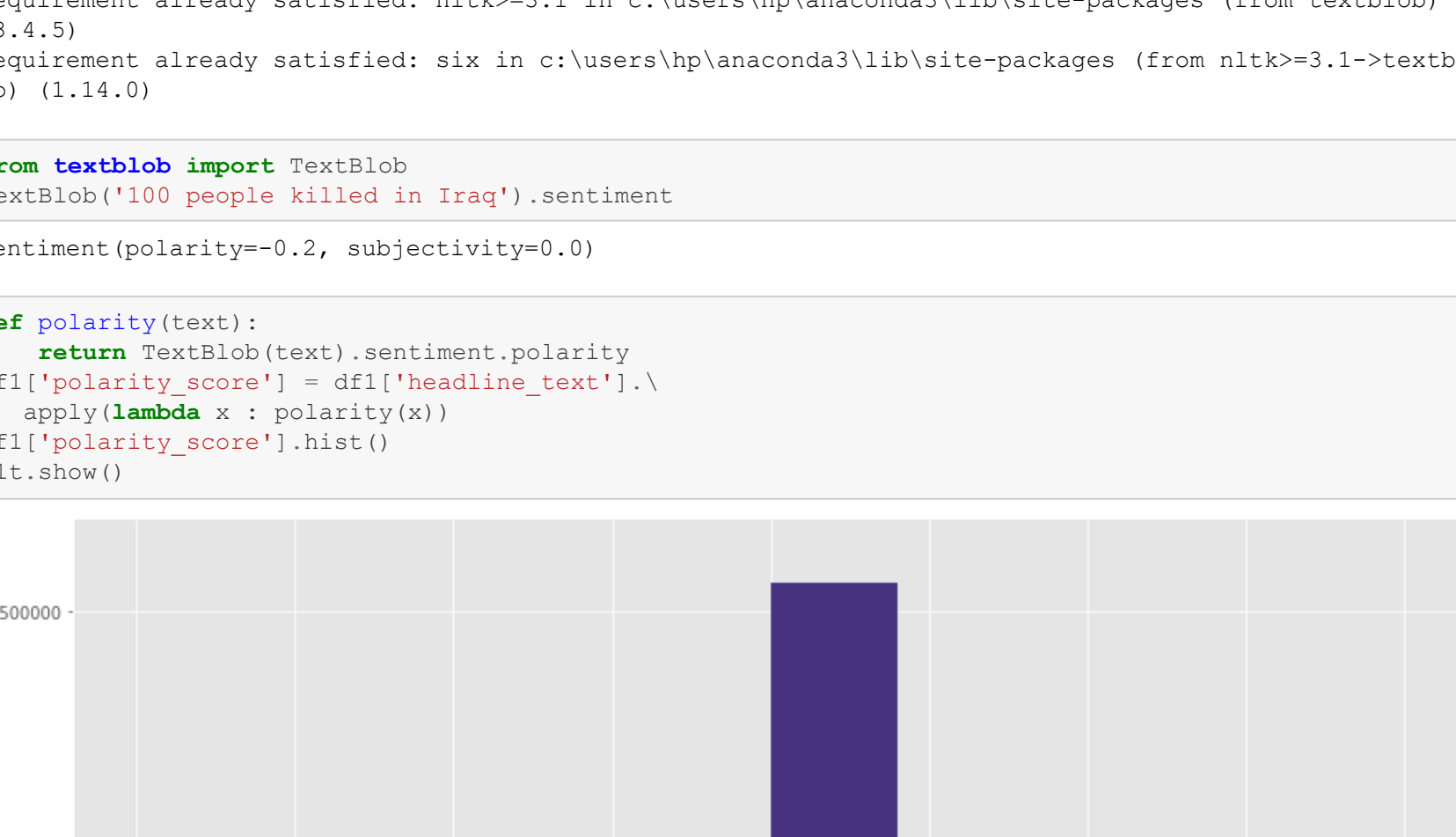
Ngram are simply contiguous sequences of n words. For example "riverbank", "The three musketeers" etc. If the number of words is 2, it is called bigram. For 3 words it is called a trigram and so on.

```
In [36]: from nltk.util import ngrams
list(ngrams(['I', 'went', 'to', 'the', 'river', 'bank'], 2))
```

```
Out[36]: [('I', 'went'),
('went', 'to'),
('to', 'the'),
('the', 'river'),
('river', 'bank')]
```

```
In [37]: from sklearn.feature_extraction.text import CountVectorizer
def get_top_ngrams(corpus, n = None):
    vec = CountVectorizer(ngram_range = (n, n)).fit(corpus)
    bag_of_words = vec.transform(corpus)
    sum_words = bag_of_words.sum(axis = 0)
    words_freq = [(word, sum_words[word_idx]) for word, word_idx in vec.vocabulary_.items()]
    words_freq = sorted(words_freq, key = lambda x : x[1], reverse = True)
    return words_freq[:10]
```

```
In [46]: top_n_bigrams = get_top_ngram(df['headline_text'], 2)[:10]
x,y = map(list, zip(*top_n_bigrams))
sns.barplot(x = y, y = x)
plt.show()
```



Textblob

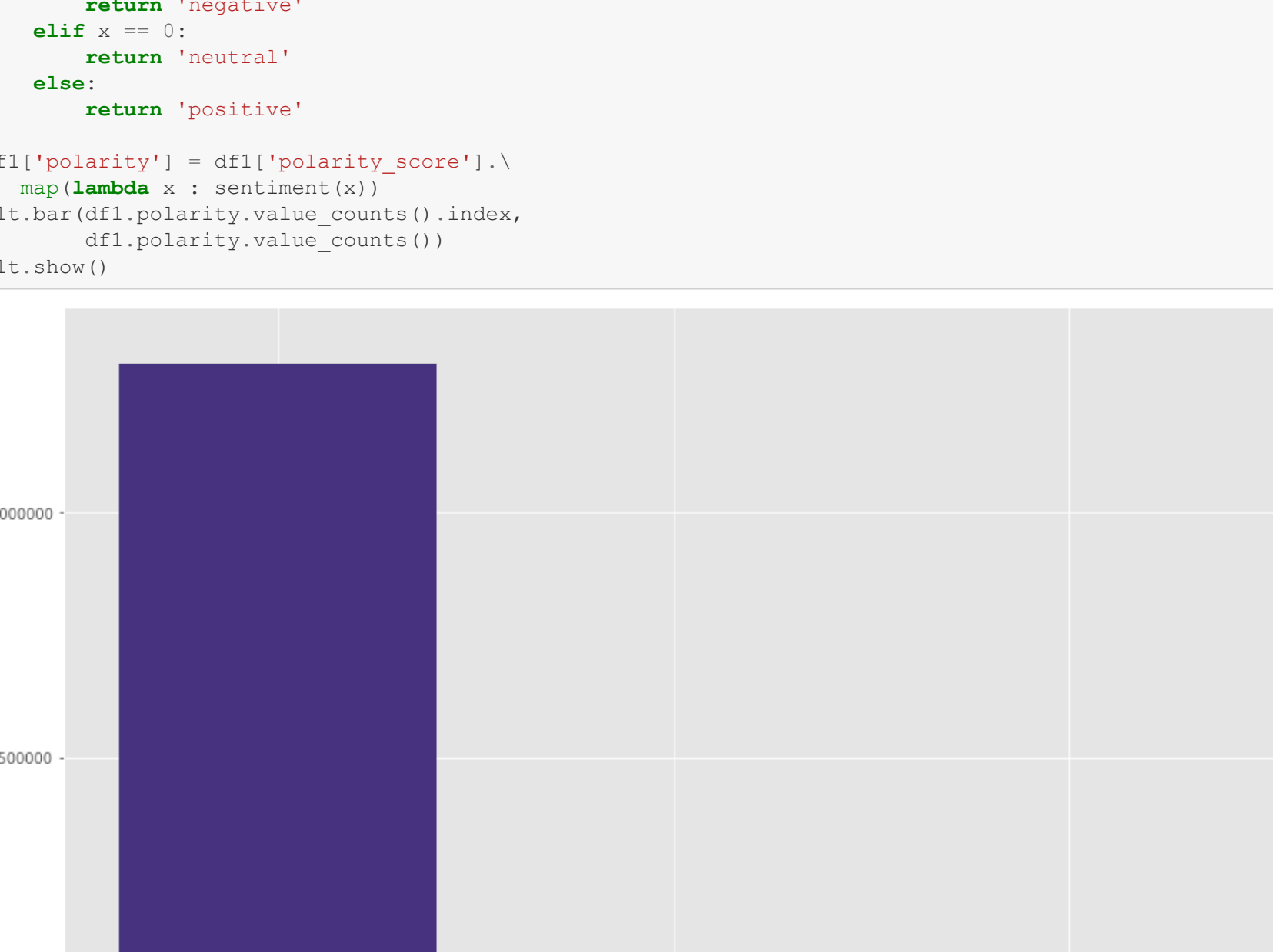
```
In [43]: !pip install textblob
```

Requirement already satisfied: textblob in c:\users\hp\anaconda3\lib\site-packages (0.15.3)
Requirement already satisfied: nltk>=3.1 in c:\users\hp\anaconda3\lib\site-packages (from textblob) (3.4.5)
Requirement already satisfied: six in c:\users\hp\anaconda3\lib\site-packages (from nltk>=3.1->textblob) (1.14.0)

```
In [40]: from textblob import TextBlob
TextBlob("100 people killed in Iraq").sentiment
```

```
Out[40]: Sentiment(polarity=0.2, subjectivity=0.0)
```

```
In [41]: def polarity(text):
def sentiment(text):
    return TextBlob(text).sentiment.polarity
df['polarity_score'] = df['headline_text'].apply(lambda x : polarity(x))
df['polarity_score'].hist()
plt.show()
```



```
In [44]: def sentiment(x):
if x < 0:
    return 'negative'
elif x == 0:
    return 'neutral'
else:
    return 'positive'

df['polarity'] = df['polarity_score'].map(lambda x : sentiment(x))
plt.bar(df['polarity'].value_counts().index, df['polarity'].value_counts())
plt.show()
```



```
In [45]: df[df['polarity'] == 'positive']['headline_text'].head()
```

```
Out[45]: 0    Win over cena satisfying but defeating underta...
5    Extra buses to clear tourist traffic
13   Will Qureshi's return really help the govt?
31   Extra buses to clear tourist traffic
39   Will Qureshi's return really help the govt?
Name: headline_text, dtype: object
```

```
In [48]: df[df['polarity'] == 'negative']['headline_text'].head()
```

```
Out[48]: 66   Destroying myths and doubts on sexuality
87   Powerless north India gropes in the dark
103  10-year-old girl missing
132   Net lottery: A winner or a sucker?
143   Mental illness can pass from parent to child
Name: headline_text, dtype: object
```