

```

1 from collections import Counter
2 from tqdm import tqdm
3 from scipy.sparse import csr_matrix
4 import math
5 import operator
6 from sklearn.preprocessing import normalize
7 import numpy

1 def fit_50(corpus):
2     #creating set for storing unique_words
3     u_w = set()
4     if isinstance(corpus, list):
5         #iterating through rows in corpus
6         for row in corpus:
7             #iterating through each words in row
8             for words in row.split():
9                 #checking for length of words
10                if len(words) < 2:
11                    continue
12                #adding each words in set
13                u_w.add(words)
14
15            #converting set into sorted list
16            u_w = list(u_w)
17
18            #creating dict with words in list as keys and enumerated index as values
19            vocab = {j:i for i, j in enumerate(u_w)}
20
21            #creating dict with keys in vocab dict as keys and its idf as values
22            vocab_idf = {word:get_idf_50(word,corpus) for word in vocab.keys()}
23
24            #gettingtop 50 values based on idf values
25            vocab_idf_50 = dict(sorted(vocab_idf.items(),key=operator.itemgetter(1),reverse=True)
26
27            #creating dict with top 50 words as keys and its enumerated index as values
28            vocab_50 = {j:i for i, j in enumerate(vocab_idf_50.keys())}
29
30            return vocab_50
31    else:
32        print("Send list of Sentences")

1 def transform_50(corpus, vocab_50):
2     #creating empty lists for rows, column, values
3     rows = []
4     columns = []
5     values = []
6     if isinstance(corpus, list):
7         #iterating through rows in corpus
8         for idx, row in enumerate(corpus):
9             #creating dict with words in row and its count as values
10            word_freq = dict(Counter(row.split()))
11            #iterating through keys in top 50 vocab dict

```

```

12     for word in vocab_50.keys():
13
14         tfidf = (word_freq.get(word, 0) / len(row.split())) * get_idf_50(word, corpus)
15
16         if tfidf != 0:
17             #appending row index(idx) into row list
18             rows.append(idx)
19
20             #getting column index from top 50 vocab dict
21             col_index = vocab_50.get(word, 0)
22             columns.append(col_index)
23
24             #appending tfidf to values list
25             values.append(tfidf)
26     #norm = normalize(values)
27
28     #return csr_matrix(norm)
29     return csr_matrix((values, (rows, columns)), shape=(len(corpus),len(vocab_50)))
30 else:
31     print("Send list of Sentences")

```

```

1 def get_idf_50(word, corpus):
2     count=0
3     #iterating through rows in corpus
4     for r in corpus:
5         #if word in that row increament count by one
6         if word in r:
7             count += 1
8     idf_key= 1 + math.log((1+len(corpus)) / (count+1))
9     return idf_key

```

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

Mounted at /content/drive

```

1 import pickle
2 with open('/content/drive/My Drive/Dataset/cleaned_strings', 'rb') as f:
3     corpus = pickle.load(f)
4
5 # printing the length of the corpus loaded
6 print("Number of documents in corpus = ",len(corpus))

```

Number of documents in corpus = 746

```

1 #printing keys in top 50 vocab dict
2 vocab_50 = fit_50(corpus)
3 print(list(vocab_50.keys()))

```

['garfield', 'tongue', 'sacrifice', 'superlative', 'cheerless', 'guys', 'letting', 's



```

1 #top 50 idf values after fit method

```

1

✓ 0s completed at 15:33

