

TF-IDF stands for Term Frequency Inverse Document Frequency of records. It can be defined as the calculation of how relevant a word in a series or corpus is to a text. The meaning increases proportionally to the number of times in the text a word appears but is compensated by the word frequency in the corpus (data-set).

In document  $d$ , the frequency represents the number of instances of a given word  $t$ . Therefore, we can see that it becomes more relevant when a word appears in the text, which is rational. Since the ordering of terms is not significant, we can use a vector to describe the text in the bag of term models. For each specific term in the paper, there is an entry with the value being the term frequency.

$tf(t,d) = \text{count of } t \text{ in } d / \text{number of words in } d$

$df(t) = \text{occurrence of } t \text{ in documents}$

$df(t) = N(t)$  where  $df(t)$  = Document frequency of a term  $t$   $N(t)$  = Number of documents containing the term  $t$

$idf(t) = N / df(t) = N / N(t)$

$tf-idf(t, d) = tf(t, d) * idf(t)$

```
from sklearn.feature_extraction.text import TfidfVectorizer

# assign documents
d0 = 'Geeks for geeks'
d1 = 'Geeks'
d2 = 'r2j'

# merge documents into a single corpus
string = [d0, d1, d2]

# create object
tfidf = TfidfVectorizer()

# get tf-df values
result = tfidf.fit_transform(string)

print('\nidf values:')
for ele1, ele2 in zip(tfidf.get_feature_names(), tfidf.idf_):
```

```
print(ele1, ': ', ele2)
```

```
idf values:
```

```
for : 1.6931471805599454
```

```
geeks : 1.2876820724517808
```

```
r2j : 1.6931471805599454
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning:
warnings.warn(msg, category=FutureWarning)
```

```
print('\nWord indexes:')
```

```
print(tfidf.vocabulary_)
```

```
# display tf-idf values
```

```
print('\ntf-idf value:')
```

```
print(result)
```

```
# in matrix form
```

```
print('\ntf-idf values in matrix form:')
```

```
print(result.toarray())
```

```
Word indexes:
```

```
{'geeks': 1, 'for': 0, 'r2j': 2}
```

```
tf-idf value:
```

```
(0, 0) 0.5493512310263033
```

```
(0, 1) 0.8355915419449176
```

```
(1, 1) 1.0
```

```
(2, 2) 1.0
```

```
tf-idf values in matrix form:
```

```
[[0.54935123 0.83559154 0.          ]
```

```
 [0.          1.          0.          ]
```

```
 [0.          0.          1.          ]]
```

---

✓ 0s completed at 12:35 PM

