2. Create representation of document by calculating Term Frequency and Inverse Document Frequency.

```
!pip install scikit-learn
```

```python
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
import math
```

```python
documentA = 'Jupiter is the largest Planet'
documentB = 'Mars is the fourth planet from the Sun'
```

```python
bagOfWordsA = documentA.split(' ')
bagOfWordsB = documentB.split(' ')
```

— + Code — + Text —

```python
uniqueWords = set(bagOfWordsA).union(set(bagOfWordsB))
```

```python
numOfWordsA = dict.fromkeys(uniqueWords, 0)
for word in bagOfWordsA:
    numOfWordsA[word] += 1

numOfWordsB = dict.fromkeys(uniqueWords, 0)
for word in bagOfWordsB:
    numOfWordsB[word] += 1
```

```python
def computeTF(wordDict, bagOfWords):
    tfDict = {}
    bagOfWordsCount = len(bagOfWords)
    for word, count in wordDict.items():
        tfDict[word] = count / float(bagOfWordsCount)
    return tfDict

tfA = computeTF(numOfWordsA, bagOfWordsA)
tfB = computeTF(numOfWordsB, bagOfWordsB)
```

```python
def computeIDF(documents):
    N = len(documents)
    idfDict = dict.fromkeys(documents[0].keys(), 0)
    for document in documents:
        for word, val in document.items():
            if val > 0:
                idfDict[word] += 1
    for word, val in idfDict.items():
        idfDict[word] = math.log(N / float(val))
    return idfDict

idfs = computeIDF([numOfWordsA, numOfWordsB])
```

```python
def computeTFIDF(tfBagOfWords, idfs):
    tfidf = {}
    for word, val in tfBagOfWords.items():
        tfidf[word] = val * idfs[word]
    return tfidf

tfidfA = computeTFIDF(tfA, idfs)
tfidfB = computeTFIDF(tfB, idfs)
```

```python
df = pd.DataFrame([tfidfA, tfidfB], index=['Document A', 'Document B'])
df
```

|  | planet | Jupiter | Mars | from | fourth | Planet | largest | is | the | Sun |
|---|---|---|---|---|---|---|---|---|---|---|
| **Document A** | 0.000000 | 0.138629 | 0.000000 | 0.000000 | 0.000000 | 0.138629 | 0.138629 | 0.0 | 0.0 | 0.000000 |
| **Document B** | 0.086643 | 0.000000 | 0.086643 | 0.086643 | 0.086643 | 0.000000 | 0.000000 | 0.0 | 0.0 | 0.086643 |

Start coding or generate with AI.