

1 a

$$\begin{aligned}\text{TF-IDF}_{\text{coffee doc1}} &= \log_{10}(25+1) \times \log_{10}(4/2) \\ &= 0.4259\end{aligned}$$

$$\begin{aligned}\text{TF-IDF}_{\text{coffee doc2}} &= \log_{10}(18+1) \times \log_{10}(4/2) \\ &= 0.3849\end{aligned}$$

$$\begin{aligned}\text{TF-IDF}_{\text{coffee doc3}} &= \log_{10}(0) \times \log_{10}(4/2) \\ &= 0\end{aligned}$$

$$\begin{aligned}\text{TF-IDF}_{\text{coffee doc4}} &= 0 \times \log_{10}(4/2) \\ &= 0\end{aligned}$$

1 b

$$\begin{aligned}\text{TF-IDF}_{\text{vacation doc1}} &= \log(20+1) \times \log(4/4) \\ &= 0\end{aligned}$$

$$\begin{aligned}\text{TF-IDF}_{\text{vacation doc2}} &= \log(23+1) \times \log(4/4) \\ &= 0\end{aligned}$$

$$\begin{aligned}\text{TF-IDF}_{\text{vacation doc3}} &= \log(18+1) \times \log(4/4) \\ &= 0\end{aligned}$$

$$\begin{aligned}\text{TF-IDF}_{\text{vacation doc4}} &= \log(13+1) \times \log(4/4) \\ &= 0\end{aligned}$$

1c

~~TF-IDF indicates that~~

1c

TF-IDF scores for coffee indicates that it is a discriminating term in doc1 and doc2 but not in doc3 and doc4, on the other hand, "vacation" is all zero in all document. This suggests that it is not relevant in any of the document.

[10]

$$\text{TFIDF}_{\text{cup}} \text{ doc1} = 0 \times \log_{10}(4/3) = 0$$

$$\text{TFIDF}_{\text{cup}} \text{ doc2} = \log_{10}(87+1) \times \log_{10}(4/3) = 0.2429$$

$$\text{TFIDF}_{\text{cup}} \text{ doc3} = \log_{10}(78+1) \times \log_{10}(4/3) = 0.2371$$

$$\text{TFIDF}_{\text{cup}} \text{ doc4} = \log_{10}(3+1) \times \log_{10}(4/3) = 0.0752$$

$$\text{TFIDF}_{\text{summer}} \text{ doc1} = \log_{10}(10+1) \times \log_{10}(4/4) = 0$$

$$\text{TFIDF}_{\text{summer}} \text{ doc2} = \log_{10}(12+1) \times \log_{10}(4/4) = 0$$

$$\text{TFIDF}_{\text{summer}} \text{ doc3} = \log_{10}(15+1) \times \log_{10}(4/4) = 0$$

$$\text{TFIDF}_{\text{summer}} \text{ doc4} = \log_{10}(80+1) \times \log_{10}(4/4) = 0$$

$$\text{TFIDF}_{\text{walk}} \text{ doc1} = \log_{10}(3+1) \times \log_{10}(4/4) = 0$$

$$\text{TFIDF}_{\text{walk}} \text{ doc2} = \log_{10}(6+1) \times \log_{10}(4/4) = 0$$

$$\text{TFIDF}_{\text{walk}} \text{ doc3} = \log_{10}(32+1) \times \log_{10}(4/4) = 0$$

$$\text{TFIDF}_{\text{walk}} \text{ doc4} = \log_{10}(20+1) \times \log_{10}(4/4) = 0$$

	Doc1	Doc2	Doc3	Doc4
cup	0	0.2429	0.2371	0.0752
summer	0	0	0	0
coffee	0.4259	0.3849	0	0
walk	0	0	0	0
vacation	0	0	0	0

$$\text{doc similarity} = \text{cos-similarity} = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \times \|\vec{y}\|}$$

$$\text{cos}[\text{doc1}, \text{doc2}] = \frac{0.163929}{0.4259 \times 0.455136} + \dots$$

$$\cos(\text{doc1}, \text{doc2}) = \frac{0.163929}{\sqrt{0.4239^2} \times \sqrt{0.2429^2 + 0.3849^2}} = 0.8456$$

$$\cos(\text{doc1}, \text{doc3}) = \frac{0}{\sqrt{0.4239^2} \times \sqrt{0.2371^2}} = 0$$

$$\cos(\text{doc1}, \text{doc4}) = 0 \quad \text{Dot product is zero}$$

$$\cos(\text{doc2}, \text{doc3}) = \frac{0.057592}{0.455136 \times 0.2371} = 0.533687$$

$$\cos(\text{doc2}, \text{doc4}) = \frac{0.018266}{0.455136 \times 0.0752} = 0.533687$$

$$\cos(\text{doc3}, \text{doc4}) = \frac{0.01983}{0.2371 \times 0.0752} = 1$$

most similar

doc3 doc4

doc1 doc2

doc2 doc3

doc2 doc4

least similar

doc1 doc3

doc1 doc4

29

	keyboard	webcam	sheep	less	racinG	count(w)
Whiteboard	1434	512	8	42	3	1999
computer	5290	2351	19	9	42	7211
mouse	6102	3425	352	4251	532	14662
horse	1	53	531	2542	3513	6640
Valetion	14	53	34	5	9	115
Cont(context)	12841	6394	944	6849	4099	31127

25

	keyboard	webcam	sheep	less	racinG	PC(w)
Whiteboard	0.0461	0.01644	0.0003	0.0013	0.0001	0.0642
computer	0.1699	0.0755	0.0006	0.0003	0.0013	0.2477
mouse	0.1960	0.1100	0.0113	0.1366	0.0171	0.4710
horse	0.0003	0.0017	0.01705	0.0817	0.1129	0.2133
Valetion	0.0004	0.0017	0.0011	0.0010	0.0002	0.0037
PC(context)	0.4125	0.2054	0.0303	0.2200	0.1316	1

2C

	keyboard	webcam	sheep	less	racinG
Whiteboard	0.7990	0.3181	0	0	0
computer	0.7336	0.5723	0	0	0
mouse	0.0126	0.18527	0	0.3986	0
horse	0	0	0	0.7875	1.988
Valetion	0	1.1643	3.2945	0.2969	0

2D

$$\text{whiteboard} = [0.7991, 0.3181, 0, 0, 0]$$

$$\text{computer} = [0.7336, 0.5723, 0, 0, 0]$$

$$\text{mouse} = [0.0126, 0.18527, 0, 0.3986, 0]$$

$$\text{horse} = [0, 0, 0, 0.7875, 1.988]$$

$$\text{vacation} = [0, 1.1643, 3.2945, 0.2969, 0]$$

~~whiteboard - mouse =~~

$$\text{cosine similarity} = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| * \|\vec{y}\|}$$

$$\text{whiteboard computer} = \frac{0.768268}{0.86086 \times 0.930428} = 0.960037$$

$$\text{whiteboard mouse} = \frac{0.069003}{0.860086 \times 0.439734} = 0.182447$$

$$\text{horse - vacation} = \frac{0.233809}{2.138294 \times 3.506775} = 0.031181$$

most similar

whiteboard computer

least similar

whiteboard mouse

horse vacation

③

```
[ ] #question a
def find_hyponym():
    house = wn.synset('house.n.01')
    hyponyms = house.hyponyms()

    if hyponyms:
        print("Hyponyms of 'house:")
        for hyponym in hyponyms:
            # Split the hyponym name by periods and keep only the first part
            hyponym_name = hyponym.name().split('.')[0]
            print(hyponym_name)
    else:
        print("No hyponyms found for 'house.n.01'.")

find_hyponym()
```

```
Hyponyms of 'house:
building
dwelling
```



```
from nltk.corpus import wordnet as wn

# Get the synsets for the words "mouse" and "horse"
mouse_synsets = wn.synsets('mouse')
horse_synsets = wn.synsets('horse')
vacation_synsets = wn.synsets('vacation')

# Get the first synset for each word
mouse_synset = mouse_synsets[0]
horse_synset = horse_synsets[0]
vacation_synset = vacation_synsets[0]

# Calculate the path similarity between the first senses of "mouse" and "horse"
mouse_horse_similarity = mouse_synset.path_similarity(horse_synset)
print("Path similarity between 'mouse' and 'horse':", mouse_horse_similarity)

# Calculate the path similarity between the first senses of "horse" and "vacation"
horse_vacation_similarity = horse_synset.path_similarity(vacation_synset)
print("Path similarity between 'horse' and 'vacation':", horse_vacation_similarity)
```




```
#question c
import numpy as np

# Load Glove word vectors
def load_glove_vectors(file_path):
    word_vectors = {}
    with open(file_path, 'r', encoding='utf-8') as f:
        for line in f:
            values = line.split()
            word = values[0]
            vector = np.array(values[1:], dtype='float32')
            word_vectors[word] = vector
    return word_vectors

# Compute cosine similarity between two words
def cosine_similarity(word1, word2, word_vectors):
    if word1 not in word_vectors or word2 not in word_vectors:
        return None

    vec1 = word_vectors[word1]
    vec2 = word_vectors[word2]

    dot_product = np.dot(vec1, vec2)
    norm1 = np.linalg.norm(vec1)
    norm2 = np.linalg.norm(vec2)

    similarity = dot_product / (norm1 * norm2)
    return similarity

# Example usage
glove_file_path = '/content/drive/MyDrive/nlp/hw5/glove.6B/glove.6B.50d.txt'
word_vectors = load_glove_vectors(glove_file_path)
```

```
[ ] #question d
word1 = 'mouse'
word2 = 'horse'
similarity = cosine_similarity(word1, word2, word_vectors)
print(f"Similarity between '{word1}' and '{word2}': {similarity:.4f}")
```

Similarity between 'mouse' and 'horse': 0.4356



```
word1 = 'horse'
word2 = 'vacation'
similarity = cosine_similarity(word1, word2, word_vectors)
print(f"Similarity between '{word1}' and '{word2}': {similarity:.4f}")
```

Similarity between 'horse' and 'vacation': 0.3102