# Quiz 2

## Yulin Hu

**Task 1:**

```python
22    def cosine(x1: Dict[str, float], x2: Dict[str, float]) -> float:
23        # TODO: to be updated
24        up = sum((s1 * x2.get(term, 0) for term, s1 in x1.items()))
25        down_i = sum((s1 ** 2 for term, s1 in x1.items()))
26        down_j = sum((s2 ** 2 for term, s2 in x2.items()))
27        cos = up/(down_i ** 0.5 * down_j ** 0.5)
28        return cos
```

For Task 1, I created 3 for-loops (up, down_i, down_j) each denoting:

up: $\sum_{\forall k}(x_{ik} \cdot x_{jk})$

down_i: $\sqrt{\sum_{\forall k}(x_{ik})^2}$

down_j: $\sqrt{\sum_{\forall k}(x_{jk})^2}$

The function cos = up/(down_i ** 0.5 * down_j ** 0.5) then returns the result of cosine similarity.

**Task 2:**

```python
36    def similar_documents(X: Dict[str, Dict[str, float]], Y: Dict[str, Dict[str, float]]) -> Dict[str, str]:
37        # Feel free to update this function
38        def most_similar_cosine(Y: Dict[str, Dict[str, float]], x: Dict[str, float]) -> str:
39            m, t = -1, None
40            for title, y in Y.items():
41                d = cosine(x, y)
42                if m < 0 or d > m:
43                    m, t = d, title
44            return t
45
46        return {k: most_similar_cosine(Y, x) for k, x in X.items()}
```

For Task 2, I only modified the similar_documents function to only use cosine similarity instead of Euclidean distance. In fact, I implemented the most_similar_cosine function within the similar_documents function for convenience. The reason why I chose to use cosine similarity is because that cosine similarity measures the vectors with respect to the origin.

In VSM, the two vector spaces are both in high dimensions. The property of cosine similarity makes it only subject to the common terms between two vectors, whereas Euclidean distance performs poorly on high dimensional data.

In the vectors we are comparing, they both have high dimensions and mainly are the same stories reconstructed in different story-telling orders. Therefore, it is the common terms between vectors that matter more than other factors, in which case cosine similarity is known to have far better performance than Euclidean distance.