

3-3

SIMILARITY

Similarity between Two Sets

- Jaccard Index

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}.$$

(If A and B are both empty, we define $J(A, B) = 1$.)

Similarity between Two Real-valued Vectors

- Euclidean distance

The **Euclidean distance** between points \mathbf{p} and \mathbf{q} is the length of the line segment connecting them ($\overline{\mathbf{pq}}$).

In Cartesian coordinates, if $\mathbf{p} = (p_1, p_2, \dots, p_n)$ and $\mathbf{q} = (q_1, q_2, \dots, q_n)$ are two points in Euclidean n -space, then the distance (d) from \mathbf{p} to \mathbf{q} , or from \mathbf{q} to \mathbf{p} is given by the Pythagorean formula:

$$\begin{aligned} d(\mathbf{p}, \mathbf{q}) &= d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2} \\ &= \sqrt{\sum_{i=1}^n (q_i - p_i)^2}. \end{aligned} \tag{1}$$

Similarity between Two Real-valued Vectors

- Cosine Similarity

Given two vectors of attributes, A and B , the cosine similarity, $\cos(\theta)$, is represented using a dot product and magnitude as

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}, \text{ where } A_i \text{ and } B_i$$

are components of vector A and B respectively.

The resulting similarity ranges from -1 meaning exactly opposite, to 1 meaning exactly the same, with 0 indicating orthogonality (decorrelation), and in-between values indicating intermediate similarity or dissimilarity.

Similarity between Two Real-valued Vectors

- Cosine Similarity

In the case of `information retrieval`, the cosine similarity of two documents will range from 0 to 1, since the term frequencies (`tf-idf` weights) cannot be negative. The angle between two term frequency vectors cannot be greater than 90° .

If the attribute vectors are normalized by subtracting the vector means (e.g., $A - \bar{A}$), the measure is called centered cosine similarity and is equivalent to the `Pearson Correlation Coefficient`.

Similarity between Two Real-valued Vectors

- Pearson Correlation Coefficient

$$W_{a,u} = \frac{\sum_{i=1}^m (r_{a,i} - \bar{r}_a)(r_{u,i} - \bar{r}_u)}{\sqrt{\sum_{i=1}^m (r_{a,i} - \bar{r}_a)^2} \sqrt{\sum_{i=1}^m (r_{u,i} - \bar{r}_u)^2}}$$

The Cosine similarity of user-mean norm
is Pearson correlation.

Similarity between Two Real-valued Vectors

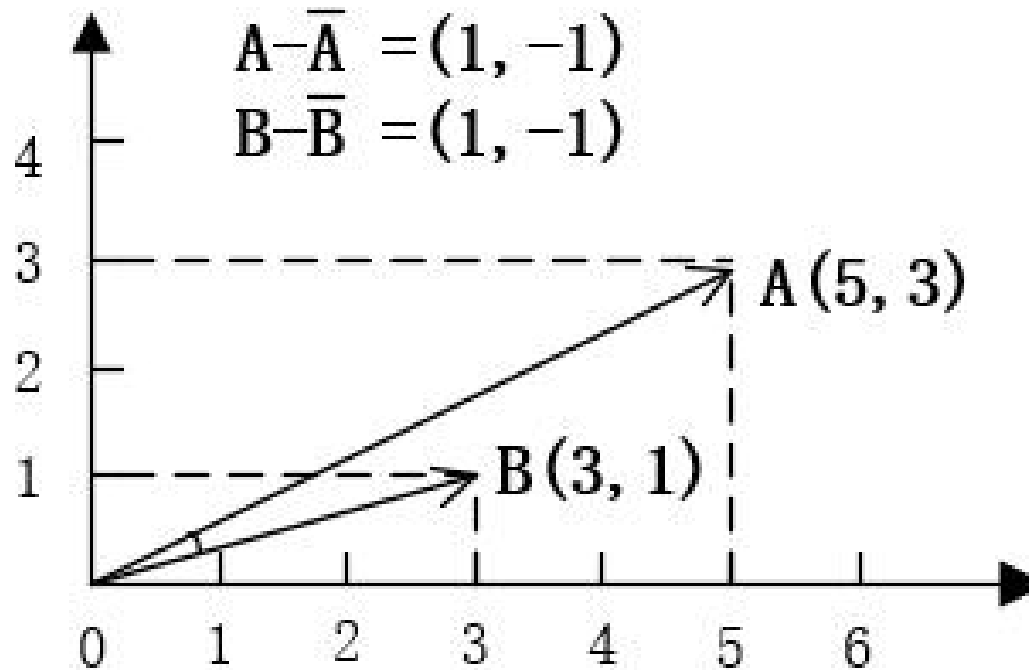
- Pearson Correlation Coefficient (expressed by z-score)

$$W_{a,u} = \frac{\sum_{i=1}^m \left(\frac{r_{a,i} - \bar{r}_a}{\sigma_a} \right) \left(\frac{r_{u,i} - \bar{r}_u}{\sigma_u} \right)}{m}$$

$$\sigma_a = \sqrt{\frac{\sum_{i=1}^m (r_{a,i} - \bar{r}_a)^2}{m}}$$

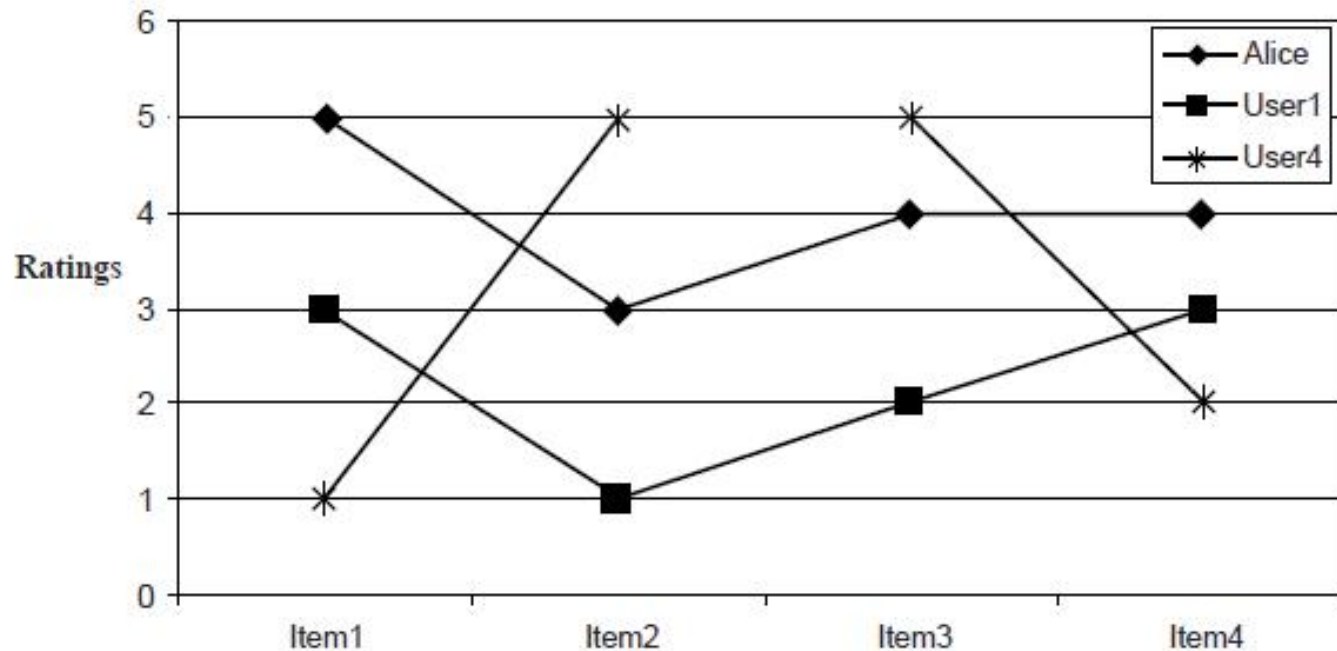
Similarity between Two Real-valued Vectors

- Pearson Correlation Coefficient vs. Cosine Similarity



Similarity between Two Real-valued Vectors

- Pearson Correlation Coefficient vs. Cosine Similarity



Similarity between Two Real-valued Vectors

- Spearman's rank correlation coefficient
 - Covert the raw scores to their ranks
 - Compute the Pearson Correlation Coefficient based on the vectors of ranks
- Examples

V1	5	10	3	8
V2	3	6	9	1

Rank(V1)	3	1	4	2
Rank(V2)	3	2	1	4