

Data Mining Homework 01

Qiu Yihang

March 2023

1 Whether Distances are Metrics

1.1 Jaccard Distance is A Metric

Proof. Recall that Jaccard Distance is

$$d(C_1, C_2) = 1 - \frac{|C_1 \cap C_2|}{|C_1 \cup C_2|}$$

Property 1. For any set A, B , $|A \cap B| \leq |A \cup B|$, i.e. $\frac{|A \cap B|}{|A \cup B|} \leq 1$. Thus, $d(A, B) \geq 0$.

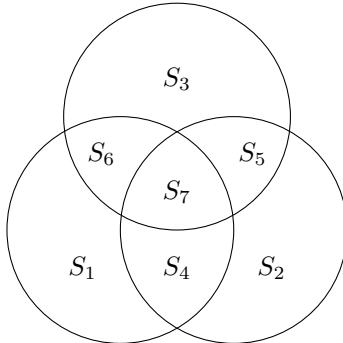
Property 2. $d(A, B) = 0$ iff. $A = B$. The proof is as follows.

$$A = B \implies |A \cup B| = |A \cap B| \implies d(A, B) = 0.$$

$$\begin{aligned} d(A, B) = 0 \implies |A \cup B| &= |A \cap B| \implies \begin{cases} \forall x \in A, x \in A \cup B \Rightarrow x \in A \cap B \Rightarrow x \in B \\ \forall x \in B, x \in A \cup B \Rightarrow x \in A \cap B \Rightarrow x \in A \end{cases} \\ &\implies A \subset B, B \subset A \implies A = B. \end{aligned}$$

Property 3. $d(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|} = 1 - \frac{|B \cap A|}{|B \cup A|} = d(B, A)$.

Property 4. For any set A, B, C , $d(A, B) \leq d(A, C) + d(C, B)$.



Let $A = \bigcup \{S_1, S_4, S_6, S_7\}, B = \bigcup \{S_3, S_5, S_6, S_7\}, C = \bigcup \{S_2, S_4, S_5, S_7\}$, where $\forall i, j \in \{1, 2, \dots, 7\}, S_i \cap S_j = \emptyset$.

First we prove the lemma that for $A, B, C \subseteq X$, it holds that

$$|A \cap C| \cdot |B \cup C| + |A \cup C| \cdot |B \cap C| \leq |C| \cdot (|A| + |B|)$$

The proof is as follows.

$$\begin{aligned} & |A \cap C| \cdot |B \cup C| + |A \cup C| \cdot |B \cap C| \leq |C| \cdot (|A| + |B|) \\ \iff & (|S_6| + |S_7|) \cdot (|S_2| + |S_4| + |S_5| + |S_7| + |S_3| + |S_6|) + \\ & (|S_4| + |S_7|) \cdot (|S_1| + |S_4| + |S_6| + |S_7| + |S_3| + |S_5|) \\ & \leq (|S_3| + |S_5| + |S_6| + |S_7|) \cdot (|S_1| + |S_4| + |S_6| + |S_7| + |S_2| + |S_4| + |S_5| + |S_7|) \\ \iff & (|S_6| + |S_7|) \cdot |S_3| + (|S_4| + |S_7|) \cdot (|S_5| + |S_6|) \\ & \leq (|S_3| + |S_5|) (|S_2| + |S_4| + |S_5| + |S_7| + |S_6|) + (|S_6| + |S_7|) (|S_1| + |S_4| + |S_7|) \\ \iff & 0 \leq |S_3| \cdot |S_2| + |S_3| \cdot |S_4| + |S_3| \cdot |S_5| + |S_2| \cdot |S_5| + |S_5|^2 + \\ & |S_5| \cdot |S_6| + |S_6| \cdot |S_1| + |S_7| \cdot |S_1| + |S_7| \cdot |S_4| + |S_7|^2 \text{ (Always Holds.)} \end{aligned}$$

Moreover, we have $|A \cup C| \cdot |B \cup C| \geq |A \cup B| \cdot |C|$. (The proof is as follows.)

$$\begin{aligned} & |A \cup C| \cdot |B \cup C| \geq |A \cup B| \cdot |C| \\ \iff & (|S_1| + |S_4| + |S_6| + |S_7| + |S_3| + |S_5|) \cdot (|S_2| + |S_4| + |S_5| + |S_7| + |S_3| + |S_6|) \\ & \geq (|S_1| + |S_4| + |S_6| + |S_7| + |S_2| + |S_5|) \cdot (|S_3| + |S_5| + |S_6| + |S_7|) \\ \iff & |S_3| \cdot (|S_2| + |S_4| + |S_5| + |S_7| + |S_3| + |S_6|) + (|S_1| + |S_4| + |S_6| + |S_7| + |S_5|) \cdot (|S_2| + |S_4|) \\ & \geq |S_2| \cdot (|S_3| + |S_5| + |S_6| + |S_7|) \\ \iff & (|S_4| + |S_5| + |S_6| + |S_7|) \cdot (|S_3| + |S_4|) + |S_3|^2 + |S_2| \cdot (|S_1| + |S_4|) \geq 0 \end{aligned}$$

Meanwhile, we have

$$\begin{aligned} |A| + |B| &= |S_1| + |S_4| + |S_6| + |S_7| + |S_2| + |S_5| \\ &= (|S_1| + |S_6| + |S_2| + |S_5|) + (|S_4| + |S_7|) = |A \cup B| + |A \cap B| \end{aligned}$$

Then we know

$$\begin{aligned} |A \cup B| \cdot (|A \cap C| \cdot |B \cup C| + |A \cup C| \cdot |B \cap C|) &\leq |A \cup B| \cdot |C| \cdot (|A| + |B|) \\ &= |A \cup B| \cdot |C| \cdot (|A \cap B| + |A \cup B|) \\ &\leq |A \cup C| \cdot |B \cup C| (|A \cap B| + |A \cup B|) \\ \text{i.e. } \frac{|A \cap C|}{|A \cup C|} + \frac{|B \cap C|}{|B \cup C|} &\leq \frac{|A \cap B|}{|A \cup B|} + 1 \iff 1 - \frac{|A \cap B|}{|A \cup B|} \leq 1 - \frac{|A \cap C|}{|A \cup C|} + 1 - \frac{|C \cap B|}{|C \cup B|} \\ &\iff d(A, B) \leq d(A, C) + d(C, B) \end{aligned}$$

Therefore, Jaccard distance is a metric. ■

1.2 Cosine Distance is Not A Metric

Disproof. Recall that cosine distance for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ is

$$d(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^d x_i y_i}{\sqrt{\sum_{i=1}^d x_i^2} \sqrt{\sum_{i=1}^d y_i^2}}$$

Obvious exists $\mathbf{x} = (1, 1, 1, 1)$ and $\mathbf{y} = (-1, -1, -1, -1)$ s.t. $d(\mathbf{x}, \mathbf{y}) = -1 < 0$.

Thus, cosine distance is not a metric. ■

1.3 Edit Distance is A Metric

Proof. Recall that edit distance (Levenshtein Distance) for two strings x and y is

$$\text{lev}(x, y) = \begin{cases} |x|, & \text{if } |y| = 0 \\ |y|, & \text{if } |x| = 0 \\ \text{lev}(x[1:], y[1:]), & \text{if } x[0] = y[0] \\ 1 + \min \{\text{lev}(x[1:], y), \text{lev}(x, y[1:]), \text{lev}(x[1:], y[1:])\} & \text{otherwise} \end{cases}$$

where $s[k:]$ means the string of all but the first k characters of s .

Property 1. By the definition of $\text{lev}(\cdot)$, obvious for any string x, y , $\text{lev}(x, y) \geq 0$.

Property 2. For any string x , there is

$$\text{lev}(x, x) = \text{lev}(x[1:], x[1:]) = \text{lev}(x[2:], x[2:]) = \dots = \text{lev}(x[|x|:], x[|x|:]) = 0.$$

Property 3. By the definition of edit distance, it is obvious that $\text{lev}(x, y) = \text{lev}(y, x)$.

Property 4. For any string x, y and z , $\text{lev}(x, y) \leq \text{lev}(x, z) + \text{lev}(z, y)$. The proof is as follows.

We prove Property 4 by contradiction.

Assume exist x, y, z s.t. $\text{lev}(x, y) > \text{lev}(x, z) + \text{lev}(z, y)$.

Consider the arguments of $\text{lev}(\cdot)$ when recursively calculate $\text{lev}(x, y), \text{lev}(x, z), \text{lev}(z, y)$. Then we get three sequences, $x \rightarrow a_1 \rightarrow a_2 \rightarrow \dots \rightarrow a_N \rightarrow y$, $x \rightarrow b_1 \rightarrow b_2 \rightarrow \dots \rightarrow b_M \rightarrow z$, $z \rightarrow c_1 \rightarrow c_2 \rightarrow \dots \rightarrow c_K \rightarrow y$, which in fact tell us how to change x to y , x to z , z to y respectively.

By the meaning of edit distance, we know $\text{lev}(x, y), \text{lev}(x, z), \text{lev}(z, y)$ are costs when we change $x \rightarrow a_1 \rightarrow a_2 \rightarrow \dots \rightarrow a_N \rightarrow y$, $x \rightarrow b_1 \rightarrow b_2 \rightarrow \dots \rightarrow b_M \rightarrow z$, $z \rightarrow c_1 \rightarrow c_2 \rightarrow \dots \rightarrow c_K \rightarrow y$ respectively. Moreover, $x \rightarrow a_1 \rightarrow a_2 \rightarrow \dots \rightarrow a_N \rightarrow y$ is the way to change x to y with the minimal cost.

Meanwhile, there exists a sequence $x \rightarrow b_1 \rightarrow b_2 \rightarrow \dots \rightarrow b_M \rightarrow z \rightarrow c_K \rightarrow \dots \rightarrow c_2 \rightarrow c_1 \rightarrow y$ with cost $\text{lev}(x, z) + \text{lev}(z, y) < \text{lev}(x, y)$. **Contradiction!**

In conclusion, edit distance is a metric. ■

1.4 Hamming Distance is A Metric

Proof. Recall that Hamming distance for two strings x and y with the same length L is

$$d(x, y) = \sum_{i=1}^L \mathbb{1}[x_i \neq y_i]$$

Property 1. Obvious for any string x, y , $d(x, y) \geq 0$.

Property 2. For any string x , $d(x, x) = \sum_{i=1}^L 0 = 0$.

Property 3. For any string x and y , $d(x, y) = \sum_{i=1}^L \mathbb{1}[x_i \neq y_i] = \sum_{i=1}^L \mathbb{1}[y_i \neq x_i] = d(y, x)$.

Property 4. For any string x, y and z , $d(x, y) \leq d(x, z) + d(z, y)$. The proof is as follows.

First we prove that $\mathbb{1}[a \neq b] \leq \mathbb{1}[a \neq c] + \mathbb{1}[c \neq b]$ for any char a, b, c by contradiction.

Assume exists a, b, c s.t. $\mathbb{1}[a \neq b] \leq \mathbb{1}[a \neq c] + \mathbb{1}[c \neq b]$.

The only possible case is that $\mathbb{1}[a \neq b] = 1$ and $\mathbb{1}[a \neq c] = \mathbb{1}[c \neq b] = 0$.

Then we get $a \neq b$ and $a = c = b$. **Contradiction!**

Further, we have

$$\begin{aligned} d(x, y) &= \sum_{i=1}^L \mathbb{1}[x_i \neq y_i] \leq \sum_{i=1}^L \mathbb{1}[x_i \neq z_i] + \mathbb{1}[z_i \neq y_i] \\ &= \sum_{i=1}^L \mathbb{1}[x_i \neq z_i] + \sum_{i=1}^L \mathbb{1}[z_i \neq y_i] = d(x, z) + d(z, y) \end{aligned}$$

Therefore, Hamming Distance is a metric. ■

2 Average Distance Between A Pair Of Points

Proof. Let the line of length L be AB . Let the pair of points be M and N . Let $x \triangleq |AM|, y \triangleq |AN|$.

Obvious $x, y \sim \mathcal{U}(0, L)$.

$$\begin{aligned} \mathbb{E}[|x - y|] &= \int_0^L \frac{1}{L^2} \cdot x \cdot \sqrt{2}(L - x) \cdot \frac{\sqrt{2}}{2} dx + \int_0^L \frac{1}{L^2} \cdot y \cdot \sqrt{2}(L - y) \cdot \frac{\sqrt{2}}{2} dy \\ &= \frac{2}{L^2} \int_0^L x(L - x) dx = \frac{2}{L^2} \left(\frac{1}{2} L^3 - \frac{1}{3} L^3 \right) \\ &= \frac{1}{3} L \end{aligned}$$

Thus, the average distance between a pair of points is $\frac{1}{3}L$. ■

3 Eckart-Young-Mirsky Theorem

Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ and $\mathbf{B} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ where \mathbf{S} = diagonal $r \times r$ matrix with $s_i = \begin{cases} \sigma_i & \text{if } i = 1, 2, \dots, k \\ 0 & \text{otherwise} \end{cases}$

Prove that \mathbf{B} is one of the best k -rank approximations to \mathbf{A} in terms of Frobenius norm error.

Proof. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$. To prove the proposition, we just need to prove that

$$\text{for any } \mathbf{C} \in \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \text{rank}(\mathbf{X}) = k\}, \min_{\mathbf{C}} \|\mathbf{A} - \mathbf{C}\|_F = \|\mathbf{A} - \mathbf{B}\|_F.$$

First we prove the following lemma.

Lemma. For any matrix $\mathbf{A} \in \mathbb{R}^{m' \times n'}$ and $\mathbf{W} \in \mathbb{R}^{p' \times n'}$ s.t. $\mathbf{W}^\top \mathbf{W} = \mathbf{I}$, $\|\mathbf{A}\mathbf{W}^\top\|_F = \|\mathbf{A}\|_F$.

For any matrix $\mathbf{A} \in \mathbb{R}^{m' \times n'}$ and $\mathbf{V} \in \mathbb{R}^{m' \times p'}$ s.t. $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$, $\|\mathbf{V}\mathbf{A}\|_F = \|\mathbf{A}\|_F$.

The proof for the lemma is as follows.

$$\begin{aligned} \|\mathbf{A}\mathbf{W}^\top\|_F &= \sqrt{\text{tr}\left(\left(\mathbf{A}\mathbf{W}^\top\right)^\top \mathbf{A}\mathbf{W}^\top\right)} = \sqrt{\text{tr}\left(\mathbf{W}\mathbf{A}^\top \mathbf{A}\mathbf{W}^\top\right)} = \sqrt{\text{tr}\left(\mathbf{A}^\top \mathbf{A}\mathbf{W}^\top \mathbf{W}\right)} \\ &= \sqrt{\text{tr}\left(\mathbf{A}^\top \mathbf{A}\mathbf{I}\right)} = \sqrt{\text{tr}\left(\mathbf{A}^\top \mathbf{A}\right)} = \|\mathbf{A}\|_F \quad \text{where } \mathbf{I} \text{ is identity matrix.} \\ \|\mathbf{V}\mathbf{A}\|_F^2 &= \sqrt{\text{tr}\left((\mathbf{V}\mathbf{A})^\top \mathbf{V}\mathbf{A}\right)} = \sqrt{\text{tr}\left(\mathbf{A}^\top \mathbf{V}^\top \mathbf{V}\mathbf{A}\right)} = \sqrt{\text{tr}\left(\mathbf{A}^\top \mathbf{I}\mathbf{A}\right)} \\ &= \sqrt{\text{tr}\left(\mathbf{A}^\top \mathbf{A}\right)} = \|\mathbf{A}\|_F \quad \text{where } \mathbf{I} \text{ is identity matrix.} \quad \square \end{aligned}$$

By the lemma, we know

$$\|\mathbf{A} - \mathbf{B}\|_F^2 = \|\mathbf{U}(\mathbf{\Sigma} - \mathbf{S})\mathbf{V}^\top\|_F^2 = \|\mathbf{\Sigma} - \mathbf{S}\|_F^2 = \sum_{i=1}^r (\sigma_i - s_i)^2 = \sum_{i=k+1}^r \sigma_i^2.$$

Now we just need to prove that

$$\text{for any } \mathbf{C} \in \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \text{rank}(\mathbf{X}) = k\}, \min_{\mathbf{C}} \|\mathbf{A} - \mathbf{C}\|_F^2 = \sum_{i=k+1}^r \sigma_i^2.$$

By SVD Theorem, we know there exist matrices $\tilde{\mathbf{U}}, \tilde{\mathbf{V}}$ and an $r \times r$ diagonal matrix $\mathbf{\Gamma} = \text{diag}(\gamma_1, \gamma_2, \dots, \gamma_k, 0, \dots, 0)$ s.t. $\mathbf{C} = \tilde{\mathbf{U}}\mathbf{\Gamma}\tilde{\mathbf{V}}^\top$. By the lemma, we know $\|\mathbf{C}\|_F = \sum_{i=1}^r \gamma_i^2$.

Let $\hat{\mathbf{U}} = \mathbf{U}^\top \tilde{\mathbf{U}} = (\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots, \hat{\mathbf{u}}_r)^\top$, $\hat{\mathbf{V}}^\top = \mathbf{V}^\top \tilde{\mathbf{V}} = (\hat{\mathbf{v}}_1, \hat{\mathbf{v}}_2, \dots, \hat{\mathbf{v}}_r)$.

Obvious $\hat{\mathbf{U}}$ is orthogonal since $\hat{\mathbf{U}}^\top \hat{\mathbf{U}} = \tilde{\mathbf{U}}^\top \mathbf{U}\mathbf{U}^\top \tilde{\mathbf{U}} = \tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} = \mathbf{I}$. Similarly, $\hat{\mathbf{V}}$ is orthogonal.

Then we have

$$\begin{aligned} \langle \mathbf{A}, \mathbf{C} \rangle &= \text{tr}\left(\mathbf{V}\mathbf{\Sigma}^\top \mathbf{U}^\top \tilde{\mathbf{U}}\mathbf{\Gamma}\tilde{\mathbf{V}}^\top\right) = \text{tr}\left(\mathbf{\Sigma}^\top \mathbf{U}^\top \tilde{\mathbf{U}}\mathbf{\Gamma}\tilde{\mathbf{V}}^\top \mathbf{V}\right) = \text{tr}\left(\mathbf{\Sigma}^\top \hat{\mathbf{U}}\mathbf{\Gamma}\hat{\mathbf{V}}^\top\right) \\ &= \sum_{i=1}^k \sigma_i \gamma_i \hat{\mathbf{u}}_i^\top \hat{\mathbf{v}}_i \leq \sum_{i=1}^k \sigma_i \gamma_i \|\hat{\mathbf{u}}_i\| \|\hat{\mathbf{v}}_i\| = \sum_{i=1}^k \sigma_i \gamma_i. \end{aligned}$$

Thus, we have

$$\begin{aligned}
\|\mathbf{A} - \mathbf{C}\|_F^2 &= \text{tr} \left((\mathbf{A} - \mathbf{C})^\top (\mathbf{A} - \mathbf{C}) \right) = \text{tr} \left(\mathbf{A}^\top \mathbf{A} - \mathbf{A}^\top \mathbf{C} - \mathbf{C}^\top \mathbf{A} + \mathbf{C}^\top \mathbf{C} \right) \\
&= \text{tr} \left(\mathbf{A}^\top \mathbf{A} \right) - 2\langle \mathbf{A}, \mathbf{C} \rangle + \text{tr} \left(\mathbf{C}^\top \mathbf{C} \right) = \|\mathbf{A}\|_F^2 - 2\langle \mathbf{A}, \mathbf{C} \rangle + \|\mathbf{C}\|_F^2 \\
&\geq \sum_{i=1}^r \sigma_i^2 - 2 \sum_{i=1}^k \sigma_i \gamma_i + \sum_{i=1}^k \gamma_i^2 = \sum_{i=k+1}^r \sigma_i^2 + \sum_{i=1}^k (\sigma_i - \gamma_i)^2 \\
&\geq \sum_{i=k+1}^r \sigma_i^2.
\end{aligned}$$

Therefore, $\mathbf{B} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ is one of the best k -rank approximations to \mathbf{A} in terms of Frobenius norm error. ■

4 Average Jaccard Similarity of Two Randomly-Sampled Sets

Solution. Let $J(A, B)$ be the Jaccard similarity of set A and B .

Totally, there are $\binom{n}{m} \cdot \binom{n}{m}$ possible pairs of (S, T) .

It is obvious that $\max\{0, 2m - n\} \leq |S \cap T| \leq m$.

When $|S \cap T| = k$, there are $\binom{n}{k} \cdot \binom{n-k}{m-k} \cdot \binom{n-m}{m-k}$ possible pairs of (S, T) .

In this case, the Jaccard similarity of S and T is $\frac{k}{2m-k}$.

Thus,

$$\mathbb{E}[J(S, T)] = \sum_{k=\max\{0, 2m-n\}}^m \frac{\binom{n}{k} \cdot \binom{n-k}{m-k} \cdot \binom{n-m}{m-k}}{\binom{n}{m}^2} \cdot \frac{k}{2m-k}.$$
■