

Course: CSC522, Automated Learning and Data Analysis

Homework 1

Student: Xusheng Xiao

1. **Problem 1 (6 points).** Please classify the following attributes as binary, discrete, or continuous. Also classify them as qualitative (nominal or ordinal) or quantitative (interval or ratio). Some cases may have more than one interpretation, so briefly indicate your reasoning then in doubt.
 - (a) Time in terms of AM or PM
Discrete and interval.
 - (b) Brightness as measured by a light meter
Continuous and ratio.
 - (c) Brightness as measured by peoples judgments
Discrete and ordinal. People may specify it like "dark, bright, very bright". Since it still has an way to order the brightness, it should be "discrete and ordinal"
 - (d) Angles as measured in degrees between 0 and 360.
Discrete and ratio.
 - (e) Bronze, Silver, and Gold medals as awarded at the Olympics
Discrete and nominal.
 - (f) Height above sea level.
Continuous and ratio.
 - (g) Number of patients in a hospital
Discrete and ratio.
 - (h) ISBN numbers for books. (Look up the format on the Web.)
Discrete and nominal.
 - (i) Ability to pass light in terms of the following values: opaque, translucent, transparent.
Discrete and ordinal.
 - (j) Military rank
Discrete and ordinal.
 - (k) Distance from the center of campus.
Continuous and ratio.
 - (l) Density of a substance in grams per cubic centimeter.
Continuous and ratio.
 - (m) Coat check number. (When you attend an event and you give your coat in exchange to a number that you can use to claim your coat when you leave.)
Discrete and nominal.

2. **Problem 2 (5 points).** For the following vectors, x and y, calculate the indicated similarity or distance measures.

(a) $X=(1,1,1,1)$, $Y=(2,2,2,2)$: cosine, correlation coefficient, Euclidean distance

cosine distance:

$$\begin{aligned}x \cdot y &= 1 * 2 + 1 * 2 + 1 * 2 + 1 * 2 = 8 \\||x|| &= \sqrt{1 * 1 + 1 * 1 + 1 * 1 + 1 * 1} = 2 \\||y|| &= \sqrt{2 * 2 + 2 * 2 + 2 * 2 + 2 * 2} = 4 \\cos(x, y) &= \frac{x \cdot y}{||x|| ||y||} = \frac{8}{2 * 4} = 1\end{aligned}$$

correlation coefficient:

$$\begin{aligned}\bar{x} &= \frac{1}{4} \sum_{k=1}^4 x_k = 1, \bar{y} = \frac{1}{4} \sum_{k=1}^4 y_k = 2 \\covariance(x, y) &= s_{xy} = \frac{1}{4-1} \sum_{k=1}^4 (x_k - \bar{x})(y_k - \bar{y}) = 0 \\standard_deviation(x) &= s_x = \sqrt{\frac{1}{4-1} \sum_{k=1}^4 (x_k - \bar{x})^2} = 0 \\standard_deviation(y) &= s_y = \sqrt{\frac{1}{4-1} \sum_{k=1}^4 (y_k - \bar{y})^2} = 0 \\corr(x, y) &= \frac{covariance(x, y)}{standard_deviation(x) * standard_deviation(y)} = \frac{s_{xy}}{s_x * s_y} = \frac{0}{0}\end{aligned}$$

$$\begin{aligned}\text{Euclidean distance} &= \sqrt{(1-2)^2 + (1-2)^2 + (1-2)^2 + (1-2)^2} \\&= 2\end{aligned}$$

(b) $X=(0,1,0,1)$, $Y=(1,0,1,0)$: cosine, correlation coefficient, Euclidean distance, Jaccard coefficient

cosine distance:

$$\begin{aligned}x \cdot y &= 0 * 1 + 1 * 0 + 0 * 1 + 1 * 0 = 0 \\||x|| &= \sqrt{0 * 0 + 1 * 1 + 0 * 0 + 1 * 1} = \sqrt{2} \\||y|| &= \sqrt{1 * 1 + 0 * 0 + 1 * 1 + 0 * 0} = \sqrt{2} \\cos(x, y) &= \frac{x \cdot y}{||x|| ||y||} = \frac{0}{\sqrt{2} * \sqrt{2}} = 0\end{aligned}$$

correlation coefficient:

$$\begin{aligned}\bar{x} &= \frac{1}{4} \sum_{k=1}^4 x_k = 0.5, \bar{y} = \frac{1}{4} \sum_{k=1}^4 y_k = 0.5 \\covariance(x, y) &= s_{xy} = \frac{1}{4-1} \sum_{k=1}^4 (x_k - \bar{x})(y_k - \bar{y}) = -\frac{1}{3} \\standard_deviation(x) &= s_x = \sqrt{\frac{1}{4-1} \sum_{k=1}^4 (x_k - \bar{x})^2} = \sqrt{\frac{1}{3}} \\standard_deviation(y) &= s_y = \sqrt{\frac{1}{4-1} \sum_{k=1}^4 (y_k - \bar{y})^2} = \sqrt{\frac{1}{3}}\end{aligned}$$

$$\text{corr}(x, y) = \frac{\text{covariance}(x, y)}{\text{standard_deviation}(x) * \text{standard_deviation}(y)} = \frac{s_{xy}}{s_x * s_y} = -1$$

$$\text{Euclidean distance} = \sqrt{(0-1)^2 + (1-0)^2 + (0-1)^2 + (1-0)^2} = 2$$

$$\text{Jaccard coefficient} = \frac{f_{11}}{f_{01} + f_{10} + f_{11}} = \frac{0}{2 + 2 + 0} = 0$$

(c) X=(0,-1,0,1), Y=(1,0,-1,0): cosine, correlation, Euclidean distance

cosine distance:

$$\begin{aligned} x \cdot y &= 0 * 1 + -1 * 0 + 0 * -1 + 1 * 0 = 0 \\ \|x\| &= \sqrt{0 * 0 + -1 * -1 + 0 * 0 + 1 * 1} = \sqrt{2} \\ \|y\| &= \sqrt{1 * 1 + 0 * 0 + -1 * -1 + 0 * 0} = \sqrt{2} \\ \cos(x, y) &= \frac{x \cdot y}{\|x\| \|y\|} = \frac{0}{2} = 0 \end{aligned}$$

correlation coefficient:

$$\begin{aligned} \bar{x} &= \frac{1}{4} \sum_{k=1}^4 x_k = 0, \bar{y} = \frac{1}{4} \sum_{k=1}^4 y_k = 0 \\ \text{covariance}(x, y) &= s_{xy} = \frac{1}{4-1} \sum_{k=1}^4 (x_k - \bar{x})(y_k - \bar{y}) = 0 \\ \text{standard_deviation}(x) &= s_x = \sqrt{\frac{1}{4-1} \sum_{k=1}^4 (x_k - \bar{x})^2} = 0 \\ \text{standard_deviation}(y) &= s_y = \sqrt{\frac{1}{4-1} \sum_{k=1}^4 (y_k - \bar{y})^2} = 0 \\ \text{corr}(x, y) &= \frac{\text{covariance}(x, y)}{\text{standard_deviation}(x) * \text{standard_deviation}(y)} = \frac{s_{xy}}{s_x * s_y} = \frac{0}{0} \end{aligned}$$

$$\text{Euclidean distance} = \sqrt{(0-1)^2 + (-1-0)^2 + (0-(-1))^2 + (1-0)^2} = 2$$

(d) X=(1,1,0,1,0,1), Y=(1,1,1,0,0,1): cosine, correlation coefficient, Jaccard coefficient

cosine distance:

$$\begin{aligned} x \cdot y &= 1 * 1 + 1 * 1 + 0 * 1 + 1 * 0 + 0 * 0 + 1 * 1 = 3 \\ \|x\| &= \sqrt{1 * 1 + 1 * 1 + 0 * 0 + 1 * 1 + 0 * 0 + 1 * 1} = 2 \\ \|y\| &= \sqrt{1 * 1 + 1 * 1 + 1 * 1 + 0 * 0 + 0 * 0 + 1 * 1} = 2 \\ \cos(x, y) &= \frac{x \cdot y}{\|x\| \|y\|} = \frac{3}{2 * 2} = \frac{3}{4} \end{aligned}$$

correlation coefficient:

$$\begin{aligned} \bar{x} &= \frac{1}{6} \sum_{k=1}^6 x_k = \frac{2}{3}, \bar{y} = \frac{1}{6} \sum_{k=1}^6 y_k = \frac{2}{3} \\ \text{covariance}(x, y) &= s_{xy} = \frac{1}{6-1} \sum_{k=1}^6 (x_k - \bar{x})(y_k - \bar{y}) = \frac{1}{15} \end{aligned}$$

$$\begin{aligned}
\text{standard_deviation}(x) &= s_x = \sqrt{\frac{1}{6-1} \sum_{k=1}^6 (x_k - \bar{x})^2} = \sqrt{\frac{4}{15}} \\
\text{standard_deviation}(y) &= s_y = \sqrt{\frac{1}{6-1} \sum_{k=1}^6 (y_k - \bar{y})^2} = \sqrt{\frac{4}{15}} \\
\text{corr}(x, y) &= \frac{\text{covariance}(x, y)}{\text{standard_deviation}(x) * \text{standard_deviation}(y)} = \frac{s_{xy}}{s_x * s_y} = \frac{1}{4}
\end{aligned}$$

$$\text{Jaccard coefficient} = \frac{f_{11}}{f_{01} + f_{10} + f_{11}} = \frac{3}{1 + 1 + 3} = \frac{3}{5}$$

(e) $X=(2,-1,0,2,0,-3)$, $Y=(-1,1,-1,0,0,-1)$: cosine, correlation coefficient

cosine distance:

$$\begin{aligned}
x \cdot y &= 2 * -1 + -1 * 1 + 0 * -1 + 2 * 0 + 0 * 0 + -3 * -1 = 0 \\
||x|| &= \sqrt{2 * 2 + -1 * -1 + 0 * 0 + 2 * 2 + 0 * 0 + -3 * -3} = \sqrt{18} \\
||y|| &= \sqrt{-1 * -1 + 1 * 1 + -1 * -1 + 0 * 0 + 0 * 0 + -1 * -1} = 2 \\
\cos(x, y) &= \frac{x \cdot y}{||x|| ||y||} = \frac{0}{\sqrt{18} * 2} = 0
\end{aligned}$$

correlation coefficient:

$$\begin{aligned}
\bar{x} &= \frac{1}{6} \sum_{k=1}^6 x_k = 0, \bar{y} = \frac{1}{6} \sum_{k=1}^6 y_k = -\frac{1}{3} \\
\text{covariance}(x, y) &= s_{xy} = \frac{1}{6-1} \sum_{k=1}^6 (x_k - \bar{x})(y_k - \bar{y}) = 0 \\
\text{standard_deviation}(x) &= s_x = \sqrt{\frac{1}{6-1} \sum_{k=1}^6 (x_k - \bar{x})^2} = \sqrt{\frac{18}{5}} \\
\text{standard_deviation}(y) &= s_y = \sqrt{\frac{1}{6-1} \sum_{k=1}^6 (y_k - \bar{y})^2} = \sqrt{\frac{2}{3}} \\
\text{corr}(x, y) &= \frac{\text{covariance}(x, y)}{\text{standard_deviation}(x) * \text{standard_deviation}(y)} = \frac{s_{xy}}{s_x * s_y} = 0
\end{aligned}$$

3. **Problem 3 (6 points).** Show that the set difference metric given by $D(A, B) := \text{size}(A - B) + \text{size}(B - A)$ satisfies the metric axioms given on pages 70/71 of our textbook. Here, A and B are sets, and $A - B$ indicates the set difference.

Theorem 1 $D(A, B) := \text{size}(A - B) + \text{size}(B - A)$ satisfies the metric axioms, where A and B are sets, and $A - B$ indicates the set difference.

Proof. To prove $D(A, B) := \text{size}(A - B) + \text{size}(B - A)$ satisfies the metric axioms, we need to show three properties of metrics hold for $D(A, B)$.

(a) **Positivity.**

- i. $D(A, A) \geq 0$ for all A . Since $A - B$ indicates the set difference, $A - A = \emptyset \Rightarrow \text{size}(A - A) = 0$.
As a result, $D(A, A) := \text{size}(A - A) + \text{size}(A - A) = 0 + 0 = 0 \Rightarrow D(A, A) \geq 0$.
- ii. $D(A, B) = 0$ only if $A = B$. Suppose $A \neq B$, then $A - B \neq \emptyset \Rightarrow \text{size}(A - B) > 0, \text{size}(B - A) > 0$.
As a result, $D(A, B) := \text{size}(A - B) + \text{size}(B - A) > 0 \Rightarrow D(A, B) \neq 0$.
Since $D(A, A) \geq 0$ for all A , $D(A, B) = 0$ only if $A = B$.
- (b) **Symmetry.**
 $D(A, B) = D(B, A)$ for all x and y . Since $D(A, B) = \text{size}(A - B) + \text{size}(B - A)$, $D(B, A) = \text{size}(B - A) + \text{size}(A - B) = D(A, B)$.
- (c) **Triangle Inequality.**
 $D(A, C) \leq D(A, B) + D(B, C)$ for all set A, B , and C .
By De Morgan, $D(A, B) = \text{size}(A - B) + \text{size}(B - A) = \text{size}(A) + \text{size}(B) - 2\text{size}(A \cap B)$
 $D(B, C) = \text{size}(B - C) + \text{size}(C - B) = \text{size}(B) + \text{size}(C) - 2\text{size}(B \cap C)$
 $D(A, C) = \text{size}(A - C) + \text{size}(C - A) = \text{size}(A) + \text{size}(C) - 2\text{size}(A \cap C)$.
 $D(A, B) + D(B, C) - D(A, C) = \text{size}(A) + \text{size}(B) - 2\text{size}(A \cap B) + \text{size}(B) + \text{size}(C) - 2\text{size}(B \cap C) - \text{size}(A) - \text{size}(C) + 2\text{size}(A \cap C)$
 $\Rightarrow D(A, B) + D(B, C) - D(A, C) = 2\text{size}(B) - 2\text{size}(A \cap B) - 2\text{size}(B \cap C) + 2\text{size}(A \cap C)$
By De Morgan, $\text{size}(B) + \text{size}(A \cap B \cap C) \geq \text{size}(A \cap B) + \text{size}(B \cap C)$.
Since $\text{size}(A \cap C) \geq \text{size}(A \cap B \cap C)$, $\text{size}(B) + \text{size}(A \cap C) \geq \text{size}(A \cap B) + \text{size}(B \cap C)$.
 $\Rightarrow 2\text{size}(B) - 2\text{size}(A \cap B) - 2\text{size}(B \cap C) + 2\text{size}(A \cap C) \geq 0$
 $\Rightarrow D(A, B) + D(B, C) - D(A, C) \geq 0 \Rightarrow D(A, C) \leq D(A, B) + D(B, C)$
Thus, $D(A, C) \leq D(A, B) + D(B, C)$ for all set A, B , and C .

Since all these three properties hold for $D(A, B) := \text{size}(A - B) + \text{size}(B - A)$, $D(A, B)$ satisfies the metric axioms.

■

4. **Problem 4 (6 points).** Describe how you would create visualizations to display information that describes
- Computer networks. Be sure to include both the static aspects of the network, such as connectivity, and the dynamic aspects, such as traffic.
 - The distribution of specific plant and animal species around the world for a specific moment in time.
 - The use of computer resources such as processor time, main memory, and disk for asset of benchmark database programs.

In your answers, please address the following issues:

- Representation: how will you map objects, attributes, and relationships to visual elements?
- Arrangement: are there special considerations that need to be taken into account with respect to how visual elements are displayed? E.g. choice of viewpoint, use of transparency, etc.
- Selection: how will you handle a large number of attributes and data objects?

5. Problem 5 (5 points).

- (a) Describe how a box plot can give information about whether the value of an attribute is symmetrically distributed. What can you say about the symmetry of the distribution of the attributes shown in Figure 3.11 on page 115 of our textbook.
- (b) Compare sepal length, sepal width, petal length, and petal width using Figure 3.12 on page 115.