

CS 513 Knowledge Disc and Data Mining

Mid Term

#1 (10 Points)

Is the following function a proper distance function? Why? Explain your answer.

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_i |x_i - y_i| \right)^3$$

Hint: Measure the distance between (0,0), (0,1) and (1,1)

Solution:

Let us assume that, $X = (0,0)$, $Y = (0,1)$ and $Z = (1,1)$

For any distance function to work the following conditions must be satisfied:

1. $d(x, y) \geq 0$ Non-negativity or separation axiom
2. $d(x, y) = 0 \Leftrightarrow x = y$ Identity of indiscernibles
3. $d(x, y) = d(y, x)$ Symmetry
4. $d(x, z) \leq d(x, y) + d(y, z)$ Subadditivity or triangle inequality

Using given distance function,

$$\begin{aligned} \text{The distance between } X(0,0) \text{ \& } Y(0,1) \Rightarrow d(x, y) &= (|0-0| + |0-1|)^3 \\ &= (0+1)^3 \\ &= (1)^3 \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{The distance between } Y(0,1) \text{ \& } X(0,0) \Rightarrow d(y, x) &= (|0-0| + |1-0|)^3 \\ &= (0+1)^3 \\ &= (1)^3 \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{The distance between } Y(0,1) \text{ \& } Z(1,1) \Rightarrow d(y, z) &= (|0-1| + |1-1|)^3 \\ &= (1+0)^3 \\ &= (1)^3 \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{The distance between } Z(1,1) \text{ \& } Y(0,1) \Rightarrow d(z, y) &= (|1-0| + |1-1|)^3 \\ &= (1+0)^3 \end{aligned}$$

$$= (1)^3$$

$$= 1$$

The distance between Z (1,1) & X (0,0) => $d(z, x)$

$$= (|1 - 0| + |1 - 0|)^3$$

$$= (1 + 1)^3$$

$$= (2)^3$$

$$= 8$$

The distance between X (0,0) & Z (1,1) => $d(x, z)$

$$= (|0 - 1| + |0 - 1|)^3$$

$$= (1 + 1)^3$$

$$= (2)^3$$

$$= 8$$

Checking validity of the distance function properties on the distance values calculated using given distance function.

1. $d(x, y) \geq 0, d(y, x) \geq 0, d(y, z) \geq 0, d(z, y) \geq 0, d(z, x) \geq 0, d(x, z) \geq 0$.

Clearly $d(x, y) \geq 0$ and $d(x, y) = 0 \Leftrightarrow x = y$ are satisfied.

2. $d(x, y) = d(y, x), d(y, z) = d(z, y), d(z, x) = d(x, z)$

Clearly $d(x, y) = d(y, x)$ is satisfied.

3. $d(x, z) = 8, d(x, y) = 1, d(y, z) = 1$

$$d(x, z) \leq d(x, y) + d(y, z)$$

$$8 \leq 1 + 1$$

$8 \leq 2$ which is false. So, condition 4 failed

$$d(z, x) = 8, d(z, y) = 1, d(y, x) = 1.$$

$$d(z, x) \leq d(z, y) + d(y, x)$$

$$8 \leq 1 + 1$$

$8 \leq 2$ which is false. So, condition 4 failed here as well.

As per above calculations and observations, given distance function satisfies the first 3 conditions but fails to meet the last condition (Triangle inequality). Therefore, given function is not a proper distance function.

5 (10 Points)

There are three major manufacturing companies that make a product: manufactures A , B , and C . Manufacture A has a 50% market share, and manufacture B has a 30% market share. 5% of A's products are defective, 6% of B's products are defective, and 8% of C's products are defective.

- a) What is the probability that a randomly selected product is defective?
 $P(\text{Defective})$?
- b) What is the probability that a randomly selected product is defective and that it came from A? $P(A \text{ and Defective})$?
- c) What is the probability that a defective product came from B? $P(B/\text{Defective})$?
- d) Are these events (being defective and coming from B) independent? Why?

Solution:

Let's assume there are 1000 items of the product in the market $\Rightarrow N = 1000$

Based on Market Share,

A has 50% of market share. $\Rightarrow N(A) = 50\% \text{ of } 1000 = 500$

B has 30% of market share. $\Rightarrow N(B) = 30\% \text{ of } 1000 = 300$

Remaining are from C $\Rightarrow N(c) = 1000 - 500 - 300 = 200$

Number of defective pieces by manufacturer are as follows:

A's defective products = $N(\text{Defective} \mid A) = 5\% \text{ of } 500 \text{ items} = 25$

B's defective products = $N(\text{Defective} \mid B) = 6\% \text{ of } 300 \text{ items} = 18$

C's defective products = $N(\text{Defective} \mid C) = 8\% \text{ of } 200 \text{ items} = 16$

$$\text{a) } P(\text{Defective}) = (N(\text{Defective} \mid A) + N(\text{Defective} \mid B) + N(\text{Defective} \mid C)) / N$$

$$= (25 + 18 + 16) / 1000 = 59 / 1000 = 0.059 = 5.9\%$$

$$\text{b) } P(A \cap \text{Defective}) = N(\text{Defective} \mid A) / N = 25 / 1000 = 0.025 = 2.5\%$$

$$\text{c) } P(B \mid \text{Defective}) = P(\text{Defective} \mid B) / P(\text{Defective}) = 18 / 59 = 0.3051 = 30.51\%$$

$$\text{d) } P(B) = 300 / 1000 = 0.3$$

$$P(\text{Defective}) = 59 / 1000 = 0.059$$

For events to be independent $\Rightarrow P(B \cap \text{Defective}) = P(B) * P(\text{Defective})$

$$P(B) * P(\text{Defective}) = 0.3 * 0.059 = 0.0177$$

$$P(B \cap \text{Defective}) = 18 / 1000 = 0.018$$

Since, $P(B \cap \text{Defective}) \neq P(B) * P(\text{Defective})$

Therefore, the events are **not independent** of each other.

6 (10 Points)

True or False:

- 1. In data mining, we usually delete all the data rows that contain a missing value to obtain a clean dataset. – False**

Explanation: Though clearing the entire row is an option, we also have other means to clean data set by filling the values with either mean of the column or mode/ most repeated value of the column or using algorithms to determine the most probable value.

- 2. Supervised data mining methods are those that use expert opinions. -- False**

- 3. Usually, low-complexity classifiers/separators need not change very much to accommodate new data points. – True**

Explanation: Low complexity separators have low variance

- 4. The optimal level of model complexity is obtained at the minimum error rate on the training dataset. -- False**

- 5. Data mining processes are autonomous, requiring little or no human oversight. -- False**

Solutions for Question 2 &3 are in R files
Solution for Question 4 is in the Excel file.

******* THE END *******