

Exam 1 Written

Benny Chen

February 22, 2023

Question 1

For each attribute given, classify its type as:

- discrete or continuous AND
- qualitative or quantitative AND
- nominal, ordinal, interval, or ratio

Indicate your reasoning if you think there may be some ambiguity in some cases.

Example: Age in years.

Answer: Discrete, quantitative, ratio.

- (a) Daily user traffic volume at YouTube.com (i.e., number of daily visitors who visited the Web site).
 - (b) Air pressure of a car/bicycle tire (in psi).
 - (c) Credit card number.
- (a) Answer: Discrete, quantitative, ratio.
- (b) Answer: Continuous, quantitative, ratio.
- (c) Answer: Discrete, qualitative, nominal.

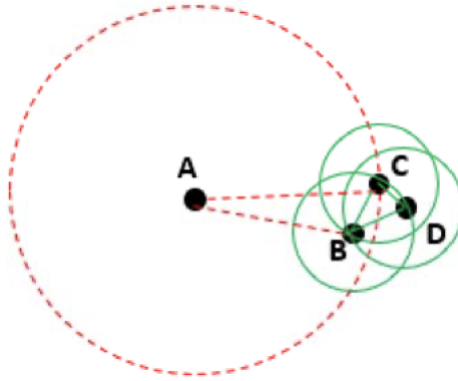
Question 2

Consider the following set of frequent 2-itemsets: $\{p, q\}$, $\{p, r\}$, $\{p, s\}$, $\{p, t\}$, $\{q, r\}$, $\{q, t\}$, $\{r, s\}$, $\{s, t\}$

- (a) List all the candidate 3-itemsets produced during the candidate generation step of the Apriori algorithm.
- (b) List all the candidate 3-itemsets that survive the pruning step of the Apriori algorithm.

- (a) Answer: $\{p, q, r\}, \{p, q, s\}, \{p, q, t\}, \{p, r, s\}, \{p, s, t\}, \{p, r, t\}, \{q, r, t\}, \{q, r, s\}, \{q, s, t\}, \{r, s, t\}$
- (b) Answer: $\{p, q, r\}, \{p, q, s\}, \{p, q, t\}, \{p, r, s\}, \{p, s, t\}, \{p, r, t\}, \{q, r, t\}, \{r, s, t\}$
 Removed: $\{q, r, s\}, \{q, s, t\}$

Question 3



This question aims at finding the local outlier factor (LOF) for the data points (a) A and (b) C from above figure. Suppose $k = 2$. We know that: B, C are the two nearest neighbors to A; B and D are the two nearest neighbors to C. We also know the given distances: $d(A, B) = 4$, $d(A, C) = 5$, $d(B, C) = 1.5$, $d(C, D) = 1$, $d(B, D) = 1.2$.

- (a) LOF (A)
- (b) LOF (C)

$$A \quad k=2$$

B and C closest

$$1 \quad d(AB) = 4$$

$$2 \quad d(AC) = 5$$

k Nearest neighbor

$$A \quad 5 = AC \quad \text{set } AB \quad AC$$

$$B \quad 1.5 = BC \quad \text{set } BD \quad BC$$

$$C \quad 1.5 = BC \quad \text{set } CD \quad BC$$

$$D \quad 1.2 = BD \quad \text{set } CD \quad BD$$

average RD

$$d_2(A) = 5$$

$$d_2(B) = 1.5$$

$$d_2(C) = 1.5$$

$$d_2(D) = 1.2$$

LRD

$$\frac{1}{5}$$

$$\frac{1}{1.5}$$

$$\frac{1}{1.5}$$

$$\frac{1}{1.2}$$

$$RD \quad AB \quad \max(1.5, 4)$$

$$RD \quad AC \quad \max(1.5, 5)$$

$$RD \quad BD \quad \max(1.2, 1.2)$$

$$RD \quad BC \quad \max(1.5, 1.5)$$

$$RD \quad CD \quad \max(1, 1.2)$$

$$RD \quad BC \quad \max(1.5, 1.5)$$

$$RD \quad CD \quad \max(1, 1.5)$$

$$RD \quad BD \quad \max(1.2, 1.5)$$

LOF:

$$\frac{\frac{LRD_2 B}{LRD_2 A} + \frac{LRD_2 C}{LRD_2 A}}{2}$$

$$= \frac{\frac{1.5}{\frac{1}{5}} + \frac{1.5}{\frac{1}{5}}}{2}$$

$$3.333333$$

$$C \quad k=2$$

B and D closest

$$1 \quad d(CD) = 1.2$$

$$2 \quad d(BC) = 1.5$$

k Nearest neighbor

$$A \quad 5 = AC \quad \text{set } AB \quad AC$$

$$B \quad 1.5 = BC \quad \text{set } BD \quad BC$$

$$C \quad 1.5 = BC \quad \text{set } CD \quad BC$$

$$D \quad 1.2 = BD \quad \text{set } CD \quad BD$$

$$RD \quad AB \quad \max(1.5, 4)$$

$$RD \quad AC \quad \max(1.5, 5)$$

$$RD \quad BD \quad \max(1.2, 1.2)$$

$$RD \quad BC \quad \max(1.5, 1.5)$$

$$RD \quad CD \quad \max(1, 1.2)$$

$$RD \quad BC \quad \max(1.5, 1.5)$$

$$RD \quad CD \quad \max(1, 1.5)$$

$$RD \quad BD \quad \max(1.2, 1.5)$$

average RD

$$d_2(A) = 5$$

$$d_2(B) = 1.5$$

$$d_2(C) = 1.5$$

$$d_2(D) = 1.2$$

LRD

$$\frac{1}{5}$$

$$\frac{1}{1.5}$$

$$\frac{1}{1.5}$$

$$\frac{1}{1.2}$$

$$\frac{\frac{1}{1.2}}{\frac{1}{1.5}} + \frac{\frac{1}{1.5}}{\frac{1}{1.5}} =$$

$$2$$

$$1.125$$