

Olabode Alamu

Guide to Engineering Data Science

1498663

Datasets

$$X = [61, 10, 32, 19, 22, 5, 100, 29, 36, 14, 49, 3]$$

$$Y = [2.3, 2.7, 1.7, 1.9, 2.1, 2.8, 1.8, 2.4, 5.9]$$

(a) Find the outliers using Tukey method

Arrange in order

$$X = [3, 5, 10, 14, 19, 22, 29, 32, 36, 49, 61, 100]$$

$$\text{median} = \frac{22 + 29}{2} = 25.5$$

first quartile of $[3, 5, 10, 14, 19, 22]$

$$Q1 = \frac{10 + 14}{2} = 12$$

Third quartile of $[29, 32, 36, 49, 61, 100]$

$$Q3 = \frac{36 + 49}{2} = 42.5$$

Interquartile Range $IQR = Q3 - Q1$

$$= 42.5 - 12 = 30.5$$

$$\begin{aligned}\text{Lower bound of box plot} &= Q1 - 1.5(IQR) \\ &= 12 - 1.5(30.5) \\ &= -33.75\end{aligned}$$

$$\begin{aligned}\text{Upper bound of box plot} &= Q3 + 1.5(IQR) \\ &= 42.5 + 1.5(30.5) \\ &= 88.25\end{aligned}$$

$$\text{Outlier present} = 100$$

$$Y = [2.3, 2.7, 1.7, 1.9, 2.1, 2.8, 1.8, 2.4, 5.9]$$

Arrange in order

$$Y = [1.7, 1.8, 1.9, 2.1, 2.3, 2.4, 2.7, 2.8, 5.9]$$

$$\text{Median} = 2.3$$

$$\text{first quartile in } [1.7, 1.8, 1.9, 2.1, 2.3]$$

$$Q1 = \underline{1.9}$$

$$\text{Third quartile in } [2.3, 2.4, 2.7, 2.8, 5.9]$$

$$Q3 = \underline{2.7}$$

$$IQR = Q3 - Q1 = \underline{0.8}$$

$$\text{lower band of box plot} = Q1 - 1.5(IQR) = 0.7$$

$$\text{Upper band of box plot} = Q3 + 1.5(IQR) = \underline{3.9}$$

$$\text{Outlier present} = \underline{5.9}$$

LOF

Data points

$$A(15.03393, 15.67469)$$

$$B(12.75117, 15.75761)$$

$$C(13.87044, 15.52042)$$

$$D(19.03499, 12.02895)$$

$$E(28.54179, 21.59978)$$

$$k=2$$

Using manhattan distance

$$\text{dist}(A, B) = 2.36568$$

$$\text{dist}(A, C) = 1.317755$$

$$\text{dist}(A, D) = 7.6468$$

$$\text{dist}(A, E) = 19.43294$$

$$\text{dist}(B, C) = 1.3564625$$

$$\text{dist}(B, D) = 10.0124826$$

$$\text{dist}(B, E) = 21.63278$$

$$\text{dist}(C, D) = 8.6566$$

$$\text{dist}(C, E) = 20.750697$$

$$\text{dist}(D, E) = 19.077625$$

$$\text{dist}_k(0)$$

$$\text{dist}_2(A) = 2.36568$$

$$\text{dist}_2(B) = 2.36568$$

$$\text{dist}_2(C) = 1.35646$$

$$\text{dist}_2(D) = 8.65602$$

$$\text{dist}_2(E) = 19.43294$$

k - distance neighborhood of o

$$N_2(A) = (B, C) = 2$$

$$N_2(B) = (A, C) = 2$$

$$N_2(C) = (B, A) = 2$$

$$N_2(D) = (A, B) = 2$$

$$N_2(E) = (B, D) = 2$$

$$\text{Lrc}_k(A) = \frac{1}{|N_2(A)|}$$

$$Lrd_2(A) = \frac{2}{2.36568 + 1.35646} = 0.537325$$

$$Lrd_2(B) = \frac{||N_2(B)||}{reach_{dist_2}(A \leftarrow B) + reach_{dist_2}(C \leftarrow B)}$$

$$reach_{dist_2}(A \leftarrow B) = \max[dist_2(A), dist(A, B)]$$

$$= \max[2.3656811, 2.3656811]$$

$$= 2.3656811$$

$$reach_{dist_2}(C \leftarrow B) = \max[dist_2(C), dist(C, B)]$$

$$= \max[1.3564625, 1.3564625]$$

$$= 1.3564625$$

$$Lrd_2(B) = 0.537325$$

$$Lrd_2(C) = \frac{||N_2(C)||}{reach_{dist_2}(B \leftarrow C) + reach_{dist_2}(A \leftarrow C)}$$

$$reach_{dist_2}(B \leftarrow C) = \max[dist_2(B), dist_2(B, C)]$$

$$= \max[2.36568, 1.3564625]$$

$$= 2.3656811$$

$$\begin{aligned} reachdist_2(A \leftarrow C) &= \max [dist_2(A), dist(A, C)] \\ &= \max [2.38568, 1.3177546] \\ &= 2.3856811 \end{aligned}$$

$$Lrd_2(C) = 0.4227112$$

$$Lrd_2(D) = \frac{|N_2(D)|}{reachdist_2(A \leftarrow D) + reachdist_2(B \leftarrow D)}$$

$$\begin{aligned} reachdist_2(A \leftarrow D) &= \max [dist_2(A), dist(A, D)] \\ &= 7.64680 \end{aligned}$$

$$\begin{aligned} reachdist_2(B \leftarrow D) &= \max [dist_2(B), dist(B, D)] \\ &= 10.01248 \end{aligned}$$

$$Lrd_2(D) = 0.11325487$$

$$Lrd_2(E) = \frac{|N_2(E)|}{reachdist_2(B \leftarrow E) + reachdist_2(D \leftarrow E)}$$

$$\begin{aligned} reachdist_2(B \leftarrow E) &= \max [dist_2(B), dist(B, E)] \\ &= 21.63278 \\ reachdist_2(D \leftarrow E) &= 19.07763 \end{aligned}$$

$$\text{Lrd}_2(E) = 0.0491275$$

$$\text{LOF}(A) = (\text{Lrd}_2(B) + \text{Lrd}_2(C)) \cdot \left[\begin{array}{c} \text{reach}_{\text{int}_2}(B \leftarrow A) \\ + \\ \text{reach}_{\text{int}_2}(C \leftarrow A) \end{array} \right]$$

$$= 3.57339$$

$$\text{LOF}(B) = \text{Lrd}_2(A) + \text{Lrd}_2(B) \cdot \left[\begin{array}{c} \text{reach}_{\text{int}_2}(A \leftarrow B) \\ + \\ \text{reach}_{\text{int}_2}(C \leftarrow B) \end{array} \right]$$

$$= 3.57339$$

and for the remaining factor,

$$\text{LOF}(E) = 26.48537$$

$$\text{LOF}(D) = 18.97755$$

$$\text{LOF}(C) = 5.08456$$

$$\text{LOF}(A) = 3.57339$$

$$\text{LOF}(B) = 3.57339$$