

Suppose that the data for analysis includes the attribute *Number*. The *Number* values for the dataset are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 52, 70. Answer the following questions.

(a) Use smoothing by bin means to smooth the above data, using a bin depth of 3. Illustrate the steps. Comment on the effect of this technique for the given data.

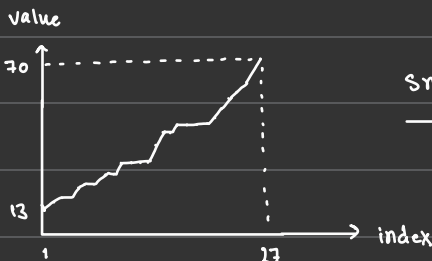
Partition into bins

1: 13, 15, 16	6: 33, 33, 35
2: 16, 19, 20	7: 35, 35, 35
3: 20, 21, 22	8: 36, 40, 45
4: 22, 25, 25	9: 46, 52, 70
5: 25, 25, 30	

Smoothing by bin means:

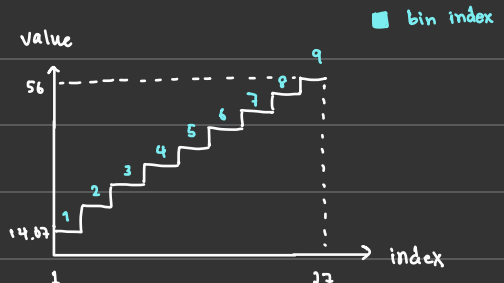
1: 14.67, 14.67, 14.67	6: 33.67, 33.67, 33.67
2: 18.33, 18.33, 18.33	7: 35, 35, 35
3: 21, 21, 21	8: 40.33, 40.33, 40.33
4: 24, 24, 24	9: 56, 56, 56
5: 26.67, 26.67, 26.67	

หลังจากทำการ smoothing data แล้วพบว่า data ที่อยู่ใน bin เดียวกันจะมีค่าเปลี่ยนกันทั้งหมด ช่วงเวลาในคอมพิวเตอร์ของ data ลดลง ลดผลกระทบของค่า outlier ที่ส่อ data สิ่งเกิดได้จาก ค่า 70 ที่ถูกแปลงเป็น 56 โดยสรุปแล้วหลังจากทำการ smoothing เป็นการลด noise ที่สามารถลด data ไปวิเคราะห์ได้จางมากขึ้น



กราฟค่าของ data
ก่อน smoothing

Smoothing
→



กราฟค่าของ data
หลัง smoothing

(b) How might you determine outliers in the data?

$$n \text{ (data size)} = 27$$

$$Q1 \text{ position} = \frac{27+1}{4} = 7$$

$$Q1 \text{ value} = 20$$

$$Q3 \text{ position} = \left(\frac{27+1}{4} \right) \times 3 = 21$$

$$Q3 \text{ value} = 35$$

$$IQR = Q3 - Q1$$

$$= 35 - 20 = 15$$

$$Q1 - 1.5 IQR = 20 - 1.5(15) = -2.5$$

$$Q3 + 1.5 IQR = 35 + 1.5(15) = 57.5$$

$$\text{จาก } Q1 - 1.5 IQR = -2.5 \text{ และ } Q3 + 1.5 IQR = 57.5$$

จะได้ค่าที่ไม้อยู่ในช่วง $[-2.5, 57.5]$ เป็น outlier

ดังนั้น มีค่า outlier ค่าเดียวคือ 70 \neq

(c) What other methods are there for data smoothing?

นอกจากวิธี bin means ก็สามารถใช้วิธีอื่นได้เช่นกัน

1) bin medians : ทุกค่าใน bin จะถูกแทนที่ด้วยค่า median ในแต่ละ bin นั้น ๆ

Partition into bins :

1: 13, 15, 16	5: 25, 25, 30
2: 16, 19, 20	6: 33, 33, 35
3: 20, 21, 22	7: 35, 35, 35
4: 22, 25, 25	8: 36, 40, 45
	9: 46, 52, 70

Smoothing by bin medians :

1: 15, 15, 15	5: 25, 25, 25
2: 19, 19, 19	6: 33, 33, 33
3: 21, 21, 21	7: 35, 35, 35
4: 25, 25, 25	8: 40, 40, 40
	9: 52, 52, 52

2) bin boundaries : ใช้ค่า min, max ของแต่ละ bin ในกรณีนั้น โดยสิ้นหลักการแทนค่าดังนี้

หากค่าที่พิจารณาใกล้กับ min มากกว่า max ก็จะใช้แทนด้วยค่า min แต่ในทางกลับกันหากมีค่าใกล้กับ max มากกว่า ก็จะใช้แทนด้วยค่า max

Partition into bins :

1: 13, 15, 16	5: 25, 25, 30
2: 16, 19, 20	6: 33, 33, 35
3: 20, 21, 22	7: 35, 35, 35
4: 22, 25, 25	8: 36, 40, 45
	9: 46, 52, 70

Smoothing by bin boundaries :

1: 13, 16, 16	5: 25, 25, 30
2: 16, 20, 20	6: 33, 33, 35
3: 20, 20, 22	7: 35, 35, 35
4: 22, 25, 25	8: 36, 36, 45
	9: 46, 46, 70

เลือกแทนด้วยค่า min

(d) Use min-max normalization to transform the value 35 for *Number* onto the range [0.0, 1.0]

$$\text{ดังนั้น } x_{i,\text{new}} = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad x_i \in x$$

พิจารณาค่า 70 ที่แปล outlier ออก

$$\text{แล้ว: } \min(x) = 13, \max(x) = 52$$

$$\text{จะได้ } 35_{\text{new}} = \frac{35 - 13}{52 - 13}$$

$$= \frac{22}{39} \approx 0.56 \quad \#$$

(e) Use z-score normalization to transform the value 35 for *Number*

$$\text{ดังนั้น } x_{i,\text{new}} = \frac{x_i - \text{mean}(x)}{\text{stdev}(x)} \quad x_i \in x$$

พิจารณาค่า 70 ที่แปล outlier ออก

$$\text{แล้ว: } \text{mean}(x) = \frac{\sum x_i}{N} \approx 28.42, \text{stdev}(x) = \frac{\sum (x_i - \mu)^2}{N} \approx 10.37$$

$$\text{จะได้ } 35_{\text{new}} = \frac{35 - 28.42}{10.37} \approx 0.63 \quad \#$$

(g) Compare min-max and z-score normalization techniques.

z-score normalization: เหมาะกับการ normalize ข้อมูลแบบ normal หรือใกล้เคียงกับแบบ normal โดยเป็นการ represent ค่าต่างๆ อยู่ห่างจากค่า mean อย่างไร (หลังจาก normalize ข้อมูลที่สัมพันธ์กับ mean จะมีค่าเท่ากับ 0) และช่วงของค่าที่ normalize แล้วไม่ได้ถูกกำหนดอะไรไว้ ไม่ได้กำหนดและค่อนข้าง sensitive กับ outlier น้อยกว่า min-max normalization

min-max normalization: เหมาะสำหรับข้อมูลที่ไม่เป็น normal ใด หลังจาก normalize แล้ว ค่าที่ได้จะมีช่วงที่ชัดเจน ในกรณีที่มีเฉพาะค่าบวกจะใช้ช่วง $[0,1]$ ส่วนกรณีที่มีข้อมูลมีค่าลบด้วยจะใช้ช่วง $[-1,1]$
และ sensitive กับ outlier มากกว่า z-score normalization เหมาะกับการนำไปใช้ Neural network, Linear regression เป็นต้น