# 1. Data cleaning

missing values, noise, inconsistencies

attributes

i) Ignore the tuple

Temp. GPS speed Auto/Man Sign

tuples

(30°C ... 50 Auto Spd

(31°e ... .. — — )

90°C

ii) fill in manually

(iii) filling in using global constant

$$\propto \longrightarrow \uparrow 999$$

(iv) filling in with mean value

(v). categorized

$25^0 - 35^0$

(vi) most-probable value

Regression

$x=10$

interpolation

$\dfrac{x}{y}$ ⊖

# Noise Elimination

meaningless data /
a deviation from normal pattern

20, 21, 25, 30, 31, (91), 34, 38, 25,

21, 20, 20

i) Binning

Bin1 : [ 20, 21, 25, 30, 31 ] → median → 23 ✓

Bin2 : 31, 34, 38, 91 → (36)

Bin3 : 30, 25, 21, 20, 20 → 20.5 ✓
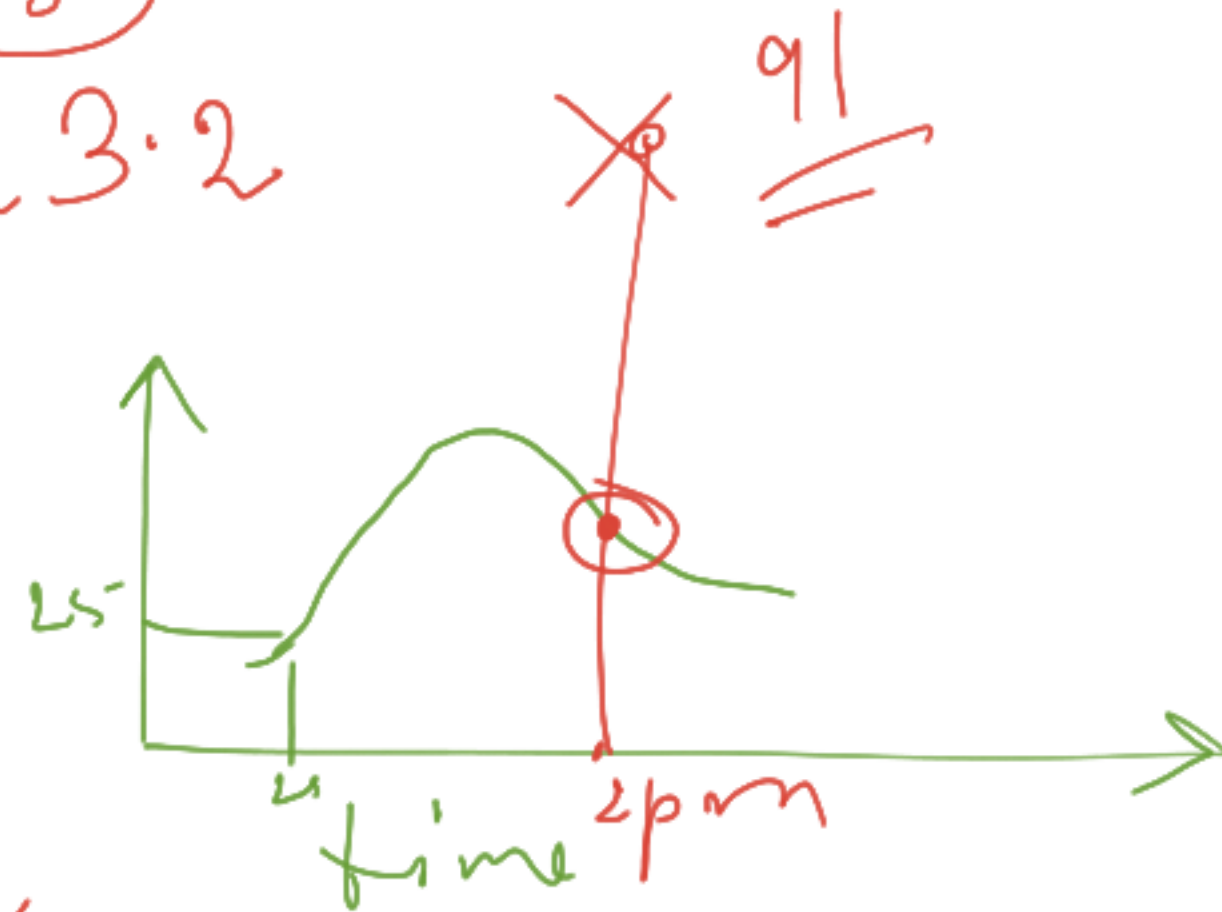
25
34
21

# Replacing by mean of each bin

Bin 1 $\longrightarrow$ 25.4

Bin 2 $\longrightarrow$ (45̄)

Bin 3 $\longrightarrow$ 23.2

ii) Regression

computationaly expensive

better solution

# iii) clustering

| | Red / Blue | |
|---|---|---|
| plan 1 | Red | X Training |
| Plan 2 Plan 3 | Blu Red | |
| | √ Test | |



P.T.O

# 2. Data Integration

i) Redundancy

ii) Conflict Resolution

iii) Data Structure mismatch

$A = 2$
$B = 4$

$A = 3$
$B = 6$

$B = 2A$

$\dfrac{A \ B \ C \ X \ Z}{\text{or}}$

$A \ C \ X \ Y \ Z$

$\dfrac{A \quad \boxed{B} \quad C}{}$

$N \begin{cases} \\ \\ \\ \end{cases}$

$\dfrac{X \quad \textcircled{Y} \quad Z}{}$

# Correlation Coefficient ⟨Pearson's Coefficient⟩

$$r_{A,B} = \frac{\sum\limits_{i=1}^{N} (a_i - \bar{A})(b_i - \bar{B})}{N \, \sigma_A \, \sigma_B}$$

$$= \frac{\sum\limits_{i=1}^{N} (a_i b_i) - N\bar{A}\bar{B}}{N \, \sigma_A \, \sigma_B}$$

$$\frac{\sum a_i}{N} = \bar{A}$$

$$\text{or, } \sum a_i = \bar{A} \cdot N$$

$$r_{A,B} = \frac{\sum\limits_{i=1}^{N} (a_i - \overline{A})(b_i - \overline{B})}{N \sigma_A \sigma_B}$$

$$\frac{\sum a_i}{N} = \overline{A}$$

$$= \frac{\sum\limits_{i=1}^{N} (a_i b_i - a_i \overline{B} - b_i \overline{A} + \overline{A}\,\overline{B})}{N \sigma_A \sigma_B}$$

$$= \frac{\sum\limits_{i=1}^{N} a_i b_i - \overline{B} \sum\limits_{i=1}^{N} a_i - \overline{A} \sum\limits_{i=1}^{N} b_i + \overline{A}\,\overline{B} \sum\limits_{i=1}^{N} 1}{N \sigma_A \sigma_B}$$

$$(1 + 1 + 1 + \cdots) \to N$$

$$= \frac{\sum\limits_{i=1}^{N} a_i b_i - N\overline{A}\,\overline{B} - N \cdot \overline{A}\,\overline{B} + N \cdot \overline{A}\,\overline{B}}{N \sigma_A \sigma_B} = \frac{\sum\limits_{i=1}^{N} a_i b_i - N\overline{A}\,\overline{B}}{N \sigma_A \sigma_B}$$

$$r_{A,B} = \frac{\sum_{i=1}^{N}(a_i - \overline{A})(b_i - \overline{B})}{N \cdot \sqrt{\frac{1}{N}\sum_{i=1}^{N}(a_i - \overline{A})} \cdot \sqrt{\frac{1}{N}\sum_{i=1}^{N}(b_i - \overline{B})}}$$

$$= \frac{\sum_{i=1}^{N}(a_i - \overline{A})(b_i - \overline{B})}{\sqrt{\sum_{i=1}^{N}(a_i - \overline{A})} \cdot \sqrt{\sum_{i=1}^{N}(b_i - \overline{B})}}$$

$\rightarrow$ cov

$\rightarrow$ SD

$-1$ to $+1$

If $\quad r_{A,B} > 0 \quad \longrightarrow \quad$ A & B are posetively correlated
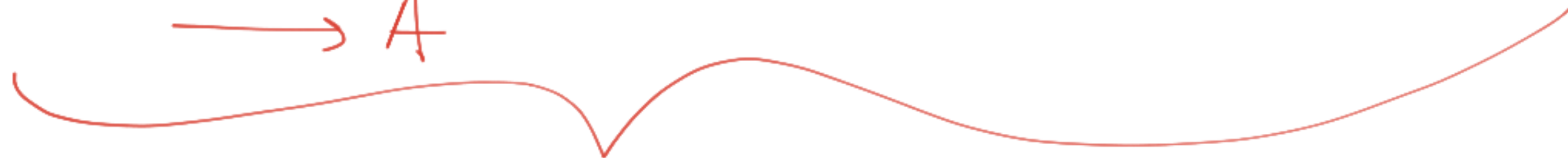
$r_{A,B} < 0 \quad \longrightarrow$ -vely correlated

$r_{A,B} = 0 \quad \longrightarrow$ A & B are independent

is able to detec only linear correlation

$B$ ↑     ⟶ $A$

linear correlation

$A = 2B$

22

| A | 2 | 3 | 4 | 5 | 8 |
|---|---|---|---|---|---|
| B | 4 | 6 | 8 | 10 | 16 |

$\overline{A} = 4.4$

$\overline{B} = 8.8$

$\sigma_A = 2.1$

$\sigma_B = 4.118$

$r_{A,B} = 13.358$

$$\sigma_B = \sqrt{\frac{1}{5} \underbrace{\left[ (4.8)^2 + (2.8)^2 + (0.8)^2 + (1.2)^2 + (7.2)^2 \right]}_{84.8}}$$

$$\sigma_A = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (a_i - \overline{A})^2}$$

$$= \sqrt{\frac{1}{5} \left[ (-2.4)^2 + (1.4)^2 + (0.4)^2 + (0.6)^2 + (3.6)^2 \right]}$$

$$\sigma_A = \sqrt{4.24}$$

$$= 2.1$$

$$r_{A,B} = \frac{(-2.4) \cdot (-4.4) + (-1.4)(-2.8) + (-0.4)(-6.8) +}{5 \times 2.1 \times 4.118}$$

$$(1.2)(0.1) + (7.2)(3.6)$$

$\longrightarrow 2.059$

$\longrightarrow 4.48$

$$= \frac{42.4}{42.39}$$
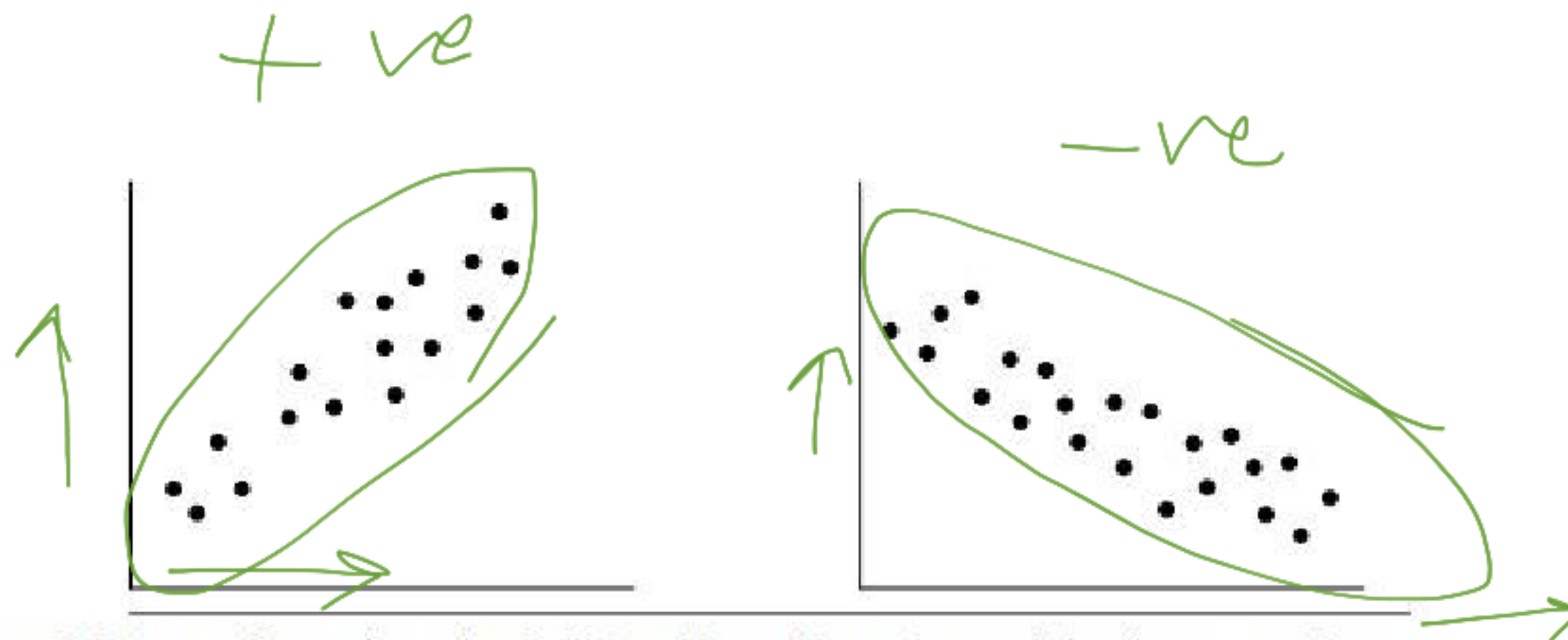
$$\approx 1$$

+ ve

−ve

**Figure 2.8** Scatter plots can be used to find (a) positive or (b) negative correlations between attributes.
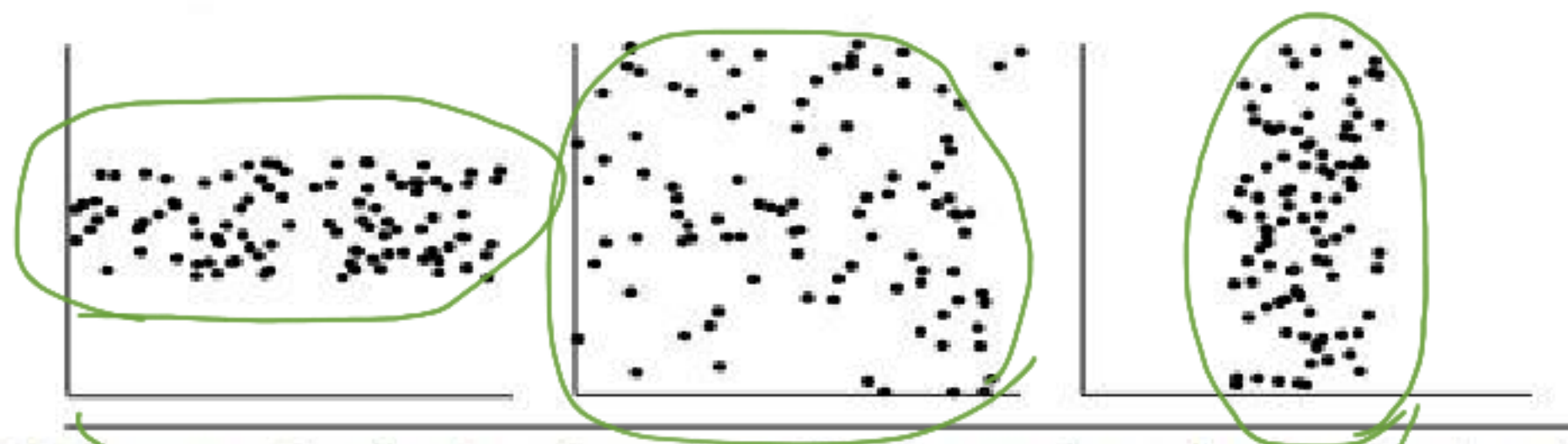
**Figure 2.9** Three cases where there is no observed correlation between the two plotted attributes in each of the data sets.

No correlation

# Next

$x^2$ Data Normalization