

Q) imagine that you have selected data from the all electronics data warehouse for analysis. the data set will be huge. the following data are sorted: 1, 1, 5, 5, 5, 5, 5, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30. (i) partition the dataset using an equal-frequency partitioning method with bin equal to 3 (ii) apply data smoothing using bin means and bin boundaries (iii) plot histogram for the above frequency divided (i) partitioning using equal frequency.

We divide the data set into 3 equal frequency bins, each containing the same number of observations. to calculate the bin boundaries, we count the number of observations in the data set and divide that by the number of bins.

In this case 3. Each bin will contain $40/3 = 13$ observations. the bin boundaries for

equal frequency partitioning method are:

Bin 1: - 1 - 12
 Bin 2: - 12 - 21
 Bin 3: - 21 - 30

ii) Data Smoothing using bin means and bin boundaries:

for data smoothing, we calculate the mean of the each bin and use that as the representative value for all observations in that bin.

Bin 1: Mean $(1+1+5+5+5+5+5+8+8+8+10+10+10)/13 = 6$

Bin 2: Mean $(10+10+10+12+14+14+14+15+15+15+15+15+15)/13 = 15$

Bin 3: Mean $(15+15+15+18+18+18+18+18+18+20+20+20+20+20+20+20+21+21+21+21+25+25+25+25+25+28+28+30+30+30)/13 = 24$

the bin boundaries for smoothed data using bin means etc. are:

Bin 1: - 6 - 12
 Bin 2: - 12 - 21
 Bin 3: - 21 - 30

iii) plotting histogram:

using bin boundaries obtained from either equal frequency (or) data smooth we can plot a histogram by creating bars of the same width that span the bin boundaries and the height of each bar is proportional and the frequency of observations in that bin. the x-axis represents the price of the item and the y-axis represents the frequency of observations.

R-program:-

Load the ggplot2 library
 library(ggplot2)

Create a vector of the Prices data $\leftarrow c(1, 1, 5, 5, 5, 5, 5, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30)$

15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 21, 25, 15, 15, 25, 25, 28, 28, 30, 30, 30

Partition the data using equal frequency partitioning with bin = 3

binned - data $\leftarrow cut(data, breaks = 3, labels = c("1-10", "20-30", "40+"), right = false)$

calculate the bin means

bin - means $\leftarrow apply(data, binned - data, mean)$

calculate the bin boundaries

bin - boundaries $\leftarrow c(-inf, 10, 30, inf)$

Apply data smoothing using bin means and bin boundaries

smoothed - data $\leftarrow cut(data, breaks = bin - boundaries, labels = bin - means, right = false)$

plot the histogram

ggplot(data.frame(smoothed - data), aes(smoothed - data))

-data)) + geom_histogram

(binwidth = 1, color = "black",
fill = "white") + labs (x = "price",
y = "frequency") + ggtitle ("Histogram
of smoothed all electronics prices")

show the plot

plot (ggplot (data.frame
(smoothed-data), aes (smoothed-
data)))

② the following table would be
plotted as (x, y) points, with
the first column being
the x values as number of
mobile phones sold and the
second column being the
 y values as money. ~~to~~ use
the scatter plot for
how many mobile phones
sold.

$x \rightarrow 4 \quad 1 \quad 5 \quad 7 \quad 10 \quad 2 \quad 50 \quad 25 \quad 90 \quad 36$

$y \rightarrow 12 \quad 5 \quad 13 \quad 19 \quad 31 \quad 7 \quad 153 \quad 72 \quad 275 \quad 110$

The scatter plot for the
given table can be plotted
as follows:

$(4, 12), (1, 5), (5, 13), (7, 19)$
 $(10, 31), (2, 7), (50, 153), (25, 72)$
 $(90, 275), (36, 110)$

