

Homework 2

Part 1:

1: Whiskers usually extend interquartile range 1.5 times I beyond 1st and 3rd quartiles / for normal distribution.

For non-normal distribution the range can result of extreme points being considered as outliers.

For instance: heavy tailed data.

~~1.5~~ 1.5, IQR can decide value of point as an outlier

2: Boxplot → symmetry → skewed data is not symmetric skewness; multiple peaks can lead to a distorted conclusions; will mask true outliers.

~~Med-couple-based skewness correction;~~

quantile-based outlier detection / less sensitive to kernel density estimation / distribution

3. mean \rightarrow sensitive to extreme values

average of all values

median \rightarrow middle value

robust to outliers

Boxplots tend to represent typical

value and the spread of the data

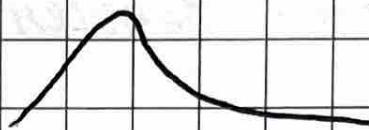
\hookrightarrow minimal impact of extreme
values

Multimodal distribution \rightarrow median can

be in a low-density region; it will
not represent actual center of a
data.

4. Strong right skewness \rightarrow

\rightarrow underlying distribution with long tail extending to the right; most data concentrated to the left



- Variance increase;
- Skewness coefficient > 0

\hookrightarrow statistical models assuming normality become invalid.

log-normal; exponential distribution

5. Boxplots efficiently compare central tendency; spread across groups.

Overlapping distributions of small sample hide multimodality or exaggerate quartile variability

6. few bins \rightarrow merge modes.

\rightarrow oversmooth the data

many bins \rightarrow create noise

\rightarrow disturb pattern

\rightarrow overfit data.

In KDE Bandwidth selection has similar role

too few \rightarrow noisy density estimate

too wide \rightarrow oversmooth data

F. Histograms \rightarrow area of each bar
represents the frequency of data
falling in the bin.

Bar chart \rightarrow height of bar shows
the frequency or value of specific
category.

Bin choice in histogram is more
crucial cause it determines how the
continuous data is grouped \Rightarrow affects
the distribution.

Bar charts show discrete category.

8. Wide Bins in case of bimodal
distribution. If bin width is larger
than separation between two modes,
may show single peak.

KDE preserve true density
structure.

9. density plots \rightarrow probability density function; kernel function

\hookrightarrow smooth representation of distribution

\hookrightarrow Kernel and band width is

crucial \rightarrow poor choice can

lead to over ; under smoothing

Sparse data \rightarrow challenge Bandwidth

small \rightarrow spikiness

large \rightarrow obscure details

10. Area under the curve represents

a probability function / tot. probability of all values

direct comparison of distribution

regardless of sample size ,

Part 2:

1. ECDF plot.

Part a

-5; -2; 0; 3; 4; 5; 5; 6; 7; 7; 8; 9; 9; 10; 12; 15

1) sort \rightarrow ✓

2) $n = 16$ data points

3) -5 1/16

-2 2/16

0 3/16

3 4/16

4 5/16

5 7/16

6 8/16

7 10/16

8 11/16

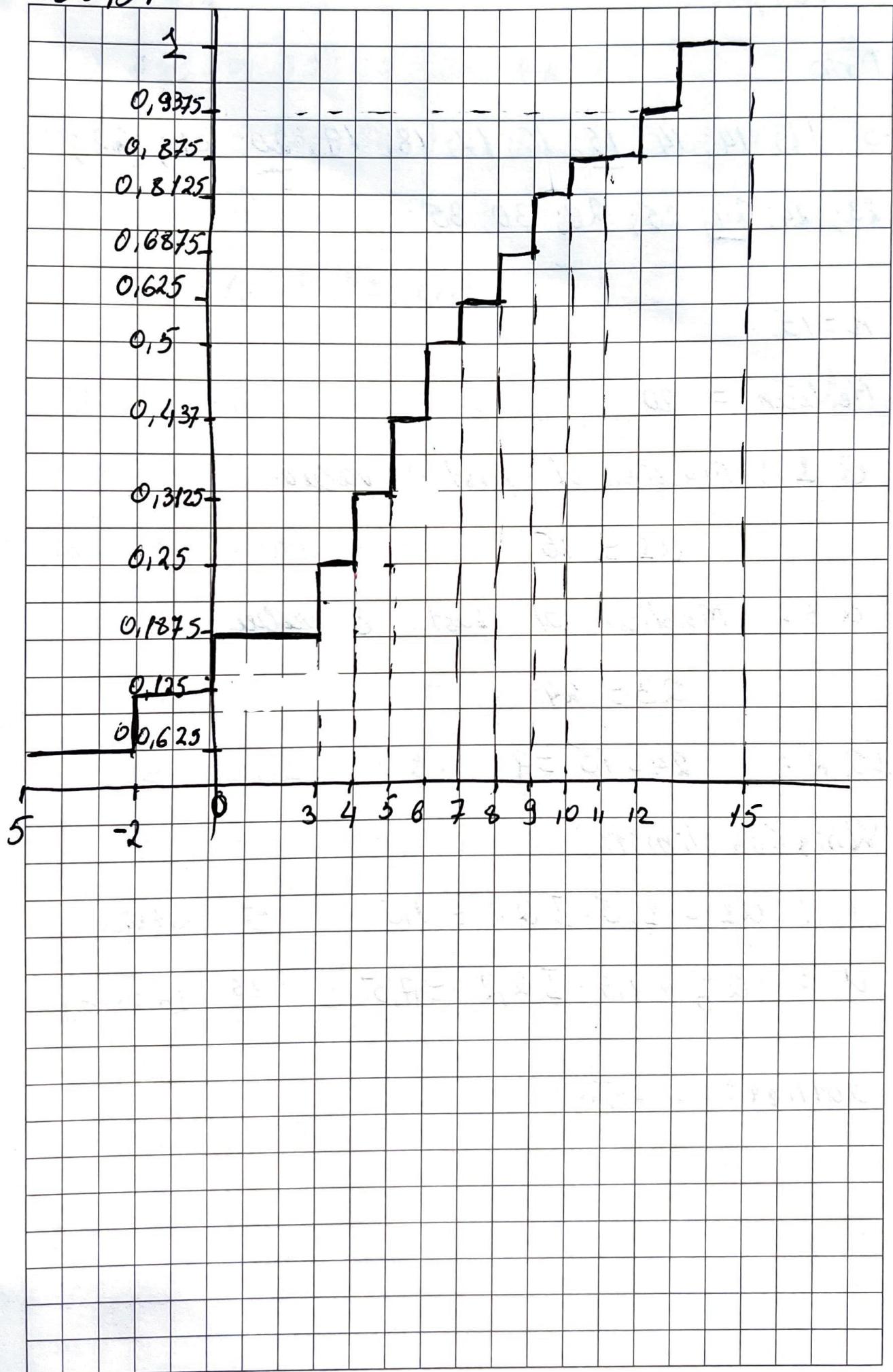
9 13/16

10 14/16

12 15/16

15 16/16

ECDF



2. Boxplot:

Data:

-5; 12; 14; 14; 15; 16; 17; 18; 19; 20; 21; 22;
23; 24; 24; 25; 29; 30; 35.

$n = 15$.

Median = 20.

Q1: Median of first 9 values.

$$Q_1 = 15$$

Q3: Median of last 9 values

$$Q_3 = 24$$

IQR: $24 - 15 = 9$

Whisker limits:

$$L: Q_1 - 1,5 \cdot IQR = 1,5$$

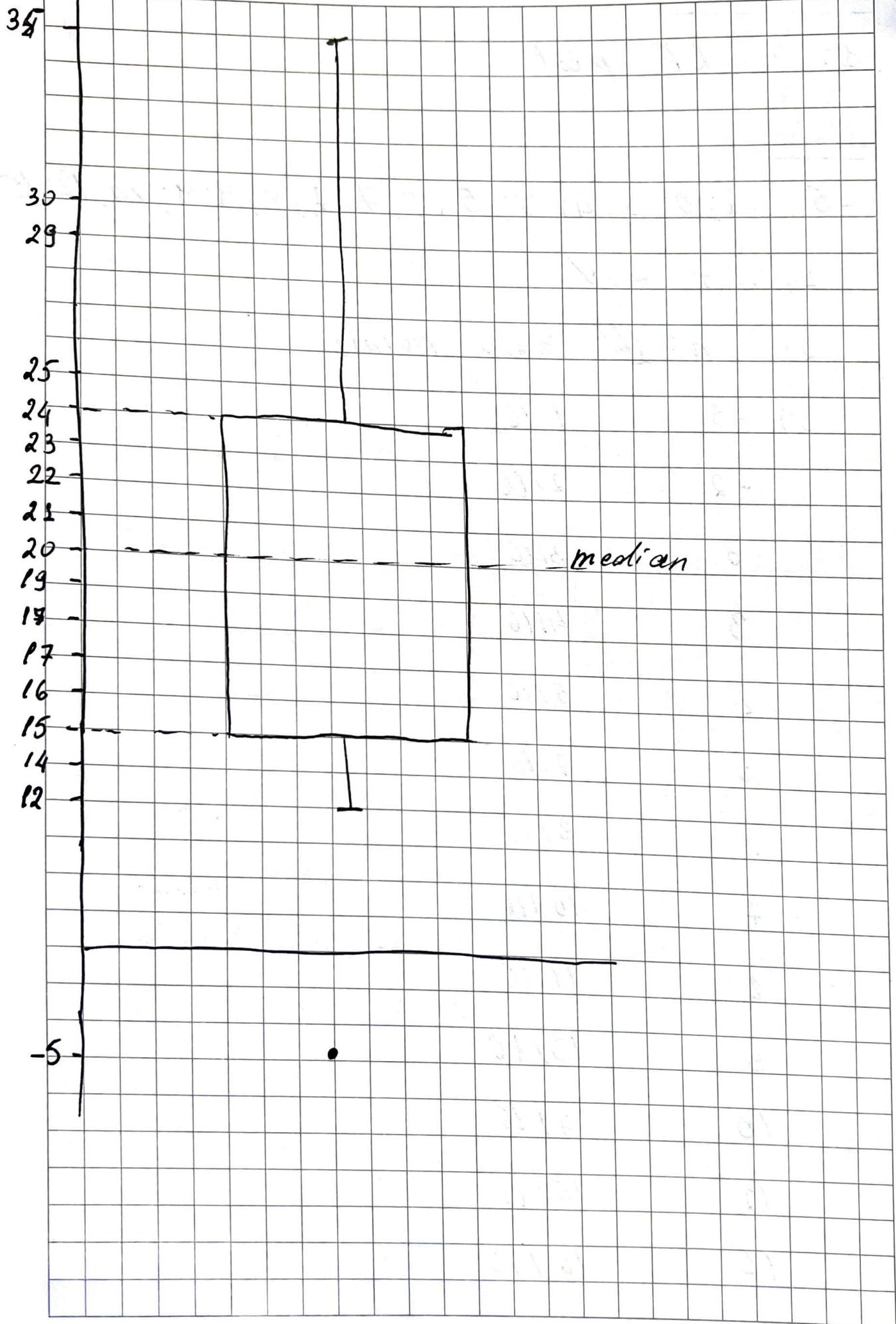
-5 outlier

$$U: Q_3 + 1,5 \cdot IQR = 37,5$$

35 in range

Outlier: $\dots -5 \dots$

Box plot



Histogram

-10; 45; 50; 55; 55; 60; 62; 65; 68; 70; 73;
74; 80; 80; 82; 85; 88; 80 ; 91; 92; 94; 97; 100;
105

$$\text{Range} : 105 - (-10) = 115$$

$$\text{Bin Width} : \frac{\text{Range}}{\text{number of bins}} = \frac{115}{5} = 23$$

frequency

Bin 1:	-10; \rightarrow 13		1
2:	13; \rightarrow 36		0
3:	36 \rightarrow 59		4
4:	59 \rightarrow 82		10
5:	82 \rightarrow 105		9

Histogram

