# Introduction

# to

# Statistics

Class Notes

# Introduction to Statistics

## Introduction

The course will be divided into 3 parts

1) Descriptive statistics
2) Probability
3) Statistical inference.

1) Descriptive Statistics

Idea: We have a set of data and we wish to capture its essence or summarize its main features.

## The sigma notation

We want a convenient way to write a sum of numbers $x_1, x_2, \ldots, x_n$

We use the sigma sign, $\Sigma$ to represent such a sum :

$$\sum_{i=m}^{n} x_i = x_m + x_{m+1} + \cdots + x_n$$

Add the numbers $x_i$, starting with $i = m$ and ending with $i = n$. $(n \geq m)$

$$\underset{i=\boxed{5}}{\overset{\boxed{10}}{\Sigma}} x_i = x_5 + x_6 + \cdots + x_{10}$$

end $\rightarrow$

start $\nearrow$

# Properties of summation $\Sigma$

1) $\sum\limits_{i=1}^{n} c x_i = c \sum\limits_{i=1}^{n} x_i$  if $c$ is a constant.

2) $\sum\limits_{i=1}^{n} (x_i + y_i) = \sum\limits_{i=1}^{n} x_i + \sum\limits_{i=1}^{n} y_i$

3) $\left( \sum\limits_{i=1}^{n} x_i \right)^2 \neq \sum\limits_{i=1}^{n} x_i^2$

$(x_1 + x_2)^2 = x_1^2 + 2 x_1 x_2 + x_2^2 = x_1^2 + x_2^2$

There are three ways of measuring the center

1) Sample mean $\overline{X}$

2) Sample median $m$

3) mode.

## Sample mean.

Given a set of quantities

$$x_1, x_2 \ldots x_n,$$

We call

$$\overline{X} = \sum_{i=1}^{n} x_i / n = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

the sample mean of the $x_i$'s

# Note:

1) $\bar{x}$ "x bar" represents the "center" of a set of numbers in some sense.

2) Note that $\bar{x}$ may not represent any number in the dataset.

$$x_1 = 1 . \quad x_2 = 3.$$

$$\bar{x} = (1+3)/2 = 2 \text{ is not equal to 1 or 3 !}$$

3) $\bar{x}$ is easily influenced by extreme observations (either very large or small)

Ex $\quad x_1 = 1, \; x_2 = 1, \; x_3 = 1, \; x_4 = 1, \; x_5 = 100$

(Income) $\quad \bar{x} = \dfrac{1+1+1+1+100}{5} = 20.8$

does not provide a picture of typical income. values pull towards 100.

## Sample median

Let $x_1, x_2, \ldots, x_n$ be a set of values.

The sample median $m$ is defined as follows.

1) If the number of observations $n$ is odd.

○   rank the observations from smallest to largest. $x_1, x_2, \ldots, x_n \Rightarrow x_1^*, x_2^*, \ldots, x_n^*$

○   then the median $m$ is the middle number after ranking

$$m = x_{\frac{n+1}{2}}^*$$

## Ex

Data: $x_1 = 1$, $x_2 = 5$ $x_3 = 4$ $x_4 = 2$ $x_5 = 3$

After ranking: $x_1^* = 1$ $x_2^* = 2$ $x_3^* = 3$ $x_4^* = 4$ $x_5^* = 5$

$$n = 5 \qquad m = X^*_{\frac{n+1}{2}} = X^*_3 = 3$$

i.e. the sample median is the observation with the property that half of the observations are $\leq$ it and half $\geq$ it.

2) When the number of observations $n$ is even

o  rank the observations as before.

o  identify two middle observations.
   the sample median $m$ is the average of these two observations.

$$m = \frac{X^*_{\frac{n}{2}} + X^*_{\frac{n}{2}+1}}{2}$$

# Ex

## Data:

$$x_1 = 5 \quad x_2 = 2 \quad x_3 = 4 \quad x_4 = 3 \quad x_5 = 1 \quad x_6 = 6$$

## After ranking:

$$x_1^* = 1 \quad x_2^* = 2 \quad x_3^* = 3 \quad x_4^* = 4 \quad x_5^* = 5 \quad x_6^* = 6$$

$$n = 6. \quad m = \frac{x_{\frac{n}{2}}^* + x_{\frac{n}{2}+1}^*}{2} = \frac{x_3^* + x_4^*}{2}$$

$$= \frac{3+4}{2} = 3$$

Note:

1) Sample median also represents a measure of centrality — different possibly from $\overline{X}$

2) It is not easily affected by extreme observations as $\overline{X}$. We say that the sample median is more **robust** than the sample mean.

Thus, if $X_1 = 1$, $X_2 = 2$, $X_3 = 1000$

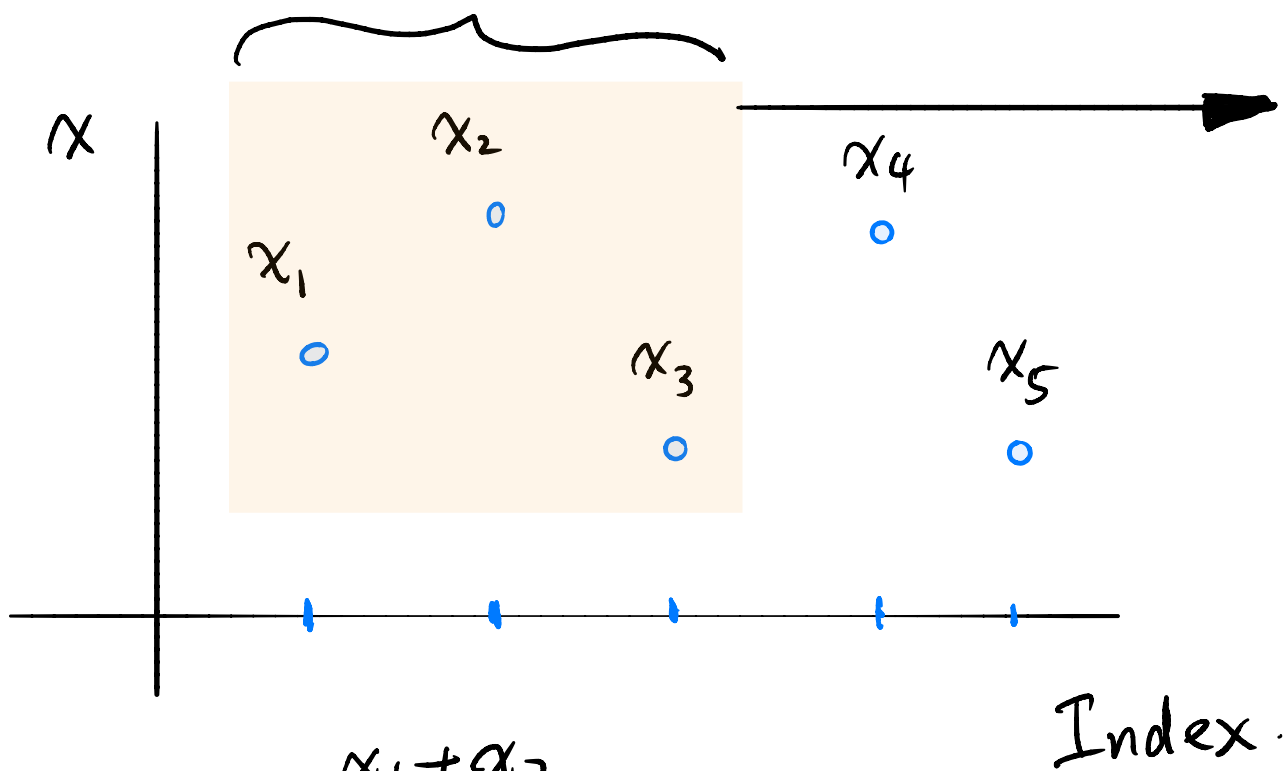$$\overline{X} = \frac{1 + 2 + 1000}{3} \approx 334$$

$m = 2$.

# Applications of Sample mean & median



50-day simple moving average

# How to compute Moving Average (MA) with a window size 3 ?

Moving window size = 3

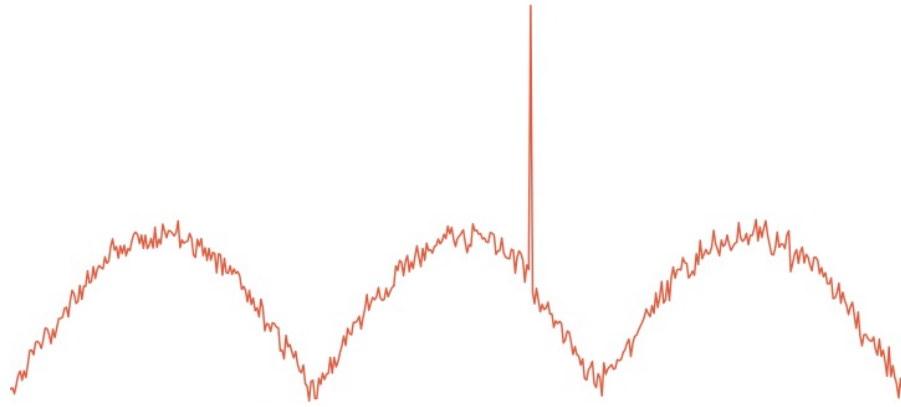$$a_1 = \frac{x_1 + x_2}{2}$$

$$a_2 = \frac{x_1 + x_2 + x_3}{3}$$

$$a_3 = \frac{x_2 + x_3 + x_4}{3}$$
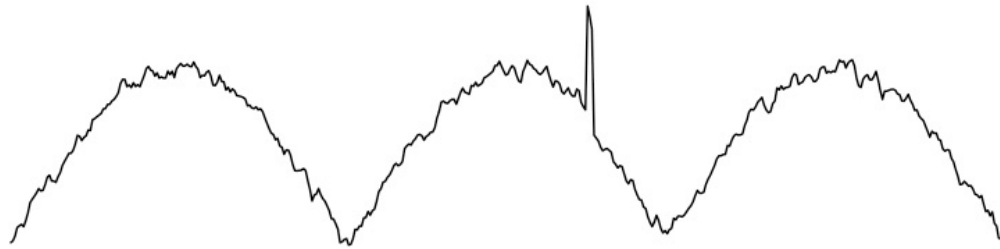
$$a_4 = \frac{x_3 + x_4 + x_5}{3}$$

$$a_5 = \frac{x_4 + x_5}{2}$$

Moving Median can be computed using a similar way
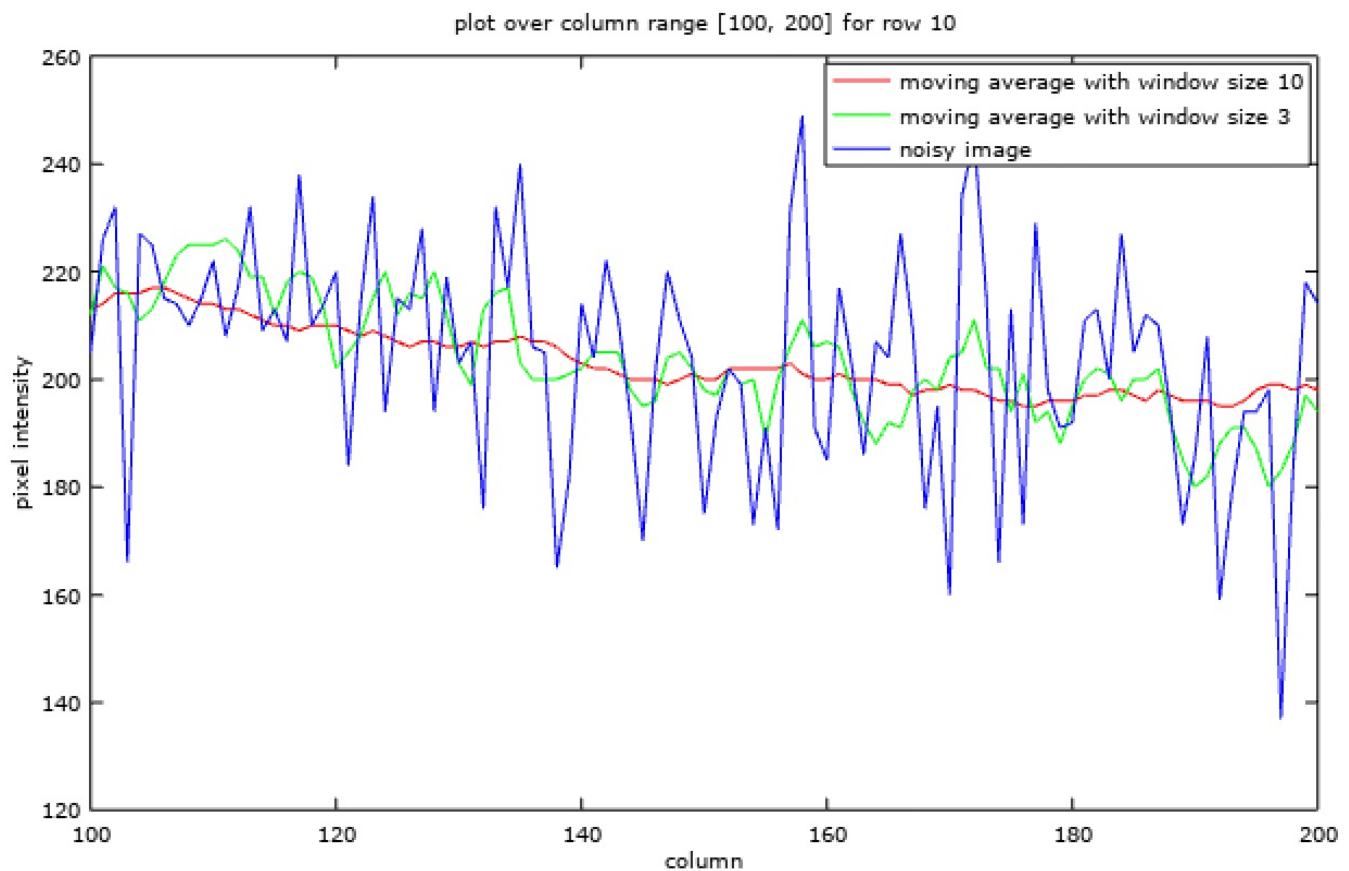
# Original Data:



# Moving Average:



# Moving Median:

Moving Median is less affected by extreme values than Moving Average is

How about the window size ?



plot over column range [100, 200] for row 10

Legend:
- moving average with window size 10
- moving average with window size 3
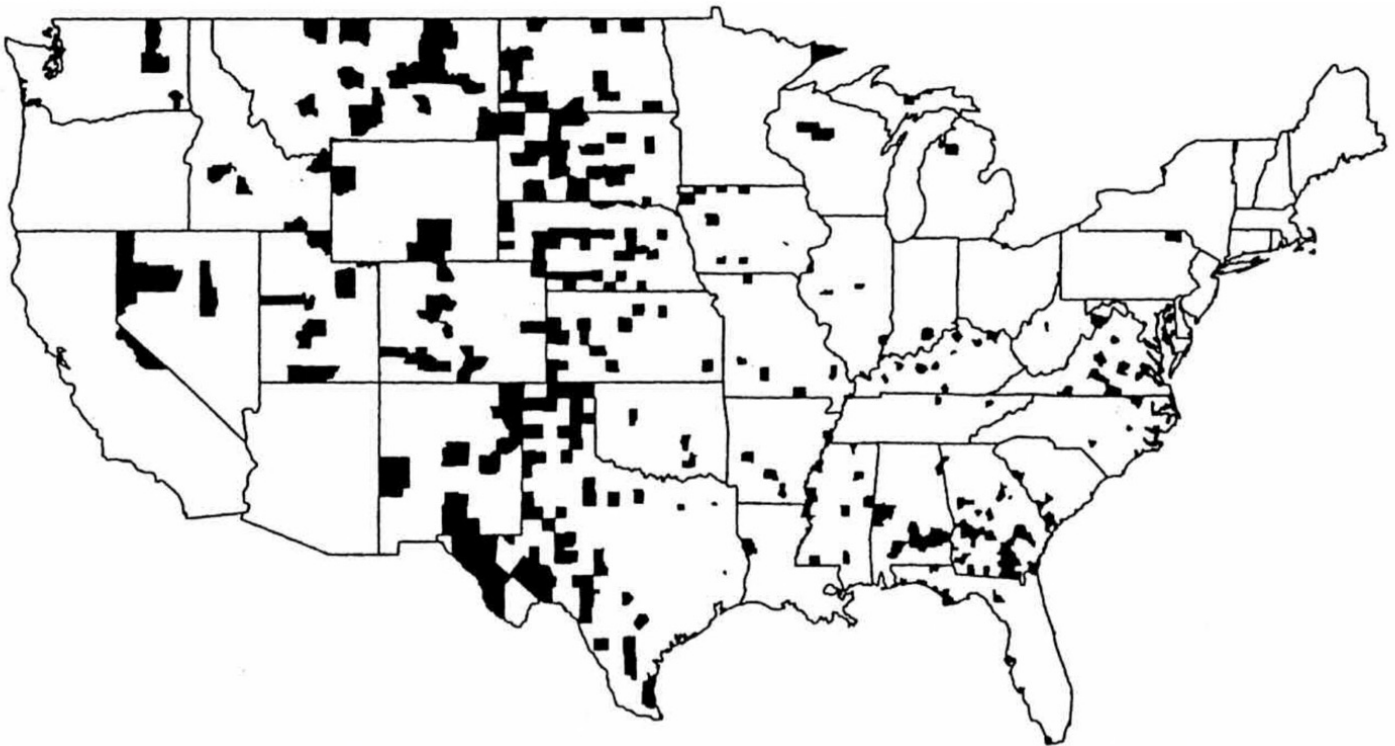- noisy image

x-axis: column
y-axis: pixel intensity

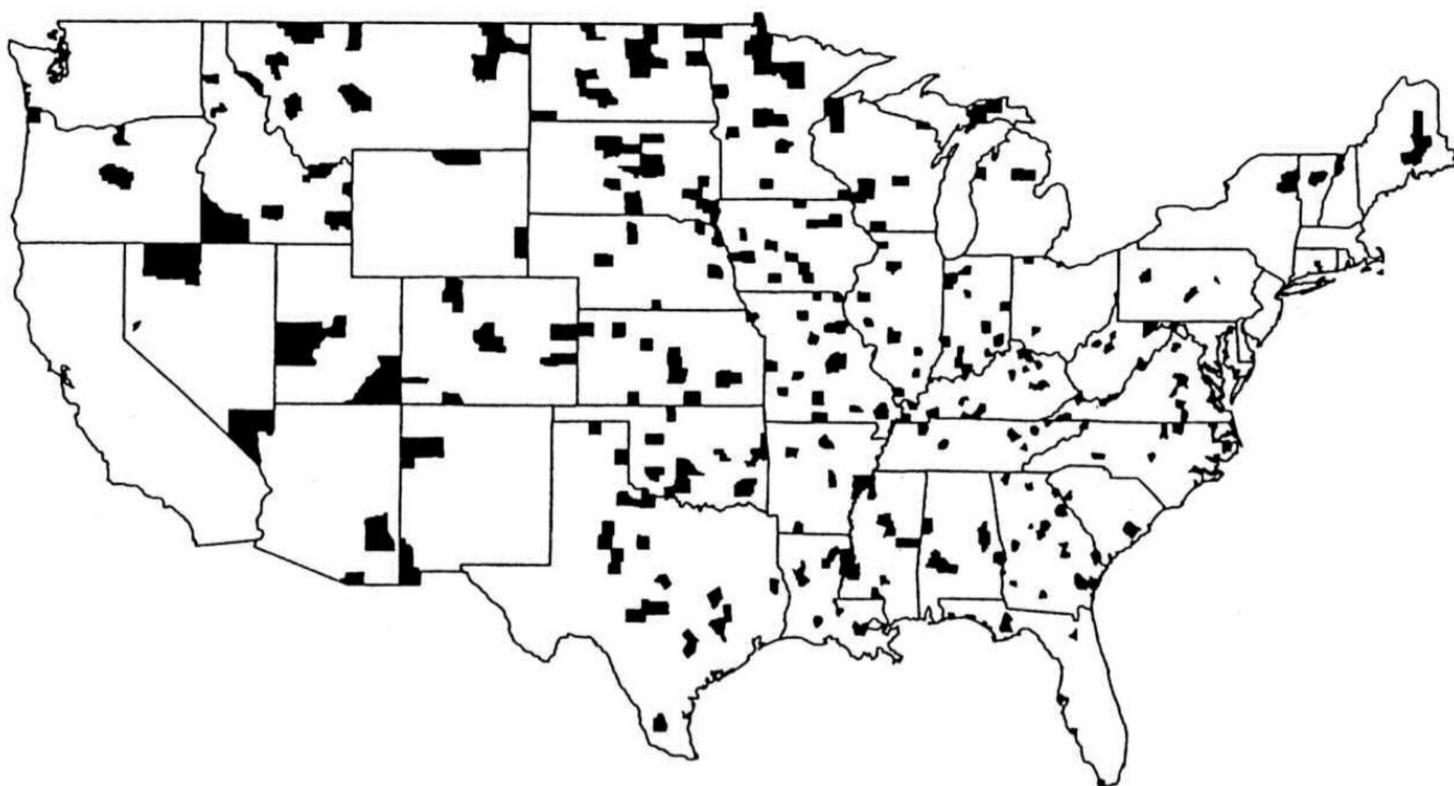The larger the window size, the smoother the curve

In general, the large the sample size, the more stable the mean (and median)
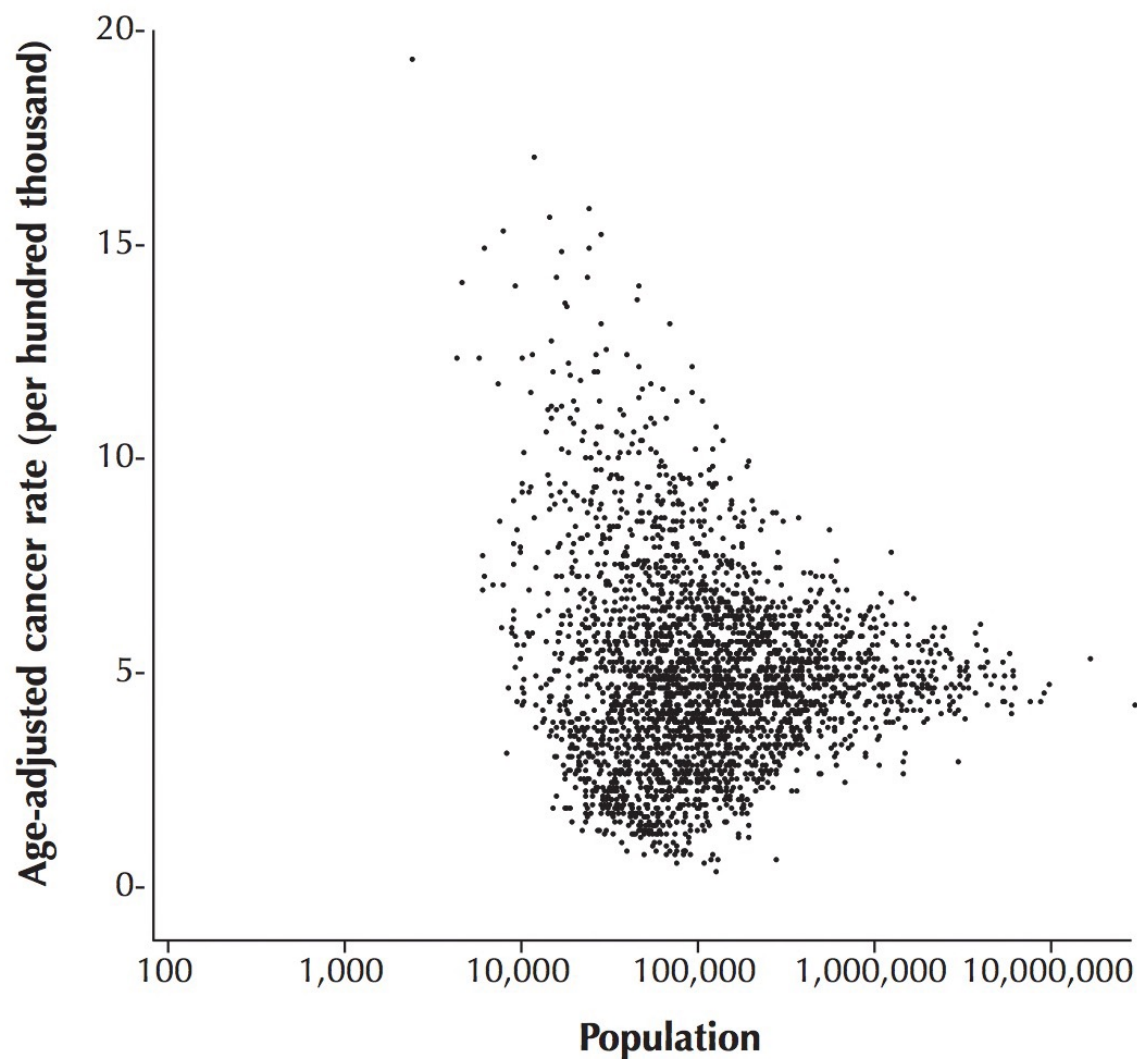
**FIGURE 1.**



The counties with the lowest 10% age-standardized death rates for cancer of the kidney/ureter for U.S. males, 1980-89. Reprinted, by permission, from Andrew Gelman and Deborah Nolan, *Teaching Statistics: A Bag of Tricks* (New York: Oxford University Press, 2002), p. 15.

**FIGURE 2.**



The counties with the highest 10% age-standardized death rates for cancer of the kidney/ureter for U.S. males, 1980-89. Reprinted, by permission, from Andrew Gelman and Deborah Nolan, *Teaching Statistics: A Bag of Tricks* (New York: Oxford University Press, 2002), p. 14.

**FIGURE 3.**



Age-adjusted death rates for cancer of the kidney/ureter for U.S. males, for all U.S. counties, 1980-89, shown as a function of the log of the county population.

## Mode

A much less common measure centrality

Defined as the most commonly occurring observation

(obs. that occurs most often)

e.g. $x_1 = -1.2$, $x_2 = 3.6$  $x_3 = 3.7$

$x_4 = -1.2$

The mode is $-1.2$

If I add $x_5 = 3.6$, then mode becomes

$-1.2$, $3.6$ (two modes)

In addition to a measure of the center, we need to measure variability of the data.

1) range

2) Sample variance

3) Sample standard deviation

## Range

$$\text{Range} = X_{largest} - X_{smallest}$$

## Sample Variance

### EX (Asset Return)

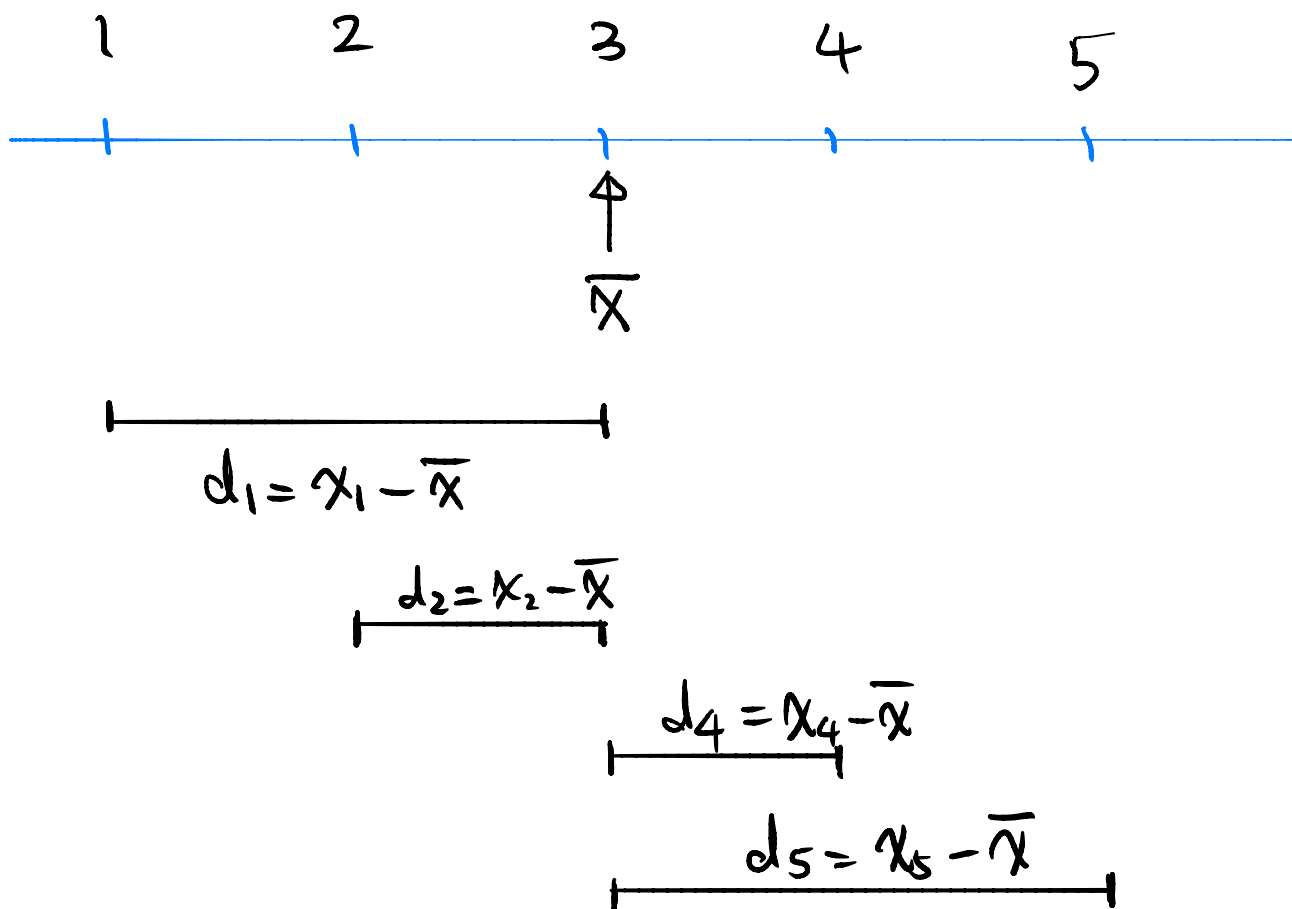| Data 1 | Data 2 |
|---|---|
| 1, 2, 3, 4, 5 | 2, 3, 3, 3, 4 |
| $\overline{x}_1 = 3$ | $\overline{x}_2 = 3$ |

Same averaged return, which one has more variation?

# Definition:

Variance: $s^2 = \dfrac{1}{n-1} \sum\limits_{i=1}^{n} (x_i - \bar{x})^2$

obs. ↑ $(x_i - \bar{x})$

sample mean ↑ $\bar{x}$

1    2    3    4    5

↑
$\bar{x}$

$d_1 = x_1 - \bar{x}$

$d_2 = x_2 - \bar{x}$

$d_4 = x_4 - \bar{x}$

$d_5 = x_5 - \bar{x}$

# Variance of Data 1

$$S_1^2 = \frac{1}{5-1}\left(d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2\right)$$

$$= \frac{1}{5-1}\left((1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2\right)$$

$$= \frac{10}{4} = 2.5$$

# Variance of Data 2

$$S_2^2 = \frac{1}{5-1}\left(d_1^2 + \cdots + d_5^2\right)$$

$$= \frac{1}{5-1}\left((2-3)^2 + (3-3)^2 + (3-3)^2 + (3-3)^2 + (4-3)^2\right)$$

$$= \frac{2}{4} = 0.5$$

Note:

1) $S^2$ is "roughly" ($n-1$ instead of $n$) the <u>averaged</u> <u>squared</u> divation of $x_i$'s from $\bar{X}$.

2) Reason of using $n-1$ will be explained later

3) Alternative formula for computing $S^2$.

$$S^2 = \frac{1}{n-1}\left[\sum_{i=1}^{n} X_i^2 - n\cdot(\bar{X})^2\right]$$

<u>proof</u>: as homework question.

Sketch:
$$S^2 = \frac{1}{n-1}\left[\sum_{i=1}^{n}(X_i - \bar{X})^2\right]$$

$$= \frac{1}{n-1}\left[\sum_{i=1}^{n}(X_i^2 - 2X_i\bar{X} + \bar{X}^2)\right]$$

$$= \frac{1}{n-1}\left[\sum_{i=1}^{n}X_i^2 - \underbrace{\sum_{i=1}^{n}2X_i\bar{X}}_{-2\bar{X}\sum_{i=1}^{n}X_i} + \underbrace{\sum_{i=1}^{n}\bar{X}^2}_{n\bar{X}^2}\right]$$

$$\underset{\parallel}{}$$
$$-2n\bar{X}^2$$

4) Alternatively, use $|x_i - \bar{x}|$ instead of $(x_i - \bar{x})^2$ to eliminate minus signs with deviations. Hence, it results in

$$L(x) = \frac{1}{n-1} \sum_{i=1}^{n} |x_i - \bar{x}|$$

## Sample Standard deviation

Definition :

Standard deviation : $S = \sqrt{S^2}$

$\uparrow$ Variance

Note that the unit of $s^2$ are the same as the units of $x_i^2$. Interpretation of $s^2$ is more difficult. Hence we consider $s = \sqrt{s^2}$ as our measure of spread. s has the

Same unit as $x_i$'s.

## Numerical Example

$x_1 = 4.$  $x_2 = 1$  $x_3 = 3$  $x_4 = 1$  $x_5 = 3$

$x_6 = 1$  $x_7 = 2$  $x_8 = 2$.

Solution:

$$\bar{x} = \frac{1}{8} \cdot \sum_{i=1}^{8} x_i = \frac{4 + 1 + \cdots + 2}{8}$$

$$= \frac{17}{8} = 2.125$$

| $x_2$ | $x_4$ | $x_6$ | $x_7$ | $x_8$ | $x_3$ | $x_5$ | $x_1$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| 1 | 1 | 1 | 2 | 2 | 3 | 3 | 4 |

$m = 2.$        $mode = 1$

$$s^2 = \frac{1}{8-1} \left[ \sum_{i=1}^{8} x_i^2 - 8\bar{x}^2 \right]$$

$$= \frac{1}{8-1} \left[ \sum_{i=1}^{8} x_i^2 - 8 \cdot (2.125)^2 \right]$$

$$= 1.26$$

$$\sum_{i=1}^{8} x_i^2 = 4^2 + 1^2 + \ldots + 2^2 = 45$$

$$\delta = \sqrt{\delta^2} = \sqrt{1.26} = 1.126.$$

range $= 4 - 1 = 3$

# Types of Data

1) Quantitative.

Definition : measurements recorded on numerical scale.

Ex: Tempreture, unemployment rate, interest rate, exam scores.

2) Qualitative.

Definition : measurements that cannot be measured on a numerical scale. only be classified into categories.

Ex: species, car type, gender

Note : Qualitative data are often coded in arbitrary numerical

values, but can not be meaningfully "+" "−" "×" "÷".

| Car Type | code |
|----------|------|
| SUV | 1 |
| sedan | 2 |
| Truck | 3 |
| .. | |

# Describing Qualitative Data

Original Data :

| Subject | Education |
|---------|-----------|
| 1 | BS |
| 2 | Master |
| 3 | PhD |
| 4 | PhD |
| 5 | BS |

Education is Qualitative, can be classified into 3 classes , BS. Master . PhD

Such data can be summarized in three ways

1) Class Frequency

the number of obs that falls in a particular class

2) Class relative frequency

$$\text{class relative freq.} = \frac{\text{class freq.}}{n}$$

n is the total number of obs.

3) Class percentage.

$$\text{class percentage} = (\text{class relative freq.}) \times 100$$
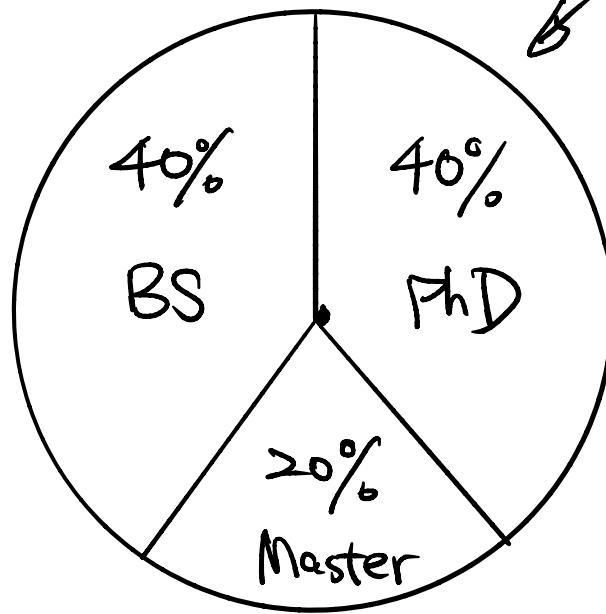
# Table:

$$n = 5$$

|         | Freq | Relative Freq. | Percentage |
|---------|------|----------------|------------|
| BS      | 2    | 2/5            | 40%        |
| Master  | 1    | 1/5            | 20%        |
| PhD     | 2    | 2/5            | 40%        |
| Total   | 5    | 1              | 100%       |

# Graphs

## Bar Graph

Freq.   or   Relative Freq.   Percentage.



BS        Master        PhD

# Pie Chart



relative Freq.
percentage.

40%
BS

40%
PhD

20%
Master

In January 1971 the Gallup poll asked: "A proposal has been made in Congress to require the U.S. government to bring home all U.S. troops before the end of this year. Would you like to have your congressman vote for or against this proposal?"

Guess the results, for respondents in each education category, and fill out this table (the two numbers in each column should add up to 100%):

| | Adults with: | | | |
| | Grade school education | High school education | College education | Total adults |
|---|---|---|---|---|
| % for withdrawal of U.S. troops (doves) | | | | 73% |
| % against withdrawal of U.S. troops (hawks) | | | | 27% |
| Total | 100% | 100% | 100% | 100% |

| $P_1$ | $P_2$ | $P_3$ |
|---|---|---|
| $1 - P_1$ | $1 - P_2$ | $1 - P_3$ |

the proportions of adults with grade school, high school, and college education are $\lambda_1$, $\lambda_2$, $\lambda_3$ respectively. They must satisfy

$$\lambda_1 + \lambda_2 + \lambda_3 = 1$$

Total proportion of adults against war is 0.73
Therefore.

$$\lambda_1 \times P_1 + \lambda_2 \times P_2 + \lambda_3 \times P_3 = 0.73$$

Assume

$$\lambda_1 = \lambda_2 = \lambda_3 = 1/3$$

Then $P_1$, $P_2$ and $P_3$ must satisfy.

$$(P_1 + P_2 + P_3)/3 = 0.73$$

In January 1971 the Gallup poll asked: "A proposal has been made in Congress to require the U.S. government to bring home all U.S. troops before the end of this year. Would you like to have your congressman vote for or against this proposal?"

Guess the results, for respondents in each education category, and fill out this table (the two numbers in each column should add up to 100%):

| | Adults with: | | | |
| | Grade school education | High school education | College education | Total adults |
|---|---|---|---|---|
| % for withdrawal of U.S. troops (doves) | 80% | 75% | 60% | 73% |
| % against withdrawal of U.S. troops (hawks) | 20% | 25% | 40% | 27% |
| Total | 100% | 100% | 100% | 100% |

► "Back in Vietnam days, the anti-war movement spread from the intelligentsia into the rest of the population, eventually paralyzing the country's will to fight." -*The Economist* (2000)

# Describing Quantitative Data

- Dot plot

- Histogram.

## Dot Plot

30.3

*rounding*

30.3 → 30.5 "to nearest half."

```
+    +    +    +    +    +
30.0  .1   .2   .3   .4   .5
```

Data:

30.5 ✎ 30.3   30.6   30.9   30.8   31.1
        31.2   31.0   31.7   31.4   32.1



```
      •
      •
      •
  •   •
  •   •   •
          •   •
——+———+———+———+———+——
  30  30.5 31.0 31.5 32.0
```

# Histogram (Discretization of quantitative data)

- Divide the range of the data into $K$ equal intervals or bins.

- Allocate data into bins

- Count the frequency of observations in bins.

## Ex

Data: 30.1  30.2  30.5  30.6

31.5  31.3

32.3  32.7  32.9

33.1

$n = 10$ obs.

| Class Interval | Freq | Rel. Freq |
|---|---|---|
| (30, 31] | 4 | 4/10 |
| (31, 32] | 2 | 2/10 |
| (32, 33] | 3 | 3/10 |
| (33, 34] | 1 | 1/10 |
| | 10 | 1. |

Freq. or Rel. Freq.

If you use the relative frequency then the total area represented by the histogram is 1.

If the size of the interval is 2 instead of 1.

| Interval | Freq | Rel. Freq. |
| --- | --- | --- |
| (30, 32] | 6 | 6/10 |
| (32, 34] | 4 | 4/10 |
| | 10 | 1 |

Note:

1) The interval size is reasonably important.

If size is too small, too detailed does not reflect distribution

If size is too large, the shape will be lost and it conveys little information.

2) Rule for selecting the interval size.

| # obs | # of intervals |
| --- | --- |
| < 25 | 5 – 6 |
| 25 – 50 | 7 – 14 |
| > 50 | 15 – 20 |

3) The scale is important.
When comparing hists. make sure the scales are the same.

Need a separate interval for the value "zero"

Zero - inflated
right - skewed.



If the insurance company charge everyone at the price of $\bar{X}$, can they cover the basic cost and won't lose money in a long run?

# 4) You can get an idea of the shape of the dist. of the data.

## Skewness

A dataset is said to be skewed if one tail of the distribution has more extreme observations than the other tail.

### Symmetric

## right skewed



## left skewed



## Detecting Skewness

- ○ Check the shape of the histogram.

- ○ or, compare the $\bar{x}$ and $M$.

  mean        median

Rule:

$$\overline{x} > M \implies \text{right skewed}$$

$$\overline{x} < M \implies \text{left skewed}$$

Ex : symmetric

4 , 4 , 4 , 5, 5 ,5. 5 , 6. 6, 6

symmetric



4  5  6

$$\overline{x} = 5$$

$$M = 5$$

$$\overline{x} = M$$

**Ex:** right skewed.

$4, 4, 4, 5.5, 5.5, 6.6, 6, 100$

right skewed



4 5 6          100

$$\overline{x} = 13.6$$

$$M = 5$$

$$\underline{\overline{x} > M}$$

**Ex:** left skewed.

$-100, 4, 4, 4, 5.5, 5.5, 6.6, 6$

left skewed



$-100$      4   5   6

$\overline{x} = -4.5$

$M = 5$

$\overline{x} < M.$

# Is this data skewed ?



Total Insurance Claim Amount (in $1000) Per Policy Year

$$M = 0$$

$$\overline{x} > 0$$

$$\overline{x} > M$$

The data is right skewed

# Measures of Relative Standing

Describing the quantitative location of a particular measurement within a dataset.

- Percentile ranking (median)

- Z - score (mean)



80 is 25th percentile

25% of students their grades ≤ 80

## Percentile

Definition:

p-th percentile ($p \in [0, 100]$) is the

number such that $p\%$ of the observations fall below the number and $(100-p)\%$ fall above it.

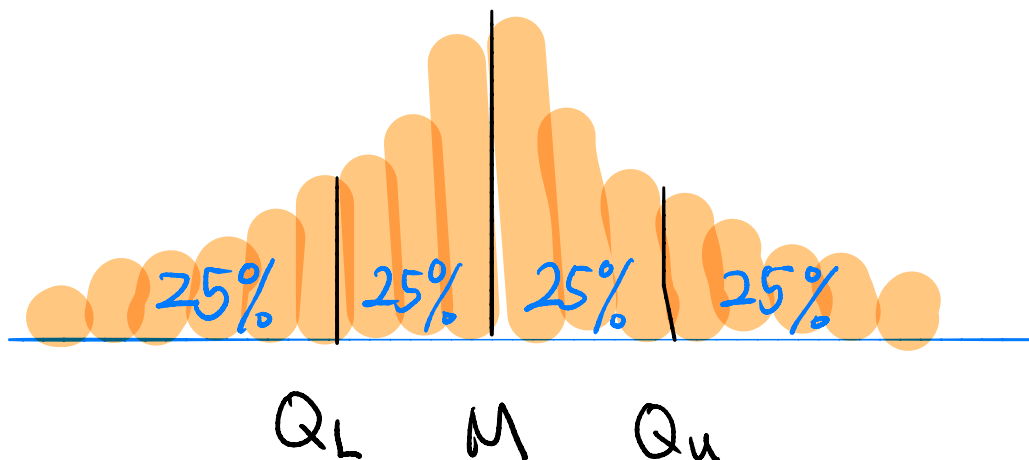## Quartile (special percentile)

○ lower quartile

$$Q_L = 25\text{-th percentile}$$

○ median

$$M = 50\text{-th percentile}.$$
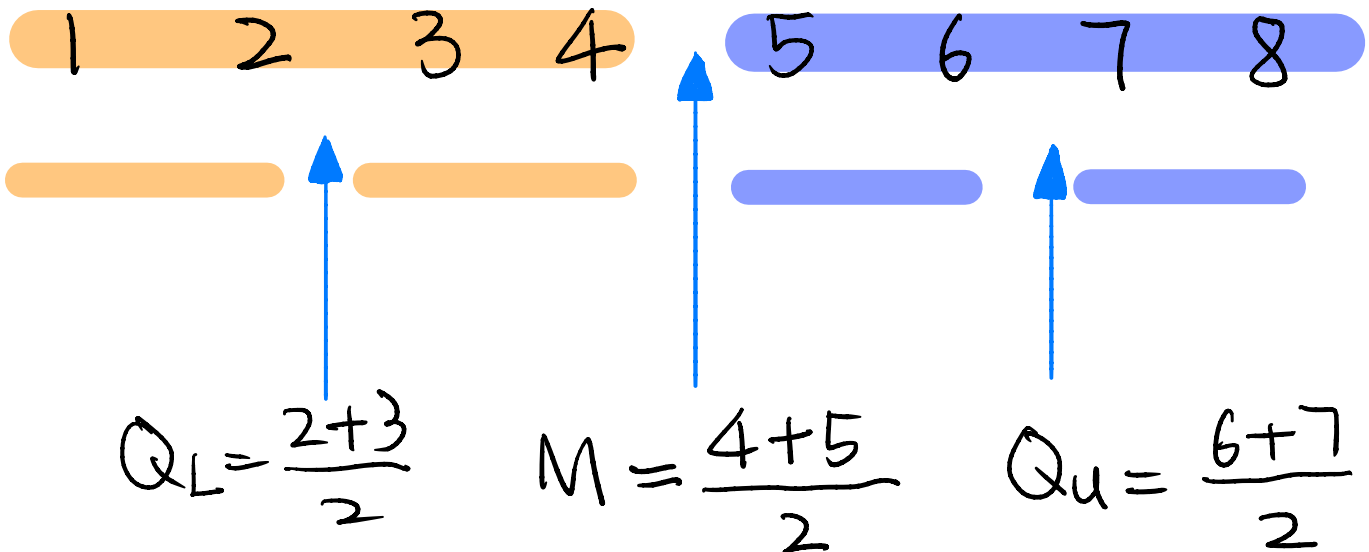
○ upper quartile

$$Q_u = 75\text{-th percentile}$$

# How to find quartiles

- Rank the observations in order (from small to large)
- Cut the order sequence into 4 equal parts
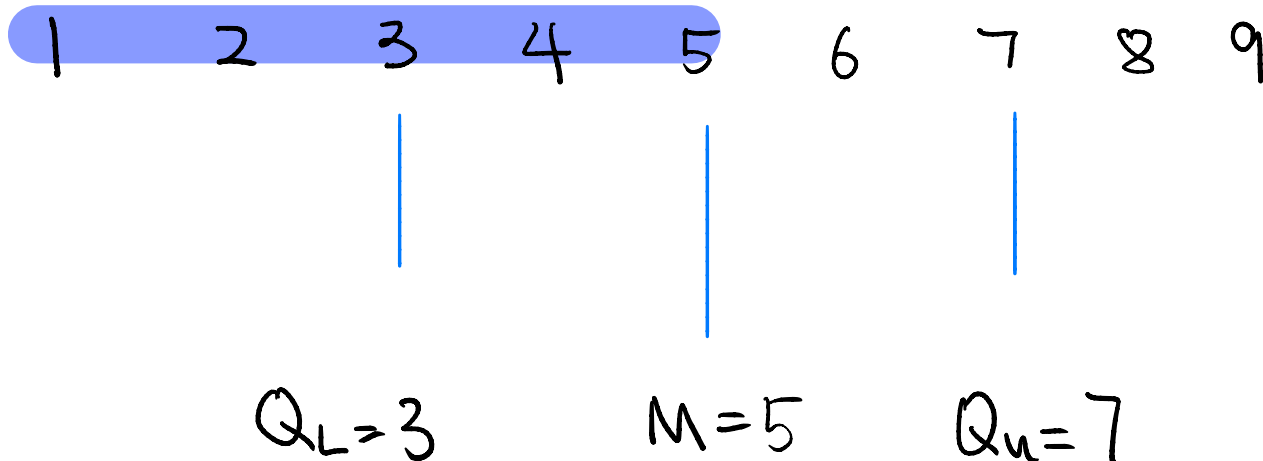- Find the quartiles at these "cuts"

## Ex

Data:

$$1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8$$

$$Q_L = \frac{2+3}{2} \qquad M = \frac{4+5}{2} \qquad Q_u = \frac{6+7}{2}$$

Data:

1 2 3 4 5 6 7 8 9

$Q_L = 3$         $M = 5$         $Q_u = 7$

## Z - score

Definition :
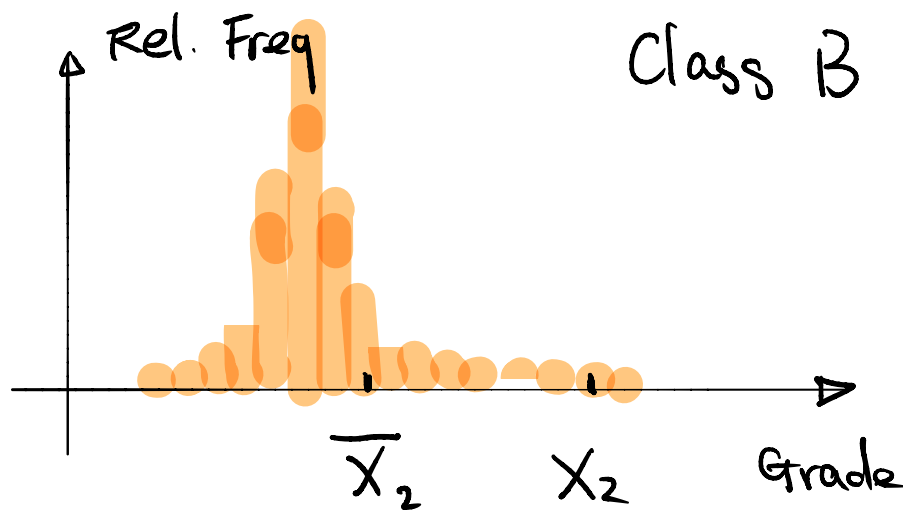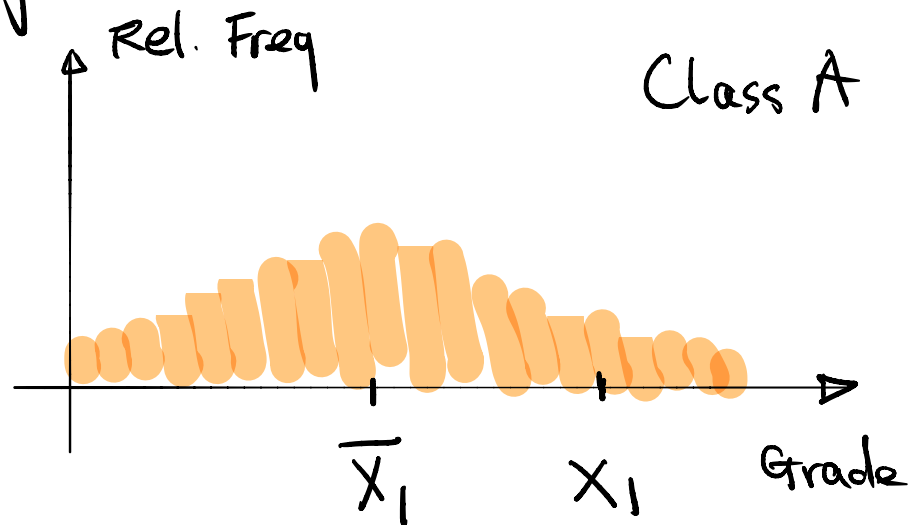
obs.

$$Z = \frac{X - \bar{X}}{S}$$

Idea:   $X - \bar{X}$ measure the deviation from the mean

$\frac{1}{S}$ is the weight of that deviation.

## Why there is a weight $\frac{1}{S}$ ?

S is standard deviation, representing the variability of the data

Histogram



Class A

$S_1$

$\overline{X_1}$   $X_1$   Grade

Class B

$S_2$

$\overline{X_2}$   $X_2$   Grade

$S_1 > S_2$   $\overline{X_1} = \overline{X_2}$

If $x_1 - \bar{x} = x_2 - \bar{x}$

but $s_1 > s_2$ , $\frac{1}{s_1} < \frac{1}{s_2}$

then

$$\frac{x_1 - \bar{X}_1}{s_1} < \frac{x_2 - \bar{X}_2}{s_2}$$

i.e. $z_1 < z_2$

Interpretation :

$z$ represents the distance between a given observation $x$ and mean $\bar{x}$, expressed by in standard deviation $s$. (adjusted by variability)

Given the same deviation $X - \bar{X}$ the higher the variability, the less extreme the observation $X$ is, relatively to the other observations.

## Ex

A sample of 2000 students scores $\bar{X} = 550$ and $s = 75$

$$X_1 = 475 \quad X_2 = 625 \quad X_3 = 700$$

Solution:

$$Z_1 = \frac{475 - \bar{X}}{s} = \frac{475 - 550}{75} = -1$$

$$z_2 = \frac{625 - 550}{75} = 1$$

$$z_3 = \frac{700 - 550}{75} = 2.$$

Student 3 has the best relative performance.

What if we change $S = 25$ ?

$$z_1' = \frac{475 - 550}{25} = -3$$

$$z_2' = \frac{625 - 550}{25} = 3$$

$$z_3' = \frac{700 - 550}{25} = 6.$$

# Empirical rule

If data dist. is bell-shape and symmetric then

1) 68% of obs will have Z-score between -1 and 1

$$68\% \text{ of } x's \in (\bar{X} - S, \bar{X} + S)$$

2) 95% of obs. will have Z-score between -2 and 2.

$$95\% \text{ of } x's \in (\bar{X} - 2S, \bar{X} + 2S)$$

$$Z = \frac{X - \bar{X}}{S} \in (-2, 2) \Rightarrow X - \bar{X} \in (-2S, S)$$

$$\Rightarrow X \in (\bar{X} - 2S, \bar{X} + 2S)$$

3) 99.7% $\cdots$ $---$

$Z \in (-3, 3)$ , $X \in (\bar{X} - 3S, \bar{X} + 3S)$

# Detecting outlier

To identify inconsistent or unusual observations in a dataset.

## Definition :

An outlier is an obs that is actually large or small observations relative to the other observations.

## Causes :

1) incorrectly recorded, wrong entry.

2) come from a different population,

3) correct entry, but rare event.

## Methods to detect outlier

○ Boxplot

- z-score

## Boxplot

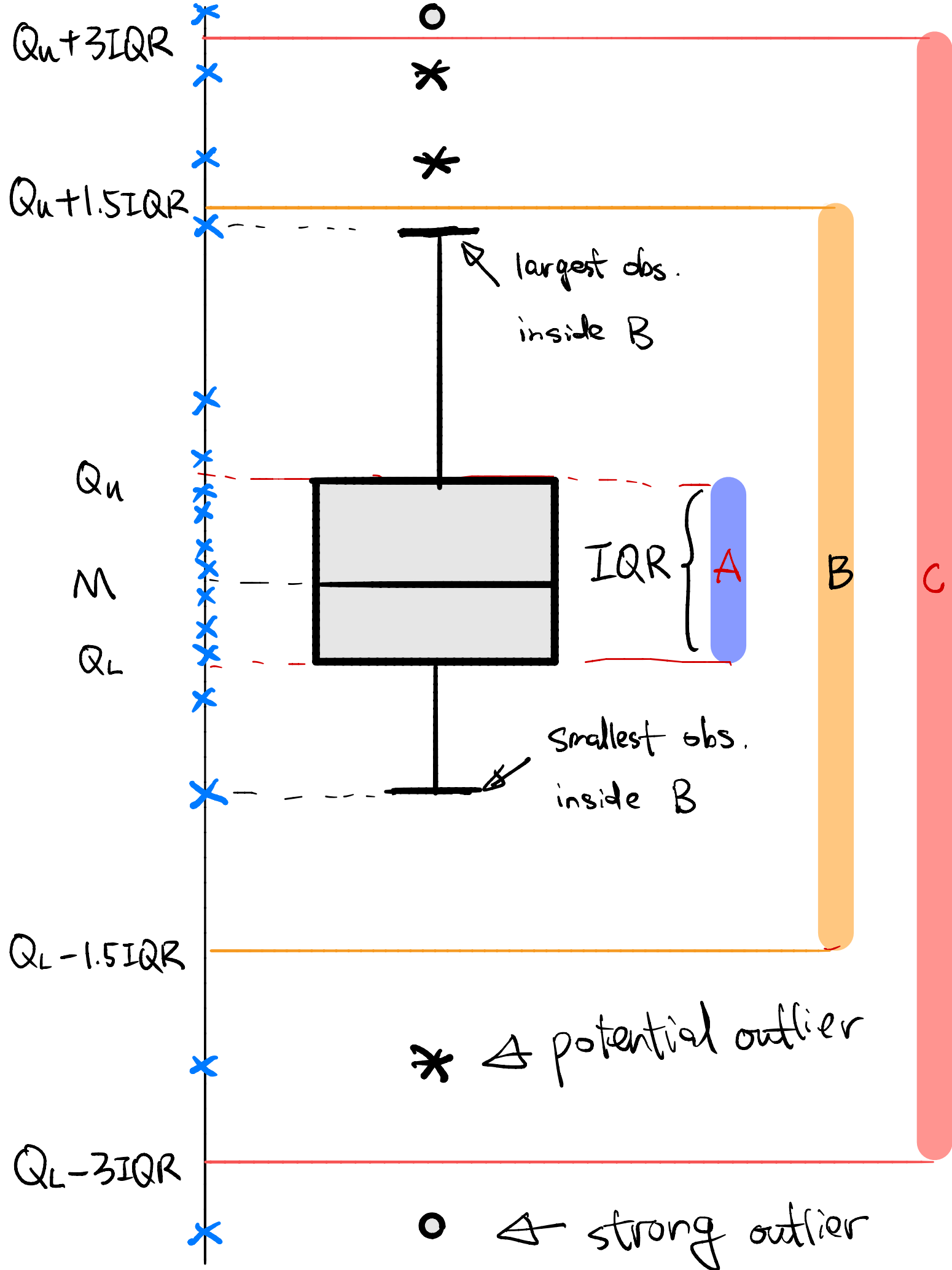**Range A.** Interquartile Range (IQR)

$$IQR = Q_u - Q_L$$

This size of IQR indicates the range of the middle 50% of the observations.

**Range B** Range of inner fences.

$$(Q_L - 1.5 \times IQR, \ Q_u + 1.5 \times IQR)$$

**Range C** Range of outer fences

$$(Q_L - 3 \times IQR, \ Q_u + 3 \times IQR)$$

$Q_u + 3IQR$

$Q_u + 1.5IQR$

largest obs.
inside B

$Q_u$

$IQR \Big\{$  A  B  C

$M$

$Q_L$

smallest obs.
inside B

$Q_L - 1.5IQR$

✳ ◁ potential outlier

$Q_L - 3IQR$

○ ◁ strong outlier

# Rule of Thumb for detecting outlier

(1)    $x \in$ Range C

     $x \notin$ Range B

     $x \in (Q_L - 3 IQR, Q_L - 1.5 IQR)$

or    $x \in (Q_u + 1.5 IQR, Q_u + 3 IQR)$

Then $x$ is a **potential** outlier

(2)    $x$ is outside range C.

     $x \in (-\infty, Q_L - 3 IQR)$

     $x \in (Q_u + 3 IQR, +\infty)$

Then $x$ is a **strong** outlier.

The largest and smallest observations in the dataset are not necessarily outliers.

## Ex.

Data: $n = 13$

0, 1, 2, 4, 5, 5, 7, 10, 10, 12, 13, 17, 39

$$M = 7$$
$$Q_L = 4$$
$$Q_u = 12$$
$$IQR = Q_u - Q_L = 8$$

B: inner fences.

$(Q_L - 1.5 \, IQR, \quad Q_u + 1.5 \, IQR)$

$(4 - 1.5 \times 8, \quad 12 + 1.5 \times 8) = (-8, 24)$

C. outer fences

$(Q_L - 3 \, IQR, \quad Q_u + 3 \, IQR)$

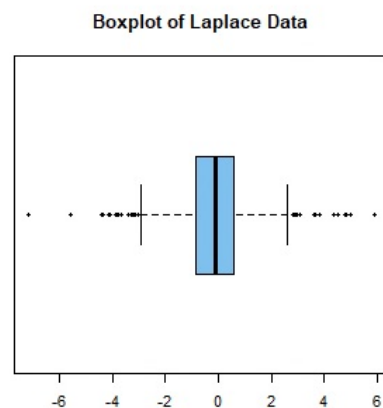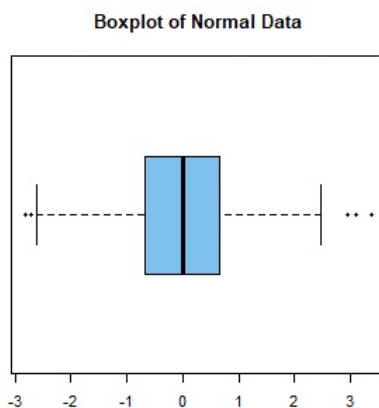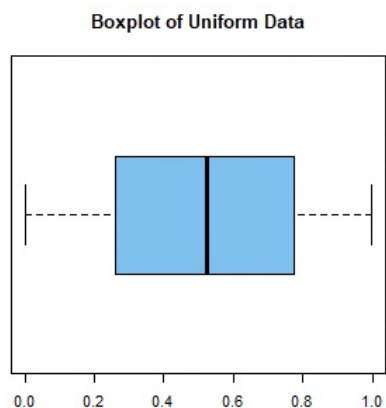$(4 - 3 \times 8, \quad 12 + 3 \times 8) = (-20, 36)$

The 39 is a strong outlier.
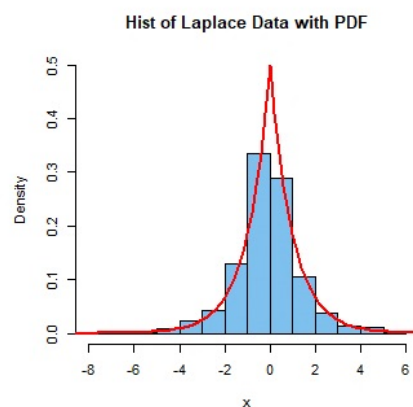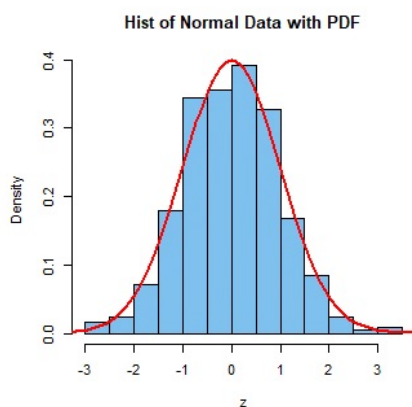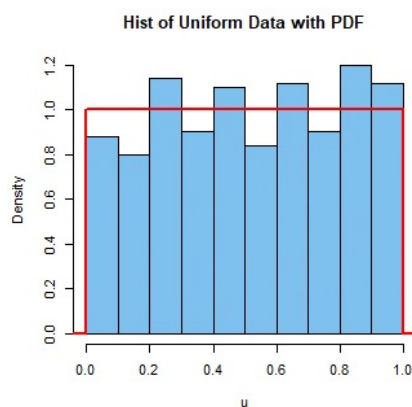
## Z - score

$|z| > 3 \implies$ outlier

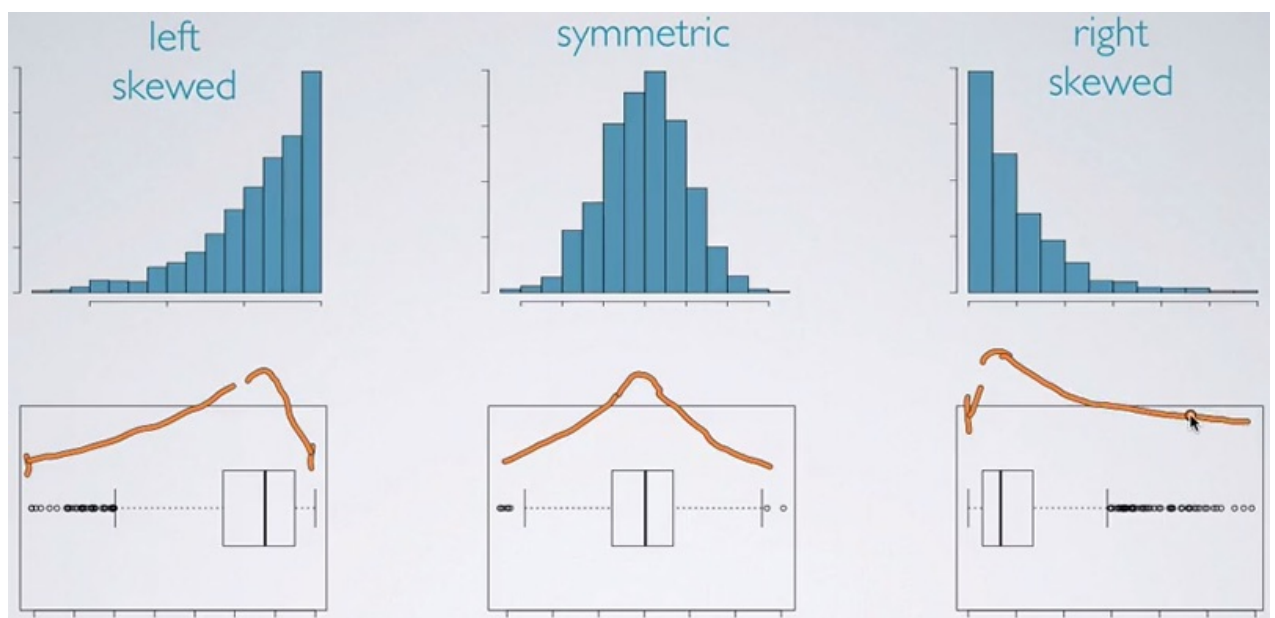$\left| \dfrac{39 - \bar{x}}{s} \right| \not> 3$

# Interpretation of boxplots

1. The line inside the boxplot is the center (median) of the distribution of the data.

2. Examine the length of the box. $(IQR = Q_u - Q_L)$ measures data variation



Hist of Uniform Data with PDF     Hist of Normal Data with PDF     Hist of Laplace Data with PDF

Boxplot of Uniform Data     Boxplot of Normal Data     Boxplot of Laplace Data

3. Visually compare the lengths of the whiskers

   Rule of thumb: the distribution of data is skewed in the direction of longer whisker.

   (the data distribution must be unimodal)



4. The boxplot is less informative than histogram
   It only gives a few pieces of information about the whole data.