## 3.3 Position of a data point within the Group

> We seek to define the location of any given
> data point with in that data Point's Group

This tells us how to compare different
data points

> **Ex** Amongst their gender who is taller?
>
> a man @ 73"      or      a woman @ 68"
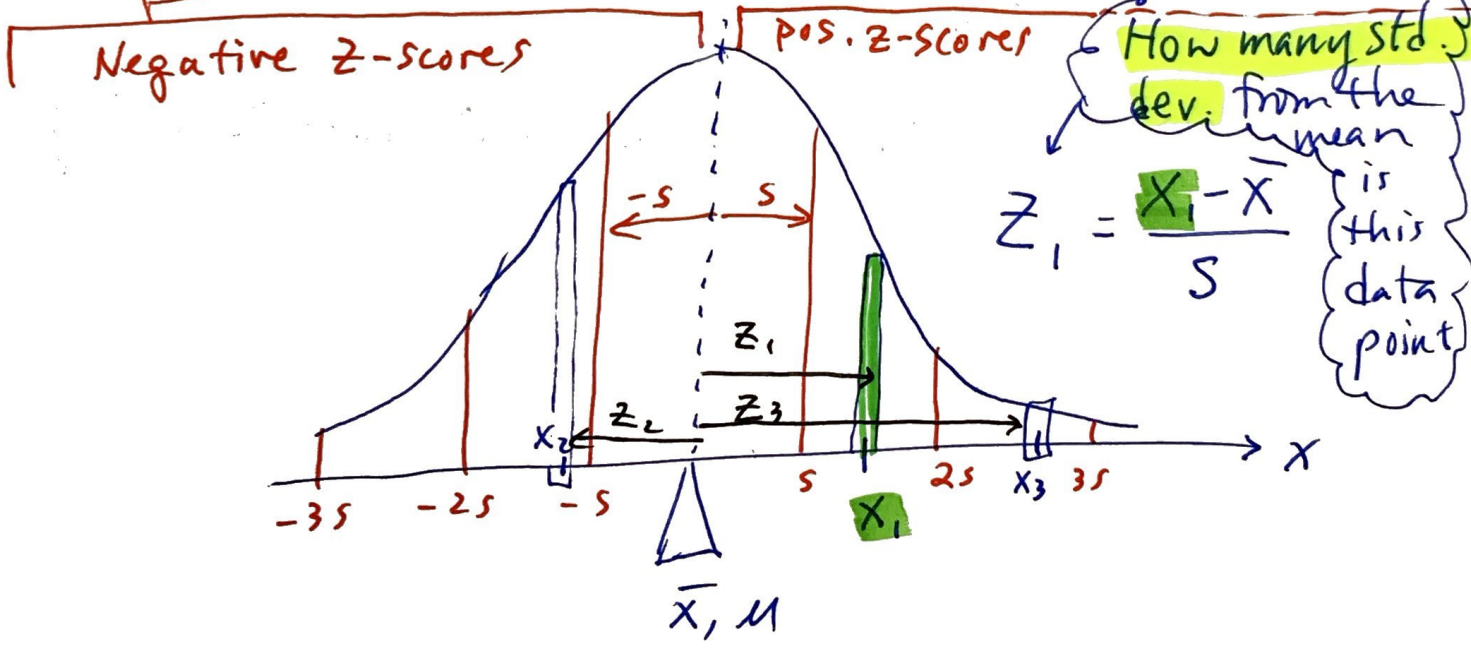
- To answer this question we introduce the

## "Z-score"

$$Z \equiv \frac{X - \mu}{\sigma}$$    and    $$Z = \frac{X - \bar{X}}{s}$$

population                    sample from a
census                       population

Negative Z-scores    |    Pos. Z-scores

How many std.
dev. from the mean
is

$$Z_1 = \frac{X_1 - \bar{X}}{s}$$   this data point



$Z_1$

$X_2$  $Z_2$  $Z_3$

$-3s$  $-2s$  $-s$    $s$    $2s$  $X_3$  $3s$

$X_1$

$\bar{X}, \mu$

**Ex**

Statisticians found that in a pop of college men the average height is 69.4 inches with a std. deviation of 3.1 inches

For college women the mean is 63.8 in and a std. dev. of 2.8 inches

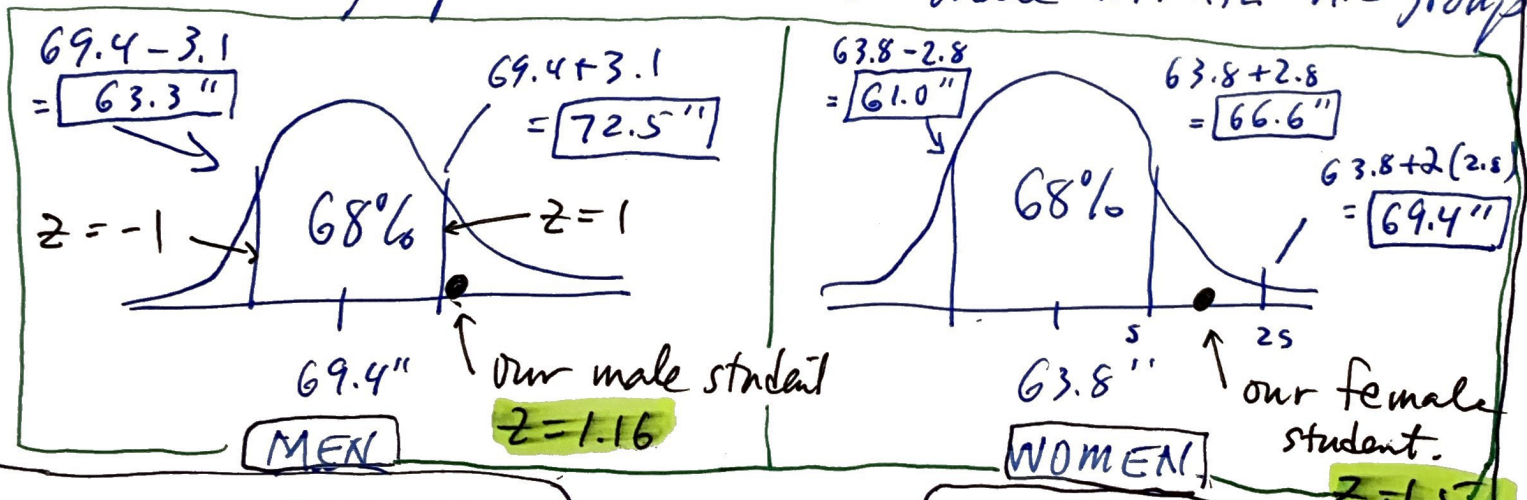Q: who is taller within their gender group: a 73" male or a 68" female?

---

• We can assume these groups has a unimodal and symmetric distribution thus allowing us to use z-scores

$$z = \frac{x - \mu}{\sigma}$$

• The 73" male: $z = \frac{73 - 69.4}{3.1} = 1.16$ std. dev. from the mean

• The 68" female: $z = \frac{68 - 63.8}{2.8} = 1.50$

**Conclusion** The 68" woman is taller within her group vs. the 73" male within his group

$69.4 - 3.1 = 63.3"$

$69.4 + 3.1 = 72.5"$

$z = -1$    68%    $z = 1$

69.4"    our male student $z = 1.16$

**MEN**

$63.8 - 2.8 = 61.0"$

$63.8 + 2.8 = 66.6"$

$63.8 + 2(2.8) = 69.4"$

68%

s    2s

63.8"    our female student $z = 1.5$

**WOMEN**

EX (Cont.)

Q: <u>Is a 75" man unusual ?</u>

i.e. is the z score of this man larger than 2 ?

$$Z = \frac{75 - 69.4}{3.1} = 1.8$$ <u>Not Unusual</u>

Q: <u>Is a 5'4" woman unusual ?</u>

$$Z = \frac{54 - 63.8}{2.8} = -3.5$$ very unusual

Unusual

Recall <u>unusual</u> is $\pm 2s$ or beyond;

this means that Z-scores beyond $\pm 2$ are unusual



data points here are unusual

data points here are unusual

Data here is usual data!

$-3s$ $-2s$ $-1s$ $\bar{X}$ $+1s$ $2s$ $3s$

$\mu$

Z-Score: $-3$ $-2$ $-1$ $0$ $1$ $2$ $3$

three S. Dev. above the mean
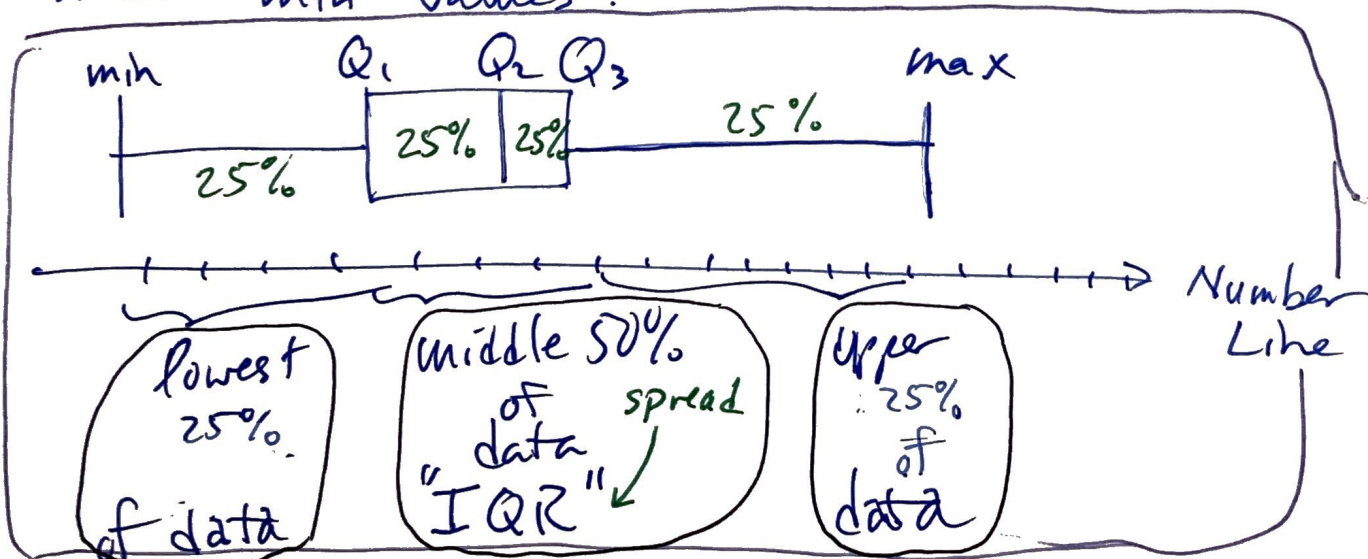
three std. Dev. away below the mean

# ✳ Quartiles

**Def**

The [1st Quartile], $Q_1$, is a data point that separates the lowest 25% of the data from the upper 75% of the Data.

The (2nd Quartile,) $Q_2$, is a data point that separates the lower 50% of data from the upper 50% of data. (aka median)

The (3rd Quartile, $Q_3$, is that data point that separates the lower 75% from the upper 25% of the data

✳ **Box- Plot** — A rectangle whoes edges are at $Q_1$ & $Q_3$ and that has "whiskers" that extend from these edges to the max and min values.



min          $Q_1$    $Q_2$  $Q_3$            max

25%  |25%          25%

25%

Number Line

lowest 25% of data

middle 50% of data "IQR"    spread

upper 25% of data

**Def!** The **IQR** is called the

**Inter Quartile Range**

$$IQR = Q_3 - Q_1$$

Contains 50% of the middle data

* **5 - number Summary!**
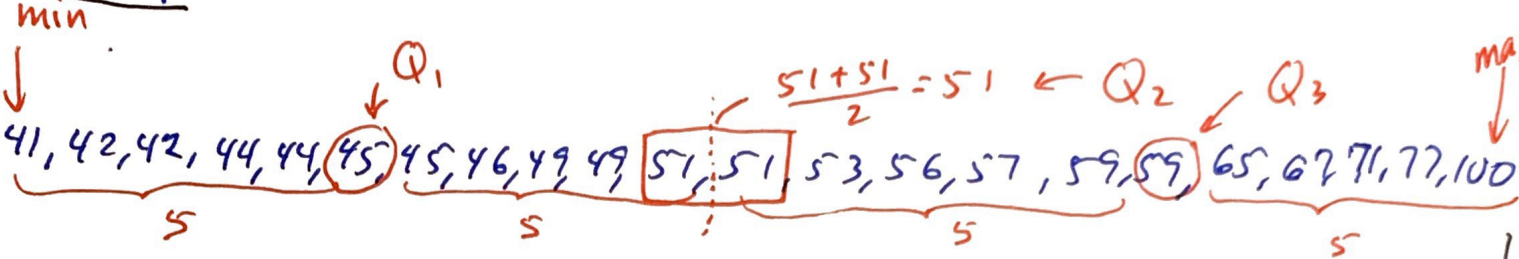
min _____

$Q_1$ _____

median _____

$Q_3$ _____

max _____

The **IQR is a measure of spread**. It is effective for not only unimodal and symmetric data distributions but also unimodal and skewed data.

EX Build a box plot for the follow data

- Raw: 65, 67, 71, 57, 51, 49, 44, 41, 59, 49, 42, 56, 45, 77
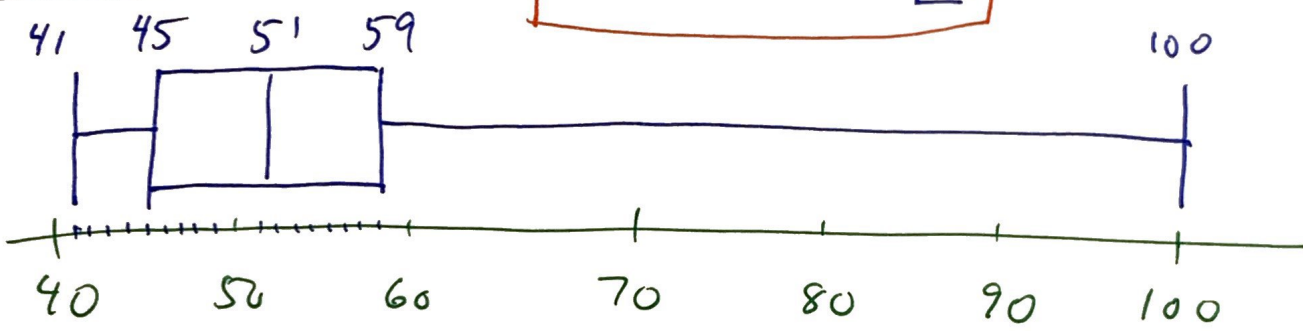  44, 42, 45, 46, 100, 59, 53, 51            N=22

- Order:

min ↓

$$\frac{51+51}{2} = 51 \leftarrow Q_2$$

41, 42, 42, 44, 44, (45) 45, 46, 49, 49 |51, 51| 53, 56, 57, 59, (59) 65, 67, 71, 77, 100

Q₁          Q₂          Q₃          max ↓

5          5          5          5

- 5 number Summary

| min | 41 |
| Q₁ | 45 |
| Q₂ | 51 |
| Q₃ | 59 |
| max | 100 |

$$IQR = 59 - 45 = \boxed{14}$$

- Box Plot

41  45  51  59                                    100

```
        40        50        60        70        80        90        100
```
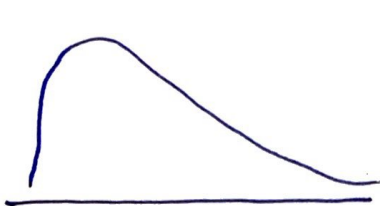
spread of data set is the IQR for non-sym data sets.

- Stat disk:    enter data → data → explore data
                                        ↘ Box - Plot
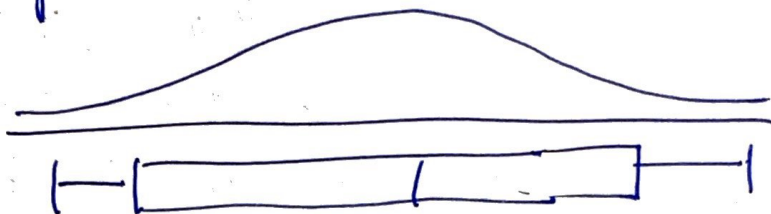
skewed right      symmetrical      skewed Left

- wide data spread



- narrow spread

EX | Consider the data set (pre-ordered)
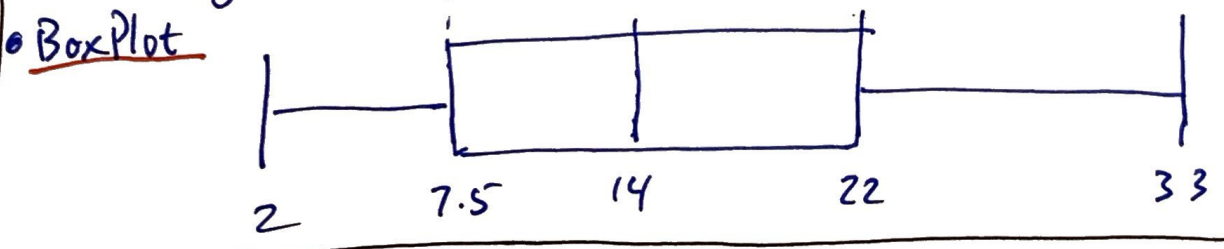
| 2 | 2 | 2 | 2 | 5 | 7 | 8 | 8 | 9 | 9 | 14 | 14 |
|---|---|---|---|---|---|---|---|---|---|----|----|
| 14| 16| 19| 20| 21| 22| 23| 24| 24| 27| 32 | 33 |

$N = 24$

## Construct a Dot Plot (see below)

- **5# summary** :

min : $\underline{\qquad 2 \qquad}$

$Q_1$ : $\dfrac{7+8}{2} = 7.5$

$Q_2$ : $\dfrac{14+14}{2} = 14$

$Q_3$ : $\dfrac{22+22}{2} = 22$

max : $\underline{\qquad 33 \qquad}$

- Dot Plot (not neci'y)

  mode = 2

  median = 14

- Ruler

0    5    10    15    20    25    30    35

- Box Plot

2    7.5    14    22    33

- statdisk
  → col in col 1 of data editor
  → data → box plot → select col 1

  → Box Plot for traditional
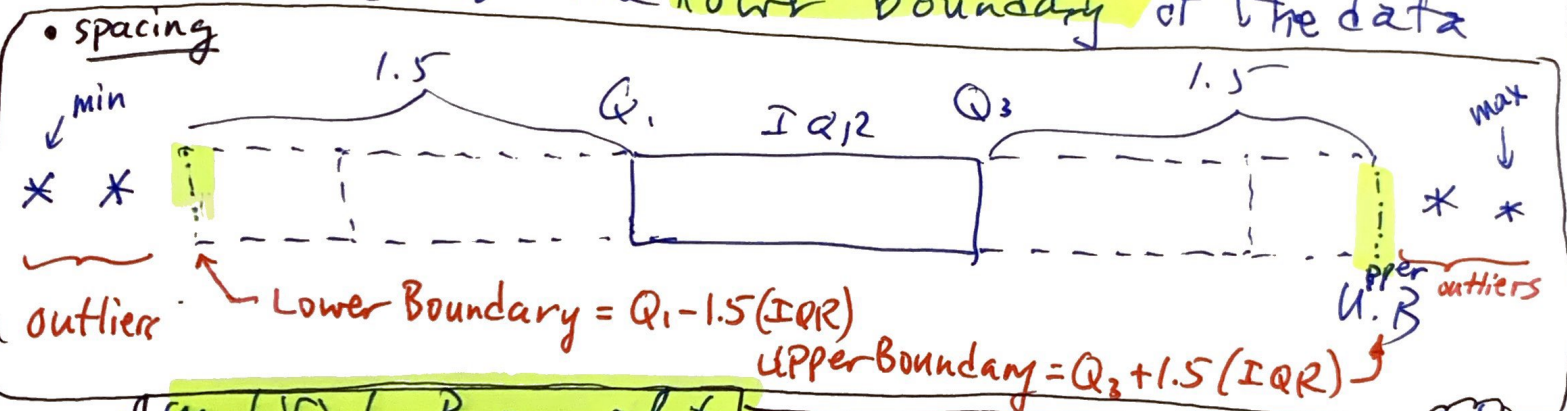  → Modified Box Plot for outlier display

**Outliers** An outlier is a data point that stands apart from the data set.

→ outliers may be erroneous: Correct or ignore

→ outliers may be a fact: explain or ignore (but tell the reader why).

- We identify outliers as being points that are **1.5 IQR's** above $Q_3$, this is called the **Upper Boundary**

— OR —

- Points that are **1.5 IQR's** below $Q_1$, this is called the **lower boundary** of the data

• spacing



Lower Boundary = $Q_1 - 1.5(IQR)$

Upper Boundary = $Q_3 + 1.5(IQR)$

- **Modified Box plot**



last data point inside the lower boundary

last data point inside U.B.

$L.B.$      $Q_1$      $Q_3$      $U.B.$

25%      $IQR$ 50%      25%

**EX** Modified Box Plot

• ordered data    $Q_1$    $Q_2$

4.1, 42, 42, 44, 45, 45, 46, 49, 49, 51, 51, 53, 56, 57, 58, 59, 65, 67, 71, 90, 100

$N = 21$

• 5 # summary

min : 41
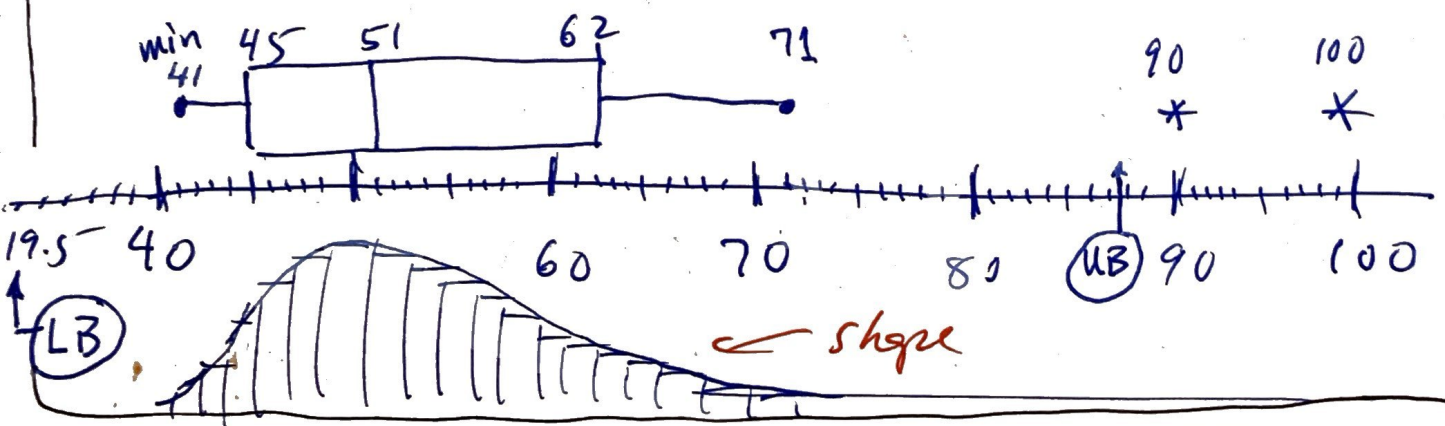
$Q_1 : (45+45)/2 = 45$

$Q_2 : 51$

$Q_3 : (59+65)/2 = 62$    $\Big\} IQR = 62 - 45 = \boxed{17}$

max : 100

• <u>Lower Boundary</u> : $Q_1 - 1.5\,IQR$

$= 45 - 1.5(17)$

$= \boxed{19.5}$

• <u>Upper Boundary</u> : $Q_3 + 1.5\,IQR$

$= 62 + 1.5(17)$

$= \boxed{87.5}$

• modified Box Plot w/ ruler

min 45    51         62              71                    90        100

41                                                                    *          *

19.5  40              60      70          80   (UB) 90     100

(LB)                              ← shape

# ✳ Percentiles

Quartiles divide the data into quarters.
we may desire more refined details

> **EX** what data point divides the upper ⅓
> from the lower ⅔?

- We use percentiles for this task.

> **Def** Given a number $p$, between 1 and 99,
> the **$p^{th}$ percentile (%)** separates the
> lowest $p\%$ from the upper $(100-p)\%$.

$p^{th}$ percentile

**Procedure** — Finding the $p^{th}$ %

1. order the data set, count the data

2. Find the **Locator L**, that point
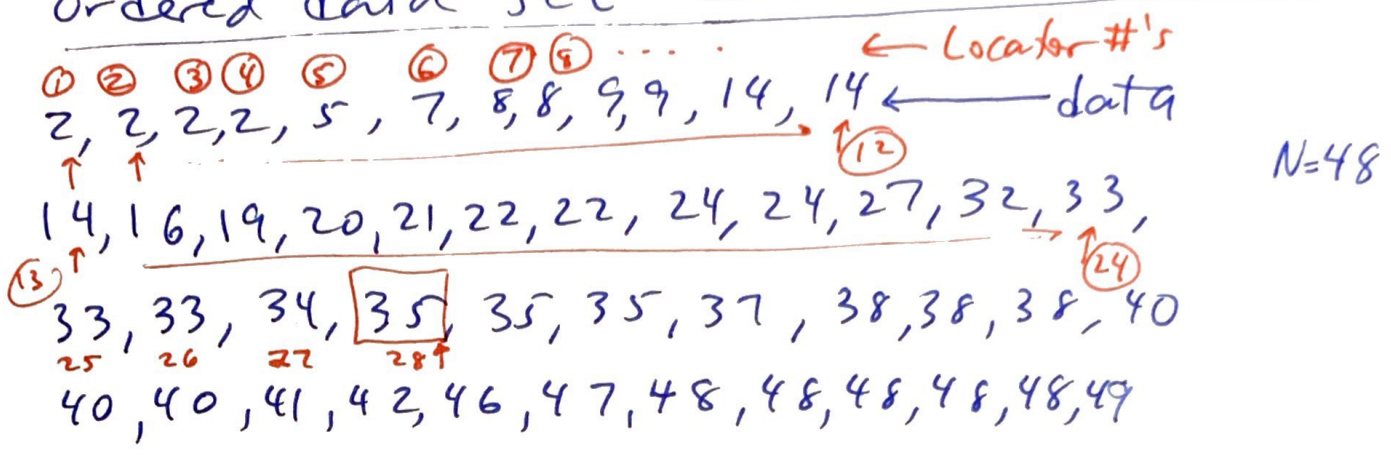   to the data value in the data set.

   $$L = \left(\frac{p}{100}\right)\cdot n$$

   Locator or, in
   Comp. Sci, **pointer**.

3. • If $L$ is a whole # then the $p^{th}$ % is the
     average of the $L$ and $L+1$ location
   • If $L$ is not a whole number then round-up
     to get the Location.

**EX** Find the 58th percentile from the ordered data set below

← Locator #'s

①②③④⑤  ⑥ ⑦⑧ · · · ·  ←———— data
2, 2, 2, 2, 5, 7, 8, 8, 9, 9, 14, 14 ⑫

N=48

14, 16, 19, 20, 21, 22, 22, 24, 24, 27, 32, 33,
⑬ ⑭

33, 33, 34, [35] 35, 35, 37, 38, 38, 38, 40
25  26  27  28↑

40, 40, 41, 42, 46, 47, 48, 48, 48, 48, 48, 49

1.  Order data   (done)

2.  Locator of 58th percentile:  $L = \left(\dfrac{58}{100}\right) 48 = 27.8$

3.  Round up to 28

4.  Count over from the lowest to the 28th data point : Here the 28th data point is [ 35 ]

5.  State Results:

" The data value 35 seperates the lower 58th % of data from the upper 42%.

(( Statdisk has no percentile support ))
• graphically (Dot Plot with 58th %)

```
                                   • 35
 :       : o : :   :           : :  • :
 :  . o : :   :   .  .        :  :  ⊙ :       .  .  :
```

[ 58% ]            ↑                [ 42% ]
                [ 58th percentile ]   of our data