

# Chapter 1

## Descriptive statistics

Statistics is a science of data. To study statistics, we have to describe data from some proper perspectives. Generally speaking, there are two ways to describe data, graphical description and numerical description. Those are what we are going to learn in this section.

Table 1.1: Example Data<sup>1</sup>

NAME	CLASS	GENDER	MID	FINAL	BASKETBALL
James	23	M	74	86	N
Andrew	23	M	74	86	Y
Jim	23	M	23	47	Y
Kim	23	M	61	78	Y
Mark	23	M	97	98	N
Owen	23	M	73	85	Y
Cook	23	M	98	99	Y
Albert	23	M	81	90	Y
Donald	23	M	70	84	N
Peter	23	M	53	72	Y
Vince	23	M	68	82	Y
Davis	23	M	83	91	Y
Alan	23	M	82	90	Y
Nick	23	M	64	80	N
Elina	23	F	72	85	N
Daisy	23	F	68	82	Y
Crystal	23	F	53	72	N
Karida	23	F	66	81	N
Linda	23	F	83	91	N
Dale	23	F	70	83	Y
Sandy	23	F	56	74	N
Emma	23	F	65	80	N
Angela	23	F	72	85	N
Katie	23	F	84	91	N
Eileen	23	F	73	85	N
Meggie	23	F	68	82	N
Jack	24	M	45	67	Y
Stan	24	M	23	47	Y
Ryan	24	M	60	77	Y
Murphy	24	M	36	60	N
Mike	24	M	82	90	Y
Antony	24	M	18	42	Y
Clare	24	M	86	93	Y
David	24	M	83	91	Y
Taylor	24	M	69	83	Y
Park	24	M	78	88	N
Gary	24	M	51	71	Y
Carson	24	M	85	92	Y
Elvis	24	M	25	49	Y
Kelly	24	F	59	76	N
Sara	24	F	77	88	N
Cherry	24	F	61	78	N
Lucy	24	F	54	73	Y
Hellen	24	F	46	68	N
Chloe	24	F	95	97	Y
Dorothy	24	F	82	90	N
Natalie	24	F	73	85	N
Vivien	24	F	76	87	N
Cathy	24	F	70	84	N
Carol	24	F	55	74	N
Bella	24	F	96	98	Y
Veronica	24	F	60	77	N

---

<sup>1</sup>MID is the scores in the midterm exam. FINAL is the scores in the final exam. BASKETBALL indicates whether a student plays basketball or not.

## 1.1 Basic concepts

- In table 1.1, each student is an **individual**.
- All students are described through perspectives of *NAME*, *CLASS*, *GENDER*, *MID*, *FINAL*, and *BASKETBALL*. Those different perspectives are called **variables**, for they may take different values for different students.
- The values of *MID* and *FINAL* can be operated on like normal numbers, such as taking average, subtraction. Those variables are called **quantitative variables**.
- The values of *NAME*, *CLASS*, *GENDER* and *BASKETBALL* only play the role of sorting individuals into different categories. Those variables are called **categorical variables**
- The way a variable takes different values is called the **distribution** of this variable.
- All the individuals we want to study is called the **population**.
- A subset of the population is called a **sample**.

Samples and populations are relative. If you take all the Chinese people as the population, people in Shanghai is a sample. If you take all the people of the whole world as the population, then Chinese people is a sample.

- The number of individuals in the sample is called the **sample size**.

A variable takes values of numbers doesn't mean it is quantitative variable. Is variable *CLASS* a quantitative or categorical variable?

## 1.2 Some basic graphs

- **Pie chart** There are 25 girls and 27 boys in table 1.1. The pie chart of the distribution of *GENDER* is given by figure 1.1.

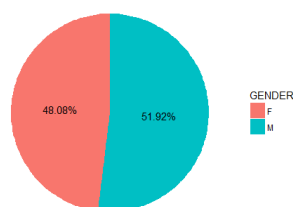


Figure 1.1: Pie chart of the distribution of the *GENDER*

- **Bar graph**

Similarly, we can draw bar graph to show the information about *GENDER*

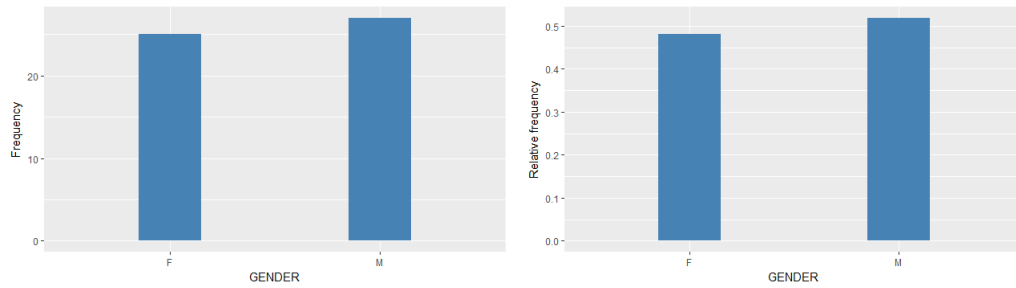


Figure 1.2: Bar graphs with percentage and frequency as vertical axis

In figure 1.2, the vertical axes are **frequency** and **percentage** or (**relative frequency**) respectively. When the sample size is too big, it is better to use relative frequency as the vertical axis.

Be sure to label the axes whenever a graph is drawn!

- **Dotplot**

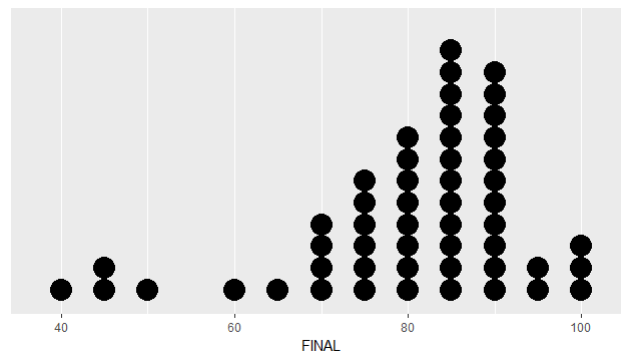


Figure 1.3: Dotplot of the distribution of the *FINAL*

In figure 1.3, the **bin width** is 5. For example, there is only one score lies in the interval  $(37.5, 42.5]$ , which is "42" from the student whose name is "Antony", and the width of this interval is  $42.5 - 37.5 = 5$ , which is the bin width. Similarly, there are eight scores lies in the interval  $(77.5, 82.5]$ .

- **Histogram**

If the dots in dotplot is replaced by bars, the graph will be histogram, as shown in figure

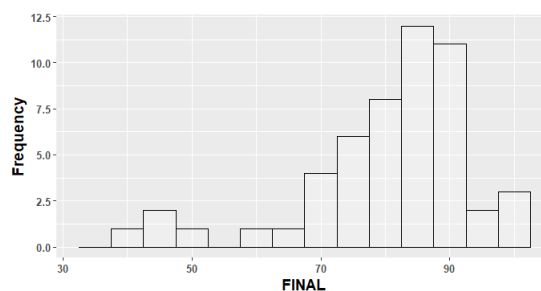


Figure 1.4: Histogram of the distribution of the *FINAL*

The vertical axis can be relative frequency as well.

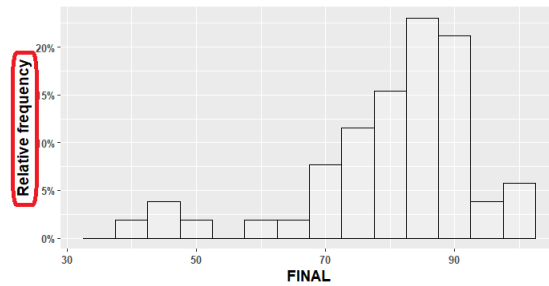


Figure 1.5: Histogram of with vertical axis **relative frequency**

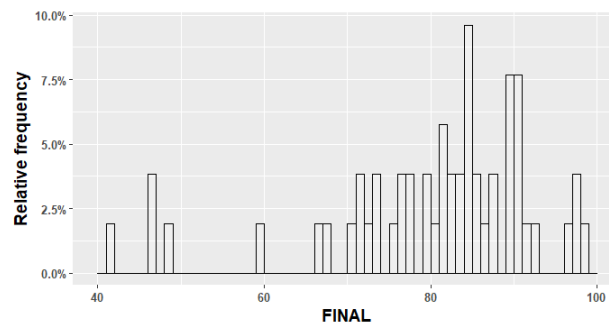
What is the difference between histogram and bar graph?

- **Density curve**

In figure 1.5, we can tell the percentage of  $FINAL \leq 50$  is approximately 8% by adding up the percentages of the first three columns. Here, the percentages are indicated by the height of the bars. (4 out of 52 students with  $FINAL \leq 50$ . They are 42, 47, 47, 49. )

Now, if we draw a histogram with bin width 1 (figure 1.6), then the percentage a bar can be calculated by

$$\text{percentage} = \text{bar height} \times \text{bin width} = \text{area of the bar}.$$



is called **density curve**. The function of the density curve is called **probability density function(pdf)**.

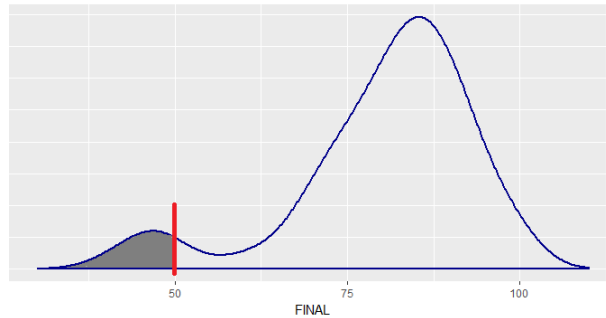


Figure 1.7: Density curve of the distribution of the *FINAL*

Figure 1.7 is a density curve of the distribution of the *FINAL*, the shaded area gives the percentage of  $FINAL \leq 50$ , which is approximately 8%.

Sometimes the vertical axis may be suppressed, because it doesn't mean too much in this book.

What is the total area under the density curve?

- **Cumulative relative frequency curve**

In figure 1.7, for each value of  $x$  there is an area to the left side of this value. Therefore we can get a function  $F$ , such that

$$F(x) = \text{Area to the left of } x.$$

If we draw a smooth graph of  $F(x)$ , it will be like figure 1.8. This curve is called the **cumulative relative frequency curve**. Function  $F(x)$  is called **cumulative density function(cdf)**.

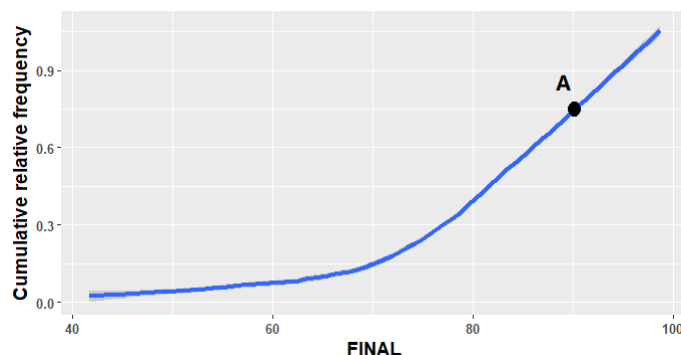


Figure 1.8: Cumulative relative frequency curve of *FINAL*

How to interpret point **A** in figure 1.8?

What is the theoretical relation between figure 1.8 and figure 1.7?

### 1.3 Some terms to describe graphs

- Shape

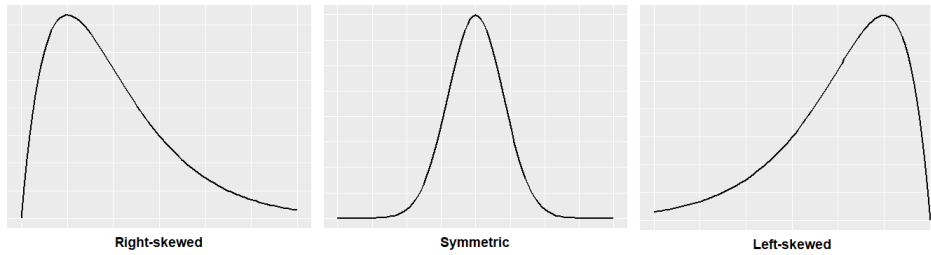


Figure 1.9: Shapes of distributions

As shown in figure 1.9, the shapes of the distributions are **right-skewed**, **symmetric** and **left-skewed** respectively.

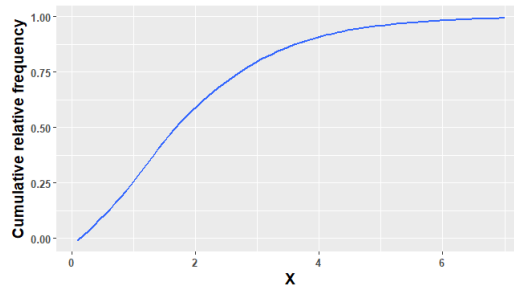


Figure 1.10: A cumulative relative frequency curve

Can you tell whether the distribution in figure 1.10 is right-skewed, left-skewed or symmetric?

- Clusters and gaps

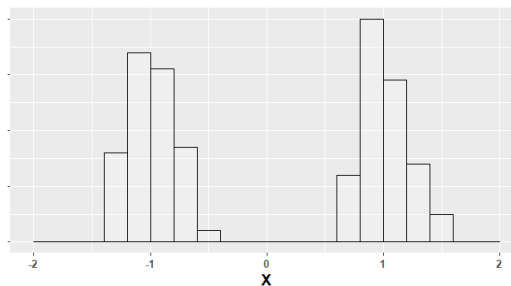


Figure 1.11: Two clusters and a gap

As shown in figure 1.11, we say the distribution has two **clusters(modes)** with a **gap**.

- **Outliers**

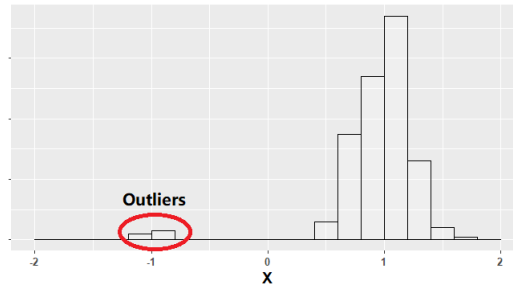


Figure 1.12: Outliers

If some values have striking departures from the pattern of the majority, those values are called **outliers**.

Outliers need special attention, for they may be generated by mistakes or some other mechanisms that are not considered.

## 1.4 Summarizing distributions

- **Center**

There are different ways to describe the center of a distribution. Here we only consider two primary ways of denoting the center: **median** and **mean**.

- **Median**

Arrange the data in increasing or decreasing order, median is the middle one or the average of the middle two.

- **Mean**

For data set  $\{x_1, x_2, \dots, x_n\}$ , the mean  $\bar{x}$  is given by

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

Sometime the notation of the mean is  $\mu$ .  $\mu$  is used for population mean,  $\bar{x}$  is for sample mean. For example, we draw a sample of 100 students from a high school, and the mean weight of those 100 students is 60kg. We use the notation  $\bar{x}$ , because the mean weight 60kg comes from the sample of 100 students. If the mean weight of all the students in this high school is 60kg, we use the notation  $\mu$ . The formulas for  $\mu$  and  $\bar{x}$  are the same.

**When to use mean and when to use median?** Let's take a look at a simple example. Say, we have a set of data  $\{1, 2, 3, 4, 5\}$ . Both the mean and the median are 3. If 5 is recorded as 500 by accident. Now, the mean is 102, while the median is still the same. That is to say, the median is not easily influenced by the extreme values. If a value is not easily influenced by extreme values, it is **resistant**.

The relationship of the mean and median can be shown by figure 1.13.



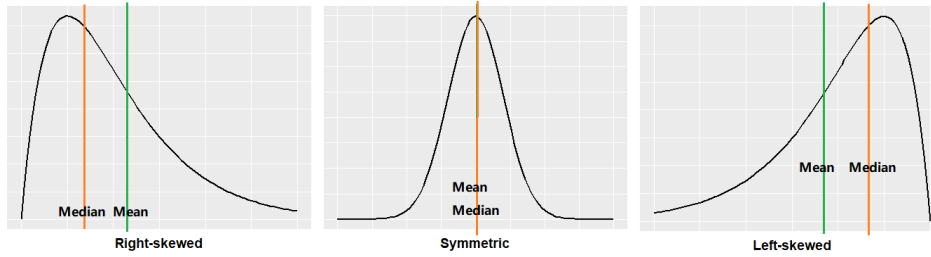


Figure 1.13: The relationship between mean and median

Generally speaking, if the distribution is symmetric, **mean** = **median**; if the distribution is right-skewed, **mean** > **median**; if the distribution is left-skewed, **mean** < **median**. A simple explanation comes as: if the distribution is right-skewed, there are more extreme values to the right side. While the mean is not so resistant as median, it can be more easily dragged to the right than the median. Thus **mean** > **median**.

Therefore, if the distribution is strongly skewed, it is better to use the **median** to describe the center of the distribution.

Is mean or median a better description of the center of the distribution of personal incomes?

- **Spread**

Spread is to measure the variability or the dispersion of the data.

- **Range**

The difference between the maximum and the minimum,

$$\text{range} = \text{maximum} - \text{minimum}$$

- **Interquartile range(IQR)**

**First quartile**( $Q_1$ ) is the value with one quarter of the data less than(or equal) to it. For data {1, 2, 3, 4, 5, 6, 7, 8}, 2 is the first quartile.

**Third quartile**( $Q_3$ ) is the value with 3/4 of the data less than(or equal) to it. For data {1, 2, 3, 4, 5, 6, 7, 8}, 6 is the third quartile.

Interquartile range is given by

$$\text{IQR} = Q_3 - Q_1.$$

For data {1, 2, 3, 4, 5, 6, 7, 8}, the interquartile range **IQR** = 6 – 2 = 4

- **Variance(Var)**

$$\text{Population variance } Var = \sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}.$$

**Sample variance**  $Var = s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$ .

Variance gives the average square of the difference between the mean and data.

– **Standard deviation**

**Population standard deviation**  $\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$ .

**Sample standard deviation**  $\bar{x} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$ .

Standard deviation gives the average distance of the data from the mean.

\* Calculate the mean and standard deviation of sample data  $\{1, 2, 3\}$  by hand

\* Calculate the mean and standard deviation of the *FINAL* of students in class 23 by calculator.

• **Location**

– **Percentile**

$n^{th}$  percentile is the value with n percent of the data smaller or equal to it. For data  $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$ , 1 is the 10<sup>th</sup> percentile, 6 is the 60<sup>th</sup> percentile.

What percentiles are  $Q_1$ , median and  $Q_3$ ?

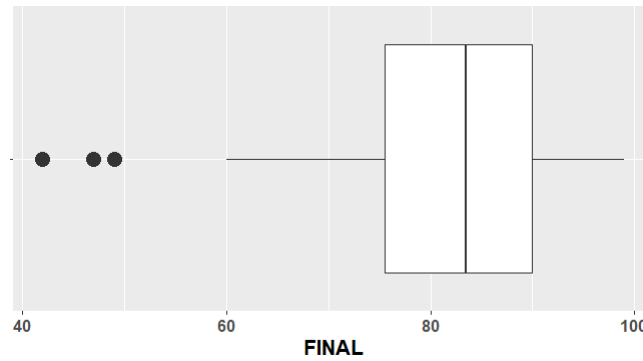
\* **Five number summary** There are five important locations for a given set of data, they are **min**,  **$Q_1$** , **median**,  **$Q_3$**  and **max**. They are called a **five number summary**. The following is a five number summary of the *FINAL*.

##	Min.	Q1.	Median	Q3.	Max.
##	42.00	75.50	83.50	90.00	99.00

\* **1.5 IQR rule** gives a simple rule to tell whether a value is an outlier or not.

$$x \notin [Q_1 - 1.5 \times IQR, Q_3 + 1.5 \times IQR] \implies x \text{ is an outlier.}$$

Find out the outliers of the *FINAL* by the 1.5 IQR rule.

\* **Boxplot**Figure 1.14: The boxplot of the *FINAL*

In figure 1.14, the three dots are outliers, the left vertical line of the box indicates the value of the  $Q_1$ , the middle vertical line indicates the median and the right vertical line indicates the  $Q_3$ .

– **z-score**

For a value  $x$ , its z-score is given by

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

The z-score gives the distance from the mean in terms of standard deviation.

Calculate and interpret the z-score of the *FINAL* of Vince.

Calculate the percentile of the *FINAL* of Vince.

Of the measurements of the spread, which are more resistant and which are not?

• **Data transformation**

Suppose we have a sample data set  $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$  with mean  $\bar{x}$  and standard deviation  $s_x$ .

- Add a constant  $c$  to the data

$$\mathbf{X} + c = \{x_1 + c, x_2 + c, \dots, x_n + c\}.$$

$$\begin{aligned} \text{The mean of } \mathbf{X} + c &= \frac{(x_1 + c) + (x_2 + c) + \dots + (x_n + c)}{n} \\ &= \frac{x_1 + x_2 + \dots + x_n}{n} + c \\ &= \bar{x} + c. \end{aligned}$$

The mean is added by the same constant  $c$ .

$$\begin{aligned} \text{The standard deviation of } \mathbf{X} + c &= \sqrt{\frac{\sum_{i=1}^n [(x_i + c) - (\bar{x} + c)]^2}{n}} \\ &= \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}} \\ &= s_x. \end{aligned}$$

The standard deviation doesn't change.

What about the **IQR** and percentiles if the data is added by a constant  $c$ .

- Multiply the data by constant  $a$

$$a\mathbf{X} = \{ax_1, ax_2, \dots, ax_n\}.$$

$$\begin{aligned} \text{The mean of } a\mathbf{X} &= \frac{(ax_1) + (ax_2) + \dots + (ax_n)}{n} \\ &= a \frac{x_1 + x_2 + \dots + x_n}{n} \\ &= a\bar{x}. \end{aligned}$$

The mean is multiplied by the same constant  $a$ .

$$\begin{aligned} \text{The standard deviation of } a\mathbf{X} &= \sqrt{\frac{\sum_{i=1}^n (ax_i - a\bar{x})^2}{n}} \\ &= |a| \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}} \\ &= |a|s_x. \end{aligned}$$

The standard deviation is multiplied by the absolute value of constant  $a$ .

Calculate the mean and standard deviation of the z-scores for any give data set.

## 1.5 Comparing distributions

- Perspectives to describe a distribution

When you are asked to describe the distribution give by a graph, you are supposed to describe the **shape**, **center** and **spread**. If there are **outliers**, **more than one clusters** and **gaps**, you are suppose to enunciate them.

For example, by referring to figure 1.14, we say the distribution of the *FINAL* is roughly symmetric, with median around 84, and IQR about 16, and there are 3 outliers.

- Graphs for comparing distributions

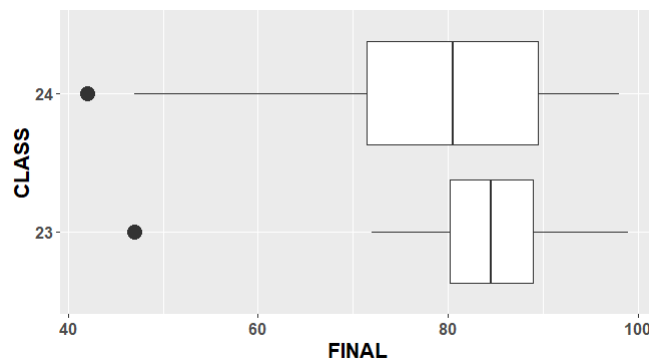


Figure 1.15: Distributions of the *FINAL* of *CLASS 23* and *CLASS 24*

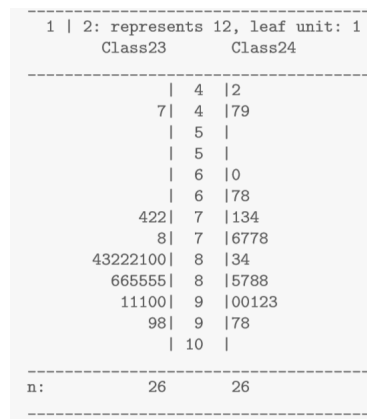


Figure 1.16: Back-to-back stem plot with splitting stems

There are many graphs to compare distributions. Figure 1.15 and figure 1.16 are about the distributions of the *FINAL* of *CLASS 23* and *CLASS 24*.

### • Compare two distributions

The distributions are compared through the same perspectives as when we describe the distributions. They are **shape**, **center** and **spread**, and **outliers**, **clusters** and **gaps** if necessary.

The terms to describe the shape are **right-skewed**, **symmetric** and **left-skewed**.

The terms to describe center are **mean** and **median**. Choose a appropriate one.

The terms to describe the spread are **range**, **IQR** and **standard deviation**. Choose a appropriate one.

For example, if figure 1.15 is given. We can say:

Both of the distributions are approximately symmetric.

The median of the *FINAL* of *CLASS 24* is approximately 80, less than the median of the *FINAL* of *CLASS 23*, which is approximately 85.

Then IQR of *CLASS 24* is approximately 20, larger than the IQR of *CLASS 23*, which is about 10. Thus the *FINAL* of *CLASS 24* is more widely spread than that of *CLASS 23* in the angle of IQR.

Both of the distributions have an outlier to the lower end.

Distribution comparison or distribution description problem always show up in AP exam.

## 1.6 The relation between two variables

- **Two categorical variables** By reading table 1.1, we may suspect that there is a relation between categorical variable *BASKETBALL* and *GENDER*. Maybe boys are more likely to play basketball than girls. How can we describe the relation between *BASKETBALL* and *GENDER*?

– **Two-way table**

		BASKETBALL		Total
		Yes	No	
GENDER	Male	21	6	27
	Female	5	20	25
Total		26	26	52

Table 1.2: Two-way table of Gender  $\times$  Basketball

Table 1.2 gives a simple description about the relation between *BASKETBALL* and *GENDER*.

\* **Conditional distribution**

The **conditional distribution** of a categorical variable is defined as the distribution of this variable while the value of the other variable is fixed.

<i>BASKETBALL</i>	Percentage of Male	Percentage of Female
Yes	$\frac{21}{26} \approx 81\%$	$\frac{5}{26} \approx 19\%$
No	$\frac{6}{26} \approx 23\%$	$\frac{20}{26} \approx 77\%$

Table 1.3: Conditional distribution

Table 1.3 gives the conditional distribution of *GENDER* conditioned on different different values of *BASKETBALL*. For example, the conditional distribution of *GENDER* among those who play basketball is that: about 81% of them are boys, 19% are girls.

If a student plays basketball, this students is more likely to be a boy than a girl. Clearly, there is some association between *GENDER* and *BASKETBALL*.

\* **Association**

If the conditional distributions of a variable are different while conditioned on different values of the other variable, we say there is an **association** between those variables. Otherwise, they are **independent**.

The conditional distributions of *GENDER* are different conditioned on different values(Yes, No) of *BASKETBALL*. Therefore, there is an association between *GENDER* and *BASKETBALL*.

\* **Marginal distribution**

If we consider the distribution of *GENDER* regardless of the *BASKETBALL*, we just look at the data at the right margin of tabel 1.2.

$$\text{Percentage of girls: } \frac{25}{52} = 48\%, \quad \text{Percentage of boys: } \frac{27}{52} = 52\%$$

Similarly, we can find out the distribution of *BASKETBALL* regardless of the *GENDER* by looking at the data at the bottom margin of tabel 1.2.

All the distributions are **marginal distribuions**, for we only consider the data at the margin of the two-way table.

Is it true that if two variables have no association(independent), the marginal distributions and the conditional distributions are the same?

– **Side-by-side bar graph**

**Side-by-side bar graph** is a graph to show the relation between two categorical variables.