

## Stem-and-Leaf Plots

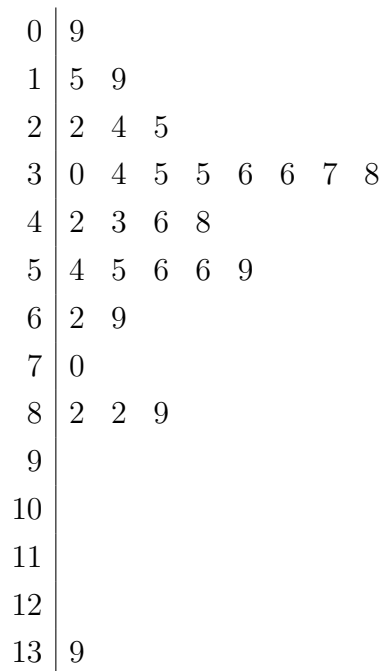
**Stem-and-leaf** plots are similar to histograms, although they carry more information. Namely, they also show how the data are distributed *within columns*. To construct a stem-and-leaf plot, we need to draw a stem and a leaf. The first one or several digits form a “stem”, and the next digit forms a “leaf”. Other digits are dropped; in other words, the numbers get rounded. For example, the number 139 can be written as

$$13 \mid 9$$

with 13 going to the stem and 9 to the leaf, or as

$$1 \mid 3$$

with 1 joining the stem, 3 joining the leaf, and the digit 9 being dropped. In the first case, the leaf unit equals 1 and the stem unit is 10, while in the second case, the leaf unit is 10 and the stem unit is  $10^2$ , showing that the (rounded) number is not 13, but 130. The stem and leaf units *must be carefully specified* for each such plot.



Stem-and-leaf plot, Example 4.1

**Example 4.1.** Consider again the data from Example 4.3 from last time, about the CPU times (in seconds) for  $N = 30$  randomly chosen jobs (sorted increasingly),

9 15 19 22 24 25 30 34 35 35  
 36 36 37 38 42 43 46 48 54 55  
 56 56 59 62 69 70 82 82 89 139

Let us draw a stem-and-leaf plot with leaf unit 1 (i.e., the last digits form a leaf). The remaining digits go to the stem. Each CPU time is then written as

$$10 \text{ “stem”} + \text{ “leaf”},$$

making the stem-and-leaf plot above.

Turning this plot by 90 degrees counterclockwise, we get a histogram with 10–unit bins (because each stem unit equals 10). Thus, all the information seen on a histogram can be obtained here too. In addition, now we can see *individual* values within each column.

Stem-and-leaf plots can also be used to compare two samples. For this purpose, one can put two leaves on the same stem.

**Example 4.2.** The following two samples represent transmission times (in seconds) of signals - known as “pings”- from two different locations.

L1: 0.0156, 0.0396, 0.0355, 0.0480, 0.0419, 0.0335, 0.0543, 0.0350,  
 0.0280, 0.0210, 0.0308, 0.0327, 0.0215, 0.0437, 0.0483,  
 L2: 0.0298, 0.0674, 0.0387, 0.0787, 0.0467, 0.0712, 0.0045, 0.0167,  
 0.0661, 0.0109, 0.0198, 0.0039.

Let us sort the two samples in increasing order.

L1: 0.0156, 0.0210, 0.0215, 0.0280, 0.0308, 0.0327, 0.0335, 0.0350  
 0.0355, 0.0396, 0.0419, 0.0437, 0.0480, 0.0483, 0.0543,  
 L2: 0.0039, 0.0045, 0.0109, 0.0167, 0.0198, 0.0298, 0.0387, 0.0467  
 0.0661, 0.0674, 0.0712, 0.0787.

Since all numbers start with 0.0..., we choose a stem unit of 0.01, a leaf unit of 0.001 and drop the last digit. We construct the following two stem-and-leaf plots (two in one), one to the left (L1) and



Is there a connection between the frequency of running antivirus software and the number of worms in the system? A scatter plot of these data is given in Figure 1(a). It clearly shows that the number of worms reduces, in general, when the antivirus is employed more frequently. This relationship, however, is not certain because no worm was detected on some “lucky” computers although the antivirus software was launched only once a week (4 times a month) on them.

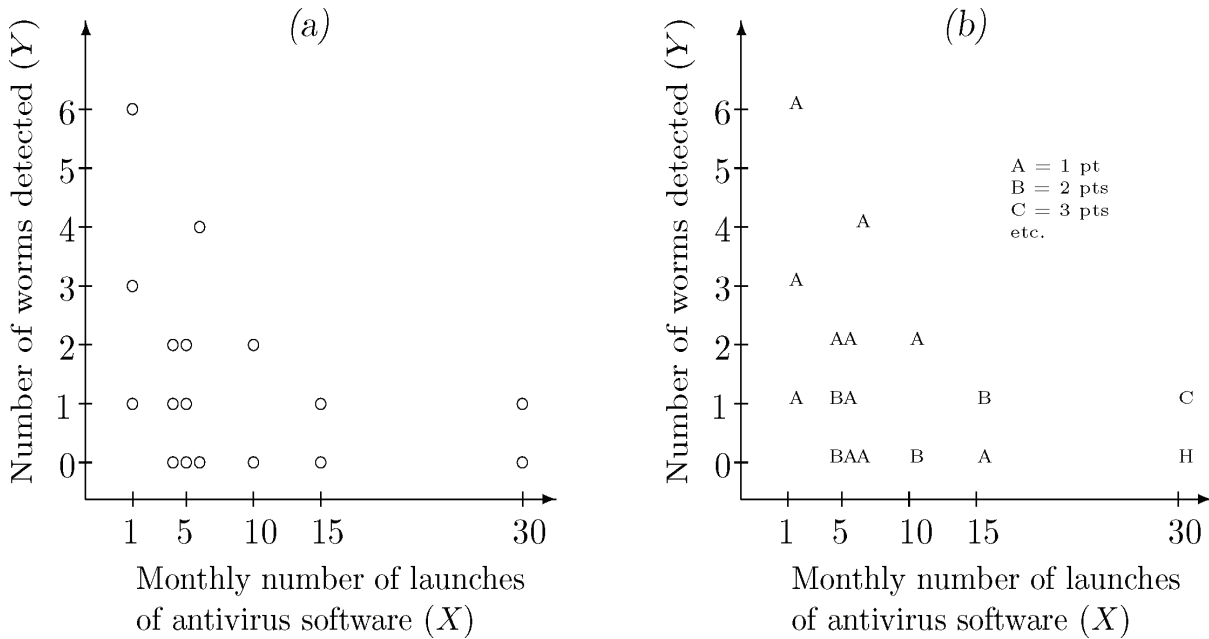


Fig. 1: Scatter plots for Example 4.3

Looking at the scatter plot in Figure 1(a), the manager realized that a portion of data is hidden there because there are identical observations. For example, no worms were detected on 8 computers where the antivirus software is used daily (30 times a month). Then, Figure 1(a) may be misleading. When the data contain identical pairs of observations, the points on a scatter plot are often depicted with either numbers or letters (e.g., “A” for 1 point, “B” for two identical points, “C” for three, ..., “H” for eight, etc.). You can see the result in Figure 1(b).

When we study time trends and development of variables over time, we use **time plots**. These are scatter plots with  $x$ -variable representing time.

**Example 4.4.** Here is how the world population increased between 1950 and 2012 (Figure 2). We can clearly see that the population increases at an almost steady rate.

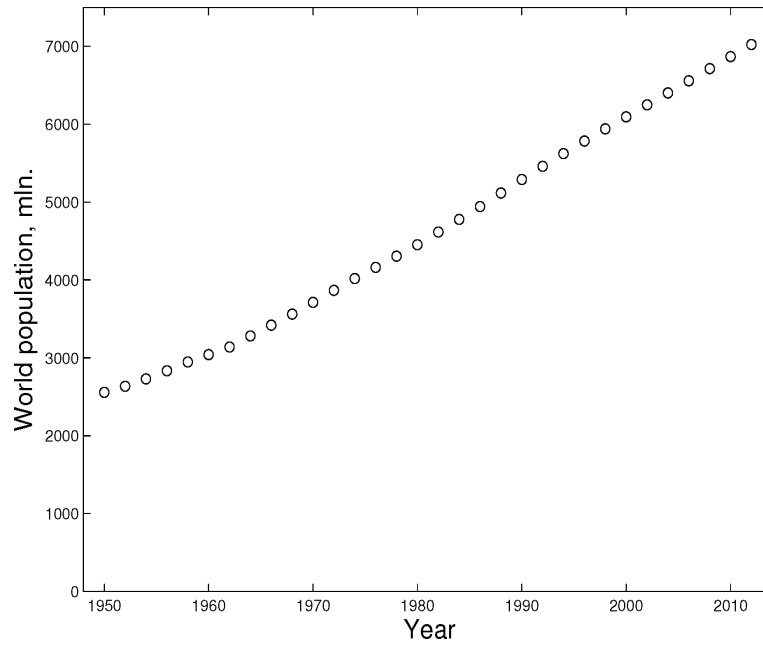


Fig. 2: Time plot of the world population in 1950–2012, Example 4.4

The actual data will be given and studied in the next chapter (Correlation and Regression). We will estimate the trends seen on time plots and scatter plots and even make forecasts for the future.

# Chapter 3. Calculative Descriptive Statistics

In the previous chapter we have considered some graphical methods for getting an idea of the shape of the density function of the population from which the sample data was drawn. Some characteristics, such as symmetry, regularity can be observed from these graphical displays of the data. Next, we consider some statistics that allow us to summarize the data set analytically. Simple **descriptive statistics** measuring the location, spread, variability and other characteristics can be computed immediately. It is hoped that these will give us some idea of the values of the parameters that characterize the entire population from which the sample was pooled. We are looking mainly at two types of statistics: *measures of central tendency*, i.e. values that locate the observations with highest frequencies (so, where most of the data values lie) and *measures of variability*, that indicate how much the values are spread out.

## 1 Measures of Central Tendency

These are values that tend to locate in some sense the “middle” of a set of data. The term “average” is often associated with these values. Each of the following measures of central tendency can be called the “average” value of a set of data.

### 1.1 Mean

**Definition 1.1.** The (*arithmetic*) *mean* ( $\overline{\text{mean}}$ ) of the data  $x_1, \dots, x_N$  is the value

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i. \quad (1.1)$$

For grouped data,  $\left( \begin{array}{c} x_i \\ f_i \end{array} \right)_{i=1, n}$ ,

$$\bar{x} = \frac{1}{N} \sum_{i=1}^n f_i x_i.$$

**Remark 1.2.** The sum of all deviations from the mean is equal to 0. Indeed,

$$\sum_{i=1}^N (x_i - \bar{x}) = \sum_{i=1}^N x_i - N\bar{x} = 0.$$

**Example 1.3.** We consider again the data from Example 4.1 (CPU times):

70 36 43 69 82 48 34 62 35 15  
59 139 46 37 42 30 55 56 36 82  
38 89 54 25 35 24 22 9 56 19

The *mean* CPU time is

$$\bar{x} = \frac{70 + 36 + \dots + 56 + 19}{30} = 48.2333 \text{ seconds.}$$

We may conclude that the mean CPU time of *all* the jobs handled by that particular processor is about the same, “near” 48.2333 seconds. In other words, we try to estimate the *population mean* by the *sample mean*. How good would that approximation be? We will learn later how to assess the accuracy of our estimates.

**Example 1.4.** Let us assume that the value  $x = 139$  (that seemed extreme, out of place, when we looked at the histogram) was *not* in this sample. Then the mean would be

$$\bar{x}_1 = 45.1034,$$

somewhat lower.

Now, in the other direction, let us suppose that the CPU time of one more job (a heavier one) is recorded and it is found to be 30 minutes = 1800 seconds. The mean of the new sample is

$$\bar{x}_2 = 104.7419 \text{ seconds,}$$

way larger than the first value!

## 1.2 Median

One disadvantage of the sample mean is its *sensitivity to extreme observations*. As we have seen in the previous example, one extreme value can significantly shift the value of the mean, to the point where it becomes almost irrelevant.

The next measure of location is the *median*, which is much less sensitive than the mean.

**Definition 1.5.** The *median* (median) is the value  $\bar{M}$  that divides a set of ordered data  $X$  into two equal parts, i.e. the value with the property that it is exceeded by at most a half of observations and is preceded by at most a half of observations.

A sample is always *discrete*, since it consists of a finite number of observations. Then, computing a sample median is similar to the case of discrete distributions. In simple random sampling, all observations are equally likely, and thus, equal probabilities on each side of a median translate into an equal number of observations. There are two cases, depending on the sample size  $N$ .

If the sorted primary data is

$$x_1 \leq \dots \leq x_N,$$

then

$$\bar{M} = \begin{cases} x_{k+1}, & \text{if } N = 2k + 1 \\ \frac{x_k + x_{k+1}}{2}, & \text{if } N = 2k \end{cases}.$$

**Remark 1.6.** The median may or may not be one of the values in the data.

**Example 1.7.** Let us find the median for the data in Example 1.3 (the CPU times).

Since there are  $N = 30$  observations, there are two middle values, the 15th and the 16th entries.

9	15	19	22	24	25	30	34	35	35
36	36	37	38	<b>42</b>	<b>43</b>	46	48	54	55
56	56	59	62	69	70	82	82	89	139

Then the median is  $\bar{M} = 42.5$ .

**Remark 1.8.** For an even number of observations, the median can be chosen to be any number between the two middle values. So in the previous example, we could say that any number in the interval  $(42, 43)$  is a median.

**Example 1.9.** Let us add again the extreme value of 30 minutes = 1800 seconds. The new sample

9	15	19	22	24	25	30	34	35	35	
36	36	37	38	42	<b>43</b>	46	48	54	55	
56	56	59	62	69	70	82	82	89	139	1800

has 31 observations, there is only one middle value (the 16th entry), so the median of the new sample is

$$\bar{M}_2 = 43.$$

Notice that the new value differs very little from the previous one and is *still relevant*, unlike the mean. So the median is a *robust* statistic, not being influenced (so much) by outliers.



### 1.3 Mode

**Definition 1.10.** A *mode*  $Mo$  of a random variable  $X$  is a value with the highest pdf, i.e., it is the point with the highest concentration of probability,  $Mo = \operatorname{argmax}\{f(x)\}$ . A *sample mode*,  $\bar{x}_{mo}$ , of a set of data is a most frequent value.

**Remark 1.11.** Notice from the wording of the definition that the mode may not be unique. A distribution can have one mode – **unimodal**, two modes – **bimodal**, three modes – **trimodal**, or more – **multimodal**.

When the pdf of a continuous distribution has multiple local maxima, it is common to refer to *all* of the local maxima as modes of the distribution.

If every value occurs only once in a sample, we say that there is **no mode**.

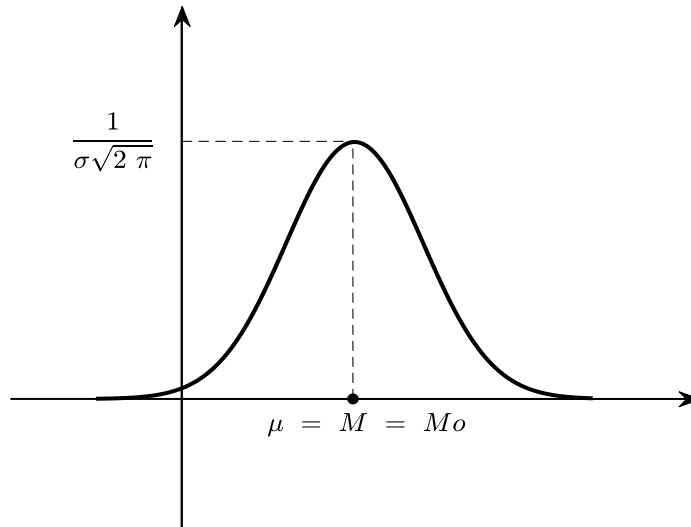


Fig. 3: Normal Distribution (unimodal)

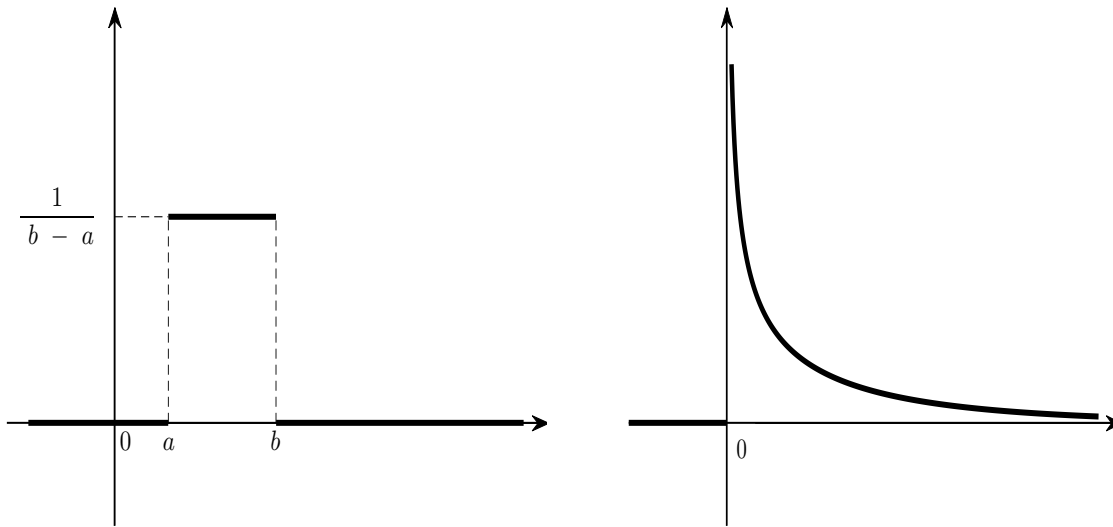
For data drawn from symmetric distributions, we have

$$\bar{x} = \bar{M} = x_{mo}.$$

This is true, for instance, for the Normal distribution which is unimodal (Figure 3). For a Uniform  $U(a, b)$  distribution, *all* values in the interval  $[a, b]$  are modes (Figure 4(a)), while the  $\chi^2(1)$  distribution (with  $\nu = 1$  degree of freedom) has no mode (Figure 4(b)).

In general,

$$x_{mo} \approx \bar{x} - 3(\bar{x} - \bar{M}).$$



(a) Uniform Distribution (multimodal)

(b)  $\chi^2$  Distribution (no mode)

Fig. 4: Multiple modes and no mode

This empirical formula was given by K. Pearson.

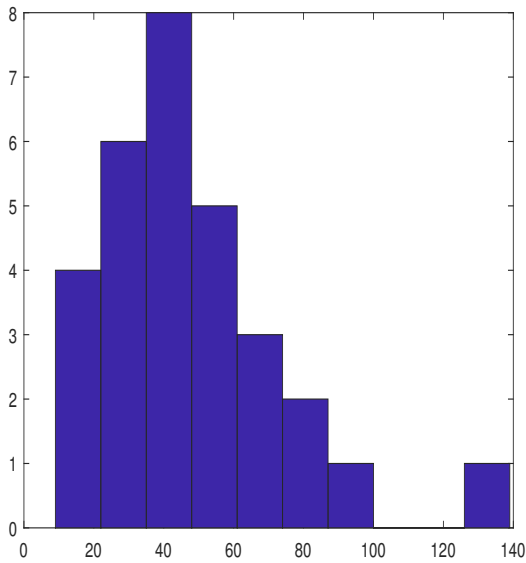
**Example 1.12.** In our example about the CPU times (Example 1.3), the values 35, 36, 56 and 82 appear twice, while all the other values have a frequency of 1. So all four are modes, this is multimodal data.

9 15 19 22 24 25 30 34 **35 35**  
**36 36** 37 38 42 43 46 48 54 55  
**56 56** 59 62 69 70 **82 82** 89 139

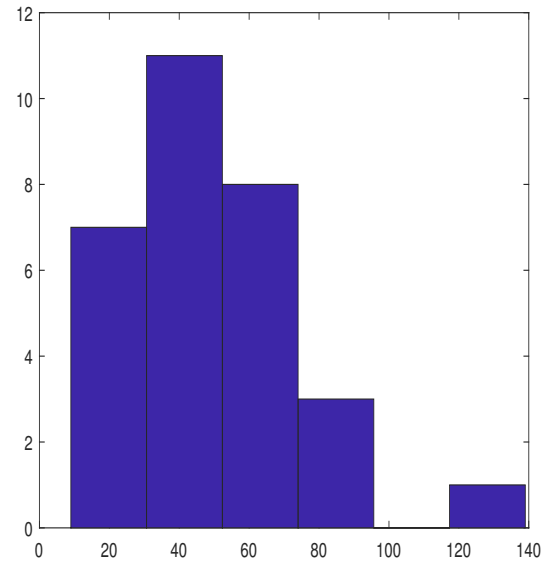
If we group the data into 10 classes, then the *modal class* is the third one,  $(35, 48]$ , with modal mark 41.5 (Figure 5(a)). If we have only 6 classes, then the second one is the modal class,  $[30.7, 52.4)$ , with mark 41.55 (Figure 5(b)).

## 2 Measures of Variability

Once we have located the central values of a set of data, it is important to measure the *variability*, whether the data values are tightly clustered or spread out. At the heart of Statistics lies variability: measuring it, reducing it, distinguishing random from real variability, identifying the various sources of real variability and making decisions in the presence of it. We need to know how “unstable” the



(a)  $n = 10$  bins



(b)  $n = 6$  bins

Fig. 5: Modal class

data is and how much the values differ from its average or from other middle values. These numbers will have small values for closely grouped data (little variation) and larger values for more widely spread out data (large variation).

The measures of variation will also help us assess the reliability of our estimates and the accuracy of our forecasts.

## 2.1 Quantiles, percentiles and quartiles

Consider the primary data  $X = \{x_1, \dots, x_N\}$ . The first two measures of variation give a very general idea of the spread in the data values.

**Definition 2.1.** The *range* ( $\boxed{\text{range}}$ ) of  $X$  is the difference

$$x_{max} - x_{min}.$$

If the values of  $X$  are sorted in increasing order, then the range is  $x_N - x_1$ .

**Definition 2.2.** The *mean absolute deviation* ( $\boxed{\text{mad}}$ ) of  $X$  is the mean of the absolute value of the

deviations from the mean, i.e. the value

$$MAD_1 = \frac{1}{N} \sum_{i=1}^N |x_i - \bar{x}|.$$

The **median absolute deviation** ( $\boxed{\text{mad}}$ ) of  $X$  is the median of the absolute value of the deviations from the median, i.e. the value

$$MAD_2 = \text{median}\{|x_i - \bar{M}|\}.$$

Like the median, the median absolute deviation is not influenced by extreme values, whereas the mean absolute deviation is.

Next, following the idea behind the definition of the median, we define values that divide the data into certain percentages. We simply replace 0.5 in its definition by some probability  $0 < p < 1$ .

**Definition 2.3.** Let  $X$  be a set of data sorted increasingly,  $p \in (0, 1)$  and  $k = 1, 2, \dots, 99$ .

- (1) A **sample  $p$ -quantile** ( $\boxed{\text{quantile}}$ ) is any number that exceeds at most  $100p\%$  of the sample and is exceeded by at most  $100(1 - p)\%$  of the sample.
- (2) A  **$k$ -percentile** ( $\boxed{\text{prctile}}$ )  $P_k$  is a  $(k/100)$ -quantile. So,  $P_k$  exceeds at most  $k\%$  and is exceeded by at most  $(100 - k)\%$  of the data
- (3) The **quartiles** of  $X$  are the values

$$Q_1 = P_{25}, \quad Q_2 = P_{50} = \bar{M} \quad \text{and} \quad Q_3 = P_{75}.$$

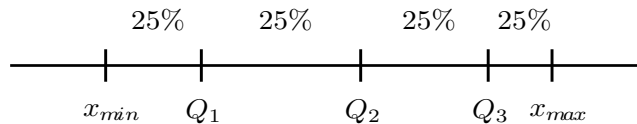


Fig. 6: Quartiles

**Definition 2.4.** Let  $X$  be a set of sorted data with quartiles  $Q_1, Q_2$  and  $Q_3$ .

(1) The **interquartile range** ( $\boxed{\text{iqr}}$ ) is the difference between the third and the first quartile

$$IQR = Q_3 - Q_1. \quad (2.2)$$

(2) The **interquartile deviation** or the **semi interquartile range** is the value

$$IQD = \frac{IQR}{2} = \frac{Q_3 - Q_1}{2}. \quad (2.3)$$

(3) The **interquartile deviation coefficient** or the **relative interquartile deviation** is the value

$$IQDC = \frac{IQD}{\bar{M}} = \frac{Q_3 - Q_1}{2Q_2}. \quad (2.4)$$

**Remark 2.5.**

1. The interquartile deviation is an absolute measure of variation and it has an important property: the range  $\bar{M} \pm IQD$  contains approximately 50% of the data.
2. The interquartile deviation coefficient  $IQDC$  varies between  $-1$  and  $1$ , taking values close to  $0$  for symmetrical distributions, with little variation and values close to  $\pm 1$  for skewed data with large variation.

**Example 2.6.** Consider again the CPU times (in seconds) for  $N = 30$  randomly chosen jobs (sorted ascendingly):

9 15 19 22 24 25 30 **34** 35 35  
36 36 37 38 42 43 46 48 54 55  
56 56 **59** 62 69 70 82 82 89 139

**Solution.** For this example, the range is

$$139 - 9 = 130 \text{ seconds}$$

and the mean and median absolute deviations are

$$MAD_1 = 19.6133,$$

$$MAD_2 = 13.5.$$

To determine the quartiles, notice that 25% of the sample equals  $30/4 = 7.5$  and 75% of the sample is  $90/4 = 22.5$  observations. From the ordered sample, we see that the 8th element, 34, has

7 observations to its left and 22 to its right, so it has *no more* than 7.5 observations to the left and *no more* than 22.5 observations to the right of it. Hence,  $Q_1 = 34$ .

Similarly, the third quartile is the 23rd smallest element,  $Q_3 = 59$ . Recall from last time that the second quartile (the median) is  $Q_2 = \bar{M} = 42.5$ . Then

$$\begin{aligned}IQR &= 59 - 34 = 25, \\IQD &= IQR/2 = 12.5, \\IQDC &= IQD/Q_2 = 0.2941.\end{aligned}$$

The interval

$$\bar{M} \pm IQD = [30, 55]$$

contains 14 observations.

The value of the *IQDC* is close neither to 0, nor to the values  $\pm 1$ . So the data doesn't show strong symmetry or strong asymmetry. This may be due to the extreme values 9 and/or 139. ■

**Example 2.7.** A computer maker sells extended warranty on the produced computers. It agrees to issue a warranty for  $x$  years if it knows that only 10% of computers will fail before the warranty expires. It is known from past experience that lifetimes of these computers have a Gamma distribution with parameters  $\alpha = 60$  and  $\lambda = 1/5$  years. Compute  $x$  and advise the company on the important decision under uncertainty about possible warranties.

**Solution.** We just need to find the tenth percentile of the specified Gamma distribution and let  $x = P_{10}$ . In Matlab, that would be computed (as the *inverse* of the cdf) by

$$x = \text{gaminv}(0.1, 60, 1/5) = 10.0624.$$

Thus, the company can issue a 10-year warranty rather safely. ■

**Remark 2.8.** For populations or very large data sets, calculating exact percentiles can be computationally very expensive since it requires sorting all the data values. Machine learning and statistical software use special algorithms (such as linear interpolation) to get an approximate percentile that can be calculated very quickly and is guaranteed to have a certain accuracy.

## Outliers

The interquartile range is also involved in another important aspect of statistical analysis, namely the detection of outliers. An *outlier*, as the name suggests, is basically an atypical value, “far away” from the rest of the data, that does not seem to belong to the distribution of the rest of the values in the data set.

We have seen how the mean is very sensitive to outliers. Other statistical procedures can be gravely affected by the presence of outliers in the data. Thus, the problem of detecting and locating an outlier is an important part of any statistical data analysis process.

How to classify a value as being “extreme”? First, we could use a simple property, known as the “ $3\sigma$  rule”. This is an application of Chebyshev’s inequality

$$P(|X - E(X)| < \varepsilon) \geq 1 - \frac{V(X)}{\varepsilon^2}, \forall \varepsilon > 0.$$

If we use the classical notations  $E(X) = \mu$ ,  $V(X) = \sigma^2$ ,  $\text{Std}(X) = \sigma$  for the mean, variance and standard deviation of  $X$  and take  $\varepsilon = 3\sigma$ , we get

$$\begin{aligned} P(|X - \mu| < 3\sigma) &\geq 1 - \frac{\sigma^2}{9\sigma^2} \\ &= \frac{8}{9} \approx .89. \end{aligned}$$

This is saying that it is *very* probable (at least 0.89 probable) that  $|X - \mu| < 3\sigma$ , or, equivalently, that  $\mu - 3\sigma < X < \mu + 3\sigma$ . In words, the  $3\sigma$  rule states that *most of the values that any random variable takes, at least 89%, lie within 3 standard deviations away from the mean*. This property is true in general, for any distribution, but especially for unimodal and symmetrical ones, where that percentage is even higher.

Based on that, one simple procedure would be to consider an outlier any value that is more than 2.5 standard deviations away from the mean, and an *extreme* outlier a value more than 3 standard deviations away from the mean.

A more general approach, that works well also for skewed data, is to consider an outlier any observation that is outside the range

$$\left[ Q_1 - \frac{3}{2}IQR, Q_3 + \frac{3}{2}IQR \right] = [Q_1 - 3IQD, Q_3 + 3IQD].$$

Also, the coefficient  $3/2$  can be replaced by some other number to decrease or enlarge the

interval of “normal” values (or, equivalently, the domain that covers the outliers):

$$[Q_1 - w \cdot IQR, Q_3 + w \cdot IQR], \quad w = 0.5, 1, 1.5.$$

For our example on CPU times of processors, we have

$$Q_1 - \frac{3}{2}IQR = -3.5,$$

$$Q_3 + \frac{3}{2}IQR = 96.5,$$

so observations outside the interval  $[-3.5, 96.5]$  are considered outliers. In this case, there is only one, the value 139.

## Boxplots

All the information we discussed above is summarized in a graphical display, called a **boxplot** (boxplot), a plot in which a rectangle is drawn to represent the second and third quartiles (so the interquartile range), with a line inside for the median value and which indicates which values are considered extreme. The “whiskers” of the boxplot are the endpoints of the interval on which normal values lie (so everything outside the whiskers is considered an outlier).

For the data in Example 2.6, the boxplot is displayed in Figure 7.

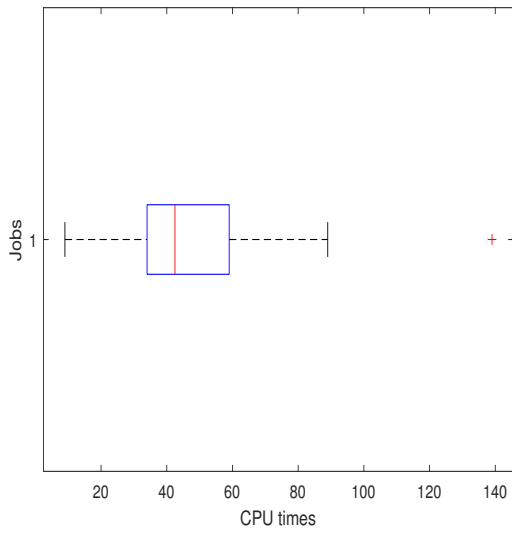
A boxplot can be displayed vertically (default) or horizontally, as in Figure 7. The box can have a “notch” (indentation) at the value of the median, as in Figure 8(a). The width of the interval of the whiskers can be changed. The interval that determines the outliers (i.e., outside of which values are considered too extreme, outliers) is

$$[Q_1 - w \cdot IQR, Q_3 + w \cdot IQR].$$

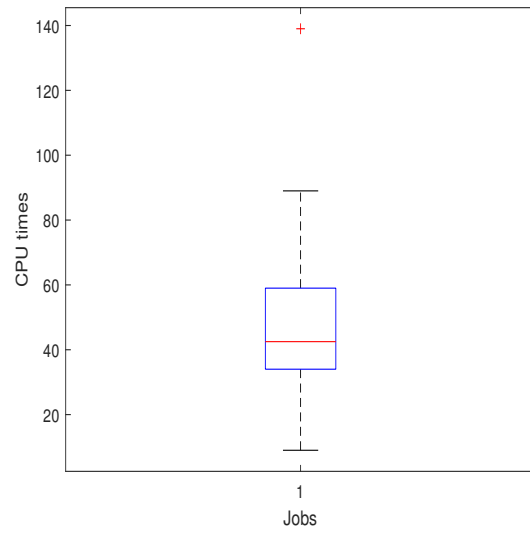
The default value is  $w = 1.5$ . With the smaller whiskers, boxplot displays more data points as outliers. In Figure 8(b), the whisker size is set to  $w = 0.5$ . Then, outliers are all the values outside the interval  $[Q_1 - 0.5 \cdot IQR, Q_3 + 0.5 \cdot IQR] = [21.5, 71.5]$ . These would be 9, 15, 19 (too small) and 82, 89, 139 (too large).

Boxplots are also very useful when we want to compare data from different samples (see Figure 9). We can compare the interquartile ranges, to examine how the data is dispersed between each sample. The longer the box, the more dispersed the data.



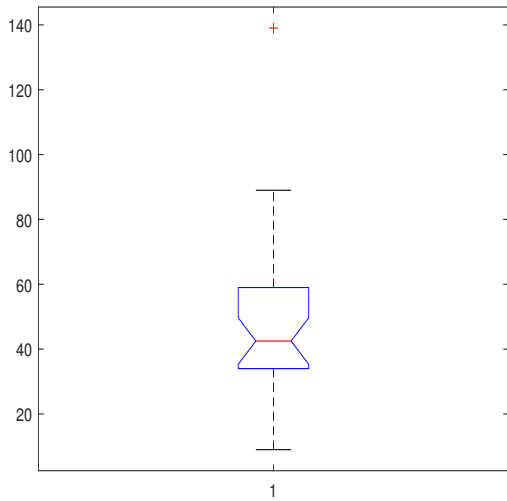


(a) horizontally

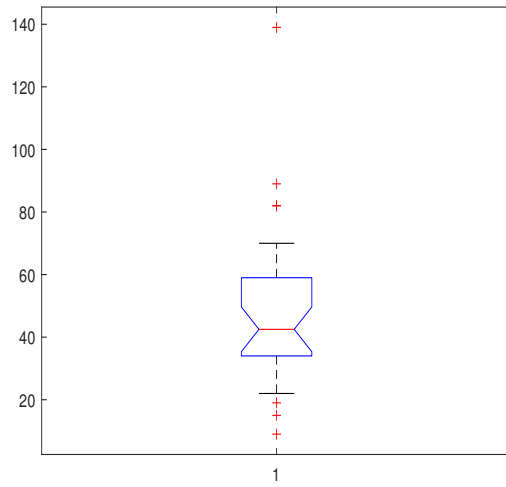


(b) vertically

Fig. 7: Quartiles, Interquartile Range, Outliers



(a) boxplot with a notch



(b) whisker  $w = 0.5$

Fig. 8: Boxplots

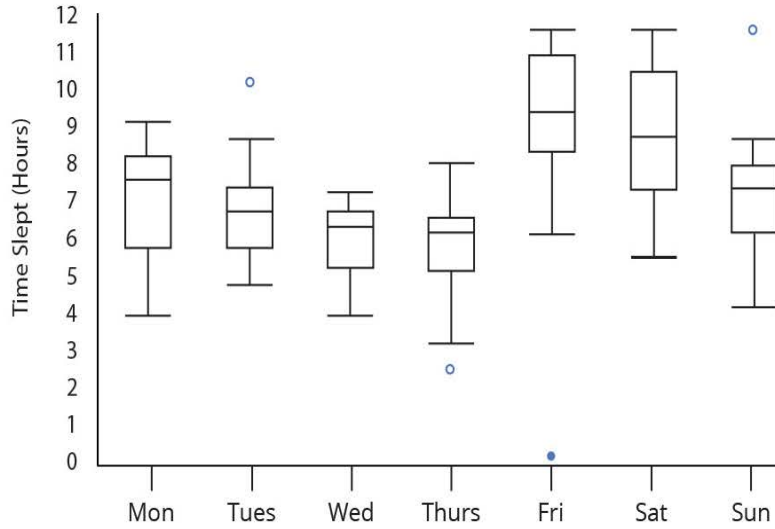


Fig. 9: Multiple boxplots

## 2.2 Moments, variance, standard deviation and coefficient of variation

The idea of the mean can be generalized, by taking various powers of the values in the data.

### Definition 2.9.

(1) The *moment of order  $k$*  is the value

$$\bar{\nu}_k = \frac{1}{N} \sum_{i=1}^N x_i^k, \quad \bar{\nu}_k = \frac{1}{N} \sum_{i=1}^n f_i x_i^k, \quad (2.5)$$

for primary and for grouped data, respectively.

(2) The *central moment of order  $k$*  (moment) is the value

$$\bar{\mu}_k = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^k, \quad \bar{\mu}_k = \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^k \quad (2.6)$$

for primary and for grouped data, respectively.

(3) The **variance** ( $\boxed{\text{var}}$ ) is the value

$$\bar{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2, \quad \bar{\sigma}^2 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^2 \quad (2.7)$$

for primary and for grouped data, respectively. The quantity  $\bar{\sigma} = \sqrt{\bar{\sigma}^2}$  is the **standard deviation** ( $\boxed{\text{std}}$ ).

**Remark 2.10.**

1. A more efficient computational formula for the variance is

$$\bar{\sigma}^2 = \frac{1}{N} \left( \sum_{i=1}^N x_i^2 - \frac{1}{N} \left( \sum_{i=1}^N x_i \right)^2 \right) = \frac{1}{N} \left( \sum_{i=1}^N x_i^2 - N\bar{x}^2 \right), \quad (2.8)$$

which follows straight from the definition.

2. We will see later that when the data represents a sample (not the entire population), a better formula is

$$\begin{aligned} s^2 &= \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 &= \frac{1}{N-1} \left( \sum_{i=1}^N x_i^2 - N\bar{x}^2 \right), \\ s^2 &= \frac{1}{N-1} \sum_{i=1}^n f_i (x_i - \bar{x})^2 &= \frac{1}{N-1} \left( \sum_{i=1}^N f_i x_i^2 - N\bar{x}^2 \right), \end{aligned} \quad (2.9)$$

for the *sample* variance for primary or grouped data. The reason the sum is divided by  $N - 1$  instead of  $N$  will have to do with the “bias” of an estimator and will be explained later on in the next chapters. To fully explain why using  $N$  leads to a biased estimate involves the notion of *degrees of freedom*, which takes into account the number of constraints in computing an estimate. The sample observations  $x_1, \dots, x_N$  are independent (by the definition of a random sample), but when computing the variance, we use the variables  $x_1 - \bar{x}, \dots, x_N - \bar{x}$ . Notice that by subtracting the sample mean  $\bar{x}$  from each observation, there exists a linear relation among the elements, namely

$$\sum_{k=1}^N (x_k - \bar{x}) = 0$$

and, thus, we lose 1 degree of freedom due to this constraint. Hence, there are only  $N - 1$  degrees of freedom. So, we will use (2.8) to compute the variance of a set of data that represents a population and (2.9) for the variance of a sample.

**Example 2.11.** Consider again our previous example on CPU times (in seconds) for  $N = 30$  ran-

domly chosen jobs:

70	36	43	69	82	48	34	62	35	15
59	139	46	37	42	30	55	56	36	82
38	89	54	25	35	24	22	9	56	19

Recall that for this data the sample mean was  $\bar{x} = 48.2333$  seconds. The sample variance is

$$s^2 = \frac{(70 - 48.2333)^2 + \dots + (19 - 48.2333)^2}{30 - 1} = \frac{20391}{29} \approx 703.1506 \text{ sec}^2.$$

Alternatively, using (2.8),

$$s^2 = \frac{70^2 + \dots + 19^2 - 30 \cdot 48.2333^2}{30 - 1} = \frac{90185 - 69794}{29} \approx 703.1506 \text{ sec}^2.$$

The sample standard deviation is

$$s = \sqrt{703.1506} \approx 26.1506 \text{ sec}.$$

By the  $3\sigma$  rule, using  $\bar{x}$  and  $s$  as estimates for the population mean  $\mu$  and population standard deviation  $\sigma$ , we may infer that at least 89% of the tasks performed by this processor require between  $\bar{x} - 3s = -30.2185$  and  $\bar{x} + 3s = 126.6851$  (so less than 126.6851) seconds of CPU time.

**Definition 2.12.** *The **coefficient of variation** is the value*

$$CV = \frac{s}{\bar{x}}.$$

**Remark 2.13.**

1. The coefficient of variation is also known as the **relative standard deviation (RSD)**.
2. It can be expressed as a ratio or as a percentage. It is useful in comparing the degrees of variation of two sets of data, even when their means are different.
2. The coefficient of variation is used in fields such as Analytical Chemistry, Engineering or Physics when doing quality assurance studies. It is also widely used in Business Statistics. For example, in the investing world, the coefficient of variation helps brokers determine how much volatility (risk) they are assuming in comparison to the amount of return they can expect from a certain investment. The lower the value of the CV, the better the risk-return trade off.