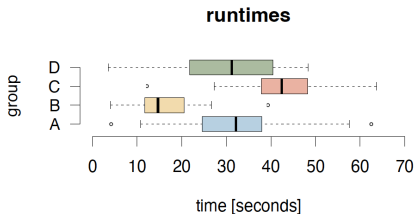


Descriptive Statistics



All examples are fictitious. All data are simulated and the graphics were created with the statistical program package R.

The materials are protected by copyright and are only provided for personal use for studies at TU Vienna. Further use is not permitted. In particular, it is not permitted to distribute the materials or make them publicly available (e.g. in social networks, on learning platforms, etc.).

Sämtliche Beispiele sind frei erfunden. Alle Daten sind simuliert und die Grafiken wurden mit statistischen Programmpaket R erstellt.

Die Materialien sind urheberrechtlich geschützt und dürfen ausschließlich für den Eigengebrauch im Rahmen des Studiums an der TU Wien genutzt werden. Eine weitere Nutzung ist nicht gestattet. Insbesondere ist es nicht gestattet, die Materialien zu verbreiten oder öffentlich zugänglich zu machen (etwa im Rahmen sozialer Netzwerke, Lernplattformen etc.).

Overview

We differentiate:

Probability theory
(Stochastics)

=

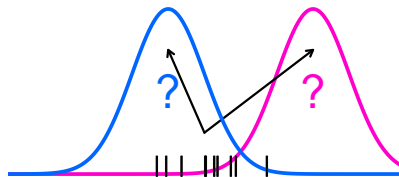
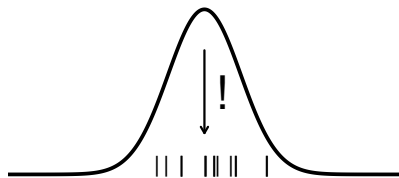
Theory of randomness

and

Statistics

=

Description of data \longrightarrow
(using stochastic **models**)



Today: Short excursion to descriptive Statistics

How do data look like? How can they be summarized?

From then on: inferential Statistics (Modelling)

How did the data occur?

Scales

We differentiate scales

- Categorical data (nominal scale, no ordering)
 - Do you drink coffee? **yes** or **no** (two categories)
 - What is the color of your hair? **blond**, **brown**, **black**, **red**, **neither** (five categories)
- Ordinal data (order, but no metric distance)
 - How much did you learn in the course? **nothing**, **few**, **much** or **very much** (four ordered categories)
 - How often do you use Tuwel? **never**, **sometimes**, **often** (three ordered categories)
- Metric data (Ratio scale, metric distance, $2*3=6$, $0=0$)
 - How large are you? **size in cm**
 - How long is the runtime of an algorithm that you implemented? **time in seconds**

(Today we stick to metric data)

Data collection

How long is the runtime of an algorithm that you implemented?

$n = 121$ students requested (same technical setup)

Results (in seconds):

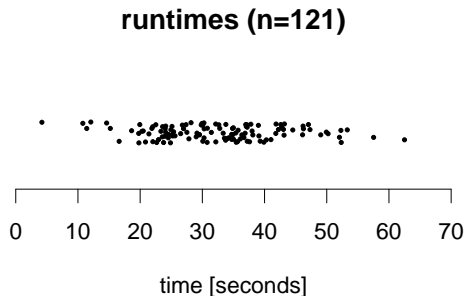
24.6, 24, 31.4, 29.9, 37.8, 19.9, 46.1, 32.8, 30.3, 29, 47.1, 27.8, 33.8, 30.1, 53.3, 23.8, 32.1, 4.2, 42.8, 25.2, 52.3, 35, 30.1, 43.2, 25.4, 62.5, 35.4, 25.2, 37.6, 37.1, 22.9, 29.5, 44.5, 34.8, 33.3, 21.9, 37.2, 24, 37, 34, 24.1, 10.8, 24.9, 37.2, 52, 30.8, 22, 18.6, 22, 26.8, 52.3, 27, 23.6, 33.5, 30.8, 20.9, 35.6, 37.2, 57.5, 46.2, 36.1, 19.8, 38.1, 36.9, 26.5, 23.6, 30.3, 49.9, 39, 50.2, 35.7, 11.4, 24.1, 27.5, 36.4, 29.8, 49, 42.6, 22.5, 32.7, 34.3, 21.4, 34.7, 47.3, 20.3, 35.4, 41.8, 24.9, 15.2, 42.2, 29.1, 25.1, 22.7, 41, 28.2, 30.3, 25.6, 41.8, 16.6, 38, 43.1, 29.5, 40.3, 20.5, 39.9, 24.5, 33.7, 14.6, 23.3, 36.7, 34.7, 34.9, 39.1, 32.2, 43, 12.1, 19.8, 27.4, 39.3, 35, 46.3

We see: n data: $x_1 = 24.6, x_2 = 24.0, \dots, x_n = 46.3$

We understand: nothing?

Thus: descriptive Statistics \rightarrow graphical representation and summary of data

Stripchart

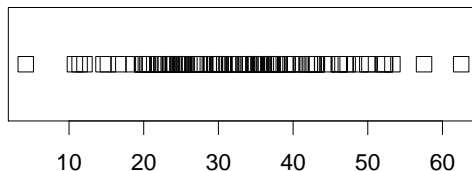


At first sight we understand how the n data distribute:

- Many data lie close to 30 (typical runtime)
- The minimum is about 5 (fastest runtime), the maximum is about 65 (slowest runtime)
- Remark.: the y -value has no meaning. The data are 'jittered' along the y -direction for a better overview.

Stripchart in R

```
#Enter data  
x <- c(24.6, 24.0, 31.4, 29.9,...,39.3, 35.0, 46.3)  
#Create stripchart  
stripchart(x)
```

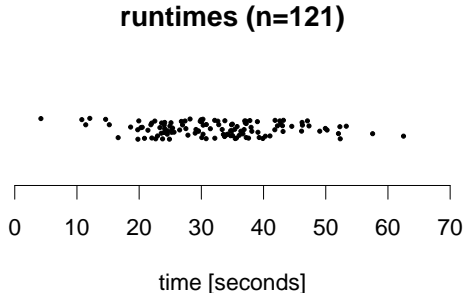


We don't understand too much - points superposed, axes annotations are missing, title is missing etc.

→ customize graphic using additional arguments or lowlevel graphics

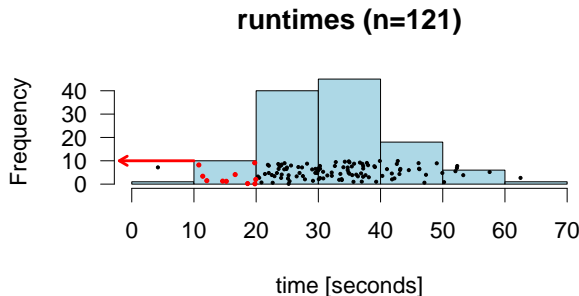
Stripchart in R

```
#Enter data
x <- c(24.6, 24.0, 31.4, 29.9,...,39.3, 35.0, 46.3)
#Create stripchart with additional arguments
stripchart(x,method="jitter",pch=19,cex=0.4,axes=FALSE,
           xlim=c(0,70),main="runtimes_(n=121)",xlab="time_[seconds]")
#add x-axis (lowlevelgraphic)
axis(1,at=seq(0,70,10))
```



Much more informative!

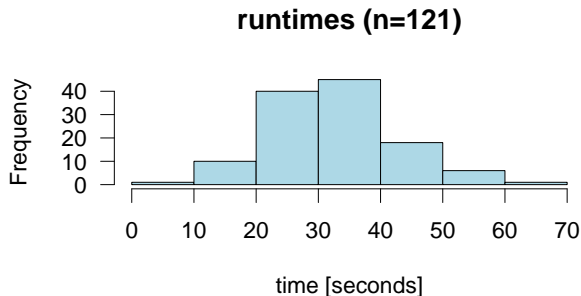
Histogram



- Description of the distribution of data
Here: approximately *bell-shaped*, i.e., unimodal and symmetric
- Absolute frequencies in the intervals $\{(10k, 10(k+1)] : k = 0, 1, \dots, 6\}$
given through the height of the bars
e.g.: **10 data** are > 10 and ≤ 20 , for short $\sum_{i=1}^n \mathbb{1}_{(10,20]}(x_i) = 10$
Consequence: The sum of the bar heights is $n = 121$

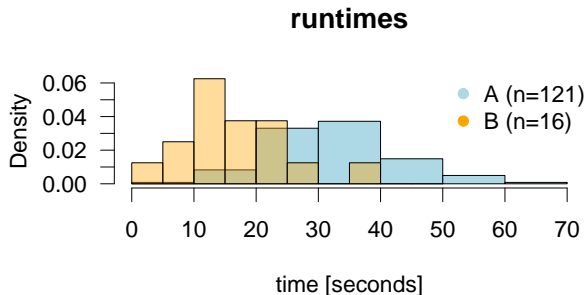
Histogram in R

```
# Histogram with additional arguments  
hist(x, las=1, xlab="time_[seconds]", ylab="Frequency",  
main="runtimes_(n=121)", col="lightblue")
```



Histogram

The same algorithm was implemented by 16 other students after they attended a certain programming course (group B)

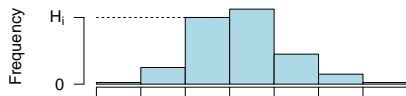


- Comparison of group A ($n_A = 121$) and group B ($n_B = 16$) inappropriate, because the sizes of the groups differ tremendously.
- Idea: Norm the areas \rightarrow total area of 1 each
The distributions are now nicely visible:
shifted against each other and about bell-shaped each.

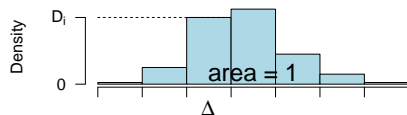
Histogram

What happens when norming?

$$\sum H_i = n$$



$$\sum D_i \times \Delta = 1$$



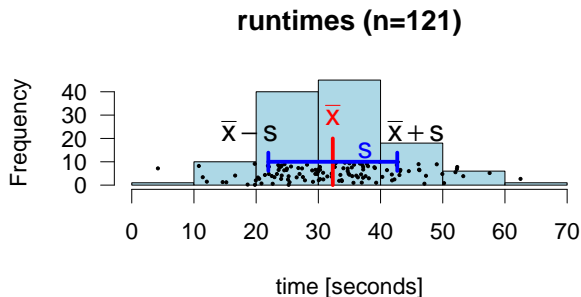
- Same 'picture', but different y -axis

Search D_i such that total area $\sum D_i \cdot \Delta \stackrel{!}{=} 1$

$$\sum H_i = n \Leftrightarrow 1 = \sum \frac{H_i}{n} = \sum \frac{H_i}{n \cdot \Delta} \cdot \Delta, \text{ hence } D_i = \frac{H_i}{n \cdot \Delta}$$

- R normes automatically via `hist(..., prob=TRUE)`

Mean and empirical standard deviation



If the data distribute approximately bell-shaped, then they can be summarized nicely by two prominent *statistics*, i.e., functions of the data:

- 1. the mean \bar{x} → where? (location)
- 2. the (empirical) standard deviation s → how variable? (dispersion)

Mean and empirical standard deviation

Data x_1, x_2, \dots, x_n

- The mean is

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$$

(center of mass of the data)

- The (empirical) variance is

$$s^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

'the mean squared deviation of the data from the mean'

- The (empirical) standard deviation is

$$s = \sqrt{s^2}$$

'the square root of the variance'

Mean and empirical standard deviation

Data x_1, x_2, \dots, x_n

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i \qquad s^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \qquad s = \sqrt{s^2}$$

Random variable X (here discrete)

$$\mathbb{E}[X] := \sum x \cdot \mathbb{P}(X = x) \qquad \mathbb{V}\text{ar}(X) := \mathbb{E}[(X - \mathbb{E}[X])^2] \qquad \sigma_X := \sqrt{\mathbb{V}\text{ar}(X)}$$

Remark:

- The factor $n - 1$ in s^2 (instead of e.g., n) has technical reasons
We speak about the *corrected* empirical variance, while for large n this correction has no practical relevance.
- Analogy to the 'universe of randomness': mean \leftrightarrow expectation

Notation

Convention:

We use *capital letters* for random variables, e.g.,

$$X_1, X_2, \dots, X_n \quad (\text{'random'})$$

and *lowercase letters* for data or realizations of the random variables

$$x_1, x_2, \dots, x_n \quad (\text{'non-random'})$$

Outlook:

The main idea of statistical modelling:

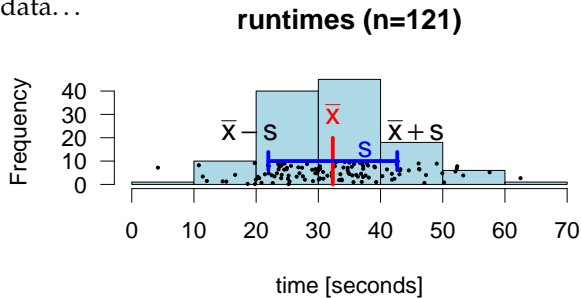
Treat data x_1, x_2, \dots, x_n ('real world')

as realizations of random variables X_1, X_2, \dots, X_n ('universe of randomness')

Note that we evaluate *statistics* either on data, e.g., $\bar{x} = (1/n) \sum^n x_i$ (\rightarrow non-random), or on random variables $\bar{X} = (1/n) \sum^n X_i$ (\rightarrow random)

Mean and empirical standard deviation

Back to the data...



Data x_1, x_2, \dots, x_n

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

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$s = \sqrt{s^2}$$

Evaluation

$$\bar{x} \approx 32.3$$

$$s^2 \approx 107.4$$

$$s \approx 10.4$$

in R via

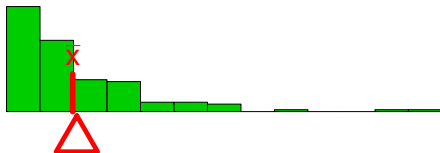
`mean(x)`

`var(x)`

`sd(x)`

Mean and empirical standard deviation

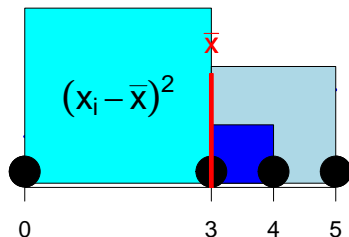
Geometrical interpretation of the mean \bar{x}



- Geometrically: Center of mass
points of same mass on a balance
Where is the **center of rotation** Δ , such that the balance is in **equilibrium**?
- Consequence: Naive estimation from graphic
Distribution not bell-shaped but *asymmetric*
few large values 'pull' \bar{x} to the right

Mean and empirical standard deviation

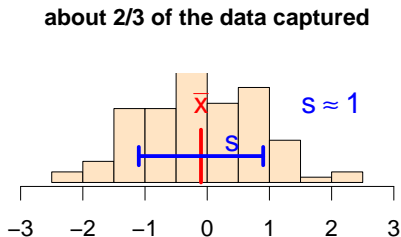
For the standard deviation s



- numerically: $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{3} (3^2 + 0^2 + 1^2 + 2^2) = \frac{14}{3} \rightarrow s = \sqrt{\frac{14}{3}}$
- Large deviations from the mean have a large impact (squaring)

Mean and empirical standard deviation

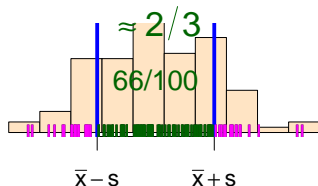
Naive estimation of s (only for bell-shaped distributions!)



- Fact: About 2/3 of the data lie in the s -neighborhood of \bar{x}
- Turn the tables
 - Estimate \bar{x} (\rightarrow balance)
 - Capture 2/3 of the data around \bar{x}
- Numerically: $\bar{x} \approx -0.1$ and $s \approx 0.94$

Mean and empirical standard deviation

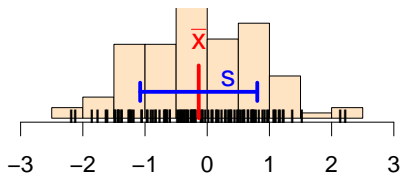
We used: For a bell-shaped distribution about $2/3$ of the data lie in the s -neighborhood of \bar{x} . But why?



- Recall: Normal distribution $N(\mu, \sigma^2)$
 - Let $X \sim N(\mu, \sigma^2)$. Then $\mathbb{P}(X \in [\mu - \sigma, \mu + \sigma]) \approx 0.68 \approx 2/3$
 X falls in the σ -neighborhood of μ with probability about $2/3$
 - Consider data $n = 100$ independent copies X_1, \dots, X_n of X
data x_1, \dots, x_n are interpreted as realizations of X_1, \dots, X_n , reasonable as data is approx bell-shaped
 - The *proportion* within $\mu \pm \sigma$ lies close to $2/3$ (\rightarrow Law of large numbers)
 - \bar{X} and S consistently *estimate* μ and σ (\rightarrow Law of large numbers)
 - Also the *proportion* within $\bar{X} \pm S$ is close to $2/3$

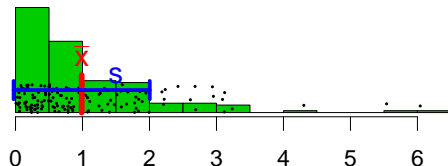
Mean and empirical standard deviation

Interpretation (only for bell-shaped distributions of data)



- \bar{x} is interpreted as a *typical observation*
- s is interpreted as the *typical deviation* of an observation (from the mean)
- These two statistics (only two!) suitably summarize the whole set of data (many!)

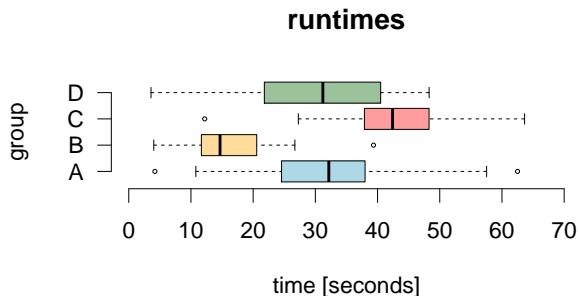
Mean and empirical standard deviation



- If the data are not distributed approximately bell-shaped, then this interpretation is not useful
- Here \bar{x} is not a typical observation. Much more data lie left of \bar{x} than right of it
- s does not describe the typical deviation of \bar{x} . Almost all of the data lie within the s -neighborhood of \bar{x} , only few outliers lie outside of it
- \bar{x} and s should not be used for the description of the location and the dispersion of the data

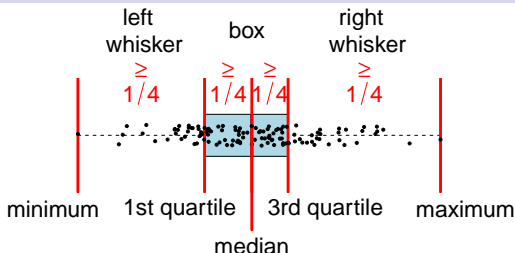
Boxplot

Comparison of four groups *A*, *B*, *C* and *D*



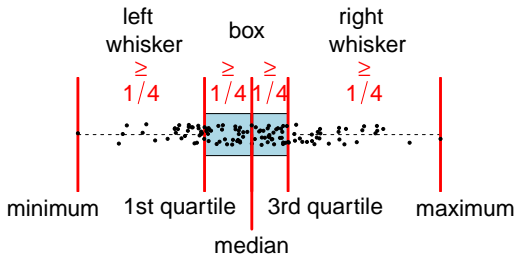
- Histograms overplotted
Could represent the data in a stripchart
Other possibility: the *box and whisker plot*, short boxplot

Boxplot



- Consists of a box and two whisker ('Schnurrhaare', meow!)
- Four sections, contain at least $1/4$ of the data
- → five statistics:
 - *Minimum*, smallest observation
 - *Maximum*, largest observation
 - *Median* (m), at least 50% of the data $\geq m$ and at least 50% are $\leq m$
 - *1st quartile* ($q_{1/4}$), at least 25% are $\leq q_{1/4}$ and at least 75% are $\geq q_{1/4}$
 - *3rd quartile* ($q_{3/4}$), at least 75% are $\leq q_{3/4}$ and at least 25% are $\geq q_{3/4}$
- Interpretation:
 - Median m is a measure for the location of the observations (→ where?)
 - Interquartile range $q_{3/4} - q_{1/4}$ (width of the box) is a measure for the dispersion of the data (→ how variable?)

Boxplot



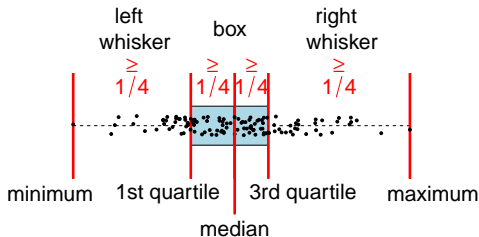
Empirical quantile (general)

- Definition: Given n data x_1, \dots, x_n . Let $p \in (0, 1)$. A number $q_p \in \mathbb{R}$ is called an (empirical) p -quantile, if
 - i. the proportion of the data that are smaller or equal q_p is at least p and
 - ii. the proportion of the data that are larger or equal q_p is at least $1 - p$.

In formulas:

$$i.: \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(-\infty, q_p]}(x_i) \geq p \quad \text{and} \quad ii.: \quad \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[q_p, \infty)}(x_i) \geq 1 - p$$

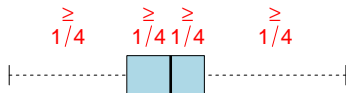
- We already know three prominent candidates (with their own name):
 - a median is a 50%-quantile ($p = 1/2$)
 - a 1st quartile is a 25%-quantile ($p = 1/4$)
 - a 3rd quartile is a 75%-quantile ($p = 3/4$)



Empirical quantile (general)

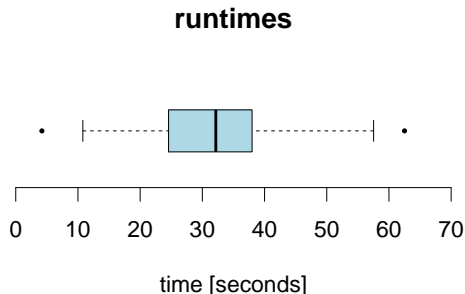
- Example: Four observations $x = (1, 2, 3, 4)^t$ superscript t denotes the transpose
 - Many medians: Every number in the interval $[2, 3]$ is a median
 - Often: Define the *unique* median as the mean value of the bounds, here 2.5
 - Analog: Every number in $[1, 2]$ is 1/4-quantile, the unique quartile is 1.5
 - Many quantiles equal: The number 2 is a p -quantile for every p of $[0.25, 0.5]$
- Remark.: These kind of 'exotic' messages may support the understanding of the definition of a quantile. The main message however is, that the boxplot appropriately summarizes many data using only five simple statistics

Take home: Many data \rightarrow at first sight: "1/4, 1/4, 1/4, 1/4"



Boxplot in R

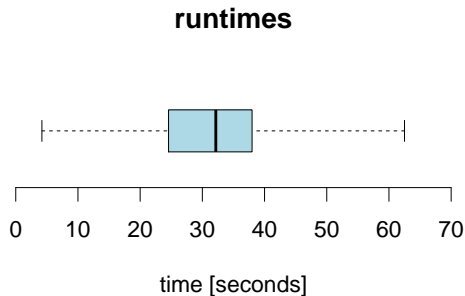
```
#Boxplot, horizontal representation  
boxplot(x, horizontal=TRUE, ...)
```



Attention: per default a whisker ranges to the observation which is most far away from the box, but does not exceed 1.5 times the interquartile range. Extreme values ('outliers') are plotted seperately.

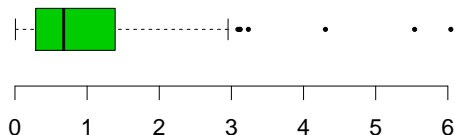
Boxplot in R

```
#Boxplot, Whisker up to the most extreme values  
boxplot(x, horizontal=TRUE, range=0, ...)
```



Through the argument `range=0` the whiskers are extended to the extreme values

Boxplot



Reminder

- due to the asymmetric distribution of the data, \bar{x} and s should not be used for the description of the location and the dispersion
- The five statistics of the boxplot are more appropriate for the description of the data

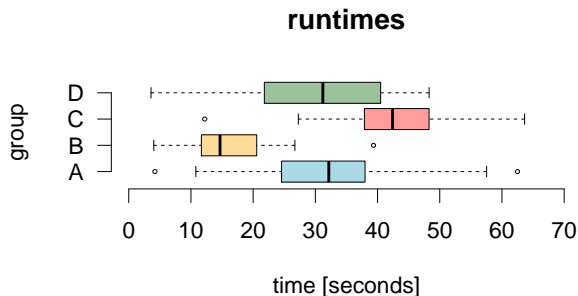
Most important message today

Always graphically visualize
your data first

(and start computing afterwards)

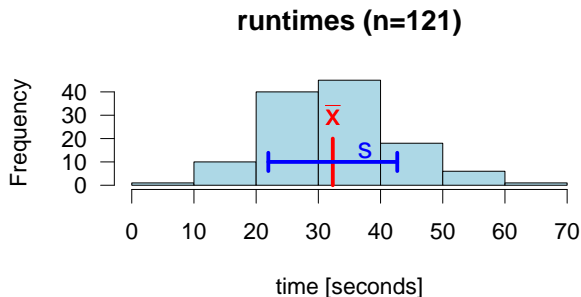
Questions

Comparison of four groups A, B, C und D



- The slowest runtime in C was about? 65
- The fastest runtime in A is about? 5
- The median runtime in D is about? 30
- What is the percentage of runtimes in group B that are smaller than 20?
about 75%
- Were 50% of the runtimes in A faster than 75% of the times in C? yes
- In group B, apart from a single runtime all others were faster than half of those of group A, half of those of C and half of those of D. Correct

Questions



- What is the mean runtime? about 32
- The standard deviation of the runtimes is about 10

Multiple-choice questions

(1) Regarding the data

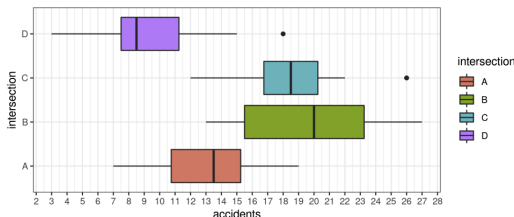
11, 21, 22, 9, 3, 5

it holds

- a. 8 is a median
- b. the set of 50%-quantiles is $[11, 21]$
- c. 22 is a $5/6$ -quantile
- d. the 1.-quartile is not unique.

Multiple-choice questions

- (2) Data on the number of yearly accidents were collected from four intersections (A-D) over a 20 year period. The corresponding boxplots are given below.



Which of the following statements is **false**.

- During at least 75% of years, intersection D had viewer fewer accidents than the lowest 25% of years at intersection A.
- During at least 15 years, fewer than 12 accidents occurred at intersection D.
- The maximum number accidents that occurred in a single intersection was 27.
- All of the accidents totals at intersection D were lower than the median number of accidents at intersection B.

Multiple-choice questions

- (3) If the standard deviation of a set of observations is zero, we can conclude
- a. that there is no relationship between the observations.
 - b. that all observations are the same value.
 - c. that the average value is 0.
 - d. that a mistake in arithmetic has been made.
- (4) Suppose the average score on a certain test is 500 with a standard deviation of 100. If each score is increased by 25, what are the new mean and standard deviation?
- a. 500, 100
 - b. 500, 125
 - c. 525, 100
 - d. 525, 125

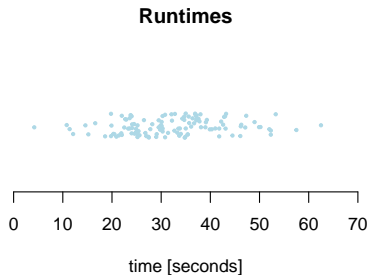
Multiple-choice questions

(5) Which one of the following is an **incorrect** statement?

- a. The sampling distribution of \bar{x} has mean equal to the population mean μ even if the population is not normally distributed.
- b. The sampling distribution of \bar{x} has standard deviation σ/\sqrt{n} even if the population is not normally distributed.
- c. When n is large, the sampling distribution of \bar{x} is approximately normal even if the population is not normally distributed.
- d. The larger the value of the sample size n , the closer the standard deviation of the sampling distribution of \bar{x} is to the standard deviation of the population.

Multiple-choice questions

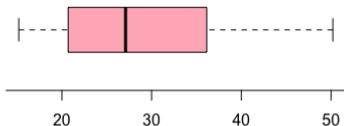
(6) Which R-command produces the following graphical representation?



- a. boxplot
- b. stripchart
- c. fivenum
- d. qqplot

Multiple-choice questions

(7) The boxplot of a given sample is given below.



Which one of the four statements is correct?

- a. The sample size is 50.
- b. The interquartile range is about 10.
- c. There are no outliers.
- d. The upper quartile is about 50.

Multiple-choice questions

(8) Which of the following are true statements?

- I If the sample has variance zero, the variance of the population is also zero.
- II If the population has variance zero, the variance of the sample is also zero.
- III If the sample has variance zero, the sample mean and the sample median are equal.

- a. I and II
- b. I and III
- c. II and II
- d. None of the above gives the complete set of true responses.

Multiple-choice questions

- (8) A financial analyst considers the following sample of book value of six companies (in Euro)

25, 7, 22, 33, 18, 15.

The sample mean and sample standard deviation are (approximately)

- a. 20 and 79.2
- b. 20 and 8.9
- c. 120 and 79.2
- d. 120 and 8.9

Multiple-choice questions

- (9) A sample of 99 distances has a mean of 24 centimeters and a median of 24.5 centimeters. Unfortunately, it has just been discovered that an observation which was recorded as “30” actually had a value of “35”. If we make this correction to the data, then
- a. the mean remains the same, but the median is increased
 - b. the mean and median are both increased
 - c. the median remains the same, but the mean is increased
 - d. we do not know how the mean and median are affected without further calculations, but the variance is increased

Multiple-choice questions

- (10) Earthquake intensities are measured using a device called a seismograph which is designed to be most sensitive for earthquakes with intensities between 4.0 and 9.0 on the open-ended Richter scale. Measurements of ten earthquakes gave the following data

4.5, L , 5.5, H , 8.7, 8.9, 6.0, H , 5.2, 7.2,

where L indicates that the earthquake had an intensity below 4.0 and a H indicates that the earthquake had an intensity above 9.0.

One measure of central tendency is the $x\%$ trimmed mean computed after trimming $x\%$ of the upper values and $x\%$ of the bottom values.

The value of the 20% trimmed mean is

- a. 6.0
- b. 6.6
- c. 6.9
- d. cannot be computed because all of the values are not known.

Thank you for your attention!