

Statistics in Computer Science

Seminar Exercises

Stanislav Katina, Mojmir Vinkler

November 26, 2016

1 Location characteristics

Exercise 1 (Code Vectorization). *Implement two versions of mean function in R - the first one is `mean_slow = function(x)` and use for loop to sum numbers. The second version is `mean_fast = function(x)` and use built-in function `sum` to sum numbers. Generate very long random vector of uniformly distributed random numbers and compare performance of both functions and also of built-in mean function.*

Hints. Use `runif(...)` for generating random numbers and `system.time({...})` for profiling.

Realizations will be denoted as x_1, x_2, \dots, x_n , **sorted realizations** as $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$. Then we define following estimations of location characteristics (sample location characteristics):

- **sample minimum** X_{\min} , with realization $x_{\min} = x_{(1)}$;
- **sample maximum** X_{\max} , with realization $x_{\max} = x_{(n)}$;
- **sample (arithmetic) mean** \bar{X} , with realization $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{j=1}^{n_j} x_j f_j, n_j \leq n$, where f_j are frequencies (counts) of x_j and $n = \sum_j f_j$;
- **sample mode** X_{mod} , with realization x_{mod} is the most common value (in case of discrete variable it is value x in which probability function has its maximum; in case of continuous variable it is value x in which density function has its maximum);
- **sample median** \tilde{X} (robust estimation of location), with realization

$$\tilde{x} = \begin{cases} x_{(\frac{n+1}{2})} & \text{if } n \text{ is odd,} \\ \frac{1}{2} \left(x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)} \right) & \text{if } n \text{ is even;} \end{cases}$$

distribution is *symmetric*, if $\bar{x} = \tilde{x} = x_{\text{mod}}$, distribution is *positively skewed* (right), if $\bar{x} > \tilde{x} > x_{\text{mod}}$ and distribution is *negatively skewed* (left), if $\bar{x} < \tilde{x} < x_{\text{mod}}$;

- **sample percentile** \tilde{X}_p (read as 100p-percentile), with realization \tilde{x}_p defined as

$$\tilde{x}_p = \begin{cases} x_{(k+1)} & \text{for } k \neq np, \\ \frac{1}{2} (x_{(k)} + x_{(k+1)}) & \text{for } k = np, \end{cases}$$

where $k = \lfloor np \rfloor$, which is floor of np ;

- **sample quartiles** there are three

- **sample first (lower) quartile** Q_1 , with realization $\tilde{x}_{0.25}$ is a value that splits off the lowest 25% of data from the highest 75%,

$$\Pr [x_{\min}, \tilde{x}_{0.25}] = \Pr [X \leq \tilde{x}_{0.25}] = \frac{1}{4}, \Pr [\tilde{x}_{0.25}, x_{\max}] = \Pr [X \geq \tilde{x}_{0.25}] = \frac{3}{4};$$

- **sample second quartile** (median) Q_2 , with realization $\tilde{x}_{0.5} = \tilde{x}$ is a value that splits off the lowest 50% of data from the highest 50%,

$$\Pr [x_{\min}, \tilde{x}_{0.5}] = \Pr [X \leq \tilde{x}_{0.5}] = \frac{1}{2}, \Pr [\tilde{x}_{0.5}, x_{\max}] = \Pr [X \geq \tilde{x}_{0.5}] = \frac{1}{2};$$

- **sample third (upper) quartile** Q_3 , with realization $\tilde{x}_{0.75}$ is a value that splits off the lowest 75% of data from the highest 25%,

$$\Pr [x_{\min}, \tilde{x}_{0.75}] = \Pr [X \leq \tilde{x}_{0.75}] = \frac{3}{4}, \Pr [\tilde{x}_{0.75}, x_{\max}] = \Pr [X \geq \tilde{x}_{0.75}] = \frac{1}{4};$$

- **sample deciles** \tilde{X}_k , with realizations \tilde{x}_k splits data to ten buckets, i.e. $k/10$ of data are lower than a decile and $(10 - k)/10$ are higher, where $k \in \{0, 1, \dots, 10\}$;
- **sample five-number summary** $(X_{\min}, Q_1, Q_2, Q_3, X_{\max})^T$, with realizations $(x_{\min}, \tilde{x}_{0.25}, \tilde{x}_{0.50}, \tilde{x}_{0.75}, x_{\max})^T$.

Robust location characteristics (resistant to outliers) are

- **sample γ -truncated arithmetic average** \bar{X}_g , with realization \bar{x}_g that is calculated as

$$\bar{x}_g = \frac{1}{n - 2g} (x_{(g+1)} + x_{(g+2)} + \dots + x_{(n-g)}),$$

where $g = \lfloor \gamma n \rfloor$, $\gamma = 0.1, 0.2$. More than $\gamma 100$ % observations must be replaced for the γ -truncated average to become large or small relative to the original [¹breakdown point \bar{x}_g is therefore γ],

¹Breakdown point represents number of observations we need to significantly change value of location characteristics. For γ -truncated and γ -winsorized arithmetic average it is γn observations, for median $n/2$ observations and for simple arithmetic average changing just one observation is enough (that's the reason we say that arithmetic average is very sensitive to outliers).

- **sample γ -winsorized arithmetic average** \bar{X}_w , with realization \bar{x}_w is defined as

$$\bar{x}_w = \frac{1}{n} \left((g+1)x_{(g+1)} + x_{(g+2)} + \dots + (g+1)x_{(n-g)} \right).$$

More than $\gamma 100\%$ must be replaced for the γ -winsorized average to become large or small relative to the original [*breakdown point* \bar{x}_w is therefore γ].

Exercise 2 (height of 10-year old girls). *Let's have $n = 12$ heights (in cm) of randomly sampled 10-year old girls sorted by height (**order** denoted as r_i for $x_{(i)}$; in case the values are equal, r_i is calculated as average of their order numbers).*

Table 1: Sorted realizations x_i and their order r_i for heights of 10-year old girls

i	1	2	3	4	5	6	7	8	9	10	11	12
$x_{(i)}$	131	132	135	141	141	141	141	142	143	146	146	151
r_i	1	2	3	5.5	5.5	5.5	5.5	8	9	10.5	10.5	12

Then $\bar{x} = 140.83$, $\tilde{x} = \frac{1}{2}(x_{(6)} + x_{(7)}) = 141$, $\tilde{x}_{0.25} = \frac{1}{2}(x_{(3)} + x_{(4)}) = 138$, where $k = \lfloor 12 \times 0.25 \rfloor = 3$, $Q_3 = \tilde{x}_{0.75} = \frac{1}{2}(x_{(9)} + x_{(10)}) = 144.5$, where $k = \lfloor 12 \times 0.75 \rfloor = 9$.

Write functions for calculation of all location characteristics. Verify your functions on characteristics above. Don't use built-in functions for location characteristics such as *mean*, *quantile*, etc. Use $\gamma = 0.1$ for truncated and winsorized arithmetic averages.

2 Spread (variability) characteristics

Then we define following estimations of spread (variability) characteristics (sample spread characteristics):

- **sample variance** S^2 , with realization

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

under linear transformation sample variance changes like this²

$$s_y^2 = s_{a+bx}^2 = b^2 s_x^2,$$

i.e.

$$\begin{aligned} s_y^2 &= s_{a+bx}^2 = \frac{1}{n-1} \sum_{i=1}^n (a + bx_i - \overline{a+bx})^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (a + bx_i - (a + b\bar{x}))^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (b(x_i - \bar{x}))^2 = b^2 s_x^2; \end{aligned}$$

²Equation tells us that variance of shifted and rescaled variable y is equal to square of scale multiplied by variance of original variable x .

- **sample standard deviation** S , with realization

$$s = s_{n-1} = s_x = \sqrt{s_x^2};$$

under linear transformation standard deviation changes like this

$$s_y = s_{a+bx} = |b| s_x,$$

- **coefficient of variation** V_k , with realization v_k represents normalized form of standard deviation (inversion of *signal-to-noise ratio*; fraction of variability to mean)

$$v_k = \frac{s_x}{\bar{x}};$$

it is usually denoted in percentage points, i.e. $100 \times (s_x/\bar{x}) \%$ and can be used only for realizations with positive values;

- **sample variance of arithmetic average** $S_{\bar{X}}^2$, with realization

$$s_{\bar{x}}^2 = \frac{s_x^2}{n};$$

- **sample standard deviation of arithmetic average (sample standard error)** $S_{\bar{X}}$, with realization

$$s_{\bar{x}} = \frac{s_x}{\sqrt{n}};$$

- **sample skewness** B_1 , with realization

$$b_1 = \frac{n^{-1} \sum_{i=1}^n (x_i - \bar{x})^3}{[n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2]^{3/2}} = \frac{\sqrt{n} \sum_{i=1}^n (x_i - \bar{x})^3}{[\sum_{i=1}^n (x_i - \bar{x})^2]^{3/2}},$$

distribution is *symmetric*, if $b_1 = 0$, *positive skewness* (density on the left side is steeper than the right side), if $b_1 > 0$ a *negative skewness* (density on the right side is steeper than the left side), if $b_1 < 0$;

- **sample kurtosis** B_2 , with realization

$$b_2 = \frac{n^{-1} \sum_{i=1}^n (x_i - \bar{x})^4}{[n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2]^2} - 3 = \frac{n \sum_{i=1}^n (x_i - \bar{x})^4}{[\sum_{i=1}^n (x_i - \bar{x})^2]^2} - 3,$$

distribution is normal (*mesokurtic*), if $b_2 = 0$, pointy (*leptokurtic*), if $b_2 > 0$ and flat (*platykurtic*), if $b_2 < 0$;

- **sample sum of squares** $SS = \sum_{i=1}^n (X_i - \bar{X})^2$, with realization

$$SS_{\text{obs}} = \sum_{i=1}^n (x_i - \bar{x})^2,$$

- **sample sum of absolute deviation** $SAD = \sum_{i=1}^n |X_i - \tilde{X}_{0.5}|$, with realization

$$SAD_{\text{obs}} = \sum_{i=1}^n |x_i - \tilde{x}_{0.5}|;$$

- **sample arithmetic average deviation** $MAD = \frac{1}{n} \sum_{i=1}^n |X_i - \tilde{X}_{0.5}|$, with realization

$$MAD_{\text{obs}} = SAD_{\text{obs}}/n;$$

- **sample range** $D = X_{\max} - X_{\min}$, with realization

$$d_{\text{obs}} = x_{\max} - x_{\min};$$

- **sample interquartile range** $D_Q = Q_3 - Q_1$, with realization

$$d_Q = \tilde{x}_{0.75} - \tilde{x}_{0.25};$$

distribution is (between quartiles) *symmetric*, if $\tilde{x}_{0.75} - \tilde{x}_{0.50} = \tilde{x}_{0.50} - \tilde{x}_{0.25}$, *positively skewed*, if $\tilde{x}_{0.75} - \tilde{x}_{0.50} > \tilde{x}_{0.50} - \tilde{x}_{0.25}$ and *negatively skewed*, if $\tilde{x}_{0.75} - \tilde{x}_{0.50} < \tilde{x}_{0.50} - \tilde{x}_{0.25}$;

- **sample decile range** $D_D = \tilde{X}_{0.9} - \tilde{X}_{0.1}$, with realization

$$d_D = \tilde{x}_{0.9} - \tilde{x}_{0.1};$$

- **sample percentile range** $D_P = \tilde{X}_{0.99} - \tilde{X}_{0.01}$, with realization

$$d_P = \tilde{x}_{0.99} - \tilde{x}_{0.01}.$$

Robust spread characteristics (variability) are

- **sample γ -truncated variance** S_{tg}^2 , with realization s_{tg}^2 calculated as

$$s_{tg}^2 = \frac{1}{n - 2g - 1} \sum_{i=g+1}^{n-g} (x_{(i)} - \bar{x}_{tg})^2;$$

more than $\gamma 100\%$ must be replaced so that γ -truncated variance changes to large or small relative to the original s^2 [*breakdown point* s_{tg}^2 is γ]; it applies that $s_{tg}^2 < s^2$ because truncating removes outliers;

- **sample γ -winsorized variance** S_{wg}^2 , with realizations s_{wg}^2 calculated as

$$s_{wg}^2 = \frac{1}{n - 1} ((g + 1)(x_{(g+1)} - \bar{x}_{wg})^2 + (x_{(g+2)} - \bar{x}_{wg})^2 + \dots + (g + 1)(x_{(n-g)} - \bar{x}_{wg})^2);$$

more than $\gamma 100\%$ must be replaced so that *gamma*-winsorized variance changes to large or small relative to the original s^2 [*breakdown point* s_{wg}^2 is γ]; it applies that $s_{wg}^2 < s^2$ because winsorization pulls outliers closer to the mean;

- **sample quartile coefficient of variation** $V_{k,Q} = (Q_3 - Q_1)/Q_2$, with realization $v_{k,Q}$ calculated as

$$v_{k,Q} = \frac{\tilde{x}_{0.75} - \tilde{x}_{0.25}}{\tilde{x}}.$$

Other robust spread characteristics characterized by modified boundaries are

- **sample robust minimum and maximum** (“inner boundaries”) $X_{\min}^* = B_L = Q_1 - 1.5D_Q$ and $X_{\max}^* = B_U = Q_3 + 1.5D_Q$, with realizations defined as

$$x_{\min}^* = b_L = \tilde{x}_{0.25} - 1.5(\tilde{x}_{0.75} - \tilde{x}_{0.25}),$$

$$x_{\max}^* = b_U = \tilde{x}_{0.75} + 1.5(\tilde{x}_{0.75} - \tilde{x}_{0.25}),$$

values outside of boundaries are considered to be *suspicious, potential outliers*;

- **sample robust minimum and maximum** (“outer boundaries”) defined as $B_U^* = Q_1 - 3(Q_3 - Q_1)$, $B_L^* = Q_3 + 3(Q_3 - Q_1)$, with realizations $b_L^* = \tilde{x}_{0.25} - 3d_Q$, $b_U^* = \tilde{x}_{0.75} + 3d_Q$;

- if there are any $x_i < b_L^* \vee x_i > b_U^*$, we call them *distant values*³,
- if $x_i \in \langle b_L^*, b_L \rangle \vee \langle b_U, b_U^* \rangle$, these are *outer values*,
- if $x_i \in \langle b_L, b_U \rangle$, these are *inner values* or *values close to median*;
- normal distribution has these properties $B_U - B_L = Q_3 + 1.5D_Q - Q_1 + 1.5D_Q = 4D_Q \doteq 4.2$; probability of $x_i \notin \langle B_L, B_U \rangle$ is then 0.04;

- **sample robust skewness** B_{1Q} and B_{1O} and their variance under asymptotic normality $B_{1\cdot}$, where $\cdot = Q$ or O , with realizations defined as

- quartile skewness

$$b_{1Q} = \frac{(\tilde{x}_{0.75} - \tilde{x}_{0.50}) - (\tilde{x}_{0.50} - \tilde{x}_{0.25})}{\tilde{x}_{0.75} - \tilde{x}_{0.25}}, \text{Var}_{as}(b_{1Q}) = 1.84,$$

- octile skewness

$$b_{1O} = \frac{(\tilde{x}_{0.875} - \tilde{x}_{0.50}) - (\tilde{x}_{0.50} - \tilde{x}_{0.125})}{\tilde{x}_{0.875} - \tilde{x}_{0.125}}, \text{Var}_{as}(b_{1O}) = 1.15.$$

Exercise 3 (height of 10-year old girls). *Calculate all spread characteristics for the sample with heights of 10-year old girls.*

³Symbol \vee means “or” and symbol \wedge means “and”.

3 Basics of Probability

Exercise 4 (Simple random sample). *In a simple random sample of size n from population of finite size N , each element has an equal probability of being chosen. If we avoid choosing any member of the population more than once, we call it **simple random sample without replacement**⁴. (Dalgaard 2008). If we put a member back to population after choosing it, we talk about **simple random sample with replacement**⁵. Let's have a set \mathcal{M} with $N = 10$ elements and we want to choose $n = 3$ elements (a) without replacement and (b) with replacement. How many combinations there are? How do these combinations look like if $\mathcal{M} = \{1, 2, \dots, 10\}$? Do the same for $N = 100$, $n = 30$ and set $\mathcal{M} = \{1, 2, \dots, 100\}$.*

Solution without  code:

(a) Number of combinations is $\binom{N}{n}$. If $N = 10$ and $n = 3$, then $\binom{N}{n} = \frac{N!}{(N-n)!n!} = \binom{10}{3} = 120$.
(b) Number of combinations with replacement is $\binom{N+n-1}{n}$. If $N = 10$ and $n = 3$, then $\binom{N+n-1}{n} = \frac{(N+n-1)!}{(N-1)!n!} = \binom{10+3-1}{3} = 220$. If $N = 100$ and $n = 30$, then $\binom{N+n-1}{n} = \binom{100+30-1}{30} = 2.009491 \times 10^{29}$.

Hints. `choose(n,k)`, `combn(n,k)`⁶, `sample(x=..., size=..., replace=...)`

Exercise 5 (Simple random sample). *A group of people are labeled by their identification numbers (ID) from 1 to 30. Choose (a) randomly 5 people out of 30 without replacement, (b) randomly 5 people out of 30 with replacement and finally (c) randomly 5 people out of 30 without replacement, where people with ID between 28 and 30 have 4× higher probability of being chosen than people with ID between 1 and 27.*

Repeat the sampling 100 times, concatenate all samples and use `table` to get frequencies and confirm that IDs 28 to 30 are indeed sampled more frequently.

Hints. `rep(1, k)`, `sample(..., prob=p)`

Exercise 6 (Normal distribution). *Let X be a random variable (it could represent for example adult height) and let's assume it is normally distributed with parameters μ (expectation or mean) and σ^2 (standard deviation) which is written as $X \sim N(\mu, \sigma^2)$, $\mu = 140.83$, $\sigma^2 = 33.79$. Normal distribution represents a probability distribution model for this random variable. Calculate probability $\Pr(a < X < b) = \Pr(X < b) - \Pr(X < a) = F_X(b) - F_X(a)$, where $a = \mu - k\sigma$, $b = \mu + k\sigma$, $k = 1, 2, 3$.⁷ Plot probability density function and fill area between points a and b and label axes x and y as shown in figure 1.*

Partial solution:

$$a = \mu - \sigma = 135.0171, b = \mu + \sigma = 146.6429,$$

$$\Pr(|X - \mu| > \sigma) = 0.3173, \Pr(|X - \mu| < \sigma) = 1 - 0.3173 = 0.6827,$$

$$a = \mu - 2\sigma = 129.2042, b = \mu + 2\sigma = 152.4558,$$

$$\Pr(|X - \mu| > 2\sigma) = 0.0455, \Pr(|X - \mu| < 2\sigma) = 1 - 0.0455 = 0.9545,$$

⁴ n -combination without replacement from N members of set \mathcal{M} .

⁵ n -combination with replacement from N members of set \mathcal{M} .

⁶requires library `utils`

⁷Note that $\Pr(a < X < b) = \Pr(a \leq X \leq b)$ because probability of a point (here a and b) is zero for continuous random variables, i.e $\Pr(a) = \Pr(b) = 0$. This does not apply to discrete random variables.

$a = \mu - 3\sigma = 123.3913$, $b = \mu + 3\sigma = 158.2687$,
 $\Pr(|X - \mu| > 3\sigma) = 0.0027$, $\Pr(|X - \mu| < 3\sigma) = 1 - 0.0027 = 0.9973$.

```

1 mu <- 0
2 sig <- 1
3 bin <- seq(mu-3*sig,mu+3*sig,by=sig)
4 pnorm(bin[7]) - pnorm(bin[1]) # 0.9973002
5 pnorm(bin[6]) - pnorm(bin[2]) # 0.9544997
6 pnorm(bin[5]) - pnorm(bin[3]) # 0.6826895
  
```

Probabilities 68.27 – 95.45 – 99.73 are called (empirical) rule or 3-sigma rule.

Hints. Similar example for exponential distribution.

```

1 # probability density curve
2 x = seq(0, 4, by=0.01)
3 y = dexp(x)
4
5 # tell R we will plot 3 graphs in one row
6 par(mfrow=c(1,3))
7
8 # plot three graphs
9 for(p in c(0.2, 0.1, 0.05)){
10   # plot probability density as a line
11   plot(x, y, type='l', xlab=paste0('Area under the curve = ', p), ylab='density')
12
13   # fill area under the curve
14   p_quantile = qexp(1-p)
15   xx = seq(4, p_quantile, by=-0.01)
16   pol_x = c(p_quantile, 4, xx)
17   pol_y = c(0, 0, dexp(xx))
18   polygon(pol_x, pol_y, col = "grey")
19 }
  
```

Exercise 7 (Normal distribution). Let $X \sim N(\mu, \sigma^2)$, where $\mu = 150$, $\sigma^2 = 6.25$. Calculate $a = \mu - x_{1-\frac{\alpha}{2}}\sigma$ and $b = \mu + x_{1-\frac{\alpha}{2}}\sigma$ so that $\Pr(a \leq X \leq b) = 1 - \alpha$ is equal to 0.90, 0.95 a 0.99. Value $x_{1-\alpha}$ is a quantile of standardized normal distribution, i.e. $\Pr(Z = \frac{X-\mu}{\sigma} < x_{1-\alpha}) = 1 - \alpha$, $Z \sim N(0, 1)$. Plot density function, fill area between points a and b and label axes x a y as shown in the following figure.

Hints. `qnorm(alpha)`

This gives us rule 90 – 95 – 99 (so called **adjusted 3-sigma rule**). We used property $\Pr(u_{\alpha/2} < Z < u_{1-\alpha/2}) = \Phi(u_{1-\alpha/2}) - \Phi(u_{\alpha/2}) = 1 - \alpha$, where Φ is cumulative distribution function of normal distribution and in general $\alpha \in (0, 1/2)$; in the exercise we used $\alpha = 0.1$, 0.05 a 0.01. See figure 2.

Exercise 8 (Interactive Normal Distribution). Create an interactive *Shiny application* showing probability density function of normal distribution $N(\mu, \sigma^2)$, where $\mu = 0$ and σ can be interactively set by user. Fill area between points a and b , where $a = \mu - k\sigma$, $b = \mu + k\sigma$, $k = 1, 2, 3$ and make k interactive too. Finally, add graph title showing probability $\Pr(a \leq X \leq b)$.

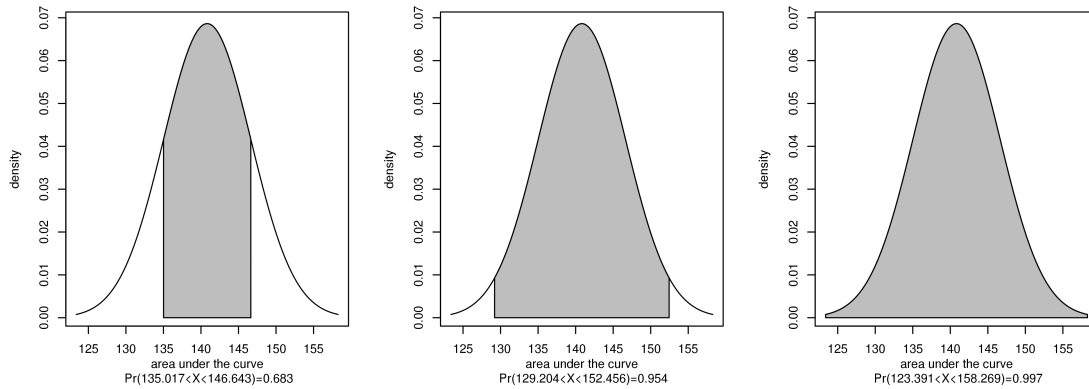


Figure 1: 3-sigma rule; density curve with colored area under this curve between corresponding quantiles on x -axis; volume of the area is equal to the probability of realization of random variable between these quantiles

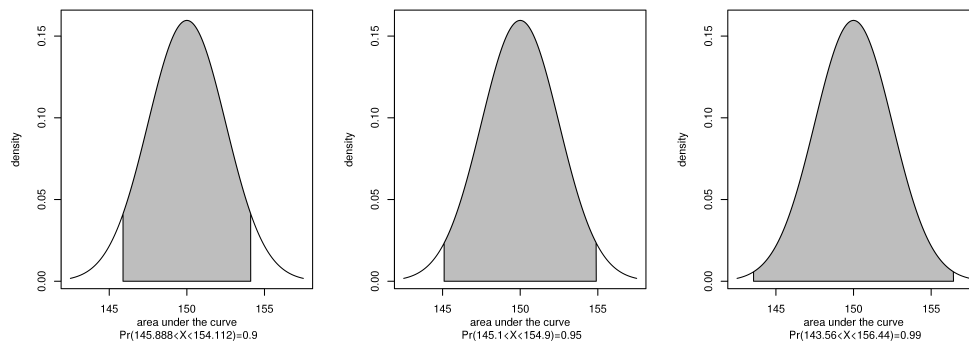


Figure 2: Adjusted 3-sigma rule; density curve with colored area under this curve between corresponding quantiles on x -axis; volume of the area is equal to the probability of realization of random variable between these quantiles

Hints. Installing shiny

To install shiny on your computer simply run `install.packages('shiny')` in R. If you're installing it on a faculty computer, you have to specify library too and also load `jsonlite` package

```

1 # create a directory called "rlib" somewhere in your profile
2 libpath = 'path_to_your_rlib_directory'
3
4 # install shiny package (need only once)
5 install.packages("shiny", lib=libpath)
6
7 # load shiny package (need to run each time you start R)
8 library("shiny", lib=libpath)
9 library("jsonlite", lib=libpath)

```

Hints. Sample app

Create a folder **normalPlot** and copy following files **ui.R** and **server.R** into it (find these files in study materials). Run it from R console with command `runApp('normalPlot', display.mode="showcase")` from parent directory of **normalPlot** (use function `setwd` to set your working directory if needed). Insert your code from previous exercises into function `renderPlot` in **server.R** to finish the exercise.

Listing 1: **ui.R**

```

1 # if missing, install with command 'install.packages('shiny')'
2 library(shiny)
3
4 # Define UI
5 shinyUI(fluidPage(
6
7   # Application title
8   titlePanel("Normal Distribution"),
9
10  # Sidebar with a slider input for the number of bins
11  sliderInput("sig",
12             "sigma:",
13             min = 0.1,
14             max = 3,
15             value = 1),
16  sliderInput("k",
17             "k:",
18             min = 1,
19             max = 3,
20             value = 1,
21             step = 1),
22
23  plotOutput('normalPlot')
24 ))

```

Listing 2: **server.R**

```

1 library(shiny)
2
3 # Define server logic
4 shinyServer(function(input, output) {
5
6   output$normalPlot <- renderPlot({
7     # insert your code from previous exercises here and modify it so that
8     # it uses values input$sig and input$k from ui.R
9     plot(c(1,2,3), c(1,2,3), main=paste('Sigma:', input$sig))
10  })
11 })

```

Exercise 9 (Interactive Normal Distribution 2). *Build similar app for Exercise 7, but let user choose if he wants to display normal distribution or exponential distribution and let him choose α too, where $\Pr(a \leq X \leq b) = 1 - \alpha$ for normal distribution and $\Pr(X \leq b) = 1 - \alpha$ for exponential distribution. Furthermore, add sliders for μ and σ if user selects normal distribution or slider λ if user selects exponential distribution.*

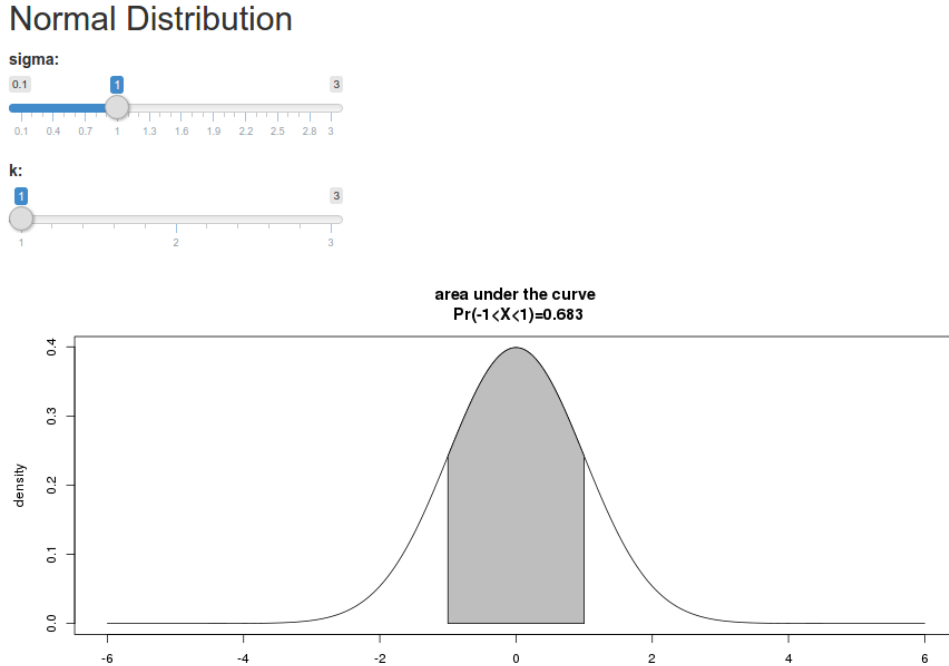


Figure 3: Screenshot from interactive Shiny application with normal distribution.

Your plot should have fixed x-axis limits (`xlim`). Also set reasonable min/max values for sliders.

Hints. Use control widgets `conditionalPanel` and `selectInput` in **ui.R**.

Exercise 10 (Binomial Distribution). Let's assume that number of people preferring treatment A over treatment B follows a binomial distribution with parameters p (probability of success) and N (number of independent trials) denoted $\text{Bin}(N, p)$, where $N = 20, p = 0.5$, i.e. people prefer both treatments equally. (a) What is the probability that 16 and more patients will prefer treatment A over treatment B? (b) What is the probability that 16 and more or 4 or less patients will prefer treatment A over treatment B?

Solution without \mathbb{R} code:

$$(a) \Pr(X \geq 16) = 1 - \Pr(X < 16) = 1 - \Pr(X \leq 15) = 1 - \sum_{i: x_i \leq 15} \Pr(X = x_i) = 1 - \sum_{i: x_i \leq 15} \binom{N}{x_i} p^{x_i} (1-p)^{N-x_i} = 1 - \sum_{i: x_i \leq 15} \binom{20}{x_i} 0.5^{x_i} (1-0.5)^{20-x_i} = 0.0059.$$

$$(b) \Pr(X \leq 4, X \geq 16) = 1 - \sum_{i: x_i \leq 15} \Pr(X = x_i) + \sum_{i: x_i \leq 4} \Pr(X = x_i) = 0.012. \text{ This probability is twice the previous one since } \text{Bin}(N, 0.5) \text{ is symmetric around } 0.5.$$

Hints. `pbinom(x, size=..., prob=...)` gives you probability $\Pr(X \leq x)$. How do you get probability $\Pr(X \geq x)$ using this function?

Exercise 11 (Binomial Distribution). Let's assume that $\Pr(\text{swirl}) = 0.533 = p_1$ is the probability of having dermatological pattern swirl on right thumb of male population and $\Pr(\text{other}) = 0.467 = p_2$ is the probability of other patterns on right thumb of the same population. Random variable X represents number of swirls and Y number of other patterns, where $X \sim \text{Bin}(N, p_1)$ and $Y \sim \text{Bin}(N, p_2)$. Calculate

1. $\Pr(X \leq 120)$ if $N = 300$
2. $\Pr(Y \leq 120)$ if $N = 300$

Compare the results with your intuition (e.g. what do you think is the probability of getting less than 120 heads if you flip a coin 300 times?)

Exercise 12 (Moments of Bernoulli distribution). Random variable X has Bernoulli distribution ($X \sim \text{Ber}(p)$) if $\Pr(X = 1) = 1 - \Pr(X = 0) = 1 - q = p$. Derive its expected value $\text{EX} = \sum_k k \Pr(X = k)$ and variance $\text{Var}X = \text{E}[(X - \text{EX})^2] = \text{E}[X^2] - (\text{EX})^2$

Exercise 13 (Moments of Binomial distribution). Random variable X has Binomial distribution ($X \sim \text{Bin}(N, p)$) if

$$\Pr(X = k) = \binom{N}{k} p^k (1 - p)^{N-k} = \frac{N!}{k!(N-k)!} p^k (1 - p)^{N-k}$$

for $k \in \{0, \dots, N\}$. Derive its expected value $\text{EX} = \sum_k k \Pr(X = k)$ and variance $\text{Var}X = \text{E}[(X - \text{EX})^2] = \text{E}[X^2] - (\text{EX})^2$

Hints. The hard way involves [Binomial theorem](#), for the easy way you need to realize relationship between Binomial and Bernoulli random variable and use results from previous exercise.

Exercise 14 (Normal approximation of binomial distribution).⁸

Let $\Pr(\text{man}) = 0.515$ be a proportion of men in population and $\Pr(\text{women}) = 0.485$ proportion of women. Let X represent number of men and Y number of women. Under the assumption of model $\text{Bin}(N, p)$ calculate

1. $\Pr(X \leq 3)$ if $N = 5$
2. $\Pr(X \leq 5)$ if $N = 10$
3. $\Pr(X \leq 25)$ if $N = 50$.

Compare these probabilities with those approximated by normal distribution $N(Np, Npq)$. Plot density function of normal distribution and superimpose it with probability distribution of binomial distribution. Plot cumulative distribution function of normal distribution and superimpose it with cumulative distribution function for binomial distribution. See figure 4.

Hints. `par(mfrow=c(2, 3)), plot(..., type='h')`

`plot` function creates a new graph, to add points or lines to an existing graph, use commands `points` or `lines`. These functions take the same arguments as `plot`.

⁸ Approximation means "similar but not exactly equal", i.e. we approximate some distribution with a different one (that has certain advantages over the approximated one) or we approximate data with some distribution (that describes data with help of easily interpretable parameters)

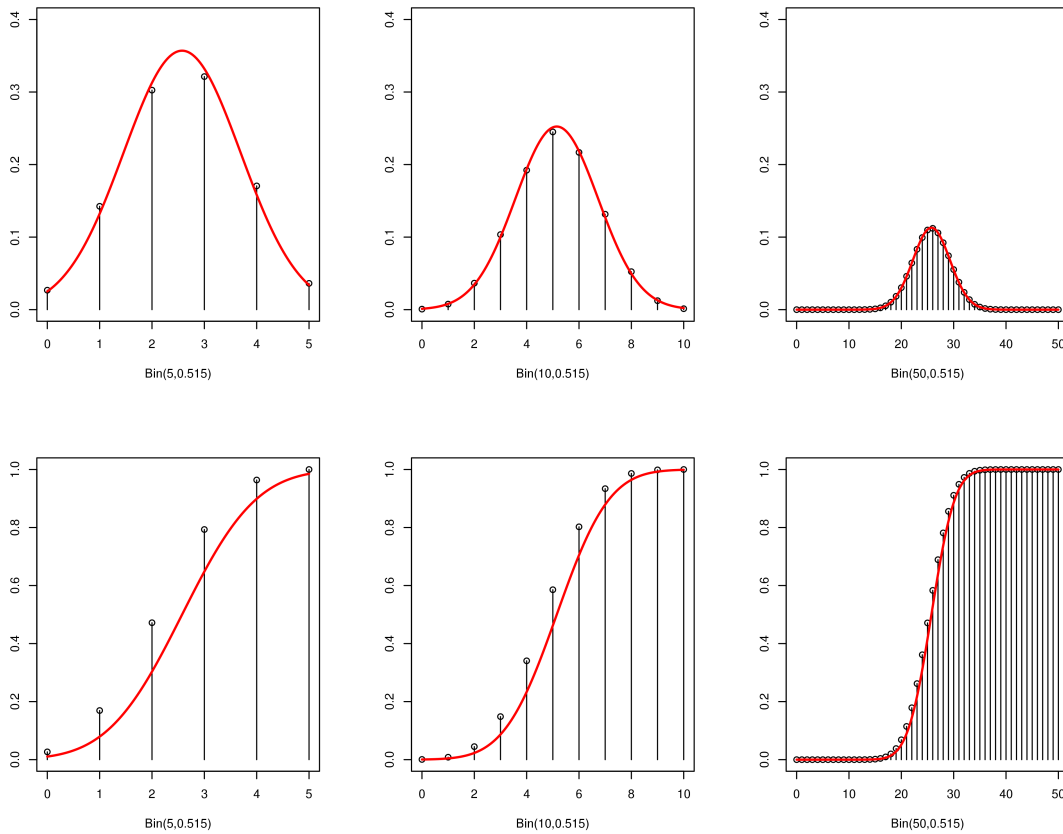


Figure 4: Approximating binomial distribution with normal distribution for $p = 0.515$, $N = 5, 10$ and 50 ; scatter plot superimposed with density (first row) and cumulative distribution function (second row).

Exercise 15 (Levelplots / Contingency tables).

Analyze **data/reputation.data** for potential security risks. Plot heatmap of Risk vs Reliability for all IPs (see figure 5).

Then plot the same risk-reliability heatmap, but grouped by SimpleType. SimpleType has the same values as Type, but all values with multiple types (multiple types are separated by ;) are replaced by "Multiples" value (see figure 6).

Try removing dominating type "Scanning Host" to get more insights into other types.

Hints. Simple heatmap

1. Load data into variable `av` and check their correctness using `head(av)` and `str(av)`.
2. Create a frequency table for heatmap

```
1 | rr.df = data.frame(table(av$Risk, av$Reliability))
2 | colnames(rr.df) <- c("Risk", "Reliability", "Freq")
```

3. Plot level plot using following command (experiment with other options)

```

1 library(lattice)
2 levelplot(Freq~Risk*Reliability,
3           data=rr.df, main="Risk~Reliability", ylab="Reliability", xlab = "Risk",
4           shrink = c(0.5, 1),
5           col.regions = colorRampPalette(c("#F5F5F5", "#01665E"))(20))

```

Hints. Heatmap grouped by type

1. Create a new column `SimpleType` that is the same as `Type` and convert it from factor variable to character with `as.character`. Then replace all values containing ; with value "Multiples" (use `grep` to select these rows).
2. Create frequency table like before and change the first argument of `levelplot` to `levelplot(Freq ~Reliability*Risk|Simpletype, ...)`.

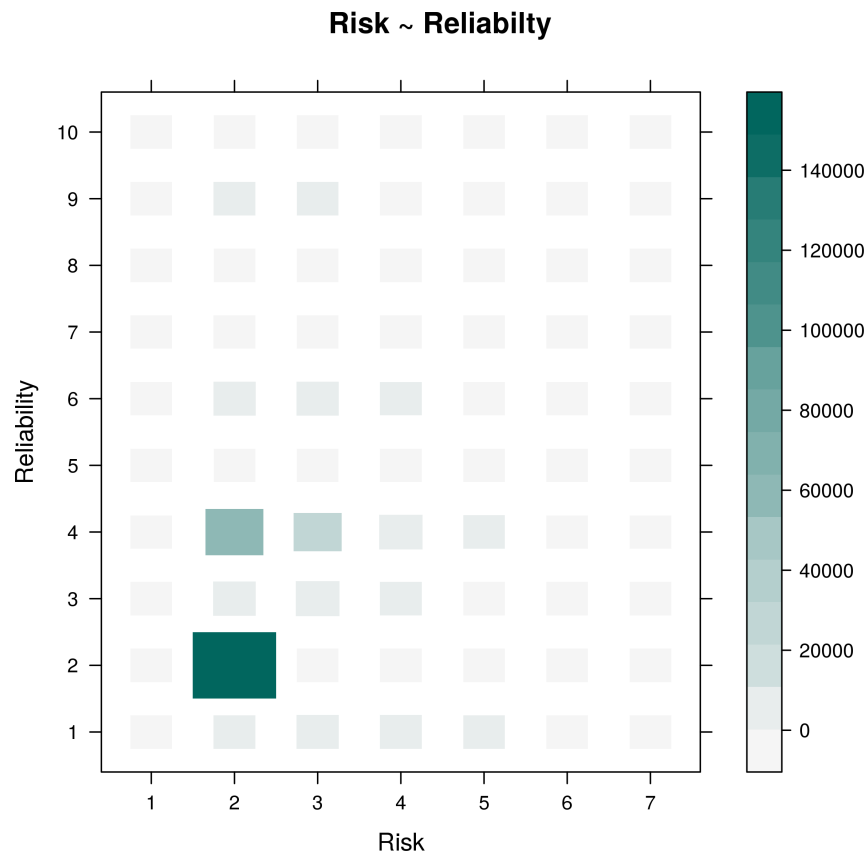


Figure 5: Heatmap for reputation data.

Exercise 16 (Bivariate normal distribution). *Random vector $(X, Y)^T$ has bivariate normal distribution*

$$N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \text{ where } \boldsymbol{\mu} = (\mu_1, \mu_2)^T \text{ and } \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix},$$

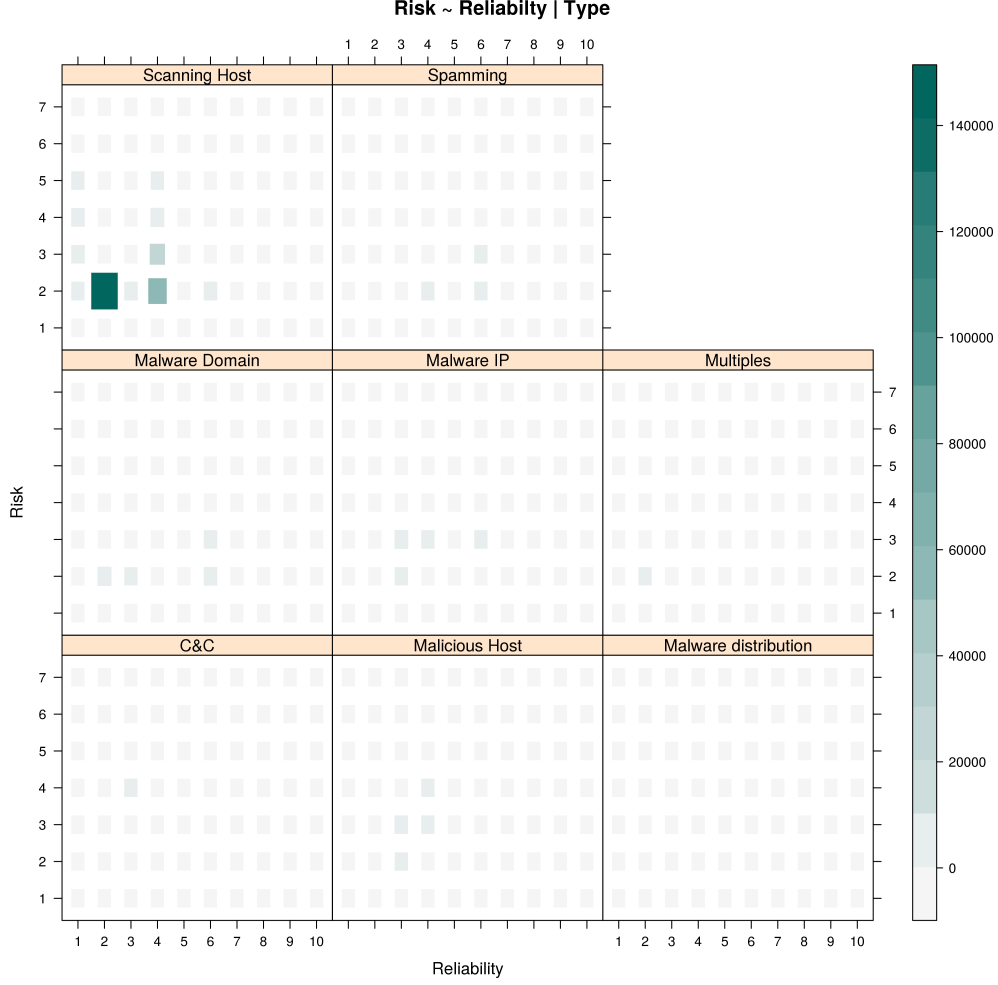


Figure 6: Heatmap for reputation data grouped by Type with all multiple types merged into one value "Multiples".

with density

$$f(x, y) = \frac{1}{2\pi\sqrt{\sigma_1^2\sigma_2^2(1-\rho^2)}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left\{ \frac{(x-\mu_1)^2}{\sigma_1^2} - 2\rho\frac{(x-\mu_1)(y-\mu_2)}{\sigma_1\sigma_2} + \frac{(y-\mu_2)^2}{\sigma_2^2} \right\} \right\},$$

where $(x, y)^T \in \mathbb{R}^2$, $\mu_j \in \mathbb{R}^1$, $\sigma_j^2 > 0$, $j = 1, 2$, $\rho \in \langle -1, 1 \rangle$ are parameters, then $\theta = (\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$. Expression in the exponent can be also written as

$$-\frac{1}{2} \begin{pmatrix} x - \mu_1 \\ y - \mu_2 \end{pmatrix}^T \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}^{-1} \begin{pmatrix} x - \mu_1 \\ y - \mu_2 \end{pmatrix},$$

marginal distributions⁹ are $X \sim N(\mu_1, \sigma_1^2)$ and $Y \sim N(\mu_2, \sigma_2^2)$, ρ is a correlation coefficient¹⁰.

⁹Marginal distribution is a distribution of marginal random variable. Marginal distribution of multivariate normal distribution is again normal, which is a very useful property.

¹⁰From this example it is clear that to sufficiently describe bivariate normal distribution we need 5 pa-

Exercise 17 (Bivariate normal distribution). (1) Create a function `bnorm(x,y,mu1,mu2,sigma1,sigma2,corr)` returning density of a bivariate normal distribution. Use function properties to test that it works correctly.

(2) Plot density of bivariate normal distribution $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ using function `image()` and superimpose it with contour graph of the same distribution using function `contour()`.

(3) Plot density of bivariate normal distribution $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ using function `persp()`. Cut density into 12 intervals, where values in these intervals correspond to colors `terrain.colors(12)`.

Use following parameters

(a) $\mu_1 = 0, \mu_2 = 0, \sigma_1 = 1, \sigma_2 = 1, \rho = 0$;

(b) $\mu_1 = 0, \mu_2 = 0, \sigma_1 = 1, \sigma_2 = 1, \rho = 0.5$;

(c) $\mu_1 = 0, \mu_2 = 0, \sigma_1 = 1, \sigma_2 = 1.2, \rho = 0.5$.

See figure 5 for correct solution.

Hints.

1. create a function `bnorm(x,y,mu1,mu2,sigma1,sigma2,corr)` returning density of bivariate normal distribution
2. to test your implementation, use properties such as symmetry, the fact that function value goes to zero as you go further away from mean, location of the maximum, etc.
3. create vectors x and y with values from -3 to 3 and length n and their cartesian product $(x_i, y_j), i = 1, \dots, n, j = 1, \dots, n$ (represented as either $n^2 \times 2$ matrix or two vectors of length n^2)
4. apply your function `bnorm` on the cartesian and reshape it to a matrix Z with dimensions $n \times n$
5. use x, y and Z in `image` function to plot density

You don't have to use Greek letters in axis labels (use `mu.1` instead of μ_1), but if you want to, look up `expression` function. For coloring of `persp` look at the last example in `persp` documentation and modify it to your needs.

Exercise 18 (Kernel density estimation). Simulate N samples from given distribution and plot them as + symbols on x -axis together with their kernel density estimation. Add theoretical normal distribution (resp. exponential distribution) with parameters estimated from data.

Do it for (a) $N = 50$ (b) $N = 1000$ and following distributions

1. $X \sim N(0, 1)$

2. $X \sim pN(-3, 1) + (1 - p)N(3, 1)$ is a mixture of normals, where $p = 0.3$

3. $X \sim \text{Exp}(\lambda)$, where $\lambda = 3$. Plot theoretical exponential distribution with estimated parameter $\hat{\lambda} = \frac{1}{\bar{X}} = \frac{N}{\sum_{i=1}^N X_i}$.

rameters, i.e. mean and variance for marginal distributions of random variables X and Y and correlation coefficient $\rho = \rho(X, Y)$ describing the strength of linear relationship of X and Y .

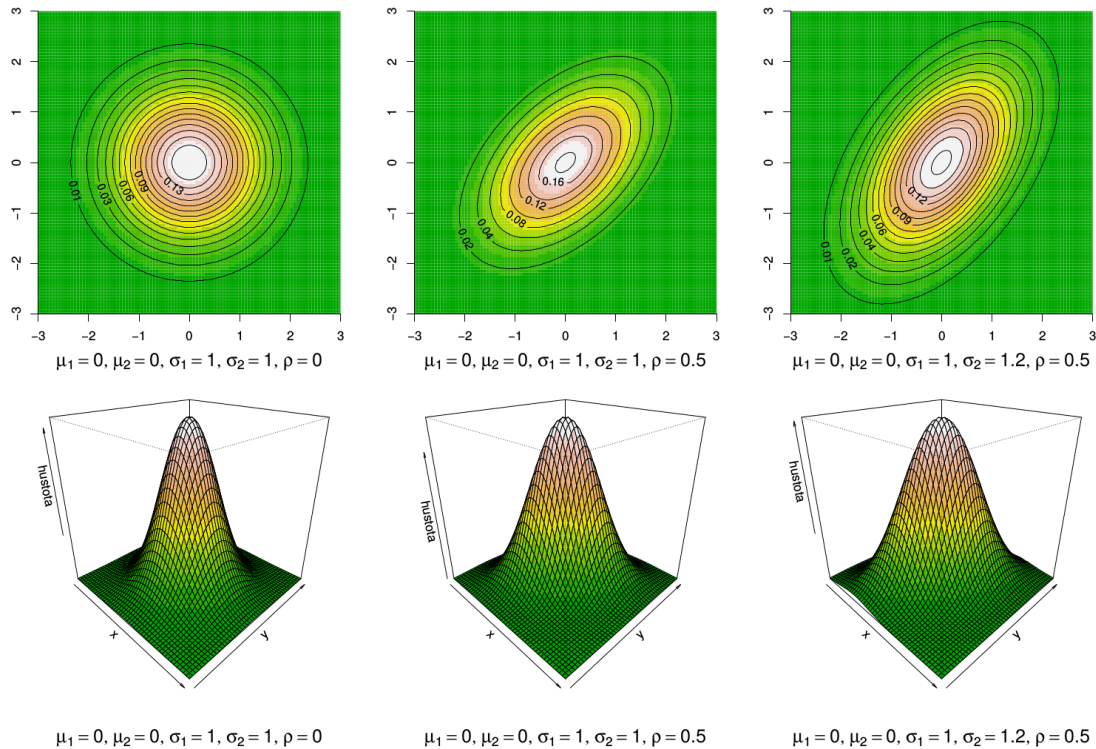


Figure 7: Density of bivariate normal distribution with different parameters (first row – contour graph, second row – perspective 3D graph displayed as surface); the larger absolute value of ρ , the more are contours different from circles (they are changing into ellipses); as the difference between σ_1 and σ_2 is getting larger, we say that difference in variability of X_1 and X_2 is increasing

See figure 8 for correct solution.

Hints. To get the same results as in the figure, call `set.seed(5)` at the beginning of your script.

Useful functions: `density(..., from=, to=, n=)`, `plot(..., lty=2)`, `points(..., pch=3)`

Exercise 19 (Bivariate normal distribution). *Simulating pseudorandom numbers from $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ can be done with \mathbb{R} in several ways:*

- 1) library `library(MASS)` and function `mvrnorm()`;
- 2) library `library(mvtnorm)` and function `rmvnorm()`;
- 3) function `rnorm()` and this algorithm – let $X_1 \sim N(0, 1)$ and $X_2 \sim N(0, 1)$; then $(Y_1, Y_2) \sim N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu} = (\mu_1, \mu_2)^T$ is vector of means and σ_1^2 , σ_2^2 and ρ are parameters of covariance matrix $\boldsymbol{\Sigma}$, where the strength of linear relationship Y_1 and Y_2 is given by magnitude and sign of ρ ; $Y_1 = \sigma_1 X_1 + \mu_1$ a $Y_2 = \sigma_2(\rho X_1 + \sqrt{1 - \rho^2} X_2) + \mu_2$.

Simulate pseudorandom numbers Y_1 and Y_2 from $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ using the first method. Calculate kernel density estimation of $(Y_1, Y_2)^T$ using function `kde2d()`. Plot it using function `image()` and superimpose it with contour graph of bivariate normal density of $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

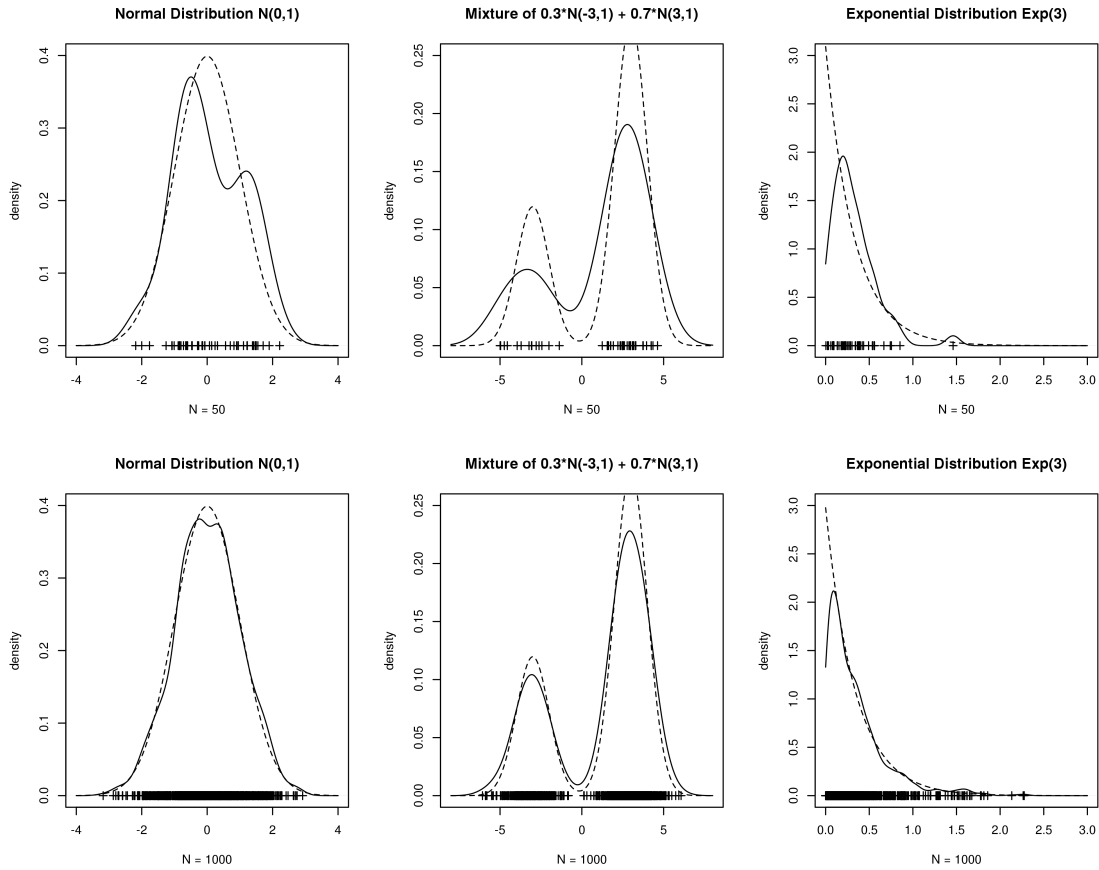


Figure 8: Kernel density estimation of several distributions (first row $n = 50$; second row $n = 1000$)

using function `contour()`. Use following parameter in simulations

- (a) $\mu_1 = 0, \mu_2 = 0, \sigma_1 = 1, \sigma_2 = 1, \rho = 0$; (1) $n = 50$ a (2) $n = 1000$;
 (b) $\mu_1 = 0, \mu_2 = 0, \sigma_1 = 1, \sigma_2 = 1, \rho = 0.5$; (1) $n = 50$ a (2) $n = 1000$;
 (c) $\mu_1 = 0, \mu_2 = 0, \sigma_1 = 1, \sigma_2 = 1.2, \rho = 0.5$; (1) $n = 50$ a (2) $n = 1000$.

See figure 11 for correct solution.

Hints. `kde2d(..., n=100, lims=c(-3, 3, -3, 3))`

Exercise 20 (Mixture of two normal distributions). Simulate pseudorandom numbers (1) from mixture of normal distributions $pN_2(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + (1-p)N_2(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ where $\boldsymbol{\theta}_1 = (\mu_{11}, \mu_{12}, \sigma_{11}^2, \sigma_{12}^2, \rho_1, \mu_{21}, \mu_{22}, \sigma_{21}^2, \sigma_{22}^2, \rho_2)^T$ and (2) from bivariate normal distribution $N_2(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where parameters represent combined vector of means a combined covariance matrix. i.e. $\boldsymbol{\theta}_2 = (\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)^T$. For (1) calculate kernel density estimation of $(X, Y)^T$ using function `kde2d()`.

(a) Plot theoretical density (2) using function `image()` and superimpose it with contour graph of theoretical density (2) using function `contour()`.

(b) Plot theoretical density (1) using function `image()` and superimpose it with contour graph

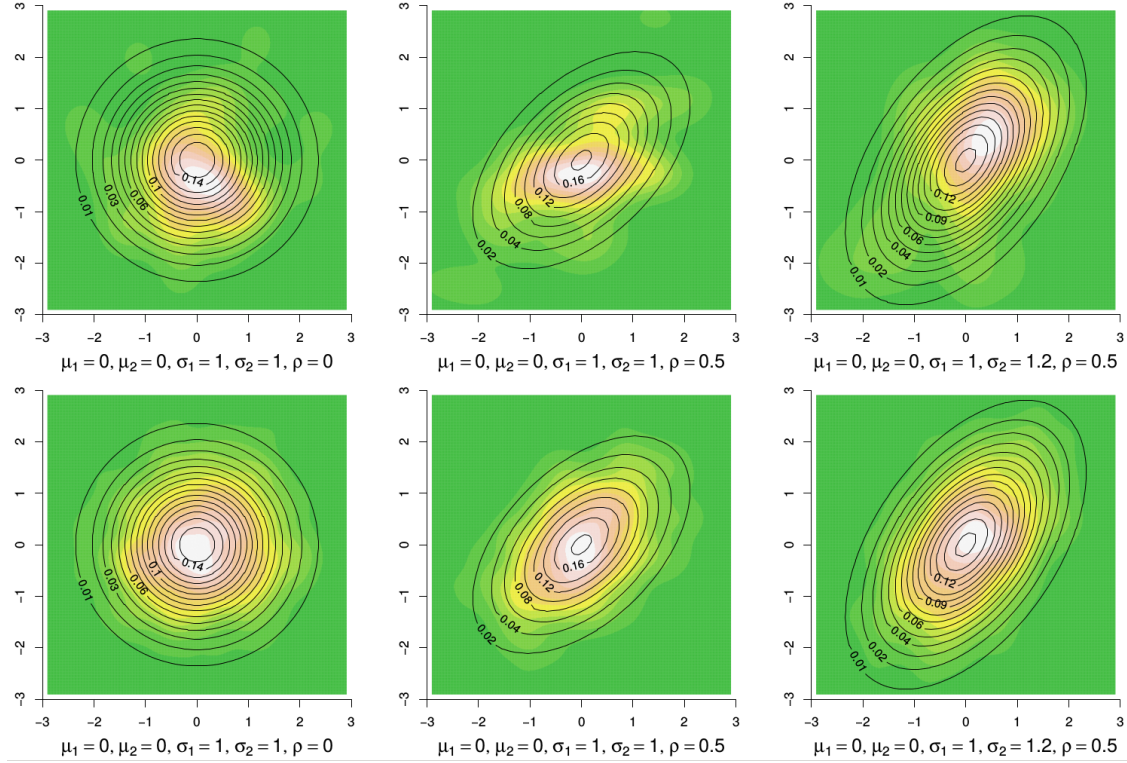


Figure 9: Density of bivariate normal distribution (first row $n = 50$; second row $n = 1000$)

of theoretical density (1) using function `contour()`.

(c) Plot kernel density estimation of (1) using function `image()` and superimpose it with contour graph of theoretical density (1) using function `contour()`.

Cut density into 12 intervals, where values in these intervals correspond to colors `terrain.colors(12)`. For simulation use these parameters $\theta_1 = (-1.2, -1.2, 1, 1, 0, 1, 1, 1, 1, 0)^T$, $n = 100$ and $p = 0.5$.

Use parameters from simulation for case (1), i.e. $\hat{\theta}_1 = \theta_1 = (\mu_{11}, \mu_{12}, \sigma_{11}^2, \sigma_{12}^2, \rho_1, \mu_{21}, \mu_{22}, \sigma_{21}^2, \sigma_{22}^2, \rho_2)^T = (-1.2, -1.2, 1, 1, 0, 1, 1, 1, 1, 0)^T$. For case (2) estimate parameters from simulated data $\hat{\theta}_2 = (\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}_1^2, \hat{\sigma}_2^2, \hat{\rho})^T$.

See figure 8 for correct solution.

Hints. For parameter estimations in case (2) use functions `mean`, `std` and `cor`.

Exercise 21 (Poisson distribution). We have a data (Greenwood and Yule 1920) with number of injuries of factory workers in the following table

n	0	1	2	3	4	≥ 5
m_n	447	132	42	21	3	2

where n is number of injuries and m_n number of workers with n injuries.

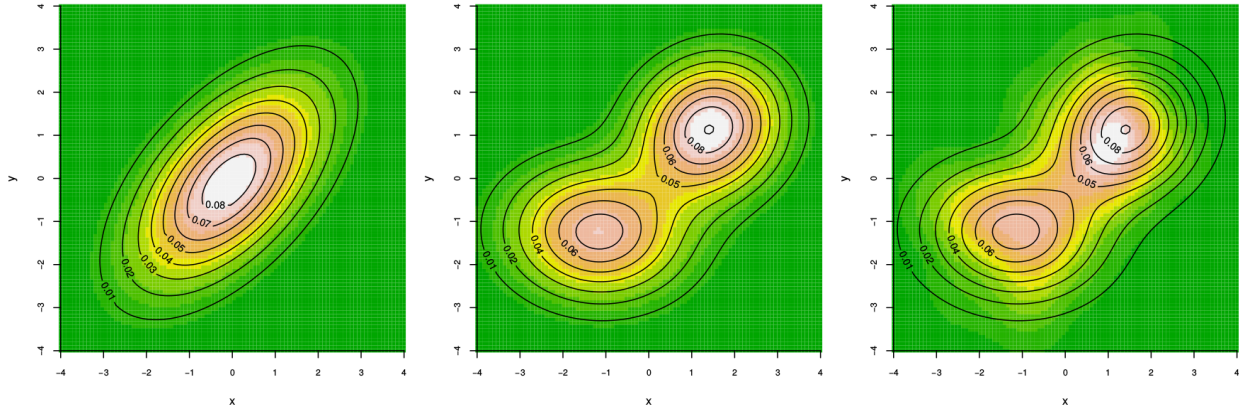


Figure 10: Combined density of bivariate normal distribution (left), mixture density of two bivariate normal distributions (middle) and kernel density estimation superimposed by mixture density of two bivariate normal distributions (right)

Calculate expected number of worker injuries under the assumption that random variable X representing injuries has Poisson distribution with parameter $\lambda = \frac{\sum_n nm_n}{\sum_n m_n}$, i.e. $X \sim \text{Poiss}(\lambda)$. Create a table with m_n and expected m_n and display it (figure 11).

Hints. Create a dataframe with two columns and proper row names, then print it. `data.frame(mn=..., expected_mn=..., row.names=...)`.

	0	1	2	3	4	>=5
mn	447	132	42	21	3	2
expected mn	406	189	44	7	1	0

Figure 11: Observed and expected frequencies of Poisson distribution.

Exercise 22 (Binomial distribution). In a study from 1889 based on medical records in Saxony professor Geissler (1889) recorded distribution of boys in families. The study included $M = 6115$ families with $N = 12$ children with following number of boys (n stands for number of boys and m_n number of families with n boys)

n	0	1	2	3	4	5	6	7	8	9	10	11	12
m_n	3	24	104	286	670	1033	1343	1112	829	478	181	45	7

Calculate expected m_n under the assumption that number of boys X in families follow binomial distribution $\pi = \frac{\sum_{n=0}^N nm_n}{NM}$ and $N = 12$ (i.e. $X \sim \text{Bin}(N, \pi)$).

Compare expected and observed frequencies - do you see any difference? Display m_n and expected m_n in a table and visualize both observed and expected frequencies in one graph (see figure 10 or use your own imagination). Calculate probability that family will have ≥ 4 and ≤ 6 boys from theoretical distribution (i.e. $\text{Pr}(4 \leq X \leq 6)$).

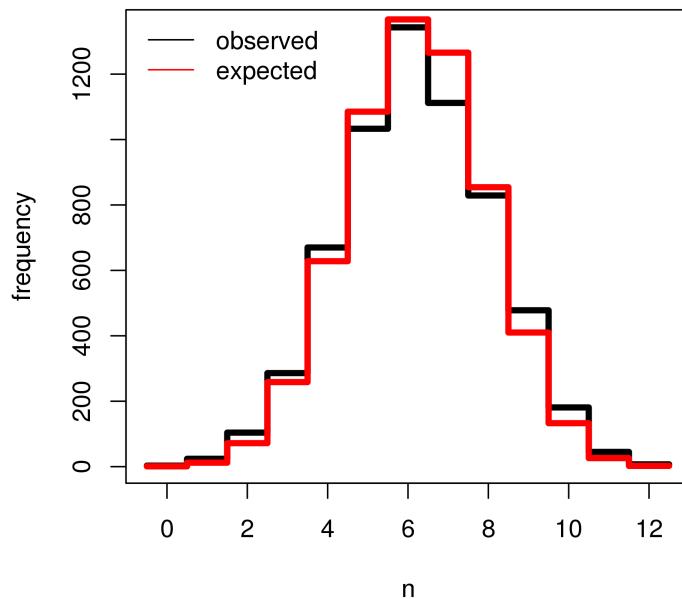


Figure 12: Observed and expected frequencies of Binomial distribution

Hints. Add legend with `legend('topleft', c('observed', 'expected'), col=c('black', 'red'), lty=1, bty='n')`. For stairs graph use argument `type='S'` or `type='s'`

Exercise 23 (Approximating Binomial distribution by Poisson). *If each of 50 million people in Italy drives a car next week (independently), then probability of dying in an accident will be 0.000002, where number of fatalities has binomial distribution $\text{Bin}(50\text{mil}, 0.000002)$. Approximate this Binomial distribution with Poisson that has the same mean. Plot both probabilities for values $\{50, 51, \dots, 150\}$.*

Exercise 24 (Contingency table). *356 people have been polled on their smoking status (Smoke) and their socioeconomic status (SES). For each person it was determined whether or not they are current smokers, former smokers, or have never smoked. Also, for each person their socioeconomic status was determined (low, middle, or high). The data file `smoker.csv` (available from study materials) contains only two columns - Smoke and SES. Load this data into \mathbb{R} and create contingency table with `table()` function.*

These observations (after converting to probability) form joint distribution $\Pr(X = x, Y = y)$, where X is a discrete random variable for smoker and Y for SES. Calculate marginal distributions $\Pr(X = x) = \sum_y \Pr(X = x, Y = y)$, $x \in \{\text{current}, \text{former}, \text{never}\}$ and $\Pr(Y = y)$, $y \in \{\text{high}, \text{low}, \text{middle}\}$. Also calculate and compare $\Pr(X = x|Y = \text{high})$ with $\Pr(X = x|Y = \text{low})$ and $\Pr(Y = y|X = \text{smoker})$.

Create a table with expected frequencies using assumption that X and Y are independent (figure). Visualize both tables (observed and expected) with `mosaicplot` function (figure).

Hints. To calculate marginal probabilities use either special `table` functions `margin.table` and `prop.table` or simply `rowSums` and `colSums`.

	High	Low	Middle
current	51	43	22
former	92	28	21
never	68	22	9

	High	Low	Middle
current	68.75281	30.30337	16.94382
former	83.57022	36.83427	20.59551
never	58.67697	25.86236	14.46067

Figure 13: Observed and expected frequencies for contingency table

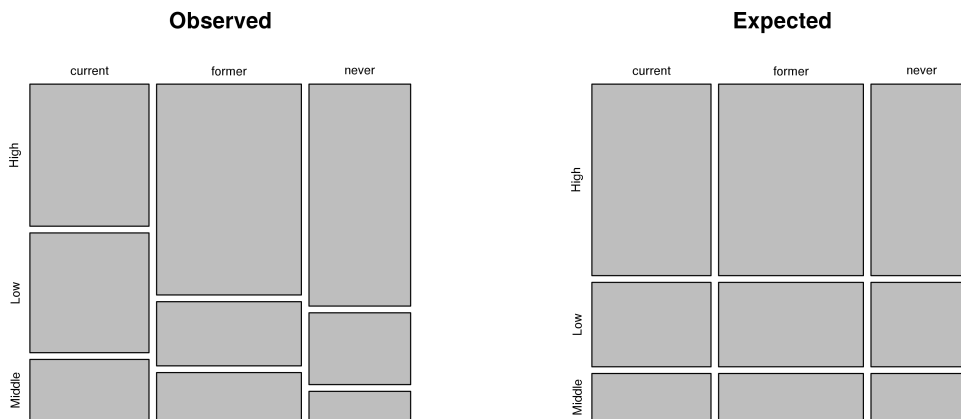


Figure 14: Mosaic plot of contingency table

Exercise 25 (Number of contingency tables). *How many contingency tables of size 4×2 there are for $N = 20$?*

Exercise 26 (Multinomial distribution). *We have following random variables (1) socioeconomic status (high - H , low - Lo), (2) political affiliation (democrat - D , republican - R) and (3) political philosophy (liberal - Li , conservatism - C). Let's denote their interactions like this: X_1 (H - D - Li), X_2 (H - D - C), X_3 (H - R - Li), X_4 (H - R - C), X_5 (Lo - D - Li), X_6 (Lo - D - C), X_7 (Lo - R - Li) a X_8 (Lo - R - C). We have random sample of size $N = 50$. Probabilities p_j are in the following table*

	D - Li	D - C	R - Li	R - C	all
H	0.12	0.12	0.04	0.12	0.4
Lo	0.18	0.18	0.06	0.18	0.6
all	0.30	0.30	0.10	0.30	1.0

Calculate $\text{Var}[X_1]$, $\text{Var}[X_3]$, $\text{Cov}[X_1, X_3]$, $\text{Cor}[X_1, X_3]$ and expected frequencies $Np_j, j = 1, 2, \dots, 8$.

Hints. $X = (X_1, X_2, \dots, X_8) \sim \text{Mult}(N, \mathbf{p})$, where $N = 50$, $\mathbf{p} = (p_1, p_2, \dots, p_8)^T$, we know that $X_j \sim \text{Bin}(N, p_j)$, p_j are given by table above $j = 1, 2, \dots, 8$.

Exercise 27 (Product-multinomial distribution). *Let's have the same probabilities as in previous exercise, but now with two separate random samples - first with size $N_1 = 30$ from group H and the second with size $N_2 = 20$ from group Lo. Denote interactions of random variables $X_{11} = X_{1|1}$ (H-D-Li), $X_{12} = X_{2|1}$ (H-D-C), $X_{13} = X_{3|1}$ (H-R-Li), $X_{14} = X_{4|1}$ (H-R-C), $X_{21} = X_{1|2}$ (Lo-D-Li), $X_{22} = X_{2|2}$ (Lo-D-C), $X_{23} = X_{3|2}$ (Lo-R-Li) and $X_{24} = X_{4|2}$ (Lo-R-C), where $\mathbf{X}_1 = (X_{11}, X_{12}, X_{13}, X_{14})^T$ and $\mathbf{X}_2 = (X_{21}, X_{22}, X_{23}, X_{24})^T$. Then $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$ has product-multinomial distribution with $K = 2$, $N_1 = 30$, $J_1 = 4$, $N_2 = 20$, $J_2 = 4$. Notation $X_{j|k}$, where $j = 1, 2, 3, 4$ and $k = 1, 2$ highlights the fact, that distribution is conditioned by socioeconomic status (high - H, low - Lo), i.e. distribution in table columns is conditioned by its row. Realizations $X_{j|k}$ are denoted as $n_{j|k} = n_{kj}$, similarly probabilities $X_{j|k} = X_{kj}$ as $p_{j|k} = p_{kj}$. Calculate conditional probabilities $p_{j|k}$, expected frequencies $N_k p_{kj}$, $\text{Var}[X_{13}]$, $\text{Cov}[X_{21}, X_{23}]$ and $\text{Cor}[X_{11}, X_{23}]$.*

Exercise 28 (Binomial distribution, simulation study). *Generate $M = 1000$ pseudo-random numbers from distribution $\text{Bin}(5, 0.5)$. Create a table with observed and theoretical probabilities (for $n = 0, 1, \dots, 5$).*

r	0	1	2	3	4	5
observed probabilities	0.031	0.158	0.302	0.324	0.161	0.024
theoretical probabilities	0.031	0.156	0.312	0.312	0.156	0.031

Superimpose histogram of simulated numbers with theoretical density function. See figure 15.

Hints. For histogram plot use `breaks` and `probability` arguments.

Exercise 29 (Normal distribution, simulation study). *If random variable $X \sim N(150, 6.25)$, what distribution has its arithmetic average \bar{X}_n ? Verify your result using simulations for $n = 30$. Make 5000 simulations (each simulation involves generating 30 realizations of X) and for each simulation calculate arithmetic average \bar{x}_m , $m = 1, 2, \dots, M$, where $M = 5000$. Superimpose their histogram (normalized) with theoretical density function for \bar{X}_n you derived. Calculate $\Pr(\bar{X}_n > 151)$ from simulated data and compare this result with theoretical (expected) probability. See figure 16.*

Exercise 30 (Normal distribution, simulation study II). *Let $X \sim N(\mu_1, \sigma_1^2)$ a $Y \sim N(\mu_2, \sigma_2^2)$. Then $\bar{X}_{n_1} - \bar{Y}_{n_2} \sim N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$. Use $\mu_1 = 100, \sigma_1 = 10, \mu_2 = 50, \sigma_2 = 9$ and (a) $n_1 = 4, n_2 = 5$, (b) $n_1 = 100, n_2 = 81$. Make $M = 1000$ simulations and for each calculate difference $\bar{x}_m - \bar{y}_m$, $m = 1, 2, \dots, M$. Superimpose histogram (normalized) of the differences with theoretical density function of difference $\bar{X}_{n_1} - \bar{Y}_{n_2}$. For both cases (a) and (b) calculate $\Pr(\bar{X}_{n_1} - \bar{Y}_{n_2}) < 52$ from simulated data and compare this result with theoretical probability.*

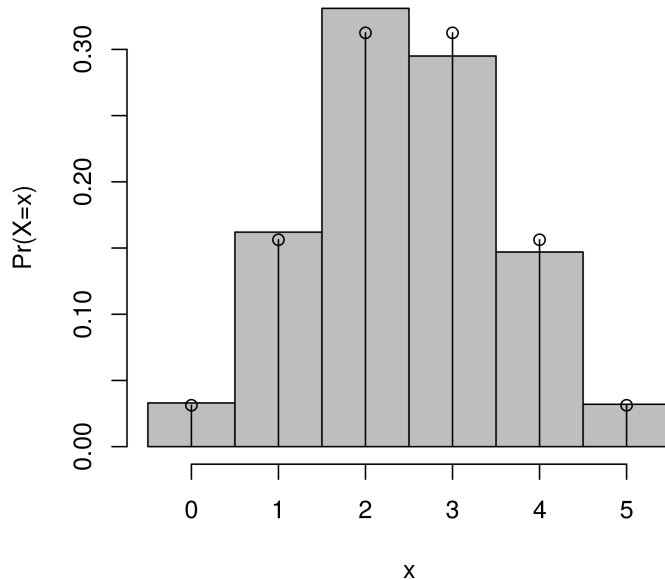


Figure 15: Histogram superimposed with theoretical distribution of $Bin(5, 0.5)$

4 Likelihood

Exercise 31 (Maximum-likelihood derivation). Let $X \sim Poiss(\lambda)$ and x_1, \dots, x_n are its i.i.d. realizations. Analytically derive its log-likelihood function $l(\lambda|\mathbf{x})$. Recall that

$$l(\lambda|\mathbf{x}) = \ln \mathcal{L}(\lambda|\mathbf{x}) = \ln \prod_{i=1}^n f(x_i) = \sum_{i=1}^n \ln f(x_i)$$

Exercise 32 (Maximum-likelihood, Poisson distribution). Let $X \sim Poiss(\lambda)$. Simulate realizations x_1, \dots, x_n . Use $n = 100$ and $\lambda = 4$. Plot log-likelihood function of Poisson distribution

$$l(\lambda|\mathbf{x}) = \sum_{i=1}^n x_i \ln \lambda - n\lambda.$$

for $\lambda \in [2, 6]$. Find maximum likelihood estimate of λ and mark it in graph with dashed line. See figure 17.

Hints. Start with writing a function `loglikelihood = function(x, lambda)` that calculates log-likelihood for **single** λ . Then in a for loop or using `sapply` function calculate log-likelihood for $\lambda \in [2, 6]$.

To find maximum of a function, either use formula for MLE from the lecture or find it empirically from vector of log-likelihoods using function `which.max(lambdas)`. To plot a vertical line use `abline(v=..., lty='dashed')`.

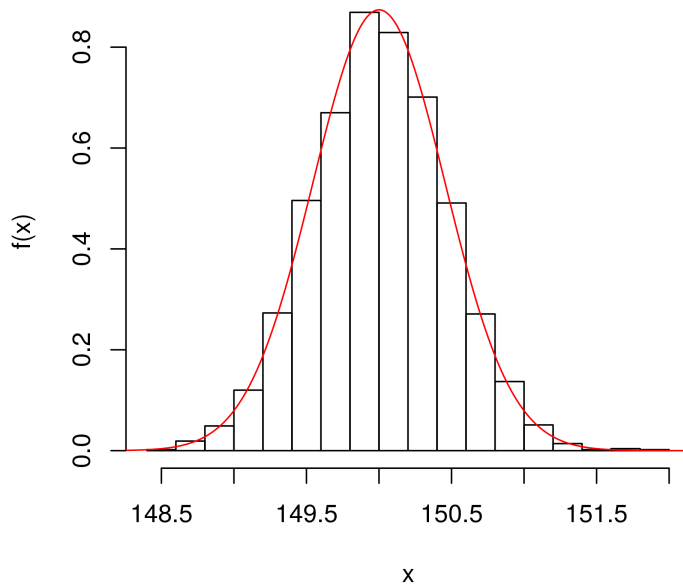


Figure 16: Histogram superimposed with theoretical distribution of $N(150, 6.25/n)$

Exercise 33 (Maximum-likelihood, Poisson distribution). Use data from exercise 21 for log-likelihood function from the previous exercise and plot it. Note that the log-likelihood function will have to be slightly modified to work with **counts** of realizations from exercise 21. See figure 17.

Exercise 34 (Maximum-likelihood, Binomial distribution). Let $X \sim \text{Bin}(N, p)$ and x its realization¹¹. Formulate its likelihood function $L(p|x)$ and log-likelihood function $l(p|x)$. Use log-likelihood function to derive maximum-likelihood estimate $\hat{p} = \underset{p}{\operatorname{argmax}} l(p|x)$ of parameter p and Fisher information $\mathcal{I}(\hat{p}) = -\frac{\partial^2}{\partial p^2} l(p|x)|_{p=\hat{p}}$.

Exercise 35 (Quadratic approximation of log-likelihood function). Let $X \sim \text{Bin}(N, p)$ and x its realization. Plot scaled log-likelihood function of binomial distribution with p on x -axis and scaled log-likelihood $l^*(p|\mathbf{x}) = l(p|\mathbf{x}) - \max(l(p|\mathbf{x}))$ on y -axis. Compare $l^*(p|\mathbf{x})$ with its quadratic approximation calculated using Taylor approximation $\ln\left(\frac{L(p|\mathbf{x})}{L(\hat{p}|\mathbf{x})}\right) \approx -\frac{1}{2}\mathcal{I}(\hat{p})(p - \hat{p})^2$. Use (a) $x = 8, N = 10$, (b) $x = 80, N = 100$ and (c) $n = 800, N = 1000$ with reasonable ranges for both axes. See figure 18.

Hints. Use `lines(..., lty="dashed")` for the dashed line.

¹¹We work with just **one** realization, so $L(p|x) = \prod_{i=1}^1 f(x|p)$. You would get the same results if you worked with N realizations with Bernoulli distribution $\text{Ber}(p)$.

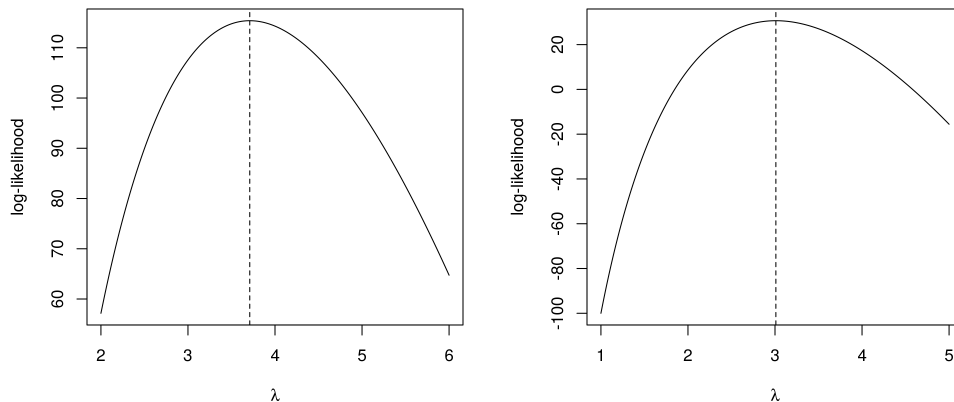


Figure 17: Left: Log-likelihood function for simulated data from $Poiss(\lambda = 4)$. Right: Log-likelihood function for number of injuries of factory workers.

5 Hypothesis Testing

Exercise 36 (Hypothesis testing). Let $X \sim Bin(N, p)$ be a random variable with $N = 1000$ and $p = 0.5$. Make 100000 simulations x_1, \dots, x_{100000} and for each simulation calculate maximum likelihood estimate of p , i.e. $\hat{p} = \frac{x}{N}$. Plot histogram of \hat{p} and add vertical line at value 0.534. Calculate percentage of \hat{p} that are higher than 0.534.

Exercise 37 (Hypothesis testing). Out of 1000 newborn children, 534 of them were boys and 466 girls. Test null hypothesis that probabilities of having a boy or a girl are equal. Calculate p -value, confidence interval, test statistic and critical region. Use Wald statistic and Likelihood statistic.

Hints. Use `uniroot` to solve inequality in confidence interval for Likelihood statistic.

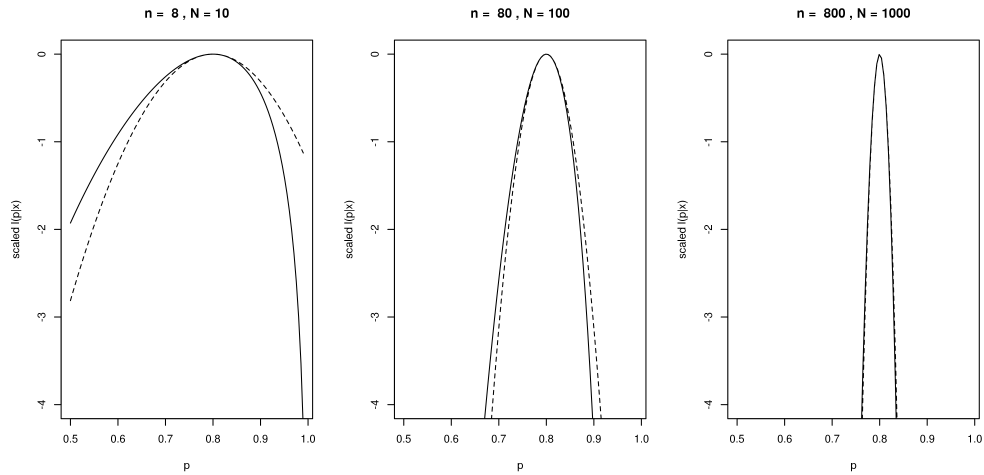
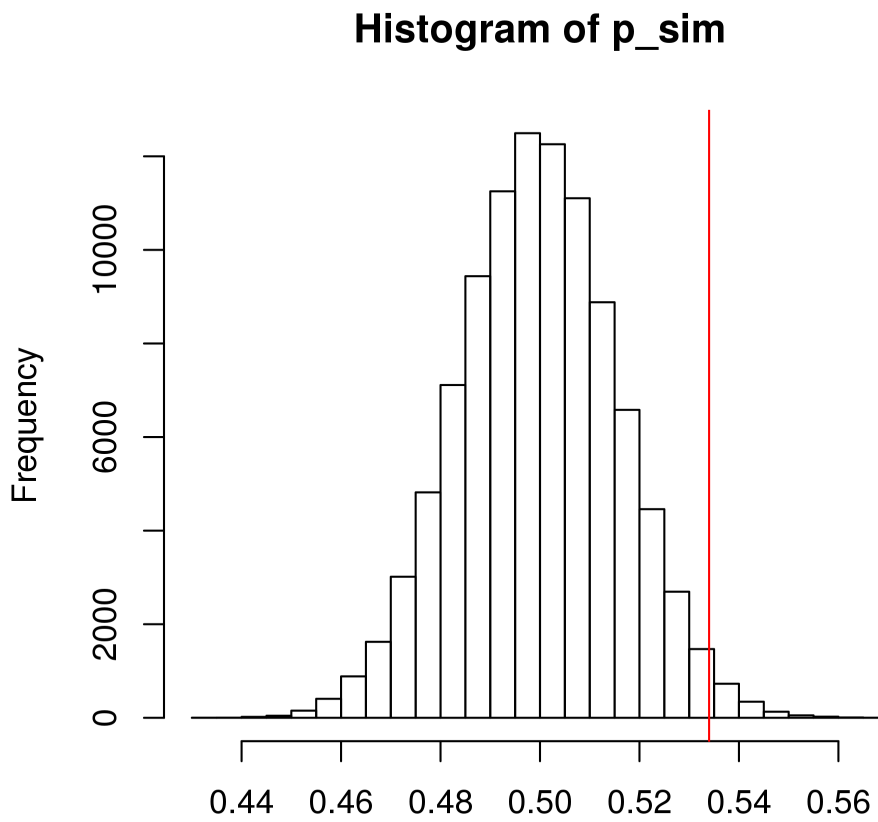


Figure 18: Compare scaled log-likelihood function (full line) with its quadratic approximation (dashed line).



Pct of time that $p_{hat} > 0.534$: 1.52%

Figure 19: Histogram of simulations of \hat{p} with line at 0.534.